

Patienten-IDentifikatoren in medizinischen Forschungsnetzen: Evaluation des Matchalgorithmus

Jutta Glock, Klaus Pommerening; Institut für Medizinische Biometrie, Epidemiologie und Informatik der Johannes-Gutenberg-Universität Mainz

Begriffe

PID: Studienübergreifender, eindeutiger Patienten-Identifikator

PID-Generator: Verarbeitet PID-Anfragen, vergleicht bestimmte Merkmale des eingegebenen Datensatzes mit der Patientenliste, liefert bei Match den PID zurück oder generiert einen neuen PID und speichert den Fall in Patientenliste

Homonymfehlerrate: Wie oft wird verschiedenen Patienten derselbe PID zugeteilt, d.h. fälschlicherweise gematched? Die Homonymfehlerrate hängt maßgeblich von den gewählten Datenfeldern ab („echte“ Homonyme vermeiden).

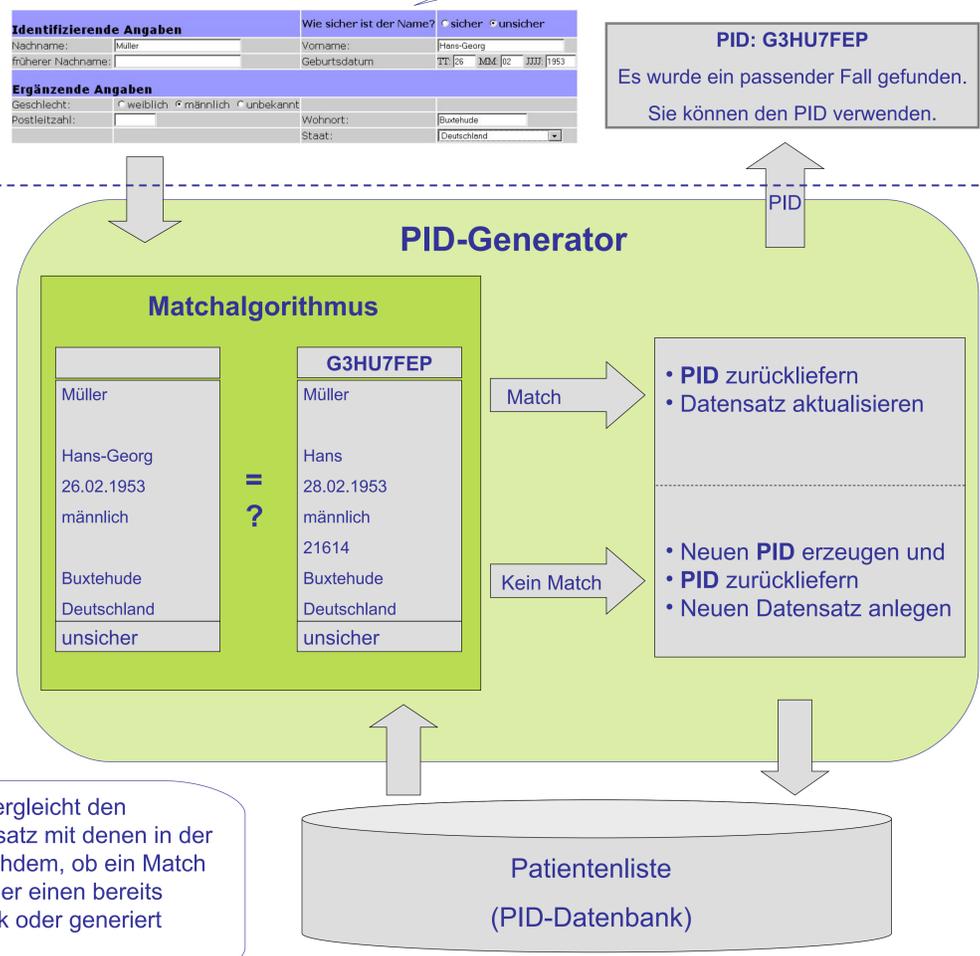
Synonymfehlerrate: Wie oft werden einem Patienten zwei oder mehrere PIDs zugeordnet, d. h. fälschlicherweise nicht gematched? Die Synonymfehlerrate hängt vor allem von der Datenqualität und organisatorischen Gegebenheiten (Häufigkeit von Mehrfacheingaben) ab und ist daher je nach Anwendung und Datenquelle unterschiedlich.

Ablauf einer PID-Anfrage

Der Anwender gibt die Patientendaten über eine Weboberfläche ein und erhält als Ergebnis einen PID oder eine Fehlermeldung.

Anwender
(Studienzentrale)

PID-Server



KPOH-Konfiguration

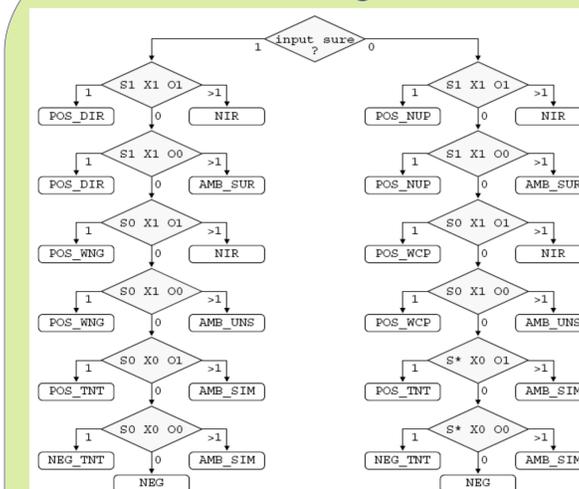
Das Matchverfahren ist insbesondere abhängig von der Wahl der Datenfelder, die zum Abgleich verwendet werden, und der Gestaltung des Entscheidungsbaums. Beides ist im PID-Generator frei konfigurierbar. Hier dargestellt sind die getesteten Spezifikationen des Kompetenznetzes POH (KPOH).

Datenfelder

Feldname	Bedeutung	Wertebereich	Pflichtfeld	Relevanz [1]
lname	Nachname	string	Ja	+
aname	Alternativer Nachname (z. B. Geburtsname)	string	Nein	+
fname	Vorname(n)	string	Ja	+
bd	Geburtstag	0-31	Ja	+
bm	Geburtsmonat	0-12	Ja	+
by	Geburtsjahr	1000-9999	Ja	+
plz	Postleitzahl	string	Nein	-
loc	Wohnort	string	Nein	*
state	Land	string	Nein	*
sex	Geschlecht	[f m n]	Nein	*

[1] + wird immer beim Matchen verwendet
* wird optional beim Matchen verwendet
- keine Matchrelevanz

Entscheidungsbaum



Legende

Rauten stellen Datenbankabfragen, abgerundete Rechtecke Resultate dar. Nach jeder Abfrage werden entweder 0 Matche, 1 Match oder mehr als 1 Match gefunden. Abhängig davon wird entweder eine neue Abfrage durchgeführt oder die Suche durch ein Resultat abgeschlossen. Nur bei Erreichen des Resultats NEG wird ein neuer PID generiert, in allen anderen Fällen wird entweder ein bereits vorhandener PID oder eine Fehlermeldung zurückgeliefert. Die Kürzel in den Rauten stellen Filter dar, die auf die Datenbank angewendet werden.

S Sureness [0 = unsicher | 1 = sicher | * = beides]
X Exactitude [0 = phonetische | 1 = exakte Übereinstimmung]
O Optionality [0 = ohne optionale Daten | mit optionalen Daten]

Tests / Ergebnisse

Funktionstest mit fiktiven Daten

- Alle Pfade des Matchbaums werden durchlaufen
- Alle Tests liefern die erwarteten Ergebnisse

Realtests mit 14.915 Datensätzen mit und ohne Verschlüsselung der Datenfelder

- 14.913 neue PIDs
- und 2 Matche erzeugt
- Verschlüsselung hat erwartungsgemäß keinen Einfluss
- 0 Homonymfehler (da beide Matche korrekt)
- 6 Synonymfehler (beruhen meist auf Fehler im Geburtsdatum)
- Realbetrieb Verpiddung des Deutschen Kinderkrebsregisters (DKKR)**
 - 44.248 Datensätze des DKKR vs. 2.579 Datensätze in Patientenliste
 - 1.569 Matche (davon 1 Duplikat)
 - 53 Mal ist Zuordnung zu unsicher (NEG_TNT)
 - 0 Homonyme innerhalb der 44.248 Datensätze (Homonyme mit „Altdaten“?)
 - < 50 Synonyme:
 - Fehlerhafte PID-Übermittlung
 - oder echte Synonyme aufgrund fehlerhaften Daten

Fiktive Testdaten (Beispiele)

Sicher?	Name	Altname	Vorname	Geburtstag	Geschlecht	Resultat	PID
Ja	Albrecht		Anton	19.01.2001	m	NEG	G1QP56LV
Ja	Alt	Albrecht	Anton-Armdt	19.01.2001	m	POS_DIR	G1QP56LV
Nein	Meier		Moritz	12.11.2002	m	NEG	4VV50DPW
Nein	Maier		Mohritz	12.11.2002	m	POS_TNT	4VV50DPW
Ja	Veit		Verena	12.08.2003	w	NEG	MAPVTF8K
Ja	Veit		Viviane	12.08.2003	w	NEG	1HUZZWLY
Nein	Veit		Verena-Viviane	12.08.2003	w	NIR	???

Stand der GPOH-Patientenliste (20.07.2005)

PIDs gesamt	45.693
Mehrfachanfragen	
4-6-fach	28
3-fach	214
2-fach	3.299
Resultate (seit 03.06.05)	
POS_DIR	9
POS_NUP	1.247
POS_WNG	28
POS_WCP	2.204
NEG_TNT	123
NEG	43.117