



دانشکده مهندسی کامپیوتر

پروژه درس

هوش مصنوعی

نام و نام خانوادگی : مهدی فقهی

بهمن ۱۳۹۹

۱.۰ فازیک

در قسمت اول مساله از ما خواسته شده بود که دو فایل را دریافت کنیم . در یکی اطلاعات مربوط به ایمیل‌های هرزنامه و در دیگری ایمیل‌های سالم را داشتیم . در اولین قدم به کمک `Pandas library` اطلاعات هر دو فایل را به تبدیل به یک `data frame` کردم و این دو تا دیتا فریم را یکی کردم . در قدم بعدی براساس نوع فایلی که از آن آمده بودند یک ستون به دیتا فریم مربوطه اضافه کردم به اسم `real` که مشخص می‌کرد ایمیل موجود سالم هست یا نه .

```

۱ def faz_one():
۲     pol_fake = pd.read_csv("pol-fake.csv")
۳     pol_real = pd.read_csv("pol-real.csv")
۴     pol_real["real"] = 1
۵     pol_fake["real"] = 0
۶
۷     result = pd.concat([pol_real, pol_fake])
۸
۹     shuffled_resutl = result.iloc[np.random.permutation(len(result))].reset_index(drop=True)

```

در قدم بعد اسم تمامی ستون‌های موجود در دیتا فریم ها را ذخیره کردم سپس یک لیست را به عنوان لیست که ستون‌هایی هستند که به کمک آن‌ها قرار هست که درخت تصمیم موردنظرمون را بسازیم را جدا کردیم . بعد دیتاهای خویش را به دو دسته دیتا که یکی برای ساخت مدل و دیگری دیتایی که برای بررسی صحت و کارآمدی مدل موجود به کار ما میاید تقسیم میکنیم . در انتها لیستی از دیتا فریم اصلی و دیتا فریم‌هایی که برای ساخت مدل اصلی و آزمایش مدل ساخته شده استفاده می‌شوند را به عنوان خروجی تابع بازمیگردانیم

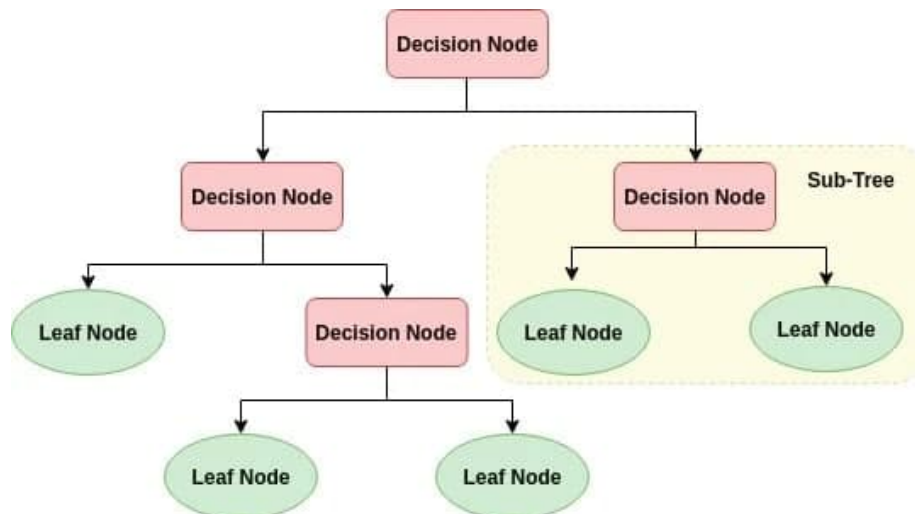
```

۱ label_name = list(shuffled_resutl.columns)
۲ feature_cols = label_name[:-1]
۳ variable target and features in dataset split #
۴ x = shuffled_resutl[feature_cols] Features #
۵ y = shuffled_resutl.real variable Target #
۶ set test and set training into dataset Split #
۷ X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=2.0,
۸                                                     random_state=1) test 20% and training 80% #
۹ return [X_train, X_test, y_train, y_test, feature_cols], shuffled_resutl
۱۰

```

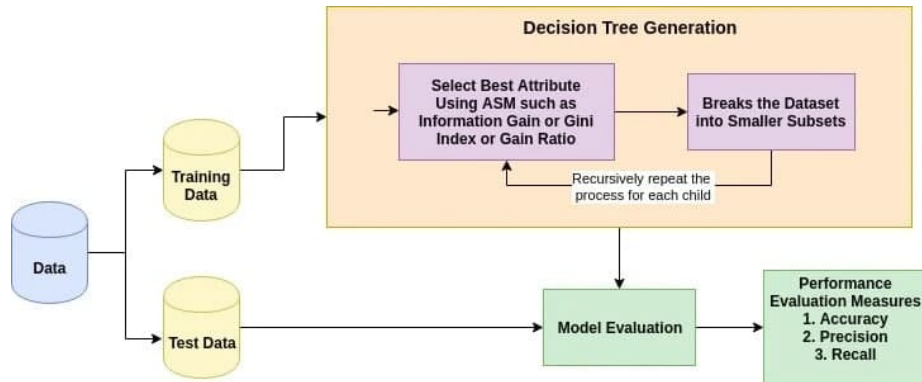
۲.۰ فاز دو

برای حل قسمت دوم پروژه ابتدا لازم هست که با مفهوم **Decision Tree** آشنا بشویم . یکی از پرکاربردترین الگوریتم‌های داده‌کاوی، الگوریتم درخت تصمیم است. در داده‌کاوی، درخت تصمیم یک مدل پیشبینی کننده است به طوری که می‌تواند برای هر دو مدل رگرسیون و طبقه‌ای مورد استفاده قرار گیرد. زمانی که درخت برای کارهای طبقه‌بندی استفاده می‌شود، به عنوان درخت طبقه‌بندی **Classification Tree** هنگامی که برای فعالیت‌های رگرسیونی به کار می‌رود درخت رگرسیون **Classification Tree** نامیده میشود. پیچیدگی زمانی درختان تصمیم تابعی از تعداد سوابق و تعداد صفات در داده های داده شده است. درخت تصمیم روشی بدون توزیع یا غیر پارامتری است که به فرضیات توزیع احتمال بستگی ندارد. درختان تصمیم می توانند داده های با ابعاد بالا را با دقت خوب اداره کنند.



۱.۲.۰ الگوریتم درخت تصمیم چگونه کار می کند؟

- با استفاده از **Attribute Selection Measures (ASM)** بهترین ویژگی را برای تقسیم سوابق انتخاب کنید
- این ویژگی را گره تصمیم گیری کرده و مجموعه داده را به زیرمجموعه های کوچکتر تقسیم می کند.
- درخت سازی را با تکرار این فرآیند به صورت بازگشتی برای هر کودک شروع می کند.



۲.۲.۰ Scikit-learn in Building Classifier Tree Decision

ابتدا مجموعه هایی که در فاز قبلی دسته بندی کرده بودیم را به عنوان ورودی تابع به تابع فاز جدید پاس می‌دهیم . با استفاده از توابعی که در لایبری فوق موجود هست مدل خودمان را براساس دیتاهای موجود در دیتافریم مربوط به به مدل سازی میسازیم و با کمک تابع زیر به ترتیب `accuracy` ، `confusion matrix` ، `precision` ، `recall` ، `F1` را بدست می‌آوریم.

```

1 def faz_two(input_list):
2     X_train, X_test, y_train, y_true, feature_cols = input_list
3     Classifier Tree Decision Train #
4     clf = DecisionTreeClassifier()
5     Classifier Tree Decision Train #
6     clf.fit(X_train, y_train)
7     dataset test for response the Predict #
8     y_pred = clf.predict(X_test)
9
10    accuracy = metrics.accuracy_score(y_true, y_pred)
11    print(" Accuracy:, accuracy)
12
13    print( matrix: "confusion\n", metrics.confusion_matrix(y_true, y_pred))
14
15    print(" precision:, metrics.precision_score(y_true, y_pred))
16
17    print(" recall:, metrics.recall_score(y_true, y_pred))
18
19    print(" F1:, metrics.f1_score(y_true, y_pred))
20
21
22    return accuracy
  
```

تعریف ۱ (Accuracy) دقت را می‌توان با مقایسه مقادیر مجموعه آزمایش واقعی و مقادیر پیش بینی شده محاسبه کرد.

تعریف ۲ (confusion-matrix) جدولی با دو ردیف و دو ستون است که تعداد مثبت کاذب ، منفی کاذب ، مثبت واقعی و منفی واقعی را گزارش می دهد. این اجازه می دهد تا تجزیه و تحلیل دقیق تر از تناسب طبقه ما به ما میدهد. به عنوان مثال ، اگر ۹۵ گربه و فقط ۵ سگ در داده ها وجود داشته باشد ، یک طبقه بندی کننده خاص ممکن است همه مشاهدات را به عنوان گربه طبقه بندی کند. دقت کلی ۹۵٪ است ، اما با جزئیات بیشتر طبقه بندی کننده دارای ۱۰۰٪ میزان تشخیص (حساسیت) برای کلاس گربه اما ۰٪ میزان تشخیص برای کلاس سگ است.

		Prediction	
		0	1
Actual	0	TN	FP
	1	FN	TP

تعریف ۳ (precision) . حاصل تقسیم tp یا همان تعداد مثبت های واقعی بر $tp+fp$ میباشد

تعریف ۴ (recall) . حاصل تقسیم tp یا همان تعداد مثبت های واقعی بر $tp+fn$ میباشد fn همان تعداد منفی های کاذب است .

تعریف ۵ ($F1$) .

$$F1 = 2 * (precision * recall) / (precision + recall)$$

$F1$ یک اندازه گیری کلی از دقت مدل است که دقت و یادآوری را با هم ترکیب می کند ، با این روش عجیب و غریب که جمع و ضرب فقط دو ویژگی قبل را مخلوط می کند تا یک ویژگی جدید به وجود بیاید . یعنی نمره خوب $F1$ به این معنی است که شما مثبت کاذب کم و منفی کاذب کم دارید ، بنابراین تهدیدهای

واقعی را به درستی شناسایی می شود. نمره $F1$ در صورت ۱ کامل در نظر گرفته می شود ، در حالی که مدل در صورت ۰ کاملاً شکست می خورد. .

۳.۲.۰ Trees Decision Visualizing

در این قسمت درخت تصمیم خود را به صورت یک گراف چاپ میکنیم . درخت موردنظر یک مدل هست که به کمک این درخت می توانیم به صورت دستی نیز با توجه به داده هایی که داریم متوجه بشویم که در انتهای درخت یعنی در برگ یک ایمیل هرزنامه هست یا نه . صرف طی کردن درخت با توجه به فیچرهایی که داریم و رسیدن به برگ موردنظر .

```

۱ from sklearn.tree import export_graphviz
۲ from sklearn.externals.six import StringIO
۳ from IPython.display import Image
۴ import pydotplus
۵
۶ dot_data = StringIO()
۷ export_graphviz(clf, out_file=dot_data,
۸                 filled=True, rounded=True,
۹                 special_characters=True, feature_names = feature_cols, class_names=['0', '1'])
۱۰ graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
۱۱ graph.write_png('faze2.png')
۱۲ Image(graph.create_png())

```

یا به شکل زیر پیاده سازی کنیم

```

۱ from sklearn.externals.six import StringIO
۲ from IPython.display import Image
۳ from sklearn.tree import export_graphviz
۴ import pydotplus
۵ dot_data = StringIO()
۶ export_graphviz(clf, out_file=dot_data,
۷                 filled=True, rounded=True,
۸                 special_characters=True, feature_names = feature_cols, class_names=['0', '1'])
۹ graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
۱۰ graph.write_png('teeFaz2.png')
۱۱ Image(graph.create_png())

```

در نهایت درخت موردنظر را به صورت عکس ذخیره میکنیم .

۳.۰ فاز سه

برای اجرای فاز سوم ابتدا باید Attribute Selection Measures آشنا بشویم. انتخاب ویژگی که به کمک آن داده ها را به بهترین شکل تقسیم کنیم، ابتکاری است. پس ASM به عنوان قوانین تقسیم شناخته می شود زیرا به ما کمک می کند تا نقاط شکست tuples را در یک گره مشخص تعیین کنیم. ASM با توضیح مجموعه داده داده شده به هر ویژگی (یا ویژگی) رتبه می دهد.

مشهورترین معیارهای انتخاب عبارتند از: Gain، Gain Ratio، همچنین Gini Index

۱.۳.۰ Measures Selection Attribute

تعریف ۶ (Gain-Information).

شانون مفهوم آنتروپی را ابداع کرد که ناخالصی مجموعه ورودی را اندازه گیری می کند. در فیزیک و ریاضیات، آنتروپی به عنوان تصادفی یا ناخالصی در سیستم گفته می شود. در نظریه اطلاعات، به ناخالصی در گروهی از داده ها اشاره دارد. کسب اطلاعات، کاهش آنتروپی است. Information Gain بین آنتروپی قبل از تقسیم و میانگین آنتروپی پس از تقسیم مجموعه داده را بر اساس مقادیر مشخصه محاسبه می کند.

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

$$\text{Info}_A(D) = \sum_{j=1}^V \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

تعریف ۷ (Ratio-Gain) Information Gain برای attribute که outcomes زیادی دارند به صورت جانبدارانه عمل میکند. به این معنی که صفت ای که values distinct بیشتری دارند را ترجیح میدهد. نسبت سود را می توان به این صورت تعریف کرد:

Y

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo_A(D)}$$

Ratio Gain

تعریف \wedge (Gini-index) $Gini\ index$ الگوریتم درخت تصمیم دیگر CART (درخت طبقه بندی و رگرسیون) از روش جینی برای ایجاد نقاط تقسیم استفاده می کند.

$$Gini(D) = 1 - \sum_{i=1}^m P_i^2$$

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

$$\Delta Gini(A) = Gini(D) - Gini_A(D).$$

ویژگی با حداقل شاخص جینی به عنوان ویژگی تقسیم انتخاب می شود.

۲.۳.۰ Implementatio

خوب در این مرحله ما درخت خودمون را به کمک information Gain پیاده سازی میکنیم . توجه داشته باشید که در مرحله قبل چون این قسمت از تابع را که در پایین می بینید مقدار دهی نکرده بودم درواقع

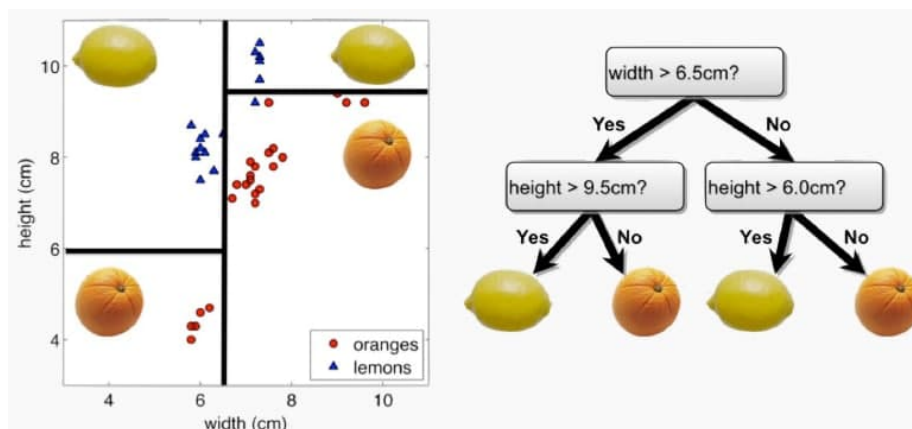
به صورت پیش فرض `Gini index` پیاده سازی شد .

```
clf = DecisionTreeClassifier(criterion="entropy")
```

بقیه موارد هم همانند فاز قبل با `Library sklearn` پیاده سازی شده است .

۴.۰ فاز چهار

برای اینکه این قسمت را بهتر متوجه بشویم ابتدا لازم هست که مفهوم درخت تصمیم رو دوباره بررسی کنیم .



یادگیری کوچکترین درخت تصمیم گیری، برای مجموعه داده های آموزش داده شده کار دشواری است. در هر گره ، ما باید پیش بینی کننده مطلوب را برای تقسیم انتخاب کنیم و مقداری بهینه برای تقسیم انتخاب کنیم. در شکل بالا ، ابتدا با بررسی عرض میوه و استفاده از آستانه ۵.۶ سانتی متر تصمیم می گیریم. ما همچنین می توانستیم با ارتفاع میوه یا انتخاب یک مقدار آستانه متفاوت برای عرض شروع کنیم. انتخاب ما در شکل مناطق درخت تأثیر دارد.

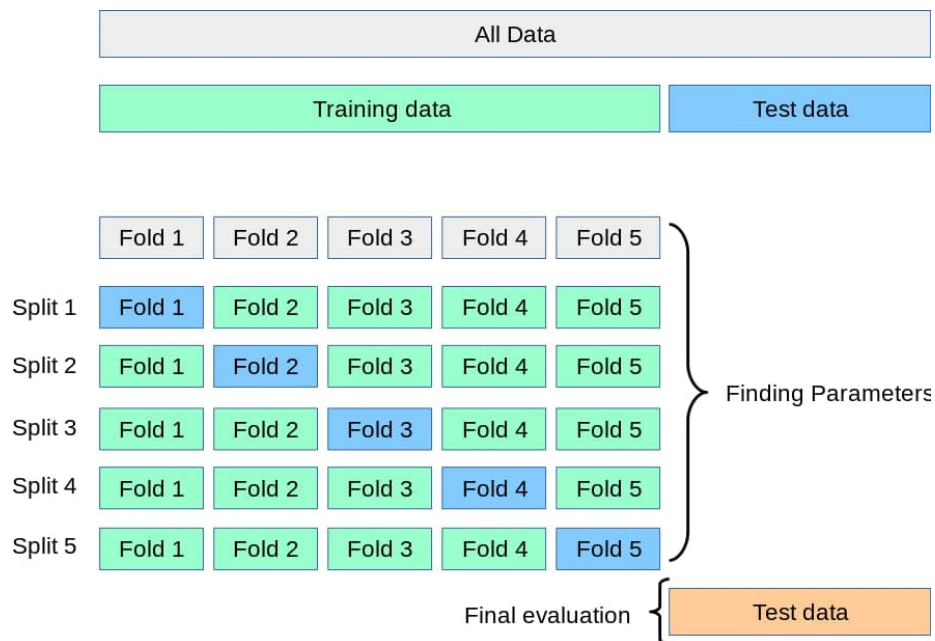
در حین تمرین ، رشد درخت ادامه خواهد یافت تا جایی که هر منطقه دقیقاً یک نقطه تمرین داشته باشد (۱۰۰٪ دقت آموزش). این منجر به یک درخت طبقه بندی کامل می شود که داده های آموزش را تقسیم می کند تا هر مرخصی شامل یک مشاهده باشد. به ترتیب ، درخت کامل بیش از حد با داده های آموزشی سازگار است.

اینجاست که باید با مفهوم `Over fitting` آشنا بشویم.

یعنی اینکه مدل ما به گونه‌ای با داده‌های آموزشی اول ست شود که با هر تغییر کوچکی در دیتای ورودی که قبلاً ندیده است پیش‌بینی جدید تغییرات بسیاری داشته باشد. عملاً مدل ما روی یک سری دیتا خوب کار خواهد کرد که قبلاً دیده است. و برای دیتاهای جدید ممکن است خیلی بد پیش‌بینی انجام بدهد. برای حل این مشکل ما از **trade of** استفاده می‌کنیم.

ما باید یک شرط توقف تعریف کنیم درواقع در این نوع مسئله توقف در سطحی از درخت که هرچه درخت واریانس داده‌ها را کم می‌کنیم بایاس داده‌ها نیز افزایش زیادی نداشته باشد. درواقع به دنبال یک سازش بین واریانس و بایاس هستیم.

در اینجا از **Cross Validation** کمک می‌گیریم. عمق مناسب را می‌توان با ارزیابی درخت از طریق یک مجموعه داده نگهدارنده از طریق **cross validation** تعیین کرد. با بازگیری مجدد داده‌ها به دفعات، تقسیم به داده‌های آموزش و داده‌های تأیید اعتبار، بررسی درختان با اندازه ارتفاع مختلف و ارزیابی دقیق طبقه‌بندی به وجود آمده نسبت به داده مورد تأیید، می‌توانیم عمق درخت را پیدا کنیم که بهترین مدل را به ما می‌دهد و عملاً سازشی بین واریانس و بایاس ایجاد می‌کنیم.



در ادامه این مسیر نیاز هست که از **K-fold CrossValidation** استفاده کنیم. K در اینجا مشخص می‌کند که داده‌های خودمون را به چند قسمت تقسیم کنیم و در هر قسمت یکی را به عنوان داده‌ای

۱۰

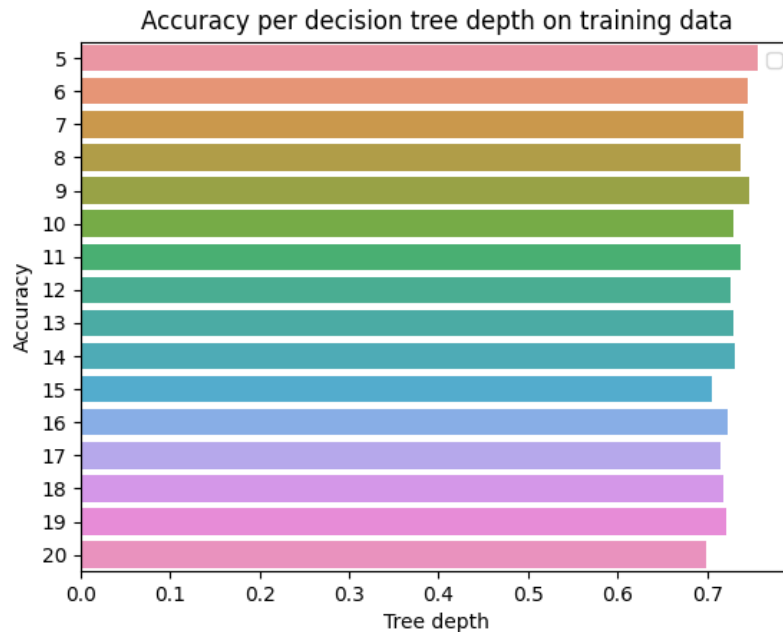
که می‌خواهیم براساس آن داده‌های موجود را تست کنیم در نظر می‌گیریم .

در خواست مسئله از ما خواسته است که آن را ۱۰ در نظر بگیریم .

```
1 def faz_four(input_list):
2
3     kf = KFold(n_splits=10, shuffle=False, random_state=None)
4     lable_name = list(input_list.columns)
5     feature_cols = lable_name[:-1]
6     X = input_list[feature_cols]   Features #
7     y = input_list.real   variable Target #
8
9     best_acuracy = 0
10    best_depth = 0
11    list_of_depth = []
12    list_of_acurrcy = []
13    for val in range(5, 21):
14        score = cross_val_score(DecisionTreeClassifier(max_depth=val, random_state=None), X, y, cv=kf,
15                                scoring="accuracy")
16
17        list_of_depth.append(str(val))
18        list_of_acurrcy.append(score.mean())
19        if score.mean() > best_acuracy:
20            best_acuracy = score.mean()
21            best_depth = val
22
23
24    make_ser = pd.Series(list_of_acurrcy, index=list_of_depth)
25    sns.barplot(x=make_ser, y=make_ser.index)
26    graph your to labels Add #
27    plt.xlabel('depth Tree')
28    plt.ylabel('Accuracy')
29    plt.title(data" training on depth tree decision per "Accuracy)
30    plt.legend()
31    plt.show()
32
33    return best_acuracy, best_depth
```

براساس عمق ۵ تا ۲۰ 10-fold CrossValidation را بدست می‌اریم و در هر عمق متوسط دقت را در نظر می‌گیریم که اگر بهتر از دقت در عمق قبلی بود همین عمق را به عنوان خروجی در نظر می‌گیریم و در غیر اینصورت عمقی که درخت را تشکیل می‌دهیم را به‌روزرسانی می‌کنیم . درنهایت بهترین عمق و دقتش را به عنوان خروجی بازمیگردانیم.

همچنین در کد قسمتی را برای نشان دادن دقت هر عمق و میزان عمق به کمک Matplotlib به شکل یک نمودار نشان می‌دهیم که شکل زیر یکی از خروجی های آن هست برای داده‌های موجود .



که نشان می‌دهد عمق پنج عمق بهتری نسبت به بقیه عمق‌ها برای پیدا کردن حاصل و جواب هست .

۵.۰ فاز نهایی

Random Forests Classifiers

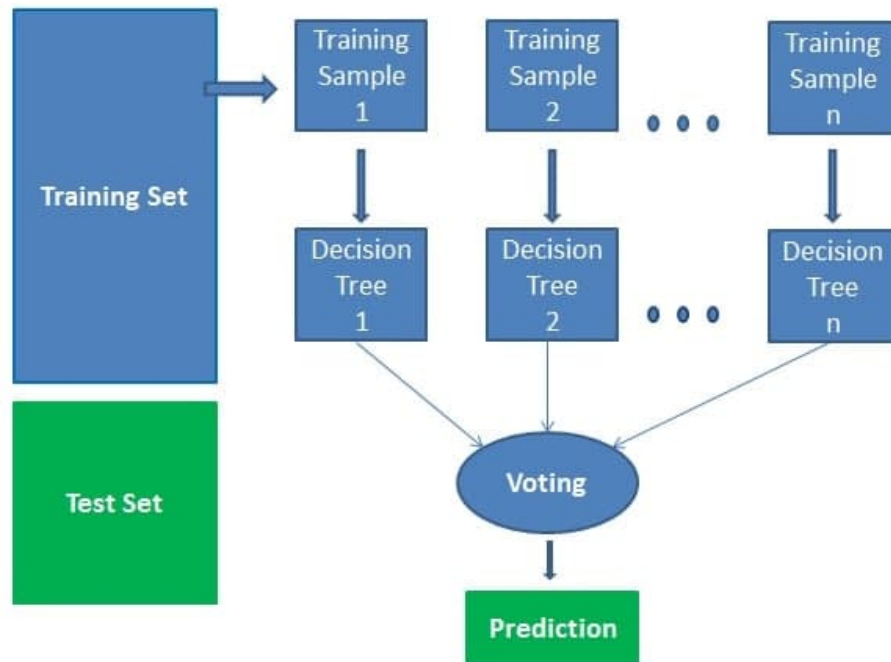
جنگل‌های تصادفی یک الگوریتم یادگیری **Supervised learning** است. هم برای طبقه‌بندی و هم برای رگرسیون قابل استفاده است. همچنین انعطاف پذیرترین و آسان‌ترین کاربرد الگوریتم است. یک جنگل از درختان تشکیل شده است. گفته می‌شود هرچه تعداد درختان آن بیشتر باشد، جنگل نیز از استحکام بیشتری برخوردار است. جنگل‌های تصادفی درختان تصمیم‌گیری را بر روی نمونه‌های داده‌ای که به طور تصادفی انتخاب شده‌اند ایجاد می‌کنند، از هر درخت پیش‌بینی می‌کنند و با استفاده از رای‌گیری بهترین راه حل را انتخاب می‌کنند. همچنین یک شاخص بسیار خوب از اهمیت ویژگی را فراهم می‌کند.

بیایید الگوریتم را از نظر عرفی درک کنیم. فرض کنید می‌خواهید به سفر بروید و دوست دارید به مکانی سفر کنید که از آن لذت بیشتری می‌برید. بنابراین برای یافتن مکانی که دوست داشته باشید چه می‌کنید؟ می‌توانید به صورت آنلاین جستجو کنید، نظرات را در وبلاگ‌ها و پورتال‌های مسافرتی بخوانید، یا می‌

توانید از دوستان خود نیز سوال کنید. فرض کنید شما تصمیم گرفته اید از دوستان خود سوال کنید و با آنها در مورد تجربه سفر گذشته آنها به مکان های مختلف صحبت کنید. از هر دوستی توصیه هایی دریافت خواهید کرد. اکنون شما باید لیستی از آن مکانهای پیشنهادی تهیه کنید. سپس ، از آنها می خواهید از لیست مکانهای پیشنهادی که ایجاد کرده اید رأی دهند (یا یکی از بهترین مکانها را برای سفر انتخاب کنید). مکانی که بیشترین تعداد آرا را داشته باشد انتخاب نهایی شما برای سفر خواهد بود. در روند تصمیم گیری فوق ، دو قسمت وجود دارد. ابتدا ، از دوستان خود در مورد تجربه سفر شخصی خود می پرسید و از چندین مکان که بازدید کرده اند ، یک توصیه دریافت می کنید. این قسمت مانند استفاده از الگوریتم درخت تصمیم است. در اینجا ، هر دوست انتخاب مکانهایی را که تاکنون بازدید کرده است ، انجام می دهد. قسمت دوم ، پس از جمع آوری تمام توصیه ها ، روش رأی گیری برای انتخاب بهترین مکان در لیست توصیه ها است. کل این فرایند دریافت توصیه ها از دوستان و رای گیری درباره آنها برای یافتن بهترین مکان به عنوان الگوریتم جنگل های تصادفی شناخته می شود. این از نظر فنی یک روش گروهی است (براساس conquer and divide) درختان تصمیم که در یک مجموعه داده تقسیم شده تصادفی تولید می شوند. این مجموعه از طبقه بندی کننده درخت تصمیم به جنگل نیز معروف است. درختان تصمیم گیری فردی با استفاده از یک شاخص انتخاب ویژگی مانند کسب اطلاعات ، نسبت سود و شاخص جینی برای هر ویژگی تولید می شوند. هر درخت به یک نمونه تصادفی مستقل بستگی دارد. در یک مسئله طبقه بندی ، هر درخت رأی می دهد و محبوب ترین کلاس به عنوان نتیجه نهایی انتخاب می شود. در حالت رگرسیون ، میانگین تمام خروجی های درخت به عنوان نتیجه نهایی در نظر گرفته می شود. در مقایسه با الگوریتم های طبقه بندی غیر خطی ساده تر و قدرتمندتر است.

Algorithm

- نمونه های تصادفی را از یک مجموعه داده مشخص انتخاب کنید.
- برای هر نمونه یک درخت تصمیم بسازید و از هر درخت تصمیم یک نتیجه پیش بینی بگیرید.
- برای هر نتیجه پیش بینی شده رأی دهید.
- نتیجه پیش بینی را با بیشترین آرا به عنوان پیش بینی نهایی انتخاب کنید.



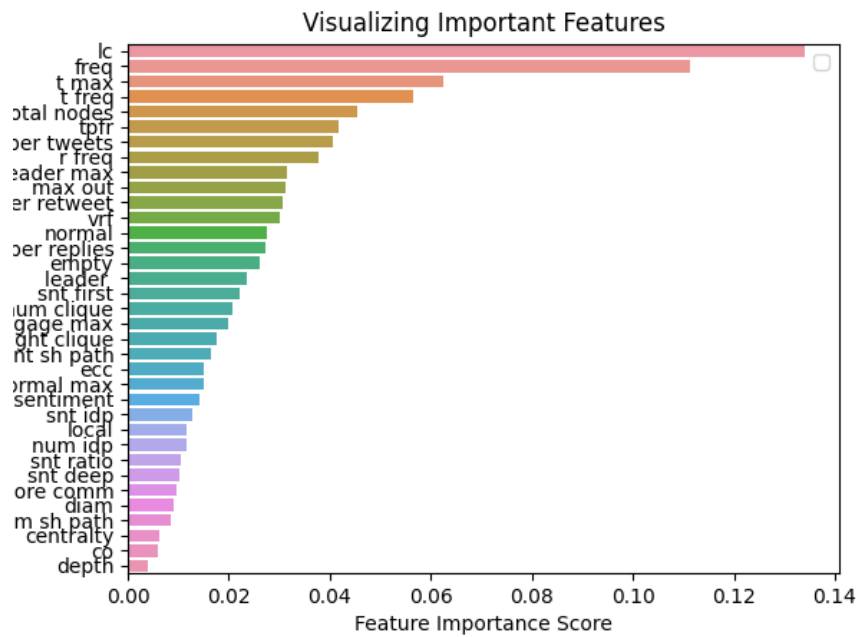
مراحل زیر در شکل زیر توسط کد زیر اجرا شده است .

```

1 def faz_final(shuffled_resutl):
2     label_name = list(shuffled_resutl.columns)
3     feature_cols = label_name[:-1]
4
5     variable target and features in dataset split #
6     x = shuffled_resutl[feature_cols] Features #
7     y = shuffled_resutl.real variable Target #
8
9     set test and set training into dataset Split #
10    X_train, X_test, y_train, y_true = train_test_split(x, y, test_size=2.0,
11                                                         random_state=1) test 20% and training 80% #
12
13    Classifier Gaussian a Create #
14    clf = RandomForestClassifier(n_estimators=100)
15
16    y_pred=clf.predict(X_test) sets training the using model the Train #
17    clf.fit(X_train, y_train)
18
19    y_pred = clf.predict(X_test)
20
21    accuracy = metrics.accuracy_score(y_true, y_pred)
22    ) ,accuracyprint("Accuracy:" #
  
```

Finding important features

یکی از موارد خوبی که این نوع پیاده‌سازی دارد این هست که به کمک این پیاده‌سازی می‌توانیم `feature` که تاثیر بیشتری بر روی نتیجه‌گیری درخت ما دارد را پیدا کنیم. در پیدا سازی پایین من این موارد رو پیدا کردم و به صورت یک نمودار نشان دادم که نمودار را در پایین می‌بینید.



در نهایت می‌توانیم مواردی که کمتر از یک مقدار خاص هستند را حذف کنیم و دقت مدل را دوباره بررسی می‌کنیم اگر بهتر شده بود که مدل جدید با `feature` حذف شده را مورد استفاده قرار می‌دهیم در غیر این صورت همان مدل قبلی با دقت و `accuracy` را برمی‌گردانیم.

```

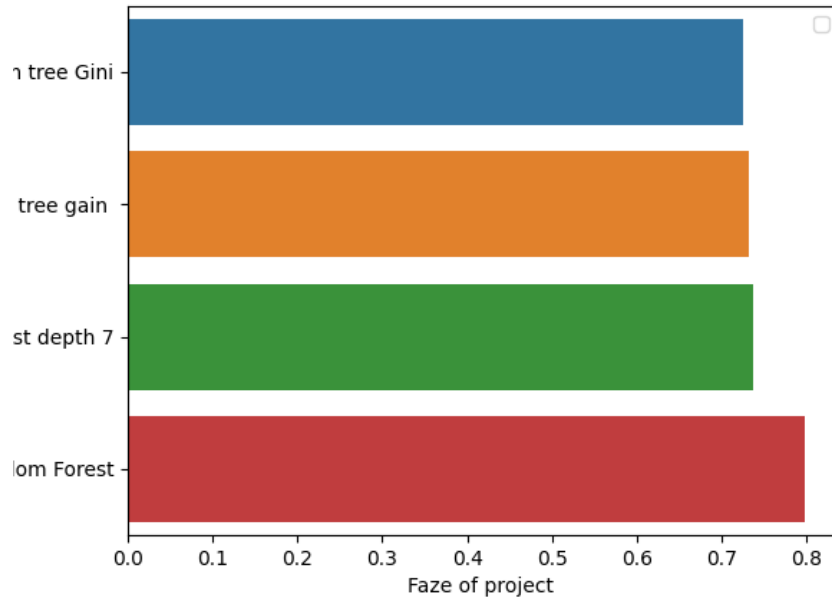
۱ feature_imp = pd.Series(clf.feature_importances_, index=feature_cols).sort_values(ascending=False)
۲
۳ sns.barplot(x=feature_imp, y=feature_imp.index)
۴ graph your to labels Add #
۵ plt.xlabel('Score Importance Feature')
۶ plt.ylabel('Features')
۷ plt.title('Features" Important "Visualizing')
۸ plt.legend()
۹ plt.show()
۱۰
۱۱ new_feature_cols = []
۱۲ number = 0
۱۳
۱۴ feature_imp_dict = feature_imp.to_dict()
۱۵
۱۶ for item in feature_imp_dict:
۱۷
۱۸     if feature_imp_dict[item] >= 0.09:
۱۹         new_feature_cols.append(item)
۲۰
۲۱ x = shuffled_resutl[new_feature_cols] Features #
۲۲
۲۳ X_train, X_test, y_train, y_true = train_test_split(x, y, test_size=2.0,
۲۴                                                     random_state=1) test 20% and training 80% #
۲۵
۲۶ clf = RandomForestClassifier(n_estimators=100)
۲۷ clf.fit(X_train, y_train)
۲۸
۲۹ y_pred = clf.predict(X_test)
۳۰
۳۱ accurecy2 = metrics.accuracy_score(y_true, y_pred)
۳۲ accurecy2) ,Accuracy:" print("new #
۳۳
۳۴ return max(accurecy, accurecy2)

```

۶.۰ نتیجه گیری

در نهایت اینجانب دقت‌های بدست آمده در هر ۳ روش قبل را باهم مقایسه کردم که شکل زیر نشان می‌دهد که به ترتیب چه روشی بهتر جواب می‌دهد.

pare the accuracy of Random Forest model with the accuracy of three former



٧.٠ مراجع

datacamp.com

satishgunjal.com

medium.com