

((به نام خدا))

نام دانشجو:

پویان حسابی

شماره دانشجویی:

۹۷۳۱۱۲۲

گزارش پروژه پایانی درس داده کاوی

استاد درس: دکتر ناظر فرد

بهار ۱۴۰۱

"پروژه پایانی درس داده کاوی"

مقدمه

در این پروژه قصد داریم با استفاده از تکنیک های classification و کتابخانه XGBoost مشخص کنیم که فردی با مشخصات خاص، دارای دیابت می باشد یا خیر. مشخصه ها و داده های مربوط به پروژه در فایل قرار دارد که در قالب زیر می باشد، هر کدام از فیلدها در تایپ های عددی، بولین (به شکل ۰ و ۱) و رشته که کتگوریکال است می باشند.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1		Diabetes_1	HighBP	HighChole	Cholesterc	BMI	Smoker	Stroke	HeartDise	PhysicalAc	Fruits	Veggies	HeavyAlcc	AnyHealth	NoDoctor	GeneralH	MentalHe	PhysicalHi	Difficulty	Sex	Age	Education	Income
2	2456	0	0	1	1	24	0	0	0	1	1	0	0	1	1	Very Low	30	30	0	male	1	Cat6	Cat6

این پروژه در ۳ فاز انجام می شود که در ادامه هر کدام از آنها توضیح داده می شوند.

مرحله پیش پردازش

قبل از اینکه پردازش اصلی و ساخت مدل را انجام دهیم، نیاز به پیش پردازش داریم تا دیتای ما به اصطلاح تمیز شود. از طرفی پیش پردازش دقت مدل را به شدت افزایش می دهد و یک گام الزامی می باشد.

در این مرحله ابتدا ستون سمت چپ که اضافه است را پاک می کنیم، سپس تمام white space ها را به کاراکتر " _ " تبدیل می کنیم تا کتابخانه XGBoost به مشکل نخورد. بعد از آن داده هایی که نال یا خالی هستند را پر می کنیم به این روش که اگر عدد یا بولین (صفر یا یک) بود مقدار میانه آن ستون را جایگزین مقدار خالی می کنیم، و اگر داده های categorized باشد مقداری در ستون که بیشترین تکرار را دارد.

بعد از آن داده هایی که به اصطلاح اشتباه می باشند یا خارج از عرف هستند اصلاح می کنیم. به طور مثال تقریباً BMI بالای ۸۰ غیر ممکن است آن را کاهش می دهیم یا اگر فیلدی در بولین ۲ باشد باید تغییر کند.

در گام بعدی داده ها را نرم می کنیم که دسته های کمتری داشته باشند، به طور مثال age را به دسته های ۱۰ تایی تقسیم کرده و BMI را در بازه های استاندارد واقعی BMI جایگزین می کنیم.

مرحله بعد دسته بندی داده های categorical است که با استفاده از one-hot-encoding آنها را به دسته های مختلف تقسیم می کنیم و جای آن دسته ها صفر و یک قرار می گیرد.

در آخر ستون Diabetes_binary را به عنوان labeled_diabetes جدا کرده و از data_frame حذف می کنیم.

"پروژه پایانی درس داده کاوی"

ساخت مدل طبقه بند

در این قسمت با استفاده از کتابخانه XGBoost دسته بندی را انجام می دهیم به این شکل که قسمتی از داده ها برای آموزش و قسمتی دیگر برای تست به نسبت ۷۵ و ۲۵ درصد در نظر می گیریم.

پارامتر های زیر را به صورت ورودی به این تابع می دهیم:

```
XGBClassifier(Learning_rate=0.1, Max_depth=4, N_estimator=200, Subsample=0.5,  
              Colsample_bytree=1, Random_seed=123, Eval_metric='auc', Verbosity=1)
```

بعد از ساخت مدل و استفاده از تابع های predict و fit، مدل را به تابع model_evaluation می دهیم تا دقت، صحت، پوشش و ماتریس درهم ریختگی را نشان دهد.

```
This is for train  
Confusion Matrix: [[7005 1783]  
 [1281 7604]]  
Accuracy: 0.8266281898941888  
precision: 0.8266281898941888  
recall: 0.8454018826937002  
This is for test  
Confusion Matrix: [[18547 8011]  
 [ 6107 20354]]  
Accuracy: 0.7337181010581112  
precision: 0.7337181010581112  
recall: 0.7522917173683784
```

تنظیم هایپر پارامتر

پارامتر های زیادی در ساخت مدل تاثیر گذار هستند، که هر کدام از آنها در عملکرد و سرعت مدل اثرگذار است. در این بخش قصد داریم بهترین مقادیر را برای پارامتر ها بیابیم به این شکل که مدل حداکثر عملکرد را داشته باشد.

در واقع با ترکیب پارامتر های زیر این کار را انجام می دهیم:

```
learning_rates = [0.02, 0.05, 0.1, 0.3], max_depths = [2, 3, 4]
```

```
n_estimators = [100, 200, 300], colsample_bytrees = [0.8, 1]
```

```
hyper_parameters = [learning_rates, max_depths, n_estimators, colsample_bytrees]
```

"پروژه پایانی درس داده کاوی"

سپس مدل و پارامترها را به تابع `GridSearchCV` داده و با تکنیک سه نقطه جداسازی و `scoring='roc_auc'` بهترین پارامترها و دقت، صحت، پوشش و ماتریس درهم ریختگی را نمایش می‌دهیم:

```
GridSearchCV(cv=StratifiedKFold(n_splits=3, random_state=7, shuffle=True),
             estimator=XGBClassifier(base_score=None, booster=None,
                                     callbacks=None, colsample_bylevel=None,
                                     colsample_bynode=None,
                                     colsample_bytree=None,
                                     early_stopping_rounds=None,
                                     enable_categorical=False,
                                     eval_metric='auc', gamma=None, gpu_id=None,
                                     grow_policy=None, importance_type=None,
                                     interaction...,
                                     max_leaves=None, min_child_weight=None,
                                     missing=None, monotone_constraints=None,
                                     n_estimators=100, n_jobs=None,
                                     num_parallel_tree=None, predictor=None,
                                     random_state=None, reg_alpha=None,
                                     reg_lambda=None, ...),
             n_jobs=-1,
             param_grid={'colsample_bytree': [0.8, 1],
                         'learning_rate': [0.02, 0.05, 0.1, 0.3],
                         'max_depth': [2, 3, 4],
                         'n_estimators': [100, 200, 300]},
             scoring='roc_auc')
```

```
Best parameters: {'colsample_bytree': 0.8, 'learning_rate': 0.05, 'max_depth': 4, 'n_estimators': 300}
```

ارزیابی مدل:

```
Best parameters: {'colsample_bytree': 0.8, 'learning_rate': 0.05, 'max_depth': 4, 'n_estimators': 300}
This is for test
Confusion Matrix: [[18491  8067]
 [ 6407 20054]]
Accuracy:  0.7270035270374772
precision: 0.7270035270374772
recall:    0.7426700939834525
This is for train
Confusion Matrix: [[6960 1828]
 [1365 7520]]
Accuracy:  0.8193289198211962
precision: 0.8193289198211962
recall:    0.836036036036036
```

مشاهده می‌شود که مقادیر Accuracy, Precision, Recall همگی در مرحله train نسبت به حالت قبل، به مراتب بهتر شده و عملاً مدل عملکرد بهتری پیدا کرده است.