

# Week 10 Quiz

Qi Meng - qm2162

Due Sun. Nov 28th 11:59pm

In this quiz, we're going to load documents from 2 categories (space, cars) in the 20newsgroups dataset.

The goal is to train a classifier that classifies documents into these 2 categories based on a term frequency representation of the documents.

We will then calculate mean cross-validation accuracy of a RandomForestClassifier using this transformation.

## Setup Environment

In [8]:

```
import numpy as np
import pandas as pd
```

## Load the Dataset

In [9]:

```
# Import fetch_20newsgroups from sklearn.datasets
from sklearn.datasets import fetch_20newsgroups

# Load the dataset using fetch_20newsgroups().
# Only fetch the two categories of interest using categories=['sci.space', 'rec.autos']
# Store in the result into newsgroups
newsgroups = fetch_20newsgroups(categories=['sci.space', 'rec.autos'])

# Store the newsgroups.data as docs, newsgroups.target as y and newsgroups.target_names as y_names
docs, y, y_names = newsgroups.data, newsgroups.target, newsgroups.target_names

# Print the number of observations by printing the length of docs
print(len(docs))
```

1187

In [10]:

```
# Print the text of the first document in docs.
# Note: try printing both with and without the print() statement
# with the print statement, linebreaks are parsed,
# without, linebreaks are printed as escape characters
print(docs[0])
print('-----')
docs[0]
```

From: prb@access.digex.com (Pat)  
 Subject: Re: Proton/Centaur?  
 Organization: Express Access Online Communications USA  
 Lines: 15  
 NNTP-Posting-Host: access.digex.net

Well thank you dennis for your as usual highly detailed and informative posting.

The question i have about the proton, is could it be handled at one of KSC's spare pads, without major malfunction, or could it be handled at kourou or Vandenberg?

Now if it uses storables, then how long would it take for the russians to equip something at cape york?

If Proton were launched from a western site, how would it compare to the T4/centaur? As i see it, it should lift very close to the T4.

pat

```
-----
Out[10]: "From: prb@access.digex.com (Pat)\nSubject: Re: Proton/Centaur?\nOrganization: Express Access Online Commu
nications USA\nLines: 15\nNNTP-Posting-Host: access.digex.net\n\n\nWell thank you dennis for your as usual
highly detailed and informative \nposting. \n\nThe question i have about the proton, is could it be ha
ndled at\none of KSC's spare pads, without major malfunction, or could it be\nhandled at kourou or Vand
enberg? \n\nNow if it uses storables, then how long would it take for the russians\nto equip something
at cape york?\n\nIf Proton were launched from a western site, how would it compare to the\nT4/centaur?
As i see it, it should lift very close to the T4.\n\npat\n"
```

```
In [11]: # Print the target value of the first document in y.
y[0]
```

```
Out[11]: 1
```

```
In [12]: # Print the target_name of the first document using y_names and y.  
y_names[y[0]]
```

```
Out[12]: 'sci.space'
```

## Use CountVectorizer to Convert To TF

```
In [13]: # Import CountVectorizer from sklearn  
from sklearn.feature_extraction.text import CountVectorizer  
  
# Initialize a CountVectorizer object. It should  
# lowercase all text,  
# include both unigrams and bigrams  
# exclude terms that occur in fewer than 10 documents  
# exclude terms that occur in more than 95% of documents  
# Store as cvect  
cvect = CountVectorizer(lowercase=True,  
                        ngram_range=(1,2),  
                        min_df=10/len(docs),  
                        max_df=0.95)  
  
# Fit cvect on docs and transform docs into their term frequency representation.  
# Store as X_tf  
X_tf = cvect.fit_transform(docs)  
  
# Print the shape of X_tf.  
# The number of rows should match the number of documents  
# and the number of columns should be in the thousands  
X_tf.shape
```

```
Out[13]: (1187, 5893)
```

```
In [15]: # Print out the last 5 terms in the learned vocabulary in cvect  
# using .get_feature_names() which returns a list of terms corresponding  
# to the order of the columns in X_tf  
# They should all be related to zoos  
cvect.get_feature_names()[-5:]
```

```
Out[15]: ['zoo', 'zoo toronto', 'zoology', 'zoology kipling', 'zoology lines']
```

```
In [16]: # The stopwords learned by cvect are stored as a set in cvect.stop_words_  
# We'd like to print out a small subset of these terms.  
# One way to get a subset of a set is to treat it as a list.  
# First, convert the stop_words_ set to a list.  
# Store as stop_words_list  
stop_words_list = list(cvect.stop_words_)  
  
# Print out the first 5 elements in stop_words_list.  
# Note that, since a set is unordered,  
# there is no meaning to the ordering of these terms and they may vary over runs.  
stop_words_list[:5]
```

```
Out[16]: ['annex3',  
          'for diesel',  
          'demonstrated that',  
          'astron astrophys',  
          'following texts']
```

## Calculate Mean CV Accuracy Using RandomForestClassifier

In [18]:

```
# Import cross_val_score from sklearn
from sklearn.model_selection import cross_val_score

# Import RandomForestClassifier from sklearn
from sklearn.ensemble import RandomForestClassifier

# Get a set of 5-fold CV scores using
#   a RandomForestClassifier with 50 trees and n_jobs=-1
#   and the full dataset X_tf and y
# Store as cv_scores
cv_scores = cross_val_score(RandomForestClassifier(n_estimators=50, n_jobs=-1),
                             X_tf,
                             Y,
                             cv=5)

# Print the mean of these cv_scores. The mean accuracy should be above .9
cv_scores.mean()
```

Out[18]: 0.9679856752827714

## Optional: Find Important Features

In [ ]:

```
# CountVectorizer stores the feature names (terms in the vocabulary) in two ways:
# 1. as a dictionary of term:column_index pairs, accessed via cvect.vocabulary_
# 2. as a list of terms, in column index order, accessed via cvect.get_feature_names()
#
# We can get the indices of the most important features by retraining a RandomForestClassifier on X_tf,y
# and accessing .feature_importances_
#
# Using some combination of the above data-structures,
# print out the top 10 terms in the vocabulary
# ranked by the feature importances learned by a RandomForestClassifier with 50 trees
#
# The terms you find will likely not be surprising given the document categories.
```

---

