

Goal and Motivation

34 mins · Stanford, CA

Goal: Predict sexual orientation from Facebook status updates.

Motivation: We want to examine the hypothesis that people with different sexual orientations express themselves differently on social media. Combining our results with our CS 221 Project, which extracted gender features from status updates, we seek to test the stereotype that male homosexuals tend to use more feminine language.

Like Comment Share

420

Data

34 mins · Stanford, CA

- We used data from myPersonality.org, with kind permission from Dr. Michal Kosinski (Stanford GSB), which contains 22M Facebook status updates and included demographic details (e.g. gender) of every user in the dataset.
- We derived the sexual orientation labels by looking at the gender of a user's partner, and comparing it to the user's gender.
- Word stemming was applied on the status updates.
- Our dataset is skewed in a 9:1 ratio. As such, our test error did not provide a meaningful sense of how our model performed, and we used alternative measures like F-1 score and ROC curves instead.

Like Comment Share

229

Features and Models

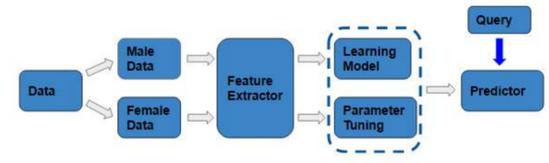
34 mins · Stanford, CA

Features:

- N-grams (tuned across a range of hyper-parameters)
- Counts of periods, exclamation marks, smileys and capital letters.

Learning Algorithms:

- Support Vector Machine
- Multinomial Naïve Bayes
- Logistic Regression
- Random Forest



Like Comment Share

1337

Results

34 mins · Stanford, CA

| Model | Males | | Females | |
|---------------------|---------|---------------------|---------|---------------------|
| | ROC AUC | F1 Score* | ROC AUC | F1 Score* |
| Logistic Regression | 0.57 | 0.92 (0.97,0.20) | 0.62 | 0.84 (0.94,0.24) |
| Naïve Bayes | 0.52 | 0.91 (0.96,0.21) | 0.58 | 0.84 (0.93,0.30) |
| SVM | 0.61 | 0.94 (0.98,0.17) | 0.62 | 0.85 (0.93,0.36) |
| Random Forest | 0.55 | - | 0.63 | - |

*For F1 Score, the figures in parentheses indicate F1-scores for heterosexuals and homosexuals respectively.

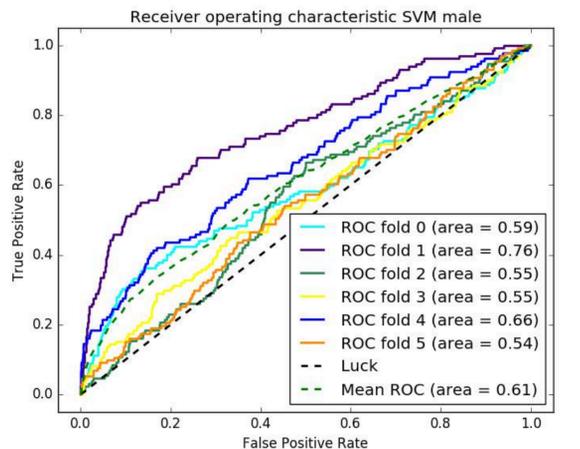
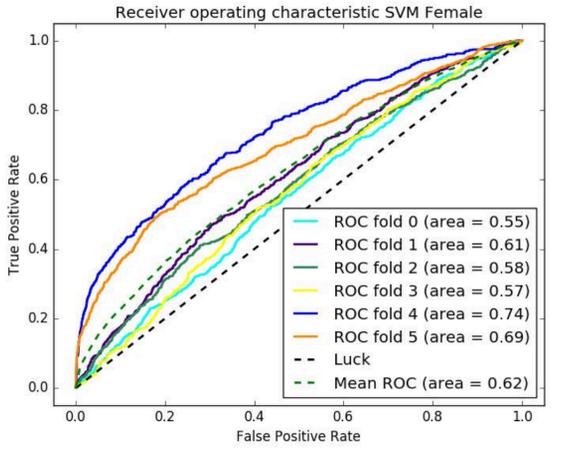
SVM Model Parameters Tuning

| | | | |
|------------------------|--------|-------|-------|
| N-gram range | (1,2) | (1,4) | (1,5) |
| Min document frequency | 1 | 0.95 | 0.9 |
| Max document frequency | 1 | 0.95 | 0.9 |
| Kernel | Linear | Poly | Rbf |

Like Comment Share

221

Receiver operating characteristic SVM

Confusion matrix, without normalization

| True Label \ Predicted Label | Heterosexual | Homosexual |
|------------------------------|--------------|------------|
| heterosexual | 51337 | 1457 |
| homosexual | 3240 | 569 |

Top Word Features

| Gender | Heterosexual | Homosexual |
|--------|--|--|
| Male | angel, mile, game, wife, drive, gotta, girlfriend, wat, texas, goal | sister, yay, okay, gay, funni-, hair, omg, b*tch, comment, sex |
| Female | pray, work, bless, church, hubbi-, clean, famili-, boyfriend, husband, today | your, random, hahaha, aint, mum, wanna, gunna, test, drunk |

Like Comment Share

952

Analysis and Future Work

34 mins · Stanford, CA

Analysis

- ROC scores for both females and males were above 60%, which told us that there were distinctions in how homosexuals expressed themselves on social media, even if the distinction was not great enough to consistently predict one's sexual orientation.
- Mentions of another partner of the opposite gender (e.g. when males mention 'wife') are strong indicators that a person is heterosexual.
- Our model showed that a male homosexual was 4 times more likely to use the word "gay". In fact, a male who mentions "gay" in a status update has a 1 in 4 chance of being homosexual.

Limitations

- Age could be a confounding factor that is driving the differences between homosexuals and heterosexuals. For example, it may be popular for young girls to declare on Facebook that there are in a "relationship" with a good friend if they are heterosexual. Also, the top word features for female homosexuals are more associated with young people (e.g. omg, sex, b*tch), whereas that for heterosexuals are more associated with older people (e.g. husband, church, work).

Applying the gender model to our data

- We had earlier built and trained models to predict gender for our CS 221 Project. We wanted to use these models to test the stereotype that male homosexuals expressed themselves in a more effeminate manner.
- We found that our gender model predicted that 45% of all male homosexuals were female. It predicted that 40% of all male heterosexuals were female. This meant that a male homosexual was 5 percentage points more likely to be predicted female.
- Our results suggest that there is slight evidence that male homosexuals express themselves more like females, as compared to male heterosexuals. However, the evidence is not strong enough to support the social stereotype.

Future Work

- Exploration of other methods of feature extraction (e.g. Word2Vec), and more nuanced feature engineering. We can also use neural networks to automatically learn features in the data.

Like Comment Share

952