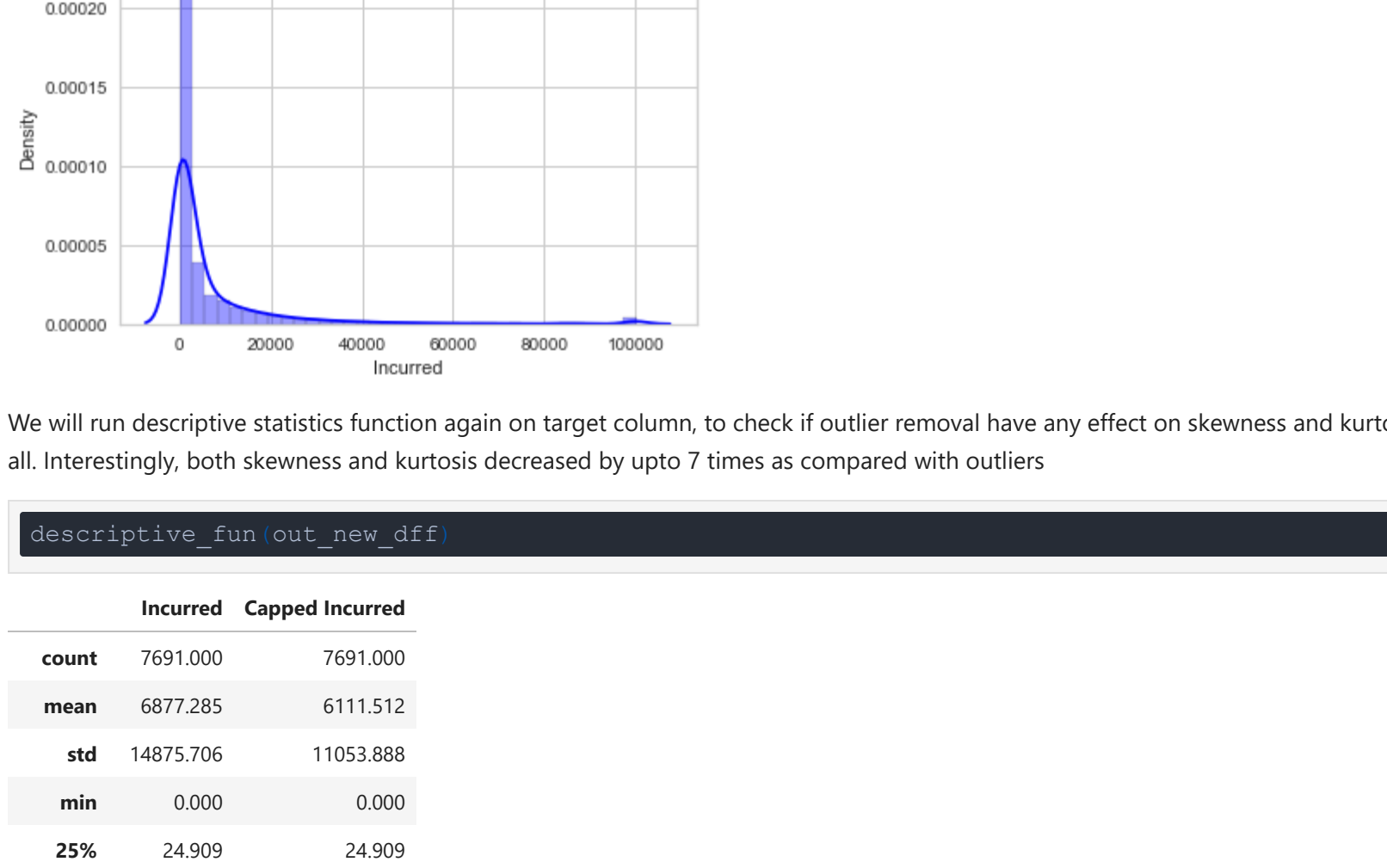


descriptive_fun() will be removed in a future version. Please adapt your code to use either descriptive_fun() with similar flexibility or histplot (an easy-to-use function for histograms).



We will run descriptive statistics function again on target column, to check if outlier removal have any effect on skewness and kurtosis at all. Interestingly, both skewness and kurtosis decreased by upto 7 times as compared with outliers

```
In [92]: descriptive_fun(out_new_def)
```

	Incurred	Capped Incurred
count	7691.000	7691.000
mean	6877.285	6115.152
std	14875.706	11053.888
min	0.000	0.000
25%	24.909	24.909
50%	1237.714	1237.714
75%	6258.404	6258.404
max	100110.251	50000.000
unique	5595.000	5474.000
zero%	23.430	23.430
skew	3.995	2.559
kurtosis	18.806	6.302

```
In [93]: # Adding back columns to main dataframe
excel_def.drop(outlier_columns, axis=, inplace=)
excel_def=pd.merge(excel_def, out_new_def, left_index=, right_index=)
```

Feature Engineering

Since no given column has strong correlation with target column therefore we will create some new features by combining already present features

- Adding Flag Columns: All TP columns are count therefore adding flag for all TP columns where count is greater than 0
- Date Columns - In algorithm, dates doesn't have any effect on its own therefore converting date_of_loss column to 'Month' and 'Weekend' flag
- Range reduction - Converted Inception, to loss to Months Since Loss (MSL), this kind of techniques is used to reduce spread of data
- Injury Flag - If count is greater than 0 for 'Wiplash', 'Traumatic', 'Fatality' or 'Unclear'
- Insured Passenger Injury - If count is greater than 0 for insured passenger at back and Injury flag is 1
- Combined Categories - Some of the categories are combined to create missing category (like n/k and n/a)

```
In [94]: flag_int_def, col_list =

    binarizer = Binarizer(threshold=, copy =)

    columns_in_col_list =
        df[columns_flag] = binarizer.fit_transform(df[columns].values.reshape(-1,))

    return df
```

```
In [95]: tp_cols = col_col_in_excel_def[ col.startwith('TP,') and excel_def.col.nunique() > ]

new_def = flag_int[excel_def, tp_cols]
```

```
In [ ]: new_def['Month'] = pd.DatetimeIndex(new_def['date_of_loss']).month
new_def['Weekend'] = np.where(pd.to_datetime(new_def['date_of_loss']).dt.dayofweek > 6, 1, 0)
new_def['MSL'] = round(new_def['Inception.to_loss']/, 1)

new_def['TP_injury_flag'] = np.where(new_def['TP_injury_wiplash'] > 0 |
new_def['TP_injury_traumatic'] > 0 | new_def['TP_injury_fatality'] > 0 |
new_def['TP_injury_unclear'] > 0, 1, 0)
new_def['TP_insured_pass_injury'] = np.where(new_def['TP_type_insured_pass_back'] > 0 &
new_def['TP_injury_flag'] > 0, 1, 0)

new_def['Vehicle_mobile'] = np.where(new_def['Vehicle_mobile'].str.lower().isin(['n/k']) ==
new, 'missing', new_def['Vehicle_mobile'].str.lower())
new_def['TP_considered_TP_at_fault'] =
np.where(new_def['TP_considered_TP_at_fault'].str.lower().isin(['n/k', 'n']), new, 'missing')
new_def['TP_considered_TP_at_fault'].str.lower())
new_def['Location_of_incident'] =
np.where(new_def['Location_of_incident'].str.lower().isin(['n/k', 'not applicable']) == new,
'missing', new_def['Location_of_incident'].str.lower())
new_def['Weather_conditions'] = np.where(new_def['Weather_conditions'].str.lower() ==
'know,low, fog', 'chilly', new_def['Weather_conditions'].str.lower())
new_def['Weather_conditions'] = np.where(new_def['Weather_conditions'].isin(['n/a', 'n/k']) |
new_def['Weather_conditions'].isnull() == new, 'missing', new_def['Weather_conditions'])
```

```
In [97]: new_def.drop(['date_of_loss', 'date_of_loss', 'date_of_loss', 'TP_type_insured_pass_back'],
['TP_type_insured_pass_back'], axis=, inplace=new, errors='ignore')
```

```
In [99]: print('Number of columns in dataset before feature engineering:', excel_def.shape[1])
print('Number of columns in dataset after feature engineering:', new_def.shape[1])
```

```
In [100]: new_def.head()
```

	Claim Number	Notifier	Notification_period	Inception.to_loss	Location_of_incident	Weather_conditions	Vehicle_mobile	Time_hour	Main_driver
0	1	PH	22	13	main road	normal	y	10	Other
1	2	CNF	1	9	main road	wet	y	18	Other
2	3	CNF	5	17	main road	wet	y	16	Y
3	4	CNF	1	23	main road	missing	y	14	Other
4	5	CNF	1	48	other	missing	n	9	Other
5	6	PH	16	23	other	missing	n	0	Other
6	7	Other	5	4	main road	normal	n	18	Other
7	8	Other	0	40	main road	wet	y	7	Other
8	9	PH	4	26	minor road	wet	y	22	Other
9	10	CNF	2	85	main road	missing	n	22	Other
10	11	Other	0	109	minor road	wet	y	7	Other
11	12	Other	1	22	main road	normal	n	23	Other
12	13	Other	0	46	main road	normal	y	11	Other
13	14	CNF	2	7	other	normal	y	16	Y
14	15	PH	0	57	minor road	normal	n	18	Other
15	16	Other	6	104	minor road	wet	y	10	Other
16	17	TP	61	55	missing	missing	missing	0	Other
17	18	PH	2	27	main road	normal	n	13	Other
18	19	PH	0	128	minor road	wet	n	8	Other
19	20	CNF	0	18	other	normal	n	11	Y

Partition dataset for Train and Test split

StratifiedShuffleSplit from sklearn is used to get similar distribution of train and test set

```
In [101]: min_value = descriptive_fun(pd.DataFrame(new_def['Incurred']).loc['min', :]).min()
max_value = descriptive_fun(pd.DataFrame(new_def['Incurred']).loc['max', :]).max()
print(min_value)
print(max_value)

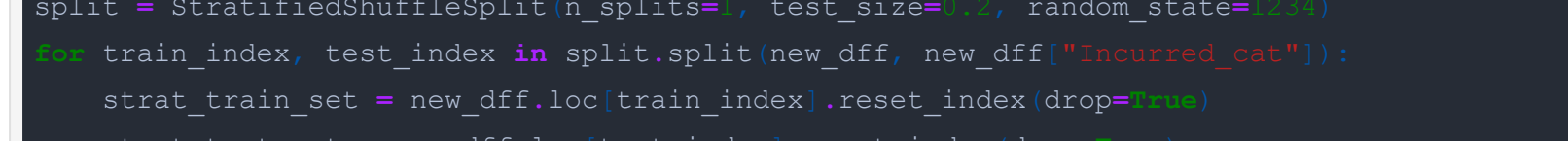
bins = np.linspace(min_value, max_value, )
print(bins)
```



```
In [102]: new_def['Incurred_cat'] = pd.cut(new_def['Incurred'],
bins = [-np.inf, , , , , , np.inf],
labels = range(, ))

print(new_def['Incurred_cat'].hist())

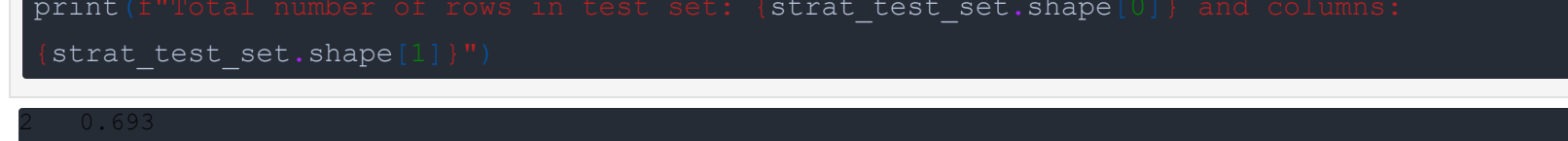
print(new_def['Incurred_cat'].value_counts()/len(new_def))
```



```
In [103]: # Split data using "Incurred_cat" column defined above
split = StratifiedShuffleSplit(n_splits=, test_size=, random_state=)
train_index, test_index = split.split(new_def, new_def['Incurred_cat'])
strat_train_set = new_def.loc[train_index, :].reset_index(drop=)
strat_test_set = new_def.loc[test_index, :].reset_index(drop=)

print(strat_test_set['Incurred_cat'].value_counts()/len(strat_test_set))

print(f"Total number of rows in train set: {strat_train_set.shape[0]} and columns: {strat_train_set.shape[1]}")
print(f"Total number of rows in test set: {strat_test_set.shape[0]} and columns: {strat_test_set.shape[1]}")
```



```
In [104]: # Separating features, target and ID column for train set
train_features = strat_train_set.drop(['Incurred', 'Claim Number', 'Incurred_cat'], axis=)
train_target = strat_train_set['Incurred']

# Drop ID Column and save it for future use
train_id_col = strat_train_set['Claim Number']
```

```
In [105]: # Separating features, target and ID column for test set
test_features = strat_test_set.drop(['Incurred', 'Claim Number', 'Incurred_cat'], axis=)
test_target = strat_test_set['Incurred']

# Drop ID Column and save it for future use
test_id_col = strat_test_set['Claim Number']
```

Separate Numeric and Categorical features

```
In [106]: # Integer continuous columns
numeric = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
int_list = train_features.select_dtypes(include=numeric).columns
int_cols = col_col_in_train_features(int_list)
int_cols = int_cols - ['TP_type_cycles', 'TP_type_pedestrian']

int_flag = col_col_in_train_features(int_list)
int_flag_remove = ['TP_type_cycles', 'TP_type_pedestrian']
int_flag = list(set(int_flag) - set(int_flag_remove))

# Categorical columns
cat_list = train_features.select_dtypes(exclude=numeric).columns
cat_cols = col_col_in_train_features(cat_list)
```

```
In [108]: descriptive_fun(train_features)
```

	Notification_period	Inception.to_loss	Time_hour	Vehicle_registration_present	Incident_details_present	Injury_details_present	TP_type
count	6152.000	6152.000	6152.000	6152.000	6152.000	6152.000	6152.000
mean	7.170	166.799	12.708	0.999	0.899	0.233	
std	37.255	104.246	5.122	0.025	0.393	0.423	
min	-18.000	0.000	0.000	0.000	0.000	0.000	
25%	0.000	75.000	9.000	1.000	1.000	0.000	
50%	1.000	161.000	13.000	1.000	1.000	0.000	
75%	2.000	252.250	17.000	1.000	1.000	0.000	
max	925.000	365.000	23.000	1.000	1.000	1.000	
unique	174.000	366.000	24.000	2.000	2.000	2.000	
zero%	44.262	0.341	4.698	0.065	19.083	76.658	
skew	12.405	0.175	-0.493	-39.189	-15.574	1.261	
kurtosis	207.285	-1.150	-0.056	1534.248	0.477	-0.411	

Target column scaling validation

```
In [113]: power = PowerTransformer(method='yeo-johnson', standardize=)
pipeline_iso = Pipeline(steps=(('p', power),))

data_trans = pipeline_iso.fit_transform(train_target.values.reshape(-1,))
```

```
In [114]: descriptive_fun(pd.DataFrame(data_trans))
```

	0
count	6152.000
mean	0.000
std	1.000
min	-1.458
25%	-0.848
50%	0.170
75%	0.759
max	2.043
unique	4489.000
zero%	0.000
skew	-0.170
kurtosis	-1.071

```
In [115]: # seaborn histogram
sns.distplot(pd.DataFrame(data_trans, columns = ['Incurred']).hist=, kde=,
bins=int (/), color = 'blue',
hist_kws={'edgecolor':'black'})

# Add labels
plt.title('Histogram of Incurred')
```

descriptive_fun() will be removed in a future version. Please adapt your code to use either descriptive_fun() with similar flexibility or histplot (an easy-to-use function for histograms).



Data Transformation

From the prior analysis we found that Incurred target feature is skewed and not show uniform distribution. To deal with this, we can transform data to Power Scale using PowerTransformer method provided by sklearn using yeo-johnson method which can scale negative and zero values. As we observed in above plot, skewness and kurtosis reduce drastically after transformation and now data looks more like normal distribution (Note: peak at beginning may suggest another distribution however with over data analysis we know that this is due to high number of 0 values so we can safely ignore it in first iteration).

For features we again use PowerTransformer method but before that we scaled data with StandardScaler method. Categorical Features are converted into binary One Hot Encoded features. Since we don't have too many unique categories

```
In [188]: X = train_features.copy()
y = train_target.copy()

# Hot encode categorical features
OH_encoder = OneHotEncoder(handle_unknown='ignore', sparse=)
OH_encoder.fit(X[cat_cols])
cat_features = pd.DataFrame(OH_encoder.transform(X[cat_cols])
cat_features.columns = OH_encoder.get_feature_names(cat_cols)

# Numerical continuous features
num_features = X[int_cols]
scaler = StandardScaler()
power = PowerTransformer(method='yeo-johnson', standardize=)
num_pipeline = Pipeline(steps=(('s', scaler), ('p', power)))
num_pipeline.fit(num_features)

num_trans_features = pd.DataFrame(num_pipeline.transform(num_features))

num_trans_features.columns = num_features.columns

# Numerical flag features
flag_features = X[int_flag]

data_df = pd.concat([num_trans_features, cat_features, flag_features], axis=)
ignore_index=, reset_index(drop=)

X_new = data_df.copy()

# scaler = MinMaxScaler()
power = PowerTransformer(method='yeo-johnson', standardize=)
target_pipeline = Pipeline(steps=(('p', power),))
target_pipeline.fit(y.values.reshape(-1,))

y_new = pd.DataFrame(target_pipeline.transform(y.values.reshape(-1,)))

print(f"Feature space increased from {X.shape[1]} to {X_new.shape[1]} after data transformation")
```

```
In [189]: # Now all features seems to follow normal distribution with mean 0 and standard deviation close to 1
descriptive_fun(X_new)
```

	Notification_period	Inception.to_loss	Time_hour	TP_type_insured_pass_back	TP_type_driver	TP_type_pass_back	TP_type_pass_front	TP_type
count	6152.000	6152.000	6152.000	6152.000	6152.000	6152.000	6152.000	
mean	0.000	-0.000	0.000	0.000	0.000	-0.000	0.000	
std	1.000	1.000	1.000	1.000	1.000	1.000	1.000	
min	-26.447	-1.701	-2.133	-0.146	-1.184	-0.177	-0.246	
25%	-0.644	-0.876	-0.782	-0.146	-1.184	-0.177	-0.246	
50%	-0.222	0.000	-0.049	-0.146	0.643	-0.177	-0.246	
75%	0.148	0.845	0.838	-0.146	0.643	-0.177	-0.246	
max	2.759	1.783	2.304	6.866	7.858	5.664	4.072	
unique	174.000	366.000	24.000	2.000	6.000	3.000	3.000	
zero%	0.000	0.000	0.000	0.000	0.000	0.000	0.000	
skew	-1.541	0.037	-0.032	6.716	0.309	5.488	3.827	
kurtosis	-80.564	-1.191	-0.454	43.120	1.307	28.130	12.649	

Model Benchmarking

Benchmark has been established on LinearRegressor model from sklearn. Training data is feed through LinearRegressor (with feature input and target) without much data preprocessing. Some of the basic preprocessing steps done on data are listed below.

- Removed feature "Claim Number" from train and validation data as it doesn't provide any value for identification. Its like index attribute to uniquely identify each claim row.
- Categorical variables are one hot encoded so that numerical values are given as input to model.

After fitting through LinearRegressor (considering all default parameters), prediction is done on test data to compute MAE.

```
In [125]: # Benchmark model, Running Linear Regression on features data

# Split train and test data in ratio of 9:1
X_train, X_dev, y_train, y_dev = train_test_split(X_new, y_new, test_size=, random_state=)

from sklearn.metrics import mean_absolute_error

# Run linear regression on input train and validation data, Compute Mean absolute error.
regr_fit = LinearModel.LinearRegression()
regr_fit.fit(X_train, y_train)

# Make predictions using the test data
y_pred = regr.predict(X_dev)
# y_pred = np.nan to num(y_pred)

# Remove high outlier predictions
for i, num in enumerate(y_pred):
    num > descriptive_fun(y_dev).loc['max', :][0]:
    y_pred[i] = y_pred[i]

# Compute MAE from prediction and truth
mae = mean_absolute_error(target_pipeline.inverse_transform(y_dev),
target_pipeline.inverse_transform(y_pred))
print(f"Mean Absolute error: {mae}%")
```

Principal Component Analysis

From the data visualization on features we concluded no evidence of strong correlation. Therefore we are using unsupervised dimensionality reduction technique to remove not so useful features. So we will set the variance explainability component to 99% and let PCA select components on its own.

After running PCA fit on Numeric and Categorical data, using explained variance ratio we got 69 components. Original dimension of data is 87

```
In [184]: # PCA on Numeric features
# =====
list = []

# Instantiate PCA with number of component = 0.99
pca = PCA(n_components=)

# Perform fit on 11 features
pca.fit(num_trans_features)

reduced_cont_feature = pca.transform(num_trans_features)

list.append(pca.explained_variance_ratio_)

print(f"Total variance explained by PCA on Numeric Features {sum(list[:11])}")
print(f"Number of components selected by PCA on Numeric Features {reduced_cont_feature.shape[1]}")
```

PCA explained variance by PCA on Numeric Features (0.9999999999999999)

```
In [185]: # PCA on Categorical features
# =====
cat_def = pd.concat([cat_features, flag_features], axis=)
ignore_index=, reset_index(drop=)
cat_def.shape

list = []

pca = PCA(n_components=)

# Perform fit on categorical features
pca.fit(cat_def)

# Perform transform
reduced_cat_feature = pca.transform(cat_def)

list.append(pca.explained_variance_ratio_)

print(f"Total variance explained by PCA on Categorical Features {sum(list[:11])}")
print(f"Number of components selected by PCA on Categorical Features {reduced_cat_feature.shape[1]}")

PCA explained variance by PCA on Categorical Features (0.9999999999999999)
```

Combine datasets

After PCA and achieving dimensionality reduction on Cat and Num features, we will combine both dataset into one and review final shape

```
In [187]: # After PCA combine Cat and Num features into single dataset
reduced_feature = np.hstack([reduced_cat_feature, reduced_cont_feature])
print(f"Number of features present in data {data_df.shape[1]}")
print(f"Number of components selected by PCA {reduced_feature.shape[1]}")

Number of features present in data 87
Number of components selected by PCA 69
```

Metrics

Mean Absolute error (MAE) is the metric we have chosen to evaluate performance of the model. MAE is measure of absolute difference between actual Incurred and predicted Incurred averaged over all samples.

This metric is considered for the following reasons:

- In this problem statement, prediction of Incurred (claim value) is performed which is numerical quantity. Predicting numerical value from input feature is the case of linear regression. Closer the predicted value to actual, lesser the error term and more accurate regression function can be approximated.
- Another metric for regression problem is Mean Squared Error (MSE). In MSE, square of the difference is taken when evaluating error function. If the error term lies between 0 and 1 then MSE is preferred over MAE as square of number (error term) further reduces. But in given problem statement predicted Incurred value can error by significant margin in hundreds or even thousands by positive or negative margin, so considering MAE is more viable as error term and it is not further amplified by taking square.

```
In [151]: # Mean absolute error performance metric

# performance metric y_true, y_predict
""" calculates and returns the performance scores between
true and predicted values based on the metric chosen. """

# TODO: Calculate the performance score between 'y_true' and 'y_predict'
mae = mean_absolute_error(y_true, y_predict)

# Return the score
return mae
```

DecisionTreeRegressor

DecisionTreeRegressor model training is done on train data with default parameters to get initial MAE estimate. As we can see resulted MAE of 3377.15 is comparatively less than benchmarked we established with LinearRegressor.

```
In [152]: #Running DecisionTreeRegressor with default parameters

# Split original data into train and test in ratio of 9:1
X_train, X_dev, y_train, y_dev = train_test_split(reduced_feature,
y_new, test_size=, random_state=)

# Instantiate class
regressor = DecisionTreeRegressor(max_depth=)

# Fit train feature and target
regressor.fit(X_train, y_train)

# Predict on Test features
y_pred = regressor.predict(X_dev)

# Print MAE from predicted and actual Loss
print(f"Mean Absolute error: {mae}")
# Absolute error: target_pipeline.inverse_transform(y_dev),
target_pipeline.inverse_transform(y_pred).reshape(-1,))
```

Model Complexity Analysis

In this process we have to find the most optimal parameter setting which gives best MAE score for DecisionTreeRegressor. I have chosen max_depth parameter with range options from 5 to 15.

max_depth = [5,6,7,8,9,10,11,12,13,14,15].

Validation curve of training and testing scores are plotted with max_depth on x-axis and scores on y-axis. Scoring function 'neg_mean_absolute_error' is used for evaluation.

```
In [153]: # Model complexity using the max_depth parameter
from sklearn.model_selection import validation_curve

# ModelComplexity(X, y):

# Vary the max_depth parameter from 1 to 10
max_depth = np.arange(, 10)

# Calculate the training and validation scores
train_scores, test_scores = validation_curve(DecisionTreeRegressor(), X, y,
param_name = 'max_depth', param_range = max_depth, cv=, scoring =
'neg_mean_absolute_error')

# Find the mean and standard deviation for smoothing
train_mean = np.mean(train_scores, axis=)
train_std = np.std(train_scores, axis=)
test_mean = np.mean(test_scores, axis=)
test_std = np.std(test_scores, axis=)

print(train_mean, test_mean)

# Plot the validation curve
plt.figure(figsize=(, ))
plt.title('Decision Tree Regressor Complexity Performance')
plt.plot(max_depth, train_mean, 'o', color = 'r', label = 'Training Score')
plt.plot(max_depth, test_mean, 'o', color = 'b', label = 'Validation Score')
plt.fill_between(max_depth, train_mean - train_std,
train_mean + train_std, alpha =, color = 'r')
plt.fill_between(max_depth, test_mean - test_std,
test_mean + test_std, alpha =, color = 'b')

# Visual aesthetics
plt.legend(loc = 'lower right')
plt.xlabel('Max depth')
plt.ylabel('Score')
# plt.ylim((-0.40,-0.60))
plt.show()
```

Observation

As we can observe from model complexity plot below:

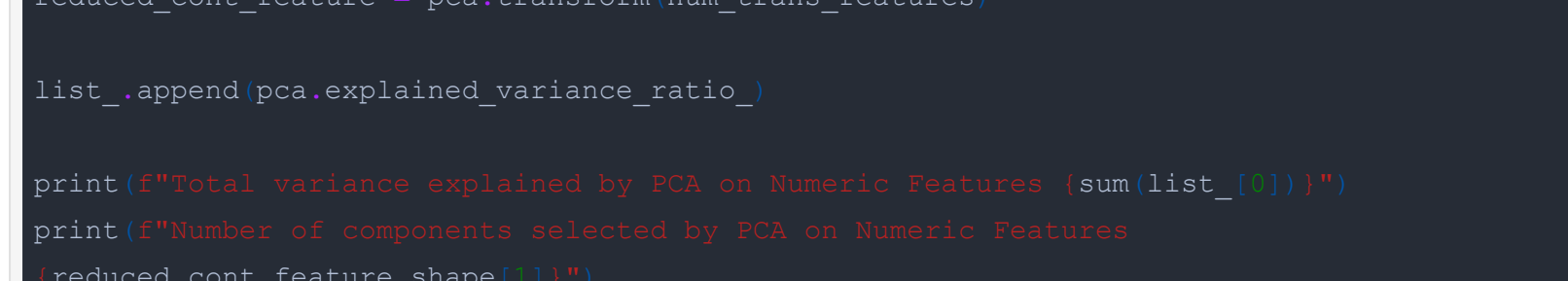
- For max_depth = [5, 6] both training and testing scores are low showing model is underfitting.
- For max_depth = [7, 15] training score kept increasing as max_depth increased on the other hand, testing score increased until max_depth 9 and then remain constant comparatively. This shows sign of high variance as model is overfitting and failing to generalize well.

Hence conclusion can be made from this plot is that we can generalize model more perfectly if we can tune all hyperparameters together to get optimal MAE

```
In [154]: ModelComplexity(X_train, y_train)

ModelComplexity(X_train, y_train)
-0.11366203 -0.087190561 -0.053268972 -0.044418771 -0.03184781 -0.02574234 -0.10950249 -0.29933626 -0.29279891
-0.28818281
-0.28148492 -0.27754036 -0.27977128 -0.272947151

Decision Tree Regressor Complexity Performance
```



Model Selection

This is a final step where two functions are constructed namely, split_train_holdout and runModel.

- The first function takes data and split it into train and evaluation. Train data is used to train model and evaluation is used as a test set
- Second function is used to train different models (new models can be easily added and tested) and test using evaluation set

