

Homework 2 Report

Part 1: Viterbi Algorithm

Below is the fully worked problem on paper:

Part 1:

Given: - Transition Probability Matrix:

	s_1	s_2
s_1	0.6	0.3
s_2	0.4	0.7

- State Probability Matrix:

s	s_1	s_2
s	0.5	0.5

- Emissions Matrix:

	s_1	s_2
A	0.3	0.1
C	0.2	0.4
G	0.3	0.1
T	0.2	0.4

Viterbi Matrix:

	$T=0$	$T=1$	$T=2$	$T=3$	$T=4$
s_1	0.1	0.018	0.0022	0.00026	0.000098
s_2	0.2	0.014	0.0039	0.00109	0.000076

Backtrace Matrix:

	*	s_1	s_1	s_1	s_2
s_1	*	s_1	s_1	s_1	s_2
s_2	s_2	s_2	s_2	s_2	s_2

Sequence: * s_2 s_2 s_2 s_1

Calculations:

$$V_0(s_1) = P(c | s_1) \cdot P(s_1) = 0.2 \cdot 0.5 = 0.1$$

$$V_0(s_2) = P(c | s_2) \cdot P(s_2) = 0.4 \cdot 0.5 = 0.2$$

$$V_1(s_1) = \max \begin{cases} V_0(s_1) \cdot P(s_1 | s_1) \cdot P(G | s_1) = 0.1 \cdot 0.6 \cdot 0.3 = 0.018 \\ V_0(s_2) \cdot P(s_1 | s_2) \cdot P(G | s_2) = 0.2 \cdot 0.3 \cdot 0.3 = 0.018 \end{cases}$$

$$V_1(s_2) = \max \begin{cases} V_0(s_1) \cdot P(s_2 | s_1) \cdot P(G | s_2) = 0.1 \cdot 0.4 \cdot 0.1 = 0.004 \\ V_0(s_2) \cdot P(s_2 | s_2) \cdot P(G | s_2) = 0.2 \cdot 0.7 \cdot 0.1 = 0.014 \end{cases}$$

$$V_2(s_1) = \max \begin{cases} V_1(s_1) \cdot P(s_1 | s_1) \cdot P(T | s_1) = 0.018 \cdot 0.6 \cdot 0.2 = 0.0022 \\ V_1(s_2) \cdot P(s_1 | s_2) \cdot P(T | s_2) = 0.014 \cdot 0.3 \cdot 0.2 = 0.0008 \end{cases}$$

$$V_2(s_2) = \max \begin{cases} V_1(s_1) \cdot P(s_2 | s_1) \cdot P(T | s_2) = 0.018 \cdot 0.4 \cdot 0.4 = 0.0029 \\ V_1(s_2) \cdot P(s_2 | s_2) \cdot P(T | s_2) = 0.014 \cdot 0.7 \cdot 0.4 = 0.0039 \end{cases}$$

$$V_3(s_1) = \max \begin{cases} V_2(s_1) \cdot P(s_1 | s_1) \cdot P(c | s_1) = 0.0022 \cdot 0.6 \cdot 0.2 = 0.00026 \\ V_2(s_2) \cdot P(s_1 | s_2) \cdot P(c | s_1) = 0.0039 \cdot 0.3 \cdot 0.2 = 0.00023 \end{cases}$$

$$V_3(s_2) = \max \begin{cases} V_2(s_1) \cdot P(s_2 | s_1) \cdot P(c | s_2) = 0.0022 \cdot 0.4 \cdot 0.4 = 0.00035 \\ V_2(s_2) \cdot P(s_2 | s_2) \cdot P(c | s_2) = 0.0039 \cdot 0.7 \cdot 0.4 = 0.00109 \end{cases}$$

$$V_4(s_1) = \max \begin{cases} V_3(s_1) \cdot P(s_1 | s_1) \cdot P(A | s_1) = 0.00026 \cdot 0.6 \cdot 0.3 = 0.000047 \\ V_3(s_2) \cdot P(s_1 | s_2) \cdot P(A | s_1) = 0.00109 \cdot 0.3 \cdot 0.3 = 0.000098 \end{cases}$$

$$V_4(s_2) = \max \begin{cases} V_3(s_1) \cdot P(s_2 | s_1) \cdot P(A | s_2) = 0.00026 \cdot 0.4 \cdot 0.1 = 0.0000104 \\ V_3(s_2) \cdot P(s_2 | s_2) \cdot P(A | s_2) = 0.00109 \cdot 0.7 \cdot 0.1 = 0.000076 \end{cases}$$

Below is the code and program output for the viterbi algorithm:

```
def Viterbi(model, observations):

    # Length of observations
    T = len(observations)

    # Number of states
    N = len(model.states)

    viterbiMatrix = np.zeros(shape=(N, T))

    translation = Translate(observations)
    backpointer = np.zeros(shape=(N, T))

    # Initialization step
    for s in range(N):
        # Initialize the starting hidden states with appropriate probabilities
        # for future states
        viterbiMatrix[s][0] = model.startProbabilities[s] *
model.emissions[translation[0]][s]
        backpointer[s][0] = -1

    # For the remaining number of observations
    for t in range(1, T):
        # For each of the hidden states, N
        for s in range(N):
            transition_probs = list()
            for sprime in range(N):
                transition_probs.append(viterbiMatrix[sprime][t-1] *
model.states[sprime][s] * model.emissions[translation[t]][s])

            viterbiMatrix[s][t] = max(transition_probs)

            backpointer[s][t] = np.argmax(viterbiMatrix[:, t-1] *
model.states[:, s] * model.emissions[translation[t], s])

    # Get the overall probability of the best path through the HMM
    bestPathProbability = 0
    bestPathPointer = list()
    TViterbi = viterbiMatrix.transpose()

    bestPathProbability = max(TViterbi[T-1][:])

    bestPathPointer = np.argmax(TViterbi[T-1][:])
```

```
bestPath = FindBestPath(viterbiMatrix, T, backpointer)

print("\n")
print("Viterbi Matrix: ")
print(viterbiMatrix)
print("\n")
print("Backtrace Matrix: ")
print(backpointer)
print("\n")
print("Best Path(from t = 0 to t = T): ", end="")

print(bestPath)
```

```
def FindBestPath(viterbiMatrix, T, backpointer):
    bestPath = np.zeros(T+1)
    TViterbi = viterbiMatrix.transpose()
    bestPathProbability = 0

    i = T-1
    bestPath[-1] = np.argmax(viterbiMatrix[:, T-1])
    for j in range(i, -1, -1):

        bestPath[j] = backpointer[int(bestPath[j+1]), j]
        bestPathProbability = max(viterbiMatrix[0, T-1], viterbiMatrix[1, T-1])
    print("Best path probability (Sumtotal for each node in the path): " +
str(bestPathProbability))
    return bestPath

def Translate(observations):

    translation = list()

    for i in observations:
        if i == 'A':
            translation.append(0)
        elif i == 'C':
            translation.append(1)
        elif i == 'G':
            translation.append(2)
        elif i == 'T':
            translation.append(3)
        else:
```

```
print("Unrecognized character. Exiting program...")  
exit()  
  
return translation
```

```
Best path probability (Sumtotal for each node in the path): 9.878399999999997e-05  
  
Viterbi Matrix:  
[[1.0000e-01 1.8000e-02 2.1600e-03 2.5920e-04 9.8784e-05]  
 [2.0000e-01 1.4000e-02 3.9200e-03 1.0976e-03 7.6832e-05]]  
  
Backtrace Matrix:  
[[-1.  0.  0.  0.  1.]  
 [-1.  1.  1.  1.  1.]]  
  
Best Path(from t = 0 to t = T): [-1.  1.  1.  1.  1.  0.]  
  
C:\Users\19255\Documents\GitHub\Natural-Language-Processing>
```

Part 2: Multinomial Naive Bayes w/t SpaCy & Sklearn

This part included two implementations. Model 1 uses pure sklearn libraries which were permitted for use by Dr. Liu on 2/22/2022. Model 2 incorporates a SpaCy pipeline for text preprocessing, then passes the processed data to Sklearn's multinomial naive bayes model.

The implementation was done in a colaboratory environment and is recommended that the included source code be run the same way.

Model 1: Pure Sklearn:

The first model used a simple pipeline consisting of three components:

1. CountVectorizer: This created a vectorized bag of words from all documents in the training section.
2. TfidfTransformer: This takes the vectorized bag of words and normalizes it. This is to reduce the negative impact of frequently occurring features within the bag. Features that appear less frequently can thus have a similar impact and the model will not be affected by biases.

3. MultinomialNB: The sklearn multinomial naive bayes model is then fed the normalized bag of words, fitted to the training set, and tested.

```
"""HW2_Part2.ipynb
```

```
Automatically generated by Colaboratory.
```

```
Original file is located at
```

```
    https://colab.research.google.com/drive/1uiDIF6UQZpaIQzQHR3dCBIvqrGdfV3Md  
"""
```

```
# NOTE: As of 2/22/2022, Dr. Liu allowed the use of sklearn preprocessing  
methods along with  
# spaCy to develop the classifier. She also allowed the use of sklearn's  
premade 20newsgroups dataset.
```

```
import numpy as np  
from sklearn import datasets  
from sklearn.naive_bayes import MultinomialNB  
from sklearn import pipeline  
from sklearn.feature_extraction.text import CountVectorizer  
from sklearn.feature_extraction.text import TfidfTransformer  
from sklearn import metrics  
import re
```

```
"""# Model 1: Pure Sklearn Setup
```

```
Sklearn and re library documentation was used in the creation of this program.  
"""
```

```
# Extract only the train and test datasets for our categories and remove  
unnecessary components  
categories = ['rec.autos', 'comp.graphics']  
train = datasets.fetch_20newsgroups(subset='train', categories=categories,  
remove=('headers', 'footers', 'quotes'), shuffle=True)  
test = datasets.fetch_20newsgroups(subset='test', categories=categories,  
remove=('headers', 'footers', 'quotes'), shuffle=True, random_state=42)  
  
# Get the number of documents in each category for train set  
cat1 = 0  
cat2 = 0  
for i in range(len(train.target)):  
    if train.target[i] == 0:  
        cat1 += 1  
    else:
```

```
cat2 += 1

# Remove all numbers and special characters from document texts
bad_patterns = "[^a-zA-Z. ]"

for doc in range(len(train.data)):
    new_doc = re.sub(bad_patterns, '', train.data[doc])
    train.data[doc] = new_doc

# Tokenize all documents in our training set and get the vocabulary.
vectorizer = CountVectorizer()
vectorizer.fit_transform(train.data)
vocabulary = vectorizer.vocabulary_
print(len(vocabulary))

# Use TfidfVectorizer() as initial pipeline to handle current setup of the
dataset.
model = pipeline.make_pipeline(CountVectorizer(), TfidfTransformer(),
MultinomialNB())

model.fit(train.data, train.target)

predicted = model.predict(test.data)

print("Number of documents in rec.autos: " + str(cat1))
print("Number of documents in comp.graphics: " + str(cat2))
print("Vocabulary Size: " + str(len(vocabulary)))
print(metrics.classification_report(test.target, predicted,
target_names=test.target_names))
```

The results of model 1 are shown below:

```
20604
Number of documents in rec.autos: 584
Number of documents in comp.graphics: 594
Vocabulary Size: 20604
```

	precision	recall	f1-score	support
comp.graphics	0.97	0.89	0.93	389
rec.autos	0.90	0.97	0.94	396
accuracy			0.94	785
macro avg	0.94	0.93	0.93	785
weighted avg	0.94	0.94	0.93	785

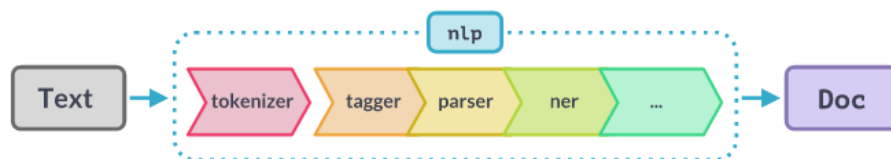
Model 2: Sklearn + SpaCy

For this model, I pre-processed all documents from both categories through a SpaCy pipeline. **Please note that the vocabulary size and number of documents within both categories does not change between models. Thus, this information was outputted only in model 1's results.**

The pipeline consisted of the following components:

1. Sentencizer: Apply sentence segmentation
2. Tokenizer: Apply word tokenization (Works in the background prior to Tagger)
3. Tagger: Assigned part-of-speech tags
4. Parser: Assigning dependency labels
5. Entity Recognizer: Detect and label entities

The below illustration details a general overview of how the pipeline works, excluding the sentence segmentation component:



The results shown at the bottom displayed very poor results overall, when in comparison to the pure sklearn model. No errors were thrown during this implementation, but the current pipeline may be insufficient or be in the incorrect order for SpaCy to properly process the documents.

```
"""# Model 2: Combo of Sklearn and SpaCy Features
```

```
Predefined sklearn pipeline is used below and then fed
```

to sklearn's CountVectorizer method.

Sklearn, re, and spaCy documentation was used in the creation of this program.

"""

```
import spacy as sp
from spacy.lang.en.stop_words import STOP_WORDS

# Re-import unedited datasets for new model
train2 = datasets.fetch_20newsgroups(subset='train', remove=('headers',
'footers', 'quotes'), categories=categories, shuffle=True)
test2 = datasets.fetch_20newsgroups(subset='test', categories=categories,
remove=('headers', 'footers', 'quotes'), shuffle=True, random_state=42)

# Remove unnecessary characters like numbers and special non-punctuation
characters
bad_patterns = "[^a-zA-Z.]"

for doc in range(len(train2.data)):
    new_doc = re.sub(bad_patterns, '', train2.data[doc])
    train2.data[doc] = new_doc

# Process each document individually using the below steps
# NOTE: Colab's spaCy library version is 2.2.4. Only version 3.0 has
# Lemmatization as a separate pipeline component. Thus, Lemmatization is
# implemented, but acts behind the scenes of the parser component.

# Import premade English processing pipeline
nlp = sp.load("en_core_web_sm")

# Add sentence segmentation
sentencizer = nlp.create_pipe("sentencizer")
nlp.add_pipe(sentencizer)

# spaCy pipeline for pre-processing
def spacy_pipeline(document):
    doc = nlp(document)
    return doc

# Create bag of words with spacy tokenizer
vectorizer = CountVectorizer(tokenizer=spacy_pipeline)

# Train the model
# model2 = pipeline.Pipeline([("bow", bag_of_words), ("classifier",
```



```
classifier)])  
model2 = pipeline.make_pipeline(vectorizer, TfidfTransformer(),  
MultinomialNB())  
model2.fit(train2.data, train2.target)  
  
predicted2 = model2.predict(test2.data)  
  
print(metrics.classification_report(test2.target, predicted2,  
target_names=test2.target_names, zero_division=1))
```

The results of model 2 are shown below:

	precision	recall	f1-score	support
comp.graphics	1.00	0.00	0.00	389
rec.autos	0.50	1.00	0.67	396
accuracy			0.50	785
macro avg	0.75	0.50	0.34	785
weighted avg	0.75	0.50	0.34	785