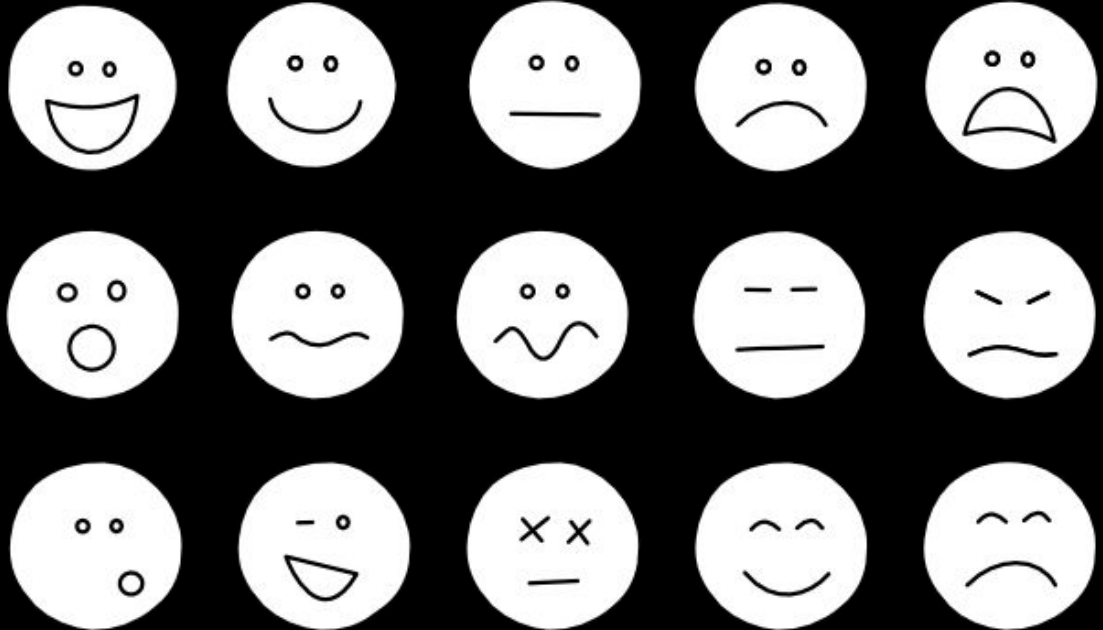


ProQuest Sentiment Analysis MDP

Taylor Murray

Honors Capstone

Dec 16th, 2020





Context

Our Sponsor

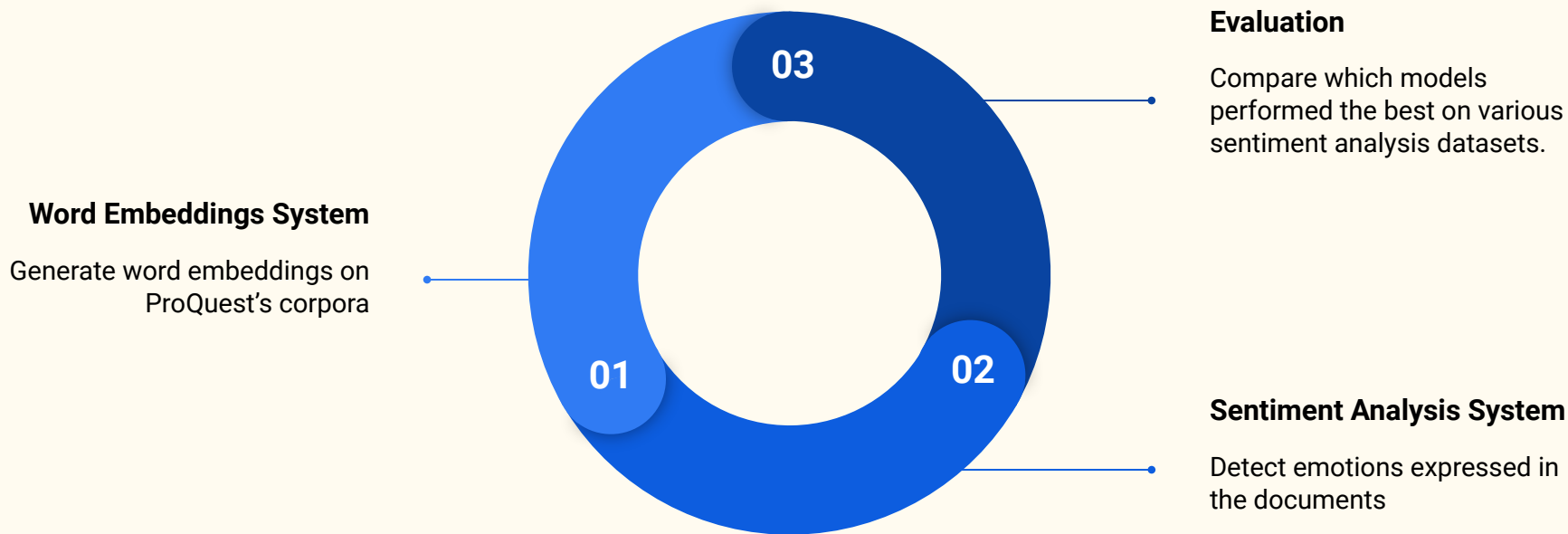
- Education technology company
- Content collection that encompasses 90,000 authoritative sources, 6 billion digital pages and spans 6 centuries

Motivations

- Enhance user experience in the TDM (text-data-mining) environment
- Improve ProQuest's internal use of data

High Level Goals

Help ProQuest users (commonly academic researchers) create word embeddings and predict sentiment on their own datasets.



Word Embeddings

Word embeddings (WEs) are representations of words in the vector space.

Dimensionality of the vector is a hyperparameter, typically chosen to be between 50-300.

Each word embedding takes up constant space.

	living being	feline	human	gender	royalty	verb	plural
<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8

Properties of Word Embeddings

Semantic property: Related to a word's *meaning*

e.g. **doctor, child, woman** → human

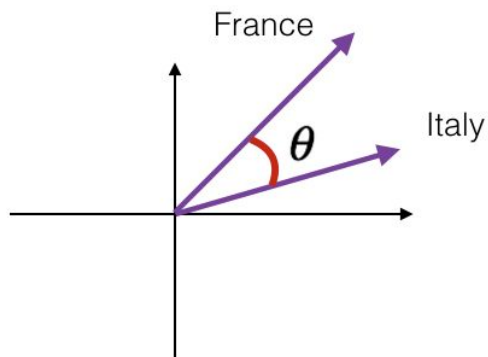
Syntactic property: Related to a word's *structure* and *grammatical features*

e.g. **geese, children, apples** → plural nouns; **ran, helped, changed** → past tense

Linear Behaviour

$\text{Vector}(\textit{King}) - \text{Vector}(\textit{Man}) + \text{Vector}(\textit{Woman}) \approx \text{Vector}(\textit{Queen})$

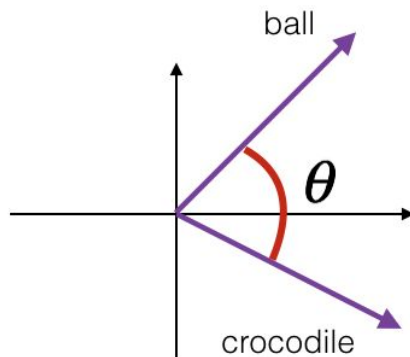
Cosine Similarity



France and Italy are quite similar

θ is close to 0°

$\cos(\theta) \approx 1$



ball and crocodile are not similar

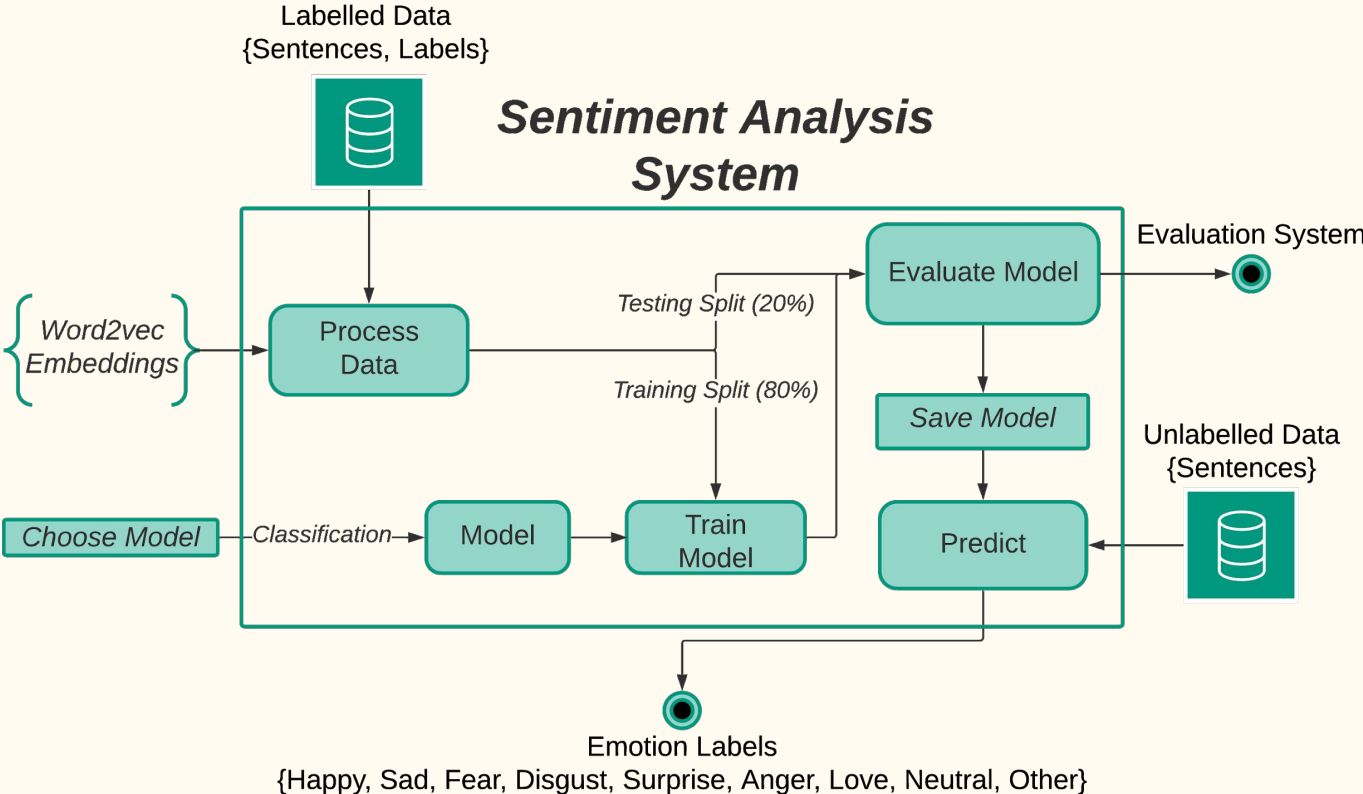
θ is close to 90°

$\cos(\theta) \approx 0$ (2)

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

Calculate similarity between different words

Sentiment Analysis System

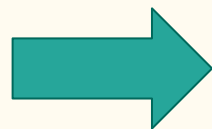


Emotion prediction

Goal: Classify text as having one of nine sentiments

- Happy
- Sad
- Fear
- Disgust
- Surprise
- Anger
- Love
- Neutral
- Other

**Ekman's 6
Basic
Emotions⁽³⁾**



Example Text:



Sentiment:

Happy

ProQuest corpora

1. New York Times
2. Book Blurbs
3. LION (Literature Online) Poems

The
New York
Times



Part of the Literature Online family

Key Question:

Can using **in-domain** corpora for word embeddings improve sentiment analysis performance?

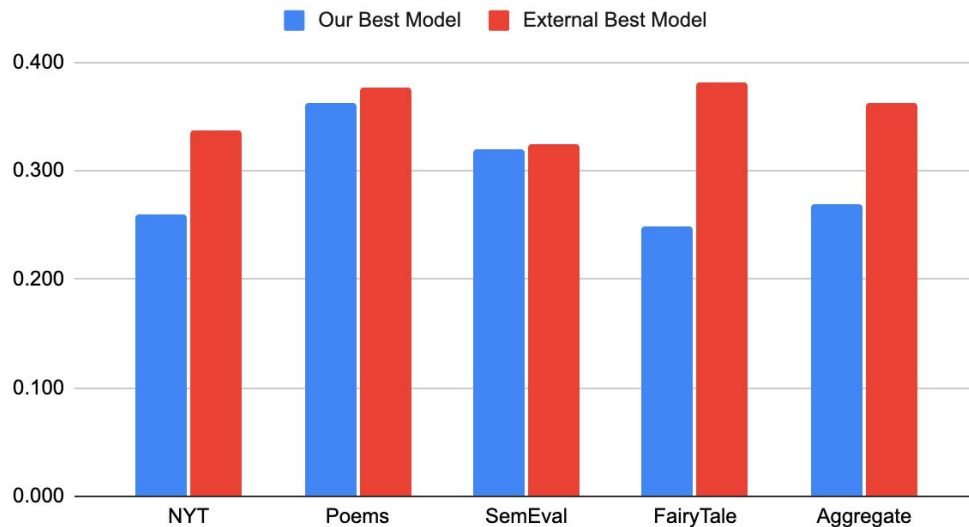
What did we discover?

We tested 18 models, 9 of which used custom word embeddings.

We used different sampling techniques, filtering methods, ML models, and sentence embeddings.

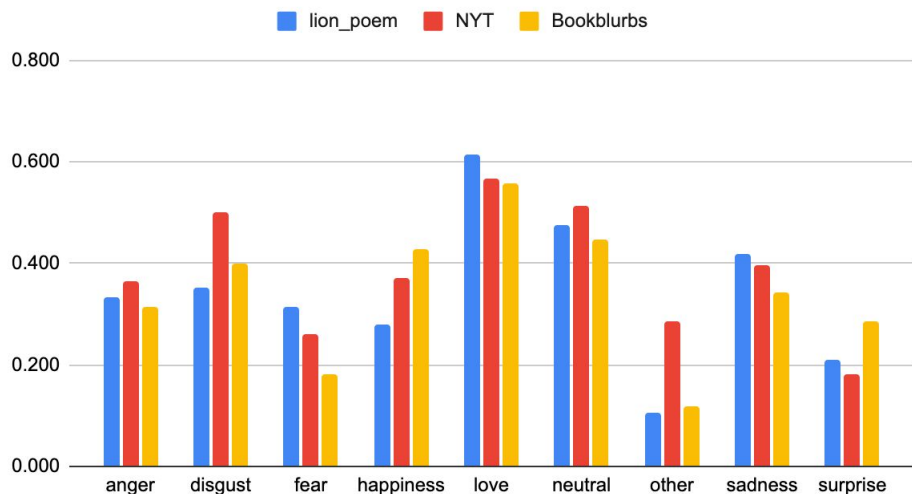
In the end, it seemed conclusive that more generally trained models performed the best.

Macro-F1 Scores for Sentiment Analysis

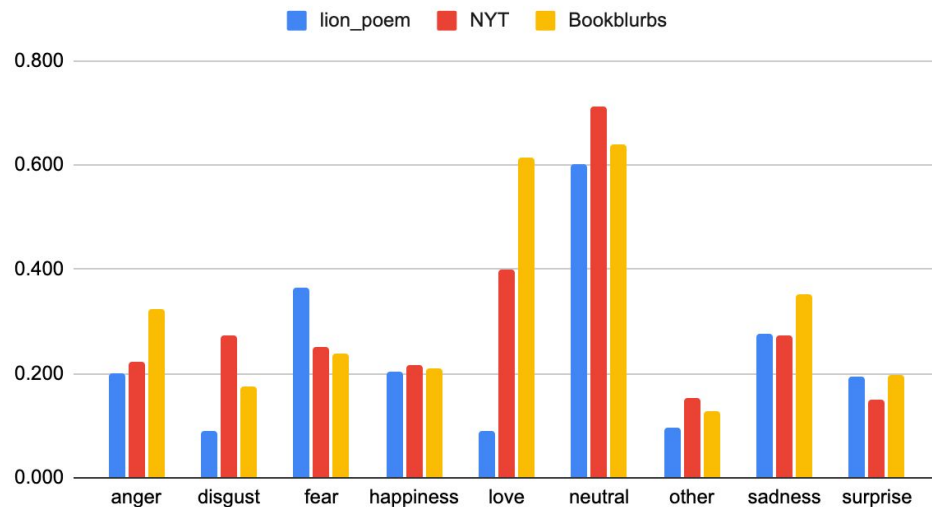


Using in-domain corpora not correlated with better performance on sentiment tasks

Lion Poems



NYT Data



End Deliverables

- Made Jupyter Notebooks for:
 - Word Embeddings Generation
 - Word Embeddings Usage
 - Predicting affective state (“happiness”, “sadness”, etc.)
 - Predicting valence state (“very negative”, “positive”, etc.)
 - Visualizations for Sentiment Analysis Results
- Set of word embeddings
- Set of sentiment analysis models
- Emotion-labelled data

Lesson: You should record **all** of your parameters and brainstorm new ones you might try later

The parameters we noted:

The parameters we failed to note:

```
parameters
-----
Note: param2 is deprecated because we haven't named all the
      feed in an array of paths to the models

param1: a list of datasets ['Stanford', 'semEval_emotions']
param2: choose one WE model ['cbow', 'sg-lion', 'sg-1m', 'sg
param3: a list of SEs ['AVG', 'SIF', 'concat', 'SBERT']
param4: choose ML models -> ['LR', 'SVM', 'Tree', 'LSTM']

return value: Dict['param1']['{param2}_{param3}']['param4']
Exceptions: when param4 == LSTM, then Dict[param1][param2][f
"""
```

- What data was trained on
- If there was over/undersampling done
- Whether or not the data trained on was stratified
- What the ML model was maximizing (F1-score or accuracy)

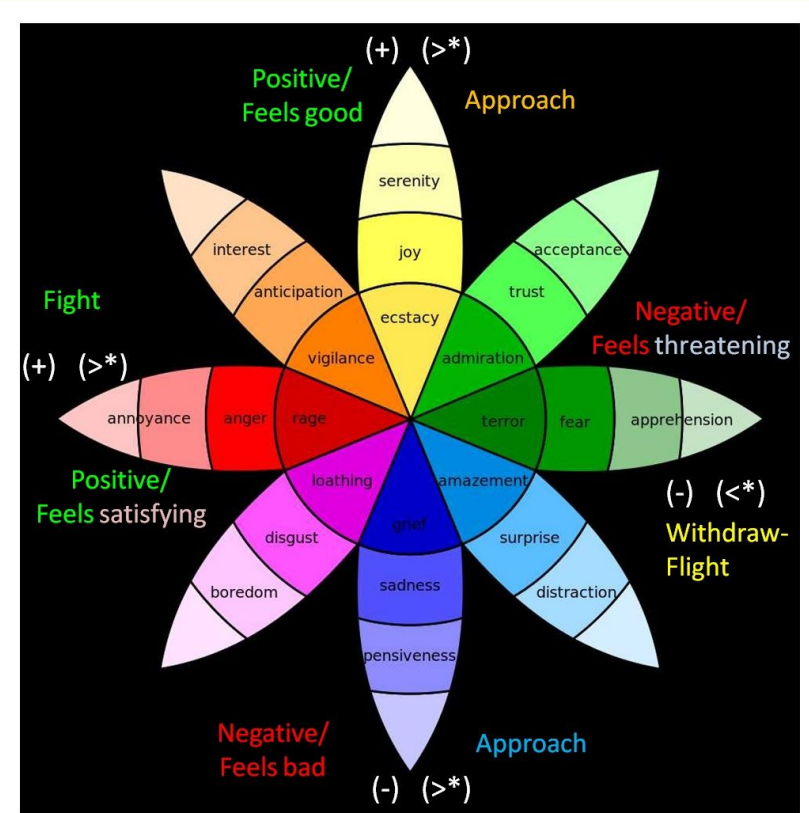
Parameter specification

How I got involved?

- Interested in machine learning
- Wanted to try my hand at research
- Enjoyed the psychological aspect

Next steps for the project

- Emotion as a vector
- Using word embeddings for optical character recognition
- Neutral filter



Thanks!

Sponsor Partners: John Dillon and Dan Hepp

Advisor: Sugih Jamin

Team: Pranay Shah, Vishnu Nair, Arun Annamalai, Eamy Mo, Rakshit Gogia,
Sebastian Jin