

CS613 Final Project

Group 9: Stroke Prediction Model

Authors: Danny Li - Tien Nguyen - Emily Wang

Part1. Preprocessing

In [1]:

```
1  # Import libraries
2  import numpy as np
3  import math
4  import csv
5  import pandas as pd
6  import random
7  from collections import defaultdict
8  from matplotlib import pyplot as plt
```

In [44]:

```

1  def count_branch(var, y):
2      br = []
3      for i in range(len(var)):
4          if var.ndim == 2:
5              num = var[0,0]
6              #print(num)
7          else:
8              num = var[i]
9              if num not in br:
10                 br.append(num)
11         #print(br)
12     class_num = []
13     for i in range(len(y)):
14         if y[i] not in class_num:
15             class_num.append(y[i])
16
17     d_br = defaultdict()
18     for item in br:
19         d_br[item] = dict.fromkeys(class_num, 0)
20     for ele in range(len(y)):
21         for i in br:
22             for j in class_num:
23                 if var[ele]==i and y[ele] == j:
24                     d_br[i][j]+=1
25     return d_br
26
27 def entropy_ind(num, total):
28     """
29     This function return the entropy value given the ratio of the attribute
30     """
31     if num == 0:
32         return 0
33     else:
34         return (-num/total)*np.log2(num/total)
35
36 # This function calculate entropy of a branch
37 #input: a dict
38 def entropy_br(d_br):
39     br_en = defaultdict()
40     keys = list(d_br.keys())
41     total = np.sum(list(d_br.values()))
42     for item in keys:
43         br_en[item] = 0
44     for i in range(len(keys)):
45         br_en[keys[i]] += entropy_ind(d_br[keys[i]],total)
46     en= np.sum(list(br_en.values()))
47     return [en, total]
48
49 # This function calculate the entropy of a feature
50 #entropy and count could be an array
51 def entropy_var(entropy, count):
52     en_var = 0
53     total = np.sum(count)
54     if total == 0:
55         return entropy
56     for i in range(len(entropy)):
57         en_var += (count[i]/total)*entropy[i]
58     return en_var
59

```

```
60 # This function return the feature that has min entropy
61 def top12features(dataset, header):
62     en_all = []
63     en_dict = defaultdict()
64     for i in range(dataset.shape[1]-1):
65         en_ls = []
66         count_ls = []
67         d_br = count_branch(dataset[:,i],dataset[:, -1])
68         for key in d_br.keys():
69             ind_en, ind_count = entropy_br(d_br[key])
70             en_ls.append(ind_en)
71             count_ls.append(ind_count)
72         en_avg = entropy_var(np.array(en_ls), np.array(count_ls))
73         en_dict[i]=en_avg
74     en_sorted = sorted(en_dict.items(), key=lambda item:item[1])
75     top12 = en_sorted[:12]
76     ls_index = sorted([item[0] for item in top12])
77     d = defaultdict()
78     for item in ls_index:
79         d[header[item]]=dataset[:,item]
80     d[header[-1]] = dataset[:, -1]
81     df = pd.DataFrame(data=d)
82     return df
```

In [45]:

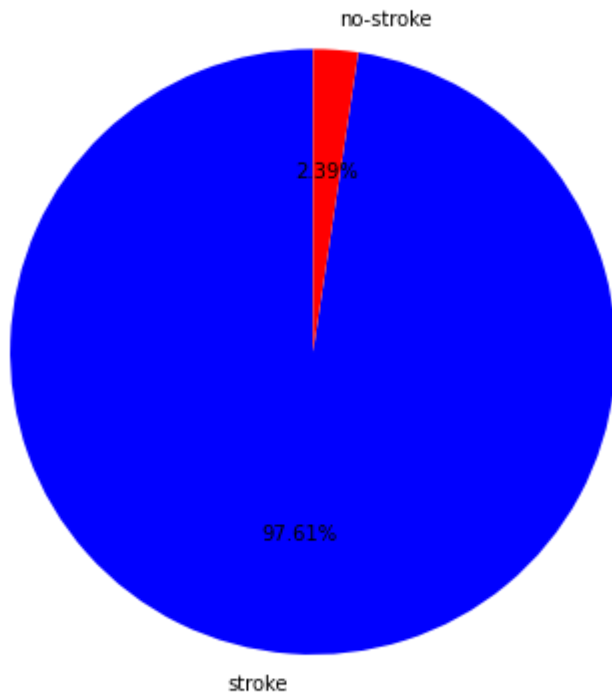
```

1  s_data = pd.read_csv('data.csv')
2  df = s_data.dropna(axis='columns')
3  f_index = list(range(2,20))
4  f_index.append(52)
5  df_initial = s_data.iloc[:,f_index]
6  df_initial = df_initial.dropna()
7  df_initial=df_initial.drop(df_initial[df_initial['RDEF1']=='C'].index)
8  df_initial=df_initial.drop(df_initial[df_initial['RDEF2']=='C'].index)
9  df_initial=df_initial.drop(df_initial[df_initial['RDEF3']=='C'].index)
10 df_initial=df_initial.drop(df_initial[df_initial['RDEF4']=='C'].index)
11 df_initial=df_initial.drop(df_initial[df_initial['RDEF5']=='C'].index)
12 df_initial=df_initial.drop(df_initial[df_initial['RDEF6']=='C'].index)
13 df_initial=df_initial.drop(df_initial[df_initial['RDEF7']=='C'].index)
14 df_initial=df_initial.drop(df_initial[df_initial['RDEF8']=='C'].index)
15 df_initial=df_initial.drop(df_initial[df_initial['DNOSTRK']=='U'].index)
16
17 df_initial.SEX = df_initial.SEX.replace({'M':0, 'F':1}).astype(np.uint8)
18 df_initial = df_initial.replace({'Y':1, 'N':0})
19 df_initial.RCONSC = df_initial.RCONSC.replace({'F':0, 'D':1, 'U':2}).astype(np.uint8)
20 df_initial['DNOSTRK'] = df_initial['DNOSTRK'].replace([0,1],[1,0])
21
22 plt.rcParams['font.sans-serif']=['SimHei']
23 plt.figure(figsize=(6,9))
24 labels = [u'stroke',u'no-stroke']
25 sizes = [12767, 312]
26 colors = ['blue', 'red']
27 explode = (0,0)
28 patches,text1,text2 = plt.pie(sizes,
29                               explode=explode,
30                               labels=labels,
31                               colors=colors,
32                               autopct = '%3.2f%%',
33                               shadow = False,
34                               startangle = 90,
35                               pctdistance = 0.6)
36 plt.axis('equal')
37 plt.show()

```

/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3146: DtypeWarning: Columns (31) have mixed types.Specify dtype option on import or set low_memory=False.

```
has_raised = await self.run_ast_nodes(code_ast.body, cell_name,
```



In [46]:

```

1 paper_features = df_initial[['AGE', 'SEX', 'RSLEEP', 'RATRIAL', 'RVISINF', 'RDEF1',
2 top12_features = top12features(np.array(df_initial), list(df_initial.columns))
3
4 p_nstrk = paper_features[paper_features['DNOSTRK']==0]
5 p_strk = (paper_features[paper_features['DNOSTRK']==1]).sample(n = len(p_nstrk))
6 p_balanced_data = pd.concat([p_strk, p_nstrk])
7
8 t_nstrk = top12_features[top12_features['DNOSTRK']==0]
9 t_strk = (top12_features[top12_features['DNOSTRK']==1]).sample(n = len(t_nstrk))
10 t_balanced_data = pd.concat([t_strk, t_nstrk])
11
12 df_initial.to_csv('18features.csv', header = True, sep = ",", index = False)
13 paper_features.to_csv('12paper_features.csv', header = True, sep = ",", index = False)
14 top12_features.to_csv('12top_features.csv', header = True, sep = ",", index = False)
15 p_balanced_data.to_csv('12paper_features_balanced.csv', header = True, sep = ",", index = False)
16 t_balanced_data.to_csv('12top_features_balanced.csv', header = True, sep = ",", index = False)

```

Part 2. Machine Learning Classifiers

2.1. Logistic Regression

In [65]:

```

1  def cal_pr(th, y_hat_val, y_val): #thread
2      # count TP, FP, FN
3      TP = 0
4      FP = 0
5      TN = 0
6      FN = 0
7
8      for i in range(len(y_hat_val)):
9          if y_hat_val[i] >= th and y_val[i] == 1:
10             TP += 1
11          elif y_hat_val[i][0] >= th and y_val[i][0] == 0:
12             FP += 1
13          elif y_hat_val[i][0] < th and y_val[i][0] == 1:
14             FN += 1
15          else:
16             TN +=1
17      if TP == 0:
18          return 0,0
19      pre = TP/(TP+FP)
20      rec = TP/(TP+FN)
21      f_measure = 2*pre*rec/(pre+rec)
22      accuracy = (TP + TN)/(TP+TN+FP+FN)
23      return pre, rec, f_measure, accuracy
24
25  def Logistic_Regression(filename):
26      data= pd.read_csv(filename)
27      spam_rand = data.sample(frac=1, random_state=0)
28
29      num_train = math.ceil(((2/3)*len(spam_rand)))
30      num_valid = len(spam_rand)-num_train
31
32      train = spam_rand.iloc[0:num_train,:]
33      val = spam_rand.iloc[num_train:,:]
34
35      X_train = np.ones(num_train).reshape(-1,1)
36      # zscore the training data
37      for j in range(len(train.columns)-1): #of features
38          vector = train.iloc[:,j].to_numpy()
39          mean = np.mean(vector)
40          std = np.std(vector,ddof=1)
41          stand = (vector-mean)/std
42          X_train = np.append(X_train,stand.reshape(-1,1),axis=1)
43
44      # add bias feature   validation part
45      X_val = np.ones(num_valid).reshape(-1,1) # // validation
46
47      # zscore the validation data
48      for j in range(len(train.columns)-1):
49          vector_val = val.iloc[:,j].to_numpy()
50          vector = train.iloc[:,j].to_numpy()
51          mean = np.mean(vector)
52          std = np.std(vector,ddof=1)
53          stand = (vector_val-mean)/std
54          X_val = np.append(X_val,stand.reshape(-1,1),axis=1)
55
56      np.random.seed(0)
57      omega = np.random.uniform(-0.01,0.01,X_train.shape[1]).reshape(-1,1)
58
59      y_hat = 1 / (1 + np.exp(-X_train.dot(omega)))

```

```

60 y_train = train.iloc[:, -1].to_numpy().reshape(-1,1)
61 j = np.mean(y_train*(np.log(y_hat))+(1-y_train)*(np.log(1-y_hat)))
62 y_val = val.iloc[:, -1].to_numpy().reshape(-1,1)
63 y_hat_val = 1 / (1 + np.exp(-X_val.dot(omega)))
64 j_val = np.mean(y_val*(np.log(y_hat_val))+(1-y_val)*(np.log(1-y_hat_val)))
65
66 # terminate until meet criteria
67 count = 0
68 value_change = 1
69
70 list1 = [j]
71 list2 = [j_val]
72
73 while count <= 1000 and value_change >= 2**-23: # number of epoch
74     #print(count)
75     omega = omega + (10**(-4))*(X_train.T.dot(y_train-y_hat))
76     y_hat = 1 / (1 + np.exp(-X_train.dot(omega)))
77     for i in range(len(y_hat)):
78         if y_hat[i][0] == 1:
79             y_hat[i][0] = y_hat[i][0] - 0.01
80         if y_hat[i][0] == 0:
81             y_hat[i][0] = y_hat[i][0] + 0.01
82     j_new = np.mean(y_train*(np.log(y_hat))+(1-y_train)*(np.log(1-y_hat)))
83     value_change = np.abs(j_new-j)
84
85     y_hat_val = 1 / (1 + np.exp(-X_val.dot(omega)))
86     for k in range(len(y_hat_val)):
87         if y_hat_val[k][0] == 1:
88             y_hat_val[k][0] = y_hat_val[k][0] - 0.01
89         if y_hat_val[k][0] == 0:
90             y_hat_val[k][0] = y_hat_val[k][0] + 0.01
91     j_val_new = np.mean(y_val*(np.log(y_hat_val))+(1-y_val)*(np.log(1-y_hat_val)))
92
93     list1.append(j_new)
94     list2.append(j_val_new)
95     j = j_new
96     count += 1
97
98 pre_train, recall_train, f_train, accuracy_train = cal_pr(0.5,y_hat, y_train)
99 print("Training")
100 print("Precision\tRecall\t\tf-Measure\tAccuracy")
101 print("{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(pre_train, recall_train, f_train, accuracy_train))
102
103 pre_valid, recall_valid, f_valid, accuracy_valid = cal_pr(0.5,y_hat_val, y_val)
104 print("Validation")
105 print("Precision\tRecall\t\tf-Measure\tAccuracy")
106 print("{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(pre_valid, recall_valid, f_valid, accuracy_valid))

```

In [66]:

```

1 # Logistic Regression for 18-Features dataset
2 LogisticRegression('18features.csv')

```

Training

Precision	Recall	f-Measure	Accuracy
0.976	1.000	0.988	0.976

Validation

Precision	Recall	f-Measure	Accuracy
0.977	1.000	0.988	0.977

In [67]:

```
1 # Naive Bayes for 12-paper-Features dataset
2 Logistic_Regression('12paper_features.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.976	1.000	0.988	0.976

Validation

Precision	Recall	f-Measure	Accuracy
0.977	1.000	0.988	0.977

In [68]:

```
1 # Logistic Regression for 12-top-Features dataset
2 Logistic_Regression('12top_features.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.976	1.000	0.988	0.976

Validation

Precision	Recall	f-Measure	Accuracy
0.977	1.000	0.988	0.977

In [69]:

```
1 # Logistic Regression for 12-paper-Features balanced dataset
2 Logistic_Regression('12paper_features_balanced.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.692	0.594	0.639	0.683

Validation

Precision	Recall	f-Measure	Accuracy
0.697	0.600	0.645	0.635

In [70]:

```
1 # Logistic Regression for 12-top-Features balanced dataset
2 Logistic_Regression('12top_features_balanced.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.709	0.619	0.661	0.700

Validation

Precision	Recall	f-Measure	Accuracy
0.683	0.487	0.569	0.591

2.2. Naive Bayes

In [36]:

```

1  # Function to z_score the training dataset
2  def zscore_train(dataset):
3      zx_train = (dataset - dataset.mean(axis = 0))/np.std(dataset,axis = 0,ddof
4      return zx_train
5  # Function to z-score the validation dataset using the mean and standard deviat
6  def zscore_valid(valid_set, train_set):
7      zx_valid = (valid_set - train_set.mean(axis = 0))/np.std(train_set,axis = 0
8      return zx_valid
9
10 # Function to count for output true positive (TP), true negative (TN), false po
11 def confMatrixPara(pred_y, y):
12     TP = 0
13     TN = 0
14     FP = 0
15     FN = 0
16     for i in range(len(pred_y)):
17         if pred_y[i] == y[i] ==1:
18             TP +=1
19         elif pred_y[i] == y[i] ==0:
20             TN +=1
21         elif (pred_y[i] == 0) and (y[i] ==1):
22             FN +=1
23         else:
24             FP +=1
25     return TP, TN, FP, FN
26
27 def count_class(array_y):
28     ls = []
29     for i in array_y:
30         if i not in ls:
31             ls.append(i)
32     return ls
33
34
35 # Function to calculate for the required statistics
36 def statistics_result(TP, TN, FP, FN, array):
37     precision =TP/(TP+FP)
38     recall = TP/(TP +FN)
39     f_measure = 2*precision*recall/(precision+recall)
40     accuracy = (1/len(array))*(TP+TN)
41     return precision, recall, f_measure, accuracy
42
43 def pdf(x, mean, std):
44     return np.exp(-0.5*((x-mean)/std)**2)/(std*np.sqrt(2*np.pi))
45
46 def probability(x_valid, prior, mean_ls, std_ls):
47     prob = np.log(prior)
48     for i in range(len(mean_ls)):
49         pdf_i = pdf(x_valid[i],mean_ls[i], std_ls[i])
50         if pdf_i ==0:
51             return 0
52         else:
53             prob = prob+ np.log(pdf_i)
54     return np.exp(prob)
55
56 # Function to apply threshold for classification
57 def apply_threshold(array, threshold):
58     array2 = np.zeros((array.shape))
59     for i in range(len(array)):

```

```

60         if array[i] >= threshold:
61             array2[i] = 1
62         else:
63             continue
64     return array2
65 def Naive_Bayes(filename):
66     data = pd.read_csv(filename, sep = ',')
67     data_list = data.values.tolist()
68
69     # Randomize the data
70     np.random.seed(0)
71     np.random.shuffle(data_list)
72     data_rand = np.array(data_list)
73
74
75     # Separate features data and rename to x and y
76     y_data = np.array(data_rand[:, -1])
77     x_data = np.array(data_rand[:, :-1])
78
79     # Selecte the first 2/3 (round up) of the data for training and the remaini
80     num_train = math.ceil(((2/3)*len(data_rand)))
81     y_train = y_data[:num_train]
82     x_train = x_data[:num_train,:]
83     y_valid = y_data[num_train:]
84     x_valid = x_data[num_train:,:]
85
86     # Zcores the features training data
87     zx_train = zscore_train(x_train)
88     zx_valid = zscore_valid(x_valid, x_train)
89
90     class_set = count_class(y_train)
91     data_set = []
92     prior_set = []
93     mean_set = []
94     std_set = []
95     for i in class_set:
96         array = zx_train[y_train==i]
97         data_set.append(array)
98         prior_ = len(array)/len(zx_train)
99         prior_set.append(prior_)
100         mean_ = np.array([np.mean(array[:, j]) for j in range(array.shape[1])])
101         std_ = np.array([np.std(array[:, j], ddof = 1) for j in range(array.shape
102         mean_ = mean_[std_ >= 1e-4]
103         std_ = std_[std_ >= 1e-4]
104         mean_set.append(mean_)
105         std_set.append(std_)
106
107     training = []
108     for i in range(len(zx_train)):
109         p_set = []
110         for j in range(len(class_set)):
111             p_ = probability(zx_train[i], prior_set[j], mean_set[j], std_set[j])
112             p_set.append(p_)
113         pred_train_y = class_set[p_set.index(max(p_set))]
114         training.append(pred_train_y)
115
116     TP_train, TN_train, FP_train, FN_train = confMatrixPara(training, y_train)
117     pre_train, recall_train, f_train, accuracy_train = statistics_result(TP_tra
118     print("Training")
119     print("Precision\tRecall\t\ttf-Measure\tAccuracy")
120     print("{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(pre_train, recall_train

```

```

121
122     prediction =[]
123     for i in range(len(zx_valid)):
124         p_set = []
125         for j in range(len(class_set)):
126             p_ = probability(zx_valid[i], prior_set[j], mean_set[j], std_set[j])
127             p_set.append(p_)
128         pred_y =class_set[p_set.index(max(p_set))]
129         prediction.append(pred_y)
130
131     TP_valid, TN_valid, FP_valid, FN_valid = confMatrixPara(prediction, y_valid)
132     pre_valid, recall_valid, f_valid, accuracy_valid = statistics_result(TP_val
133     print("Validation")
134     print("Precision\tRecall\t\tf-Measure\tAccuracy")
135     print("{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(pre_valid, recall_valid

```

In [37]:

```

1 # Naive Bayes for 18-Features dataset
2 Naive_Bayes('18features.csv')

```

Training			
Precision	Recall	f-Measure	Accuracy
0.980	0.943	0.961	0.925
Validation			
Precision	Recall	f-Measure	Accuracy
0.980	0.940	0.960	0.923

In [40]:

```

1 # Naive Bayes for 12-paper-features dataset
2 Naive_Bayes('12paper_features.csv')

```

Training			
Precision	Recall	f-Measure	Accuracy
0.979	0.963	0.971	0.943
Validation			
Precision	Recall	f-Measure	Accuracy
0.979	0.960	0.969	0.941

In [41]:

```

1 # Naive Bayes for 12-top-features dataset
2 Naive_Bayes('12top_features.csv')

```

Training			
Precision	Recall	f-Measure	Accuracy
0.979	0.951	0.965	0.932
Validation			
Precision	Recall	f-Measure	Accuracy
0.980	0.946	0.962	0.928

In [47]:

```
1 # Naive Bayes for 12-top-features dataset
2 Naive_Bayes('12paper_features_balanced.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.722	0.528	0.610	0.680

Validation

Precision	Recall	f-Measure	Accuracy
0.730	0.470	0.571	0.611

In [48]:

```
1 # Naive Bayes for 12-top-features dataset
2 Naive_Bayes('12top_features_balanced.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.727	0.594	0.654	0.702

Validation

Precision	Recall	f-Measure	Accuracy
0.675	0.487	0.566	0.587

2.3. Decision Tree

In [51]:

```

1  # This function return the feature that has min entropy
2  def find_node(dataset):
3      en_all = []
4      for i in range(dataset.shape[1]-1):
5          en_ls = []
6          count_ls = []
7          d_br = count_branch(dataset[:,i],dataset[:, -1])
8          for key in d_br.keys():
9              ind_en, ind_count = entropy_br(d_br[key])
10             en_ls.append(ind_en)
11             count_ls.append(ind_count)
12             en_avg = entropy_var(np.array(en_ls), np.array(count_ls))
13             en_all.append(en_avg)
14             node = en_all.index(min(en_all))
15             return node
16
17 # array is a dictionary
18 def leaf_node(array):
19     d = dict()
20     for key in array.keys():
21         num = 0
22         for item in array[key]:
23             if num <= array[key][item]:
24                 num = array[key][item]
25                 d[key]=item
26     return d
27
28
29 def create_set(pre_set, node):
30     count_br = set(pre_set[:,node])
31     d = {}
32     data = np.delete(pre_set, obj = node,axis =1)
33     for i in range(len(count_br)):
34         array_ls = []
35         for j in range(len(pre_set)):
36             if pre_set[:,node][j]==(list(count_br))[i]:
37                 array_ls.append(data[j,:])
38             d[(list(count_br))[i]] = np.array(array_ls)
39
40     return d
41
42 def DTL(data,root, d):
43     if len(data)<=1:
44         children=[]
45         for key, array in data.items():
46             n = array
47             shape_n = (np.array(n)).shape
48             ls_y =[]
49             for i in array[:, -1]:
50                 if i not in ls_y:
51                     ls_y.append(i)
52             if shape_n[0]<=1:
53                 return n[0][ -1]
54             elif shape_n[1] <=2:
55                 if len(ls_y)<=1:
56                     return ls_y[0]
57
58                 num_br0 = count_branch(array[:, -1],array[:, -1])
59                 if len(num_br0)<=1:

```

```

60         dict_ = list(num_br0.values())
61         keys_ = list(dict_[0].keys())
62         values_ = list(dict_[0].values())
63         a = keys_[values_.index(max(values_))]
64         return a
65     else:
66         d = leaf_node(num_br0)
67         return d
68     elif len(ls_y)<=1:
69         return ls_y[0]
70     else:
71         node = find_node(array)
72         child_data = create_set(array,node)
73         children.append({key:DTL(child_data,node,d)})
74 else:
75     children = []
76     for value, array in data.items():
77         if array.ndim ==1:
78             continue
79         elif array.shape[1] <=2:
80             num_br0 = count_branch(array[:, :-1], array[:, -1])
81             if np.isscalar(num_br0):
82                 d = num_br0
83             else:
84                 d = leaf_node(num_br0)
85             return d
86         else:
87             node = find_node(array)
88             child_data = create_set(array,node)
89             children.append({value:DTL(child_data,node,d)})
90
91     return {root:children}
92
93 def testing_tree(valid_sample, tree):
94     if tree is None:
95         return 1
96     elif len(tree)>0:
97         node = list(tree.keys())
98         values = list(tree.values()) #list dict type
99         if np.isscalar(values[0]):
100             return values
101         else:
102             ans = valid_sample[node]
103             for i in range(len(values[0])):
104                 d = values[0][i]
105                 if ans != list(d.keys()):
106                     continue
107                 else:
108                     new_tree = list(d.values())[0]
109                     if np.isscalar(new_tree):
110                         return new_tree
111                     new_array = del_col(valid_sample, node)
112                     result = testing_tree(new_array, new_tree)
113             return result
114
115     else:
116         d = list(tree.values())
117         return d
118
119 def del_col(valid_sample, node):
120     new_array = np.delete(valid_sample, obj = node)

```

```

121     return new_array
122
123 def decision_tree(filename):
124     # Read in the data from csv file and save it to a list
125     data = pd.read_csv(filename, sep = ',')
126     data_list = data.values.tolist()
127
128     # Randomize the data
129     np.random.seed(0)
130     np.random.shuffle(data_list)
131     data_rand = np.array(data_list)
132
133     # Separate features data and rename to x and y
134     y_data = np.array(data_rand[:, -1])
135     x_data = np.array(data_rand[:, :-1])
136
137     # Selecte the first 2/3 (round up) of the data for training and the remaini
138     num_train = math.ceil(((2/3)*len(data_rand)))
139     y_train = y_data[:num_train]
140     x_train = x_data[:num_train,:]
141     y_valid = y_data[num_train:]
142     x_valid = x_data[num_train:,:]
143
144     # Zcores the features training data
145     zx_train = zscore_train(x_train)
146     zx_valid = zscore_valid(x_valid, x_train)
147
148     mean_train = np.array([np.mean(zx_train[:, i]) for i in range(zx_train.shape
149     mean_valid = np.array([np.mean(zx_valid[:, i]) for i in range(zx_train.shape
150
151     zx_trainb = np.zeros((zx_train.shape))
152     for i in range(zx_train.shape[1]):
153         for j in range(len(zx_train)):
154             if zx_train[j, i] >= mean_train[i]:
155                 zx_trainb[j,i] = 1
156             else:
157                 continue
158     zx_validb = np.zeros((zx_valid.shape))
159     for i in range(zx_valid.shape[1]):
160         for j in range(len(zx_valid)):
161             if zx_valid[j, i] >= mean_valid[i]:
162                 zx_validb[j,i] = 1
163             else:
164                 continue
165
166     z_train = np.concatenate((zx_trainb, np.reshape(y_train, (len(y_train), 1))
167
168     data = z_train[:, :]
169     tree = defaultdict()
170     root = find_node(data)
171     dict_ = create_set(data, root)
172
173     decision_tree = DTL(dict_,root, None)
174     training = []
175     for i in range(len(zx_trainb)):
176         pre = testing_tree(zx_trainb[i], decision_tree)
177         training.append(pre)
178
179     TP_train, TN_train, FP_train, FN_train = confMatrixPara(training, y_train)
180     pre_train, recall_train, f_train, accuracy_train = statistics_result(TP_tra
181     print("Training")

```

```

182     print("Precision\tRecall\t\tf-Measure\tAccuracy")
183     print("{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(pre_train, recall_train,
184
185     pre_y = []
186     for i in range(len(zx_validb)):
187         pre = testing_tree(zx_validb[i], decision_tree)
188         pre_y.append(pre)
189
190     TP_valid, TN_valid, FP_valid, FN_valid = confMatrixPara(pre_y, y_valid)
191     pre_valid, recall_valid, f_valid, accuracy_valid = statistics_result(TP_val
192     print("Validation")
193     print("Precision\tRecall\t\tf-Measure\tAccuracy")
194     print("{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(pre_valid, recall_valid

```

In [52]:

```

1  # Decision Tree for 18-Features dataset
2  decision_tree('18features.csv')

```

Training

Precision	Recall	f-Measure	Accuracy
0.985	1.000	0.992	0.985

Validation

Precision	Recall	f-Measure	Accuracy
0.977	0.985	0.981	0.963

In [53]:

```

1  # Decision Tree for 12-paper-Features dataset
2  decision_tree('12paper_features.csv')

```

Training

Precision	Recall	f-Measure	Accuracy
0.976	1.000	0.988	0.975

Validation

Precision	Recall	f-Measure	Accuracy
0.975	0.999	0.987	0.973

In [54]:

```

1  # Decision Tree for 12top-Features dataset
2  decision_tree('12top_features.csv')

```

Training

Precision	Recall	f-Measure	Accuracy
0.977	1.000	0.988	0.976

Validation

Precision	Recall	f-Measure	Accuracy
0.976	0.997	0.986	0.973

In [55]:

```
1 # Decision Tree for 12-paper-Features balanced dataset
2 decision_tree('12paper_features_balanced.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.716	0.906	0.800	0.815

Validation

Precision	Recall	f-Measure	Accuracy
0.504	0.652	0.569	0.562

In [56]:

```
1 # Decision Tree for 12-top-Features balanced dataset
2 decision_tree('12top_features_balanced.csv')
```

Training

Precision	Recall	f-Measure	Accuracy
0.702	0.952	0.808	0.817

Validation

Precision	Recall	f-Measure	Accuracy
0.540	0.698	0.609	0.587

2.4. Random Forest

In [57]:

```

1  def random_forest(train, valid, sample_size, tree_numbers):
2      y_valid = valid[:, -1]
3      x_valid = valid[:, :-1]
4
5      # sample size
6      count = 0
7      tree_ls = []
8      while count < tree_numbers:
9          train_sample = np.array(random.choices(train, k = sample_size))
10         y_train = train_sample[:, -1]
11         x_train = train_sample[:, :-1]
12
13         # Zcores the features training data
14         zx_train = zscore_train(x_train)
15         zx_valid = zscore_valid(x_valid, x_train)
16
17         mean_train = np.array([np.mean(zx_train[:, i]) for i in range(zx_train.shape[1])])
18         mean_valid = np.array([np.mean(zx_valid[:, i]) for i in range(zx_train.shape[1])])
19
20         zx_trainb = np.zeros((zx_train.shape))
21         for i in range(zx_train.shape[1]):
22             for j in range(len(zx_train)):
23                 if zx_train[j, i] >= mean_train[i]:
24                     zx_trainb[j, i] = 1
25             else:
26                 continue
27         zx_validb = np.zeros((zx_valid.shape))
28         for i in range(zx_valid.shape[1]):
29             for j in range(len(zx_valid)):
30                 if zx_valid[j, i] >= mean_valid[i]:
31                     zx_validb[j, i] = 1
32             else:
33                 continue
34
35         z_train = np.concatenate((zx_trainb, np.reshape(y_train, (len(y_train), 1))))
36
37         data = z_train[:, :]
38
39         root = find_node(data)
40         dict_ = create_set(data, root)
41         decision_tree = DTL(dict_, root, None)
42         #print(decision_tree)
43         tree_ls.append(decision_tree)
44         count += 1
45
46     predictions = []
47
48     for i in range(len(zx_validb)):
49         pre_ls = []
50         for j in range(len(tree_ls)):
51             pre = testing_tree(zx_validb[i], tree_ls[j])
52             pre_ls.append(pre)
53         for item in range(len(pre_ls)):
54             if isinstance(pre_ls[item], list):
55                 ls = pre_ls[item]
56                 pre_ls[item] = ls[0]
57         pre_y = max(set(pre_ls), key = pre_ls.count)
58         predictions.append(pre_y)
59

```

```

60     TP_valid, TN_valid, FP_valid, FN_valid = confMatrixPara(predictions, y_valid)
61     pre_valid, recall_valid, f_valid, accuracy_valid = statistics_result(TP_valid, TN_valid, FP_valid, FN_valid)
62
63     return pre_valid, recall_valid, f_valid, accuracy_valid
64
65 def random_forest_results(filename, tree_numbers, sample_size):
66     # Read in the data from csv file and save it to a list
67     df = pd.read_csv(filename, sep = ',')
68     data_list = df.values.tolist()
69
70     # Randomize the data
71     np.random.seed(0)
72     np.random.shuffle(data_list)
73
74
75     # Selecte the first 2/3 (round up) of the data for training and the remainin
76     num_train = int(np.ceil((2/3)*len(data_list)))
77     train = data_list[:num_train]
78     valid = np.array(data_list[num_train:])
79
80     print("#tree\tPrecision\tRecall\t\tf-Measure\tAccuracy")
81     for i in tree_numbers:
82         pre_valid, recall_valid, f_valid, accuracy_valid = random_forest(train,
83             print("{:.2f}\t{:.3f}\t\t{:.3f}\t\t{:.3f}\t\t{:.3f}".format(i, pre_valid, recall_valid, f_valid, accuracy_valid))

```

In [58]:

```

1 # Random Forest for 18-Features dataset with sample size = 8000
2 random_forest_results('18features.csv', [1,10],8000)

```

#tree	Precision	Recall	f-Measure	Accuracy
1.00	0.978	0.988	0.983	0.967
10.00	0.977	0.994	0.986	0.972

In [61]:

```

1 # Random Forest for 12-paper-Features dataset with sample size = 8000
2 random_forest_results('12paper_features.csv', [1,10,100,200,300,500],8000)

```

#tree	Precision	Recall	f-Measure	Accuracy
1.00	0.977	0.991	0.984	0.968
10.00	0.977	0.998	0.987	0.975
100.00	0.977	0.999	0.988	0.976
200.00	0.977	1.000	0.988	0.977
300.00	0.977	0.999	0.988	0.976
500.00	0.977	0.999	0.988	0.976

In [62]:

```

1 # Random Forest for 12-paper-Features dataset with sample size = 6000
2 random_forest_results('12paper_features.csv', [1,10],6000)

```

#tree	Precision	Recall	f-Measure	Accuracy
1.00	0.977	0.988	0.983	0.966
10.00	0.977	0.998	0.987	0.975

In [64]:

```
1 # Random Forest for 12-top-features dataset with sample size = 8000
2 random_forest_results('12top_features.csv', [1,10],8000)
```

#tree	Precision	Recall	f-Measure	Accuracy
1.00	0.977	0.993	0.985	0.970
10.00	0.977	0.999	0.988	0.975

Part 3 Cross Validation on Balanced Datasets

3.1 Logistic Regression

In [76]:

```

1  def cross_valid_LR(filename, k):
2      sample = pd.read_csv(filename,header=0)
3      sample.shape
4      from scipy import stats
5      from sklearn.model_selection import KFold
6      for item in k:
7          kf = KFold(n_splits=item)
8          trains = []
9          tests = []
10         for train,test in kf.split(sample):
11             trains.append(train)
12             tests.append(test)
13
14         mse = []
15         ls = []
16         for num in range(item):
17             cv_train = sample.iloc[trains[num]]
18             cv_test = sample.iloc[tests[num]]
19             train_X = np.ones(cv_train.shape[0]).reshape(-1,1)
20             test_X = np.ones(cv_test.shape[0]).reshape(-1,1)
21             train_y = cv_train.iloc[:,12].to_numpy().reshape(-1,1)
22             test_y = cv_test.iloc[:,12].to_numpy().reshape(-1,1)
23
24             for i in range (12):
25                 cv_train_temp = cv_train.iloc[:,i].to_numpy()
26                 cv_train_temp_mean = cv_train_temp.mean()
27                 cv_train_temp_std = cv_train_temp.std()
28                 cv_train_temp_stand= (cv_train_temp-cv_train_temp_mean)/cv_train_temp_std
29
30                 train_X = np.append(train_X,cv_train_temp_stand.reshape(-1,1),axis=1)
31
32                 cv_test_temp = cv_test.iloc[:,i].to_numpy()
33                 cv_test_temp_stand = (cv_test_temp-cv_train_temp_mean)/cv_train_temp_std
34
35                 test_X = np.append(test_X,cv_test_temp_stand.reshape(-1,1),axis=1)
36
37             omega = np.random.uniform(-0.01,0.01,train_X.shape[1]).reshape(-1,1)
38
39             X_train = train_X
40             y_hat = 1 / (1 + np.exp(-X_train.dot(omega)))
41             y_train = train_y
42             j = np.mean(y_train*(np.log(y_hat))+(1-y_train)*(np.log(1-y_hat)))
43             y_val = test_y
44             X_val = test_X
45             y_hat_val = 1 / (1 + np.exp(-X_val.dot(omega)))
46             j_val = np.mean(y_val*(np.log(y_hat_val))+(1-y_val)*(np.log(1-y_hat_val)))
47
48             # terminate until meet criteria
49             count = 0
50             value_change = 1
51
52             list1 = [j]
53             list2 = [j_val]
54
55             while count <= 1000 and value_change >= 2**-23:
56                 #print(count)
57                 omega = omega + (10**(-4))*(X_train.T.dot(y_train-y_hat))
58                 y_hat = 1 / (1 + np.exp(-X_train.dot(omega)))
59                 for i in range(len(y_hat)):

```

```

60         if y_hat[i][0] == 1:
61             y_hat[i][0] = y_hat[i][0] - 0.01
62         if y_hat[i][0] == 0:
63             y_hat[i][0] = y_hat[i][0] + 0.01
64     j_new = np.mean(y_train*(np.log(y_hat))+(1-y_train)*(np.log(1-y_hat)))
65     value_change = np.abs(j_new-j)
66
67     y_hat_val = 1 / (1 + np.exp(-X_val.dot(omega)))
68     for k in range(len(y_hat_val)):
69         if y_hat_val[k][0] == 1:
70             y_hat_val[k][0] = y_hat_val[k][0] - 0.01
71         if y_hat_val[k][0] == 0:
72             y_hat_val[k][0] = y_hat_val[k][0] + 0.01
73     j_val_new = np.mean(y_val*(np.log(y_hat_val))+(1-y_val)*(np.log(1-y_hat_val)))
74
75     list1.append(j_new)
76     list2.append(j_val_new)
77     j = j_new
78     count +=1
79
80     n = 0
81     #y_hat_val = 1 / (1 + np.exp(-X_val.dot(omega)))
82     for i in range(len(y_val)):
83         if y_hat_val[i][0] >= 0.5 and y_val[i][0] == 1:
84             n += 1
85         elif y_hat_val[i][0] < 0.5 and y_val[i][0] == 0:
86             n += 1
87
88     acc = n/len(y_val)
89     ls.append(acc)
90     a = sum(ls)/len(ls)
91     print("k\tAccuracy")
92     print("{:.2f}\t{:.3f}".format(item, a))

```

In [78]:

```

1  # Cross-validation Logistic Regression on 12-paper-features-balanced dataset
2  cross_valid_LR("12paper_features_balanced.csv", [5,10,50,100,500])

```

k	Accuracy
5.00	0.484
10.00	0.645
50.00	0.667
100.00	0.667
500.00	1.000

3.2 Naive Bayes

In [87]:

```

1  def cross_valid_NB(filename, k_ls):
2      from sklearn.model_selection import KFold
3      sample = pd.read_csv(filename, header=0)
4      sample.shape
5
6      for item in k_ls:
7          kf = KFold(n_splits=item)
8          trains = []
9          tests = []
10         for train, test in kf.split(sample):
11             trains.append(train)
12             tests.append(test)
13
14         mse = []
15         ls = []
16         for num in range(item):
17             cv_train = sample.iloc[train[num]]
18             cv_test = sample.iloc[test[num]]
19             train_X = np.ones(cv_train.shape[0]).reshape(-1, 1)
20             test_X = np.ones(cv_test.shape[0]).reshape(-1, 1)
21             train_y = cv_train.iloc[:, 12].to_numpy().reshape(-1, 1)
22             test_y = cv_test.iloc[:, 12].to_numpy().reshape(-1, 1)
23
24             for i in range(12):
25
26                 cv_train_temp = cv_train.iloc[:, i].to_numpy()
27                 cv_train_temp_mean = cv_train_temp.mean()
28                 cv_train_temp_std = cv_train_temp.std()
29                 cv_train_temp_stand = (cv_train_temp - cv_train_temp_mean) / cv_train_temp_std
30
31                 train_X = np.append(train_X, cv_train_temp_stand.reshape(-1, 1), axis=1)
32
33                 cv_test_temp = cv_test.iloc[:, i].to_numpy()
34                 cv_test_temp_stand = (cv_test_temp - cv_train_temp_mean) / cv_train_temp_std
35
36                 test_X = np.append(test_X, cv_test_temp_stand.reshape(-1, 1), axis=1)
37
38                 X_and_lastcol = np.append(train_X, train_y.reshape(-1, 1), axis=1)
39                 X_df = pd.DataFrame(X_and_lastcol)
40
41
42                 train_spam = X_df[X_df[13] == 1]
43
44                 # split into Non-Spam samples
45                 train_nospam = X_df[X_df[13] == 0]
46
47                 dict_spam = {}
48                 for i in range(13):
49                     arr = train_spam.iloc[:, i].to_numpy()
50                     dict_spam[i] = []
51                     dict_spam[i].append(np.mean(arr))
52                     dict_spam[i].append(np.var(arr, ddof=1))
53                     dict_spam[i].append(np.std(arr, ddof=1))
54
55
56                 # create a dictionary to store the mean, variance, std for each feature
57                 dict_nospam = {}
58                 for i in range(13):
59                     arr = train_nospam.iloc[:, i].to_numpy()

```

```

60         dict_nospam[i] = []
61         dict_nospam[i].append(np.mean(arr))
62         dict_nospam[i].append(np.var(arr, ddof=1))
63         dict_nospam[i].append(np.std(arr, ddof=1))
64
65     import math
66     # calculate p(y=Spam)
67     pro_spam = len(train_spam)/len(train_X)
68
69     # calculate p(y=Non-Spam)
70     pro_nospam = len(train_nospam)/len(train_X)
71
72     # create a function to calculate guassian distribution
73     def P_spam(x, mean, std):
74         var = float(std)**2+10**-100
75         denom = (2*math.pi*var)**.5
76         num = math.exp(-(float(x)-float(mean))**2/(2*var))
77         return num/denom
78
79
80     # calcualte all the probability in spam class
81     lst_spam = []
82     for row in range(len(test_X)):
83         lst_s = []
84         for col in range(12):
85             pro = P_spam(test_X[row,:][col], dict_spam[col][0], dict_spam[col][1])
86             lst_s.append(pro+10**-100)
87         a = np.array(lst_s)
88         product = np.log(pro_spam) + np.sum(np.log(a))
89         lst_spam.append(product)
90
91     # calcualte all the probability in non-spam class
92     lst_nospam = []
93
94     for row in range(len(test_X)):
95         lst_nons = []
96         for col in range(12):
97             pro = P_spam(test_X[row,:][col], dict_nospam[col][0], dict_nospam[col][1])
98             lst_nons.append(pro+10**-300)
99         a = np.array(lst_nons)
100        product = np.log(pro_nospam)+np.sum(np.log(a))
101        lst_nospam.append(product)
102
103    # compare the probability for these two class
104    lst = []
105    for i in range(len(test_X)):
106        if lst_spam[i] >= lst_nospam[i]:
107            lst.append(1)
108        else:
109            lst.append(0)
110
111    # calculate Accuracy
112    count = 0
113    for i in range(len(test_X)):
114        if test_y[i] == lst[i]:
115            count += 1
116    Accuracy = count/len(test_X)
117    ls.append(Accuracy)
118    a = sum(ls)/len(ls)
119    print("k\tAccuracy")
120    print("{:.2f}\t{:.3f}".format(item, a))

```


In [88]:

```
1 # Cross-validation Naive on 12-paper-features-balanced dataset
2 cross_valid_NB("12top_features_balanced.csv", [5,10,50,100,500])
```

k	Accuracy
5.00	0.571
k	Accuracy
10.00	0.616
k	Accuracy
50.00	0.660
k	Accuracy
100.00	0.657
k	Accuracy
500.00	0.671

3.3 Decision Tree

In [89]:

```

1  class GadId3Classifier:
2      def fit(self, input, output):
3          data = input.copy()
4          data[output.name] = output
5          self.tree = self.decision_tree(data, data, input.columns, output.name)
6
7      def predict(self, input):
8          # convert input data into a dictionary of samples
9          samples = input.to_dict(orient='records')
10         predictions = []
11
12         # make a prediction for every sample
13         for sample in samples:
14             predictions.append(self.make_prediction(sample, self.tree, 1.0))
15
16         return predictions
17
18     def entropy(self, attribute_column):
19         # find unique values and their frequency counts for the given attribute
20         values, counts = np.unique(attribute_column, return_counts=True)
21
22         # calculate entropy for each unique value
23         entropy_list = []
24
25         for i in range(len(values)):
26             probability = counts[i]/np.sum(counts)
27             entropy_list.append(-probability*np.log2(probability))
28
29         # calculate sum of individual entropy values
30         total_entropy = np.sum(entropy_list)
31
32         return total_entropy
33
34     def information_gain(self, data, feature_attribute_name, target_attribute_name):
35         # find total entropy of given subset
36         total_entropy = self.entropy(data[target_attribute_name])
37
38         # find unique values and their frequency counts for the attribute to be split
39         values, counts = np.unique(data[feature_attribute_name], return_counts=True)
40
41         # calculate weighted entropy of subset
42         weighted_entropy_list = []
43
44         for i in range(len(values)):
45             subset_probability = counts[i]/np.sum(counts)
46             subset_entropy = self.entropy(data.where(data[feature_attribute_name]==values[i]))
47             weighted_entropy_list.append(subset_probability*subset_entropy)
48
49         total_weighted_entropy = np.sum(weighted_entropy_list)
50
51         # calculate information gain
52         information_gain = total_entropy - total_weighted_entropy
53
54         return information_gain
55
56     def decision_tree(self, data, original_data, feature_attribute_names, target_attribute_name):
57         # base cases:
58         # if data is pure, return the majority class of subset
59         unique_classes = np.unique(data[target_attribute_name])

```

```

60     if len(unique_classes) <= 1:
61         return unique_classes[0]
62     # if subset is empty, ie. no samples, return majority class of original data
63     elif len(data) == 0:
64         majority_class_index = np.argmax(np.unique(original_data[target_attribute_name])[0])
65         return np.unique(original_data[target_attribute_name])[majority_class_index]
66     # if data set contains no features to train with, return parent node class
67     elif len(feature_attribute_names) == 0:
68         return parent_node_class
69     # if none of the above are true, construct a branch:
70     else:
71         # determine parent node class of current branch
72         majority_class_index = np.argmax(np.unique(data[target_attribute_name], return_counts=True)[1])
73         parent_node_class = unique_classes[majority_class_index]
74
75         # determine information gain values for each feature
76         # choose feature which best splits the data, ie. highest value
77         ig_values = [self.information_gain(data, feature, target_attribute_name) for feature in feature_attribute_names]
78         best_feature_index = np.argmax(ig_values)
79         best_feature = feature_attribute_names[best_feature_index]
80
81         # create tree structure, empty at first
82         tree = {best_feature: {}}
83
84         # remove best feature from available features, it will become the parent node
85         feature_attribute_names = [i for i in feature_attribute_names if i != best_feature]
86
87         # create nodes under parent node
88         parent_attribute_values = np.unique(data[best_feature])
89         for value in parent_attribute_values:
90             sub_data = data.where(data[best_feature] == value).dropna()
91
92             # call the algorithm recursively
93             subtree = self.decision_tree(sub_data, original_data, feature_attribute_names)
94
95             # add subtree to original tree
96             tree[best_feature][value] = subtree
97
98         return tree
99
100     def make_prediction(self, sample, tree, default=1):
101         # map sample data to tree
102         for attribute in list(sample.keys()):
103             # check if feature exists in tree
104             if attribute in list(tree.keys()):
105                 try:
106                     result = tree[attribute][sample[attribute]]
107                 except:
108                     return default
109
110             result = tree[attribute][sample[attribute]]
111
112             # if more attributes exist within result, recursively find best result
113             if isinstance(result, dict):
114                 return self.make_prediction(sample, result)
115             else:
116                 return result
117
118     def cross_valid_DT(filename, k_ls):
119         sample = pd.read_csv(filename, header=0)
120
121         for item in k_ls:

```

```

121 from sklearn.model_selection import KFold
122 kf = KFold(n_splits=item)
123 trains = []
124 tests = []
125 for train,test in kf.split(sample):
126     trains.append(train)
127     tests.append(test)
128
129 mse = []
130 ls = []
131 for num in range(item):
132     cv_train = sample.iloc[trains[num]]
133     cv_test = sample.iloc[tests[num]]
134     train_X = np.ones(cv_train.shape[0]).reshape(-1,1)
135     test_X = np.ones(cv_test.shape[0]).reshape(-1,1)
136     train_y = cv_train.iloc[:,12].to_numpy().reshape(-1,1)
137     test_y = cv_test.iloc[:,12].to_numpy().reshape(-1,1)
138
139     for i in range (12):
140
141         cv_train_temp = cv_train.iloc[:,i].to_numpy()
142         cv_train_temp_mean = cv_train.iloc[:,i].mean()
143         cv_train_temp_std = cv_train.iloc[:,i].std()
144         cv_train_temp_stand= (cv_train_temp-cv_train_temp_mean)/cv_train
145
146         train_X = np.append(train_X,cv_train_temp_stand.reshape(-1,1),a
147
148         cv_test_temp = cv_test.iloc[:,i].to_numpy()
149         cv_test_temp_stand = (cv_test_temp-cv_train_temp_mean)/cv_train
150
151         test_X = np.append(test_X,cv_test_temp_stand.reshape(-1,1),axis
152
153     train_X = train_X[:,1:13]
154     test_X = test_X[:,1:13]
155
156
157     train_X = pd.DataFrame(train_X)
158     train_y = pd.DataFrame(train_y)
159     test_X = pd.DataFrame(test_X)
160     test_y = pd.DataFrame(test_y)
161     train_X.name = 'train_X'
162     train_y.name = 'train_y'
163     test_X.name = 'test_X'
164     test_y.name = 'test_y'
165
166
167     from sklearn.metrics import accuracy_score
168     model = GadId3Classifier()
169     model.fit(train_X, train_y)
170
171     y_hat = model.predict(test_X)
172
173     # print(accuracy_score(test_y, y_hat))
174     ls.append(accuracy_score(test_y, y_hat))
175     a = sum(ls)/len(ls)
176
177     print("k\tAccuracy")
178     print("{:.2f}\t{:.3f}".format(item, a))

```

In [90]:

```
1 # Cross-validation Decision Tree on 12-paper-features-balanced dataset  
2 cross_valid_DT("12paper_features_balanced.csv", [5,10,50,100,500])
```

k	Accuracy
5.00	0.505
k	Accuracy
10.00	0.534
k	Accuracy
50.00	0.567