

# Causal Record Linkage

Kayané Robach

Joint work with Michel Hof, Mark van de Wiel, Stéphanie van der Pas



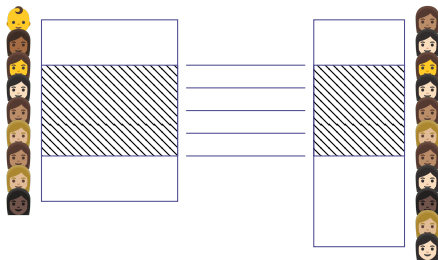
EPIDEMIOLOGY AND  
DATA SCIENCE



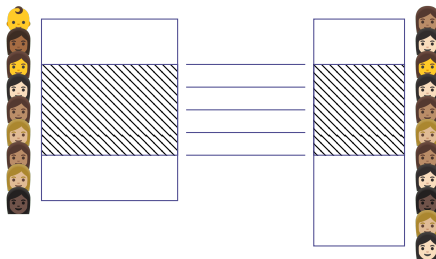
Bigstatistics

How to combine individuals records across heterogeneous data sets?

How to combine individuals records across heterogeneous data sets? *Record Linkage*

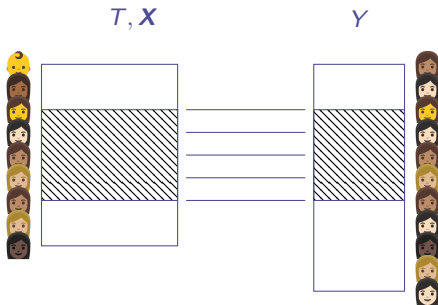


How to combine individuals records across heterogeneous data sets? *Record Linkage*



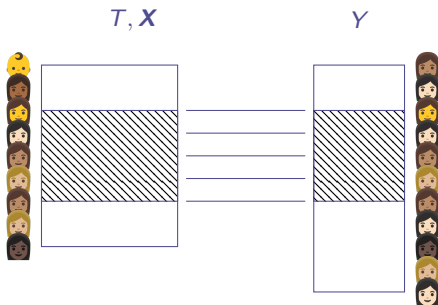
Connect information on the same individuals

How to combine individuals records across heterogeneous data sets? *Record Linkage*



baseline information and exposure  $\longleftrightarrow$  outcomes

How to combine individuals records across heterogeneous data sets? *Record Linkage*



Perform causal inference

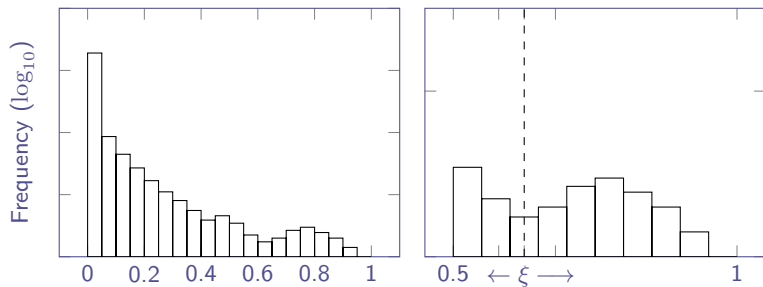
Without unique identifiers?

Without unique identifiers → use partially identifying information  
(birth year, ethnicity, postal code) ⚠ also causally relevant covariates **X**



Probabilistically recover complete profiles  $(T, \mathbf{X}, Y)$  with *Record Linkage*

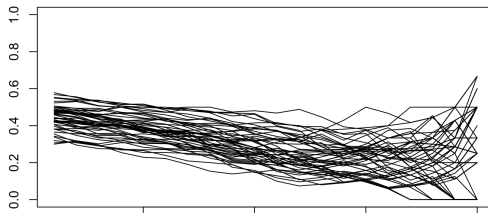
Rigorously select the profiles with linkage score  $> \xi$



How to choose this threshold  $\xi$ ?

We select the set of linked records, trying to minimise:

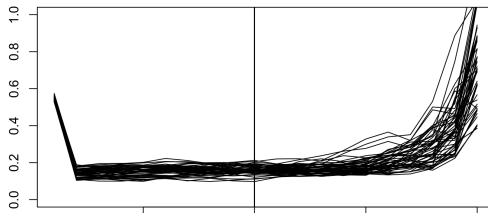
→ the estimated False Discovery Proportion



How to choose this threshold  $\xi$ ?

We select the set of linked records, trying to minimise:

→ the Maximum Mean Discrepancy



Stringency on the selection of linked records



Low ( $\xi$  closer to 0.5)

High ( $\xi$  closer to 1)

## Stringency on the selection of linked records

Low ( $\xi$  closer to 0.5)

High ( $\xi$  closer to 1)

Positivity 😊

## Stringency on the selection of linked records

Low ( $\xi$  closer to 0.5)

High ( $\xi$  closer to 1)

Positivity 😊

Conditional exchangeability 😬

## Stringency on the selection of linked records

Low ( $\xi$  closer to 0.5)

High ( $\xi$  closer to 1)

Positivity 😊

Conditional exchangeability 😬

Consistency 🤖



## Stringency on the selection of linked records

Low ( $\xi$  closer to 0.5)

Positivity 😊

Conditional exchangeability 😬

Consistency 🤪

High ( $\xi$  closer to 1)

Consistency 😊

## Stringency on the selection of linked records

---

Low ( $\xi$  closer to 0.5)

Positivity 😊

Conditional exchangeability 😞

Consistency 🙄

High ( $\xi$  closer to 1)

Conditional exchangeability 😊

Consistency 😊

## Stringency on the selection of linked records

---

Low ( $\xi$  closer to 0.5)

Positivity 😊

Conditional exchangeability 😞

Consistency 🙄

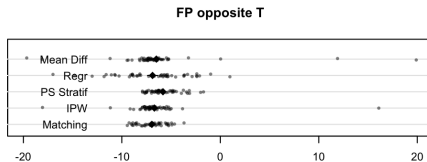
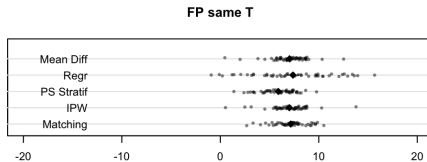
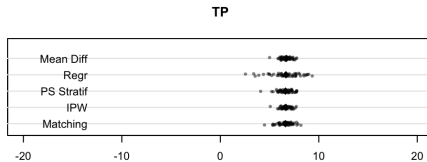
High ( $\xi$  closer to 1)

Positivity 🤖

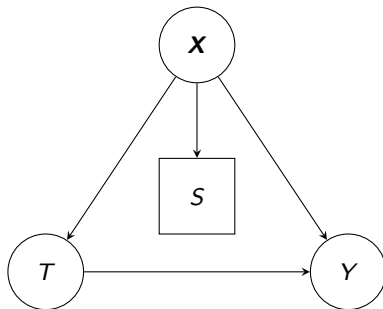
Conditional exchangeability 😊

Consistency 😊

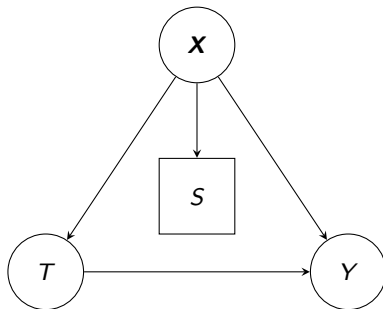
Atypical profiles over-represented!



What is happening?

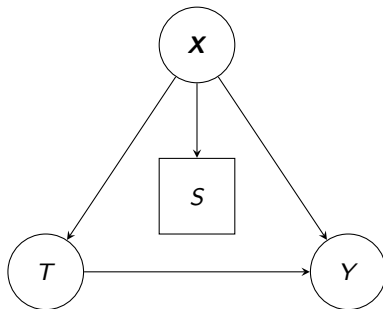


What is happening?



Selection  $S$  due to linkage

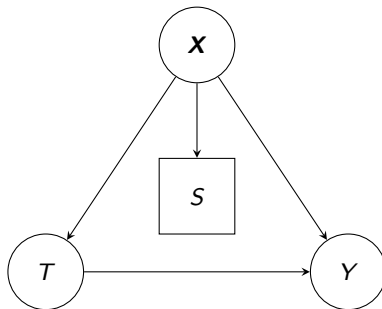
What is happening?



Selection  $S$  due to linkage

*(selection backdoor criterion, effect is transportable from selected to baseline population)*

What is happening?



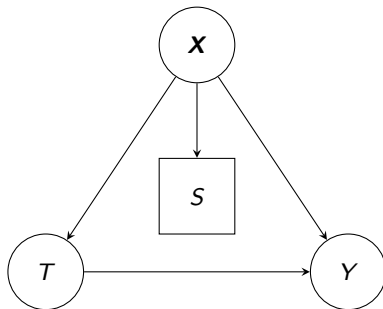
Selection  $S$  due to linkage

(*selection backdoor criterion, effect is transportable from selected to baseline population*)

'Recoverability from sampling bias'



What is happening?



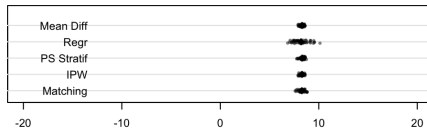
Selection  $S$  due to linkage

(*selection backdoor criterion, effect is transportable from selected to baseline population*)

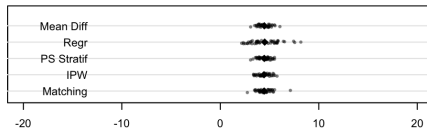
'Recoverability from sampling bias'

→ Use G-methods

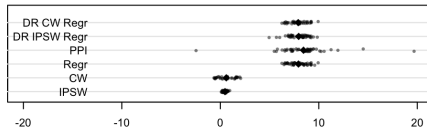
True Links



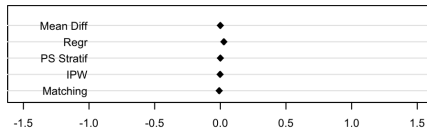
Linked



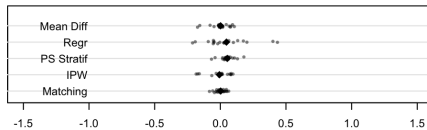
G-Linked



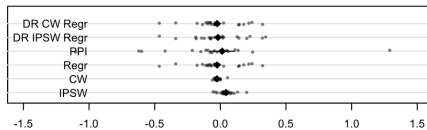
### True Links



### Linked



### G-Linked



Thank You!