# Web APIs and NLP Classification

# Objective & Aim

- Objective

- A beverage company wants to know if people in an area posts more on social media about coffee or tea.

- Aim

- A classification model is needed to classify posts into either coffee or tea.

# Content

- Data Gathering
- Data Cleaning
- Naïve Bayes Model
- Alternative Models
- Conclusions and Recommendations

# Data Collection

- Using Pushshift's API, posts about coffee and tea were collected.

- Empty, removed or deleted posts were removed.
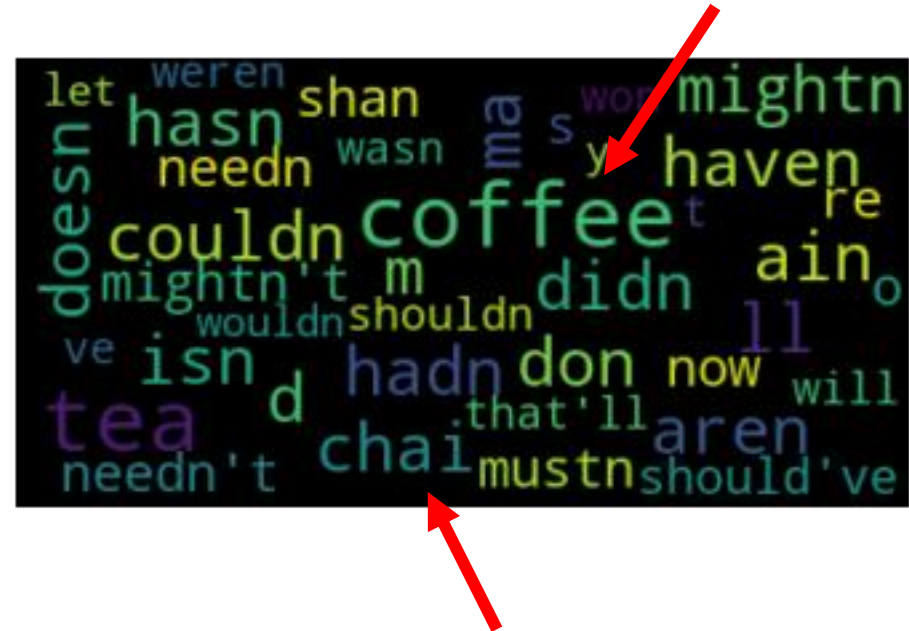
- Duplicates were removed.

# Data Cleaning

**Original text**



**Removal of stop words**

# Data Cleaning

**Original text**

Get a fine italian expresso, melt a little square of dark chocolate in it (at least 70% cocoa) and a dash of cinemon. Nothing more!\n\nAnd then, Like Homer said, "It\'s like a party in my mouth and everyone\'s invited!"Just for real cofee lovers!

**After data cleaning**

get fine italian melt little square dark chocolate least cocoa dash cinemon nothing like homer said like party mouth everyone invited real cofee lovers

# Lemmatization and Stemming

- **Cleaned sample**

- using starbucks beans grinding size tamping twice get oz secondshelp brevle dual temp machine

- **Lemmatization**

- using starbucks bean grinding size tamping twice get oz secondshelp brevle dual temp machine

- **Stemming**

- use starbuck bean grind size tamp twice get oz secondshelp brevl dual temp machin

# Naïve Bayes Model

- <span style="color:red">MultinomialNB</span>

- It is easy to implement.

- One of the most popular supervised learning classification.


- <span style="color:red">Metrics:</span>

- Accuracy score – emphasis on correct positive predictions.

- f1-score – emphasis on Type 1 and Type 2 errors.

# Naïve Bayes Model

- **No differences** in the scores when using cleaned, lemmatized and stemmed datasets.

| Dataset | Train score | Accuracy | f1-score |
|---------|-------------|----------|----------|
| Cleaned | 0.893 | 0.871 | 0.879 |
| Lemmatized | 0.895 | 0.872 | 0.879 |
| Stemmed | 0.897 | 0.865 | 0.870 |

# Naïve Bayes Model – Feature Importance

- <span style="color:red">'bean'</span>
- 'grind', 'burr', 'roaster',

- <span style="color:red">'machines'</span>
- 'breville', 'baratza', 'jx', 'zpresso'

- <span style="color:red">Method</span>
- 'moka', 'drip', 'barista'

# Naïve Bayes Model – Feature Importance

- **Types**
- 'chinese', 'sencha', 'yunnan', 'oolong', 'chamomile', 'pu', 'erh', 'gong', 'fu'

- **Device**
- 'gaiwan'

- **Other**
- 'leaf', 'sourcing',
- 'worry', 'stories', 'topics', 'august', 'gal'

# Spam !!

- Not dropped during cleaning.
- They were properly classified as 'tea'.



drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September

drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September

drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September

drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September
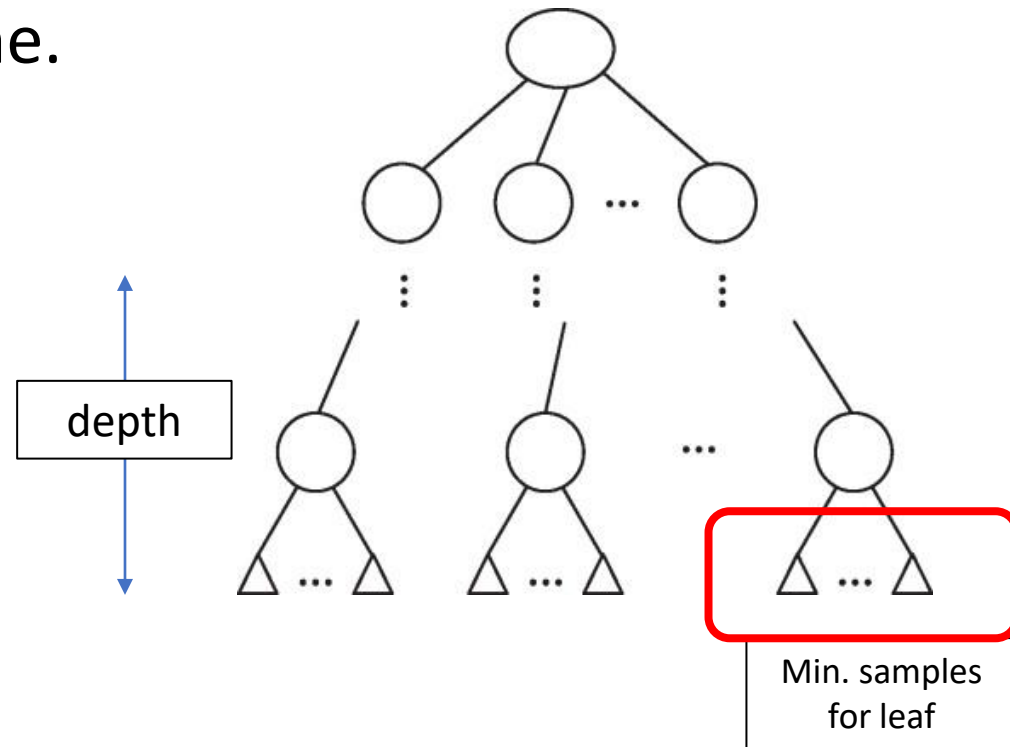
drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September

drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September

drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk non related topics maybe want advice guy gal talk life general cup daily discussion questions stories September

drinking today questions mind stories share worry one make fun drink questions ask also talk anything else mind specific routine making oolong kick lately feel free link pictures well even talk

# Alternative Models – Decision Tree Classifier

- One decision tree, one outcome.
- Not recommended.
- Tendency to overfit.
- Hyperparameter tuning:
  - Tree depth
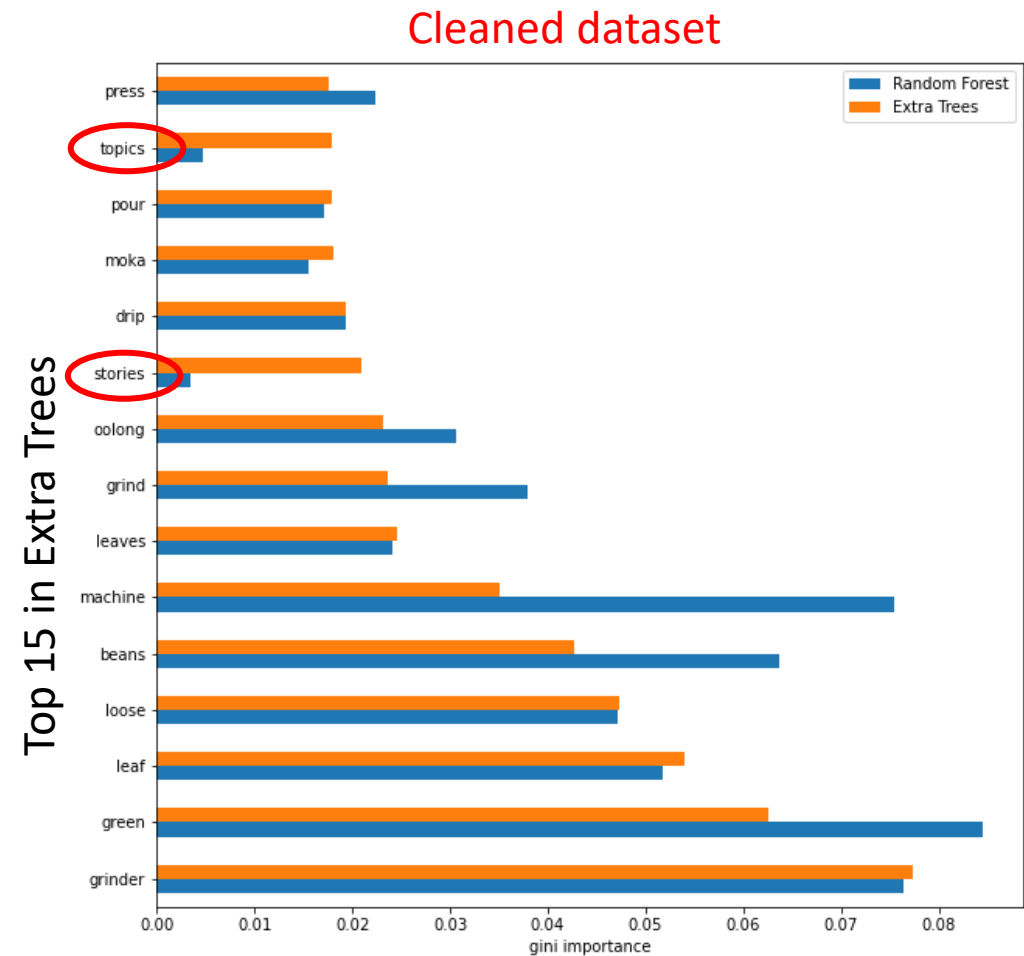  - Min sample for split
  - Min samples for leaf

# Random Forest vs Extra Trees

- Hyperparameter tuning:
  - Number of trees
  - Tree depth
  - Min sample for split
  - Min samples for leaf

- Which is better?

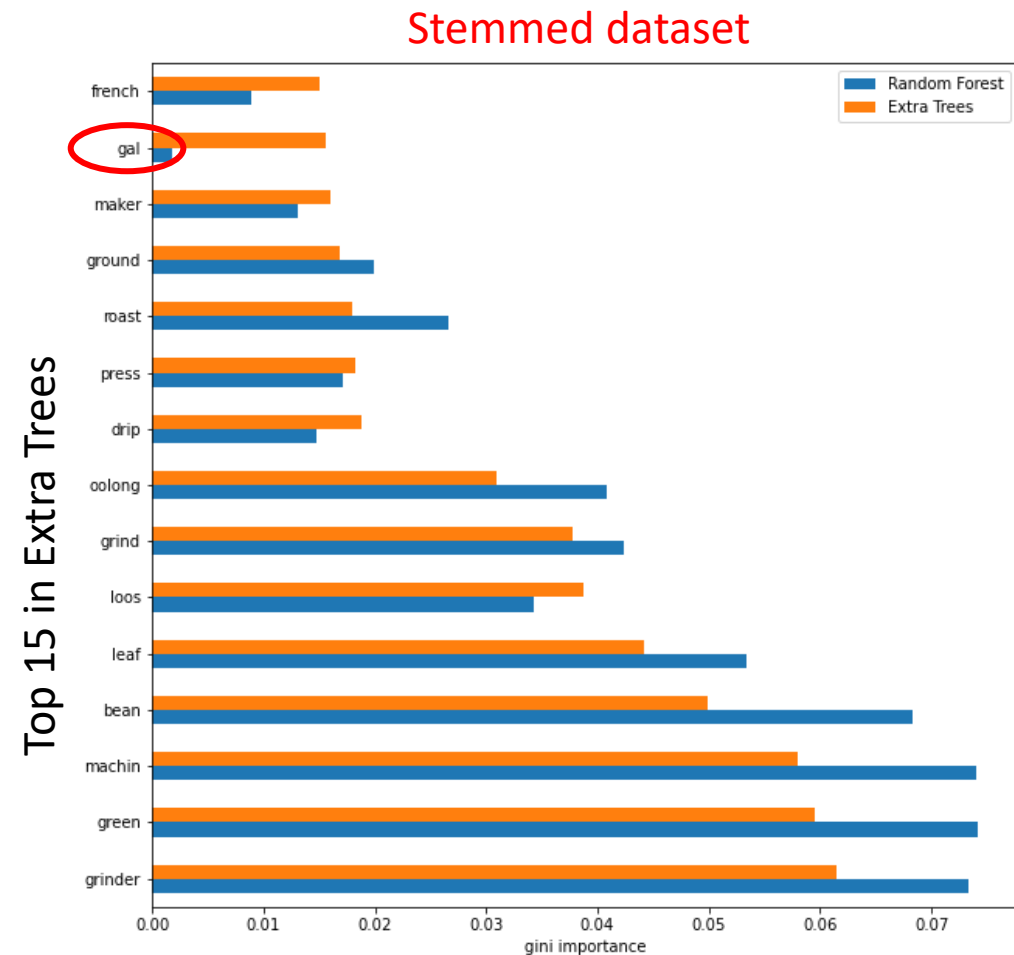| Dataset | Accuracy scores | |
|---|---|---|
| | Random Forest | Extra Trees |
| Cleaned | 0.832 | 0.812 |
| Lemmatized | 0.850 | 0.833 |
| Stemmed | 0.862 | 0.850 |

# Random Forest vs Extra Trees

- Words in spam posts has higher importance in Extra Trees.

- Not found in Random Forest.



Cleaned dataset

# Random Forest vs Extra Trees

- Stemmed features?

- 'gal' in spam posts has higher importance in Extra Trees.

- Not found in Random Forest.



Stemmed dataset

# Random Forest vs Extra Trees

- <span style="color:red">Metrics:</span>
- Random Forest has higher accuracy and f1 scores than Extra Trees.

- Both models performed better for stemmed dataset.
- If there were many irrelevant features, stemming might be better at consolidating features.

# Conclusions

- Naïve Bayes model the model of choice.
  - It provided good accuracy and f1-scores, compared to other models.
  - It reported good scores for all datasets.
  - The model performed well *despite* the presence of spam.


- Recommendations:
  - Better cleaning
  - Validity longer term