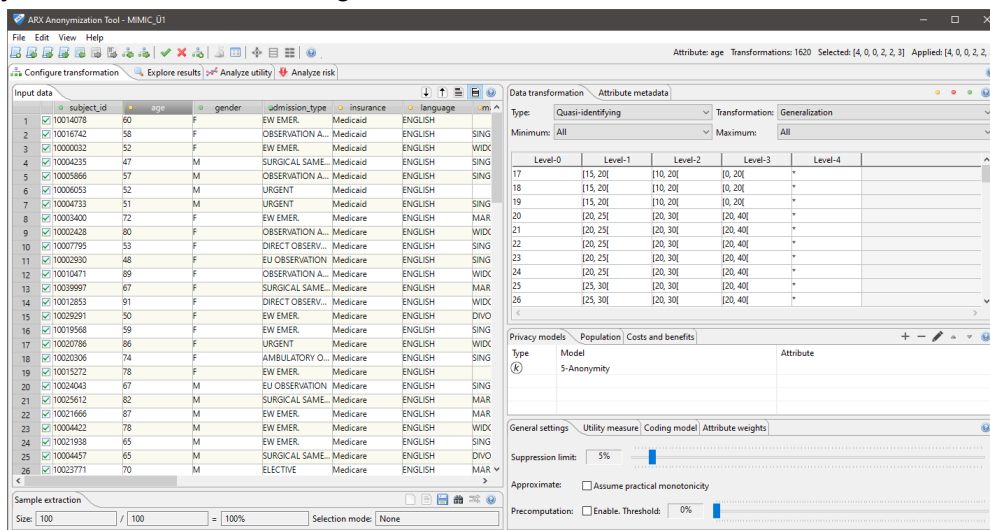


ARX MIMIC Übung

Basiert auf der ARX Übung von Dr. Fabian Prasser. [Musterlösung](#)

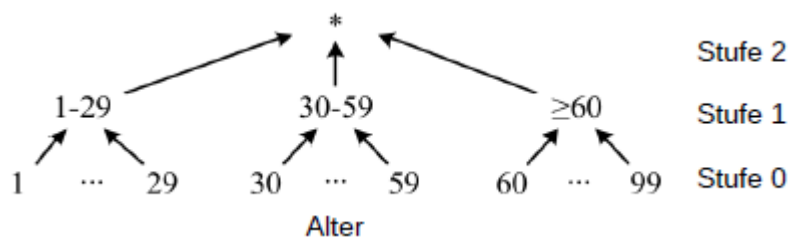
Aufgabe 1)

1. Öffnen Sie ARX und erstellen Sie ein neues Projekt. Importieren Sie den Datensatz "MIMIC IV Demo full Ü3.csv".
2. Wählen Sie die Spalten "age", "insurance", "language", "marital_status", "race" und "weight" nacheinander aus und wählen Sie jeweils den Attributtyp "Quasi-identifying". Importieren Sie nun unter „File“ → „Import hierarchy“ die entsprechende Hierarchie für jedes Attribut das Sie ausgewählt haben.



Hintergrund

- Generalisierungshierarchien werden zur Transformation von Werten eingesetzt



- In ARX werden Generalisierungshierarchien meist als Tabellen dargestellt
- Jede Zeile einer solchen Tabelle enthält die Abbildungsregel für einen Attributwert
- Für die hier gezeigte Hierarchie ergibt sich bspw. folgende Darstellung:

Stufe 0	Stufe 1	Stufe 2
1	1-29	*
...
29	1-29	*
30	30-59	*
...

- **Hinweis:** Für die sichere Durchführung von Anonymisierungsverfahren mit ARX ist es zwingend notwendig, dass jedes in einem Datensatz abgebildete Individuum durch genau einen Eintrag repräsentiert ist!
3. Schützen Sie nun die Einträge des Datensatzes vor Re-Identifikation. Spezifizieren Sie ein maximales Re-Identifikationsrisiko von 20% für jeden Eintrag im Datensatz. Verwenden Sie Generalisierung von Attributen und Unterdrückung von Einträgen zum Schutz des Datensatzes. Stellen Sie sicher, dass nicht mehr als 5% der Einträge entfernt werden.
- Wählen Sie unter “General Settings” ein Unterdrückungs Limit von 5 %. Das ist die Begrenzung des Anteils von Einträgen die entfernt werden können.

General settings | Utility measure | Coding model | Attribute weights

Suppression limit: 5%

Approximate: ☐ Assume practical monotonicity

Precomputation: ☐ Enable. Threshold: 0%

Spezifizieren Sie nun die Risikogrenzwerte / Datenschutzmodelle. Fügen Sie unter der Liste aller aktiver Schutzmodelle ein neues Schutzmodell hinzu. Wählen Sie k-Anonymity mit $k=5$ (Risikogrenzwert $1/5=20\%$).

Add a new privacy model

Please select a privacy model which will be applied to the data set

Type	Model	Attribute
(ϵ, δ)	(ϵ, δ) -Differential privacy	
(k)	k-Anonymity	
(k)	k-Map	
(δ)	δ -Presence	
(ϵ)	Profitability	
(r)	Average-reidentification-risk	
(r)	Population-uniqueness	
(r)	Sample-uniqueness	

Configuration

K: 5

Note: you can also enter values by double-clicking the control knobs

OK Cancel

Privacy models | Population | Costs and benefits

Type	Model	Attribute
(k)	5-Anonymity	

4. Anonymisieren Sie nun den Datensatz über die Menüleiste (✓). Nach erfolgreicher Anonymisierung werden in der rechten oberen Ecke von ARX Informationen über den Prozess angezeigt

Hintergrund

- Eine Transformation der Eingabedaten verringert die Eindeutigkeit von Wertekombinationen
- Dadurch entstehen Gruppen von Einträgen, auch Klassen genannt

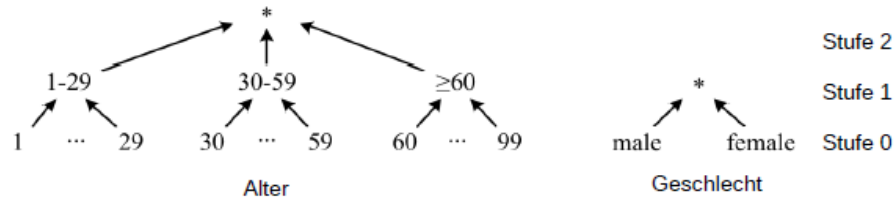
- Die Größe der Gruppe in der sich der zu einem Individuum gehörende Datensatz befindet bestimmt die Wahrscheinlichkeit mit der der Eintrag dem Individuum zugeordnet werden kann
- Befindet sich ein Eintrag beispielsweise in einer Gruppe der Größe k , so kann dieser nur mit einer Wahrscheinlichkeit von $1/k$ dem Individuum richtig zugeordnet werden
- Man nennt diese Wahrscheinlichkeit das Re-Identifikationsrisiko
- Die Attribute die zur Gruppenbildung herangezogen werden nennt man Quasi-Identifikatoren
- Das Datenschutzmodell k-Anonymity fordert für jeden Eintrag eine Klassengröße von mindestens k
- Das bedeutet, dass k-Anonymity das Re-Identifikationsrisiko jedes Eintrags auf maximal $1/k$ beschränkt
- Die Transformation von Eingabedaten erfolgt meist durch das Generalisieren von Attributwerten und das Entfernen (auch "Unterdrücken" genannt) von Einträgen
- Üblicherweise wird eine Obergrenze für die Anzahl der Einträge die entfernt werden dürfen festgelegt („Suppression limit“)

Aufgabe 2)

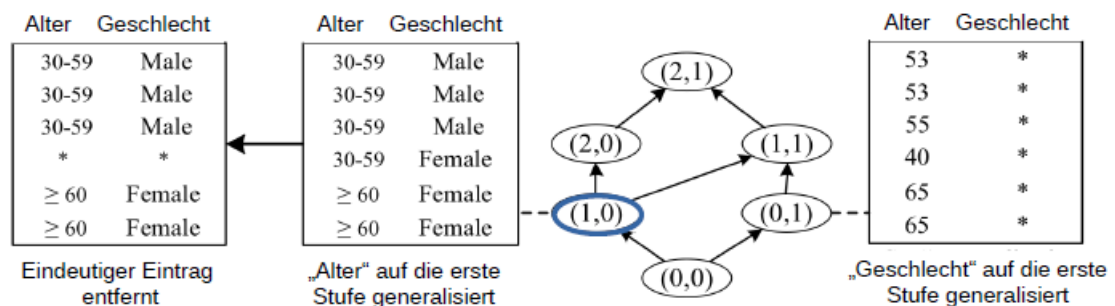
Untersuchen Sie den Lösungsraum. Finden Sie alternative Lösungen mit einem nahezu optimalen Informationsverlust. Untersuchen Sie Unterschiede zur vorhergehenden Lösung.

Hintergrund

- Generalisierungshierarchien



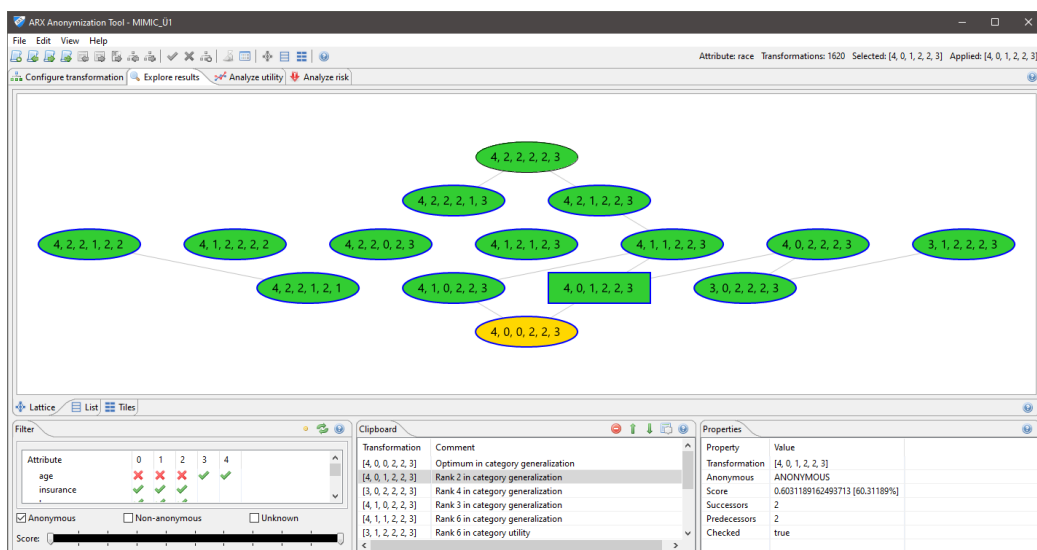
- Lösungsraum



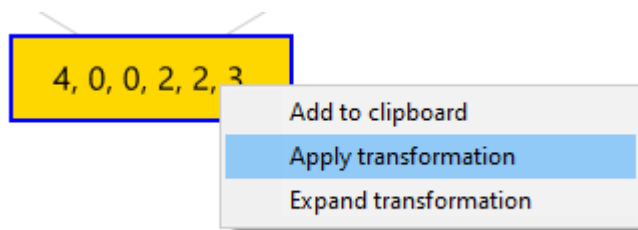
- Jede Transformation (im Bild markiert) im Lösungsraum definiert eine eindeutige Kombination von Generalisierungsstufen, eine für jeden Quasi-Identifikator
- Zusätzlich werden oftmals Einträge automatisch entfernt („suppression“)

1. Verwenden Sie die Explorationsperspektive von ARX:

- Grün: anonyme Transformation
- Gelb: optimale Lösung
- Rechteck: ausgewählte Transformation



Hinweis: Mit einem Rechtsklick auf eine Transformation können Sie „Apply transformation“ wählen, um eine ausgewählte Transformation auf den Datensatz anzuwenden



2. Wechseln Sie nun in die Qualitätsanalyseperspektive von ARX:

- Links: Eingabe Daten
- Rechts: Ausgabe der angewandten Transformation

ARX Anonymization Tool - MIMIC_U1

File Edit View Help

Attribute: age Transformations: 1620 Selected: [4, 0, 0, 2, 3] Applied: [4, 0, 0, 2, 3]

Configure transformation Explore results Analyze utility Analyze risk

Input data Classification performance Quality models

	subject_id	age	gender	admission_type	insurance	language	rm
1	10014078	60	F	EW EMER.	Medicaid	ENGLISH	SING
2	10016742	58	F	OBSERVATION A...	Medicaid	ENGLISH	SING
3	10000032	52	F	EW EMER.	Medicaid	ENGLISH	WIDC
4	10004235	47	M	SURGICAL SAME...	Medicaid	ENGLISH	SING
5	10005866	57	M	OBSERVATION A...	Medicaid	ENGLISH	SING
6	10006053	52	M	URGENT	Medicaid	ENGLISH	SING
7	10004733	51	M	URGENT	Medicaid	ENGLISH	SING
8	10003400	72	F	EW EMER.	Medicare	ENGLISH	MAR
9	10002428	80	F	OBSERVATION A...	Medicare	ENGLISH	WIDC
10	10007795	53	F	DIRECT OBSERV...	Medicare	ENGLISH	SING
11	10002930	48	F	EU OBSERVATION	Medicare	ENGLISH	SING
12	10010471	89	F	OBSERVATION A...	Medicare	ENGLISH	WIDC
13	10039997	67	F	SURGICAL SAME...	Medicare	ENGLISH	MAR
14	10012853	91	F	DIRECT OBSERV...	Medicare	ENGLISH	WIDC
15	10029291	50	F	EW EMER.	Medicare	ENGLISH	DIVO
16	10019568	59	F	EW EMER.	Medicare	ENGLISH	SING
17	10020786	86	F	URGENT	Medicare	ENGLISH	WIDC
18	10020306	74	F	AMBULATORY O...	Medicare	ENGLISH	SING
19	10015272	78	F	EW EMER.	Medicare	ENGLISH	SING
20	10024043	67	M	EU OBSERVATION	Medicare	ENGLISH	SING
21	10025612	82	M	SURGICAL SAME...	Medicare	ENGLISH	MAR

Output data Classification performance Quality models

	subject_id	age	gender	admission_type	insurance	language	rm
1	10014078	*	F	EW EMER.	Medicaid	ENGLISH	*
2	10016742	*	F	OBSERVATION A...	Medicaid	ENGLISH	*
3	10000032	*	F	EW EMER.	Medicaid	ENGLISH	*
4	10004235	*	M	SURGICAL SAME...	Medicaid	ENGLISH	*
5	10005866	*	M	OBSERVATION A...	Medicaid	ENGLISH	*
6	10006053	*	M	URGENT	Medicaid	ENGLISH	*
7	10004733	*	M	URGENT	Medicaid	ENGLISH	*
8	10003400	*	F	EW EMER.	Medicare	ENGLISH	*
9	10002428	*	F	OBSERVATION A...	Medicare	ENGLISH	*
10	10007795	*	F	DIRECT OBSERV...	Medicare	ENGLISH	*
11	10002930	*	F	EU OBSERVATION	Medicare	ENGLISH	*
12	10010471	*	F	OBSERVATION A...	Medicare	ENGLISH	*
13	10039997	*	F	SURGICAL SAME...	Medicare	ENGLISH	*
14	10012853	*	F	DIRECT OBSERV...	Medicare	ENGLISH	*
15	10029291	*	F	EW EMER.	Medicare	ENGLISH	*
16	10019568	*	F	EW EMER.	Medicare	ENGLISH	*
17	10020786	*	F	URGENT	Medicare	ENGLISH	*
18	10020306	*	F	AMBULATORY O...	Medicare	ENGLISH	*
19	10015272	*	F	EW EMER.	Medicare	ENGLISH	*
20	10024043	*	M	EU OBSERVATION	Medicare	ENGLISH	*
21	10025612	*	M	SURGICAL SAME...	Medicare	ENGLISH	*

Summary statistics Distribution Contingency Class sizes Properties Classification models

Measure	Value (incl. suppressed)	Value (excl. suppressed)
Average class size	1.0101 (1.0101%)	1.0101 (1.0101%)
Maximal class size	2 (2%)	2 (2%)
Minimal class size	1 (1%)	1 (1%)
Suppressed records	0 (0%)	0
Number of classes	99	99
Number of records	100	100

Summary statistics Distribution Contingency Class sizes Properties Classification models

Measure	Value (incl. suppressed)	Value (excl. suppressed)
Average class size	24.25 (24.25%)	24.25 (25%)
Maximal class size	49 (49%)	49 (50.51546%)
Minimal class size	5 (5%)	5 (5.15464%)
Suppressed records	3 (3%)	0
Number of classes	4	4
Number of records	100	97

3. Schauen Sie sich den Reiter „Class sizes“ an und beantworten Sie die folgenden Fragen:

Füllen Sie aus:

• Optimale Transformation

• [_, _, _, _, _, _, _, _]

• Eigenschaften

• Vollständig entfernte(s) Attribut(e): ____ (_____)

• Unterdrückte Einträge: ____ (____, ____%)

• Informationsverlust: ____,%

• Durchschn. Klassengröße / Average class size: __, __

• Ihre Auswahl

• [_, _, _, _, _, _, _, _]

• Eigenschaften

• Vollständig entfernte(s) Attribut(e): ____ (_____)

• Unterdrückte Einträge: ____ (____, ____%)

• Informationsverlust: ____,%

• Durchschn. Klassengröße / Average class size: __, __

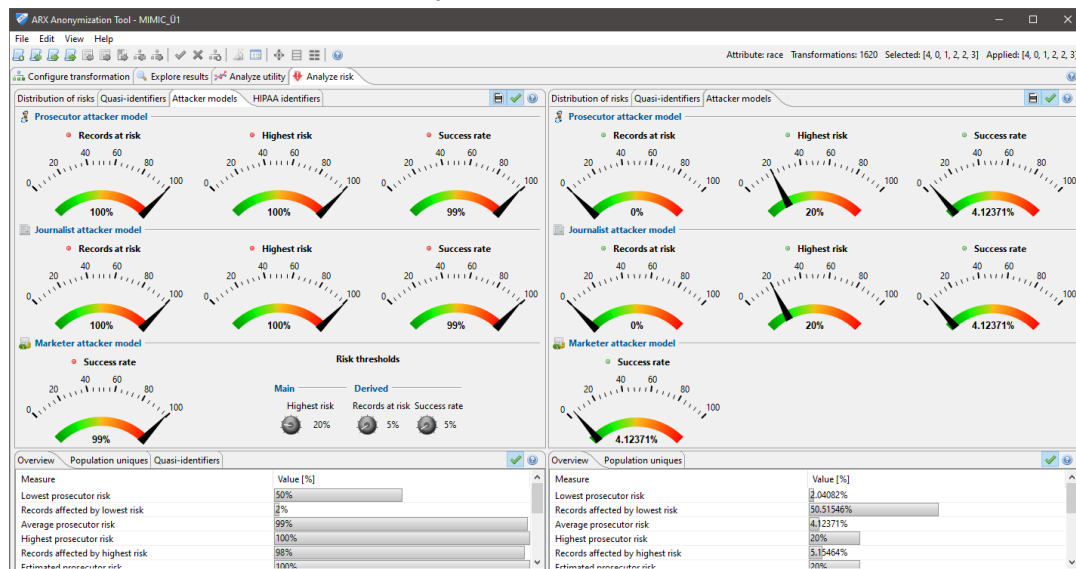
Aufgabe 3)

Analysieren Sie Re-Identifikationsrisiken von Einträgen in einem 5-anonymen Datensatz. Hier entspricht 20% 5-anonymity. Untersuchen Sie wie sich Re-Identifikationsrisiken durch die Anonymisierung verändert haben.

Hintergrund:

- Die Größe der Gruppe in der sich ein Eintrag befindet bestimmt die Wahrscheinlichkeit mit der er dem zugehörigen Individuum richtig zugeordnet werden kann, d.h. das Re-Identifikationsrisiko
- Selbst bei einem k-anonymen Datensatz muss davon ausgegangen werden, dass durch zufälliges Abbilden von Individuen in passende Gruppen ein nicht geringer Anteil der Einträge richtig zugeordnet werden kann
- Diese Betrachtungen setzen voraus, dass der Angreifer bereits weiß, dass Daten über die Individuen die er versucht abzubilden im Datensatz vorhanden sind
- Man nennt dieses Worst-Case Szenario das Prosecutor Modell

1. Wechseln Sie in den Reiter „Analyze Risk“ und wählen Sie „Attacker models“:



Füllen Sie aus:

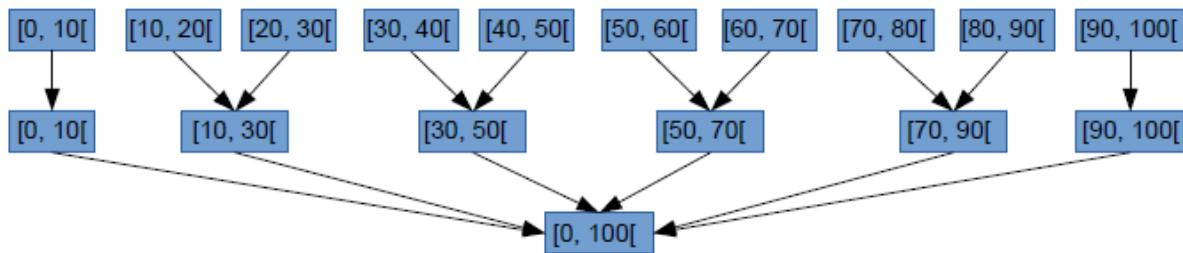
- Prosecutor Re-Identifikationsrisiko für den Eingabedatensatz
 - __, __% der Einträge sind einem Risiko von mehr als 20% ausgesetzt
 - Die zu erwartende relative Anzahl an korrekt re-identifizierten Einträgen ist __, __%
 - Das niedrigste Risiko ist __%, das höchste Risiko ist __%
 - Das höchste Risiko betrifft __% der Datensätze
- Prosecutor Re-Identifikationsrisiko für den Ausgabedatensatz
 - __, __% der Einträge sind einem Risiko von mehr als 20% ausgesetzt
 - Die zu erwartende relative Anzahl an korrekt re-identifizierten Einträgen ist __, __%
 - Das niedrigste Risiko ist __%, das höchste Risiko ist __%
 - Das höchste Risiko betrifft __% der Datensätze

Feststellung 1: _____

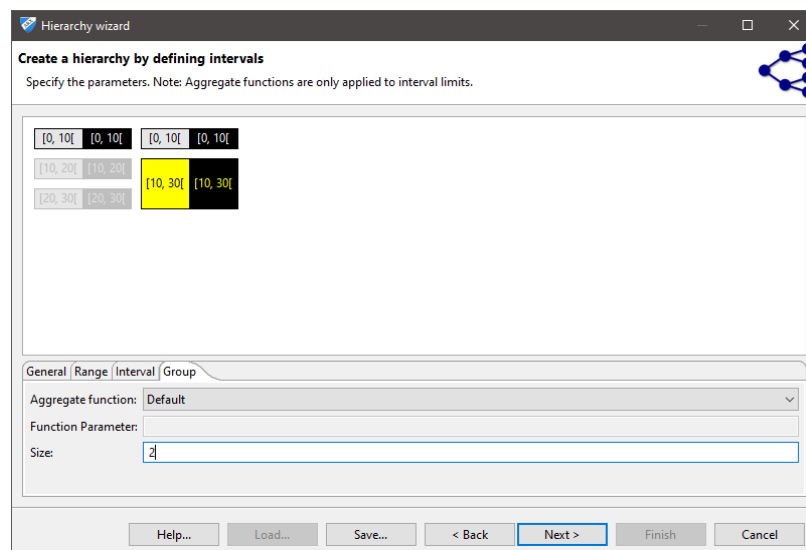
Feststellung 2: _____

Aufgabe 4)

Erstellen Sie eine Generalisierungshierarchie für das Attribut „age“ die folgende Struktur besitzt:



1. Wählen Sie „Edit → Create hierarchy...“ aus dem Anwendungsmenü.
Mit ARX können Sie drei Arten von Hierarchien erzeugen:
 - Basierend auf Intervallen
Beispiel: Alter, Laborwerte
 - Basierend auf einer Ordnung der Attributwerte
Beispiel: Kodierte Diagnosen
 - Durch Maskierung
Beispiel: Postleitzahl
2. Wählen Sie „Use intervals“ und klicken Sie „Next“. Wählen Sie nun den Reiter "Range":
Setzen Sie die untere Schranke auf "0" und die obere Schranke auf "100".
3. Klicken Sie anschließend auf das eine angezeigte Intervall. Definieren Sie ein Basisintervall der Größe 10:
Setzen Sie hier die untere Schranke auf "0" und die obere Schranke auf "10".
4. Fügen Sie nun eine weitere Generalisierungsstufe ein.
 - 4.1. Machen Sie einen Rechtsklick auf das Basisintervall und wählen Sie „Add new level“.
 - 4.2. Dann wählen Sie mit einem Rechtsklick auf den neu entstandenen Knoten „Add after“ um einen Geschwisterknoten auf der gleichen Ebene einzufügen.
 - 4.3. Um das Intervall]10, 30] zu definieren, müssen zwei Intervalle der darunterliegenden Ebene zusammengefasst werden. Definieren Sie demnach eine Größe von 2:



5. Erzeugen Sie anschließend die Intervalle [30, 50[, [50,70[, [70, 90[, [90, 100[.

Hierarchy wizard

Review the hierarchy
Overview of groups and values

#Groups

Level-0	Level-1	Level-2	Level-3
21	[21, 30[[21, 30[*
26	[21, 30[[21, 30[*
28	[21, 30[[21, 30[*
29	[21, 30[[21, 30[*
34	[30, 40[[30, 50[*
37	[30, 40[[30, 50[*
38	[30, 40[[30, 50[*
40	[40, 50[[30, 50[*
43	[40, 50[[30, 50[*
44	[40, 50[[30, 50[*
45	[40, 50[[30, 50[*
46	[40, 50[[30, 50[*
47	[40, 50[[30, 50[*
48	[40, 50[[30, 50[*
50	[50, 60[[50, 70[*

Help... Load... Save... < Back Next > Finish Cancel

Aufgabe 5)

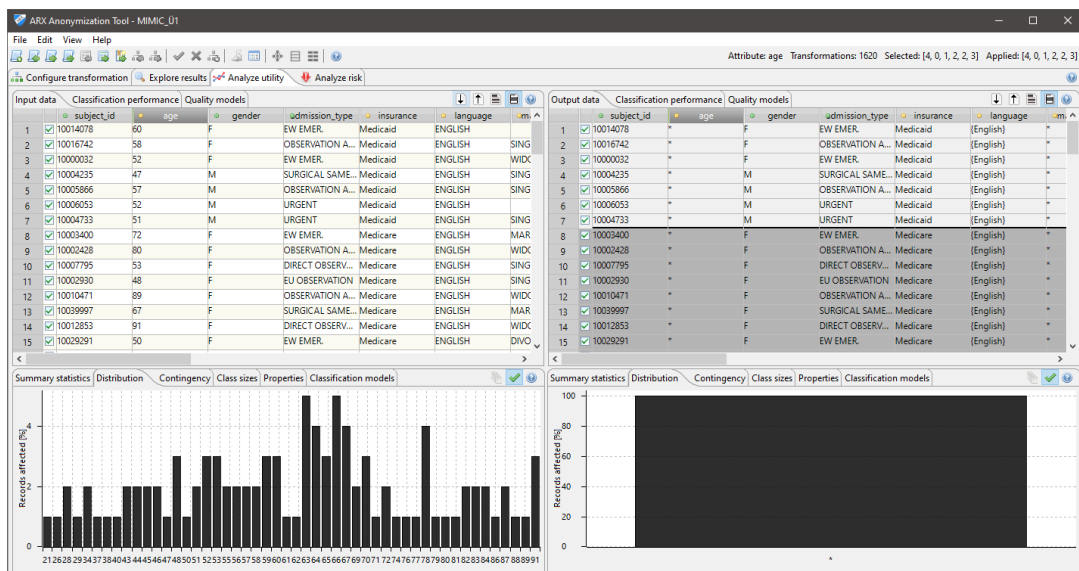
Machen Sie sich weiter mit ARX vertraut. Untersuchen Sie...

- ... die Verteilung des Attributs „age“
- ... statistische Parameter der Verteilung des Attributs „age“
- ... die Kontingenz zwischen „age“ und „marital-status“
- ... die Auswirkung des Anonymisierungsprozesses auf die Vorhersagbarkeit des Attributs „marital-status“

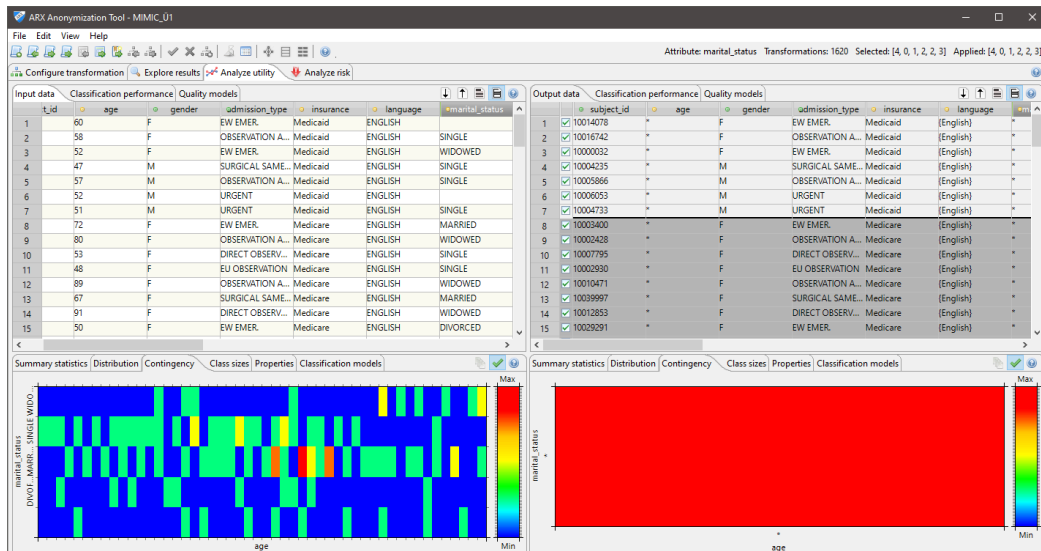
Hintergrund:

- Da es oft nur schwer möglich ist, automatisiert eine angemessene Lösung für ein Anonymisierungsproblem zu finden, werden von ARX verschiedene Funktionalitäten zur Analyse der Nützlichkeit der Ausgabedaten bereitgestellt
- Zur Untersuchung der Auswirkungen des Anonymisierungsprozesses auf ein häufiges Anwendungsszenario implementiert ARX eine Methode aus dem Bereich maschinelles Lernen
- Hierzu wird ein logistisches Regressionsmodell so trainiert, dass es auf Basis einer vorgegebenen Menge von Feature-Attributen den Wert eines Klassenattributs für Einträge im Eingabedatensatz vorhersagen kann
- Die Auswirkung des Anonymisierungsprozesses zeigen sich in einem Rückgang der Genauigkeit des Klassifikators wenn er statt mit den Eingabedaten, mit den anonymisierten Ausgabedaten trainiert wurde
- Die von ARX angezeigten Werte sind als Richtwerte zu sehen

1. Wechseln Sie in die Qualitätsanalyseperspektive von ARX. Selektieren Sie das Attribut „age“ durch einen Mausklick auf die entsprechende Spalte und dann den Reiter „Distribution“.

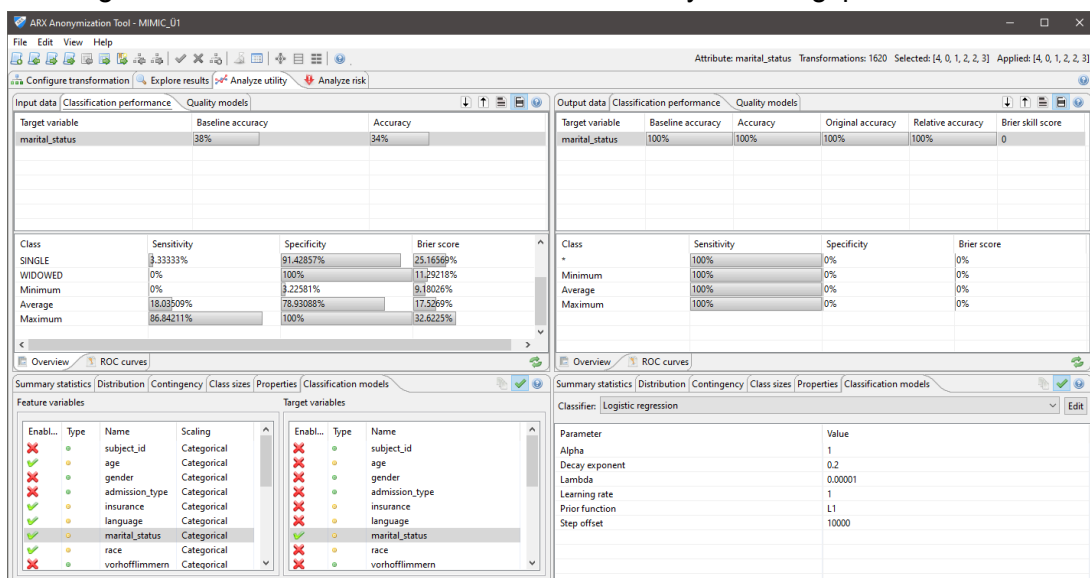


2. Selektieren Sie den Reiter „Summary statistics“ und vergleichen Sie die „Scale of measure“. Selektieren Sie erst das Attribut „age“, dann das Attribut „marital-status“ und wählen Sie dann den Reiter „Contingency“.



Die Kontingenztabelle ist als Heatmap dargestellt. Vergrößerung durch Generalisierung von Attributwerten. Verzerrung durch Unterdrückung von Einträgen. Alle Einträge von marital_status wurden entfernt.

- Wählen Sie nun den Reiter „Classification performance“. Der Klassifikator sagt Werte der Class-Attribute auf Basis der Feature-Attribute vorher. Vergleichen Sie nun die Genauigkeit des Klassifikators vor & nach dem Anonymisierungsprozess.



Hier ist die Accuracy (Relative Genauigkeit) ist von 35 % auf 100% gestiegen.

Aufgabe 6)

Analysieren Sie die Auswirkung von Unterdrückung („Suppression“) auf den Informationsverlust. Anonymisieren Sie den Datensatz mit unterschiedlichen Grenzwerten für die Anzahl der maximal zu entfernenden Einträge („Suppression limit“).

Hintergrund:

- Generalisierung von Attributwerten reduziert die Eindeutigkeit von Einträgen
- Häufig gibt es Einträge in einem Datensatz, die nur durch die Verwendung einer sehr hohen Generalisierungsstufe an Eindeutigkeit verlieren (Ausreisser)
- Um zu verhindern, dass der gesamte Datensatz stark generalisiert werden muss, können diese Einträge automatisch entfernt werden
- Dadurch kann der Verlust an Information reduziert werden
- **Hinweis:** Hier zeigt sich eine grundlegende Eigenschaft vieler Anonymisierungsverfahren. Meist werden häufig auftretende „Muster“ extrahiert.

1. Wechseln Sie in die Konfigurationsperspektive von ARX

- Bestimmen Sie den maximalen Grenzwert der sich zu untersuchen lohnt. Anonymisieren Sie den Datensatz dazu mit einem "suppression limit" von 100%

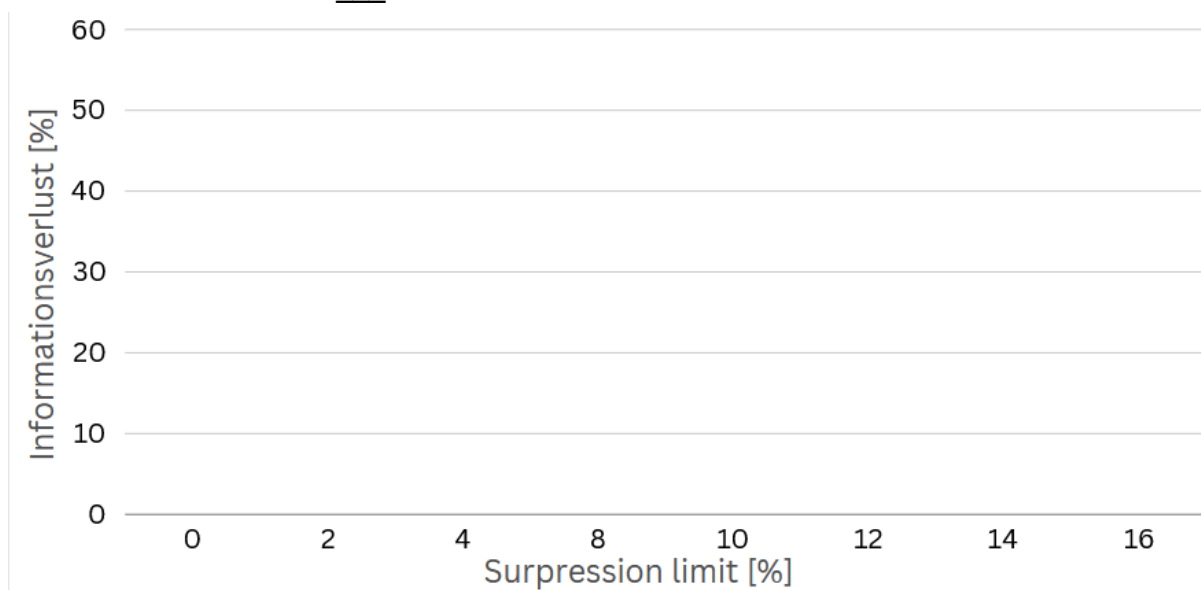
Wechseln Sie in die Qualitätsanalyseperspektive von ARX

- Der Grenzwert entspricht der relativen Anzahl der unterdrückten Einträge (aufgerundet). Sie finden diese Kennzahl im Reiter „Class sizes“.

2. Anonymisieren sie den Datensatz anschließend mit „suppression limits“ von 0%, 2%, ..., bis einschließlich dem bestimmten Grenzwert (aufgerundet).

Füllen Sie aus:

- Der Grenzwert ist ____%



Feststellung 1: _____

Feststellung 2: _____

Aufgabe 7)

Analysieren Sie die Auswirkung von Mikroaggregation auf die Ausgabedaten. Vergleichen Sie dazu Mikroaggregation und Generalisierung. Betrachten Sie das Attribut „age“.

Hintergrund:

- Generalisierung und Mikroaggregation verringern die Eindeutigkeit von Werten
- Bei Generalisierung wird auf einen gemeinsamen allgemeineren Wert abgebildet:
 $50 \rightarrow [50,70]$
 $60 \rightarrow [50,70]$
 $70 \rightarrow [50,70]$
- Mikroaggregation ist ein Verfahren für numerische Attribute, bei dem Werte durch ein Aggregat ersetzt werden:
 $50 \rightarrow \text{Avg}(50,60,70) = 60$
 $60 \rightarrow \text{Avg}(50,60,70) = 60$
 $70 \rightarrow \text{Avg}(50,60,70) = 60$

1. Anonymisieren Sie den Datensatz mit verschiedenen Transformationsmethoden für "age"
 1. Generalisierung
 2. Mikroaggregation mit Modus („Mode“)Wechseln Sie nach jedem Anonymisierungsvorgang in die Qualitätsanalyseperspektive. Notieren Sie die Parameter aus den Reitern „Summary statistics“ und „Class sizes“.

Füllen Sie aus:

	Input	Generalization	Microaggregation
Scale of measure			
Mode			
Median			
Min			
Max			
Suppressed records			

Feststellung 1: _____

Feststellung 2: _____