## ✓ XGBoosting

```
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  %matplotlib inline
6  from sklearn.model_selection import GridSearchCV
7  from __future__ import print_function
8  import os
9
10 from google.colab import drive
11 drive.mount('/content/drive')
12
13 filepath = '/content/drive/My Drive/Colab Notebooks/Dataset/Walmart.csv'
14 xg = pd.read_csv(filepath)
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=Tr
```

```
1  xg.head()
```

|   | Store | Date | Weekly_Sales | Holiday_Flag | Temperature | Fuel_Price | CPI | Unemployment |
|---|-------|------|--------------|--------------|-------------|------------|-----|--------------|
| **0** | 1 | 05-02-2010 | 1643690.90 | 0 | 42.31 | 2.572 | 211.096358 | 8.106 |
| **1** | 1 | 12-02-2010 | 1641957.44 | 1 | 38.51 | 2.548 | 211.242170 | 8.106 |
| **2** | 1 | 19-02-2010 | 1611968.17 | 0 | 39.93 | 2.514 | 211.289143 | 8.106 |
| **3** | 1 | 26-02-2010 | 1409727.59 | 0 | 46.63 | 2.561 | 211.319643 | 8.106 |
| **4** | 1 | 05-03-2010 | 1554806.68 | 0 | 46.50 | 2.625 | 211.350143 | 8.106 |

```
1  xg.shape
```

```
(6435, 8)
```

```
1  xg.dtypes.value_counts()
```

```
float64    5
int64      2
object     1
dtype: int64
```

```
1  xg = xg.drop('Date',axis=1)
```

```
1  float_columns = (xg.dtypes == np.float)
2  print( (xg.loc[:,float_columns].max()==1.0).all() )
3  print( (xg.loc[:,float_columns].min()==-1.0).all() )
```

```
False
False
<ipython-input-7-c6b27d4c13c4>:2: DeprecationWarning: `np.float` is a deprecated alias for the builtin `float`. To silence t
Deprecated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
  float_columns = (xg.dtypes == np.float)
```

```
1  from sklearn.preprocessing import LabelEncoder
2
3  le = LabelEncoder()
4
5  xg['Holiday_Flag'] = le.fit_transform(xg['Holiday_Flag'])
6
7  le.classes_
```

```
array([0, 1])
```

```
1  xg.Holiday_Flag.unique()
```

```
array([0, 1])
```

```
1  from sklearn.model_selection import train_test_split
2
3  feature_columns = [x for x in xg.columns if x != 'Holiday_Flag']
4
```

```
5 X_train, X_test, y_train, y_test = train_test_split(xg[feature_columns], xg['Holiday_Flag'],
6                     test_size=0.3, random_state=42)
```

```
1 X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

```
((4504, 6), (4504,), (1931, 6), (1931,))
```

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.metrics import accuracy_score
3
4 error_list = list()
5
6 tree_list = [15, 50, 100, 200, 400]
7 for n_trees in tree_list:
8
9     GBC = GradientBoostingClassifier(n_estimators=n_trees,
10                                     subsample=0.5,
11                                     max_features=4,
12                                     random_state=42)
13
14     GBC.fit(X_train.values, y_train.values)
15     y_pred = GBC.predict(X_test)
16
17     error = 1. - accuracy_score(y_test, y_pred)
18
19     error_list.append(pd.Series({'n_trees': n_trees, 'error': error}))
20
21 error_xg = pd.concat(error_list, axis=1).T.set_index('n_trees')
22
23 error_xg
```

```
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:443: UserWarning: X has feature names, but GradientBoostingClassifier
  warnings.warn(
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:443: UserWarning: X has feature names, but GradientBoostingClassifier
  warnings.warn(
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:443: UserWarning: X has feature names, but GradientBoostingClassifier
  warnings.warn(
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:443: UserWarning: X has feature names, but GradientBoostingClassifier
  warnings.warn(
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:443: UserWarning: X has feature names, but GradientBoostingClassifier
  warnings.warn(
```
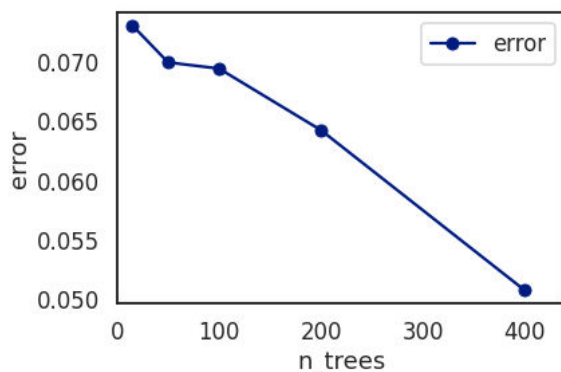
|         | error    |
| ------- | -------- |
| n_trees |          |
| 15.0    | 0.073019 |
| 50.0    | 0.069912 |
| 100.0   | 0.069394 |
| 200.0   | 0.064215 |
| 400.0   | 0.050751 |

```
1 sns.set_context('talk')
2 sns.set_style('white')
3 sns.set_palette('dark')
4
5 ax = error_xg.plot(marker='o')
6
7 ax.set(xlabel='n_trees', ylabel='error')
8 ax.set_xlim(0, max(error_xg.index)*1.1);
```



Number of trees vs error graph

```
1 param_grid = {'n_estimators': [200, 400],
2               'learning_rate': [0.1, 0.01]}
3
4 GV_GBC = GridSearchCV(GradientBoostingClassifier(subsample=0.5,
5                                                  max_features=4,
6                                                  random_state=42),
7                       param_grid=param_grid,
8                       scoring='accuracy',
9                       n_jobs=-1)
10
11 GV_GBC = GV_GBC.fit(X_train, y_train)
```

```
1 GV_GBC.best_estimator_
```

```
GradientBoostingClassifier(max_features=4, n_estimators=400, random_state=42,
                           subsample=0.5)
```

## ∨ Accuracy, Classification Report, ROC AUC Score, Confusion Matrix for XGBoosting

```
1 from sklearn.metrics import classification_report
2
3 y_pred = GV_GBC.predict(X_test)
4 print(classification_report(y_pred, y_test))
```

```
              precision    recall  f1-score   support

           0       1.00      0.95      0.97      1880
           1       0.33      0.92      0.49        51

    accuracy                           0.95      1931
   macro avg       0.67      0.94      0.73      1931
weighted avg       0.98      0.95      0.96      1931
```
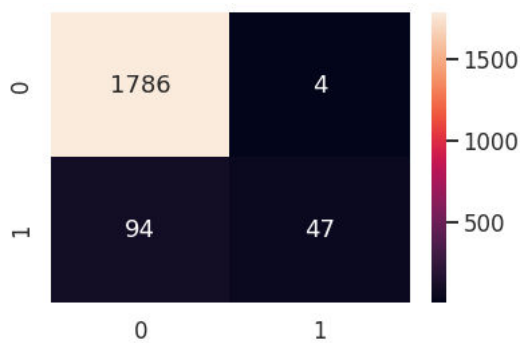
- Based from the classification report, the team got an accuracy of 0.95 and with 1.00, 0.95 and 0.97 for precision, recall and f1-score, respectively.

```
1 from sklearn.metrics import confusion_matrix
2
3 sns.set_context('talk')
4 cm = confusion_matrix(y_test, y_pred)
5 ax = sns.heatmap(cm, annot=True, fmt='d')
```



```
1 from scipy.stats import uniform, randint
2 import xgboost as xgb
3 from sklearn.metrics import classification_report,confusion_matrix,accuracy_score, roc_auc_score
```

```
1 xgb_model = xgb.XGBClassifier(objective="binary:logistic", voting='hard', random_state=42)
2 xgb_model.fit(X_train, y_train)
3
4 y_pred = xgb_model.predict(X_test)
5
6 print(confusion_matrix(y_test, y_pred))
```

```
[[1789    1]
 [ 130   11]]
```

```
1 y_pred = xgb_model.predict(X_test)
2 print(classification_report(y_pred, y_test))
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 1.00 | 0.93 | 0.96 | 1919 |
| 1 | 0.08 | 0.92 | 0.14 | 12 |
| accuracy |  |  | 0.93 | 1931 |
| macro avg | 0.54 | 0.92 | 0.55 | 1931 |
| weighted avg | 0.99 | 0.93 | 0.96 | 1931 |

```
1 print(accuracy_score(y_test, y_pred))
```

```
0.9321595028482651
```

```
1 print(roc_auc_score(y_pred,y_test))
```
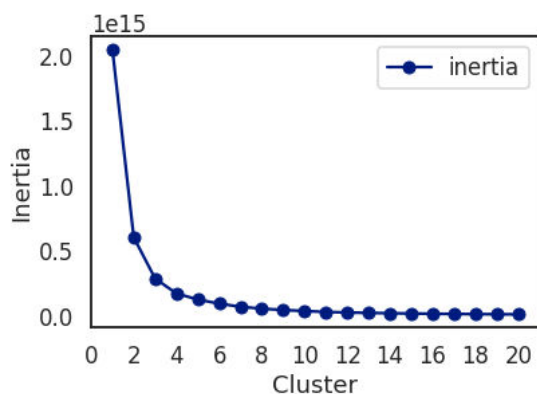
```
0.9244615250998783
```

After importing the xgboost as xgb and applying the code above, the classification report has slightly changed as well as as the accuracy. It has reduced to 0.93 now which is still a good value and the team got 0.924 as the roc auc score as well.

## ⌄ Agglomerative Clustering

```
1 from sklearn.cluster import KMeans
2
3 float_columns = [x for x in xg.columns if x not in ['Holiday_Flag', 'Unemployment']]
4
5 km = KMeans(n_clusters=2, random_state=42)
6 km = km.fit(xg[float_columns])
7
8 xg['kmeans'] = km.predict(xg[float_columns])
```

```
1 km_list = list()
2
3 for clust in range(1,21):
4     km = KMeans(n_clusters=clust, random_state=42)
5     km = km.fit(xg[float_columns])
6
7     km_list.append(pd.Series({'clusters': clust,
8                               'inertia': km.inertia_,
9                               'model': km}))
```

```
1 plot_xg = (pd.concat(km_list, axis=1)
2             .T
3             [['clusters','inertia']]
4             .set_index('clusters'))
5
6 ax = plot_xg.plot(marker='o',ls='-')
7 ax.set_xticks(range(0,21,2))
8 ax.set_xlim(0,21)
9 ax.set(xlabel='Cluster', ylabel='Inertia');
```



Inertia vs. Cluster graph

```
1 from sklearn.cluster import AgglomerativeClustering
2 from scipy.cluster import hierarchy
3 from matplotlib import colors
4
5 ag = AgglomerativeClustering(n_clusters=2, linkage='ward', compute_full_tree=True)
```
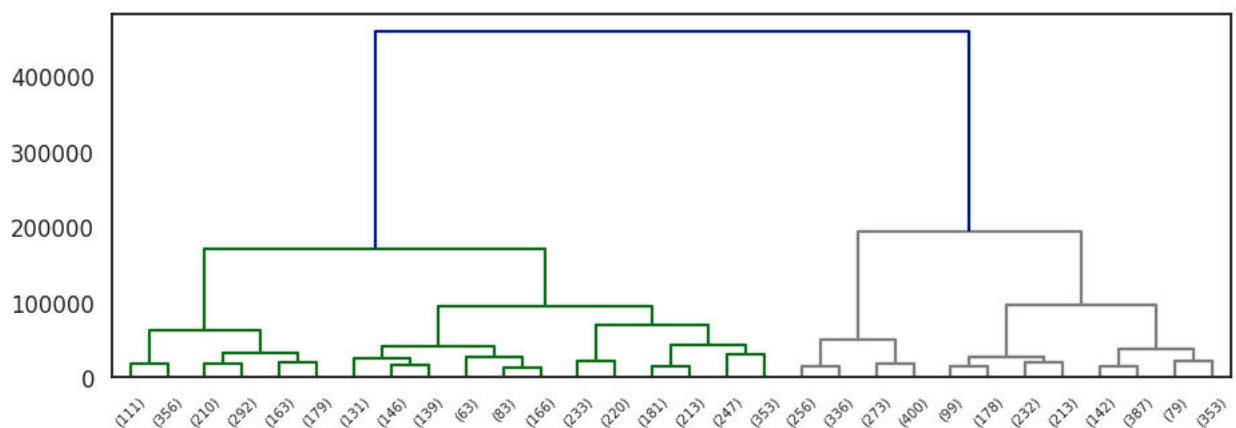
```
6 ag = ag.fit(xg[float_columns])
7 xg['agglom'] = ag.fit_predict(xg[float_columns])
```

```
1 (xg[['Holiday_Flag','agglom','kmeans']]
2 .groupby(['Holiday_Flag','agglom','kmeans'])
3 .size()
4 .to_frame()
5 .rename(columns={0:'number'}))
```

|  |  |  | number |
| --- | --- | --- | --- |
| Holiday_Flag | agglom | kmeans |  |
| 0 | 0 | 0 | 1384 |
|  |  | 1 | 2375 |
|  | 1 | 0 | 2226 |
| 1 | 0 | 0 | 95 |
|  |  | 1 | 197 |
|  | 1 | 0 | 158 |

## Dendogram produced by agglomerative clustering

```
 1 from scipy.cluster import hierarchy
 2 from matplotlib import colors
 3
 4 Z = hierarchy.linkage(ag.children_, method='ward')
 5
 6 fig, ax = plt.subplots(figsize=(15,5))
 7
 8 # Some color customization
 9 dark_palette = sns.color_palette()
10 red = colors.to_hex(dark_palette[2])
11 blue = colors.to_hex(dark_palette[0])
12
13 hierarchy.set_link_color_palette([red, 'gray'])
14
15 den = hierarchy.dendrogram(Z, orientation='top',
16                            p=30, truncate_mode='lastp',
17                            show_leaf_counts=True, ax=ax,
18                            above_threshold_color=blue)
```



```
 1 from sklearn.ensemble import RandomForestClassifier
 2 from sklearn.metrics import classification_report, roc_auc_score
 3 from sklearn.model_selection import StratifiedShuffleSplit
 4
 5
 6 y = (xg['Unemployment'] > 7).astype(int)
 7 X_with_kmeans = xg.drop(['agglom', 'Holiday_Flag', 'Unemployment'], axis=1)
 8 X_without_kmeans = X_with_kmeans.drop('kmeans', axis=1)
 9 sss = StratifiedShuffleSplit(n_splits=10, random_state=6532)
10
11
12 def get_avg_roc_10splits(estimator, X, y):
13     roc_auc_list = []
14     for train_index, test_index in sss.split(X, y):
15         X_train, X_test = X.iloc[train_index], X.iloc[test_index]
16         y_train, y_test = y.iloc[train_index], y.iloc[test_index]
```

```
17        estimator.fit(X_train, y_train)
18        y_predicted = estimator.predict(X_test)
19        y_scored = estimator.predict_proba(X_test)[:, 1]
20        roc_auc_list.append(roc_auc_score(y_test, y_scored))
21     return np.mean(roc_auc_list)
22 # return classification_report(y_test, y_predicted)
23
24
25 estimator = RandomForestClassifier()
26 roc_with_kmeans = get_avg_roc_10splits(estimator, X_with_kmeans, y)
27 roc_without_kmeans = get_avg_roc_10splits(estimator, X_without_kmeans, y)
28 print("Without kmeans cluster as input to Random Forest, roc-auc is \"{0}\"".format(roc_without_kmeans))
29 print("Using kmeans cluster as input to Random Forest, roc-auc is \"{0}\"".format(roc_with_kmeans))
30
```

```
Without kmeans cluster as input to Random Forest, roc-auc is "0.9994285576505761"
Using kmeans cluster as input to Random Forest, roc-auc is "0.9992308434199739"
```

## ⌄ Pipeline

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6
7 from google.colab import drive
8 drive.mount('/content/drive')
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=Tr
```

```
1 filepath = '/content/drive/My Drive/Colab Notebooks/Dataset/Walmart.csv'
2 xg = pd.read_csv(filepath)
```

```
1 xg
```

|      | Store | Date       | Weekly_Sales | Holiday_Flag | Temperature | Fuel_Price | CPI        | Unemployment |
|------|-------|------------|--------------|--------------|-------------|------------|------------|--------------|
| 0    | 1     | 05-02-2010 | 1643690.90   | 0            | 42.31       | 2.572      | 211.096358 | 8.106        |
| 1    | 1     | 12-02-2010 | 1641957.44   | 1            | 38.51       | 2.548      | 211.242170 | 8.106        |
| 2    | 1     | 19-02-2010 | 1611968.17   | 0            | 39.93       | 2.514      | 211.289143 | 8.106        |
| 3    | 1     | 26-02-2010 | 1409727.59   | 0            | 46.63       | 2.561      | 211.319643 | 8.106        |
| 4    | 1     | 05-03-2010 | 1554806.68   | 0            | 46.50       | 2.625      | 211.350143 | 8.106        |
| ...  | ...   | ...        | ...          | ...          | ...         | ...        | ...        | ...          |
| 6430 | 45    | 28-09-2012 | 713173.95    | 0            | 64.88       | 3.997      | 192.013558 | 8.684        |
| 6431 | 45    | 05-10-2012 | 733455.07    | 0            | 64.89       | 3.985      | 192.170412 | 8.667        |
| 6432 | 45    | 12-10-2012 | 734464.36    | 0            | 54.47       | 4.000      | 192.327265 | 8.667        |
| 6433 | 45    | 19-10-2012 | 718125.53    | 0            | 56.47       | 3.969      | 192.330854 | 8.667        |
| 6434 | 45    | 26-10-2012 | 760281.43    | 0            | 58.85       | 3.882      | 192.308899 | 8.667        |

6435 rows × 8 columns

```
1 xg = xg.drop('Date',axis=1)
```

```
1 for col in xg.columns:
2     xg[col] = xg[col].astype(np.float)
```

```
<ipython-input-56-6dbce3b12efe>:2: DeprecationWarning: `np.float` is a deprecated alias for the builtin `float`. To silence
Deprecated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
  xg[col] = xg[col].astype(np.float)
```

```
1 xg_orig = xg.copy()
```

```
1 corr_mat = xg.corr()
2
3 # Strip the diagonal for future examination
4 for x in range(corr_mat.shape[0]):
5     corr_mat.iloc[x,x] = 0.0
```

```
6
7 corr_mat
```

|  | Store | Weekly_Sales | Holiday_Flag | Temperature | Fuel_Price | CPI | Unemployment |
|---|---|---|---|---|---|---|---|
| **Store** | 0.000000e+00 | -0.335332 | -4.386841e-16 | -0.022659 | 0.060023 | -0.209492 | 0.223531 |
| **Weekly_Sales** | -3.353320e-01 | 0.000000 | 3.689097e-02 | -0.063810 | 0.009464 | -0.072634 | -0.106176 |
| **Holiday_Flag** | -4.386841e-16 | 0.036891 | 0.000000e+00 | -0.155091 | -0.078347 | -0.002162 | 0.010960 |
| **Temperature** | -2.265908e-02 | -0.063810 | -1.550913e-01 | 0.000000 | 0.144982 | 0.176888 | 0.101158 |
| **Fuel_Price** | 6.002295e-02 | 0.009464 | -7.834652e-02 | 0.144982 | 0.000000 | -0.170642 | -0.034684 |
| **CPI** | -2.094919e-01 | -0.072634 | -2.162091e-03 | 0.176888 | -0.170642 | 0.000000 | -0.302020 |
| **Unemployment** | 2.235313e-01 | -0.106176 | 1.096028e-02 | 0.101158 | -0.034684 | -0.302020 | 0.000000 |

```
1 corr_mat.abs().idxmax()
```

```
Store            Weekly_Sales
Weekly_Sales             Store
Holiday_Flag      Temperature
Temperature               CPI
Fuel_Price                CPI
CPI              Unemployment
Unemployment              CPI
dtype: object
```

```
1 log_columns = xg.skew().sort_values(ascending=False)
2 log_columns = log_columns.loc[log_columns > 0.75]
3
4 log_columns
```

```
Holiday_Flag     3.373499
Unemployment     1.188144
dtype: float64
```

```
1 for col in log_columns.index:
2     xg[col] = np.log1p(xg[col])
```

```
1 from sklearn.preprocessing import MinMaxScaler
2
3 mms = MinMaxScaler()
4
5 for col in xg.columns:
6     xg[col] = mms.fit_transform(xg[[col]]).squeeze()
```

```
1 from sklearn.preprocessing import FunctionTransformer
2 from sklearn.pipeline import Pipeline
3
4 # The custom NumPy log transformer
5 log_transformer = FunctionTransformer(np.log1p)
6
7 # The pipeline
8 estimators = [('log1p', log_transformer), ('minmaxscale', MinMaxScaler())]
9 pipeline = Pipeline(estimators)
10
11 # Convert the original data
12 xg_pipe = pipeline.fit_transform(xg_orig)
```

```
1 np.allclose(xg_pipe, xg)
```

```
False
```

```
1 from sklearn.decomposition import PCA
2
3 pca_list = list()
4 feature_weight_list = list()
5
6 # Fit a range of PCA models
7
8 for n in range(1, 6):
9
10     # Create and fit the model
11     PCAmod = PCA(n_components=n)
12     PCAmod.fit(xg)
13
14     # Store the model and variance
15     pca_list.append(pd.Series({'n':n, 'model':PCAmod,
```

```
16                                    'var': PCAmod.explained_variance_ratio_.sum()}))
17
18     # Calculate and store feature importances
19     abs_feature_values = np.abs(PCAmod.components_).sum(axis=0)
20     feature_weight_list.append(pd.DataFrame({'n':n,
21                                              'features': xg.columns,
22                                              'values':abs_feature_values/abs_feature_values.sum()}))
23
24 pca_df = pd.concat(pca_list, axis=1).T.set_index('n')
25 pca_df
```

|   | model | var |
|---|-------|-----|
| **n** | | |
| **1** | PCA(n_components=1) | 0.37273 |
| **2** | PCA(n_components=2) | 0.563031 |
| **3** | PCA(n_components=3) | 0.717861 |
| **4** | PCA(n_components=4) | 0.832671 |
| **5** | PCA(n_components=5) | 0.90656 |

```
1 features_df = (pd.concat(feature_weight_list)
2                 .pivot(index='n', columns='features', values='values'))
3
4 features_df
```
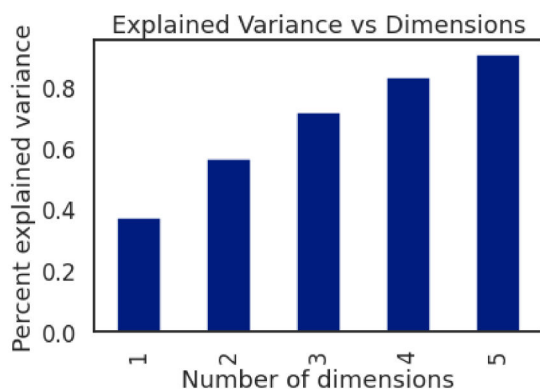
| features | CPI | Fuel_Price | Holiday_Flag | Store | Temperature | Unemployment | Weekly_Sales |
|----------|-----|-----------|--------------|-------|-------------|--------------|--------------|
| **n** | | | | | | | |
| **1** | 0.577673 | 0.080874 | 0.001512 | 0.196542 | 0.050184 | 0.088784 | 0.004432 |
| **2** | 0.372387 | 0.047145 | 0.012918 | 0.364854 | 0.045827 | 0.075438 | 0.081430 |
| **3** | 0.252680 | 0.101460 | 0.193040 | 0.257552 | 0.079229 | 0.057383 | 0.058657 |
| **4** | 0.201596 | 0.209181 | 0.201394 | 0.193538 | 0.085942 | 0.063331 | 0.045017 |
| **5** | 0.160937 | 0.175849 | 0.173813 | 0.168016 | 0.153904 | 0.123667 | 0.043813 |

```
1 sns.set_context('talk')
2
3 ax = pca_df['var'].plot(kind='bar')
4
5 ax.set(xlabel='Number of dimensions',
6        ylabel='Percent explained variance',
7        title='Explained Variance vs Dimensions');
```
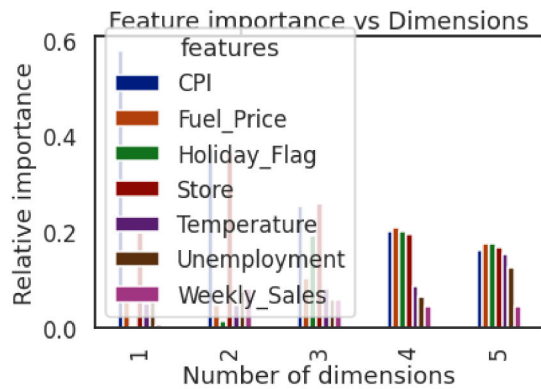


This graph is the variance vs dimensions. As seen above, it is increasing consistently.

```
1 ax = features_df.plot(kind='bar')
2
3 ax.set(xlabel='Number of dimensions',
4        ylabel='Relative importance',
5        title='Feature importance vs Dimensions');
```
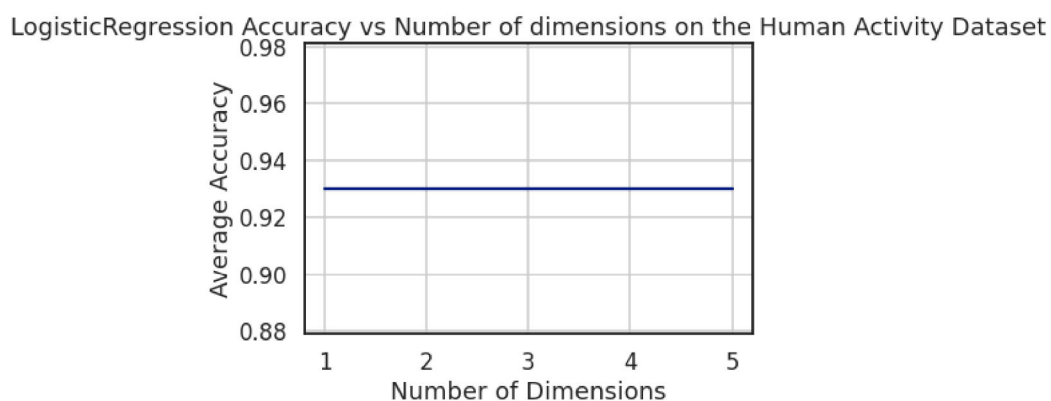
This graph shows the features as well as the number of dimensions vs the importance of the features with the features having different colors for the legends of the graph.

```python
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import StratifiedShuffleSplit
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import MinMaxScaler


X = xg.drop('Holiday_Flag', axis=1)
y = xg.Holiday_Flag
sss = StratifiedShuffleSplit(n_splits=10, random_state=42)

def get_avg_score(n):
    pipe = [
        ('scaler', MinMaxScaler()),
        ('pca', PCA(n_components=n)),
        ('estimator', LogisticRegression())
    ]
    pipe = Pipeline(pipe)
    scores = []
    for train_index, test_index in sss.split(X, y):
        X_train, X_test = X.loc[train_index], X.loc[test_index]
        y_train, y_test = y.loc[train_index], y.loc[test_index]
        pipe.fit(X_train, y_train)
        scores.append(accuracy_score(y_test, pipe.predict(X_test)))
    return np.mean(scores)


ns = [1,2,3,4,5]
score_list = [get_avg_score(n) for n in ns]
```

```python
sns.set_context('talk')

ax = plt.axes()
ax.plot(ns, score_list)
ax.set(xlabel='Number of Dimensions',
       ylabel='Average Accuracy',
       title='LogisticRegression Accuracy vs Number of dimensions on the Human Activity Dataset')
ax.grid(True)
```



```python
Start coding or generate with AI.
```

The graph shows the average accuracy of our data vs the number of dimensions.

```
1 xg.head()
```

|   | Store | Weekly_Sales | Holiday_Flag | Temperature | Fuel_Price | CPI | Unemployment |
|---|-------|--------------|--------------|-------------|------------|-----|--------------|
| **0** | 0.0 | 0.397291 | 0.0 | 0.434149 | 0.050100 | 0.840500 | 0.545562 |
| **1** | 0.0 | 0.396811 | 1.0 | 0.396967 | 0.038076 | 0.841941 | 0.545562 |
| **2** | 0.0 | 0.388501 | 0.0 | 0.410861 | 0.021042 | 0.842405 | 0.545562 |
| **3** | 0.0 | 0.332458 | 0.0 | 0.476419 | 0.044589 | 0.842707 | 0.545562 |
| **4** | 0.0 | 0.372661 | 0.0 | 0.475147 | 0.076653 | 0.843008 | 0.545562 |