# Is the E-Mail a SPAM or a HAM?

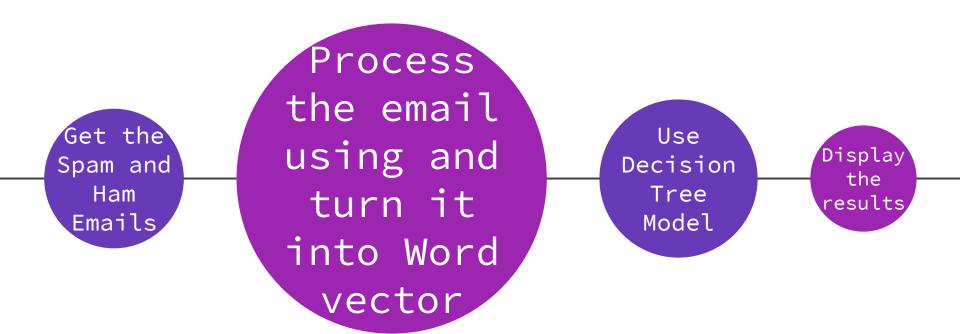
A simple Decision Tree Based classifier to identify if a given mail is spam or ham

-BY Vikraant Pai Roll No. D005 SAP ID: 70271019005

#### **About The Problem**

- The problem that is being solved here is finding the right decision tree model that can exam whether the mail is spam or a ham
- 2. The method that is used to solve this problem is to first extract information from the mail to find out the words in the mail and then fit a decision tree model
- 3. Using this model it is tested on a test sample for precision and recall and compared with logistic regression

#### Steps Taken to build the model

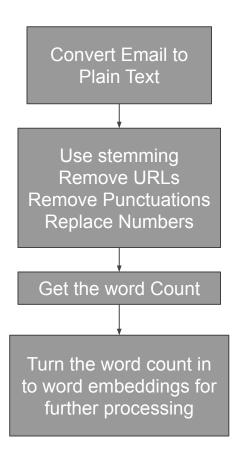


#### Extracting the features from the mail

Email is converted to plain text and then word vector using Natural language Toolkit for stemming and based on the size of the vocabulary for the corpus, it is converted to

word embeddings for further

processing



# Using the custom transformer to transform the words into Word embeddings

## Using the custom transformer to transform the words into Word embeddings

```
[112] from sklearn.tree import DecisionTreeClassifier
   dtc = DecisionTreeClassifier( max depth=5, max leaf nodes=5, criterion="gini")
   dtc.fit(X train transformed, y train)
   score = cross val score(dtc, X train transformed, y train, cv=3, verbose=3)
   [CV] .....
   [CV] ....., score=0.931, total= 0.0s
    [CV] .....
   [CV] ....., score=0.940, total= 0.0s
   [CV] .....
   [CV] ....., score=0.929, total= 0.0s
   [Parallel(n jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
   [Parallel(n jobs=1)]: Done 1 out of 1 | elapsed: 0.0s remaining:
                                                        0.05
   [Parallel(n jobs=1)]: Done 2 out of 2 | elapsed:
                                          0.0s remaining:
                                                        0.05
                                          0.1s finished
   [Parallel(n jobs=1)]: Done 3 out of 3 | elapsed:
   print(f"Mean cross validation score {score.mean()}")
```

Mean cross validation score 0.9333333333333333

#### Hyperparameters for the model

```
dtc.get_params()
{'ccp alpha': 0.0,
  'class weight': None,
  'criterion': 'gini',
  'max depth': 5,
  'max features': None,
 'max leaf nodes': 5,
  'min impurity decrease': 0.0,
 'min impurity_split': None,
  'min samples leaf': 1,
 'min samples_split': 2,
  'min weight fraction leaf': 0.0,
  'presort': 'deprecated',
  'random state': None,
  'splitter': 'best'}
```

#### Some emails for testing

print(email to text(X\_test[50])[:1000]) print(email\_to\_text(X\_test[1])[:1000]) Hi, Never mind, there was some cron thing doing rpm -qf ??? seems fine now. print(email to text(X test[20] -- On Wednesday, February 06, 2002 07:37:44 +1300 Mark Derricutt Red Hat 8.0 is released tomor Whore eructed: RPMs of GStreamer for it. <mark@talios.com> wrote: >--lIt's an amusing anecdote. >-- certainly nothing here sup All of them (core, plugins ar >--1"Status: False". > Fetched 88.1kB in 2m31s (581B/s) repository in a new "redhat-{ > error: cannot get exclusive lock on /var/lib/rpm/Packages >So thats the trick, just let > error: cannot open Packages index using db3 - Operation not permitted (1) >true [...] The repository for dependenci > E: could not open RPM database:cannot open Packages index using db3 the gstreamer rpms is "redhat > Operation not permitted (1) Exsqueeze me, but what part of contains all the necessary pa did you fail to grok? I perso heard of Bush and Chirac going A screenshot of Red Hat 8.0 r http://thomas.apestaart.org/c -- \m/ --Next time I hear a joke, I pro "...if I seem super human I have been misunderstood." (c) Dream Theater out primary sources for confir we have you around to keep us mark@talios.com - ICQ: 1934853 JID: talios@myjabber.net Here are some known issues w inundating us with such erudit "fight the powers that be, fro a) gstreamer-nautilus isn't ł other dippy bromides. have a -devel package for it RPM-List mailing list < RPM-List@freshrpms.net> http://lists.freshrpms.net/mailman/listinfo/rpm-list

http://xent.com/mailman/listin

## Decoding the word embeddings for better understanding

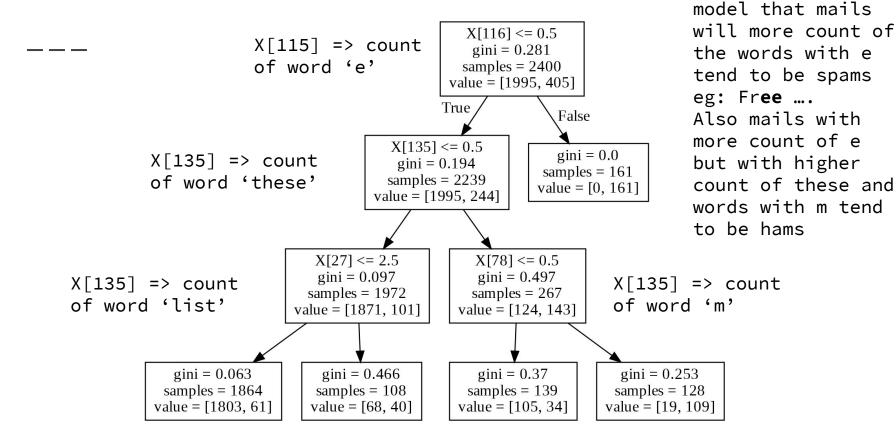
We need to find the values for X[116], X[135], X[27], X[78]

We do it using the word to vector transform in the pipeline

```
[153] preprocess_pipeline["wordcount_to_vector"].most_common_[116][0]
  [166] preprocess pipeline["wordcount to vector"].most common [135][0]
       'these'
  [167] preprocess pipeline["wordcount to vector"].most common [27][0]
   「→ 'list'
         preprocess pipeline["wordcount to vector"].most common [78][0]
```

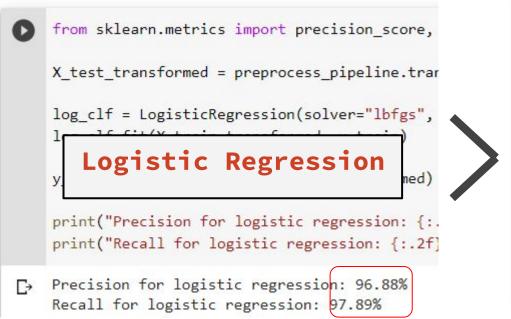
#### Decision Tree Generated By the Model

It seems from the



### Comparison with Logistic Regression on Test Set

\_\_\_\_ Logistic Regression Performs Better on the test set. However we also see that even though Decision Tree performs poorly in comparison to Log Reg, it is more explainable



from sklearn.metrics import precision score, rec X test transformed = preprocess pipeline.transfo **Decision Tree** print("Precision for Decision Tree Classifier: { print("Recall for Decision Tree Classifier : {:. Precision for Decision Tree Classifier: 88.16% Recall for Decision Tree Classifier: 70.53%

#### **Example Predictions by the model**

Actual: Ham
Predicted: Ham

y\_pred[1], y\_test[0]

[→ (0, 0)

print(email\_to\_text(X\_test[1])[:1

, Нi,

Red Hat 8.0 is released tomorrow (monday). I took some time out to make RPMs of GStreamer for it.

All of them (core, plugins and player) have been uploaded to the apt repository in a new "redhat-80-i386" directory.

The repository for dependencies is again called "deps", and the one for the gstreamer rpms is "redhat", because this time around the base distro contains all the necessary packages.

A screenshot of Red Hat 8.0 running gst-player is up at <a href="http://thomas.apestaart.org/download/screenshots/redhat-80-gst-player.png">http://thomas.apestaart.org/download/screenshots/redhat-80-gst-player.png</a>

Here are some known issues with the resulting rpms :

a) gstreamer-nautilus isn't built, the package got renamed and I don't have a -devel package for it yet

Actual: Spam
Predicted: Spam

[199] y\_pred[4], y\_test[4]

[ (1, 1)

print(email\_to\_text(X\_test[4])[:1000])

A POWERHOUSE GIFTING PROGRAM You Don't Want To Miss!

GET IN WITH THE FOUNDERS! The MAJOR PLAYERS are on This ONE For ONCE be where the PlayerS are This is YOUR Private Invitation

EXPERTS ARE CALLING THIS THE FASTEST WAY
TO HUGE CASH FLOW EVER CONCEIVED
Leverage \$1,000 into \$50,000 Over and Over Again

THE QUESTION HERE IS:
YOU EITHER WANT TO BE WEALTHY
OR YOU DON'T!!!
WHICH ONE ARE YOU?
I am tossing you a financial lifeline and for your sake I
Hope you GRAB onto it and hold on tight For the Ride of youR life!

Testimonials

Hear what average people are doing their first few days:

- ♦We've received 8,000 in 1 day and we are doing that over and over again!'
- ♦I'm a single mother in FL and I've received 12,000 in the last 4 days.♦
- ♦I was not sure about this when I sent off my \$1,000 pledge, but I got bac
- lacktriangleI didn't have the money, so I found myself a partner to work this with. W

#### **Confusion Matrix for Decision Tree Model**

