

הפקולטה להנדסה ע"ש איבי ואלדר פליישמן  
המחלקה להנדסת תעשייה



## מבוא ללמידת מכונה – תש"ף (2019-20) פרויקט מסכם

**מרצה:** ד"ר אייל קולמן

**מתרגלת:** שי עובר

**מגישים:** אולגה סולדטנקו – 342480308

יליזבטה גולדברג – 337762462

קבוצה 03

## **תקציר מנהלים:**

בפרויקט זה עסקנו בבעיית Binary Classification בה סיווגנו רשומות לשתי קטגוריות על סמך מספר פיצ'רים בדאטה סט. ראשית חקרנו את הנתונים – אופי התפלגות הפיצ'רים, התנהגות קורלטיבית בין הפיצ'רים, נתונים סטטיסטיים. כמו כן בנינו מספר גרפים שונים המשמשים לוויזואליזציה של נתונים ושיפור הבנתם. לאחר מכן ביצענו שלב של עיבוד מקדים – הסרת נתונים חריגים, נרמול נתונים, השלמת נתונים חסרים, הקטנת ממדיות של הבעיה ובניית פיצ'רים חדשים. כל זה על מנת להכין את הדאטה לשלב של הרצת מודלים. בנינו 4 מודלים – 2 בסיסיים (Logistic Regression, KNN) ו-2 מתקדמים (ANN, Adaptive Boosting). על חלק מהדאטה הרצנו את המודלים (train) ועל חלק שני (validation) ביצענו פרדיקציה וערכנו טיב המודלים. קיבלנו כי חודל Logistic Regression בעל מדדים הכי טובים ולכן השתמשנו בו לסיווג הנתונים הסופי.

## **חלק ראשון – אקספלורציה**

בחלק זה ביצענו בדיקת הנתונים. בעיה בעלת 22161 דגימות ו-25 פיצ'רים. מהם 6 קטגוריאליים והשער מספריים. משתנים קטגוריאליים – מופעים בפיצ'רים מס' 5, 6, 9, 14, 18, 19. פיצ'ר 13 על פניו נראה קטגוריאלי אך בפועל הוא בשלמים. פיצ'ר 14 גם הפיך לפיצ'ר נומרי בקלות על ידי הסרת 'mm' מכל הדגימות.

ראינו שאין העתקים של אותה דגימה.

בנינו את Heat map וראינו שקיימת קורלציה חיובית גבוהה (יותר מ-0.85) בין פיצ'רים מסוימים: [0,1,2], [7,16], [7,17], [8,17], [16,17]. לא קיימים פיצ'רים בעלי קורלציה שלילית גבוהה. בנוסף לחלק מהפיצ'רים בעלי קורלציה גבוהה בנינו scatter plots כדי לראות את האופן תלות ביניהם.

עבור הפיצ'רים הנומריים בוצעו היסטוגרמות ו-Boxplots לפי כל קלס כדי להבין איך הערכים מתפלגים. לפי גרפים ניתן לראות כי הנתונים לא מנורמלים - עבור פיצ'רים נדרש לשנות את הסקלה של הנתונים (טווחי ערכים שונים בפיצ'רים שונים) כמו כן קיימים פיצ'רים עם נתונים חריגים.

עבור פיצ'רים קטגוריאליים בנינו גרפים שכיחויות.

דוגמאות של הגרפים ניתן לראות בנספח 1.

## **חלק שני – עיבוד מקדים:**

- **נתונים חריגים** – באקספלורציה ראינו שרוב הפיצ'רים הנומריים הם בעלי פילוג שיש בו מרכז וההסתברות יורדת ככל שמתרחקים ממנו וקיימים נתונים חריגים, לכן בחרנו להסיר כל הדגימות אשר עולות על 3 סטיות התקן.

- נרמול נתונים – לנרמול השתמשנו ב z-score normalization. נרמול חשוב לעבודה טובה של רוב האלגוריתמים של למידת מכונה - הם משתמשים במרחקים אוקלידיים.
- טיפול בנתונים חסרים – החלטנו בכלל לא להשתמש בפיצ'רים שחסרים להם יותר מ 5% מהנתונים. כבר בשלב זה הסרנו פיצ'רים 5 ו-15. לכל שאר הפיצ'רים השלמנו דגימות חסרות בממוצה (mean). בפיצ'ר 13 מכיוון שערכיו שלמים הוחלט להסיר את הדגימות הלא ידועות (פחות מ 5% מה dataset).
- הקטנת ממדיות – ממדיות הבעיה לא גדולה מידי, אחרי הורדת פיצ'רים קטגוריאליים נשארו רק 19. זה טוב לגודל הדאטה סט שיש לנו. אין סכנה של curse of dimensionality. בכל זאת ניסינו קצת להקטין את הממדיות לריצה יותר מהירה וחסכון של זיכרון. הורדנו פיצ'רים בעלי קורלציה גבוהה אשר גילינו בשלב האקספלורציה (פיצ'ר 1 ירד).
- בזמן אימון המודלים וניסיונות השתמשנו ב-PCA, לקחנו רק את הפיצ'רים הראשונים אשר ביחד מסבירים את 85% מהנתונים (expected variance), אבל ראינו ששימוש ב-PCA מוריד את הביצועים לכן במודלים סופיים הוחלט לא להשתמש בו.
- בניית פיצ'רים חדשים – הוחלט לבנות פיצ'רים חדשים מפיצ'רים שכבר קיימים. עשינו סכום של פיצ'רים 8 ו-17 ופיצ'רים 0 ו-2 (פיצ'רים בעלי קורלציה גבוהה בניהם). בנוסף השתמשנו ב-KFold Clustering כדי לשייך דגימות לקבוצות שונות. ראינו שהדאטה מתחלק הכי טוב ל 3 קלסטרים (לפי שיטת Elbow plot). הוספנו לדאטה של train שורה עם מספר קלסטר שאליו דגימה השתייכה ('clustering').
- החלת העיבוד המקדים על סט ה-test – בשביל סיווג נכון סט הטסט צריך להיראות כמו סט שעליו אומנו המודלים. בשביל זה עושים נרמול, PCA (אם מחליטים להשתמש בו), מסירים את הפיצ'רים ובונים את הפיצ'רים החדשים (קומבינציות ו-Clustering) בדיוק כמו בסט האימון.

לכל הבדיקות והרצות השתמשנו באובייקט Pipeline שבנינו בעצמנו. הוא מאפשר להריץ גרסאות שונות של הבעיה ללא הרצה ידנית ולשמור את כל הנתונים בטבלאות נוחות. הקוד אפשר לראות בנספח 2.

## **חלק שלישי ורביעי – הרצת המודלים והערכתם:**

- עבור מודלים ראשוניים בחרנו לא להשתמש ב- Naïve Bayes Classifier משום שמודל זה מניח אי-תלות של הפיצ'רים, אשר ככל הנראה לא נכון במקרה שלנו. לכן ממשנו את Logistic Regression. ו- KNN.
- במודלים מתקדמים בחרנו לממש Adaptive Boosting ורשתות נוירונים (ANN) אשר לדעתנו מתאימות הכי טוב למבנה הנתונים שלנו.
- הערכת מודלים התבצעה על ידי השוואת מדד AUC ובדיקת overfitting. על מנת להגדיל יכולת הכללה של המודלים השתמשנו ב-K-Fold Cross Validation כאשר מספר K-Folds האופטימלי הוא 5. הגדנו כי מודל נחשב overfitted כאשר פערי ביצוע על ה-Train ועל ה-Validation עולים מעל 5%. גרפים של K-Fold ניתן לראות בנספח 3.

בנוסף עבור כל מודל חשבנו מדד דיוק (Accuracy) משוקלל. מדד זה צריך להתייחס לעלות סיווג שגוי אשר חמורה פי 5 לדגימה בעלת תיוג "1":  
 $(w_{tp}, w_{fp}, w_{tn}, w_{fn}) = (1, 1, 1, 5)$

$$\frac{w_{tp}TP + w_{tn}TN}{w_{tp}TP + w_{tn}TN + w_{fp}FP + w_{fn}FN}$$

במקרים קצה (כל הנבואות נכונות, כל הנבואות שגויות) מדד זה מקבל ערכים 1 ו-0 בדומה למדד הדיוק רגיל. במקרה ויש יותר סיווגים שגויים של מדדים בעלי סיווג 1 מדד חדש מקבל ארך יותר קטן.

- KNN - למציאת פרמטרים למודל השתמשנו ב-Grid Search.  
 פרמטרים שהתקבלו: 'metric='euclidean', n\_neighbors = 100, weights='uniform'.  
 מדד AUC שהתקבל: 0.84. קיבלנו הבדל של 1.5% בין Train ל-Validation ולכן ניתן להסיג שלא קיים overfitting. מדד דיוק משוקלל - 0.54. מדד זה יצאה נמוך מכיוון שכמות הטעויות מסוג false negative יותר גדולה מ-false positive.
- Logistic Regression - למציאת פרמטרים למודל השתמשנו ב-Grid Search.  
 הפרמטרים שהתקבלו: C=0.5, max\_iter=300, penalty='l1', 'solver='liblinear', tol=0.0001.  
 מדד AUC שהתקבל: 0.862. קיבלנו הבדל של 0.09% בין Train ל-Validation ולכן ניתן להסיג שלא קיים overfitting. מדד דיוק משוקלל - 0.6. מדד זה יצאה נמוך מכיוון שכמות הטעויות מסוג false negative יותר גדולה מ-false positive.
- ANN - למציאת פרמטרים למודל השתמשנו ב-Random Search.  
 הפרמטרים שהתקבלו: max\_iter = 1000, learning\_rate\_init = 0.01, 'hidden\_layer\_sizes = (100,)', batch\_size = 10, activation = 'relu'.  
 מדד AUC שהתקבל: 0.856. קיבלנו הבדל של 2.2% בין Train ל-Validation ולכן ניתן להסיג שלא קיים overfitting. מדד דיוק משוקלל - 0.6. מדד זה יצאה נמוך מכיוון שכמות הטעויות מסוג false negative יותר גדולה מ-false positive.
- Adaptive Boosting - למציאת פרמטרים למודל השתמשנו ב-Grid Search.  
 הפרמטרים שהתקבלו: learning\_rate=0.1, n\_estimators=300, random\_state=0.  
 מדד AUC שהתקבל: 0.857. קיבלנו הבדל של 0.41% בין Train ל-Validation ולכן ניתן להסיג שלא קיים overfitting. מדד דיוק משוקלל - 0.6. מדד זה יצאה נמוך מכיוון שכמות הטעויות מסוג false negative יותר גדולה מ-false positive.
- Confusion Matrix - בעזרתה ניתן לעריך איכות הפלט של המודל.

לדוגמה מופיע מטריצה עבור מודל Logistic Regression בנספח 3.  
לאלן הסבר על מרכיבי המטריצה:

TP 379 תצפיות היו 1 ומודל חזה 1	FP 159 תצפיות היו 0 ומודל חזה 1
FN 517 תצפיות היו 1 ומודל חזה 0	TN 3918 תצפיות היו 0 ומודל חזה 0

מהמטריצה ניתן להסיג כביצועי המודל טובים אך יש מספר די גדול של טעויות  
מסוג false negative.

### **חלק חמישי – ביצוע פרדיקציה:**

בחרנו להשתמש במודל Logistic Regression לפרדיקציה מכיוון שהיא בעלת AUC הכי גבוה.

### **סיכום:**

בפרויקט זה בנינו 4 מודלים. ראינו שכל המודלים נותנות תוצאות טובות, אין overfitting. למרות זאת מדד דיוק משוקלל יחסית נמוך, מודלים מחזירים הרבה שגיאות false negative. מכיוון שנדרשנו לבחור מודל לפי ROC החלטנו לא לנסות לשפר מדד זה.

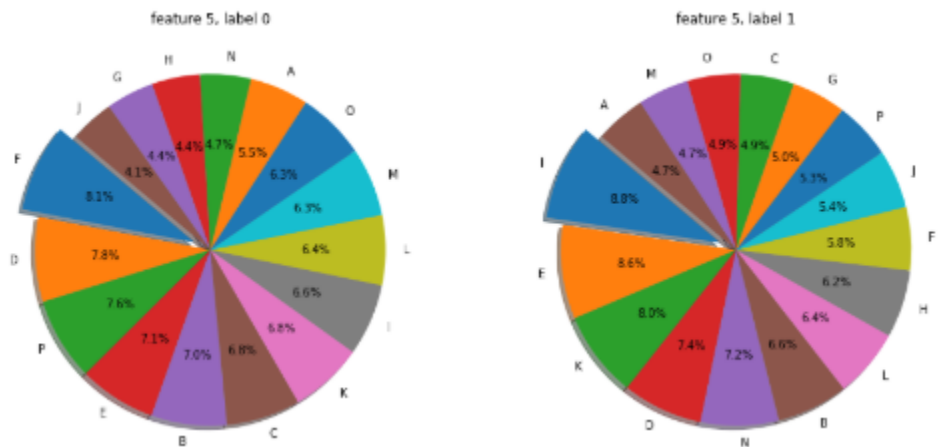
במהלך ניתוח נתונים ראינו כי היה נירש לבצע הרבה פעולות בעיבוד מקדים, אך לא קיבלנו ידע מספק לגבי פיצ'רים ודאטה ולכן ניתן לומר שלא נעשה עיבוד מקדים בצורה המיטבית ולכן תוצאות המודלים ניתנים לשיפור.

קיבלנו 3 מודלים (Adaptive Boosting, Logistic Regression, ANN) בעלי מדדי ביצוע מאוד דומים: AUC קצת יותר טוב ל-Logistic Regression, אך שווה בערך ל-0.86 עבור שלושת המודלים; מדד דיוק משוקלל אומד על 0.6 עבור שלושת המודלים; פערי ביצוע בין test ו-train הקטן ביותר עבור מודל Logistic Regression.

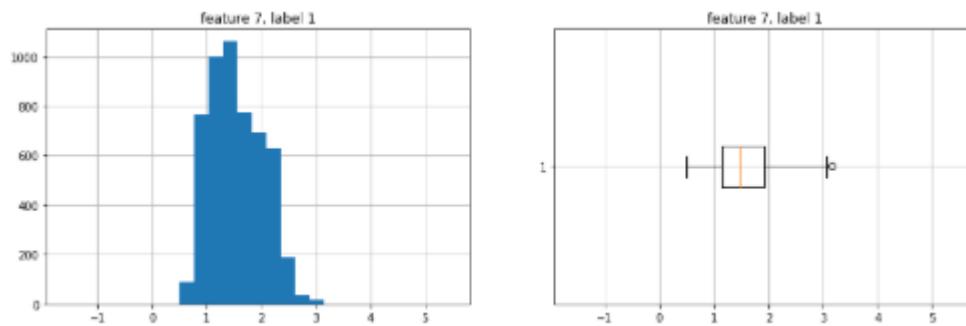
לבסוף בחרנו במודל Logistic Regression כמסווג הכי טוב לפי מטריקת AUC ובשכלול מדדים נוספים.

## נספחים

### נספח 1: אקספלורציה – דוגמאות לגרפים שונים



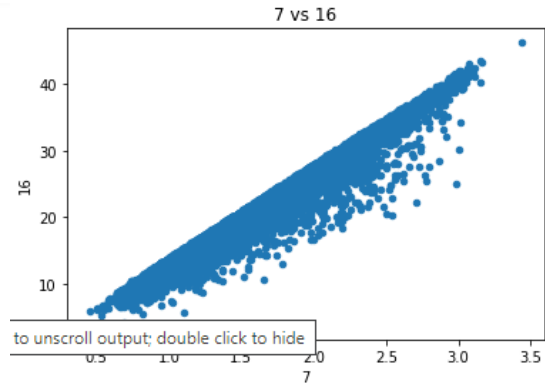
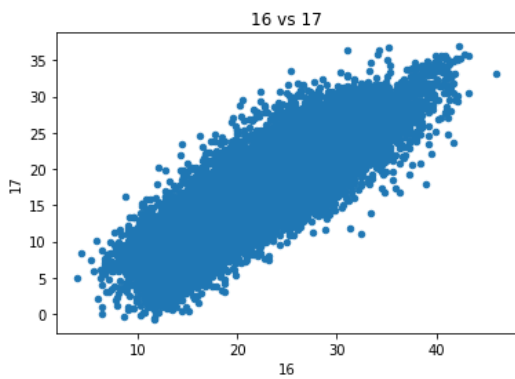
איור 1: pieplot עבור פיצ'ר קטגוריאלי (5)



איור 2: histogram ו-boxplot עבור פיצ'ר נומרי (7)



איור 3: Heatmap עבור סט נתונים train



איור 4: גרף עבור פיצ'רים בעלי קורלציה גבוהה (7 ו-16, 16 ו-17)

## נספח 2: עיבוד מקדים

### אובייקט Pipeline

```

1. class Pipeline:
2.     def __init__(self, normalization = 1):
3.         self.one_hot_spot = 0
4.         self.combine_features = 0
5.         self.drop_features = 0
6.         self.anova = 0
7.         self.normalization = normalization
8.         self.model = 'KNN'
9.         self.outliers = 3
10.        self.pca_n_components = 0
11.        self.combine_features = 0
12.        self.clustering = 0
13.
14.    def set_combine_features(self, features):
15.        temp = copy.copy(self)
16.        temp.combine_features = features
17.        return temp
18.
19.    def set_drop_features(self, features):
20.        temp = copy.copy(self)
21.        temp.drop_features = features
22.        return temp
23.
24.    def set_remove_outliers(self, z):
25.        temp = copy.copy(self)
26.        temp.outliers = z
27.        return temp
28.
29.    def set_clustering(self, n_clusters = 3):
30.        temp = copy.copy(self)
31.        temp.clustering = n_clusters
32.        return temp
33.

```

```

34.     def set_PCA(self, n_components):
35.         temp = copy.copy(self)
36.         temp.pca_n_components = n_components
37.         return temp
38.
39.     def set_model(self, model):
40.         temp = copy.copy(self)
41.         temp.model = model
42.         return temp
43.
44.
45.     def run_pipeline(self, X,y,i,lenth, start):
46.         print('Starting ', i , 'out of ',lenth)
47.
48.         if self.combine_features:
49.             for feature in self.combine_features:
50.                 X = combine_features(X,feature[0],feature[1])
51.                 print('Combining features - done')
52.
53.         if self.drop_features:
54.             X = X.drop(self.drop_features,axis = 1)
55.             print('Dropping features - done')
56.
57.         if self.outliers:
58.             X,y = Outliers(X,y)
59.             print('Outliers - done')
60.
61.         X = Filling_the_missing_data(X)
62.         print('Filling_the_missing_data - done')
63.
64.         if self.clustering:
65.             X = K_means(X)
66.             print('Clustering - done')
67.
68.         if self.normalization:
69.             X, std = Normalization(X)
70.             print('Normalization - done')
71.
72.         if self.pca_n_components:
73.             X, pca_clf = our_PCA(X,self.pca_n_components)
74.             print('PCA - done')
75.
76.         X_train, X_valid, y_train, y_valid = split_test(X,y)
77.
78.         if self.model == 'KNN':
79.             print('Entered KNN')
80.             res_KNN = KNN(X_train,y_train)
81.             self.best_estimator = res_KNN.best_estimator_
82.             self.auc = res_KNN.best_score_
83.
84.         elif self.model == 'LogReg':
85.             print('Entered LogReg')
86.             res_LogReg = LogReg(X_train,y_train)
87.             self.best_estimator = res_LogReg.best_estimator_
88.             self.auc = res_LogReg.best_score_

```



```

89.
90.     elif self.model == 'AdaBoost':
91.         print('EnteredAdaBoost')
92.         res_AdaBoost = AdaBoost(X_train,y_train)
93.         self.best_estimator = res_AdaBoost.best_estimator_
94.         self.auc = res_AdaBoost.best_score_
95.
96.     elif self.model == 'AdaBoost_SVC':
97.         print('EnteredAdaBoost_SVC')
98.         res_AdaBoost_SVC = AdaBoost_SVC(X_train,y_train)
99.         self.best_estimator = res_AdaBoost_SVC.best_estimator_
100.        self.auc = res_AdaBoost_SVC.best_score_
101.
102.        elif self.model == 'ANN':
103.            print('ANN')
104.            res_ANN = ANN(X_train,y_train)
105.            self.best_estimator = res_ANN.best_estimator_
106.            self.auc = res_ANN.best_score_
107.
108.            stop = timeit.default_timer()
109.            print('Time: ', (stop - start)/60)
110.            self.execution_differences,self.auc_train =
            execution_differences(X_train, y_train,X_valid, y_valid,
            self.best_estimator)

```

דוגמה של שימוש ב Pipeline:

```

111.     #defining parameters
112.     Params_combine_features_1 = [[('8','17'),('0','2')],[('0','2')]]
113.     Params_drop_features_1 = [['16'],['12']]
114.     Params_PCA = [0,0.9]
115.
116.     #creating list of pipelines
117.     Pipelines = [Pipeline().set_model('LogReg')]
118.     Pipelines = [P.set_combine_features(Param) for P,Param in
        product(Pipelines,Params_combine_features_1) ]
119.     Pipelines = [P.set_drop_features(Param) for P,Param in
        product(Pipelines,Params_drop_features_1) ]
120.     Pipelines = [P.set_PCA(Param) for P,Param in
        product(Pipelines,Params_PCA) ]
121.
122.     #control time
123.     start = timeit.default_timer()
124.     i=1
125.     lenth = len(Pipelines)
126.
127.     #run the pipelines
128.     for P in Pipelines:
129.         dat = P.run_pipeline(X,y,i,lenth,start)
130.         i+=1
131.         Results_df = Results_df.append(P.__dict__, ignore_index=True)
132.
133.     stop = timeit.default_timer()

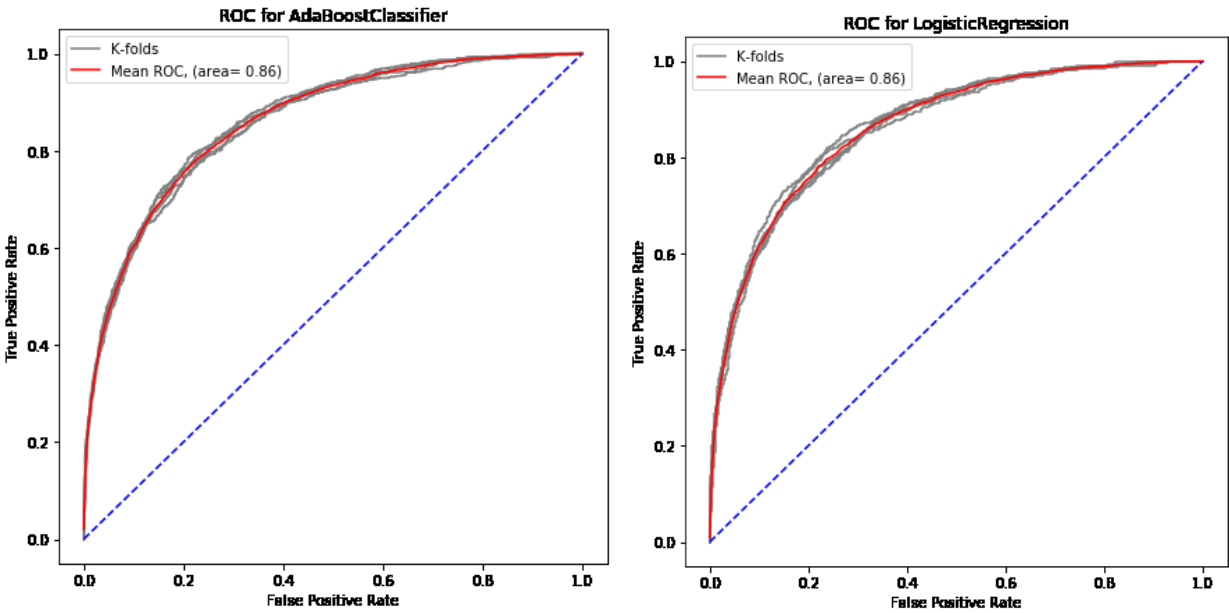
```

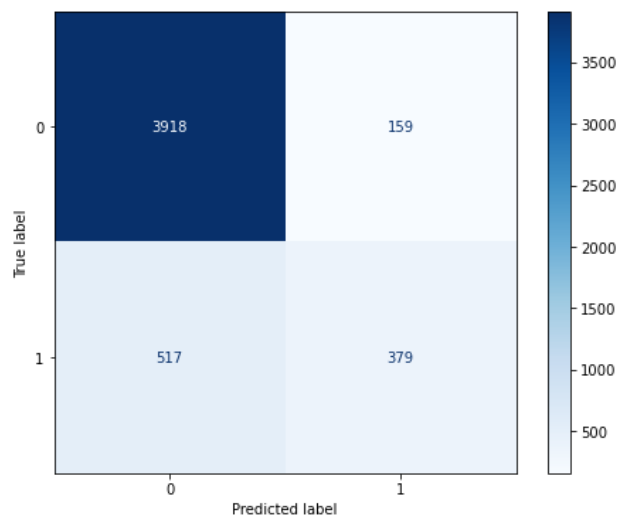
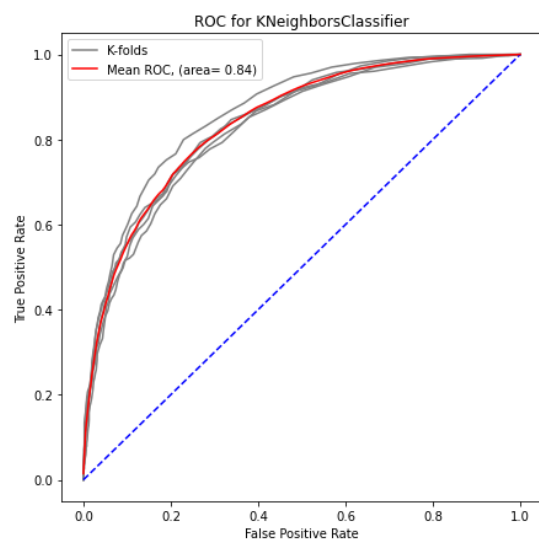
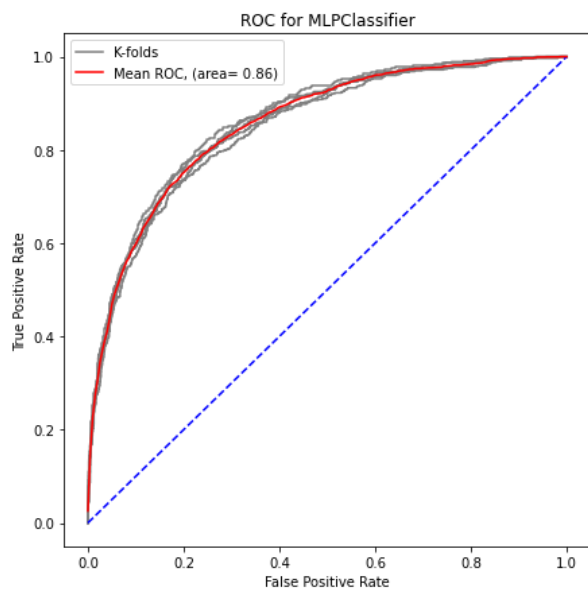
```
134.     print('Time: ', (stop - start)/60)
```

דוגמה של טבלה עם תוצאות של Pipeline:

	drop_features	Combine_features	normalizatio n	model	outliers	pc a	clusterin g	auc
0	['16']	[('2' , '0') , ('17' , '8')]	1	LogReg	3	0	0	0.8635
1	['16']	[('2' , '0') , ('17' , '8')]	1	LogReg	3	0.9	0	0.8313
2	['16']	[('2' , '0') , ('17' , '8')]	1	LogReg	3	0	3	0.8634
3	['16']	[('2' , '0') , ('17' , '8')]	1	LogReg	3	0.9	3	0.8310
4	['12']	[('2' , '0')]	1	LogReg	3	0	0	0.8601
5	['12']	[('2' , '0')]	1	LogReg	3	0.9	0	0.8389
6	['12']	[('2' , '0') , ('17' , '8')]	1	LogReg	3	0	3	0.8601
7	['12']	[('2' , '0') , ('17' , '8')]	1	LogReg	3	0.9	3	0.8388

נספח 3: הערכת המודלים





איור 5: Confusion matrix for logistic regression