

**STA130H1F**

**Class #12**

**Prof. Nathan Taback**

**2018-12-03**

# Today's Class

- Conduct during final exams at UofT
- Finish confounding from last class
- Ethical Issues in Data Science Research

# A Classic Example: Treatment for kidney stones

Source of data: *British Medical Journal (Clinical Research Edition)* March 29, 1986

- Observations are patients being treated for kidney stones.
- treatment is one of 2 treatments (open or Invasive)
- outcome is success or failure of the treatment
- Doctors get to choose the treatment, depending on the patient
- What might influence how a doctor chooses a treatment for their patient?

# Kidney stone example

```
tab <- table(kidney_stones$outcome,  
             kidney_stones$treatment, deparse.level = 2)  
addmargins(prop.table(tab))
```

```
##                kidney_stones$treatment  
## kidney_stones$outcome    Invasive      Open      Sum  
##                failure 0.08714286 0.11000000 0.19714286  
##                success 0.41285714 0.39000000 0.80285714  
##                Sum      0.50000000 0.50000000 1.00000000
```

# Kidney stones come in various sizes

```
kidney_stones %>%  
  count(size, treatment, outcome) %>%  
  group_by(size, treatment) %>%  
  mutate(per_success = n / sum(n))
```

```
## # A tibble: 8 x 5  
## # Groups:   size, treatment [4]  
##   size treatment outcome      n per_success  
##   <chr> <chr>      <chr>   <int>      <dbl>  
## 1 large Invasive failure    25      0.312  
## 2 large Invasive success    55      0.688  
## 3 large Open      failure    71      0.270  
## 4 large Open      success   192      0.730  
## 5 small Invasive failure    36      0.133  
## 6 small Invasive success   234      0.867  
## 7 small Open      failure     6      0.0690  
## 8 small Open      success    81      0.931
```

Column percentages (conditional distribution of success given treatment):

```
prop.table(table(kidney_stones$outcome, kidney_stones$treatment), margin = 2)
```

```
##  
##           Invasive      Open  
##  failure 0.1742857 0.2200000  
##  success 0.8257143 0.7800000
```

```
large <- kidney_stones %>% filter(size == "large")  
prop.table(table(large$outcome, large$treatment), margin = 2)
```

```
##  
##           Invasive      Open  
##  failure 0.312500 0.269962  
##  success 0.687500 0.730038
```

```
small <- kidney_stones %>% filter(size == "small")  
prop.table(table(small$outcome, small$treatment), margin = 2)
```

```
##  
##           Invasive      Open  
##  failure 0.1333333 0.06896552  
##  success 0.8666667 0.93103448
```

*Which treatment is better?*

This example is another case of **Simpson's paradox**.

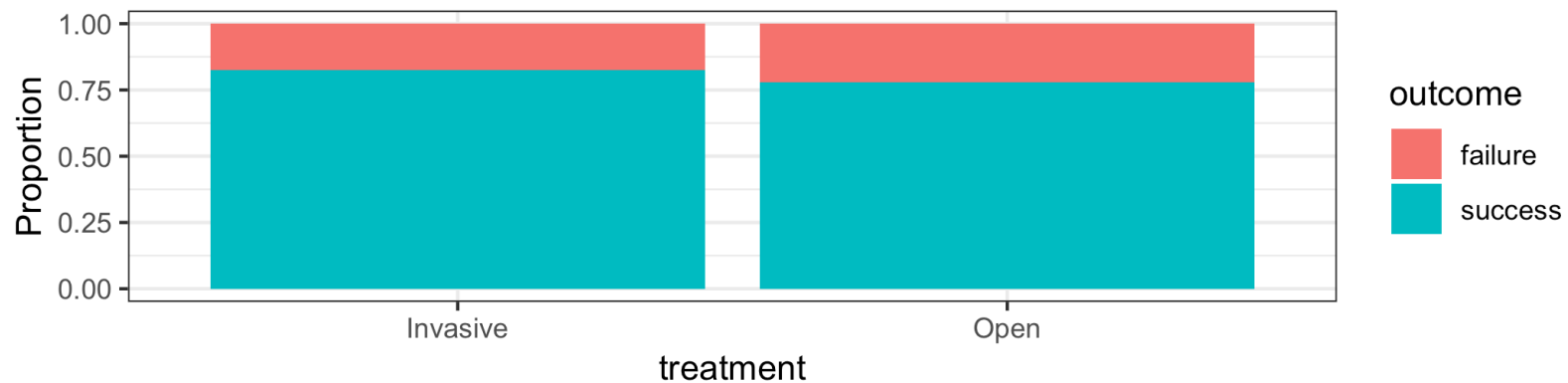
## Moral of the story:

Be careful drawing conclusions from data!

It's important to understand how the data were collected and what other factors might have an affect.

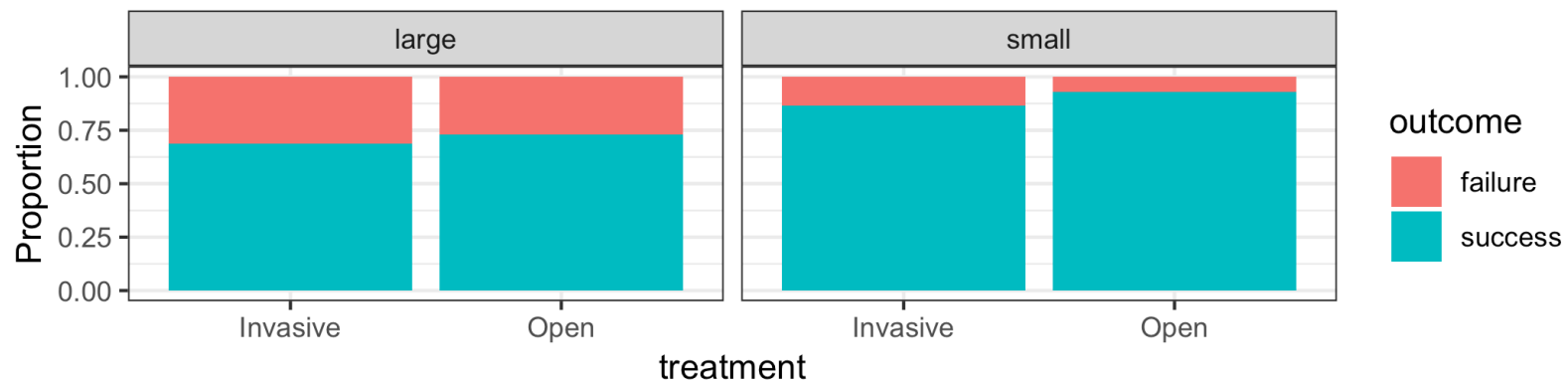
## Visualizing the kidney stone data: treatment and outcome

```
ggplot(kidney_stones, aes(x=treatment, fill=outcome)) +  
  geom_bar(position = "fill") +  
  labs(y = "Proportion") + theme_bw()
```



## Visualizing the kidney stone data: treatment and outcome by size

```
ggplot(kidney_stones, aes(x=treatment, fill=outcome)) +  
  geom_bar(position = "fill") + labs(y = "Proportion") +  
  facet_grid(. ~ size) +  
  theme_bw()
```





# Confounding

# What is a confounding variable?

- When examining the relationship between two variables in observational studies, it is important to consider the possible effects of other variables.

# What is a confounding variable?

- When examining the relationship between two variables in observational studies, it is important to consider the possible effects of other variables.
- A third variable is a **confounding variable** if it affects the nature of the relationship between two other variables, so that it is impossible to know if one variable causes another, or if the observed relationship is due to the third variable.

# What is a confounding variable?

- When examining the relationship between two variables in observational studies, it is important to consider the possible effects of other variables.
- A third variable is a **confounding variable** if it affects the nature of the relationship between two other variables, so that it is impossible to know if one variable causes another, or if the observed relationship is due to the third variable.
- The possible presence of confounding variables means we must be cautious when interpreting relationships.

# Examples of confounding?

- A 2012 [study](#) showed that heavy use of marijuana in adolescence can negatively affect IQ.

*Is it possible that there are other variables, such as socioeconomic status, that is associated with both marijuana use and IQ?*

# Examples of confounding?

- A 2012 [study](#) showed that heavy use of marijuana in adolescence can negatively affect IQ.

*Is it possible that there are other variables, such as socioeconomic status, that is associated with both marijuana use and IQ?*

- Another 2012 [study](#) showed that coffee drinking was inversely related to mortality.

*Should we all drink more coffee so we will live longer? Or is it possible that healthy people, who will live longer because they are healthy, are also more likely to drink coffee than unhealthy people?*

# Examples of confounding?

- A 2012 [study](#) showed that heavy use of marijuana in adolescence can negatively affect IQ.

*Is it possible that there are other variables, such as socioeconomic status, that is associated with both marijuana use and IQ?*

- Another 2012 [study](#) showed that coffee drinking was inversely related to mortality.

*Should we all drink more coffee so we will live longer? Or is it possible that healthy people, who will live longer because they are healthy, are also more likely to drink coffee than unhealthy people?*

- Many nutrition studies.

*Are people who are likely to stick to a diet different than those who won't in important ways?*

# How can confounding be avoided?

- Data can be collected through *experiments* or *observational studies*.



# How can confounding be avoided?

- Data can be collected through *experiments* or *observational studies*.
- In **observational studies**, data are collected without intervention. The data are measurements of existing characteristics of the individuals being measured.

# How can confounding be avoided?

- Data can be collected through *experiments* or *observational studies*.
- In **observational studies**, data are collected without intervention. The data are measurements of existing characteristics of the individuals being measured.
- In **experiments**, an investigator imposes an intervention on the individuals being studied, randomly assigning some individuals to one treatment and randomly assigning other individuals to another treatment (sometimes this other treatment is a *control*).

# How can confounding be avoided?

- Data can be collected through *experiments* or *observational studies*.
- In **observational studies**, data are collected without intervention. The data are measurements of existing characteristics of the individuals being measured.
- In **experiments**, an investigator imposes an intervention on the individuals being studied, randomly assigning some individuals to one treatment and randomly assigning other individuals to another treatment (sometimes this other treatment is a *control*).
- Randomized experiments are often used when we want to be able to say a treatment **causes** a change in a measurement.

# How can confounding be avoided?

- Data can be collected through *experiments* or *observational studies*.
- In **observational studies**, data are collected without intervention. The data are measurements of existing characteristics of the individuals being measured.
- In **experiments**, an investigator imposes an intervention on the individuals being studied, randomly assigning some individuals to one treatment and randomly assigning other individuals to another treatment (sometimes this other treatment is a *control*).
- Randomized experiments are often used when we want to be able to say a treatment **causes** a change in a measurement.
- Other than the difference in treatment received, any differences between the individuals in the treatment and control groups are just due to random chance in their group assignment.

# How can confounding be avoided?

- In a randomized experiment, if there is a difference in our measurement of interest, we *maybe* able to conclude it was caused by the treatment, and not due to some other systematic difference that can confound our interpretation of the effect of the treatment.

# How can confounding be avoided?

- In a randomized experiment, if there is a difference in our measurement of interest, we *maybe* able to conclude it was caused by the treatment, and not due to some other systematic difference that can confound our interpretation of the effect of the treatment.
- Example experiment from Week 5 lecture:  
Students were randomly assigned to be sleep-deprived or to have unrestricted sleep and how they learned a visual discrimination task was compared between these two groups.

# How can confounding be avoided?

- In a randomized experiment, if there is a difference in our measurement of interest, we *maybe* be able to conclude it was caused by the treatment, and not due to some other systematic difference that can confound our interpretation of the effect of the treatment.
- Example experiment from Week 5 lecture:  
Students were randomly assigned to be sleep-deprived or to have unrestricted sleep and how they learned a visual discrimination task was compared between these two groups.
- It's not always practical or ethical to carry out an experiment. For example, it would be considered unethical to randomly assign people to smoke marijuana.

# Ethical Issues in Data Science

Identification of ethical considerations involving research where data is collected:

- History of Ethical Codes: Nuremberg Code; and Declaration of Helsinki.
- Tuskegee Syphilis Study
- Informed consent
- Ethical issues in Data Science Research
- Ethical issues using Public Data
- Bias and Inclusion in AI Systems



## Class Survey

### **STA130 Quick Survey**

To complete the survey, go to  
**[Pollev.com/nathantaback](https://Pollev.com/nathantaback)**

# Research Ethics

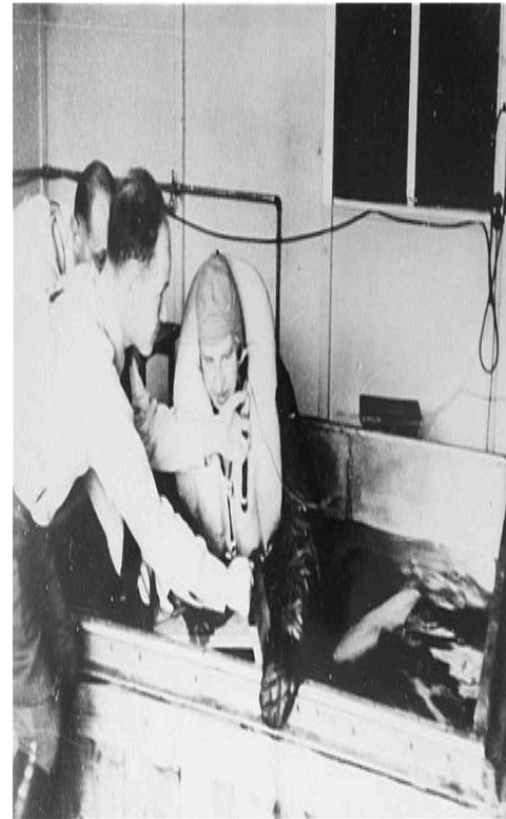
- Medical and scientific research involving humans usually require approval and monitoring by an ethics board.
- UofT and other universities have several ethics review boards.
- For example, if a survey is given to STA130 students then unless we are given permission, by the UofT ethics board, and students in the course then we would not be able to publish the findings (even the summary statistics) of the survey.
- Why are institutions such as Universities so cautious?

# History of Ethical Codes

# Nazi Medical Experiments

Experiments on prisoners included:

- Injecting dye into the eyes of twins to study conjoined twins.
- Removal of bones, muscles, nerves without anesthesia to study bone, muscle, and nerve regeneration.
- Infection with malaria then treated with various drugs to test efficacy.



A victim of a Nazi medical experiment is immersed in icy water at the Dachau concentration camp. SS doctor Sigmund Rascher oversees the experiment. Germany, 1942.

— Bildarchiv Preussischer Kulturbesitz

# Government Misuse of Data

- Systematic killing of several million Jews in Europe required extensive planning, organization, and coordination.

# Government Misuse of Data

- Systematic killing of several million Jews in Europe required extensive planning, organization, and coordination.
- Government statistical agencies conducted special censuses or population registries in certain places.

# Government Misuse of Data

- Systematic killing of several million Jews in Europe required extensive planning, organization, and coordination.
- Government statistical agencies conducted special censuses or population registries in certain places.
- The completed forms from these data collections were used to provide names and addresses of Jews to be included in transports to concentration camps and extermination camps.

# Government Misuse of Data

- Systematic killing of several million Jews in Europe required extensive planning, organization, and coordination.
- Government statistical agencies conducted special censuses or population registries in certain places.
- The completed forms from these data collections were used to provide names and addresses of Jews to be included in transports to concentration camps and extermination camps.
- The 1930 Dutch census were one of several data sources used to identify high density of Jewish population.



# Government Misuse of Data

- Systematic killing of several million Jews in Europe required extensive planning, organization, and coordination.
- Government statistical agencies conducted special censuses or population registries in certain places.
- The completed forms from these data collections were used to provide names and addresses of Jews to be included in transports to concentration camps and extermination camps.
- The 1930 Dutch census were one of several data sources used to identify high density of Jewish population.
- 1940 U.S. population census collected information on ancestry.

# Government Misuse of Data

- Systematic killing of several million Jews in Europe required extensive planning, organization, and coordination.
- Government statistical agencies conducted special censuses or population registries in certain places.
- The completed forms from these data collections were used to provide names and addresses of Jews to be included in transports to concentration camps and extermination camps.
- The 1930 Dutch census were one of several data sources used to identify high density of Jewish population.
- 1940 U.S. population census collected information on ancestry.
- This data was used in internment of Japanese-Americans at the outbreak of World War II.

(Sletzer, 1998)

# Nuremberg Code

- Ethical codes often emerge out of crisis events.
- The Nuremberg code was formulated in August 1947, in Nuremberg, Germany, by American judges sitting in judgment of Nazi doctors accused of conducting murderous and torturous human experiments in the concentration camps.
- The judges at Nuremberg realized the importance of Hippocratic code (do no harm) was not sufficient.



The defendants listen as the prosecution begins introducing documents at the International Military Tribunal trial of war criminals at Nuremberg. November 22, 1945.

— National Archives and Records Administration, College Park, Md.

# Nuremberg Code

The Nuremberg code codified many of our standard principles of ethical research today including:

- research must appropriately balance risk and potential reward (e.g., clinical equipoise),
- researchers must be well versed in their discipline and ground human experiments in animal trials.

# Ethical Scandals lead to U.S. Law

- Ethical codes did not carry the weight of law in the U.S. until after a series of scandals in the 1960s and 1970s.
- This led to the 1974 National Research Act, which established the National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research.
- A notable result of the commission was establishing institutional **ethics review boards** (also known as IRB or REB) which act as independent panels that review research proposals to assess possible harms to human subjects.
- This gives research institutions the power and responsibility to self-regulate through these boards.

# Ethical Scandals: Tuskegee Syphilis Study

- The Tuskegee syphilis study was one of the most notable scandals.
- In 1932 US government scientists enrolled 400 African American males from Alabama, known to be infected with syphilis (a sexually transmitted infection that can cause serious health problems).
- Study participants were followed-up to examine long-term effects of syphilis.

# Ethical Scandals: Tuskegee Syphilis Study

- Study participants told they were being treated for “bad blood” but received no medical intervention.
- The subjects were never told they had syphilis. Subjects were denied access to treatment, even for years after penicillin came into use in 1947.
- Study lasted for four decades; exposed in 1972.

# Ethical Scandals: Tuskegee Syphilis Study

By 1972,

- 28 study participants had died due to syphilis
- 100 had died from conditions related to syphilis
- 40 wives were infected
- 19 infants infected

On May 16 1997, President Clinton apologized for the role of US government in this study.





# Ethical Scandals: Tuskegee Syphilis Study

- Subjects were not given appropriate information about the risks and benefits of the study.
- After subjects informed about the risks and benefits they did not give their permission to participate in the study.

# Informed Consent

Three elements of informed consent:

1. Information
2. Comprehension
3. Voluntariness

# Informed Consent

**Information:** the research procedure, their purposes, risks and anticipated benefits, alternative procedures (where therapy is involved), and a statement offering the subject the opportunity to ask questions and to withdraw at any time from the research.

# Informed Consent

**Information:** the research procedure, their purposes, risks and anticipated benefits, alternative procedures (where therapy is involved), and a statement offering the subject the opportunity to ask questions and to withdraw at any time from the research.

# Informed Consent

**Information:** the research procedure, their purposes, risks and anticipated benefits, alternative procedures (where therapy is involved), and a statement offering the subject the opportunity to ask questions and to withdraw at any time from the research.

A physician researcher asks her patients to enrol in a study she is conducting. She slowly explains the study including the risks and benefits, and gives her patients ample time to consider being a research subject. Since research is an important part of her position at the hospital she informs patients that if they don't enrol in her study that she can no longer act as their physician. Which element of informed consent is not present?

 Respond at **PollEv.com/nanthantaback**  Text **NATHANTABACK** to **37607** once to join, then **A, B, C, or D**

Information	<b>A</b>
Comprehension	<b>B</b>
Voluntariness	<b>C</b>
All of the elements of informed consent are present.	<b>D</b>

# Ethics in Data Science Research

# Ethics in Data Science Research

- Regulations built around the research/practice distinction in medicine are a method for signaling and negotiating temporary changes to physician-patient relationship.
- A patient must be informed, and consent to, situations in which a physician may no longer be making or be able to make decisions in the best interest of the patient.
- In a research context, a physician has the best interest of the social collective as an explicit competing interest to the well-being of the patient.
- There is no easy analogue for the physician-researcher in data science.

(Metcalfe and Crawford, 2016)



# Ethics in Data Science Research

- What are the ethical obligations data scientists have for the well-being of human subjects in data science research? How do we assess that those obligations are being met?
- Currently there are no laws that address these questions.
- Several organizations (e.g., [Statistical Society of Canada](#), [American Statistical Association](#), [Association for Computing Machinery](#)) have developed detailed statements on topics such as professionalism, integrity of data and methods, responsibilities to stakeholders, conflicts of interest, and the response to allegations of misconduct.

# Ethics in Data Science Research

- One challenge in developing an analogue of medical research ethics in the data science framework is the criteria for human-subjects' protections.
- The criteria depends on an unstated assumption that the risk to research subjects depends on what kind of data is obtained and how it is obtained, not what is done with the data after it is obtained.
- This assumption is based on the idea that data which is public poses no new risks for human subjects.
- Data science drives significant changes to "how we know" by creating new knowledge through tying together previously disconnected data sets.

(Metcalfe and Crawford, 2016)

# Ethics in Data Science Research

- Most ethics boards exempt research of existing data, documents, records, and specimens if that data is publicly available.
- This means that most non-medical data science will receive very little review. (Metcalf and Crawford, 2016)

## Principles to Determine Exemptions from Research Ethics Review

### Preamble

According to the Tri-Council Policy Statement (TCPS), the mandate of a Research Ethics Board (REB) is to provide **Research Ethics Review (RER)** for all **research involving human subjects**, as defined in Article 1.1. It is not within the mandate of the REB to review research activities outside of this definition. Items c) and d) of this article give some guidance as to the types of activities which may be **exempt**: research based on publicly available data or individuals in the public arena, research-like activities that fall into the areas of quality assurance or performance reviews and non-research activities such as testing within normal educational requirements.

# New York City Taxi & Limousine Commission

- In 2013, the New York City Taxi & Limousine Commission released a data set of 173 million individual cab rides.
- The data set included the pickup and drop off times, locations, fare and tip amounts.
- The taxi drivers' medallion numbers were anonymized

(Metcalf and Crawford, 2016)

# New York City Taxi & Limousine Commission

- Researchers were able to de-anonymized the data to reveal sensitive information such as any driver's annual income and enabling researchers to infer their home address (Franceschi-Bicchierai, 2015).
- A data scientist at Neustar Research showed that by combining this data set with other forms of public information like celebrity blogs you could track well-known actors, and predict likely home addresses of people who frequented strip clubs (Tockar, 2014).
- Another researcher demonstrated how the taxi data set could be used to predict which drivers were devout Muslims by observing which drivers stopped at prayer times (Franceschi-Bicchierai, 2015).

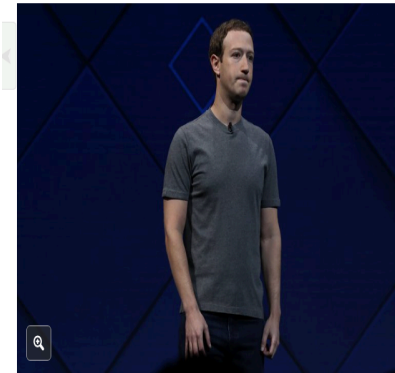
(Metcalf and Crawford, 2016)

# Cambridge Analytica

- News reports earlier this year claimed that Cambridge Analytica bought 50 million Facebook profiles from a researcher.
- Only 270,000 users consented to having their data used by the researcher.
- Cambridge Analytica combined The Facebook data with other databases to build profiles of these users.




## *Zuckerberg Takes Steps to Calm Facebook Employees*

By SHEERA FRENKEL MARCH 23, 2018



Mark Zuckerberg, Facebook's chief executive, has been on an apology tour of sorts this week for his company's mishandling of data privacy. On Friday, he spoke to employees. Jim Wilson/The New York Times

### RELATED COVERAGE

-  Zuckerberg, Facing Facebook's Worst Crisis Yet, Pledges Better Privacy  
MARCH 21, 2018
-  Facebook's Role in Data Misuse Sets Off Storms on Two Continents  
MARCH 18, 2018
-  How Trump Consultants Exploited the Facebook Data of Millions  
MARCH 17, 2018

- This allowed the company to use the data to target users with specific ads.
- What ethical obligations should a private company have in protecting users' data?

# Ethics of Public Data

- Users of data services, such as social media platforms, often know very little about how their private data will be used in research.
- Should Facebook users now expect that their social media activities could affect their ability to get a loan or influence how they vote?
- Is it reasonable to assume that social behavior on Facebook is the same as social relationships outside of Facebook? Could this assumption cause economic harm to individuals and communities?

(Metcalfe and Crawford, 2016)

# Ethics of Public Data

- If human subjects research regulations assume that public data sets are harmless, it will be nearly impossible to review the consequences of the people affected.
- These "data subjects" may have no knowledge that they are part of a study that might affect their future opportunities or well-being.

(Metcalfe and Crawford, 2016)

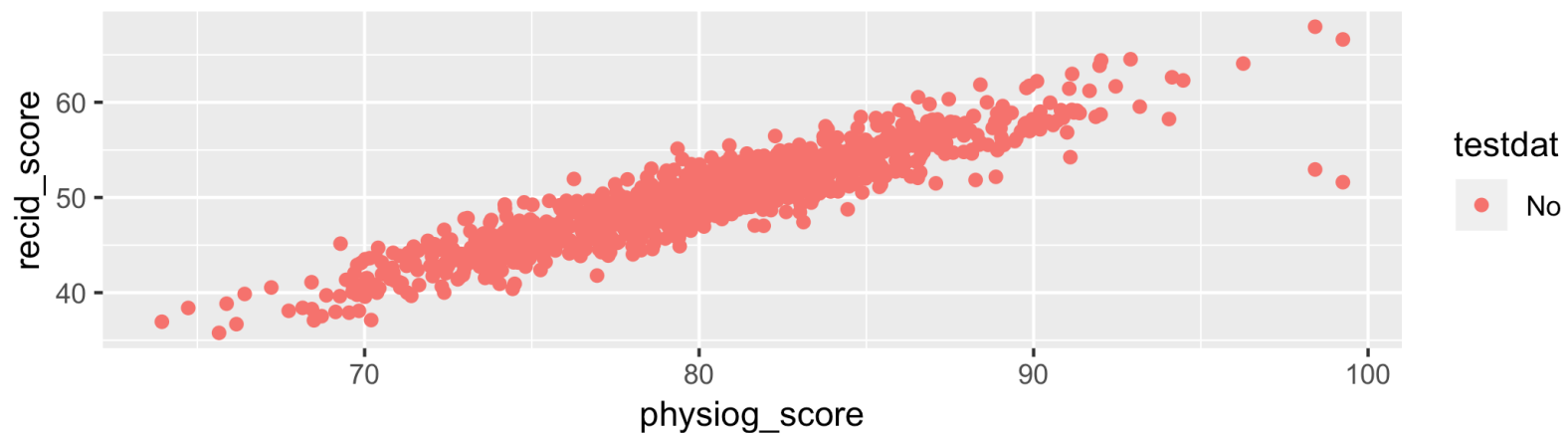


# Experimentation on Data Subjects

- An A/B test randomizes subjects to two conditions A and B.
- Randomized experiments can be used to establish causation not just correlation.
- Facebook randomized 700,000 users to two different types of newsfeeds to investigate if it could manipulate their emotions.
- Since Facebook is a private company it was not required to use any independent review process to approve the research.

# Predicting Recidivism

- Consider a simple Artificial Intelligence system designed to predict recidivism based on a person's facial features.
- A data set of 1,001 adult males convicted of crimes were assessed for personality traits associated with their facial features.
- Each person in the data set was given a recidivism score (extremely unlikely to reoffend 0 - 100 extremely likely to reoffend), and physiognomy score (facial features atypical of criminal 0 - 100 facial features typical of a criminal)



# Predicting Recidivism

- Linear regression was used to build a prediction model.

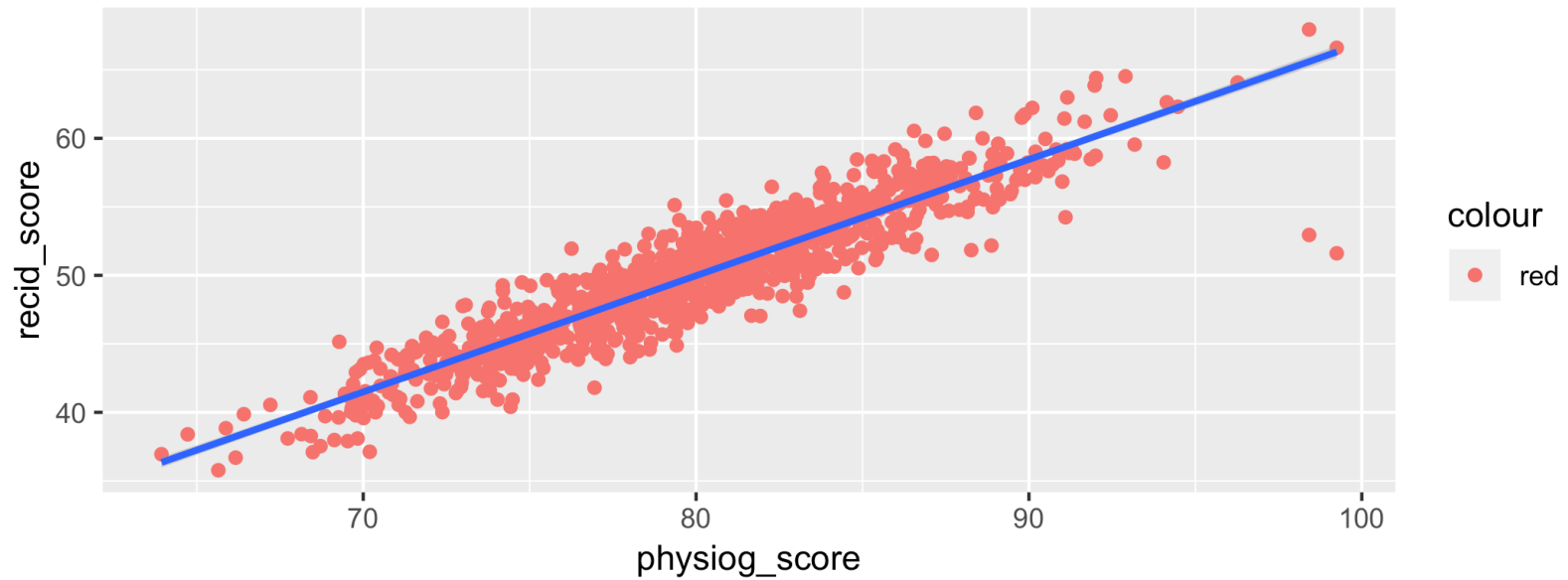
```
set.seed(10)
train <- dat_crime %>% sample_frac(size = 0.8)
test <- dat_crime %>% anti_join(train, by = 'id')
reg_mod <- lm(recid_score ~ physiog_score, data = train)
summary(reg_mod)$r.squared
```

```
## [1] 0.854069
```

```
yhat <- predict(reg_mod, newdata = test)
y <- test$recid_score
sqrt(sum((y - yhat)^2) / length(test$recid_score))
```

```
## [1] 2.077955
```

```
dat_crime %>%  
  ggplot(aes(x = physiog_score, y = recid_score)) +  
  geom_point(aes(colour = "red")) +  
  geom_smooth(method = "lm")
```



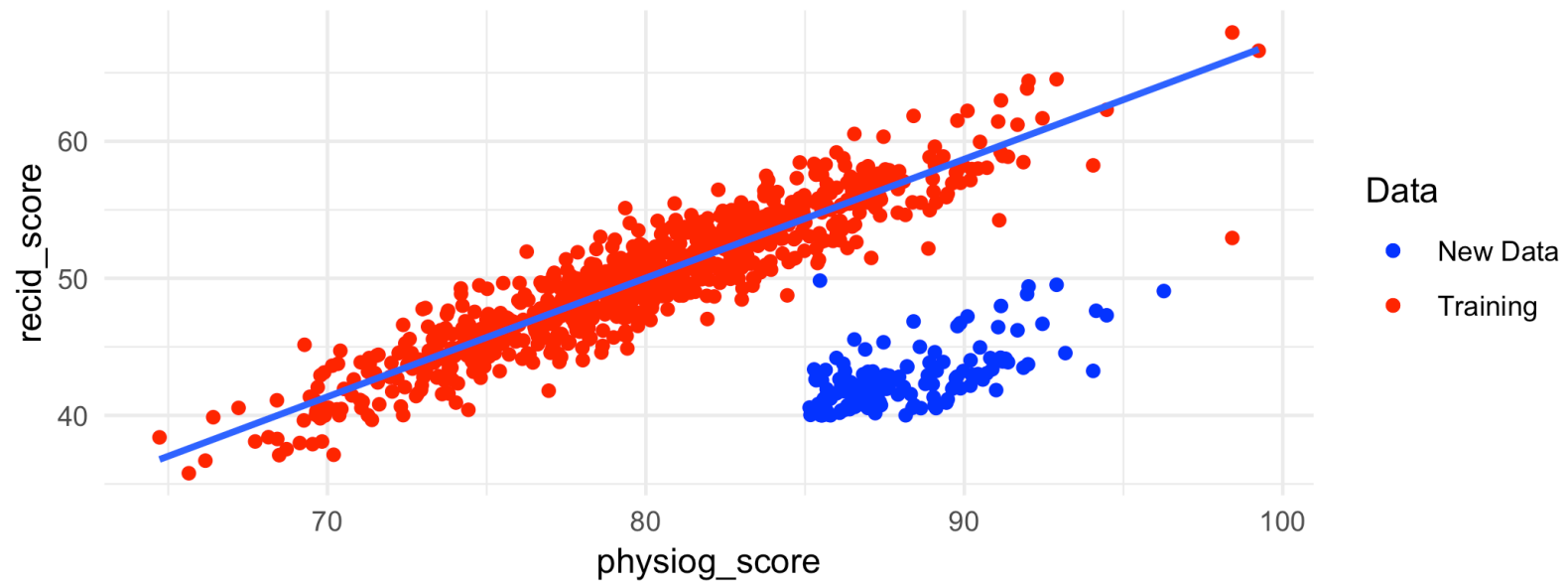
# Predicting Recidivism

- The company now claims that the model will work in predicting recidivism in another country.
- The software is used a few cases at a time.
- After 140 cases the police force that is using the software decide to evaluate the accuracy of the predictions.

```
yhat <- predict(reg_mod, newdata = dat_new)
y <- dat_new$recid_score
sqrt(sum((y - yhat)^2) / length(dat_new$recid_score))
```

```
## [1] 14.48995
```

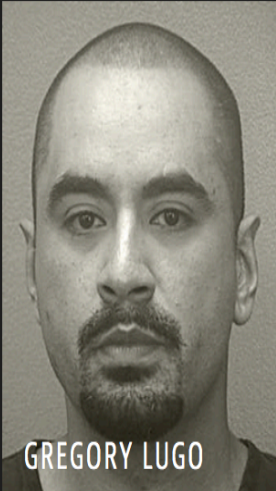

# Predicting Recidivism



# Predicting Future Crime - Pro Publica

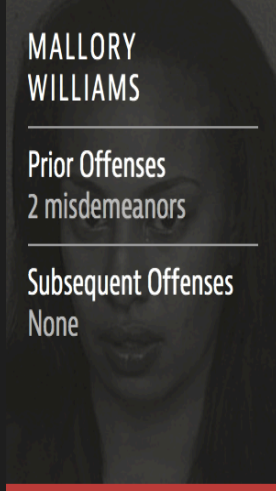
## Investigation of COMPAS

Two DUI Arrests

Gregory Lugo	Mallory Williams
	
LOW RISK 1	MEDIUM RISK 6

Lugo crashed his Lincoln Navigator into a Toyota Camry while drunk. He was rated as a low risk of reoffending despite the fact that it was at least his fourth DUI.

Two DUI Arrests

Gregory Lugo	Mallory Williams
	
Prior Offenses 3 DUIs, 1 battery	Prior Offenses 2 misdemeanors
Subsequent Offenses 1 domestic violence battery	Subsequent Offenses None
LOW RISK 1	MEDIUM RISK 6

Lugo crashed his Lincoln Navigator into a Toyota Camry while drunk. He was rated as a low risk of reoffending despite the fact that it was at least his fourth DUI.

# Predicting Future Crime - Pro Publica Investigation of COMPAS

- Pro Publica investigated an AI system used by courts and judges to predict recidivism.
- Pro Publica obtained the risk scores assigned to more than 7,000 people arrested in Broward County, Florida, in 2013 and 2014.
- They checked to see how many were charged with new crimes over the next two years, the same benchmark used by the creators of the algorithm.
- 20% of the people predicted to commit violent crimes actually went on to do so.



# Bias and Inclusion in AI Systems and ML Algorithms

- AI systems are taught what they “know” from training data.
- Training data can be incomplete, biased, or skewed. This is sometimes referred to as "algorithmic bias".
- Training data may come from poorly defined non-representative samples of a population.
- These problems with training data may not be obvious if the data set construction is non-transparent.

# Ethical Concerns in AI and ML Algorithms

- Should AI systems be used in sensitive or high-stakes contexts?
- Who gets to make these decisions?
- What is the proper degree of human involvement in various types of decision-making?