

```
1 nltk.download()
```

```

|   Unzipping grammars/spanish_grammars.zip.
| Downloading package state_union to /root/nltk_data...
|   Unzipping corpora/state_union.zip.
| Downloading package stopwords to /root/nltk_data...
|   Unzipping corpora/stopwords.zip.
| Downloading package subjectivity to /root/nltk_data...
|   Unzipping corpora/subjectivity.zip.
| Downloading package swadesh to /root/nltk_data...
|   Unzipping corpora/swadesh.zip.
|
| Downloading package switchboard to /root/nltk_data...
|   Unzipping corpora/switchboard.zip.
| Downloading package tagsets to /root/nltk_data...
|   Unzipping help/tagsets.zip.
| Downloading package timit to /root/nltk_data...
|   Unzipping corpora/timit.zip.
| Downloading package toolbox to /root/nltk_data...
|   Unzipping corpora/toolbox.zip.
| Downloading package treebank to /root/nltk_data...
|   Unzipping corpora/treebank.zip.
| Downloading package twitter_samples to /root/nltk_data...
|   Unzipping corpora/twitter_samples.zip.
| Downloading package udhr to /root/nltk_data...
|   Unzipping corpora/udhr.zip.
| Downloading package udhr2 to /root/nltk_data...
|   Unzipping corpora/udhr2.zip.
| Downloading package unicode_samples to /root/nltk_data...
|   Unzipping corpora/unicode_samples.zip.
| Downloading package universal_tagset to /root/nltk_data...
|   Unzipping taggers/universal_tagset.zip.
| Downloading package universal_treebanks_v20 to
|   /root/nltk_data...
| Downloading package vader_lexicon to /root/nltk_data...
| Downloading package verbnet to /root/nltk_data...
|   Unzipping corpora/verbnet.zip.
| Downloading package verbnet3 to /root/nltk_data...
|   Unzipping corpora/verbnet3.zip.
| Downloading package webtext to /root/nltk_data...
|   Unzipping corpora/webtext.zip.
| Downloading package wmt15_eval to /root/nltk_data...
|   Unzipping models/wmt15_eval.zip.
| Downloading package word2vec_sample to /root/nltk_data...
|   Unzipping models/word2vec_sample.zip.
| Downloading package wordnet to /root/nltk_data...
| Downloading package wordnet2021 to /root/nltk_data...
| Downloading package wordnet31 to /root/nltk_data...
| Downloading package wordnet_ic to /root/nltk_data...
|   Unzipping corpora/wordnet_ic.zip.
| Downloading package words to /root/nltk_data...
|   Unzipping corpora/words.zip.
| Downloading package ycoe to /root/nltk_data...
|   Unzipping corpora/ycoe.zip.

```

Done downloading collection all

d) Download l) List u) Update c) Config h) Help q) Quit

Downloader> q
True

```
1 import csv
2 import re
3 import nltk
4 import string
5 import math
6 import numpy as np
7 import sklearn
8 from nltk.corpus import stopwords
9 from pprint import pprint
10 from termcolor import colored
11 from sklearn.metrics import precision_recall_fscore_support
```

```
1 #tempaltes for printing
2 printReverseGreen = lambda x : print(colored(x,'green',attrs=['reverse','bold'])
3 printReverseRed = lambda x : print(colored(x,'red',attrs=['reverse','bold']))
```

```
1 def printDataset(data,row):
2     for i in range(row):
3         printReverseRed('Row '+str(i))
4         for attr in data[i]:
5             print(attr)
```

```
1 def createDataset(fname):
2     #reading dataset
3     lol=list(csv.reader(open(fname), delimiter=':'))
4     tarClass=list(set([e[0] for e in lol]))
5     tarClass.sort()
6     tarClass_map=dict(zip(tarClass,range(len(tarClass))))
7     printReverseGreen("Mapping for target class")
8     print(tarClass_map)
9     #Mapping target class to numeric value
10    data=[[tarClass_map[row[0]]," ".join(row[1].split()[1:]),None,None,None,None]
11    #removing punctuation and numeric value from question and generating vocab
12    vocab=[]
13    for row in data:
14        questn=[word.lower() for word in re.sub('[^a-zA-Z]', ' ',row[1]).split()]
15        row[2]=questn
16        row[3]=len(questn)
17        vocab.extend(questn)
18    #creating ngram
19    ngram=list(nltk.ngrams(vocab,1))
20    #counting a ngram
21    ngramFreq={}
22    for e in ngram:
23        if " ".join(e) in ngramFreq:ngramFreq[" ".join(e)]+=1
24        else:ngramFreq[" ".join(e)]=1
25    #sorting in decreasing order
26    ngramFreq=sorted(ngramFreq.items(),key=lambda x:x[1],reverse=True)
27    #selecting top 500 from sorted list
28    moFreq500=[e[0] for e in ngramFreq[:500]]
29    #converting ngramFreq again back to dictionary from list of tuples
```

```

29 #converting ngramFreq again back to dictionary from list of tuples
30 ngramFreq=dict(ngramFreq)
31 printReverseGreen("Dataset Attribute")
32 print("-"*120)
33 print("| Targte class | Raw Question | Question After preprocessing | length")
34 print("-"*120)
35 printReverseGreen("first 5 rows from dataset after processing the question")
36 printDataset(data,5)
37 #extracting lexical and syntactic data
38 for row in data:
39     lexical=[]
40     syntax=[]
41     for word in row[2]:
42         if word in moFreq500:
43             lexical.append(word)
44             syntax.append(nltk.pos_tag([word])[0][1])
45     row[4]=lexical
46     row[5]=syntax
47 printReverseGreen("first 5 rows from dataset after extracting lexical and syntactic data")
48 printDataset(data,5)
49 return tarClass_map,ngramFreq,moFreq500,data

```

```

1 def probWordGivenTarClass(data,word,column_no,tarClass,vocab_length):
2     tarClass_rows=[questn for questn in data if questn[0]==tarClass]
3     count_word_given_tarClass=sum([row[column_no].count(word) for row in tarClass_rows])
4     return (count_word_given_tarClass+1)/(len(tarClass_rows)+vocab_length)

```

```

1 def probWordsGivenTarClasses(tarClasses,column_no,data,words):
2     map_probWordGivenTarClass={}
3     for word in words:
4         map={}
5         for tarClass in tarClasses:
6             map[tarClass]=probWordGivenTarClass(data,word,column_no,tarClass,vocab_length)
7         map_probWordGivenTarClass[word]=map
8     return map_probWordGivenTarClass

```

```

1 def trainCalculateColProbs(data,mpwgtc,mppgtc):
2     for row in data:
3         row[4]=np.product([mpwgtc[word][row[0]] for word in row[4]])
4         row[5]=np.product([mppgtc[word][row[0]] for word in row[5]])
5     printReverseGreen("first 5 rows from training dataset after probability calculation")
6     printDataset(data,5)
7     return data

```

```

1 def testCalculateColProbs(data,mpwgtc,mppgtc):
2     for row in data:
3         row[4]=np.product([max(mpwgtc[word].values()) if word in mpwgtc else 1 for word in row[4]])
4         row[5]=np.product([max(mppgtc[word].values()) if word in mppgtc else 1 for word in row[5]])
5     printReverseGreen("first 5 rows from testing dataset after probability calculation")
6     printDataset(data,5)
7     return data
8

```

```

1 def bestSplit(data,algorithm):
2     transpose_data=np.array(data).T.tolist()
3     impurity_attribute=[]
4     #print(len(data))
5     for i in range(1,len(transpose_data)):
6         data=transpose_data[i]
7         data_target=list(zip(data,transpose_data[0]))
8         unique_data=list(set(data))
9         unique_data.sort()
10        #print("unique_data",unique_data)
11        split_position=[(unique_data[i]+unique_data[i+1])/2 for i in range(len(unique_data)-1)]
12        if len(unique_data)==1:
13            split_position=[unique_data[0]]
14        #print("split positions",split_position)
15        impurity=[]
16        minGiniSplit=1
17        for position in split_position:
18            le=[]
19            le_target_frequency={}
20            gt=[]
21            gt_target_frequency={}
22            for row in data_target:
23                if row[0]<=position:le.append(row)
24                else:gt.append(row)
25            for target in tarClass_map.values():
26                le_target_frequency[target]=0
27                gt_target_frequency[target]=0
28            for row in le:
29                le_target_frequency[row[1]]+=1
30            for row in gt:
31                gt_target_frequency[row[1]]+=1
32            le.append(None)
33            gt.append(None)
34            #gini
35            if algorithm=="gini":
36                algorithm_le=1-sum([(e[1]/len(le))**2 for e in le_target_frequency.items()])
37                algorithm_gt=1-sum([(e[1]/len(gt))**2 for e in gt_target_frequency.items()])
38                algorithm_split=(len(le)/(len(le)+len(gt)))*algorithm_le+(len(gt)/(len(le)+len(gt)))*algorithm_gt
39            #entropy
40            if algorithm=="entropy":
41                algorithm_le=sum([-1*(e[1]/len(le))*math.log2((e[1]+1)/len(le)) for e in le_target_frequency.items()])
42                algorithm_gt=sum([-1*(e[1]/len(gt))*math.log2((e[1]+1)/len(gt)) for e in gt_target_frequency.items()])
43                algorithm_split=(len(le)/(len(le)+len(gt)))*algorithm_le+(len(gt)/(len(le)+len(gt)))*algorithm_gt
44            #classification error
45            if algorithm=="ce":
46                algorithm_le=1-max([e[1]/len(le) for e in le_target_frequency.items()])
47                algorithm_gt=1-max([e[1]/len(gt) for e in gt_target_frequency.items()])
48                algorithm_split=(len(le)/(len(le)+len(gt)))*algorithm_le+(len(gt)/(len(le)+len(gt)))*algorithm_gt
49            tarClass=None
50            '''
51            if int(algorithm_split)==0:
52                for e in tarClass_map.values():
53                    if le_target_frequency[e]+gt_target_frequency[e]==len(data_target):
54                        tarClass=e

```

```

55         ...
56         target_frequency=[e,le_target_frequency[e]+gt_target_frequency[e]]
57         tarClass=sorted(target_frequency,key=lambda x:x[1],reverse=True)[0]
58         impurity.append([i,position,algorithm_split,tarClass])
59         #print("impurity",impurity)
60         #print("tarClass",tarClass)
61         impurity_attribute.append(sorted(impurity,key=lambda x:x[2])[0])
62         #print("impurity attribute",impurity_attribute)
63         #if tarClass:break
64     return sorted(impurity_attribute,key=lambda x:x[2])[0]

```

```

1 class node:
2     def __init__(self,data):
3         self.attr=None
4         self.val=None
5         self.impurity=None
6         self.data=data
7         self.left=None
8         self.right=None
9         self.tarClass=None

```

```

1 def bulidTree(train_data,algorithm):
2     start=node(train_data)
3     level=[start]
4     cnt=0
5     while level:
6         cur=level[0]
7         #if len(cur.data)>1:
8         decision=bestSplit(cur.data,algorithm)
9         cur.attr=decision[0]
10        cur.val=decision[1]
11        cur.impurity=decision[2]
12        cur.tarClass=decision[3]
13        #if len(cur.data)==1:
14        #cur.tarClass=cur.data[0][0]
15        #cur.impurity=0
16        if cur.impurity!=0:
17            left=[row for row in cur.data if row[cur.attr]<=cur.val]
18            right=[row for row in cur.data if row[cur.attr]>cur.val]
19            if left:
20                cur.left=node(left)
21                level.append(cur.left)
22            if right:
23                cur.right=node(right)
24                level.append(cur.right)
25        level.pop(0)
26        cnt+=1
27        if cnt==20000:
28            break
29    return start

```

```

1 def decisionTreeClassifier(start,attributes):

```

```

2     cur=start
3     tarClass=None
4     while cur and cur.attr:
5         #print("current attribute",cur.attr)
6         if attributes[cur.attr-1]<=cur.val:
7             tarClass=cur.tarClass
8             cur=cur.left
9         else:
10            tarClass=cur.tarClass
11            cur=cur.right
12     return tarClass

```

```

1 #creation of train dataset
2 tarClass_map,train_ngramFreq,train_moFreq500,train_data=createDataset("dt_train

```

```

None
Row 2
1
How can I find a list of celebrities ' real names ?
['how', 'can', 'find', 'list', 'of', 'celebrities', 'real', 'names']
8
None
None
Row 3
2
What fowl grabs the spotlight after the Chinese Year of the Monkey ?
['what', 'fowl', 'grabs', 'the', 'spotlight', 'after', 'the', 'chinese', 'y
12
None
None
Row 4
0
What is the full form of .com ?
['what', 'is', 'the', 'full', 'form', 'of', 'com']
7
None
None
first 5 rows from dataset after extracting lexical and syntactic data

Row 0
1
How did serfdom develop in and then leave Russia ?
['how', 'did', 'serfdom', 'develop', 'in', 'and', 'then', 'leave', 'russia'
9
['how', 'did', 'in', 'and', 'then']
['WRB', 'VBD', 'IN', 'CC', 'RB']
Row 1
2
What films featured the character Popeye Doyle ?
['what', 'films', 'featured', 'the', 'character', 'popeye', 'doyle']
7
['what', 'the', 'character']
['WP', 'DT', 'NN']
Row 2
1
How can I find a list of celebrities ' real names ?
['how', 'can', 'find', 'list', 'of', 'celebrities', 'real', 'names']
8
['how', 'can', 'find', 'list', 'of', 'real', 'names']
['how', 'can', 'find', 'list', 'of', 'real', 'names']

```

```

1 WRD , MD , VB , NN , IN , JJ , NNS ]
Row 3
2
What fowl grabs the spotlight after the Chinese Year of the Monkey ?
['what', 'fowl', 'grabs', 'the', 'spotlight', 'after', 'the', 'chinese', 'y
12
['what', 'the', 'after', 'the', 'chinese', 'year', 'of', 'the']
['WP', 'DT', 'IN', 'DT', 'JJ', 'NN', 'IN', 'DT']
Row 4
0
What is the full form of .com ?
['what', 'is', 'the', 'full', 'form', 'of', 'com']
7
['what', 'is', 'the', 'full', 'form', 'of', 'com']

```

```

1 #uniques pos from train data
2 unique_pos=set([e for row in train_data for e in row[-1]])

```

```

1 #mapper for word given all target classes
2 map_probWordGivenTarClass=probWordsGivenTarClasses(sorted(tarClass_map.values())

```

```

1 #mapper for pos given all target classes
2 map_prob_pos_given_tarClass=probWordsGivenTarClasses(sorted(tarClass_map.values

```

```

1 #calculate the probabilities for lexical and syntactical attribute
2 train_data_after_prob_cal=trainCalculateColProbs(train_data,map_probWordGivenTa

```

```

first 5 rows from training dataset after probability calculations
Row 0
1
How did serfdom develop in and then leave Russia ?
['how', 'did', 'serfdom', 'develop', 'in', 'and', 'then', 'leave', 'russia']
9
7.734196724847904e-08
0.00013266783798205007
Row 1
2
What films featured the character Popeye Doyle ?
['what', 'films', 'featured', 'the', 'character', 'popeye', 'doyle']
7
0.0013101844897959183
0.7824579360530675
Row 2
1
How can I find a list of celebrities ' real names ?
['how', 'can', 'find', 'list', 'of', 'celebrities', 'real', 'names']
8
2.7404778811929026e-13
2.4588954000112153e-05
Row 3
2
What fowl grabs the spotlight after the Chinese Year of the Monkey ?
['what', 'fowl', 'grabs', 'the', 'spotlight', 'after', 'the', 'chinese', 'yea
12
4.671329086249706e-10
0.08633452119520049

```

```

Row 4
0
What is the full form of .com ?
['what', 'is', 'the', 'full', 'form', 'of', 'com']
7
1.75915223862936e-12
0.02025786748867755

```

```

1 #selecting only target class,length,lexical,syntax attribute
2 train_data=[[row[0],row[3],row[4],row[5]] for row in train_data_after_prob_cal]

```

```

1 #creation of test dataset
2 tarClass_map,test_ngramFreq,test_moFreq500,test_data=createDataset("dt_test.csv")

```

```

None
Row 2
3
Who was Galileo ?
['who', 'was', 'galileo']
3
None
None
Row 3
1
What is an atom ?
['what', 'is', 'an', 'atom']
4
None
None
Row 4
5
When did Hawaii become a state ?
['when', 'did', 'hawaii', 'become', 'state']
5
None
None
first 5 rows from dataset after extracting lexical and syntactic data
Row 0
5
How far is it from Denver to Aspen ?
['how', 'far', 'is', 'it', 'from', 'denver', 'to', 'aspen']
8
['how', 'far', 'is', 'it', 'from', 'denver', 'to', 'aspen']
['WRB', 'RB', 'VBZ', 'PRP', 'IN', 'NN', 'TO', 'VB']
Row 1
4
What county is Modesto , California in ?
['what', 'county', 'is', 'modesto', 'california', 'in']
6
['what', 'county', 'is', 'modesto', 'california', 'in']
['WP', 'NN', 'VBZ', 'NN', 'NN', 'IN']
Row 2
3
Who was Galileo ?
['who', 'was', 'galileo']
3
['who', 'was', 'galileo']

```



```

1 WP , VBD , NN ]
Row 3
1
What is an atom ?
['what', 'is', 'an', 'atom']
4
['what', 'is', 'an', 'atom']
['WP', 'VBZ', 'DT', 'NN']
Row 4
5
When did Hawaii become a state ?
['when', 'did', 'hawaii', 'become', 'state']
5
['when', 'did', 'hawaii', 'become', 'state']
['WDD', 'VDD', 'NN', 'NN', 'NN']

```

```

1 #unique pos from test data
2 unique_pos=set([e for row in test_data for e in row[-1]])

```

```

1 #calculate the probabilities for lexical and syntactical attribute
2 test_data_after_prob_cal=testCalculateColProbs(test_data,map_probWordGivenTarCl
3 actual_target=np.array(test_data_after_prob_cal).T.tolist()[0]
4 test_predicted_algorithm=[]

```

first 5 rows from testing dataset after probability calculations

```

Row 0
5
How far is it from Denver to Aspen ?
['how', 'far', 'is', 'it', 'from', 'denver', 'to', 'aspen']
8
5.4853999922335294e-08
0.0008350159563586373
Row 1
4
What county is Modesto , California in ?
['what', 'county', 'is', 'modesto', 'california', 'in']
6
1.1856409585819994e-06
1.5584426076095899
Row 2
3
Who was Galileo ?
['who', 'was', 'galileo']
3
0.058518982365854205
0.469523748848022
Row 3
1
What is an atom ?
['what', 'is', 'an', 'atom']
4
0.005470750385762961
0.7515745108399025
Row 4
5
When did Hawaii become a state ?
['when', 'did', 'hawaii', 'become', 'state']
5

```

2.907471246516451e-06
0.6002377132007415

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: VisibleDeprecationWarning: This is separate from the ipykernel package so we can avoid doing imports u

```
1 for algorithm in ["gini","entropy","ce"]:
2     printReverseRed(algorithm)
3     #build tree
4     printReverseGreen("Building Tree ...")
5     start=bulidTree(train_data,algorithm)
6 #     prediction
7     printReverseGreen("Predicting test data ...")
8     test_predicted=[]
9     for questn in test_data_after_prob_cal:
10         test_predicted.append(decisionTreeClassifier(start,[questn[3],questn[4]
11     cnt=0
12     for i in range(len(actual_target)):
13         if actual_target[i]==test_predicted[i]:
14             cnt+=1
15     print("Accuracy",cnt/len(test_predicted)*100)
16     #calculating precision,recall,fscore
17     metric=precision_recall_fscore_support(actual_target,test_predicted,labels=
18     printReverseGreen(" "*10+"| ".join([e[0]+" "+str(e[1]) for e in tarClass_ma
19     metric_name=['precision','recall','f-score','support']
20     for i in range(len(metric)):
21         print(metric_name[i],metric[i].tolist())
22     test_predicted_algorithm.append(test_predicted)
```

gini

Building Tree ...

Predicting test data ...

Accuracy 15.0

ABBR 0| DESC 1| ENTY 2| HUM 3| LOC 4| NUM 5

precision [0.0, 0.18, 0.11920529801324503, 0.14213197969543148, 0.17142857142

recall [0.0, 0.06521739130434782, 0.19148936170212766, 0.4307692307692308, 0.

f-score [0.0, 0.09574468085106382, 0.1469387755102041, 0.2137404580152672, 0.

support [9, 138, 94, 65, 81, 113]

entropy

Building Tree ...

Predicting test data ...

Accuracy 16.6

ABBR 0| DESC 1| ENTY 2| HUM 3| LOC 4| NUM 5

precision [0.0, 0.23255813953488372, 0.14507772020725387, 0.1568627450980392,

recall [0.0, 0.07246376811594203, 0.2978723404255319, 0.36923076923076925, 0.

f-score [0.0, 0.11049723756906078, 0.19512195121951217, 0.22018348623853212,

support [9, 138, 94, 65, 81, 113]

ce

Building Tree ...

Predicting test data ...

Accuracy 18.6

ABBR 0| DESC 1| ENTY 2| HUM 3| LOC 4| NUM 5

precision [0.015625, 0.35135135135135137, 0.19402985074626866, 0.15, 0.202898

recall [0.11111111111111111, 0.18840579710144928, 0.2765957446808511, 0.276923

f-score [0.0273972602739726, 0.24528301886792453, 0.2280701754385965, 0.19459

support [9, 138, 94, 65, 81, 113]

```
1 #Observe how many samples are mis-classified using gini index based
2 #model but correctly classified by mis-classification error and
3 #cross-entropy based model.
4 false_gini_true_entropy=0
5 false_gini_true_ce=0
```

```
1 for i in range(len(test_predicted)):
2     if test_predicted_algo[0][i]!=actual_target[i] and test_predicted_algo[1][i]
3         false_gini_true_entropy+=1
4     if test_predicted_algo[0][i]!=actual_target[i] and test_predicted_algo[2][i]
5         false_gini_true_ce+=1
```

```
1 print("samples misclassified by Gini but correctly classified by entropy",false_g
2 print("samples misclassified by Gini but correctly classified by Classification Error")
```

```
samples misclassified by Gini but correctly classified by entropy 38
samples misclassified by Gini but correctly classified by Classification Error
```



[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 23:37

