# Comparing NLP Algorithms

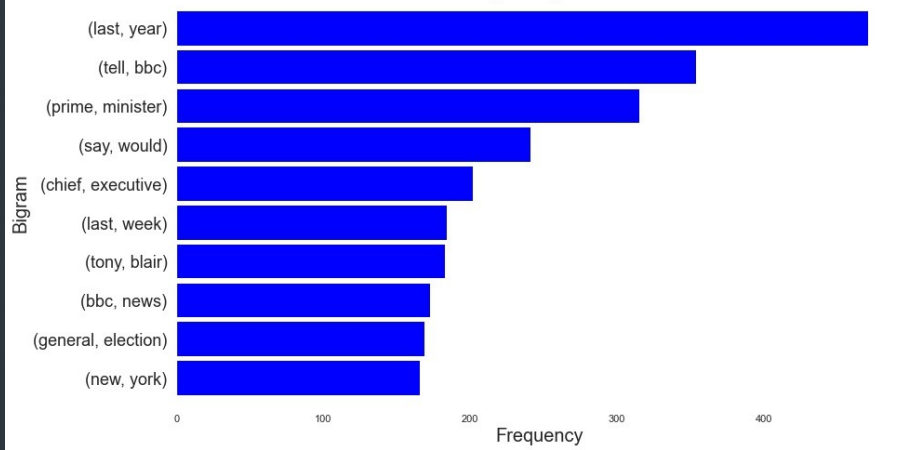Multinomial Naïve Bayes, SVM, LDA, NER

# Goals

- Train NLP models for use in production data pipeline for unstructured scientific and technical data

  - Automate pre-processing

  - Text Classification

  - Topic Modeling

  - N.E.R. for metadata tagging

- Determine feasibility of LDA for Topic Modeling in production pipeline

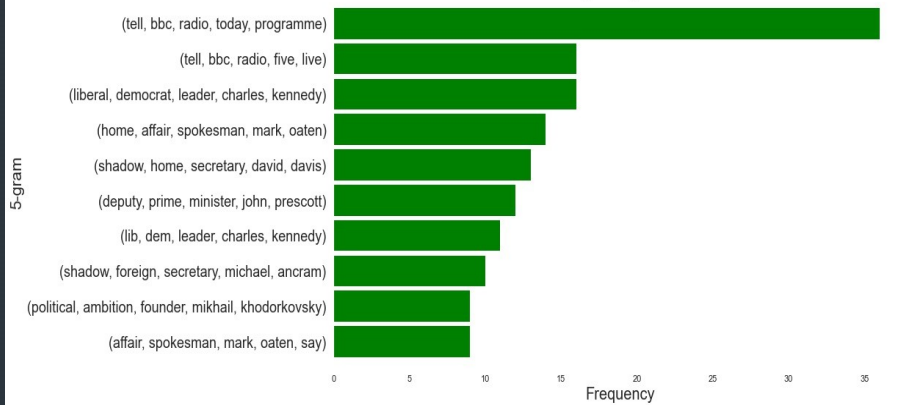- Use Named Entity Recognition to extract metadata

# The Data Set

- Consists of 2225 documents from the BBC news website corresponding to stories in five topical areas from 2004-2005.Natural Classes: 5 (business, entertainment, politics, sport, tech)

- Data set came from an academic research project that I can't even begin to comprehend

  - D. Greene and P. Cunningham. "Practical Solutions to the Problem of Diagonal Dominance in Kernel Document Clustering", Proc. ICML 2006.

- Series of text documents in folders grouped by topic

- Required very little wrangling (read in text, strip 1st line as title, split filepath to get topic, store in dataframe)

- Data Set stats (after trimming some outliers)

```
count    2200.000000
mean      198.202273
std        86.017350
min        46.000000
25%       132.000000
50%       178.000000
75%       251.000000
max       499.000000
```

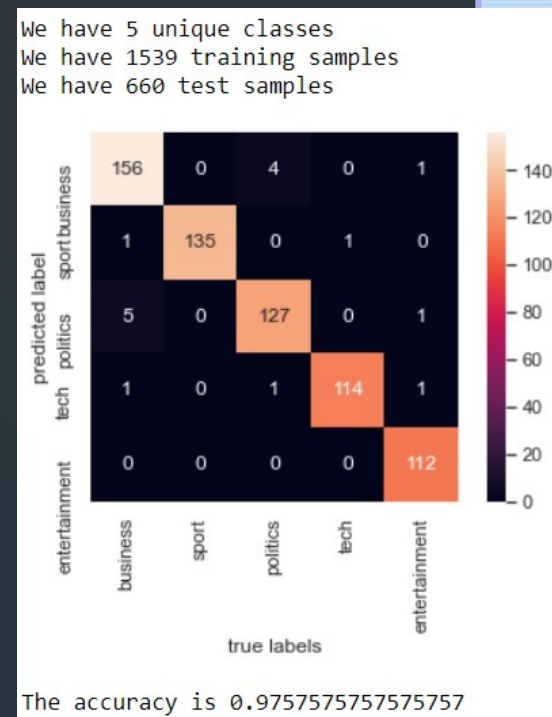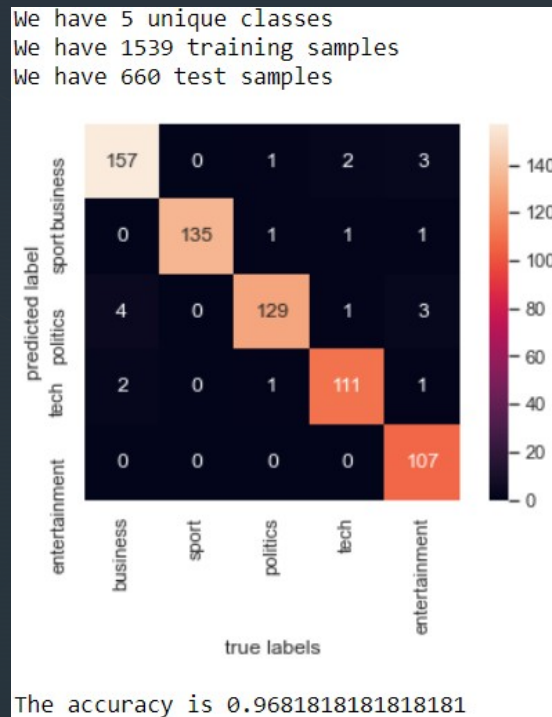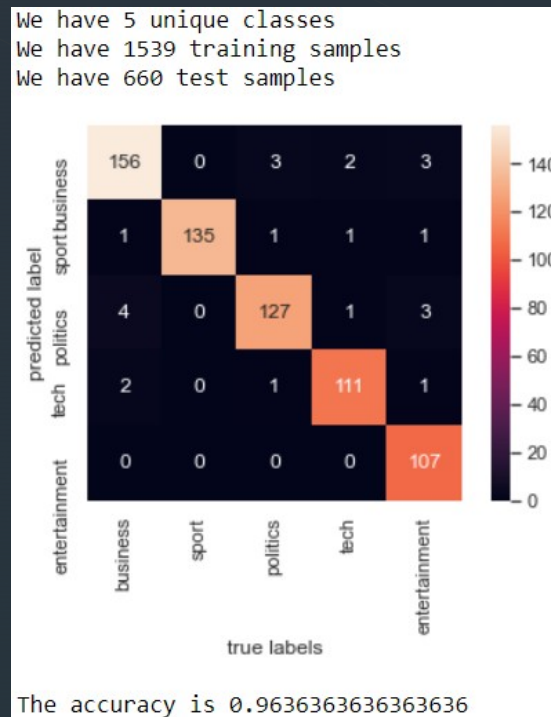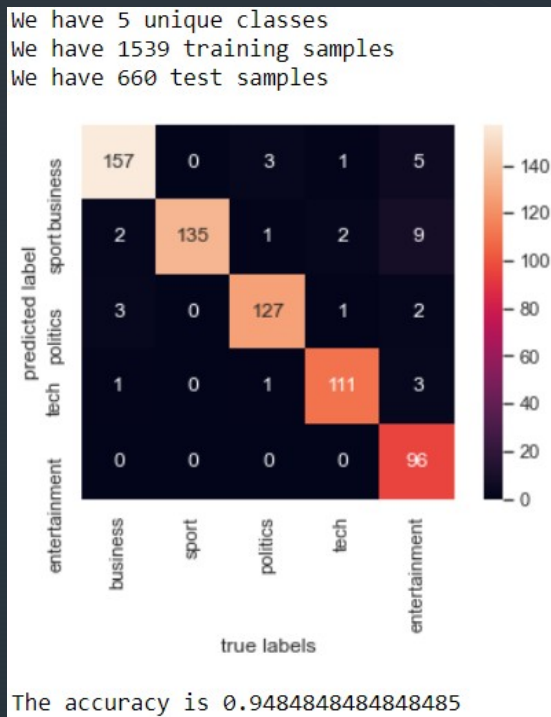| | |
|---|---|
| business | 510 |
| sport | 505 |
| politics | 413 |
| tech | 392 |
| entertainment | 380 |

# Text Classification – MNB & SVM

- Text classification works…too well?

- Pre-processing options:

  - Manual, TfidfVectorizer, PorterStemmer, Lemmatizer, others.

- Multinomial Naïve Bayes: 94.8% - 96.8%

- SVM: 97.5%

# Topic Modeling – LDA

- Can you be a library purist?
  - Leveraging what works in the moment
  - nltk, gensim, spaCy, sklearn, etc…
- What is LDA good for anyway?
  - People with a ton of time and patience…
  - Or for determining topics from bulk unstructured data
  - Or for clustering classified documents to gain new insights
  - Or … Or … Or …
- Can LDA make predictions?
  - Why, yes. Yes, it can.
    - Just not as easy as Text Classification…

# Named Entity Recognition

- Not much to say here.

- I used spaCy, but I'm open to Stanford's stuff too

- NER just works for metadata extraction

  - CAVEAT – manual process

- The pre-trained model is pretty fantastic

- I really need to learn how train a custom NER model

{'DATE': ['fourth quarter',
        '2005',
        'the preceding three quarters',
        '2003',
        'Quarterly',
        'the full-year',
        'Friday',
        'our full-year',
        '2000',
        'the fourth quarter',
        'year-earlier',
        'the three months to December'],
 'GPE': ['US'],
 'NORP': ['German'],
 'ORG': ['AOL Europe',
        'Rings',
        'Time Warner',
        'Bertelsmann',
        'the US Securities Exchange Commission',
        'TimeWarner',
        'AOL',
        'Google',
        'SEC',
        'Warner Bros'],
 'PERSON': ['Alexander', 'Richard Parsons',
'Catwoman']}

# Outcomes

- Text classification models work very well
  - MNB vs SVM vs ?? … How good is too good
  - Need more (diverse) data
- LDA looks like it could be useful with a large enough data set and good analytic insight…and patience
- NER is a beautiful thing for metadata extraction
  - Need to customize it for specific mission sets
- Pipeline and make_pipeline are … beautiful
- Holy crap I'm glad I learned about pickle
  - Pickling trained models!?!? Yes please and thank you

# Next Steps

- Read/write to MongoDB vs Pandas Dataframe

- Customized stop words list to remove things like 'say' and 'bbc'

- Incorporate foreign language handling

  - Machine Translation (AWS Translate, Google Cloud Translation, etc)

  - NLP in foreign languages

- Customize N.E.R. for tailored results

- Incorporate sentiment analysis for social media data

- Web application for data interaction

  - Django / Flask / ??

  - Visualize existing data

    - word cloud, histograms, ngrams, sentiment spikes over time, heatmaps, etc

  - Upload new document(s)

    - pre-process, predict, update metadata, correct errors, add to training set, retrain model