

https://colab.research.google.com/drive/1OoNTJPnoXJWwPW2QGQEdUeNDVLYqsUxOi#scrollTo=1uCZO7P_0695

گزارش:

مرحله اولیه و بسیار اساسی در فرایند تحلیل داده، EDA یا تحلیل اکتشافی داده است. EDA ابزاری قدرتمندی است که به ما این امکان را می‌دهد تا اطلاعات اولیه از داده‌ها به دست آوریم و الگوها، ارتباطات و نقاط قوت و ضعف داده‌ها را شناسایی کنیم و در ادامه بتوانیم از مدل‌های یادگیری ماشین برای دسته‌بندی بهینه استفاده کنیم.

در ابتدا تمامی کتابخانه‌های مورد نیاز را وارد می‌کنیم. در ادامه نیز هر کجا متوجه شدیم نیاز به وارد کردن کتابخانه‌ی جدید است، آن را نیز اضافه می‌کنیم.

در ادامه، مسیر فایل با فرمت CSV که می‌خواهیم اطلاعات آن را بخوانیم را در متغیر file_path ذخیره می‌کنیم. سپس با استفاده از تابع read_csv از کتابخانه pandas، فایل را می‌خوانیم. در ادامه این داده‌ها را در یک ساختار داده به نام data ذخیره می‌کنیم. این ساختار می‌تواند به عنوان یک جدول دوبعدی از داده‌ها در نظر گرفته شود. سپس تابع head() سطرهای ابتدایی ساختار داده data را نمایش می‌دهد.

در ادامه 3 خط را به طور تصادفی و بدون امکان تکرار چاپ می‌کنیم.

با تابع np.random.choice()، اندیس‌های تصادفی برای انتخاب ردیف‌ها ایجاد می‌شود. پارامتر size تعیین کننده تعداد اندیس‌های تصادفی است و با استفاده از replace=False ردیف‌ها غیرقابل تکرار خواهند بود. سپس با تابع iloc محتوای ردیف‌های انتخاب شده دسترسی پیدا کرده و آن‌ها را چاپ می‌کنیم.

سپس اطلاعات کلی دیتاها را مشاهده می‌کنیم که به طور مثال شامل نوع داده و تعداد داده‌های غیر null در هر ستون می‌باشد.

در بخش بعدی ویژگی‌هایی شامل تعداد، میانگین، انحراف معیار، حداقل، چارک‌ها و حداکثر را بررسی می‌کنیم. به این صورت اطلاعات دقیق‌تری از دیتاست پیدا کرده و میتوانیم حدود داده‌ها و داده‌های پرت را بررسی کنیم.

سپس برای بررسی توزیع در داده هدف از یک نمودار countplot استفاده می‌کنیم، زیرا توازن یا ناتوازن کلاس‌ها می‌تواند بر توانایی مدل در یادگیری و پیش‌بینی تأثیر بگذارد.

این نمودار که تعداد yes و no را در ستون هدف بررسی می کند نشان می دهد تعداد no ها بیشتر است. نمودار هیستوگرام برای ویژگی های مختلف است که اطلاعاتی در مورد توزیع و توازن ویژگی ها در دیتاست ارائه می دهد. هر هیستوگرام نمایانگر توزیع های مختلفی است که در داده وجود دارد.

سپس از نمودار جعبه ای (Box Plot) استفاده می کنیم تا توزیع دما را در زمان مشخصی (ساعت 9 صبح) بررسی کنیم. این نوع نمودار به ما اطلاعات مفیدی از قبیل میانه، چارک ها، و نقاط دورافتاده (outliers) می دهد. در ادامه همین نمودار را برای ساعت 3 بعدازظهر بررسی می کنیم. متوجه میشویم که دما در زمان دوم بیشتر است. (می توان از طریق میانه مقایسه کرد)

در نمودار میله ای (Bar Plot) برای نمایش میانگین دمای حداکثر در ایستگاه های مختلف هواشناسی استفاده می کنیم. این نوع نمودار به ما کمک می کند تا مقایسه ای بین میانگین دماهای مختلف در ایستگاه های مختلف انجام دهیم.

در نمودار بعد، تغییرات دما در بازه های زمانی مختلف را بررسی میکنیم. برای جلوگیری از تراکم نمودار نیز می توان بازه های زمانی ماهانه یا سه ماهه را در نظر گرفت.

سپس به سراغ داده های گمشده می رویم. ابتدا تعداد آن ها را در هر ستون بررسی می کنیم.

سپس ستون های بلا استفاده را حذف می کنیم و اقدام به پر کردن مقادیر گمشده می نماییم.

داده های عددی با میانگین مقدارشان و داده های دسته ای با بیشترین مقدار (mode) تکمیل می شوند. در نهایت، تعداد مقادیر نال در هر ستون چاپ می شود تا اطمینان حاصل شود که داده ها به درستی تکمیل شده اند.

کد بعدی برای تبدیل مقادیر غیر عددی در DataFrame به مقادیر عددی با استفاده از LabelEncoder از کتابخانه scikit-learn است. این داده های دسته ای (categorical) را برای استفاده در مدل های ماشین لرنینگ آماده می کند. بخش های مختلف:

تبدیل ستون تاریخ به نوع تاریخ و تنظیم آن به عنوان ایندکس: این کار برای اطمینان از این است که داده ها به درستی نمونه گیری و تحلیل شوند.

تکمیل داده های ناقص: این فرآیند تضمین می کند که هیچ داده ی ناقصی در DataFrame وجود ندارد.

استفاده از LabelEncoder: این کار داده های دسته ای را به مقادیر عددی تبدیل می کند که برای مدل های یادگیری ماشین مناسب هستند.

سپس چند سطر اول از داده ها نمایش داده می شود تا بتوان تغییر مقادیر را مشاهده کرد.

(در پیوست یک تفاوت دو روش برای تبدیل داده های غیر عددی به داده های عددی بررسی می شود.)

حال به سراغ حذف داده های پرت می رویم و از روش IQR استفاده می کنیم. برای شهود بهتر، یک بار نمودار جعبه ای را قبل از حذف داده های پرت و یک بار بعد از حذف آن ها رسم می کنیم . به وضوح تغییر در داده های پرت مشاهده میشود.

سپس ماتریس کرویشن را رسم می کنیم.

ماتریس همبستگی نشان دهنده ی تاثیر ویژگی های مختلف روی یکدیگر است و باعث درک ارتباط بین آن ها می شود. این اطلاعات می توانند در انتخاب ویژگی ها یا حتی تفسیری از رفتار مدل مفید باشند.

در این ماتریس مقادیر نزدیک به ۱ یا -1 نشان دهنده ارتباط قوی مثبت یا منفی است.

اگر عدد نشان داده شده در هر بخش مربعی، عددی مثبت و نزدیک به 1 باشد ، نشان دهنده یک ارتباط مثبت بزرگ است. یعنی این که دو ویژگی با هم همبستگی بسیار زیادی دارند و اگر که یکی از آنها افزایش پیدا کند، دیگری نیز افزایش پیدا می کند و برعکس.

اگر منفی و نزدیک به -1 باشد ،نشان دهنده یک ارتباط منفی بزرگ است. این به معنای این است که این دو ویژگی نیز با هم همبستگی بسیار زیادی دارند ولی این بار هنگامی که یکی از آنها افزایش می یابد، دیگری کاهش می یابد و برعکس، یعنی مخالف هم هستند.

اگر عدد، مثبت یا منفی و نزدیک به 0 باشد به معنای این است که دو ویژگی با هم همبستگی دارند ولی روی هم تاثیر کمی میگذارند و اگر عدد مثبت باشد تاثیر مستقیم و اگر منفی باشد تاثیرش عکس است.

بدیهی است که در سلول متناظر هر متغیر با خودش عدد 1 را می بینیم، یعنی بیشترین تاثیر گذاری .

در ماتریس اول رسم شده ، چند سطر و ستون کاملاً سفید می بینیم . سه ردیف خالی (سفید) در ماتریس همبستگی به این دلیل است که این ویژگی ها در داده ها به احتمال زیاد شامل مقادیر ثابت زیاد یا دارای واریانس صفر هستند. این ویژگی ها نمی توانند همبستگی مناسبی با سایر ویژگی ها داشته باشند زیرا تنوع داده ها برای محاسبه همبستگی وجود ندارد.

برای حل این مشکل قبل از محاسبه ماتریس همبستگی، میتوان ویژگی هایی که واریانس صفر یا بسیار کم دارند را حذف کرد. این کار را می توان با استفاده از تابع `VarianceThreshold` از کتابخانه `scikit-learn` انجام داد.

بعد از انجام این کار مجدداً ماتریس را رسم می کنیم و ویژگی هایی که همبستگی بالایی با یکدیگر دارند (با تعیین یک آستانه همبستگی) را حذف می کنیم. باید دقت کنیم که در محاسبات، همبستگی یک ستون با خودش در نظر گرفته نشود.

همبستگی بالا بین ویژگی ها می تواند منجر به افزونگی (redundancy) شود و مدل ها را ناپایدار یا کم دقت کند.

استفاده از `set` برای ویژگی های حذف شده تضمین می کند که هر ویژگی تنها یک بار حذف شود.

آستانه‌ی همبستگی یک هاپیرپارامتر است که با توجه به ویژگی‌های مسئله و دیتاست باید تعیین شود و مقدار 0.9 برای همبستگی، مقدار بالایی می‌باشد.

سپس به سراغ پیاده‌سازی مدل‌ها می‌رویم.

ابتدا ستون هدف را جدا می‌کنیم.

سپس برای توازن داده‌های نامتوازن از تکنیک SMOTE (Synthetic Minority Over-sampling Technique) استفاده می‌کنیم.

پارامتر `sampling_strategy='auto'` به این معنی است که نمونه‌برداری برای توازن تمام کلاس‌ها انجام خواهد شد.

در نهایت نیز تعداد نمونه‌های هر کلاس پس از اعمال SMOTE نمایش داده می‌شود. بعد از اجرای این کد، داده‌ها توازن یافته و تعداد نمونه‌های هر کلاس برابر خواهد شد.

همانطور که در ابتدا گفته شد داده‌های نامتوازن مشکلی رایج در بسیاری از مسائل یادگیری ماشین است که می‌تواند عملکرد مدل را تحت تأثیر قرار دهد. SMOTE یکی از تکنیک‌های `oversampling` است که با تولید نمونه‌های مصنوعی از داده‌های اقلیت به توازن داده‌ها کمک می‌کند و باعث بهبود عملکرد مدل‌های یادگیری ماشین می‌شود.

در ادامه داده‌ها را استاندارد می‌کنیم.

استانداردسازی یکی از مراحل پیش‌پردازش مهم در مسائل یادگیری ماشین است. این تکنیک مقادیر ویژگی‌ها را به گونه‌ای تبدیل می‌کند که میانگین آن‌ها برابر با 0 و انحراف معیار آن‌ها برابر با 1 شود. این کار به مدل‌های یادگیری ماشین کمک می‌کند تا عملکرد بهتری داشته باشند، به ویژه برای مدل‌هایی که مبتنی بر گرادیان هستند (مانند رگرسیون لجستیک، شبکه‌های عصبی و ...).

متد `fit_transform` برای اعمال استانداردسازی به داده‌های `X_res` استفاده می‌شود. این متد دو عملیات را به صورت متوالی انجام می‌دهد:

`fit`: محاسبه‌ی میانگین و انحراف معیار برای هر ویژگی در داده‌های `X_res`.

`transform`: استانداردسازی داده‌ها با استفاده از میانگین و انحراف معیار محاسبه شده.

بعد از اجرای این کد، داده‌های `X_res` که قبلاً با استفاده از تکنیک SMOTE توازن یافته‌اند، استانداردسازی می‌شوند و در متغیر `X_res_scaled` ذخیره می‌گردند.

سپس کاهش بعد داده‌ها با استفاده از روش تحلیل مولفه‌های اصلی (PCA) را انجام می‌دهیم.

تحلیل مولفه‌های اصلی (PCA) یک تکنیک کاهش بعد است که به کاهش پیچیدگی داده‌ها و افزایش کارایی مدل‌های یادگیری ماشین کمک می‌کند. این تکنیک با پیدا کردن مولفه‌های اصلی، ابعاد داده‌ها را کاهش می‌دهد و اطلاعات اصلی را حفظ می‌کند. برخی از مزایای استفاده از PCA عبارتند از:

کاهش ابعاد داده‌ها و در نتیجه کاهش زمان محاسباتی مدل‌ها.

حذف نویز و کاهش همبستگی بین ویژگی‌ها.

افزایش دقت مدل‌های یادگیری ماشین با کاهش overfitting.

ابتدا یک شیء PCA ایجاد می‌کنیم و تعداد مولفه‌هایی که باید حفظ شوند را تعیین می‌کنیم. پارامتر `n_components` مشخص می‌کند که داده‌ها به چند مولفه اصلی کاهش داده شوند.

متد `fit_transform` برای اعمال PCA به داده‌های `X_res_scaled` استفاده می‌شود. این متد دو عملیات را به صورت متوالی انجام می‌دهد:

محاسبه‌ی مولفه‌های اصلی از داده‌ها `X_res_scaled`.

تبدیل داده‌ها به فضای مولفه‌های اصلی جدید.

بعد از اجرای این کد، داده‌های استانداردسازی شده‌ی `X_res_scaled` به فضای مولفه‌های اصلی کاهش داده می‌شوند و در متغیر `X_res_pca` ذخیره می‌گردند. این داده‌های کاهش بعد یافته می‌توانند به عنوان ورودی به مدل‌های یادگیری ماشین استفاده شوند تا عملکرد بهتری حاصل شود.

هایپرپارامتر `n_components` در PCA:

`n_components` در PCA تعیین می‌کند که چند مؤلفه اصلی (Principal Components) باید حفظ شوند. انتخاب مقدار مناسب برای این هایپرپارامتر بسیار مهم است، زیرا:

مقادیر کوچک: انتخاب تعداد کمی مؤلفه ممکن است به از دست رفتن اطلاعات مهم و کاهش عملکرد مدل منجر شود.

مقادیر بزرگ: انتخاب تعداد زیادی مؤلفه ممکن است منجر به حفظ نویز و افزایش پیچیدگی مدل شود، که می‌تواند به overfitting منجر شود.

روش‌های مختلفی برای تنظیم `n_components` وجود دارد:

انتخاب یک عدد ثابت: مانند 15. این روش ساده است ولی ممکن است بهترین نتیجه را ندهد.

نسبت واریانس: انتخاب `n_components` به گونه‌ای که یک درصد مشخصی از واریانس کل داده‌ها را توضیح دهد.

جستجوی هایپرپارامتر :استفاده از روش هایی مانند جستجوی شبکه ای یا جستجوی تصادفی برای پیدا کردن بهترین مقدار.

به عنوان اولین مدل ، داده ها را به مجموعه های آموزشی، اعتبارسنجی و تست تقسیم می کنیم و یک مدل SVM (Support Vector Machine) را پیاده سازی می کنیم.

ابتدا داده ها به دو مجموعه تقسیم می شوند: 70٪ برای آموزش و 30٪ برای مجموعه موقت. (X_{temp} , y_{temp}) سپس مجموعه موقت به دو قسمت مساوی تقسیم می شود: 15٪ برای اعتبارسنجی و 15٪ برای تست.

یک مدل SVM با کرنل خطی ایجاد و آموزش داده می شود.

انتخاب کرنل خطی برای SVM (Support Vector Machine) بستگی به ماهیت داده ها و مسئله دارد.

سادگی و سرعت

کرنل خطی نسبت به کرنل های پیچیده تر مانند RBF یا پلی نومیک، ساده تر و سریع تر است.

داده های قابل جداسازی خطی

اگر داده ها به طور تقریبی به صورت خطی قابل جداسازی باشند، کرنل خطی می تواند عملکرد بسیار خوبی داشته باشد. در این صورت نیازی به استفاده از کرنل های غیرخطی نیست.

تفسیرپذیری

مدل های خطی تفسیرپذیری بهتری دارند. وزن های مدل خطی مستقیماً نشان می دهند که کدام ویژگی ها در تصمیم گیری مهم تر هستند.

ابعاد بالا

در مسائل با تعداد ویژگی های بسیار زیاد (به عنوان مثال، در مسائل متن کاوی و ژنومیک)، کرنل خطی معمولاً عملکرد خوبی دارد. این به دلیل خاصیت کرنل خطی در فضاها با ابعاد بالا است که باعث می شود داده ها به طور طبیعی جدا شوند.

جلوگیری از overfitting

زمانی که داده های آموزشی کم باشند استفاده از کرنل های پیچیده تر می تواند منجر به overfitting شود.

تقسیم بندی داده ها به سه بخش برای جلوگیری از overfitting میتواند مفید باشد. با استفاده از مجموعه اعتبارسنجی می توان مدل را بهینه سازی کرد و با استفاده از مجموعه تست عملکرد نهایی مدل را ارزیابی نمود.

سپس ارزیابی مدل را انجام می‌دهیم . به این منظور ماتریس درهم ریختگی و معیارهای ارزیابی (accuracy, precision, recall, F1) را با استفاده از کتابخانه های آماده محاسبه میکنیم.

در معیار های ارزیابی به دو عبارت زیر برمیخوریم:

macro avg: میانگین دقت، بازخوانی و F1-Score برای هر کلاس به صورت جداگانه محاسبه شده و سپس میانگین گرفته شده است. این معیار برابر با میانگین اعداد به صورت میانگین وزن دار نیست.

weighted avg: مشابه میانگین ماکرو است، اما وزن دهی به این میانگین ها بر اساس تعداد نمونه ها انجام می شود. این معیار نشان می دهد که مدل چقدر کلاس های مختلف در مجموعه داده را درست تشخیص داده است.

KNN

برای به کار بردن KNN باید ابتدا مقدار بهینه k را بدست آوریم. در این پروژه اعداد 1 تا 10 را برای k انتخاب می کنیم و هر کدام از k ها که دقت مدل در آن بیشتر بود آن مقدار را در نظر می گیریم. برای این کار از تابع cross_val_score استفاده می کنیم که داده ها و برچسب ها را می گیرد و آرایه ای از دقت ها را با توجه به انتخاب های مختلف از داده های تست و ترین خروجی می دهد. در پارامتر cv=k_folds تعداد بخش های که داده ها به آنها تقسیم می شوند را نشان می دهیم که ما عدد 5 را انتخاب کردیم. وقتی این پارامتر 5 باشد این تابع داده ها را به 5 قسمت مساوی تقسیم می کند و سپس هر بار یکی از این 5 قسمت را به عنوان داده های تست در نظر می گیرد و بقیه را به عنوان داده های ترین در نظر می گیرد سپس دقت های بدست آمده را در یک آرایه خروجی می دهد.

```
k_folds = 5

# که می‌خواهیم امتحان کنیم K لیستی از مقادیر برای پارامتر
k_values = list(range(1, 10))

# K لیستی برای نگهداری دقت‌های هر مقدار
accuracies = []

# انجام اعتبارسنجی متقاطع برای هر مقدار
for k in k_values:
    knn = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn, X_res_scaled, y_res, cv=k_folds, scoring='accuracy')
    accuracies.append(scores.mean())
```

در هر تکرار از این for میانگین دقت های بدست آمده را به آرایه accuracies اضافه می کنیم.

```
# بر اساس دقت میانگین K یافتن بهترین مقدار
best_k = k_values[np.argmax(accuracies)]
best_accuracy = max(accuracies)

print('K بهترین مقدار: ', best_k)
print('دقت متناظر: {:.2f}%'.format(best_accuracy * 100))
```

سپس بهترین مقدار k همراه با دقت حاصل آن را بدست می آوریم.

```
X_train, X_test, y_train, y_test = train_test_split(X_res_scaled, y_res, test_size=0.2, random_state=42)
Knn = KNeighborsClassifier(n_neighbors=2)
knn.fit(X_train, y_train)
```

سپس همانطور که در تصویر مشاهده می کنیم داده های متوازن و استاندارد شده را به داده های تست و ترین با ضریب 0.2 تقسیم می کنیم و مدل را با بهترین k (2) ترین می کنیم.

```
accuracy = accuracy_score(y_test, y_pred)
print('دقت مدل: {:.2f}%'.format(accuracy * 100))

print("Train Classification Report:")
print(classification_report(y_test, y_pred))

conf_matrix = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
plt.xlabel('predicted')
plt.ylabel('actual')
plt.title('confusion matrix')
plt.show()
```

سپس دقت مدل بدست آمده را چاپ می کنیم و همچنین معیار های precision و recall و f1-score را با استفاده از تابع classification_report که داده های y_test و y_pred که توسط مدل بدست آمده است را با هم مقایسه می کند و سپس ماتریس سردرگمی را برای داده های تست رسم می کنیم. دقت در این مدل 84 درصد است و معیار های recall و precision و f1-score به ترتیب برای کلاس 1 94 و 79 و 86 است که بالا بودن معیار recall قابل توجه است.

درخت تصمیم

```
dt = DecisionTreeClassifier()

# آموزش مدل بر روی داده های آموزشی
dt.fit(X_train, y_train)

# پیش‌بینی بر روی داده های آزمون
y_pred = dt.predict(X_test)

# محاسبه دقت مدل
accuracy = accuracy_score(y_test, y_pred)
print('دقت مدل درخت تصمیم: {:.2f}%'.format(accuracy * 100))
```

یک مدل درخت تصمیم با پرامترهای پیش فرض را در نظر می گیریم. البته درخت تصمیم با پرامتر $\text{min_samples_split} = 10, 20, 100$ و $\text{min_samples_leaf} = 10, 20, 100$ را امتحان کردیم که در همه حالت ها از حالت پیش فرض دقت کاهش می یابد پس همان حالت پیش فرض را در نظر گرفتیم. در ادامه نیز معیارهای recall و precision و f1-score را برای هر دو کلاس چاپ می کنیم را برای کلاس 1 به ترتیب 88 و 86 و 87 خروجی داده می شود. می بینیم که با توجه به KNN این مدل تقریباً f1-score مشابه ای دارد اما این در این مدل هر دو معیار precision و recall متعادل تر هستند و نزدیک به هم هستند. البته هر دو مدل با توجه به svm نتایج بهتری را کسب می کنند(شاید به دلیل اینکه در svm داده های اعتبارسنجی هم داریم)

پیوست 1-

روش‌هایی مانند Label Encoding و One-Hot Encoding برای تبدیل مقادیر دسته‌بندی (categorical) به مقادیر عددی استفاده می‌شوند. با این حال، این دو روش در نحوه تبدیل داده‌ها و کاربردهایشان تفاوت دارند.

Label Encoding

مقادیر دسته‌بندی را به مقادیر عددی تبدیل می‌کند. در این روش، هر دسته به یک عدد منحصر به فرد نگاشت می‌شود. به عنوان مثال، دسته‌ها به این صورت ممکن است نگاشت شوند:

Red -> 0

Green -> 1

Blue -> 2

این روش برای مدل‌هایی که می‌توانند با مقادیر ترتیبی کار کنند، مناسب است. اما در برخی موارد ممکن است باعث بروز مشکل شود، زیرا مدل ممکن است رابطه ترتیبی بین دسته‌ها را فرض کند (مثلاً فرض کند که 2 بزرگ‌تر از 1 است).

One-Hot Encoding

هر دسته را به یک بردار باینری تبدیل می‌کند. در این روش، هر دسته به یک بردار که فقط یک عنصر آن برابر با 1 و بقیه عناصر برابر با 0 هستند، نگاشت می‌شود. به عنوان مثال:

Red -> [1, 0, 0]

Green -> [0, 1, 0]

Blue -> [0, 0, 1]

این روش برای مدل‌هایی مناسب است که با مقادیر ترتیبی کار نمی‌کنند و نمی‌خواهید مدل رابطه ترتیبی بین دسته‌ها فرض کند.

نتیجه‌گیری:

Label Encoding برای مدل‌هایی که با مقادیر ترتیبی به درستی کار می‌کنند مانند درخت‌های تصمیم، جنگل‌های تصادفی و GBM مناسب است.

One-Hot Encoding برای مدل‌هایی که به مقادیر ترتیبی حساس هستند مانند رگرسیون لجستیک، SVM، KNN، Naive Bayes و شبکه‌های عصبی مناسب است.