## Stochastische Populationsmodelle

Notizen zu einer Vorlesung an der Johannes-Gutenberg-Universität Mainz, Sommer 2024

Matthias Birkner

Vorläufige Version, 18. Juli 2024

Kommentare, Korrekturvorschläge, Hinweise auf (Tipp-)fehler gerne per Email an birkner@mathematik.uni-mainz.de senden

# Inhaltsverzeichnis

I	Wright-Fisher-Modell und Kingman-Koaleszent				
	I.I	Wright-Fisher-Modell: Fundamentalmodell für genetische Drift	2		
	1.2	Genealogien und Kingmans Koaleszent	8		
	1.3	Moran-Modell	13		
	I.4	Dualität	14		
	1.5	Beispiel: Die Beobachtungen von Dorit et al, 1995	18		
2	Mutationen und der markierte Koaleszent				
	<b>2.</b> I	Infinitely-many-alleles-Modell (IMA)	24		
	2.2	Infinitely-many-sites-Modell (IMS)	35		
3	Selektion				
	3.1	Vorbemerkung: Modelle für den diploiden Fall	57		
	3.2	Vorbemerkung: Deterministische Dynamik	59		
	3.3	(2-Typ) Moran-Modell mit (gerichteter) Selektion	63		
	3.4	(2-Typ) Moran-Modell mit (gerichteter) Selektion und Mutation	74		
4	Räu	mliche Struktur	82		
A	Anhang				
	А.1	Ein Exkurs zum Poissonprozess und zu zeitkontinuierlichen Markovketten	91		
	A.2	(Weitere) Eigenschaften des Kingman-Koaleszenten	95		
	A.3	Die Verteilung der Summe unabhängiger, exponentialverteilter Zufallsvariablen	97		
	A.4	Erwartete Fixationszeit im Wright-Fisher-Modell: exakte Rechnung	98		
	A.5	Ein Steilkurs über Martingale in diskreter Zeit	103		

# Kapitel 1

# Wright-Fisher-Modell und Kingman-Koaleszent

## 1.1 Wright-Fisher-Modell: Fundamentalmodell für genetische Drift

Das klassische *Wright-Fisher-Modell*<sup>1,2</sup> ist ein grundlegendes Modell der mathematischen Populationsgenetik zur Beschreibung des Phänomens der sogenannten Gendrift: genetische Typenhäufigkeiten in Populationen verändern sich im Lauf der Zeit aufgrund von Zufälligkeiten im Reproduktionserfolg, auch wenn keine systematischen Unterschiede der Typen oder äußere Einflüsse einwirken. Es ist ein (sehr) idealisiertes Populationsmodell, wir treffen folgende Annahmen:

- feste Populationsgröße:  $N \in \mathbb{N}$  Individuen in jeder Generation
- die Population entwickelt sich in diskreten, nicht-überlappenden Generationen
- jedes Individuum hat nur ein "Elter"<sup>3</sup>
- es gibt Zufälligkeit bezüglich der Anzahl der Nachkommen, dabei aber keine systematischen Vorteile einzelner Individuen: die (gemeinsame) Kinderzahlverteilung ist "symmetrisch"
- es gibt verschiedene genetische Typen, die von Elter zu Kind vererbt werden

Der Einfachheit halber nehmen wir hier (zunächst) an, dass es nur zwei verschiedene Typen, bezeichnet *a* und *A* gibt (zudem: keine "Kopierfehler", sog. Mutationen<sup>4</sup>, bei der Vererbung).

<sup>&</sup>lt;sup>1</sup>Sewall Green Wright, 1889 – 1988, amerikanischer Genetiker

<sup>&</sup>lt;sup>2</sup>Ronald Aylmer Fisher, 1890 – 1962, britischer Statistiker und Genetiker.

<sup>&</sup>lt;sup>3</sup>Im Jargon der Genetik sind die Individuen "haploid" – wörtlich angemessen z.B. für Bakterien, mitochondriale genetische Typen, Y-Chromosom. Viele Spezies sind "diploid", besitzen also zwei Kopien jedes Chromosoms [ggfs. mit Ausnahme der Geschlechtschromosomen], manche Pflanzen sind "polyploid". Asymptotisch, mit Ersetzung  $N \rightsquigarrow 2N$ , ist das Modell aber auch für Gene in diploiden Populationen passend.

<sup>&</sup>lt;sup>4</sup>Eine sehr schöne Einführung in die Grundlagen der Genetik findet sich beispielsweise auf der Webseite DNA from the beginning https://www.dnaftb.org/ des Cold Spring Harbor Laboratory.

Nehmen wir an, jedes Individuum einer gegebenen Generation hat (unabhängig) eine zufällige, Poisson-verteilte Anzahl Nachkommen mit Mittelwert 1, sei  $M_i$  = Anzahl Nachkommen von Individuum  $i, 1 \le i \le N$ ; angesichts der konstanten Gesamtpopulationsgröße müssen wir auf  $\{M_1 + M_2 + \dots + M_N = N\}$  bedingen. Für  $m_1, \dots, m_N \in \mathbb{N}_0$  (mit  $m_1 + \dots + m_N = N$ ) ist dann (beachte: die Faltungeigenschaft der Poissonverteilung liefert  $M_1 + M_2 + \dots + M_N = d$  Pois(N))

$$\mathbb{P}\Big(M_1 = m_1, \dots, M_N = m_N \Big| \sum_{i=1}^N M_i = N\Big) = \frac{\prod_{i=1}^N e^{-1} \frac{1^{m_i}}{m_i!}}{e^{-N} \frac{N^N}{N!}} = \frac{N!}{m_1! m_2! \cdots m_N!} \Big(\frac{1}{N}\Big)^N$$

Somit: die gemeinsame Verteilung der Nachkommenszahlen der Individuen einer gegebenen Generation ist Multinom $(N, \frac{1}{N}, \dots, \frac{1}{N})$ -verteilt. (Wir nehmen zudem Unabhängigkeit über die verschiedenen Generationen an.)

Alternative Interpretation: jedem Individuum in Generation r wird unabhängig ein uniform aus allen Individuen der Generation r - 1 gewähltes Individuum als Elter zugeordnet.

Wir können daraus die Dynamik des Typenzählprozesses  $(X_r^{(N)})_{r \ge r_0}$  ablesen: Sei

 $X_r^{(N)}$  = Anzahl Typ *A*-Individuen in Generation *r* 

(demnach:  $N - X_r^{(N)}$  Typ *a*-Individuen in Generation *r*).

Für  $x, y \in \{0, 1, \dots, N\}$  gilt somit

$$\mathbb{P}(X_{r+1}^{(N)} = y | X_r^{(N)} = x) = {\binom{N}{y}} (\frac{x}{N})^y (1 - \frac{x}{N})^{N-y},$$

d.h. gegeben  $X_r^{(N)} = x$  ist  $X_{r+1}^{(N)} \sim Bin(N, x/N)$ . Insbesondere gilt

$$\mathbb{E}\left[X_{r+1}^{(N)} \mid X_r^{(N)} = x\right] = x, \quad \operatorname{Var}\left[X_{r+1}^{(N)} \mid X_r^{(N)} = x\right] = N\frac{x}{N}\left(1 - \frac{x}{N}\right), \qquad x = 0, 1, \dots, N \quad (\mathbf{I}.\mathbf{I})$$

Abkürzend werden wir die folgende Notation verwenden:  $\mathbb{P}_x$  für das W'maß in der Situation, dass  $X_0^{(N)} = x$ .

Sei  $T_{\text{fix}} := \inf\{r \in \mathbb{N}_0 : X_r^{(N)} = 0 \text{ oder } X_r^{(N)} = N\}$  Zeitpunkt, zu dem einer der beiden Typen verloren geht (angesichts  $\min_{1 \le x \le N-1} \mathbb{P}_x(X_{r+1}^{(N)} \in \{0, N\}) > 0$  ist offenbar  $\mathbb{P}_x(T_{\text{fix}} < \infty) = 1$  für jedes x = 0, 1, ..., N).

#### Wie wahrscheinlich ist es, dass sich Typ A durchsetzt?

$$h(x) := \mathbb{P}_x(X_{T_{\text{fix}}}^{(N)} = N), \quad x = 0, 1, \dots, N$$

ist die eindeutige Lösung des Gleichungssystems

$$h(x) = \sum_{y=0}^{N} {N \choose y} \left(\frac{x}{N}\right)^{y} \left(1 - \frac{x}{N}\right)^{N-y} h(y), \quad x \in \{1, 2, \dots, N-1\},$$
  
$$h(0) = 0, \qquad h(N) = 1$$

(Dies ist eine Instanz der allgemeinen Beziehung zwischen Auftreffwahrscheinlichkeiten von Markovketten und diskreten Dirichlet-Problemen, siehe beispielsweise [Geo15, Kap. 6.2], [Bir24, Kap. 7.1].)

Da der Erwartungswert einer Bin(N, x/N)-verteilten Zufallsvariable gerade  $N\frac{x}{N} = x$  ist, sieht man, dass der Ansatz h(x) = x/N die eindeutige Lösung liefert, d.h.

$$\mathbb{P}_x(X_{T_{\mathrm{fix}}}^{(N)} = N) = \frac{x}{N}$$

Die Fixationswahrscheinlichkeit entspricht also genau dem Startanteil.

**Wie lange wird es typischerweise dauern, bis einer der beiden Typen verschwunden ist?** Dazu betrachten wir (zunächst) die erwartete *Stichprobenheterozygotie* 

$$\mathbb{E}_{x}\left[2\frac{X_{r}^{(N)}}{N}\left(1-\frac{X_{r}^{(N)}}{N}\right)\right]$$

Dies die Wahrscheinlichkeit, in einer zufälligen Stichprobe der Größe zwei (mit Zurücklegen) zwei unterschiedliche genetischen Typen vorzufinden.

Nach (I.I) ist für  $x \in \{0, 1, ..., N\}$ 

$$\mathbb{E}\left[2\frac{X_{r}^{(N)}}{N}\left(1-\frac{X_{r}^{(N)}}{N}\right)\middle|X_{r-1}^{(N)}=x\right]$$

$$=\frac{2}{N}\mathbb{E}\left[X_{r}^{(N)}\middle|X_{r-1}^{(N)}=x\right]-\frac{2}{N^{2}}\mathbb{E}\left[\left(X_{r}^{(N)}\right)^{2}\middle|X_{r-1}^{(N)}=x\right]$$

$$=\frac{2x}{N}-\frac{2}{N^{2}}\left(\mathbb{E}\left[X_{r}^{(N)}\middle|X_{r-1}^{(N)}=x\right]+\left(\mathbb{E}\left[X_{r}^{(N)}\middle|X_{r-1}^{(N)}=x\right]\right)^{2}\right)$$

$$=\frac{2x}{N}-\frac{2}{N^{2}}\left(x\left(1-\frac{x}{N}\right)+x^{2}\right)=\frac{2x}{N}\left(1-\frac{x}{N}\right)-\frac{2x}{N^{2}}\left(1-\frac{x}{N}\right)$$

$$=\left(1-\frac{1}{N}\right)2\frac{x}{N}\left(1-\frac{x}{N}\right)$$

somit

$$\mathbb{E}_{x}\left[2\frac{X_{r}^{(N)}}{N}\left(1-\frac{X_{r}^{(N)}}{N}\right)\right] = \mathbb{E}_{x}\left[\mathbb{E}_{x}\left[2\frac{X_{r}^{(N)}}{N}\left(1-\frac{X_{r}^{(N)}}{N}\right)\middle|X_{r-1}^{(N)}\right]\right]$$
$$= \left(1-\frac{1}{N}\right)\mathbb{E}_{x}\left[2\frac{X_{r-1}^{(N)}}{N}\left(1-\frac{X_{-1}^{(N)}}{N}\right)\right]$$

und iterativ

$$\mathbb{E}_{x}\left[2\frac{X_{r}^{(N)}}{N}\left(1-\frac{X_{r}^{(N)}}{N}\right)\right] = \left(1-\frac{1}{N}\right)^{r}2\frac{x}{N}\left(1-\frac{x}{N}\right) \tag{1.2}$$

Wir sehen: um bei großem N eine substantielle Änderung über r Generationen zu sehen, muss  $r \propto N$  sein, denn für  $r = \lfloor tN \rfloor$  und  $x = x^{(N)} = \lfloor yN \rfloor$  mit  $t \in (0, \infty)$  und  $y \in (0, 1)$  ergibt sich

$$\mathbb{E}_{\lfloor yN \rfloor} \left[ 2 \frac{X_{\lfloor tN \rfloor}^{(N)}}{N} \left( 1 - \frac{X_{\lfloor tN \rfloor}^{(N)}}{N} \right) \right] \underset{N \to \infty}{\longrightarrow} e^{-t} 2y (1-y)$$
(1.3)



Abbildung 1.1: Beobachtungen aus 105 *Drosophila melanogaster*-Populationsexperimenten (aus P. Buri, *loc. cit.*, Table 14, S. 387). Für die Generationen 0, 2, 4, ..., 18, 19 ist jeweils die empirische Verteilung der  $bw^{75}$ -Anteile über die 105 Populationen als Histogramm aufgetragen.

#### **Das Buri-Beispiel**

Peter Buri, Gene frequency in small populations of mutant Drosophila, *Evolution* 10, 367–402 (1956) berichtet ein Experiment in "künstlicher Evolution":

- 105 Populationen von jeweils konstant<sup>5</sup> 16 Taufliegen (8 weibl., 8 männl.) wurden für 19 Generationen (1 Gen. ≈ 14d) unter konstanten Bedingungen gehalten.
- 2 Allele: bw und  $bw^{75}$ , die 3 Genotypen bw/bw,  $bw/bw^{75}$ ,  $bw^{75}/bw^{75}$  sind anhand der Augenfarbe unterscheidbar
- Vorexperimente legten nahe, dass diese Genotypen keinen Einfluss auf den erwarteten Reproduktionserfolg haben.
- Die Anzahl bw<sup>75</sup>-Chromosomen in jeder Population und Generation wurde beobachtet, s.a. Abb. 1.1.

In dieser Laborsituation kann man die Wirkung der Gendrift direkt beobachten, siehe auch Abbildung 1.1. Beispielsweise passt der beobachtete Abfall der (empirischen) Heterozygotie, gemittelt über die 105 Populationen, recht gut zum mittels (1.2) theoretisch vorhergesagten geometrischen Abfallen der erwartete Stichprobenheterozygotie, allerdings muss der "reale" Populationsgrößenparameter 2N = 32 durch die "effektive" Populationsgröße 2N = 23 ersetzt werden, siehe Abb. 1.2.

<sup>&</sup>lt;sup>5</sup>Die konstante Populationsgröße wurde jeweils "von Hand" beim Umsetzen der nächsten Generation in ein neues Glas erzwungen, zwischenzeitlich waren die Populationen natürlich angewachsen.



Abbildung 1.2: Beobachtete empirische Heterozygotie und nach (1.2) angepasste Kurven für die (theoretische) erwartete Stichprobenheterozygotie für zwei Parameterwahlen (2*N* = 32 und 2*N* = 23)

**Erwartete Zeit bis zur Fixierung** Anhand der Rechnung rund um die erwartete Stichprobenheterozygotie haben wir in (1.3) bereits gesehen, dass für das Modell mit Populationsgröße *N* die Zeitskala *N* relevant ist.

Es gilt

$$\mathbb{E}_{x}[T_{\text{fix}}] = 1 + \sum_{y=0}^{N} \mathbb{P}_{x}(X_{1} = y) \mathbb{E}_{y}[T_{\text{fix}}], \quad x = 1, 2, \dots, N-1$$
(1.4)

mit Randwerten  $\mathbb{E}_0[T_{\text{fix}}] = \mathbb{E}_N[T_{\text{fix}}] = 0$  (dies verwendet Zerlegung nach dem ersten Schritt und die Markoveigenschaft, siehe z.B. [Bir24, Kap. 7.1]). Allerdings ist das resultierende lineare Gleichungssystem in N - 1 Unbekannten voll besetzt und nicht explizit lösbar.

Ein prinzipiell gangbarer Weg liegt in der Heuristik, dass  $\mathbb{E}_x[T_{\text{fix}}]$  für großes N in "genügend glatter" Weise von  $p \coloneqq \frac{x}{N}$  abhängen sollte, und dann eine Taylorenwicklung von  $\mathbb{E}_{pN}[T_{\text{fix}}]$  als Funktion von p anzusetzen.

**Satz 1.1.** Set  $x_N \in \mathbb{N}$  mit  $x_N/N \to p \in [0,1]$  für  $N \to \infty$ . Dann gilt

$$\lim_{N \to \infty} c_N \frac{\mathbb{E}_{x_N}[T_{\text{fix}}^{(N)}]}{2N} = H(p)$$
(1.5)

 $mit H(p) = -p \log(p) - (1-p) \log(1-p).$ 

Wir betrachten hier nur eine Beweisheuristik, siehe Abschnitt A.4 für das volle Argument.

*Beweisskizze.* Nehmen wir an, es gibt eine genügend glatte Funktion  $f_N : [0,1] \rightarrow \mathbb{R}_+$  mit

$$f_N\left(\frac{x}{N}\right) = \mathbb{E}_x[T_{\text{fix}}^{(N)}],$$

so gilt gemäß (1.4) mit Taylor-Entwicklung von  $f_N$  um  $\frac{x}{N}$ 

$$f_{N}\left(\frac{x}{N}\right) = 1 + \sum_{y=0}^{N} \mathbb{P}_{x}(X_{1} = y)f_{N}\left(\frac{y}{N}\right)$$
  
$$= 1 + \sum_{y=0}^{N} \mathbb{P}_{x}(X_{1} = y)\left[f_{N}\left(\frac{x}{N}\right) + \left(\frac{y-x}{N}\right)f_{N}'\left(\frac{x}{N}\right) + \frac{1}{2}\left(\frac{y-x}{N}\right)^{2}f_{N}''\left(\frac{x}{N}\right)\right] + R(N, x)$$
  
$$= 1 + f_{N}\left(\frac{x}{N}\right) + \frac{1}{N}f_{N}'\left(\frac{x}{N}\right)\mathbb{E}_{x}[X_{1} - x] + \frac{1}{2}\frac{1}{N^{2}}f_{N}''\left(\frac{x}{N}\right)\operatorname{Var}_{x}[X_{1} - x] + R(N, x)$$
  
$$= 1 + f_{N}\left(\frac{x}{N}\right) + \frac{1}{2N}\frac{x(N-x)}{N^{2}}f_{N}''\left(\frac{x}{N}\right) + R(N, x)$$

mit Restterm

$$R(N,x) = \sum_{y=0}^{N} \mathbb{P}_x(X_1 = y) \frac{1}{6} \left(\frac{y-x}{N}\right)^3 f_N^{\prime\prime\prime}(\zeta(\frac{x}{N}, \frac{y}{N}))$$

 $(\zeta(\frac{x}{N}, \frac{y}{N})$  ist eine Zahl zwischen  $\frac{x}{N}$  und  $\frac{y}{N}$ ).

Nun ist

$$\sum_{y=0}^{N} \mathbb{P}_{x}(X_{1} = y) \left(\frac{y - x}{N}\right)^{3} = \frac{1}{N^{3}} \mathbb{E}\left[\left(Y_{N, x/N} - x\right)^{3}\right]$$

mit  $Y_{N,x/N} \sim Bin(N, x/N)$  und es gilt

$$\max_{x=0,1,...,N} \mathbb{E}\Big[ |Y_{N,x/N} - x|^3 \Big] \le N^{3/2}$$

(siehe z.B. Lemma A.7 in Anhang A.4). Daher ist

$$R(N,x) = o(1/N)$$

zumindest plausibel.

Schreibe  $\frac{x}{N} = p$ , also erfüllt  $f_N$  näherungsweise

$$f_N''(p) = -2N \frac{1}{p(1-p)}, \quad 0 
(1.6)$$

mit den Randbedingungen  $f_N(0) = f_N(1) = 0$ . Man sieht nun leicht, dass eine explizite Lösung von (1.6) für  $p \in (0, 1)$  gegeben ist durch

$$f_N(p) = -2N(p\log(p) + (1-p)\log(1-p)),$$

denn

$$f'_N(p) = 2N(\log(p) + 1 - \log(1 - p) - 1) = -\frac{2}{c_N}(\log(p) - \log(1 - p)),$$

und

$$f_N''(p) = \left(-2N(\log(p) - \log(1-p))\right)' = -2N\left(\frac{1}{p} + \frac{1}{1-p}\right).$$

Der Satz zeigt insbesondere, dass für  $X_0 = N/2$  die erwartete Zeit bis zur Absorption von entweder a oder A gegeben ist durch

$$\mathbb{E}_{N/2}[T_{\text{fix}}] \approx -2N(1/2\log(1/2) + 1/2\log(1/2)) = 2\log(2) \cdot N \approx 1,39 \cdot N$$

Generationen.

## 1.2 Genealogien und Kingmans Koaleszent

**Genealogischer Blickpunkt** Das Wright-Fisher-Modell genealogisch ausgesprochen: Wir nummerieren die Individuen jeder Generation  $r \in \mathbb{Z}$  mit i = 1, 2, ..., N durch. Sei

$$A_{r,i}^{(N)} :=$$
Nr. des Vorfahren (in Gen.  $r - 1$ ) von Ind. Nr.  $i$  in Generation  $r$ , (1.7)

aus der Modellannahme: die  $A_{r,i}^{(N)}$ ,  $r \in \mathbb{Z}$ ,  $i \in [N]$  sind u.i.v. uniform auf  $[N] \coloneqq \{1, 2, \dots, N\}$ .

**Ahnenverhältnisse** Sei  $A_{r,i}^{(N)}[k]$  die Nummer des Ahnen vor k Generationen von Individuum Nr. i in Generation r [dieser Ahne lebte in Generation r - k], aus (1.7) ist diese

rekursiv bestimmt durch  $A_{r,i}^{(N)}[1] = A_{r,i}^{(N)}$  und  $A_{r,i}^{(N)}[k+1] = A_{r-k,A_{r,i}^{(N)}[k]}^{(N)}$  für  $k \in \mathbb{N}$ .

Wir betrachte eine Stichprobe von n verschiedenen (zufällig gezogenen) Individuen aus Generation r = 0, sagen wir die Individuen Nr.  $J_1, \ldots, J_n$  mit  $\mathbb{P}(J_1 = j_1, \ldots, J_n = j_n) = 1/(N)_{n\downarrow}$  für paarweise verschiedene  $j_1, \ldots, j_n \in [N]$ . (Wir notieren fallende Faktorielle als  $(x)_{k\downarrow} := x(x-1)(x-2)\cdots(x-k+1)$  für  $x \in \mathbb{R}, k \in \mathbb{N}$  mit Setzung  $(x)_{0\downarrow} = 1$ .)

Die Verwandtschaftsverhältnisse innerhalb der Stichprobe kodieren wir durch

 $R_k^{(N,n)}$ , eine (zufällige) Äquivalenzrelation,

gegeben durch  $i \sim_k j$   $(i, j \in [n], k = 0, 1, ...)$ , wenn  $A_{0,J_i}^{(N)}[k] = A_{0,J_j}^{(N)}[k]$  gilt, d.h. Stichproben i und j haben denselben Ahnen vor k Generationen.

Sei

 $\mathcal{E}_n \coloneqq \{ \ddot{\mathsf{A}}$ quivalenzrelationen auf  $[n] \}$ 

wir notieren  $\xi \in \mathcal{E}_n$  etwa durch eine (ungeordnete) Liste der Äquivalenzklassen (z.B.  $\xi = \{\{1\}, \{2, 3\}\} \in \mathcal{E}_3$  bedeutet  $2 \sim_{\xi} 3, 1 \neq_{\xi} 2, 1 \neq_{\xi} 3$ ).

Wir schreiben  $\xi \leq \eta$ , falls

$$i \sim_{\xi} j \implies i \sim_{\eta} j \qquad \text{gilt},$$

d.h.  $\eta$  entsteht aus  $\xi$  durch Vereinigung einiger Klassen, ggfs. in mehreren Gruppen (beispielsweise ist  $\{\{1\}, \{2\}, \{3,4\}, \{5\}, \{6,7\}, \{8\}\} \le \{\{1,2,6,7\}, \{3,4,5\}, \{8\}\}$ ).

Offensichtlich ist  $i \sim_0 j \iff i = j$ , d.h.  $R_0^{(N,n)} = \{\{1\}, \{2\}, \dots, \{n\}\}$  und es gilt stets  $R_k^{(N,n)} \leq R_{k+1}^{(N,n)}$ .

Betrachten wir zunächst den Fall n = 2: Für eine Stichprobe der Größe n = 2 ist die "korrekte" Zeitskala der Genealogie [Vielfache von] N, denn die Paarverschmelzungsw'keit ist

$$p^{(N,2)}(\{\{1\},\{2\}\},\{1,2\}) = \frac{1}{N}$$

und somit die Zeit

$$\tau_1^{(N,2)} \coloneqq \inf \left\{ k \in \mathbb{N} : R_k^{(N,n)} \{ \{1,2\} \} \right\}$$

bis die Stichprobe ihren ersten gemeinsamen Vorfahren findet (gemessen in Generationen), ~ geom(1/N), d.h.  $\mathbb{P}(\tau_1^{(N,2)} > 0) = (1 - 1/N)^k$ . Somit für  $t \in \mathbb{R}_+$ 

$$\mathbb{P}\left(\frac{\tau_1^{(N,2)}}{N} > t\right) = \left(1 - \frac{1}{N}\right)^{\left[Nt\right]} \xrightarrow[N \to \infty]{} e^{-t}$$

d.h. für große N ist  $\tau_1^{(N,2)}/N$  ungefähr Exp(1)-verteilt. (Das passt auch zur Beobachtung aus (1.3), dass die "relevante" Zeitskala N ist.)

Für n = 3 sieht die Übergangsmatrix  $p^{(N,3)}(\cdot, \cdot)$  von  $R^{(N,3)}$  folgendermaßen aus:

	$\{\{1\},\{2\},\{3\}\}$	$\{\{1,2\},\{3\}\}$	$\{\{1,3\},\{2\}\}$	$\{\{1\},\{2,3\}\}$	$\{\{1,2,3\}\}$
$\{\{1\},\{2\},\{3\}\}$	$1 - 3\frac{1}{N} + 2\frac{1}{N^2}$	$\frac{1}{N}(1-\frac{1}{N})$	$\frac{1}{N}(1-\frac{1}{N})$	$\frac{1}{N}(1-\frac{1}{N})$	$\frac{1}{N^2}$
$\{\{1,2\},\{3\}\}$	0	$1 - \frac{1}{N}$	0	0	$\frac{1}{N}$
$\{\{1,3\},\{2\}\}$	0	0	$1 - \frac{1}{N}$	0	$\frac{1}{N}$
$\{\{1\},\{2,3\}\}$	0	0	0	$1 - \frac{1}{N}$	$\frac{1}{N}$
$\{\{1, 2, 3\}\}$	0	0	0	0	1

d.h.

**Lemma 1.2.** Für festes  $N \ge n$  ist  $(R_k^{(N,n)})_{k \in \mathbb{N}_0}$  eine Markovkette mit Werten in  $\mathcal{E}_n$ . Die Übergangswahrscheinlichkeiten sind gegeben durch

$$p^{(N,n)}(\xi,\eta) := \mathbb{P}(R_{k+1}^{(N,n)} = \eta | R_k^{(N,n)} = \xi) = \frac{(N)_{a\downarrow}}{N^b}$$

sofern  $\xi \leq \eta$ , wobei  $\eta$  aus  $|\eta| = a$  Klassen besteht,  $\xi$  aus  $b = |\xi| = b_1 + \dots + b_a$  Klassen besteht und  $\eta$  aus  $\xi$  durch Verschmelzen von a Gruppen von Klassen in Gruppengrößen  $b_1, \dots, b_a$  entsteht (d.h.  $\eta = \{C_1, \dots, C_a\}$  und  $\xi = \{C_{\alpha\beta} : 1 \leq \alpha \leq a, 1 \leq \beta \leq b_\alpha\}$  mit  $C_\alpha = \cup_{\beta=1}^{b_\alpha} C_{\alpha\beta}$  für  $\alpha = 1, \dots, a$ ).

*Beweis.*  $R_k^{(N,n)} = \xi$  bedeutet, dass es (k Generationen vor der Gegenwart) b verschiedene "aktive Ahnenlinien" geben muss, d.h. es gibt b paarweise verschiedene Zahlen  $i_{\alpha,\beta} \in [N], \beta = 1, ..., b_{\alpha}, \alpha = 1, ..., a$ , so dass

$$A_{0,J_i}^{(N)}[k] = i_{\alpha,\beta} \text{ für } i \in C_{\alpha\beta}, \ \beta = 1, \dots, b_{\alpha}, \ \alpha = 1, \dots, a$$

gilt. Gegeben dies tritt das Ereignis  $\{R_{k+1}^{(N,n)} = \eta\}$  genau dann ein, wenn es paarweise verschiedene  $j_1, j_2, \ldots, j_a \in [N]$  gibt mit

$$A_{-k,i_{\alpha,\beta}}^{(N)} = j_{\alpha} \text{ für } \beta = 1, \dots, b_{\alpha}, \ \alpha = 1, \dots, a.$$
(1.8)

Nach Konstruktion des Wright-Fisher-Modells gilt für jede solche Wahl

$$\mathbb{P}\left(A_{-k,i_{\alpha,\beta}}^{(N)}=j_{\alpha} \text{ für } \beta=1,\ldots,b_{\alpha}, \ \alpha=1,\ldots,a\right)=\prod_{\alpha=1}^{a}\prod_{\beta=1}^{b_{\alpha}}\frac{1}{N}=\frac{1}{N^{b}}$$

und es gibt  $N \cdot (N-1) \cdot (N-2) \cdots (N-a+1) = (N)_{a\downarrow}$  viele mögliche Wahlen.

**Beobachtung 1.3.** Für  $\xi, \eta \in \mathcal{E}_n \min \xi \leq \eta$ , wobei  $|\eta| = a$ ,  $|\xi| = b = b_1 + b_2 + \dots + b_a \min b_1, \dots, b_a \geq 1$  zeigt Lemma 1.2

I. 
$$p^{(N,n)}(\xi,\xi) = \frac{(N)_{b\downarrow}}{N^b} = \prod_{i=0}^{b-1} \left(1 - \frac{i}{N}\right) = 1 - \frac{1 + 2 + \dots + (b-1)}{N} + O\left(\frac{1}{N^2}\right)$$
$$= 1 - \frac{1}{N} {b \choose 2} + O\left(\frac{1}{N^2}\right)$$

2. Falls a = b - 1, d.h.  $\eta$  entsteht aus  $\xi$  durch Verschmelzung genau eines Paars von Klassen,

$$p^{(N,n)}(\xi,\eta) = \frac{(N)_{(b-1)\downarrow}}{N^b} = \frac{1}{N} \prod_{i=0}^{b-2} \left(1 - \frac{i}{N}\right) = \frac{1}{N} + O\left(\frac{1}{N^2}\right)$$

3. Falls  $a \leq b - 2$ , d.h. mehr als zwei Klassen sind an Verschmelzung(en) beteiligt,

$$p^{(N,n)}(\xi,\eta) = O\left(\frac{1}{N^2}\right)$$

Dies legt nahe, den zeitreskalierten Prozess der Ahnenverhältnisse  $(R_{\lfloor Nt \rfloor}^{(N,n)})_{t \ge 0}$  zu betrachten. Es stellt sich heraus, dass dieser gegen eine zeitkontinuierliche Markovkette konvergiert. Für dazu notwendige Techniken siehe den Exkurs in Abschnitt A.1.

Die allgemeine Struktur des Grenzwerts (für beliebige Stichprobengröße n) ist folgende:

**Definition 1.4.** Die zeitkontinuierliche Markovkette  $(R_t^{(n)})_{t\geq 0}$  auf  $\mathcal{E}_n$  mit Sprungratenmatrix

$$q_{\xi\eta} = \begin{cases} 1 & \text{falls } \eta \text{ aus } \xi \text{ durch Verschmelzung von genau zwei Klassen entsteht,} \\ -\binom{|\xi|}{2} & \text{falls } \eta = \xi, \\ 0 & \text{sonst} \end{cases}$$
(1.9)

heißt Kingmans<sup>6</sup> (*n*-)Koaleszent.

Zumeist betrachten wir den Startzustand  $R_0^{(n)} = \{\{1\}, \{2\}, \dots, \{n\}\}$ . Wir können den Pfad  $(R_t^{(n)})_{t\geq 0}$  als Baum interpretieren, dessen Blätter mit  $1, \ldots, n$  markiert sind: Zu den Zeitpunkten  $0 = \tau_n^{(n)} < \tau_{n-1}^{(n)} < \cdots < \tau_2^{(n)} < \tau_1^{(n)}$ , wo

$$\tau_k^{(n)} \coloneqq \inf\{t \ge 0 : |R_t^{(n)}| \le k\}$$

verschmelzen jeweils zwei Zweige. Für Stichproben  $i, j \in [n]$  können wir den genealogischen Abstand von i und  $j, \inf\{t \ge 0: i \sim_{R_{\star}^{(n)}} j\}$  aus dem Baum ablesen.

Satz 1.5. Es gilt

$$(R_{\lfloor Nt \rfloor}^{(N,n)})_{t\geq 0} \longrightarrow (R_t^{(n)})_{t\geq 0} \quad f \ddot{u}r N \to \infty.$$

Wir beweisen die in Satz 1.5 formulierte Konvergenz im Sinne der endlich-dimensionalen Verteilungen; tatsächlich gilt auch Konvergenz in Verteilung auf dem Pfadraum  $D([0,\infty),\mathcal{E}_n)$ .

*Beweis von Satz 1.5.* Fixiere *n*. Gemäß Lemma A.2 müssen wir zeigen, dass für  $\xi, \eta \in \mathcal{E}_n$  gilt

$$p^{(N,n)}(\xi,\eta) = \delta_{\xi,\eta} + \frac{1}{N}q_{\xi\eta} + o\left(\frac{1}{N}\right)$$
(1.10)

 $[da |\mathcal{E}_n| < \infty \text{ ist dann der Fehler gleichmäßig klein, d.h. wir zeigen, dass } \lim_{N \to \infty} \max_{\xi, \eta \in \mathcal{E}_n} N |p^{(N,n)}(\xi, \eta) - N| = 0$  $\left|\delta_{\xi,\eta} - (1/N)q_{\xi\eta}\right| = 0 \text{ gilt}].$ 

Dies folgt aus Beobachtung 1.3 und der Form der Sprungraten aus (1.9).

<sup>&</sup>lt;sup>6</sup>J.F.C. Kingman, The coalescent, Stochastic Process. Appl. 13 (1982), no. 3, 235–248.

**Beobachtung 1.6.** 1. (Die Zeit bis zum jüngsten gemeinsamen Vorfahren) Aus der Struktur der Sprungratenmatrix (1.9) folgt

$$\tau_1^{(n)} = \left(\tau_{n-1}^{(n)} - \tau_n^{(n)}\right) + \left(\tau_{n-2}^{(n)} - \tau_{n-1}^{(n)}\right) + \dots + \left(\tau_1^{(n)} - \tau_2^{(n)}\right) \stackrel{d}{=} S_n + S_{n-1} + \dots + S_2, \qquad n \ge 2,$$

wobei die  $S_k$  unabhängige exponentialverteilte Zufallsvariablen mit Parameter  $\binom{k}{2}$  sind, somit

$$\mathbb{E}[\tau_1^{(n)}] = \sum_{k=2}^n E[S_k] = \sum_{k=2}^n \frac{2}{k(k-1)} = 2\sum_{k=2}^n \left(\frac{1}{k-1} - \frac{1}{k}\right) = 2\left(1 - \frac{1}{n}\right)$$

und 1 =  $\mathbb{E}[\tau_1^{(2)}] \leq \mathbb{E}[\tau_1^{(n)}] < \lim_{m \to \infty} \mathbb{E}[\tau_1^{(m)}] = 2.$ 

Für ein Populationsmodell (aus der von uns betrachteten Schar) mit Populationsgröße N bedeutet dies, das der jüngste gemeinsame Vorfahre der heute lebenden Population im Mittel vor etwa 2N Generationen gelebt hat.

Weiter ist

$$\operatorname{Var}[\tau_1^{(n)}] = \sum_{k=2}^n \operatorname{Var}[S_k] = \sum_{k=2}^n {\binom{k}{2}}^{-2} = \sum_{k=2}^n \frac{4}{k^2(k-1)^2} = \sum_{k=2}^n \left\{ 4\left(\frac{1}{k^2} + \frac{1}{(k-1)^2}\right) + 8\left(\frac{1}{k} - \frac{1}{k-1}\right) \right\}$$
$$= \left\{ 8\sum_{k=1}^{n-1} \frac{1}{k^2} \right\} - 4 + \frac{4}{n^2} + \frac{8}{n} - 8 = \left\{ 8\sum_{k=1}^{n-1} \frac{1}{k^2} \right\} - 4\left(1 - \frac{1}{n}\right)\left(3 + \frac{1}{n}\right),$$

insbesondere

$$1 = \operatorname{Var}[\tau_1^{(2)}] \le \operatorname{Var}[\tau_1^{(n)}] < \lim_{n \to \infty} \operatorname{Var}[\tau_1^{(n)}] = 8\frac{\pi^2}{6} - 12 \approx 1.16$$

Der wesentliche Beitrag zur Gesamtvarianz kommt also von der letzten Verschmelzungszeit  $S_2$ .

2. (Teilstichproben-Konsistenz) Sei  $\pi_{n,n-1} : \mathcal{E}_n \to \mathcal{E}_{n-1}$  die Einschränkung aller Äquivalenzklassen auf [n-1], so gilt

$$(\pi_{n,n-1}(R_t^{(n)}))_{t\geq 0} \stackrel{d}{=} (R_t^{(n-1)})_{t\geq 0}.$$

Dies folgt aus der Form der Sprungraten oder auch aus der Tatsache, dass für die Approximanten (wie in Satz 1.5) nach Konstruktion  $\pi_{n,n-1}(R_k^{(N,n)}) = R_k^{(N,n-1)}$  (realisierungsweise) gilt.

3. (Invarianz der Verteilung unter Permutation der Stichprobennummern) Für eine Permutation  $\sigma$ von [n] und  $\xi = \{C_1, \ldots, C_a\} \in \mathcal{E}_n$  sei  $\sigma(\xi) = \{\sigma(C_1), \ldots, \sigma(C_a)\}$  die Äquivalenzrelation, die man erhält, indem man die Elemente der Blöcke von  $\xi$  gemäß  $\sigma$  umnummeriert. Es gilt

$$\left(\sigma(R_t^{(n)})\right)_{t\geq 0} \stackrel{d}{=} \left(R_t^{(n)}\right)_{t\geq 0}.$$

Dies folgt aus der Symmetrie der Sprungraten oder auch aus der Tatsache, dass für die Approximanten (wie in Satz 1.5) nach Konstruktion  $(\sigma(R_k^{(N,n)}))_{k\in\mathbb{N}_0} \stackrel{d}{=} (R_k^{(N,n)})_{k\in\mathbb{N}_0}$  gilt. [Man sagt auch, dass  $R_t^{(n)}$  eine austauschbare zufällige Äquivalenzrelation ist.]

**Bericht.** Mit Beob. 1.6, 2. und Kolmogorovs Erweiterungssatz ist es möglich, den Kingman-Koaleszenten  $(R_t)_{t\geq 0}$  mit Stichprobengröße  $n = \infty$  als Markovprozess auf  $\mathcal{E} := {\text{Äquivalenzrelationen auf } \mathbb{N}}$  mit



Abbildung 1.3: Zwei Realisierungen des Kingman-100-Koaleszenten (wobei die Blätter jeweils so sortiert wurden, dass der Baum überschneidungsfrei zu zeichnen ist)

Startwert  $R_0 = \{\{1\}, \{2\}, \dots\}$  zu definieren mit der Eigenschaft  $(\pi_{\infty,n}(R_t))_{t\geq 0} \stackrel{d}{=} (\sigma(R_t^{(n)}))_{t\geq 0}$  für jedes  $n \in \mathbb{N}$ .

Siehe auch Konstruktion 1.9 unten für eine explizite Kopplung aller (Kingman-)*n*-Koaleszenten via "look down", die dies ohne (allzugroßen) Theorie-Aufwand leistet.

Beob. 1.6, 1. zeigt, dass  $\mathbb{E}[\tau_1^{(\infty)}] = 2 < \infty$ , d.h. auch eine "unendlich große Stichprobe" findet f.s. in endlicher Zeit ihren ersten gemeinsamen Vorfahren. Obwohl  $|R_0| = \infty$  ist, gilt  $|R_t| < \infty$  für jedes t > 0 fast sicher. Man sagt auch, dass der Kingman-Koaleszent "aus dem Unendlichen herabsteigt." Siehe auch die Simulationsbilder in Abbildung 1.3 für einen Eindruck dieses Phänomens.

### 1.3 Moran-Modell

Das Moran-Modell ist gewissermaßen das "zeitkontinuierliche Analogon" zum Wright-Fisher-Modell, es ist (ebenfalls) eines der fundamentalen Modelle der mathematischen Populationsgenetik.

**Definition 1.7** ((Neutrales 2 Typ-)Moran-Modell<sup>7</sup>). Man betrachtet eine Population von konstant N (haploiden) Individuen, jedes Individuum besitzt eine unabhängige, Exp(1)-verteilte Lebenszeit und wird am Ende seiner Lebenszeit durch den Nachkommen eines rein zufällig aus der Population gezogenen Individuums ersetzt (es gibt nur ein Elter und, sagen wir, man kann durch sein eigenes Kind ersetzt werden).

Wir nehmen zusätzlich an, dass es zwei Typen A und a gibt, die ohne Mutation vererbt werden. Sei

 $X_t^{(N)}$  = Anzahl Typ *A*-Ind. zur Zeit *t*.

<sup>&</sup>lt;sup>7</sup>Nach Patrick Alfred Pierce Moran, 1917–1988 benannt

Angesichts der Gedächtnislosigkeit der Exponentialverteilung ist  $(X_t^{(N)})_{t\geq 0}$  eine zeitkontinuierliche Markovkette mit Werten in  $\{0, 1, \ldots, N\}$  und Sprungraten

$$q_{i,i+1} = i \frac{(N-i)}{N} = (N-i) \frac{i}{N} = q_{i,i-1}, \quad q_{i,i} = -2 \frac{i(N-i)}{N}$$

(die übrigen Einträge der Sprungratenmatrix  $Q = (q_{i,j})$  sind = 0).

#### **Graphische Konstruktion**

Für jedes geordnete Paar  $(i, j), i, j \in \{1, ..., N\}, i \neq j$  sei  $(N_t^{(i,j)})_{t\geq 0}$  ein Poissonprozess auf  $\mathbb{R}_+$  mit Rate  $\frac{1}{N}$ , u.a. für verschiedene Paare. Zu den Sprungzeiten von  $(N_t^{(i,j)})_{t\geq 0}$  stirbt Individuum j und wird durch einen Nachkommen von Individuum i ersetzt (s.a. Abb. 1.4).

[Bild an der Tafel]

Abbildung 1.4: Im Bild: *N* Kopien der Zeitachse, gerichtete Pfeile zwischen ihnen zu den Sprungzeitpunkten von u.a. Poissonprozessen; das Individuum an der Pfeilspitze stirbt jeweils und wird durch einen Nachkommen des Individuums am Pfeilschaft ersetzt.

Sei

 $X_t(i)$  = Typ von Individuum *i* zur Zeit *t*.

Die Dynamik des Prozesses  $(X_t(1), X_t(2), \dots, X_t(N))_{t \ge 0}$ , der über die Typen der Individuen in der Population Buch führt (und nicht nur über die Anzahlen) ist somit folgende:

Ersetze zu jedem Sprungzeitpunkt t von  $N^{(i,j)}$  den Typ  $X_{t-}(j)$  durch  $X_t(j) = X_{t(-)}(i)$ .

Dies ist wohldefiniert, da unabhängige Poissonprozesse f.s. keine gemeinsamen Sprungzeitpunkte besitzen. Diese Konstruktion ist ein Spezialfall einer sogenannten Harris-Konstruktion<sup>8</sup>, ein in der Theorie der interagierenden Teilchensysteme übliches (und nützliches) Werkzeug.

#### 1.4 Dualität

**Bemerkung 1.8** (Ablesen der Genealogie und der Typen aus der graphischen Konstruktion). Für  $t > 0, i \in [N]$  sei

$$A_s^{(i,t)} =$$
Nr. des Ahnenindividuums zur Zeit  $t - s$  von Ind.  $i$  zur Zeit  $t$  (für  $0 \le s \le t$ , Werte in  $[N]$ )

Zur Konstruktion von  $A^{(i,t)} = (A_s^{(i,t)})_{0 \le s \le t}$  verfolgen wir die derzeitige "Zeitachse" rückwärts und folgen den Pfeilen jeweils in entgegengesetzter Richtung, vgl. auch Abb. 1.4.

<sup>&</sup>lt;sup>8</sup>nach Theodore Edward Harris, 1919–2005 benannt

In Formeln können wir den Pfad von  $A^{(i,t)}$  beispielsweise folgendermaßen fassen (wir schreiben  $N^{(j,i)}([a,b))$  für die Anzahl Sprünge des Poissonprozesses  $N^{(j,i)}$  im Zeitintervall [a,b)):

Sei  $T_0^{(i,t)} := 0$ ,  $\widetilde{A}_0^{(i,t)} := A_0^{(i,t)} := i$ , für  $k \in \mathbb{N}$  setzen wir

$$T_{k}^{(i,t)} \coloneqq \inf \left\{ u > T_{k-1}^{(i,t)} \colon \text{es gibt ein } j \neq i \text{ mit } N^{(j,\widetilde{A}_{k-1}^{(i,t)})} ([t-u,t-T_{k-1}^{(i,t)})) = 1 \right\}$$

bzw.  $T_k^{(i,t)} \coloneqq t$ , falls es kein solches u gibt. Falls  $T_k^{(i,t)} = t$  gilt, so setzen wir  $M^{(i,t)} \coloneqq k$  und wir brechen die Konstruktion hier ab, andernfalls sei

$$\widetilde{A}_{k}^{(i,t)} \operatorname{das}(\mathrm{f.s.}) \operatorname{eindeutig} \operatorname{bestimmte} j \operatorname{mit} N^{(j,\widetilde{A}_{k-1}^{(i,t)})} \left( \left[ t - T_{k}^{(i,t)}, t - T_{k-1}^{(i,t)} \right) \right) = 1$$

und wir setzen die Konstruktion fort. Da die endlich vielen Poissonprozesse  $N^{(j,i)}$  f.s. keine Häufungspunkte in [0, t] besitzen, bricht die Konstruktion mit Wahrscheinlichkeit 1 nach endlich vielen Schritten ab und wir setzen dann für  $0 < s \le t$ 

$$A^{(i,t)}_s \coloneqq \widetilde{A}^{(i,t)}_\ell \, \text{ falls } T^{(i,t)}_\ell \leq s < T^{(i,t)}_{\ell+1} \text{ für } 0 \leq \ell < M^{(i,t)}$$

bzw.  $A_t^{(i,t)} \coloneqq \widetilde{A}_{M^{(i,t)}-1}^{(i,t)}$ .

 $(A_s^{(i,t)})_{0 \le s \le t}$  ist eine zeitkontinuierliche Markovkette mit (vollkommen symmetrischen) Sprungraten

$$q_{jk} = \begin{cases} \frac{1}{N}, & k \neq j, \\ -\frac{N-1}{N}, & k = j, \end{cases}$$
(1.11)

man nennt eine solche Kette auch eine (zeitkontinuierliche) "Irrfahrt auf dem vollständigen Graphen  $V_N$  der Ordnung N".

Für  $i_1 \neq i_2$  bewegen sich  $A^{(i_1,t)}$  und  $A^{(i_2,t)}$  unabhängig bis zum "Verschmelzungszeitpunkt"

$$\tau_{i_1,i_2} \coloneqq \inf\{s \in [0,t] : A_s^{(i_1,t)} = A_s^{(i_2,t)}\},\$$

ab dann, d.h. für  $u \ge \tau_{i_1,i_2}$ , gilt  $A_u^{(i_1,t)} = A_u^{(i_2,t)}$ .

Für paarweise verschiedene  $i_1, i_2, \ldots i_n (\leq N)$  bilden

 $A^{(i_1,t)},\ldots,A^{(i_n,t)}$  ein System verschmelzender Irrfahrten auf  $V_n$ 

und mit

$$k \sim_{s,N} \ell : \iff A_s^{(i_k,t)} = A_s^{(i_\ell,t)}, \quad 1 \le k, \ell \le n$$

ist

$$\mathcal{R}_{s}^{(n,N)} \coloneqq \ddot{\mathsf{A}}$$
quivalenzklassen bezüglich  $\sim_{s,N}, s \in [0, t]$ 

ein (zeittransformierter) Kingman-*n*-Koaleszent. ( $(\mathcal{R}_{Ns/2}^{(n,N)})_{s\geq 0}$  wäre wörtlich ein Koaleszent, wenn wir die Zeitachsen in der graphischen Konstruktion "bis – $\infty$  fortsetzten.")

Aus der Konstruktion ergibt sich folgende (realisierungsweise Form) der "Dualität":

$$X_t(i) = X_0(A_t^{(i,t)}) \quad \text{für } 1 \le i \le N, t > 0.$$
(I.12)

*Beweisskizze.* Die Tatsache, dass  $A^{(i,t)}$  eine zeitkontinuierliche Markovkette ist, folgt anschaulich gesehen aus der Unabhängigkeit der Zuwächse der "treibenden" Poissonprozesse  $N^{(j,k)}$ , die symmetrische Form der Sprungratenmatrix (1.11) stammt daher, dass alle Poissonprozesse dieselbe Rate 1/Nhaben. Wenn aktuell  $A_s^{(i,t)} = j$ , so gibt es für  $0 < h \ll 1$  und jedes  $j' \neq j$  mit Wahrscheinlichkeit  $\approx h/N$  einen Sprung von  $N^{(j',j)}$  im Zeitintervall [t-s-h, t-s) und dann springt  $A^{(i,t)}$  von j nach j'.

Etwas formaler: Sei

$$\mathcal{F}_u^t \coloneqq \sigma \left( N^{(j,k)}([a,b)) : j \neq k, t - u \le a < b \le t \right)$$

die  $\sigma$ -Algebra, die die Informationen über alle Sprünge der  $N^{(j,k)}$  zwischen t - u und t enthält. Offenbar kann man  $A_s^{(i,t)}$  für  $s \le u$  anhand der Pfade der  $N^{(j,k)}$  zwischen t - u und t rekonstruieren (d.h.  $A_s^{(i,t)}$  ist  $\mathcal{F}_u^t$ -messbar für  $s \le u$ ) und für  $s < t, j \ne j' \in [N]$  ist auf dem Ereignis  $\{A_s^{(i,t)} = j\}$ 

$$\frac{1}{h} \mathbb{P} \left( A_{s+h}^{(i,t)} = j' \,|\, \mathcal{F}_s^t \right) \\
= \frac{1}{h} \mathbb{P} \left( N^{(j',j)} ([t-s-h,t-s)) = 1, N^{(j'',j)} ([t-s-h,t-s)) = 0 \text{ für } j'' \neq j' \right) + \frac{1}{h} R_h \\
= \frac{1}{h} \cdot e^{-h/N} \frac{h/N}{1!} \cdot \left( e^{-h/N} \right)^{N-2} + R_h = \frac{1}{N} + o(1)$$

für  $h \downarrow 0$ , wobei der Resterm

$$|R_h| \le \mathbb{P}\left(\sum_{k\neq\ell}^N N^{(k,\ell)}([t-s-h,t-s)) \ge 2\right) = O(h^2)$$

erfüllt.

Wir können die Formel (1.12) und den dazugehörigen Gedankengang verwenden, um (faktorielle) Momente des Typenanteilsprozesses im Moran-Modell zu berechnen: Betrachten wir ein Moran-Modell mit Populationsgröße N und Startanzahl  $X_0^{(N)} = x_0^{(N)}$  von Typ A-Individuen. Für  $t \ge 0$ und  $n \in \mathbb{N}$  ist

$$\mathbb{E}_{x_0^{(N)}}\left[\frac{X_t^{(N)}(X_t^{(N)}-1)\cdots(X_t^{(N)}-n+1)}{N(N-1)\cdots(N-n+1)}\right]$$

die Wahrscheinlichkeit, bei n Zügen ohne Zurücklegen aus der Population zur Zeit t jedesmal Typ A zu ziehen. Andererseits seien  $J_1, \ldots, J_n$  mit  $\mathbb{P}(J_1 = j_1, \ldots, J_n = j_n) = 1/(N)_{n\downarrow}$  für paarweise verschiedene  $j_1, \ldots, j_n \in [N]$  die Nummern der n zur Zeit t gezogenen Individuen. Obige Wahrscheinlichkeit ist

$$\mathbb{P}_{x_0^{(N)}} \Big( X_t(J_1) = X_t(J_2) = \dots = X_t(J_n) = A \Big) \\ = \mathbb{P}_{x_0^{(N)}} \Big( X_0(A_t^{(J_1,t)}) = X_0(A_t^{(J_2,t)}) = \dots = X_0(A_t^{(J_2,t)}) = A \Big) \\ = \mathbb{E} \Big[ \prod_{i=1}^{\#\{A_t^{(J_1,t)},\dots,A_t^{(J_n,t)}\}} \frac{x_0^{(N)} - i + 1}{N - i + 1} \Big]$$

wobei wir für die erste Gleichung (1.12) verwenden und für die zweite Gleichung beobachten, dass die Nummern  $A_t^{(J_1,t)}, \ldots, A_t^{(J_n,t)}$  der Ahnenindividuen der Stichprobe (es kann in dieser Liste Mehrfacheinträge geben) gerade  $\#\{A_t^{(J_1,t)}, \ldots, A_t^{(J_n,t)}\}$  Zügen ohne Zurücklegen aus [N] entsprechen. Da  $\#\{A_t^{(J_1,t)}, \ldots, A_t^{(n,N)}\} = d \#\mathcal{R}_t^{(n,N)}$  gilt, folgt die "Stichprobendualitätsformel"

$$\mathbb{E}_{x_0^{(N)}}\left[\frac{X_t^{(N)}(X_t^{(N)}-1)\cdots(X_t^{(N)}-n+1)}{N(N-1)\cdots(N-n+1)}\right] = \mathbb{E}\left[\prod_{i=1}^{\#\mathcal{R}_t^{(n,N)}}\frac{x_0^{(N)}-i+1}{N-i+1}\right]$$
(1.13)

Die Beobachtung, dass  $\mathcal{R}_{Nt/2}^{(n,N)} = d R_t^{(n)}$ , der Kingman *n*-Koaleszent zur Zeit *t*, liefert damit für  $z \in [0,1]$  und  $n \in \mathbb{N}_0$ 

$$\lim_{N \to \infty} \mathbb{E}_{\lfloor Nz \rfloor} \left[ \left( X_{Nt/2}^{(N)} / N \right)^n \right] = \mathbb{E}_n \left[ z^{\#R_t} \right]$$
(1.14)

Dies ist zumindest ein Indiz, dass der reskalierte Typenanteilsprozess  $(X_{Nt/2}^{(N)}/N)_{t\geq 0}$  ein Limesobjekt besitzt (die Wright-Fisher-Diffusion).

**Bemerkung.** Eine analoge Überlegung greift für das Wright-Fisher-Modell aus Abschnitt 1.1 unter Verwendung der Ahneninformationen aus Abschnitt 1.2. Man findet für das Wright-Fisher-Modell

$$\lim_{N \to \infty} \mathbb{E}_{\lfloor Nz \rfloor} \left[ \left( X_{\lfloor Nt \rfloor}^{(N)} / N \right)^n \right] = \mathbb{E}_n \left[ z^{\#R_t} \right]$$

Der Limesprozess ist tatsächlich derselbe wie beim Moran-Modell.

**Konstruktion 1.9** (eine explizite Kopplung aller (Kingman-)*n*-Koaleszenten via "look down"). Für  $1 \le i < j$  seien  $(L_{j,i}(t))_{t\ge 0}$  unabhängige Poissonprozesse auf  $\mathbb{R}_+$  mit Rate 1. Für  $k \in \mathbb{N}$  sei  $(A_k(t))_{t\ge 0}$  gegeben als die Lösung von

$$A_k(t) = k - \sum_{1 \le i < j \le k} \int_0^t (j-i) \mathbf{1} (A_k(s-) = j) L_{j,i}(ds), \quad t \ge 0$$

(wobei wir  $L_{j,i}$  als die Verteilungsfunktion des – zufälligen – Zählmaßes auffassen, das jeweils an den Sprungstellen von  $L_{j,i}$  Atome der Masse 1 besitzt). Zur Veranschaulichung des Systems der Lösungen  $A_k$  folgendes Bild: Für jedes k = 1, 2, ... betrachte eine Kopie der Zeitachse auf Niveau k, für i < jzeichne zu den Sprungzeitpunkten des Prozesses  $L_{j,i}$  einen Pfeil von Niveau j nach Niveau i.



Weiter sei  $Q = (q_{\ell,m})_{\ell,m\in\mathbb{N}}$  mit

$$q_{\ell,m} = \begin{cases} 1, & 1 \le m < \ell \\ -(\ell - 1), & m = \ell, \\ 0, & \text{sonst} \end{cases}$$

Dann gilt:

a) Für  $k \in \mathbb{N}$  ist  $A_k = (A_k(t))_{t \ge 0}$  Markovkette auf  $[k] \coloneqq \{1, \ldots, k\}$ , deren Sprungratenmatrix durch die Einschränkung von Q auf [k] gegeben ist.

b) Die Prozesse  $A_k, k \in \mathbb{N}$  bilden ein System verschmelzender Markovketten, d.h. für  $k \neq \ell$  und  $t \ge 0$  gilt

$$A_k(t) = A_\ell(t) \implies A_k(t+s) = A_\ell(t+s)$$
 für all  $s > 0$ 

c) Für  $n \in \mathbb{N}$  und  $t \ge 0$  definieren wir eine Äquivalenzrelation  $R_t^{(n)}$  of [n] via

$$k \sim_{R_{\ell}^{(n)}} \ell \iff A_k(t) = A_\ell(t)$$

 $(R_t^{(n)})_{t\geq 0}$  ist verteilt wie Kingmans-*n*-Koaleszent.

d)  $R_t^{(n)}$  entsteht aus  $R_t^{(n+1)}$  durch Einschränkung der Klassen von  $R_t^{(n+1)}$  auf [n] (eine etwaige leere Klasse  $\emptyset = \{n + 1\} \cap [n]$  wird stillschweigend entfernt).

Insbesondere ist für  $t \ge 0$  die zufällige Aquivalenzrelation  $R_t$  of  $\mathbb{N}$  via

$$k \sim_{R_t} \ell \iff k \sim_{R^{(n)}} \ell \text{ für } n \ge \max\{k, \ell\}$$

wohldefiniert. Der Prozess  $(R_t)_{t\geq 0}$  beschreibt den Kingman-Koaleszenten, der mit unendlich vielen Blättern startet.

## 1.5 Beispiel: Die Beobachtungen von Dorit et al, 1995

Robert L. Dorit, Hiroshi Akashi und Walter Gilbert berichten in Absence of Polymorphism at the ZFY Locus on the Human Y Chromosome, *Science* 268, 1183–1185 (1995) die Ergebnisse einer genetischen Studie:

- Weltweite<sup>9</sup> Stichprobe von 38 Männern (*homo sapiens*)
- Ein 729 Basenpaare langes, nicht-kodierendes Stück des Y-Chromosoms (das 3. Intron des ZFY-Gens) wurde für jede Stichprobe sequenziert
- Es wurden keinerlei Mutationen gefunden: Alle 38 Stichproben identisch
- Inter-spezies-Vergleich mit Schimpanse, Gorilla, Orang-Utan (und Pavian als "outgroup") zeigt, dass am betrachteten Lokus Mutationen vorkommen können

<sup>&</sup>lt;sup>9</sup> Loc. cit., S. 1184: "Human DNA samples were obtained from male volunteers who donated hair follicle samples or from cell lines provided by L. L. Cavalli-Sforza and K. K. Kidd. Geographic origins were determined by interview. Whenever possible, geographic origins of parents and grandparents were also ascertained. The samples are grouped by continent of origin, and the number of individuals is given in parentheses. Africa: Nigeria\* (1), Ivory Coast (1), Tanzania (1), Southern Africa (2), Algeria (1), Central African Republic\* (2), African American (2); Americas: Mexico (2), Guatemala (1), Peru\* (1), Argentina (1), Native American (2); Asia: China\* (2), Korea (1), Japan\* (2), Taiwan (2), Indonesia (1), India (1); Europe/Middle East: Ireland\* (1), Belgium (1), Italy\* (1), Spain (1), Russia\* (2), Poland\* (1), Saudi Arabia\* (1), Turkey (1); South Pacific: Melanesia (1), New Guinea\* (1), Australia\* (1). (\*) Indicates samples where the 3<sup>2</sup>-most zinc-finger exon was also sequenced."

 Molekulare Uhr-Annahme und auf Fossilien beruhende Annahmen über die Zeit seit der Aufspaltung von der Vorfahren von Mensch und Schimpanse bzw. Orang-Utan ergeben geschätzte Rate von (fixierten) Mutationen

 $1,35 \times 10^{-3}$  Mutationen pro Basenpaar pro Million Jahre

Was können wir angesichts dieser Beobachtungen über die Zeit bis zum jüngsten gemeinsamen Vorfahren der gezogenen 38 Y-Chromosomen (und damit implizit auch über den jgV aller heute lebenden Männer) sagen?

Wir verwenden den Kingman-Koaleszenten als Modell der Genealogie.

A-priori-Verteilung Ohne Berücksichtigung der Beobachtungen würden wir annehmen, dass

$$T_{jgV} \stackrel{d}{=} S_{38} + S_{37} + \dots + S_2$$

wo  $T_{jgV}$  die Zeit (in Koaleszenten-Zeiteinheiten) bis zum jüngsten gemeinsamen Vorfahren der 38 gezogenen Männer, die  $S_k$  unabhängig mit  $S_k \sim Exp(\binom{k}{2})$ ,

1 Koaleszenten-Zeiteinheit  $\widehat{=} N_{\text{eff}} \times g$  Jahre

mit  $N_{\text{eff}}$  ... effektive Populationsgröße (für Männer), g ... Generationslänge (in Jahren), also

a-priori-Verteilung:  $\mathscr{L}(T_{jgV}) = \underset{k=2}{\overset{38}{*}} \operatorname{Exp}(\binom{k}{2})$ , d.h. mit Lemma A.6 ist die Dichte

$$f_{\text{a-pri}} = \sum_{i=2}^{38} {\binom{i}{2}} \exp\left(-\binom{i}{2}t\right) \prod_{j=2, j \neq i}^{38} \frac{\binom{j}{2}}{\binom{j}{2} - \binom{i}{2}}$$

und

$$\mathbb{E}[T_{jgV}] = \sum_{k=2}^{38} \frac{2}{k(k-1)} = 2\left(1 - \frac{1}{38}\right) = \frac{37}{19} \doteq 1,947,$$

5%-Quantil von  $T_{jgV}$ :  $q_{0,05} = 0,744,95\%$ -Quantil:  $q_{0,95} = 4,041$ .

Mit Annahmen  $N_{\text{eff}} = 5.000$ , g = 20 Jahre übersetzt sich dies zu MW = 195.000 Jahre,  $q_{0.05} = 74.000$  Jahre,  $q_{0.95} = 404.000$  Jahre.



**A-posteriori-Verteilung** Sei  $S_k$  die Länge des Zeitintervalls (in Koaleszenten-Zeiteinheiten) während dessen die Genealogie der Stichprobe aus k Linien bestand,  $M_k$  die Anzahl Mutationen, die während dieses Intervalls in der Genealogie auftreten.

Gegeben

$$S_{k} = t \text{ ist } M_{k} \text{ Poisson-verteilt mit Parameter } tk\frac{\theta}{2}, \text{ d.h.}$$

$$\mathbb{P}(M_{k} = m \mid S_{k} = t) = \exp\left(-tk\frac{\theta}{2}\right)\frac{\left(tk\frac{\theta}{2}\right)^{m}}{m!} \text{ wobei } \theta = 2N_{\text{eff}} \times g \times \mu$$
(1.15)

mit  $N_{\text{eff}}$  effektive Populationsgröße, g Generationslänge (in Jahren),  $\mu$  Mutationsrate der betrachteten Region im Genom (pro Jahr) (und gegeben  $S_2, \ldots, S_n$  sind  $M_2, \ldots, M_n$  unabhängig).

Eine heuristische Begründung für (1.15) ist folgende (wir werden dies im weiteren Verlauf der Vorlesung noch genauer betrachten): Angesichts Satz 1.5 entsprechen t Koaleszenten-Zeiteinheiten im Populationsmodell mit Populationsgröße N etwa $t/c_N$  Generationen und somit etwa $tg/c_N = tgN_{\text{eff}}$ Jahren. Wenn wir annehmen, dass pro Jahr (unabhängig von allem anderen) eine Mutation mit der (sehr kleinen) Wahrscheinlichkeit  $\mu$  auftritt, so ist die Verteilung der Anzahl Mutationen, die wir auf einem Stück der Genealogie dieser Länge sehen,

$$\operatorname{Bin}(tgN_{\operatorname{eff}},\mu) \approx \operatorname{Poi}(tgN_{\operatorname{eff}}\mu) = \operatorname{Poi}(t\theta/2).$$

Da gegeben  $S_k = t$  der Teil der Genealogie, währenddessen k Linien existieren, aus k Stücken von je t Koaleszenten-Zeiteinheiten besteht, ist

$$\mathbb{P}(M_k = m \mid S_k = t) = \underbrace{\operatorname{Poi}(t\theta/2) * \cdots * \operatorname{Poi}(t\theta/2)}_{k \text{ mal}} = \operatorname{Poi}(tk\theta/2).$$

(Die Normierung des Mutationsparameters als  $\theta/2$  hat historische Gründe und sorgt auch dafür, dass manche Formeln "schöner" aussehen.)

Frage: Wie ist  $S_{38} + \dots + S_2$  verteilt, gegeben dass  $M_{38} + \dots + M_2 = 0$ ?  $S_k \sim \text{Exp}(\binom{k}{2})$ ,  $\mathscr{L}(M_k | S_k = t) = \text{Poi}(tk\theta/2)$ , dann ist

$$\mathbb{P}(M_k = m) = \int_0^\infty {\binom{k}{2}} \exp\left(-\binom{k}{2}t\right) e^{-tk\theta/2} \frac{(tk\theta/2)^m}{m!} dt$$
$$= {\binom{k}{2}} \frac{(k\theta/2)^m}{m!} \int_0^\infty t^m \exp\left(-\binom{k}{2} + k\theta/2\right) t dt = \frac{k-1}{k-1+\theta} \left(\frac{\theta}{k-1+\theta}\right)^m,$$

(wir substituieren  $u = (\binom{k}{2} + k\theta/2)t$  und nutzen  $\int_0^\infty u^m e^{-u} du = \Gamma(m+1) = m!$ ) d.h.  $\mathscr{L}(M_k) = \operatorname{Geom}(\frac{k-1}{k-1+\theta})$  — man könnte dies alternativ auch über ein "konkurrierende Raten"-Argument einsehen. Weiter ist damit

$$\mathbb{P}(S_k \le t \mid M_k = 0) = \frac{k - 1 + \theta}{k - 1} \int_0^t {\binom{k}{2}} \exp\left(-{\binom{k}{2}s}\right) e^{-sk\theta/2} ds.$$
$$\mathscr{L}(S_k \mid M_k = 0) = \exp\left(\frac{k(k - 1 + \theta)}{2}\right).$$

Bedingt auf  $M_2 = \cdots = M_{38} = 0$  sind  $S_2, \ldots, S_{38}$  (weiterhin) unabhängig.

Demnach: Verteilung der Zeit bis zum jüngsten gemeinsamen Vorfahren (in Koaleszenten-Zeiteinheiten), bedingt auf  $M := M_2 + \dots + M_{38} = 0$  ist

$$T_{jgV}|_{\{M=0\}} \stackrel{d}{=} S'_{38} + S'_{37} + \dots + S'_{22}$$

mit  $S'_k$  u.a.,  $S'_k \sim \operatorname{Exp}\left(\frac{k(k-1+\theta)}{2}\right)$ , d.h.  $\mathscr{L}(T_{jgV}|M=0) = \overset{38}{*} \operatorname{Exp}\left(\frac{k(k-1+\theta)}{2}\right)$ . Die Dichte von

 $\mathscr{L}(T_{jgV} | M = 0)$  (and damit auch den Erwartungswert und die Verteilungsfunktion) können wir wiederum mit Lemma A.6 bestimmen.

Wir fixieren g = 20a,  $\mu = 729 \times 1.35 \cdot 10^{-9}a^{-1} \doteq 0.98 \cdot 10^{-6}a^{-1}$  (aus Dorit et al (1995), diese Werte waren auch in der Literatur unstrittig), so hängt die bedingte Verteilung von  $T_{jgV}$  (und nicht nur ihre "Übersetzung in Realzeit") vom Parameter  $N_{eff}$  ab.

$N_{\rm eff}$	EW	$q_{0,05}$	$q_{0,95}$
2.500	91.519	35.851	187.369
5.000	173.007	69.263	349.909
10.000	313.234	130.095	620.279
20.000	532.785	233.853	1.020.819

Dichte von  $T_{jgV}$  (Realzeit)





Wir sehen insbesondere: Die Verteilung von  $T_{igV}$  hängt in nicht-linearer Weise von  $N_{eff}$  ab.

**Literaturbericht** Die ursprüngliche Studie erschien in Robert L. Dorit, Hiroshi Akashi und Walter Gilbert, Absence of Polymorphism at the ZFY Locus on the Human Y Chromosome, *Science* 268, 1183–1185 (1995), siehe auch die Diskussionsbeiträge ("technical comments") dazu in *Science* 272, 1356– 1362 (1996) von Y.-X. Fu und W.-H. Li, von P. Donnelly, S. Tavaré, D.J. Balding und R.C. Griffiths, von G. Weiss und A. von Haeseler und von J. Rogers, P.B. Samollow und A.G. Comuzzie, die insbesondere eine etwas unpräzise Anwendung der Koaleszenten-Theorie von Dorit et al korrigierten. Die Darstellung fußt in weiten Teilen auf Wakeley [Wako9, Ch. 8.1]

# Kapitel 2

## Mutationen und der markierte Koaleszent

Wir betrachten in diesem Kapitel die Situation, dass die Individuen der Population verschiedene gentische Typen (abstrakt: aus einer Menge E möglicher Typen) haben und dass – im Gegensatz zur Situation in Kapitel I – Kinder mit einer gewissen (typischerweise kleinen) Wahrscheinlichkeit einen anderen Typ als ihr Elter haben (sogenannte "Mutationen"). Konkret stellen wir uns das Wright-Fisher-Modell aus Abschnitt I.I mit Populationsgröße N durch folgenden Mechanismus ergänzt vor: Jedes Kind ist unabhängig mit Wahrscheinlichkeit

$$\mu_N = \frac{\theta}{2N} \tag{2.1}$$

ein "mutiertes" Kind, wobei  $\theta \in (0, \infty)$ . (Welche Anderung des Typs diese Mutation bewirkt, bleibt noch zu spezifizieren, wir betrachten in den folgenden Abschnitten konkrete Wahlen.)

Annahme (2.1) führt dazu, dass für  $t \in (0, \infty)$  auf einer einzelnen Ahnenlinie über  $\lfloor Nt \rfloor$  Generationen  $Bin(\lfloor Nt \rfloor, \theta/2N) \approx Pois(t\theta/2)$  viele Mutationen auftreten. In der Skalierung der Genealogie einer *n*-Stichprobe wie in Satz 1.5 bedeutet dies: Das Limesobjekt ist ein Kingman-Koaleszent mit *n* Blättern, längs dessen Ästen Mutationen gemäß einem Poissonprozess mit Rate  $\theta/2$  auftreten (vgl. auch die Diskussion in Abschnitt A.1).

**Bemerkung.** Die Annahme (2.1), die Populationsgröße und Mutationswahrscheinlichkeiten aneinander koppelt, mag auf den ersten Blick unnatürlich erscheinen: Warum sollte die Mutationswahrscheinlichkeit, die sich aus der Effektivität der biochemischen Kopier- und Reparaturmechanismen innerhalb der Zellen eines Individuums, ggfs. im Zusammenspiel mit diversen Umwelteinflüssen, ergibt, irgend etwas mit der Populationsgröße zu tun haben?

Annahme (2.1) bedeutet, dass Mutation und Gendrift "auf derselben Zeitskala" wirken. Wenn  $\mu_N \ll 1/N$ , so wirkt die Gendrift (wie wir wissen, über Zeiten der Größenordnung O(N)), bevor Mutationen irgendeinen merklichen Einfluss auf die Zusammensetzung der Population haben, wenn  $\mu_N \gg 1/N$ , so stellt sich "reines Mutationsgleichgewicht" ein, an dem Gendrift nichts ändert.

Anders gewendet: Für eine gegebene Population (mit endlichem, aber sehr großem N) ist

$$2N\mu_N = \theta$$

der "entscheidende" Parameter, um die Zeitentwicklung des Typenanteils zu beschreiben.

## 2.1 Infinitely-many-alleles-Modell (IMA)

**Definition 2.1** (Infinitely-many-alleles-Mutationsmechanismus). Wir treffen die Modellannahme, dass jede Mutation einen völlig neuen Typ erzeugt (und die Mutationen sind "neutral", d.h. sie beeinflussen den Fortpflanzungserfolg nicht). In der Literatur ist auch der Name "infinite alleles model" üblich.

Mathematisch realisieren wir dies durch die Wahl E = [0, 1] als Typenmenge, bei jedem Mutationsereignis wird (unabhängig) der neue Typ uniform([0, 1])-verteilt gewählt.

Wenn Mutationen sich mit Rate  $\theta/2 > 0$  ereignen, ist die Vorwärtsentwicklung des Typs längs einer Abstammungslinie demnach beschrieben durch den Markov-Prozess mit Generator

$$Bf(x) = \frac{\theta}{2} \int_0^1 f(u) - f(x) \, du, \ x \in [0, 1]$$

für  $f : [0, 1] \rightarrow \mathbb{R}$  beschränkt und messbar.

Betrachte Stichprobe der Größe n: Beobachtete genetische Variabilität modelliert durch n-Koaleszent, längs dessen Kanten sich mit Rate  $\frac{\theta}{2}$  Mutationen gem. IMA-Modell ereignen.

```
[Bild an der Tafel]
```

Offenbar ist nur der Teil der Genealogie jeweils "oberhalb" der jüngsten Mutation relevant, für die Beobachtungen an den Blättern können wir also folgende äquivalente Dynamik betrachten:

**Definition 2.2** ("Getöteter *n*-Koaleszent"). • Beginne mit  $\{\{1\}, \{2\}, \ldots, \{n\}\}$  (alle aktiv).

- Jedes Paar von aktiven Klassen verschmilzt mit Rate 1.
- Jede Klasse wird mit Rate  $\frac{\theta}{2}$  getötet/inaktiviert: allen Elementen wird derselbe, unif([0, 1])verteilte Typ zugeordnet und die Klasse wird inaktiviert

(sie hat ihre "definierende Mutation" getroffen).

• Ende, wenn keine aktiven Klassen mehr übrig.

Analog zu Def. 1.4 (Kingman-Koaleszent) kann man dies formal als zeitkontinuierliche Markovkette auf

 $\widetilde{\mathcal{E}}_n = \{ \ddot{A}$ quivalenzrelationen auf [n], deren Klassen als aktiv/inaktiv markiert sind  $\}$ 

ausformulieren (Übung: man stelle die Sprungratenmatrix auf).

**Bemerkung 2.3.** Angesichts der Symmetrien des Koaleszenten ist die eigentlich relevante Information das *Typenhäufigkeitsspektrum*  $(B_1, \ldots, B_n)$ , wobei

 $B_i$  = #Typen, die *i*-mal in der Stichprobe vorkommen, i = 1, ..., n

(offenbar ist  $\sum_{i=1}^{n} i B_i = n$ ).

Gegeben  $(B_1, \ldots, B_n) = (b_1, \ldots, b_n)$  mit  $\sum_{i=1}^n ib_i = n$  und

$$\sum_{i=1}^{n} b_i = k$$

sind die beobachteten Typen in der Stichprobe folgendermaßen verteilt: Zerlege  $\{1, ..., n\}$  uniform in k Teilmengen der Größen

$$\underbrace{1,\ldots,1}_{b_1},\underbrace{2,\ldots,2}_{b_2},\ldots,\underbrace{n}_{b_n},$$

ordne jeder Teilmenge u.a. einen unif([0, 1])-verteilten Typ zu.

Seien  $E_n, \ldots, E_1$  ZVn mit Werten in {koal, mut},

 $E_k$  = Typ des Ereignisses, das die Anz. aktiver Klassen von k auf k - 1 reduziert.

Es gilt

$$\mathbb{P}(E_k = \text{coal}) = \frac{\binom{k}{2}}{\binom{k}{2} + k\frac{\theta}{2}} = \frac{k-1}{k-1+\theta},$$
$$\mathbb{P}(E_k = \text{mut}) = \frac{k\frac{\theta}{2}}{\binom{k}{2} + k\frac{\theta}{2}} = \frac{\theta}{k-1+\theta},$$

und  $E_n, \ldots, E_1$  sind unabängig (verwende "konkurrierende-Raten"-Argument und Symmetrien der Sprungraten). Also gilt für  $e_n, e_{n-1}, \ldots, e_1 \in \{\text{coal}, \text{mut}\}$ 

$$\mathbb{P}(E_n = e_n, \dots, E_1 = e_1) = \frac{\prod_{k=1}^n \left(\theta \mathbf{1}(e_k = \text{mut}) + (k-1)\mathbf{1}(e_k = \text{coal})\right)}{\prod_{k=1}^n (k-1+\theta)}.$$
 (2.2)

**Definition 2.4** (Hoppe<sup>1</sup>-Urne). Urne enthält eine schwarze Kugel ("Mutationskugel") mit Masse  $\theta$ , farbige Kugeln jew. mit Masse 1.

- Zu Beginn: Urne enthält nur die schwarze Kugel.
- In jedem Schritt: Ziehe eine Kugel mit W'keit proportional zu ihrer Masse.
- Falls farbige Kugel gezogen: Lege zurück zusammen mit einer weiteren Kugel derselben Farbe.
- Falls schwarze Kugel gezogen: Lege zurück zusammen mit einer weiteren Kugel einer völlig neuen Farbe.

**Lemma 2.5.** Die von den n nicht-schwarzen Kugeln erzeugte Verteilung der Familiengrößen (Typenhäufigkeitsspektrum) nach n Zügen der Hoppe-Urne entspricht der des getöteten n-Koaleszenten (aus Def. 2.2).

<sup>&</sup>lt;sup>1</sup>Fred M. Hoppe, Pólya-like urns and the Ewens' sampling formula, J. Math. Biol. 20, no. 1, 91–94, (1984).

Beweisskizze. Lese (2.2) "rückwärts."

Sei

 $K_n$  = #verschiedene Typen in *n*-Stichprobe.

**Satz 2.6.** Im IMA-Modell ( $\theta > 0$  fest) gilt für  $n \to \infty$ 

$$\mathbb{E}[K_n] = \sum_{i=1}^n \frac{\theta}{\theta + i - 1} \sim \theta \log n, \quad \operatorname{Var}[K_n] = \sum_{i=1}^n \frac{\theta}{\theta + i - 1} \cdot \frac{i - 1}{\theta + i - 1} \sim \theta \log n,$$
  
und  $\frac{K_n - \mathbb{E}[K_n]}{\sqrt{\operatorname{Var}(K_n)}} \xrightarrow{d} \mathcal{N}(0, 1).$ 

Beweis. Verwende Hoppe-Urne, schreibe

$$K_n = A_1 + \dots + A_n$$

mit

 $A_i = 1$  (im *i*-ten Zug wurde die schwarze Kugel gezogen).

Nach Konstruktion sind  $A_1, \ldots, A_n$  u.a. mit

$$\mathbb{P}(A_i=1) = \theta/(\theta+i-1).$$

Für die Asymptotik vergleiche mit Riemann-Integral, für asymptotische Normalität bilde ein (unabängiges, zentriertes, normiertes) Dreiecksschema

$$X_{ni} = \frac{A_i - \frac{\theta}{\theta + i - 1}}{\sqrt{\operatorname{Var}[K_n]}},$$

dies erfüllt (trivialerweise) die Lindeberg-Bedingung: Es gilt für

$$S_n = X_{n1} + \dots + X_{nn}$$

 $\operatorname{Var}[S_n] = 1$ , für jedes  $\varepsilon > 0$  gilt wegen  $\operatorname{Var}[K_n] \to \infty$ , dass  $\mathbb{E}[X_{ni}^2 \mathbf{1}(X_{ni}^2 > \varepsilon)] = 0$  für n genügend groß, insbesondere

$$L_n(\varepsilon) \coloneqq \sum_{i=1}^n \mathbb{E}[X_{ni}^2 \mathbf{1}(X_{ni}^2 > \varepsilon)] \to 0.$$

**Bemerkung 2.7** (Hoppes Urne und zufällige Permutationen<sup>2</sup>). n Züge aus der Hoppe-Urne generieren sukzessive eine zufällige Permutation  $\Pi_n$  von  $\{1, \ldots, n\}$  (in Zyklen-Darstellung) folgendermaßen:

Nummeriere die farbigen Kugeln in der Reihenfolge des Erscheinens, wenn im k-ten Zug

<sup>&</sup>lt;sup>2</sup>Eine Beobachtung aus Paul Joyce, Simon Tavaré, Cycles, permutations and the structure of the Yule process with immigration, *Stochastic Process. Appl.* 25, no. 2, 309–314, (1987).

- schwarze Kugel gezogen : füge neuen Zyklus (k) hinzu,
   (insbesondere: nach dem ersten Zug entsteht die Identität (1))
- farbige Kugel  $j_k$  gezogen : füge im Zyklus, der  $j_k$  enthält, links von  $j_k$  ein.

Für jede Permutation  $\pi$  mit k Zyklen gilt dann

$$\mathbb{P}(\Pi_n = \pi) = \frac{\theta^k}{\theta(\theta+1)\cdots(\theta+n-1)},$$

denn  $\pi$  legt die Reihenfolge der Ereignisse i.d. Urne fest, wenn im k-ten Zug schwarze Kugel gezogen: Faktor  $\frac{\theta}{\theta+k-1}$ , wenn farbige gezogen: Faktor  $\frac{1}{\theta+k-1}$ .

Man kann dies als eine Version des sogenannten China-Restaurant-Prozesses auffassen, siehe z.B. [Kle20], Kap. 24.3 (und speziell S. 523f dort).

**Satz 2.8** (Ewens'sche Stichprobenformel<sup>3</sup>). Seien  $b_1, \ldots, b_n \in \mathbb{N}_0$  mit

$$\sum_{j=1}^{n} b_j = k \le n \quad und \quad \sum_{j=1}^{n} jb_j = n$$

gegeben. Die Wahrscheinlichkeit, in einer n-Stichprobe (im IMA-Modell) jeweils  $b_j$  Typen mit genau j Repräsentanten (für j = 1, ..., n) zu beobachten, ist

$$\frac{n!}{1^{b_1}2^{b_2}\cdots n^{b_n}} \cdot \frac{1}{b_1!b_2!\cdots b_n!} \cdot \frac{\theta^k}{\theta(\theta+1)\cdots(\theta+n-1)}.$$
(2.3)

Man kann (2.3) auch schreiben als

$$C(n,\theta) \times \prod_{j=1}^{n} e^{-\theta/j} \frac{(\theta/j)^{b_j}}{b_j!}$$
(2.4)

mit

$$C(n,\theta) = n! \exp\left(\theta \sum_{j=1}^{n} 1/j\right) / (\theta(\theta+1)\cdots(\theta+n-1)),$$

d.h. die Verteilung des Typenhäufigkeitsspektrums  $(B_1, \ldots, B_n)$  in einer *n*-Stichprobe ist  $\bigotimes_{j=1}^n \text{Poi}(\theta/j)$ , bedingt auf  $\sum_{j=1}^n jB_j = n$ .

*Beweis.* Man kann dies per Induktion beweisen, indem man nach dem "jüngsten" Ereignis im markierten Koaleszenten zerlegt. Wir betrachten hier ein direktes, kombinatorisches Argument aus dem Artikel Robert C. Griffiths, Sabin Lessard, Ewens' sampling formula and related formulae: combinatorial proofs, extensions to variable population size and applications to ages of alleles, *Theoretical Population Biology* 68, no. 3, 167–177, (2005).

Seien  $E_n, E_{n-1}, \ldots E_1$  die "Elementarübergänge" in der Historie des getöteten Koaleszenten aus Def. 2.2 (vgl. (2.2) und Diskussion vor der Hoppe-Urne, Def. 2.4),

 $E_m$  beschreibt das Ereignis, das die Anzahl aktiver Linien von m auf m-1 reduziert.

<sup>&</sup>lt;sup>3</sup>Warren J. Ewens, The sampling theory of selectively neutral alleles, *Theoretical Population Biology* 3, 87–112, (1972) und S. Karlin, J. McGregor, Addendum to a paper of W. Ewens, *Theoretical Population Biology* 3, 113–116, (1972).

Wir nummerieren die Linien mit  $1, \ldots, n$  und führen (genauer) Buch, welche Linie von einer Mutation getroffen wird bzw. welche Linie mit welcher Linie verschmilzt, mögliche Werte der Elementarübergänge sind also

mut(i), Linie *i* trifft eine Mutation, und

koal $(i \rightarrow j)$ , Linie *i* verschmilzt in Linie  $j (\neq i)$ .

Wir denken bei den Verschmelzungereignissen an "gerichtete" Verschmelzungen, d.h. für jedes aktuell noch aktive Paar von Linien i und j

verschmilzt Linie i in Linie j mit Rate  $\frac{1}{2}$ .

Eine Liste  $e_n, e_{n-1}, \ldots, e_1$  solcher möglicher Elementarübergänge nennen wir ein *Feinprotokoll*. Sei für  $m \ge 2$ 

$$p_m(e_m) = \begin{cases} \frac{1/2}{\frac{1}{2}m(m-1) + \frac{\theta}{2}m} = \frac{1}{m(m-1+\theta)}, & \text{wenn } e_m \text{ eine Verschmelzung,} \\ \frac{\theta/2}{\frac{1}{2}m(m-1) + \frac{\theta}{2}m} = \frac{\theta}{m(m-1+\theta)}, & \text{wenn } e_m \text{ eine Mutation,} \end{cases}$$

 $p_1(e_1) = 1$ , wenn  $e_1$  eine Mutation ist, und  $p_1(e_1) = 0$  für eine (dann unmögliche) Verschmelzung.

Für ein gegebenes mögliches Feinprotokoll  $e_n, e_{n-1}, \ldots, e_2, e_1$  mit k Mutationsereignissen (und somit k Typen) ist

$$\mathbb{P}(E_n = e_n, \dots, E_1 = e_1) = p_n(e_n)p_{n-1}(e_{n-1})\cdots p_1(e_1) = \frac{\theta^k}{\prod_{m=1}^n m(m-1+\theta)} = \frac{\theta^k}{n!(\theta)_{n\uparrow}} \quad (2.5)$$

(Produkt der Übergangswahrscheinlichkeiten der Skelettkette des getöteten Kingman-n-Koaleszenten).

Für 
$$b_1, ..., b_n \in \mathbb{N}_0$$
 mit  $\sum_{j=1}^n b_j = k$  (und  $\sum_{j=1}^n jb_j = n$ ) gibt es
$$\frac{(n!)^2}{\prod_{j=1}^n (b_j! j^{b_j})}$$
(2.6)

mögliche Feinprotokolle, die auf dieses Typenhäufigkeitsspektrum  $(b_1, \ldots, b_n)$  führen. Das Produkt von (2.5) und (2.6) liefert (2.3).

Zu (2.6): Stellen wir uns für den Moment die k Typen ("künstlich") nummeriert vor, mit

Typenhäufigkeitsvektor  $(n_1, n_2, \ldots, n_k)$ ,

d.h.  $n_\ell$  Stich<br/>proben sind vom  $\ell\text{-ten}$  Typ, und es gilt

$$|\{\ell: n_\ell = j\}| = b_j, \quad j = 1, 2, \dots, n.$$

Es gibt

$$n! \times \binom{n}{n_1 \dots n_k} \times (n_1 - 1)! \dots (n_k - 1)! = \frac{(n!)^2}{\prod_{\ell=1}^k n_\ell} = \frac{(n!)^2}{\prod_{j=1}^k j^{b_j}}$$
(2.7)

verschiedene Feinprotokolle mit k nummerierten Typen, die auf diesen Typenhäufigkeitsvektor  $(n_1, n_2, ..., n_k)$  führen:

- 1. *n*! Möglichkeiten für die Reihenfolge, in der die Linien inaktiv werden,
- 2.  $\binom{n}{n_1 \dots n_k}$  Möglichkeiten, die *n* Linien auf die *k* Typen aufzuteilen (*n* nummerierte Kugeln in *k* Schachteln legen),
- 3. für  $\ell = 1, ..., k$  gibt es  $(n_{\ell} 1)!$  viele Arten, innerhalb von Typ  $\ell$  die "Verschmelzungsziele" festzulegen.

(Nehmen wir an, wir haben in Schritt 1 und 2 festgelegt, dass Linien  $i_1, i_2, \ldots, i_{n_\ell}$  vom Typ  $\ell$  sind und dass diese in der Reihenfolge  $i_1, i_2, \ldots$  inaktiviert werden. Dann gibt es  $n_\ell - 1$  viele Wahlen für das Verschmelzungsziel von  $i_1, n_\ell - 2$  viele Wahlen für das Verschmelzungsziel von  $i_2$ , u.s.w.)

Schließlich führen

$$b_1! \cdot b_2! \cdot \dots \cdot b_n! \tag{2.8}$$

verschiedene Feinprotokolle mit nummerierten Typen auf dasselbe Feinprotokoll ohne Typennummerierung:

Typen 
$$\ell$$
 und  $\ell'$  können vertauscht werden, sofern  $n_{\ell} = n_{\ell'}$ . (2.9)

Der Quotient von (2.7) und (2.8) liefert (2.6).

Bemerkung 2.9. Nach Satz 2.6 ist

$$\widehat{\theta}_{\text{naiv}} \coloneqq \frac{K_n}{\log n}$$

ein (schwach) konsistenter Schätzer für die Mutationsrate  $\theta$ , d.h. für jedes  $\theta \in (0, \infty)$  gilt

 $\widehat{\theta}_{\text{naiv}} \xrightarrow[n \to \infty]{} \theta$  stochastisch bzw. in Verteilung.

Satz 2.6 liefert auch asymptotische Normalität von  $\widehat{\theta}_{naiv}$ . Allerdings ist

$$\operatorname{Var}_{\theta}[\widehat{\theta}_{\operatorname{naiv}}] \sim \frac{\theta}{\log n}.$$

(Dies ist allerdings "deprimierend langsam": z.B. müsste  $n = e^{100} \approx 2.7 \cdot 10^{43}$  sein, damit die Streuung  $\approx 0.1\sqrt{\theta}$  ist.)

**Beobachtung 2.10.** Im IMA-Modell ist  $K_n$  suffizient für  $\theta$ , d.h. gegeben  $K_n = k$  hängt die Verteilung der beobachteten Typen nicht von  $\theta$  ab.

Beweis. Sei

$$C_{n,k} := \sum_{(b_1,...,b_n)}' \frac{n!}{\prod_{i=1}^n i^{b_i} b_i!},$$

 $(\sum' \text{ bezeichnet die Summe über } (b_1, \ldots, b_n) \in \mathbb{N}_0^n \text{ mit } \sum b_i = k, \sum i b_i = n).$ 

Satz 2.8 (Ewens-Formel) liefert

$$\mathbb{P}_{\theta}(K_n = k) = C_{k,n} \frac{\theta^k}{\theta(\theta + 1)\cdots(\theta + n - 1)},$$

für  $b_1, \ldots, b_n$  mit  $b_1 + \cdots + b_n = k$  (und  $\sum i b_i = n$ ) also

$$\mathbb{P}_{\theta}(B_1 = b_1, \dots, B_n = b_n | K_n = k) = \frac{1}{C_{k,n}} \frac{n!}{1^{b_1} 2^{b_2} \cdots n^{b_n}} \cdot \frac{1}{b_1! b_2! \cdots b_n!}.$$

Sei  $K_n = k$  beobachtet. Der Maximum-Likelihood-Schätzer  $\widehat{\theta}_{ML}$  ist dasjenige  $\theta \ge 0$ , das die Likelihood

$$L_n(\theta, k) = \mathbb{P}_{\theta}(K_n = k) = C_{k,n} \frac{\theta^k}{\theta(\theta + 1)\cdots(\theta + n - 1)}$$

(als Funktion von  $\theta$ ) maximiert.

Es ist

$$\frac{\partial}{\partial \theta} \log L_n(\theta, k) = \frac{\partial}{\partial \theta} \left( k \log \theta - \sum_{i=0}^{n-1} \log(\theta + i) \right) = \frac{k}{\theta} - \sum_{i=0}^{n-1} \frac{1}{\theta + i} = \frac{1}{\theta} \left( k - \sum_{i=0}^{n-1} \frac{\theta}{\theta + i} \right),$$

also ist  $\widehat{\theta}_{\rm ML}$  Lösung von

$$k = \sum_{i=0}^{n-1} \frac{\widehat{\theta}_{\mathrm{ML}}}{\widehat{\theta}_{\mathrm{ML}} + i} \quad \left( \dots = \mathbb{E}_{\widehat{\theta}_{\mathrm{ML}}} [K_n]^{\mathsf{``}} \right)$$

(d.h. derjenige  $\theta$ -Wert, unter dem erwartet=beobachtet).

Die Fisher-Information ist

$$I(\theta) = \mathbb{E}_{\theta} \left[ \left( \frac{\partial}{\partial \theta} \log L_n(\theta, K_n) \right)^2 \right] = \frac{1}{\theta^2} \mathbb{E}_{\theta} \left[ \left( K_n - \sum_{i=0}^{n-1} \frac{\theta}{\theta+i} \right)^2 \right] = \frac{1}{\theta^2} \operatorname{Var}_{\theta}(K_n)$$

 $(\sim \frac{1}{\theta} \log n \text{ nach Satz 2.6}).$ 

#### 2.1.1 Die GEM-Verteilung

Betrachte die Hoppe-Urne (für festes n alternativ: den Koaleszent mit Mutationen gemäß IMA-Modell), die Typen/Familien seien in Reihenfolge des Erscheinens nummeriert ("age order")

 $X_k(n) :=$  Größe der *k*-ten Familie nach dem *n*-ten Zug aus der Hoppe-Urne,

(offenbar  $X_1(1) = 1, X_k(n) = 0$  für k > n).

Frage:

$$\left(\frac{1}{n}X_1(n), \frac{1}{n}X_2(n), \frac{1}{n}X_3(n), \dots\right) \underset{n \to \infty}{\longrightarrow} ?$$
(2.10)

Beobachtung:  $(n + \theta)^{-1}X_1(n)$ , n = 2, 3, ..., ist ein (beschränktes) Martingal:

$$\mathbb{E}\Big[\frac{1}{n+1+\theta}X_1(n+1)\,\Big|\,\mathcal{F}_n\Big] = \frac{X_1(n)+1}{n+1+\theta}\frac{X_1(n)}{n+\theta} + \frac{X_1(n)}{n+1+\theta}\frac{\theta+n-X_1(n)}{n+\theta} = \frac{X_1(n)}{n+\theta}$$

(mit  $\mathcal{F}_n = \sigma$  (Beobachtungen bis einschließlich *n*-tem Zug)), konvergiert also f.s. Analog ist für  $k \ge 2$ mit  $\alpha_k$  = Zeitpunkt des ersten Auftretens von Typ k

$$(n + \alpha_k + \theta)^{-1} X_k(n + \alpha_k), \quad n = 1, 2, \dots$$

ein (beschr.) Martingal.

Demnach: (2.10) konvergiert f.s. (zumindest koordinaten-weise).

**Satz 2.11** (GEM-Verteilung<sup>4</sup>). Seien  $B_1, B_2, \ldots$  u.i.v. Beta $(1, \theta)$ , d.h. sie besitzen die Dichte  $\theta(1 - b)^{\theta-1}$  auf [0, 1]. Die Verteilung des Grenzwerts in (2.10) ist gegeben durch

$$(B_1, (1-B_1)B_2, (1-B_1)(1-B_2)B_3, (1-B_1)(1-B_2)(1-B_3)B_4, \dots).$$

**Definition 2.12** (Yule<sup>5</sup>-Prozess). Ein zeitkontinuierlicher reiner Geburtsprozess (jedes Individuum erzeugt u.a. mit Rate 1 ein neues Individuum) heißt ein Yule-Prozess.

Es handelt sich also um eine zeitkontinuierliche Markovkette auf ℕ mit Sprungraten

$$q_{n,n+1} = n = -q_{n,n}, n \in \mathbb{N} \quad (q_{n,m} = 0 \text{ für } m \neq n, n+1).$$

Es ist [auch] ein Spezialfall eines zeitkont. (Galton-Watson-)Verzweigungsprozesses.

**Lemma 2.13.** Sei  $(Y_t)_{t\geq 0}$  ein Yule-Prozess (mit Geburtsrate 1 pro Individuum) und Startwert  $Y_0 = 1$ . Es gilt  $\mathscr{L}(Y_t) = \text{geom}(e^{-t})$  und  $(e^{-t}Y_t)_{t\geq 0}$  ist ein L<sup>2</sup>-beschränktes Martingal mit

$$\lim_{t\to\infty} e^{-t} Y_t \stackrel{d}{=} \operatorname{Exp}(1).$$

Beweis. Sei

 $T_i := |\{t : Y_t = i\}|$  die Länge des Zeitintervalls, in dem *i* Ind. leben.

Die Form der Raten zeigt:

$$T_1, T_2, \ldots$$
 sind u.a.,  $T_i \sim \operatorname{Exp}(i)$ .

Somit

$$\mathbb{P}(Y_t > n) = \mathbb{P}(T_1 + \dots + T_n < t) = \mathbb{P}(\max_{i=1,\dots,n} \tau_i < t) = (1 - e^{-t})^n, \quad n = 0, 1, 2, \dots$$

mit  $\tau_i$  u.i.v.  $\operatorname{Exp}(1)$ , d.h.  $\mathscr{L}(Y_t) = \operatorname{Geom}(e^{-t})$ .

(Alternativ beachte, dass die Lösung der Vorwärtsgleichung

$$\frac{\partial}{\partial t}\mathbb{P}_1(Y_t=n) = (n-1)\mathbb{P}_1(Y_t=n-1) - n\mathbb{P}_1(Y_t=n), \quad \mathbb{P}_1(Y_0=n) = \delta_{1n}$$

<sup>&</sup>lt;sup>4</sup>Nach Bob Griffiths, Steinar Engen und John William Thomas Mccloskey benannt

<sup>&</sup>lt;sup>5</sup>nach George Udny Yule, 1871–1951

gegeben ist durch  $\mathbb{P}_1(Y_t = n) = e^{-t}(1 - e^{-t})^{n-1})$ 

Zusammen mit der Verzweigungseigenschaft

$$\mathscr{L}(Y_t|Y_0 = k+j) = \mathscr{L}(Y_t|Y_0 = k) * \mathscr{L}(Y_t|Y_0 = j)$$

folgt:  $\mathbb{E}[Y_{t+h}|Y_t] = e^h Y_t$ , d.h.  $(e^{-t}Y_t)_{t\geq 0}$  ist Martingal.

Weiter folgt leicht:  $\sup_{t\geq 0} \mathbb{E}[(e^{-t}Y_t)^2] < \infty$  und  $e^{-t}Y_t \to d \operatorname{Exp}(1)$ .

**Lemma 2.14.** Seien  $G_1$  und  $G_2$  u.a.,  $G_i \sim \text{Gamma}(\theta_i)$  (d.b. Dichte  $(\Gamma(\theta_i))^{-1}g^{\theta_i-1}e^{-g}$  auf  $\mathbb{R}_+$ ). Dann ist

$$\mathscr{L}\left(G_1+G_2,\frac{G_1}{G_1+G_2}\right) = \operatorname{Gamma}(\theta_1+\theta_2) \otimes \operatorname{Beta}(\theta_1,\theta_2).$$

*Beweis.*  $G := G_1 + G_2 (G \sim \text{Gamma}(\theta_1 + \theta_2))$ . Die gemeinsame Dichte von  $(G_1, G)$  ist

$$f_{(G_1,G)}(g_1,g) = c\mathbf{1}(0 \le g_1 \le g)g_1^{\theta_1 - 1}e^{-g_1}(g - g_1)^{\theta_2 - 1}e^{-(g - g_1)} = ce^{-g}\mathbf{1}(0 \le g_1 \le g)g_1^{\theta_1 - 1}(g - g_1)^{\theta_2 - 1},$$

demnach ist die bedingte Dichte von  $G_1$ , gegeben G = g

$$f_{G_1|G=g}(g_1) = c(g)\mathbf{1}(0 \le g_1 \le g)g_1^{\theta_1-1}(g-g_1)^{\theta_2-1},$$

und somit die bedingte Dichte von  $G_1/G$ , gegeben G = g

$$f_{(G_1/G)|G=g}(b) = \tilde{c}(g)\mathbf{1}(0 \le b \le 1)b^{\theta_1 - 1}(1 - b)^{\theta_2 - 1}$$

 $Da \int_0^1 f_{(G_1/G)|G=g}(b) db = 1$  gilt, muss

$$\tilde{c}(g) = \Gamma(\theta_1 + \theta_2) / (\Gamma(\theta_1)\Gamma(\theta_2))$$

für jedes g > 0 gelten.

Beweis von Satz 2.11. Darstellung via Yule-Prozess mit Immigration:

Seien  $0 < T_1 < T_2 < \cdots$  die Sprungzeitpunkte eines Poissonprozesses auf  $[0, \infty)$  mit Rate  $\theta$ .

Der *i*-te Immigrant erscheint zum Zeitpunkt  $T_i$ , gründet *i*-te Familie, diese wächst ab dann als Yule-Prozess (vgl. Def. 2.12) unabhängig von allen anderen.

[Bild an der Tafel]

Seien

$$Z_i(t) \coloneqq \text{Größe der } i\text{-ten Familie zur Zeit } t \text{ (wir setzen } Z_i(t) = 0 \text{ für } t < T_i, Z_i(T_i) = 1\text{)},$$
  

$$S(t) \coloneqq \sum_{i=1}^{\infty} Z_i(t) \quad \text{die Gesamtgröße der Population zur Zeit } t,$$
  

$$\tau_n \coloneqq \min\{t : S(t) = n\} \quad \text{der Zeitpunkt, zu dem die Population auf } n \text{ anwächst.}$$

Es gilt

$$\left(\frac{1}{n}Z_{1}(\tau_{n}), \frac{1}{n}Z_{2}(\tau_{n}), \frac{1}{n}Z_{3}(\tau_{n}), \dots\right)_{n=1,2,\dots} \stackrel{d}{=} \left(\frac{1}{n}X_{1}(n), \frac{1}{n}X_{2}(n), \frac{1}{n}X_{3}(n), \dots\right)_{n=1,2,\dots}$$
(2.11)

Dazu Vergleich der Sprungraten: Es gebe aktuell

k Familien d. Größen  $j_1, j_2, \ldots, j_k$  mit  $j_1 + \cdots + j_k = n$ .

- S(t) springt nach n + 1 mit Rate  $n + \theta$ ,
- der Zuwachs
  - betrifft *i*-te Familie mit W'keit  $j_i/(n + \theta)$ ,
  - erzeugt neue Familie mit W'keit  $\theta/(n+\theta)$ .

Also: Skelettkette des Yule-Prozesses mit Immigration = Hoppe-Urne.

Lemma 2.13 zeigt:

$$\left(e^{-t}Z_1(t), e^{-t}Z_2(t), e^{-t}Z_3(t), \dots\right) \rightarrow \left(e^{-T_1}A_1, e^{-T_2}A_2, e^{-T_3}A_3, \dots\right)$$
 f.s., (2.12)

(koordinaten-weise) mit  $A_1, A_2, \ldots$  u.i.v. Exp(1).

Daraus folgt

$$e^{-t}S(t) \to \sum_{n=1}^{\infty} e^{-T_n} A_n$$
 f.s. (2.13)

(Zur Rechtfertigung der Grenzwertvertauschung beachte, dass  $M_i := \sup_{t\geq 0} e^{-t} Z_i(T_i + t)$  u.i.v. sind mit  $\mathbb{E}M_1 < \infty$ , also  $\limsup M_n/n = 0$  gemäß Borel-Cantelli-Lemma  $(\sum_{n=1}^{\infty} \mathbb{P}(M_i > \epsilon n) < \infty$  für jedes  $\epsilon > 0$ ).

Weiterhin gilt  $T_n/n \rightarrow \theta^{-1}$  f.s. gem. dem starken Gesetz der großen Zahlen, somit ist für  $m \ge N_0$ 

$$\sup_{t \ge 0} \sum_{n=m}^{\infty} e^{-T_n} e^{-(t-T_n)} Z_n(t) \le \sum_{n=m}^{\infty} e^{-T_n} M_n \le \sum_{n=m}^{\infty} n e^{-2n/\theta}.$$
 )

Weiterhin ist

$$\mathscr{L}\left(\sum_{n=1}^{\infty} e^{-T_n} A_n\right) = \operatorname{Gamma}(\theta), \tag{2.14}$$

denn  $\sum_i \delta_{(A_i,T_i)}$  ist ein Poissonscher Punktprozess auf  $\mathbb{R}_+ \times \mathbb{R}_+$  mit Intensitätsmaß  $\theta dt \otimes e^{-x} dx$ (und dies ist das Lévy-Maß des Gamma-Prozess/der Gamma-Verteilung, siehe z.B. Klenke [Kle20], Bsp. 16.15)

Für die Form des Intensitätsmaßes: sei  $h: (0, \infty) \to \mathbb{R}_+$ , sagen wir, stetig mit kompaktem Träger, so ist

$$\int_0^\infty \int_0^\infty h(e^{-t}a)\,\theta dt\,e^{-a}da = \int_0^\infty \int_0^a h(r)\theta \frac{dr}{r}e^{-a}da = \int_0^\infty h(r)\int_r^\infty e^{-a}\,da\,\theta \frac{dr}{r} = \int_0^\infty h(r)\theta e^{-r}\frac{dr}{r}.$$

(Die allgemeine Beobachtung dahinter ist folgende: Wenn  $\Pi = \sum \delta_{a_i}$  ein PPP auf E mit Intensitätsmaß  $\nu$  ist und  $f : E \to E'$ , dann ist  $\tilde{\Pi} = \sum \delta_{f(a_i)}$  ein PPP auf E' mit Intensitätsmaß  $\tilde{\nu} = \nu \circ f^{-1}$ .)

Das Argument für (2.14) zeigt auch

$$\mathscr{L}(G,T) = \operatorname{Gamma}(1+\theta) \otimes \operatorname{Exp}(1) \Rightarrow \mathscr{L}(e^{-T/\theta}G) = \operatorname{Gamma}(\theta),$$
 (2.15)

denn  $\sum_{n=1}^{\infty} e^{-T_n} A_n = e^{-T_1} \left( A_1 + \sum_{n=2}^{\infty} e^{-(T_n - T_1)} A_n \right).$ 

(Alternativ beachte man, dass  $e^{-T/\theta} \sim \text{Beta}(\theta, 1)$  gilt und verwende Lemma 2.14.)

Somit gilt

$$\frac{Z_1(t)}{S(t)} = \frac{e^{T_1 - t} Z_1(t)}{e^{T_1 - t} Z_1(t) + \sum_{i=2}^{\infty} e^{T_1 - t} Z_i(t)} \to \frac{A_1}{A_1 + \sum_{i=2}^{\infty} e^{-(T_i - T_1)} A_i} =: B_1 \quad \text{f.s.},$$

und  $\mathscr{L}(B_1) = \text{Beta}(1,\theta)$ , wobei  $B_1$  und  $A_1 + \sum_{i=2}^{\infty} e^{-(T_i - T_1)} A_i$  u.a. (Lemma 2.14).

Sei

$$C_n \coloneqq A_n + \sum_{i=n+1}^{\infty} e^{-(T_i - T_n)} A_i, \quad B_n \coloneqq \frac{A_n}{C_n}.$$

Zeige induktiv:

$$\mathscr{L}(C_1, B_1, B_2, \dots, B_n) = \operatorname{Gamma}(1+\theta) \otimes \operatorname{Beta}(1, \theta)^{\otimes n} \quad \text{für } n \in \mathbb{N}.$$
 (2.16)

Der Fall n = 1 stimmt nach obigem.

Für den Schluss von  $n \rightarrow n + 1$ : I.V. und Stationarität sowie Unabhängigkeit der Poisson-Zuwächse liefern

$$\mathscr{L}(C_2, B_2, B_3, \dots, B_{n+1}) = \operatorname{Gamma}(1+\theta) \otimes \operatorname{Beta}(1,\theta)^{\otimes n}$$

zudem sind  $(C_2, B_2, B_3, ..., B_{n+1})$  und  $(A_1, T_2 - T_1)$  unabhängig.

Nach Def. ist

$$C_1 = A_1 + e^{-(T_2 - T_1)}C_2, \ B_1 = \frac{A_1}{C_1} = \frac{A_1}{A_1 + e^{-(T_2 - T_1)}C_2}$$

Es ist  $e^{-(T_2-T_1)}C_2 \sim \text{Gamma}(\theta)$  nach (2.15) und  $A_1 \sim \text{Exp}(1)$  u.a. von  $e^{-(T_2-T_1)}C_2$ , d.h.

$$(C_1, B_1) \sim \text{Gamma}(\theta) \otimes \text{Beta}(1, \theta)$$

nach Lemma 2.14, dies liefert den Induktionssschluss.

Schließlich gilt

$$\frac{e^{-t}Z_{n}(t)}{e^{-t}S(t)} \rightarrow \frac{e^{-T_{n}}A_{n}}{\sum_{i=1}^{\infty}e^{-T_{i}}A_{i}} = \frac{\sum_{i=2}^{\infty}e^{-T_{i}}A_{i}}{\sum_{i=1}^{\infty}e^{-T_{i}}A_{i}} \times \dots \times \frac{\sum_{i=n}^{\infty}e^{-T_{i}}A_{i}}{\sum_{i=n-1}^{\infty}e^{-T_{i}}A_{i}} \times \frac{e^{-T_{n}}A_{n}}{\sum_{i=n}^{\infty}e^{-T_{i}}A_{i}} \\
= (1 - B_{1}) \times \dots \times (1 - B_{n}) \times \frac{A_{n}}{A_{n} + \sum_{i=n+1}^{\infty}e^{-(T_{i} - T_{n})}A_{i}} \\
= (1 - B_{1}) \times \dots \times (1 - B_{n-1})B_{n}.$$

**Beobachtung 2.15** (Poisson-Dirichlet-Verteilung). Seien  $1 \ge V_1 > V_2 > \cdots$  die (der Größe) nach sortierten Einträge (=Typenhäufigkeiten) aus GEM-verteiltem

$$(B_1, (1-B_1)B_2, (1-B_1)(1-B_2)B_3, (1-B_1)(1-B_2)(1-B_3)B_4, \dots).$$

Sei  $\Pi = \sum_i \delta_{X_i}$  PPP auf  $\mathbb{R}_+$  mit Intensitätsmaß  $(\theta/x)e^{-x}dx$  ( $\Pi$  beschreibt die Sprünge eines Standard-Gamma-Subordinators bis zur Zeit  $\theta$ ) und  $S := \sum X_i$ , seien  $X_{[1]} > X_{[2]} > \cdots$  die Ordnungsstatistik der  $X_i$ s. Dann ist

$$(X_{[1]}/S, X_{[2]}/S, \dots) \stackrel{d}{=} (V_1, V_2, \dots).$$

Dies folgt aus dem Beweis von Satz 2.11. Diese Verteilung (ein W'maß auf  $\{(x_1, x_2, ...) \in [0, 1]^{\mathbb{N}} : x_1 + x_2 + \cdots = 1\}$ ), heißt die Poisson-Dirichlet-Verteilung (mit Parameter  $\theta$ ).

**Bericht 2.16** (Endlich viele Typen mit elternunabhängiger Mutation). Betrachte Mutationsmodell mit d neutralen Typen (Typenmenge  $E = \{1, \ldots, d\}$ ), jede Linie mutiert mit Rate  $\theta/2$ , Typ nach Mutation ist j mit W'keit  $\pi_j$  (> 0) (( $\pi_1, \ldots, \pi_d$ ) Ws-Gewichte auf  $\{1, \ldots, d\}$ ), u.a. vom vorigen Typ.

Die Typenverteilung in einer unendlichen Population im Gleichgewicht ist dann

$$\mathscr{L}(Z_1(\infty),\ldots,Z_d(\infty)) = \text{Dirichlet}(\theta\pi_1,\ldots,\theta\pi_d),$$
 (2.17)

d.h. die Dichte ist

$$\frac{\Gamma(\theta)}{\Gamma(\theta\pi_1)\cdots\Gamma(\theta\pi_d)} x_1^{\theta\pi_1-1} \cdots x_k^{\theta\pi_d-1}$$

bezüglich dem Lebesgue-Maß auf  $\{(x_1, \ldots, x_d) : 0 \le x_i \le 1, x_1 + \cdots + x_d = 1\}$ .

Man kann den allgemeinen Fall aus Beob. 2.15 herleiten: Wenn man jeden Sprung des Gamma-Subordinators unabhängig gemäß  $\pi$  mit einer "Farbe" aus  $\{1, \ldots, d\}$  einfärbt, so bilden die Sprünge jeder Farbe für sich jeweils unabhängige Gamma-Subordinatoren (mit entsprechend verkleinerter Intensität), somit:  $Y_i \sim \text{Gamma}(\theta \pi_i)$  und  $Y_1, \ldots, Y_d$  unabhängig, so ist

$$\left(Z_1(\infty),\ldots,Z_d(\infty)\right) \stackrel{d}{=} \left(\frac{Y_1}{Y_1+\cdots+Y_k},\ldots,\frac{Y_k}{Y_1+\cdots+Y_k}\right)$$

und die rechte Seite ist Dirichlet ( $\theta \pi_1, \ldots, \theta \pi_d$ )-verteilt (dies ist eine multivariate Verallgemeinerung von Lemma 2.14, siehe z.B. Ch. 40.5 in Norman L. Johnson, Samuel Kotz, *Distributions in statistics: continuous multivariate distributions*, Wiley, 1972).

## 2.2 Infinitely-many-sites-Modell (IMS)

**Definition 2.17** (Infinitely-many-sites-Modell<sup>6</sup>). Man nimmt an, dass jede Mutation eine neue, bisher noch nie mutierte Position am betrachteten Lokus betrifft.

Mathematisch realisiert man dies z.B. folgendermaßen: Die betrachtete Stelle im Genom (eine gewisse Abfolge von Nukleotiden im DNS-Doppelstrang eines Chromosoms) entspricht [0, 1], jede Mutation erhält eine neue, uniform aus [0, 1] gewählte "Position", der Typ eines Individuums ist ein (einfaches) Zählmaß auf [0, 1] (bzw. alternativ eine Teilmenge von [0, 1]), der Typ eines Individuums gibt an, wo dieses relativ zu einen "Referenztyp" (oder "Wildtyp") mutiert ist.

Das infinitely-many-sites-Modell (IMS-Modell, in der Literatur auch infinite-sites-Modell genannt) ist für viele praktische Zwecke eine angemessene Approximation für die Beschreibung von Mutationen auf dem Niveau der DNS-Sequenz : Wenn die Mutationsrate pro Basenpaar sehr klein und die betrachtete Stelle im Genom (der sog. Lokus) nicht "zu lang" ist, ist es plausibel, die Möglichkeit der Mehrfachmutation einer Stelle (und andere Effekte, die im IMS-Modell nicht berücksichtigt werden, etwa Rekombination oder Insertionen/Deletionen längerer Stücke im Genom) zu vernachlässigen. (Man denke an eine Abfolge von  $L \gg 1$  Basenpaaren, bei Mutation wird eine der zufällig gewählte der L Positionen zufällig modifiziert.)

<sup>&</sup>lt;sup>6</sup>Eingeführt in G. A. Watterson, On the number of segregating sites in genetical models without recombination, Theoretical Population Biology 7 (2), 256–276, (1975).
**Beispiel 2.18.** John Parsch, Colin D. Meiklejohn, and Daniel L. Hartl, Patterns of DNA sequence variation suggest the recent action of positive selection in the janus-ocnus region of drosophila simulans, *Genetics* 159:647–657, (2001) berichten (u.a.) genetische Variabilität in einem ca. 1.700 Basenpaare langen Stück des Chromosoms 3 in einer (weltweiten) Stichprobe von 8 *Drosophila simulans* und einer Stichprobe von *Drosophila melanogaster*, zwei verwandten Arten von Taufliegen. An insgesamt 31 Stellen sind Unterschiede zwischen den Individuen sichtbar (siehe Tabelle 2.1, es sind nur die sogenannten variablen oder segregierenden Positionen aufgeführt).

Die Sequenzinformation von *Drosophila melanogaster* — diese Stichprobe bildet bezüglich der 8 Stichproben von *Drosophila simulans* eine "outgroup" — gestattet (zusammen mit den IMS-Modellannahmen) an jeder Position zu entscheiden, welche Base die anzestrale und welche die mutierte ist.

Position I I 2 2 3 5 5 5 6 6 6 6 6 6 7 8 I 2 3 3 3 4 6 6 6 6 I 9 9 6 4 8 4 5 5 5 5 6 6 o 4 3 9 0 2 8 3 3 5 I 7 4 2 8 9 5 4 2 7 8 7 7 29 I 0 2 8 2 I 5 3 gctcgataagccga ST с с а a ta t а a а t с а • . . g **S2** . . . . . . с а t g с с с g g g g a t c t a t c c t c t g t t \$3 g . S٢ g s6 **S**7 g **s**8 g тı . . . c a t g c c c a g . . . . t . . c c . . t g . t g c a

Tabelle 2.1: Beobachtete genetische Variabilität in einer Region in Chromosom 3 aus einer Stichprobe von 8 *Drosophila simulans* (Zeilen s1–s8) und einer Stichprobe von *Drosophila melanogaster* (Zeile m1) aus Parsch et al, *Genetics* 159:647–657, (2001). Siehe Figure 2 dort, wir betrachten hier nur den Teil der Sequenz, der die Gene *janA* und *janB* umfasst.

Zur Beschreibung von beobachteten Sequenzdaten in einer Stichprobe der Größe n im Kontext des IMS-Modells betrachten wir folgende Modellvorstellung: Die n-Stichprobe entsteht aus einem n-Koaleszent, längs dessen Ästen sich mit Rate  $\frac{\theta}{2}$  Mutationen ereignen (und jede trifft eine völlig neue Position).

Die Anzahl segregierender Stellen ist

 $S_n$  = # verschiedene Mutationen, die in *n*-Stichprobe vorkommen

(im Sinne von: Positionen, an denen sich mindestens zwei Stichproben unterscheiden). Wenn  $S_n = s$ , so entsprechen die Beobachtungen einer  $n \times s$ -Datenmatrix  $(D_{ik})_{i=1,...,n;k=1,...,s}$ 

 $D_{ik} = 1$ (Stichprobe *i* ist an *k*-ter segregierender Stelle mutiert).

Beispielsweise sehen wir in Abbildung 2.1 eine Realisierung mit  $S_4 = 3$ .



Abbildung 2.1: Ein 4-Koaleszent, längs dessen Kanten Mutationen gemäß IMS-Modell auftreten. Wir registrieren für jede Mutation die mutierte Position (in [0, 1]) und an den Blättern (den Stichproben) sämtliche Mutationen, die sich auf dem kürzesten Weg vom jeweiligen Blatt zur Wurzel befinden.

**Bemerkung 2.19** (Unbekannter Wildtyp). Wenn man "nur" die Stichprobe sieht und keine externen Zusatzinformationen (z.B. eine "outgroup" durch inter-Spezies-Vergleich wie in Bsp. 2.18) besitzt, kann man an den segregierenden Stellen nicht entscheiden, welcher Typ der Wildtyp und welcher die Mutante ist (im Genetik-Jargon: die Mutationen sind "unpolarisiert"). In dieser Situation ist obige Datenmatrix nur bis auf "Umklappen" von Spalten definiert, d.h. die eigentliche Information ist  $S_n$ und

 $\Delta_{i,j}(k) = \begin{cases} 1, & \text{Stichproben } i \text{ und } j \text{ an } k \text{-ter segregierender Stelle verschieden,} \\ 0, & \text{sonst} \end{cases}$ 

für  $k = 1, ..., S_n$ .

### Beobachtung 2.20. Es gilt

$$\mathbb{E}_{\theta}[S_n] = \theta h_n, \quad \operatorname{Var}_{\theta}[S_n] = \theta h_n + \theta^2 \sum_{i=1}^{n-1} \frac{1}{i^2}$$

mit  $h_n = \sum_{i=1}^{n-1} \frac{1}{i}$ .

$$\widehat{\theta}_W \coloneqq \frac{S_n}{h_n}$$

ist ein erwartungstreuer Schätzer für  $\theta$  (der sogenannte Watterson-Schätzer) mit

$$\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right] \sim \frac{\theta}{\log n} \quad \text{und} \quad \frac{\widehat{\theta}_{W} - \theta}{\sqrt{\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]}} \xrightarrow{d} \mathcal{N}(0, 1) \qquad \text{für } n \to \infty.$$

Beweis. Schreibe

$$S_n = S_{n,n} + S_{n,n-1} + \dots + S_{n,2}$$

mit

$$S_{n,j}$$
 = # Mutationen, während die Genealogie aus j Linien besteht.

Wir hatten in der Diskussion in Kapitel 1.5 gesehen, dass  $S_{n,j} \sim \text{geom}(\frac{j-1}{\theta+j-1})$  (denn gegeben  $T_j$ , die Zeit, während der j Linien in der Genealogie existieren, ist  $S_{n,j} \sim \text{Pois}(\frac{\theta}{2}jT_j)$ ) und  $S_{n,n}, \ldots, S_{n,2}$  sind unabhängig. Somit

$$\mathbb{E}_{\theta}[S_{n,j}] = \frac{\theta + j - 1}{j - 1} - 1 = \theta / (j - 1), \quad \operatorname{Var}_{\theta}[S_{n,j}] = \left(\frac{\theta}{\theta + j - 1}\right) / \left(\frac{j - 1}{\theta + j - 1}\right)^2 = \frac{\theta(\theta + j - 1)}{(j - 1)^2} = \frac{\theta}{j - 1} + \frac{\theta^2}{(j - 1)^2}$$

was die Formeln für Erwartungswert und Varianz von  $S_n$  beweist.

Zur asymptotischen Normalität von  $\widehat{\theta}_W$ : Schreibe

$$X_{n,j} \coloneqq \frac{S_{n,j} - \theta/(j-1)}{\sqrt{\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]}}, \quad j = 2, 3, \dots, n$$

Die  $X_{n,j}$  bilden ein unabhängiges, zentriertes und normiertes Dreiecksschema (für festes  $n \operatorname{sind} X_{n,2}, \ldots, X_{n,n}$ unabhängig mit  $\mathbb{E}_{\theta}[X_{n,j}] = 0$  und  $\sum_{j=2}^{n} \operatorname{Var}_{\theta}[X_{n,j}] = 1$ ), für  $\varepsilon > 0$  und  $n \operatorname{so} \operatorname{groß}$ , dass  $\varepsilon \operatorname{Var}_{\theta}[\widehat{\theta}_{W}] > \theta^{2}$  gilt, ist

$$\mathbb{E}\left[X_{n,j}^{2}\mathbf{1}\left(X_{n,j}^{2} > \varepsilon\right)\right] = \frac{1}{\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]} \mathbb{E}\left[\left(S_{n,j} - \theta/(j-1)\right)^{2}\mathbf{1}\left(S_{n,j} > \theta/(j-1) + \sqrt{\varepsilon}\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]\right)\right] \\ \leq \frac{1}{\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]} \mathbb{E}\left[S_{n,j}^{2}\mathbf{1}\left(S_{n,j} > \theta/(j-1) + \sqrt{\varepsilon}\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]\right)\right].$$

Demnach erfüllt das Schema die Lindeberg-Bedingung

$$\lim_{n \to \infty} \mathbb{E} \Big[ X_{n,j}^2 \mathbf{1} \big( X_{n,j}^2 > \varepsilon \big) \Big] = 0$$
(2.18)

und daher ist das renormierte  $\widehat{\theta}_W$  asymptotisch normalverteilt (siehe z.B. [Kle20, Satz 15.43]).

(beachte: Für  $X \sim \text{geom}(p)$  gilt

$$\mathbb{E}\left[X^{2}\mathbf{1}(X \ge m)\right] = \sum_{k=m}^{\infty} k^{2}p(1-p)^{k-1} \le \sum_{k=m}^{\infty} (k+1)kp(1-p)^{k-1} = p\sum_{k=m}^{\infty} \left[\frac{d^{2}}{dy^{2}}(1-y)^{k+1}\right]_{y=p}$$
$$= p\left[\frac{d^{2}}{dy^{2}}\sum_{k=m}^{\infty} (1-y)^{k+1}\right]_{y=p} = p\left[\frac{d^{2}}{dy^{2}}\frac{(1-y)^{m+1}}{y}\right]_{y=p}$$
$$= p\frac{(1-p)^{m-1}}{p}\left((m+1)m+2(m+1)\frac{1-p}{p}+2\frac{(1-p)^{2}}{p^{2}}\right)$$
$$\le \left((m+1)(m+2)+2\right)\frac{(1-p)^{m-1}}{p^{2}},$$

demnach für  $X = S_{n,j}$  mit  $p = p_j = \frac{j-1}{\theta+j-1}$  und z.B.  $m = \sqrt{\frac{1}{2}\varepsilon\theta\log n}$  ist

$$\frac{1}{\operatorname{Var}_{\theta}\left[\widehat{\theta}_{W}\right]} \mathbb{E}\left[S_{n,j}^{2} \mathbf{1}\left(S_{n,j} > \sqrt{\frac{1}{2}\varepsilon\theta \log n}\right)\right] \leq C_{\theta}\left(\frac{\theta}{\theta+j-1}\right) \sqrt{\frac{1}{2}\varepsilon\theta \log n}$$

für ein  $C_{\theta} < \infty$ , d.h (2.18) gilt.)

**Bemerkung 2.21** (Alternativer Zugang zu Beob. 2.20). Wir könnten Erwartungswert und Varianz von  $S_n$  auch folgendermaßen berechnen: Gegeben die Gesamtlänge  $L_{\text{ges}}$  des Koaleszenten ist  $S_n$  Poi $((\theta/2)L_{\text{ges}})$ -verteilt.

$$L_{\rm ges} \stackrel{d}{=} \sum_{j=2}^{n} jT_j,$$

mit  $T_n, T_{n-1}, \ldots, T_2$  u.a.,  $\mathscr{L}(T_j) = \operatorname{Exp}(\binom{j}{2})$ , somit

$$\mathbb{E}_{\theta}[S_n] = \frac{\theta}{2} \sum_{j=2}^n j \binom{j}{2} = \theta \sum_{i=1}^{n-1} \frac{1}{i},$$

$$\operatorname{Var}_{\theta}[S_n] = \mathbb{E}_{\theta}[\operatorname{Var}_{\theta}[S_n | L_{\operatorname{res}}]] + \operatorname{Var}_{\theta}[\mathbb{E}_{\theta}[S_n | L_{\operatorname{res}}]]$$
(2.19)

$$\begin{aligned} \operatorname{tr}_{\theta}[S_{n}] &= \mathbb{E}_{\theta}[\operatorname{Var}_{\theta}[S_{n} | L_{\operatorname{ges}}]] + \operatorname{Var}_{\theta}[\mathbb{E}_{\theta}[S_{n} | L_{\operatorname{ges}}]] \\ &= \mathbb{E}_{\theta}\left[\frac{\theta}{2}L_{\operatorname{ges}}\right] + \operatorname{Var}_{\theta}\left[\frac{\theta}{2}L_{\operatorname{ges}}\right] = \theta \sum_{i=1}^{n-1} \frac{1}{i} + \theta^{2} \sum_{i=1}^{n-1} \frac{1}{i^{2}}. \end{aligned}$$
(2.20)

Die Konvergenzordnung  $O(1/\log(n))$  der Varianz des Watterson-Schätzers ist zwar von Standpunkt der Statistik gesehen "frustrierend langsam" (zumal im Vergleich zur klassischen Situation, in der man einen Parameter basierend auf *n unabhängigen* Beobachtungen schätzt, dort hat man typischerweise Abfall der Varianz O(1/n) für plausible Schätzer). Andererseits haben Y. X. Fu and W. H. Li, Maximum Likelihood Estimation of Population Parameters, *Genetics* 134 (4), 1261–1270, (1993) gezeigt, dass es (zumindest asymptotisch) auch nicht besser möglich ist:

**Satz 2.22.** Jeder erwartungstreue Schätzer für  $\theta$  im IMS-Modell

hat unter 
$$\mathbb{P}_{\theta}$$
 mindestens Varianz  $\theta / \sum_{k=1}^{n-1} \frac{1}{\theta + k} \quad (\sim \theta / \log n \, f \ddot{u} r \, n \to \infty)$ 

*Beweis.* Nehmen wir (zunächst) an, wir könnten  $S_{n,2} = s_{n,2}, S_{n,3} = s_{n,3}, \dots, S_{n,n} = s_{n,n}$  beobachten (was anhand von Sequenzdaten an den Blättern des Koaleszenten nicht immer möglich ist):

Die Likelihoodfunktion (die Verteilungsgewichte von  $(S_{n,n}, \ldots, S_{n,2})$ , aufgefasst als Funktion des Parameters  $\theta$ ) ist

$$L_n(s_{n,2},...,s_{n,n};\theta) = \prod_{j=2}^n \frac{j-1}{\theta+j-1} \left(\frac{\theta}{\theta+j-1}\right)^{s_{n,j}}$$
  
=  $(n-1)!\theta^{s_n} \prod_{j=2}^n (\theta+j-1)^{-(s_{n,j}+1)}$ 

mit  $s_n = s_{n,2} + \dots + s_{n,n}$ , also

$$\frac{\partial}{\partial \theta} \log L_n(s_{n,2},\ldots,s_{n,n};\theta) = \frac{s_n}{\theta} - \sum_{j=2}^n \frac{s_{n,j}+1}{\theta+j-1}$$

d.h.  $\widehat{\theta}_{ML,hyp}$ , der Maximum-Likelihood-Schätzer für  $\theta$  basierend auf  $(S_{n,n}, \ldots, S_{n,2})$ , ist die Lösung (in  $\theta$ ) von  $s_n = \theta \sum_{j=2}^n \frac{s_{n,j+1}}{\theta+j-1}$ .

Weiter ist

$$\frac{\partial^2}{\partial \theta^2} \log L_n(s_{n,2},\ldots,s_{n,n};\theta) = -\frac{s_n}{\theta^2} + \sum_{j=2}^n \frac{s_{n,j}+1}{(\theta+j-1)^2},$$

die Fisher-Information ist somit

$$I(\theta) = -\mathbb{E}_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log L_n(s_{n,2}, \dots, s_{n,n}; \theta) \right]$$
  
$$= \mathbb{E}_{\theta} \left[ \frac{S_n}{\theta^2} \right] - \sum_{j=2}^n \mathbb{E}_{\theta} \left[ \frac{S_{n,j} + 1}{(\theta + j - 1)^2} \right] = \frac{\sum_{k=1}^{n-1} 1/k}{\theta} - \sum_{j=2}^n \frac{\theta + j - 1}{(j - 1)(\theta + j - 1)^2}$$
  
$$= \frac{\sum_{k=1}^{n-1} 1/k}{\theta} - \sum_{j=2}^n \frac{1}{(j - 1)(\theta + j - 1)} = \frac{1}{\theta} \sum_{k=1}^{n-1} \left( \frac{1}{k} - \frac{\theta}{k(\theta + k)} \right) = \frac{1}{\theta} \sum_{k=1}^{n-1} \frac{1}{\theta + k}$$

Gemäß der Cramér-Rao-Ungleichung (siehe z.B. John A. Rice, *Mathematical statistics and data analysis*, Duxbury Press, 1995, Ch. 8.6 oder die knappe Diskussion unten) gilt für jeden erwartungstreuen Schätzer

$$T = T(S_{n,n},\ldots,S_{n,2})$$

für  $\theta$  (d.h. der Schätzer wird durch eine Funktion  $T : \mathbb{N}_0^{n-1} \to (0, \infty)$  mit  $\mathbb{E}_{\theta} [T(S_{n,n}, \dots, S_{n,2})] = \theta$  für alle  $\theta > 0$  dargestellt)

$$\operatorname{Var}_{\theta}[T(S_{n,n},\ldots,S_{n,2})] \geq \frac{1}{I(\theta)} = \frac{\theta}{\sum_{k=1}^{n-1} \frac{1}{\theta+k}}$$

d.h. die Behauptung.

Nachträge:

1. Heuristik zur Cramér-Rao-Ungleichung:

Betrachten wir die allgemeine Situation, dass die Beobachtungen X ein Zufallsvektor (mit Werten in einer geeigneten Teilmenge von  $\mathbb{R}^d$  für ein d) sind, die Dichte-/Likelihood-Funktion  $f(x; \theta)$  (im diskreten Fall: die Gewichte) sei genügend glatt, so dass die folgenden Vertauschungen von Ableitung und Integral gerechtfertigt sind.

Sei  $V(X) := \frac{\partial}{\partial \theta} \log f(X; \theta) = \frac{1}{f(X; \theta)} \frac{\partial}{\partial \theta} f(X; \theta)$  die "Score-Funktion", es ist  $\mathbb{E}_{\theta}[V(X)] = 0$  stets, denn

$$\int f(x;\theta) \frac{1}{f(x;\theta)} \frac{\partial}{\partial \theta} f(x;\theta) \, dx = \int \frac{\partial}{\partial \theta} f(x;\theta) \, dx = \frac{\partial}{\partial \theta} \int f(x;\theta) \, dx = \frac{\partial}{\partial \theta} 1 = 0.$$
$$I(\theta) \coloneqq \operatorname{Var}_{\theta} [V(X)] = \mathbb{E}_{\theta} [V(X)^2]$$

heißt die Fisher-Information, beachte auch

$$I(\theta) = -\mathbb{E}_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log f(X; \theta) \right]$$

(d.h. wir können die Fisher-Information als (erwartete) Krümmung der Likelihood-Funktion an den beobachteten Daten interpretieren), denn

$$\frac{\partial^2}{\partial\theta^2}\log f(x;\theta) = \frac{\frac{\partial^2}{\partial\theta^2}f(x;\theta)}{f(x;\theta)} - \left(\frac{\frac{\partial}{\partial\theta}f(x;\theta)}{f(x;\theta)}\right)^2 = \frac{\frac{\partial^2}{\partial\theta^2}f(x;\theta)}{f(x;\theta)} - \left(\frac{\partial}{\partial\theta}\log f(x;\theta)\right)^2$$

und

$$\mathbb{E}_{\theta}\left[\frac{\frac{\partial^2}{\partial\theta^2}f(X;\theta)}{f(X;\theta)}\right] = \int \frac{\frac{\partial^2}{\partial\theta^2}f(x;\theta)}{f(x;\theta)}f(x;\theta)\,dx = \int \frac{\partial^2}{\partial\theta^2}f(x;\theta)\,dx = \frac{\partial^2}{\partial\theta^2}\int f(x;\theta)\,dx = \frac{\partial^2}{\partial\theta^2}1 = 0$$

Sei nun T(X) irgendein Schätzer für  $\theta$  (d.h. eine Funktion der Beobachtungen X mit Werten in  $[0, \infty)$  (und  $\mathbb{E}_{\theta}[(T(X))^2] < \infty$ ), dann ist

$$Cov_{\theta}[T(X), V(X)] = \mathbb{E}_{\theta}[T(X)V(X)] = \int T(x)\frac{1}{f(x;\theta)}\frac{\partial}{\partial\theta}f(x;\theta)f(x;\theta)dx$$
$$= \int T(x)\frac{\partial}{\partial\theta}f(x;\theta)dx = \frac{\partial}{\partial\theta}\mathbb{E}_{\theta}[T(X)].$$

Die Cauchy-Schwarz-Ungleichung liefert

$$\sqrt{\operatorname{Var}_{\theta}[T(X)]\operatorname{Var}_{\theta}[V(X)]} \ge |\operatorname{Cov}_{\theta}[T(X), V(X)]| = \left|\frac{\partial}{\partial \theta} \mathbb{E}_{\theta}[T(X)]\right|$$

und somit gilt

$$\operatorname{Var}_{\theta}[T(X)] \geq \frac{\left|\frac{\partial}{\partial \theta} \mathbb{E}_{\theta}[T(X)]\right|^{2}}{\operatorname{Var}_{\theta}[V(X)]} = \frac{\left|\frac{\partial}{\partial \theta} \mathbb{E}_{\theta}[T(X)]\right|^{2}}{I(\theta)}$$

falls T(X) ein unverzerrter Schätzer für  $\theta$  ist, d.h.  $\mathbb{E}_{\theta}[T(X)] = \theta$  für alle  $\theta$ , so ist natürlich  $\frac{\partial}{\partial \theta} \mathbb{E}_{\theta}[T(X)] = 1$  (und dies ist die Form der Cramér-Rao-Ungleichung, die wir oben verwendet haben).

2. Wir hatten die Schranke unter der Annahme hergeleitet, dass wir  $(S_{n,n}, \ldots, S_{n,2})$  tatsächlich beobachten könnten, was anhand der Daten i.A. nicht möglich ist. Intuitiv erscheint es zumindest sehr plausibel, dass jeder "reale" erwartungstreue Schätzer für  $\theta$  (d.h. jedes  $\widetilde{T} = \widetilde{T}(D)$ , das eine Funktion der Datenmatrix D ist, das also weniger Informationen verwenden darf, als wir oben angenommen hatten), ebenfalls mindestens Varianz  $1/I(\theta)$  hat.

Diese Intuition kann man durch den Begriff der (statistischen) Suffizienz folgendermaßen formalisieren: Sei X die "volle" Information, die die Genealogie der *n*-Stichprobe und die darauf vorkommenden Mutationen beschreibt, d.h. X enthält die "topologische" Information, in welcher Reihenfolge die Verschmelzungen der Linien stattfinden, und für jede Kante im Baum die Information, welche Mutationen auf dieser liegen (wir verzichten hier darauf, dies in Formeln zu fassen). Offenbar kann man aus X die Datenmatrix D ablesen und somit kann  $\widetilde{T} = \widetilde{T}(D(X))$  als eine Funktion von X interpretiert werden.

Die entscheidende Beobachtung ist, dass  $Y \coloneqq (S_{n,2}, S_{n,3}, \dots, S_{n,n})$  suffizient für  $\theta$  ist, d.h. die bedingte Verteilung  $\mathscr{L}_{\theta}(X | Y)$  hängt nicht von  $\theta$  ab — gegeben  $S_{n,2} = s_{n,2}, \dots, S_{n,n} = s_{n,2}$ entsteht X, indem man für  $j = n, n-1, \dots, 2$  auf den j Kanten in "Niveau" j des Koaleszenten  $s_{n,j}$  Mutationen uniform verteilt und unter allen aktuell möglichen Verschmelzungen uniform eine auswählt; demnach enthalten die Gewichte von  $\mathscr{L}_{\theta}(X | Y)$  nur kombinatorische Terme, aber keine  $\theta$ -Abhängigkeit.

Nun ist (beachte, dass wir den bedingten Erwartungswert bilden können, ohne  $\theta$  zu kennen)

$$\widehat{T} \coloneqq \widehat{T}(Y) \coloneqq \mathbb{E}[\widetilde{T} \mid Y]$$

ebenfalls ein erwartungstreuer Schätzer für  $\theta$  und nach Konstruktion ist  $\widehat{T}$  eine gewisse Funktion von  $Y = (S_{n,2}, S_{n,3}, \dots, S_{n,n})$ , d.h. nach obigem ist  $\operatorname{Var}_{\theta}[\widehat{T}] \ge 1/I(\theta)$  und folglich auch

$$\operatorname{Var}_{\theta}[\widetilde{T}] = \mathbb{E}_{\theta}[\underbrace{\operatorname{Var}_{\theta}[\widetilde{T} \mid Y]}_{\geq 0}] + \operatorname{Var}_{\theta}[\underbrace{\mathbb{E}_{\theta}[\widetilde{T} \mid Y]}_{=\widehat{T}}] \geq \operatorname{Var}_{\theta}[\widehat{T}] \geq \frac{1}{I(\theta)}.$$

Definition 2.23 (Frequenzspektrum). Sei

 $\xi_i^{(n)} := #$  Mutationen, die in genau *i* der *n* Stichproben vorkommen, i = 1, ..., n - 1.

(Wir nehmen dabei an, dass an jeder Position der anzestrale oder "Wildtyp" bekannt ist, z.B. durch Interspezies-Vergleich.) Der Vektor

$$\xi^{(n)} = \left(\xi_1^{(n)}, \xi_2^{(n)}, \dots, \xi_{n-1}^{(n)}\right)$$

heißt das Frequenzspektrum (der segregrierenden Stellen).

Wenn der anzestrale Typ nicht bekannt ist, betrachtet man stattdessen das gefaltete Frequenzspektrum  $(\eta_1^{(n)}, \eta_2^{(n)}, \dots, \eta_{\lfloor n/2 \rfloor}^{(n)})$  mit

$$\eta_i^{(n)} := \xi_i^{(n)} + \xi_{n-i}^{(n)} \mathbf{1}_{i \neq n/2}, \quad 1 \le i \le \lfloor n/2 \rfloor.$$

**Satz 2.24.** *Es gilt* 

$$\mathbb{E}_{\theta}\left[\xi_{i}^{(n)}\right] = \frac{\theta}{i}, \quad \operatorname{Cov}_{\theta}\left[\xi_{i}^{(n)}, \xi_{j}^{(n)}\right] = \mathbf{1}_{i=j}\frac{\theta}{i} + \theta^{2}\sigma_{ij}, \quad 1 \le i \le j \le n$$

*mit*  $h_n \coloneqq \sum_{i=1}^{n-1} \frac{1}{i}, \beta_n(i) \coloneqq \frac{2n}{(n-i+1)(n-i)} (h_{n+1} - h_i) - \frac{2}{n-i}$ 

$$\sigma_{ii} = \begin{cases} \beta_n(i+1), & i < \frac{n}{2}, \\ 2\frac{h_n - h_i}{n-i} - \frac{1}{i^2}, & i = \frac{n}{2}, \\ \beta_n(i) - \frac{1}{i^2}, & i > \frac{n}{2}, \end{cases} \quad f \ddot{u} r \, i > j \, ist \quad \sigma_{ij} = \begin{cases} \frac{\beta_n(i+1) - \beta_n(i)}{2}, & i + j < n, \\ \frac{h_n - h_i}{n-i} + \frac{h_n - h_j}{n-j} & -\frac{\beta_n(i) + \beta_n(j+1)}{2} - \frac{1}{ij}, & i + j = n, \\ \frac{\beta_n(j) - \beta_n(j+1)}{2} - \frac{1}{ij}, & i + j > n \end{cases}$$

(und  $\sigma_{ij} = \sigma_{ji}$ ).

Die Diagonaleinträge  $\sigma_{i,i}$  (d.h.  $\operatorname{Var}_{\theta}[\xi_i^{(n)}]$ ) dominieren die Kovarianzmatrix  $(\sigma_{i,j})$ : Abb. 2.2 zeigt die Diagonaleinträge  $\sigma_{i,i}$  und die Antidiagonaleinträge  $\sigma_{i,n-i}$  für n = 25 und  $\theta = 1$ , Abb. 2.3 zeigt eine dreidimensionale Darstellung von  $(\sigma_{i,j})$  für n = 25 und  $\theta = 1$ , Abb. 2.3 zeigt  $(-\sigma_{i,j})$ , wobei der besseren Sichtbarkeit wegen die (auch betragsmäßig) deutlich größeren Diagonal- und Antidiagonal- einträge auf 0 gesetzt wurden.

*Beweis (der Formel für den Erwartungswert).* Wir denken uns die Kanten des *n*-Koaleszenten auf jedem Niveau (u.a. zufällig) nummeriert.



Abbildung 2.2: Diagonaleinträge  $\sigma_{i,i}$  und Antidiagonaleinträge  $\sigma_{i,n-i}$  der Kovarianzmatrix von  $\xi^{(n)}$  für  $n = 25, \theta = 1$ 



Abbildung 2.3: Kovarianzmatrix von  $\xi^{(n)}$  für  $n = 25, \theta = 1$ 

Sei

 $u_{k,\ell} \coloneqq \#$  Mutationen auf  $\ell$ -ter Kante auf Niveau k,  $J_{k,\ell} \coloneqq \#$  Blätter oberhalb  $\ell$ -ter Kante auf Niveau k,

damit ist

$$\xi_i^{(n)} = \sum_{k=2}^{n-i+1} \sum_{\ell=1}^k \nu_{k,\ell} \mathbf{1}(J_{k,\ell} = i).$$
(2.21)

Somit gilt

$$\begin{split} \mathbb{E}_{\theta}\left[\xi_{i}^{(n)}\right] &= \sum_{k=2}^{n-i+1} \sum_{\ell=1}^{k} \mathbb{E}_{\theta}\left[\nu_{k,\ell} \mathbf{1}(J_{k,\ell}=i)\right] = \sum_{k=2}^{n-i+1} \sum_{\ell=1}^{k} \mathbb{E}_{\theta}\left[\nu_{k,\ell}\right] \mathbb{P}_{\theta}(J_{k,\ell}=i) \\ &= \sum_{k=2}^{n-i+1} \sum_{\ell=1}^{k} \frac{\theta}{k(k-1)} \frac{\binom{n-i-1}{k-2}}{\binom{n-1}{k-1}} = \sum_{k=2}^{n-i+1} k \frac{\theta}{k(k-1)} \frac{\binom{n-i-1}{k-2}}{\binom{n-1}{k-1}} \\ &= \theta \sum_{k=2}^{n-i+1} \frac{1}{k-1} \frac{(n-i-1)!}{(k-2)!(n-i-k+1)!} \frac{(k-1)!(n-k)!}{(n-1)!} \times \frac{i!}{i(i-1)!} \\ &= \frac{\theta}{i} \frac{1}{\binom{n-1}{i}} \sum_{k=2}^{n-i+1} \binom{n-k}{i-1} = \frac{\theta}{i} \end{split}$$



Abbildung 2.4: (-1)×Kovarianzmatrix von  $\xi^{(n)}$  für n = 25,  $\theta = 1$ , wobei Diagonal- und Antidiagonaleinträge auf 0 gesetzt wurden

Wir verwenden hierbei in der ersten Zeile, dass

$$\mathbb{E}_{\theta} \big[ \nu_{k,\ell} \big] = \mathbb{E}_{\theta} \big[ \mathbb{E}_{\theta} \big[ \nu_{k,\ell} \big| T_k \big] \big] = \mathbb{E}_{\theta} \big[ \frac{\theta}{2} T_k \big] = \frac{\theta}{2} \frac{2}{k(k-1)} = \frac{\theta}{k(k-1)}$$

gilt und dass  $\nu_{k,\ell}$  (das nur vom Poissonprozess der Mutationen abhängt) und  $J_{k,\ell}$  (das nur die Kombinatorik der Abstammungsverhältnisse widerspiegelt) unabhängig sind.

In der zweiten Zeile ersetzen wir

$$\mathbb{P}_{\theta}(J_{k,\ell}=i)=\frac{\binom{n-i-1}{k-2}}{\binom{n-1}{k-1}},$$

denn mit Korollar A.4 ist die Aufteilung in k (zufällig nummerierte) Familiengrößen uniform auf allen  $\{(m_1, \ldots, m_i) \in \mathbb{N}^i : m_1 + \cdots + m_i = n\}$ : Es gibt  $\binom{n-1}{k-1}$  viele Möglichkeiten, bei  $\binom{n-i-1}{k-2}$  davon ist  $J_{k,\ell} = i$ .

Schließlich beachte in der letzten Zeile  $\sum_{k=2}^{n-i+1} \binom{n-k}{i-1} = \binom{n-1}{i}$ , denn es gibt  $\binom{n-1-(k-1)}{i-1} = \binom{n-k}{i-1}$  viele Teilmengen von  $\{1, \ldots, n-1\}$  der Größe *i*, deren kleinstes Element k-1 ist, und insgesamt  $\binom{n-1}{i}$  Teilmengen von  $\{1, \ldots, n-1\}$  der Größe *i*.

Um  $\mathbb{E}_{\theta}[\xi_{i}^{(n)}\xi_{j}^{(n)}]$  zu bestimmen kann man die Darstellungen (2.21) für *i* und für *j* miteinander multiplizieren und erhält analog zu oben eine Darstellung via eine Doppelsumme über Paare von Kanten im Koaleszenten-Baum. Mittels einer Verfeinerung von Korollar A.4 kann man den kombinatorischen Ausdruck  $\mathbb{P}_{\theta}(J_{k,\ell} = i, J_{k',\ell'} = j)$  bestimmen (man unterscheidet verschiedene Fälle, je nachdem ob die betrachtete Kante  $(k', \ell')$  ein Nachfahre der Kante  $(k, \ell)$  im Baum ist oder nicht) und erhält nach recht umfangreichen Umformungen die oben angegebenen Ausdrücke für  $\operatorname{Cov}_{\theta}[\xi_{i}^{(n)}, \xi_{j}^{(n)}]$ , für Details siehe den Artikel von Yun-Xin Fu, Statistical Properties of Segregating Sites, *Theor. Pop. Biol.* 48, 172–197 (1995), in dem dieser Satz bewiesen wurde.

#### Tajimas7 Test

Betrachte eine *n*-Stichprobe (im IMS-Modell), für  $1 \le i < j \le n$  sei

 $\Delta_{i,j} :=$  Anzahl segregierende Stellen, an denen sich Stichproben *i* und *j* unterscheiden.

Die mittlere Anzahl paarweiser Unterschiede,

$$\widehat{\theta}_{\pi} \coloneqq \frac{1}{\binom{n}{2}} \sum_{1 \le i < j \le n} \Delta_{ij}$$

("Tajimas  $\hat{\theta}_{\pi}$ "), ist ein (auf den beobachteten Sequenzen basierender) Schätzer für die Mutationsrate  $\theta$ .

Beobachtung 2.25. Es gebe s segregierende Stellen.

$$\widehat{\theta}_{\pi} = \frac{1}{\binom{n}{2}} \sum_{1 \le i < j \le n} \sum_{m=1}^{s} \mathbf{1} (\text{Stichpr. } i \text{ und } j \text{ unterschiedl. an } m \text{-ter segr. Stelle}) = \frac{1}{\binom{n}{2}} \sum_{k=1}^{n-1} \xi_k^{(n)} k(n-k)$$

<sup>&</sup>lt;sup>7</sup>Fumio Tajima, Statistical Method for Testing the Neutral Mutation Hypothesis by DNA Polymorphism, *Genetics* 123, 585–595, (1989)



Abbildung 2.5: Formen eines *n*-Koaleszenten für n = 3 (bis auf Umnummerierung der Bläter nur eine Möglichkeit) und n = 4 (zwei Möglichkeiten)

(mit  $\xi_k^{(n)} = \#$  Mut., die in k Stichpr. vorkommen, aus Def. 2.23), d.h.  $\hat{\theta}_{\pi}$  ist eine (lineare) Funktion des Frequenzspektrums.

(Darüberhinaus kann  $\hat{\theta}_{\pi}$  ebenso wie  $S_n$  und  $\hat{\theta}_W$  als Funktion des gefalteten Frequenzspektrums aufgefasst werden, d.h. wir können dies auch dann bestimmen, wenn wir die anzestralen Typen nicht kennen.)

Proposition 2.26. Es gilt

$$\mathbb{E}_{\theta}\left[\widehat{\theta}_{\pi}\right] = \theta, \quad \operatorname{Var}_{\theta}\left(\widehat{\theta}_{\pi}\right) = \frac{n+1}{3(n-1)}\theta + \frac{2(n^2+n+3)}{9n(n-1)}\theta^2$$

Insbesondere:  $\hat{\theta}_{\pi}$  ist erwartungstreuer Schätzer für  $\theta$ , allerdings ist es nicht konsistent:

$$\lim_{n \to \infty} \operatorname{Var}_{\theta} \left( \widehat{\theta}_{\pi} \right) = \frac{1}{3} \theta + \frac{2}{9} \theta^2 > 0.$$

**Bemerkung.**  $\widehat{\theta}_{\pi} = \frac{1}{\binom{n}{2}} \sum_{1 \le i < j \le n} \Delta_{ij}$  wird in der Literatur auch mit  $\pi$  bezeichnet und die (empirische) "Nukleotid-Diversität" (engl. "nucleotide diversity") genannt.

 $(\mathbb{E}_{\theta}[\pi] = \theta$  ist einer der Gründe für die Parametrisierung, dass Mutationen mit Rate  $\theta/2$  längs der Genealogie erscheinen.)

*Beweis.* Betrachte zunächst eine Stichprobe der Größe n = 2: Es ist

$$\mathbb{E}_{\theta}[\Delta_{1,2}] = \mathbb{E}_{\theta}\left[\mathbb{E}_{\theta}[\Delta_{1,2} | T_2]\right] = \mathbb{E}_{\theta}[\theta T_2] = \theta \mathbb{E}_{\theta}[T_2] = \theta$$

(mit  $T_2$  = Zeit, währenddessen die Genealogie aus 2 Linien besteht = Zeit bis zum jgV der beiden Stichproben,  $T_2 \sim \text{Exp}(1)$ ) und

$$\operatorname{Var}_{\theta}[\Delta_{1,2}] = \operatorname{Var}_{\theta}\left[\mathbb{E}_{\theta}[\Delta_{1,2} \mid T_{2}]\right] + \mathbb{E}_{\theta}\left[\operatorname{Var}_{\theta}[\Delta_{1,2} \mid T_{2}]\right]$$
$$= \operatorname{Var}_{\theta}[\theta T_{2}] + \mathbb{E}_{\theta}[\theta T_{2}] = \theta^{2} \operatorname{Var}_{\theta}[T_{2}] + \theta \mathbb{E}_{\theta}[T_{2}] = \theta^{2} + \theta.$$

Betrachte nun eine Stichprobe der Größe n = 3, es bezeichne  $\eta_a$  die Anzahl Mutationen auf Kante a, etc., siehe Abbildung 2.5. Jede Kante kommt in 2 von 3 paarweisen Vergleichen vor, also ist

$$\widehat{\theta}_{\pi,n=3} = \frac{1}{3} \left( \Delta_{1,2} + \Delta_{1,3} + \Delta_{2,3} \right) = \frac{2}{3} \left( \eta_a + \eta_b + \eta_c + \eta_d \right).$$

Sei  $T_j$  die Länge der Zeitspanne, währenddessen der Koaleszent aus j Linien besteht. Nach Definition sind  $\eta_a, \eta_b, \eta_c, \eta_d$  unabhängig, gegeben  $T_3$  und  $T_2$ , und  $\eta_a, \eta_b \sim \text{Poi}(\frac{\theta}{2}T_3), \eta_c \sim \text{Poi}(\frac{\theta}{2}T_2), \eta_d \sim \text{Poi}(\frac{\theta}{2}(T_3 + T_2)), \text{ d.h. } \eta_a + \eta_b + \eta_c + \eta_d \sim \text{Pois}(\theta L_3/2), \text{ wobei } L_3 = 3T_3 + 2T_2 \text{ die Gesamtlänge des Baums ist. Daher ist}$ 

$$\begin{aligned} \operatorname{Var}_{\theta} \left[ \widehat{\theta}_{\pi, n=3} \right] &= \frac{4}{9} \operatorname{Var}_{\theta} \left[ \eta_{a} + \eta_{b} + \eta_{c} + \eta_{d} \right] \\ &= \frac{4}{9} \operatorname{Var}_{\theta} \left[ \mathbb{E}_{\theta} \left[ \eta_{a} + \eta_{b} + \eta_{c} + \eta_{d} \mid L_{3} \right] \right] + \frac{4}{9} \mathbb{E}_{\theta} \left[ \operatorname{Var}_{\theta} \left[ \eta_{a} + \eta_{b} + \eta_{c} + \eta_{d} \mid L_{3} \right] \right] \\ &= \frac{4}{9} \operatorname{Var}_{\theta} \left[ \frac{\theta}{2} L_{3} \right] + \frac{4}{9} \mathbb{E}_{\theta} \left[ \frac{\theta}{2} L_{3} \right] = \frac{4}{9} \frac{\theta^{2}}{4} \left( 9 \cdot \frac{1}{3^{2}} + 4 \cdot 1 \right) + \frac{4}{9} \frac{\theta}{2} \left( 3 \cdot \frac{1}{3} + 2 \cdot 1 \right) = \frac{5}{9} \theta^{2} + \frac{2}{3} \theta. \end{aligned}$$

Andererseits ist wegen der Symmetrien der Verteilung des Koaleszenten  $\text{Cov}_{\theta}[\Delta_{1,2}, \Delta_{1,3}] = \text{Cov}_{\theta}[\Delta_{1,2}, \Delta_{2,3}]$ , etc. und somit

$$\operatorname{Var}_{\theta}\left[\widehat{\theta}_{\pi,n=3}\right] = \frac{1}{9} \cdot 3 \cdot \operatorname{Var}_{\theta}\left[\Delta_{1,2}\right] + \frac{1}{9} \cdot 6 \cdot \operatorname{Cov}_{\theta}\left[\Delta_{1,2}, \Delta_{1,3}\right] = \frac{1}{3}(\theta^{2} + \theta) + \frac{2}{3}\operatorname{Cov}_{\theta}\left[\Delta_{1,2}, \Delta_{1,3}\right],$$

folglich

$$\operatorname{Cov}_{\theta}[\Delta_{1,2}, \Delta_{1,3}] = \frac{3}{2} \operatorname{Var}_{\theta}[\widehat{\theta}_{\pi,n=3}] - \frac{1}{2} (\theta^2 + \theta) = \frac{1}{3} \theta^2 + \frac{1}{2} \theta.$$
(2.22)

Betrachte nun eine Stichprobe der Größe n = 4: Es gibt 2 mögliche Baumtopologien (siehe Abb. 2.5).

Wir untersuchen zunächst Fall I (das mittlere Diagramm in Abb. 2.5). Gegeben  $T_4, T_3, T_2$  sind hier  $\eta_a \sim \text{Poi}(\frac{\theta}{2}T_4), \eta_b \sim \text{Poi}(\frac{\theta}{2}T_4), \eta_c \sim \text{Poi}(\frac{\theta}{2}(T_4 + T_3)), \eta_d \sim \text{Poi}(\frac{\theta}{2}(T_4 + T_3 + T_2)), \eta_e \sim \text{Poi}(\frac{\theta}{2}T_3), \eta_f \sim \text{Poi}(\frac{\theta}{2}T_2)$  und unabhängig. Weiter ist in diesem Fall

$$\widehat{\theta}_{\pi,n=4} = \Delta(1) = \frac{1}{\binom{4}{2}} \left( 3\eta_a + 3\eta_b + 3\eta_c + 3\eta_d + 4\eta_e + 3\eta_f \right) =: \frac{1}{6} X_1$$

(beachte: wenn oberhalb einer Kante  $\ell$  Blätter liegen, so tritt sie in  $\ell \cdot (n - \ell)$  paarweisen Vergleichen auf), somit ist

$$\begin{split} \mathbb{E}_{\theta} \Big[ \widehat{\theta}_{\pi,n=4} \, \big| \, \text{Top.=I} \Big] &= \frac{1}{6} \frac{\theta}{2} \mathbb{E} \Big[ 3T_4 + 3T_4 + 3(T_4 + T_3) + 3(T_4 + T_3 + T_2) + 4T_3 + 3T_2 \Big] \\ &= \frac{\theta}{12} \Big( \frac{3}{\binom{4}{2}} + \frac{3}{\binom{4}{2}} + 3\Big( \frac{1}{\binom{4}{2}} + \frac{1}{\binom{3}{2}} \Big) + \frac{4}{\binom{3}{2}} + 3\Big( \frac{1}{\binom{4}{2}} + \frac{1}{\binom{3}{2}} + 1 \Big) + 3 \cdot 1 \Big) = \frac{17}{18} \theta, \\ \text{Var}_{\theta} \Big[ \widehat{\theta}_{\pi,n=4} \, \big| \, \text{Top.=I} \Big] &= \frac{1}{36} \mathbb{E}_{\theta} \Big[ \, \text{Var}_{\theta} \big[ X_1 \, \big| \, T_4, T_3, T_2 \big] \Big] + \frac{1}{36} \, \text{Var}_{\theta} \Big[ \mathbb{E}_{\theta} \big[ X_1 \, \big| \, T_4, T_3, T_2 \big] \Big] \\ &= \frac{1}{36} \mathbb{E}_{\theta} \Big[ \frac{\theta}{2} \big( 9T_4 + 9T_4 + 9\big(T_3 + T_4\big) + 9\big(T_4 + T_3 + T_2\big) + 16T_3 + 9T_2 \big) \Big] \\ &+ \frac{1}{36} \, \text{Var}_{\theta} \Big[ \frac{\theta}{2} \big( 3T_4 + 3T_4 + 3\big(T_4 + T_3\big) + 3\big(T_4 + T_3 + T_2\big) + 4T_3 + 3T_2 \big) \Big] \\ &= \frac{\theta}{72} \Big( \frac{\theta}{6} + \frac{\theta}{6} + 9\big(\frac{1}{6} + \frac{1}{3}\big) + 9\big(\frac{1}{6} + \frac{1}{3} + 1\big) + \frac{16}{3} + 9 \cdot 1 \big) \\ &+ \frac{\theta^2}{144} \, \text{Var}_{\theta} \Big[ 12T_4 + 10T_3 + 6T_2 \Big] \\ &= \frac{53}{108} \theta + \frac{\theta^2}{144} \Big( \frac{12^2}{6^2} + \frac{10^2}{3^2} + \frac{6^2}{1^2} \Big) = \frac{53}{108} \theta + \frac{115}{324} \theta^2. \end{split}$$

Untersuchen wir nun Fall 2 (das rechte Diagramm in Abb. 2.5). Gegeben  $T_4, T_3, T_2$  sind hier  $\eta_a \sim \operatorname{Poi}(\frac{\theta}{2}T_4), \eta_b \sim \operatorname{Poi}(\frac{\theta}{2}T_4), \eta_c \sim \operatorname{Poi}(\frac{\theta}{2}(T_4 + T_3)), \eta_d \sim \operatorname{Poi}(\frac{\theta}{2}(T_4 + T_3)), \eta_e \sim \operatorname{Poi}(\frac{\theta}{2}(T_3 + T_2)), \eta_f \sim \operatorname{Poi}(\frac{\theta}{2}T_2)$  und unabhängig, weiter ist in diesem Fall (mit Argumentation analog zu Fall 1)

$$\widehat{\theta}_{\pi,n=4} = \Delta(2) = \frac{1}{\binom{4}{2}} \left( 3\eta_a + 3\eta_b + 3\eta_c + 3\eta_d + 4\eta_e + 4\eta_f \right) =: \frac{1}{6} X_2,$$

somit ergibt sich

$$\begin{split} \mathbb{E}_{\theta} \Big[ \widehat{\theta}_{\pi,n=4} \, \big| \, \text{Top.=2} \Big] &= \frac{1}{6} \frac{\theta}{2} \mathbb{E} \Big[ 3T_4 + 3T_4 + 3(T_4 + T_3) + 3(T_4 + T_3) + 4(T_3 + T_2) + 4T_2 \Big] \\ &= \frac{\theta}{12} \mathbb{E} \Big[ 12T_4 + 10T_3 + 8T_2 \Big] = \frac{\theta}{12} \Big( \frac{12}{6} + \frac{10}{3} + 8 \Big) = \frac{10}{9} \theta, \\ \text{Var}_{\theta} \Big[ \widehat{\theta}_{\pi,n=4} \, \big| \, \text{Top.=2} \Big] &= \frac{1}{36} \mathbb{E}_{\theta} \Big[ \, \text{Var}_{\theta} \big[ X_2 \, \big| \, T_4, T_3, T_2 \big] \Big] + \frac{1}{36} \, \text{Var}_{\theta} \Big[ \mathbb{E}_{\theta} \big[ X_2 \, \big| \, T_4, T_3, T_2 \big] \Big] \\ &= \frac{1}{36} \mathbb{E}_{\theta} \Big[ \frac{\theta}{2} \big( 9T_4 + 9T_4 + 9(T_4 + T_3) + 9(T_4 + T_3) + 16(T_3 + T_2) + 16T_2 \big) \Big] \\ &\quad + \frac{1}{36} \, \text{Var}_{\theta} \Big[ \frac{\theta}{2} \big( 3T_4 + 3T_4 + 3(T_4 + T_3) + 3(T_4 + T_3) + 4(T_3 + T_2) + 4T_2 \big) \Big] \\ &= \frac{\theta}{72} \mathbb{E} \Big[ 36T_4 + 34T_3 + 32T_2 \Big] + \frac{\theta^2}{144} \, \text{Var}_{\theta} \Big[ 12T_4 + 10T_3 + 8T_2 \Big] \\ &= \frac{\theta}{72} \Big( \frac{36}{6} + \frac{34}{3} + 32 \cdot 1 \Big) + \frac{\theta^2}{144} \Big( \frac{12^2}{6^2} + \frac{10^2}{3^2} + \frac{8^2}{1^2} \Big) = \frac{37}{54} \theta + \frac{89}{162} \theta^2. \end{split}$$

Insgesamt ist mit  $\mathbb{P}(\text{Top.=1}) = 2/3 = 1 - \mathbb{P}(\text{Top.=2})$  (denn damit der 2. Fall für die Baumtopologie eintritt, muss die zweitjüngste Verschmelzung das Paar von Linien betreffen, das bis dahin noch an keiner Verschmelzung teilgenommen hat, dies ist dann 1 von 3 Möglichkeiten)

$$\mathbb{E}_{\theta}\left[\widehat{\theta}_{\pi,n=4}\right] = \frac{2}{3} \cdot \frac{17}{18}\theta + \frac{1}{3} \cdot \frac{10}{9}\theta = \theta$$

und

$$\begin{aligned} \operatorname{Var}_{\theta} \left[ \widehat{\theta}_{\pi, n=4} \right] &= \mathbb{E}_{\theta} \left[ \operatorname{Var}_{\theta} \left[ \widehat{\theta}_{\pi, n=4} \, | \, \operatorname{Top.} \right] \right] + \, \operatorname{Var}_{\theta} \left[ \mathbb{E}_{\theta} \left[ \widehat{\theta}_{\pi, n=4} \, | \, \operatorname{Top.} \right] \right] \\ &= \frac{2}{3} \left( \frac{53}{108} \theta + \frac{115}{324} \theta^2 \right) + \frac{1}{3} \left( \frac{37}{54} \theta + \frac{89}{162} \theta^2 \right) + \frac{2}{3} \left( \frac{17}{18} \theta - \theta \right)^2 + \frac{1}{3} \left( \frac{10}{9} \theta - \theta \right)^2 \\ &= \frac{23}{54} \theta^2 + \frac{5}{9} \theta. \end{aligned}$$

Andererseits ist wie oben wegen der Symmetrien der Verteilung des Koaleszenten

$$\begin{aligned} \operatorname{Var}_{\theta} \left[ \widehat{\theta}_{\pi, n=4} \right] &= \frac{1}{36} \operatorname{Cov}_{\theta} \left[ \sum_{1 \le i < j \le 4} \Delta_{i, j}, \sum_{1 \le k < \ell \le 4} \Delta_{k, \ell} \right] \\ &= \frac{1}{36} \Big( 6 \operatorname{Var}_{\theta} \left[ \Delta_{1, 2} \right] + 6 \cdot 2 \cdot 2 \operatorname{Cov}_{\theta} \left[ \Delta_{1, 2}, \Delta_{1, 3} \right] + 6 \operatorname{Cov}_{\theta} \left[ \Delta_{1, 2}, \Delta_{3, 4} \right] \Big) \\ &= \frac{1}{6} \big( \theta^{2} + \theta \big) + \frac{2}{3} \big( \frac{1}{3} \theta^{2} + \frac{1}{2} \theta \big) + \frac{1}{6} \operatorname{Cov}_{\theta} \big[ \Delta_{1, 2}, \Delta_{3, 4} \big] \end{aligned}$$

und somit

$$\operatorname{Cov}_{\theta}[\Delta_{1,2}, \Delta_{3,4}] = 6\operatorname{Var}_{\theta}[\widehat{\theta}_{\pi,n=4}] - (\theta^2 + \theta) - 4(\frac{1}{3}\theta^2 + \frac{1}{2}\theta) = \frac{2}{9}\theta^2 + \frac{1}{3}\theta.$$
(2.23)

Schließlich betrachten wir den allgemeinen Fall einer Stichprobe der Größe n:

$$\begin{split} \mathbb{E}_{\theta} \Big[ \widehat{\theta}_{\pi} \Big] &= \frac{1}{\binom{n}{2}} \sum_{1 \le i < j \le n} \mathbb{E}_{\theta} \Big[ \Delta_{i,j} \Big] = \frac{1}{\binom{n}{2}} \binom{n}{2} \mathbb{E}_{\theta} \Big[ \Delta_{1,2} \Big] = \theta, \\ \operatorname{Var}_{\theta} \Big[ \widehat{\theta}_{\pi} \Big] &= \frac{1}{\binom{n}{2}} \operatorname{Cov}_{\theta} \Big[ \sum_{1 \le i < j \le n} \Delta_{i,j}, \sum_{1 \le k < \ell \le n} \Delta_{k,\ell} \Big] \\ &= \frac{1}{\binom{n}{2}} \Big( \binom{n}{2} \operatorname{Var}_{\theta} \Big[ \Delta_{1,2} \Big] + \binom{n}{2} 2 (n-2) \operatorname{Cov}_{\theta} \Big[ \Delta_{1,2}, \Delta_{1,3} \Big] \\ &+ \binom{n}{2} \binom{n-2}{2} \operatorname{Cov}_{\theta} \Big[ \Delta_{1,2}, \Delta_{3,4} \Big] \Big) \\ &= \frac{1}{\binom{n}{2}} \Big( \theta^{2} + \theta + 2(n-2) \Big( \frac{1}{3} \theta^{2} + \frac{1}{2} \theta \Big) + \binom{n-2}{2} \Big( \frac{2}{9} \theta^{2} + \frac{1}{3} \theta \Big) \Big) \\ &= \frac{n+1}{3(n-1)} \theta + \frac{2(n^{2}+n+3)}{9n(n-1)} \theta^{2}. \end{split}$$

Passen beobachtete Sequenzdaten zum Modell? Unser wahrscheinlichkeitstheoretisches Modell beschreibt die Verteilung von n beobachteten Sequenzen, die wir an den Blättern eines Kingmann-Koaleszenten ablesen, auf dessen Kanten gemäß einem Poissonprozess mit einer gewissen Rate  $\theta/2$ Mutationen liegen, die den Typ jeweils gemäß dem IMS-Modell ändern. Angesichts Satz 1.5 ist die biologische Interpretation, dass wir n Stichproben aus einer "panmiktischen" Population konstanter Größe sehen und dass die genetische Variabilität (am betrachteten Ort im Genom) "neutral" ist (und dass die Annahmen des IMS-Modells wenigstens approximativ zutreffen).

Die Tatsache, dass sowohl  $\hat{\theta}_W$  als auch  $\hat{\theta}_{\pi}$  in diesem Modell erwartungstreue Schätzer für (das unbekannte)  $\theta$  sind, gestattet es, für die Nullhypothese "das Modell beschreibt die Daten zutreffend" einen statistischen Test zu formulieren. Wenn das Modell zutrifft, sollte nämlich

$$\widehat{\theta}_{\pi} - \widehat{\theta}_{W} \approx 0$$

bis auf "zufällige Fluktuationen" gelten. Diese Idee geht auf F. Tajima zurück, siehe den in Fußnote 7 auf S. 46 zitierten Artikel.

Um einzuschätzen, wie groß die "typischen" Fluktuationen sind, sollten wir (zumindest)  $\operatorname{Var}_{\theta} \left[ \widehat{\theta}_{\pi} - \widehat{\theta}_{W} \right]$  bestimmen können.

**Bericht 2.27.** Es gilt  $\operatorname{Cov}_{\theta} \left[ S_n, \widehat{\theta}_{\pi} \right] = \theta + \left( \frac{1}{2} + \frac{1}{n} \right) \theta^2$ , also  $\operatorname{Cov}_{\theta} \left[ \widehat{\theta}_W, \widehat{\theta}_{\pi} \right] = \frac{\theta}{h_n} + \left( \frac{1}{2} + \frac{1}{n} \right) \frac{\theta^2}{h_n}$  und somit  $\operatorname{Vor}_{\theta} \left[ \widehat{\theta}_{-} - \widehat{\theta}_{++} \right] = \left( -\frac{n+1}{2} - \frac{1}{2} \right) \theta + \left( \frac{2(n^2+n+3)}{2} - \frac{n+2}{2} + \frac{g_n}{2} \right) \theta^2$ 

$$\operatorname{Var}_{\theta}\left[\widehat{\theta}_{\pi} - \widehat{\theta}_{W}\right] = \left(\frac{n+1}{3(n-1)} - \frac{1}{h_{n}}\right)\theta + \left(\frac{2(n^{2}+n+3)}{9n(n-1)} - \frac{n+2}{nh_{n}} + \frac{g_{n}}{h_{n}^{2}}\right)\theta^{2}$$

(mit  $h_n = \sum_{i=1}^{n-1} \frac{1}{i}, g_n = \sum_{i=1}^{n-1} \frac{1}{i^2}$ ).

Weiterhin ist

$$\widehat{V} \coloneqq \alpha_1 S_n + \alpha_2 S_n (S_n - 1) \tag{2.24}$$

mit

$$\alpha_1 = \left(\frac{n+1}{3(n-1)} - \frac{1}{h_n}\right) / h_n, \quad \alpha_2 = \left(\frac{2(n^2 + n + 3)}{9n(n-1)} - \frac{n+2}{nh_n} + \frac{g_n}{h_n^2}\right) / \left(h_n^2 + g_n\right)$$
(2.25)

ein erwartungstreuer Schätzer für  $\operatorname{Var}_{\theta} \left[ \widehat{\theta}_{\pi} - \widehat{\theta}_{W} \right]$ .

Die Formel für  $\operatorname{Cov}_{\theta}[S_n, \widehat{\theta}_{\pi}]$  kann man mittels einer ähnlichen Zerlegung wie im Beweis von Proposition 2.26 beweisen, siehe F. Tajima, a.a.O. Zusammen mit Beobachtung 2.20 und Proposition 2.26 ergibt sich daraus die Formel für  $\operatorname{Var}_{\theta}[\widehat{\theta}_{\pi} - \widehat{\theta}_{W}]$ .

Aus Beobachtung 2.20 folgt auch

$$\mathbb{E}_{\theta}[S_n] = \theta h_n \quad \text{und} \quad \mathbb{E}_{\theta}[S_n(S_n - 1)] = \operatorname{Var}_{\theta}[S_n] + \left(\mathbb{E}_{\theta}[S_n]\right)^2 - \mathbb{E}_{\theta}[S_n] = \theta^2 h_n + \theta^2 g_n$$
  
d.h.  $\mathbb{E}_{\theta}[\widehat{V}] = \operatorname{Var}_{\theta}[\widehat{\theta}_{\pi} - \widehat{\theta}_W].$ 

**Definition 2.28** (Tajimas *D*).  $D \coloneqq \frac{\widehat{\theta}_{\pi} - \widehat{\theta}_{W}}{\sqrt{\widehat{V}}}$  mit  $\widehat{V}$  aus (2.24) heißt Tajimas *D*.

Die Teststatistik D erfüllt  $\mathbb{E}_{\theta}[D] \approx 0$ ,  $\operatorname{Var}_{\theta}(D) \approx 1$  (die Erwartung ist nicht exakt = 0, da Zähler und Nenner nicht unabhängig sind, die Varianz ist nicht exakt = 1, da  $\widehat{V}$  nur ein Schätzer für die Varianz des Zählers ist). Die Formulierung ist (beispielsweise) durch den klassischen *t*-Test inspiriert: Dort normiert man einen empirischen Mittelwert von *n* Beobachtungswerten mit dem Standardfehler, einem Schätzer für die Streuung.

Um anhand von D einen statistischen Test zu formulieren, benötigen wir (für ein vorgegebenes Signifikanzniveau  $\alpha$ ) sogenannte kritische Werte, d.h. geeignete Quantile von D unter der Nullhypothese.

Auf dem Ereignis  $\{S_n = s\}$  gilt

$$\widehat{\theta}_W = s/h_n, \quad \widehat{V} = \alpha_1 s + \alpha_2 s(s-1),$$

der kleinste möglicher Wert von  $\widehat{\theta}_{\pi}$  ist dann

$$\frac{1}{\binom{n}{2}}s(n-1) = \frac{2s}{n}$$

Dies geschieht, wenn  $\xi_1^{(n)} + \xi_{n-1}^{(n)} = n$ ,  $\xi_i^{(n)} = 0$  für  $2 \le i \le n - 2$  gilt (insbesondere, wenn alle Mutationen auf sogenannten externen Kanten – die direkt zu einem Blatt führen – liegen). Der kleinste mögliche Wert von D ist dann somit

$$\frac{2s/n - s/h_n}{\sqrt{\alpha_1 s + \alpha_2 s(s-1)}} \xrightarrow[s \to \infty]{} \frac{2/n - 1/h_n}{\sqrt{\alpha_2}} =: d_{\min} \ \left( = d_{\min}(n) \right).$$



Abbildung 2.6: Ein "sternförmiger" (links) und ein "Hühnerbein" - (rechts) Koaleszentenbaum

Während ein so kleiner Wert von *D* unter dem Modell, in dem Mutationen auf den Kingman-Koaleszenten fallen, eher untypisch ist, wäre dies ist in einer "sternförmigen" Genealogie, in der die externen Äste die Gesamtlänge des Baumes dominieren (siehe Abbildung 2.6), typisch.

Andererseits ist auf  $\{S_n = s\}$  der größte mögliche Wert von  $\widehat{\theta}_{\pi}$ 

$$\frac{1}{\binom{n}{2}}s\lceil n/2\rceil\lfloor n/2\rfloor = 2s\frac{\lceil n/2\rceil\lfloor n/2\rfloor}{n(n-1)}.$$

Dies geschieht, wenn  $\xi_{\lceil n/2 \rceil}^{(n)} = n$ ,  $\xi_i^{(n)} = 0$  für  $i \neq \lceil n/2 \rceil$  (d.h. wenn alle Mutationen auf sehr "balanzierten" Kanten liegen, die die Blätter in genau zwei Hälften teilen). Der größte mögliche Wert von D ist dann

$$\frac{\frac{2s[n/2][n/2]}{n(n-1)} - s/h_n}{\sqrt{\alpha_1 s + \alpha_2 s(s-1)}} \xrightarrow[s \to \infty]{} \xrightarrow[s \to \infty]{\frac{2[n/2][n/2]}{n(n-1)} - 1/h_n}{\sqrt{\alpha_2}} =: d_{\max} \ \left( = d_{\max}(n) \right)$$

Während ein so kleiner Wert von *D* unter dem Kingman-Koaleszenten untypisch wäre, wäre dies ist in einer (sehr balanzierten) "Hühnerbein-artigen"-Genealogie, in der zwei innere Äste die Gesamtlänge des Baums dominieren (siehe Abbildung 2.6), typisch.

Im Gegensatz etwa zum klassischen t-Test ist die Verteilung von D unter der Nullhypothese

"die Beobachtungen entstehen durch die Typen an den Blättern eines  

$$n$$
-Koaleszenten, längs dessen Kanten sich mit Rate  $\theta/2$  Mutationen gemäß (2.26)  
IMS-Modell ereignen"

nicht explizit bekannt und hängt von dem unbekannten  $\theta$  ab.

Tajimas pragmatisch-heuristische Lösung: Approximiere die Verteilung von D durch eine skalierte Beta-Verteilung, so dass der Träger =  $[d_{\min}, d_{\max}]$ , EW= 0 und Var= 1 gilt (was recht plausibel passt, siehe Abbildung 2.7): Verwende die approximative Dichte

$$f_{\rm appr}(d) = \frac{\Gamma(u+v)(d-d_{\rm min})^{u-1}(d_{\rm max}-d)^{v-1}}{\Gamma(u)\Gamma(v)(d_{\rm max}-d_{\rm min})^{u+v-1}}, \quad d_{\rm min} < d < d_{\rm max}$$
(2.27)

mit

$$u = \frac{(1 + d_{\max}d_{\min})d_{\min}}{d_{\max} - d_{\min}}, \quad v = -\frac{(1 + d_{\max}d_{\min})d_{\max}}{d_{\max} - d_{\min}}.$$
 (2.28)

(beachte  $d_{\min} < 0 < d_{\max}$ ).

Diese Formeln entspringen dem Ansatz

$$D \approx (d_{\max} - d_{\min})B + d_{\min}$$
 mit  $B \sim \text{Beta}(u, v)$ .

Beta(u, v) hat EW  $\frac{u}{u+v}$  und Var  $\frac{uv}{(u+v)^2(u+v+1)}$ , aus dem Ansatz und den geforderten Normierungen ergibt sich

$$\frac{u}{u+v} = \frac{-d_{\min}}{d_{\max} - d_{\min}} \implies v = u \frac{d_{\max} - d_{\min}}{-d_{\min}} - u = u \frac{d_{\max}}{-d_{\min}},$$

$$\frac{uv}{(u+v)^2(u+v+1)} = \frac{u^2 \frac{d_{\max}}{-d_{\min}}}{u^2 \left(1 + \frac{d_{\max}}{-d_{\min}}\right)^2 \left(u(1 + \frac{d_{\max}}{-d_{\min}}) + 1\right)} = \frac{-d_{\max}d_{\min}}{(d_{\max} - d_{\min})^2 \left(u(1 + \frac{d_{\max}}{-d_{\min}}) + 1\right)}$$

$$= \frac{1}{(d_{\max} - d_{\min})^2}$$

$$\implies u = \frac{-d_{\max}d_{\min} - 1}{1 - \frac{d_{\max}}{d_{\min}}} = \frac{d_{\min}(1 + d_{\max}d_{\min})}{d_{\max} - d_{\min}}, \quad v = -\frac{d_{\max}(1 + d_{\max}d_{\min})}{d_{\max} - d_{\min}},$$

woraus sich (2.28) ergibt.

**Definition 2.29** (Tajimas Test). Sei  $\alpha \in (0, 1)$ ,  $q_{\text{Beta}(u,v)}(\alpha/2)$ ,  $q_{\text{Beta}(u,v)}(1 - \alpha/2)$  das  $\alpha/2$ - bzw.  $(1 - \alpha/2)$ -Quantil der Beta(u, v)-Verteilung mit angepassten Parametern u, v aus (2.28).

Lehne  $H_0$ : (2.26) ab, wenn

$$D < (d_{\max} - d_{\min})q_{\text{Beta}(u,v)}(\alpha/2) + d_{\min} \quad \text{oder}$$
$$D > (d_{\max} - d_{\min})q_{\text{Beta}(u,v)}(1 - \alpha/2) + d_{\min}.$$

Dieser Test hält (zumindest approximativ) das Signifikanzniveau  $\alpha$  ein.

**Beispiel.** Für die Daten aus Bsp. 2.18 ergibt sich  $n = 8, s = 31, \xi_1^{(8)} = 13, \xi_2^{(8)} = 1, \xi_7^{(8)} = 17$ , somit  $\widehat{\theta}_{\pi} \doteq 7.93, \widehat{\theta}_W \doteq 11.96, D \doteq -1.79$ 

Tajimas Approximation liefert ein 95%-Konfidenzintervall für D unter dem Standard-Kingman-Koaleszenten von [-1.663, 1.975], d.h. die Abweichung von 0 ist auf dem 5%-Niveau signifikant (s. Tajima, a.a.O., Table 2, S 592).

**Diskussion.** In der biologischen Interpretation nennt man Tajimas Test gelegentlich etwas salopp einen "Test auf Neutralität", da die Nullhypothese (2.26) aus einem Modell ohne Selektion stammt.

Signifikante Abweichungen von  $D \approx 0$  legen Alternativhypothesen nahe, unter denen der Baum, der die Stichproben verbindet, eher nicht wie ein "typischer" Koaleszent aussicht.

Ein signifikant negatives D < 0 passt eher zu einem Baum, in dem externe Aste dominieren (Abbildung 2.6, links). Biologische Szenarien, in denen solche Genealogien typisch sind, wären beispielsweise gerichtete Selektion am betrachteten Ort im Genom (oder in dessen "Nähe", ein sogenannter selektiver "sweep") oder eine stark wachsende Population.



Abbildung 2.7: Simulation der Verteilung von D für n = 25 und  $\theta = 10$  (links) bzw.  $\theta = 2$  (rechts) unter dem Kingman-Koaleszenten mit IMS-Mutationen sowie angepasste skalierte Beta-Dichte aus (2.27). Histogramm jeweils basierend auf 100.000 simulierten Datensätzen.

Ein signifikant positives D > 0 passt eher zu einem Baum, in dem wenige interne Aste dominieren (Abbildung 2.6, rechts). Populationsszenarien, in denen solche Genealogien typisch sind, sind beispielsweise (räumlich stark) strukturierte Populationen oder sogenannte balanzierende Selektion (bei der selektive Kräfte gewissermaßen eine genetische Substruktur in der Population aufrecht erhalten).

### Eine "exakte" Version von Tajimas Test

K. L. Simonsen, G. A. Churchill und C. F. Aquadro haben in dem Artikel Properties of statistical tests of neutrality for DNA polymorphism data, *Genetics* 141:413–429, (1995) eine Version von Tajimas Test vorgeschlagen, die ohne die (nicht wörtlich gerechtfertigte) Approximation von *D* durch eine Beta-Verteilung auskommt<sup>8</sup>.

Das unbekannte  $\theta$  wird dabei als "Störparameter" (engl. "nuisance parameter") aufgefasst. Wir wählen  $\beta > 0$  (und typischerweise klein) und konstruieren zunächst ein Konfidenzintervall für  $\theta$ zum Irrtumsniveau  $\beta$ :

Sei für  $s \in \mathbb{N}_0$ 

$$\widehat{\theta}_L(s) = \min\left\{\theta > 0 : \mathbb{P}_{\theta}(S_n \ge s) > \beta/2\right\}, \quad \widehat{\theta}_R(s) = \max\left\{\theta > 0 : \mathbb{P}_{\theta}(S_n \le s) > \beta/2\right\}$$

<sup>&</sup>lt;sup>8</sup>Die Konstruktion verwendet ein allgemeines statistisches Prinzip, siehe R. L. Berger und D. D. Boos, *P* values maximized over a confidence set for the nuisance parameter, *Journal of the American Statistical Association* 89, No. 427, 1012–1016, (1994).

Dies ist mittels der Verteilungsfunktion von  $S_n$  unter  $\mathbb{P}_{\theta}$  aus Lemma 2.30 unten zumindest numerisch möglich; da diese als Funktion von  $\theta$  stetig ist, gilt tatsächlich  $\mathbb{P}_{\widehat{\theta}_L(s)}(S_n \ge s) = \beta/2$  und  $\mathbb{P}_{\widehat{\theta}_R(s)}(S_n \le s) = \beta/2$  für  $s \in \mathbb{N}_0$ .

Da  $\theta \mapsto \mathbb{P}_{\theta}(S_n \ge s)$  monoton wachsend in  $\theta$  ist, gilt

$$\widehat{\theta}_L(s) > \theta \iff \mathbb{P}_{\theta}(S_n \ge s) \le \beta/2 \quad \text{und} \quad \widehat{\theta}_R(s) < \theta \iff \mathbb{P}_{\theta}(S_n \le s) \le \beta/2.$$

Damit gilt

$$\forall \theta > 0 : \mathbb{P}_{\theta} ( [\widehat{\theta}_L(S_n), \widehat{\theta}_R(S_n)] \ni \theta ) \ge 1 - \beta$$

denn für $\theta > 0$ ist

$$\mathbb{P}_{\theta}\left(\left[\widehat{\theta}_{L}(S_{n}),\widehat{\theta}_{R}(S_{n})\right] \not = \mathbb{P}_{\theta}\left(\theta < \widehat{\theta}_{L}(S_{n})\right) + \mathbb{P}_{\theta}\left(\theta > \widehat{\theta}_{R}(S_{n})\right)$$
$$= \mathbb{P}_{\theta}\left(S_{n} \in \left\{s : \theta < \widehat{\theta}_{L}(s)\right\}\right) + \mathbb{P}_{\theta}\left(S_{n} \in \left\{s : \theta > \widehat{\theta}_{R}(s)\right\}\right)$$
$$= \sum_{s : \mathbb{P}_{\theta}(S_{n} \ge s) \le \beta/2} \mathbb{P}_{\theta}\left(S_{n} = s\right) + \sum_{s : \mathbb{P}_{\theta}(S_{n} \le s) \le \beta/2} \mathbb{P}_{\theta}\left(S_{n} = s\right) \le \frac{\beta}{2} + \frac{\beta}{2}.$$

Dann bestimmt man bei beobachtetem Wert von  $S_n$  (mittels Simulation, für  $\theta$ -Werte aus einem geeignet feinen Gitter in  $[\widehat{\theta}_L(S_n), \widehat{\theta}_R(S_n)]$ )

$$D_{L}^{*} = \min\left\{\frac{\alpha}{2}\text{-Quantil von }\mathcal{L}_{\theta}(D) : \theta \in [\widehat{\theta}_{L}(S_{n}), \widehat{\theta}_{R}(S_{n})]\right\},\$$
$$D_{L}^{*} = \max\left\{\left(1 - \frac{\alpha}{2}\right)\text{-Quantil von }\mathcal{L}_{\theta}(D) : \theta \in [\widehat{\theta}_{L}(S_{n}), \widehat{\theta}_{R}(S_{n})]\right\}$$

Somit gilt

$$\forall \theta > 0 : \mathbb{P}_{\theta} \left( D \notin [D_L^*, D_R^*] \right) \le \alpha + \beta,$$

d.h. der Test

lehne 
$$H_0$$
: (2.26) ab, wenn  $D < D_L^*$  oder  $D > D_R^*$ 

hält Niveau  $\alpha + \beta$  ein (zumindest theoretisch, wenn man die Quantile im 2. Schritt exakt bestimmen könnte).

**Beispiel.** Für die Daten aus Bsp. 2.18 ( $n = 8, D \doteq -1.79$ ) berichten Simonsen, Churchill und Aquadro, a.a.O., Table 3 gemäß diesem Ansatz ein 95%-Konfidenzintervall für D unter dem Standard-Kingman-Koaleszenten von [-1.80, 1.83] (für  $n = 10, S_n \in [27, 41]$ ), d.h. die Abweichung ist "gerade so" nicht signifikant auf dem 5%-Niveau.

**Lemma 2.30** (Explizite Verteilung von  $S_n$ ). Es gilt

$$\mathbb{P}_{\theta} \left( S_n = m \right) = \frac{n-1}{\theta} \sum_{k=1}^{n-1} (-1)^{k-1} {\binom{n-2}{k-1}} \left( \frac{\theta}{\theta+k} \right)^{m+1}, \quad m \in \mathbb{N}_0,$$
$$\mathbb{P}_{\theta} \left( S_n \le s \right) = 1 - \sum_{k=1}^{n-1} (-1)^{k-1} {\binom{n-1}{k}} \left( \frac{\theta}{\theta+k} \right)^{s+1}, \quad s \in \mathbb{N}_0.$$

**Bemerkung.** Dies ist eine Version von Lemma A.6 (Dichte der Faltung exponentieller ZVn) für den diskreten Fall (Faltung geometrischer ZVn).

Beweis.  $S_n = S_{n,n} + S_{n,n-1} + \dots + S_{n,2} \text{ mit } S_{n,j} \sim \text{geom}(\frac{i-1}{\theta+i-1})$  u.a. Sei  $u \in [0, 1]$ : Es ist

$$\mathbb{E}\left[u^{S_{n,j}}\right] = \sum_{\ell=0}^{\infty} u^{\ell} \frac{j-1}{\theta+j-1} \left(\frac{\theta}{\theta+j-1}\right)^{\ell} = \frac{j-1}{\theta+j-1} \frac{1}{1-u\frac{\theta}{\theta+j-1}} = \frac{j-1}{j-1+\theta(1-u)},$$

somit

$$\mathbb{E}\left[u^{S_n}\right] = \prod_{j=2}^n \mathbb{E}\left[u^{S_{n,j}}\right] = \prod_{k=1}^{n-1} \frac{k}{k + \theta(1-u)}$$

Weiter ist

$$\prod_{k=1}^{n-1} \frac{k}{k+z} = \sum_{k=1}^{n-1} \frac{a_{n,k}}{k+z} \quad (z \in \mathbb{C} \smallsetminus -\mathbb{N}) \quad \text{mit } a_{n,k} = \frac{(n-1)!}{\prod_{j \neq k}^{n-1}(j-k)} = (-1)^k (n-1) \binom{n-2}{k-1},$$

also

$$\mathbb{E}_{\theta}\left[u^{S_n}\right] = \sum_{m=0}^{\infty} u^m \mathbb{P}_{\theta}\left(S_n = m\right) = \sum_{k=1}^{n-1} a_{n,k} \sum_{m=0}^{\infty} \left(\frac{\theta}{\theta+k}\right)^m u^m = \sum_{m=0}^{\infty} u^m \sum_{k=1}^{n-1} a_{n,k} \left(\frac{\theta}{\theta+k}\right)^m$$

und

$$\begin{aligned} \mathbb{P}_{\theta} \left( S_n \leq s \right) &= \sum_{m=0}^{s} \mathbb{P}_{\theta} \left( S_n = m \right) = \sum_{m=0}^{s} \frac{n-1}{\theta} \sum_{k=1}^{n-1} (-1)^{k-1} \binom{n-2}{k-1} \binom{\theta}{\theta+k}^{m+1} \\ &= \frac{n-1}{\theta} \sum_{k=1}^{n-1} (-1)^{k-1} \binom{n-2}{k-1} \sum_{m=0}^{s} \left( \frac{\theta}{\theta+k} \right)^{m+1} \\ &= \frac{n-1}{\theta} \sum_{k=1}^{n-1} (-1)^{k-1} \binom{n-2}{k-1} \frac{\theta}{\theta+k} \frac{1 - \left(\frac{\theta}{\theta+k}\right)^{s+1}}{1 - \left(\frac{\theta}{\theta+k}\right)} \\ &= \frac{n-1}{\theta} \sum_{k=1}^{n-1} (-1)^{k-1} \binom{n-2}{k-1} \frac{\theta}{k} \left( 1 - \left( \frac{\theta}{\theta+k} \right)^{s+1} \right) \\ &= \sum_{k=1}^{n-1} (-1)^{k-1} \binom{n-1}{k} \left( 1 - \left( \frac{\theta}{\theta+k} \right)^{s+1} \right) \\ &= 1 - \sum_{k=1}^{n-1} (-1)^{k-1} \binom{n-1}{k} \left( \frac{\theta}{\theta+k} \right)^{s+1} \end{aligned}$$

 $(\operatorname{denn} - \sum_{k=1}^{n-1} (-1)^k {n-1 \choose k} = 1 - (1-1)^{n-1} = 1).$ 

# Kapitel 3

## Selektion

## 3.1 Vorbemerkung: Modelle für den diploiden Fall

Viele Spezies sind *diploid*, besitzen also zwei Kopien jedes (autosomalen) Chromosoms, und typischerweise hat jedes Individuum zwei (verschiedene) Eltern. Gemäß den Mendelschen Regeln erbt ein Kind von jedem Elter jeweils eine Kopie eines der beiden Chromosomen dieses Elters (welche, wird im Idealfall rein zufällig ausgewählt), siehe Abb. 3.1 für ein schematisches Diagramm. Die bis-



Abbildung 3.1: Schematisches Bild zur Vermehrung im diploiden Fall

her betrachteten Populationsmodelle beschreiben den *haploiden* Fall, in dem jedes Individuum nur ein Elter besitzt. Wir diskutieren hier kurz, wie die Modelle und die Ergebnisse an den diploiden Fall angepasst werden können.

Wir betrachten eine Population von N diploiden Individuuen pro Generation und interessieren uns für einen gewissen Ort im Genom, auf einem gewissen Chromosom. Wir nehmen der Einfachheit halber an, dass es nur ein Geschlecht gibt (die Individuuen sind Hermaphroditen) und dass Selbstbefruchtung prinzipiell möglich ist. Ein einfaches stochastisches Modell, das in dieser Situation Zufälligkeit in der Nachkommensverteilung beschreibt, ist das diploide Wright-Fisher-Modell:

**Definition 3.1** (diploides Wright-Fisher-Modell). Die Population besteht aus 2N Chromosomen pro Generation (die in N diploide Individuen aufgeteilt sind, z.B. Chrom. 1 und 2 in Ind. 1, Chrom. 3 und 4 in Ind. 2, etc.). Es findet Zufallspaarung in folgendem Sinn statt: Für jedes der N Kinder unabhängig werden zwei Individuen der Elterngeneration rein zufällig (sagen wir, mit Zurücklegen) als Eltern ausgewählt. Für die beiden Chromosomen in einem Kind bedeutet dies, dass jedes eine

Kopie eines rein zufällig aus den 2N Chromosomen der Elterngeneration gezogenen Chromosoms ist.

Wir ersetzen also N durch 2N im "haploiden Wright-Fisher-Modell" (wie in Abschnitt 1.1) und ändern die Interpretation des Begriffs "Individuum" entsprechend.

Demnach (nach Satz 1.5) konvergiert die Genealogie einer zufälligen Stichprobe von n Chromosomen im diploiden Wright-Fisher-Modell (z.B. aus n/2 diploiden Individuen), wenn man die Zeit mit 2N reskaliert, gegen den Kingman-Koaleszenten (eine Koaleszenten-Zeiteinheit entspricht im Modell mit N Individuen nun etwa 2N Generationen).

**Bemerkung 3.2** (Hardy-Weinberg-Gleichgewicht<sup>1</sup>). Zwei Typen von Chromosomen (im Jargon der Genetik: zwei "Allele") bedeuten, dass es drei mögliche (diploide) Genotypen der Individuen gibt

$$AA$$
,  $Aa$ ,  $aa$ 

(Typ A-Homozygot, Heterozygot, Typ a-Homozygot).

Anhand der Information  $X_r^{(N)} = k$  allein (mit  $X_r^{(N)} =$  Anz. Typ A-Chromosomen in Generation r wie in Abschnitt 1.1) ist die Aufteilung in diploide Individuen nicht festgelegt. Seien

$$(Y_r^{(N)}(AA), Y_r^{(N)}(Aa), Y_r^{(N)}(aa))$$

die Anzahlen diploider Ind. (der drei möglichen Genotypen) in Generation r. Gegeben  $X_r^{(N)} = k = \lfloor 2Np \rfloor$  (mit  $p \in [0, 1]$ ) ist

$$\left(Y_{r+1}^{(N)}(AA), Y_{r+1}^{(N)}(Aa), Y_{r+1}^{(N)}(aa)\right) \sim \text{Multinom}\left(N, p^2, 2p(1-p), (1-p)^2\right)$$

Mit dem Gesetz der großen Zahlen und dem multivariaten zentraler Grenzwertsatz folgt: Gegeben  $\frac{1}{2N}X_r^{(N)} = p$  ist

$$\left(\frac{1}{N}Y_{r+1}^{(N)}(AA), \frac{1}{N}Y_{r+1}^{(N)}(Aa), \frac{1}{N}Y_{r+1}^{(N)}(aa)\right) = (p^2, 2p(1-p), (1-p)^2) + O_{\mathbb{P}}\left(\frac{1}{\sqrt{N}}\right)$$

(wobei  $O_{\mathbb{P}}(1/\sqrt{N})$  einen (zufälligen) Korrekturterm beschreibt, dessen "typische" Größe höchstens  $C/\sqrt{N}$  für eine Konstante  $C < \infty$  ist).

Die Aufteilung in Genotypen (AA, Aa, aa) gemäß  $p^2 : 2p(1-p) : (1-p)^2$  heißt *Hardy-Weinberg-Gleichgewicht*. Wir sehen: In einer sehr großen Population mit "Zufallspaarungen" und Startanteil p von Allel 1 stellt sich nach nur einer Generation für die Verteilung der diploiden Genotypen nahezu das Hardy-Weinberg-Gleichgewicht ein (in der mathematischen Idealisierung im Grenzwert  $N \rightarrow \infty$  stellt es sich exakt ein). Andererseits braucht es Zeit  $\approx \Theta(2N)$ , bis eine spürbare Änderung des Allelanteils auftritt (vgl. die Diskussion am Ende von Abschnitt 1.4), im Vergleich dazu stellt sich HW-Glgw. also "instantan" ein. Dies ist eine Begründung, warum man sich in Populationsgenetik-Modellierung oft auf den haploiden Fall beschränkt (und beschränken darf).

<sup>&</sup>lt;sup>1</sup>Nach Godfrey Harold Hardy (1877–1947) und Wilhelm Weinberg (1862–1937) benannt. G.H. Hardy, Mendelian proportions in a mixed population, *Science* 28 (706), 49–50, (1908); W. Weinberg, Über den Nachweis der Vererbung beim Menschen, *Jahreshefte des Vereins für Vaterländische Naturkunde in Württemberg*, 64, 369–382, (1908).

## 3.2 Vorbemerkung: Deterministische Dynamik

Wir betrachten eine sehr große Population aus diploiden, (der mathematischen Einfachheit halber) hermaphroditischen Individuen und nehmen an, dass es bezüglich eines gewissen Orts im Genom zwei verschiedene Typen, sogenannte Allele, sagen wir A und a, gibt. Jedes Individuum besitzt also zwei Chromosomenkopien, von denen es jeweils eine von jedem seiner beiden Eltern geerbt hat (gemäß den Mendelschen Regeln, d.h. die vererbte Chromosomenkopie wird rein zufällig ausgewählt); es gibt aber keine expliziten Geschlechter, jedes Individuum könnte sich prinzipiell mit jedem anderen paaren. Die zeitliche Entwicklung laufe in diskreten Generationen ab. Sei die Populationsgröße konstant N und

$$(Y_r^{(N)}(AA), Y_r^{(N)}(Aa), Y_r^{(N)}(aa))$$

seien die Anzahlen diploider Individuen (der drei möglichen diploiden Genotypen AA, Aa, aa) in Generation r,

$$X_r^{(N)} = 2Y_r^{(N)}(AA) + Y_r^{(N)}(Aa)$$

die Anzahl A-Chromosomen (unter den insgesamt 2N Chromosomen) in Generation r.

**Erinnerung.** In Abschnitt 3.1 hatten wir diese Situation im Kontext des diploiden Wright-Fisher-Modells (Def. 3.1) betrachtet. Dort hatten wir vorausgesetzt, dass die genetischen Typen für den Fortpflanzungserfolg irrelevant sind und im Modell angenommen, dass die N Kinder der Nachfolgegeneration aus unabhängigen Zügen aus der Elterngeneration (wörtlich: mit Zurücklegen entstehen). Demnach ist gegeben  $X_r^{(N)} = k = 2Np$  (mit einem  $p \in [0, 1] \cap \frac{1}{2N}\mathbb{Z}$ )

$$\left(Y_{r+1}^{(N)}(AA), Y_{r+1}^{(N)}(Aa), Y_{r+1}^{(N)}(aa)\right) \sim \text{Multinom}\left(N, p^2, 2p(1-p), (1-p)^2\right)$$

und  $\mathscr{L}(X_{r+1}^{(N)} | X_r^{(N)} = k) = \operatorname{Bin}(2N, \frac{k}{2N}).$ 

Wenn  $\frac{1}{2N}X_0^{(N)} \Rightarrow p$  für  $N \to \infty$  gilt, so folgt mit dem Gesetz der großen Zahlen

$$\left(\frac{1}{N}Y_1^{(N)}(AA), \frac{1}{N}Y_1^{(N)}(Aa), \frac{1}{N}Y_1^{(N)}(aa)\right) \to \left(p^2, 2p(1-p), (1-p)^2\right)$$

(in Wahrscheinlichkeit) und iterativ ergibt dann sich für jedes  $T \in \mathbb{N}$ ,  $\varepsilon > 0$  (beachte  $p = p^2 + \frac{1}{2}2p(1-p))$ 

$$\mathbb{P}\Big(\Big|\frac{1}{2N}X_{r}^{(N)} - p\Big| < \varepsilon \text{ für } r = 1, 2, \dots, T \Big| \frac{1}{2N}X_{0}^{(N)} = p_{N}\Big) \underset{N \to \infty}{\longrightarrow} 1 \text{ und}$$

$$\mathbb{P}\Big(\Big|\Big(\frac{1}{N}Y_{r}^{(N)}(AA), \frac{1}{N}Y_{r}^{(N)}(Aa), \frac{1}{N}Y_{r}^{(N)}(aa)\Big) - \Big(p^{2}, 2p(1-p), (1-p)^{2}\Big)\Big|\Big| < \varepsilon$$

$$\text{für } r = 1, 2, \dots, T \Big| \frac{1}{2N}X_{0}^{(N)} = p_{N}\Big) \underset{N \to \infty}{\longrightarrow} 1,$$

d.h. es stellt sich Hardy-Weinberg-Gleichgewicht ein (vgl. Bem. 3.2) und die Allelanteile ändern sich (auf dieser Zeitskala) überhaupt nicht.

**(Diploide) Selektion** Wir betrachten nun die allgemeine Situation, in der der diploide Genotyp eines Individuums seinen Überlebens- und Fortpflanzungserfolg beeinflusst. Dazu seien  $w_{AA}, w_{Aa}, w_{aa} \ge 0$  gegeben und wir nehmen an, dass die Chance eines Typ AA-Kinds, bis zum Reproduktionsalter zu überleben und somit als potentielles Elter in Frage zu kommen, proportional zu  $w_{AA}$  ist, etc. Der Vektor  $(w_{AA}, w_{Aa}, w_{aa})$  gibt die relative Fitness der drei Genotypen an.

Wir betrachten wieder eine feste Populationgröße N und bezeichnen mit

$$(Y_r^{(N)}(AA), Y_r^{(N)}(Aa), Y_r^{(N)}(aa))$$

die Anzahlen diploider Individuen sowie mit  $X_r^{(N)} = 2Y_r^{(N)}(AA) + Y_r^{(N)}(Aa)$  die Anzahl A-Chromosomen in Generation r. Für das diploide Wright-Fisher-Modell mit Selektion ist gegeben  $X_r^{(N)} = k = 2Np$  (mit einem  $p \in [0, 1] \cap \frac{1}{2N}\mathbb{Z}$ )

$$\left(Y_{r+1}^{(N)}(AA), Y_{r+1}^{(N)}(Aa), Y_{r+1}^{(N)}(aa)\right) \sim \text{Multinom}\left(N, \frac{p^2 w_{AA}}{w_{\text{ges}}(p)}, \frac{2p(1-p)w_{Aa}}{w_{\text{ges}}(p)}, \frac{(1-p)^2 w_{aa}}{w_{\text{ges}}(p)}\right)$$
(3.1)

mit

$$w_{\text{ges}}(p) \coloneqq p^2 w_{AA} + 2p(1-p)w_{Aa} + (1-p)^2 w_{aa}$$

der "Gesamtfitness" der Population (bei A-Anteil p). Offenbar bildet

$$\left(Y_{r+1}^{(N)}(AA), Y_{r+1}^{(N)}(Aa), Y_{r+1}^{(N)}(aa)\right)_{r \in \mathbb{N}_0}$$

mit Setzung (3.1) eine Markovkette.

Wir können (3.1) folgendermaßen interpretieren (siehe auch Abbildung 3.3): Zunächst wird ein großes Reservoir von  $\gg N$  "Juvenilen" gemäß Zufallspaarung wie im neutralen diploiden Wright-Fisher-Modell gebildet, d.h. wenn der Anteil A-Allele in der Elterngeneration p ist, verhalten sich die Anteile der Genotypen AA, Aa, aa in diesem Reservoir wie  $p^2 : 2p(1-p) : (1-p)^2$ . Dann werden aus diesem Reservoir N Juvenile herausgezogen, wobei ein Typ AA-Juveniles mit Wahrscheinlichkeit proportional zu  $w_{AA}$  gewählt wird, etc., und diese N Herausgezogenen reifen zu den N Erwachsenenindividuen der Folgegeneration heran. (Eine denkbare Assoziation sind Pflanzen, die viele Früchte produzieren, von denen aber wegen Ressourcenbeschränkungen nur wenige tatsächlich keimen.)

Alternativ können wir (3.1) folgendermaßen generieren: o. E. dürfen wir  $w_{AA}, w_{Aa}, w_{aa} \in [0, 1]$ annehmen (es kommt nur auf die Verhältnisse an). Wenn im Modell ein Kind erzeugt werden soll, werden zwei Eltern rein zufällig gewählt und in diesen jeweils rein zufällig ein Chromosom, dessen Typ kopiert wird; wenn dabei ein AA-Kind entsteht, überlebt es mit Wahrscheinlichkeit  $w_{AA}$ , etc. Die Wahrscheinlichkeit, dass ein Kind überlebt, ist dann  $p^2w_{AA} + 2p(1-p)w_{Aa} + (1-p)^2w_{aa} = w_{ges}(p)$  und der Typ eines überlebenden Kinds ist somit

$$AA ext{ mit W'keit } \frac{p^2 w_{AA}}{w_{\text{ges}}(p)}, \ Aa ext{ mit W'keit } \frac{2p(1-p)w_{Aa}}{w_{\text{ges}}(p)}, \ aa ext{ mit W'keit } \frac{(1-p)^2 w_{aa}}{w_{\text{ges}}(p)},$$

wenn der Anteil A-Allele in der Elterngeneration p ist.

Wir wiederholen diesen Mechanismus unabhängig so lange, bis N überlebende Kinder erzeugt wurden und registrieren deren Typen, dies ergibt (3.1).

Analog zum neutralen Fall folgt aus (3.1) mit dem Gesetz der großen Zahlen, dass, sofern  $\frac{1}{2N}X_0^{(N)} =: p_N \rightarrow p$  für  $N \rightarrow \infty$  gilt,

$$\left(\frac{1}{N}Y_{1}^{(N)}(AA), \frac{1}{N}Y_{1}^{(N)}(Aa), \frac{1}{N}Y_{1}^{(N)}(aa)\right) \xrightarrow[N \to \infty]{} \left(\frac{p^{2}w_{AA}}{w_{\text{ges}}(p)}, \frac{2p(1-p)w_{Aa}}{w_{\text{ges}}(p)}, \frac{(1-p)^{2}w_{aa}}{w_{\text{ges}}(p)}\right)$$

und daher auch

$$\frac{1}{2N}X_1^{(N)} \xrightarrow[N \to \infty]{} \frac{1}{2} \frac{2p^2 w_{AA} + 2p(1-p)w_{Aa}}{w_{\text{ges}}(p)} = \frac{p^2 w_{AA} + p(1-p)w_{Aa}}{w_{\text{ges}}(p)}$$

(in Wahrscheinlichkeit).

Sei

$$f(p) \coloneqq \frac{p^2 w_{AA} + p(1-p) w_{Aa}}{w_{\text{ges}}(p)},$$
(3.2)

für  $x_0 \in [0, 1]$  definiert die Funktionsiteration

$$x_{n+1} \coloneqq f(x_n), \quad n \in \mathbb{N}_0 \tag{3.3}$$

eine Folge in [0, 1].

Durch Iteration der Argumente oben ergibt sich (analog zum neutralen Fall), sofern  $\frac{1}{2N}X_0^{(N)} \rightarrow x_0$  gilt, für jedes  $T \in \mathbb{N}, \varepsilon > 0$ 

$$\mathbb{P}\left(\left|\frac{1}{2N}X_{r}^{(N)}-x_{r}\right| < \varepsilon \text{ für } r \leq T \left|\frac{1}{2N}X_{0}^{(N)}=p_{N}\right) \underset{N \to \infty}{\longrightarrow} 1 \text{ und} \tag{3.4}\right) \\
\mathbb{P}\left(\left|\left(\frac{1}{N}Y_{r}^{(N)}(AA), \frac{1}{N}Y_{r}^{(N)}(Aa), \frac{1}{N}Y_{r}^{(N)}(aa)\right) - \left(\frac{x_{r}^{2}w_{AA}}{w_{\text{ges}}(x_{r})}, \frac{2x_{r}(1-x_{r})w_{Aa}}{w_{\text{ges}}(x_{r})}, \frac{(1-x_{r})^{2}w_{aa}}{w_{\text{ges}}(x_{r})}\right)\right| < \varepsilon \text{ für } r \leq T \left|\frac{1}{2N}X_{0}^{(N)}=p_{N}\right) \underset{N \to \infty}{\longrightarrow} 1, \tag{3.5}$$

d.h. der Prozess des A-Anteils konvergiert gegen die deterministische, durch (3.3) beschriebene Folge. Auf dieser (unskalierten) Zeitskala spielen somit für eine große Population Fluktuationen, die aus der Zufälligkeit des Reproduktionsprozesses stammen, nahezu keine Rolle.

Da nur Quotienten der relativen Fitnesswerte  $w_{AA}$ ,  $w_{Aa}$ ,  $w_{aa}$  die Dynamik bestimmen, parametrisiert man (3.3) zur leichteren Interpretation o.E. folgendermaßen um:

$$w_{AA} = 1, \ w_{Aa} = 1 - hs, \ w_{aa} = 1 - s$$
 (3.6)

mit  $s \in [0, 1]$  und  $h \in \mathbb{R}$ . s misst den selektiven Nachteil des Homozygoten aa im Vergleich zu AA (wir nehmen o.E. an, dass AA mindestens so fit ist wie aa, sonst vertausche die Rollen von A und a), h heißt der (Koeffizient des) "Heterozygoteneffekt(s)".

Die biologische Interpretation von (3.2), (3.3) variiert je nach dem Wert von h:

 $\begin{array}{lll} h=1 & : & \text{Allel } A \text{ ist rezessiv} \\ h=0 & : & \text{Allel } A \text{ ist dominant} \\ 0 < h < 1 & : & \text{unvollständige Dominanz} \\ h < 0 & : & \text{Überdominanz} (Aa ,, am fittesten") \\ h > 1 & : & \text{Unterdominanz} \end{array}$ 

(der Fall h = 1/2: "additiver Fitnesseffekt" ist mathematisch besonders angenehm und häufig für  $s \approx 0$  gerechtfertigt)

In der Parametrisierung (3.6) lautet (3.2)

$$f(p) = \frac{p^2 + p(1-p)(1-hs)}{p^2 + 2p(1-p)(1-hs) + (1-p)^2(1-s)}$$

und die Änderung über eine Generation ist

$$d(p) = f(p) - p = \frac{sp(1-p)(ph + (1-p)(1-h))}{p^2 + 2p(1-p)(1-hs) + (1-p)^2(1-s)}.$$

**Langzeitverhalten** Verändern sich die Allelhäufigkeiten im Laufe der Zeit? Wenn ja, wie? Wie sieht es nach sehr langer Zeit aus? Wird sich das *A*-Allel durchsetzen?

Das Langzeitverhalten von  $(x_r)_{r \in \mathbb{N}_0}$  aus (3.3) wird (hauptsächlich) von h bestimmt (wir nehmen s > 0 an, sonst ist  $x_r \equiv x_0$ ). Offensichtlich gilt d(0) = d(1) = 0 und, sofern  $h \neq 1/2$ , auch  $d(p_*) = 0$  mit  $p_* = \frac{1-h}{1-2h}$ .

Falls  $0 \le h \le 1$ : Es gilt d(p) > 0 für  $0 , somit <math>x_r \nearrow 1$  für  $r \to \infty$  sobald  $x_0 > 0$ . Man spricht von "gerichteter Selektion": Das "fittere" Allel A verdrängt a.

Falls h < 0: Es gilt  $0 < p_* < 1$  und d(p) > 0 für 0 , <math>d(p) < 0 für  $p_* . Daher gilt <math>x_r \nearrow p_*$  für  $r \rightarrow \infty$  falls  $0 < x_0 \le p_*$  und  $x_r \searrow p_*$  falls  $p_* < x_0 < 1$ . Man spricht von "balancierender Selektion": Beide Allele bleiben langfristig in der Population erhalten, das genaue Verhältnis hängt von h ab (das hier die Stärke der Überdominanz misst).

Falls h > 1: Es gilt  $0 < p_* < 1$  und d(p) < 0 für 0 , <math>d(p) > 0 für  $p_* . Daher gilt <math>x_r \ge 0$  für  $r \to \infty$  falls  $0 < x_0 < p_*$  und  $x_r \nearrow 1$  falls  $p_* < x_0 < 1$ . Man spricht von "disruptiver Selektion": Langfristig setzt sich einer der beiden Typen durch (es sei denn, man beginnt mit genau  $x_0 = p_*$ ), welcher, hängt von der Startbedingung ab. Zwei Populationen, deren Startanteil an A sehr ähnlich ist, aber einer ober- und einer unterhalb von  $p_*$ , werden sich auf lange Sicht auseinander entwickeln.

**Beispiel.** Das *medionigra*-Allel (*a*) wurde in einer Population von *Callimorpha dominula* (ein Nachtfalter, deutsch Schönbär) in der Umgebung von Oxford recht intensiv studiert (*a* verändert im Vergleich zum "Wildtyp" die Flügelfärbung). Abbildung 3.2 zeigt den beobachteten *a*-Anteil für die Jahre 1939–1972 und eine angepasste deterministische Folge ( $x_r$ ), erzeugt aus (3.3), (3.6) mit s = 0,1, h = 0,5 (d.h.  $w_{AA} = 1$ ,  $w_{Aa} = 0,95$ ,  $w_{aa} = 0,9$ ). Die Werte sind Kap. 3 des Buches John H. Gillespie, *Population genetics : a concise guide*, Johns Hopkins Univ. Press, 1998, entnommen. Die Modellkurve passt – zumindest dem Augenschein nach – recht gut.

(Das beweist allerdings nicht, dass tatsächlich gerichtete Selektion für die beobachteten Änderungen des *a*-Anteils verantwortlich ist — es könnte andere Effekte geben, die z.T. kontrovers in der Literatur diskutiert wurden, siehe z.B. E.B. Ford, P.M. Sheppard, The medionigra polymorphism of Panaxia dominula. *Heredity* 24, 112—134, 1969.)

**Beispiel.** Ein "klassisches Lehrbuchbeispiel" für den Effekt balancierender Selektion ist die Verbreitung der Sichelzellenanämie in Gebieten mit endemischer Malaria. Das *a*-Allel (in unserer Notation) ruft in der homozygoten Form eine Verformung der roten Blutkörperchen hervor, die zum Krankheitsbild der Sichelzellenanämie führt. *Aa*-Individuen haben nur eine milde Form der Krankheit und



Abbildung 3.2: Beobachteter Anteil des *a*-Allels 1939–1972 in *Callimorpha dominula* und Modellvorhersage



Abbildung 3.3: Schematische Darstellung eines Lebenszyklus im Modell mit Selektion

sind zugleich vor gewissen Formen der Malaria geschützt. Sie haben daher in Regionen, in denen Malaria weit verbreitet ist – zumal, wenn keine Therapien zur Verfügung stehen – gegenüber AA-Homozygoten effektiv einen Vorteil. J.H. Gillespie, a.a.O., Kap. 3.3 berichtet, dass die Wahl s = 1, h = -0.17 und somit  $p_* \approx 0.87$ , d.h. ein Anteil a von  $1 - p_* \approx 0.13$  relativ gut zu Beobachtungswerten in West- und Zentralafrika passt.

## 3.3 (2-Typ) Moran-Modell mit (gerichteter) Selektion

Wir schränken uns (aus Zeit- und Platzgründen) im Folgenden bei der detaillierteren Diskussion von Modellen, die sowohl Selektion als auch Gendrift enthalten, auf das Moran-Modell (mit gerichteter Selektion) ein, da hier (im Ggs. zum Wright-Fisher-Modell) wegen der einfacheren Struktur der Übergangsdynamik verschiedene explizite Rechnungen möglich sind.

**Definition 3.3** ((2 Typ-)Moran-Modell mit gerichteter Selektion). Sei  $s \ge 0$ . Man betrachtet eine Population von konstant N (haploiden) Individuen, die zwei mögliche Typen A und a haben. Typ A-Ind. vermehren sich mit Rate 1 + s, Typ a-Ind. vermehren sich mit Rate 1 und bei jedem Vermehrungsereignis wird ein rein zufällig gezogenes Individuum durch den gerade erzeugten Nachkommen ersetzt (der den Typ des Elters erbt und, sagen wir, man kann durch sein eigenes Kind ersetzt werden). Sei

$$X_t^{(N)} = \#$$
Typ A-Ind. zur Zeit  $t_s$ 



Abbildung 3.4: Zeitliche Entwicklung des A-Anteils  $x_r$  bei verschiedenen Startwerten. Links s=0,1, h=0,5 (gerichtete Selektion), rechts: s=1, h=-0,17 (balancierende Selektion)

 $(X_t^{(N)})_{t \ge 0}$  ist eine zeitkontinuierliche Markovkette mit Werten in  $\{0, 1, \dots, N\}$  und Sprungratenmatrix

$$q_{ij} = \begin{cases} (1+s)i\frac{N-i}{N}, & j = i+1, \\ (N-i)\frac{i}{N}, & j = i-1, \\ -(2+s)\frac{i(N-i)}{N}, & j = i, \\ 0, & \text{sonst.} \end{cases}$$

**Bemerkung.** 1. Wörtlich unterscheiden sich in der Formulierung von Def. 3.3 die Typen bezüglich der Vermehrungsgeschwindigkeit (man nennt dies auch "fecundity selection").

Wir könnten dieselben Sprungraten allerdings auch als Unterschiede in der Überlebensfähigkeit ("viability selection") interpretieren: Nehmen wir an, Typ A-Individuen sterben mit Rate 1 (d.h. ihre Lebenszeit ist Exp(1)-verteilt), Typ a-Individuen sterben mit Rate 1 + s (d.h. ihre Lebenszeit ist  $\sim Exp(1 + s)$ ) und das gestorbene Ind. wird instantan durch den Nachkommen eines rein zufällig gezogenen Ind. ersetzt, so liefert dies dieselbe Dynamik von  $(X_t^{(N)})_{t\geq 0}$ .

**Bemerkung 3.4** (Graphische Konstruktion des Moran-Modells mit Selektion). Man kann die graphische Konstruktion aus Abschnitt 1.3 auf den Fall mit Selektion übertragen:

Für jedes geordnete Paar  $(i, j), i, j \in \{1, ..., N\}, i \neq j$  sei  $(N_t^{(i,j)})_{t\geq 0}$  ein Poissonprozess auf  $\mathbb{R}_+$  mit Rate  $\frac{1}{N}$ , u.a. für verschiedene Paare. Zusätzlich sei für jedes Paar  $(S_t^{(i,j)})_{t\geq 0}$  ein weiterer, unabhängiger Poissonprozess auf  $\mathbb{R}_+$  mit Rate s/N.

Wie vorher erzeugt zu den Sprungzeiten von  $(N_t^{(i,j)})_{t\geq 0}$  Individuum *i* einen Nachkommen, der Individuum *j* ersetzt, wir legen in der Konstruktion zu einem solchen Zeitpunkt einen Pfeil von *i* nach *j*. Zu den Sprungzeiten von  $(S_t^{(i,j)})_{t\geq 0}$  legen wir "s-Pfeile", zu einem solchen Zeitpunkt erzeugt Individuum *i* einen Nachkommen, der *j* verdrängt, aber nur, wenn *i* aktuell Typ *A* hat (s-Pfeile sind also nur "wirksam", wenn der Typ am "Pfeilschaft" *A* ist).

[Bild an der Tafel]

Offenbar sind 0 und N (wiederum) absorbierende Zustände von  $X^{(N)}$ . Die Fixierungswahrscheinlichkeiten sind für s > 0 (für  $s \downarrow 0$  ergibt sich z.B. mit der Regel von l'Hospital i/N) wie folgt:

**Satz 3.5.** Set  $\tau_k = \inf\{t \ge 0 : X_t^{(N)} = k\}$ . Es gilt

$$\mathbb{P}_i(\tau_N < \tau_0) = \frac{1 - (1 + s)^{-i}}{1 - (1 + s)^{-N}}$$

(Wir schreiben  $\mathbb{P}_i(\cdot)$  für  $\mathbb{P}(\cdot | X_0^{(N)} = i)$  und unterdrücken die Abhängigkeit von N in der Notation von  $\tau_k$ .)

**Beobachtung 3.6.** Für s > 0 (fest),  $N \gg 1$  ist demnach

$$h(1) \approx 1 - \frac{1}{1+s} = \frac{s}{1+s} \quad (\approx s \text{ für } s \text{ klein});$$

für  $N \gg 1$ ,  $s \ll 1$  mit  $\sigma \coloneqq Ns \in (0, \infty)$ ,  $x \in (0, 1)$  ist

$$h(\lfloor Nx \rfloor) = \frac{1 - (1 + \sigma/N)^{-\lfloor Nx \rfloor}}{1 - (1 + \sigma/N)^{-N}} \approx \frac{1 - e^{-\sigma x}}{1 - e^{-\sigma}}$$

*Beweis von Satz 3.5.* Sei  $h(i) := \mathbb{P}_i(\tau_N < \tau_0)$ , offenbar gilt h(0) = 0, h(N) = 1, Zerlegung gemäß dem ersten Sprung zeigt

$$h(i) = \frac{1}{(2+s)\frac{i(N-i)}{N}} \left( (1+s)\frac{i(N-i)}{N}h(i+1) + \frac{i(N-i)}{N}h(i-1) \right)$$
  
=  $\frac{1+s}{2+s}h(i+1) + \frac{1}{2+s}h(i-1)$  für  $1 \le i < N$ .

Dies ist ein homogenes, lineares Differenzengleichungssystem 2. Ordnung, demnach gilt

$$h(i) = c_1 u_1^i + c_2 u_2^i$$

mit  $u_{1/2}$  Lösungen von  $(1 + s)u^2 - (2 + s)u + 1 = 0$ , d.h.

$$u_{1/2} = \frac{(2+s) \pm \sqrt{(2+s)^2 - 4(1+s) \cdot 1}}{2(1+s)} = \frac{(2+s) \pm s}{2(1+s)} = \{1, \frac{1}{1+s}\}$$

und Koeffizienten  $c_1, c_2 \in \mathbb{R}$ , deren Werte  $c_1 = 1/(1 - (1 + s)^{-N}), c_2 = -1/(1 - (1 + s)^{-N})$  sich aus den Randbedingungen ergeben.

Stellen wir uns vor, in einer bisher homogenen *a*-Population ist gerade eine vorteilhafte Mutation *A* aufgetreten (weitere Mutationen sieht unser Modell zunächst nicht vor). Obwohl *A* selektiv bevorzugt ist, könnte der Typ *A* aufgrund der zufälligen Fluktuationen im Reproduktionsprozess trotzdem wieder verschwinden. Beobachtung 3.6 gibt die Wahrscheinlichkeit an, dass *A* tatsächlich fixiert. Der folgende Satz geht der Frage nach, wie lange es dauert, bis eine solche "frisch aufgetretene" vorteilhafte Mutation fixiert (gegeben, dass dies tatsächlich passiert). **Satz 3.7.** Set s > 0 fest, starte mit  $X_0^{(N)} = 1$ . Für  $N \to \infty$  gilt

$$\frac{s}{2\log N}\tau_N \to 1 \quad in \ Verteilung \ unter \mathbb{P}(\cdot \mid \tau_N < \tau_0),$$

 $d.b. \mathbb{P}(\tau_N \leq b \log N \mid \tau_N < \tau_0) \underset{N \to \infty}{\longrightarrow} \begin{cases} 0, & \textit{falls } b < 2/s, \\ 1, & \textit{falls } b > 2/s. \end{cases}$ 

*Bew.idee für Satz 3.7:* Wir zerlegen den Pfad (auf dem Weg von  $X_0^{(N)} = 1$  nach  $X_{\tau_N}^{(N)} = N$ ) in drei Phasen.

1. Solange  $X^{(N)}$  klein ist (sagen wir,  $\leq \varepsilon N$ ), geschehen die Sprünge

$$x \to x + 1$$
 mit Rate  $(1+s)i\frac{N-i}{N} \approx (1+s)i$ ,  
 $x \to x - 1$  mit Rate  $i\frac{N-i}{N} \approx i$ ,

d.h.  $X^{(N)}$  verhält sich "beinahe" wie ein superkritischer binärer Galton-Watson(-Verzweigungs)-Prozess  $(Y_t)$  (vgl. Def. 3.9 unten). Da

 $Y_t \approx$  eine zufällige Konstante  $\times e^{st}$ 

(siehe Proposition 3.10 unten), dauert es etwa Zeit  $\approx \frac{1}{s} \log N$ , bis  $X^{(N)}$  auf  $\varepsilon N$  angewachsen ist.

2. Sobald  $X^{(N)}$  "mittelgroß" ist (sagen wir,  $\varepsilon N < X_t^{(N)} < (1 - \varepsilon)N$ ) folgt der Pfad des Anteilsprozesses  $\frac{1}{N}X_t^{(N)}$  nahezu einer deterministischen, "makroskopischen" Dynamik (vgl. auch Diskussion in Abschnitt 3.2).

Diese Phase ist daher "kurz" im Vergleich zu Phasen 1 und 3 (ihre Länge divergiert nicht mit  $N \rightarrow \infty$ )

 Sobald X<sup>(N)</sup> "groß" ist (sagen wir, ≥ (1-ε)N) können wir analog zu Phase I mit einem subkritischem Galton-Watson-Prozess vergleichen (oder wir vertauschen Rollen von a und A und kehren Zeit um, um wörtlich auf Beweisschritte für Phase I zurückzugreifen).

Die Phase dauert ebenfalls  $\approx \frac{1}{s} \log N$ .

Offensichtlich benötigen wir für den Beweis von Satz 3.7 Informationen über das Verhalten von Markovketten, die wir auf das Erreichen gewisser Zustände bedingen. Eine allgemeine Antwort gibt die Doob-Transformation.

Sei X zeitkontinuierliche Markovkette auf E (E endl. oder abz.b. unendliche Menge) mit (konservativer) Ratenmatrix  $Q = (q_{ij})$  (und wir nehmen an  $\sup_x -q_{xx} < \infty$ , d.h. die Gesamtsprungrate ist global beschränkt), sei  $E_0 \subset E$  eine endliche Menge von absorbierenden Zuständen, Q sei irreduzibel auf  $E \smallsetminus E_0$  und jeder Punkt  $z \in E_0$  mit pos. W'keit erreichbar von  $E \searrow E_0$  aus, für

$$\tau \coloneqq \inf\{t \ge 0 : X_t \in E_0\}$$

gelte  $\mathbb{P}_x(\tau < \infty) = 1$  für alle  $x \in E$ .

Sei  $z_0 \in E_0$  und  $h : E \rightarrow [0, \infty)$  (beschränkt) mit  $h(z_0) = 1$ , h(z) = 0 für  $z \in E_0 \setminus \{z_0\}$ , h(Q)-harmonisch in  $E \setminus E_0$ , d.h.

$$\sum_{y} q_{xy} h(y) = 0 \quad \text{für } x \in E \setminus E_0.$$
(3.7)

**Lemma 3.8** (Doob-Transformation). Sei  $X_0 \in E \setminus E_0$  (f.s).  $(X)_{t \ge 0}$  bedingt auf  $\{X_\tau = z_0\}$  ist verteilt wie die zeitkontinuierliche Markovkette  $\widetilde{X}$  mit Sprungraten  $\widetilde{q}_{ij} = \frac{h(j)}{h(i)}q_{ij}$ .

(Beachte:  $\sum_{j \neq i} \frac{h(j)}{h(i)} q_{ij} = \frac{1}{h(i)} (-h(i)q_{ii}) = -q_{ii}$  nach Voraussetzung, d.h.  $(\widetilde{q}_{ij})$  ist ebenfalls eine (konservative) Ratenmatrix und  $\widetilde{X}$  und X haben in jedem Punkt dieselbe Gesamtsprungrate).

*Beweis.* Sei  $p_{ij} \coloneqq \frac{q_{ij}}{-q_{ii}}$  für  $i \neq j \in E$ ,  $p_{ii} \coloneqq 0$  die Übergangsmatrix der (zeitdiskreten) Skelettkette X' von X (d.h. im Zustand i verbringt X eine  $Exp(-q_{ii})$ -verteilte Wartezeit und springt dann wie X' von i aus gemäß  $p_{ij}, j \in E$  in einen neuen Zustand j) und analog  $\widetilde{p}_{ij} \coloneqq \frac{\widetilde{q}_{ij}}{-\widetilde{q}_{ii}}$  für  $i \neq j \in E$ ,  $\widetilde{p}_{ii} \coloneqq 0$  die Übergangsmatrix der (zeitdiskreten) Skelettkette  $\widetilde{X}'$  von  $\widetilde{X}$ . Dann gilt auch

$$\widetilde{p}'_{ij} = \frac{h(j)}{h(i)} p'_{ij}.$$

Es ist  $\{X_{\tau} = z_0\} = \{X'_{\tau'} = z_0\}$  mit  $\tau' := \min\{k \in \mathbb{N}_0 : X'_k \in E_0\}$  und nach Definition folgt aus (3.7) für  $x \in E \smallsetminus E_0$ 

$$-q_{xx}h(x) = \sum_{y \neq x} q_{xy}h(y), \quad \text{ somit } h(x) = \sum_{y} p_{xy}h(y),$$

d.h. h ist auch p-harmonisch in  $E \\ E_0$ . Das Standard-Argument der Zerlegung nach dem ersten Sprung zeigt, dass auch  $i \mapsto \mathbb{P}_i(X'_{\tau'} = z_0) p$ -harmonisch in  $E \\ E_0$  ist, auf  $E_0$  stimmt dies natürlich mit h überein, daher gilt

$$h(i) = \mathbb{P}_i(X'_{\tau'} = z_0) = \mathbb{P}_i(X_\tau = z_0).$$

Siehe z.B. Klenke, [Kle20, Satz 19.7] für die Eindeutigkeit des Dirichlet-Problems (wörtlich dort im Fall endlichen Zustandsraums).

Für  $\ell \in \mathbb{N}, x_0, x_1, \dots, x_{\ell-1} \in E \smallsetminus E_0, x_\ell \coloneqq z_0$  ist

$$\mathbb{P}_{x_0} \Big( X'_0 = x_0, X'_1 = x_1, \dots, X'_{\ell} = x_{\ell} \, \Big| \, X'_{\tau'} = z_0 \Big) = \frac{1}{h(x_0)} \prod_{i=1}^{\ell} p'_{x_{i-1}x_i} = \frac{1}{h(x_0)} \prod_{i=1}^{\ell} \frac{h(x_{i-1}) \widetilde{p}'_{x_{i-1}x_i}}{h(x_i)} \\ = \frac{1}{h(x_{\ell})} \prod_{i=1}^{\ell} \widetilde{p}'_{x_{i-1}x_i} = \mathbb{P}_{x_0} \Big( \widetilde{X}'_0 = x_0, \widetilde{X}'_1 = x_1, \dots, \widetilde{X}'_{\ell} = x_{\ell} \Big),$$

denn  $h(x_{\ell}) = h(z_0) = 1$ . Demnach gilt die Behauptung für die Skelettketten, da X und  $\widetilde{X}$  in jedem Punkt dieselbe Gesamtsprungrate haben, gilt sie auch für X und  $\widetilde{X}$ . Für  $X^{(N)}$ , das Moran-Modell mit gerichteter Selektion (Def. 3.3) und  $E_0 = \{0, N\}, z_0 = N$  ergibt sich mit Satz 3.5 aus Lemma 3.8:

$$\widetilde{q}_{i,i+1} = \frac{1 - (1+s)^{-(i+1)}}{1 - (1+s)^{-i}} i \frac{N-i}{N} (1+s) = i \left(1 - \frac{i}{N}\right) \left(1 + \frac{s}{1 - (1+s)^{-i}}\right),$$
  
$$\widetilde{q}_{i,i-1} = \frac{1 - (1+s)^{-(i-1)}}{1 - (1+s)^{-i}} \left(N - i\right) \frac{i}{N} = i \left(1 - \frac{i}{N}\right) \left(1 + \frac{s(1+s)^{-i+1}}{1 - (1+s)^{-i}}\right)$$

Für unser Vergleichsargument benötigen wir (gewisse) Verzweigungsprozesse (in Def. 2.12 hatten wir bereits den Yule-Prozess betrachtet, einen Verzweigungsprozess, in dem Individuen niemals sterben).

**Definition 3.9** (Binärer (zeitkontinuierlicher) Galton-Watson-Prozess). Wir betrachten eine Population von Individuen, die sich jeweils unabhängig mit Rate  $\lambda > 0$  verdoppeln und mit Rate  $\mu > 0$ sterben (es gibt keine Restriktionen an die Größe der Population); bezeichne  $Y_t$  die Anzahl Individuen zur Zeit t.

 $(Y_t)_{t\geq 0}$  ist eine zeitkontinuierliche Markovkette mit Sprungratenmatrix

$$q_{ij}^{\rm GW} = \begin{cases} \lambda i, & j = i+1, \\ \mu i, & j = i-1, \\ -(\lambda + \mu)i, & j = i, \\ 0, & \text{sonst.} \end{cases}$$

Y heißt ein (binärer zeitkontinuierlicher) Galton-Watson-Prozess.

Man sieht aus der Form der Sprungraten, dass Y die Verzweigungseigenschaft

$$\mathscr{L}(Y_t | Y_0 = k + j) = \mathscr{L}(Y_t | Y_0 = k) * \mathscr{L}(Y_t | Y_0 = j)$$

besitzt (anschaulich: wenn man die Startpopulation in zwei Gruppen aufteilt, dann kann man die beiden Gruppen sich für Zeit *t* unabhängig getrennt entwickeln lassen und die sich daraus ergebenden Populationen dann wieder zusammenfügen, dies ändert die Verteilung der Gesamtgröße nicht).

#### **Proposition 3.10.** *Es gilt*

$$\mathbb{E}_{k}[Y_{t}] = e^{(\lambda-\mu)t}k, \quad \operatorname{Var}_{k}[Y_{t}] = \frac{\lambda+\mu}{\lambda-\mu} \left(e^{2(\lambda-\mu)t} - e^{(\lambda-\mu)t}\right)k \tag{3.8}$$

(die Formel für die Varianz gilt im Fall  $\lambda \neq \mu$ , falls  $\lambda = \mu$  gilt, ist  $\operatorname{Var}_k[Y_t] = 2\lambda kt$ ). Mit  $W_t := e^{-(\lambda-\mu)t}Y_t$  ist  $(W_t)_{t\geq 0}$  ein nicht-negatives Martingal, demnach

$$W_t \to W_\infty$$
 f.s. für ein  $W_\infty \ge 0$ .

Falls  $\lambda > \mu$  gilt, ist  $(W_t)_{t \ge 0} L^2$ -beschränkt, insbesondere gleichgradig integrierbar, und es gilt  $\{W_{\infty} > 0\} = \{Y_t \neq 0 \ \forall t \ge 0\}$  f.s. (d.h. falls der Prozess überlebt, wächst er auch exponentiell mit Rate  $\lambda - \mu$ ) sowie

$$h(k) := \mathbb{P}_k(Y_t = 0 \text{ schließlich}) = (\mu/\lambda)^k, \quad y \in \mathbb{N}.$$

**Bemerkung 3.11.** Ein binärer zeitkontinuierlicher Galton-Watson-Prozess  $Y = (Y_t)_{t\geq 0}$  mit  $\lambda > \mu$ heißt *superkritisch*, Proposition 3.10 zeigt, dass Y dann mit positiver Wahrscheinlichkeit überlebt und in diesem Fall unbeschränkt (exponentiell) wächst.

Falls  $\lambda < \mu$  gilt, heißt Y subkritisch, in diesem Fall gilt  $\mathbb{P}_k(Y_t > 0) \leq \mathbb{E}_k[Y_t] = ke^{(\lambda - \mu)t} \to 0$  für  $t \to \infty$  und wegen  $\{Y_t = 0\} \subset \{Y_u = 0 \text{ für alle } u \geq t\}$  also  $\mathbb{P}_k(Y_t = 0 \text{ schließlich}) = 1$ .

Im Fall  $\lambda = \mu$  heißt *Y kritisch*, man kann mit etwas mehr Aufwand zeigen, dass auch dann gilt  $\mathbb{P}_k(Y_t = 0 \text{ schließlich}) = 1.$ 

*Beweisskizze für Proposition 3.10.* Für das 1. Moment beachten wir mit f(y) = y, dass der Generator  $Lf(y) = \lambda y \cdot (y+1) + \mu y \cdot (y-1) - (\lambda + u)y = (\lambda - \mu)y$  erfüllt, also

$$\frac{d}{dt}\mathbb{E}_k[Y_t] = (\lambda - \mu)\mathbb{E}_k[Y_t], \quad \mathbb{E}_k[Y_0] = k$$

gemäß Kolmogorovs Rückwärtsgleichung mit Lösung  $\mathbb{E}_k[Y_t] = e^{(\lambda - \mu)t}k$ .

Zusammen mit der Markoveigenschaft ergibt sich

$$\mathbb{E}_k \Big[ Y_{t+h} \, \big| \, \sigma(Y_u, u \le t) \Big] = e^{(\lambda - \mu)h} Y_t,$$

d.h.  $W_t = e^{-(\lambda - \mu)t} Y_t$  ist ein Martingal.

Für das 2. Moment bzw. die Varianz betrachten wir  $f(k) = k^2$ ,

$$Lf(k) = k(\lambda(k+1)^{2} + \mu(k-1)^{2} - (\lambda+\mu)k^{2}) = k(2\lambda k + \lambda - 2\mu k + \mu) = 2(\lambda-\mu)k^{2} + (\lambda+\mu)k^{2}$$

Damit ist

$$\frac{d}{dt}\mathbb{E}_1[Y_t^2] = 2(\lambda - \mu)\mathbb{E}_1[Y_t^2] + (\lambda + \mu)\mathbb{E}_1[Y_t] = 2(\lambda - \mu)\mathbb{E}_1[Y_t^2] + (\lambda + \mu)e^{(\lambda - \mu)t}$$

(und  $\mathbb{E}_1[Y_0^2] = 1$ ), Variation der Konstanten liefert

$$\begin{split} \mathbb{E}_{1}[Y_{t}^{2}] &= e^{2(\lambda-\mu)t} \mathbb{E}_{1}[Y_{0}^{2}] + \int_{0}^{t} e^{2(\lambda-\mu)(t-u)} (\lambda+\mu) e^{(\lambda-\mu)u} \, du \\ &= e^{2(\lambda-\mu)t} + (\lambda+\mu) e^{2(\lambda-\mu)t} \int_{0}^{t} e^{-(\lambda-\mu)u} \, du = e^{2(\lambda-\mu)t} + (\lambda+\mu) e^{2(\lambda-\mu)t} \frac{1}{\lambda-\mu} (1-e^{-(\lambda-\mu)t}) \\ &= e^{2(\lambda-\mu)t} + \frac{\lambda+\mu}{\lambda-\mu} (e^{2(\lambda-\mu)t} - e^{(\lambda-\mu)t}) = \frac{2\lambda}{\lambda-\mu} e^{2(\lambda-\mu)t} - \frac{\lambda+\mu}{\lambda-\mu} e^{(\lambda-\mu)t}. \end{split}$$

Somit

$$\operatorname{Var}_{1}[Y_{t}] = \mathbb{E}_{1}[Y_{t}^{2}] - \left(\mathbb{E}_{1}[Y_{t}]\right)^{2} = \left(\frac{2\lambda}{\lambda-\mu} - 1\right)e^{2(\lambda-\mu)t} - \frac{\lambda+\mu}{\lambda-\mu}e^{(\lambda-\mu)t} = \frac{\lambda+\mu}{\lambda-\mu}\left(e^{2(\lambda-\mu)t} - e^{(\lambda-\mu)t}\right),$$

die Verzweigungseigenschaft liefert  $\operatorname{Var}_{k}[Y_{t}] = k \operatorname{Var}_{1}[Y_{t}] (\operatorname{denn} \mathscr{L}(Y_{t} | Y_{0} = 0) =^{d} Y_{t}^{(1)} + \dots + Y_{t}^{(k)},$ wobei  $Y^{(i)}$  unabhängige Kopien jeweils mit Startwert  $Y^{(i)} = 1$  sind).

Die Formel für den Fall  $\lambda = \mu$  kann man beispielsweise erhalten, indem man in obigem  $\lambda \searrow \mu$  betrachtet (oder die entsprechende Differentialgleichung für  $\frac{d}{dt}\mathbb{E}_1[Y_t^2] = 2\lambda\mathbb{E}_1[Y_t]$  direkt löst).

Im Fall  $\lambda > \mu$  zeigt die Formel für das 2. Moment, dass  $\sup_{t\geq 0} \mathbb{E}\left[(e^{-(\lambda+\mu)t}Y_t)^2\right] < \infty$  gilt.

Zur Aussterbewahrscheinlichkeit:  $h(y) := \mathbb{P}_y(Y_t = 0 \text{ schließlich})$  löst (zerlege gemäß dem ersten Sprung, benutze die starke Markov-Eigenschaft)

$$h(y) = \frac{\lambda y}{(\lambda + \mu)y} h(y+1) + \frac{\mu y}{(\lambda + \mu)y} h(y-1) = \frac{\lambda}{\lambda + \mu} h(y+1) + \frac{\mu}{\lambda + \mu} h(y-1), \quad y \in \mathbb{N}$$

mit Randbedingung h(0) = 1 (und  $\lim_{y\to\infty} h(y) = 0$ ), die (eind.) Lösung ist  $h(y) = (\mu/\lambda)^y$ .

Offensichtlich gilt

$$\{Y_t = 0 \text{ schließlich}\} \subset \{W = 0\} \text{ stets},\$$

um die Behauptung

$$\{W_{\infty} > 0\} = \{Y_t \neq 0 \ \forall \ t \ge 0\}$$
 f.s

zu zeigen, müsste man etwa zeigen, dass  $\mathbb{P}_k(W_{\infty} = 0) = \mathbb{P}_k(Y_t = 0 \text{ schließlich}) = (\mu/\lambda)^k$  gilt.

Siehe dazu beispielsweise Thm. III.2 in K.B. Athreya, P. Ney, Branching processes, Springer, 1972 oder Prop. 5.6 und Cor. 5.7 in Russell Lyons, Yuval Peres, Probability on trees and networks, Cambridge University Press, 2016+ (auch elektronisch unter http://mypage.iu.edu/~rdlyons/prbtree/ prbtree.html).

Mit Lemma 3.8 (Doob-Transformation) ist die Sprungratenmatrix von Y (für  $\lambda = 1 + s, \mu = 1$ ) bedingt auf  $\{Y_t \neq 0 \forall t \ge 0\}$  (=  $\{Y_t \rightarrow \infty\}$  =  $\{W_\infty > 0\}$ ):

$$\widetilde{q}_{ij}^{\rm GW} = \frac{1 - (1 + s)^{-j}}{1 - (1 + s)^{-i}} q_{ij}^{\rm GW}$$

(Wörtlich wäre hier noch ein kleines Approximationsargument notwendig, da " $\infty$ " eigentlich kein Punkt des Zustandsraums von Y ist: Bedinge zunächst darauf,  $N \gg 1$  vor 0 zu treffen, dann lasse  $N \rightarrow \infty$ .)

**Beobachtung 3.12.** Sei  $\varepsilon > 0$ ,  $\tilde{\tau}_{\varepsilon N} := \inf\{t \ge 0 : Y_t \ge \varepsilon N\}$ , es gilt

$$\frac{s}{\log N}\tilde{\tau}_{\varepsilon N}\underset{N\to\infty}{\longrightarrow}1\quad\text{f.s. auf }\{W_{\infty}>0\}$$

Intuitiv/anschaulich klar:  $Y_t \approx e^{st} W_\infty$ , also sollte  $\tilde{\tau}_{\varepsilon N} \approx \frac{1}{s} (\log N + \log \varepsilon - \log W_\infty)$  sein.

*Beweis.* Sei  $T_{\delta} := \sup\{t : \left|\frac{Y_t}{e^{st}W_{\infty}} - 1\right| > \delta\}$  (bzw.  $T_{\delta} = \infty$  auf  $\{W_{\infty} = 0\}$ ), es gilt  $\{T_{\delta} < \infty\}$  (f.s.) auf  $\{W_{\infty} > 0\}$ .

Wähle  $\delta < \varepsilon$ . Es gilt

$$\big\{ \widetilde{\tau}_{\varepsilon N} > \tfrac{1+\delta}{s} \log N \big\} \subset \big\{ Y_{\tfrac{1+\delta}{s} \log N} < \varepsilon N \big\} \subset \big\{ T_{\delta} \geq \tfrac{1+\delta}{s} \log N \big\} \cup \big\{ W_{\infty} \leq N^{-\delta} \big\},$$

 $\operatorname{denn} \operatorname{auf} \left\{ T_{\delta} < \frac{1+\delta}{s} \log N \right\} \cap \left\{ W_{\infty} > N^{-\delta} \right\} \operatorname{ist} Y_{\frac{1+\delta}{s} \log N} \ge (1-\delta) e^{(1+\delta) \log N} W_{\infty} \ge (1-\delta) N > \varepsilon N.$  Weiter ist

 $\big\{\tilde{\tau}_{\varepsilon N} < \tfrac{1-\delta}{s}\log N\big\} \subset \Big\{\sup_{t \leq \log\log N} Y_t \geq \varepsilon N\Big\} \cup \big\{T_\delta \geq \log\log N\big\} \cup \big\{W_\infty \geq \tfrac{\varepsilon}{2(1+\delta)}N^\delta\big\},$ 

denn auf  $\{T_{\delta} < \log \log N\} \cap \{W_{\infty} < \frac{\varepsilon}{2(1+\delta)}N^{\delta}\}$  ist

$$\sup_{\log \log N \le t \le \frac{1-\delta}{s} \log N} Y_t \le \sup_{\log \log N \le t \le \frac{1-\delta}{s} \log N} (1+\delta) e^{st} W_{\infty} \le (1+\delta) e^{(1-\delta) \log N} \frac{\varepsilon}{2(1+\delta)} N^{\delta} = \frac{\varepsilon}{2} N.$$

Wegen

$$\limsup_{N \to \infty} \left\{ \sup_{t \le \log \log N} Y_t \ge \varepsilon N \right\} \subset \left\{ W_\infty = \infty \right\}$$

(und  $\mathbb{P}(W_{\infty} = \infty) = 0)$  gilt somit

$$\mathbb{P}\Big(\liminf_{N\to\infty}\Big\{\tfrac{1-\delta}{s}\log N\leq \tilde{\tau}_{\varepsilon N}\leq \tfrac{1+\delta}{s}\log N\Big\}\Big)=1.$$

**Bericht.** Es gilt auch  $\mathbb{E}_1[\tilde{\tau}_{\varepsilon N} | Y_t \neq 0 \forall t \ge 0] \sim \frac{1}{s} \log N.$ )

Für den Beweis von Satz 3.7 vergleichen wir verschiedene zeitkontinuierliche Markovketten (nämlich das Moran-Modell und einen Verzweigungsprozess), indem wir sie mittels Zeittransformation ineinander überführen. Die allgemeine Situation beschreibt das folgende Lemma.

**Lemma 3.13.** Sei  $(X_t)_{t\geq 0}$  zeitkontinuierliche Markovkette auf E mit Q-Matrix  $(q_{ij})$ , sei  $\varphi : E \rightarrow (0, \infty)$  (sagen wir, beschr. und glm. positiv), setze

$$T_u \coloneqq \int_0^u \varphi(X_v) \, dv, \quad u \ge 0,$$

 $((T_u) \text{ ist ein sogenanntes ,} additives Funktional" von X).$  $u \mapsto T_u$  hat stetige, strikt wachsende Pfade, die Inverse ist

$$T_t^{-1} := \inf \{ u \ge 0 : T_u > t \} \quad f \ddot{u} r t \ge 0.$$

Sei  $\widehat{X}_t := X_{T_t^{-1}}, t \ge 0$ .  $\widehat{X} = (\widehat{X}_t)_{t\ge 0}$  ist zeitkontinuierliche Markovkette auf E mit Sprungratenmatrix  $\widehat{Q} = (\widehat{q}_{ij})$ , wobei

$$\widehat{q}_{ij} = \frac{q_{ij}}{\varphi(i)}$$

Beweisskizze. Sei  $X_0 = i$ ,

 $\tau_1 \coloneqq \text{erster Sprungzeitpkt. von } X (\sim \exp(-q_{ii})),$ 

es ist  $T_{\tau_1}$  =  $\tau_1 \varphi(i)$  und

$$t < T_{\tau_1} \iff T_t^{-1} < T_{T_{\tau_1}}^{-1} = \tau_1$$

d.h.  $\widehat{\tau}_1$  = erster Sprungzeitpkt. von  $\widehat{X} = T_{\tau_1} = \tau_1 \varphi(i)$  (~ Exp $(-q_{ii}/\varphi(i))$ ).

Dann verwendet man die starke Markov-Eigenschaft von X und die Tatsache, dass die Skelettketten von X und von  $\hat{X}$  nach Konstruktion übereinstimmen.
*Beweis von Satz 3.7.* Wähle  $\varepsilon > 0$ .

1. Phase:

Wende Lemma 3.13 an auf  $X^{(N)}$  (startend von  $X_0^{(N)} = 1$ ) mit  $\varphi(i) = 1 - \frac{i}{N}$ , also

$$T_u = \int_0^u 1 - \frac{X_v^{(N)}}{N} \, dv, \quad T_t^{-1} \coloneqq \inf\{u \ge 0 : T_u > t\}.$$

Sei $\widetilde{Y}_t\coloneqq X_{T_t^{-1}}^{(N)}$  und

$$\widetilde{\tau}_{\varepsilon N} = \inf\{t \ge 0 : \widetilde{Y}_t \ge \varepsilon N\}.$$

 $(\widetilde{Y}_t)_{t\geq 0}$  ist gemäß Lemma 3.13 verteilt wie ein superkritischer Galton-Watson-Prozess mit  $\lambda = 1 + s$ ,  $\mu = 1$  (der gestoppt wird, sobald N Ind. erreicht).

Solange  $\widetilde{X}_t^{(N)} \leq \varepsilon N$  gilt, ist  $(1 - \varepsilon)u \leq T_u \leq u$ , also

$$t \leq T_t^{-1} \leq \frac{1}{1-\varepsilon}t \quad \text{und daher} \quad \tau_{\varepsilon N} = T_{\widetilde{\tau}_{\varepsilon N}}^{-1} \in \left[\widetilde{\tau}_{\varepsilon N}, \frac{1}{1-\varepsilon}\widetilde{\tau}_{\varepsilon N}\right]$$

Mit Beob. 3.12 gilt

$$\mathbb{P}_1\left(\frac{s}{\log N}\widetilde{\tau}_{\varepsilon N}\in \left(1-\varepsilon,\frac{1}{1-\varepsilon}\right)\middle|\widetilde{\tau}_N<\widetilde{\tau}_0\right)\to 1\quad\text{für }N\to\infty.$$

Beachte:  $\{\widetilde{Y} \text{ überlebt}\} \subset \{\widetilde{\tau}_N < \widetilde{\tau}_0\} \text{ mit } \mathbb{P}_1(\widetilde{Y} \text{ überlebt}) > 0 \text{ und we$ gen

$$\mathbb{P}_1(\{\widetilde{\tau}_N < \widetilde{\tau}_0\} \cap \{\widetilde{Y} \text{ überlebt}\}^c) = \mathbb{P}_N(\widetilde{Y} \text{ stirbt aus}) = 1/(1+s)^N \to 0 \quad \text{für } N \to \infty$$

spielt es keine Rolle, ob wir auf  $\{\widetilde{\tau}_N < \widetilde{\tau}_0\}$  oder auf  $\{\widetilde{Y} \text{ überlebt}\} = \{\lim_{t \to \infty} e^{-st}\widetilde{Y}_t > 0\}$  bedingen.

Daher gilt auch

$$\mathbb{P}_1\left(\frac{s}{\log N}\tau_{\varepsilon N} \in \left((1-\varepsilon)^2, \frac{1}{(1-\varepsilon)^2}\right) \middle| \tau_N < \tau_0\right) \to 1 \quad \text{für } N \to \infty.$$

2. Phase:

Zeige

$$\mathscr{L}(\tau_{(1-\varepsilon)N} - \tau_{\varepsilon N} | \tau_N < \tau_0), \ N \in \mathbb{N}, \ \text{ ist straff.}$$

Intuitiv ist plausibel (zumindest wenn man ohne die Bedingung  $\tau_N < \tau_0$  argumentiert), dass  $(\frac{1}{N}X_t^{(N)})_{t\geq 0}$  startend von  $X_0^{(N)} = [Nx_0]$  mit  $x_0 \in (0, 1)$  für  $N \to \infty$  gegen die Lösung der logistischen Differentialgleichung

$$\frac{d}{dt}y(t) = sy(t)(1 - y(t)), \quad y(0) = x_0$$

konvergiert (beachte: wir transformieren hier nicht die Zeitachse).

Startend von  $y(0) = \varepsilon$  erreicht die Kurve  $(y(t))_{t\geq 0}$  den Wert  $1 - \varepsilon$  in endlicher Zeit.

Vgl. auch die Diskussion in Abschnitt 3.2, folgende heuristische Rechnung zeigt, dass die 2. Momente von Inkrementen trivial werden: Für  $h \downarrow 0$  ist

$$\frac{1}{h} \mathbb{E} \Big[ \frac{1}{N} X_{t+h}^{(N)} - \frac{1}{N} X_t^{(N)} \Big| \frac{1}{N} X_t^{(N)} = x \Big] \\ = \frac{1}{h} \frac{1}{N} \Big( h(1+s) N x \frac{N-Nx}{N} \cdot (+1) + h(N-Nx) \frac{Nx}{N} \cdot (-1) + o(h) \Big) = sx(1-x) + o(1)$$

und

$$\begin{split} &\frac{1}{h} \mathbb{E} \Big[ \Big( \frac{1}{N} X_{t+h}^{(N)} - \frac{1}{N} X_t^{(N)} \Big)^2 \, \big| \, \frac{1}{N} X_t^{(N)} = x \Big] \\ &= \frac{1}{h} \frac{1}{N^2} \Big( h(1+s) N x \frac{N-Nx}{N} \cdot (+1)^2 + h(N-Nx) \frac{Nx}{N} \cdot (-1)^2 + o(h) \Big) \\ &= \frac{1}{N} \Big( (2+s) x (1-x) + o(1) \Big), \end{split}$$

was im nahelegt, dass  $X^{(N)}$  gegen die Lösung einer deterministischen Differentialgleichung konvergiert.

Man kann dieses Argument präzise machen, für unsere Zwecke genügt hier eine gröbere Abschätzung via einem (erneuten) Vergleich mit einem superkritischen Galton-Watson-Prozess:

Für t genügend groß  $(t > \frac{1}{s} \log \frac{1-\varepsilon}{\varepsilon})$  ist

$$\lim_{N \to \infty} \mathbb{P}_{[\varepsilon N]}(Y_t \ge (1 - \varepsilon)N \,|\, \mathbf{\ddot{U}berleben}) = 1$$

(verwende  $\mathscr{L}_{[\varepsilon N]}(Y_t) = \mathscr{L}_1(Y_t)^{*[\varepsilon N]}$  und das Ges.d.gr.Z.)

Solange  $X_u^{(N)} \leq (1 - \varepsilon)N$  gilt, ist  $\varepsilon u \leq T_u \leq u$ , also  $t \leq T_t^{-1} \leq \frac{1}{\varepsilon}t$ , somit

$$\tau_{(1-\varepsilon)N} - \tau_{\varepsilon N} \in \left[\widetilde{\tau}_{(1-\varepsilon)N} - \widetilde{\tau}_{\varepsilon N}, \frac{1}{\varepsilon} (\widetilde{\tau}_{(1-\varepsilon)N} - \widetilde{\tau}_{\varepsilon N})\right]$$

(sofern  $\tau_{\varepsilon N} < \infty$ ) und demnach gilt insbesondere

$$\mathbb{P}\Big(\tau_{(1-\varepsilon)N} - \tau_{\varepsilon N} \leq \frac{1}{s} \frac{1}{\varepsilon} \log \frac{1-\varepsilon}{\varepsilon}\Big) \underset{N \to \infty}{\longrightarrow} 1$$

(und man überlege, dass dies auch unter der Bedingung  $\tau_N < \infty$  richtig bleibt: gegeben  $\tau_{\varepsilon N} < \infty$  hat  $\{\tau_N < \infty\}$  W'keit  $\geq 1 - (1/1 + s)^{\varepsilon N}$ ).

3. Phase:

Sobald  $X^{(N)}$  den Wert  $(1 - \varepsilon)N$  erreicht, vertauschen wir die Rollen der Typen A und a:  $\overline{X}_{t}^{(N)} \coloneqq N - X_{t}^{(N)}$  hat Sprungraten

$$q_{ij} = \begin{cases} i\frac{N-i}{N}, & j = i+1, \\ (1+s)i\frac{N-i}{N}, & j = i-1, \\ -(2+s)i\frac{N-i}{N}, & j = i, \\ 0, & \text{sonst.} \end{cases}$$

Wir starten in  $\overline{X}_{0}^{(N)} = \varepsilon N$ , die 3. Phase endet, sobald  $\overline{X}_{t}^{(N)} = 0$  gilt (wörtlich müssten wir darauf bedingen, dass  $\overline{X}^{(N)}$  niemals den Wert N übersteigt, dies fällt asymptotisch nicht ins Gewicht, da  $\mathbb{P}_{\varepsilon N}(\overline{X}^{(N)}$  erreicht N) exponentiell klein in N wird).

Wir vergleichen mit einem subkritischen Galton-Watson-Prozess ( $\overline{Y}_t$ ) (Individuen sterben mit Rate 1 + s, verdoppeln sich mit Rate 1) mit Sprungraten

$$\overline{q}_{ij} = \begin{cases} i, & j = i+1, \\ (1+s)i, & j = i-1, \\ -(2+s)i, & j = i, \\ 0, & \text{sonst.} \end{cases}$$

Prop. 3.10 liefert

$$\mathbb{E}_k[\overline{Y}_t] = ke^{-st}, \quad \operatorname{Var}_k[\overline{Y}_t] = k\frac{2+s}{s}(e^{-st} - e^{-2st}).$$

Somit gilt für  $\delta > 0$ :

$$\mathbb{P}_{[\varepsilon N]}(\overline{Y}_{\frac{1+\delta}{s}\log N} > 0) \le \mathbb{E}_{[\varepsilon N]}[\overline{Y}_{\frac{1+\delta}{s}\log N}] = [\varepsilon N] \exp(-(1+\delta)\log N) \sim \varepsilon N^{-\delta} \to 0,$$

and ererseits ist (beachte  $\mathbb{E}_{[\varepsilon N]}[\overline{Y}_{\frac{1-\delta}{s}\log N}] = [\varepsilon N]N^{\delta-1})$ 

$$\begin{aligned} \mathbb{P}_{[\varepsilon N]}\left(\overline{Y}_{\frac{1-\delta}{s}\log N} = 0\right) &\leq \mathbb{P}_{[\varepsilon N]}\left(|\overline{Y}_{\frac{1-\delta}{s}\log N} - \mathbb{E}_{[\varepsilon N]}[\overline{Y}_{\frac{1-\delta}{s}\log N}]| \geq \varepsilon N^{\delta}\right) \\ &\leq \frac{1}{\varepsilon^2 N^{2\delta}} \operatorname{Var}_{[\varepsilon N]}\left(\overline{Y}_{\frac{1-\delta}{s}\log N}\right) \\ &= \frac{1}{\varepsilon^2 N^{2\delta}}\left[\varepsilon N\right] \frac{2+s}{s} \left(\frac{1}{N^{1+\delta}} - \frac{1}{N^{2(1+\delta)}}\right) \to 0 \end{aligned}$$

Demnach gilt

$$\mathbb{P}_{[\varepsilon N]}\left(\frac{1-\delta}{s}\log N < \inf\left\{t \ge 0: \overline{Y}_t = 0\right\} \le \frac{1+\delta}{s}\log N\right) \underset{N \to \infty}{\longrightarrow} 1$$

und Argumente wie für Phase 1 zeigen, dass Analoges für  $\overline{X}^{(N)}$  gilt.

Schließlich folgt die Behauptung mit  $\delta \downarrow 0$ , dann  $\varepsilon \downarrow 0$ .

## 3.4 (2-Typ) Moran-Modell mit (gerichteter) Selektion und Mutation

Wir nehmen nun zusätzlich zur Dynamik aus Abschnitt 3.3 (vgl. Def. 3.3) die Möglichkeit der *Mutation* an:

- Jedes Typ a-Individuum mutiert mit Rate  $m_A$  zu Typ A,
- jedes Typ A-Ind. mutiert mit Rate  $m_a$  zu Typ a

mit gewissen  $m_A, m_a \ge 0$ .

**Definition 3.14** ((2 Typ-)Moran-Modell mit gerichteter Selektion und Mutation). Die zeitkontinuierliche Markovkette  $X^{(N)}$  auf  $\{0, 1, 2, ..., N\}$  mit Sprungratenmatrix

$$q_{ij} = \begin{cases} (1+s)i\frac{N-i}{N} + (N-i)m_A, & j = i+1, \\ (N-i)\frac{i}{N} + im_a, & j = i-1, \\ -(2+s)\frac{i(N-i)}{N} - (N-i)m_A - im_a, & j = i, \\ 0, & \text{sonst.} \end{cases}$$

interpretieren wir als

$$X_t^{(N)} = \#$$
Typ A-Individuen zur Zeit  $t, \quad t \ge 0$ 

in einer Population der Größe N wie oben beschrieben.  $X^{(N)}$  heißt (Typenzählprozess des) 2 Typ-Moran-Modell(s) mit gerichteter Selektion und Mutation.

**Bemerkung.** Für gegebenes  $N \in \mathbb{N}$ ,  $s \ge 0$  und  $m_a, m_A > 0$  ist  $X^{(N)}$  offensichtlich irreduzibel, demnach existiert ein eindeutiges Gleichgewicht.

**Beobachtung 3.15** (Gleichgewichte von Geburts- und Todesprozessen). Betrachte eine zeitkontinuierliche Markovkette X auf  $\{0, 1, ..., N\}$  mit Sprungraten

$$q_{i,i+1} = \lambda_i, \quad q_{i,i-1} = \mu_i, \quad q_{i,i} = -(\lambda_i + \mu_i)$$

mit  $\lambda_i > 0$  für  $0 \le i < N$ ,  $\mu_i > 0$  für  $0 < i \le N$ ,  $\lambda_N = \mu_0 = 0$  (ein solches X heißt auch ein Geburtsund Todesprozess).

X besitzt eine reversible Gleichgewichtsverteilung  $(\pi_k)_{k=0,1,\dots,N}$ , d.h.  $(\pi_k)$  erfüllt

$$\pi_{k+1}\mu_{k+1} = \pi_k \lambda_k \quad \text{für } k = 0, \dots, N-1 \tag{3.9}$$

("detaillierte Balance"). Folglich gilt

$$\pi_k = \frac{\lambda_{k-1}}{\mu_k} \pi_{k-1} = \dots = \pi_0 \prod_{j=1}^k \frac{\lambda_{j-1}}{\mu_j}$$

und somit

$$\pi_k = \frac{\varphi_k}{\sum_{j=0}^N \varphi_j} \quad \text{mit} \quad \varphi_k \coloneqq \prod_{j=1}^k \frac{\lambda_{j-1}}{\mu_j} \quad (\text{mit } \varphi_0 = 1).$$

Aus der detaillierten Balancegleichung (3.9) folgt natürlich die Gleichgewichtsbedingung

$$\pi_k(\lambda_k + \mu_k) = \pi_{k+1}\mu_{k+1} + \pi_{k-1}\lambda_{k-1}$$

(für k = 0, ..., N, mit  $\pi_{N+1} \coloneqq \pi_{-1} \coloneqq 0$ ), d.h. die (wegen Irreduzibilität eindeutige) Gleichgewichtsverteilung hat tatsächlich diese Gestalt.

Geburts- und Todesprozesse sind also stets reversibel, dies ist eine Spezialität der 1-dim. Situation.

Aus dem Zusammenspiel der verschiedenen evolutionären Kräfte ergibt sich nun ein Mutations-Selektions-Drift-Gleichgewicht:

**Satz 3.16** (Mutations-Selektions-Drift-Gleichgewicht im Moran-Modell). Seien  $s \ge 0$ ,  $m_A$ ,  $m_a > 0$ ,

 $\pi^{(N)}(k) = \lim_{t \to \infty} \mathbb{P}(X_t^{(N)} = k), \quad k = 0, \dots, N$ 

die Gleichgewichts-Verteilung des 2 Typ-Moran-Modells mit gerichteter Selektion und Mutation (aus Def. 3.14).

a) Es ist

$$\pi^{(N)}(k) = \frac{1}{C_N(m_A, m_a, s)} (1+s)^k {\binom{N}{k}} \frac{\left(\frac{Nm_A}{1+s}\right)_{k\uparrow} \left(Nm_a\right)_{(N-k)\uparrow}}{\left(N\frac{m_A+(1+s)m_a}{1+s}\right)_{N\uparrow}}$$

mit einer Normierungskonstante  $C_N(m_A, m_a, s)$ .

Die Normierungskonstante erfüllt

$$C_N(m_A, m_a, s) = \mathbb{E}\left[(1+s\xi)^N\right] \quad mit\,\xi \sim \text{Beta}\left(\frac{Nm_A}{1+s}, Nm_a\right).$$

b) Mit  $\theta_A, \theta_a > 0, \sigma \ge 0$  gelte

$$s = \sigma/N, m_A = \theta_A/N, m_a = \theta_a/N, \tag{3.10}$$

dann konvergiert  $\sum_{k=0}^{N} \pi^{(N)}(k) \delta_{k/N}$  für  $N \to \infty$  (schwach als W-Maß auf [0,1]) gegen das Wahrscheinlichkeitsmaß mit Dichte

$$\frac{1}{C_{\theta_A,\theta_a,\sigma}} x^{\theta_A - 1} (1 - x)^{\theta_a - 1} e^{\sigma x}, \quad x \in (0, 1),$$
(3.11)

wobei  $C_{\theta_A,\theta_a,\sigma} = \int_0^1 x^{\theta_A - 1} (1 - x)^{\theta_a - 1} e^{\sigma x} dx.$ 

**Bemerkung.** Annahme (3.10) bewirkt, dass die "evolutionären Kräfte" Selektion, Mutation und Gendrift (im Modell) auf vergleichbarer Zeitskala wirken. Anders gewendet: Das Modell passt zu einer gegebenen Situation, wenn die Population groß und Mutationsraten und Selektionsvorteil (dazu quantitativ passend) klein sind – s.a. die entsprechende Diskussion in Kapitel 2, speziell die Bemerkung auf Seite 23 zur analogen Annahme (2.1) im neutralen Fall.

Beweis von Satz 3.16. a)

$$\begin{split} \prod_{i=0}^{k-1} \lambda_i &= \prod_{i=0}^{k-1} \left( \underbrace{(1+s)i\frac{N-i}{N} + (N-i)m_A}_{=\frac{1}{N}(1+s)(N-i)\left(i+\frac{Nm_A}{1+s}\right)} \right) = \frac{1}{N^k} (1+s)^k (N)_{k\downarrow} \prod_{i=0}^{k-1} \left( \frac{Nm_A}{1+s} + i \right) \\ &= \frac{1}{N^k} (1+s)^k (N)_{k\downarrow} \left( \frac{Nm_A}{1+s} \right)_{k\uparrow}, \\ \prod_{i=1}^k \mu_i &= \prod_{i=1}^k \left( \underbrace{(N-i)\frac{i}{N} + im_a}_{=\frac{i}{N}(N-i+Nm_a)} \right) = \frac{k!}{N^k} \prod_{i=1}^k (N-i+Nm_a) = \frac{k!}{N^k} \frac{(Nm_a)_{N\uparrow}}{(Nm_a)_{(N-k)\uparrow}}, \end{split}$$

also

$$\pi^{(N)}(k) \propto \frac{\prod_{i=0}^{k-1} \lambda_i}{\prod_{i=1}^k \mu_i} \propto (1+s)^k \binom{N}{k} \frac{\left(\frac{Nm_A}{1+s}\right)_{k\uparrow} \left(Nm_a\right)_{(N-k)\uparrow}}{\left(N\frac{m_A+(1+s)m_a}{1+s}\right)_{N\uparrow}}$$

Für die Normierung beachte :  $\xi \sim \text{Beta}(\alpha_1, \alpha_2)$  erfüllt

$$\mathbb{E}[\xi^k(1-\xi)^\ell] = \frac{\text{Beta}(\alpha_1+k,\alpha_2+\ell)}{\text{Beta}(\alpha_1,\alpha_2)} = \frac{\Gamma(\alpha_1+\alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \frac{\Gamma(\alpha_1+k)\Gamma(\alpha_2+\ell)}{\Gamma(\alpha_1+\alpha_2+k+\ell)}$$

und  $\Gamma(\alpha + k)/\Gamma(\alpha) = (\alpha)_{k\uparrow}$ , (wende an mit  $\alpha_1 = \frac{Nm_A}{1+s}$ ,  $\alpha_2 = Nm_a$ ), erhalte

$$\mathbb{E}\left[(1+s\xi)^{N}\right] = \mathbb{E}\left[\left((1+s)\xi + (1-\xi)\right)^{N}\right] = \sum_{k=0}^{N} (1+s)^{k} \binom{N}{k} \mathbb{E}\left[\xi^{k}(1-\xi)^{N-k}\right].$$

b) Schreibe

$$\lambda_{i} = N \frac{i}{N} \left( 1 - \frac{i}{N} \right) \left( 1 + s + \frac{m_{A}}{i/N} \right) \quad (i \neq 0, \ \lambda_{0} = N m_{A}),$$
  
$$\mu_{i} = N \frac{i}{N} \left( 1 - \frac{i}{N} \right) \left( 1 + \frac{m_{a}}{1 - i/N} \right) \quad (i \neq N, \ \mu_{N} = N m_{a}).$$

$$\begin{split} \varphi_k^{(N)} &= \frac{\prod_0^{k-1} \lambda_i}{\prod_1^k \mu_i} = \frac{Nm_A}{N\frac{k}{N} \left(1 - \frac{k}{N}\right) \left(1 + \frac{m_a}{1 - k/N}\right)} \prod_{i=1}^{k-1} \frac{1 + s + \frac{m_A}{i/N}}{1 + \frac{m_a}{1 - i/N}} \\ &= \frac{1}{N} \frac{\theta_A}{\frac{k}{N} \left(1 - \frac{k}{N}\right) \left(1 + \frac{1}{N} \frac{\theta_a}{1 - k/N}\right)} \prod_{i=1}^{k-1} \frac{1 + \frac{\sigma}{N} + \frac{\theta_A}{i}}{1 + \frac{\theta_a}{N - i}} \end{split}$$

Sei  $k = \lfloor Nx \rfloor$  mit  $x \in [\varepsilon, 1 - \varepsilon]$  ( $\varepsilon$  zunächst fest)

$$\log \prod_{i=1}^{k-1} \frac{1 + \frac{\sigma}{N} + \frac{\theta_A}{i}}{1 + \frac{\theta_a}{N-i}} = \sum_{i=1}^{k-1} \log \left(1 + \frac{\sigma}{N} + \frac{\theta_A}{i}\right) - \log \left(1 + \frac{\theta_a}{N-i}\right)$$
$$= \sum_{i=1}^{k-1} \frac{\sigma}{N} + \frac{\theta_A}{i} - \frac{\theta_a}{N-i} + R_{N,1}(k)$$

mit  $|R_{N,1}(k) - \tilde{R}_{\varepsilon}| \le C/N^2$  (beachte  $\log(1 + x) = x + O(x^2)$ ), dann verwende<sup>2</sup>

$$\sum_{i=1}^{k} \frac{1}{i} = \log k + c_{\gamma} + o(1)$$
 mit  $c_{\gamma} \approx 0.57721...$ , die Euler-Mascheroni-Konstante

<sup>&</sup>lt;sup>2</sup>Für die explizite Asymptotik siehe z.B. M. Abramowitz und I.A. Stegun, *Handbook of mathematical functions*, Dover publications, 9. Aufl., 1970, 6.3, 18 beachte  $\psi(z) = \Gamma'(z)/\Gamma(z) \left(=\frac{d}{dz}\log(\Gamma(z))\right)$ , die Digamma-Funktion, erfüllt  $\psi(n) + \gamma = \sum_{k=1}^{n-1} k^{-1}$  nach 6.3, 2); für ein Argument "von Hand" beachte (vgl. z.B. Heuser, *Lehrbuch der Analysis I*, Aufg. 88, 5)  $d_m := \sum_{k=1}^m \frac{1}{k} - \int_1^m \frac{1}{x} dx$  erfüllt  $d_m \ge 0$  (Integralvergleich) und  $d_m - d_{m+1} = -\frac{1}{m+1} + \int_m^{m+1} \frac{1}{x} dx = \int_m^{m+1} \frac{m+1-x}{x(m+1)} dx \in (0, \frac{1}{m(m+1)})$ , d.h.  $d_m \searrow \gamma$  (die Euler-Mascheroni-Konstante) und  $d_m - \gamma \le \sum_{\ell=m}^{\infty} (d_\ell - d_{\ell-1}) \le C/m$ ]

Somit ist

$$\begin{split} \varphi_k^{(N)} &= \frac{1}{N} \frac{\theta_A}{\frac{k}{N} \left(1 - \frac{k}{N}\right) \left(1 + \frac{1}{N} \frac{\theta_a}{1 - k/N}\right)} \\ &\times \exp\left(\frac{k - 1}{N} \sigma + \theta_A \log(k - 1) - \theta_a \log(N - 1) + \theta_a \log(N - k + 1)\right) \\ &\times \exp\left(\tilde{R}_{\varepsilon} + O(1/N^2)\right) \\ &= \frac{C}{N} x^{\theta_A - 1} (1 - x)^{\theta_a - 1} e^{\sigma x} \left(1 + O(1/N^2)\right). \end{split}$$

**Bericht 3.17** (Konvergenz gegen die Wright-Fisher-Diffusion mit Mutation und Selektion). Unter (3.10) konvergiert

$$\left(\frac{1}{N}X_{tN/2}^{(N)}\right)_{t\geq 0} \to X = (X_t)_{t\geq 0} \quad \text{für } N \to \infty$$

(in Vert. auf  $D([0,\infty);[0,1])$ ), wobei X Lösung von

$$dX_t = \left(\frac{\sigma}{2}X_t(1-X_t) + \frac{1}{2}(\theta_A - (\theta_A + \theta_a)X_t)\right)dt + \sqrt{X_t(1-X_t)}\,dB_t$$

mit  $(B_t)$  Standard-Brownbewegung. X heißt die (2-Typ) Wright-Fisher-Diffusion mit (gerichteter) Selektion und Mutation, sie ist ein Markovprozess auf [0, 1] mit Generator

$$Lf(x) = \left(\frac{\sigma}{2}x(1-x) + \frac{1}{2}(\theta_A - (\theta_A + \theta_a)x)\right)f'(x) + \frac{1}{2}x(1-x)f''(x) \quad \text{für } f \in C^2([0,1]).$$

Heuristisch können wir hier zumindest beobachten, dass für  $h \downarrow 0$  gilt

$$\mathbb{E}\left[\frac{1}{N}\left(X_{(t+h)N/2}^{(N)} - X_{tN/2}^{(N)}\right) \middle| X_{tN/2}^{(N)} = k\right]$$
  
=  $\frac{1}{N}\left(h\frac{N}{2}\left((1 + \frac{\sigma}{N})k\frac{N-k}{N} + (N-k)\frac{\theta_A}{N}\right)(+1) + h\frac{N}{2}\left(k\frac{N-k}{N} + k\frac{\theta_a}{N}\right)(-1) + O(h^2)\right)$   
=  $h\left(\frac{\sigma}{2}\frac{k}{N}\left(1 - \frac{k}{N}\right) + \left(1 - \frac{k}{N}\right)\frac{\theta_A}{2} - \frac{k}{N}\frac{\theta_a}{2}\right) + O(h^2)$ 

und

$$\begin{split} & \mathbb{E}\left[\frac{1}{N^2} \left(X_{(t+h)N/2}^{(N)} - X_{tN/2}^{(N)}\right)^2 \left|X_{tN/2}^{(N)} = k\right] \\ &= \frac{1}{N^2} \left(\frac{N}{2} h\left((2 + \frac{\sigma}{N})k\frac{N-k}{N} + \left(1 - \frac{k}{N}\right)\theta_A + \frac{k}{N}\theta_a\right)\right) + O(h^2) \\ &= h\frac{k}{N} \left(1 - \frac{k}{N}\right) + O(h/N) + O(h^2), \end{split}$$

was zu obigem Generator führt.

Die Dichte (3.11)

$$\frac{1}{C_{\theta_A,\theta_a,\sigma}} x^{\theta_A-1} (1-x)^{\theta_a-1} e^{\sigma x}$$

aus Satz 3.16 ist die Gleichgewichtsdichte der Wright-Fisher-Diffusion mit Selektion und Mutation.

### 3.4.1 Graphische Konstruktion und anzestraler Selektionsgraph

Für jedes geordnete Paar  $(i, j), i \neq j$  sei  $(N_t^{(i,j)})_{t\geq 0}$  ein Poissonprozess mit Rate  $1/N, (S_t^{(i,j)})_{t\geq 0}$  ein Poissonprozess mit Rate s/N, für jedes i sei  $(M_t^{(a,i)})_{t\geq 0}$  ein Poissonprozess mit Rate  $m_a$  und  $(M_t^{(A,i)})_{t\geq 0}$  ein Poissonprozess mit Rate  $m_A$ .

Zusätzlich zu den Pfeilen wie in Bem. 3.4 (für den Fall ohne Mutationen) legen wir nun noch "Mutationsereignisse" auf die "Zeitachsen": ein Sprung von  $M^{(a,i)}$  ändert den Typ von Individuum *i* auf *a*, ein Sprung von  $M^{(A,i)}$  ändert ihn auf *A* (wobei "triviale Wechsel" erlaubt sind, z.B. bleibt ein Typ *A*-Individuum *i* nach einem Sprung von  $M^{(A,i)}$  vom Typ *A*).

[Skizze an der Tafel]

**Beobachtung 3.18** (Ablesen der "Genealogie" und anzestraler Selektionsgraph (ASG)). 1. Eine *n*-Stichprobe und ihre Genealogie können in einem zweistufigen Prozess gewonnen werden

- zuerst "rückwärts" (in der Zeit): Linien verschmelzen oder spalten (wenn von der Spitze eines s-Pfeils getroffen), verfolge bis # Linien = 1 (der "ultimate ancestor" ist erreicht)
- dann "vorwärts" : lege Mutationen (die die Typen bestimmen), löse damit potentielle Verzweigungsereignisse auf

(sobald wir wissen, ob ein gewisses Individuum zum Zeitpunkt eines s-Pfeils Typ A oder a hat, können entscheiden, ob der s-Pfeil benutzt wurde).

2. Alternative Rückwärtsdynamik in "einem Schritt": Sei  $(A_t)_{t\geq 0}$  zeitkontinuierliche Markovkette mit Werten in  $2^{\{1,\ldots,N\}} \cup \{\partial\}$ , wenn aktuell  $A_{t-} = B \subset \{1,\ldots,N\}$ :

für 
$$i \in B$$
: springe von  $A_{t-} = B$  nach  $A_t = B \setminus \{i\}$  mit Rate  $\frac{|B| - 1}{N} + m_a$ ,  
für  $j \notin B$ : springe von  $A_{t-} = B$  nach  $A_t = B \cup \{j\}$  mit Rate  $\frac{s}{N}|B|$ ,  
springe nach  $\partial$  mit Rate  $m_A|B|$  (der Prozess wird "getötet")

Set 
$$t > 0, n \le N, i_1, ..., i_n \in \{1, ..., N\}$$
, starte in  $A_0 = \{i_1, ..., i_n\}$ . Es ist

$$\{X_t(i_1) = a, \dots, X_t(i_n) = a\} = \{A_t \text{ nicht getötet}\} \cap \{X_0(j) = a, \text{ für alle } j \in A_t\},$$
(3.12)

wie wir von der graphischen Konstruktion ablesen.

Der Zählprozess  $(Y_t^{(N)})$  mit

$$Y_t^{(N)} = |A_t|$$
, aktuelle Anzahl "potentieller Ahnenlinien",

hat dann (offenbar) Sprungratenmatrix

$$q_{ij}^{Y,N} = \begin{cases} s\frac{i(N-i)}{N}, & j = i+1, \\ \frac{i(i-1)}{N} + im_a, & j = i-1, \\ im_A, & j = \partial, \\ -(1+s)\frac{i(N-i)}{N} - i(m_a + m_A), & j = i \end{cases}$$

Nehmen wir an  $X_0^{(N)}(1), \ldots, X_0^{(N)}(N)$  sind austauschbar verteilt (und u.a. von den PPPen der graphischen Konstruktion)

Sei  $n \le N, t \ge 0$ , ziehe *n*-mal ohne Zurücklegen aus der Population zur Zeit *t*. Die W'keit, dann *n*-mal Typ *a* zu sehen, ist somit

$$\mathbb{E}\left[\frac{(N-X_t^{(N)})(N-X_t^{(N)}-1)\cdots(N-X_t^{(N)}-n+1)}{N(N-1)\cdot(N-n+1)}\right] = \frac{1}{(N)_{n\downarrow}}\mathbb{E}\left[\left(N-X_t^{(N)}\right)_{n\downarrow}\right]$$
$$= \frac{1}{(N)_{n\downarrow}}\mathbb{E}\left[\left(N-X_0^{(N)}\right)_{Y_t^{(N)}\downarrow}\mathbf{1}(Y_t^{(N)} \text{ nicht getötet})\right]$$

(durch Ablesen von der graphischen Konstruktion, vgl. (3.12)).

**Beobachtung 3.19.**  $(Y_{tN/2}^{(N)})_{t\geq 0} \to Y = (Y_t)_{t\geq 0}$  (in Vert. auf  $D([0,\infty); \mathbb{N}_0 \cup \{\partial\})$ ), Y zeitk. MK auf  $\mathbb{N}_0$  mit Sprungraten

$$q_{ij}^{Y} = \begin{cases} \frac{\sigma}{2}i, & j = i+1, \\ \binom{i}{2} + i\frac{\theta_{a}}{2}, & j = i-1, \\ i\frac{\theta_{A}}{2}, & j = \partial, \end{cases}$$

Die Kette *Y* ist dual zur Wright-Fisher-Diffusion mit Mutation und Selektion *X* aus Bericht 3.17:

$$\mathbb{E}\left[(1-X_t)^n \left| X_t = x_0\right] = \mathbb{E}\left[(1-x_0)^{Y_t} e^{-\frac{\theta_A}{2} \int_0^t Y_u \, du} \left| Y_0 = n\right]$$
(3.13)

Betrachte  $f(x, n) \coloneqq (1 - x)^n$ , es ist

$$(L^{X}f(\cdot,n))(x) = \left(\frac{\sigma}{2}x(1-x) + \frac{1}{2}(\theta_{A} - (\theta_{A} + \theta_{a})x)\right)n(1-x)^{n-1}(-1) + \frac{1}{2}x(1-x)n(n-1)(1-x)^{n-2} = \frac{\sigma}{2}n\left((1-x)^{n+1} - (1-x)^{n}\right) + \frac{n(n-1)}{2}\left((1-x)^{n-1} - (1-x)^{n}\right) + \frac{\theta_{a}}{2}n\left((1-x)^{n-1} - (1-x)^{n}\right) - \frac{\theta_{A}}{2}n(1-x)^{n} = \left(L^{Y}f(x,\cdot)\right)(n) - \frac{\theta_{A}}{2}nf(x,n)$$

Interpretation von (3.13) via ancestral selection graph :

- Beginne mit *n* Linien,
- jede Linie verzweigt mit Rate  $\sigma/2$ ,
- jedes Paar Linien verschmilzt mit Rate 1,
- jede Linie "endet" mit Rate  $\theta_a/2$ ,

- der gesamte Prozess wird getötet mit Rate aktuelle #Linien  $\cdot \theta_A/2$ .
- Nach Zeit t, sofern der Prozess noch nicht getötet wurde, wählen wir für jede der dann existierenden Linien u.a. einen Typ (A mit W'keit x<sub>0</sub>, a mit W'keit 1 – x<sub>0</sub>)

Die rechte Seite von (3.13) ist die Wahrscheinlichkeit, dass Prozess nicht getötet wurde und am Ende alle Linien Typ a erhalten haben, dies ist zugleich die Wahrscheinlichkeit, in einer n-Stichprobe aus einer Population, deren Typenanteils-Evolution von der Wright-Fisher-Diffusion mit Mutation und Selektion beschrieben wird, n-mal den (selektiv benachteiligten) Typ a zu sehen.

Mit  $t \to \infty$  kann man auf diese Weise die Momente der Gleichgewichtsverteilung von X beschreiben. Für den Prozess Y genügt es dabei zu warten, bis der sogenannte "ultimate ancestor" erreicht ist (und sein Typ durch eine Mutation festgelegt ist, wörtlich also bis  $Y_t = 0$  gilt oder der Prozess getötet wird).

# Kapitel 4

# Räumliche Struktur

Wir betrachten gekoppelte bzw. interagierende Wright-Fisher-Modelle (der Einfachheit halber: nur der haploide Fall), d.h. wir behandeln folgendes Modell:

- S Menge von Kolonien,  $(p(x, y))_{x,y \in S}$  stochastische Matrix, lokale Populationsgröße  $N \in \mathbb{N}$
- Individuen einer gewissen Population leben in diskreten Generationen verteilt auf "Kolonien", die mit S (endl. oder abzählbare Menge, z.B. S = Z) indiziert sind, jeweils N Individuen pro Kolonie



- Typenmenge E, Mutationswahrscheinlichkeit μ: wir betrachten (der Einfachheit halber hier nur) Typen und Mutationen gemäß Infinitely-many-alleles-Modell (mathematisch realisiert etwa via E = [0, 1] und bei jedem Mutationsereignis wird unabhängig der neue Typ uniform([0, 1])verteilt gewählt, vgl. Abschnitt 2.1).
- Dynamik: Jedes Individuum in Generation r in Kolonie x wählt mit W'keit p(x, y) ein uniform aus Kolonie y in Generation r 1 gewähltes Individuum als "Elter", kopiert dessen Typ (mit Wahrscheinlichkeit 1-μ) bzw. erhält "mutierten Typ" (mit Wahrscheinlichkeit μ), jeweils unabhängig über die Individuen und die Kolonien

Wir betrachten später speziell den Fall  $p(x, y) = (1 - \nu)\delta_{xy} + \nu q(x, y)$  mit  $\nu \in (0, 1), q(\cdot, \cdot)$ stochastische Matrix mit 0-Einträgen auf der Diagonale (und wir denken an  $\nu \ll 1$ ).



Wie beeinflusst die räumliche Struktur die genetische Zusammensetzung der Gesamtpopulation? Ziehe zufällig zwei Individuen aus Generation n, eines uniform aus Kolonie x und eines uniform aus Kolonie y (falls y = x, so denken wir an Ziehen ohne Zurücklegen), sei

$$\psi_n(x,y) \coloneqq \mathbb{P} \begin{pmatrix} \text{zwei zufällig aus Kolonien } x \text{ und } y \text{ gezogene } \\ \text{Individuen haben denselben Typ} \end{pmatrix}$$

In der Literatur wird diese Größe (im hier diskutierten Kontext) auch als die Wahrscheinlichkeit der "identity by descent", IBD, bezeichnet.

Zerlegung gemäß der Position der Vorfahren der beiden zufällig gezogenen Individuen liefert

$$\begin{split} \psi_n(x,y) &= (1-\mu)^2 \sum_{\substack{x',y' \in S \\ x' \neq y'}} p(x,x') p(y,y') \psi_{n-1}(x',y') \\ &+ (1-\mu)^2 \sum_{x' \in S} p(x,x') p(y,x') \Big( \frac{1}{N} + \frac{N-1}{N} \psi_{n-1}(x',x') \Big) \\ &= (1-\mu)^2 \sum_{\substack{x',y' \in S}} p(x,x') p(y,y') \Big( \psi_{n-1}(x',y') + \mathbf{1}(y'=x') \frac{1-\psi_{n-1}(x',x')}{N} \Big) \end{split}$$

alternativ geschrieben als

$$\psi_n(x,y) = (1-\mu)^2 \sum_{x',y' \in S} p(x,x') p(y,y') \Big( \psi_{n-1}(x',y') \frac{N-\mathbf{1}(y'=x')}{N} + \frac{\mathbf{1}(y'=x')}{N} \Big)$$
(4.1)

Iterieren von (4.1) liefert (mit Notation  $x_0 = x, y_0 = y$ )

$$\begin{split} \psi_{n}(x,y) \\ &= \sum_{k=1}^{n} \frac{1}{N} (1-\mu)^{2k} \sum_{\substack{x_{1},x_{2},\dots,x_{k} \in S \\ y_{1},y_{2},\dots,y_{k} \in S \\ x_{k} = y_{k}}} \left( \prod_{i=1}^{k} p(x_{i-1},x_{i}) \right) \left( \prod_{i=1}^{k} p(y_{i-1},y_{i}) \right) \left( 1-\frac{1}{N} \right)^{\#\{1 \le i \le k-1: x_{i} = y_{i}\}-1} \\ &+ (1-\mu)^{2n} \sum_{\substack{x_{1},x_{2},\dots,x_{n} \in S \\ y_{1},y_{2},\dots,y_{n} \in S}} \left( \prod_{i=1}^{k} p(x_{i-1},x_{i}) \right) \left( \prod_{i=1}^{k} p(y_{i-1},y_{i}) \right) \left( 1-\frac{1}{N} \right)^{\#\{1 \le i \le n: x_{i} = y_{i}\}} \psi_{0}(x_{n},y_{n}) \end{split}$$

(der *k*-te Summand in der zweiten Zeile entsteht, indem man man beim Iterieren von (4.1) k - 1 mal den ersten Term in der großen Klammer rechts in (4.1) expandiert und dann den zweiten Term in der großen Klammer rechts in (4.1) verwendet, die dritte Zeile entsteht, wenn beim Iterieren stets den ersten Term in der großen Klammer rechts in (4.1) expandiert).

Mit  $(X_k^{(1)})_{k \in \mathbb{N}_0}$ ,  $(X_k^{(2)})_{k \in \mathbb{N}_0}$  unabhängigen *p*-Markovketten auf *S* (Interpretation: räumliche Einbettung der beiden Ahnenlinien) lässt sich dies folgendermaßen formulieren: Sei  $J_n := \#\{1 \le i \le n : X_i^{(1)} = X_i^{(2)}\}$  (Anzahl Kollisionen bis Zeit *n*), dann ist

$$\psi_n(x,y) = \sum_{k=1}^n (1-\mu)^{2k} \mathbb{E}_{(x,y)} \left[ \frac{1}{N} \left( 1 - \frac{1}{N} \right)^{J_k - 1} \mathbf{1} \left( X_k^{(1)} = X_k^{(2)} \right) \right] \\ + (1-\mu)^{2n} \mathbb{E}_{(x,y)} \left[ \left( 1 - \frac{1}{N} \right)^{J_n} \psi_0 \left( X_n^{(1)}, X_n^{(2)} \right) \right]$$

mit  $\tau_{\text{coal}}$  = Anzahl Schritte bis zur Verschmelzung (verschmelze zu jedem Zeitpunkt j, zu dem  $X_j^{(1)}$  =  $X_j^{(2)}$  gilt, mit W'keit 1/N; in Formeln z.B.  $\tau_{\text{coal}} = \inf\{j \ge 1 : X_j^{(1)} = X_j^{(2)}, W_j = 1\}$  mit  $W_j, j \in \mathbb{N}$  u.i.v. Ber(1/N)-ZVn, unabhängig von  $X^{(1)}, X^{(2)}$ ), d.h.

$$\psi_n(x,y) = \mathbb{E}_{(x,y)} \Big[ (1-2\mu)^{\tau_{\text{coal}}} \mathbf{1}(\tau_{\text{coal}} \le n) \Big] + \mathbb{E}_{(x,y)} \Big[ (1-2\mu)^n \mathbf{1}(\tau_{\text{coal}} > n) \psi_0 \Big( X_n^{(1)}, X_n^{(2)} \Big) \Big]$$



Abbildung 4.1: Illustration von Ahnenlinien als verzögert verschmelzende Irrfahrten

Mit  $n \to \infty$  (in Summendarstellung oder/und in Erwartungswertdarstellung) finden wir: Betrachte Population "im Gleichgewicht"<sup>1</sup>, ziehe zufällig zwei Individuen aus der Population im Gleichgewicht, eines uniform aus Kolonie x und eines uniform aus Kolonie y (falls y = x, so denken wir an Ziehen ohne Zurücklegen)

$$\psi(x,y) \coloneqq \mathbb{P}\left( \begin{array}{c} \text{zwei zufällig aus Kolonien } x \text{ und } y \text{ im Gleichgewicht} \\ \text{gezogene Individuen haben denselben Typ} \end{array} \right)$$

Es ist  $\psi(x,y) = \lim_{n \to \infty} \psi_n(x,y)$  und

$$\psi(x,y) = \mathbb{E}_{(x,y)}\left[(1-\mu)^{2\tau_{\text{coal}}}\right]$$
(4.2)

(siehe auch Abbildung 4.1). Zerlegung nach dem ersten Schritt (in die Vergangenheit) wie oben (oder der Grenzwert  $n \rightarrow \infty$  in (4.1)) liefert

$$\psi(x,y) = (1-\mu)^2 \sum_{x',y'\in S} p(x,x')p(y,y') \Big(\psi(x',y')\frac{N-\mathbf{1}(y'=x')}{N} + \frac{\mathbf{1}(y'=x')}{N}\Big)$$
(4.3)

**Räumlich homogener Fall** Wir nehmen ab hier an, dass S eine abelsche Gruppe ist (z.B.  $\mathbb{Z}, \mathbb{Z}^d$ ,  $Z/(L\mathbb{Z}), \mathbb{Z}^d/(L\mathbb{Z}^d), ...$ ) und p(x, y) = p(y-x) hängt nur von der Differenz der Argumente ab. Wir notieren p(z) := p(0, z) und analog  $p^k(z) = p^k(0, z)$  für die k-Schritt-Übergangswahrscheinlichkeiten.

<sup>&</sup>lt;sup>1</sup>Wir beweisen hier nicht die (zutreffende) Aussage, dass die Populationskonfigurationen tatsächllich in Verteilung gegen Gleichgewicht der interagierenden Wright-Fisher-Modelle mit Mutation konvergieren. Wir sehen aber zumindest: Die hier betrachtete Fragestellung ist im Grenzwert  $n \rightarrow \infty$  sinnvoll.

Insbesondere sind dann  $(X_k^{(1)})_{k \in \mathbb{N}_0}$  und  $(X_k^{(2)})_{k \in \mathbb{N}_0}$  u.a. Kopien einer Irrfahrt. Auch  $\psi(x, y) = \phi(y - x)$  hängt dann nur von y - x ab.

Im räumlich homogenen Fall lässt sich der Erwartungswert (4.2) folgendermaßen expliziter beschreiben. Zur Darstellung von  $\tau_{coal}$  betrachten wir die Kollisionszeiten

$$\tau_1 \coloneqq \inf\{k \ge 1 : X_k^{(1)} = X_k^{(2)}\}, \quad \tau_m \coloneqq \inf\{k \ge \tau_{m-1} : X_k^{(1)} = X_k^{(2)}\}, \ m = 2, 3, \dots$$

(wegen räumlicher Homogenität sind die  $\tau_{m-1} - \tau_m$ , m = 1, 2, ... u.i.v., verteilt wie die Rückkehrzeit zur 0 der Differenzirrfahrt), V davon u.a., ~ Geom(1/N) (hier:  $\mathbb{P}(V = m) = (1/N)(1 - 1/N)^{m-1}$ , m = 1, 2, ...)

$$\tau_{\text{coal}} \stackrel{d}{=} \tau_V = \tau_1 + \mathbf{1}(V \ge 2) \sum_{j=2}^{V} (\tau_j - \tau_{j-1})$$
(4.4)

**Beobachtung 4.1.** Sei  $Z_k = X_k^{(1)} - X_k^{(2)}, k \in \mathbb{N}_0$ .  $(Z_k)_k$  ist Irrfahrt mit Inkrementverteilung  $\mathbb{P}(Z_k - Z_k - z_k) = \widetilde{p}(0, z_k) = \sum p(0, z_k) p(0, z_k - z_k)$ 

$$\mathbb{P}(Z_k - Z_{k-1} = z) = \widetilde{p}(0, z) = \sum_y p(0, y) p(0, y - z)$$

falls p symmetrisch:  $\widetilde{p}(\cdot, \cdot) = p^2(\cdot, \cdot)$ ).

Setze  $\tau_y = \inf\{k \ge 1 : Z_k = y\}$  (Eintritts-/Rückkehrzeit nach y). Für  $\lambda \in [0, 1]$  sei

$$H_{\lambda}(x,y) = \mathbb{E}_{x}[\lambda^{\tau_{y}}]$$
$$G_{\lambda}(x,y) = \mathbb{E}_{x}\left[\sum_{m=1}^{\infty} \lambda^{m} \mathbf{1}(Z_{m} = y)\right] = \sum_{m=1}^{\infty} \lambda^{m} \mathbb{P}_{x}(Z_{m} = y)$$

(und  $H_{\lambda}(x, y) = H_{\lambda}(0, y-x), G_{\lambda}(x, y) = G_{\lambda}(0, y-x)$  wegen Homogenität,  $H_{\lambda}(x, y) = H_{\lambda}(y, x),$  $G_{\lambda}(x, y) = G_{\lambda}(y, x)$  da Z nach Konstruktion symmetrische Inkremente hat). Zerlegung nach dem ersten Besuchszeitpunkt in y zusammen mit der (starken) Markoveigenschaft liefert ("Erneuerungsdarstellung der Besuchswahrscheinlichkeiten in y")

$$G_{\lambda}(x,y) = \sum_{m=1}^{\infty} \lambda^{m} \sum_{k=1}^{m} \mathbb{P}_{x}(\tau_{y} = k, Z_{m} = y) = \sum_{m=1}^{\infty} \sum_{k=1}^{m} \lambda^{m} \mathbb{P}_{x}(\tau_{y} = k) \mathbb{P}_{y}(Z_{m-k} = y)$$
$$= \sum_{k=1}^{\infty} \lambda^{k} \mathbb{P}_{x}(\tau_{y} = k) \sum_{m=k}^{\infty} \lambda^{m-k} \mathbb{P}_{y}(Z_{m-k} = y) = \sum_{k=1}^{\infty} \lambda^{k} \mathbb{P}_{x}(\tau_{y} = k) \Big(1 + \sum_{\ell=1}^{\infty} \lambda^{\ell} \mathbb{P}_{y}(Z_{\ell} = y)\Big)$$
$$= H_{\lambda}(x, y) \Big(1 + G_{\lambda}(y, y)\Big)$$

d.h.

$$H_{\lambda}(x,y) = \frac{G_{\lambda}(x,y)}{1+G_{\lambda}(y,y)} = \frac{G_{\lambda}(0,y-x)}{1+G_{\lambda}(0,0)}$$

Damit liefert (4.4):

$$\begin{split} \mathbb{E}_{(x,y)}[\lambda^{\tau_{\text{coal}}}] &= \mathbb{E}_{(x,y)}[\lambda^{\tau_{1}+1(V\geq 2)\sum_{j=2}^{V}(\tau_{j}-\tau_{j-1})}] = \mathbb{E}_{(x,y)}[\lambda^{\tau_{1}}]\mathbb{E}_{(x,y)}[\lambda^{1(V\geq 2)\sum_{j=2}^{V}(\tau_{j}-\tau_{j-1})}] \\ &= \mathbb{E}_{(x,y)}[\lambda^{\tau_{1}}]\sum_{w=1}^{\infty} \mathbb{P}(V=w) \left(\mathbb{E}_{(0,0)}[\lambda^{\tau_{1}}]\right)^{w-1} \\ &= H_{\lambda}(y-x,0)\sum_{w=1}^{\infty}\frac{1}{N} \left(1-\frac{1}{N}\right)^{w-1}H_{\lambda}(0,0)^{w-1} = \frac{H_{\lambda}(x-y,0)}{N}\frac{1}{1-(1-1/N)H_{\lambda}(0,0)} \\ &= \frac{H_{\lambda}(x-y,0)}{N-(N-1)H_{\lambda}(0,0)} = \frac{G_{\lambda}(x-y,0)}{1+G_{\lambda}(0,0)}\frac{1}{N-(N-1)\frac{G_{\lambda}(0,0)}{1+G_{\lambda}(0,0)}} = \frac{G_{\lambda}(x-y,0)}{N+G_{\lambda}(0,0)} \end{split}$$

und somit

$$\psi(x,y) = \mathbb{E}_{(x,y)} \Big[ (1-\mu)^{2\tau_{\text{coal}}} \Big] = \frac{G_{(1-\mu)^2}(x-y,0)}{N+G_{(1-\mu)^2}(0,0)} = \frac{G_{(1-\mu)^2}(x,y)}{N+G_{(1-\mu)^2}(0,0)}$$

Mit  $1 - \psi(0,0) = N/(N + G_{(1-\mu)^2}(0,0))$  kann man dies als

$$\psi(x,y) = \frac{1 - \psi(0,0)}{N} G_{(1-\mu)^2}(x,y)$$
(4.5)

schreiben.

**Symmetrisches** p Wir nehmen ab hier zusätzlich an, dass p symmetrisch ist, d.h. p(x, y) = p(y, x). Dann ist

$$\widetilde{p}(x,y) = \sum_{z} p(x,z)p(y,z) = \sum_{z} p(x,z)p(z,y) = p^{2}(x,y)$$

d.h.  $\widetilde{p}$  entspricht der zweiten (Matrix-)Potenz von p. In dieser Situation ist

$$G_{\lambda}(x,y) = \sum_{n=1}^{\infty} \lambda^n p^{2n}(x,y)$$
(4.6)

und somit

$$\psi(x,y) = \frac{1 - \psi(0,0)}{N} \sum_{n=1}^{\infty} (1 - \mu)^{2n} p^{2n}(x,y)$$
(4.7)

Im räumlich homogenen Fall ist  $\phi(z) \coloneqq \psi(0, z) = \psi(x, x + z), \phi$  erfüllt

$$\phi(z) = \frac{1 - \phi(0)}{N} \sum_{n=1}^{\infty} (1 - \mu)^{2n} p^{2n}(z)$$
(4.8)

**Bericht 4.2.** Sei X eine  $\mathbb{Z}^d$ -wertige Zufallsvariable mit charakteristischer Funktion  $\varphi_X = \mathbb{E}[e^{-i\langle t, X \rangle}]$ . Dann hat  $\varphi_X$  Periode  $2\pi\mathbb{Z}^d$ , d.h.  $\varphi_X(t + 2\pi m) = \varphi_X(t)$  für  $t \in \mathbb{R}^d$ ,  $m \in \mathbb{Z}^d$ . Es gilt die Fourier-Inversionsformel

$$\mathbb{P}(X=x) = \frac{1}{(2\pi)^d} \int_{(-\pi,\pi]^d} e^{-i\langle x,t\rangle} \varphi_X(t) \, dt, \quad x \in \mathbb{Z}^d.$$

(Idee: für  $z \in \mathbb{Z}^d$  gilt  $\int_{(-\pi,\pi]^d} e^{-i\langle z,t \rangle} dt = (2\pi)^d \mathbf{1}_{\{0\}}(z)$  und somit  $\int_{(-\pi,\pi]^d} e^{-i\langle x,t \rangle} \varphi_X(t) dt = \int_{(-\pi,\pi]^d} e^{-i\langle x,t \rangle} \sum_{y \in \mathbb{Z}^d} e^{i\langle y,t \rangle} dt.$ ) Für  $t = (t_1, \dots, t_d) \in (-\pi, \pi]^d$  ist

$$\widehat{p}(t) = \sum_{x \in \mathbb{Z}^d} e^{i\langle t, x \rangle} p(z)$$

und für  $k \in \mathbb{N}$ 

$$\widehat{p^{k}}(t) = \sum_{x \in \mathbb{Z}^{d}} e^{i\langle t, x \rangle} p^{k}(z) = (\widehat{p}(t))^{k}$$

(Zudem: da wir hier symmetrisches p betrachten, ist  $\widehat{p}(t)$  reell.)

(4.7) liefert

$$\begin{split} \widehat{\phi}(t) &= \frac{1 - \phi(0)}{N} \sum_{n=1}^{\infty} (1 - \mu)^{2n} \widehat{p^{2n}}(t) = \frac{1 - \phi(0)}{N} \sum_{n=1}^{\infty} (1 - \mu)^{2n} (\widehat{p}(t))^{2n} \\ &= \frac{1 - \phi(0)}{N} \frac{(1 - \mu)^2 (\widehat{p}(t))^2}{1 - (1 - \mu)^2 (\widehat{p}(t))^2} \end{split}$$

Fourier-Inversion liefert

$$\phi(x) = \frac{1}{(2\pi)^d} \frac{1 - \phi(0)}{N} \int_{(-\pi,\pi]^d} e^{-i\langle x,t\rangle} \frac{(1-\mu)^2(\widehat{p}(t))^2}{1 - (1-\mu)^2(\widehat{p}(t))^2} dt$$
(4.9)

Für x = 0 liefert (4.9) eine Gleichung für  $\phi(0)$  mit Lösung

$$\phi(0) = (1 + 1/I)^{-1}$$

wobei

$$I = \frac{1}{N} \frac{1}{(2\pi)^d} \int_{(-\pi,\pi]^d} \frac{(1-\mu)^2 (\widehat{p}(t))^2}{1-(1-\mu)^2 (\widehat{p}(t))^2} dt$$
(4.10)

#### "Kleine" Sprungwahrscheinlichkeiten. Betrachte

$$p(x,y) = (1-\nu)\delta_{xy} + \nu q(x,y)$$
(4.11)

mit  $\nu \in (0,1)$ , wobei  $q(\cdot, \cdot)$  eine (homogene: q(x,y) = q(0, y - x) =: q(y - x)) symmetrische stochastische Matrix auf  $\mathbb{Z}^d$  mit 0-Einträgen auf der Diagonale ist (und wir denken an  $\nu \ll 1$ ). Wegen Symmetrie ist  $\sum_z q(z)z = 0$ , wir nehmen zudem an, dass  $\sum_z q(z)||z||^3 < \infty$ .

**Satz 4.3.** Set  $S = \mathbb{Z}$ , p habe die Form (4.11). Für  $0 < \mu, \nu \ll 1$  mit  $\mu \ll \nu$  ist

$$\phi(0) \approx \frac{1}{1 + 2N\sqrt{2\mu\nu\sigma^2}} \tag{4.12}$$

$$\phi(x) \approx \phi(0) \exp\left(-((2\mu)/(\nu\sigma^2))^{1/2}|x|\right)$$
 (4.13)

wobei  $\sigma^2 = \sum_{z \in \mathbb{Z}} q(z) z^2$  die Varianz von q ist. (Wir fassen die Bedeutung des Approximationssymbols  $\approx$  in der folgenden Beweisskizze präziser.)

Bemerkung 4.4. 1. Zum Vergleich: Wenn es nur eine Kolonie gäbe, so wäre

$$\phi(0) = \frac{1/N}{(1/N) + (2\mu - \mu^2)} = \frac{1}{1 + N(2\mu - \mu^2)} \approx \frac{1}{1 + 2N\mu}$$

(verwende z.B. ein "konkurrierende Münzen"-Argument).

2. Man sieht an (4.13), dass

$$L = \sqrt{\nu \sigma^2 / (2\mu)}$$

hier die "charakteristische Längenskala" ist: Wenn man die Stichproben im Abstand  $|x| \gg L$  zieht, ist  $\phi(x) \approx 0$ , für  $|x| \ll L$  ist  $\phi(x) \approx \phi(0)$ .

3. Im Fall der nächste-Nachbar-Irrfahrt in d = 1, d.h. q(-1) = q(1) = 1/2 ist  $\widehat{q}(t) = (e^{it} + e^{-it})/2 = \cos(t)$ ,  $t \in (-\pi, \pi]$  und man könnte das Integral in (4.9) prinzipiell explizit bestimmen.

Beweisskizze für Satz 4.3. Annahme (4.11) liefert

$$\widehat{p}(t) = (1 - \nu) + \nu \widehat{q}(t)$$

zudem ist nach Voraussetzung (mit  $\sigma^2 = \sum_{x \in \mathbb{Z}} x^2 q(x)$ )

$$\widehat{q}(t) = 1 - \frac{\sigma^2}{2}t^2 + r(t)$$

mit "Restterm"  $|r(t)| \leq C|t|^3$  für eine Konstante  $C < \infty$ , somit

$$\frac{(1-\mu)^2(\widehat{p}(t))^2}{1-(1-\mu)^2(\widehat{p}(t))^2} = \frac{(1-\mu)^2(1-\nu(1-\widehat{q}(t)))^2}{1-(1-\mu)^2(1-\nu(1-\widehat{q}(t)))^2} = \frac{(1-\mu)^2(1-\nu\sigma^2t^2/2+\nu r(t))^2}{1-(1-\mu)^2(1-\nu\sigma^2t^2/2+\nu r(t))^2}$$

Für  $\nu \ll 1, \mu \ll 1$  ist

$$(1-\mu)^{2}(1-\nu\sigma^{2}t^{2}/2+\nu r(t))^{2}$$
  
=  $(1-2\mu+\mu^{2})(1-2(\nu\sigma^{2}t^{2}/2-\nu r(t))+(-\nu\sigma^{2}t^{2}/2+\nu r(t))^{2})$   
=  $1-2\mu-\nu\sigma^{2}t^{2}+O(\mu^{2}+\mu\nu+\nu^{2}+\nu r(t))$ 

und somit

$$\frac{(1-\mu)^2(\widehat{p}(t))^2}{1-(1-\mu)^2(\widehat{p}(t))^2} = \frac{1+O(\mu+\nu)}{2\mu+\nu\sigma^2t^2+O(\mu^2+\nu^2+\nu r(t))} = \frac{1}{2\mu+\nu\sigma^2t^2} \cdot \left(1+O(\mu+\nu)\right)$$
(4.14)

und (für das erste Gleichheitszeichen substituiere t = ((2 $\mu$ )/( $\nu\sigma^2$ ))<sup>1/2</sup>u)

$$\int_{-\pi}^{\pi} \frac{1}{2\mu + \nu\sigma^2 t^2} dt = \int_{-\pi((2\mu)/(\nu\sigma^2))^{-1/2}}^{\pi((2\mu)/(\nu\sigma^2))^{-1/2}} \frac{1}{2\mu + 2\mu u^2} \left(\frac{2\mu}{\nu\sigma^2}\right)^{1/2} du$$
$$= \frac{1}{\sqrt{2\mu\nu\sigma^2}} \int_{-\pi((2\mu)/(\nu\sigma^2))^{-1/2}}^{\pi((2\mu)/(\nu\sigma^2))^{-1/2}} \frac{1}{1 + u^2} du \sim \frac{1}{\sqrt{2\mu\nu\sigma^2}} \int_{-\infty}^{\infty} \frac{1}{1 + u^2} du = \frac{\pi}{\sqrt{2\mu\nu\sigma^2}}$$

wenn  $\mu, \nu \to 0 \text{ mit } \mu/\nu \to 0$  (beachte: die Dichte der Standard-Cauchyverteilung ist  $\frac{1}{\pi(1+x^2)}$ ). Mit (4.10) folgt, dass

$$I \sim \frac{1}{2N\sqrt{2\mu\nu\sigma^2}}$$

und damit

$$\phi(0) = \frac{1}{1+1/I} \approx \frac{1}{1+2N\sqrt{2\mu\nu\sigma^2}}$$

(für  $0 < \mu, \nu \ll 1 \text{ mit } \mu \ll \nu$ ).

Für den Ausdruck für  $\phi(x)$  in (4.9) betrachte

$$J = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-itx} \frac{(1-\mu)^2 (\widehat{p}(t))^2}{1-(1-\mu)^2 (\widehat{p}(t))^2} dt$$
(4.15)

Mit (4.14) ergibt sich wie oben (substituiere wieder  $t = ((2\mu)/(\nu\sigma^2))^{1/2}u)$ 

$$J \sim \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{e^{-itx}}{2\mu + \nu\sigma^2 t^2} dt = \frac{1}{2\pi} \int_{-\pi((2\mu)/(\nu\sigma^2))^{-1/2}}^{\pi((2\mu)/(\nu\sigma^2))^{-1/2}} \frac{\exp(-ix((2\mu)/(\nu\sigma^2))^{1/2}u)}{2\mu + 2\mu u^2} \left(\frac{2\mu}{\nu\sigma^2}\right)^{1/2} du$$
$$= \frac{1}{2\pi} \frac{1}{\sqrt{2\mu\nu\sigma^2}} \int_{-\pi((2\mu)/(\nu\sigma^2))^{-1/2}}^{\pi((2\mu)/(\nu\sigma^2))^{-1/2}} \frac{\exp(-ix((2\mu)/(\nu\sigma^2))^{1/2}u)}{1 + u^2} du$$
$$\sim \frac{1}{2\sqrt{2\mu\nu\sigma^2}} \int_{-\infty}^{\infty} \frac{\exp(-ix((2\mu)/(\nu\sigma^2))^{1/2}u)}{\pi(1 + u^2)} du$$
$$= \frac{1}{2\sqrt{2\mu\nu\sigma^2}} \exp(-((2\mu)/(\nu\sigma^2))^{1/2}|x|)$$
(4.16)

(für  $0 < \mu, \nu \ll 1$  mit  $\mu \ll \nu$ ).

Mit (4.9) dann

$$\begin{split} \phi(x) &= \frac{1 - \phi(0)}{N} J \\ &\sim \frac{1}{N} \frac{2N\sqrt{2\mu\nu\sigma^2}}{1 + 2N\sqrt{2\mu\nu\sigma^2}} \frac{1}{2\sqrt{2\mu\nu\sigma^2}} \exp(-((2\mu)/(\nu\sigma^2))^{1/2}|x|) \\ &= \frac{1}{1 + 2N\sqrt{2\mu\nu\sigma^2}} \exp(-((2\mu)/(\nu\sigma^2))^{1/2}|x|) \sim \phi(0) \exp\left(-((2\mu)/(\nu\sigma^2))^{1/2}|x|\right) (4.17) \end{split}$$

Für den Fall  $S = \mathbb{Z}^2$  siehe z.B. Theorem 5.7 in [Duro2, Chapter 5.2].

## Literaturverzeichnis

- [Bir24] Matthias Birkner, Einführung in die Stochastik, https://www.staff.uni-mainz.de/ birkner/GrundlStoch\_2324/Stochastik-Einfuehrung\_WS23\_24.pdf, 2024, Vorlesungsnotizen, JGU Mainz.
- [Dur02] Rick Durrett, *Probability models for DNA sequence evolution*, Probability and its Applications (New York), Springer-Verlag, New York, 2002.
- [Geo15] Hans-Otto Georgii, *Stochastik. Einführung in die Wahrscheinlichkeitstheorie und Statistik*, 5th, revised and expanded ed. ed., Berlin: De Gruyter, 2015 (German).
- [Kle20] Achim Klenke, *Wahrscheinlichkeitstheorie*, 4th revised and supplemented edition ed., Masterclass, Berlin: Springer Spektrum, 2020 (German).
- [KW10] Götz Kersting and Anton Wakolbinger, *Elementare Stochastik.*, 2nd revised ed. ed., Basel: Birkhäuser, 2010 (German).
- [Nor97] J. R. Norris, *Markov chains*, Camb. Ser. Stat. Probab. Math., vol. 2, Cambridge: Cambridge University Press, 1997 (English).
- [Wako9] John Wakeley, *Coalescent theory: an introduction*, vol. 1, Roberts & Company Publishers Greenwood Village, Colorado, 2009.

## Anhang A

# Anhang

## A.1 Ein Exkurs zum Poissonprozess und zu zeitkontinuierlichen Markovketten

Wir diskutieren hier knapp Poissonprozesse auf  $\mathbb{R}_+$  und zeitkontinuierliche Markovketten. Wesentlich mehr dazu findet man z.B. in den Büchern von J. Norris [Nor97] oder von A. Klenke [Kle20, Kap. 17.3].

Sei  $c_N \longrightarrow_{N \to \infty} 0$  eine Nullfolge,  $\lambda > 0$ ,  $Z_i^{(N)}$ ,  $i \in \mathbb{N}$  u.i.v. ~ Ber $(\lambda c_N)$  (wir betrachten o.E. nur so große N, dass  $\lambda c_N \leq 1$ ), seien

$$T_0^{(N)} \coloneqq 0, \quad T_\ell^{(N)} \coloneqq \inf \left\{ i > T_{\ell-1}^{(N)} : Z_i^{(N)} = 1 \right\}, \ \ell \in \mathbb{N}$$

 $(T_{\ell}^{(N)})$  ist der Zeitpunkt des  $\ell$ -ten Erfolgs in der Münzwurffolge  $(Z_i^{(N)})_{i\in\mathbb{N}}$ ), dann sind

$$\tau_{\ell}^{(N)} \coloneqq T_{\ell}^{(N)} - T_{\ell-1}^{(N)}, \ \ell \in \mathbb{N}$$

u.i.v.,  $\tau_{\ell}^{(N)} \sim \operatorname{geom}(\lambda c_N)$ , d.h.  $\mathbb{P}(\tau_{\ell}^{(N)} = j) = c_N \lambda (1 - c_N \lambda)^{j-1}$  für  $j \in \mathbb{N}$  und für  $x \ge 0$  gilt  $\mathbb{P}(c_N \tau_{\ell}^{(N)} > x) = \mathbb{P}(\tau_{\ell}^{(N)} > \lfloor \frac{x}{c_N} \rfloor) = (1 - c_N \lambda)^{\lfloor x/c_N \rfloor} \xrightarrow[N \to \infty]{} e^{-\lambda x},$ 

d.h.  $c_N \tau_{\ell}^{(N)} \xrightarrow{d}_{N \to \infty} \operatorname{Exp}(\lambda)$  (Übung: Beweisen Sie diese Aussagen).

Sei weiter

$$M_k^{(N)} \coloneqq \left| \{ 1 \le i \le k : Z_i^{(N)} = 1 \} \right| = \max\{\ell \in \mathbb{N}_0 : T_\ell^{(N)} \le k \},\$$

offenbar gilt für  $0 \leq k_0 < k_1 < \dots < k_m$ 

$$M_{k_1}^{(N)} - M_{k_0}^{(N)}, M_{k_2}^{(N)} - M_{k_1}^{(N)}, \dots, M_{k_m}^{(N)} - M_{k_{m-1}}^{(N)}$$
 sind unabhängig

und für  $0 \le k < k'$  ist  $M_{k'}^{(N)} - M_k^{(N)} \sim Bin(k' - k, c_N \lambda)$ , somit gilt für  $0 \le t < t'$ 

$$M_{\lfloor t'/c_N \rfloor}^{(N)} - M_{\lfloor t/c_N \rfloor}^{(N)} \xrightarrow{d}_{N \to \infty} \operatorname{Pois}(\lambda(t'-t)).$$

(Übung: Beweisen Sie diese Aussagen).

Dies lädt ein, folgendes Limesobjekt zu betrachten: Sei  $\tau_1, \tau_2, \ldots$  u.i.v.,  $\tau_{\ell} \sim \text{Exp}(\lambda), T_0 := 0, T_{\ell} := \tau_1 + \cdots + \tau_{\ell}, \ell \in \mathbb{N},$ 

$$M_t \coloneqq \max\{i \in \mathbb{N}_0 : T_i \le t\}, \quad t \in [0, \infty)$$

der stochastische Prozess  $(M_t)_{t\geq 0}$  heißt *Poissonprozess* mit Rate  $\lambda$ . (Beachte: die Definition ist so eingerichtet, dass  $t \mapsto M_t$  rechtsstetig ist, man sagt auch:  $(M_t)_t$  hat rechtsstetige Pfade.)

Aus obigen Überlegungen folgt für jedes  $m \in \mathbb{N}$ 

$$(c_N \tau_1^{(N)}, \dots, c_N \tau_m^{(N)}) \xrightarrow{d}_{N \to \infty} (\tau_1, \dots, \tau_m),$$

$$(c_N T_1^{(N)}, \dots, c_N T_m^{(N)}) \xrightarrow{d}_{N \to \infty} (T_1, \dots, T_m)$$

somit ergibt sich für  $t_1 < t_2 < t_m, k_1, \ldots, k_m \in \mathbb{N}_0$ 

$$\mathbb{P}\left(M_{\lfloor t_{1}/c_{N} \rfloor}^{(N)} = k_{1}, \dots, M_{\lfloor t_{m}/c_{N} \rfloor}^{(N)} = k_{m}\right) = \mathbb{P}\left(T_{k_{1}}^{(N)} \le \lfloor t_{1}/c_{N} \rfloor < T_{k_{1}+1}^{(N)}, \dots, T_{k_{m}}^{(N)} \le \lfloor t_{m}/c_{N} \rfloor < T_{k_{m}+1}^{(N)}\right)$$
  
$$\xrightarrow{d}_{N \to \infty} \mathbb{P}\left(T_{k_{1}} \le t_{1} < T_{k_{1}+1}, \dots, T_{k_{m}} \le t_{m} < T_{k_{m}+1}\right) = \mathbb{P}\left(M_{t_{1}} = k_{1}, \dots, M_{t_{m}} = k_{m}\right),$$

d.h. die Folge von stochastischen Prozessen  $(M_{\lfloor t/c_N \rfloor}^{(N)})_{t \ge 0}$  konvergiert gegen den Prozess  $(M_t)_{t \ge 0}$  im Sinne der endlich-dimensionalen Verteilungen.

Aus diesen Beobachtungen folgt

$$\begin{split} & \mathbb{P}\Big(M_{t_1} - M_{t_0} = j_1, M_{t_2} - M_{t_1} = j_2 \dots, M_{t_m} - M_{t_{m-1}} = j_m\Big) \\ &= \lim_{N \to \infty} \mathbb{P}\Big(M_{\lfloor t_1/c_N \rfloor}^{(N)} - M_{\lfloor t_0/c_N \rfloor}^{(N)} = j_1, \dots, M_{\lfloor t_m/c_N \rfloor}^{(N)} - M_{\lfloor t_{m-1}/c_N \rfloor}^{(N)} = j_m\Big) \\ &= \lim_{N \to \infty} \mathbb{P}\Big(M_{\lfloor t_1/c_N \rfloor}^{(N)} - M_{\lfloor t_0/c_N \rfloor}^{(N)} = j_1\Big) \times \dots \times \mathbb{P}\Big(M_{\lfloor t_m/c_N \rfloor}^{(N)} - M_{\lfloor t_{m-1}/c_N \rfloor}^{(N)} = j_m\Big) \\ &= \prod_{i=1}^m e^{-\lambda(t_i - t_{i-1})} \frac{(\lambda(t_i - t_{i-1}))^{j_i}}{j_i!}, \end{split}$$

d.h. die Inkremente eines Poissonprozesses  $(M_t)$  sind Poissonverteilt [der Parameter ist  $\lambda \times$  die Länge des betrachteten Zeitintervalls] und Inkremente über jeweils disjunkte Zeitintervalle sind unabhängig. Diese beiden Eigenschaften charakterisieren den Poissonprozess [ggfs. mit Forderung der Rechtsstetigkeit].

Der Parameter  $\lambda$  kann als Sprungrate interpretiert werden in dem Sinne, dass für ein (kurzes) Zeitintervall (t, t + h] die Wahrscheinlichkeit, einen Sprung in diesem Zeitintervall zu sehen,  $\approx \lambda \times$ Intervalllänge ist, genauer

$$\mathbb{P}(M_{t+h} = k+1 \mid M_t = k) = \mathbb{P}(M_{t+h} - M_t = 1) = e^{-\lambda h} \frac{\lambda h}{1!} = \lambda h + O(h^2) \quad \text{für } h \downarrow 0.$$

**Zu allgemeinen zeitkontinuierlichen Markovketten** Sei E endliche Menge,  $\widehat{p} = (\widehat{p}(x, y))_{x,y\in E}$ stochastische Matrix (d.h.  $\widehat{p}(x, y) \ge 0$ ,  $\sum_{y\in E} \widehat{p}(x, y) = 1$  für alle  $x \in E$ ),  $\widehat{X} = (\widehat{X}_n)_{n\in\mathbb{N}_0}$  (zeitdiskrete, homogene)  $\widehat{p}$ -Markovkette (d.h.  $\mathbb{P}(\widehat{X}_0 = x_0, \widehat{X}_1 = x_1, \dots, \widehat{X}_n = x_n) = \mathbb{P}(\widehat{X}_0 = x_0)\widehat{p}(x_0, x_1) \times \dots \times \widehat{p}(x_{n-1}, x_n)$  für  $x_0, x_1, \dots, x_n \in E$ ). Sei  $(M_t)_{t\geq 0}$  Poissonprozess mit Rate  $\lambda > 0$ , unabhängig von  $\widehat{X}$ ,

$$X_t \coloneqq \widehat{X}_{M_t}, \quad t \in [0, \infty),$$

so ist  $[\widehat{p}^m$  bezeichne die *m*-te Potenz von  $\widehat{p}$ , *I* die  $E \times E$ -Identitätsmatrix]

$$p_{t}(x,y) \coloneqq \mathbb{P}(X_{t} = y \mid X_{0} = x) = \sum_{m=0}^{\infty} \mathbb{P}(M_{t} = m, X_{t} = y \mid X_{0} = x)$$

$$= \sum_{m=0}^{\infty} e^{-\lambda t} \frac{(\lambda t)^{m}}{m!} \widehat{p}^{m}(x,y) = \sum_{m=0}^{\infty} \sum_{n=0}^{\infty} \frac{(\lambda t)^{n}}{n!} \frac{(\lambda t)^{m}}{m!} ((-I)^{n} \widehat{p}^{m})(x,y)$$

$$= \sum_{\ell=0}^{\infty} \frac{t^{\ell}}{\ell!} \lambda^{\ell} \sum_{m=0}^{\ell} {\ell \choose m} (\widehat{p}^{m}(-I)^{\ell-m})(x,y) = \sum_{\ell=0}^{\infty} \frac{t^{\ell}}{\ell!} (\lambda(\widehat{p}-I))^{\ell}(x,y) = \sum_{\ell=0}^{\infty} \frac{t^{\ell}}{\ell!} Q^{\ell}(x,y) = (e^{tQ})(x,y)$$

[wobei die Matrix  $Q = (q_{x,y})_{x,y\in E}$  Einträge  $q_{x,y} = \lambda(\widehat{p}(x,y) - \delta_{xy})$  besitzt; obige Reihe konvergiert, denn  $\max_{x,y\in E} |Q_{xy}^n| \le (|E| \max_{x,y\in E} |Q_{x,y}|)^n$ ; beachte auch, dass  $\widehat{p}^m$  und  $I^n$  kommutieren] und analoge Rechnungen, die die Unabhängigkeit der Zuwächse von  $(M_t)_{t\geq 0}$  ausnutzen, zeigen

$$\mathbb{P}(X_{t_1} = x_1, \dots, X_{t_n} = x_n | X_0 = x_0) = \prod_{i=1}^n p_{t_i - t_{i-i}}(x_{i-1}, x_i)$$

für  $0 = t_0 < t_1 < \dots < t_n, x_0, x_1, \dots, x_n \in E$ .

Die Matrix Q heißt die Sprungratenmatrix (auch: Ratenmatrix oder Q-Matrix) der zeitkontinuierlichen Markovkette X, sie hat die Eigenschaften

$$q_{x,y} \ge 0 \quad \text{für } x \neq y, \qquad \sum_{y \in E, y \neq x} q_{x,y} = -q_{x,x}.$$

Zur Interpretation der Einträge von Q als Sprungraten: Für  $x \neq y$  und  $h \downarrow 0$  ist

$$\mathbb{P}(X_{t+h} = y \mid X_t = x) = (e^{hQ})(x, y) = Q^0(x, y) + hQ^1(x, y) + O(h^2) = hq_{x,y} + O(h^2).$$

Kolmogorovs Vorwärts- und Rückwärtsgleichungen für zeitkontinuierliche Markovketten Für die Sprungratenmatrix  $Q = (q_{x,y})_{x,y\in E}$  einer zeitkontinuierlichen Markovkette auf der endlichen Menge E löst  $\exp(tQ) = \sum_{n=0}^{\infty} \frac{t^n}{n!} Q^n$  für  $t \ge 0$  das System von Differentialgleichungen

$$\frac{\partial}{\partial t} \exp(tQ) = Q \exp(tQ) = \exp(tQ)Q \tag{A.I}$$

demnach für  $t \ge 0$  und  $P_t = e^{tQ} = (p_t(x, y))_{x,y \in E}$ 

$$\frac{\partial}{\partial t}P_t = QP_t, \text{ d.h. } \forall x, y \in E : \quad \frac{\partial}{\partial t}p_t(x, y) = \sum_z q_{x,z}p_t(z, y) = \sum_z q_{x,z} \left( p_t(z, y) - p_t(x, y) \right),$$
(A.2)

$$\frac{\partial}{\partial t}P_t = P_t Q, \text{ d.h. } \forall x, y \in E : \quad \frac{\partial}{\partial t} p_t(x, y) = \sum_z p_t(x, z) q_{z,y} = p_t(x, y) q_{y,y} + \sum_{z \neq y} p_t(x, z) q_{z,y}$$
(A.3)

(beachte: gliedweise Differentiation der Exponentialreihe ist hier erlaubt).

Die Gleichungen (A.2) und (A.3) haben eine stochastische Interpretation. (A.2) heißt Kolmogorovs *Rückwärtsgleichung*, denn es gilt (wir schreiben  $\mathbb{P}_x(\cdot)$  für  $\mathbb{P}(\cdot|X_0 = x)$ )

$$\frac{p_{t+h}(x,y) - p_t(x,y)}{h} = \frac{1}{h} \Big( \mathbb{P}_x (X_{t+h} = y) - \mathbb{P}_x (X_t = y) \Big) \\ = \frac{1}{h} \Big( \sum_z \mathbb{P}_x (X_{t+h} = y | X_h = z) \mathbb{P}_x (X_h = z) - \mathbb{P}_x (X_t = y) \Big) \\ = \frac{1}{h} \Big( \sum_z p_t(z,y) \Big( \mathbb{1}_{\{x=z\}} + hq_{x,z} + o(h) \Big) - p_t(x,y) \Big) = \sum_z q_{x,z} p_t(z,y) + o(1),$$

man leitet sie also aus her durch "Rückwärtszerlegung" des Prozesses X im Intervall [0, t + h] gemäß dem Verhalten am Anfang des Intervalls. Analog heißt (A.3) Kolmogorovs *Vorwärtsgleichung*, sie entsteht durch Zerlegung gemäß dem Wert bei t:

$$\frac{p_{t+h}(x,y) - p_t(x,y)}{h} = \frac{1}{h} \Big( \sum_z \mathbb{P}_x (X_{t+h} = y | X_t = z) \mathbb{P}_x (X_t = z) - \mathbb{P}_x (X_t = y) \Big) \\ = \frac{1}{h} \Big( \sum_z p_t(x,z) \big( \mathbb{1}_{\{z=y\}} + hq_{y,z} + o(h) \big) - p_t(x,y) \Big) = \sum_z q_{x,z} p_t(z,y) + o(1).$$

Sowohl (A.2) als auch (A.3) sind (im Fall  $|E| < \infty$ ) eindeutig lösbar und beide bestimmen die Halbgruppe von Übergangsmatrizen  $(P_t)_{t\geq 0}$  – es sind beides endliche Systeme linearer Differentialgleichungen mit konstanten Koeffizienten.

**Beispiel A.i.** a) Sei  $E = \{0, 1\}, Q = {\binom{-a \ a}{b \ -b}} \text{ mit } a, b > 0$ . Für  $x, y \in \{0, 1\}$  gilt  $p_t(x, y) = \delta_{x,y}e^{-(a+b)t} + (1 - e^{-(a+b)t})\mu(y) \text{ mit } \mu(0) = b/(a+b), \mu(1) = a/(a+b).$ b) Sei  $E = \{0, 1, 2, 3\}$  (oder auch  $E = \{A, G, C, T\}$ ) und

$$Q = \begin{pmatrix} -3 & 1 & 1 & 1 \\ 1 & -3 & 1 & 1 \\ 1 & 1 & -3 & 1 \\ 1 & 1 & 1 & -3 \end{pmatrix}.$$

Für  $x, y \in E$  gilt  $p_t(x, y) = \delta_{x,y} e^{-4t} + \frac{1}{4}(1 - e^{-4t}).$ 

**Lemma A.2.** E endliche Menge,  $Q = (q_{xy})_{x,y \in E}$  Sprungratenmatrix,  $X^{(N)}$ ,  $N \in \mathbb{N}$  zeitdiskrete Ewertige Markovketten mit Übergangsmatrix

$$p^{(N)}(x,y) = \delta_{x,y} + c_N q_{xy} + o(c_N), \quad x, y \in E,$$

wo  $c_N \to 0$  für  $N \to \infty$  und  $X_0^{(N)} = x_0 \in E$ . Dann konvergieren die (zeitlich reskalierten) Prozesse  $(X_{\lfloor t/c_N \rfloor}^{(N)})_{t\geq 0}$  für  $N \to \infty$  gegen die zeitkontinuierliche Markovkette X mit Sprungratenmatrix Q (im Sinne der endlich-dimensionalen Verteilungen).

*Beweis.* Wir schreiben die Übergangsmatrix von  $X^{(N)}$  als

$$p^{(N)} = I + c_N Q_N$$

mit  $Q_N \coloneqq c_N^{-1}(p^{(N)} - I)$ , somit nach Voraussetzung  $Q_N \xrightarrow[N \to \infty]{} Q$  (eintrags-weise).

$$\begin{split} (I+c_NQ_N)^{\lfloor c_N^{-1}t\rfloor} &= \sum_{k=0}^{\lfloor c_N^{-1}t\rfloor} \binom{\lfloor c_N^{-1}t\rfloor}{k} c_N^k Q_N^k = \sum_{k=0}^{\lfloor c_N^{-1}t\rfloor} c_N^k \frac{\lfloor c_N^{-1}t\rfloor(\lfloor c_N^{-1}t\rfloor-1)\cdots(\lfloor c_N^{-1}t\rfloor-k+1)}{k!} Q_N^k \\ &\to \sum_{k=0}^{\infty} \frac{t^k}{k!} Q^k = e^{tQ} \qquad \text{für } N \to \infty. \end{split}$$

Beachte: Für  $k \in \mathbb{N}$  gilt

$$c_N^k \frac{\lfloor c_N^{-1}t \rfloor (\lfloor c_N^{-1}t \rfloor - 1) \cdots (\lfloor c_N^{-1}t \rfloor - k + 1)}{k!} Q_N^k \underset{N \to \infty}{\longrightarrow} \frac{t^k}{k!} Q^k$$

(eintrags-weise) und die Beträge der Einträge der Matrix auf der linken Seite sind für genügend großes ${\cal N}$ 

$$\leq t^{k} (2 \max\{|Q(x,y)| : x, y \in E\})^{k} / k!,$$

so dass Grenzwert und Summation vertauscht werden können. Somit

$$\mathbb{P}\left(X_{\lfloor c_N^{-1}t_1 \rfloor}^{(N)} = x_1, \dots, X_{\lfloor c_N^{-1}t_n \rfloor}^{(N)} = x_n\right) = \prod_{j=1}^n (I + c_N Q_N)^{\lfloor c_N^{-1}(t_j - t_{j-1}) \rfloor} (x_{j-1}, x_j)$$
$$\xrightarrow[N \to \infty]{} \prod_{j=1}^n (e^{(t_j - t_{j-1}Q)})(x_{j-1}, x_j) = \mathbb{P}\left(X_{t_1} = x_1, \dots, X_{t_n} = x_n\right).$$

### A.2 (Weitere) Eigenschaften des Kingman-Koaleszenten

Sei  $\xi_i^{(n)}$ , i = n, n-1, ..., 1 der Zustand des *n*-Koaleszenten zum ersten Zeitpunkt, zu dem *i* Klassen existieren, d.h.  $\xi_i^{(n)} = R_{\tau_i^{(n)}}^{(n)}$  (mit  $\tau_i^{(n)} \coloneqq \inf\{t \ge 0 \colon |R_t^{(n)}| \le i\}$  wie oben).  $[(\xi_i^{(n)})_{i=n,n-1,...,1}$  heißt die *Skelettkette* des Kingman-*n*-Koaleszenten, sie ist (offenbar) eine Markovkette.]

**Proposition A.3.** Für  $\xi \in \mathcal{E}_n$  mit i Klassen der Größen  $\lambda_1, \ldots, \lambda_i \in \mathbb{N}$  (mit  $\lambda_1 + \cdots + \lambda_i = n$ ) gilt

$$\mathbb{P}(\xi_i^{(n)} = \xi) = c_{n,i}w(\xi) \quad mit \ w(\xi) = \lambda_1! \cdots \lambda_i!, \ c_{n,i} = \frac{i!}{n!} \frac{(n-i)!(i-1)!}{(n-1)!}.$$
(A.4)

**Beispiel.** Betrachte n = 9, i = 3, es ist  $c_{9,3} = \frac{3!}{9!} \frac{6!2!}{8!} = 1/1693440$ .

Wir sehen: die Verteilung hat mehr Gewicht auf "unbalanzierten Aufteilungen."

Beweis von Prop. A.3. Rückwärtsinduktion über i: Für i = n gilt  $\mathbb{P}(\xi_n^{(n)} = \{\{1\}, \dots, \{n\}\}) = 1$  mit  $\lambda_1 = \dots = \lambda_n = 1$ , und  $c_{n,n} = w(1, \dots, 1) = 1$ .  $i \to i - 1$ : Es ist

$$\mathbb{P}(\xi_{i-1}^{(n)} = \eta \,|\, \xi_i^{(n)} = \xi) = \begin{cases} \frac{1}{\binom{i}{2}}, & \text{falls } \eta \text{ aus } \xi \text{ durch Verschmelzung eines Paars von Klassen} \\ 0, & \text{entsteht,} \\ 0, & \text{sonst.} \end{cases}$$

Sei  $\eta \in \mathcal{E}_n, |\eta| = i - 1$ , Klassengrößen  $\tilde{\lambda}_1, \ldots, \tilde{\lambda}_{i-1}$ .

$$\mathbb{P}(\xi_{i-1}^{(n)} = \eta) = \frac{2}{i(i-1)} \sum_{\xi:\xi < \eta} \mathbb{P}(\xi_{i}^{(n)} = \xi) \\
= \frac{2}{i(i-1)} \sum_{\ell=1}^{i-1} \sum_{m=1}^{\tilde{\lambda}_{\ell}-1} \frac{1}{2} {\tilde{\lambda}_{\ell} \choose m} c_{n,i} \tilde{\lambda}_{1}! \cdots \tilde{\lambda}_{\ell-1}! m! (\tilde{\lambda}_{\ell} - m)! \tilde{\lambda}_{\ell+1}! \cdots \tilde{\lambda}_{i-1}! \\
= \frac{c_{n,i}}{i(i-1)} \sum_{\ell=1}^{i-1} \sum_{m=1}^{\tilde{\lambda}_{\ell}-1} w(\eta) = \frac{c_{n,i}w(\eta)}{i(i-1)} \sum_{\ell=1}^{i-1} \sum_{m=1}^{\tilde{\lambda}_{\ell}-1} 1 \\
= \frac{c_{n,i}w(\eta)}{i(i-1)} (n - (i-1)) = c_{n,i-1}w(\eta)$$

Für das 2. Gleichheitszeichen verwenden wir die Induktionsannahme und zerlegen gemäß der "aufgespaltenen" Klasse: die  $\ell$ -te Klasse hat  $\tilde{\lambda}_{\ell}$  Elemente, zerlege in 2 Teilmengen d. Größen m und  $\tilde{\lambda}_{\ell} - m$ , es gibt  $\frac{1}{2} {\tilde{\lambda}_{\ell} \choose m}$  mögliche Wahlen; der Faktor  $\frac{1}{2}$  entsteht, weil die Klassen in  $\eta$  als ungeordnet aufgefasst werden.

**Korollar A.4.** I. Sei  $\sigma$  eine uniform verteilte Permutation von  $\{1, \ldots, i\}$ , u.a. von  $\xi^{(n)}$ ,  $M_i = |C_{i,\sigma(i)}^{(n)}|$ mit  $\xi_i^{(n)} = \{C_{i,1}^{(n)}, \ldots, C_{i,i}^{(n)}\}$ . Dann ist

 $(M_1,\ldots,M_i)$  uniform verteilt auf  $\{(m_1,\ldots,m_i)\in\mathbb{N}^i:m_1+\cdots+m_i=n\}$ .

2. Sei  $\xi = \{A_1, \ldots, A_{i-1}\} \in \mathcal{E}_n$  mit  $|A_j| = \lambda_j$ .  $\mathscr{L}(\xi_i^{(n)} | \xi_{i-1}^{(n)} = \xi)$  kann folgendermaßen beschrieben werden:

- wähle  $A_j$  mit W'keit  $\frac{\lambda_j-1}{n-i+1}$  ( $j \in \{1, \ldots, i-1\}$ ), dann wähle k uniform aus  $\{1, \ldots, \lambda_j-1\}$
- spalte  $A_j$  uniform in zwei Teile der Größen k und  $\lambda_j k$ .

**Bemerkung.** Für i = 2 zeigt Kor. A.4 die "uniforme Aufspaltung" der Stichprobe in zwei älteste Familien.

*Beweis von Kor. A.4.* I. Seien  $m_1, \ldots, m_i$  mit  $m_1 + \cdots + m_i = n$  gegeben. Nach Prop. A.3 hat jede Realisierung des (zufällig) geordneten Vektors  $(C_{i,\sigma(1)}^{(n)}, \ldots, C_{i,\sigma(i)}^{(n)})$ , die mit den geforderten Größen  $m_j$  verträglich ist, die W'keit  $\frac{1}{i!}c_{n,i}m_1!\cdots m_i!$ , somit

$$\mathbb{P}((M_1, \dots, M_i) = (m_1, \dots, m_i)) = \binom{n}{m_1 \dots m_i} \frac{1}{i!} c_{n,i} m_1! \cdots m_i!$$
$$= \frac{(n-i)!(i-1)!}{(n-1)!} = \frac{1}{\binom{n-1}{i-1}},$$

denn es gibt  $\binom{n}{m_1 \dots m_i}$  Partitionen, die bezgl. der Größe der Klassen in Frage kommen,

jede hat n. Prop. A.3 dieselbe W'keit  $c_{n,i}m_1!\cdots m_i!$ ,

die W'keit, dass zuf. Perm.  $\sigma$  geg. Ordnung liefert, ergibt nochmals einen Faktor  $\frac{1}{i!}$ .

(Beachte auch  $\#\{\{(m_1, \ldots, m_i) \in \mathbb{N}^i : m_1 + \cdots + m_i = n\}\} = \binom{n-1}{i-1}: n$  Kugeln in *i* (nummerierte) Schachteln legen, so dass keine Schachtel leer ist: n - 1 mögl. Plätze für i - 1 "Trennwände".) 2.  $\xi$  entstehe aus  $\eta$  durch Aufteilen von  $A_j$  in 2 Teile der Größen k und  $\lambda_j - k$ .

$$\mathbb{P}(\xi_{i}^{(n)} = \xi | \xi_{i-1}^{(n)} = \eta) = \frac{\mathbb{P}(\xi_{i-1}^{(n)} = \eta | \xi_{i}^{(n)} = \xi) \mathbb{P}(\xi_{i}^{(n)} = \xi)}{\mathbb{P}(\xi_{i-1}^{(n)} = \eta)} \\
= \frac{\frac{1}{\binom{i}{2}} c_{n,i} \lambda_{1}! \cdots \lambda_{j-1}! k! (\lambda_{j} - k)! \lambda_{j+1}! \cdots \lambda_{i-1}!}{c_{n,i-1} \lambda_{1}! \cdots \lambda_{j-1}! \lambda_{j}! \lambda_{j+1}! \cdots \lambda_{i-1}!} = \frac{1}{\binom{i}{2}} \cdot \frac{i(i-1)}{n-i+1} \frac{1}{\binom{\lambda_{j}}{k}} \\
= \frac{\lambda_{j} - 1}{n-i+1} \cdot \frac{1}{\lambda_{j} - 1} \cdot 2\frac{1}{\binom{\lambda_{j}}{k}}$$

 $\left(\frac{\lambda_j-1}{n-i+1} \cong \text{Wahl von } A_j; \frac{1}{\lambda_j-1} \cong \text{Wahl von } k; 2\frac{1}{\binom{\lambda_j}{k}} \cong \text{Wahl der Zerlegung von } A_j - \text{beachte: Faktor 2, da}$ die Klassen "ungeordnet" angegeben werden)

**Korollar A.5.** Betrachte eine Teilstichprobe der Grösse n in einem Kingman-m-Koaleszenten, mit m > n. Dann erfüllt die Wahrscheinlichkeit des Ereignisses  $E_{m,n}$ , dass der jüngste gemeinsame Vorfahre (jgV) der n-Stichprobe mit der Wurzel des m-Koaleszenten übereinstimmt,

$$\mathbb{P}(E_{m,n}) \to \frac{n-1}{n+1} \quad f \ddot{u} r m \to \infty.$$

*Beweis.* Wir betrachten  $\xi_2^{(n)}$ , die erste Aufspaltung des Koaleszenten von der Wurzel aus betrachtet. Diese resultiert in einer Aufspaltung der trivialen Partition  $\{\{1, \ldots, m\}\}$  in eine Äquivalenzrelation mit genau zwei Klassen der Grössen m - X und X, wobei X nach Kor. A.4 auf [m - 1] uniform verteilt ist. Falls der jgV der n-Stichprobe nicht mit der Wurzel übereinstimmt, so müssen die Ahnenlinien aller n Individuen der Stichprobe alle entweder in dem Block der Grösse m - X oder in dem Block der Grösse X liegen.

Das erste Ereignis hat Wahrscheinlichkeit  $\frac{(m-X)_{n\downarrow}}{(m)_{n\downarrow}}$  und das zweite  $\frac{(X)_{n\downarrow}}{(m)_{n\downarrow}}$ . Wir erhalten

$$\mathbb{P}(E_{mn,}) = 1 - \mathbb{P}((E_{mn,})^{c}) = 1 - \sum_{k=1}^{m-1} \left[ \frac{(m-k)_{n\downarrow}}{(m)_{n\downarrow}} + \frac{(k)_{n\downarrow}}{(m)_{n\downarrow}} \right] \underbrace{\mathbb{P}(X=k)}_{\frac{1}{m-1}}$$
$$\xrightarrow{m\to\infty} 1 - \int_{0}^{1} (x^{n} + (1-x)^{n}) \, dx$$
$$= 1 - \left[ \frac{1}{n+1} x^{n+1} \right]_{0}^{1} - \left[ \frac{1}{n+1} (1-x)^{n+1} (-1) \right]_{0}^{1} = 1 - \frac{2}{n+1}$$

für  $m \rightarrow \infty$  durch Konvergenz der Riemann-Summe gegen das Riemann-Integral.

### A.3 Die Verteilung der Summe unabhängiger, exponentialverteilter Zufallsvariablen

**Lemma A.6.** Seien  $X_1, X_2, \ldots, X_n$  unabhängig,  $X_i$  sei exponentialverteilt mit Parameter  $\lambda_i$  und  $\lambda_1 < \lambda_2 < \cdots < \lambda_n$ . Dann hat  $X \coloneqq X_1 + \cdots + X_n$  die Dichte

$$f_X(t) = \sum_{j=1}^n \lambda_j \exp(-\lambda_j t) \prod_{k=1, k\neq j}^n \frac{\lambda_k}{\lambda_k - \lambda_j},$$

insbesondere ist (mit  $a_j := \prod_{k=1, k\neq j}^n \frac{\lambda_k}{\lambda_k - \lambda_j}$ )

$$\mathbb{P}(X > t) = \sum_{j=1}^{n} a_j \exp(-\lambda_j t), \ t \ge 0.$$

*Beweisskizze.* Die Formel für die Dichte kann man beispielsweise per Induktion durch sukzessive Faltung mit der Exponentialdichte beweisen, für den Induktionsschritt beachten wir

$$\int_{0}^{t} \sum_{j=1}^{n-1} \lambda_{j} \exp(-\lambda_{j}s) \prod_{k=1,k\neq j}^{n-1} \frac{\lambda_{k}}{\lambda_{k} - \lambda_{j}} \times \lambda_{n} \exp(-\lambda_{n}(t-s)) ds$$

$$= \sum_{j=1}^{n-1} \lambda_{j} \lambda_{n} \prod_{k=1,k\neq j}^{n-1} \frac{\lambda_{k}}{\lambda_{k} - \lambda_{j}} \times e^{-\lambda_{n}t} \int_{0}^{t} e^{(\lambda_{n} - \lambda_{j})s} ds = \sum_{j=1}^{n-1} \lambda_{j} \lambda_{n} \prod_{k=1,k\neq j}^{n-1} \frac{\lambda_{k}}{\lambda_{k} - \lambda_{j}} \times \frac{e^{-\lambda_{j}t} - e^{-\lambda_{n}t}}{\lambda_{n} - \lambda_{j}}$$

$$= \sum_{j=1}^{n-1} \lambda_{j} e^{-\lambda_{j}t} \prod_{k=1,k\neq j}^{n} \frac{\lambda_{k}}{\lambda_{k} - \lambda_{j}} - \lambda_{n} e^{-\lambda_{n}t} \sum_{j=1}^{n-1} \frac{\lambda_{j}}{\lambda_{n} - \lambda_{j}} \prod_{k=1,k\neq j}^{n-1} \frac{\lambda_{k}}{\lambda_{k} - \lambda_{j}}.$$

Dann verwenden wir die Identität

$$\sum_{j=1}^{n-1} \frac{\lambda_j}{\lambda_n - \lambda_j} \prod_{k=1, k \neq j}^{n-1} \frac{\lambda_k}{\lambda_k - \lambda_j} = -\prod_{k=1}^{n-1} \frac{\lambda_k}{\lambda_k - \lambda_n},$$

die (man dividiere beide Seiten durch  $\lambda_1 \lambda_2 \cdots \lambda_{n-1}$  und sortiere Terme) äquivalent ist zu

$$\sum_{j=1}^{n} \prod_{k=1, k\neq j}^{n} \frac{1}{\lambda_k - \lambda_j} = 0.$$
 (A.5)

Sei  $\ell_j(x) \coloneqq \prod_{k=1,k\neq j}^n \frac{x-\lambda_k}{\lambda_j-\lambda_k}$  (das *j*-te Lagrange-Polynom zu  $\lambda_1, \ldots, \lambda_n$ ),  $\ell_1(x) + \ell_2(x) + \cdots + \ell_n(x)$  ist ein Polynom in x vom Grad n-1, das (mindestens) an den n verschiedenen Stellen  $\lambda_1, \ldots, \lambda_n$  den Wert 1 annimmt (denn  $\ell_j(\lambda_i) = \delta_{ji}$ ), daher gilt  $\ell_1(x) + \ell_2(x) + \cdots + \ell_n(x) \equiv 1$ . Die linke Seite von (A.5) ist  $(-1)^{n-1}$  mal der Koeffizient von  $x^{n-1}$  in diesem Polynom.

Alternativ beachte man, dass für  $\zeta \in \mathbb{R}$  ist  $\mathbb{E}[e^{i\zeta X_j}] = \int_0^\infty e^{i\zeta x} \lambda_j e^{-\lambda_j x} dx = \frac{\lambda_j}{\lambda_j - i\zeta}$ , also  $\mathbb{E}[e^{i\zeta X}] = \prod_{j=1}^n \frac{\lambda_j}{\lambda_j - i\zeta} =: \varphi_1(\zeta)$  gilt, während  $\int_0^\infty e^{i\zeta x} \sum_{j=1}^n a_j \lambda_j e^{-\lambda_j x} dx = \sum_{j=1}^n \frac{a_j \lambda_j}{\lambda_j - i\zeta} =: \varphi_2(z)$  und es *ist*  $\varphi_2 = \varphi_1(\varphi_2)$  ist die Partialbruchzerlegungs-Darstellung von  $\varphi_1$ ), denn beide sind meromorph auf  $\mathbb{C}$  mit jeweils einfachen Polen bei  $\zeta = -i\lambda_1, \ldots, -i\lambda_n$  und  $\lim_{|z|\to\infty} \varphi_{1/2}(z) = 0$ ,  $\lim_{z\to -i\lambda_j} \frac{\varphi_1(z)}{\lambda_j - iz} = \lambda_j \prod_{k=1, k\neq j}^n \frac{\lambda_k}{\lambda_k - \lambda_j} = \lim_{z\to -i\lambda_j} \frac{\varphi_2(z)}{\lambda_j - iz}$ .

Die Formel für den Verteilungsschwanz von X ergibt sich durch entsprechendes Integrieren der Dichte.

### A.4 Erwartete Fixationszeit im Wright-Fisher-Modell: exakte Rechnung

In diesem Kapitel führen wir den vollständigen Beweis von Satz 1.1 aus, d.h. für die Fixationszeit  $T_{\text{fix}}^{(N)} = \inf\{r \ge 0 : X_r^{(N)} = 0 \text{ oder } X_r^{(N)} = N\}$  im neutralen 2-Typ-Wright-Fisher-Modell mt

Populationsgröße N gilt

$$\lim_{N \to \infty} c_N \frac{\mathbb{E}_{x_N}[T_{\text{fix}}^{(N)}]}{2N} = H(p) = -p \log(p) - (1-p) \log(1-p)$$
(1.5)

falls  $x_N/N \rightarrow p \in [0, 1]$ .

Wir schreiben im Folgenden  $p_r := X_r^{(N)}/N$  für den Anteil von Typ A-Individuen in der r-ten Generation (und unterdrücken in der Notation, dass dessen Verteilung natürlich auch von N abhängt), sowie  $\mathbb{P}_p$  bzw.  $\mathbb{E}_p$  für Wahrscheinlichkeiten bzw. Erwartungswerte in der Sitation, dass der Startanteil  $p_0 = X_0^{(N)}/N = p$  beträgt.

Wir zeigen zunächst die obere Schranke

$$\mathbb{E}_p\left[T_{\text{fix}}^{(N)}\right] \le 2NH(p) \tag{A.6}$$

Für  $p \in \{0, 1\}$  gilt die Aussage trivialerweise. Für  $p \in \{\frac{1}{N}, \dots, \frac{N-1}{N}\}$  liefert die "Ein-Schritt-Analyse" (1.4)

$$\mathbb{E}_p[T_{\text{fix}}^{(N)}] = 1 + \sum_{q \in \{0, \frac{1}{N}, \dots, 1\}} \mathbb{P}_p\{p_1 = q\} \mathbb{E}_q[T_{\text{fix}}^{(N)}].$$
(A.7)

Unser Ziel ist nun zu zeigen, dass für alle  $p \in I_N = \{0, \frac{1}{N}, \dots, 1\}$  und  $N \in \mathbb{N}$  die Funktion

$$f_N(p) \coloneqq 2NH(p) - \mathbb{E}_p[T_{\text{fix}}^{(N)}],$$

die Ungleichung

$$\sum_{q \in I_N} \mathbb{P}_p(p_1 = q) \big[ f_N(q) - f_N(p) \big] \le 0, \tag{A.8}$$

mit  $f_N(0) = f_N(1) = 0$  erfüllt. Denn damit ist  $f_N$  superharmonisch und  $\{f_N(p_r)\}_{r\geq 0}$  ein  $\mathcal{F}^N$ -Supermartingal [KW10, Kapitel 2.7]. Mit dem Doobschen Stoppsatz [KW10, Satz 1.10] folgt für jedes für  $p \in I_N$ , dass

$$f_N(p) \ge \mathbb{E}_p\left[f_N(p_{T_{\text{fix}}^{(N)}})\right] = 0$$

und damit gilt

$$f_N(p) = 2NH(p) - \mathbb{E}_p[T_{\text{fix}}^{(N)}] \ge 0$$

und dies ist gerade die Aussage (A.6).

Es bleibt also noch (A.8) zu zeigen, was für Randfälle  $p \in \{0, 1\}$  offensichtlich gilt. Angesichts von (A.7) ist (A.8) für  $p \in \{\frac{1}{N}, \dots, \frac{N-1}{N}\}$  äquivalent zu folgender Ungleichung:

$$2N\sum_{q\in I_N} \mathbb{P}_p(p_1=q) \Big[ H(q) - H(p) \Big] \le -1 \tag{A.9}$$

Da die Funktion  $H(\cdot)$  symmetrisch bezüglich Spiegelung um p = 1/2 ist, d.h. invariant unter der Ersetzung  $p \mapsto 1 - p$ , gilt

$$\sum_{q \in I_N} \mathbb{P}_{1-p}(p_1 = q) \Big[ H(q) - H(1-p) \Big] = \sum_{q \in I_N} \mathbb{P}_{1-p} \{ p_1 = 1-q \} \Big[ H(1-q) - H(1-p) \Big]$$
$$= \sum_{q \in I_N} \mathbb{P}_p(p_1 = q) \Big[ H(q) - H(p) \Big]$$

Daher genügt es, den Fall 0 zu betrachten.*H*ist im Innern von <math>[0, 1] beliebig oft stetig differenzierbar; Taylor-Entwicklung bis zur ersten Ordnung mit Restglied in Integralform liefert

$$H(q) = H(p) + (q-p)H'(p) + \int_{p}^{q} (q-t)H''(t) dt$$
  
=  $H(p) + (q-p)H'(p) + (q-p)^{2} \int_{0}^{1} H''(p+u(q-p))(1-u) du, \quad q \in I_{N}$ 

wobei wir in der zweiten Zeile t = p + u(q - p) substituiert haben. Wir definieren ein Wahrscheinlichkeitsmaß  $\tilde{\mu}_{N,p}$  auf (0, 1) mittels

$$\int_{(0,1)} f(x) \,\tilde{\mu}_{N,p}(dx) \coloneqq \frac{N}{p(1-p)} \int_0^1 \mathbb{E}_p \Big[ f \Big( p + u(p_1 - p) \Big) (p_1 - p)^2 \Big] 2(1-u) \, du$$

(siehe Lemma A.7 (a)).

Damit schreibt sich

$$\sum_{q \in I_N} \mathbb{P}_p(p_1 = q) \Big[ H(q) - H(p) \Big] = \frac{p(1-p)}{2N} \int_{(0,1)} H''(x) \,\tilde{\mu}_{N,p}(dx)$$
$$= -\frac{p(1-p)}{2N} \int_{(0,1)} \varphi(x) \,\tilde{\mu}_{N,p}(dx) \le -\frac{p(1-p)}{2N} \,\varphi\bigg(\int_{(0,1)} x \,\tilde{\mu}_{N,p}(dx)\bigg)$$

mit  $\varphi(x) = \frac{1}{x(1-x)}$ , wobei das Ungleichungszeichen aus der Konvexität der Funktion  $\varphi$  auf (0, 1) und der Jensenschen Ungleichung folgt. Weiter ist

$$\int_{(0,1)} x \,\tilde{\mu}_{N,p}(dx) = \frac{N}{2p(1-p)} \int_0^1 \mathbb{E}_p \Big[ \big( p + u(p_1 - p) \big) (p_1 - p)^2 \Big] 2(1-u) \, du$$
$$= \frac{N}{2p(1-p)} \Big( p \mathbb{E}_p \Big[ (p_1 - p)^2 \Big] + \frac{1}{3} \mathbb{E}_p \Big[ (p_1 - p)^3 \Big] \Big)$$

wobei wir den Satz von Fubini sowie  $\int_0^1 2u(1-u) \, du = 1/3$  verwenden. Es ist

$$\mathbb{E}_p[(p_1 - p)^3] = \frac{1}{N^3} \mathbb{E}[(Y_{N,p} - Np)^3] = \frac{1}{N^2} p(1 - p)(1 - 2p)$$

mit  $Y_{N,p} \sim \operatorname{Bin}(N,p)$  und  $\mathbb{E}_p[(p_1 - p)^2] = p(1-p)/N$  (siehe Lemma A.7 (b)).

Damit finden wir

$$\int_{(0,1)} x \,\tilde{\mu}_{N,p}(dx) = \frac{p}{2} + \frac{1-2p}{6N}$$

und daher

$$\sum_{q \in I_N} \mathbb{P}_p(p_1 = q) \Big[ H(q) - H(p) \Big] \\ \leq -\frac{p(1-p)}{2N} \frac{1}{\left( p/2 + (1-2p)/(6N) \right) \left( 1 - p/2 - (1-2p)/(6N) \right)} \leq -\frac{1}{2N}$$

(beachte  $\frac{p(1-p)}{(p/2+(1-2p)/(6N))(1-p/2-(1-2p)/(6N))} \ge 1$  für  $p \le 1/2$ ). Dies komplettiert den Beweis von (A.9) und damit auch der Aussage (A.6).

Um die Asymptotik von  $\mathbb{E}_p[T_{\text{fix}}^{(N)}]$  aus (1.5) zu zeigen, betrachten wir nun die Taylor-Entwicklung von *H* bis zur zweiten Ordnung. Sei zunächst  $\varepsilon \in (0, \frac{1}{2})$  und setze

$$T_{\varepsilon}^{(N)} := \inf\{r : p_r < \varepsilon \text{ oder } p_r > 1 - \varepsilon\}.$$

Für  $\zeta \in (\varepsilon/2, 1 - \varepsilon/2)$  ist

$$H'''(\zeta) = \frac{1 - 2\zeta}{\zeta^2 (1 - \zeta^2)}$$

gleichmäßig beschränkt. Taylor-Entwicklung bis zur zweiten Ordnung mit Restglied in Integralform liefert

$$H(q) = H(p) + (q-p)H'(p) + \frac{(q-p)^2}{2}H''(p) + \frac{(q-p)^3}{2}\int_0^1 H'''(p+u(q-p))(1-u)^2 du$$

Weiter ist

$$H'''(p+u(q-p))(1-u)^{2}| \leq \frac{(1-u)^{2}}{((1-u)p+uq)^{2}(1-(1-u)p-uq)^{2}} \leq \frac{1}{p^{2}(1-p)^{2}} \leq \varepsilon^{-4}$$

Für das zweite Ungleichungszeichen betrachten wir hier eine die Fallunterscheidung: Falls q < p, so ist der erste Term in der zweiten Zeile höchstens

$$(1-u)^2/(((1-u)p+0)(1-(1-u)p-up))^2 = 1/(p(1-p))^2$$

falls  $q \ge p$  gilt, so ist er höchstens  $(1-u)^2/(p(1-(1-u)p-u))^2 = 1/(p(1-p))^2$ .

Somit gilt

$$\sum_{\substack{q \in \{0, \frac{1}{N}, \dots, 1\}\\q \in \{0, \frac{1}{N}, \dots, 1\}}} \mathbb{P}_p(p_1 = q) \Big[ H(q) - H(p) \Big]$$
  
= 
$$\sum_{\substack{q \in \{0, \frac{1}{N}, \dots, 1\}\\q \in \{0, \frac{1}{N}, \dots, 1\}}} \mathbb{P}_p(p_1 = q) \Big[ (q - p)H'(p) + \frac{1}{2}(q - p)^2 H''(p) + \frac{1}{6}(q - p)^3 H''(\zeta(q, p)) \Big]$$
  
= 
$$-\frac{1}{2} + \frac{1}{6} \sum_{\substack{q \in \{0, \frac{1}{N}, \dots, 1\}\\q \in \{0, \frac{1}{N}, \dots, 1\}}} \mathbb{P}_p(p_1 = q)(q - p)^3 H'''(\zeta(q, p))$$
  
= 
$$-\frac{1}{2} + R_N(p)$$

wobei das Restglied die Abschätzung

$$\max_{p \in \{0,\frac{1}{N},\dots,1\} \cap (\varepsilon,1-\varepsilon)} |R_N(p)| \le \varepsilon^{-4} \max_{p \in \{0,\frac{1}{N},\dots,1\}} \mathbb{E}_p\left[|p_1-p|^3\right] \le \frac{1}{\varepsilon^4 N^{3/2}}$$

erfüllt: Für  $p \in [0, 1]$  ist mit der Jensenschen Ungleichung und Lemma A.7 (b)

$$\mathbb{E}_p[|p_1 - p|^3] \le \left(\mathbb{E}_p[|p_1 - p|^4]\right)^{3/4} \le \left(N^{-2}\right)^{3/4}$$

Die übliche Ein-Schritt-Analyse zusammen mit obigem zeigt, dass

$$f_{N,\varepsilon}(p) \coloneqq 2NH(p) - \mathbb{E}_p[T_{\varepsilon}^{(N)}]$$

für  $p \in \{0, \frac{1}{N}, \dots, 1\} \cap (\varepsilon, 1 - \varepsilon)$  die Gleichung

$$\sum_{q \in \{0,\frac{1}{N},\dots,1\}} \mathbb{P}_p(p_1 = q) \left[ f_{N,\varepsilon}(q) - f_{N,\varepsilon}(p) \right] = 2NR_N(p)$$

erfüllt. Zudem gilt für  $p \in \{0, \frac{1}{N}, \dots, 1\} \cap ([0, \varepsilon] \cap [1 - \varepsilon, 1])$ 

$$0 \le f_{N,\varepsilon}(p) \le 2N \Big[ \varepsilon \log \frac{1}{\varepsilon} + (1-\varepsilon) \log \frac{1}{1-\varepsilon} \Big]$$

Somit erhalten wir [Referenz für geeigneten Stoppsatz heraussuchen...]

$$\begin{aligned} \left| \frac{1}{2N} f_{N,\varepsilon}(p) \right| &= \left| \frac{1}{2N} \mathbb{E}_p \Big[ f_{N,\varepsilon}(p_{T_{\varepsilon}^{(N)}}) + \sum_{j=0}^{T_{\varepsilon}^{(N)}-1} 2NR_N(p_j) \Big] \right| \\ &\leq \Big[ \varepsilon \log \frac{1}{\varepsilon} + (1-\varepsilon) \log \frac{1}{1-\varepsilon} \Big] + \frac{1}{\varepsilon^4 N^{3/2}} \mathbb{E}_p \Big[ T_{\varepsilon}^{(N)} \Big] \end{aligned}$$

 $\operatorname{Da} \mathbb{E}_p \left[ T_{\varepsilon}^{(N)} \right] \leq \mathbb{E}_p \left[ T_{\operatorname{fix}}^{(N)} \right] \leq 2NH(p) \operatorname{gemäß} (A.6) \operatorname{gilt}, \operatorname{folgt}$ 

$$\limsup_{N \to \infty} \max_{p \in \{0, \frac{1}{N}, \dots, 1\}} \left| \frac{1}{2N} f_{N, \varepsilon}(p) \right| \le H(\varepsilon)$$

mit  $H(\varepsilon) \to 0$  für  $\varepsilon \downarrow 0$ . Offenbar ist stets  $T_{\text{fix}}^{(N)} - T_{\varepsilon}^{(N)} \ge 0$ . Die starke Markov-Eigenschaft (und ein offensichliches Kopplungsargument, sowie Spiegelungssymmetrie um p = 1/2) zeigt, dass

$$0 \leq \mathbb{E}_p \Big[ T_{\mathrm{fix}}^{(N)} - T_{\varepsilon}^{(N)} \Big] \leq \sup_{0 < p' \leq \varepsilon} \mathbb{E}_{p'} \Big[ T_{\mathrm{fix}}^{(N)} \Big] \leq 2NH(\varepsilon).$$

Insgesamt ergibt sich

$$\limsup_{N \to \infty} \left| \frac{\mathbb{E}_p \left[ T_{\text{fix}}^{(N)} \right]}{2N} - H(p) \right| \le \limsup_{N \to \infty} \left| \frac{\mathbb{E}_p \left[ T_{\varepsilon}^{(N)} \right]}{2N} - H(p) \right| + \limsup_{N \to \infty} \left| \frac{\mathbb{E}_p \left[ T_{\text{fix}}^{(N)} - T_{\varepsilon}^{(N)} \right]}{2N} \right| \le H(\varepsilon) + H(\varepsilon)$$

und mit  $\varepsilon \downarrow 0$  folgt die Behauptung.

Wir tragen einige kleine Details aus dem Beweis von Satz 1.1 in folgendem Lemma zusammen:

**Lemma A.7.** (a) Für  $0 und <math>N \in \mathbb{N}$  definiert die Formel

$$\int_{(0,1)} f(x) \,\tilde{\mu}_{N,p}(dx) \coloneqq \frac{N}{p(1-p)} \int_0^1 \mathbb{E}_p \Big[ f \Big( p + u(p_1 - p) \Big) (p_1 - p)^2 \Big] 2(1-u) \, du \quad (A.10)$$

für beschränktes (und messbares)  $f : [0,1] \rightarrow \mathbb{R}$  ein Wahrscheinlichkeitsmaß  $\tilde{\mu}_{N,p}$  auf [0,1].

(b) Für  $Y_{N,p} \sim Bin(N,p)$  gilt

$$\mathbb{E}[(Y_{N,p} - Np)^3] = Np(1-p)(1-2p),$$
  
$$\mathbb{E}[(Y_{N,p} - Np)^4] = Np(1-p)(1-3p(1-p)) + 3N(N-1)p^2(1-p)^2 \le N^2$$

*Beweisksizze.* (a) Offensichtlich definiert die rechte Seite von (A.10) ein Maß. Um zu prüfen, dass es sich tatsächlich um ein Wahrscheinlichkeitsmaß handelt, setzen wir  $f(x) \equiv 1$  ein: wegen  $\mathbb{E}_p[(p_1 - p)^2] = p(1-p)/N$  und  $\int_0^1 2(1-u) du = 1$  ist  $\tilde{\mu}_{N,p}([0,1]) = 1$ .

(b) Man kann  $Y_{N,p} - Np = \sum_{i=1}^{N} (A_i - p)$  mit  $A_1, \ldots, A_n$  u.i.v. Bernoulli(*p*) darstellen und dann die Linearität des Erwartungswerts (und geeignete Gedanken zur Kombinatorik) ausnutzen. Alternativ könnte man die momentenerzeugende Funktion viermal ableiten.

Insbesondere liefert die Jensen-Ungleichung ( $\mathbb{R}_+ \ni x \mapsto x^{3/4}$  is konkav) dann

$$\mathbb{E}[|Y_{N,p} - Np|^3] = \mathbb{E}[(|Y_{N,p} - Np|^4)^{3/4}] \le (\mathbb{E}[|Y_{N,p} - Np|^4])^{3/4} \le (N^2)^{3/4} = N^{3/2}$$

### A.5 Ein Steilkurs über Martingale in diskreter Zeit

Dieses Kapitel ist eine Einladung, sich (in sehr knapper Form) mit der Theorie der (zeitdiskreten) Martingale zu beschäftigen. Eine wesentlich gründlichere Behandlung findet sich beispielsweise bei Klenke [Kle20], speziell Kapitel 8–11 (das auch Lesern ans Herz gelegt sei, die an den im Text erwähnten "Übungen" verzweifeln).

**Beispiel A.8.** Die symmetrische gewöhnliche Irrfahrt  $S_n = X_1 + \dots + X_n$ ,  $X_i$  u.i.v.,  $\mathbb{P}(X_i = \pm 1) = 1/2$ ( $S_0 \coloneqq 0$ ) wird uns hier als "Leib-und-Magen-Beispiel" dienen.

#### A.5.1 Filtrationen

Zur Erinnerung: Das "übliche" Grundobjekt der Wahrscheinlichkeitstheorie ist ein Wahrscheinlichkeitsraum  $(\Omega, \mathscr{A}, \mathbb{P})$ , bestehend aus Grundraum  $\Omega, \sigma$ -Algebra  $\mathscr{A}$  auf  $\Omega$  und einem Wahrscheinlichkeitsmaß  $\mathbb{P}$  auf  $\mathscr{A}$ . Eine *Filtration*  $(\mathscr{F}_n)_n$  ist eine aufsteigend geordnete Familie von Teil- $\sigma$ -Algebren, d.h.  $\mathscr{F}_n \subset \mathscr{F}_{n+1} \subset \mathscr{A}$  für  $n = 0, 1, 2, \ldots$  Eine naheliegende (und nützliche) Interpretation ist,  $\mathscr{F}_n$  als die Menge der bis zur Zeit n entschiedenen Ereignisse aufzufassen.

**Beispiel A.9.** Eine Folge von Zufallsvariablen  $(X_n)$  definiert via  $\mathscr{F}_n \coloneqq \sigma(X_0, X_1, \dots, X_n)$  eine Filtration (Übung: Überzeugen Sie sich davon).

### A.5.2 Bedingte Erwartung

 $\mathscr{G} \subset \mathscr{A}$  eine (Teil-) $\sigma$ -Algebra. Wenn  $\mathscr{G}$  endlich viele Atome  $A_1, \ldots, A_\ell$  hat, liegt es nahe, die "bedingte Erwartung von X, gegeben die Information aus  $\mathscr{G}$ " folgendermaßen zu definieren:

$$\mathbb{E}[X|\mathscr{G}](\omega) = \frac{1}{\mathbb{P}(A_i)} \mathbb{E}[X1_{A_i}] \quad \text{für } \omega \in A_i.$$
(A.II)

Man verallgemeinert (A.11) folgendermaßen: Eine reellwertige ZV Z heißt bedingte Erwartung von X gegeben  $\mathscr{G}$  (schreibe  $\mathbb{E}[X|\mathscr{G}]$ ), wenn gilt

- I. Z ist  $\mathscr{G}$ -messbar, d.h.  $\{Z \in B\} \in \mathscr{G}$  für jede messbare Teilmenge  $B \subset \mathbb{R}$ ,
- 2.  $\mathbb{E}[HZ] = \mathbb{E}[HX]$  für alle beschränkten  $\mathscr{G}$ -messbaren ZVn H.

(Übung: Überzeugen Sie sich, dass im Fall  $|\Omega| < \infty$  die Version aus (A.11) diese Definition erfüllt.)

**Bericht A.10.** Für integrierbares X existiert die bedingte Erwartung  $\mathbb{E}[X|\mathscr{G}]$  und ist bis auf f.s.-Gleichheit eindeutig bestimmt (Existenz beispielsweise via Projektion auf den Unterraum der (quadratintegrablen)  $\mathscr{G}$ -messbaren ZVn). Neben den "üblichen" Eigenschaften von Erwartungwerten (Linearität, Positivität) sind zwei wichtige Eigenschaften

- I.  $\mathbb{E}[\mathbb{E}[X|\mathcal{G}]|\mathcal{G}'] = \mathbb{E}[X|\mathcal{G}']$  wenn  $\mathcal{G}' \subset \mathcal{G}$  ("Turmeigenschaft"),
- 2.  $\mathbb{E}[YX|\mathcal{G}] = YE[X|\mathcal{G}]$  sofern  $Y \mathcal{G}$ -messbar ist (und  $\mathbb{E}[|XY|] < \infty$ ).

### A.5.3 Martingale

Eine Folge integrierbarer Zufallsvariablen  $(M_n)$  (so dass  $M_n \mathscr{F}_n$ -messbar ist für n = 0, 1, ...) heißt ein *Martingal* (bezüglich der Filtration  $(\mathscr{F}_n)$ ), wenn

$$\mathbb{E}[M_n | \mathscr{F}_{n-1}] = M_{n-1} \text{ (f.s.) für } n = 1, 2, \dots$$
(A.12)

gilt. Es gilt dann auch  $\mathbb{E}[M_n|\mathscr{F}_m] = M_m$  (f.s.) für  $m \leq n$  (Übung).

Die symmetrische Irrfahrt (Beispiel A.8) ist ein Martingal (Übung).

**Bericht A.11.** Wenn in (A.12) das  $_{=}$  durch  $_{>}$  ersetzt wird, spricht man von einem *Submartingal*, wenn es durch  $_{=}$  ersetzt wird, von einem *Supermartingal*.

#### Prävisible Prozesse, "Spielstrategien", Gewinnprozesse als Martingale

 $H_1, H_2, \ldots$  eine Folge (individuell) beschränkter Zufallsvariablen, so dass  $H_i \mathscr{F}_{i-1}$ -messbar ist für  $i = 1, 2, \ldots$  (man nennt dann  $(H_i)_{i\geq 1}$  auch *prävisibel*),  $(M_n)$  ein Martingal. Dann ist auch die Folge  $(Y_n)$ , definiert durch  $Y_0 \coloneqq 0$ ,

$$Y_n := \sum_{k=1}^n H_k (M_k - M_{k-1}), \quad n = 1, 2, \dots$$
 (A.13)

ein Martingal (Übung). Wenn man  $(M_n)$  als den Gewinnprozess eines Spielers, der in jeder Runde einen "Einheitseinsatz" in einem fairen Spiel wettet, interpretiert, so ergibt dies für  $(Y_n)$  folgende Interpretation: Dies ist der Gewinnprozess eines Spielers, der jeweils vor der *i*-ten Runde den  $H_i$ -fachen Einheitseinsatz setzt. Die Bedingung, dass  $H_i \mathscr{F}_{i-1}$ -messbar sein muss, beschreibt einen Spieler ohne hellseherische Fähigkeiten: die Höhe des Einsatzes muss vor der Kenntnis des Ausgangs der *i*-ten Runde festgelegt werden.

### A.5.4 Stoppzeiten

Eine Zufallsvariable  $\tau$  mit Werten in  $\{0, 1, ...\}$  mit der Eigenschaft

$$\{\tau \le n\} \in \mathscr{F}_n, \quad n = 0, 1, 2, \dots$$
 (A.14)

heißt eine *Stoppzeit* (strenggenommen:  $(\mathscr{F}_n)$ -Stoppzeit). (A.14) läßt sich folgendermaßen interpretieren: Man kann zu jedem Zeitpunkt n entscheiden, ob  $\tau$  "bereits eingetreten ist". Äquivalent kann man fordern, dass { $\tau = n$ }  $\in \mathscr{F}_n$  für alle n (Übung).

Für eine Stoppzeit  $\tau$  ist die  $\tau$ -Vergangenheit  $\mathscr{F}_{\tau}$  gegeben durch  $A \in \mathscr{F}_{\tau} : \iff A \cap \{\tau \leq n\} \in \mathscr{F}_n$  für jedes n ( $\mathscr{F}_{\tau}$  ist eine  $\sigma$ -Algebra, Übung).

Eine wichtige Klasse von Stoppzeiten erhält man mittels  $\tau_A := \min\{k \in \mathbb{Z}_+ : X_k \in A\}$ , wenn  $(X_n)$  eine  $(\mathscr{F}_n)$ -adaptierte Folge (sagen wir, reellwertiger) Zufallsvariablen und  $A \subset \mathbb{R}$  (Überzeugen Sie sich, dass  $\tau_A$  eine Stoppzeit ist). Sind  $\tau_1, \tau_2$  Stoppzeiten, so auch  $\tau_1 \wedge \tau_2$  und  $\tau_1 \vee \tau_2$  (Übung). Warum ist mit  $\tau$  stets auch  $\tau + 5$  eine Stoppzeit,  $\tau - 5$  aber im Allgemeinen nicht?

### A.5.5 Optionales Stoppen

 $(M_n)$  Martingal,  $\tau$  beschränkte Stoppzeit (d.h. es gibt eine Konstante T mit der Eigenschaft  $\mathbb{P}(\tau \leq T) = 1$ ). Dann gilt

**Satz A.12** (Optional sampling-Satz, Basisversion).  $\mathbb{E}[M_{\tau}] = \mathbb{E}[M_0]$ , all geneiner  $\mathbb{E}[M_{\tau}|\mathscr{F}_n] = M_{\tau \wedge n} f \ddot{u} r n = 0, 1, \dots$ 

Zum Beweis argumentieren Sie beispielsweise folgendermaßen: Überprüfen Sie zunächst, dass

$$M_{\tau} = \mathbb{E} \Big[ M_T \big| \mathscr{F}_{\tau} \Big] \quad \text{fast sicher} \tag{A.15}$$

gilt. Tatsächlich gilt für  $A \in \mathscr{F}_{\tau}$ 

$$\mathbb{E}[M_{\tau}\mathbf{1}_{A}] = \mathbb{E}\left[\sum_{k=0}^{T} M_{k}\mathbf{1}(\tau=k)\mathbf{1}_{A}\right] = \sum_{k=0}^{T} \mathbb{E}[M_{k}\mathbf{1}_{A\cap\{\tau=k\}}] = \sum_{k=0}^{T} \mathbb{E}\left[\mathbb{E}[M_{T}|\mathscr{F}_{k}]\mathbf{1}_{A\cap\{\tau=k\}}\right]$$
$$= \sum_{k=0}^{T} \mathbb{E}\left[\mathbb{E}[M_{T}\mathbf{1}_{A\cap\{\tau=k\}}|\mathscr{F}_{k}]\right] = \sum_{k=0}^{T} \mathbb{E}[M_{T}\mathbf{1}_{A\cap\{\tau=k\}}] = \mathbb{E}[M_{T}\mathbf{1}_{A}],$$

wobei an geeigneter Stelle (wo?) die Martingaleigenschaft  $M_k = \mathbb{E}[M_T|\mathscr{F}_k]$ , die Tatsache  $A \cap \{\tau = k\} \in \mathscr{F}_k$  und  $\mathbb{E}[\mathbb{E}[\cdot|\mathscr{F}_k]] = \mathbb{E}[\cdot]$  ausgenutzt werden. Aus (A.15) ergibt sich sofort die erste Behauptung (warum?).

Für die zweite Behauptung benutzen Sie die Turmeigenschaft der bedingten Erwartung beispielsweise folgendermaßen:  $\tau \wedge n$  ist ebenfalls eine Stoppzeit, die offenbar  $\tau \wedge n \leq \tau (\leq T)$  erfüllt. Überlegen Sie sich, dass dies  $\mathscr{F}_{\tau \wedge n} \subset \mathscr{F}_{\tau}$  impliziert (ist das anschaulich einsichtig?). Demnach gilt mit (A.15)

$$M_{\tau \wedge n} = \mathbb{E}\big[M_T \big| \mathscr{F}_{\tau \wedge n}\big] = \mathbb{E}\big[\mathbb{E}\big[M_T \big| \mathscr{F}_{\tau}\big]\big| \mathscr{F}_{\tau \wedge n}\big] = \mathbb{E}\big[M_\tau \big| \mathscr{F}_{\tau \wedge n}\big].$$

**Bericht A.13.** Man kann in Satz A.12 die Bedingung, dass  $\tau$  beschränkt sein muss, fallen lassen. Technisch ist dann die entscheidende Bedingung, dass die Familie  $(M_n)$  gleichgradig integrierbar sein muss (Siehe [Kle20, Abschn. 10.3]). Ganz ohne Bedingungen kann Satz A.12 aber nicht richtig sein, wie die gewöhnliche Irrfahrt (Beispiel A.8) mit  $\tau_{\{1\}} := \min\{n : S_n = 1\}$  zeigt: Wegen der Rekurrenz von  $(S_n)$  ist  $\tau_{\{1\}} < \infty$  f.s., also  $S_{\tau_{\{1\}}} = 1$ , somit  $\mathbb{E}[S_{\tau_{\{1\}}}] = 1 \neq \mathbb{E}[S_0] = 0$ . (Für die Glücksspielinterpretation bedeutet dies: Man kann – im Prinzip – aus einem fairen Spiel sicheren Gewinn ziehen, wenn man ggfs. beliebig lange spielen und dabei beliebig hohe Schulden ansammeln darf.)

**Bemerkung A.14.** Aus Satz A.12 folgt, dass das *gestoppte* Martingal  $(M_{\tau \wedge n})_n$  ebenfalls ein Martingal ist, wenn  $\tau$  eine (beschränkte) Stoppzeit und  $(M_n)$  ein Martingal ist.

#### A.5.6 Konvergenz

Unter ("leichten") Bedingungen konvergiert ein Martingal  $(M_n)$  fast sicher. Die auf Joseph Doob zurückgehende Beweisidee ist folgende: Wäre dies nicht der Fall, so gäbe es a < b, so dass  $(M_n)$ unendlich oft zwischen (unterhalb) a und (oberhalb) b oszilliert. Dann könnte man mit folgender Strategie beliebig großen Gewinn erzielen:

- Steige ein, sobald  $M_n$  unter a fällt,
- halte, bis  $M_n$  über b steigt.
- Erziele mindestens Gewinn b a > 0 aus jeder solchen "Aufkreuzung".

Das widerspricht allerdings den Beobachtungen aus Abschnitt (A.5.3).

Wir wollen diese Idee nun präzisieren. Sei  $(M_n)$  ein nach unten beschränktes Martingal, o.E.  $M_n \ge 0$  für alle *n*. (Warum ist die Annahme  $\ge 0$  keine zusätzliche Einschränkung?)

Seien  $0 \le a < b < \infty$ . Setzen Sie  $\sigma_0 \coloneqq 0$ ,

$$\tau_k := \inf\{n > \sigma_{k-1} : M_n \le a\}, \quad k = 1, 2, \dots, \\ \sigma_k := \inf\{n > \tau_k : M_n \ge b\}, \quad k = 1, 2, \dots$$

(Mit Verabredung  $\tau_k = \infty$  bzw.  $\sigma_k = \infty$ , wenn es kein passendes *n* mehr gibt.)

Überzeugen Sie sich, dass die  $\tau_k$  und  $\sigma_k$  Stoppzeiten sind. Betrachten Sie beispielsweise eine Skizze, um sich zu vergewissern, dass  $(M_n)$ 

im Zeitintervall { $\tau_k, \tau_{k+1}, \dots, \sigma_k$ } die k-te Aufkreuzung von (unterhalb) a nach (oberhalb) b ausführt (A.16)

(sofern  $\tau_k, \sigma_k < \infty$ ). Sei

$$U_n^{(a,b)} \coloneqq \max\{k : \sigma_k \le n\}$$

die Anzahl abgeschlossener solcher Aufkreuzungen bis zum Zeitpunkt *n*.

Sei  $I_0 \coloneqq 0$ , für  $n \ge 1$ 

$$I_n := \sum_{i=0}^{n-1} \mathbf{1} (\exists k : \tau_k \le i < \sigma_k) (M_{i+1} - M_i),$$

d.h. nur die Inkremente von  $(M_n)$  innerhalb der Aufkreuzungsintervalle zählen für  $(I_n)$ . Verifizieren Sie, dass

$$\mathbb{E}\big[I_n\big|\mathscr{F}_{n-1}\big] = I_{n-1}$$

gilt, d.h.  $(I_n)$  ist (ebenfalls) ein Martingal.

Warum gilt

$$I_n \ge (b-a)U_n^{(a,b)} + \left(M_n - M_{\tau_{U_n^{(a,b)}+1} \land n}\right) \ge (b-a)U_n^{(a,b)} + (0-a)?$$
(A.17)

(Hinweis: Für jedes k ist  $\sum_{i=\tau_k}^{\sigma_k-1} (M_{i+1} - M_i) = M_{\sigma_l} - M_{\tau_k} \ge (b-a).$ )

Lemma A.15 (Aufkreuzungsungleichung). Es gilt für jedes n

$$\mathbb{E}\left[U_n^{(a,b)}\right] \le \frac{a}{b-a}.$$
(A.18)

Offenbar  $U_n^{(a,b)} \leq U_{n+1}^{(a,b)}$  für alle n, d.h. die Folge von Zufallsvariablen  $(U_n^{(a,b)} : n \in \mathbb{N})$  konvergiert monoton gegen ein  $U_{\infty}^{(a,b)}$ , also auch

$$\mathbb{E}\left[U_{\infty}^{(a,b)}\right] = \lim_{n \to \infty} \mathbb{E}\left[U_{n}^{(a,b)}\right] \le \frac{a}{b-a} < \infty$$

(benutzen Sie den Satz von der monotonen Konvergenz für das Gleichheitszeichen und dann (A.18) für die Abschätzung), insbesondere  $\mathbb{P}(U_{\infty}^{(a,b)} < \infty) = 1$ .

Betrachten Sie nun Ereignisse (mit  $0 \le a < b, a, b \in \mathbb{Q}$ , sagen wir)

$$C^{(a,b)} := \left\{ \liminf_{n \to \infty} X_n < a \right\} \cap \left\{ \limsup_{n \to \infty} X_n > b \right\}.$$

Argumentieren Sie, dass  $C^{(a,b)} \subset \{U^{(a,b)}_{\infty}\}$ , folglich  $P(C^{(a,b)}) = 0$  nach obigem, und daher auch

$$\mathbb{P}\left(\bigcup_{\substack{0\le a< b\\a,b\in\mathbb{Q}}} C^{(a,b)}\right) = 0 \tag{A.19}$$

gilt. Warum haben Sie damit folgende Version des Martingalkonvergenzsatzes bewiesen?

Satz A.16. Ein nach unten beschränktes Martingal konvergiert mit Wahrscheinlichkeit 1.

**Bericht A.17.** Die Konvergenz  $M_n \to M_\infty$  f.s. muss i.A. nicht die Konvergenz der Erwartungwerte implizieren: Betrachten Sie beispielsweise die Irrfahrt aus Beispiel A.8, die beim Auftreffen auf -1 gestoppt wird. Für gleichgradig integrierbare Martingale gilt allerdings auch  $\mathbb{E}[M_n] \to \mathbb{E}[M_\infty]$ .

### A.5.7 Doobsche Ungleichung

Im Allgemeinen ist es sehr schwierig, aus der Verteilung eines stochastischen Prozesses zu festen Zeiten Informationen über das Pfadverhalten wie beispielsweise das laufende Maximum abzuleiten. Im Fall von Martingalen sind die Verhältnisse übersichtlicher:
**Satz A.18** (Doobs  $L^2$ -Ungleichung). Sei  $(M_n)$  Martingal mit  $M_0 \ge 0$  und  $\mathbb{E}[M_n^2] < \infty$  für alle n,  $M_n^* := \max_{0 \le k \le n} M_k$ . Dann gilt

$$\mathbb{E}\big[(M_n^*)^2\big] \le 4\mathbb{E}\big[M_n^2\big].$$

Für festes  $\lambda > 0$  gilt

$$\lambda \mathbb{P}(M_n^* \ge \lambda) \le \mathbb{E}\left[M_n \mathbf{1}(M_n^* \ge \lambda)\right] \left( \le \mathbb{E}\left[|M_n| \mathbf{1}(M_n^* \ge \lambda)\right] \right).$$
(A.20)

Argumentieren Sie beispielsweise folgendermaßen:  $\tau := \inf\{k : M_k \ge \lambda\} \land n$  ist eine (durch n) beschränkte Stoppzeit, also

$$\mathbb{E}[M_n] = \mathbb{E}[M_{\tau}] = \mathbb{E}[M_{\tau}\mathbf{1}(M_n^* \ge \lambda)] + \mathbb{E}[M_{\tau}\mathbf{1}(M_n^* < \lambda)] = \mathbb{E}[M_{\tau}\mathbf{1}(M_n^* \ge \lambda)] + \mathbb{E}[M_n\mathbf{1}(M_n^* < \lambda)]$$
  
$$\ge \lambda \mathbb{P}(M_n^* \ge \lambda) + \mathbb{E}[M_n\mathbf{1}(M_n^* < \lambda)].$$

Nun substrahiere  $\mathbb{E}[M_n \mathbf{1}(M_n^* < \lambda)]$  auf beiden Seiten.

Stets gilt

$$(M_n^*)^2 = \int_0^{M_n^*} 2\lambda \, d\lambda,$$

also (wegen  $(M_n^*)^2 \le M_0^2 + M_1^2 + \dots + M_n^2$  ist der Erwartungswert endlich)

$$\mathbb{E}[(M_n^*)^2] = \mathbb{E}\left[\int_0^{M_n^*} 2\lambda \, d\lambda\right] = \mathbb{E}\left[\int_0^\infty 2\lambda \mathbf{1}(M_n^* \ge \lambda) \, d\lambda\right] = 2\int_0^\infty \lambda \mathbb{P}(M_n^* \ge \lambda) \, d\lambda$$
$$\leq 2\int_0^\infty \mathbb{E}[|M_n|\mathbf{1}(M_n^* \ge \lambda)] \, d\lambda = 2\mathbb{E}\left[|M_n|\int_0^{M_n^*} d\lambda\right] = 2\mathbb{E}[|M_n|M_n^*]$$

Folgern Sie mit der Cauchy-Schwarz-Ungleichung:

$$\mathbb{E}\left[(M_n^*)^2\right] \le 2\sqrt{E\left[|M_n|^2\right]}\sqrt{\mathbb{E}\left[(M_n^*)^2\right]}.$$

**Bericht A.19.** Die Ungleichung gilt wörtlich auch für Submartingale. Die Annahme  $M_0 \ge 0$  ist eigentlich nicht notwendig (vereinfacht hier nur die Argumentation ein wenig). Es gilt eine analoge Aussage für jedes p > 1 statt p = 2 (Doobs  $L^p$ -Ungleichung).

## A.5.8 Die symmetrische gewöhnliche Irrfahrt auf einem Intervall

Seien  $a, b \in \mathbb{Z}$ , a < x < b,  $Z^{(x,a,b)}$  die symmetrische gewöhnliche Irrfahrt startend in  $Z_0^{(x,a,b)} = x$ , gestoppt, sobald  $Z_0^{(x,a,b)} \in \{a, b\}$ . Prüfen Sie:  $((b - Z_n^{(x,a,b)})/(b-a))_n$  und  $((b - Z_n^{(x,a,b)})(Z_n^{(x,a,b)} - a) - n)_n$  sind Martingale. Können Sie diese Information benutzen, um die Wahrscheinlichkeit, dass der obere Rand getroffen wird, sowie die erwartete Zeit bis zum Treffen des Rands zu berechnen?