# Post-Prediction Inference on Political Twitter

Alicia Gunawan, Dylan Haar, Luis Ledezma-ramos

## Abstract

Having observed data seems to be a necessary requirement to conduct inference, but what happens when observed outcomes cannot easily be obtained? The simplest practice seems to proceed with using predicted outcomes, but without any corrections this can result in issues like bias and incorrect standard errors. Our project studies a correction method for inference conducted on predicted, not observed outcomes—called post-prediction inference—through the lens of political data. We are investigating the kinds of phrases or words in a tweet that will most strongly indicate a person's political alignment to US politics. We have discovered that these correction techniques are promising in their ability to correct for post-prediction inference in the field of political science.

## 1 Intro

Machine learning is a modern task in data science that uses observed data values to model and predict data. It takes advantage of having observed data available, but what should be done when observed data cannot be obtained? A common practice is to use predicted values when observed values are unavailable, but without any corrections we inevitably run into issues such as incorrect standard errors, bias, and inflated false positive rates.

Wang et al. proposes a method to correct inference done on predicted outcomes–which they name post-prediction inference, or postpi–in *Methods for correcting inference based on outcomes predicted by machine learning*. This statistical technique takes advantage of the standard structure for machine learning and uses bootstrapping to correct statistical error when using predicted values in place of observed values.

We are exploring the applicability of Wang et al.'s postpi bootstrapping technique on political data–that is, on political twitter posts. Our project will be investigating what kinds of phrases or words in a tweet will strongly indicate a person's political alignment, in the context of US politics. By doing so, we can simultaneously test how the bootstrap post-prediction inference approach interacts with text data and how this method can be generally applicable towards analyses in political science.

## 2 Methodology

The postpi bootstrap approach by Wang et al. is a method that aims to correct inference in studies that use predicted outcomes in lieu of observed outcomes. It is effective due to its simplicity–this approach is not dependent on deriving the first principles of the prediction model, so we are free to focus on accuracy without worrying about the impact of the complexity of the model on the

bootstrap approach. The reason why it is not dependent is because this approach utilizes an easily generalizable and low-dimensional relationship between observed and predicted outcomes.

There are four assumptions that the postpi bootstrap approach rests on:

1. The training, testing, and validation dataset must all arise from the same distribution, and the training and testing set must have observed outcomes to train the prediction and relationship model.
2. There must be a simple, low-dimensional relationship between observed and predicted outcomes.
3. The relationship model arising from the relationship above must be consistent for the validation set and for any future data.
4. The features used when conducting inference must be present in the training and testing data, and used to train the prediction model.

An implementation of this algorithm is provided below:

---
**Algorithm** Bootstrap-based Postpi Correction
---
**Require:** Observations $\{x_{(tr)}, y_{(tr)}, x_{(te)}, y_{(te)}, x_{(val)}\}$, where $y \in \mathbb{R}^n$ and $x \in \mathbb{R}^{n \times m}$.
**Require:** Prediction model $\hat{f}(\cdot)$, relationship model $k(\cdot)$, and inference model of the form $g[E(y|X)] = X\beta$.
1: Use $(x_{(tr)}, y_{(tr)})$ to fit the prediction model s.t. $y_p = \hat{f}(x)$.
2: Get test set predicted outcomes: $y_{p(te)} = \hat{f}(x_{(te)})$.
3: Use $(y_{(te)}, y_{p(te)})$ to fit the relationship model s.t. $y = k(y_p)$.
4: Get validation set predicted outcomes: $y_{p(val)} = \hat{f}(x_{(val)})$
5: Use $(x_{(val)}, y_{p(val)})$ to bootstrap
6: **for** $b \in \{1, ..., B\}$ **do**
7:      Sample predicted outcomes and covariates $(x_{i(val)}^b, y_{pi(val)}^b)$ with replacement for $i = 1, ..., n$.
8:      Simulate values from the relationship model $\bar{y}_i^b = k(y_{p(val)}^b)$.
9:      Fit the inference model $g[E(\bar{y}^b | X_{(val)}^b)] = X_{(val)}^b \beta^b$.
10:      Extract coefficient estimator $\beta^b$ from the fitted inference model.
11:      Extract the SE of the estimator $se(\hat{\beta}^b)$ from the fitted inference model.
12: Estimate the inference model coefficient using a median function $\hat{\beta}^{boot} = median(\hat{\beta}^1, \hat{\beta}^2, ..., \hat{\beta}^B)$.
13: Estimate the inference model SE:
14:      The parametric method $\hat{SE}^{boot,par} = median(\hat{SE}(\hat{\beta}^1), \hat{SE}(\hat{\beta}^2), ..., \hat{SE}(\hat{\beta}^B))$.
15:      The nonparametric method $\hat{SE}^{boot,non-par} = SD(\hat{\beta}^1, \hat{\beta}^2, ..., \hat{\beta}^B)$.

---

## 3 Data

### 3.1 Data Collection

We collected our data by scraping tweets from US politicians from Twitter. Specifically, we took the Twitter handles of the President, Vice President, and all the members of US Congress except Representatives Chris Smith (R-NJ) and Jefferson Van Drew (R-NJ), as they have both deleted their Twitter accounts. These Twitter handles were compiled and provided by the UCSD library, and outdated names or Twitter handles were updated manually by ourselves. Additionally, the

two Independent members of Congress–Senators Bernie Sanders (I-VT) and Angus King (I-ME)–will be considered Democratic politicians for our purposes, as they caucus with Democrats.

Using these Twitter handles, we scraped approximately 100 tweets from each politician, although the exact number of tweets pulled from each individual will fluctuate as not all members of Congress use Twitter with the same frequency as their colleagues. Our final dataset consists of 44,328 tweets for an average of 82 tweets per politician. Of these tweets, 22,653 tweets are from Democrats, 21,478 tweets are from Republicans, and 197 tweets are from Independents (converted to Democrats).

**3.2 Cleaning Text Data**

To prepare our data for prediction and feature selection, we cleaned the tweets by expanding all contractions, converted all text into lowercase format, and removed urls, punctuation, and unicode characters. Additionally, we also removed stopwords, using the dictionary of stopwords provided by NLTK to do so.

**3.3 Exploratory Data Analysis**

Our data consists of a relatively equal number of tweets leaning either Democratic or Republican. As said earlier, with Independent politicians counting as Democrats, there are a total of 44,328 tweets–22,850 are classified as tweets from Democrats, while 21,478 are classified as tweets from Republicans.

We look at Figure 1 for a first glance at the data. Figure 1 is an overlaid histogram plotting the number of words in tweets from Democrats and Republicans. While both histograms are clearly skewed to the left, we can see that the distribution of the length of tweets for Democrats has a higher peak than the distribution for Republicans, which tells us that tweets from Democrats average more words compared to their counterparts on the opposite aisle. This could imply that the prediction model will utilize more vocabulary from Democrat-classified tweets than Republican, which might have interesting effects on the prediction model and thus the bootstrap algorithm and inference.
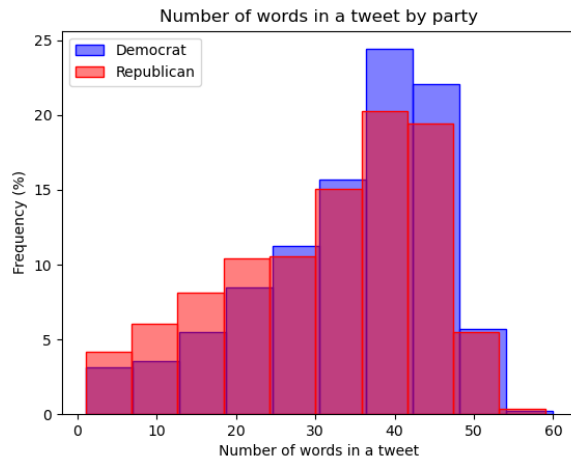
**Figure 1**: A histogram depicting the number of words in a tweet by party. We can see that Democrats generally have longer tweets compared to Republicans.

We take a deeper dive into each party in Figure 2 below, which lists the 10 most frequent words used by Democrats and Republicans, excluding stopwords. There are very few commonalities between either party–only two words are commonly used by both parties: 'today' and 'year'.

Democrats seem to focus on policy issues as suggested by 'act' and 'infrastructure', but otherwise their attentions are spread across a multitude of topics as no single unifying issue seems to be able to group together their most frequently used words. On the other hand, Republicans seem to focus more on their political opponents–words such as 'biden', 'democrats', and 'president' seem to suggest that–and on the American people. There is notably a significant reference to 'biden', with the President's name being used approximately 3500 times, almost double the frequency of the second most popular word. As such, Figure 2 shows us that Republican-classified tweets may revolve more strongly around certain themes, such as their opponents, compared to Democrat-classified tweets. Again, this may influence the prediction model and in turn the inference conducted on our features.
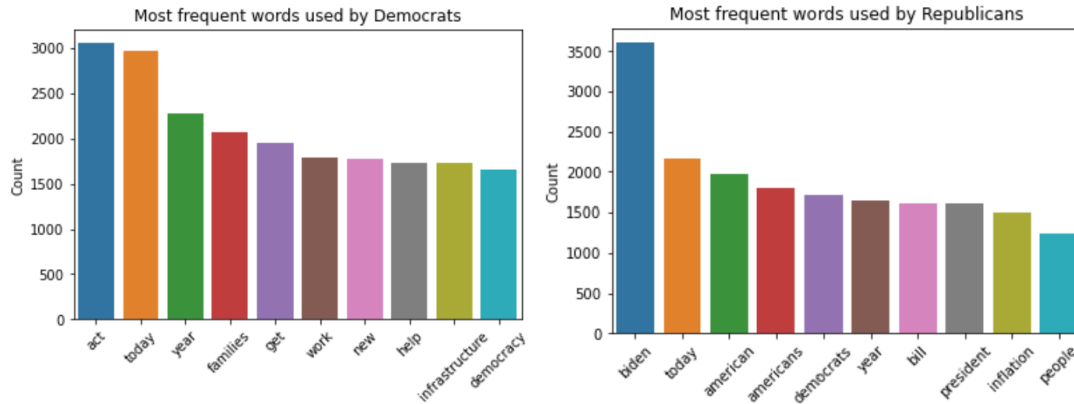
**Figure 2**: Bar plots depicting the most frequent words used by either party. We can also see a significant difference in the most frequent words used by either party–only 'today' and 'year' is a word that both parties use in common.

## 4 Methods

### 4.1 Prediction and Relationship Model

During this stage of our project, we worked on maximizing the accuracy of our prediction model. We compared several different prediction models in the process of coming up with our final model, trying other classification algorithms such as logistic regression and ridge regression (regularized). After determining which model performed the best, we tuned hyperparameters on the final model to further improve its performance. In the end, we used a TF-IDF vectorization model with 200,000 features and 1-3 words per feature, and an SVC model for prediction, with a linear kernel and C=1.5.

Following the method that Wang et al. used to prepare the prediction data for the bootstrap postpi method, we used our prediction model to generate the probability distribution for each tweet–the probability of it being classified as Democratic or Republican-leaning–and used this data and the observed outcomes from the test dataset to build a relationship model. We used a K-NN machine learning model for this as we found it to describe the relationship between the predicted and observed outcomes well compared to other models like logistic regression.

### 4.2 Feature Selection for Inference

We reviewed relevant literature in political science to develop a criteria for choosing our features.

In *Twitter Language Use Reflects Psychological Differences between Democrats and Republicans*, Sylwester and Purver discuss the differences between Democrats and Republicans in the context of previous findings and their own discoveries. For example, Haidt's Moral Foundations model, which identifies "harm, fairness, liberty, ingroup, authority, and purity" as

the pillars of morality, has been used to distinguish between liberals and conservatives. It was found that liberals prioritized the harm and fairness aspects of morality, while conservatives focused more on liberty, ingroup, authority, and purity. Sylwester and Purver also found differences between Democratic and Republican-aligned people when it came to what kinds of topics they discussed and emotions they expressed–Republicans focused more on topics such as "religion…, national identity…, government and law…, and their opponents" while Democrats were focused on emphasizing their uniqueness and generally expressed more anxiety and emotion. These findings are somewhat in line with our own observations made through the data–as stated before, we found that Republican tweets made references to their opponents on a much larger scale than Democrats, and also made mention of the American people–their national identity–plenty of times as well.

We also reviewed Chen et al.'s study, *#Election2020: the first public Twitter dataset on the 2020 US Presidential election*. Chen et al. found that more conservative Twitter users tended to share more topics related to conspiracy theories and "public health and voting misinformation" compared to liberal Twitter users.

Taking these two sources into consideration, our criteria for selecting features was whether or not they would fall into either liberal or conservative tendencies as discovered by either source. If a feature implied a discussion of harm or fairness, or was an expression of uniqueness, anxiety, or emotion, then we anticipated that this feature would connect more to Democratic-aligned tweets. On the other hand, if a feature discussed liberty, purity, religion, national identity, government and law, or Republican opponents, or implied that the topic at hand was associated with public health or voting misinformation, said feature may be connected to Republican-aligned tweets.

We ended up selecting 5 features to conduct inference, which are border, illegal, god, defund, and happy. We hypothesized that the first three would be strong indicators for a Republican-classified tweet as they allude to national identity, law, and religion, while the last two would indicate a Democratic-classified tweet as they allude to concepts of harm and fairness, as well as emotion.

## 5 Results

After conducting inference using the bootstrap postpi algorithm, we found that the parametric bootstrap method worked best to correct for inference. As such, for the inference we interpret below we will only be considering the corrections made using the parametric method, and not the non-parametric bootstrap method.

## 5.1 Inference on 'border'

The table below shows the results of conducting inference on the word 'border'. The bootstrap postpi algorithm corrects coefficients, SEs, and t-statistics as mentioned above and the results below shows that the algorithm works as intended. The true beta coefficient has a value of 7.491, but in the case that we didn't have the observed values, using the bootstrap postpi algorithm would correct the coefficient to 7.498. The corrected value is a better estimate for the coefficient compared to the no correction approach value of 8.272. The coefficient was corrected by an absolute difference of 0.007. The SE for a no correction approach results in an absolute difference of 0.018 to the true value, but after correction, the absolute difference decreases to 0.011. The t-statistic for the no correction approach results in an absolute difference of 0.72 while the corrected approach resulted in an absolute difference of 0.125. These results are meaningful because the smaller differences would suggest that we have a good bootstrap model that corrects inference using predicted values instead of observed values.

A positive coefficient for the word 'border' implies that this feature is a good predictor for the Republican party. To compare how much better the correction was on the coefficient we can compute the odds ratio. Using the actual value, Republicans are approximately $e^{7.491} = 1791.843$ times more likely to use the word 'border' than Democrats. The odds when using the coefficient with no correction would tell us that Republicans are approximately $e^{8.272} = 3912.767$ times more likely to use the word 'border' than Democrats, which is over 2000 times more than the actual odds. The odds when we use the corrected coefficient would tell us that Republicans are approximately $e^{7.498} = 1804.430$ times more likely to use the word 'border' than Democrats, which is relatively close to the odds of the true coefficient.

To test whether the feature is a statistically significant predictor we must evaluate the t-statistic. If the null hypothesis was true–that there is no significant difference between Republicans and Democrats in their use of the word 'border'–then we would expect a sample with no difference. Since the corrected t-statistic of ~ 8.966 is greater than 2, we have 95% confidence that there is a positive difference between our sample data and the null hypothesis. This implies that the word 'border' is a good predictor for the Republican party.

| Feature: border | Actual Values | No Correction | Non-Parametric | Parametric |
|---|---|---|---|---|
| Coefficient | 7.491279693899731 | 8.272380684108418 | 7.497679809201015 | 7.497679809201015 |
| SE | 0.8473168827373054 | 0.8652117033394918 | 0.38892459925697676 | 0.8361886115904971 |
| T-Stat | 8.841178367293622 | 9.561105856727531 | 19.27797784846986 | 8.966493570080837 |

## 5.2 Inference on 'illegal'

The table below shows the results of conducting inference on the word 'illegal'. The true beta coefficient has a value of 5.790, but in the case that we did not have the observed values, using the bootstrap postpi algorithm would correct the coefficient to 5.832. The corrected value is a better estimate for the coefficient compared to the no correction approach value of 6.392. The SE for the no correction approach results in an absolute difference of 0.004 but after running the bootstrap postpi algorithm, the absolute difference decreased to 0.003. The t-statistic for a no correction approach results in an absolute difference of 0.529 while the corrected absolute difference resulted in 0.051. These results are meaningful because the smaller differences would suggest that we have a good bootstrap model that corrects inference using predicted values instead of observed values.

A positive coefficient for the word 'illegal' implies that this feature is a good predictor for the Republican party. To compare how much better the correction was on the coefficient we can compute the odds ratio. Using the actual value, Republicans are approximately $e^{5.790} = 327.013$ times more likely to use the word 'illegal' than Democrats. The odds when using the coefficient with no correction would tell us that Republicans are approximately $e^{6.392} = 597.050$ times more likely to use the word 'illegal' than Democrats, which is over 200 times more than the actual odds. The odds when we use the corrected coefficient would tell us that Republicans are approximately $e^{5.832} = 341.040$ times more likely to use the word 'illegal' than Democrats, which is relatively close to the odds ratio calculated from the true coefficient.

To test whether the feature is a statistically significant predictor we must evaluate the t-statistic. If the null hypothesis was true–that there is no significant difference between Republicans and Democrats in their use of the word 'illegal'– then we would expect a sample with no difference. Since the corrected t-statistic of ~ 5.335 is greater than 2, we have 95% confidence that there is a positive difference between our sample data and the null hypothesis. This implies that the word 'illegal' is a good predictor for the Republican party.

| Feature: illegal | Actual Values | No Correction | Non-Parametric | Parametric |
|---|---|---|---|---|
| Coefficient | 5.790370678304617 | 6.3923191188455535 | 5.832335983568017 | 5.832335983568017 |
| SE | 1.0957382919323564 | 1.0996851060116388 | 0.36702547247366996 | 1.0931272276718014 |
| T-Stat | 5.2844467706729334 | 5.812863231392987 | 15.890820722222278 | 5.33545943777288 |

### 5.3 Inference on 'god'

The table below shows the results of conducting inference on the word 'god'. The true beta coefficient has a value of 4.897, but in the case that we didn't have the observed values, using the bootstrap postpi algorithm corrects the coefficient to 4.779. The corrected value is a better estimate for the coefficient compared to the no correction approach value of 5.447. The coefficient was corrected by an absolute difference of 0.55. The SE for a no correction approach results in an absolute difference of 0.007 to the true value, but after correction, the absolute difference increased to 0.018. The t-statistic for the no correction approach results in an absolute difference of 0.493 while the corrected approach resulted in an absolute difference of 0.036. These results are meaningful because the smaller differences would suggest that we have a good bootstrap model that corrects inference using predicted values instead of observed values.

A positive coefficient for the word 'god' implies that this feature is a good predictor for the Republican party. To compare how much better the correction was on the coefficient we can compute the odds ratio. Using the actual value, Republicans are approximately $e^{4.897} = 133.888$ times more likely to use the word 'god' than Democrats. The odds calculated using the coefficient with no correction would tell us that Republicans are approximately $e^{5.447} = 232.061$ times more likely to use the word 'god' than Democrats, which is more than 100 times the actual odds. The odds when we use the corrected coefficient would tell us that Republicans are approximately $e^{4.779} = 118.985$ times more likely to use the word' god' than Democrats, which is relatively close to the odds of the true coefficient.

To test whether the feature is a statistically significant predictor we must evaluate the t-statistic. If the null hypothesis was true–that there is no significant difference between Republicans and Democrats in their use of the word 'god'– then we would expect a sample with no difference. Since the corrected t-statistic of ~ 4.720 is greater than 2, we have 95% confidence that there is a positive difference between our sample data and the null hypothesis. This implies that the word 'god' is a good predictor for the Republican party.

| Feature: god | Actual Values | No Correction | Non-Parametric | Parametric |
|---|---|---|---|---|
| Coefficient | 4.896751845067567 | 5.446532227871352 | 4.77884444860249 | 4.77884444860249 |
| SE | 1.0296865349554185 | 1.0376801856497464 | 0.37360565349958674 | 1.0125472708370837 |
| T-Stat | 4.755575292901716 | 5.248758050112512 | 12.791145968586838 | 4.719626022646595 |

**5.4 Inference on 'defund'**

The table below shows the results of conducting inference on the word 'defund'. The true beta coefficient has a value of 1.181, but in the case that we didn't have the observed values, using the bootstrap postpi algorithm would correct the coefficient to 1.076. The corrected value is a better estimate for the coefficient compared to the no correction approach value of 1.511. The coefficient was corrected by an absolute difference of 0.105. The SE for a no correction approach results in an absolute difference of 0.002 but after running the bootstrap postpi algorithm, the absolute difference decreased to 0.0001. The T-Statistic for a no correction approach results in an absolute difference of 0.173 while the corrected absolute difference resulted in 0.055. These results are meaningful because the smaller differences would suggest that we have a good bootstrap model that corrects inference using predicted values instead of observed values.

To compare how much better the correction was on the coefficient we can compute the odds ratio. Using the actual value, Republicans are approximately $e^{1.181} = 3.26$ times more likely to use the word 'defund' than Democrats. The odds when using the coefficient with no correction would tell us that, Republicans are approximately $e^{1.511} = 4.53$ times more likely to use the word 'defund' than Democrats, which is close to the actual odds. The odds when we use the corrected coefficient would tell us that, Republicans are approximately $e^{1.076} = 2.933$ times more likely to use the word 'defund' than Democrats, which is closer to the odds calculated using the true coefficient.

Interestingly, conducting inference on the feature 'defund' yielded a positive coefficient, which implies that this feature is a good predictor for the Republican party, and not the Democratic party contrary to our hypothesis.

To test whether the feature is a statistically significant predictor we must evaluate the t-statistic. If the null hypothesis was true–that there is no significant difference between Republicans and Democrats in their use of the word 'defund'– then we would expect a sample with no difference. Since the corrected T-Statistic of ~ 0.560 is less than 2 and greater than -2, we have 95% confidence that there is not a positive difference between our sample data and the null hypothesis. This implies that the word 'defund' is not a good predictor for the Republican party.

| Feature: defund | Actual Values | No Correction | Non-Parametric | Parametric |
|---|---|---|---|---|
| Coefficient | 1.1809699288304498 | 1.510615350109516 | 1.0756616269881318 | 1.0756616269881318 |
| SE | 1.9195712586934353 | 1.9179793839463273 | 0.39136574186049267 | 1.919645058248869 |
| T-Stat | 0.6152258862399733 | 0.787607709839591 | 2.7484818213127227 | 0.560344018997641 |

## 5.5 Inference on 'happy'

The table below shows the results of conducting inference on the word 'happy'. The true beta coefficient has a value of 0.935, but in the case that we didn't have the observed values, using the bootstrap postpi algorithm would correct the coefficient to 0.959. The corrected value is a better estimate for the coefficient compared to the no correction approach value of 1.137. The SE for a no correction approach results in an absolute difference of 0.001 but after running the bootstrap postpi algorithm, the absolute difference increased to 0.0045. The T-Statistic for a no correction approach results in an absolute difference of 0.406 while the corrected absolute difference resulted in 0.0321. These results are meaningful because the smaller differences would suggest that we have a good bootstrap model that corrects inference using predicted values instead of observed values.

To compare how much better the correction was on the coefficient we can compute the odds ratio. Using the actual value, republicans are approximately $e^{0.935} = 2.547$ times more likely to use the word 'happy' than Democrats. The odds when using the coefficient with no correction would tell us that, Republicans are approximately $e^{1.137} = 3.117$ times more likely to use the word 'happy' than Democrats, which is close to the actual odds. The odds when we use the corrected coefficient would tell us that, Republicans are approximately $e^{0.959} = 2.609$ times more likely to use the word 'happy' than Democrats, which is closer to the odds of the true coefficient.

Once again, we find that inference on the feature 'happy' also yielded a positive coefficient, which tells us that this feature is a good predictor for the Republican party, and not the Democratic party. This is, again, contrary to what we hypothesized would be the case.

To test whether the feature is a statistically significant predictor we must evaluate the t-statistic. If the null hypothesis was true–that there is no significant difference between Republicans and Democrats in their use of the word 'happy'– then we would expect a sample with no difference. Since the corrected T-Statistic of ~ 1.920 is less than 2 and greater than -2, we have 95% confidence that there is no positive difference between our sample data and the null hypothesis. This implies that the word 'happy' is not a good predictor for the Republican party.

| Feature: happy | Actual Values | No Correction | Non-Parametric | Parametric |
|---|---|---|---|---|
| Coefficient | 0.9349812363166592 | 1.1370802106872728 | 0.9594110548987034 | 0.9594110548987034 |
| SE | 0.49531091025219376 | 0.4956741362008665 | 0.43422431723594523 | 0.49975429846889297 |
| T-Stat | 1.8876653369912684 | 2.2940075498038164 | 2.209482557323906 | 1.9197654884371576 |

# 6 Conclusion

We can conclude from the results above that the bootstrap postpi algorithm is promising for providing post-prediction inference correction for text data and for the field of political science. The correction is small but keep in mind that the goal of the bootstrap postpi algorithm is to correct inference using predicted values, and that has succeeded.

We tested whether the features 'border', 'illegal', and 'god' would be strong indicators for a Republican-classified tweet as they allude to national identity, law, and religion. Similarly we tested that the other two features, 'defund' and 'happy', would indicate a Democratic-classified tweet as they allude to concepts of harm and fairness, as well as emotion.

Based on the results, 'border', 'illegal', and 'god' are strong predictors towards classifying a republican tweet. On the contrary, 'defund' and 'happy' were not strong predictors towards classifying a democratic tweet or republican tweet. Thus we can argue that the hypotheses we derived from literature review were half correct. There is statistical evidence that Republicans tend to use language that alludes to national identity, law, and religion but there is not enough evidence to argue that Democratic tweets tend to discuss concepts of harm, fairness and emotion.

These findings are applicable in several ways. For example, politicians can use our discoveries to appeal to Republicans' values during a campaign by using language that alludes to national identity, law, and religion. The ability to resonate with people is vital for a politician to win elections and to advance their agenda. Minute changes in language to appeal to their audience can be the difference between a winning or losing campaign.

That being said, there are some limitations. One limitation of Wang et al.'s bootstrap algorithm is that the feature has to be present in the training/testing set, thus inference on features that are not present in the training or testing dataset is not able to be done. Another limitation is that the bootstrap algorithm may have unexpected results in text data, owing to the fact that some selected features may not appear much at all in the data. It is also important to consider how one chooses to pre-process text data; the removal of certain stem words or suffixes can lead to unexpected results.

In conclusion, we have demonstrated that the bootstrap postpi algorithm first developed by Wang et al. is shown to correct predicted outcomes when observed outcomes are not available on political data. In such a field where collecting observed outcomes can be exceedingly time-consuming and expensive to collect, this is a significant finding that may open doors to some studies that may otherwise be too difficult to conduct.

# 7 References

Chen, E., Deb, A. & Ferrara, E. #Election2020: the first public Twitter dataset on the 2020 US Presidential election. *J Comput Soc Sc* (2021). https://doi.org/10.1007/s42001-021-00117-9

Sylwester K, Purver M (2015) Twitter Language Use Reflects Psychological Differences between Democrats and Republicans. PLOS ONE 10(9): e0137422. https://doi.org/10.1371/journal.pone.0137422

Wang, Siruo, Tyler H. McCormick, and Jeffrey T. Leek. "Methods for correcting inference based on outcomes predicted by machine learning." *Proceedings of the National Academy of Sciences* 117.48 (2020): 30266-30275.