# Political Belief Polarization

## What it is, Why it exists, How we can help

Zhiping (Patricia) Xiao
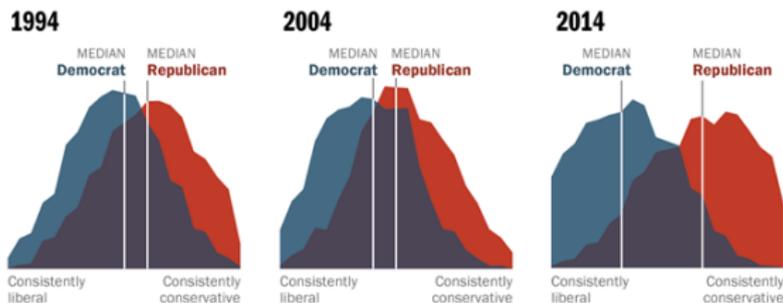
University of California, Los Angeles

# Outline

UCLA

# Background Introduction

**Democrats and Republicans More Ideologically Divided than in the Past**

*Distribution of Democrats and Republicans on a 10-item scale of political values*

Source: 2014 Political Polarization in the American Public
Notes: Ideological consistency based on a scale of 10 political values questions (see Appendix A). The blue area in this chart represents the ideological distribution of Democrats; the red area of Republicans. The overlap of these two distributions is shaded purple. Republicans include Republican-leaning independents; Democrats include Democratic-leaning independents (see Appendix B).

**PEW RESEARCH CENTER**

Polarization: ideological distance between political parties and candidates.

Pernicious Polarization: when (1) society divide into separated camps (2) disagreement on ground-truth.

UCLA

[1]Image from: https://www.pewresearch.org/

Divergence in beliefs about vaccines:

- ▶ It was in year 1998 when The Lancet, a medical journal, published a study linking the MMR (measles, mumps, and rubella) vaccine to autism.

- ▶ Although this study has since been retracted and its results refuted, it is still a rallying cry for the modern anti-vaccination movements.

- ▶ The long-lasting influence still exists today.

UCLA

- ▶ Studying Political Bias via Word Embeddings (WWW'20)
- ▶ Political audience diversity and news reliability in algorithmic ranking (Nature Human Behaviour, 2022)
- ▶ Belief polarization in a complex world: A learning theory perspective (PNAS May 11, 2021) (with Appendix)
- ▶ Encouraging Moderation: Clues from a Simple Model of Ideological Conflict (Phys. Rev. Lett, 2012)
- ▶ Reasoning about Political Bias in Content Moderation (AAAI'20 Special Programs Section, Sister Conference Track)

UCLA

# Analyzing the Polarity

Studying Political Bias via Word Embeddings

▶ polarization is reflected in vocabulary

Political audience diversity and news reliability in algorithmic ranking

▶ high quality news attract diverse audiences

Belief polarization in a complex world: A learning theory perspective

▶ exposing to similar content can still end up in polarization

▶ preference towards simplicity encourages polarization

▶ in theory, carefully-designed bias ($\tilde{\mathcal{D}} = (1 - \alpha)\mathcal{D} + \alpha\mathcal{P}$) on data distribution can help reduce polarization

Studying Political Bias via Word Embeddings

- ► Motivation: word-level political bias.
- ► Related Work: Modeling gender bias (e.g. GN-GloVe)
- ► Unique challenges: lack of definitional word pairs (e.g. "he" versus "she", "waiter" versus "waitress" in gender bias subspace) to compute a political bias subspace
- ► Solution: identify word pairs from Republicans and Democrats speech documents' frequent words

UCLA

# Bias in Word Embedding: Gender Bias

A famous way of modeling gender bias ("Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings"): [2]

- ▶ Gender subspace defined on gender word pairs (she-he, her-his, woman-man, mary-john, herself-himself, daughterson, mother-father, gal-guy, girl-boy, female-male)

- ▶ Using these pairs, compute a vector $w \in \mathbb{R}^{768}$ which captures the "gender direction"; male words will be on one end of the space and female words will be on the other end about this direction.

- ▶ Challenge in Political Bias: no ground-truth pairs

[2] https://arxiv.org/abs/1607.06520

UCLA

$$\overrightarrow{\text{grandmother}} = \overrightarrow{\text{wise}} + \overrightarrow{\text{gal}}$$

$$\overrightarrow{\text{grandfather}} = \overrightarrow{\text{wise}} + \overrightarrow{\text{guy}}$$

$$\overrightarrow{\text{grandmother}} - \overrightarrow{\text{grandfather}} = \overrightarrow{\text{gal}} - \overrightarrow{\text{guy}} = g\,,$$

$$\text{DirectBias}_c = \frac{1}{|N|} \sum_{w \in N} |\cos(\overrightarrow{\text{w}}, g)|^c\,,$$

where $N$ are the gender-neutral words, and $c$ is the parameter controls the strictness of measuring bias (when $c$ is close to 0 then the score is 0 only when $\overrightarrow{\text{w}}, g$ have no overlap).

UCLA

Another way measures the **indirect bias**. A given word vector $w \in \mathbb{R}^d$ normalized to unit length can be decomposed as:

$$w = w_g + w_\perp \qquad w_g = (w \cdot g)g\,,$$

where $w_g$ is the contribution from gender, $w_\perp$ is the remainder. We measure the similarity
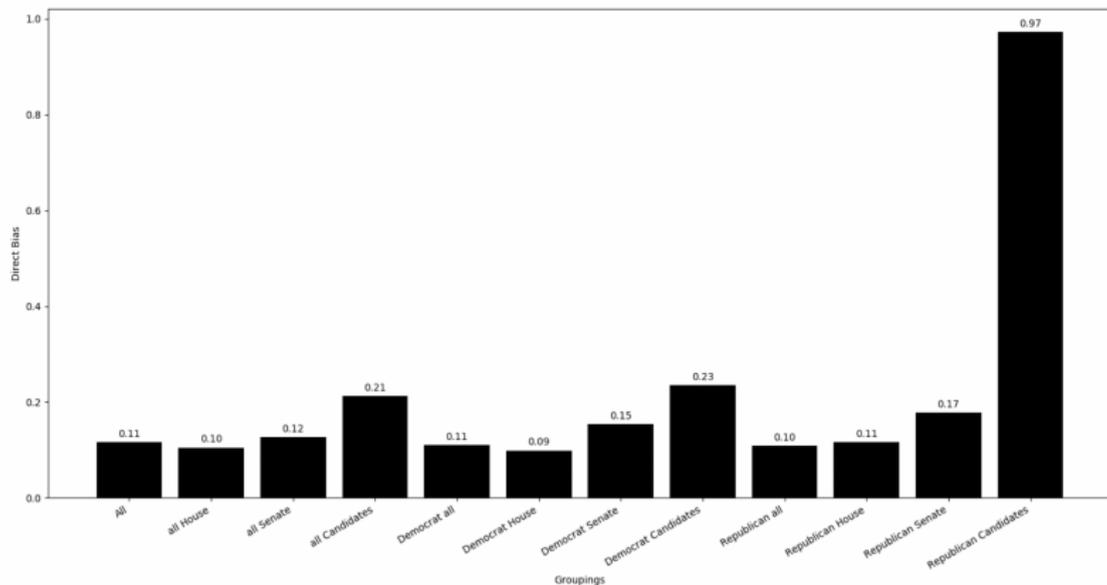
$$\beta(w, v) = \left( w \cdot v - \frac{w_\perp \cdot v_\perp}{\|w_\perp\|_2 \|v_\perp\|_2} \right) / w \cdot v\,,$$

where $\frac{w_\perp \cdot v_\perp}{\|w_\perp\|_2 \|v_\perp\|_2}$ is the re-normalized inner product after projecting out the gender space. e.g. Considering $g = \overrightarrow{\text{softball}} - \overrightarrow{\text{football}}$, *receptionist*, *waitress* are closer to *softball* than *football*, the $\beta$ scores between these words and *softball* are 67% and 35% respectively.

**UCLA**

Political-Pair Identification Steps:

1. Collect documents characteristic of the political bias. e.g. Set of tweets from Republican versus Democratic politicians in the US.

2. Extract lists of the most commonly used words.

3. From the candidate words, look for related or corresponding pairs. i.e. Different words that Republicans and Democrats might use to talk about the same idea.

Note: These words are **not** chosen to be direct antonyms, but describe the same concept / a parallel concept.

UCLA

Trump's tweets (Republican candidates) has the highest polarity score 0.97. In general, tweets from presidential candidates have higher bias than for other politicians. *Maybe more extreme views can be effective in driving political momentum.*

UCLA

Political audience diversity and news reliability in algorithmic ranking

- ▶ Motivation: political bias in news audience.
- ▶ Related Work: Collaborative Filtering
- ▶ show that popularity does not predict news reliability
- ▶ show that websites with more extreme and less politically diverse audiences have lower journalistic standards
- ▶ using the political diversity of a website's audience as a quality signal
- ▶ incorporate audience diversity into a standard collaborative filtering framework to increases trustworthiness of newsfeed algorithms
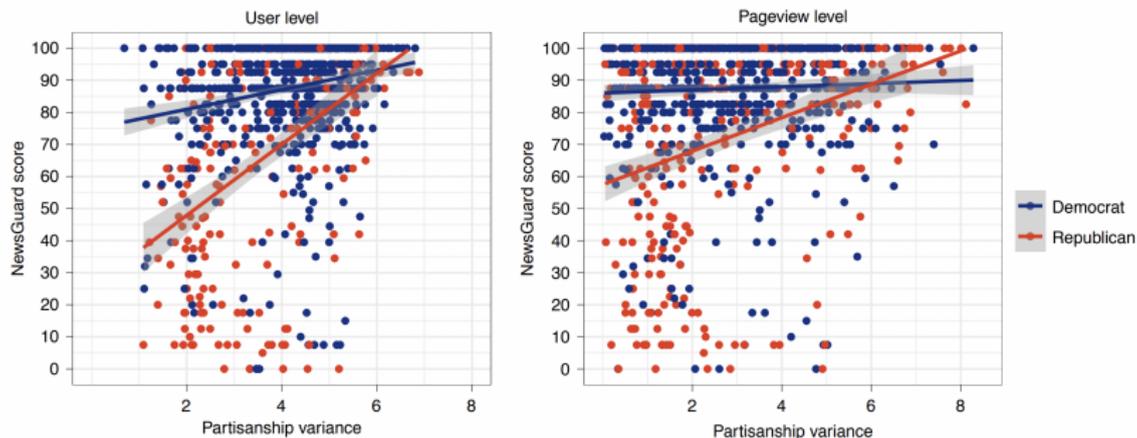
UCLA

- ► A comprehensive data set of web traffic history from $6,890$ US residents, collected along with surveys of self-reported partisan information from respondents in the YouGov Pulse survey panel;
  - ► Used as ground-truth pageview data and user political-tendency data.
- ► A dataset of $3,765$ news source reliability scores compiled by trained experts in journalism and provided by NewsGuard.
  - ► Used as ground-truth news source reliability scores.

UCLA

Current recommendation systems rely on popularity and engagement. Potential problem: domain pageviews are not associated with overall news reliability.

- ▶ Show the concern by measuring correlation between user / pageview (visit) data and reliability scores.
- ▶ Use coefficient of partial correlation between NewsGuard reliability scores and the **variance** of audience partisanship.
- ▶ Conclusion: measures of audience partisan diversity correlate with news reliability better than popularity does.

Better newsfeed algorithm that considers reliability?

- ▶ Solution: incorporating audience partisan diversity into algorithmic ranking decisions.

**UCLA**

Audiences' divergence serves as a reliable measurement of news agency's quality.

Starting from classic user-based CF (Collaborative Filtering) algorithm.

▶ Provide the best recommendations for users by learning from others with similar preferences on items.

▶ "Item" here refers to "news source domain".

Build the user-domain matrix $V \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{D}|}$ from user set $\mathcal{U}$ and domain set $\mathcal{D}$:

$$v_{u,d} = \frac{\pi_{u,d}}{\sum_h \pi_{u,h}} \left( \frac{\pi}{\sum_u \pi_{u,d}} \right)$$

where $\pi_{u,d} \in \mathbb{Z}^+$ counts how many times a user $u$ has visited domain $d$, and $\pi = \sum_u \sum_v \pi_{u,d}$ is the total number of visits.

UCLA

Then coefficient of similarity between users $u$ and $u'$ is:

$$\text{sim}(u, u') = \frac{\tau(V_u, V_{u'}) + 1}{2},$$

where $V_u \in \mathbb{R}^{1 \times |\mathcal{D}|}$ denotes a row vector of $V$. Here, $\tau$ can be correlation coefficient, such as the Kendall rank correlation coefficient, or the Pearson's correlation coefficient.

These similarity coefficients are used to calculate the predicted ratings. Part of the user-domain visit data is held-out for the prediction task.

UCLA

Standard user-based CF calculates the predicted rating of a user $u$ for a domain $d$ as:

$$\hat{v}_{u,d}^{\text{CF}} = \overline{v}_u + \frac{\sum_{u' \in N_{u_d}} \text{sim}(u, u')(v_{u',d} - \overline{v}_{u'})}{\sum_{u' \in N_{u_d}} \text{sim}(u, u')} \, ,$$

where $N_{u_d}$ is the set of $n$ most similar users to $u$ who have rated $d$, $\overline{v}_u$ is the average ratings of $u$ across all domains they visited (average of a row in $V$).
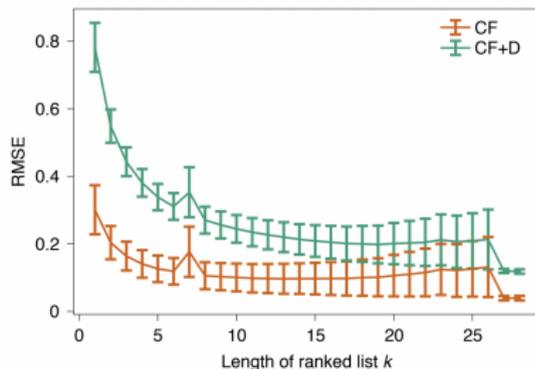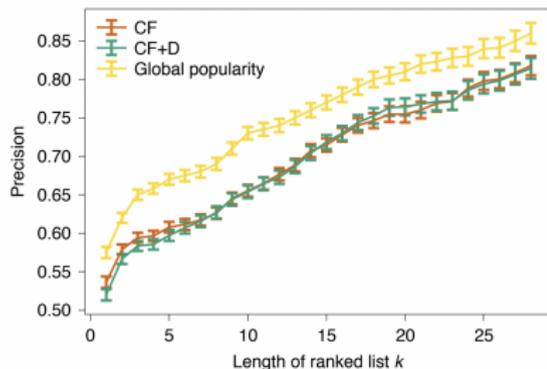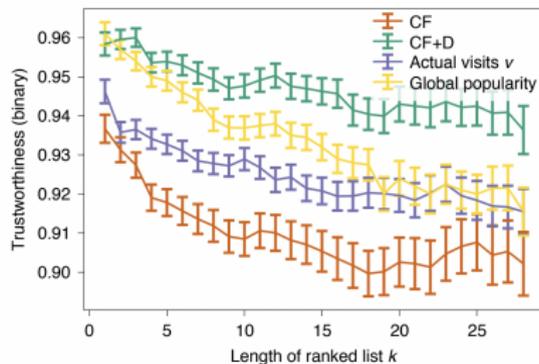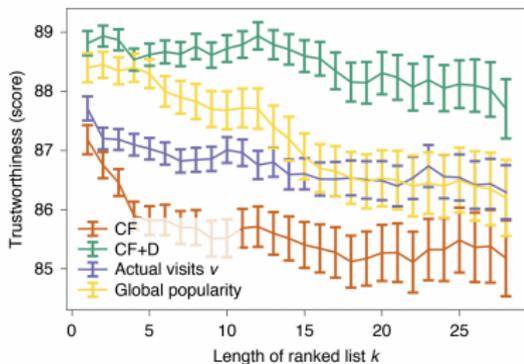
UCLA

The variant CF+D (collaborative filtering + diversity) calculates the predicted rating of a user $u$ for a domain $d$ as:

$$\hat{v}_{u,d}^{\text{CF+D}} = \hat{v}_{u,d}^{\text{CF}} + g(\delta_d) \, ,$$

where $g(\delta_d)$ is obtained by plugging the self-reported audience partisan diversity $\delta_d$ of each domain $d$ into a standard logistic function:

$$g(\delta_d) = \frac{a}{1 + \exp(-(\delta_d - t)/\psi)}$$

with parameters $a, \psi, t$ generalizing the upper asymptote, inverse growth rate, and location of the standard logistic function, respectively. $t$ is empirically estimated as $t = \overline{\delta}$, the average audience partisan diversity across all domains.

UCLA

Belief polarization in a complex world: A learning theory perspective

- ▶ Related Work: Learning Theory
- ▶ analyzing the problem of belief polarization in a purely mathematical way

UCLA

$\mathcal{D}$: distribution over domain $\mathcal{X} \times \mathcal{Y}$ where the set of labels $\mathcal{Y} = \{-1, 1\}$

$\mathcal{D} \downarrow \mathcal{X}$: the marginal distribution of $\mathcal{D}$ on $\mathcal{X}$

$\mathcal{P}, \mathcal{P}'$: any distribution over domain $\mathcal{X}$, with total variation distance denoted as supremum of the element-wise distances:

$$\text{TV}(\mathcal{P}, \mathcal{P}') = \sup_{X \in \mathcal{X}} |\mathcal{P}(X) - \mathcal{P}'(X)|$$

For ease of exposition, $L_1$ distance is used instead, knowing that

$$\|\mathcal{P} - \mathcal{P}'\|_1 = 2\text{TV}(\mathcal{P}, \mathcal{P}')$$

UCLA

For $\mathcal{D}$ and $\mathcal{D}'$ over $\mathcal{X} \times \mathcal{Y}$, we say that their conditional label distribution match each other when:

$$\forall x \in \mathcal{X}, \Pr_{(x,y)\sim\mathcal{D}}[y|x] = \Pr_{(x,y)\sim\mathcal{D}'}[y|x] \, ,$$

and when they match, we use them interchangeably:

$$\|D - D'\| \Leftrightarrow \|\mathcal{D} \downarrow \mathcal{X} - \mathcal{D}' \downarrow \mathcal{X}\|$$

Belief function class $\mathcal{F} : \mathcal{X} \to \mathcal{Y}$

The error of belief function $f \in \mathcal{F}$ over $\mathcal{D}$ is described as

$$\mathrm{err}_{\mathcal{D}}(f) = \forall x \in \mathcal{X}, \mathrm{Pr}_{(x,y)\sim\mathcal{D}}[f(x) \neq y]$$

$\mathcal{D}$ is realizable iff $\exists f \in \mathcal{F}, \ \mathrm{err}_{\mathcal{D}}(f) = 0$.

The empirical error of $f$ over a training set $S$ of $m$ labeled data points $S = \{(x^i, y^i)\}_{i\in[m]}$ is denoted by

$$\mathrm{err}_S(f) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\big(f(x^i) \neq y^i\big)$$

UCLA

For two belief functions $f, f' \in \mathcal{F}$, the disagreement of them on $\mathcal{D}$ is denoted by

$$\Delta_{\mathcal{D}}(f, f') = \Pr_{x \sim \mathcal{D} \downarrow \mathcal{X}}[f(x) \neq f'(x)]$$

For a set $\mathcal{H}$ of belief functions, we denote its diameter by

$$\text{diam}_{\mathcal{D}}(\mathcal{H}) = \max_{f, f' \in \mathcal{H}} \Delta_{\mathcal{D}}(f, f')$$

as the largest disagreement between two belief functions in this class.

Disagreement between two belief functions and the diameter of a belief function class do not depend on the labels of $\mathcal{D}$; with a slight abuse of notation, we use $\mathcal{D}$ in place of $\mathcal{D} \downarrow \mathcal{X}$ in the notation for diameter and disagreement at times.

UCLA

Learning setting is polarizing if agents learn functions whose **disagreement** is disproportionately larger than the **difference between the distributions** to which they were exposed.

Settings: two agents learn functions $f_1$ and $f_2$ from distributions $\mathcal{D}_1$ and $\mathcal{D}_2$ that have the (1) same marginal distribution $\mathcal{D} \downarrow \mathcal{X}$ and (2) either the **matched** or very similar conditional label distribution, yet $\Delta_{\mathcal{D}}(f_1, f_2)$ is large.

In this view, polarization is the lack of consensus between agents' beliefs when exposed to similar sets of information, independently of how inaccurate these beliefs may be.

UCLA

Consider realizable $\mathcal{D}_1, \mathcal{D}_2$ over $\mathcal{X} \times \mathcal{Y}$ with shared marginal distribution $\mathcal{D}$, consistent with belief functions $f_1$, $f_2$ respectively. i.e.

$$\mathrm{err}_{\mathcal{D}_1}(f_1) = \mathrm{err}_{\mathcal{D}_2}(f_2) = 0$$

while $\Delta_{\mathcal{D}}(f_1, f_2)$ is large.

We consider the two agents attempts to see the world from the other's perspective (perhaps through communication), end up observing training sets from almost identical mixtures of these two distributions and learn belief functions $\tilde{f}_1$ and $\tilde{f}_2$; we ask whether $\Delta_{\mathcal{D}}(\tilde{f}_1, \tilde{f}_2)$ will be significant.

UCLA

- Distribution (after communication): $\tilde{\mathcal{D}}_1 = (1-\alpha)\mathcal{D}_1 + \alpha\mathcal{D}_2$, $\tilde{\mathcal{D}}_2 = (1-\alpha)\mathcal{D}_2 + \alpha\mathcal{D}_1$, with $\alpha \in (0, \frac{1}{2})$

- Assume $\tilde{f}_1$ and $\tilde{f}_2$ achieve optimal accuracy on training sets $\tilde{S}_1$ and $\tilde{S}_2$ drawn from the above distributions

- Prove the possibility of large disagreement between $\tilde{f}_1$ and $\tilde{f}_2$ when $\alpha \to \frac{1}{2}$ (i.e. distribution almost identical), by proving that *the optimal belief functions will converge to the original beliefs, s.t. differ from each other by as much as the originals.*

- Show that by prob at least $1 - \delta$, $\Delta_{\mathcal{D}}(\tilde{f}_1, \tilde{f}_2) \geq \Delta_{\mathcal{D}}(f_1, f_2) - \epsilon$, $\Delta_{\mathcal{D}}(\tilde{f}_i, f_i) \leq \epsilon/4$ $(i = 1, 2)$

- Proved with help of Vapnik–Chervonenkis dimension (VC dimension) $\mathrm{VCD}(\mathcal{F})$. (bound of $m$ expressed as function of $\delta, \epsilon, \mathrm{VCD}(\mathcal{F})$)

**UCLA**

This is the case where both agents sample from the same realizable $\mathcal{D}$ consistent with $f^* \in \mathcal{F}$, but prefer simplicity by consider:

$$\text{cost}_{\mathcal{D}}(f) = \text{err}_{\mathcal{D}}(f) + \phi(f), \qquad \text{cost}_S(f) = \text{err}_S(f) + \phi(f).$$

Concrete examples are used for proof. Example 1, for $\mathcal{X} = \mathbb{R}^d$, $\mathcal{F} = \{f_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x})\}_{\mathbf{w} \in \mathbb{R}^d}$, a concrete example of the complexity cost:

$$\phi(f_{\mathbf{w}}) = h(\|\mathbf{w}\|_0),$$

where $h(0) = 0$, $h(1) \leq \frac{1}{2d}$, $h(d) \geq \frac{1}{2}$.

UCLA

- ▶ Construct an Example
- ▶ Prove the example following its concrete settings
- ▶ Show that there are at least two optimal belief functions $f_1, f_2 \in \arg\min_{f \in \mathcal{F}} \text{cost}_{\mathcal{D}}(f)$ with nontrivial disagreement.
- ▶ The intuition: polarization can easily arise after the setting involves many dimensions and complexity leads people to choose some subset; the chance that they coordinate on exactly the same dimensions is only high if there aren't other dimensions of similar importance. Otherwise, different observers can easily favor different dimensions in their belief functions.
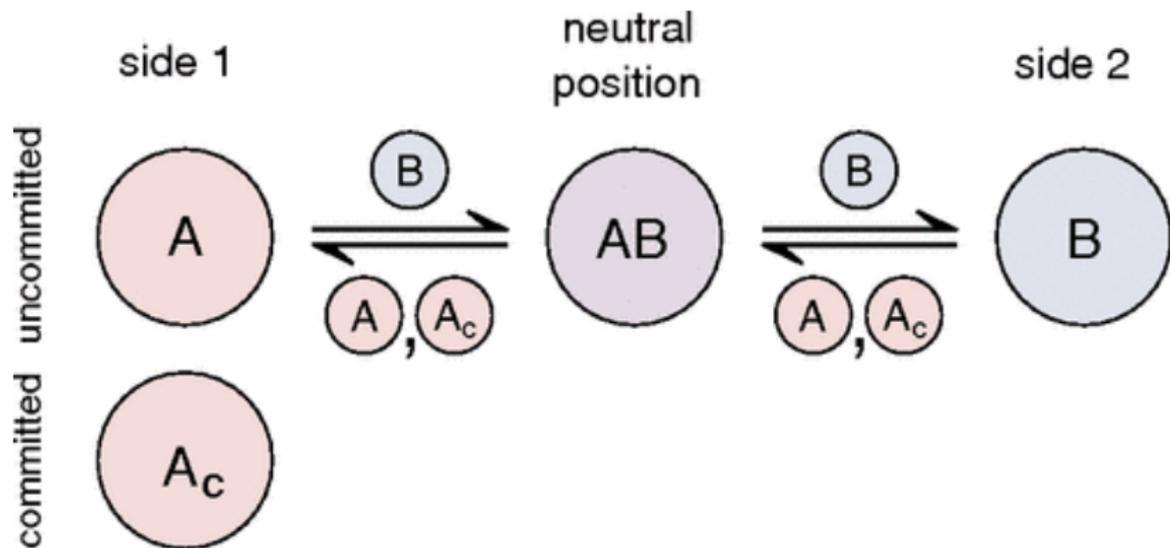
UCLA

# Reducing the Polarity

Encouraging Moderation: Clues from a Simple Model of Ideological Conflict

▶ When the fraction of zealots exceeding a threshold, they might change the ideology of the rest.

▶ Encouraging people to stay at their own position is not a good idea for moderation.

▶ Non-social stimulus is needed for moderation.

Reasoning about Political Bias in Content Moderation

▶ Moderation is hard in practice, easily biased

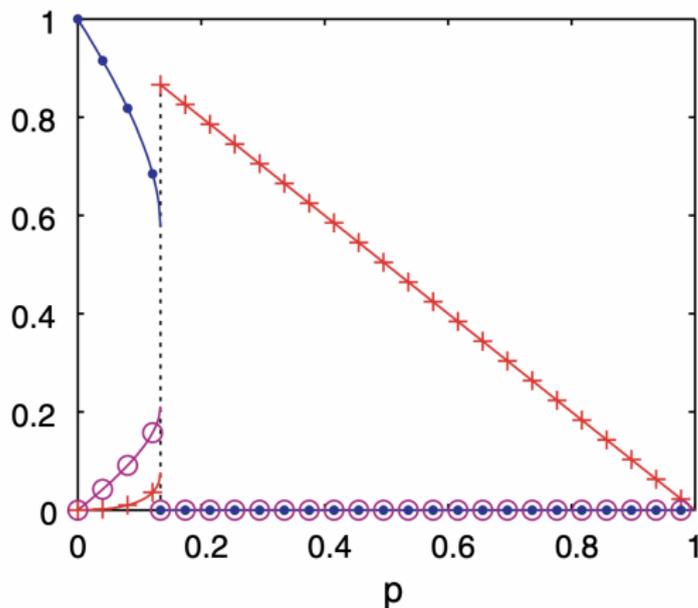| Speaker | Listener preinteraction | Listener post-interaction |
|---|:---:|:---:|
| $A, A_c$ | $B$ | $AB$ |
| | $AB$ | $A$ |
| $B$ | $A$ | $AB$ |
| | $AB$ | $B$ |

The dynamics of the basic model are deterministic, continuous and mean-field. When there is a pair of *speaker* and *listener*, the position of *listener* might change. The **zealots** $A_c$ subpopulation is **constant**, never changed.

UCLA

Let $p$ denote the constant fraction zealots $A_c$ subpopulation, and $n_A$, $n_B$, and $n_{AB}$ denote the expected fractions of the total population of $N$ individuals corresponding to the uncommitted $A$, $B$, and $AB$.

$$n_A + n_B + n_{AB} + p = 1$$

If we select a listener and a speaker uniformly at random:

$$\Delta n_A = (p + n_A)n_{AB} - n_A n_{AB}$$
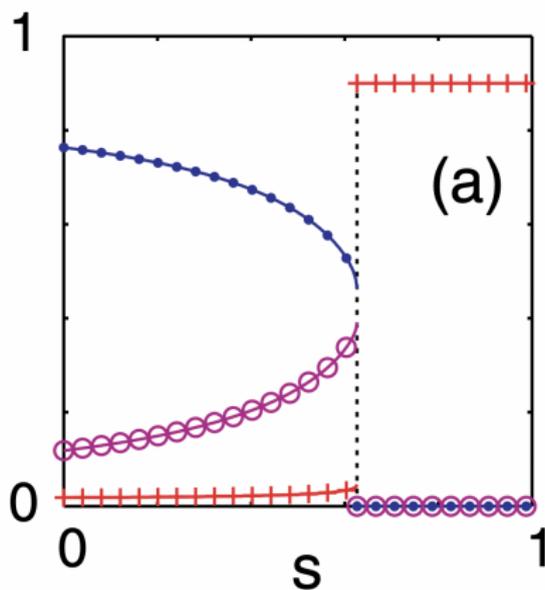$$\Delta n_B = n_B n_{AB} - (p + n_A)n_B \tag{1}$$

Equation (1) with assumed initial population:
$n_A = 0, n_B = 1 - p$. The dashed line indicates a critical point
$p_c = 1 - \frac{\sqrt{3}}{2} \approx 0.134$. $n_A$ for red plus signs, $n_B$ blue dots, and
$n_{AB}$ the magenta open circles.

UCLA

$$\Delta n_A = (1-s)(p + n_A)n_{AB} - n_A n_{AB}$$
$$\Delta n_B = (1-s)n_B n_{AB} - (p + n_A)n_B$$

(2)

where the stubbornness parameter $s$ indicates how likely a moderate is to remain moderate after listening to an extremist. The $s = 0$ special case is identical to the basic model.
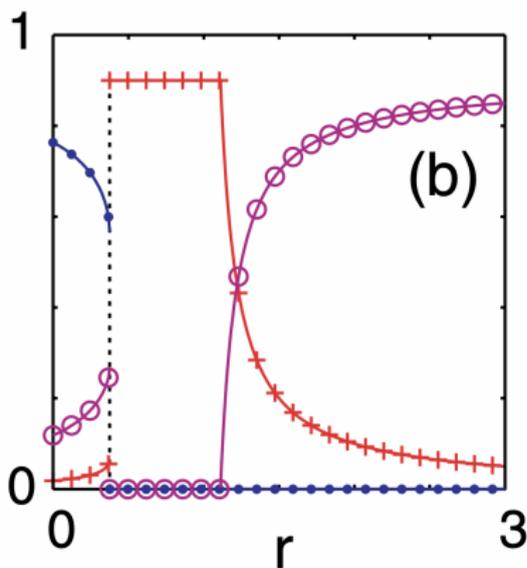
UCLA

In this example, $p = 0.1$ is fixed, and $n_A = 0, n_B = 1 - p$. The model fails to guarantee moderation.

UCLA

Increasing $s$ not only reduces the flow of $AB$ to $A$, but also reduces the flow of $AB$ to $B$, thereby depleting uncommitted subpopulations on both sides.

With competition from $B$ extremists over the $AB$ subpopulation weakened as a result, it takes fewer $A$ zealots (hence lower $p_c$) to convert the moderates to the $A$ camp.

UCLA

$$\Delta n_A = (p + n_A)n_{AB} - n_A n_{AB} - rn_A n_{AB}$$
$$\Delta n_B = n_B nAB - (p + n_A)n_B - rn_B n_{AB}$$

$$(3)$$

where the new parameter $r$ is a nonnegative real number that reflects the intensity of the moderates' evangelism. It means that there's a tendency of $AB$ members trying to convince $A$ and $B$ into $AB$.
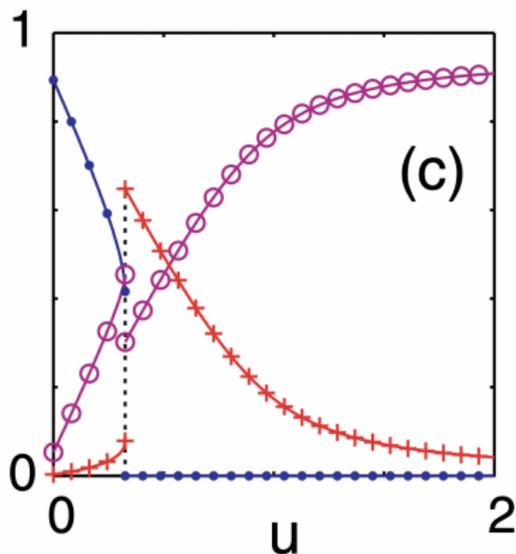
UCLA

In this example, $p = 0.1$ is fixed, and $n_A = 0, n_B = 1 - p$.
Again, the model fails to guarantee moderation.

If the moderates' campaign of persuasion is sufficiently successful from the start (i.e. $r$ starts and stays large enough) then the moderates do in fact maintain a large, robust equilibrium population.

Otherwise $AB$'s evangelistic efforts can instigate their own extinction.

$$\Delta n_A = (p + n_A)n_{AB} - n_A n_{AB} - u n_A$$
$$\Delta n_B = n_B nAB - (p + n_A)n_B - u n_B$$

(4)

where $u$ is a nonnegative parameter representing the rate at which the radicals abandon their radical position in response to the **nonsocial stimulus**.
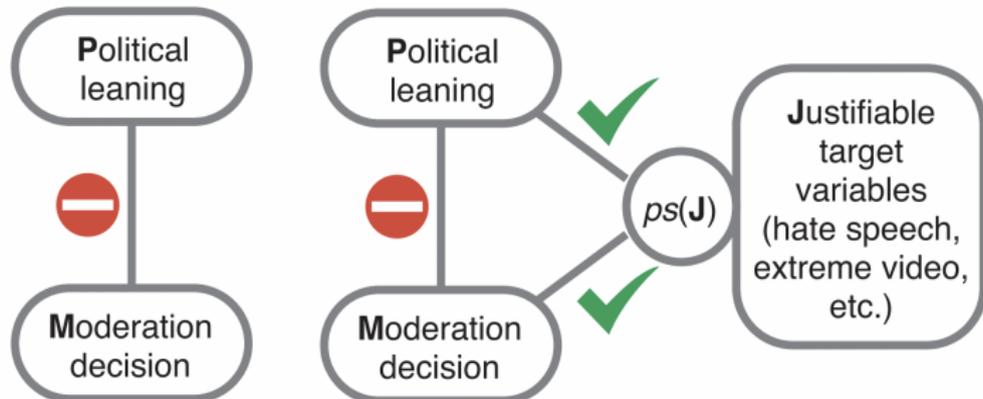
UCLA

(c)

In this example, $p = 0.05$ is fixed, and $n_A = 0, n_B = 1 - p$. This time, moderation is almost guaranteed. We still see the influence of $A_c$.

UCLA

They didn't consider this situation:

As time goes by, $B$ zealots $B_c$ form their community.

There's a trend in many movements that last for a long time (e.g. last for year): one extreme triggers its opponent, and then they become increasingly extreme.

Reasoning about Political Bias in Content Moderation



(a) Independence. 1st null hypothesis $\mathbf{H}_0^{ind}$: M ⊥⊥ P.

(b) Separation. Propensity scoring function $ps(\mathrm{J})$ is used to summarize J to a scala, hence 2nd null hypothesis $\mathbf{H}_0^{sep}$: M ⊥⊥ P $\mid ps(\mathrm{J})$.

UCLA

M: moderated / alive, P: left / right, J: justifiable target variables e.g. hate-speech.
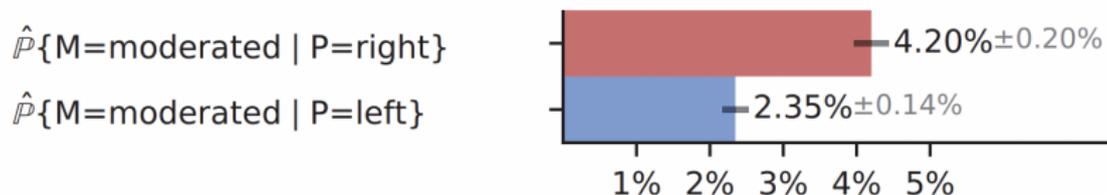
Independence measurement:
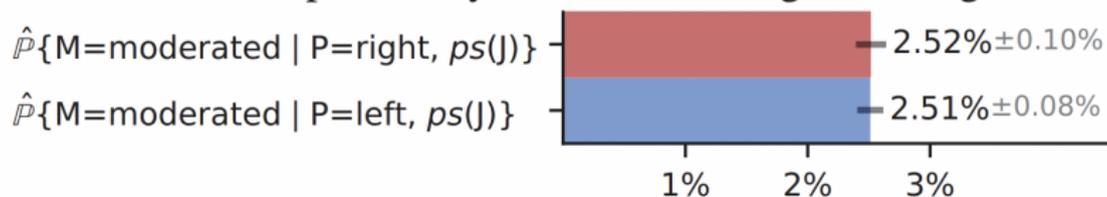
$$\mathbb{P}\{M|P = left\} = P\{M|P = right\}$$

Separation measurement:

$$\mathbb{P}\{M|P = left, J\} = P\{M|P = right, J\},$$

where in practice we usually measure J by *propensity scoring* $ps(J)$ where $ps : \mathbb{R}^{|J|} \to \mathbb{R}$.

**UCLA**

(a) $\mathbf{H}_0^{ind}$ is rejected. There is significant difference between comment moderation probability under left- and right- leaning videos.



(b) $\mathbf{H}_0^{sep}$ holds. There is no significant difference between comment moderation probability under left- and right- leaning videos with propensity scored justifiable variables.

**UCLA**

# Conclusion

Data
- ▶ Incompleteness: Platforms & Users will delete
- ▶ Lack of Labels: we humans are complex, hard to label us in nature
- ▶ Always Biased: preference of using certain platforms / not using social networks etc.
- ▶ Noisy: bot, cyborg, misinformation, disinformation...

Methodology
- ▶ Majority $\neq$ Reliability: Human Expert needed
- ▶ Theory $\neq$ Reliability: Field Studies needed

UCLA

There is an "unchanged" belief in the online & offline world. The materialism / rationalism, no matter what we call it.

Because we live in the material world, and we do care about objective reality.

It is easy to trigger polarization. Do it by telling lies, by spreading misinformation and disinformation. It is hard to moderate the polarization.

Reality itself can be complex. It is challenging to convey them in understandable ways. But writing one-sided biased articles are easy.

We all hold that "I am right, he/she is wrong". Sometimes it is true, sometimes we are on the different aspects of a big picture.

We are all victims of information explosion. The information we are exposed to has to be selected. Selecting process can easily build information cocoons.

UCLA

- ▶ Polarization and the Global Crisis of Democracy: Common Patterns, Dynamics, and Pernicious Consequences for Democratic Polities
  - ▶ a summary of the common patterns of political polarization, e.g. its relational nature
- ▶ Visual Mis- and Disinformation, Social Media, and Democracy
  - ▶ a survey of disinformation and misinformation on multi-modality data source
- ▶ Opinion dynamics and bounded confidence: models, analysis and simulation
  - ▶ simulated network bridges the theories and applications

UCLA