



Einfluss von Varianten des anonymen Record Linkage auf Gewichtsverteilung und Klassifikation

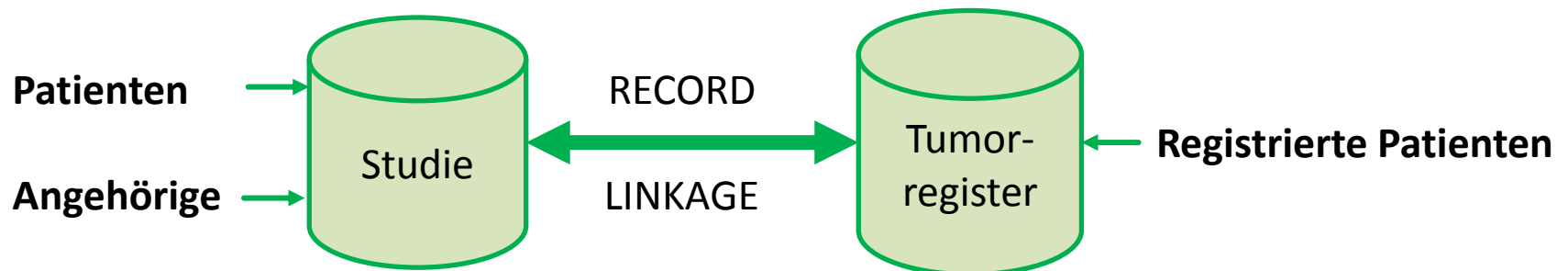
Daniel Nasseh
Jürgen Stausberg

Institut für medizinische Informationsverarbeitung, **Biometrie** und **Epidemiologie**

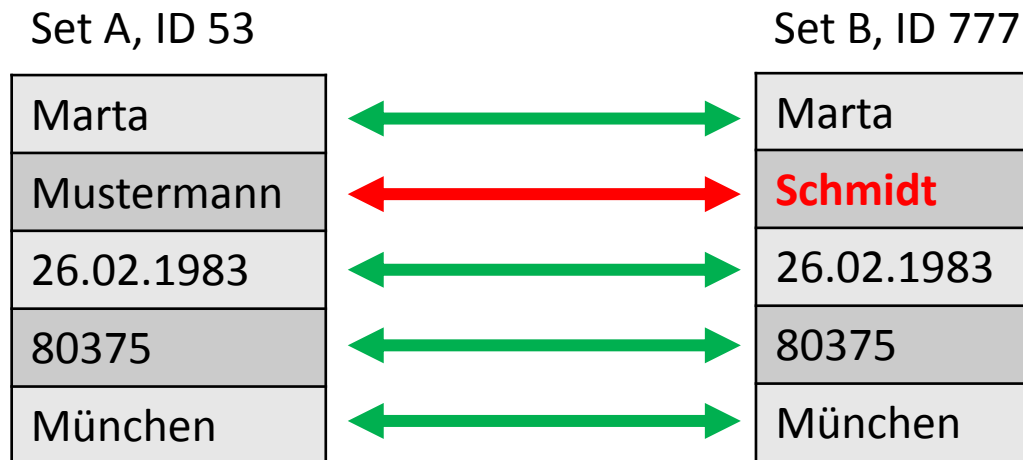
Kontext: Studie zu familiärem Darmkrebs



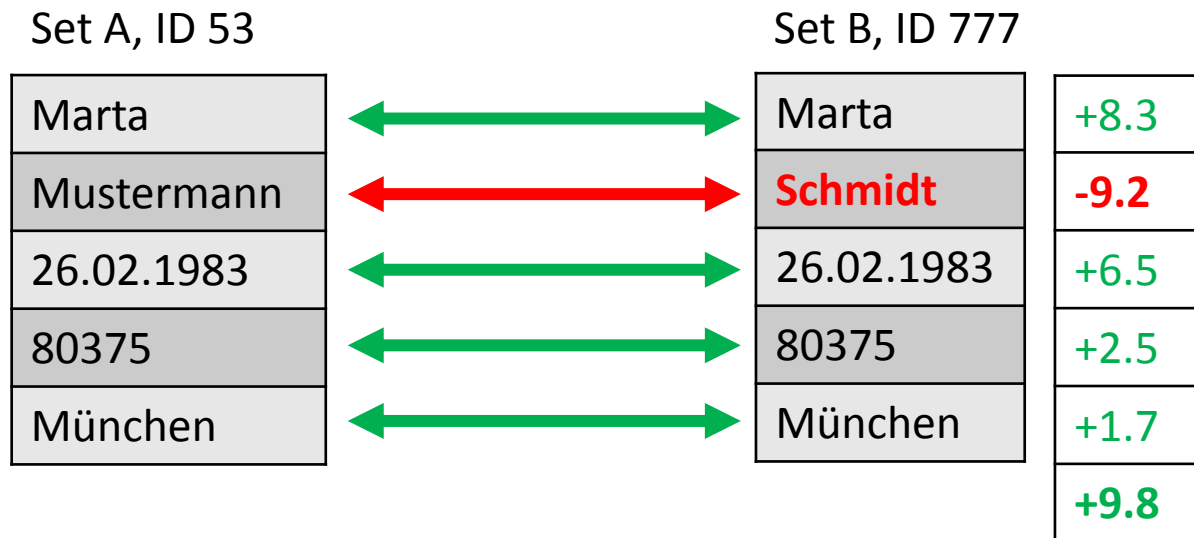
<http://www.darmkrebs-familienstudie.de/>



Probabilistisches Record Linkage



Probabilistisches Record Linkage

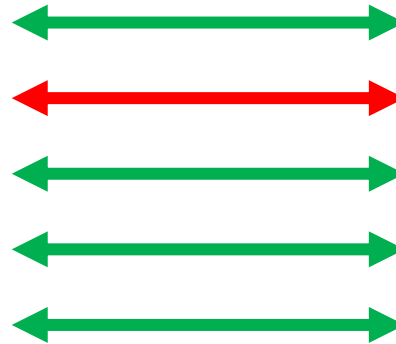


„Anonymes“ Record Linkage



Set A, ID 53

23B1EE4D
A0E2222E
3790A0DF
3FB64B52
855789B4



Set B, ID 777

23B1EE4D
920A3C5V
3790A0DF
3FB64B52
855789B4

+8.3
-9.2
+6.5
+2.5
+1.7
+9.8



Gewichtsdaten



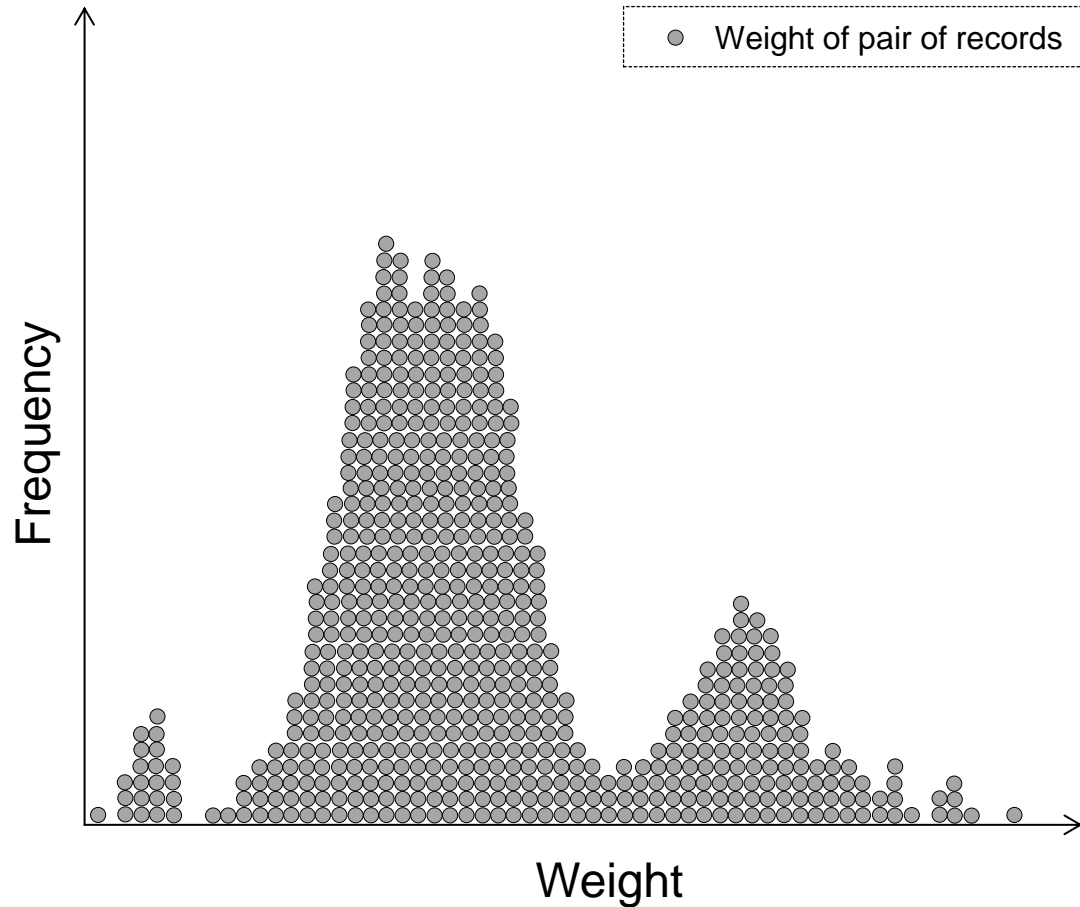
ID A	ID B	Gewicht
3452	45	54,33
2	294	34,22
36	5833	22,19
277	921	17,33
67	513	8,95
794	23	-5,33
...

Gewichtsdaten

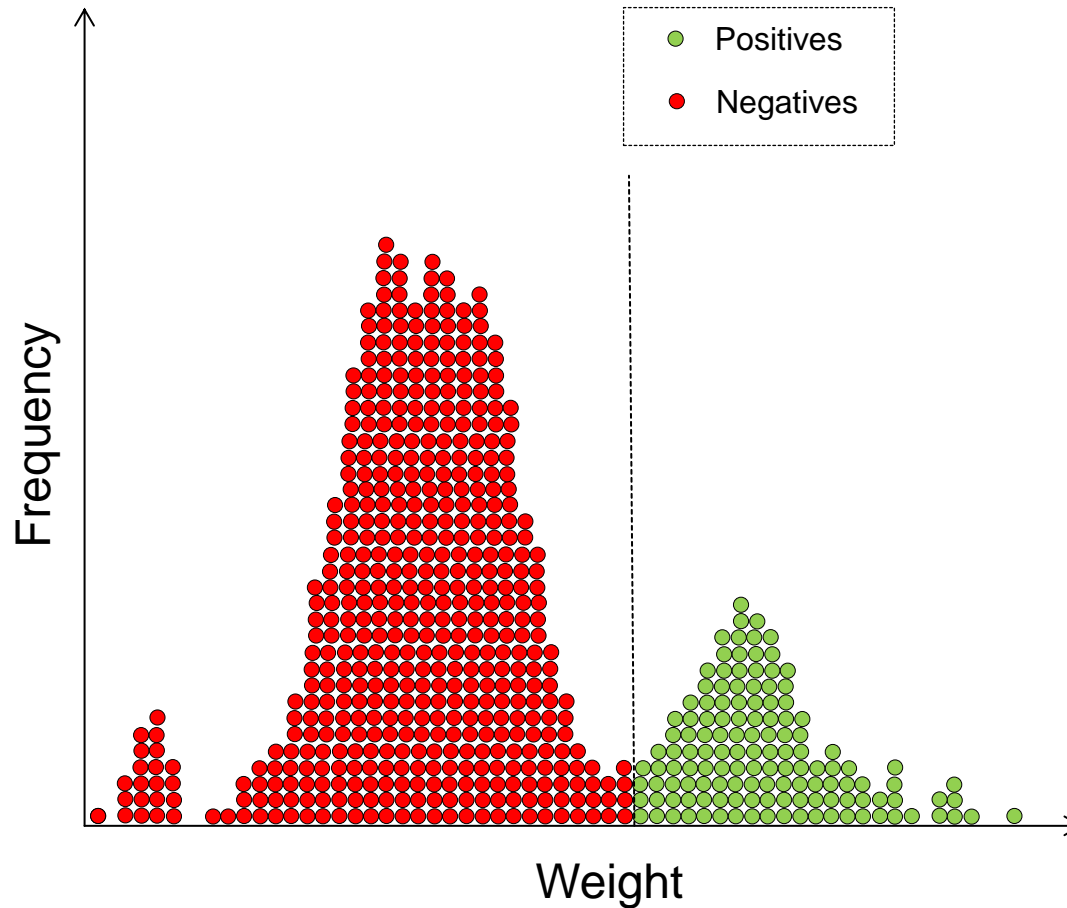


ID A	ID B	Gewicht
3452	45	54,33
2	294	34,22
36	5833	22,19
277	921	17,33
67	513	8,95
794	23	-5,33
...

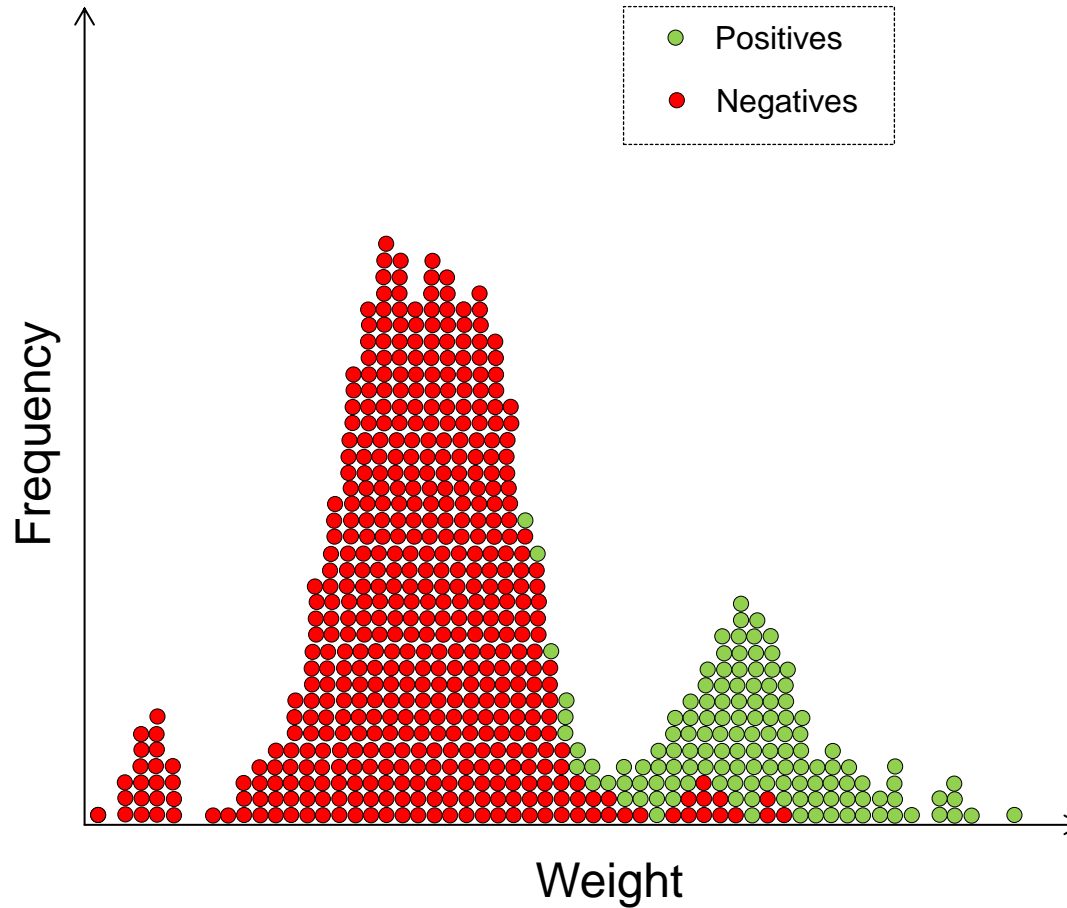
Histogramm der Gewichte paarweiser Vergleiche



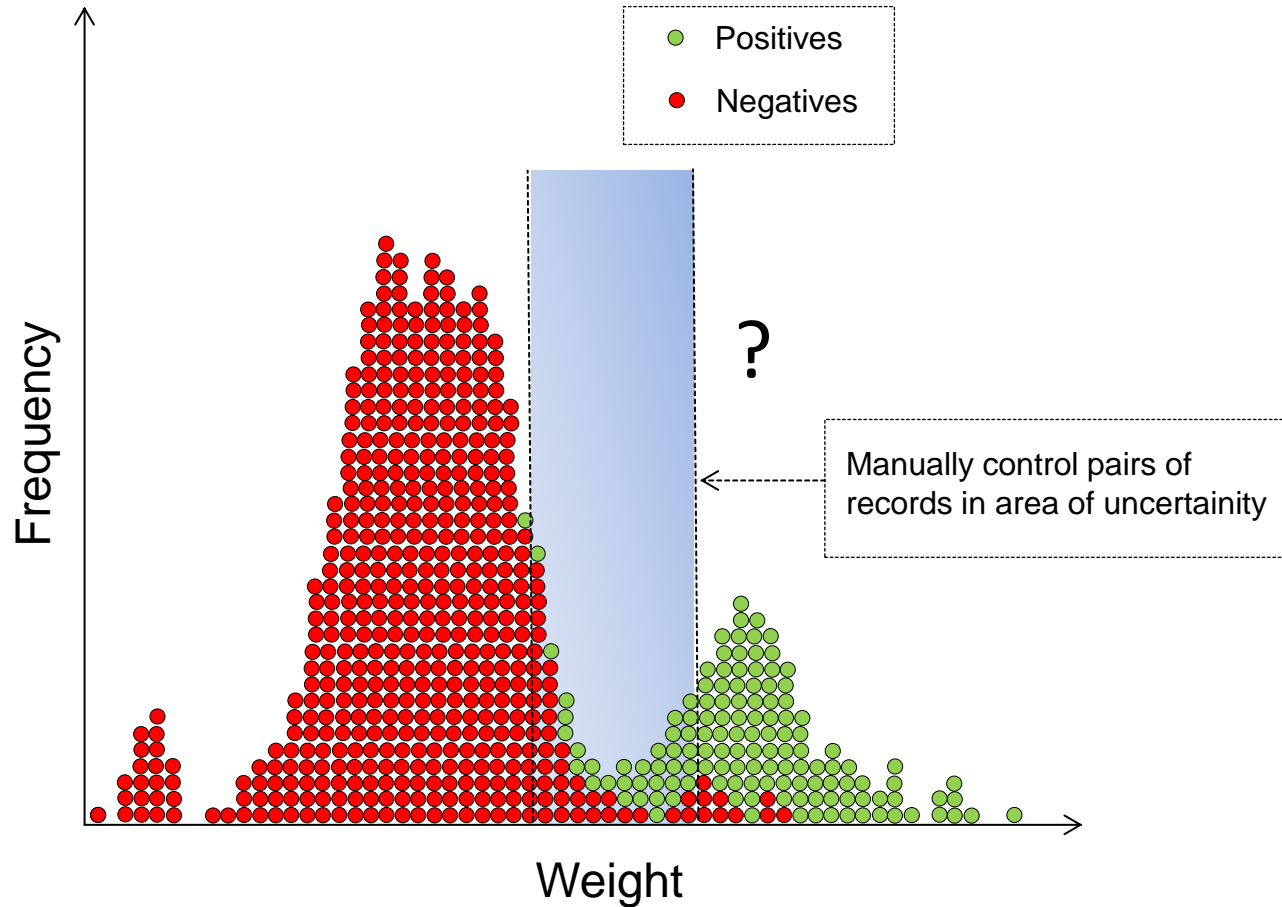
Histogramm der Gewichte paarweiser Vergleiche



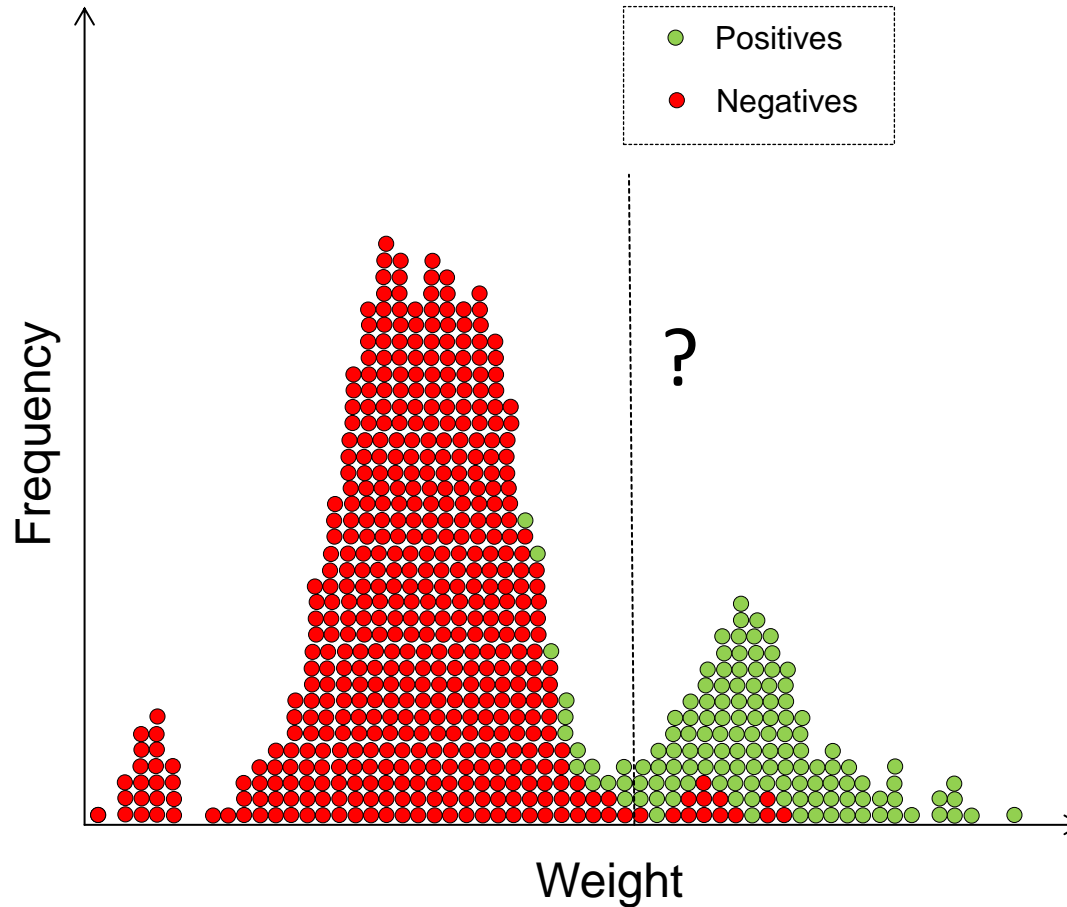
Histogramm der Gewichte paarweiser Vergleiche



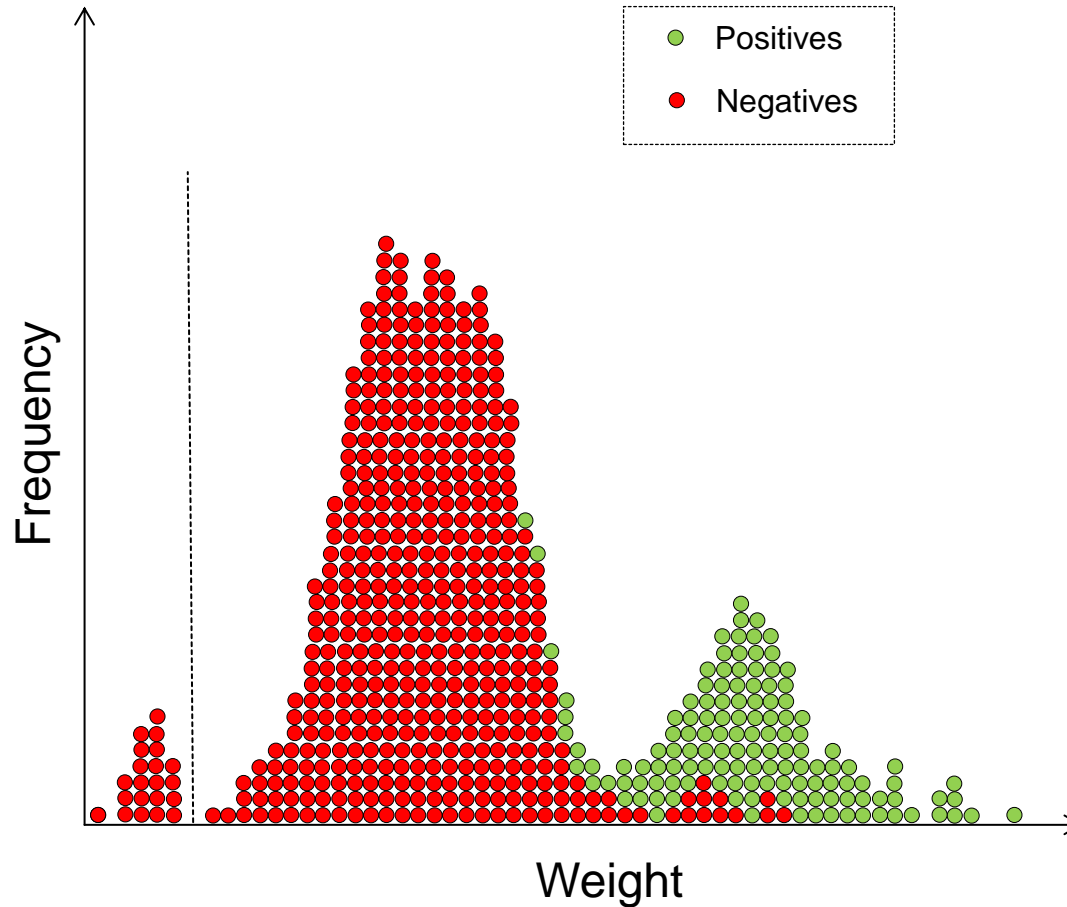
Histogramm der Gewichte paarweiser Vergleiche



Histogramm der Gewichte paarweiser Vergleiche



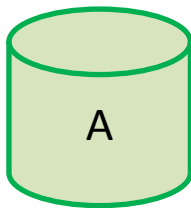
Histogramm der Gewichte paarweiser Vergleiche



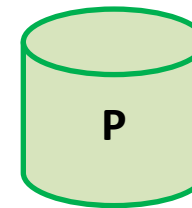
Experimentelles Setup



Künstliches Datenset



Öffentliches Datenset [*]

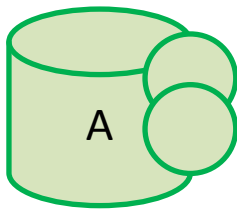


* <http://secondstring.cvs.sourceforge.net/viewvc/secondstring/secondstring/>.

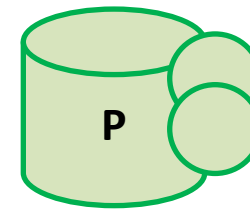
Experimentelles Setup



Künstliches Datenset



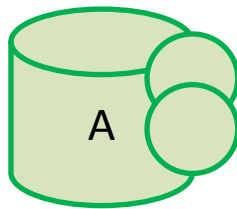
Öffentliches Datenset



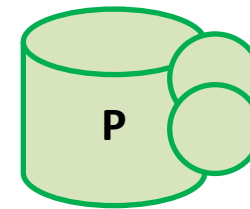
Experimentelles Setup



Künstliches Datenset



Öffentliches Datenset



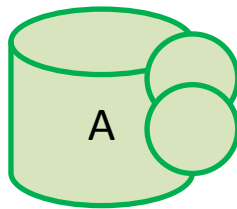
Record Linkage



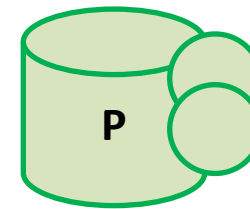
Experimentelles Setup



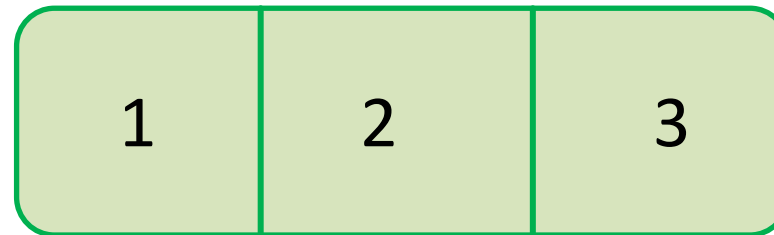
Künstliches Datenset



Öffentliches Datenset



Record Linkage

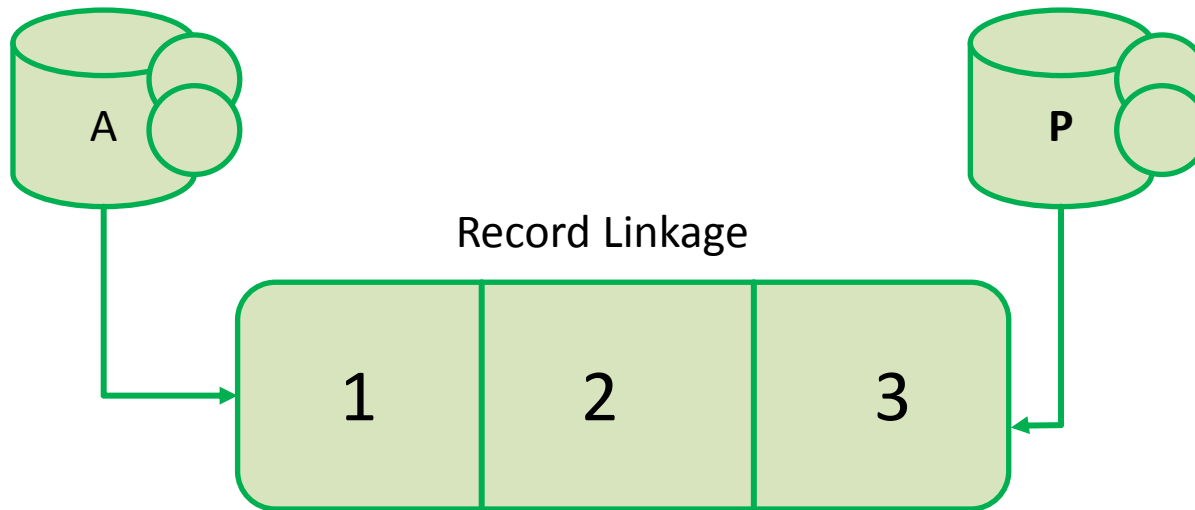


Experimentelles Setup



Künstliches Datenset

Öffentliches Datenset

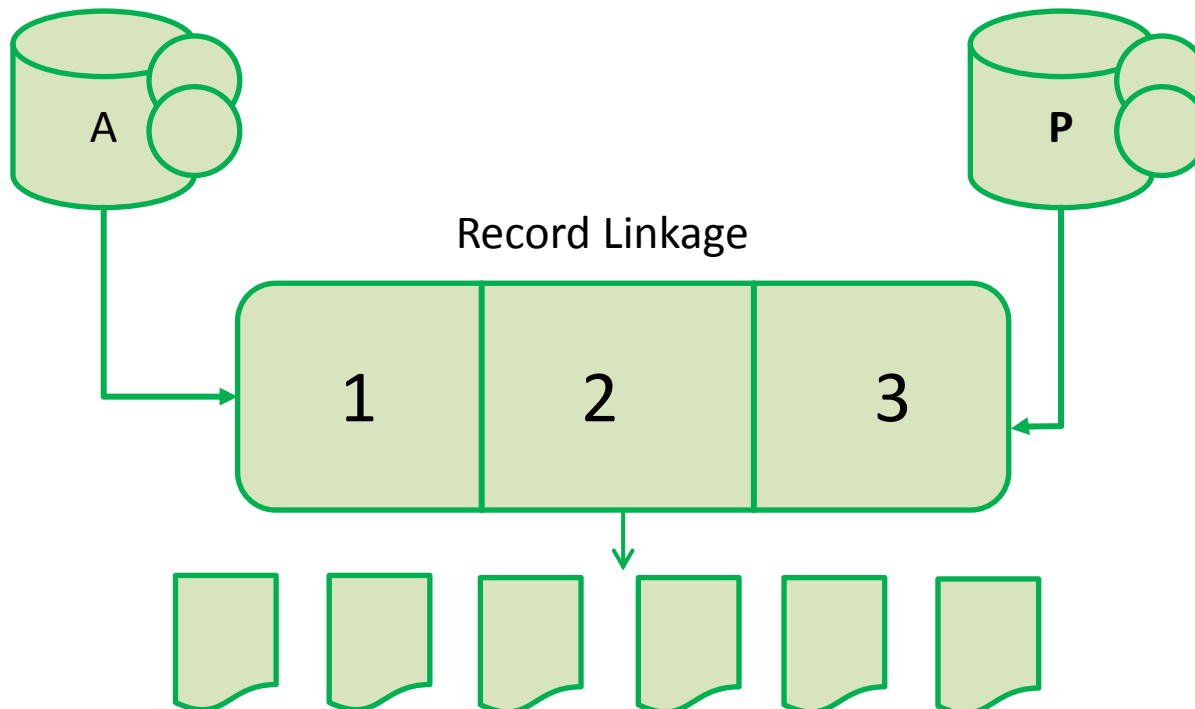


Experimentelles Setup



Künstliches Datenset

Öffentliches Datenset



Konfigurationen

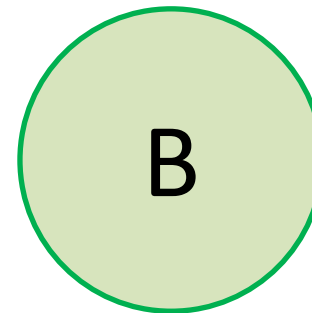
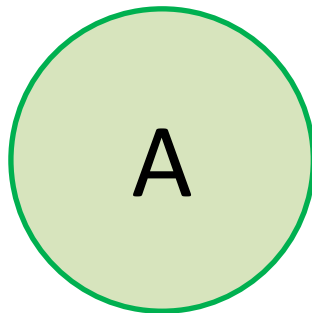


1. Blocking (naiv / mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

Konfigurationen



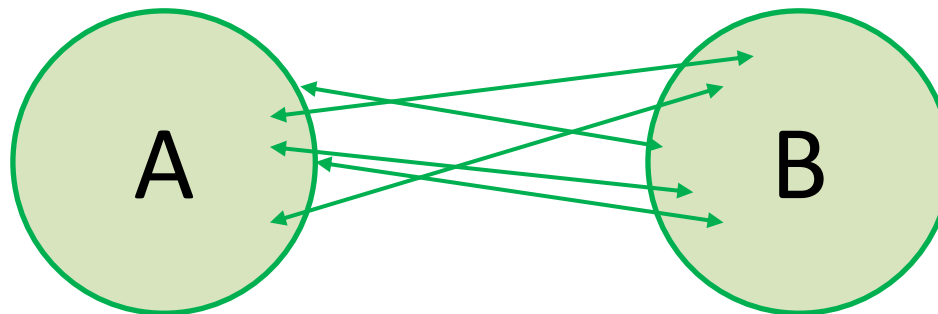
1. Blocking (mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning



Konfigurationen



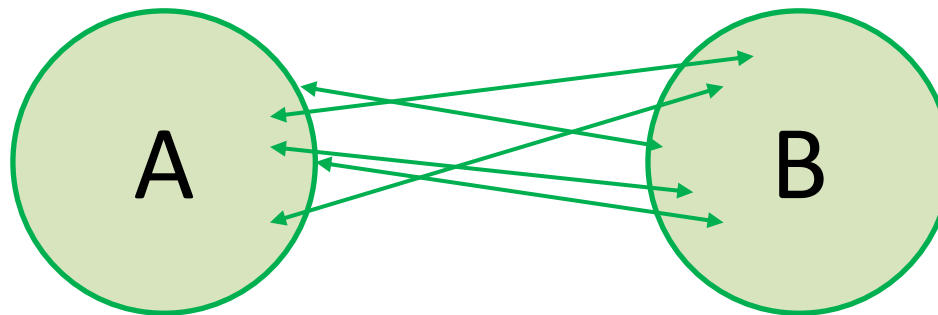
-
1. Blocking (mit mehrfach auftretenden Links)
 2. Blocking (mit „unique“ Links)
 3. Multi-Link-Cleaning



Konfigurationen



1. Blocking (mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

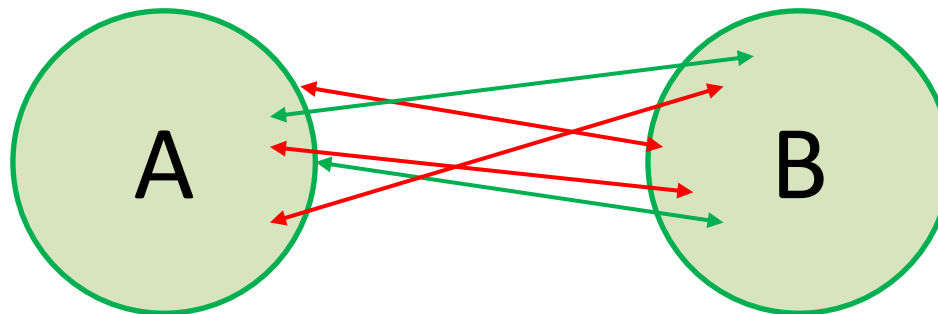


Bei $10.000 * 10.000 = 100.000.000$ Berechnungen

Konfigurationen



1. Blocking (mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

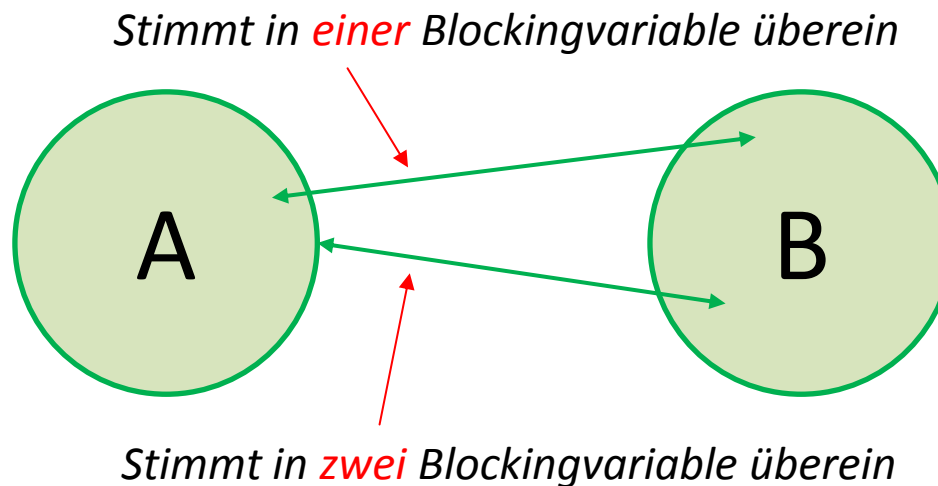


Record Linkage nur noch auf Vergleichen
die in Blocking-Variablen übereinstimmen!

Konfigurationen



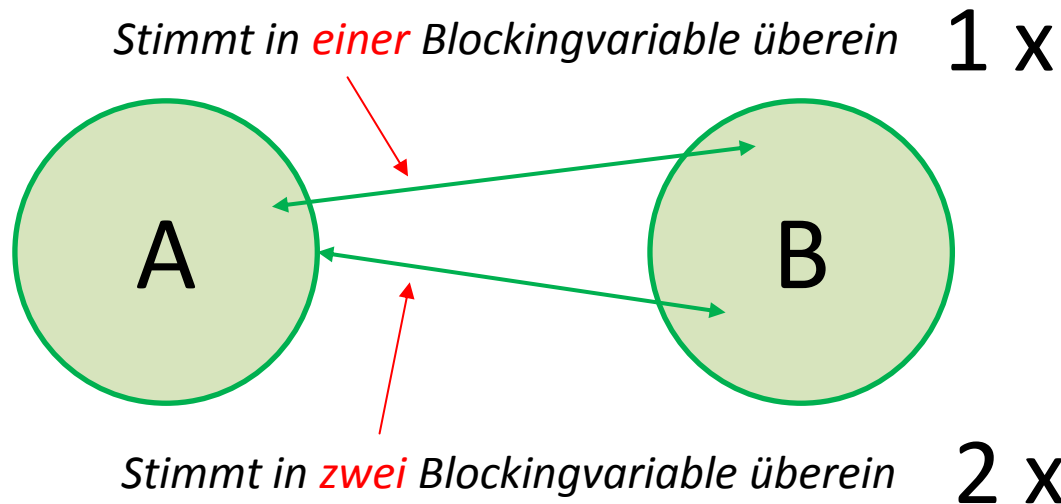
1. Blocking (mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning



Konfigurationen



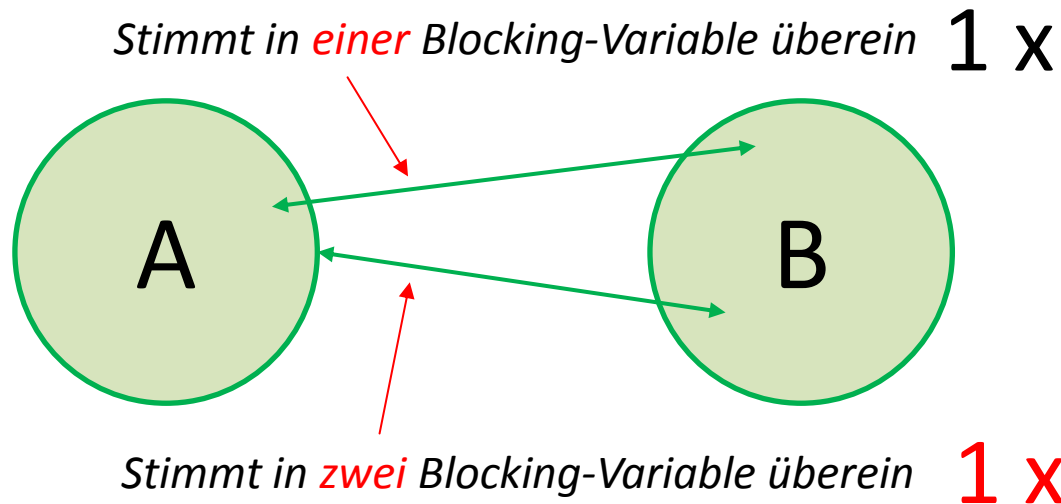
1. Blocking (mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning



Konfigurationen



1. Blocking (mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning



Konfigurationen



1. Blocking (naiv / mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

Konfigurationen



1. Blocking (naiv / mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

ID A	ID B	Gewicht
55	21	+48.8
27	33	+21.4
55	84	+19.3

Konfigurationen



1. Blocking (naiv / mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

ID A	ID B	Gewicht
55	21	+48.8
27	33	+21.4
55	84	+19.3

Konfigurationen



1. Blocking (naiv / mit mehrfach auftretenden Links)
2. Blocking (mit „unique“ Links)
3. Multi-Link-Cleaning

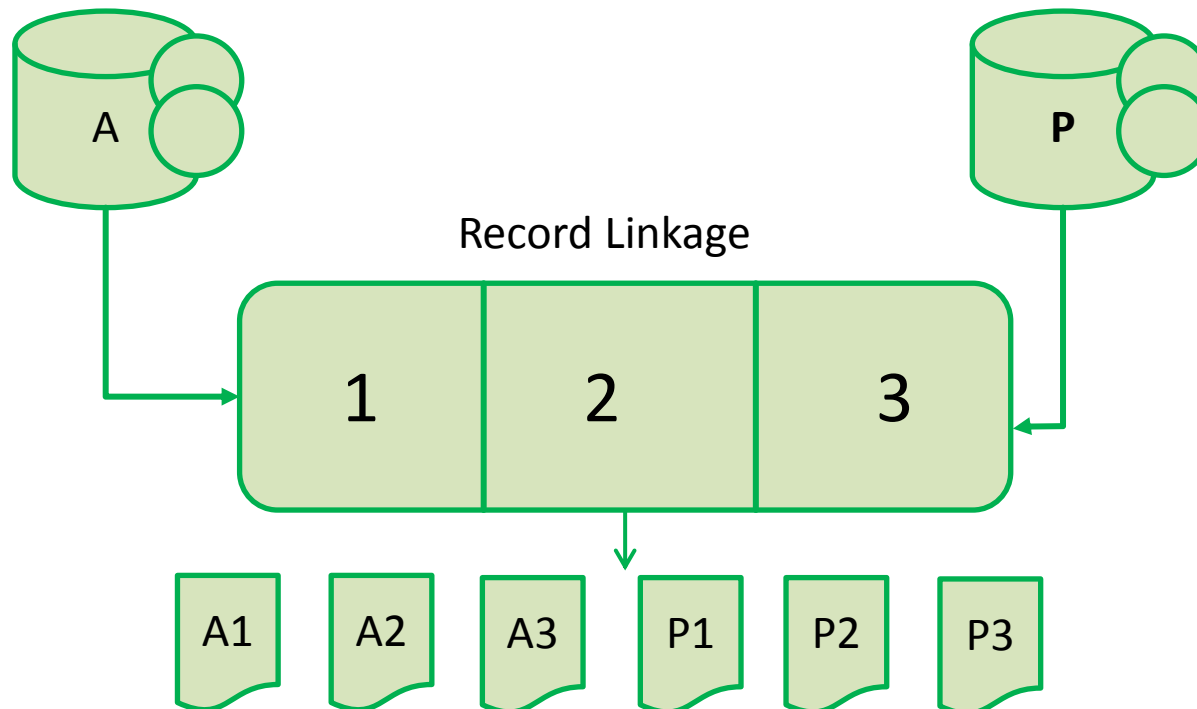
ID A	ID B	Gewicht
55	21	+48.8
27	33	+21.4
55	84	+19.3

Experimentelles Setup



Künstliches Datenset

Öffentliches Datenset

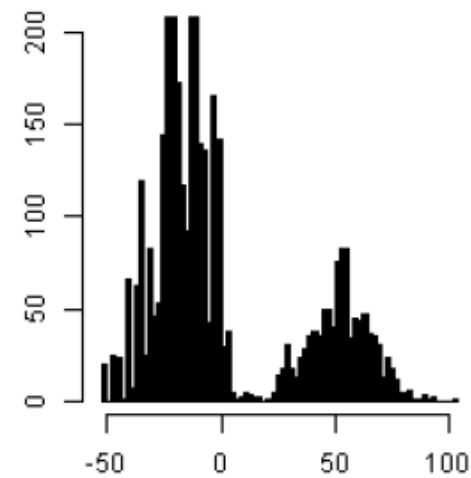
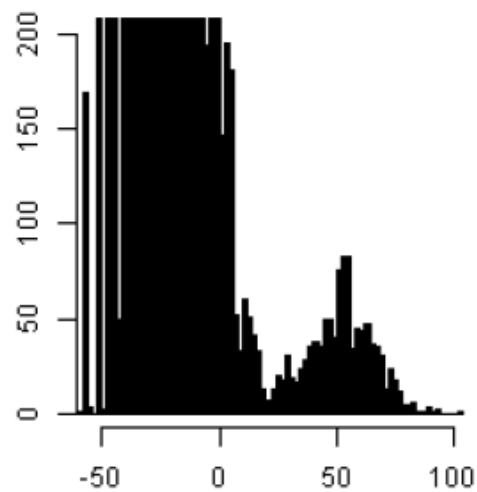
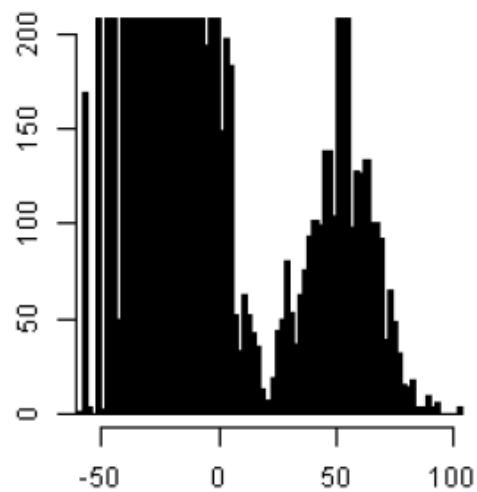


Including duplicates

Without duplicates

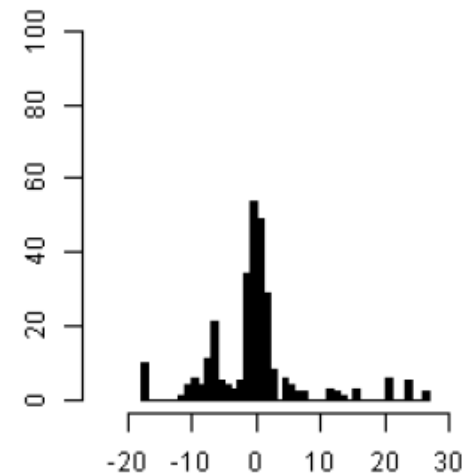
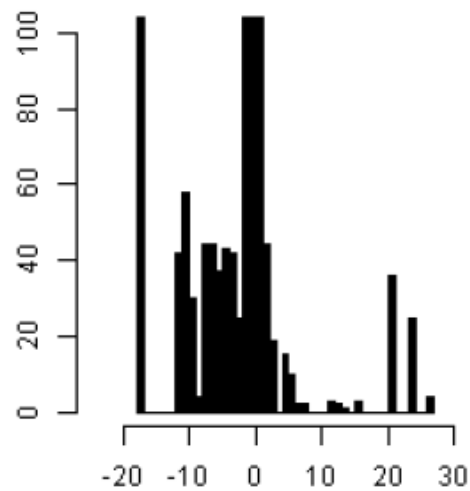
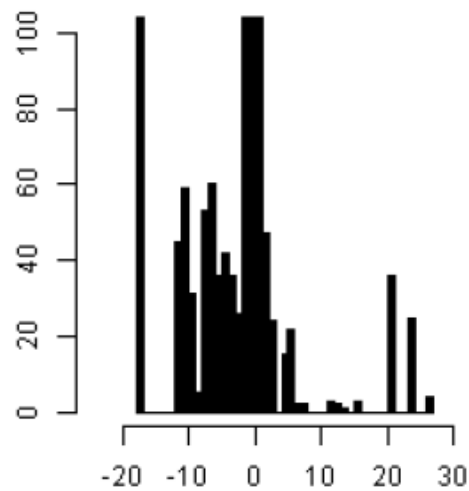
Multi-link-cleaning

Artificial dataset



Census dataset

frequency →



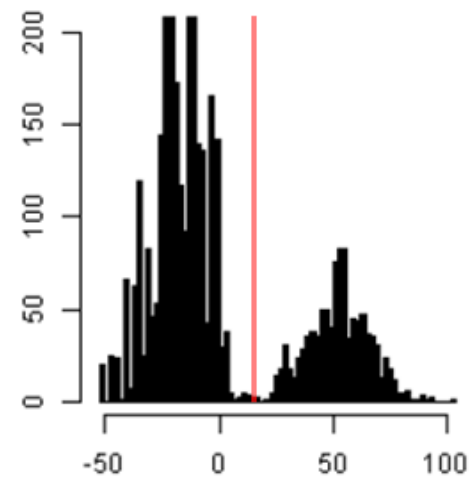
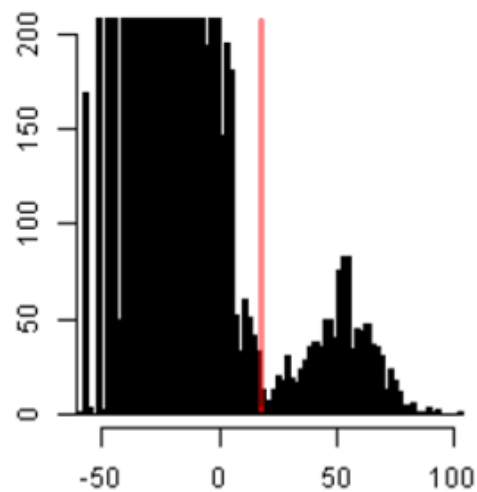
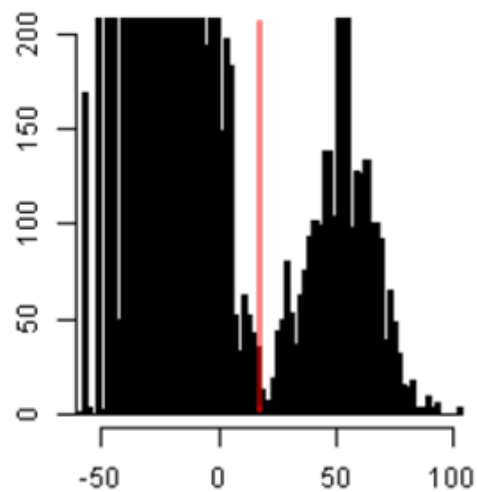
weight →

Including duplicates

Without duplicates

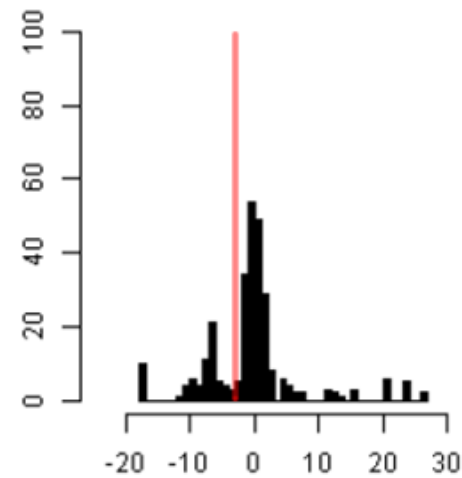
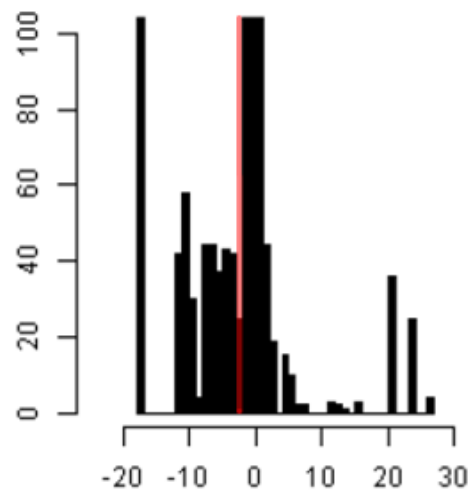
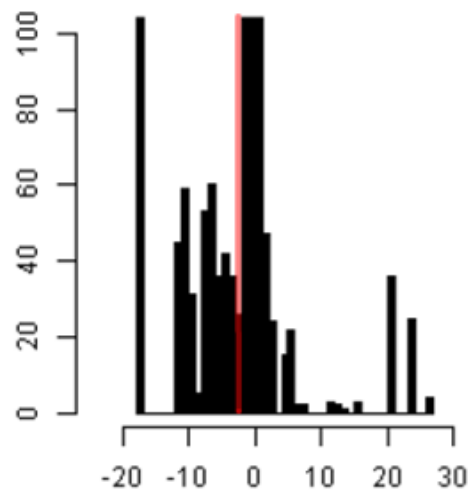
Multi-link-cleaning

Artificial dataset



Census dataset

frequency →



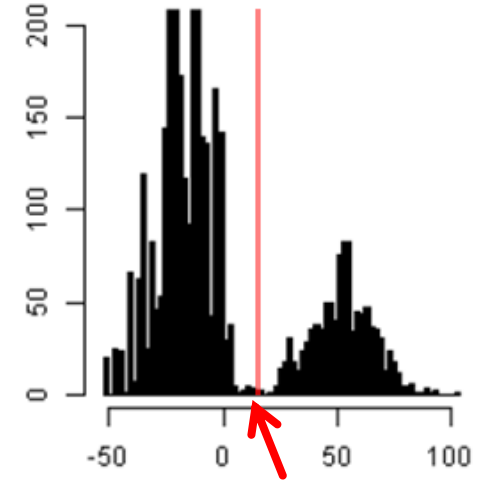
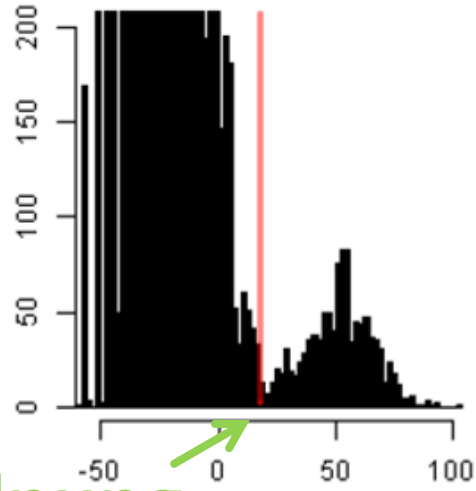
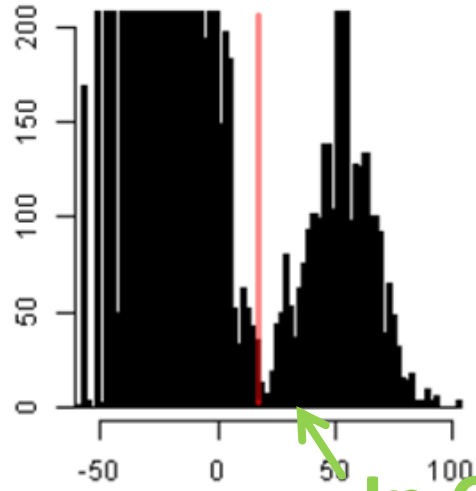
weight →

Including duplicates

Without duplicates

Multi-link-cleaning

Artificial dataset



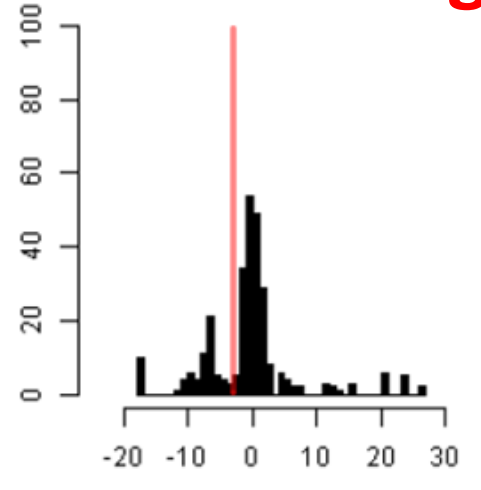
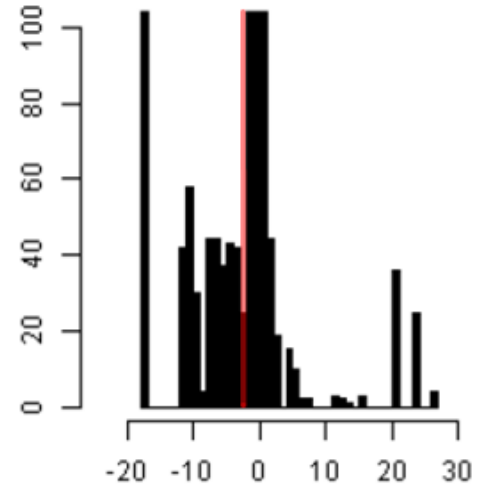
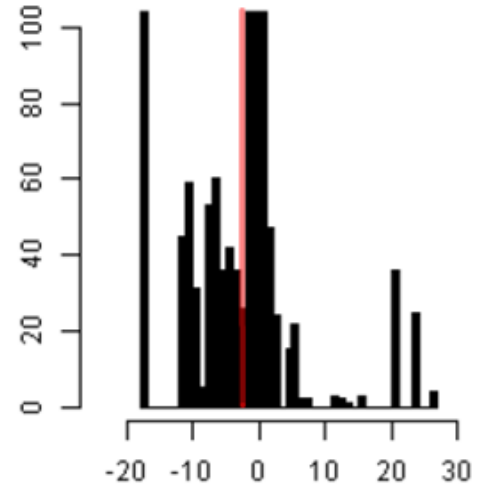
In Ordnung

uneindeutig

Census dataset

frequency →

weight →

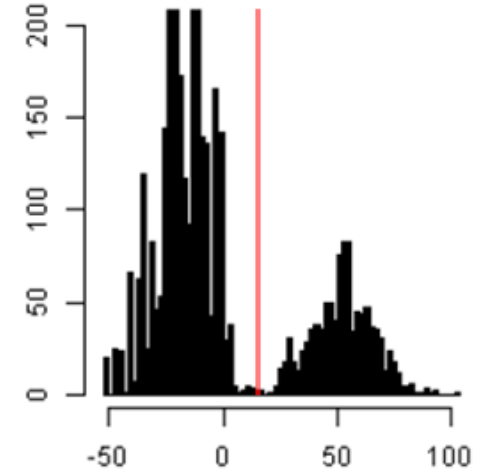
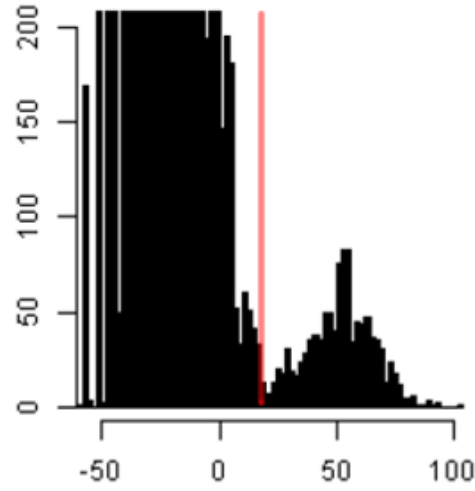
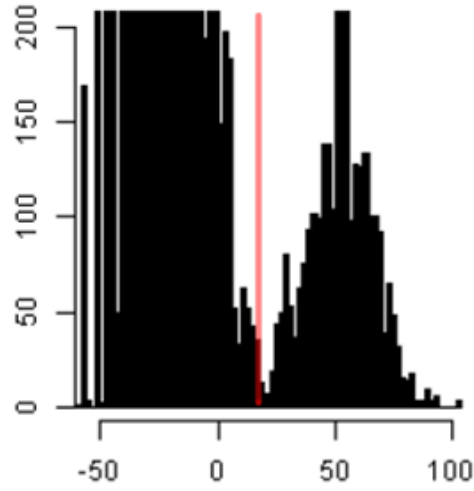


Including duplicates

Without duplicates

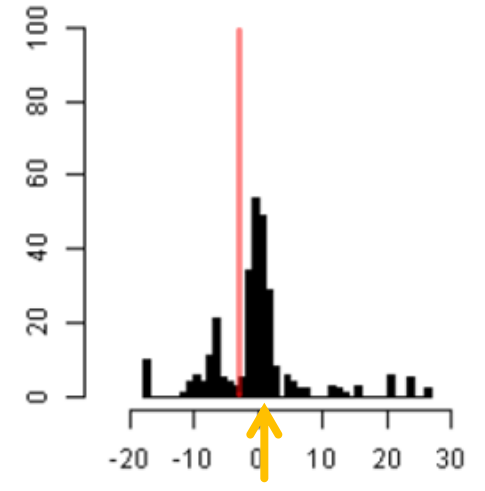
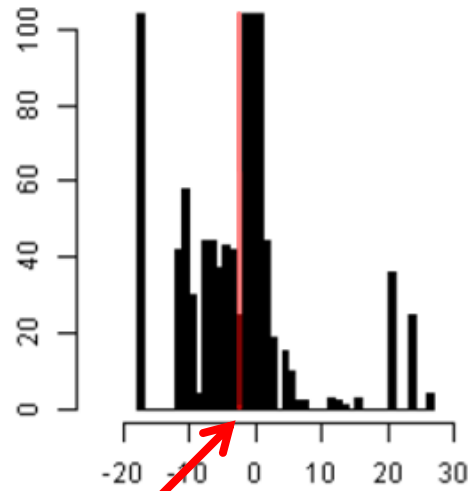
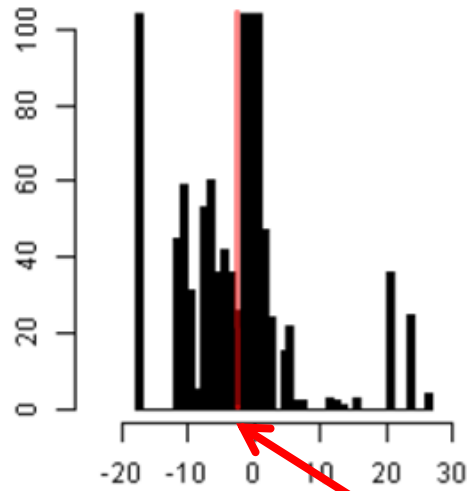
Multi-link-cleaning

Artificial dataset



Census dataset

frequency →



weight →

uneindeutig

erahnbär

Haupterkenntnisse



- Bei Verwendung von unüberwachten Klassifikationssystemen ist die Gefahr einer kompletten Fehlklassifikation vorhanden
- Klassifikation ist nicht automatisierbar
- Es kann hilfreich sein verschiedene Konfigurationen zu verwenden um die korrekte Schranke zu ermitteln
- Bei schmutzigen oder kleinen Datensets bieten sich eher überwachte Klassifikationssysteme an (Boosting, Bagging).

Haupterkenntnisse



- Bei Verwendung von unüberwachten Klassifikationssystemen ist die Gefahr einer kompletten Fehlklassifikation vorhanden
- Klassifikation ist nicht automatisierbar
- Es kann hilfreich sein verschiedene Konfigurationen zu verwenden um die korrekte Schranke zu ermitteln
- Bei schmutzigen oder kleinen Datensets bieten sich eher überwachte Klassifikationssysteme an (Boosting, Bagging).

Haupterkenntnisse



- Bei Verwendung von unüberwachten Klassifikationssystemen ist die Gefahr einer kompletten Fehlklassifikation vorhanden
- Klassifikation ist nicht automatisierbar
- Es kann hilfreich sein verschiedene Konfigurationen zu verwenden um die korrekte Schranke zu ermitteln
- Bei schmutzigen oder kleinen Datensets bieten sich eher überwachte Klassifikationssysteme an (Boosting, Bagging).

Haupterkenntnisse



- Bei Verwendung von unüberwachten Klassifikationssystemen ist die Gefahr einer kompletten Fehlklassifikation vorhanden
- Klassifikation ist nicht automatisierbar
- Es kann hilfreich sein verschiedene Konfigurationen zu verwenden um die korrekte Schranke zu ermitteln
- Bei schmutzigen oder kleinen Datensets bieten sich eher überwachte Klassifikationssysteme an (Boosting, Bagging). => *Basieren auf Trainingssets*



Vielen Dank für Ihre
Aufmerksamkeit!

Fragen?