## ⌄ *Filter Based Methods - Mutual Information* - 0.578273098930927

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.datasets import fetch_california_housing
6 from sklearn.metrics import r2_score, mean_squared_error
7 from sklearn.linear_model import LinearRegression, Lasso
```

```
1 dataset = fetch_california_housing()
2 n_features = dataset.data.shape[1]
3 print(n_features)
4 feature_names = dataset.feature_names
5 print(feature_names)
6 #print(dataset.DESCR)
```

```
8
    ['MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population', 'AveOccup', 'Latitude', 'Longitude']
```

```
1 x=dataset['data']
2 y=dataset['target']
3 feature_names = dataset['feature_names']
```
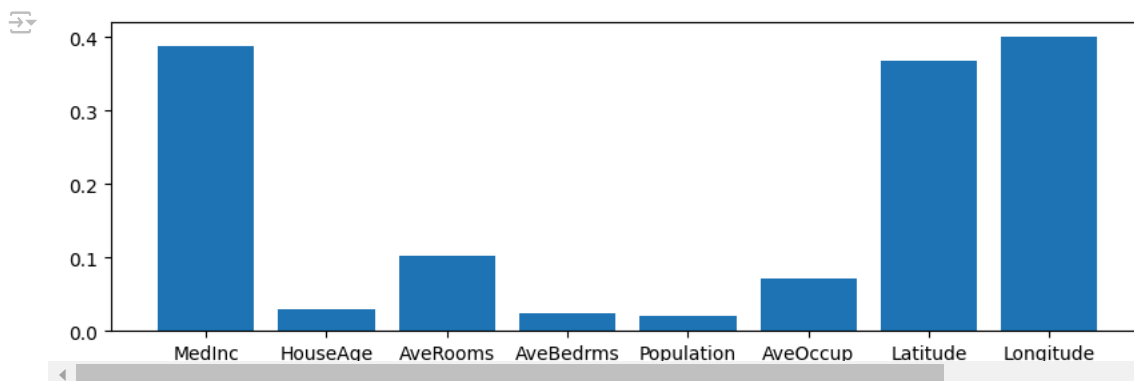
```
1 from sklearn.feature_selection import mutual_info_regression, SelectPercentile,SelectKBest
```

```
1 mi = mutual_info_regression(x,y)
2 #print(mi)
3 # Create a dictionary of feature names and their importance scores
4 feature_importance = dict(zip(dataset.feature_names, mi))
5 # Print in a formatted way
6 for feature, importance in feature_importance.items():
7     print(f"{feature}: {importance:.4f}")
```

```
MedInc: 0.3875
HouseAge: 0.0294
AveRooms: 0.1032
AveBedrms: 0.0242
Population: 0.0213
AveOccup: 0.0722
Latitude: 0.3688
Longitude: 0.4003
```

Visualise feature Selection

```
1 plt.figure(figsize=(10,3))
2 plt.bar(feature_names,mi)
3 plt.show()
```



**How to pick the feature**

*Option 1*

```
1 x_new = SelectKBest(mutual_info_regression,k=5).fit_transform(x,y)
2 print(x_new.shape)
3 # Get scores for all features
4 selector = SelectKBest(mutual_info_regression,k=5).fit(x,y)
5 scores = selector.scores_
6 # Create a dictionary of feature names and their scores
7 feature_scores = dict(zip(dataset.feature_names, scores))
8 # Sort features by score
9 sorted_features = sorted(feature_scores.items(), key=lambda x: x[1], reverse=True)[:5]
10
11 print("Top 5 features with their scores:")
12 for feature, score in sorted_features:
13     print(f"{feature}: {score:.4f}")
```

```
(20640, 5)
Top 5 features with their scores:
Longitude: 0.4024
MedInc: 0.3880
Latitude: 0.3678
AveRooms: 0.1032
AveOccup: 0.0728
```

*Option 2*

```
1 x_new = SelectPercentile(mutual_info_regression,percentile=50).fit_transform(x,y)
2 print(x_new.shape)
3 # Get scores for all features
4 selector = SelectKBest(mutual_info_regression,k=5).fit(x,y)
5 scores = selector.scores_
6 # Create a dictionary of feature names and their scores
7 feature_scores = dict(zip(dataset.feature_names, scores))
8 # Sort features by score
9 sorted_features = sorted(feature_scores.items(), key=lambda x: x[1], reverse=True)[:5]
10
11 print("Top 5 features with their scores:")
12 for feature, score in sorted_features:
13     print(f"{feature}: {score:.4f}")
```

```
(20640, 4)
Top 5 features with their scores:
Longitude: 0.4009
MedInc: 0.3874
Latitude: 0.3713
AveRooms: 0.1036
AveOccup: 0.0720
```

```
1 x_train,x_test,y_train,y_test=train_test_split(x_new,y, test_size=0.2)
2 model = LinearRegression()
3 model.fit(x_train,y_train)
4 y_pred = model.predict(x_test)
5 print(r2_score(y_test,y_pred))
```

```
0.578273098930927
```

## ⌄ *Filter Based Methods - Chi Squared Cannot be used for Continuous dataset*

Double-click (or enter) to edit

## ⌄ *Filter Based Methods - Pearson Correlation * - 0.5173065395604509

f_regression: A scoring function that computes F-statistics between each feature and target

SelectKBest: A feature selector that selects top k features based on a scoring function

```
1 from sklearn.feature_selection import f_regression, SelectKBest
```

```
1 #x_new = SelectKBest(f_regression,k=8).fit_transform(x,y)
2 #print(x_new.shape)
3
```

```
4 #Create the selector
5 selector = SelectKBest(f_regression, k=6)
6
7 # Fit the selector (without transforming)
8 selector.fit(x, y)
9
10 # Get the selected feature mask (True for selected features)
11 selected_mask = selector.get_support()
12
13 # Get the names of selected features using the mask
14 #selected_features = dataset.feature_names[selected_mask]
15 selected_features = [feature for feature, is_selected in zip(dataset.feature_names, selected_mask) if is_selected]
16 """
17 zip(dataset.feature_names, selected_mask) pairs each feature name with its corresponding boolean value (True/False) from the selected_mas
18 for feature, is_selected in ... loops through each pair, where:
19
20 feature gets the feature name from dataset.feature_names
21 is_selected gets the corresponding True/False value from selected_mask
22
23
24 if is_selected only keeps features where the mask value is True
25 [feature for ...] creates a list containing only the feature names that passed the condition
26 """
27 # Print selected features and their scores
28 for feature, score in zip(selected_features, selector.scores_):
29     print(f"{feature}: {score:.4f}")
```

```
MedInc: 18556.5716
HouseAge: 232.8415
AveRooms: 487.7575
AveBedrms: 45.1086
Latitude: 12.5474
Longitude: 11.6353
```

```
1 x_train,x_test,y_train,y_test  = train_test_split(x_new,y,test_size=0.2)
```

```
1 model=LinearRegression()
2 model.fit(x_train,y_train)
3 y_pred=model.predict(x_test)
4 print(r2_score(y_test,y_pred))
```

```
0.5979581682306951
```

To drop feature manually

```
1 import pandas as pd
2 x_pd = pd.DataFrame(x,columns=feature_names)
3 x_pd.head(10)
4 x_pd.corr()
```

| | MedInc | HouseAge | AveRooms | AveBedrms | Population | AveOccup | Latitude | Longitude |
|---|---|---|---|---|---|---|---|---|
| MedInc | 1.000000 | -0.119034 | 0.326895 | -0.062040 | 0.004834 | 0.018766 | -0.079809 | -0.015176 |
| HouseAge | -0.119034 | 1.000000 | -0.153277 | -0.077747 | -0.296244 | 0.013191 | 0.011173 | -0.108197 |
| AveRooms | 0.326895 | -0.153277 | 1.000000 | 0.847621 | -0.072213 | -0.004852 | 0.106389 | -0.027540 |
| AveBedrms | -0.062040 | -0.077747 | 0.847621 | 1.000000 | -0.066197 | -0.006181 | 0.069721 | 0.013344 |
| Population | 0.004834 | -0.296244 | -0.072213 | -0.066197 | 1.000000 | 0.069863 | -0.108785 | 0.099773 |
| AveOccup | 0.018766 | 0.013191 | -0.004852 | -0.006181 | 0.069863 | 1.000000 | 0.002366 | 0.002476 |
| Latitude | -0.079809 | 0.011173 | 0.106389 | 0.069721 | -0.108785 | 0.002366 | 1.000000 | -0.924664 |
| Longitude | -0.015176 | -0.108197 | -0.027540 | 0.013344 | 0.099773 | 0.002476 | -0.924664 | 1.000000 |

```
1 y_series = pd.Series(y)
2 x_pd.corrwith(y_series).sort_values(ascending=False)
```

|          | 0         |
|----------|-----------|
| MedInc   | 0.688075  |
| AveRooms | 0.151948  |
| HouseAge | 0.105623  |
| AveOccup | -0.023737 |
| Population | -0.024650 |
| Longitude | -0.045967 |
| AveBedrms | -0.046701 |
| Latitude | -0.144160 |

dtype: float64

```
1 import seaborn as sns
```

```
1 f,ax=plt.subplots(figsize=(8,8))
2 sns.heatmap(x_pd.corr(),annot=True,linewidths=0.5,fmt='.1f',ax=ax)
3 plt.show()
```



Latitude and Longitude have the strongest negative correlation (-0.9), meaning as latitude increases, longitude tends to decrease significantly. This makes sense geographically. AveRooms and AveBedrms have a strong positive correlation (0.8), which is logical since houses with more total rooms tend to have more bedrooms.

Before dropping these features, you should consider: Check their importance scores (using mutual_info_regression or feature_importance_ from your model) to see which features are most predictive of your target variable

```
1 mi = mutual_info_regression(x, y)
2 feature_importance = dict(zip(dataset.feature_names, mi))
3 sorted_importance = dict(sorted(feature_importance.items(), key=lambda x: x[1], reverse=True))
```

```
4
5 for feature, importance in sorted_importance.items():
6     print(f"{feature}: {importance:.4f}")
```

```
Longitude: 0.4013
MedInc: 0.3868
Latitude: 0.3718
AveRooms: 0.1032
AveOccup: 0.0727
HouseAge: 0.0331
AveBedrms: 0.0244
Population: 0.0212
```

```
1 x_new = x_pd.drop(['Longitude','AveRooms'],axis=1)
2 x_new.head()
```

|   | MedInc | HouseAge | AveBedrms | Population | AveOccup | Latitude |
|---|--------|----------|-----------|------------|----------|----------|
| 0 | 8.3252 | 41.0     | 1.023810  | 322.0      | 2.555556 | 37.88    |
| 1 | 8.3014 | 21.0     | 0.971880  | 2401.0     | 2.109842 | 37.86    |
| 2 | 7.2574 | 52.0     | 1.073446  | 496.0      | 2.802260 | 37.85    |
| 3 | 5.6431 | 52.0     | 1.073059  | 558.0      | 2.547945 | 37.85    |
| 4 | 3.8462 | 52.0     | 1.081081  | 565.0      | 2.181467 | 37.85    |

Next steps:   [ Generate code with x_new ]   [ ○ View recommended plots ]   [ New interactive sheet ]

```
1 x_train,x_test,y_train,y_test = train_test_split(x_new,y,test_size=0.2)
2 model = LinearRegression()
3 model.fit(x_train,y_train)
4 y_pred=model.predict(x_test)
5 print(r2_score(y_test,y_pred))
```

```
0.5173065395604509
```

## Wrapper Based Methods - Recursive Feature Elimination (RFE) - 0.5896765399567889

```
1 from sklearn.feature_selection import RFE
2 from sklearn.linear_model import Lasso
```

```
1 estimator = Lasso()
2 #Creates a Lasso (Least Absolute Shrinkage and Selection Operator) regression model which will be used as the base estimator.
3 #Lasso is a type of linear regression that includes L1 regularization
4 selector = RFE(estimator,n_features_to_select=5,step=1).fit(x,y)
5 print(selector.ranking_)
6
7 #Print features with their rankings
8 for feature, rank in zip(dataset.feature_names, selector.ranking_):
9     print(f"{feature}: {rank}")
10
11 # Or just print selected features
12 selected_features = [feature for feature, is_selected in zip(dataset.feature_names, selector.support_) if is_selected]
13 print("Selected features:", selected_features)
```

```
[1 1 4 3 1 2 1 1]
MedInc: 1
HouseAge: 1
AveRooms: 4
AveBedrms: 3
Population: 1
AveOccup: 2
Latitude: 1
Longitude: 1
Selected features: ['MedInc', 'HouseAge', 'Population', 'Latitude', 'Longitude']
```

```
1 x_new = selector.transform(x)
2 print(x_new.shape)
```

```
(20640, 5)
```

```
1 x_train,x_test,y_train,y_test=train_test_split(x_new,y, test_size=0.2)
```

```
2 model = LinearRegression()
3 model.fit(x_train,y_train)
4 y_pred = model.predict(x_test)
5 print(r2_score(y_test,y_pred))
```

→ 0.5896765399567889

```
2 model = LinearRegression()
3 model.fit(x_train,y_train)
4 y_pred = model.predict(x_test)
5 print(r2_score(y_test,y_pred))
```