# What makes TED talks worth spreading?

Exploration of the linguistic and text features of TED video transcripts

● ● ●

Ozge Yasar

# Goals of the project

- Using linguistic & non-linguistic features such as film date, inspiring and

  persuasive votes etc. to:

  - Build a classifier in order to predict the gender of the presenters

  - Predict the number of views for each video

- Create a classification model to accurately classify the multiple topic tags of each

  video from text (bag of words)

# Data Set

Shared on data.world - Words of Persuasion: Text Predictors of Persuasive TED Talks*

Initial Size:

- 2406 TED Talk scripts scraped from TED website.
- 187 features

Video features included:

- Number of views, comments, votes (e.i., Persuasive, Unconvincing),
- Video transcripts
- Tags, occupations of the speakers

Linguistic Inquiry Word Count (LIWC**) features:

- Linguistic Dimensions (use of pronouns, punctuation, etc.)
- Other Grammar (verbs, adjectives, etc.)
- Psychological Processes (Affective processes, Cognitive processes, Time orientations , etc.)
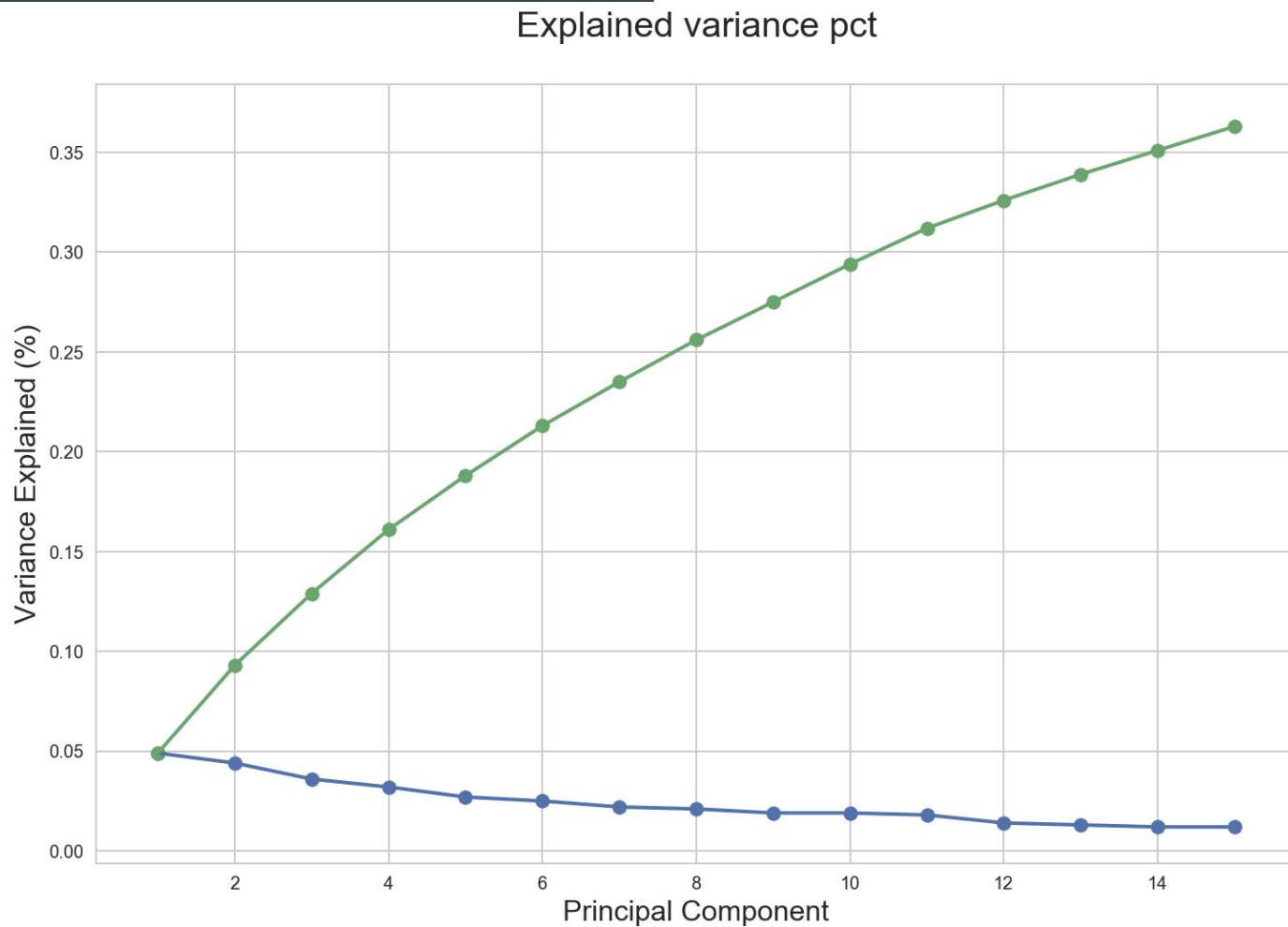
Morality features***

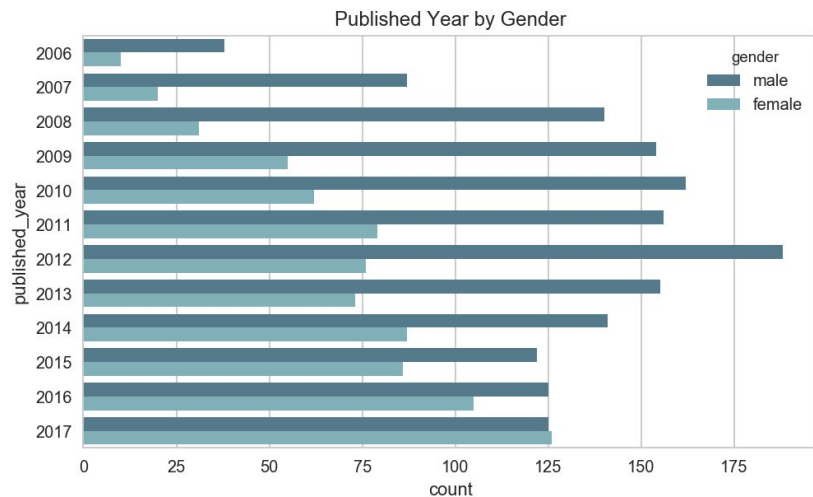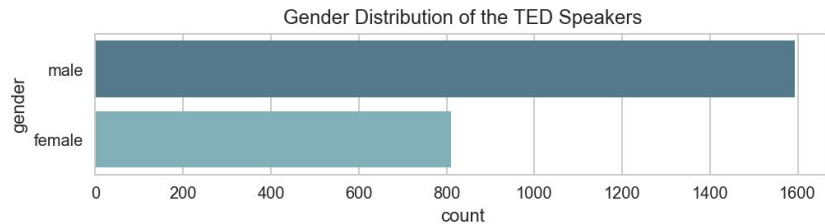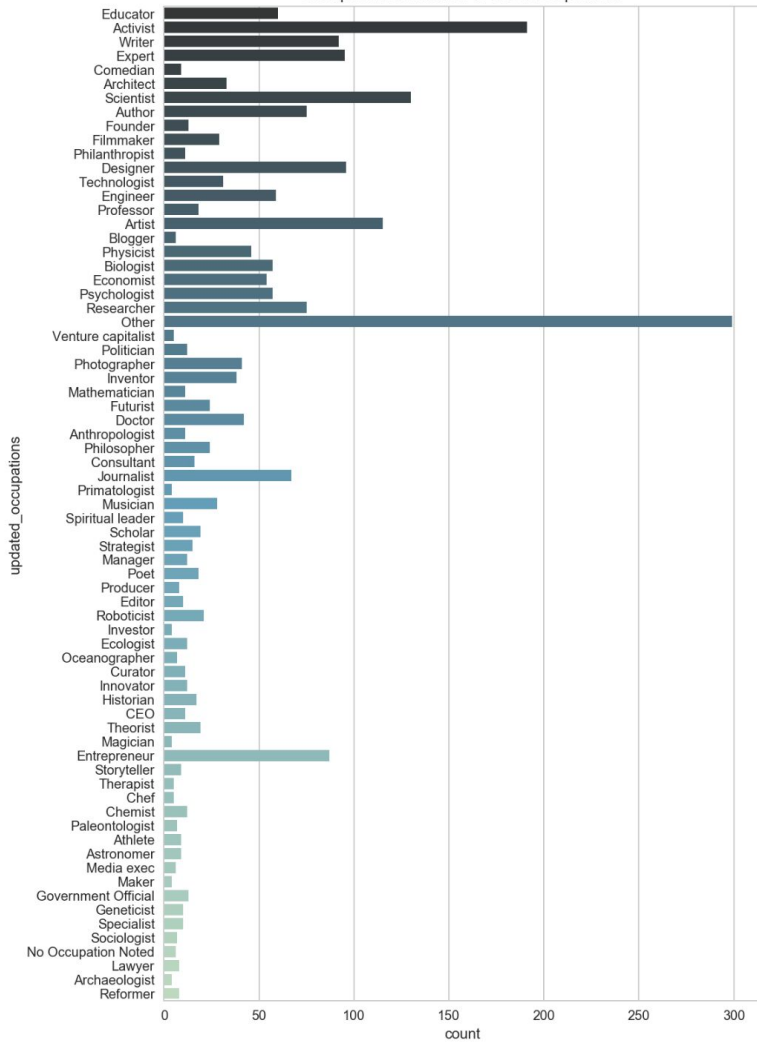- Harm. Fairness, Purity & sub categories

Explained variance pct

# Feature Distributions

# Gender Distributions

Occupation Distribution of the TED Speakers

Occupation Distribution of the TED Speakers by Gender

# Video Tags - I

# Video Tags - II

# Video Tags - III

# Distribution of Views

# Cleaning Process

Highly correlated continuous variables were removed. For example:

- **Fairness (general)**
- Fairness virtue, Fairness Vice

Log transformed outcome variable: **views**

Outliers were removed

Outside +/- 10 Std range

# Models

# View Prediction

# Predicting Number of Views

Models used:

- Random Forest Regressor
- LinearSVR
- Linear Regression
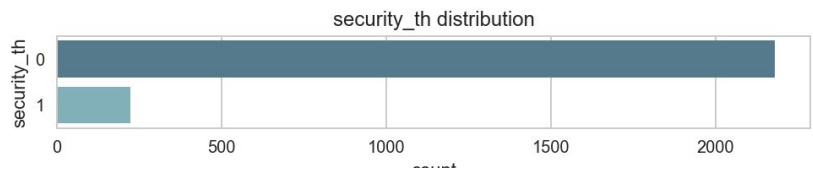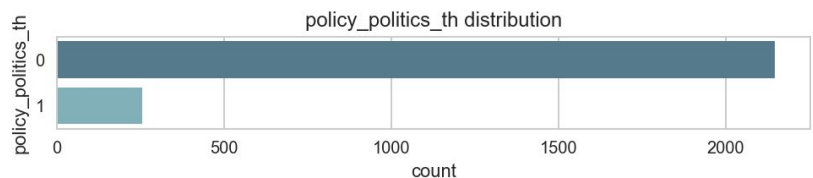- LassoCV
- Ridge CV
- MLPRegressor

| Importance | Variable |
|:---:|:---:|
| 0.32 | Languages |
| 0.086 | Film date |
| 0.083 | Normalized unconvincing votes |
| 0.053 | Word count |
| 0.032 | Words that indicates differentiation |
| 0.02 | Number of positive emotion words |

# Gender Prediction

# Predicting

Models Used:

- Random
- KNN C
- Logistic

y = 0.72

## Receiver operating characteristic: Salary over median



|  | Precision | Recall | f1-score | Support |
|---|---|---|---|---|
| **female** | 0.61 | 0.39 | 0.48 | 179 |
| **male** | 0.75 | 0.88 | 0.81 | 375 |
| **avg/total** | 0.71 | 0.72 | 0.70 | 554 |

# Topic Classification

# TED Words

# Text Processing

- TextBlob
- Words Lemmatized
- CountVectorizer
- Tf-idf transformer

# Multilabel Classification

Models used:

- OneVSRestClassifier(KNN)

- Jaccard_similarity_score (accuracy) = 0.37
- Baseline (mean of y_train matrix) = 0.19

- Model Micro-Average Precision Score = 0.39
- Baseline Micro-Average Precision Score = 0.19

# Individual Label Accuracy

| | Base Line | Model |
|---|---|---|
| art_and_design_th | 0.351367 | 0.785021 |
| business_th | 0.247919 | 0.822469 |
| continents_cities_countries_th | 0.103448 | 0.897365 |
| educ_learn_youth_th | 0.097503 | 0.898752 |
| emotions_psychology_th | 0.124257 | 0.901526 |
| entertainment_th | 0.144471 | 0.886269 |
| gender_matters_th | 0.068966 | 0.955617 |
| health_and_medicine_th | 0.258026 | 0.819695 |
| human_th | 0.313317 | 0.719834 |
| literature_th | 0.082640 | 0.938974 |
| living_spaces_cult_th | 0.294887 | 0.744799 |
| media_th | 0.077289 | 0.925104 |
| music_th | 0.024376 | 0.980583 |
| nature_environment_th | 0.383472 | 0.753121 |
| policy_politics_th | 0.102854 | 0.918169 |
| religions_spirituality_th | 0.048751 | 0.955617 |
| science_and_research_th | 0.374554 | 0.757282 |
| security_th | 0.088585 | 0.918169 |
| society_th | 0.353151 | 0.643551 |
| sports | 0.015458 | 0.988904 |
| technology_th | 0.480975 | 0.679612 |
| tedmem_awardwin_th | 0.272889 | 0.721221 |

# Precision-Recall Curves



Extension of Precision-Recall curve to multi-class

Art & Design

Nature & Environment

Technology

Science & Research

— Precision-recall for class technology_th (area = 0.63)
— Precision-recall for class society_th (area = 0.43)
— Precision-recall for class policy_politics_th (area = 0.37)
— Precision-recall for class security_th (area = 0.29)
— Precision-recall for class science_and_research_th (area = 0.60)
— Precision-recall for class nature_environment_th (area = 0.59)
— Precision-recall for class living_spaces_cult_th (area = 0.39)
— Precision-recall for class religions_spirituality_th (area = 0.23)
— Precision-recall for class tedmem_awardwin_th (area = 0.25)
— Precision-recall for class emotions_psychology_th (area = 0.27)
— Precision-recall for class business_th (area = 0.40)
— Precision-recall for class entertainment_th (area = 0.33)
— Precision-recall for class media_th (area = 0.13)
— Precision-recall for class music_th (area = 0.16)
— Precision-recall for class health_and_medicine_th (area = 0.51)
— Precision-recall for class art_and_design_th (area = 0.58)
— Precision-recall for class gender_matters_th (area = 0.48)
— Precision-recall for class educ_learn_youth_th (area = 0.20)
— Precision-recall for class sports (area = 0.01)
— Precision-recall for class literature_th (area = 0.19)
— Precision-recall for class human_th (area = 0.35)
— Precision-recall for class continents_cities_countries_th (area = 0.24)

# Next Steps

- Processing the text with Gensim Doc2Vec, Word2Vec

- Other algorithms to tackle multi-classification problem: Classifier Chain, Label Powerset from skmultilearn