

```

1 from sklearn.metrics import roc_auc_score
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import ShuffleSplit
4 import os
5 from sklearn import metrics
6 import pandas as pd
7 import numpy as np
8 from sklearn.model_selection import cross_val_score
9 from sklearn.model_selection import KFold
10 import pandas
11 from keras.models import Sequential
12 from keras.layers import Dense
13 from keras.wrappers.scikit_learn import KerasClassifier
14 from keras.utils import np_utils
15 from sklearn.model_selection import cross_val_score
16 from sklearn.model_selection import KFold
17 from sklearn.preprocessing import LabelEncoder
18 from sklearn.pipeline import Pipeline

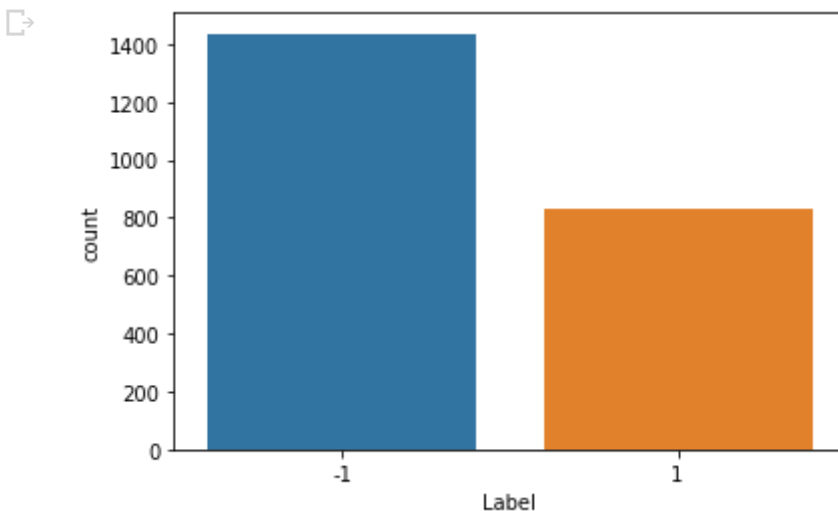
```

Load Dataset

```

1 xdata=pd.read_csv("./trainset.data")
2 ydata=pd.read_csv("./testset.dat")
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 sns.countplot(xdata['Label'],label="Sum")
6 plt.show()
7 xlab=xdata['Label']
8 xid=xdata['Sequence']
9 yid=ydata['ID']
10
11 pdA=xdata.drop(['Label'], axis = 1)
12 pdB=ydata.drop(['ID'], axis = 1)
13 pdB=pdB.rename(columns = {" Sequence": "Sequence"})

```



Composition of k-spaced amino acid pairs (CKSAAP)

```

1 def returnCKSAAP(pdA):
2     AA = 'ACDEFGHIKLMNPQRSTVWY'
3     gap=1
4     encodings = []
5     aaPairs = []
6     for aa1 in AA:
7         for aa2 in AA:
8             aaPairs.append(aa1 + aa2)
9     for i in pdA['Sequence']:
10        name, sequence = i,i
11        code = []
12        for g in range(gap+1):
13            myDict = {}
14            for pair in aaPairs:
15                myDict[pair] = 0
16            sum = 0
17            for index1 in range(len(sequence)):
18                index2 = index1 + g + 1
19                if index1 < len(sequence) and index2 < len(sequence) and sequence[index1] != sequence[index2]:
20                    myDict[sequence[index1] + sequence[index2]] = myDict[sequence[index1] + sequence[index2]] + 1
21                sum = sum + 1
22            for pair in aaPairs:
23                code.append(myDict[pair] / sum)
24        encodings.append(code)
25    return encodings

```

Dipeptide composition (DPC)

```

1 def returnDPC(pdA):
2     AA = 'ACDEFGHIKLMNPQRSTVWY'
3     encodings = []
4     diPeptides = [aa1 + aa2 for aa1 in AA for aa2 in AA]
5
6
7     AADict = {}
8     for i in range(len(AA)):
9         AADict[AA[i]] = i
10
11    for i in pdA['Sequence']:
12        name, sequence = i,i
13        code = []
14        tmpCode = [0] * 400
15        for j in range(len(sequence) - 2 + 1):
16            tmpCode[AADict[sequence[j]] * 20 + AADict[sequence[j+1]]] = tmpCode[AADict[sequence[j]] * 20 + AADict[sequence[j+1]]] + 1
17        if sum(tmpCode) != 0:
18            tmpCode = [i/sum(tmpCode) for i in tmpCode]
19        code = code + tmpCode
20        encodings.append(code)
21    return encodings

```

800 features from CKSAAP and 400 features from DPC

```

1 train_features1=np.array(returnCKSAAP(pdA))
2 print(train_features1.shape)
3 unique_d1 = [list(x) for x in set(tuple(x) for x in train_features1)]
4 print(len(unique_d1))
5
6 test_features1=np.array(returnCKSAAP(pdB))
7 print(test_features1.shape)
8 unique_d1 = [list(x) for x in set(tuple(x) for x in test_features1)]
9 print(len(unique_d1))
10
11 train_features2=np.array(returnDPC(pdA))
12 print(train_features2.shape)
13 unique_d1 = [list(x) for x in set(tuple(x) for x in train_features2)]
14 print(len(unique_d1))
15
16 test_features2=np.array(returnDPC(pdB))
17 print(test_features2.shape)
18 unique_d1 = [list(x) for x in set(tuple(x) for x in test_features2)]
19 print(len(unique_d1))
20
21
22 train_features=np.concatenate((train_features1,train_features2), axis=1)
23 print(train_features.shape)
24 unique_d1 = [list(x) for x in set(tuple(x) for x in train_features)]
25 print(len(unique_d1))
26
27
28 test_features=np.concatenate((test_features1,test_features2), axis=1)
29 print(test_features.shape)
30 unique_d1 = [list(x) for x in set(tuple(x) for x in test_features)]
31 print(len(unique_d1))

```

```

↳ (2270, 800)
2270
(567, 800)
567
(2270, 400)
2267
(567, 400)
567
(2270, 1200)
2270
(567, 1200)
567

```

I. Best: K-Neighbour (CKSAAP+DPC)

```

1 from sklearn.neighbors import KNeighborsClassifier
2 neigh = KNeighborsClassifier(n_neighbors=2)
3 ss=ShuffleSplit(n_splits=20, test_size=0.5, random_state=10)
4 scores = cross_val_score(neigh, train_features,xlab, cv=ss,n_jobs=-1, verbose=1)
5 print(np.mean(scores))
6 print(scores)

```

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
0.8038766519823788
[0.79118943 0.78942731 0.81321586 0.80352423 0.82643172 0.80176211
 0.82555066 0.80528634 0.79471366 0.80352423 0.80264317 0.80528634
 0.79295154 0.81409692 0.80440529 0.8          0.79295154 0.80881057
 0.80528634 0.79647577]
[Parallel(n_jobs=-1)]: Done 20 out of 20 | elapsed: 47.9s finished

```

```

1 neigh.fit(train_features, xlab)
2 pred= neigh.predict(test_features)
3 ids= np.arange(1001, 1568)
4 df= pd.DataFrame(ids, columns= ["ID"])
5 df['Label']= pred
6 df.to_csv('./submnk.csv', index= False)

```

II Best:Keras Classifier(CKSAAP Features)

```

1 def baseline_model():
2     # create model
3     model = Sequential()
4     model.add(Dense(1200, input_dim=800, activation='relu'))
5     model.add(Dense(2, activation='softmax'))
6     model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
7     return model
8 estimator = KerasClassifier(build_fn=baseline_model, epochs=200, batch_size=5,verbose=0)
9 kfold = KFold(n_splits=5, shuffle=True)
10 results = cross_val_score(estimator,train_features1,xlab, cv=kfold)
11 print(np.mean(results))
12 print(results)

```

```

1 estimator.fit(train_features1, xlab)
2 pred= estimator.predict(test_features1)
3 ids= np.arange(1001, 1568)
4 df= pd.DataFrame(ids, columns= ["ID"])
5 df['Label']= pred
6 df.to_csv('./submkc.csv', index= False)

```

```
1
```

