

گزارش – پروژه Ensemble Learning

در بخش اول از این تمرین، باید دیتاست را لود کنیم. چون از محیط Google Colab برای اجرای برنامه استفاده میکنم، از دستور زیر استفاده میکنیم.

```
1 from google.colab import drive
2 drive.mount('/content/gdrive')
```

در مرحله بعدی، از فایل‌های pairsDevTrain.txt و pairsDevTest.txt که مشخص کننده ی داده های آموزشی و تست هستند، باید عکسهای آموزشی و تست را پیدا کنیم. پس لازم است این دو فایل را بصورت خط به خط بخوانیم. مشخص است که باید علائمی نظیر '[' و ']' و '\t' و '\n' را از خطوط حذف کنیم:

```
1 lines = []
2 with open('/content/gdrive/My Drive/payanterm_hoosh/pairsDevTrain.txt', 'r') as f:
3     lines = f.readlines()
```

```
1 lines_test = []
2 with open('/content/gdrive/My Drive/payanterm_hoosh/pairsDevTest.txt', 'r') as f:
3     lines_test = f.readlines()
```

```
1 lin = []
2 for line in lines:
3     line = line.replace("\n", "")
4     x = line.split("\t")
5     lin.append(x)
```

```
1 lin_test = []
2 for line in lines_test:
3     line = line.replace("\n", "")
4     x = line.split("\t")
5     lin_test.append(x)
```

حالا باید داده های آموزشی و تست را از پوشه افراد داده شده بخوانیم. برای این کار ابتدا سراغ داده هایی میرویم که مربوط به یک نفر هستند. با توجه به اسم شخص و شماره ای که جلوی نام او نوشته شده، درون پوشه ها به جستجو میپردازیم.

```

1 Xtrain_first = []
2 Xtrain_second = []
3 Ytrain = []
4
5 for i in range(int(lin[0][0])):
6     Ytrain.append(1)
7     if int(lin[i+1][1]) < 10:
8         st = f'000{lin[i+1][1]}'
9     elif int(lin[i+1][1]) < 100:
10        st = f'00{lin[i+1][1]}'
11    elif int(lin[i+1][1]) < 1000:
12        st = f'0{lin[i+1][1]}'
13    else:
14        st = lin[i+1][1]
15
16    if int(lin[i+1][2]) < 10:
17        st1 = f'000{lin[i+1][2]}'
18    elif int(lin[i+1][2]) < 100:
19        st1 = f'00{lin[i+1][2]}'
20    elif int(lin[i+1][2]) < 1000:
21        st1 = f'0{lin[i+1][2]}'
22    else:
23        st1 = lin[i+1][2]
24
25    x = open(f'/content/gdrive/My Drive/payanterm_hoosh/lfw/{lin[i+1][0]}/{lin[i+1][0]}_{st}.txt', 'r')
26    y = open(f'/content/gdrive/My Drive/payanterm_hoosh/lfw/{lin[i+1][0]}/{lin[i+1][0]}_{st1}.txt', 'r')
27

```

این کار را برای داده هایی که متعلق به دو نفر متفاوت هستند نیز تکرار میکنیم.

```

1 for i in range(int(lin[0][0]), 2*int(lin[0][0])):
2     Ytrain.append(0)
3     if int(lin[i+1][1]) < 10:
4         st = f'000{lin[i+1][1]}'
5     elif int(lin[i+1][1]) < 100:
6         st = f'00{lin[i+1][1]}'
7     elif int(lin[i+1][1]) < 1000:
8         st = f'0{lin[i+1][1]}'
9     else:
10        st = lin[i+1][1]
11
12    if int(lin[i+1][3]) < 10:
13        st1 = f'000{lin[i+1][3]}'
14    elif int(lin[i+1][3]) < 100:
15        st1 = f'00{lin[i+1][3]}'
16    elif int(lin[i+1][3]) < 1000:
17        st1 = f'0{lin[i+1][3]}'
18    else:
19        st1 = lin[i+1][3]
20
21    x = open(f'/content/gdrive/My Drive/payanterm_hoosh/lfw/{lin[i+1][0]}/{lin[i+1][0]}_{st}.txt', 'r')
22    y = open(f'/content/gdrive/My Drive/payanterm_hoosh/lfw/{lin[i+1][2]}/{lin[i+1][2]}_{st1}.txt', 'r')
23

```

همین عملیات را برای داده های تست نیز تکرار میکنیم تا داده های تست نیز کامل شوند.

اما هم داده های تست هم داده های آموزشی دارای ابعاد مساوی نیستند. مثلاً بعضی از داده ها دارای 4 سطر و بعضی دیگر دارای 5 سطر هستند. پس باید آنها را به داده های سطری تبدیل کرد. چون همه دادگان دارای 512 ویژگی (پیکسل) هستند، پس میتوان هر داده را به یک بردار 512 تایی تبدیل کرد که همه داده ها از نظر ابعاد مساوی شوند.

برای اینکه داده ها بصورت آرایه های دوبعدی تبدیل شوند، سعی میکنیم 512 پیکسل را بعنوان بعد دوم قرار دهیم که در تصویر زیر این تبدیلات مشخص شده اند.

```
1 flatten_list = lambda y:[x for a in y for x in flatten_list(a)] if type(y) is list else [y]
2 temp = flatten_list(Xtrain_first)
3 temp2 = flatten_list(Xtrain_second)
```

```
1 flatten_list = lambda y:[x for a in y for x in flatten_list(a)] if type(y) is list else [y]
2 temp3 = flatten_list(Xtest_first)
3 temp4 = flatten_list(Xtest_second)
```

```
1 import numpy as np
2 Xtrain_first = np.array(temp).reshape(-1, 512)
3 Xtrain_second = np.array(temp2).reshape(-1, 512)
```

```
1 Xtest_first = np.array(temp3).reshape(-1, 512)
2 Xtest_second = np.array(temp4).reshape(-1, 512)
```

حالا باید تصمیم بگیریم که داده های دسته اول که مربوط به شخص اول هستند، چگونه به داده های دسته دوم که مربوط به شخص دوم هستند، مرتبط شوند. راه ابتدایی که به ذهن من رسید این بود که داده های دسته اول را به داده های دسته دوم concatenate کنیم. در تصویر زیر نمایش آنرا میبینیم که از طریق `axis = 1` این کار انجام شده است. بنابراین عکسهایی که قرار است باهم مقایسه شوند، در کنار هم قرار میگیرند.

```
1 import pandas as pd
2
3 df1 = pd.DataFrame(Xtrain_first)
4 df2 = pd.DataFrame(Xtrain_second)
5 print(df1.shape)
6 print(df2.shape)
7
8 df_contact = pd.concat([df1, df2], axis=1)
```

```
1 df3 = pd.DataFrame(Xtest_first)
2 df4 = pd.DataFrame(Xtest_second)
3 print(df3.shape)
4 print(df4.shape)
5
6 df_contact_test = pd.concat([df3, df4], axis=1)
```

قبل از اینکه به پیاده سازی الگوریتمها بپردازیم، میتوانیم از یک روش برای بهبود عملکردها استفاده کنیم. اگر داده های آموزشی و target آنها را باهم concatenate کنیم و سپس کل مجموعه را shuffle کنیم، باعث بهبود الگوریتم خواهد شد. چون داده ها بصورت پیش فرض به این صورت است که همه داده هایی که $y = 1$ هستند پشت سر هم است و همه داده هایی که $y = 0$ است پشت سر هم. این باعث عدم یادگیری صحیح خواهد شد. پس:

```
1 train = pd.concat([df_contact, df_y], ignore_index=True, sort=False, axis=1)
2 test = pd.concat([df_contact_test, df_y_test], ignore_index=True, sort=False, axis=1)

1 train = train.sample(frac = 1)
2 test = test.sample(frac = 1)

1 df_contact = train.loc[:, :1023]
2 df_y = train.loc[:, 1024]
3 df_contact_test = test.loc[:, :1023]
4 df_y_test = test.loc[:, 1024]
```

اکنون میخواهیم به بررسی الگوریتمهای مختلف روی این دیتاست بپردازیم. ابتدا از الگوریتم Bagging آغاز میکنیم. این الگوریتم سه پارامتر خیلی مهم دارد. میتوانیم آرگومان `base_estimator` را برابر `DecisionTreeClassifier` قرار دهیم. بنابراین پیش بینی کننده ی پایه ای درخت تصمیم خواهد بود. بسته به میزان مصرف زمان و حافظه، تعداد درختان را عوض میکنیم. در اینجا بهترین حالت که یک مصالحه ای برقرار کند را برابر 250 قرار میدهیم. عدد `random_state` هم برابر 7 یا اعداد دیگر میتوانیم قرار دهیم. پس تصویر زیر نمایشی از آن است:

```
1 from sklearn import model_selection
2 from sklearn.ensemble import BaggingClassifier
3 from sklearn.tree import DecisionTreeClassifier
4
5 model = BaggingClassifier(base_estimator = DecisionTreeClassifier(), n_estimators = 250, random_state=7)
6 model.fit(df_contact, df_y.ravel())
7 yhat = model.predict(df_contact_test)
```

```
1 from sklearn.metrics import accuracy_score
2 accuracy_score(df_y_test, yhat)
```

0.689

مشخص است که دقت حدود 69 درصد برای این الگوریتم بدست آمده است.

سپس به الگوریتم AdaBoost میپردازیم. این مدل هم دو آرگومان مهم دارد. تعداد پیش بینی کننده ها را در این بخش برابر 1000 قرار میدهیم. نتیجه در تصویر مشخص است. اما دقت این الگوریتم، دقت بالایی نیست.

```

1 import pandas
2 from sklearn import model_selection
3 from sklearn.ensemble import AdaBoostClassifier
4
5 model = AdaBoostClassifier(n_estimators=1000, random_state=7)
6 model.fit(df_contact, df_y.ravel())
7 yhat = model.predict(df_contact_test)

```

```

1 from sklearn.metrics import accuracy_score
2 df_y_test = np.array(df_y_test)
3 accuracy_score(df_y_test, yhat)

```

0.522

سومین الگوریتم بسیار مهم در حوزه ensemble، الگوریتم random forest است. این الگوریتم زمان اجرای زیادی دارد. با تست کردن مقادیر مختلف برای تعداد پیش بینی کننده ها، از 100 تا 10000، به این نتیجه رسیدیم که اگر تعداد پیش بینی کننده ها برابر 5000 باشد، بهترین نتیجه با صرف زمانی کمتر بدست می آید.

```

1 import pandas
2 from sklearn import model_selection
3 from sklearn.ensemble import RandomForestClassifier
4
5 model = RandomForestClassifier(n_estimators=10000, max_features=128)
6 model.fit(df_contact, df_y.ravel())
7 yhat = model.predict(df_contact_test)

```

```

1 from sklearn.metrics import accuracy_score
2 df_y_test_ = df_y_test
3 accuracy_score(df_y_test_, yhat)

```

0.726

اما اصلی ترین بخش برای الگوریتم های ensemble، الگوریتم هایی است که از آن بعنوان voting یاد میشود. این الگوریتم ها از نظر چند classifier استفاده میکنند. من در این بخش، از دو نوع الگوریتم voting استفاده خواهم کرد. اولین مدل، از classifier های logisticRegression، DecisionTree و SVC استفاده میکند.

دومین مدل، از classifier های DecisionTree، KNeighbors و SVC استفاده میکند. در دومین مدل، برای هر Classifier یک وزن تعیین کرده ام. این وزن‌ها با تکرار چندباره الگوریتم مشخص شده اند، طوری که هر طبقه بندی که نتیجه بهتری بدهد، همان طبقه بندی، وزن بهتری دارد.

```
1 from sklearn import model_selection
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.svm import SVC
5 from sklearn.ensemble import VotingClassifier
6
7 estimators = []
8 model1 = LogisticRegression()
9 estimators.append(('logistic', model1))
10 model2 = DecisionTreeClassifier()
11 estimators.append(('cart', model2))
12 model3 = SVC()
13 estimators.append(('svm', model3))
14
15 ensemble = VotingClassifier(estimators)
16 model.fit(df_contact, df_y.ravel())
17 yhat = model.predict(df_contact_test)
```

```
1 from sklearn.metrics import accuracy_score
2 accuracy_score(df_y_test, yhat)
```

0.522

```

1 from sklearn import datasets
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.svm import SVC
5 from itertools import product
6 from sklearn.ensemble import VotingClassifier
7
8 clf1 = DecisionTreeClassifier(max_depth=4)
9 clf2 = KNeighborsClassifier(n_neighbors=7)
10 clf3 = SVC(kernel='rbf', probability=True)
11 eclf = VotingClassifier(estimators=[('dt', clf1), ('knn', clf2), ('svc', clf3)], voting='soft', weights=[1, 3, 5])
12
13 clf1.fit(df_contact, df_y.ravel())
14 yhat_1 = clf1.predict(df_contact_test)
15 clf2.fit(df_contact, df_y.ravel())
16 yhat_2 = clf2.predict(df_contact_test)
17 clf3.fit(df_contact, df_y.ravel())
18 yhat_3 = clf3.predict(df_contact_test)
19 eclf.fit(df_contact, df_y.ravel())
20 yhat_4 = eclf.predict(df_contact_test)

```

```

1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_3))
6 print(accuracy_score(df_y_test, yhat_4))

```

```

0.523
0.682
0.732
0.735

```

پس تا اینجا کار، بیشترین دقت برای الگوریتم SVC با الگوریتم کرنل rbf ثبت شد که نتیجه ی آن هم در الگوریتم ensemble مشهود است.

الگوریتم boosting دیگری که میتوانیم برای classification استفاده کنیم، الگوریتم Hist Gradient Tree Boosting میباشد. این الگوریتم هم با توجه به اینکه زمان اجرای پایینی دارد و دقت بالایی ارائه میدهد، مدل مناسبی ایجاد میکند.

```

1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.ensemble import HistGradientBoostingClassifier
3
4 clf = HistGradientBoostingClassifier(loss='binary_crossentropy', learning_rate=0.2).fit(df_contact, df_y.ravel())
5 clf.score(df_contact_test, df_y_test)

```

```

0.734

```

در آخرین مرحله تصمیم گرفتم یک شبکه عصبی یک لایه انتخاب کنم و آنرا بعنوان یکی از classifier های الگوریتم voting استفاده کنم. بهترین دقت تا اینجا کار از شبکه عصبی بدست آمد. همانطور که مشخص است از adam بعنوان optimizer استفاده شده و نرخ یادگیری با سعی و خطا بدست آمده است.

```
1 from sklearn.neural_network import MLPClassifier
2 clf = MLPClassifier(solver='adam', alpha=1e-2, hidden_layer_sizes=(128,), random_state=1)
3 clf.fit(df_contact, df_y.ravel())
4 yhat = clf.predict(df_contact_test)
5 accuracy_score(df_y_test, yhat)
```

0.774

استفاده در مدل voting:

```
1 from sklearn import datasets
2 from sklearn.svm import SVC
3 from sklearn.ensemble import VotingClassifier
4
5 clf1 = MLPClassifier(solver='adam', alpha=1e-2, hidden_layer_sizes=(128,), random_state=1)
6 clf2 = KNeighborsClassifier(n_neighbors=7)
7 clf3 = SVC(kernel='rbf', probability=True)
8 eclf = VotingClassifier(estimators=[('mlp', clf1), ('knn', clf2), ('svc', clf3)], voting='soft', weights=[5, 3, 4])
9
10 clf1.fit(df_contact, df_y.ravel())
11 yhat_1 = clf1.predict(df_contact_test)
12 clf2.fit(df_contact, df_y.ravel())
13 yhat_2 = clf2.predict(df_contact_test)
14 clf3.fit(df_contact, df_y.ravel())
15 yhat_3 = clf3.predict(df_contact_test)
16 eclf.fit(df_contact, df_y.ravel())
17 yhat_4 = eclf.predict(df_contact_test)
```

```
1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_3))
6 print(accuracy_score(df_y_test, yhat_4))
```

0.774
0.682
0.732
0.786

در این مدل voting که از شبکه عصبی، KNeighbors و SVC استفاده کردیم، دقت به 78.6 درصد رسید.

اما مسئله در اینجا ختم نشد. من تلاش زیادی کردم که بتوانم درصد دقت را از این بالاتر ببرم. بنابراین نحوه تشکیل آرایه های تست و آموزش را مورد بازبینی قرار دادم. به این منظور، داده ی مجموعه اول که متعلق به نفر اول است را از داده ی مجموعه دوم که متعلق به نفر دوم است کم کردم. یعنی داده ی 512 تایی از مجموعه اول را از داده ی 512 تایی مجموعه دوم کم کردم و یک بردار 512 تایی بدست آمد که حاصل تفریق آن دو بود. این شکل از دادگان دو مزیت دارد.

اولا مفهوم distance در آن دیده میشود. پس مدل با داشتن فاصله دو تصویر، سعی میکند الگوریتمی برای یافتن پاسخ بیابد. ثانيا تعداد ویژگی ها نسبت به حالت قبل نصف شده است. این به بهبود نتیجه کمک فراوانی خواهد کرد.

```
1 fd_concat_distance = np.zeros((df_contact.shape[0], int(df_contact.shape[1]/2)))
2 c = 0
3 for row in df_contact:
4     for element in range(row.shape[0]):
5         if element < 512:
6             fd_concat_distance[c][element] = row[element] - row[element + 512]
7         c = c + 1
8
9 fd_concat_distance.shape
```

(2200, 512)

```
1 fd_concat_distance_test = np.zeros((df_contact_test.shape[0], int(df_contact_test.shape[1]/2)))
2 c = 0
3 for row in df_contact_test:
4     for element in range(row.shape[0]):
5         if element < 512:
6             fd_concat_distance_test[c][element] = row[element] - row[element + 512]
7         c = c + 1
8
9 fd_concat_distance_test.shape
```

(1000, 512)

در این راستا، میتوانیم از روشهای دیگری نیز استفاده کنیم. مثلا میتوانیم سه ویژگی به داده هایمان اضافه کنیم که فاصله اقلیدسی، فاصله منهتن و فاصله کسینوسی است. در تصاویر زیر، پیاده سازی ان انجام شده است.

```
1 # =====
2 from scipy.spatial import distance
3 num_row = 0
4 new_feature = np.zeros((df_contact.shape[0], 1))
5 for row in df_contact:
6     a = row[: 512]
7     b = row[512 :]
8     dst = distance.euclidean(a, b)
9     new_feature[num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance = np.append(fd_concat_distance, new_feature, axis = 1)
13 fd_concat_distance.shape
```

(2200, 513)

```
[103] 1 # =====
2 from scipy.spatial import distance
3 num_row = 0
4 new_feature = np.zeros((df_contact_test.shape[0], 1))
5 for row in df_contact_test:
6     a = row[: 512]
7     b = row[512 :]
8     dst = distance.euclidean(a, b)
9     new_feature[num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance_test = np.append(fd_concat_distance_test, new_feature, axis = 1)
13 fd_concat_distance_test.shape
```

(1000, 513)

```
1 #####
2 from sklearn.metrics.pairwise import manhattan_distances
3 num_row = 0
4 new_feature = np.zeros([df_contact.shape[0], 1])
5 for row in df_contact:
6     a = row[ : 512].reshape(1, -1)
7     b = row[512 :].reshape(1, -1)
8     dst = manhattan_distances(a, b)
9     new_feature[num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance = np.append(fd_concat_distance, new_feature, axis = 1)
13 fd_concat_distance.shape

(2200, 513)
```

```
[125] 1 #####
2 from sklearn.metrics.pairwise import manhattan_distances
3 num_row = 0
4 new_feature = np.zeros([df_contact_test.shape[0], 1])
5 for row in df_contact_test:
6     a = row[ : 512].reshape(1, -1)
7     b = row[512 :].reshape(1, -1)
8     dst = manhattan_distances(a, b)
9     new_feature[num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance_test = np.append(fd_concat_distance_test, new_feature, axis = 1)
13 fd_concat_distance_test.shape

(1000, 513)
```

```
1 #####
2 from sklearn.metrics.pairwise import cosine_similarity
3 num_row = 0
4 new_feature = np.zeros([df_contact.shape[0], 1])
5 for row in df_contact:
6     a = row[ : 512].reshape(1, -1)
7     b = row[512 :].reshape(1, -1)
8     dst = manhattan_distances(a, b)
9     new_feature[num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance = np.append(fd_concat_distance, new_feature, axis = 1)
13 fd_concat_distance.shape

(2200, 514)
```

```
[127] 1 #####
2 from sklearn.metrics.pairwise import cosine_similarity
3 num_row = 0
4 new_feature = np.zeros([df_contact_test.shape[0], 1])
5 for row in df_contact_test:
6     a = row[ : 512].reshape(1, -1)
7     b = row[512 :].reshape(1, -1)
8     dst = manhattan_distances(a, b)
9     new_feature[num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance_test = np.append(fd_concat_distance_test, new_feature, axis = 1)
13 fd_concat_distance_test.shape

(1000, 514)
```

البته انجام این تغییرات تأثیری در خروجی نداشت. من جستجو کردم و فهمیدم که تابع Euclidean که در کتابخانه distance تعریف شده، دلیل این کم بودن دقت است. بنابراین از کتابخانه sklearn استفاده کردم و نتیجه به بهترین حالت رسید.

```

1 # #####
2 from sklearn.metrics.pairwise import euclidean_distances
3 num_row = 0
4 new_feature = np.zeros([df_contact.shape[0], 1])
5 for row in df_contact:
6     a = row[ : 512].reshape(1, -1)
7     b = row[512 :].reshape(1, -1)
8     dst = euclidean_distances(a, b)
9     new_feature [num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance = np.append(fd_concat_distance, new_feature, axis = 1)
13 fd_concat_distance.shape

```

(2200, 513)

```

] 1 # #####
2 from sklearn.metrics.pairwise import euclidean_distances
3 num_row = 0
4 new_feature = np.zeros([df_contact_test.shape[0], 1])
5 for row in df_contact_test:
6     a = row[ : 512].reshape(1, -1)
7     b = row[512 :].reshape(1, -1)
8     dst = euclidean_distances(a, b)
9     new_feature [num_row] = dst
10    num_row = num_row + 1
11
12 fd_concat_distance_test = np.append(fd_concat_distance_test, new_feature, axis = 1)
13 fd_concat_distance_test.shape

```

(1000, 513)

یک مرحله دیگر که حتما باید بر روی دادگان انجام شود، بحث استاندارد سازی است. این فرایند را توسط تابع standardScaler انجام می‌دهیم.

▼ Standard

✓
0s

```

1 from sklearn.preprocessing import StandardScaler
2 sc = StandardScaler()
3 fd_concat_distance = sc.fit_transform(fd_concat_distance)
4 fd_concat_distance_test = sc.transform (fd_concat_distance_test)

```

حالا به نتایج تست مدل‌های مختلف روی این دادگان می‌پردازیم.

```

1 from sklearn.neural_network import MLPClassifier
2 clf = MLPClassifier(solver='adam', alpha=1e-1, hidden_layer_sizes=(300,), random_state=1)
3 clf.fit(fd_concat_distance, df_y.ravel())
4 yhat = clf.predict(fd_concat_distance_test)
5 accuracy_score(df_y_test, yhat)

0.88

```

میبینیم که مدل شبکه عصبی با نرخ یادگیری 0.1 و تعداد 300 نورون، 88 درصد خروجی میدهد که بهترین دقت در بین مدلهایی است که تاکنون ران شده است.

حالا از مدل voting استفاده میکنیم که مدل ensemble است.

```

1 from sklearn import datasets
2 from sklearn.svm import SVC
3 from sklearn.ensemble import VotingClassifier
4
5 clf1 = MLPClassifier(solver='adam', alpha=1e-1, hidden_layer_sizes=(300,), random_state=1)
6 clf2 = SVC(kernel='rbf', probability=True)
7 eclf = VotingClassifier(estimators=[('mlp', clf1), ('svc', clf2)], voting='soft', weights=[5, 4])
8
9 clf1.fit(fd_concat_distance, df_y.ravel())
10 yhat_1 = clf1.predict(fd_concat_distance_test)
11 clf2.fit(fd_concat_distance, df_y.ravel())
12 yhat_2 = clf2.predict(fd_concat_distance_test)
13 eclf.fit(fd_concat_distance, df_y.ravel())
14 yhat_4 = eclf.predict(fd_concat_distance_test)

1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_4))

0.88
0.859
0.88

```

برای بهبود کارایی الگوریتم یک ایده ای که میتوان مطرح کرد، اضافه کردن داده های آموزشی است. برای این منظور، فرض میکنیم تصاویر آموزشی تاریک تر باشند، یا روشن تر باشند. بنابراین میتوان به ازای هر تصویر، دو تصویر جدید بوجود آورد که تاریکتر یا روشنتر از تصویر اولیه باشند. برای این منظور، بخش darker و lighter کد زده شده اند. من بعنوان نمونه، همه پیکسلهای این تصاویر را 0.1 تیره تر یا روشنتر کرده ام.

▼ Darker

✓
26s



```
1 print(df_norm_.shape)
2 num_row = 0
3 for row in fd_concat_distance:
4     darker = np.empty([0,0])
5     for elmnt in range(row.shape[0]-3):
6         temp = row[elmnt] + 0.05
7         if temp > 1:
8             temp = 1
9         darker = np.append(darker, temp)
10
11 darker = np.append(darker, row[512])
12 darker = np.append(darker, row[513])
13 darker = np.append(darker, row[514])
14 darker = np.expand_dims(darker, axis=1)
15 df_norm_ = np.concatenate((df_norm_, darker.T), axis=0)
16 y_ = np.append(y_, y_[num_row])
17 num_row = num_row + 1
18 print(df_norm_.shape)
```

☞ (2200, 515)
(4400, 515)

▼ Lighter

✓
34s



```
1 print(df_norm_.shape)
2 num_row = 0
3 for row in fd_concat_distance:
4     lighter = np.empty([0,0])
5     for elmnt in range(row.shape[0] - 3):
6         temp = row[elmnt] - 0.05
7         if temp < 0:
8             temp = 0
9         lighter = np.append(lighter, temp)
10
11 lighter = np.append(lighter, row[512])
12 lighter = np.append(lighter, row[513])
13 lighter = np.append(lighter, row[514])
14 lighter = np.expand_dims(lighter, axis=1)
15 df_norm_ = np.concatenate((df_norm_, lighter.T), axis=0)
16 y_ = np.append(y_, y_[num_row])
17 num_row = num_row + 1
18 print(df_norm_.shape)
```

☞ (4400, 515)
(6600, 515)

پس از آن، الگوی rotate با 180 درجه زاویه را پیاده سازی کردم ولی نتیجه چندان جالب نشد. تصویر آن بصورت زیر است:

▼ Rotate



```
1 print(df_norm_.shape)
2 num_row = 0
3 for row in fd_concat_distance:
4
5     swaper = np.empty([0,0])
6     for elmnt in range(row.shape[0]):
7         temp = row[row.shape[0] - elmnt - 1]
8         swaper = np.append(swaper, temp)
9
10    swaper = np.expand_dims(swaper, axis=1)
11    df_norm_ = np.concatenate((df_norm_, swaper.T), axis=0)
12    df_y = np.append(df_y, df_y[num_row])
13    num_row = num_row + 1
14 print(df_norm_.shape)
```

اما فیلتر salt and pepper مناسبی برای تصاویر است. در این فیلتر، بصورت رندم بعضی از پیکسلها را سیاه و بعضی را سفید میکنیم. پس یعنی بعضی از پیکسلها 0 و بعضی دیگر 1 خواهند شد. در این کد 32 پیکسل بصورت رندم برای هر بردار 512 تایی انتخاب میشود.

```

1 print(df_norm_.shape)
2 num_row = 0
3 for row in fd_concat_distance:
4
5     salt_pepper = np.empty([0,0])
6     randm_dark = np.random.randint(512, size=32)
7     randm_light = np.random.randint(512, size=32)
8
9     for elmnt in range(row.shape[0]):
10         if elmnt in randm_dark and elmnt not in [512, 513, 514]:
11             temp = df_norm_.min()
12             salt_pepper = np.append(salt_pepper, temp)
13
14         elif elmnt in randm_light and elmnt not in [512, 513, 514]:
15             temp = df_norm_.max()
16             salt_pepper = np.append(salt_pepper, temp)
17
18         else:
19             temp = row[elmnt]
20             salt_pepper = np.append(salt_pepper, temp)
21
22     salt_pepper = np.expand_dims(salt_pepper, axis=1)
23     df_norm_ = np.concatenate((df_norm_, salt_pepper.T), axis=0)
24     y_ = np.append(y_, y_[num_row])
25     num_row = num_row + 1
26
27 print(df_norm_.shape)

```

```

↳ (2200, 515)
   (4400, 515)

```

اگر هرکدام از فیلترها را بصورت جداگانه پیاده کنیم، نتایج متفاوتی خواهیم دید. اما برای مثال، پس از پیاده سازی salt and pepper نتیجه ی زیر بدست می آید.


```
Test

1 from sklearn import datasets
2 from sklearn.svm import SVC
3 from sklearn.ensemble import VotingClassifier
4
5 clf1 = MLPClassifier(solver='adam', alpha=1e-1, hidden_layer_sizes=(300,), random_state=1)
6 clf2 = SVC(kernel='rbf', probability=True)
7 eclf = VotingClassifier(estimators=[('mlp', clf1), ('svc', clf2)], voting='soft', weights=[5, 4])
8
9 clf1.fit(df_norm_, y_.ravel())
10 yhat_1 = clf1.predict(fd_concat_distance_test)
11 clf2.fit(df_norm_, y_.ravel())
12 yhat_2 = clf2.predict(fd_concat_distance_test)
13 eclf.fit(df_norm_, y_.ravel())
14 yhat_4 = eclf.predict(fd_concat_distance_test)

[284] 1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_4))

0.867
0.839
0.871
```

برای مثال دیگر از درخت تصمیم استفاده میکنیم. مشاهده میکنیم که دقت آن با توجه به پردازش های صورت گرفته بشدت بالا رفته است.

```
1 from sklearn.tree import DecisionTreeClassifier
2 clf = DecisionTreeClassifier(random_state=0, criterion='entropy', class_weight='balanced')
3 clf.fit(fd_concat_distance, df_y.ravel())
4 yhat_dt = clf.predict(fd_concat_distance_test)
5 print(accuracy_score(df_y_test, yhat_dt))

0.808
```

```

1 import numpy as np
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.ensemble import RandomForestClassifier, VotingClassifier
5 from sklearn.svm import SVC
6
7 clf1 = DecisionTreeClassifier(random_state=0, criterion='entropy', class_weight='balanced')
8 clf2 = RandomForestClassifier(n_estimators=500, random_state=1)
9 clf3 = GaussianNB()
10 clf1.fit(fd_concat_distance, df_y.ravel())
11 clf2.fit(fd_concat_distance, df_y.ravel())
12 clf3.fit(fd_concat_distance, df_y.ravel())
13
14 yhat_1 = clf1.predict(fd_concat_distance_test)
15 yhat_2 = clf2.predict(fd_concat_distance_test)
16 yhat_3 = clf3.predict(fd_concat_distance_test)

```

```

90] 1 from sklearn.metrics import accuracy_score
2 print(accuracy_score(df_y_test, yhat_1))
3 print(accuracy_score(df_y_test, yhat_2))
4 print(accuracy_score(df_y_test, yhat_3))
5
6 eclf1 = VotingClassifier(estimators=[('DecisionTreeClassifier', clf1), ('RandomForestClassifier', clf2), ('GaussianNB', clf3)],
7                             voting='soft', weights=[5, 5, 4])
8 eclf1 = eclf1.fit(fd_concat_distance, df_y.ravel())
9 yhat_res = eclf1.predict(fd_concat_distance_test)
10 print(accuracy_score(df_y_test, yhat_1))

```

```

0.808
0.861
0.8
0.808

```

در مرحله بعدی میتوانیم از لایه pooling در شبکه CNN الهام بگیریم. برای این منظور، هر چهار پیکسل کنار هم را در نظر بگیرید که چهار خانه مجاور در آرایه میباشد. فرض کنیم اندیس آنها بین i تا $i+3$ است. از این چهار عدد میانگین میگیریم و آن را در خانه i ام بردار جدیدی قرار میدهیم. سپس i را یکی زیاد میکنیم تا به انتهای آرایه برسیم. این بردارها را به انتهای آرایه داده های آموزشی اضافه میکنیم. کد آن بصورت زیر است و نتیجه متعاقبا آورده شده است.

```

1 print(df_norm.shape)
2 num_row = 0
3 for row in fd_concat_distance:
4     mean_row = np.empty([0,0])
5     elmnt = 0
6     while elmnt < df_norm.shape[1]:
7         if elmnt + 3 < 512:
8             li = [row[elmnt], row[elmnt+1], row[elmnt+2], row[elmnt+3]]
9             temp = sum(li) / 4
10            mean_row = np.append(mean_row, temp)
11            elmnt = elmnt + 1
12        else:
13            temp = row[elmnt]
14            mean_row = np.append(mean_row, temp)
15            elmnt = elmnt + 1
16
17    mean_row = np.expand_dims(mean_row, axis=1)
18    df_norm = np.concatenate((df_norm, mean_row.T), axis=0)
19    y_ = np.append(y_, y_[num_row])
20    num_row = num_row + 1
21 print(df_norm.shape)

```

```

(2200, 515)
(4400, 515)

```

```

1 from sklearn import datasets
2 from sklearn.svm import SVC
3 from sklearn.ensemble import VotingClassifier
4
5 clf1 = MLPClassifier(solver='adam', alpha=1e-1, hidden_layer_sizes=(300,), random_state=1)
6 clf2 = SVC(kernel='rbf', probability=True)
7 eclf = VotingClassifier(estimators=[('mlp', clf1), ('svc', clf2)], voting='soft', weights=[5, 4])
8
9 clf1.fit(df_norm, y_.ravel())
10 yhat_1 = clf1.predict(df_norm_test)
11 clf2.fit(df_norm, y_.ravel())
12 yhat_2 = clf2.predict(df_norm_test)
13 eclf.fit(df_norm, y_.ravel())
14 yhat_4 = eclf.predict(df_norm_test)

```

```

1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_4))

```

```

0.865
0.869
0.872

```

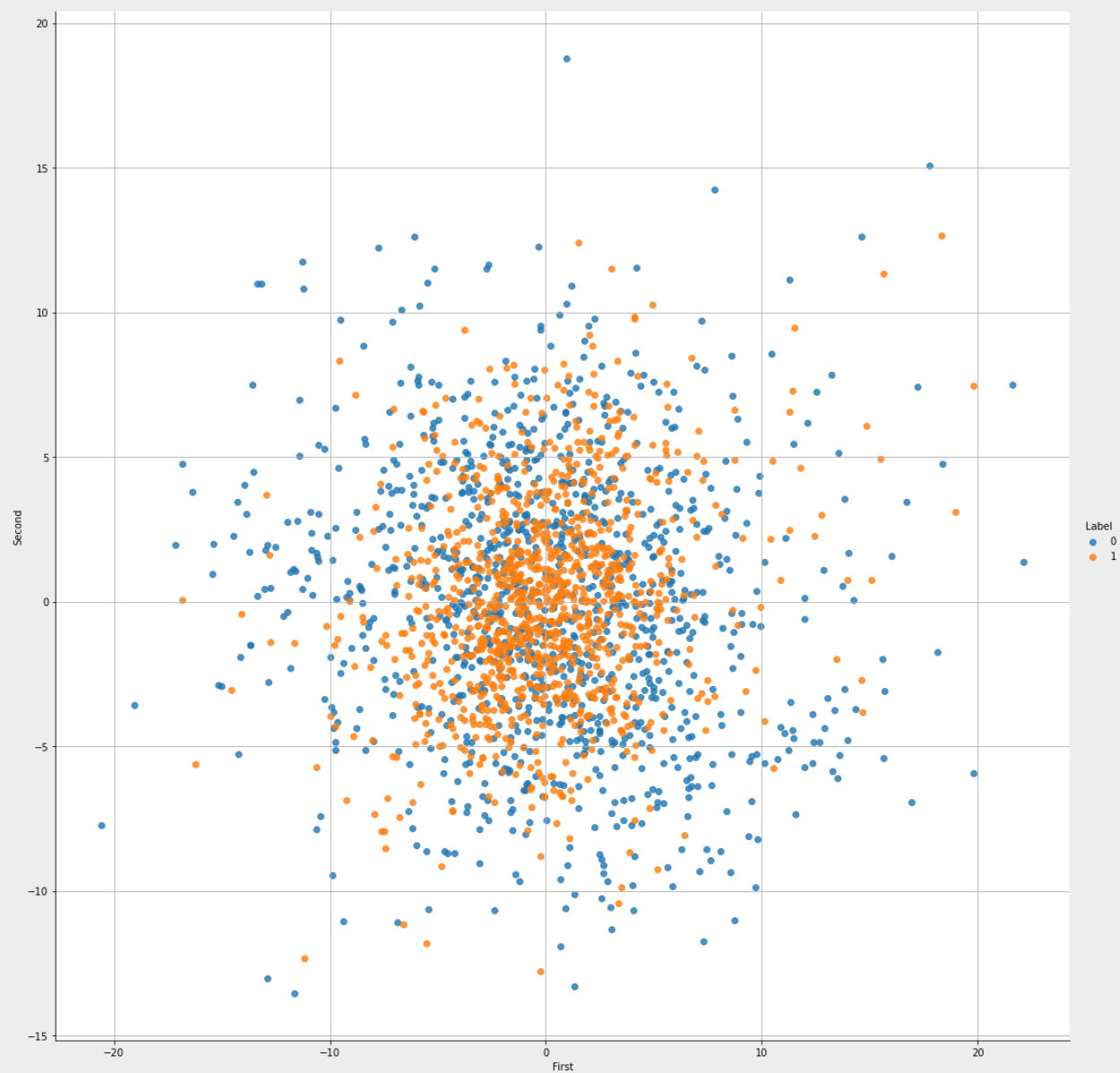
در این گام، تصمیم گرفتیم که به نحوی تعداد فیچرها را کاهش دهیم. برای اینکار از PCA استفاده کردیم. کد آن بصورت زیر است.

```
1 from sklearn.decomposition import PCA
2 from sklearn.decomposition import KernelPCA
3 pca = PCA(n_components= 256, random_state=7, svd_solver = 'arpack')
4 X_PCA = pca.fit_transform(df_norm)
5 X_PCA = pd.DataFrame(data = X_PCA)
6
7 X_PCA_test = pca.fit_transform(df_norm_test)
8 X_PCA_test = pd.DataFrame(data = X_PCA_test)
```

```
1 X_PCA.shape
```

```
(2200, 256)
```

تصویر زیر، نموداری از دو فیچر اول داده های انتقال داده شده است.



مشاهده میشود که تعداد فیچرها را برابر 256 یعنی نصف حالت قبل قرار دادیم. با این توضیح، آن را تست میکنیم. اما متأسفانه الگوریتم بسیار بد عمل میکند.

```

1 from sklearn import datasets
2 from sklearn.svm import SVC
3 from sklearn.ensemble import VotingClassifier
4
5 clf1 = MLPClassifier(solver='adam', alpha=1e-4, hidden_layer_sizes=(128,), random_state=1)
6 clf2 = SVC(kernel='rbf', probability=True)
7 eclf = VotingClassifier(estimators=[('mlp', clf1), ('svc', clf2)], voting='soft', weights=[5, 4])
8
9 clf1.fit(X_PCA, y_.ravel())
10 yhat_1 = clf1.predict(X_PCA_test)
11 clf2.fit(X_PCA, y_.ravel())
12 yhat_2 = clf2.predict(X_PCA_test)
13 eclf.fit(X_PCA, y_.ravel())
14 yhat_4 = eclf.predict(X_PCA_test)

```

```

1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_4))

```

```

0.538
0.703
0.566

```

یک راه دیگر که برای کم کردن تعداد فیچرها میشد استفاده کرد این بود که دو پیکسل از تصویر اول را با دو پیکسل از تصویر دوم که قرار بود با هم مقایسه شوند در نظر بگیریم. اکنون دو بردار دو عضوی داریم. سپس با فاصله اقلیدسی آنها فاصله آنها را حساب کنیم. این کار را ادامه میدهیم و سپس یک آرایه ای با نصف فیچرهای قبلی خواهیم داشت. کد آن بصورت زیر است. این کار روی داده های تست نیز انجام میشود.

```

1 from sklearn.metrics.pairwise import euclidean_distances
2 df_norm_ = np.zeros((df_contact.shape[0], int(df_contact.shape[1]/4)))
3 c = 0
4 for row in df_contact:
5     p = 0
6     it = 0
7     while p < (row.shape[0]-3)/2:
8         a = np.array((row[p], row[p+1])).reshape(1, -1)
9         b = np.array((row[p+512], row[p+513])).reshape(1, -1)
10        dst = euclidean_distances(a, b)
11        df_norm_[c][it] = dst
12        it = it + 1
13        p = p + 2
14
15    c = c + 1
16
17 df_norm_.shape

```

(2200, 256)

اما این روش کاهش ویژگی هم به خوبی موفق نبود.

```

1 from sklearn import datasets
2 from sklearn.svm import SVC
3 from sklearn.ensemble import VotingClassifier
4
5 clf1 = MLPClassifier(solver='adam', alpha=1e-2, hidden_layer_sizes=(128,), random_state=1)
6 clf2 = SVC(kernel='rbf', probability=True)
7 eclf = VotingClassifier(estimators=[('mlp', clf1), ('svc', clf2)], voting='soft', weights=[5, 5])
8
9 clf1.fit(df_norm_, y_.ravel())
10 yhat_1 = clf1.predict(df_norm_test_)
11 clf2.fit(df_norm_, y_.ravel())
12 yhat_2 = clf2.predict(df_norm_test_)
13 eclf.fit(df_norm_, y_.ravel())
14 yhat_4 = eclf.predict(df_norm_test_)

```

```

1 from sklearn.metrics import accuracy_score
2
3 print(accuracy_score(df_y_test, yhat_1))
4 print(accuracy_score(df_y_test, yhat_2))
5 print(accuracy_score(df_y_test, yhat_4))

```

0.704
0.838
0.768

بعنوان جمع بندی، بهترین دقت 88 درصد ثبت شد.

کامیار نصیری