



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Causal Inference: Text and NLP

Logistics

- HW 3 on causal inference has been released
 - Deadline is Friday
- Midterm Exam
 - In class next Wednesday
 - Includes all material through Wednesday 3/4 (including homeworks)
 - Sample problems released on Piazza
 - Review session Monday 3/9

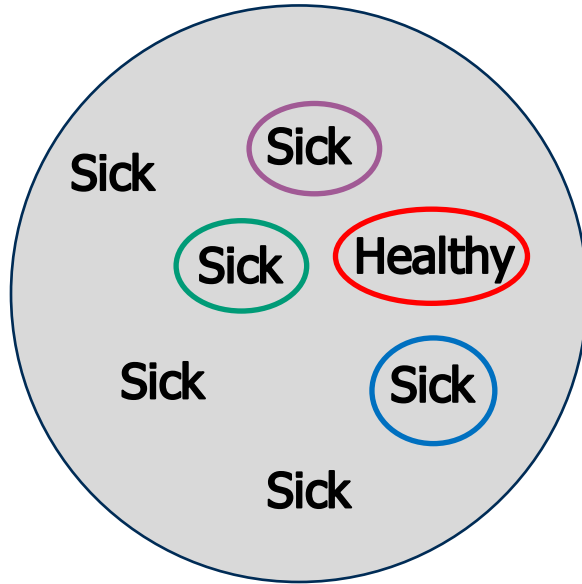
Recap

- Methods for adjusting for confounders
 - Regression
 - Matching
 - Propensity scores (matching, weighting, and stratification)

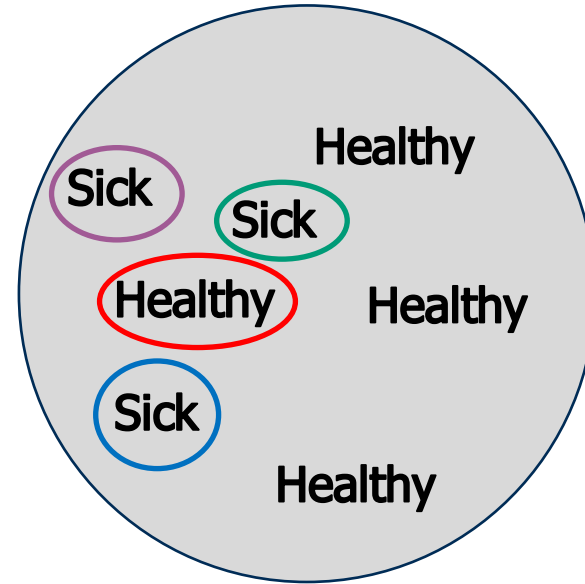
Today:

- Additional notes about when to do adjustments (confounder vs colliders vs mediators)
- Causal inference with text
 - Overview
 - Adjusting for text as confounders (or mediators)
 - Drawing from causal inference to improve NLP models

Direct Matching



Took Medicine



Didn't Take Medicine

Propensity Score

- X might be high dimensional, is it necessary to match on (or more generally adjust for) all of X?
- Define the *propensity score* as the probability of receiving treatment, given confounders:
 - $e(X) = P(T = 1 \mid X = x)$

Propensity Score Theorem

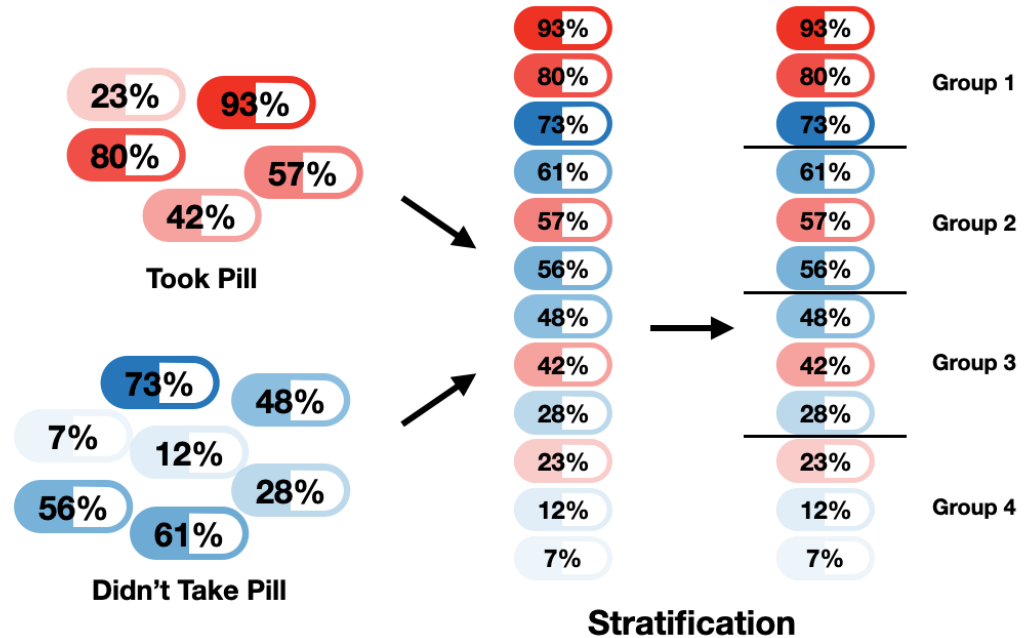
- Given positivity, unconfoundedness given X implies unconfoundedness given the propensity score $e(X)$

$$Y(1), Y(0) \perp T \mid X \Rightarrow Y(1), Y(0) \perp T \mid e(X)$$

- When we are adjusting for X , we can swap in $e(X)$ instead
- We don't typically actually know $e(X)$ but we can estimate it from the data (e.g., training a model to predict T from X)

Propensity Stratification

- Stratify (bucket) individuals into mutually exclusive subsets with the same propensity score
- 5 subsets (quintiles) is a common choice
- Compute estimand for each strata and then pool them (weighted equally for quintiles)



Rosenbaum P.R., Rubin D.B. Reducing bias in observational studies using subclassification on the propensity score. Journal of the American Statistical Association. 1984;79:516–524
Image: <https://towardsdatascience.com/propensity-score-5c29c480130c>

IPW (Inverse probability weighting)

$$w_i = \frac{T_i}{e(X_i)} + \frac{1 - T_i}{1 - e(X_i)}$$

- Define weight: inverse estimate of the probability of the treatment that the individual actually received

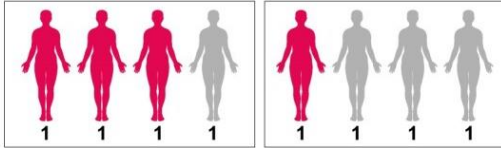
ATE = weighted avg. of treated individuals – weighted avg. of untreated individuals

- [Also called IPTW, inverse probability of treatment weighting]

IPW (Inverse probability weighting)

treatment control

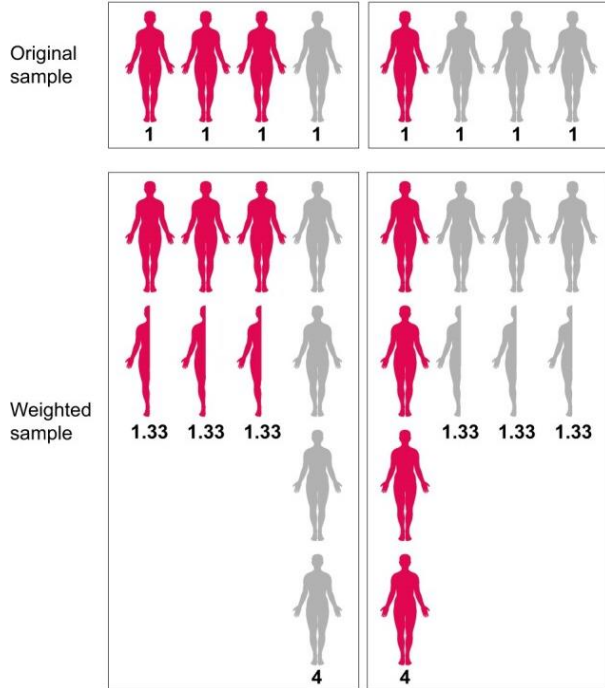
Original
sample



- Setup: Red = felt sick
 - $\frac{3}{4}$ people who felt sick took medicine
 - $P(\text{taking medicine} \mid \text{feel sick}) = 0.75$
 - $P(\text{no medicine} \mid \text{feel sick}) = 0.25$
- Weights:
 - Took medicine, felt sick: $1/0.75 = 1.333$
 - No medicine, felt sick: $1/.25 = 4$
 - [similarly calculate weights for people who didn't feel sick]
- When we apply weights, we've balanced feeling sick with not feeling sick

IPW (Inverse probability weighting)

treatment **control**



- We're creating "pseudeopopulations"
- Similar concept: when collecting survey data, you may upweight respondents of particular demographics to match population statistics

How do propensity adjustment methods compare?

- Often choice depends on what model is best suited to data and analysis
- Several studies have demonstrated that propensity score **matching** eliminates a greater proportion of the systematic differences than **stratification** (Austin, 2009a; Austin, Grootendorst, & Anderson, 2007; Austin & Mamdani, 2006)
- In some settings propensity score **matching** and **IPTW** were shown to be comparable; in others propensity score matching was slightly better (Austin, 2009a)

Regression vs. Matching?

- “matching methods should not be seen in conflict with regression adjustment and in fact the two methods are complementary and best used in combination”
 - E.g. you could stratify based on propensity scores and then use regression adjustment with each stratum to adjust for lingering differences
- “matching methods highlight areas of the covariate distribution where there is not sufficient overlap between the treatment and control groups, such that the resulting treatment effect estimates would rely heavily on extrapolation”
- “methods such as linear regression adjustment can actually increase bias in the estimated treatment effect when the true relationship between the covariate and outcome is even moderately non-linear”



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

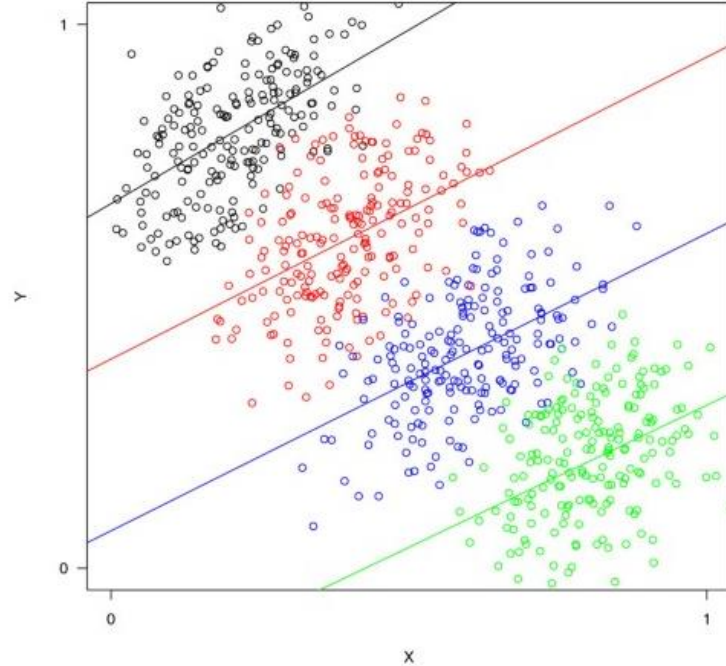
Some additional notes

Double Machine Learning

- General framework for estimating causal effects using ML (random forests, lasso or post-lasso, neural nets, boosted regression trees, and various hybrids and ensembles of these methods)
- Available in Python and R packages:
 - <https://github.com/DoubleML>

Mixed Effects Regression Models

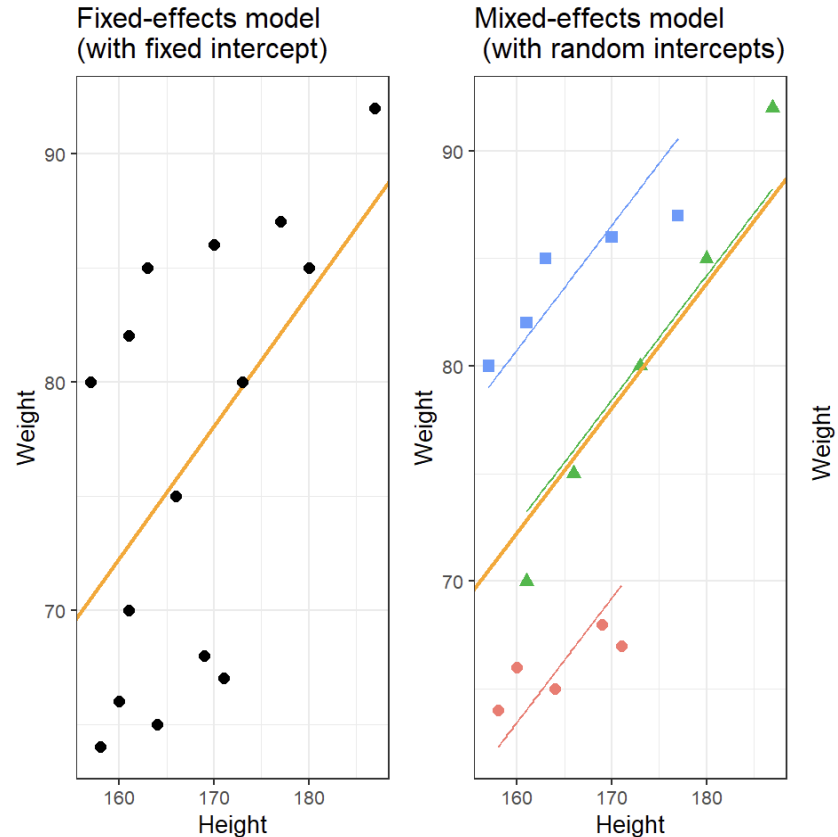
- We discussed regression adjustment for confounders
- When data is hierarchical / non-independent we need a better regression model
- E.g. you examine if dosage of medicine affects fevers
- Your data is from hospitals in different countries where underlying health conditions that affect baseline health
- Recall Simpson's Paradox



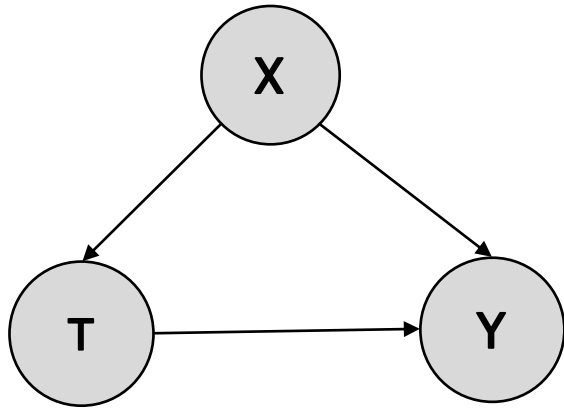
- Data looks negatively correlated overall
- Subsetting data shows positive correlations 16

Mixed Effects Regression Models

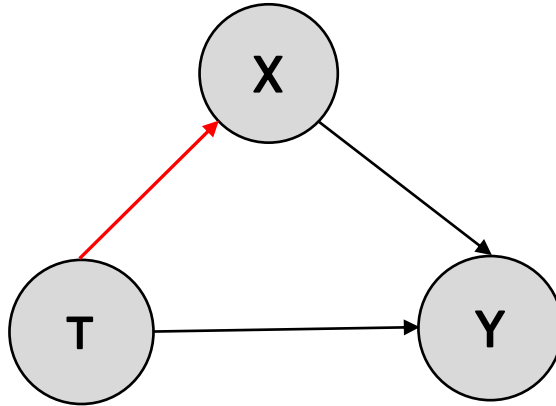
- We can account for differences across subgroups by allowing subgroups to have different parameters (e.g. different intercepts in linear regression)
- Subgroup is a *random* effect
- Dosage is a *fixed* effect



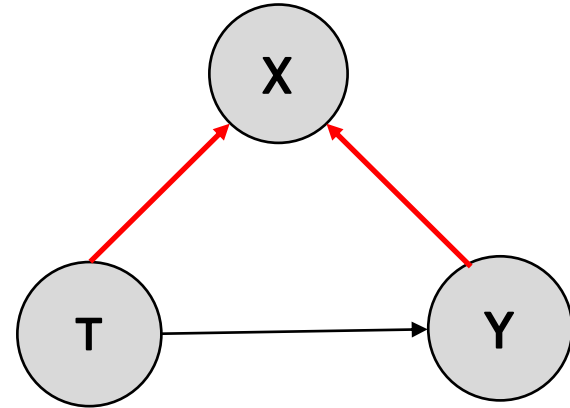
Confounders vs. Mediators vs. Colliders



confounder

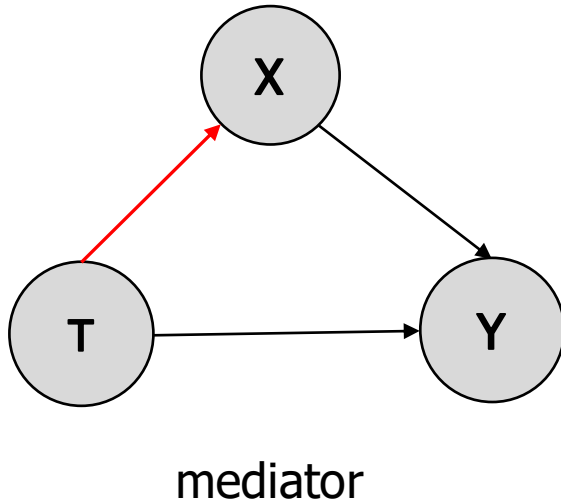


mediator



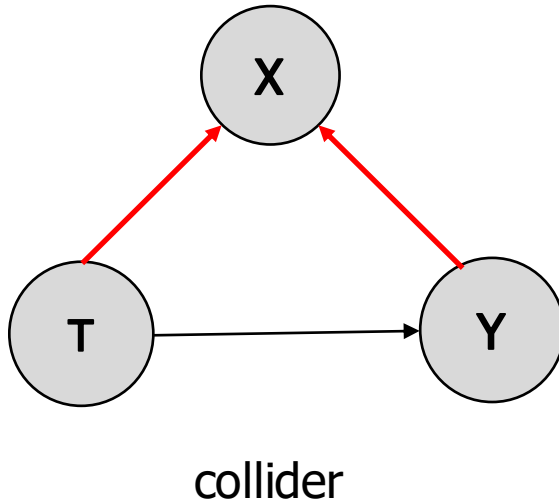
collider

Confounders vs. Mediators vs. Colliders



- Example:
 - Estimating if gender has an effect on social media likes
 - Gender (T) influences the topic of posts (X)
 - Topic of posts (X) and gender (T) influence number of likes (Y)
- If we adjust for X, we may be removing some of the effect
- We may still choose to adjust for X if we specifically want to capture the direct effect and not the indirect effect
- We may want to separate out direct and indirect effects in a mediation analysis

Confounders vs. Mediators vs. Colliders



- Example:
 - Studying if cold medicine reduces fever
 - Medicine (T) influences fever (Y) and if you went out for ice cream (X) [because medicine tastes bad]
 - People who don't have fevers (Y) eat ice cream (X)
 - If you condition on X (e.g. restrict data to people who ate ice cream), you're selecting for people who were feeling better in your control group → you find that medicine causes fever
- If we adjust for X, we are adding bias to our estimator!

Confounders vs. Mediators vs. Colliders

i	T	Y	X
1	1	0	Ice cream
2	1	fever	Ice cream
3	1	0	Ice cream
4	1	0	Ice cream
5	0	0	Ice cream
6	0	0	Ice cream
7	0	fever	no
8	0	fever	no

- Example:
 - Studying if cold medicine reduces fever
 - Medicine (T) influences fever (Y) and if you went out for ice cream (X) [because medicine tastes bad]
- 1/4 of medicine takers still has a fever
- 0/2 of non-medicine takers still has fever



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Causal Inference in text: Overview

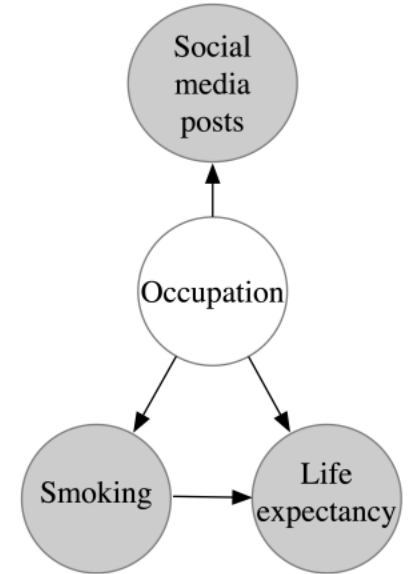
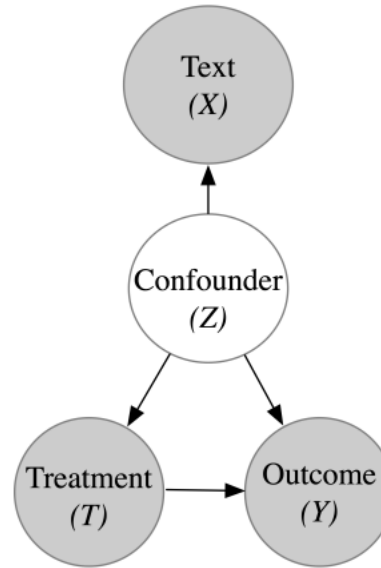
What characteristics distinguish text from other data types?

- Text is high dimensional
 - Overfitting, violations of positivity
- Compared to other high dimensional data:
 - Text can be read and evaluated by humans
 - Designing meaningful representations of text is an open problem

Text as confounders

- Text data could either:
 - (a) serve as a surrogate for potential confounders
 - (b) the language of text itself could be a confounder

Example: the linguistic content of social media posts (confounder) could influence censorship (treatment) and future posting rates (outcome)



Text as treatment or outcome

- Do Wikipedia articles contain gender bias?
 - Treatment: Perceived gender
 - Outcome: Article text
 - Confounders/Mediators: Perceived characteristics other than gender

- Does a celebrity's social media posts cause them to gain followers?
 - Treatment: The social media posts
 - Outcome: Follower counts
 - Confounders/Mediators: Changes in social media usage, current events



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Adjusting for text as confounders

Two similar approaches

- Topic Inverse Regression Matching
 - Roberts, Margaret E., Brandon M. Stewart, and Richard A. Nielsen. "Adjusting for confounding with text matching." *American Journal of Political Science* 64.4 (2020): 887-903.
- "Causally sufficient" embeddings
 - Veitch, Victor, Dhanya Sridhar, and David Blei. "Adapting text embeddings for causal inference." *Conference on Uncertainty in Artificial Intelligence*. PMLR, 2020.

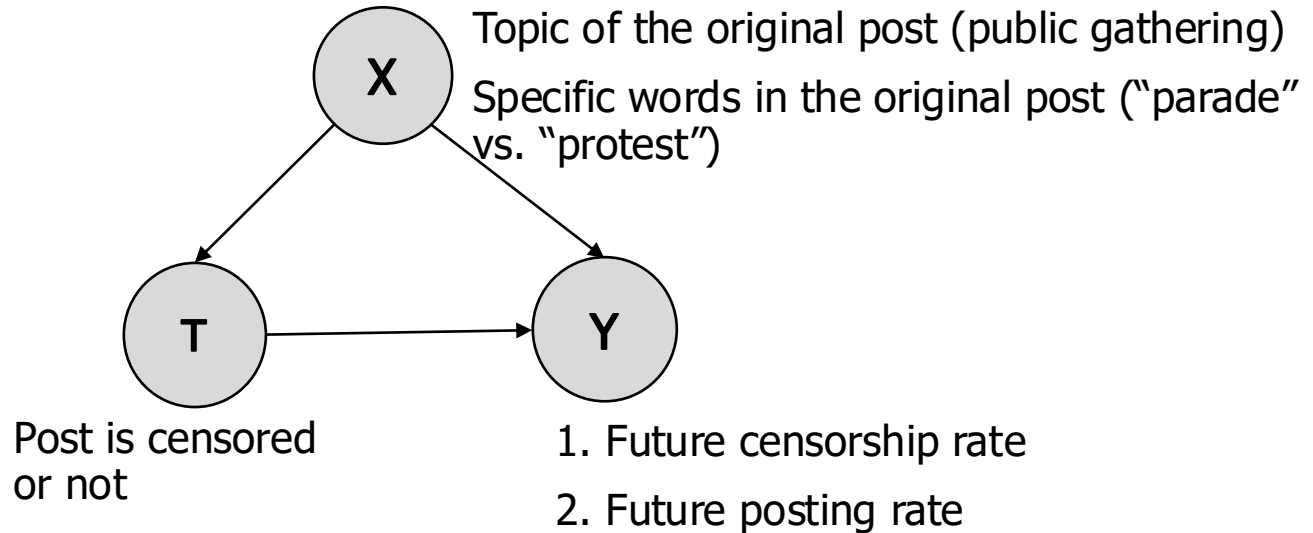
Adjusting for text as confounders: Topic Inverse Regression Matching

- Key ideas:
 - Matching (remember: direct or propensity) is a good approach for adjusting for text as confounder because analysts can manually evaluate the quality of the adjustment by comparing the matched treatment and control text
 - Most use cases what we need to match on are topics (as opposed to sentiment, punctuation, word order, etc). We also may care about individual words
 - We need to match on aspects of the text that are predictive of treatment (definition of confounders)

Example application: Effects of censorship in Chinese social media

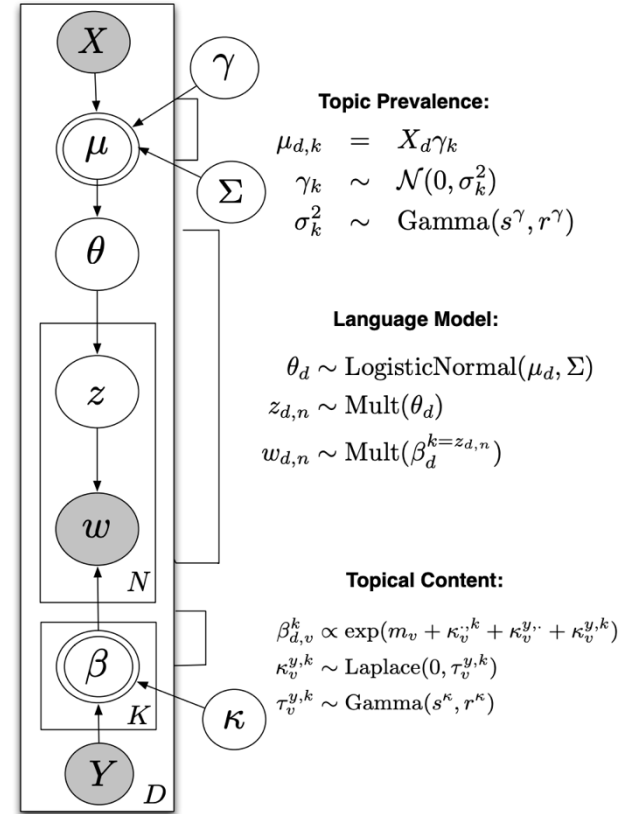
- Research questions:
 - 1. "Is censorship completely determined by the text of a particular post, or does censorship become more targeted toward users based on their previous censorship history?"
 - 2. Does having a post censored cause people to post less in the future?

Example application: Effects of censorship in Chinese social media



Topic Inverse Regression Matching using STM

- For bag-of-words representation W , define a function $g(W)$ to create a low-dimensional estimate that captures topic and word differences that relate to treatment assignment
- Primary model for text representations: *structured topic model (STM)*
- LDA-style topic model that allows flexible inclusion of covariates



Step	Rationale
1. Estimate a structural topic model including the treatment vector as a content covariate.	Reduces dimension of the text
2. Extract each document's topics calculated as though treated (part of $g(\mathbf{W})$).	Ensures semantic similarity of matched texts
3. Extract each document's projection onto the treatment variable (part of $g(\mathbf{W})$).	Ensures similar treatment probability of matched texts
4. Use a low-dimensional matching method to match on $g(\mathbf{W})$ and estimate treatment effects using the matched sample.	Standardizes matching

Example application: Effects of Censorship on Chinese social media

- Research questions:
 - 1. “Is censorship completely determined by the text of a particular post, or does censorship become more targeted toward users based on their previous censorship history?”
 - 2. Does having a post censored cause people to post less in the future?
- Methods:
 - Use TIRM to identify pairs of nearly identical social media posts written by nearly identical users, where one is censored and the other is not
 - Examine subsequent posting and censorship rates of each user

Example application: Effects of Censorship on Chinese social media

- Results:
 - Having a post censored increases the probability of future censorship significantly
 - It does not decrease number of future posts by the censored user
- Conclusions:
 - Option 1: algorithmic targeting of censorship, where social media users are more likely to be censored after censorship because they are flagged by the censors
 - Option 2: social media users may chafe against censorship and respond by posting increasingly sensitive content that is more likely to be censored

A different method: develop “causally sufficient” text embeddings

- Text is high dimensional and data is finite: difficult to fit models directly to text
- Instead, “reduce the text to a low-dimensional representation that suffices for causal identification and enables efficient estimation from finite data.”
- Two key ideas:
 - Supervised dimensionality reduction: we don’t need the full text, causal inference only requires the parts of text that are predictive of the treatment and outcome
 - Efficient language modeling: design representations of text to dispose of “linguistically irrelevant information”, presumed to also be “causally irrelevant”

General approach: develop “causally sufficient” text embeddings

- Start with a language model (BERT) and modify it to produce 3 outputs:
 - 1) document-level embeddings
 - 2) a map from the embeddings to treatment probability
 - 3) a map from the embeddings to expected outcomes for the treated and untreated
 - [(2) and (3) are small added neural networks on the original model]
- [They also do a variant based on a topic model]

General approach: develop “causally sufficient” text embeddings

- Train model to predict outcome, treatment, and with language-modeling objective (e.g. to learn meaningful text representations)

$$\begin{aligned} L(\mathbf{w}_i; \xi, \gamma) &= (y_i - \tilde{Q}(t_i, \lambda_i; \gamma))^2 \longrightarrow \text{Outcome} \\ &+ \text{CrossEnt}(t_i, \tilde{g}(\lambda_i; \gamma)) \longrightarrow \text{Treatment} \\ &+ L_U(\mathbf{w}_i; \xi, \gamma). \longrightarrow \text{Language modeling} \end{aligned}$$

- To compute average treatment effect, plug estimated embeddings, propensity scores, and conditional outcomes into a downstream estimator

Evaluation

- Two settings:
 - Peer-reviewed journal articles: Causal effect of including a theorem on paper acceptance.
 - Treatment: the word “theorem” occurs in the paper
 - Confounder: article abstract (subject of the paper)
 - Outcome: accept/reject
 - Effect of gender on Reddit popularity
 - Treatment: “male” label
 - Mediator: Post text (topic or style)
 - Outcome: Popularity score

How can we use this data for *evaluation* rather than analysis?

Evaluations

- Simulated data:
 - Use real confounders and treatments
 - Simulate outcomes (so we know the “true” causal effect)
- Their findings:
 - 1) Yes, language modeling helps recover simulated effects
 - 2) Yes, supervised dimensionality helps
 - 3) Their proposed models C-BERT and C-ATM outperform alternatives



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Drawing from Causal Inference to Improve NLP models

Drawing from Causal Inference to Improve NLP models

- ML in general typically captures associates, not causal effects
- Models are prone to overfitting, exploit spurious correlations in the data
 - E.g. train a model to identify photos of dogs from cats; Model learns that dogs always have collars



→ "DOG"

Drawing from Causal Inference to Improve NLP models

- ML in general typically captures associates, not causal effects
- Models are prone to overfitting, exploit spurious correlations in the data
 - E.g. train a model to identify photos of dogs from cats; Model learns that dogs always have collars
- Maybe by drawing from causal inference we can train models to ignore these spurious correlations, especially for tasks where it's hard to collect good training data
- Case study: drawing from causal inference to detect *subtle gender bias*

Need to develop new models

- Our goal: detect subtle gender biases like microaggressions, objectifications, and condescension in 2nd-person text
 - “Oh, you work at an office? I bet you’re a secretary”
 - “Total tangent I know, but you’re gorgeous”
- Current classifiers that detect hate speech, offensive language, or negative sentiment cannot detect these comments
- [Note: focus on binary gender]

Naive Approach: Supervised Classification



I like Bob, but you're hot, so kick his butt

Like · Reply ·



Thanks so much **Ma'am!**

Like · Reply ·



I'd vote for you if I lived in **Massachusetts**

Like · Reply ·



...a good way to celebrate **Title IX**, too!

Like · Reply ·



Naive Approach: Supervised Classification

I like Bob, but you're hot, so kick his butt

Like · Reply ·

Thanks so much **Ma'am!**

Like · Reply ·

I'd vote for you if I lived in **Massachusetts**

Like · Reply ·

...a good way to celebrate **Title IX**, too!

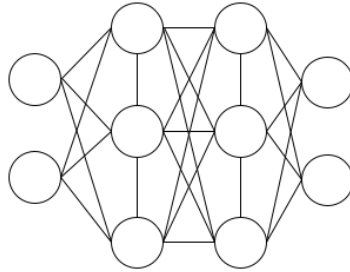
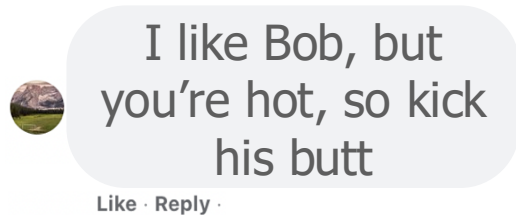
Like · Reply ·



Problem: Biases are *subtle, implicit, and context-dependent*

Proposed approach: Comments contain gender bias if they are highly predictive of gender

- Train a classifier that predicts the gender of the person the text is addressed to
- If the classifier makes a prediction with high confidence, the text likely contains bias



→ Addressed to **Man**

→ Addressed to **Woman**

If a comment is very likely to be addressed to a woman, and is very unlikely to be addressed to a man, it probably contains gender bias.

Challenge: Text main contain *confounds* that are predictive of gender, but not indicative of gender bias



I like Bob, but you're hot, so kick his butt

Like · Reply ·



Thanks so much
Ma'am!

Like · Reply ·



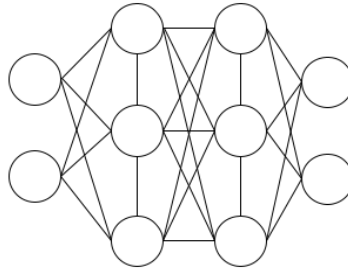
I'd vote for you if I lived in **Massachusetts**

Like · Reply ·



...a good way to celebrate **Title IX**, too!

Like · Reply ·



→ Addressed to **Woman**

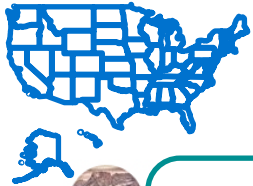
→ Addressed to **Woman**

→ Addressed to **Woman**

→ Addressed to **Woman**

Challenge: Text main contain *confounds* that are predictive of gender, but not indicative of gender bias

- Overtly gendered words
- Preceding context in the conversation
- Traits of people (other than gender) in the conversation



Saturday is the 40th anniversary of **Title IX**...

Like · Reply ·



...a good way to celebrate Title IX, too!

Like · Reply ·



I'd vote for you if I lived in Massachusetts

Like · Reply ·



Bob and I join Bill Hemmer on America's Newsroom to discuss whether or not...

Like · Reply ·



I like Bob, but you're hot, so kick his butt

Like · Reply ·



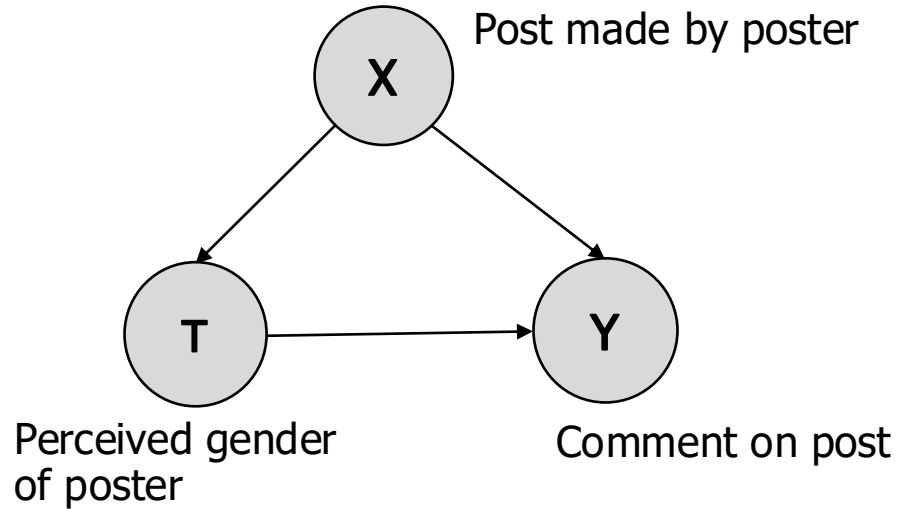
Thanks so much Ma'am!

Like · Reply ·

A note on causal set-up

- We're not really doing causal inference: we are trying to build a classifier to detect microaggressions, not draw conclusions about the state of the world:
 - "confounds": spurious correlations in our data (not necessarily "confounders")
- Some of these factors that we don't want the model to learn are confounding variables

A note on causal set-up



[Note: we have text as an outcome and as a confounder]

Preceding context is an *observed* confounding variables

Writer_Gender: F



Saturday is the 40th anniversary of **Title IX**! I'm celebrating with a Sat morning run - ladies please respond below if you want to join

Like · Reply ·



Wish I could ! Already have plans for a bike ride and breakfast with some awesome ladies - a good way to celebrate **Title IX**, too!

Like · Reply ·



Would love to!

Like · Reply ·

Writer_Gender: M



Any deal with **Iran** — a nation that the United States cut off diplomatic ties with 35 years ago — must protect America's interests at home and abroad.

Like · Reply ·



Iran might be a free, democratic nation today, if not for decades of American interference.

Like · Reply ·

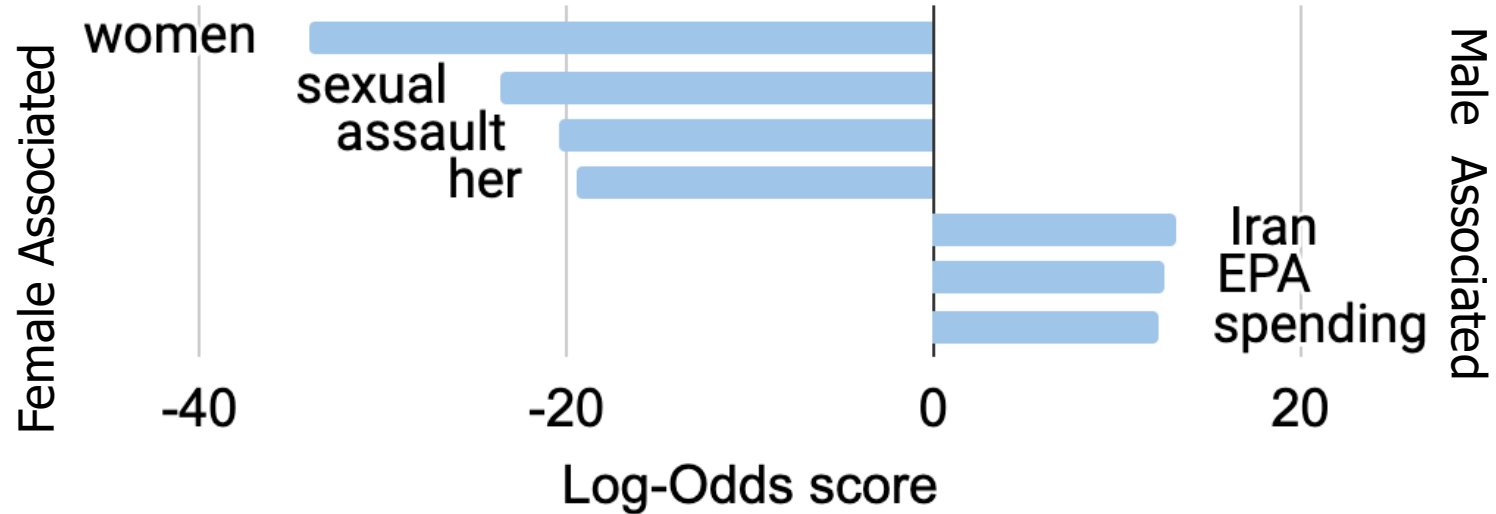


That's for sure! Worst deal he could make! We can't trust **Iran** and America knows it !!!!!

Like · Reply ·

Key problem: Men and women post different content, which is reflected in their replies

Preceding context is an *observed* confounding variables



Propensity matching for *observed* confounding variables

~~Writer_Gender: F~~

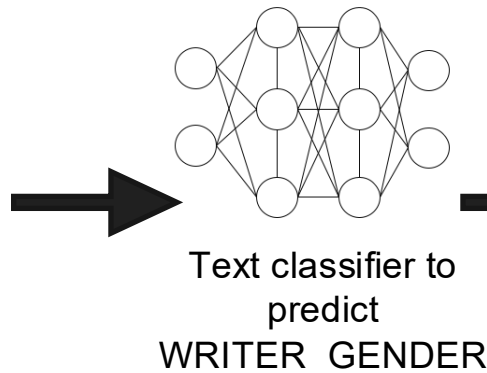
~~Saturday is the 40th anniversary of Title IX! I'm celebrating with a Sat morning run - ladies please respond below if you want to join.~~

Writer_Gender: M

Any deal with Iran — a nation that the United States cut off diplomatic ties with 35 years ago — must protect America's interests at home and abroad.

Writer_Gender: F

My overriding concern is whether or not the agreement is in the national security interest of the United States. Iran must be blocked from proceeding any further towards developing a nuclear weapon.



$$|e_i - e_l| \geq c \forall l$$

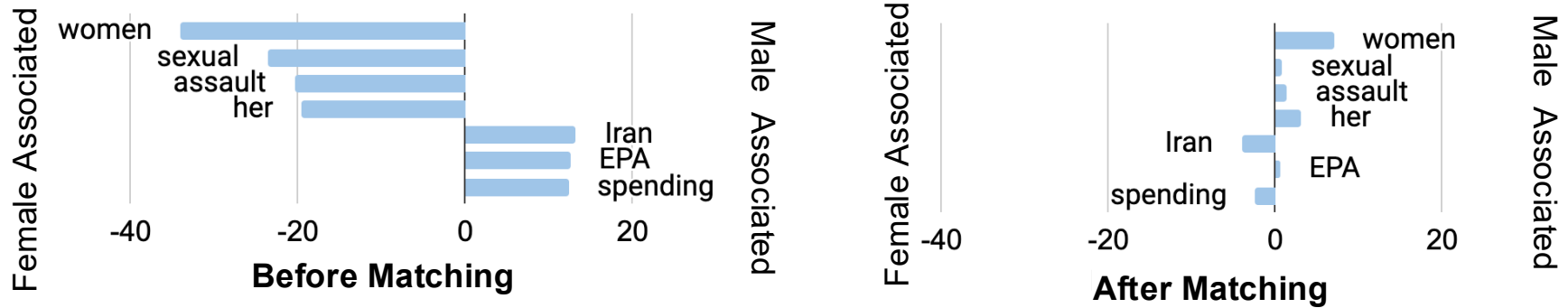
$$e_i = P(W.Gender_i = F | Post_i) \approx 0.91$$

$$e_j = P(W.Gender_j = F | Post_j) \approx 0.33$$

$$e_k = P(W.Gender_k = F | Post_k) \approx 0.32$$

$$\operatorname{argmin}_j |e_k - e_j|$$

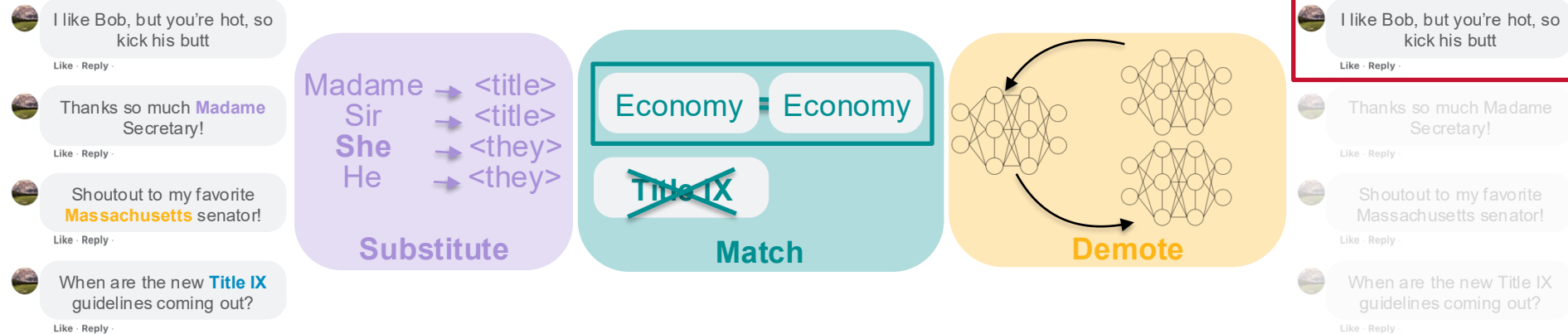
Propensity matching for *observed* confounding variables



Propensity matching breaks associations between gender and context in the training data

Proposed Model: Comments contain bias if they are highly predictive of gender *despite confound control*

- Substitute overt indicators
- Balance observed confounders through propensity matching
- Demote latent confounders through adversarial training



Self-reported microaggressions

	Public Figs		Politicians	
	F1	Acc.	F1	Acc.
base	61.3	57.3	48.1	64.2
+demotion	62.2	57.9	53.7	61.5
+match	38.9	55.9	46.9	50.7
+match+dem.	50.9	57.0	56.9	49.9
Random	46.0	49.8	-	-
Class Random	42.1	48.3	-	-

- Models are not trained at all for this task; they are only trained for gender-of-addressee prediction, but they still perform better than chance

Findings: characteristics of bias against women politicians

- Influential words:
 - Competence and domesticity
 - 'Force', 'situation', 'spouse', 'family', 'love'
- Examples:
 - "DINO I hope another real Democrat challenges you next election"
 - "I did not vote for you and have no clue why anyone should have. You do not belong in politics"

Findings: characteristics of bias against women

- Influential words:
 - Appearance and sexualization
 - 'beautiful', 'love', 'sexo'
- Examples:
 - "Total tangent I know but, you're gorgeous."
 - "I like Bob, but you're hot, so kick his butt."

Recap

- Overview:
 - Text as confounders, treatment, or outcome
- Text as confounders
 - Topic modeling and language modeling to adjust for text
- Drawing from causal inference to improve NLP models
 - Applying ideas from causal inference to model development
- Next class:
 - Network Analysis