



دانشگاه علم و صنعت ایران
دانشکده مهندسی کامپیوتر

عنوان: پروژه چهارم درس داده کاوی
رتبه بندی ویژگی های مؤثر در دسته بندی مجموعه داده ها

نام و نام خانوادگی: آیلین نائب زاده

شماره دانشجویی: ۹۹۵۲۲۱۸۵

نیم سال تحصیلی: پائیز ۱۴۰۲

مدرس: دکتر حسین رحمانی

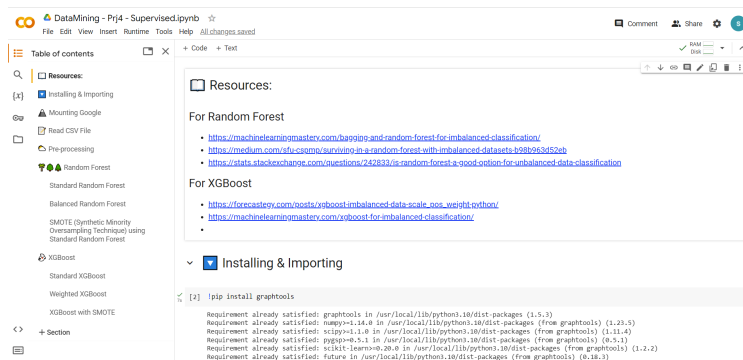
فهرست مطالب

۲	۱ گام اول
۳	۲ گام دوم
۴	۳ گام سوم
۵	۴ گام چهارم
۵	۱.۴ بخش الف)
۶	۲.۴ بخش ب)
۹	۳.۴ بخش ج)
۱۰	۵ گام پنجم
۱۱	۶ گام ششم
۱۲	۷ نتیجه گیری
	۱.۷ آیا ویژگی‌هایی وجود دارند که در هر دو نوع یادگیری Supervised و Unsupervised از نظر تاثیرگذاری هم‌پوشانی داشته باشند؟
۱۲	۲.۷ به نظر شما کدام یک از رتبه‌بندی‌ها منطقی‌تر است؟
۱۲	۳.۷ اگر اختلاف زیادی در معیارهای ارزیابی روش‌ها مشاهده شد، دلیل آن را بنویسید.

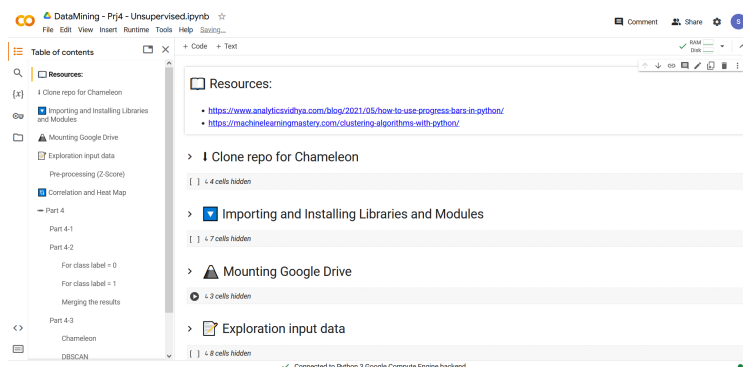
۱ گام اول

در این پروژه هدف این است تا باتوجه به مجموعه داده برچسب گذاری شده، ویژگی موجود را براساس میزان تاثیرگذاری رتبه‌بندی کنیم. فایل ورودی دارای ۸۳ ستون و حدود ۳۰۰۰۰ ردیف می‌باشد. همچنین یک ستون بیانگر شماره ردیف با شروع از ۱ و آخرین ستون نیز بیانگر برچسب نمونه‌ها می‌باشد که می‌تواند مقدار ۰ یا ۱ را داشته باشد.

بدلیل کمبود زمان و پیش بردن سریع‌تر مراحل خواسته‌شده به‌صورت موازی، گام‌های مورد نیاز برای حل پروژه، در دو فایل جداگانه به نام‌های DataMining - Prj4 - Supervised.ipynb و DataMining - Prj4 - Unsupervised.ipynb پیاده‌سازی شده‌اند.



شکل ۱: نمای کلی از پروژه با ناظر



شکل ۲: نمای کلی از پروژه بدون ناظر

در مرحله اول بااستفاده از کتابخانه pandas و تابع read_csv، فایل‌های ورودی را به متغیرهایی از نوع dataframe تبدیل می‌کنیم تا در قدم‌های بعدی راحت‌تر بتوانیم تحلیل‌های مورد نیاز را انجام بدهیم. همچنین بااستفاده از توابع head()، info() و describe(). اطلاعات بیشتری نسبت به داده‌های موجود کسب می‌کنیم.

* خروجی تمامی مراحل ذکر شده بطور کامل در فایل‌های ipynb موجود می‌باشند.

۲ گام دوم

طبق توضیحات ذکر شده در فایل پروژه به‌منظور پیش پردازش داده‌ها تنها کافی است عملیات نرمالایز را بر روی داده‌ها انجام دهیم. البته همان‌گونه که در خروجی تابع `info()` می‌توانید مشاهده کنید، هیچ یک از ستون‌ها شامل مقادیر `None` نمی‌باشد و همگی دارای مقدار غیر تهی می‌باشند.

حال به‌منظور انجام عملیات `normalization` از تابع آماده `StandardScaler()` موجود در کتابخانه `sklearn.preprocessing` استفاده می‌کنیم. درواقع این تابع عملیات `Z-Score` را بر روی ویژگی‌های گوناگون به‌طور جداگانه انجام می‌دهد. البته به این نکته دقت شود که نیازی نیست که بر روی ستون‌های اول و آخر مجموعه داده‌های ورودی نیز این عملیات انجام شود.

$$Z = \frac{x - \mu}{\sigma}$$

Score
Mean

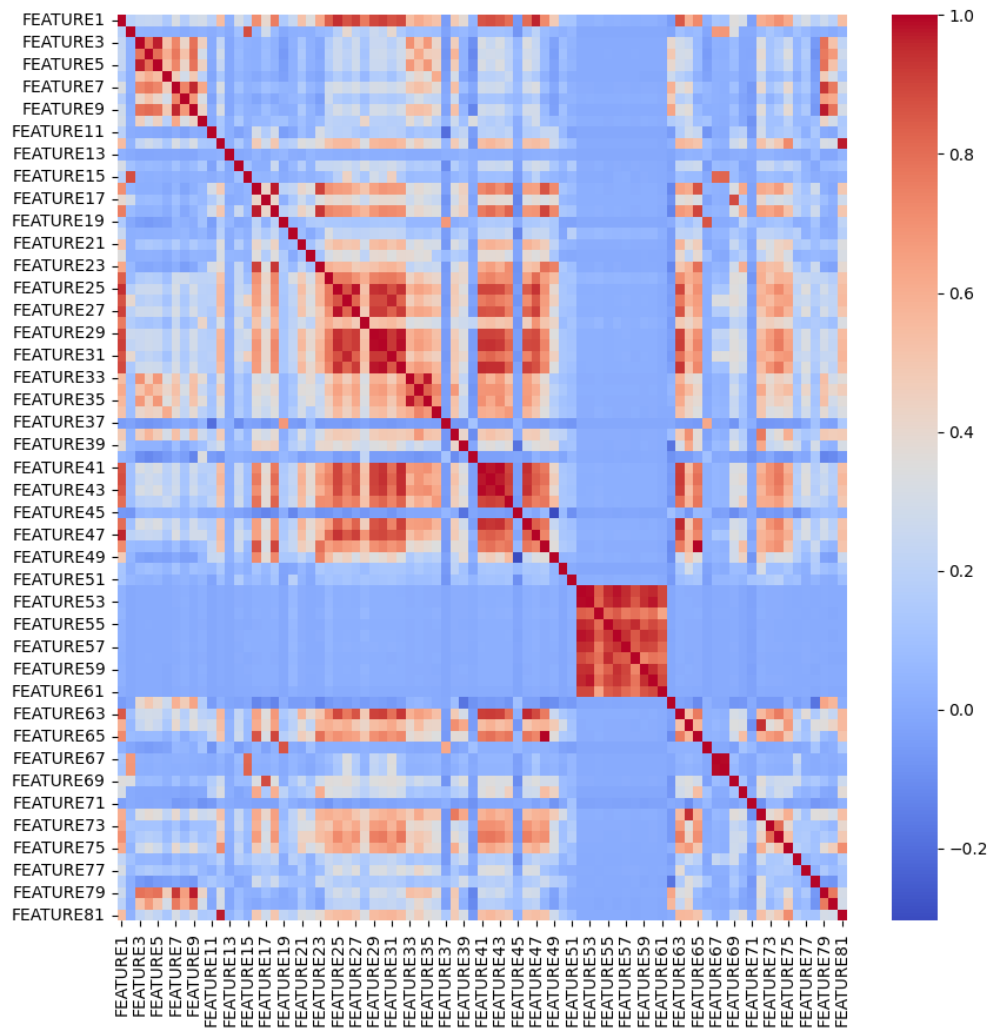
↪
↩

↪ SD

شکل ۳: رابطه Z-Score

۳ گام سوم

در این مرحله همان طور که خواسته شده است، ماتریس همبستگی میان ویژگی های موجود را بدست می آوریم. برای محاسبه این ماتریس می توانیم از تابع آماده `corr()` که در کتابخانه `pandas` موجود می باشد استفاده کنیم. تصویر ماتریس نهایی را در شکل زیر می توانید مشاهده کنید.



شکل ۴: ماتریس همبستگی میان ویژگی ها

۴ گام چهارم

این مرحله درواقع مربوط به استفاده از الگوریتم یادگیری بدون ناظر می باشد. در دو بخش اول از الگوریتم های K-Means و Mini Batch K-Means استفاده می کنیم. و در بخش سوم از الگوریتم BIRCH به همراه الگوریتم PCA به منظور کاهش ابعاد داده ها استفاده خواهیم کرد.

در این مرحله به دلیل حجم بالای داده ها با چالش های زیادی مواجه شدم. بطور مثال در زیر بخش سوم با ارور Session crashed. All the RAM was used. در پروژه Colab چندین بار مواجه شدم که مجبور به کاهش ابعاد مدل شدم. و همچنین بدلیل تعداد بالای نمونه ها، ارزیابی مدل های این مرحله با استفاده از الگوریتم Silhouette به زمان زیادی احتیاج داشتند.

۱.۴ بخش الف)

در این مرحله ابتدا با استفاده از تابع `find_optimal_k_batch()` مقدار بهینه K را برای الگوریتم K-Means که بر روی تمام داده ها اجرا خواهد شد، پیدا خواهیم کرد. نحوه ارزیابی مقدار بهینه K نیز براساس بیشینه مقدار Silhouette می باشد که همانطور که می دانیم هرچه به عدد ۱ نزدیک تر باشد یعنی خوشه های بهتری تشکیل شده اند. طبق نتایج بدست آمده، مقدار $K=3$ بهترین خروجی را برای ما خواهد داشت. (برای این بخش مقادیر K از ۲ تا ۱۰ مورد بررسی قرار گرفته اند). در نهایت الگوریتم K-Means با $k=3$ بر روی تمامی داده ها را مجدداً اجرا می کنیم. در نهایت در یک حلقه با تعداد دفعات تکرار ۸۱ (به تعداد ویژگی ها) و محاسبه نسبت بیشینه مسافت درون خوشه و خارج از آن، میزان اهمیت و تاثیرگذاری ویژگی ها را محاسبه می کنیم. بخشی از کد مربوط به این بخش را در زیر به همراه تعدادی از ویژگی ها به همراه اهمیت شان می توانید مشاهده کنید.

* در فایل اصلی تمامی ویژگی ها همراه با میزان اهمیت شان آورده شده اند.

```

1 # calculate the feature importance for each cluster
2 importance = []
3 for i in range(81):
4     imp = 1 - np.std(centroids[:, i]) / (np.max(centroids[:, i]) - np.min(
5         centroids[:, i]))
6     importance.append(imp)
7
8 # sort the features by importance in descending order
9 features = X_scaled_df.columns
10 sorted_features = sorted(zip(features, importance), key=lambda x: x[1],
11     reverse=True)
12
13 # print the feature importance
14 print("Feature importance for k-means clustering for all the samples:")
15 for f, i in sorted_features:
16     print(f"{f}: {i:.4f}")
    
```

Feature importance for k-means clustering for all the samples:

```
FEATURE64: 0.5918
FEATURE45: 0.5917
FEATURE38: 0.5917
FEATURE72: 0.5917
FEATURE11: 0.5917
FEATURE7: 0.5917
FEATURE76: 0.5915
FEATURE9: 0.5914
FEATURE4: 0.5909
FEATURE79: 0.5904
FEATURE50: 0.5903
FEATURE8: 0.5901
FEATURE71: 0.5898
FEATURE62: 0.5895
FEATURE78: 0.5873
FEATURE10: 0.5860
FEATURE6: 0.5858
FEATURE39: 0.5855
FEATURE5: 0.5847
FEATURE37: 0.5843
FEATURE3: 0.5837
FEATURE51: 0.5819
FEATURE61: 0.5819
FEATURE69: 0.5768
FEATURE80: 0.5763
FEATURE17: 0.5751
FEATURE75: 0.5739
FEATURE81: 0.5730
FEATURE12: 0.5725
FEATURE49: 0.5711
```

شکل ۵: خروجی بخش ۴-الف)

۲.۴ بخش ب)

در این بخش نیز مراحل بخش الف) را تکرار می‌کنیم ولی با این تفاوت که نیاز است الگوریتم K-Means را یک بار بر روی نمونه‌هایی با برچسب ۱ و بار دیگر بر روی نمونه‌هایی با برچسب ۰ جداگانه اجرا کنیم. در هر بار مقدار بهینه K را با استفاده از معیار Silhouette جداگانه محاسبه کردم. همچنین در نهایت با استفاده از یک میانگین وزن دار نتایج اهمیت ویژگی‌ها را با یکدیگر ادغام کردم. بخشی از کد مربوط به این بخش را در زیر به همراه تعدادی از ویژگی‌ها به همراه اهمیت‌شان می‌توانید مشاهده کنید.

* در فایل اصلی تمامی ویژگی‌ها همراه با میزان اهمیت‌شان آورده شده‌اند.

Feature importance for k-means clustering based on the weighted average:

FEATURE18: 0.5916
 FEATURE67: 0.5916
 FEATURE36: 0.5916
 FEATURE27: 0.5916
 FEATURE31: 0.5916
 FEATURE29: 0.5916
 FEATURE26: 0.5915
 FEATURE16: 0.5915
 FEATURE30: 0.5915
 FEATURE41: 0.5915
 FEATURE43: 0.5914
 FEATURE32: 0.5914
 FEATURE25: 0.5914
 FEATURE44: 0.5914
 FEATURE39: 0.5913
 FEATURE40: 0.5912
 FEATURE69: 0.5912
 FEATURE42: 0.5911
 FEATURE33: 0.5911
 FEATURE20: 0.5910
 FEATURE35: 0.5905
 FEATURE7: 0.5900
 FEATURE21: 0.5899
 FEATURE34: 0.5899
 FEATURE46: 0.5898
 FEATURE47: 0.5897
 FEATURE68: 0.5896
 FEATURE73: 0.5894
 FEATURE72: 0.5889
 FEATURE63: 0.5888
 FEATURE76: 0.5875
 FEATURE38: 0.5869
 FEATURE22: 0.5869

شکل ۶: خروجی بخش ۴-ب)

```

1 importance_0 = []
2 importance_1 = []
3 features = X.columns
4 for i in range(len(features)):
5     imp_0 = 1 - np.std(centroids_0[:, i]) / (np.max(centroids_0[:, i]) - np.min(
6         centroids_0[:, i]))
7     importance_0.append(imp_0)
8     imp_1 = 1 - np.std(centroids_1[:, i]) / (np.max(centroids_1[:, i]) - np.min(
9         centroids_1[:, i]))
10    importance_1.append(imp_1)
11 sorted_features_0 = sorted(zip(features, importance_0), key=lambda x: x[1],
12                             reverse=True)

```



```
11 sorted_features_1 = sorted(zip(features, importance_1), key=lambda x: x[1],
    reverse=True)
12 def weighted_average(x, y, w):
13     return [w * a + (1 - w) * b for a, b in zip(x, y)]
14 w_0 = len(X_0) / len(X)
15 w_1 = len(X_1) / len(X)
16 importance_avg = weighted_average(importance_0, importance_1, w_0)
17 sorted_features_avg = sorted(zip(features, importance_avg), key=lambda x: x
    [1], reverse=True)
18 print("Feature importance for k-means clustering based on the weighted average
    :")
19 for f, i in sorted_features_avg:
20     print(f"{f}: {i:.4f}")
```

۳.۴ بخش ج)

در این مرحله نیز از ما خواسته شده است تا با یک الگوریتم خوشه‌بندی دیگر میزان تاثیرگذاری ویژگی‌ها را محاسبه کنیم. ابتدا الگوریتم Chameleon را تست کردم. برای این کار طبق توضیحات فرستاده شده، سعی کردم از کتابخانه تعریف شده یکی از دوستان استفاده کنم. در هنگام استفاده از این الگوریتم به چالش‌هایی همانند نصب کتابخانه‌های مورد نیاز و هماهنگی آن‌ها با سیستم عامل مواجه شدم. همچنین پس از رفع مشکلات گفته شده به زمانی در حدود ۲۰۰ ساعت نیاز بود تا اینکه داده‌های ورودی را خوشه‌بندی کند.

به عنوان جایگزین به سراغ استفاده از الگوریتم DBSCAN می‌روم. در هنگام استفاده از این الگوریتم حتی با وجود کاهش ابعاد داده‌های ورودی، ولی خطای Session is crashed. را در محیط پروژه Colab دریافت کردم. پس از چندین تلاش، در نهایت به سراغ استفاده از الگوریتم BIRCH رفتم. برای استفاده از این الگوریتم نیز در صورتی که همه داده‌های ورودی را به عنوان ورودی استفاده کنیم، خطایی مشابه حالت DBSCAN مشاهده خواهیم کرد. در نتیجه در ابتدا با استفاده از الگوریتم PCA، ابعاد داده‌ها را کاهش می‌دهیم. خروجی نهایی مربوط به این الگوریتم و تاثیر دو ویژگی اصلی را در تصاویر زیر می‌توانید مشاهده کنید.



شکل ۷: خوشه‌های حاصل از بخش ۴-ج)

Feature importance for BIRCH clustering:

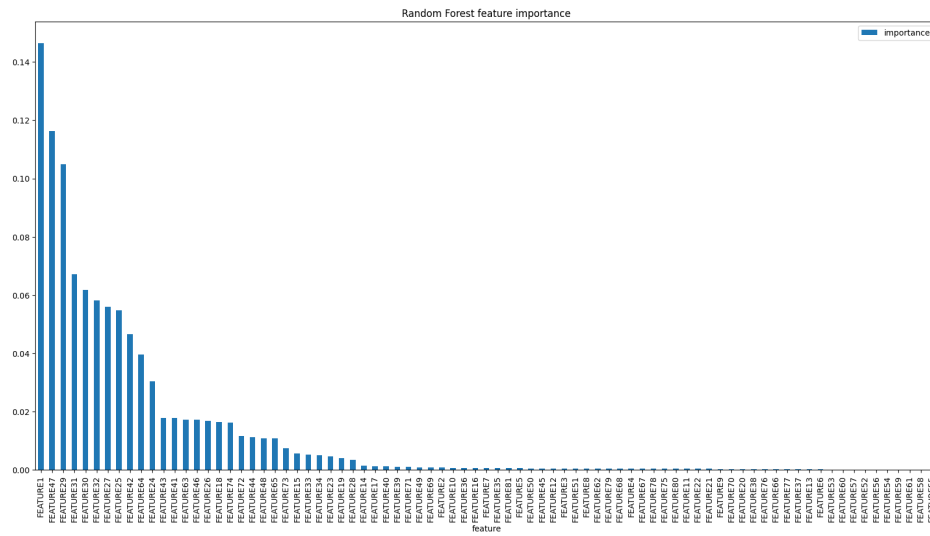
PC1: 0.1663

PC2: 0.0002

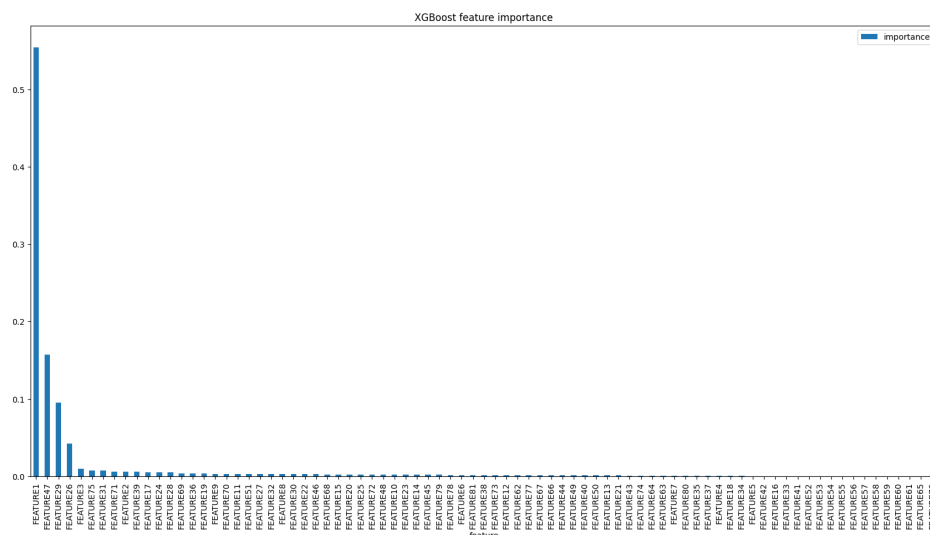
شکل ۸: خروجی بخش ۴-ج)

۵ گام پنجم

در این مرحله باتوجه به خواسته موجود در صورت پروژه الگوریتم Random Forest و XGBoost جهت یادگیری با نظارت استفاده می‌کنیم. همچنین باتوجه به توزیع نامتقارن داده‌ها با برچسب‌های ۰ و ۱ در هنگام استفاده از هر یک از الگوریتم‌های ذکر شده، از دو روش برای بهبود نتایج استفاده می‌کنیم. باتوجه به اینکه بهترین نتایج در حالت Random Forest with Oversampling و Standard XGBoost بدست می‌آیند، تنها برای این دو حالت تاثیرگذاری و میزان اهمیت هر یک از ویژگی‌ها را محاسبه می‌کنیم. نتایج گفته‌شده را در تصاویر زیر می‌توانید مشاهده کنید.



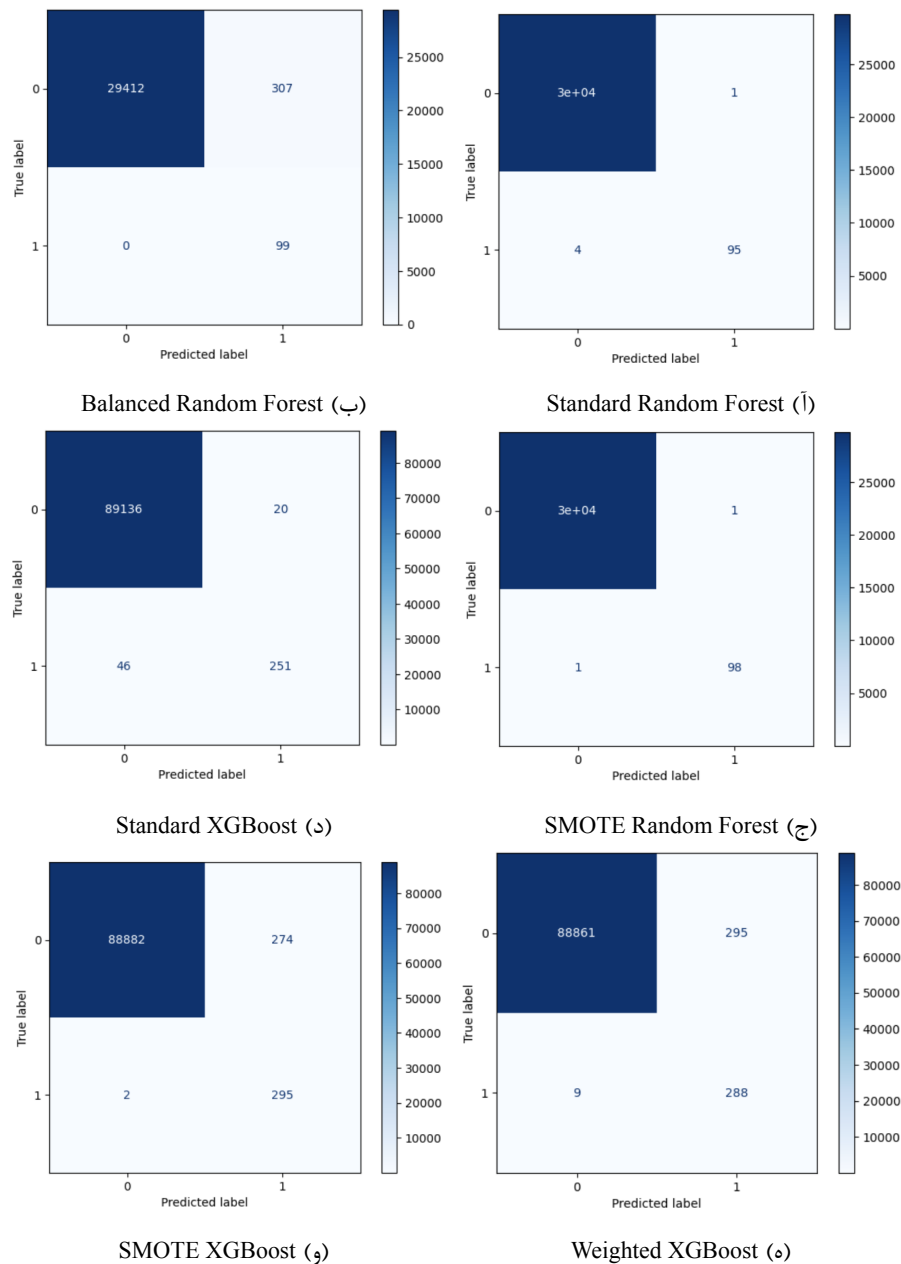
شکل ۹: میزان اهمیت و تاثیرگذاری ویژگی‌ها بر روی الگوریتم Random Forest with SMOTE



شکل ۱۰: میزان اهمیت و تاثیرگذاری ویژگی‌ها بر روی الگوریتم Standard XGBoost

۶ گام ششم

نتایج حاصل از ارزیابی معیارهای گوناگون بر روی مدل‌های استفاده‌شده را در فایل اکسل ارسال شده می‌توانید مشاهده کنید. همچنین در تصاویر زیر نیز می‌توانید Confusion Matrix الگوریتم‌های یادگیری با نظارت استفاده‌شده در حالت‌های مختلف را مشاهده کنید.



۷ نتیجه گیری

۱.۷ آیا ویژگی‌هایی وجود دارند که در هر دو نوع یادگیری **Supervised** و **Unsupervised** از نظر تاثیرگذاری هم‌پوشانی داشته باشند؟

در حالت یادگیری بدون ناظر هم‌پوشانی مشاهده نمی‌شود. ولی در یادگیری باناظر همان‌گونه که در تصاویر گام پنجم توضیح داده شده‌است، اولین ویژگی بیشترین تاثیر در حالت‌های محاسبه‌شده را داشته‌است و بین سایر ویژگی‌ها نیز هم‌پوشانی و شباهت دیده می‌شود.

۲.۷ به نظر شما کدام یک از رتبه‌بندی‌ها منطقی‌تر است؟

باتوجه به نتایج خروجی و همچنین چالش‌های حل مسئله، الگوریتم با ناظر عملکرد بهتری داشته‌اند. باتوجه به حجم بالای داده‌ها در هنگام استفاده از الگوریتم‌های بدون ناظر ممکن است با چالش‌های سخت‌افزاری مواجه شویم. ولی با تغییرات جزئی در الگوریتم‌های با ناظر می‌توانیم به دقت و صحت بالایی دست پیدا کنیم. همچنین یکی از چالش‌های دیگر در هنگام استفاده از الگوریتم بدون ناظر تعیین مقادیر برای پارامترهای مورد نیاز آن‌ها می‌باشد که نتایج نهایی را تحت تاثیر قرار می‌دهند.

۳.۷ اگر اختلاف زیادی در معیارهای ارزیابی روش‌ها مشاهده شد، دلیل آن را بنویسید.

در الگوریتم‌های با ناظر مقادیر accuracy, precision, recall و F1-score شباهت زیادی با یکدیگر دارند. در حالت استفاده از الگوریتم‌های بی‌ناظر همان‌طور که در بخش چهار اشاره شده‌است، مقادیر Silhouette شباهت زیادی به یکدیگر دارند.