# CS224C: NLP for CSS

# Computational Basics

Diyi Yang

Stanford CS

# Announcements

Please sign up for Presentation by this coming Tuesday (Apr 9th)

# Data Access

**Twitter's new data access rules will make social media research harder**
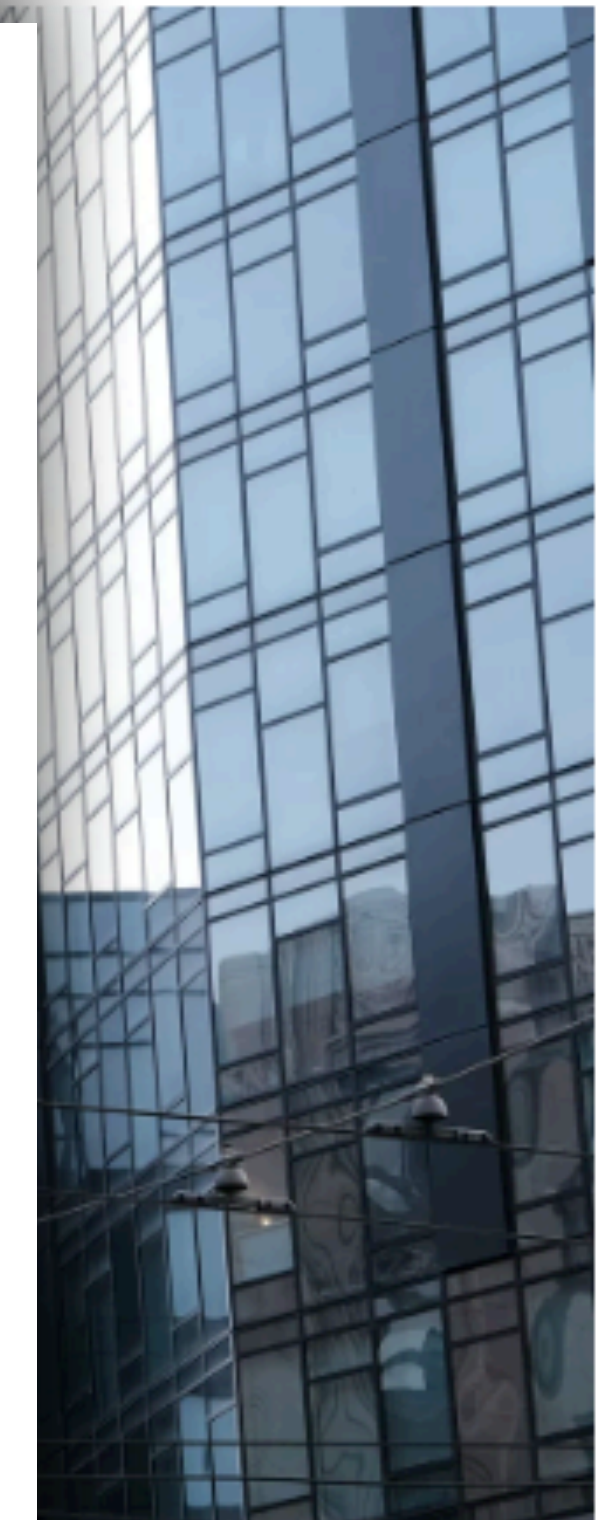
FEBRUARY 9, 2023 · 7:00 AM ET

By Huo Jingnan

# Without access to social media platform data, we risk being left in the dark

**Significance:**

Social media data are essential for studying human behaviour and understanding potential systemic risks. Social media platforms have, however, begun to remove access to these data. In response, other countries and regions have implemented legislation that compels platforms to provide researchers with data access. In South Africa, we have lagged behind the Global North when it comes to using platform data in our research and, given the recent access restrictions, we risk being left behind. In this Commentary, I call attention to this critical issue and initiate a conversation about access to social media data in South Africa.
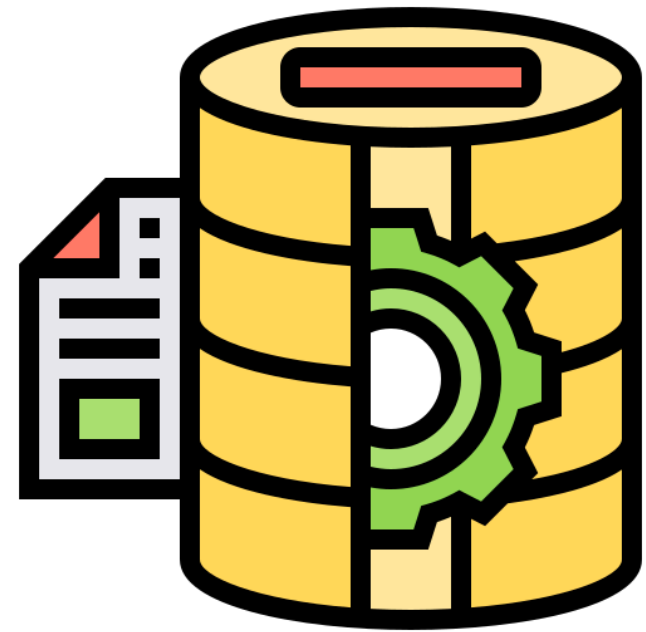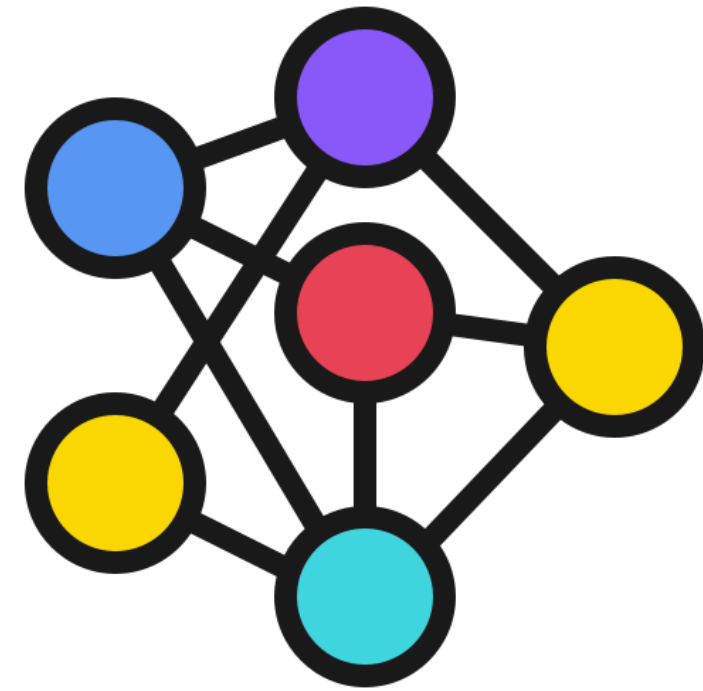
rder to researchers to

# Challenge Summary

1. The complexity of the theoretical issues confronting social science

2. The difficulty in obtaining the relevant observational data

3. The difficulty of manipulating large scale social organizationals experimentally

4. **The complexity and difficulty in computationally, scientifically and rigorously modeling such problems and data**
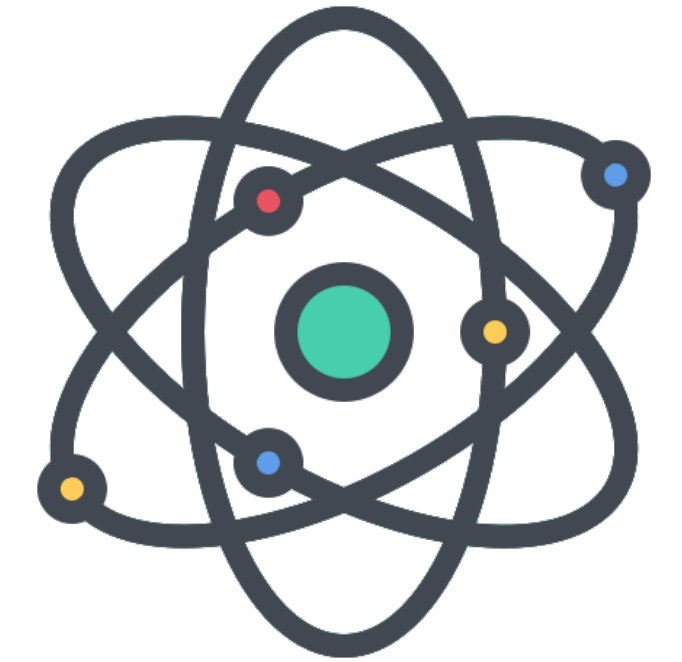
# Computational Social Science in a nutshell

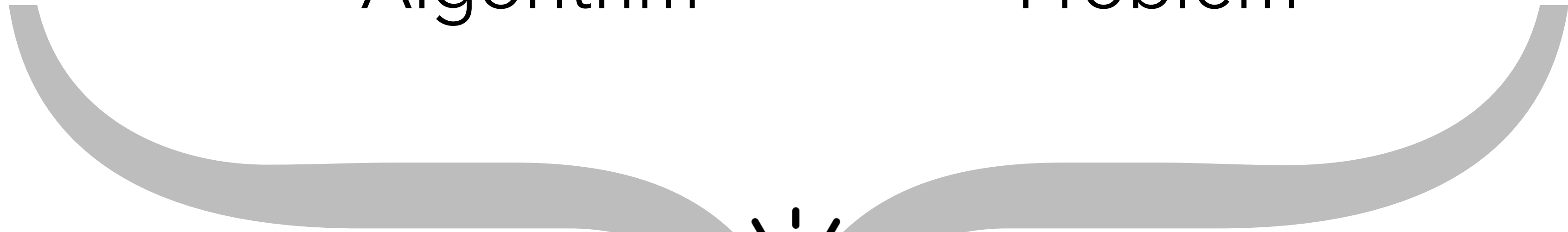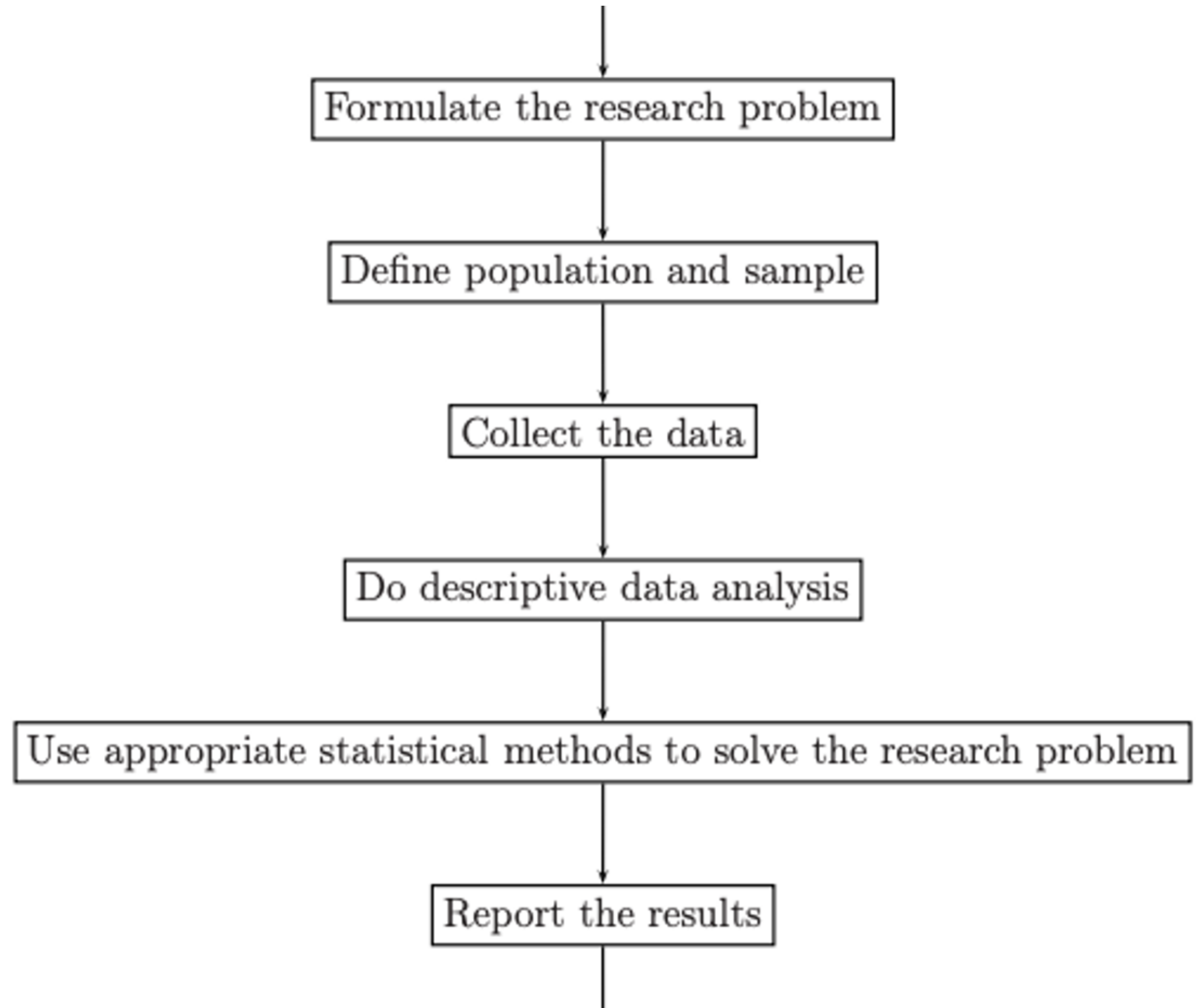

Data       Algorithm       Problem       Theory

Knowledge
Impact

# Lecture Overview

- ✦ Classification
- ✦ Regression
- ✦ Clustering
- ✦ Big Data + CSS

# Computational Framework



Formulate the research problem

↓

Define population and sample

↓

Collect the data

↓

Do descriptive data analysis

↓

Use appropriate statistical methods to solve the research problem

↓

Report the results

# Use of Classification or Regression

Two major uses of supervised classification/regression
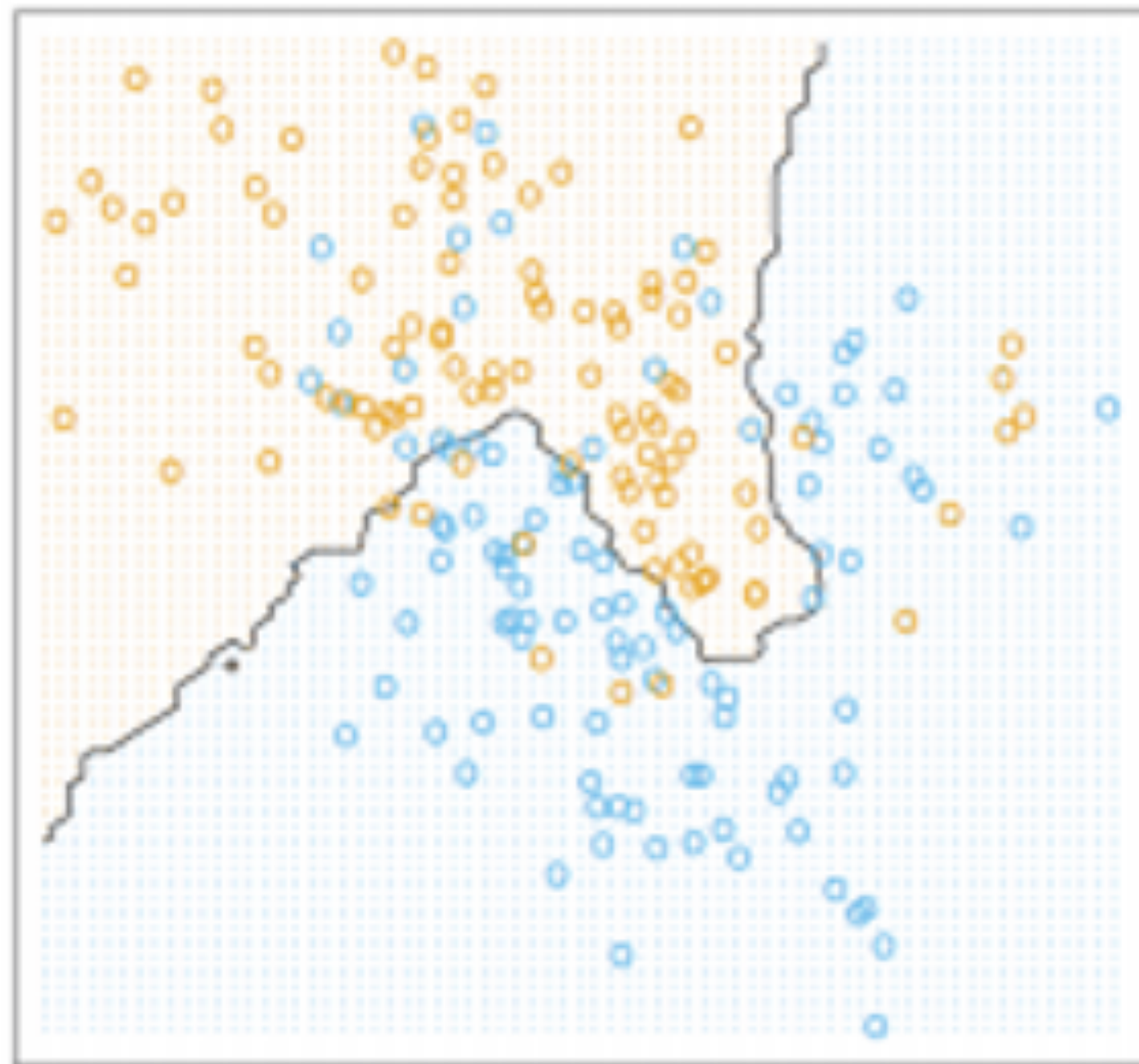
**Prediction:**

Train a model on a sample of data $(x, y)$ to predict for some new data $x'$

In zero-shot or few-shot prompting, no training is needed!

**Interpretation or Explanation:**

Train a model on a sample of data $(x, y)$ to understand the relationship between $x$ and $y$

# Common Methods



Classification                Clustering                Regression

# Classification

A ***mapping*** $h$ from input data $x$ (drawn from instance space $X$) to a label $y$ from some enumerable output space $Y$

$X$ = set of all documents

$Y$ = {English, Mandarin, Greek, …}

$x$ = a single document

$y$ = ancient Greek

# Reviews and Ratings

Reviewed: October 24, 2022

**Lovely little spot to spend some time. Very grateful for the clean, stylish room after travelling.**

8.0

🙂 · Beautiful, spacious room - after 22 hours of travel and a botched flight, I cried with happiness when I arrive.
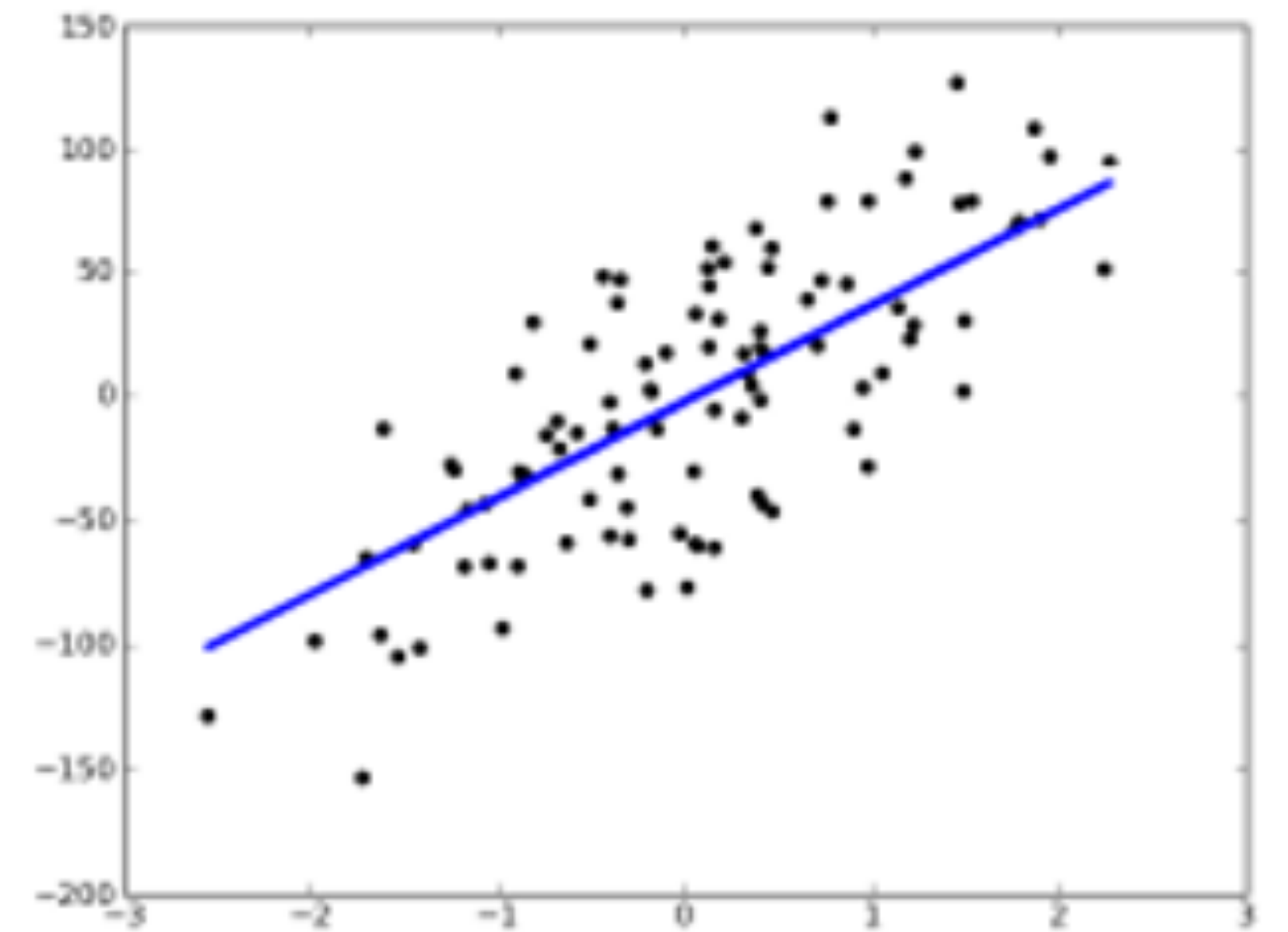Large, comfy bed - I didn't want to get out.
Friendly, helpful staff - especially the bar staff.
Accessible, enjoyable bar - open to late with a large selection of drinks.
Lots of room to sit by the pool. With a spa too. As well as on the foreshore in hammocks. Very laidback, enjoyable environment. I happily spent the day here, relaxing before a friends wedding. I loved walking along the foreshore footpath, following the shoreline and walking past the other hotels.

☹ · The beach wasn't an inviting swim, though a beautiful backdrop - which is not a fault of the hotel's. But flagging in case you're romanticising a beach swim; the hotel pool is better.
Watch: hidden costs. This might be normal/acceptable in non-Australian countries but I was caught off guard. There's the room cost, then there's the taxes (which booking.com tends to include in their final price), and THEN the hotel has a 'resort fee'. Which allows for 'amenities' access - which I find a bit "on the nose", the 'amenities' is what you automatically have access to when you book a room... but i guess some countres/states prefer a staggered bill...

## IMDb Charts

## IMDb Top 250 Movies
IMDb Top 250 as rated by regular IMDb voters.

Showing 250 Titles

Sort by: Ranking

| Rank & Title | IMDb Rating | Your Rating | |
|---|---|---|---|
| 1. The Shawshank Redemption (1994) | ⭐ 9.2 | ☆ | + |
| 2. The Godfather (1972) | ⭐ 9.2 | ☆ | + |
| 3. The Dark Knight (2008) | ⭐ 9.0 | ☆ | + |
| 4. The Godfather Part II (1974) | ⭐ 9.0 | ☆ | + |
| 5. 12 Angry Men (1957) | ⭐ 9.0 | ☆ | + |
| 6. Schindler's List (1993) | ⭐ 8.9 | ☆ | + |

# Some Text Classification Applications

| Task | x | y |
|---|---|---|
| Language identification | text | {English, Mandarin, Greek, …} |
| Spam classification | email | {spam, not spam} |
| Authorship attribution | text | {jk rowling, james joyce, …} |
| Genre classification | novel | {detective, romance, gothic, …} |
| Sentiment classification | text | {positive, negative, neutral, mixed} |

# Lots of Model Choices

Decision Trees

Logistic Regression

Support Vector Machine

Random Forests

Neural Nets (e.g., BERT)

Prompting LLMs

# Model Differences

Binary Classification

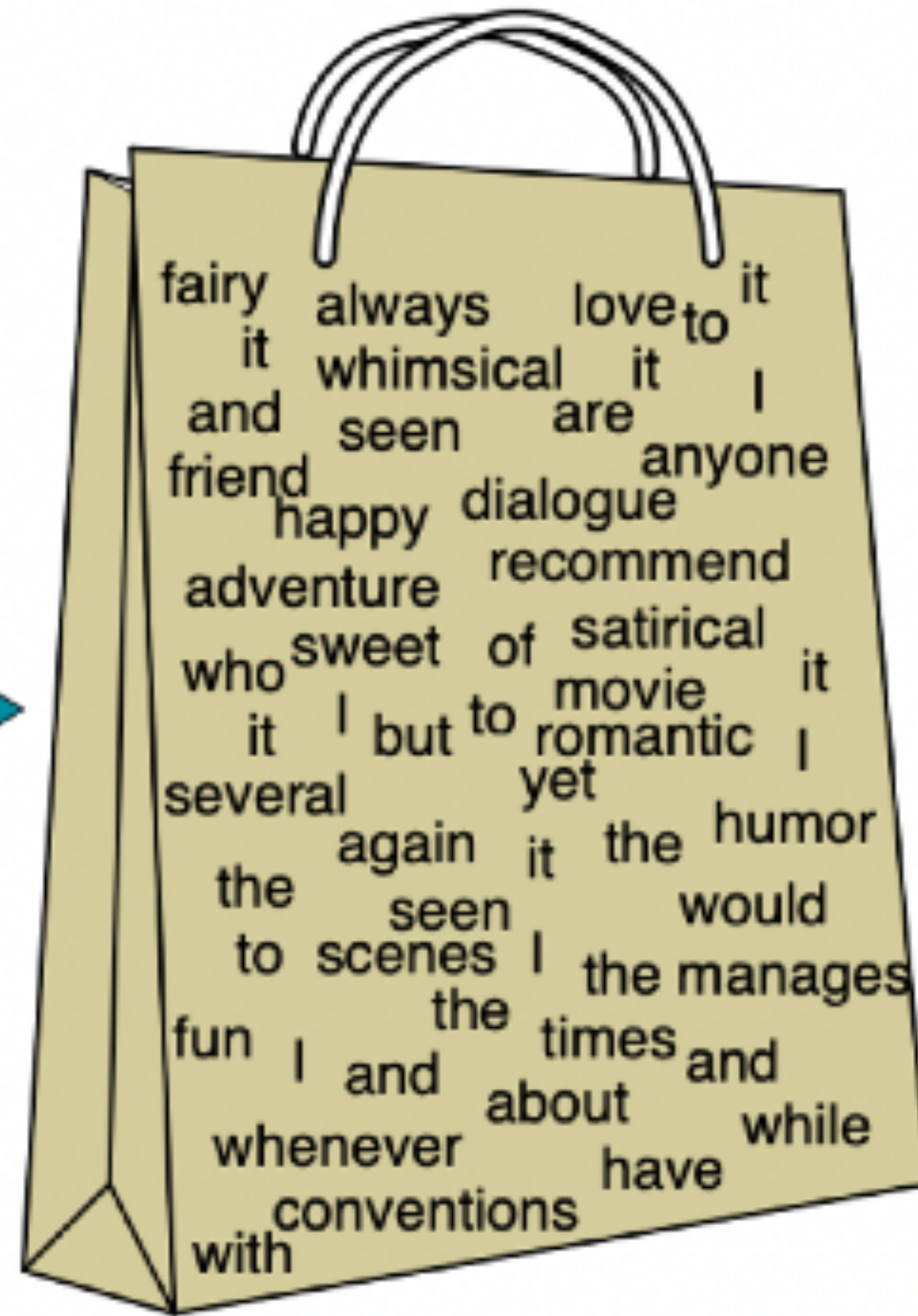*One out of 2 labels applied to a given x*

Multiclass Classification

*One out of N labels applied to a given x*

Multilabel Classification

*Multiple labels apply to a given x*

# Bag of Words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

fairy always love to it
it whimsical it I
and seen are
friend anyone
happy dialogue
adventure recommend
who sweet of satirical it
it I but to movie
several romantic I
yet
the again it the humor
seen would
to scenes I the manages
fun the times and
I and about
whenever while
conventions have
with

| it | 6 |
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| … | … |

# Beyond the Bag of Words

Some linguistic phenomena require going beyond the bag-of-words:

- *That's not bad for the first day*

- *This is not the worst thing that can happen*

- *It would be nice if you acted like you understood*

- *This film should be brilliant. The actors are first grade. Stallone plays a happy, wonderful man. His sweet wife is beautiful and adores him. He has a fascinating gift for living life fully. It sounds like a great plot, however, the film is a failure.*

# Signals in Text Beyond Words

Emojis: 🙏🤷🏽‍♂️🥰

Special characters: ' } { [ ] # @ ! * < > ~

Out of vocabulary words: icebucketchallenge, wowwwww

URLs: *https://www.nytimes.com/*

Typos or spelling errors: *typs, tpos, …*

Social media features: *@user, RT, #hashtags*

Slang words: *chill, slay, sick …*

# Signals in Text Beyond English

Language identification has very high accuracy for long texts, but struggles with social media (short informal) text

Code switching:   *I have 2 friends* *due estudiaron la contabilidad*

# Applying Text Classification

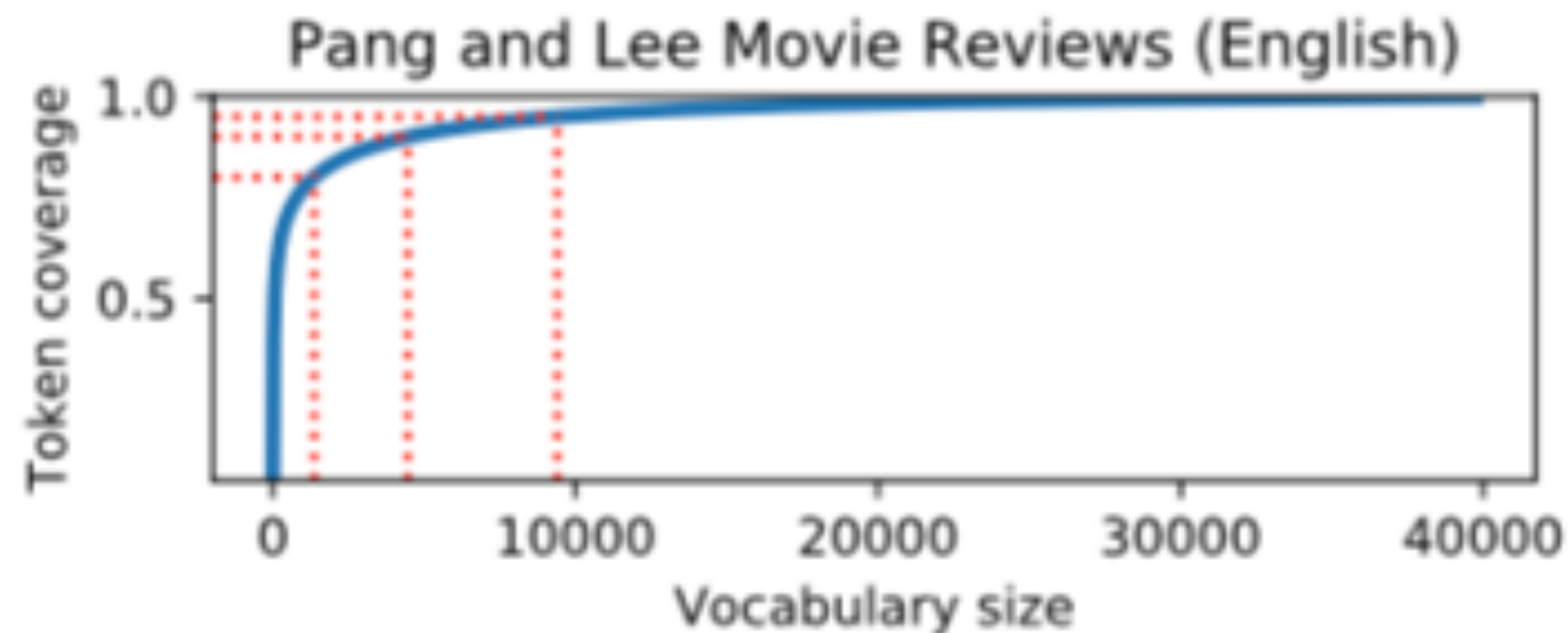The "raw" form of text is usually a sequence of characters

We do not cover this in class. Check out this after class if needed

Converting this into a meaningful feature vector $x$ requires a series of design decisions, such as tokenization, normalization, and filtering

# Vocabulary Size Filtering

A small number of word types accounts for the majority of word tokens.



The number of parameters in a classifier usually grows linearly with the size of the vocabulary. It can be useful to limit the vocabulary, e.g., to word types appearing at least x times, or in at least y% of documents.

# Experiment Setup (Train/Test Split)

|  | **Training** | **Development Or Validation** | **Testing** |
|---|---|---|---|
| **Size** | 80% | 10% | 10% |
| **Purpose** | Training Models | Model Selection *(e.g., parameters)* | Evaluation *(You should never look at it until the very end)* |

# Evaluating Your Classifier

Goal is to predict future performance, on unseen data.

It is hard to predict the future.

Do not evaluate on data that was already used …

*For training*

*For hyperparameter selection*

*For selecting the classification model or model structure*

*For making preprocessing decisions, such as vocabulary selection.*

# Beyond Right and Wrong

For any label, there are two ways to be wrong:

**False positive:** the system incorrectly predicts the label

**False negative:** the system incorrectly fails to predict the label.

Similarly, there are two ways to be right:

**True positive:** the system correctly predicts the label

**True negative:** the system correctly predicts that the label does not apply to it.

# Accuracy

$$accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

The problem with accuracy is rare labels.

Consider a system for detecting tweets written in Telugu.

0.3% of Tweets are written in Telugu.

A system that always says "Not Telugu" is 99.7% accurate.

# Recall

$$Recall = \frac{TP}{TP + FN}$$

Recall is the fraction of positive instances which were correctly classified.

The "never Telugu" classifier has zero recall.

An "always Telugu" classifier would have perfect recall.

# Precision

$$\text{Precision} = \frac{TP}{TP + FP}$$

Precision is the fraction of positive predictions that were correct.

The "never Telugu" classifier has precision 0/0.

An "always Telugu" classifier would have precision p=0.003, which is the rate of Telugu tweets in the dataset.

# Combining Recall and Precision

In binary classification, there is an inherent tradeoff btw recall and precision.

The correct navigation of this tradeoff is problem-specific!

*For a preliminary medical diagnosis, we might prefer high recall. False positives can be screened out later.*

*The "beyond a reasonable doubt" standard of U.S. criminal law implies a preference for high precision.*

# Combining Recall and Precision

In binary classification, there is an inherent tradeoff btw recall and precision.

The correct navigation of this tradeoff is problem-specific!

If recall and precision are weighted equally, they can be combined into a single number called F-measure

$$F = \frac{2 \times Recall \times Precision}{Recall + Precision}$$

# Evaluating Multi-Class Classification

Recall and precision imply binary classification: each instance is either positive or negative.

In multi-class classification, each instance is positive for one class, and negative for all other classes.

# Evaluating Multi-Class Classification

Two ways to combine performance across classes:

**Macro F-measure:** compute the F-measure per class, and average across all classes. This treats all classes equally, regardless of their frequency.

**Micro F-measure:** compute the total number of true positives, false positives, and false negatives across all classes, and compute a single F-measure. This emphasizes performance on high-frequency classes.

# Comparing Classifiers

Suppose you and your friend build classifiers to solve a problem:

You classifier $C_1$ get 82% accuracy

You friend's classifier $C_2$ get 73% accuracy

Will $C_1$ be more accurate in the future?

What is the test set had 10000 examples?

What is the test set had 11 examples?

# Getting Labels

Text classification relies on large datasets of labeled examples. There are two main ways to get labels:

Metadata sometimes tell us exactly what we want to know: Did the Senator vote for a bill? How many stars did the reviewer give? Was the request for free pizza accepted?

Other times, the labels must be annotated, by experts or by "crow-workers"

# Let's build a hate speech classifier for X

**What do we need?** 🤔

1. **Why** do we need to set it up as a classification task?
2. **What** counts as a hate speech?
3. How can we get the **ground truth** for a tweet?
4. How many **data points** do we need?
5. How can we **evaluate** this system?
6. Does the system really **work**?

# Let's build a hate speech classifier for X

A tweet is offensive if it

1. uses a sexist or racial slur.
2. attacks a minority.
3. seeks to silence a minority.
4. criticizes a minority (without a well founded argument).
5. promotes, but does not directly use, hate speech or violent crime.
6. criticizes a minority and uses a straw man argument.
7. blatantly misrepresents truth or seeks to distort views on a minority with unfounded claims.
8. shows support of problematic hash tags. E.g. "#BanIslam", "#whoriental", "#whitegenocide"
9. negatively stereotypes a minority.
10. defends xenophobia or sexism.
11. contains a screen name that is offensive, as per the previous criteria, the tweet is ambiguous (at best), and the tweet is on a topic that satisfies any of the above criteria.

**Hateful Symbols or Hateful People?**
**Predictive Features for Hate Speech Detection on Twitter**

**Zeerak Waseem**
University of Copenhagen
Copenhagen, Denmark
csp265@alumni.ku.dk

**Dirk Hovy**
University of Copenhagen
Copenhagen, Denmark
dirk.hovy@hum.ku.dk

"Annotated 16,914 tweets,
3,383 of that for sexist content sent by 613 users,
1,972 for racist content sent by 9 users
11,559 for neither sexist or racist, sent by 614 users"

|  | char $n$-grams | +gender | +gender +loc | word $n$-grams |
|---|---|---|---|---|
| F1 | 73.89 | 73.93 | 73.62* | 64.58 |
| Precision | 72.87% | 72.93% | 72.58% | 64.39% |
| Recall | 77.75% | 77.74% | 77.43% | 71.93% |

35

# Let's build a hate speech classifier for X
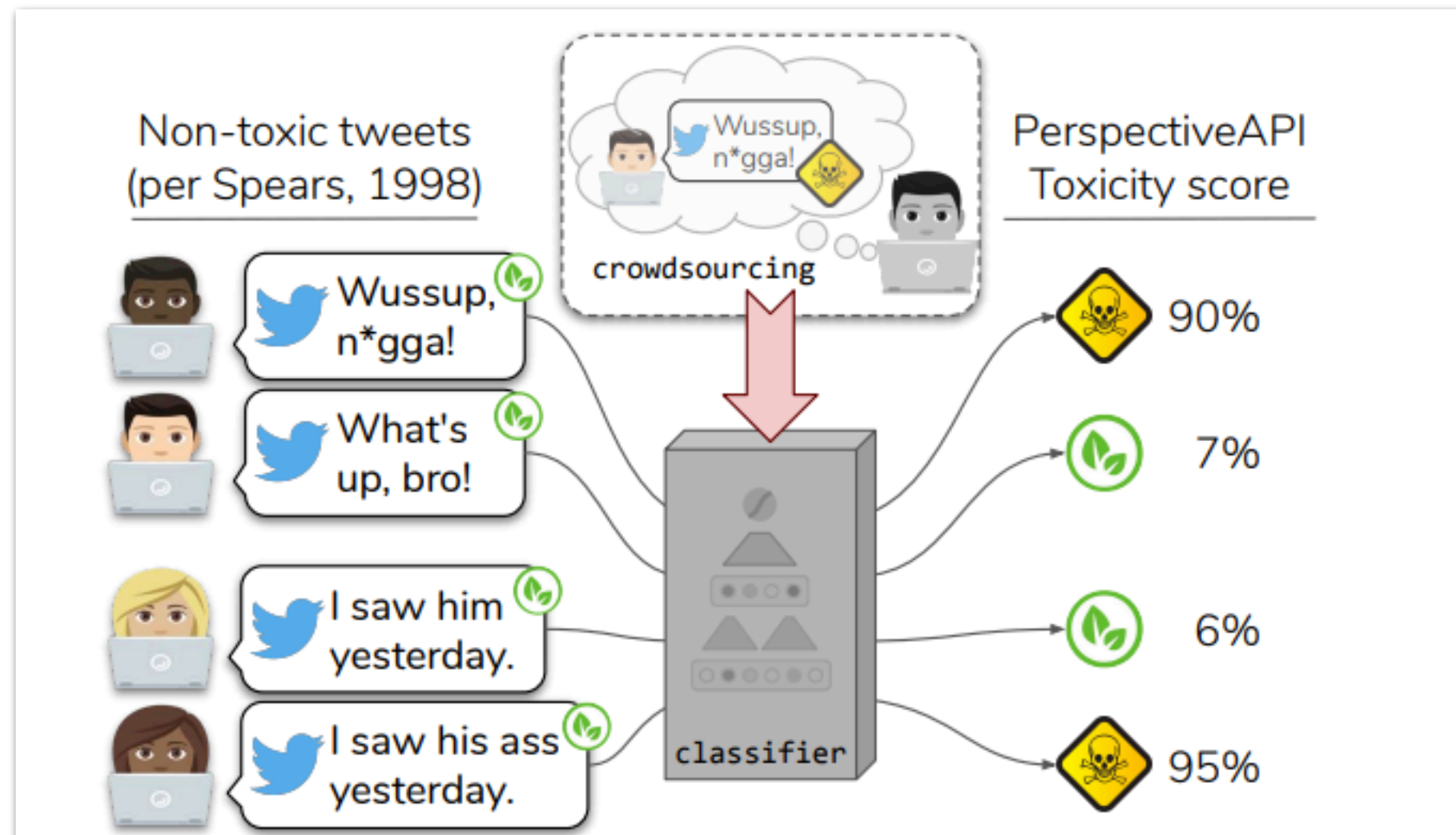
Is this classifier good?



Figure 1: Phrases in African American English (AAE), their non-AAE equivalents (from Spears, 1998), and toxicity scores from PerspectiveAPI.com. Perspective is a tool from Jigsaw/Alphabet that uses a convolutional neural network to detect toxic language, trained on crowdsourced data where annotators were asked to label the toxicity of text without metadata.

## Racial Bias in Hate Speech and Abusive Language Detection Datasets

**Thomas Davidson**
Department of Sociology
Cornell University
trd54@cornell.edu

**Debasmita Bhattacharya**
Department of
Computer Science
Cornell University
db758@cornell.edu

**Ingmar Weber**
Qatar Computer
Research Institute
iweber@hbku.edu.qa

## The Risk of Racial Bias in Hate Speech Detection

Maarten Sap◇   Dallas Card♣   Saadia Gabriel◇   Yejin Choi◇♡   Noah A. Smith◇♡
◇Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, USA
♣Machine Learning Department, Carnegie Mellon University, Pittsburgh, USA
♡Allen Institute for Artificial Intelligence, Seattle, USA

## Exploring the Role of Grammar and Word Choice in Bias Toward African American English (AAE) in Hate Speech Classification

Camille Harris
Georgia Institute of Technology
Atlanta, GA, USA
charris320@gatech.edu

Matan Halevy
Georgia Institute of Technology
Atlanta, GA, USA
matan@gatech.edu

Ayanna Howard
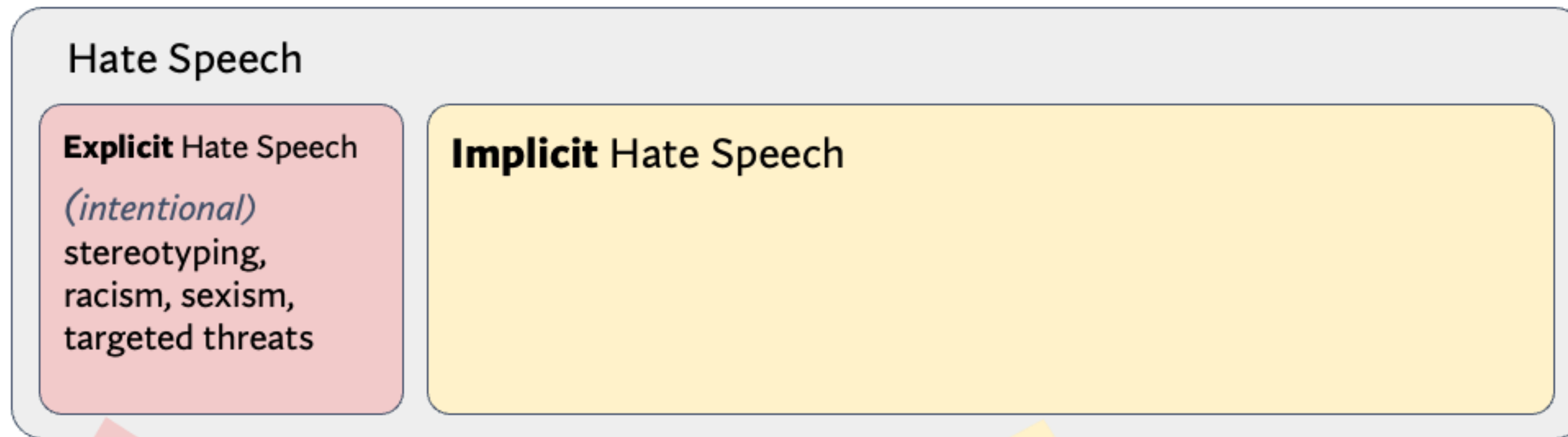The Ohio State University
OH, USA

Amy Bruckman
Georgia Institute of Technology
Atlanta, GA, USA

Diyi Yang
Georgia Institute of Technology
Atlanta, GA, USA

# Let's build a hate speech classifier for X

## Hate speech goes beyond explicit usage of keywords



Hate Speech

**Explicit** Hate Speech

*(intentional)* stereotyping, racism, sexism, targeted threats

**Implicit** Hate Speech

"[IDENTITY] destroy everything they touch"

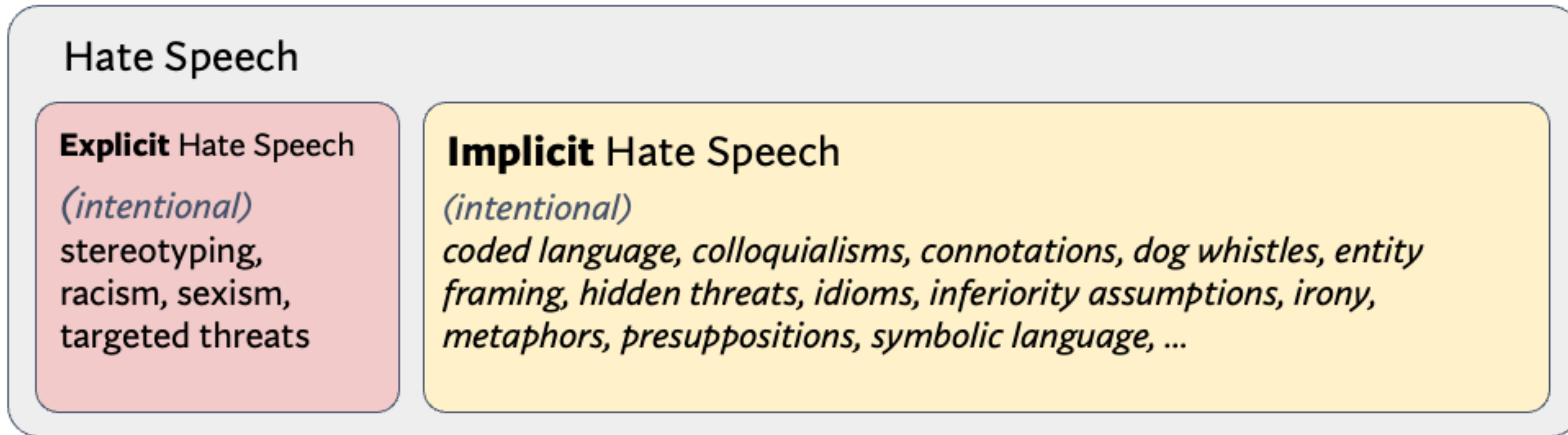*Explicit Hate Speech Example*

"White revolution is the only solution."

*Implicit Hate Speech Example*

Mai ElSherief [*◦]     Caleb Ziems [*†]     David Muchlinski[†]     Vaishnavi Anupindi[†]
Jordyn Seybolt[†]     Munmun De Choudhury[†]     Diyi Yang[†]
[◦]UC San Diego, [†]Georgia Institute of Technology
melsherief@ucsd.edu
{cziems, dmuchlinski3, vanupindi3}@gatech.edu
{jseybolt3, munmund, dyang888}@gatech.edu

# Let's build a hate speech classifier for X

**Hate speech goes beyond explicit usage of keywords**

## Hate Speech

**Explicit** Hate Speech

*(intentional)*
stereotyping,
racism, sexism,
targeted threats

**Implicit** Hate Speech

*(intentional)*
*coded language, colloquialisms, connotations, dog whistles, entity framing, hidden threats, idioms, inferiority assumptions, irony, metaphors, presuppositions, symbolic language, …*

# Let's build an implicit hate speech classifier for X

Incitement of Violence (Somerville, 2011; Assembly, 1966)

Inferiority Language (Nielsen, 2002; Kennedy et al., 2018)

Irony (Waseem and Hovy, 2016; Justo et al., 2014)

Stereotypes (Warner and Hirschberg, 2012)

Threats and Intimidation (Sanguinetti et al., 2018)

White Grievance (Berbrier, 2000; Bloch et al., 2020; Miller-Driss, 2020)

# Let's build an implicit hate speech classifier for X



**5M** tweets from prominent hate groups

Filter out tweets with explicit hate

Implicit hate categories, targets & implied statement

Annotation with good agreement Fleiss' $\kappa = 0.6$

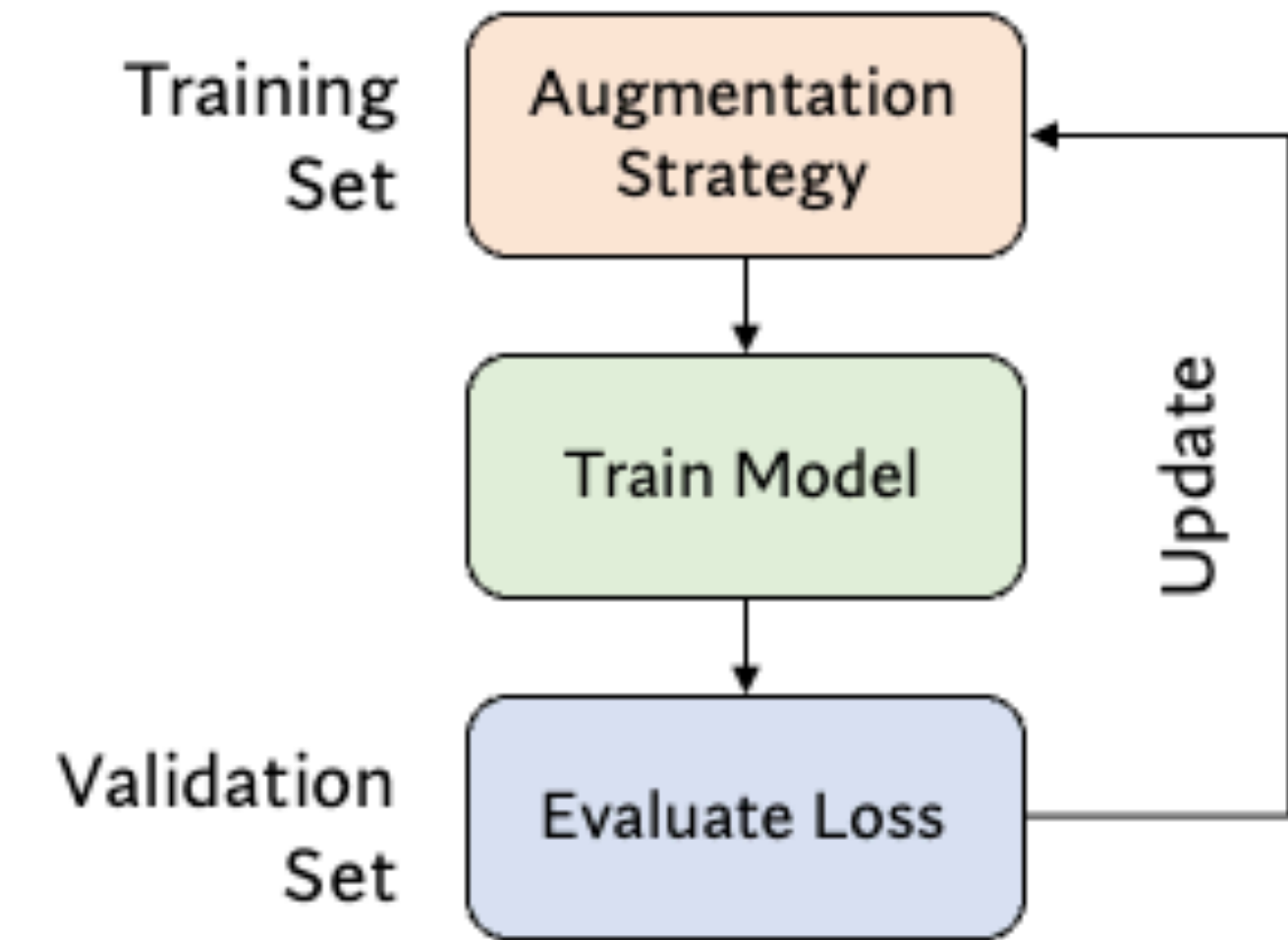# Let's build an implicit hate speech classifier for X
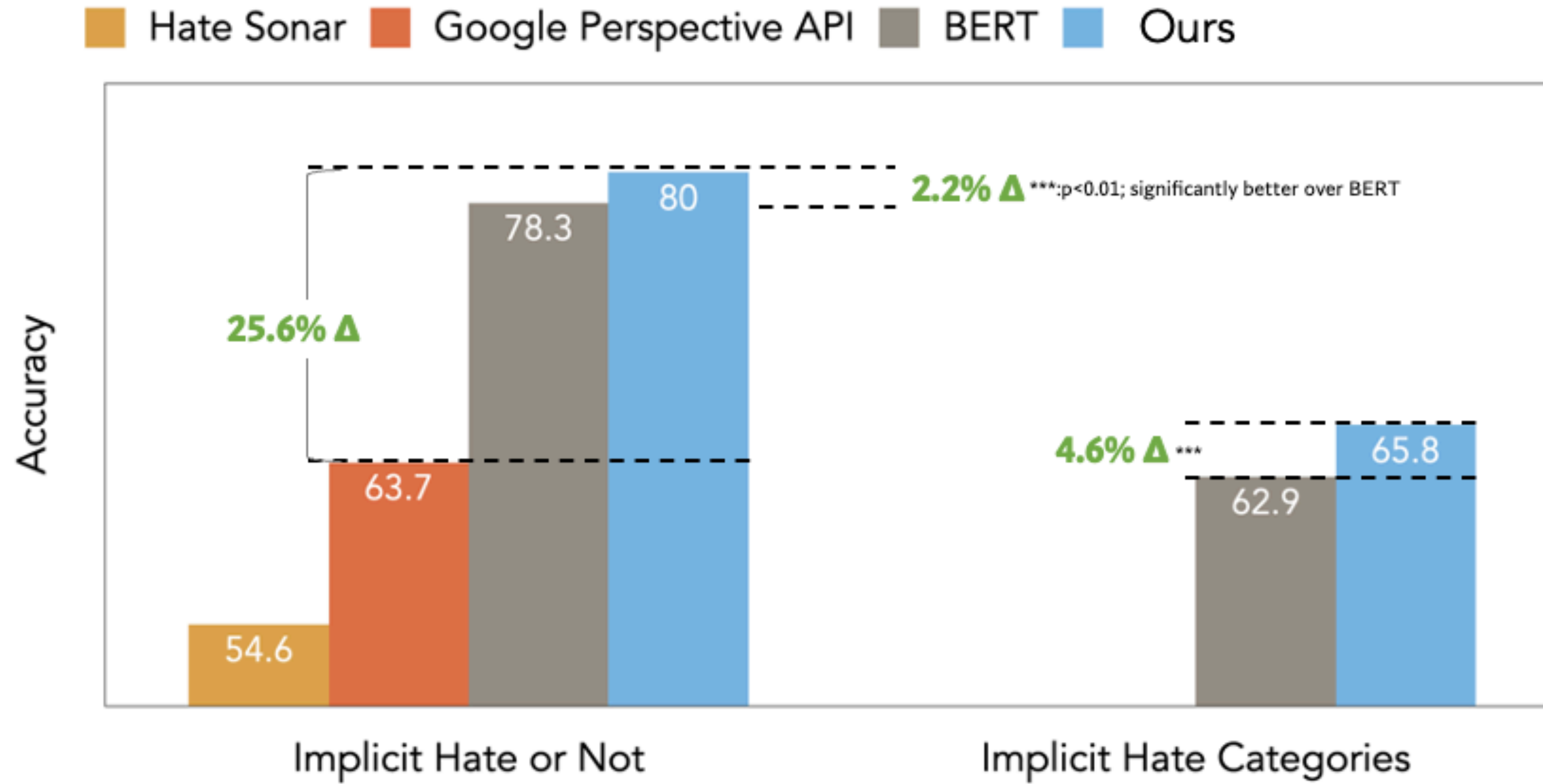


BERT_BASE

(Devlin et al., 2018)

| Warner and Hirschberg (2012) | Burnap and Williams (2014) | Djuric et al. (2015) | Waseem and Hovy (2016) |
| Gao and Huang (2017) | Davidson et al. (2017) | de Gibert et al. (2018) | Kennedy et al. (2018) |
| Founta et al. (2018) | Zampieri et al. (2019) | Basile et al. (2019) | Sap et al. (2020) |

1ˢᵗ stage fine-tune on existing hate speech datasets via *multi-task learning*

Training Set → Augmentation Strategy → Train Model → Validation Set / Evaluate Loss → Update

2ⁿᵈ stage fine-tune on Implicit Hate, with auto-augmentation

41

# Let's build an implicit hate speech classifier for X

# Recognizing A Classification Problem and Its Complexities

Can you formulate your question as a choice among some possible classes?

Can you create (or find) labeled data that marks that choice for a bunch of examples? Can you make that choice?

Can you create features that might help in distinguishing those classes?

Is the classification enough to capture the nuances?

# Regression

# Regression

A mapping from input data $x$ (drawn from instance space $X$) to a point $y$ in $R$

$R$: the set of real numbers
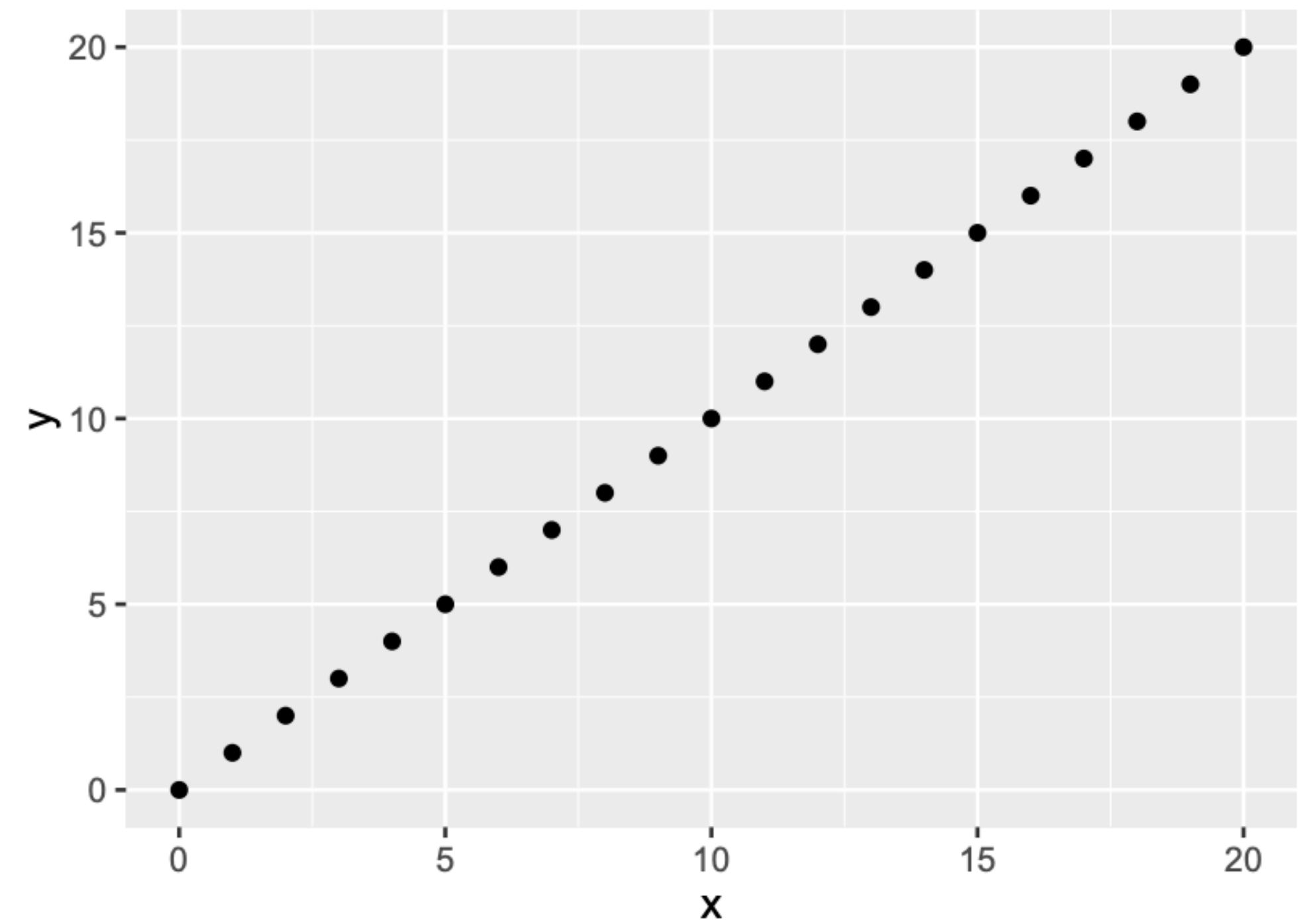
$x$ = the empire state building

$y$ = 17444.5625″

45

# Linear Regression

Suppose we have $n$ data points. For each data point $i$, we observe

$$(x_1, y_1), (x_2, y_2), (x_3, y_3), \ldots, (x_n, y_n)$$

Linear regression states that $\hat{y}_i = \sum_{i=1}^{F} x_i \beta_i$



Slide content credit to David Bamman
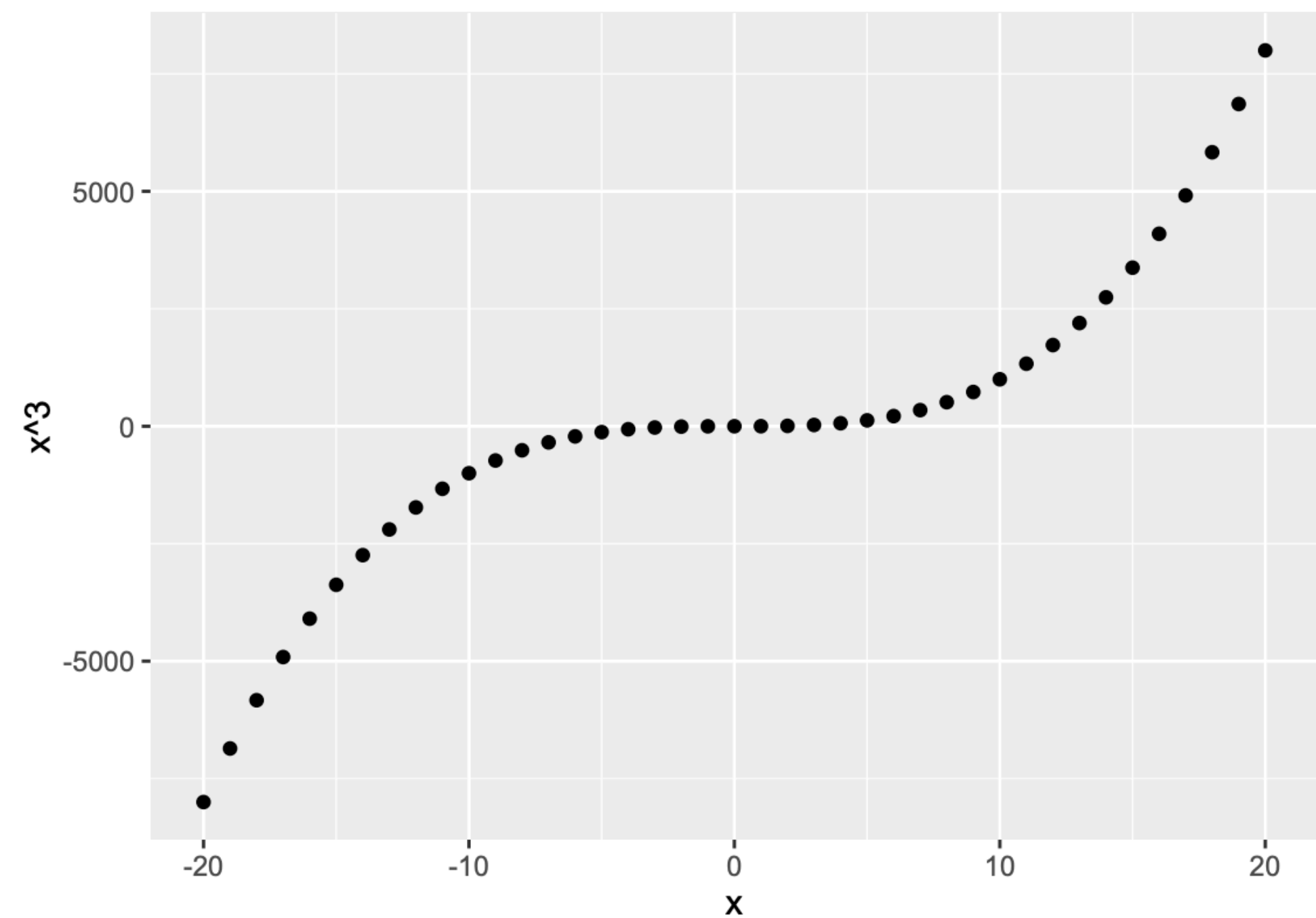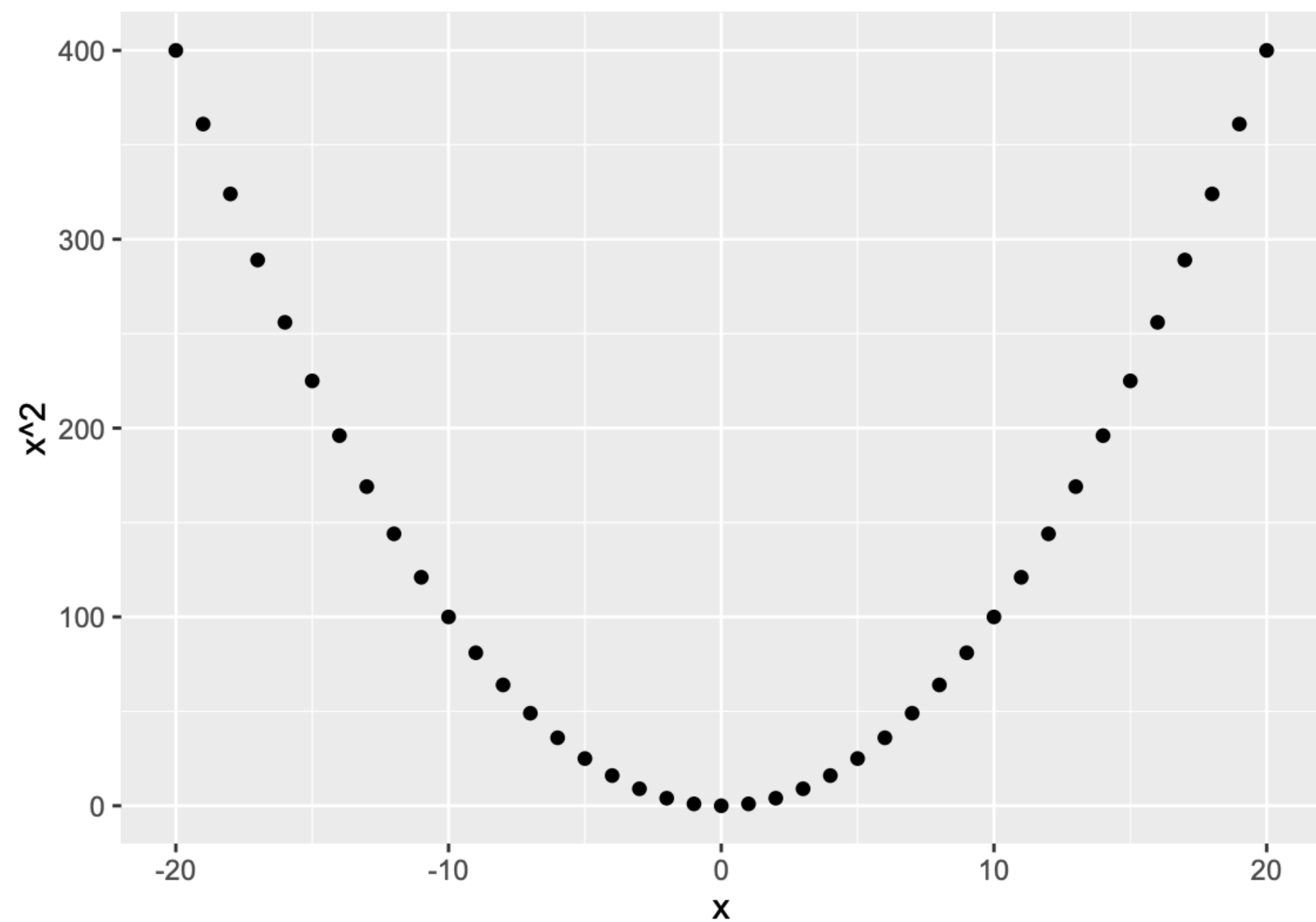
46

# Regression for Social Sciences

# Polynomial Regression

$$\hat{y}_i = \sum_{i=1}^{F} x_i \beta_{a,i} + \sum_{i=1}^{F} x_i^2 \beta_{b,i}$$

$$\hat{y}_i = \sum_{i=1}^{F} x_i \beta_{a,i} + \sum_{i=1}^{F} x_i^2 \beta_{b,i} + \sum_{i=1}^{F} x_i^3 \beta_{c,i}$$



Slide content credit to David Bamman

48

# Nonlinear Regression

Support vector machines (regression)

Neural Networks

…

# Number of Parameters

$$\hat{y}_i = \sum_{i=1}^{F} x_i \beta_{a,i}$$

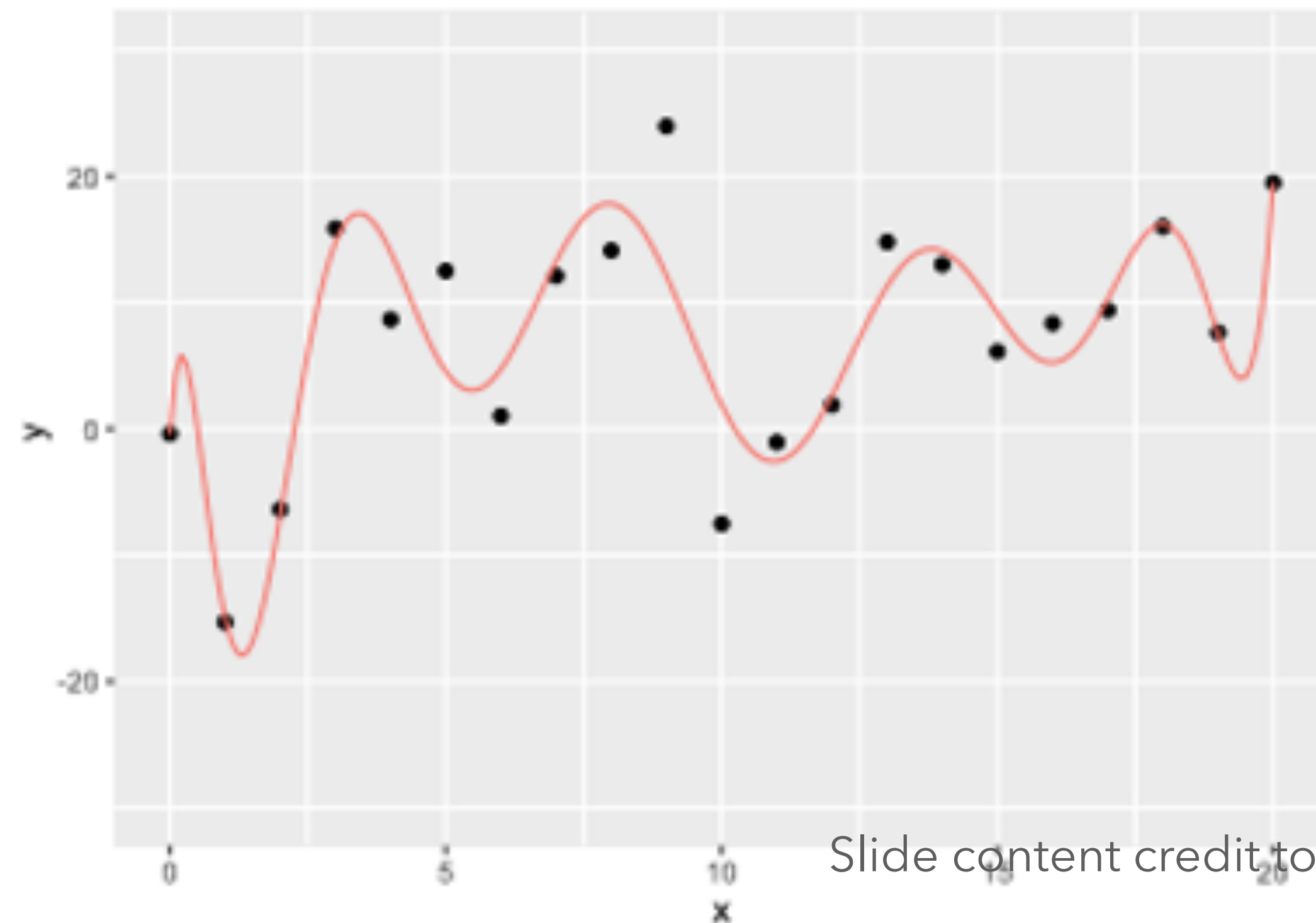$$\hat{y}_i = \sum_{i=1}^{F} x_i \beta_{a,i} + \sum_{i=1}^{F} x_i^2 \beta_{b,i}$$
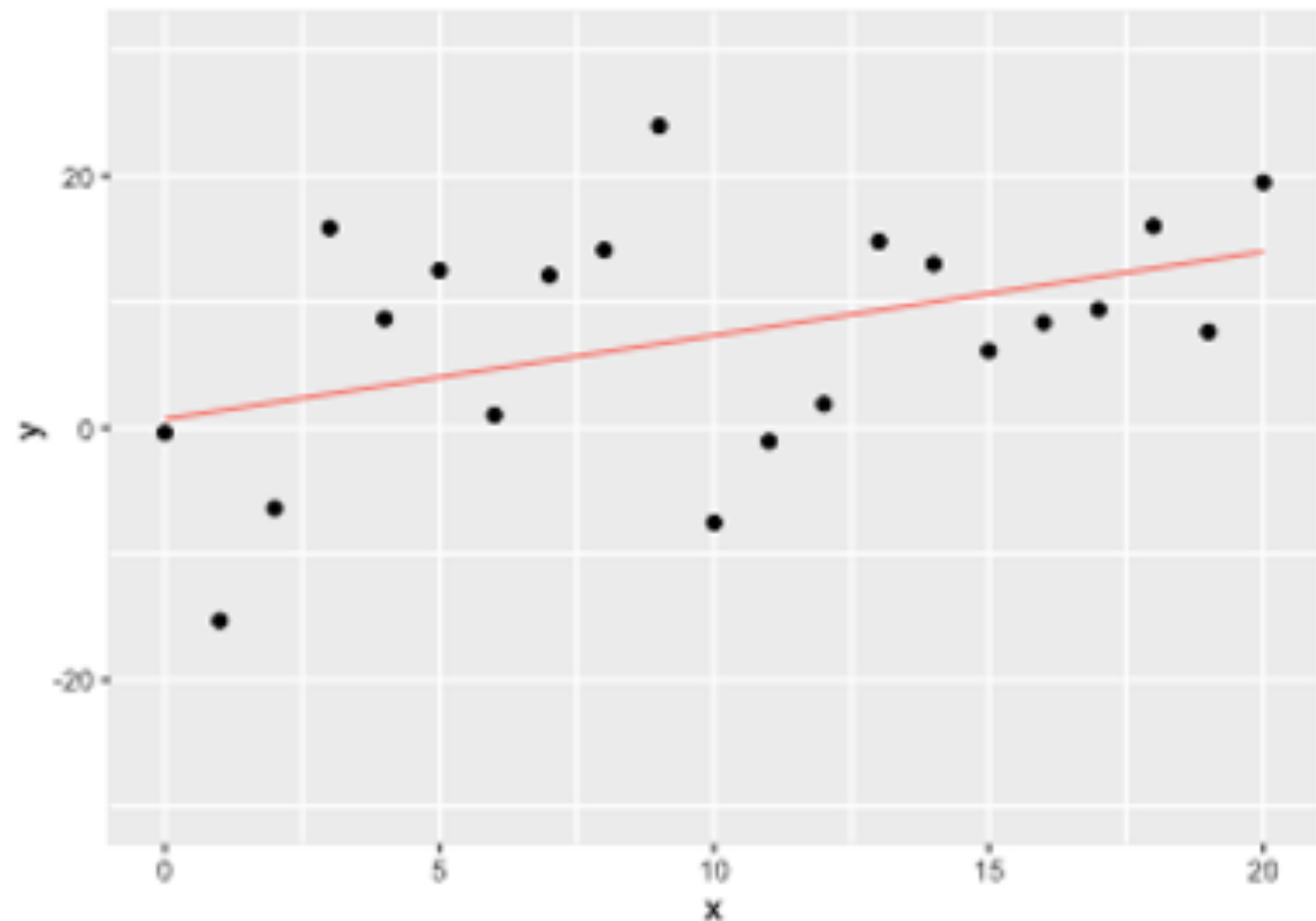
$$\hat{y}_i = \sum_{i=1}^{F} x_i \beta_{a,i} + \sum_{i=1}^{F} x_i^2 \beta_{b,i} + \sum_{i=1}^{F} x_i^3 \beta_{c,i}$$

Slide content credit to David Bamman

# Overfitting

Memorizing the nuances (and noise) of the training data that prevents generalizing to unseen data

51

# Sources of Error

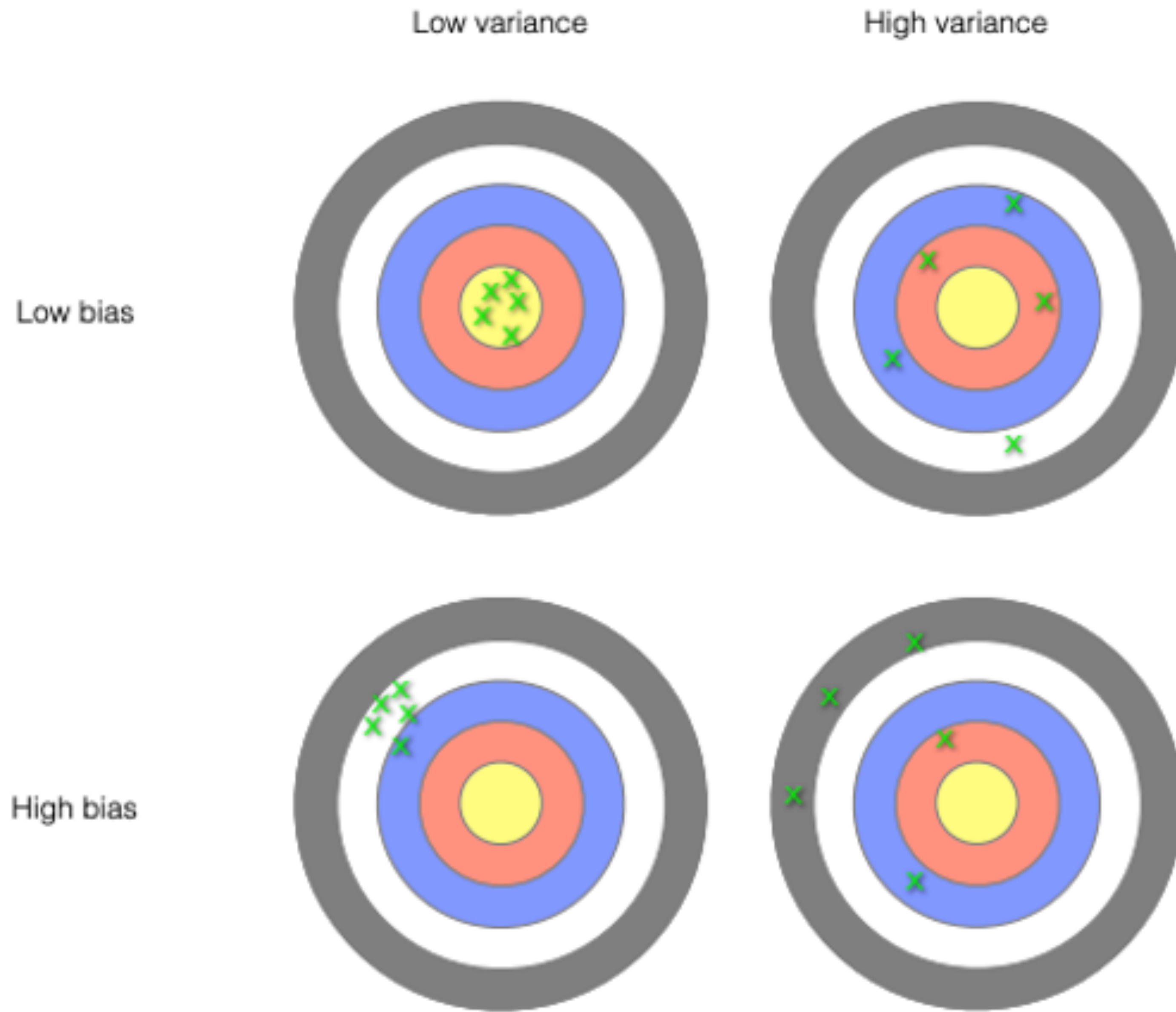**Bias:** Error due to mis-specifying the relationship between input and output

*Too few parameters, or the wrong kinds*

**Variance:** Error due to sensitivity to random fluctuations in the training data. If you train on different data, do you get radically different predictions?

*Too many parameters*

Low variance     High variance

Low bias

High bias

Slide content credit to David Bamman

Image from Flach 2012

# Regression for Social Sciences

Regression analysis is a very useful tool for social sciences

+ Understand the relationship between variables, adjusting for other potential confounders

+ Predict the value of one variable based on others

# In Other Terminology

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

Intersect

| Dependent Variable | $=$ | Independent Variable | $+$ | Independent Variable |
|---|---|---|---|---|

# How good is the Fit?

Mean squared error (MSE)   $\dfrac{1}{N}\displaystyle\sum_{i=1}^{N}(\hat{y}_i - y_i)^2$

Mean absolute error (MAE)   $\dfrac{1}{N}\displaystyle\sum_{i=1}^{N}|\hat{y}_i - y_i|$

# How good is the "fit"?

Sum of the squares total (SST): total variability about the mean

$$\sum (Y - \bar{Y})^2$$

Sum of the squared error (SSE): variability about the regression line

$$\sum (Y - \hat{Y})^2$$

Sum of the squares due to regression (SSR): total variability that is explained by the model

$$\sum (\hat{Y} - \bar{Y})^2$$

# Coefficient of Determination $r^2$

The proportion of the variability explained by regression model

$$r^2 = \frac{\text{SSR}}{\text{SST}}$$

# Recommendations for Building Regression Models

A high $r^2$ is desired with a reasonable set of variables

When more variables get added to the model, $r^2$ usually increases.

Thus, adjusted $r^2$ is often used to account for the number of variables
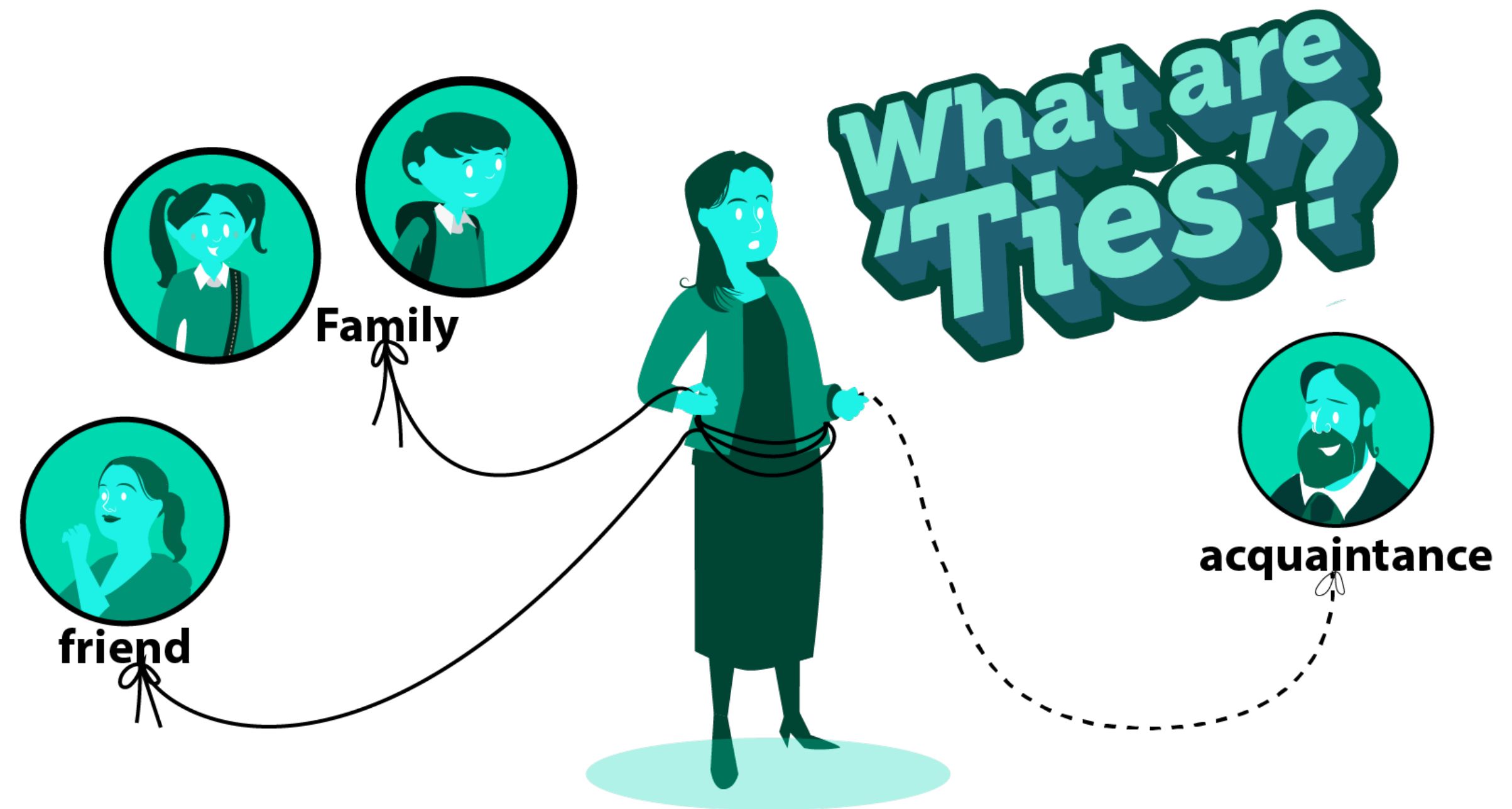
Independent variables might contain **duplicated** information

*Colinear* if two variables are correlated

*Multicolinearity* if more than two variables are correlated - this will make the interpretation of regression coefficient problematic

# Let's predict tie strength on Facebook

1. **Why** is this a regression task?
2. **What** is tie strength?
3. How can we get the **ground truth**?
4. How to get **data**?
5. How can we **evaluate** it?
6. Does the system really **work**?



https://murraydare.co.uk/marketing-theory/strong-weak-ties
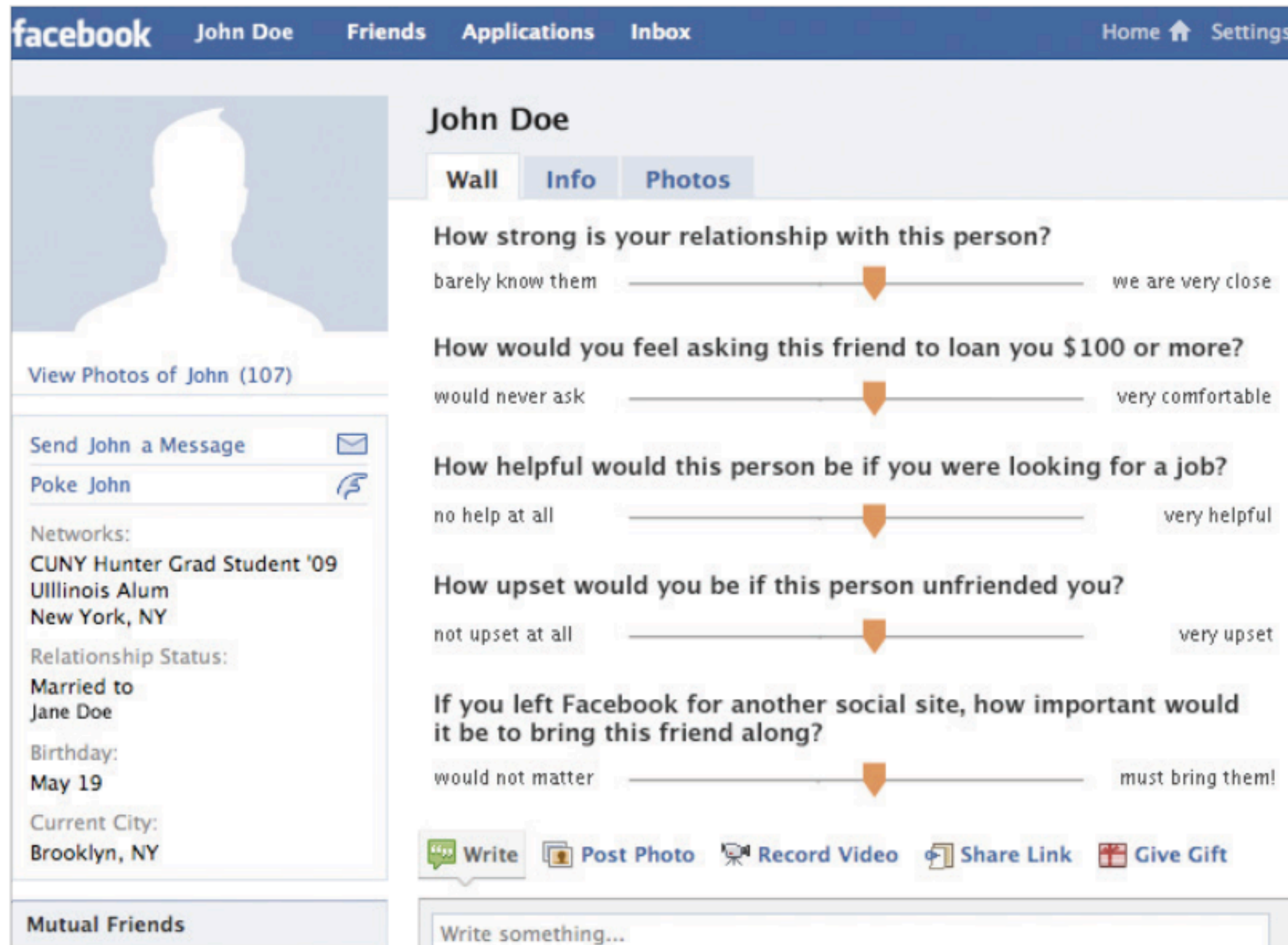
# Let's predict tie strength on Facebook

Mark Granovetter introduced the concept of **tie strength** in1973
       "**The Strength of Weak Ties**"

The strength of a tie is a (probably linear) combination of the amount of
time, the emotional intensity, the intimacy (mutual confiding), and the
reciprocal services which characterize the tie

Gilbert, Eric, and Karrie Karahalios. "Predicting tie strength with social media." In Proceedings of the SIGCHI conference on human factors in computing systems, pp. 211-220. 2009.

# Let's predict tie strength on Facebook

# Let's predict tie strength on Facebook

What features we could use to predict self-reported tie strength?

| Predictive Intensity Variables | Distribution | Max |
|---|---|---|
| Wall words exchanged | | 9549 |
| Participant-initiated wall posts | | 55 |
| Friend-initiated wall posts | | 47 |
| Inbox messages exchanged | | 9 |
| Inbox thread depth | | 31 |
| Participant's status updates | | 80 |
| Friend's status updates | | 200 |
| Friend's photo comments | | 1352 |

| Intimacy Variables | | |
|---|---|---|
| Participant's number of friends | | 729 |
| Friend's number of friends | | 2050 |
| Days since last communication | | 1115 |
| Wall intimacy words | | 148 |
| Inbox intimacy words | | 137 |
| Appearances together in photo | | 73 |
| Participant's appearances in photo | | 897 |
| Distance between hometowns (mi) | | 8182 |
| Friend's relationship status | 6% engaged 30% single | 32% married 30% in relationship |

| Duration Variable | | |
|---|---|---|
| Days since first communication | | 1328 |

| Reciprocal Services Variables | | |
|---|---|---|
| Links exchanged by wall post | | 688 |
| Applications in common | | 18 |

| Structural Variables | | |
|---|---|---|
| Number of mutual friends | | 206 |
| Groups in common | | 12 |
| Norm. TF-IDF of *interests* and *about* | | 73 |

| Emotional Support Variables | | |
|---|---|---|
| Wall & inbox positive emotion words | | 197 |
| Wall & inbox negative emotion words | | 51 |

| Social Distance Variables | | |
|---|---|---|
| Age difference (days) | | 5995 |
| Number of occupations difference | | 8 |
| Educational difference (degrees) | | 3 |
| Overlapping words in *religion* | | 2 |
| Political difference (scale) | | 4 |

# Let's predict tie strength on Facebook



**How strong?**
- +R$_i$: 0.37
- +D$_i$: 0.5
- +N(i): 0.53

**Loan $100?**
- +R$_i$: 0.35
- +D$_i$: 0.52
- +N(i): 0.54

**Helpful for job?**
- +R$_i$: 0.24
- +D$_i$: 0.38
- +N(i): 0.39

**Upset if unfriended?**
- +R$_i$: 0.27
- +D$_i$: 0.4
- +N(i): 0.42

**Bring friend to new site?**
- +R$_i$: 0.35
- +D$_i$: 0.46
- +N(i): 0.48

The model's Adjusted R2 values for all five dependent variables, broken down by the model's three main terms.

Modeling interactions between tie strength dimensions results in a substantial performance boost.

The model performs best on Loan $100? and How strong?, the most general question

64

# Let's predict tie strength on Facebook

| Top 15 Predictive Variables | β | F | p-value |
|---|---|---|---|
| Days since last communication | -0.76 | 453 | < 0.001 |
| Days since first communication | 0.755 | 7.55 | < 0.001 |
| Intimacy × Structural | 0.4 | 12.37 | < 0.001 |
| Wall words exchanged | 0.299 | 11.51 | < 0.001 |
| Mean strength of mutual friends | 0.257 | 188.2 | < 0.001 |
| Educational difference | -0.22 | 29.72 | < 0.001 |
| Structural × Structural | 0.195 | 12.41 | < 0.001 |
| Reciprocal Serv. × Reciprocal Serv. | -0.19 | 14.4 | < 0.001 |
| Participant-initiated wall posts | 0.146 | 119.7 | < 0.001 |
| Inbox thread depth | -0.14 | 1.09 | 0.29 |
| Participant's number of friends | -0.14 | 30.34 | < 0.001 |
| Inbox positive emotion words | 0.135 | 3.64 | 0.05 |
| Social Distance × Structural | 0.13 | 34 | < 0.001 |
| Participant's number of apps | -0.12 | 2.32 | 0.12 |
| Wall intimacy words | 0.111 | 18.15 | < 0.001 |

The fifteen predictive variables with highest standardized beta coefficients.

The two Days since variables have large coefficients because of the difference between never communicating and communicating once.

The utility distribution of the predictive variables forms a power-law distribution: **with only these fifteen variables, the model has over half of the information it needs to predict tie strength.**

# Let's predict tie strength on Facebook

Don't forget error analysis

**rating: 0.96; prediction: 0.47**

This friend is very special. He and I attended the same high school, we interacted a lot over 3 years and we are very very close. We trust each other. My friend are I are still interacting in ways other than Facebook such as IM, emails, phones. Unfortunately, that friend and I rarely interact through Facebook so I guess your predictor doesn't have enough information to be accurate.

**rating: 0; prediction: 0.44**

I don't know why he friended me. But I'm easy on Facebook, because I feel like I'm somehow building (at least a miniscule amount of) social capital, even when I don't know the person. We went to the same high school and have a few dozen common friends. We've never interacted with each other on Facebook aside from the friending.

**rating: 0.6; prediction: 0.11**

Ah yes. This friend is an old ex. We haven't really spoken to each other in about 6 years, but we ended up friending each other on Facebook when I first joined. But he's still important to me. We were best friends for seven years before we dated. So I rated it where I did (I was actually even thinking of rating it higher) because I am optimistically hoping we'll recover some of our "best friend"-ness after a while. Hasn't happened yet, though.