

## # Spam Email Detection Project

numpy: good for numpy arrays

pandas: good for data-frames

re: regular expression for texts

Stopwords: words that aren't too valuable ('a', 'the')

PorterStemmer: Stem our words, gives root word, removes prefix and suffix

TfidfVectorizer: Converts text to vectors

logistic Regression model: Binary Classification

### Data Cleansing:

Before cleaning, features like punctuation or capitalization

that might be helpful  
are retained.

```
SPECIALSYMBOLS = set("@%&*#^+=")
def styling(content):
    #for letters only
    letters = [c for c in content if c.isalpha()]
    upper_ratio = sum(1 for c in letters if c.isupper())/len(letters) if letters else 0
    #counts:
    exclamation_count = content.count("!")
    question_count = content.count("?")
    #special signs
    special_sign = sum(c in SPECIALSYMBOLS for c in content)
    #URL count
    url_count = len(re.findall(r"http://", content))

    return pd.Series({
        "upper_ratio": upper_ratio,
        "exclamation count ": exclamation_count,
        "question count": question_count,
        "special count": special_sign,
        "url count": url_count
    })
```

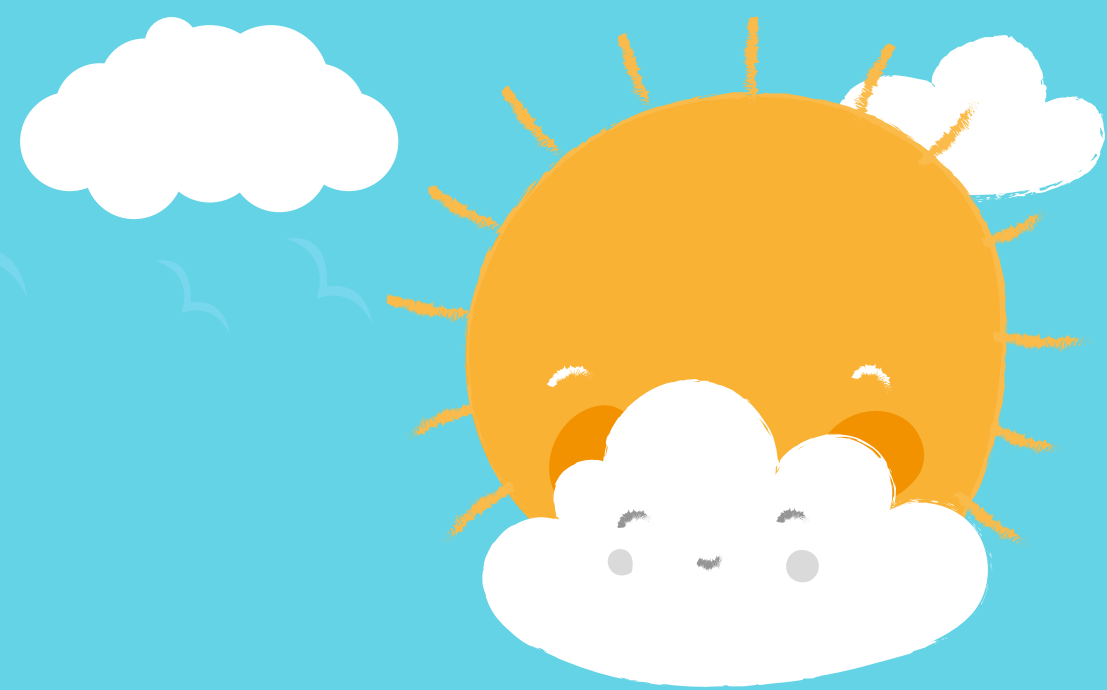
```
pstem = PorterStemmer()
def cleanse(content):
    new_content = re.sub(r"[^a-z]", '', content)
    new_content = new_content.lower()
    new_content = new_content.split()
    new_content = [word for word in new_content if word]
    new_content = [pstem.stem(word) for word in new_content if not word in stopwords.words('english')]
    new_content = ' '.join(new_content)
    return new_content
```

→ Only keeps letters

→ splits to linked list

→ only keeps root word

→ does not take insignificant words



TF ID Vectorizer: Term frequency - Inverse Document frequency  
↳ scores words based on how important they are

$$TF = \frac{\text{no. of occurrence of } t}{\text{no. of terms in doc}}$$

$$IDF = \log\left(\frac{\text{Total no. of docs}}{\text{no. of docs with } t}\right)$$

$$TFID(t) = TF(t) \times IDF(t)$$

↳ more if occurs more in one doc

↳ less if found in more docs

→ max\_range keeps vocab to a fixed set to dull noise

→ ngram\_range determines how many "bunches" to consider

```
# print(email_data.columns)
vectorizer = TfidfVectorizer(max_features = 3500, ngram_range=(1,2))
X_vec = vectorizer.fit_transform(email_data['cleaned']) #a sparse matrix
```

top 3500 words  
accepts "totaly free")

✓ 0.2s

Python

get the numeric data :

```
x_num = email_data[["upper ratio", "exclamation count", "question count", "special c"]
x_all = sp.hstack([X_vec, x_num])
```

→ stacks horizontally first we have the feature matrix and then numeric cols

we had  $n$  numeric columns, and suppose we get a vectorized matrix (no. of rows, 3500) and  $x\_num$  (no. of rows, 5)

So the numbers get stacked to form one frame.

Splitting the Dataset to training and testing models

10 % for testing

```
X_train, X_test, Y_train, Y_test = train_test_split(x_all,y,test_size=0.1, stratify=y, random_state=2)
```

properly distributed with regards to y

Python

