

Polarization Measures

Dora Demszky

May 18, 2019

1 Definitions

Let $|U| = |D| + |R|$ be the total number of users tweeting about a particular event, where D and R denote the set of Democratic and Republican users, respectively.

Let \mathbf{c}_i be the vector of token counts for user i and $m_i = \sum_j c_{ij}$ the total amount of tokens used by user i . Let $\hat{\mathbf{q}}_i = \mathbf{c}_i/m_i$ be the empirical phrase frequencies for user i . Let $\hat{\mathbf{q}}_i^P = \sum_{i \in P} \mathbf{c}_i / \sum_{i \in P} m_i$ be the empirical phrase frequencies for party P .

The definition of $\hat{\boldsymbol{\rho}}$ varies depending on the type of features (or “distance metric”, as referred to in Gentzkow et al. (Forthcoming)) we use, such as posterior, mutual information or χ^2 , and I will define $\hat{\boldsymbol{\rho}}$ in Section 3.

2 Estimators

The estimator are the same as they are defined in Gentzkow et al. (Forthcoming). In my interpretation, the non-model-based estimates can be thought of as weighted averages of token-level features, where the features are weighted by the distribution of token use for users and then the party-level averages are averaged across the two parties.

2.1 MLE

$$\hat{\pi}^{MLE} = \frac{1}{2}(\hat{\mathbf{q}}_i^R \cdot \hat{\boldsymbol{\rho}} + \hat{\mathbf{q}}_i^D \cdot (1 - \hat{\boldsymbol{\rho}})) \quad (1)$$

2.2 Leave-out Estimator

$$\hat{\pi}^{LO} = \frac{1}{2} \left(\frac{1}{|R|} \sum_{i \in R} \hat{\mathbf{q}}_i^R \cdot \hat{\boldsymbol{\rho}}_{-i} + \frac{1}{|D|} \sum_{i \in D} \hat{\mathbf{q}}_i^D \cdot (1 - \hat{\boldsymbol{\rho}}_{-i}) \right) \quad (2)$$

3 Feature Types ($\hat{\boldsymbol{\rho}}$)

3.1 Posterior Probability

$$\hat{\boldsymbol{\rho}}^{PROB} = \frac{\hat{\mathbf{q}}^R}{(\hat{\mathbf{q}}^R + \hat{\mathbf{q}}^D)} \quad (3)$$

3.2 Mutual Information

Let u_j^P be the number of users from party P mentioning the token j and let \mathbf{u}^P be the vector representing the number of users from party P mentioning each token. Also, let $\mathbf{u} = \mathbf{u}^R + \mathbf{u}^D$. Analogically, let u_{*j}^P be the number of users from party P *not* mentioning the token j and let \mathbf{u}_*^P be the vector representing the number of users from party P not mentioning each token, and let $\mathbf{u}_* = \mathbf{u}_*^R + \mathbf{u}_*^D$. I used add-one smoothing to calculate \mathbf{u}^P and \mathbf{u}_*^P as is commonly done in the case of mutual information.

The following four equations represent the components of the mutual information formula as defined here¹, when probabilities are MLE.

$$\mathbf{v}^D = \mathbf{u}^D \otimes \log_2\left(\frac{|U|}{|D|} \mathbf{u}^D \oslash \mathbf{u}\right) \quad (4)$$

$$\mathbf{v}_*^D = \mathbf{u}_*^D \otimes \log_2\left(\frac{|U|}{|D|} \mathbf{u}_*^D \oslash \mathbf{u}_*\right) \quad (5)$$

$$\mathbf{v}^R = \mathbf{u}^R \otimes \log_2\left(\frac{|U|}{|R|} \mathbf{u}^R \oslash \mathbf{u}\right) \quad (6)$$

$$\mathbf{v}_*^R = \mathbf{u}_*^R \otimes \log_2\left(\frac{|U|}{|R|} \mathbf{u}_*^R \oslash \mathbf{u}_*\right) \quad (7)$$

where \otimes and \oslash denote element-wise multiplication and division, respectively.

The mutual information of each feature (token) is estimated via

$$\hat{\rho}^{MI} = \frac{1}{|U|} (\mathbf{v}^D + \mathbf{v}^R + \mathbf{v}_*^D + \mathbf{v}_*^R) \quad (8)$$

3.3 Chi-Square (χ^2)

The χ^2 value of each feature (token) is estimated via

$$\hat{\rho}^{\chi^2} = \frac{|U|(\mathbf{u}^D \otimes \mathbf{u}_*^R - \mathbf{u}_*^D \otimes \mathbf{u}^R)^2}{\mathbf{u} \otimes \mathbf{u}_* \otimes (\mathbf{u}^D + \mathbf{u}_*^D) \otimes (\mathbf{u}^R + \mathbf{u}_*^R)} \quad (9)$$

¹<https://nlp.stanford.edu/IR-book/html/htmledition/mutual-information-1.html>