

# Week 6: Data Analysis

Charlotte Hadley

# Topics for today

1. Why is data anonymisation important?
2. What are the risks of de-anonymisation of health data?
3. Anonymity measures: k-anonymity and l-diversity
4. Case studies of de-anonymisation

... and why anonymity measures are often not

5. R packages for working with anonymous data

Why is data annotation important?

# Data isn't always a

Mostly during this course we've been talking explicitly being collected - and participating

But sometimes not the case.

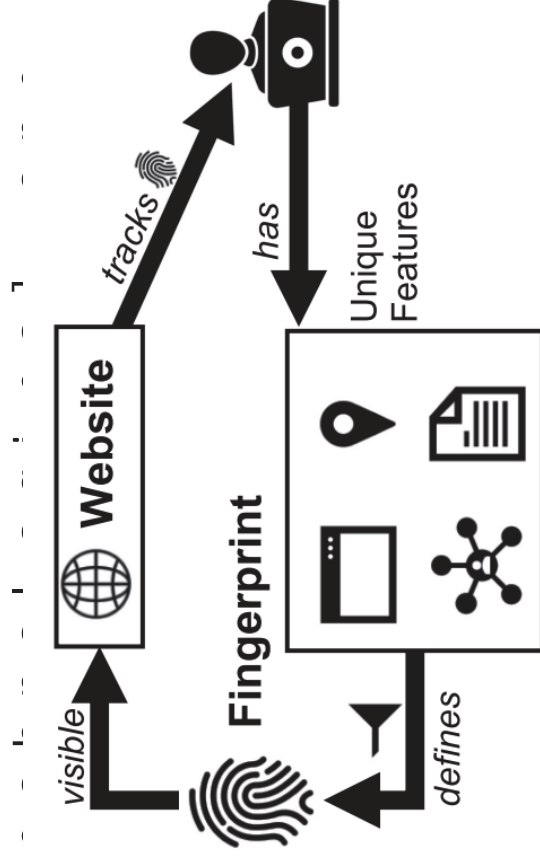
Data is collected continuously about individuals often without implicit consent.

# Data tracking across

Cookies [and similar technologies] present in the mobile

They allow websites to track what you do on their website with websites.

... but what about all the other devices you use?



(b) A representation of how device fingerprinting enables the tracking of users on the web. The specific fingerprint may vary across different websites, since it can include different features unique to a particular device, which are requested by the website.

Source: Kretschmer et

# GDPR and Cookies

The GDPR<sup>2</sup> mandated "cookie notices" everywhere.

For a common sense the policy [gdpr.eu/rec<sup>3</sup>](#)

The policy features in the cookie availability

The policy was well met necessary

In 2016 it was demonstrated the top million websites of tracking

However, ~~there~~ evidence policy hasn't materialized privacy<sup>1</sup>

# Cookies: Fingerprinting, patterns

These are the primary ways websites circumv

- The biggest issue with third-party cookies is that they are harder to track.
- User-agent cookies by default select the user's preferred language from the browser's settings.
- Dark pattern employed to prevent users from deleting cookies.

Overall, the policy has probably made things

GDPR impacted smaller advertiser segment compared to brands, such as Google and Facebook, leading to these companies, which, in turn, may increase decrease in Source: <sup>1</sup> Kretschmer et al. 2021

# Why do websites want to track you?

There are simple un-  
this:

- Selling tracking data  
networks
- Using tracking data  
users

But there are inter-

- Anticipatory shipping  
purchases
  - Amazon's patented  
in 2013 thoroughly  
Eva-Maria Nyckel
- Used in the agro-

Can we get some exam-  
ples they track?



Wired: Ma dehnukmai tt aar a\manan: My inden

I strongly recommend we a deidng oaluk / a fr t h i l . s e / a m r .

- The article is from 2014, but as we've discussed made this situation worse.
- The article also provides a wealth of other examples for your assistance
- The quote below is from the article's author

an woman  
**Madhumita Venkataramanan: My  
 identity for sale**

Even as you're reading this, your smartphone can reveal  
 your location

































































































































































































































































































































# Data tracking during

These A&E waiting sent to folks that This survey is sent into it.

The Section 251 of 2006 provides for confidential patient consent for purposes of HRA or the Secreta Health and Social

That an exception to Act<sup>11</sup>

This means that the at this data and de

This exception was COVID-19 pandemic surveying hospital

Note that the NHS may opt out once you've

However, always case, patients/ser given the opportunity

---

Data anonymisation is  
because data is coll  
all of the time

... how does that match up with GDPR and the

# Individual rights for

The ADPP provides 8 rights for ~~we~~ don't know that individuals: collected - or by which

individual rights are informed

- The right of access

This opens up lots of

- The right to rectify

We'll discuss these in

- The right to erasure about data ethics.

- The right to restrict processing

- The right to data portability

- The right to object specific risks to deanonymisation

- Rights in relation to automated decision making

What are the  
anonymisation

# Who can beat risk anonymisation?

## Individuals Organisation

It's the dangers to individuals that we are going to focus on. It should primarily be a process of

There are risks to individuals of not being able to access their data. It's a risk to the individual's ability to access their data. It's a risk to the individual's ability to access their data.

- Organisation of the data is difficult, especially

Let's focus on the individual's data for the use of data protection risk <sup>12</sup> designed for organisations

# Specific risks of de-anonymisation

First and foremost, the techniques used to de-anonymise data are often self-

- Information about the publication of the data is often not disclosed, but more about the data itself. Individual assessment, or an assessment of the sensitive information domain.

There is a lack of information in the domain that:

- Individual assessment of the data is not possible
- Attack and/or injury
- Individual assessment of the data is not possible



# Specific risks of deanonymisation

Sensitive information • Although the health care data third-party organisation has lost, the information is not lost. The organisation can still use the data to improve its services and to provide better care.

- Employers could potentially use the information in employment decisions.

Remember that the GDPR is a law, not a goal. It is not a goal to have perfect data protection. It is a goal to have a system that is secure and that can be trusted.

# Specific risks of de-anonymisation

These are the "protected characteristics" defined in the Equality Act 2010 against which groups.

- Age
- Disability
- Gender reassignment
- Marriage and Civil Partnership
- Pregnancy and Maternity
- Race
- Religion
- Sex
- Sexual Orientation.

This prejudice can be actions, from hate groups, to institutional policies

De-anonymisation of health data is not only a breach of confidentiality but also a violation of individual rights

All individual data is private, but health data is the de-anonymisation of health care data.



# A potter's history anonymist's story

# Early evidence for d

In the late 90s there was a rapid conceptual large, public datasets.

Latanya Sweeney<sup>14</sup> showed using public data and voting list (n = 54, 805) as syntique identifiers only 2 pieces of information

birth date alone	12%
birth date and gender	29%
birth date and 5-digit ZIP code	69%
birth date and full postal code	97%

**Table 3. Uniqueness of Demographic Fields in Cambridge, Massachusetts, Voter List.**

Source: Latanya<sup>14</sup> Sweeney

# Early evidence for d

In the late 90s there was a rapid conceptual large, public datasets.

Three years later in 2000 we demonstrated

that it is possible to identify a person from. Here are some surprising results using only three fields of information, even though typical data releases contain many more fields. It was found that 87% (216 million of 248 million) of the population in the United States had reported characteristics that likely made them unique based only on {5-digit ZIP, gender, date of birth}. About half of the U.S. population (132 million of 248 million or 53%) are likely to be uniquely identified by only {place, gender, date of birth}, where place is basically the city, town, or municipality in which the person resides. And even at the county level, {county, gender, date of birth} are likely to uniquely identify 18% of the U.S. population. In general, few characteristics are needed to uniquely identify a person.

87% of the US population can be uniquely identified with only ZIP, gender, and date of birth. Source: Latanya Sweeney<sup>15</sup>

# k - a n o n y m i t y : A M o d e l P r i v a c y

# k - a n o n y m i f i n y : t i A o n d ( l )

k - a n o n y m i t y i s a p r o p e r t y o f a n a n o n y m i t y a b e t m a t h e r  
h a s b e e n s u b j e c t t o a n a n o n y m i t y s k a - t a i n o n n . y m i t y i s  
p e o p l e d o n ' t d e c l a r e  
c a l c u l a t i n g k v a l u e s .

k - a n o n y m i t y e g s r a t h a l t u e  
g u a r a n t e e s u m a l q u e n e s s o f  
i n d i v i d u a k s - a n a n o n y m i t y d u a l s .



# k - a n o n y m i f i c a t i o n ( I I )

Let's consider a simple pre

name	regional_disease
Saini	Engl, 20-30 Heart
Enio	Engl, 20-30 Heart
Dauri	Engl, 20-30 Heart
Alphur	Engl, 20-30 Pancre
Baliar	Engl, 20-30 Pancre
Kenye	Engl, 20-30 Pancre
Graci	Wale, 40-50 Liver
Aliye	Wale, 40-50 Liver
Kadin	Wale, 40-50 Liver
Asaph	Wale, 40-50 Liver

We consider each column identifier

These attributes can be categorized into two types of identifiers:

- Unique identifiers

- These attributes that identify individuals are removed for anonymization.

- Quasi-identifiers

- These attributes could be used to identify individuals, even after anonymization.

# k - a n o n y m i f i c a t i o n a n d e t h i c s

Now we've thrown away the utility of the data. We need to decide whether it's worth attributing sensitivity to the data.

Sensitive attributes are medical data that we need to protect in an anonymisation process.

Non-sensitive	Sensitive
Age	Diagnosis
Engl 20-30	Heart
Engl 20-30	Heart
Engl 20-30	Heart
Engl 20-30	Pancr
Engl 20-30	Pancr
Engl 20-30	Pancr
Wale:40-50	Liver
Wale:40-50	Liver
Wale:40-50	Liver
Wale:40-50	Liver

the healthcare

# k - a n o n y m i f i c a t i o n ( I V )

There are often cases about how we calculate k - anonymous data.

- Combining together all non attributes each is a separate attribute.
- Combining together all attributes
- For each individual attribute

Let's go through each of these

Non - sens	Sensi
regi age _ r a d i sea	
Engl 20 - 30	Heart
Engl 20 - 30	Heart
Engl 20 - 30	Heart
Engl 20 - 30	Pancr
Engl 20 - 30	Pancr
Engl 20 - 30	Pancr
Wale:40 - 50	Liver
Wale:40 - 50	Liver
Wale:40 - 50	Liver
Wale:40 - 50	Liver

# k - a n o n y m i f i c a t i o n ( V )

Combining together all attributes compared to sensitive attributes.

In toy example like this manually any count how many belong to each group.

Using this measure the data anonymous as all individuals guaranteed anonymity among (like)1.

disease:age_race_income			
Heart	English	20-30	div 6
Heart	Engl	20-30	6
Heart	Engl	20-30	6
Pancr	Engl	20-30	6
Pancr	Engl	20-30	6
Pancr	Engl	20-30	6
Liver	Wale	40-50	4
Liver	Wale	40-50	4
Liver	Wale	40-50	4
Liver	Wale	40-50	4

# k - a n o n y m i t y : t i A n d ( V I )

di sea:re gi age \_r in \_i n \_g

Heart	Engl	20-30	3
Heart	Engl	20-30	3
Heart	Engl	20-30	3
Pancr	Engl	20-30	3
Pancr	Engl	20-30	3
Pancr	Engl	20-30	3
Liver	Wale	40-50	4
Liver	Wale	40-50	4
Liver	Wale	40-50	4
Liver	Wale	40-50	4

Combining together all

When we measure target

the k-anonymity of the data

Using this metric, the data

anonymity

# k - a n o n y m i t y : t i A n d V I I

For each individual attribute

When we measure the anonymous individual variable the data anonymity

We always measure the difference of our difference.

disease:regi_age_r		
Heart	Engl	20 - 30
Heart	Engl	20 - 30
Heart	Engl	20 - 30
Pancr	Engl	20 - 30
Pancr	Engl	20 - 30
Pancr	Engl	20 - 30
Liver	Wale	40 - 50
Liver	Wale	40 - 50
Liver	Wale	40 - 50
Liver	Wale	40 - 50

# known ymi fny:ti And Vll

As we've seen, each of these give separate measures of the data.

1. Combining together all non-attributes ~~is~~ <sup>can</sup> be seen as a combination of attributes.

2. Combining together all at

3. For each individual attribute

Frustratingly rare for authors to explicitly state combination of attributes they use.

The methods are listed roughly in terms of the frequency that I've seen them in the literature.

Non-sens	Sensi
regi age_radi sea	
Engl 20-30	Heart
Engl 20-30	Heart
Engl 20-30	Heart
Engl 20-30	Pancr
Engl 20-30	Pancr
Engl 20-30	Pancr
Wale:40-50	Liver
Wale:40-50	Liver
Wale:40-50	Liver
Wale:40-50	Liver

# k - a n o n y m i f i c a t i o n : t i a n d e ( l x )

The first diet ion l' ve fig ii ø ð n  
a n d p r e c i s e .

B u t b e a w a r e t h a t e y a o u ' <sup>16</sup>  
m o r e t e c h n i c a l n o t a t i o n .

## Lemma.

Let  $RT(A_1, \dots, A_n)$  be a table,  $Q_{RT} = (A_1, \dots, A_n)$  be the quasi-identifier associated with  $RT$ ,  $A_1, \dots, A_n \subseteq A_1, \dots, A_n$ , and  $RT$  satisfy  $k$ -anonymity. Then, each sequence of values in  $RT[A_x]$  appears with at least  $k$  occurrences in  $RT[Q_{RT}]$  for  $x=1, \dots, n$ .

## Example 4. $k$ occurrences of each value under $k$ -anonymity

Table  $T$  in Figure 2 adheres to  $k$ -anonymity, where  $Q_T = \{Race, Birth, Gender, ZIP\}$  and  $k=2$ . Therefore, each value that appears in a value associated with an attribute of  $Q_T$  in  $T$  appears at least  $k$  times.  $|T[Race = "black"]| = 6$ ,  $|T[Race = "white"]| = 5$ ,  $|T[Birth = "1964"]| = 5$ ,  $|T[Birth = "1965"]| = 4$ ,  $|T[Birth = "1967"]| = 2$ ,  $|T[Gender = "m"]| = 6$ ,  $|T[Gender = "f"]| = 5$ ,  $|T[ZIP = "0213*"]| = 9$ . And,  $|T[ZIP = "0214*"]| = 2$ .

Source: Formal definition for  $k$ -anonymity from Sweeney 2002<sup>16</sup>



# k - a n o n y m i t y :   H o w   i s

We are responsible for manipulating our data

Even though a minimum value of 3 is often  
in practice is to ensure it is not too high  
see ravraet iaot

We have tools available to us:



# Task: Set up our project

## SLIDE 1 OF 1

1. Create a new `webpack` project for
2. Create a new `React` `ReactDOM` `ReactDOM` `ReactDOM`
3. Install `webpack` `webpack` `webpack`

# { wakefield }

The wakefield package is very useful for creating variables.

The package provides a first class interface to the

age	dice	hair	mltsex_income
animaldna	height	months	smokes
answerdob	income	name	speed
area	dummy	internet	normstate
car	educatiq	polittstring	
childremploylanguage	race	upper	
coin	eye	level	religvalid
color	grade	likert	sat year
date_sgrade_lorem_ipsum	ez	ip	coc
death	group	marital	sex

# walkthrough: random data frame

We generate data as follows:

```
1 r_data_frame(10,  
2   id,  
3   name,  
4   dob,  
5   income,  
6   smokes,  
7   death)
```

```
# A tibble: 10 × 6  
  ID   Name      DOB      Income Smokes Death  
  <chr> <chr> <date> <dbl> <lg> <lg>  
1 01 Dreniyah 2008-05-24 42462. TRUE FALSE  
2 02 Alexiana 2008-11-06 60950. FALSE TRUE  
3 03 Ashauria 2008-10-16 17069. FALSE TRUE  
4 04 Azelea 2007-12-11 19028. FALSE FALSE  
5 05 Krystof 2009-09-11 49067. TRUE FALSE  
6 06 Loudes 2008-10-24 65993. FALSE FALSE  
7 07 Rebeckah 2008-07-21 50860. FALSE FALSE  
8 08 Jarek 2009-01-26 21266. FALSE FALSE  
9 09 Gavrielle 2008-08-03 50405. FALSE TRUE  
10 10 Tarajah 2008-05-07 28798. FALSE FALSE
```

Can you explain why our machine

# Pseudorandomness

When programming we use pseudorandom numbers.

These are algorithms that take a seed as input and produce a sequence of numbers. We can then use these numbers for many different purposes.

```
1 set.seed(1)
2 r_data_frame(10,
3   id,
4   name,
5   income,
6   dna,
7   smokes,
8   death)
```

```
# A tibble: 10 × 6
   ID Name      Income DNA      Smokes Death
  <chr> <chr>      <dbl> <fct>      <lgl> <lgl>
1 01 Donaldeen 48108. Cytosine FALSE FALSE
2 02 Martiqua 36496. Guanine FALSE FALSE
3 03 Juliaann 12130. Thymine FALSE TRUE
4 04 Poyraz 53488. Cytosine FALSE TRUE
5 05 Bolslaus 46311. Adenine FALSE FALSE
6 06 Duc 56733. Adenine FALSE FALSE
7 07 Hadeer 52217. Thymine TRUE TRUE
8 08 Camliya 31854. Thymine FALSE TRUE
9 09 Ashlay 77153. Adenine TRUE TRUE
10 10 Dutch 62490. Adenine FALSE FALSE
```

# k - a n o n y m i t y f o r o u r

Lest 'pretend our dataset for information doesn't exist

When thinking about making anonymous ...

- Are there any variables suppressed
- How could we generalise variables?

LCName	DOB	IncOSmok	Deat
01Donal	c2009-	C31854	FALSIFALS
02Marti	c2009-	C77153	FALSIFALS
03Julia	c2008-	1162490	TRUETRUE
04Poyra	c2009-	C15462	TRUEFALS
05Bolesl	c2008-	C19430	FALSITRUE
06Duc	2008-	C24504	FALSIFALS
07Hadeer	2009-	C27535	TRUETRUE
08Camill	c2008-	116200	FALSIFALS
09Ashla	c2008-	C5750	FALSIFALS
1CDutch	2008-	C59224	FALSITRUE

[1] I'm not interested by the optimal solution by {wakefel} There is NO correlation between the two variables. You can tell because I don't have a citation



# Task: k-anonymity case

## SLIDE 1 OF 2

1. Use this code to create a dataset:

```
1 library(wakefield)
2 library(tidyverse)
3 library(lubridate)
4
5 set.seed(1)
6 smoke_data <- r_data_frame(
7   50000,
8   id,
9   name,
10  dob(start = ymd("1950-01-01"),
11     k = abs(as.integer(days(ymd("1950-01-01") - Sys.Date())) - 365*18))),
12  income,
13  smokes,
14  death)
```

2. Suppress inappropriate columns from the



# Task: k-anonymity

## SLIDE 1 OF 2

1. Generalise the remaining variables as follows

- Extract year of birth
- Split income into 4 categories

- " $< \text{£}30,000$ "
- " $\text{£}30,000 - \text{£}70,000$ "
- " $\text{£}70,000 - \text{£}100,000$ "
- " $\text{£}100,000 +$ "

2. Calculate the k-anonymity of the dataset



# How to attack - anon

There are 3 known attacks for attempt to de

- Unsorted matching attack
  - If an anonymised dataset is ordered in t  
recorded this is a. pDoetpeenntdi eanl t aotnt atchke vaetctta  
either an addi ffeiro noarl (qwuoarnstt- icdaesnet)i a uni qu  
▪ It's a good practice to randomise the orde

# How to attack - anon (III)

There are 3 known attacks for attempt to de

- Unsorted mat
- Subsequent release attacks
  - Large healthcare datasets might be used multiple k-anonymised releases.
  - Temporal attacks are possible by fam always i( ntervention) or the removal of data (eg
  - Protecting against these attacks require consider the k-anonymity of combined rel

# How to attack - anon (IIII)

There are 3 known attacks for attempt to de

- Unsorted matc
- Subsequent re

# I - diversity

# I - diversity (I)

I - diversity is a more measure of the anonymity variables in an anonymous introduced by Machana

2006<sup>21</sup>

This method depends on model of background kn

We are going to cover Bayesian statistics. It is a bit different from the one you're interested in.

assignments  $\psi$  compatible with the background knowledge such that  $\psi(X) = s$  can be calculated as follows.  $X$  is assigned

the **sensitive** value  $s$ . Since  $X(Q) = q$ , out of the remaining  $N_q - 1$  individuals having the **non-sensitive** value  $q$ ,  $N_{(q,s)} - 1$  of them are assigned  $s$ . For every other **sensitive** value  $s'$ ,  $N_{(q,s')}$  out of the  $N_q - 1$  individuals are assigned  $s'$ . For every  $q' \neq q$  and every  $s'$ , some  $N_{(q',s')}$  out of the  $N_{q'}$  individuals having the **non-sensitive** value  $q'$  are assigned  $s'$ . The number of these assignments is

$$\begin{aligned} & \frac{(N_q - 1)!}{(N_{(q,s)} - 1)!} \prod_{s' \neq s} \frac{N_{(q,s')!}}{\prod_{q' \neq q} N_{(q',s')!}} \\ &= \frac{N_{(q,s)}}{N_q} \prod_{q' \in Q} \frac{N_{q'}!}{\prod_{s' \in S} N_{(q',s')!}} \end{aligned} \quad (2)$$

For each mapping  $\psi$  such that  $\psi(X) = s$ , we count the number of  $Z_n$ 's such that  $(\psi, Z_n) \vdash (T^*, X)$  as follows. Let  $q^*$  be the generalized value of  $q = X(Q)$ .  $X$ 's record will appear as  $t_X^* = (q^*, s)$  in the table  $T^*$ . Apart from  $t_X^*$ ,  $T^*$  contains  $n_{(q^*,s)} - 1$  other tuples of the form  $(q^*, s)$ . Hence, apart from  $X$ ,  $Z_n$  should contain  $n_{(q^*,s)} - 1$  other individuals  $\omega$  with  $\psi(\omega) = s$  and  $\omega(Q) = q^*$  where  $q^*$  generalizes to  $q^*$ . For all other  $(q^{*,s'})$  such that  $q^* \neq q^*$  or  $s' \neq s$ ,  $Z_n$  should contain  $n_{(q^{*,s'})}$  individuals  $\omega'$  where  $\psi(\omega') = s'$  and  $q^{*}$  is the generalized value of  $\omega(Q)$ . The number of  $Z_n$ 's is given by

$$\begin{aligned} & \left( \frac{N_{(q^*,s)} - 1}{n_{(q^*,s)} - 1} \right) \prod_{(q^{*,s'}) \in (Q^* \times S) \setminus \{(q^*,s)\}} \prod_{(q^{*,s'})} \left( \frac{N_{(q^{*,s'})}}{n_{(q^{*,s'})}} \right) \\ &= \frac{n_{q^*,s}}{N_{(q^*,s)}} \prod_{(q^{*,s'}) \in Q^* \times S} \left( \frac{N_{(q^{*,s'})}}{n_{(q^{*,s'})}} \right) \end{aligned} \quad (3)$$

The cardinality of  $T_{(X,s)}^*$  is therefore the product of Equations 2 and 3 and can be expressed as

$$\begin{aligned} |T_{(X,s)}^*| &= \frac{N_{(q,s)}}{N_q} \prod_{q' \in Q} \frac{N_{q'}!}{\prod_{s' \in S} N_{(q',s')!}} \times \frac{n_{q^*,s}}{N_{(q^*,s)}} \prod_{(q^{*,s'}) \in Q^* \times S} \left( \frac{N_{(q^{*,s'})}}{n_{(q^{*,s'})}} \right) \\ &= \frac{N_{(q,s)}}{n_{(q^*,s)} N_{(q^*,s)}} \times \frac{1}{N_q} \prod_{q' \in Q} \frac{N_{q'}!}{\prod_{s' \in S} N_{(q',s')!}} \times \prod_{(q^{*,s'}) \in Q^* \times S} \left( \frac{N_{(q^{*,s'})}}{n_{(q^{*,s'})}} \right) \\ &= \frac{N_{(q,s)}}{n_{(q^*,s)} N_{(q^*,s)}} \times \varepsilon \end{aligned}$$

The expression  $\varepsilon$  is the same for all  $s' \in S$ . Hence, the expression for the observed belief is

$$\begin{aligned} \beta_{(q,s),T^*} &= \frac{|T_{(X,s)}^*|}{\sum_{s' \in S} |T_{(X,s')}^*|} \\ &= \frac{n_{(q,s)} \frac{N_{(q,s)}}{n_{(q^*,s)} N_{(q^*,s)}}}{\sum_{s' \in S} n_{(q^*,s')} \frac{N_{(q,s')}}{n_{(q^*,s')} N_{(q^*,s')}}} \end{aligned}$$

Using the substitutions  $f(q,s) = N_{(q,s)}/N$  and  $f(q^*,s) = N_{(q^*,s)}/N$ , we get the required expression.

$$\begin{aligned} \beta_{(q,s),T^*} &= \frac{n_{(q^*,s)} \frac{f(q,s)}{f(q^*,s)}}{\sum_{s' \in S} n_{(q^*,s')} \frac{f(q,s')}{f(q^*,s')}} \\ &= \frac{n_{(q^*,s)} \frac{f(s|q)}{f(s|q^*)}}{\sum_{s' \in S} n_{(q^*,s')} \frac{f(s'|q^*)}{f(s'|q^*)}} \end{aligned}$$

Note that in the special case when  $S$  and  $Q$  are independent, The expression for the observed belief simplifies to



# I - diversity (II)

In order to estimate I-diversity once again consider the standard dataset.

We split our data into:

- Sensitive variables - the variables.
- Keys - the non-sensitive

A dataset is I-diverse if combination of key attributes leads "well-represented" values for each sensitive variable.

Non-sensitive Keys		Sensitive diseases
region_age_range		
England20-30	Heart	Heart
England20-30	Heart	
England20-30	Heart	
England20-30	Pancreas	Pancreas
England20-30	Pancreas	
England20-30	Pancreas	
Wales40-50	Liver	Liver
Wales40-50	Liver	
Wales40-50	Liver	
Wales40-50	Liver	Liver
Wales40-50	Liver	
Wales40-50	Liver	

# I - diversity (III)





# Task: Calculating di

## SLIDE 1 OF 2

1. Install the `tidyverse` package
2. Add this dataset to your

```
1 data_di_seases <- tibble(  
2   name = c("Sai ndhavi", "Eni o", "Daur y", "Al phus", "Bal ian",  
3           "Kenyea", "Graci el ynn", "Al i ye", "Kadi nce", "Asaph"),  
4   region = c(rep("England", 6), rep("Wales", 4)),  
5   age_range = c(rep("20-30", 6), rep("40-50", 4)),  
6   di_sease = c(rep("Heart", 3), rep("Pancreatic", 3), rep("Liver", 4))  
7 )  
8 data_di_seases
```

```
# A tibble: 10 x 4  
  name      region age_range di_sease  
   <chr>    <chr>    <chr>    <chr>  
1 Sai ndhavi  England 20-30    Heart  
2 Eni o      England 20-30    Heart  
3 Daur y     England 20-30    Heart  
4 Al phus    England 20-30    Pancreatic  
5 Bal ian    England 20-30    Pancreatic  
6 Kenyea     England 20-30    Pancreatic  
7 Graci el ynn Wales    40-50    Liver  
8 Al i ye    Wales    40-50    Liver  
9 Kadi nce   Wales    40-50    Liver  
10 Asaph     Wales    40-50    Liver
```



# Task: Calculate l - di

## SLIDE 2 OF 2

1. Compute diversity of the dataset

```
1 l_diseases <- data_diseases %>%  
2 mutate(di_sease = as_factor(di_sease)) %>%  
3 createSdcObj(keyVars = c("region", "age_range")) %>%  
4 ldiversity(ldiv_index = "di_sease")
```

1. Extract the l-diversity value

```
1 l_diseases@risk$ldiversity
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.0	1.0	2.0	1.6	2.0	2.0

# Summarising k - a n o diversity

# Summarizing $k$ -anonymity diversity

$k$ -anonymity has well known limitations for high-dimensional datasets [20]. Furthermore,  $k$ -anonymization completely fails on high-dimensional datasets [2], such as the Netflix Prize dataset and most real-world datasets of individual recommendations and purchases.

at least 3-dimensional Source: Narayan and datasets containing sensitive attributes.

l-diversity is a much more sophisticated tool that provides stronger privacy guarantees for top-level recommendations in a diversified manner. It is designed to protect sensitive information by ensuring that the data is not only  $k$ -anonymous but also  $l$ -diverse. This means that for any group of  $k$  records, there must be at least  $l$  distinct values for the sensitive attribute. This helps to prevent the identification of individuals in the dataset.

However, diversity does not guarantee against the coarsening of data might also degrade the usability of released data.

Case Study: Patient

# Netflix Prize Dataset:

In October 2006, Netflix competition with the their recommendation algorithm that would improve the algorithm.

Let's get into this a little bit more.

use movie\_date\_of\_g

13 Ali 2006-01-5

13 Ali 2006-01-4

It will be over 3 million most-recent ratings from those same subscribers over the same set of movies as a competition equal to the

Netflix provided over 100 million ratings (and their dates) from over 480 thousand random anonymous subscribers on nearly 8 thousand movies. The data were collected between October, 1998 and December, 2005 and reflect a massive and ongoing effort received by Netflix during this period. The ratings are on a scale from 1 to 5 (integral stars).

# Netflix Prize Dataset:

Netflix had an algorithm which attempted to predict movie ratings based on the subset of movie ratings.

The intention was to provide more personalized recommendations that average ratings for users.

Netflix chose to use the squared error (RMSE) and all other user ratings as accuracy

The main goal was 10% more accurate than the user rating

The competition challenged further improve this an additional 10%.

# Netflix Prize Dataset:

The competition was and it wasn't won!

In fact, it's a dramatic winning entries, one another in 200,

I'd recommend [this article](#) by [Dan Jacobson](#).

However

16 days after [Narayan](#) [Vitaly Shmatikov](#) demonstrated to re-identify [users](#)

This draft [play script](#) is from 2008 and we'll look in [this definition](#) was possible



# Netflix Prize Dataset:

With the announcement of a 2nd competition in 2009 a class-action lawsuit was filed.

The suit described the as the biggest "voluntary" data breach in history (FTC) also got involved

Unfortunatly the case was settled privately so we don't know the damages. The best I can get is this quote from a deleted<sup>25</sup> blog post

---

In the past few months, the Federal Trade Commission (FTC) asked us how a Netflix Prize sequel might affect Netflix members' privacy, and a lawsuit was filed by Kamber Law LLC pertaining to the sequel. While both the FTC and the plaintiffs' lawyers, we've had very productive discussions centered on our commitment to protect our members' privacy. We have reached an understanding with the FTC and

have settled the lawsuit with Netflix. The resolution to both matters involves certain parameters for how we use Netflix data in any future research programs.

In light of all this, we have decided to not pursue the Netflix Request sequel that we announced on August 6, 2009." - Neil Hunt, Chief Product Officer @ Netflix

# Netflix Prize Data Set: Overview (I)

The exact mechanics of the algorithm behind Arvind Narayanan's attack on a large dataset is beyond the scope of this course. We will discuss the attack in more detail in the next lecture.

We're going to walk through the original paper, which is the source of the attack. We'll be using the dataset from the Netflix Prize competition, which is a large dataset of movie ratings.

It's important to identify this paper in the context of statistical de-anonymization of large sparse datasets. That is not what we're doing here.

# Netflix Prize Dataset: Overview (III)

The de-anonymisation of the Netflix Prize dataset has led to a large number of new data sources: movie ratings.

- IMDb Ratings

- Date of movie release

An assumed similarity between IMDb ratings profiles is a natural attack vector against the dataset.

## Netflix Prize

The paper provides a detailed analysis of the IMDb dataset.

Despite this, the dataset is positively cross-matched between the datasets.

# Netflix Prize Data Set: O ( I I )

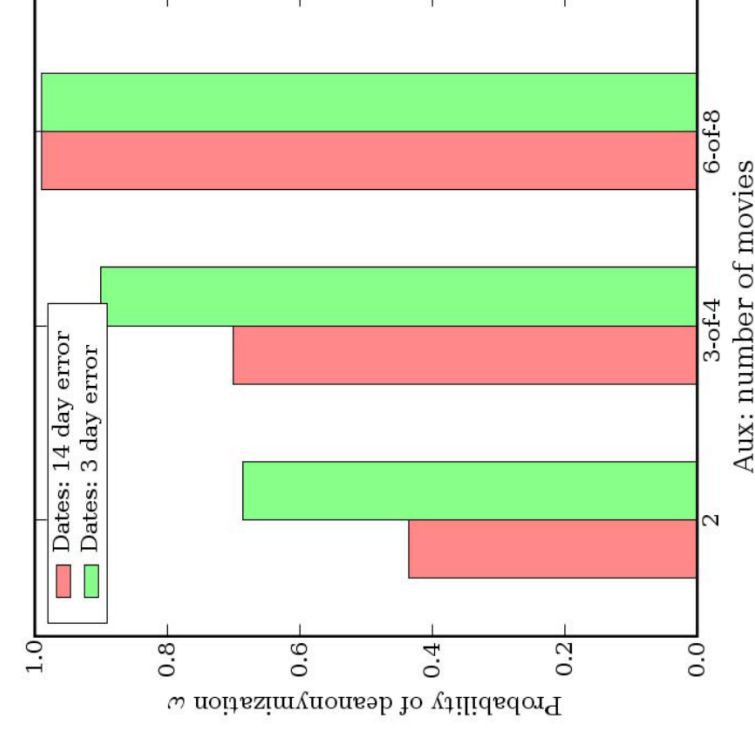
Netflix x Prize Data set: O  
( I I )

The de-anonymization attack we readt in g by two data sources: additional attack vec

- IMDb R:

- Date of movie review account creation date

The authors **heavily** use **AUC<sub>F</sub>**. F



# Selected other ca





# Case Study: Facebook

Facebook has a long a nBde aacvofnu lwahsi sdteosri yg noefd t privacy breaches and gpuoesstt i pounracbhl aes e s F a c y b o o activity feeds.

Facebook Beacon is on eT hoof ctohnep aonl yd eosrti gi n a l l examples, all the way sbearcvki cfer owra s2 0" 0o7p.t - i n " , evidence this was not

It's a rare example where Mark Zuckerberg open talks Hæbvœuyte aiftt ears car i t i c i s m mi s t<sup>27</sup>ake provided an opt-out. to t f u r t h e f s e r v i c e .

Discuss why Facebook Beacon  
breaches user privacy

A class action suit i  
\$9.5 million.

# Case Study: Google B

Google Buzz was a v Google automa public  
networking tool: Google P rpages using  
service.

- Launched: February <sup>28</sup>
- Discontinued: Dec <sup>29</sup>ce

The pages disclo  
user most frequentl  
via email or chat w

It shut down explic  
violations and Goog  
\$8.5m <sup>30</sup> will t h in n on e f m d  
service launching.

This was as designe  
Google designed thi

When you first enter Google Buzz, to make the startup experi re ease may automati at l se l et  
peopl for you to fol w based on the peopl you emai and chat wi th most Si in H ay we may al suggest  
to others that they automati at l fol w you. You can revi w and edi the l to f peopl you fol w and bl ck  
peopl from fol wi ng you.

peopl fol wi ng you and who you' er fol wi ng on your profil .e

. You may opt out of di p r ying the l to f

# Case Study: In-store

Do you know what these machines are for?



# Case Study: In-store

These devices have embedded sensors that can track customer behavior and their attention.

The video feed from the cameras is processed by software algorithms to estimate

- Gender
- Age
- Mood

Their preferences in the store are used to make targeted marketing.



"Quividi software employs advanced facial detection algorithms to recognize individual visitors marketing

# Anonymous notation

# Anonymisation software

Table 5. Comparison of the off-the-shelf privacy model-based data anonymization tools in terms of available development options, anonymization functionality and risk metrics.

Tool	Development support				Anonymiza- tion	Risk as- sessment
	Last release	Open source	Public API <sup>a</sup>	Extensibility	Cross-plat- form	Programming language
AEX	November 2019	✓ <sup>b</sup>	✓	✓	✓	Java
Amnesia	October 2019	✓	✓	✓	✓	Java
µ-ANT <sup>c</sup>	August 2019	✓	✓	✓	✓	Java
Anonimatron	August 2019	✓	✓	✓	✓	Java
SECRET <sup>a,d</sup>	June 2019					C++
sdMicro	May 2019	✓	✓	Poorly support- ed	✓	R
Airleak Insights	April 2019					Ruby
NLM <sup>e</sup> Scribber	April 2019			✓	✓	Perl
Anonymizer	March 2019	✓	✓	✓	✓	Ruby
Shiny Anonymizer	February 2019	✓	✓	✓	✓	R
µ-ARGUS	March 2018					C++
UTD <sup>f</sup> Toolbox	April 2010	✓		Poorly support- ed	✓	Java
OpenPseudonymizer	November 2011	✓		✓	✓	Java
TIAMAT <sup>g</sup>	2009			✓	✓	Java
Cornell Toolkit	2009	✓		Poorly support- ed	✓	C++

In 2021 <sup>32</sup> we refer to a system de-anonymisation tools in

There we felt we should build with R

- { s d c M i c r o }

- This package contains a measuring/exploring the via a {shiny} app.

- { ShinyAnonymizer }

- This package provides a anonymising healthcare d

# Anonymisation software

However, there is no fit & all but we sneezed to be vigilant and methodology or guarantee of privacy through anonymisation.

- Preparing data for release
- k-anonymity and l-diversity
- Privacy metrics and provide some assurance of privacy
- But background knowledge attacks might undermine these.



# Simulating fake data

# Simulating fake data



We've used `faker` package to  
fake datasets.

The `faker` package is a more  
package for simulating data  
`DeBruijne`

# Useful resources

# Useful resources (I)

Classes start on Monday August 26, 2019 (CMU Academic Calendar)

Office hours and Hangouts can be found in the [Google Calendar](#)

All homework is due 10 minutes before lecture/recitation start.

Date	Topic	Reading	Notes
Part 0: Introduction			
Week 1			
Monday August 26	Course Overview	<ul style="list-style-type: none"><li>• <a href="#">CMU Computing Policy</a></li><li>• <a href="#">CMU Policy on Academic Integrity</a></li></ul>	<a href="#">slides</a>
Wednesday August 28	Conceptual Framework for Understanding Privacy	<ul style="list-style-type: none"><li>• <a href="#">Fair Information Principles</a></li><li>• (optional) <a href="#">Privacy in Context</a></li><li>• (optional) <a href="#">Overview Article in Stanford Encyclopedia of Philosophy</a></li><li>• (optional) <a href="#">Privacy as Contextual Integrity</a></li><li>• (optional) <a href="#">A Taxonomy of Privacy</a></li></ul>	<a href="#">slides</a>
Friday August 30	No recitation this week		
Homework 1 out			
Part I: Privacy through Accountability: Formalization and Detection			
Week 2			
Monday September 2	No Class: Labor Day		
Wednesday September 4	Enforcing Purpose Restrictions through Audit	<ul style="list-style-type: none"><li>• <a href="#">Formalizing and Enforcing Purpose Restrictions in Privacy Policies</a></li><li>• <a href="#">Purpose Restrictions on Information Use</a></li><li>• <a href="#">Privacy and Contextual Integrity: Framework and Applications</a></li><li>• <a href="#">Summary of the HIPAA Privacy Rule</a> (<a href="#">Permitted Uses and Disclosures</a>, <a href="#">Authorized Uses and Disclosures</a>)</li><li>• <a href="#">A Formalization of HIPAA for a Medical Messaging System</a></li><li>• <a href="#">Experiences in the Logical Specification of the HIPAA and GLBA Privacy Laws</a></li></ul>	<a href="#">slides</a> <a href="#">Notes on MDPs</a>
Friday September 6	Recitation on Docker (Sruti)		

This is a real world tutorial on privacy from Carnegie University

You find all lectures and even exercises <https://course.ece>

This quickly becomes technical.

# Useful resources (I)

## Privacy in a Mobile-Social World

*CompSci 590.03*

*Instructor: Ashwin Machanavajjhala*

Ashwin Machanavajjhala  
excellent course a  
mobile world -  
<https://courses.cs.duke.edu/590.03/>

This quickly becomes  
technical.

Lecture 1 : 590.03 Fall 13



# Assessment

# Assessment

In this section you must explain:

- What is an open dataset?
- What is "health data"?
- Why is it important to government, industry available?
- What are some ethical considerations and groups
- What is data anonymisation and why is it important?
- What are the dangers to individuals and groups de-anonymised?
- What are some steps that can be taken to reduce

---

In answering these questions, you must include as much about data de-anonymisation. Include as much as you can about how the data was de-anonymised.

# References

1. Kretschmer, M., Pennekamp, J. & Wehr, K. **Cooki Banners and Privacy Policy: Measuring the Impact of the GDPR on the Web**. ACM Transactions on the Web 15, 1–42 (2021).
2. European Union Regulation (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL. (2016).
3. Cookies and the GDPR, and the ePrivacy Directive & GDPR.eu (2019).
4. Engelhardt, S. & Narayanan, A. Online Tracking: A 1-milestone Measurement and Analysis. Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, 1388–1401 (Association for Computing Machinery, 2016). doi [10.1145/2976749.2978313](https://doi.org/10.1145/2976749.2978313).
5. Habi, H., Li, M., Young, E. & Cranor, L. "Okay, whatever": An Evaluation of Cookie Consent Interfaces. In Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems 1–27 (Association for Computing Machinery, 2022). doi [10.1145/3491102.3501985](https://doi.org/10.1145/3491102.3501985).
6. Spiegel, J. R., McKenna, M. T., Lakshman, G. S. & Nordstrom, P. G. Method and system for anti-phishing package sharing (2013).
7. Nyckel, E.-M. Ahead of Time: The Infrastructure of Amazon's Anti-phishing System. In Medical Infrastructures and the Politics of Digital Time (eds. Volpert, A. & Stieglitz, K.) 263–278 (Amsterdam University Press, 2021). doi [10.2307/j.ctv1xcxr3n.18](https://doi.org/10.2307/j.ctv1xcxr3n.18).
8. Vintner, N. Q., Behdani, B. & Blumenthal, J. **Data-driven process redesign: Anti-phishing supply chain**. International Journal of Production Research 58, 1302–1318 (2020).
9. Venkataraman, M. Madhumi at Venkataraman: My interview for sales. LinkedIn UK (2014).
10. UK Government National Health Service Act 2006. (2006).
11. UK Government Data Protection Act 2018. (2018).
12. Information Commissioner's Office. Anonymity and the Management of data protection risk code of practice (2012).
13. UK Government Equalities Act 2010. Information.gov.uk (2010).
14. Sweeney, L. **Weaving Technology and Policy Together to Make it Confidential**. The Journal of Law, Medicine & Ethics 25, 98–110 (1997).
15. Sweeney, L. **Siempre**. Demographic Software People. Bytes (2000) doi [10.1184/R1/6625769.V1](https://doi.org/10.1184/R1/6625769.V1).
16. Sweeney, L. **K-ANONYMITY: A MODEL FOR PROTECTING PRIVACY**. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 10, 557–570 (2002).
17. Sweeney, L. **Online You, Your Doctor, and Many Others May Know: Technology Science** (2015).



18. Samarati P. & Sweeney, L. Protecting Privacy when Disclosing Information in K-Anonymity and Its Enforcement through Generalization and Suppression (1998).
19. El Emam, K. & Dankar, F. K. **Protecting Privacy: Using Privacy in the Journal of the American Medical Association on Security and Privacy** (2008) 111–125 (2008). doi [10.1109/SP.2008.33](https://doi.org/10.1109/SP.2008.33).
20. Narayanan, A. & Shmatik, D. Robust De-anonymization of Large Sparse Datasets. in 2008 IEEE Symposium on Security and Privacy (sp 2008) 111–125 (2008). doi [10.1109/SP.2008.33](https://doi.org/10.1109/SP.2008.33).
21. Machanavajjhale, A., Gehrke, J., Kifer, D. & Venkatasubramanian, M. L-diversity: Privacy beyond k-anonymity. in 22nd International Conference on Data Engineering (ICDE' 06) 24–24 (2006). doi [10.1109/ICDE.2006.1](https://doi.org/10.1109/ICDE.2006.1).
22. Bennett, J. & Lanni, S. The Netflix Prize (2007).
23. Jackson, D. The Netflix Prize: How a \$1 Million Contest Changed Big Data. *Wired* (2013).
24. Narayanan, A. & Shmatik, D. V. How To Break Anonymity of the Netflix Prize Dataset (2006) doi <https://arxiv.org/abs/cs/0610105v1>.
25. Waxman, S. Whoops! Netflix Gets Caught by FTC, Cancels Contest (2010).
26. Narayanan, A. & Shmatik, D. V. "How to Break Anonymity of the Netflix Prize Dataset" - FAQ. (2007).
27. Zuckerberg, M. Our Commitment to the Facebook Community (2011).
28. Ho, E. Google Buzz is Shutting Down. Official Google Blog (2010).
29. Google. Blog A full sweep. Official Google Blog (2011).
30. BuzzClicks. Google Buzz User Privacy Policy (2010).
31. Qui videtur CONSUMER Privacy - Qui videtur Insignificant data with full anonymous measurements. Qui videtur (2022).
32. Zuo, Z. et al **Data Anonymity: A Review of Pervasive Healthcare: Systematic Literature Mapping Study**. JMIR Medical Informatics 9, e29871 (2021).