



به نام خدا



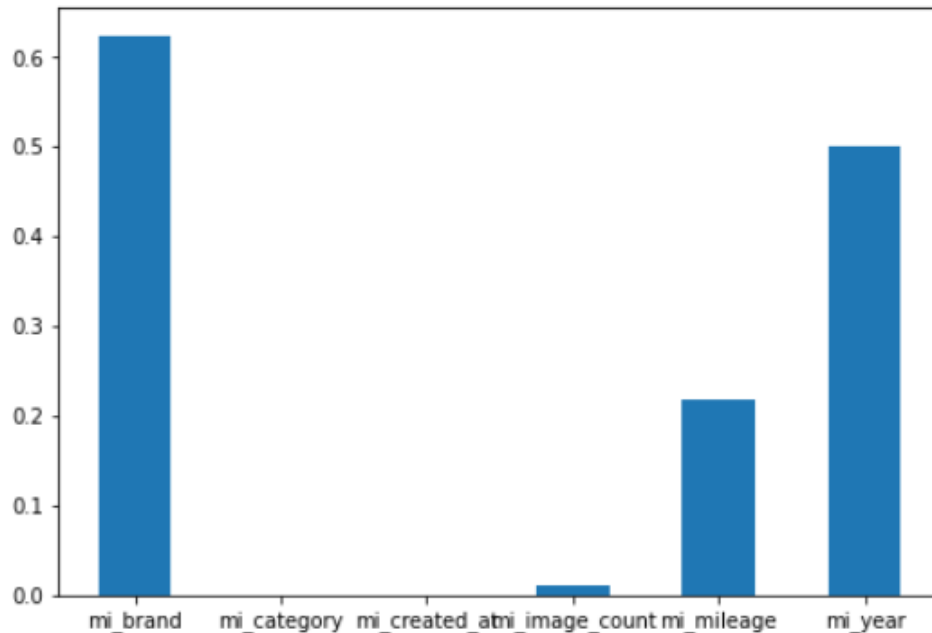
دانشگاه تهران
دانشکده مهندسی برق و کامپیوتر
هوش مصنوعی

تمرین کامپیوتری چهارم

نام و نام خانوادگی	مهسا تاجیک
شماره دانشجویی	810198126
تاریخ ارسال گزارش	14 خرداد

کاهش دهیم و در زمان و حافظه‌ی مصرفی صرفه جویی کنیم. در واقع وابستگی بین هر دو ویژگی را می‌توان محاسبه کرد و اگر مقدارش صفر شد یعنی دو ویژگی کاملاً از هم مستقل هستند

[0.62419425 0. 0. 0.01149709 0.21863502 0.5001221]



پرسش ۲) با دستور data.dtypes تایپ ستون‌ها را می‌توانیم مشاهده کنیم :

```
brand          object
category       object
created_at     object
description     object
image_count    int64
mileage        float64
price          int64
title          object
year           object
```

ستون‌هایی که تایپ object دارند می‌توانیم ابتدا به تایپ category تبدیل کرده سپس با دستور cat.codes به مقادیر عددی کد کنیم :

```
data['brand'] = data['brand'].astype('category').cat.codes
data['category'] = data['category'].astype('category').cat.codes
data['year'] = data['year'].astype('category').cat.codes
data['created_at'] = data['created_at'].map(lambda x: x[: -5])
data['created_at'] = data['created_at'].astype('category').cat.codes
```

پرسش ۳) برای شمارش تعداد کلمات در ستون‌های متنی title , description از countVectorizer استفاده کردیم که متدهایی بنام min_df و max_df دارد که با تنظیم آن‌ها میتوان کلماتی که تعداد

تکرار خیلی کم یا تعداد تکرار زیادی دارند و در اکثر داکيومنت ها تکرار شدند یا همان stop word ها را حذف کرد. همچنین متدی بنام ngram_range دارد که می توانیم مشخص کنیم که چند کلمه پشت هم را در نظر بگیرد چون همانطور که میدانیم بعضی کلمات در جاهای مختلف معانی متفاوتی دارند زمانیکه کلمات قبل یا بعد آن ها را در نظر میگیریم به فهم معنی کلمه در آن جایگاه کمک می کند. بعد از چندین بار تغییر این مقادیر، برای مقادیر زیر به جواب بهینه رسیدیم و آن ها را فیکس کردیم و در کل 36 ویژگی خواهیم داشت:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2
3 vectorizer1 = CountVectorizer(min_df=0.1,max_df=0.9,ngram_range=(1,3))
4 desc_tokens = vectorizer1.fit_transform(data_description_arr)
5
6 vectorizer2 = CountVectorizer(min_df=0.1,max_df=0.9,ngram_range=(1,3))
7 title_tokens = vectorizer2.fit_transform(data_title_arr)
8
9 desc_df = pd.DataFrame(desc_tokens.toarray(),columns=vectorizer1.get_feature_names())
10 title_df = pd.DataFrame(title_tokens.toarray(),columns=vectorizer2.get_feature_names())
11
12 data = pd.concat([data,desc_df,title_df])
13 data = data.fillna(0)
14 data
```

پرسش ۴) یکی از روش های برخورد با مقادیر از دست رفته ، حذف آن هاست که ما هم از همین روش استفاده کردیم و سطرهایی که حتی یک مقدار nan در آن ها بود ، حذف کردیم. مزیت این روش این است که یک مدل robust داریم اما با حذف این مقادیر ، اطلاعاتی را از دست خواهیم داد و این روش زمانی بهتر است استفاده شود که سائز دیتاست به نسبت مقادیر از دست رفته به اندازه کافی بزرگ باشد.

روش دیگر جایگزین کردن یک مقدار ثابت بجای مقادیر از دست رفته است مثلاً می توانیم تمام مقادیر nan را با 0 ، 1 یا -1 پر کنیم. یا همچنین میتوانیم با میانگین یا میانه ی آن سطر پر کنیم یا اینکه روش روش های backwardfill یا forwardfill استفاده کنیم و با مقدار قبلی یا بعدی در سطر یا ستون ، آن خانه از دیتافریم را پر کنیم.

فاز دوم : پیش بینی قیمت

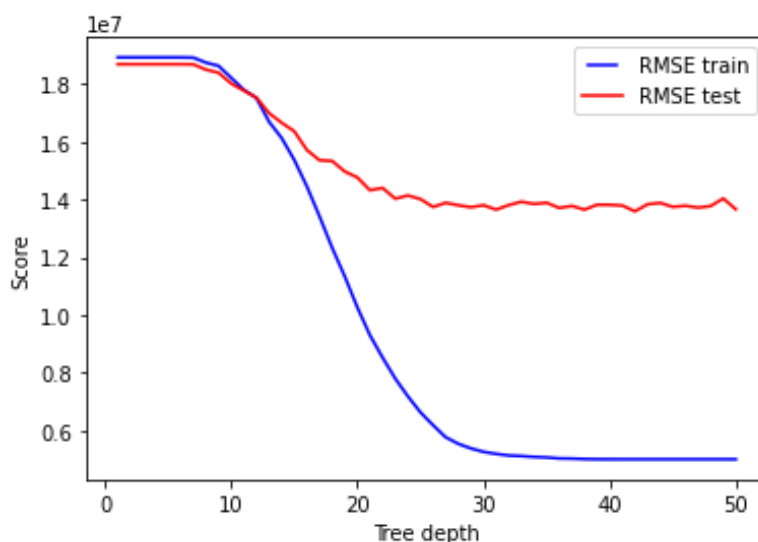
پارامتر n_neighbors مشخص می کند که برای طبقه بندی هر دیتا چندتا همسایه اطراف آن را باید برای تصمیم گیری در مورد کلاس این دیتا در نظر بگیریم که برای مقادیر 10 و 15 و 20 مقدار خطای rmse بصورت زیر است :

KNN

n_neighbors	RMSE Score
20	17415425.85978448
15	17193196.469603904
10	17045764.18711257

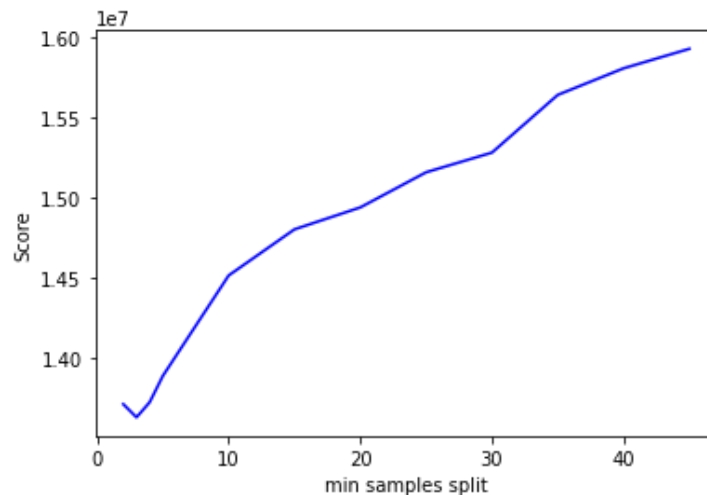
DECISION TREE

- مقدار پارامتر `min_samples_split` را در حالت دیفالت یعنی 2 قرار دادیم و مقدار پارامتر `max_depth` را از 1 تا 50 تغییر دادیم و مقادیر خطای RMSE بصورت زیر می باشد. کمترین مقدار خطای RMSE حدودا برابر است با **13,600,000**



با دستور `tree._max_depth` حالت دیفالت را برای پارامتر `max_depth` میتوان پیدا کرد که برای درخت ما این مقدار برابر است با 41.

- مقدار `max_depth` را روی 42 فیکس میکنیم و مقدار پارامتر `min_samples_split` را برای مقادیر `[2, 3, 4, 5, 10, 15, 20, 25, 30, 35, 40, 45]` تغییر می دهیم و مقادیر خطای `mse, rmse` را محاسبه میکنیم که تغییرات `rmse` بصورت زیر است:



بنابراین هرچه مقدار پارامتر `min_samples_split` کمتر باشد خطای کمتری داریم و کمترین مقدار در 3 اتفاق می افتد.

بصورت کلی که با تغییر مقادیر پارامترها خطای RMSE از حدود 13 میلیون تا 18 میلیون تغییر داشت. (در بعضی اجراها به خطای حدود 6 میلیون هم رسیدیم ولی در اجرای بعدی این مقدار مجدداً افزایش پیدا کرد).

Linear regression

خطای این مدل حدوداً 13 میلیون شد:

```
1 from sklearn.linear_model import LinearRegression
2 lr = LinearRegression().fit(X_train, y_train)
3 coef = lr.coef_
4 intercept = lr.intercept_
5 y_pred = lr.predict(X_test)
6 lr.score(X_train, y_train)
7 mse = mean_squared_error(y_test, y_pred)
8 rmse = mean_squared_error(y_test, y_pred, squared=False)
9 rmse
```

13735635.807216046

پرسش ۵) به طور معمول 15 تا 20 درصد داده‌ها برای تست و بقیه برای آموزش در نظر گرفته می شود.

تقسیم بندی مختلف داده های آموزش و تست برای **Linear Regression** :

نتیجه برای 80 درصد داده ها برای آموزش و 20 درصد برای تست :

12845970.622448027

نتیجه برای 98 درصد داده ها برای آموزش و 2 درصد برای تست :

13735635.807216046

نتیجه برای 40 درصد داده ها برای آموزش و 60 درصد برای تست :

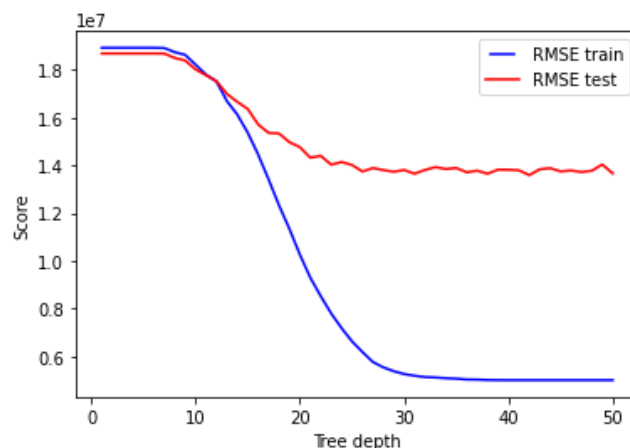
13860931.545627901

زمانیکه دیتای کافی برای آموزش نداریم و 40 درصد را برای آموزش در نظر میگیریم مدل به خوبی یاد نمیگیرد و underfitting اتفاق می افتد زمانیکه که 98 درصد داده ها را برای آموزش و دو درصد را برای تست میگیریم ، مدل بیش از حد آموزش میبندد و قابلیت generalization را از دست می دهد و overfitting اتفاق می افتد در هر دو حالت خطا بیشتر از زمانیست که 80 درصد برای آموزش و 20 درصد برای تست در نظر گرفتیم.

پرسش ۶)

بیشترین عمق ، برای decision tree بین 1 تا 50 تغییر دادیم. نمودار خطای داده های آموزش و تست بصورت زیر است:

Decision tree



فاز سوم : استفاده از مدل های تجمیعی

پرسش ۷) مقادیر خطای ترین و تست برای random forest به ازای max_depth های 2 ، 5 و 10 بصورت زیر است:

```
Rmse_train Score: 18925575.094937544  
Rmse_test Score: 18689445.033343326  
Rmse_train Score: 18925575.230869036  
Rmse_test Score: 18689445.172420606  
Rmse_train Score: 18922040.329567824  
Rmse_test Score: 18689379.429148108
```