

```
[70]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')

In [71]: from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score,r2_score,classification_report,confusion_matrix
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.feature_selection import SelectKBest,SelectPercentile,mutual_info_classif,mutual_info_regression

In [72]: data = pd.read_csv('CarPrice_Assignment.csv')

In [73]: data.head()
```

	car_ID	symboling	CarName	fueltpe	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	...	enginesize	fuelsy
0	1	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	
1	2	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	
2	3	1	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	
3	4	2	audi100/l	gas	std	four	sedan	fwd	front	99.8	...	109	
4	5	2	audi100/s	gas	std	four	sedan	4wd	front	99.4	...	136	

5 rows × 26 columns

```
In [74]: data.tail()
```

	car_ID	symboling	CarName	fueltpe	aspiration	doornumber	carbody	drivewheel	enginelocation	wheelbase	...	enginesize	fuelsy
200	201	-1	volvo 145e (sw)	gas	std	four	sedan	rwd	front	109.1	...	141	m
201	202	-1	volvo 145e	gas	turbo	four	sedan	rwd	front	109.1	...	141	m
202	203	-1	volvo 144e	gas	std	four	sedan	rwd	front	109.1	...	173	m
203	204	-1	volvo 246i	diesel	turbo	four	sedan	rwd	front	109.1	...	145	
204	205	-1	volvo 264g	gas	turbo	four	sedan	rwd	front	109.1	...	141	m

5 rows × 26 columns

```
In [75]: data.columns

Out[75]: Index(['car_ID', 'symboling', 'CarName', 'fueltpe', 'aspiration', 'doornumber', 'carbody', 'drivewheel', 'enginelocation', 'wheelbase', 'carlength', 'carwidth', 'carheight', 'curbweight', 'enginetype', 'cylindernumber', 'enginesize', 'fuelsystem', 'boreratio', 'stroke', 'compressionratio', 'horsepower', 'peakrpm', 'citympg', 'highwaympg', 'price'],
dtype='object')

In [76]: data.isnull().sum()
0

Out[76]: car_ID          0
symboling         6
CarName          147
fueltpe           2
aspiration         2
doornumber        2
carbody           5
drivewheel        3
enginelocation    2
wheelbase        53
carlength         75
carwidth          44
carheight         49
curbweight        171
enginetype         7
cylindernumber    7
enginesize        44
fuelsystem         8
boreratio         38
stroke            37
compressionratio   32
horsepower        59
peakrpm           23
citympg           29
highwaympg        30
price            189
dtype: int64

In [77]: data.nunique()

Out[77]: car_ID          205
symboling         6
CarName          147
fueltpe           2
aspiration         2
doornumber        2
carbody           5
drivewheel        3
enginelocation    2
wheelbase        53
carlength         75
carwidth          44
carheight         49
curbweight        171
enginetype         7
cylindernumber    7
enginesize        44
fuelsystem         8
boreratio         38
stroke            37
compressionratio   32
horsepower        59
peakrpm           23
citympg           29
highwaympg        30
price            189
dtype: int64

In [78]: data.iloc[:,20:30]

Out[78]:      compressionratio  horsepower  peakrpm  citympg  highwaympg  price
0              32.000000      5000.00    21.00    27.00    13499.00
1              32.000000      5000.00    21.00    27.00    16500.00
2              32.000000      5000.00    21.00    26.00    16500.00
3              32.000000      5500.00    24.00    30.00    13950.00
4              32.000000      5500.00    18.00    22.00    17450.00
...              ...              ...
200             32.000000      5400.00    23.00    28.00    16645.00
201             32.000000      5300.00    19.00    25.00    19045.00
202             32.000000      5500.00    18.00    23.00    21485.00
203             32.000000      4800.00    26.00    27.00    22470.00
204             32.000000      5400.00    19.00    25.00    22625.00

205 rows × 6 columns

In [79]: data.drop(columns=['car_ID'],inplace=True)

In [80]: data['CarName'].value_counts().sort_values(ascending=False).head(10)

Out[80]: toyota corolla      6
peugeot 504              6
subaru dl                4
mitsubishi outlander    3
masda 626                3
toyota mark ii          3
mitsubishi mirage g4    3
mitsubishi g4           3
honda civic             3
Name: CarName, dtype: int64

In [81]: label = LabelEncoder()

In [82]: data['CarName'] = label.fit_transform(data['CarName'])

In [83]: data['fueltpe'].value_counts().sort_values(ascending=False).head(10)

Out[83]: gas      185
diesel    20
Name: fueltpe, dtype: int64

In [84]: data['fueltpe'] = label.fit_transform(data['fueltpe'])

In [85]: data['aspiration'].value_counts().sort_values(ascending=False).head(10)

Out[85]: std      168
turbo     37
Name: aspiration, dtype: int64

In [86]: data['aspiration'] = label.fit_transform(data['aspiration'])

In [87]: data['doornumber'].value_counts().sort_values(ascending=False).head(10)

Out[87]: four      115
two       90
Name: doornumber, dtype: int64

In [88]: data['doornumber'] = label.fit_transform(data['doornumber'])

In [89]: data['drivewheel'].value_counts().sort_values(ascending=False).head(10)

Out[89]: fwd      120
rwd       76
4wd        9
Name: drivewheel, dtype: int64

In [90]: df = pd.get_dummies(data['drivewheel'],drop_first=True)

In [91]: data['enginelocation'].value_counts().sort_values(ascending=False).head(10)

Out[91]: front  202
rear    3
Name: enginelocation, dtype: int64

In [92]: data['enginelocation'] = label.fit_transform(data['enginelocation'])

In [93]: data['enginetype'].value_counts().sort_values(ascending=False).head(10)

Out[93]: ohc      148
ohcf     15
ohcv     13
l        12
dohc     12
rotor     4
dohcv     1
Name: enginetype, dtype: int64

In [94]: df2 = pd.get_dummies(data['enginetype'],drop_first=True)

In [95]: data['cylindernumber'].value_counts().sort_values(ascending=False).head(10)

Out[95]: four      159
six       24
five      11
eight      3
two        4
three      1
twelve     1
Name: cylindernumber, dtype: int64

In [96]: data['cylindernumber'] = label.fit_transform(data['cylindernumber'])

In [97]: data['fuelsystem'].value_counts().sort_values(ascending=False).head(10)

Out[97]: mpfi      94
2bbl      6
idl       20
lbbbl     11
spdi      9
4bbl       3
spfi       1
mfi        1
Name: fuelsystem, dtype: int64

In [98]: data['fuelsystem'] = label.fit_transform(data['fuelsystem'])

In [99]: dataal = pd.concat([data,df],axis=1)

In [100]: dataal.drop(columns=['drivewheel'],inplace=True)

In [101]: dataal.drop(columns=['enginetype'],inplace=True)

In [102]: dataal['carbody'] = label.fit_transform(dataal['carbody'])

In [103]: len(dataal.columns)

Out[103]: 25

In [104]: dataal.head()
```

	symboling	CarName	fueltpe	aspiration	doornumber	carbody	enginelocation	wheelbase	carlength	carwidth	...	boreratio	stroke	co
0	3	2	1	0	1	0	0	88.6	168.8	64.1	...	3.47	2.68	
1	3	3	1	0	1	0	0	88.6	168.8	64.1	...	3.47	2.68	
2	1	1	1	0	1	2	0	94.5	176.2	65.5	...	2.68	3.47	
3	2	4	1	0	0	3	0	99.8	171.6	66.2	...	3.19	3.40	
4	2	5	1	0	0	3	0	99.4	176.6	66.4	...	3.19	3.40	

5 rows × 25 columns

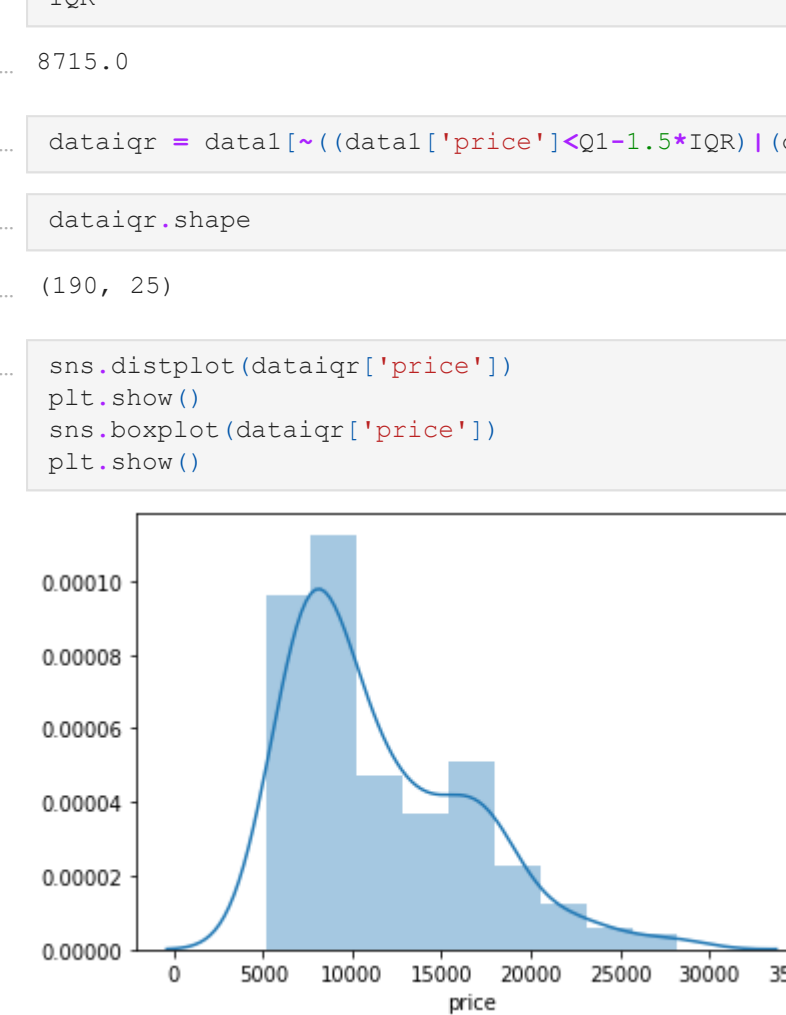
```
In [108]: data.tail()

Out[108]:      symboling  CarName  fueltpe  aspiration  doornumber  carbody  enginelocation  wheelbase  carlength  carwidth  ...  boreratio  stroke
200          -1    139         1          0          0          3          0    109.1    188.8    68.9  ...      3.78    3.15
201          -1    148         1          0          0          3          0    109.1    188.8    68.8  ...      3.78    3.15
202           1    130         1          0          0          3          0    109.1    188.8    68.9  ...      3.58    2.87
203          -1    142         0          1          0          3          0    109.1    188.8    68.9  ...      3.01    3.40
204          -1    143         1          1          0          3          0    109.1    188.8    68.9  ...      3.78    3.15

5 rows × 25 columns
```

Feature selection on the basis of Mutual Information Gain

```
In [109]: sns.distplot(dataal['price'])
plt.show()
sns.boxplot(dataal['price'])
plt.show()
```



```
In [121]: Q1 = dataal['price'].quantile(0.25)
Q3 = dataal['price'].quantile(0.75)
IQR = Q3-Q1
IQR

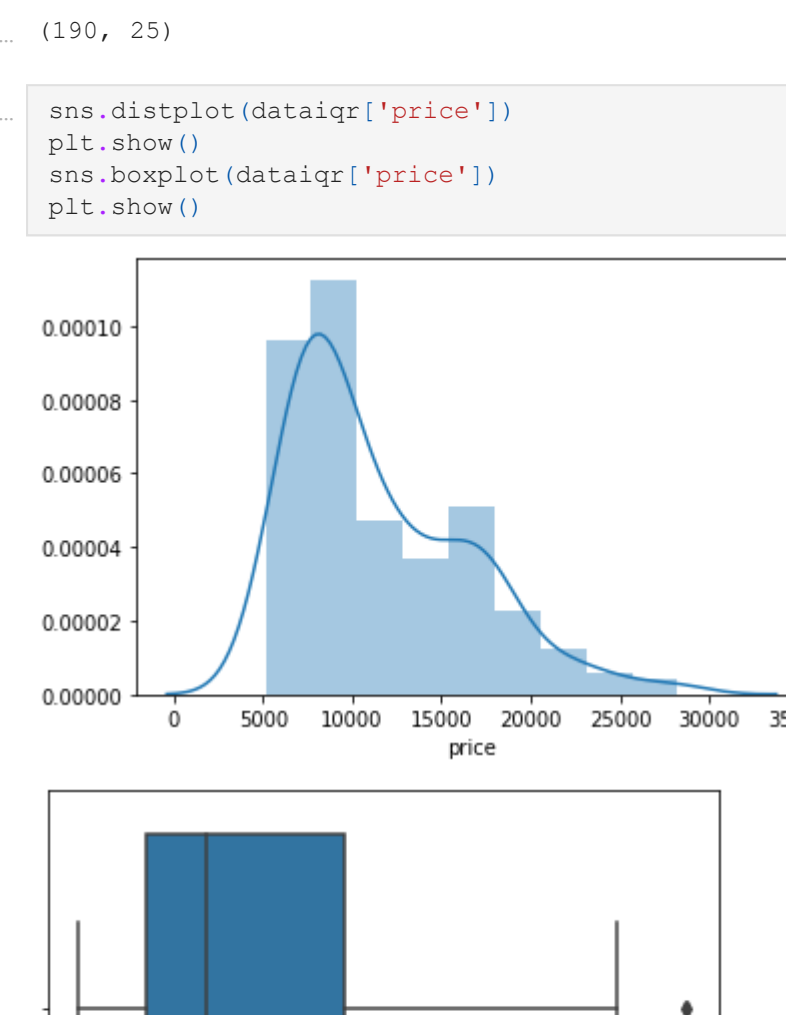
Out[121]: 8715.0

In [122]: dataalqr = dataal[~((dataal['price']<Q1|1.5*IQR| (dataal['price']>Q3|1.5*IQR)))]

In [123]: dataalqr.shape

Out[123]: (190, 25)

In [124]: sns.distplot(dataalqr['price'])
plt.show()
sns.boxplot(dataalqr['price'])
plt.show()
```



Implement the IG

```
In [125]: x = dataal.drop(columns=['price'])
y = dataalqr['price']

In [126]: x.shape,y.shape

Out[126]: ((190, 24), (190,))

In [127]: x_train,x_test,y_train,y_test = train_test_split(x,y,train_size=0.80,random_state=42)

In [128]: x_train.shape,x_test.shape

Out[128]: ((152, 24), (38, 24))

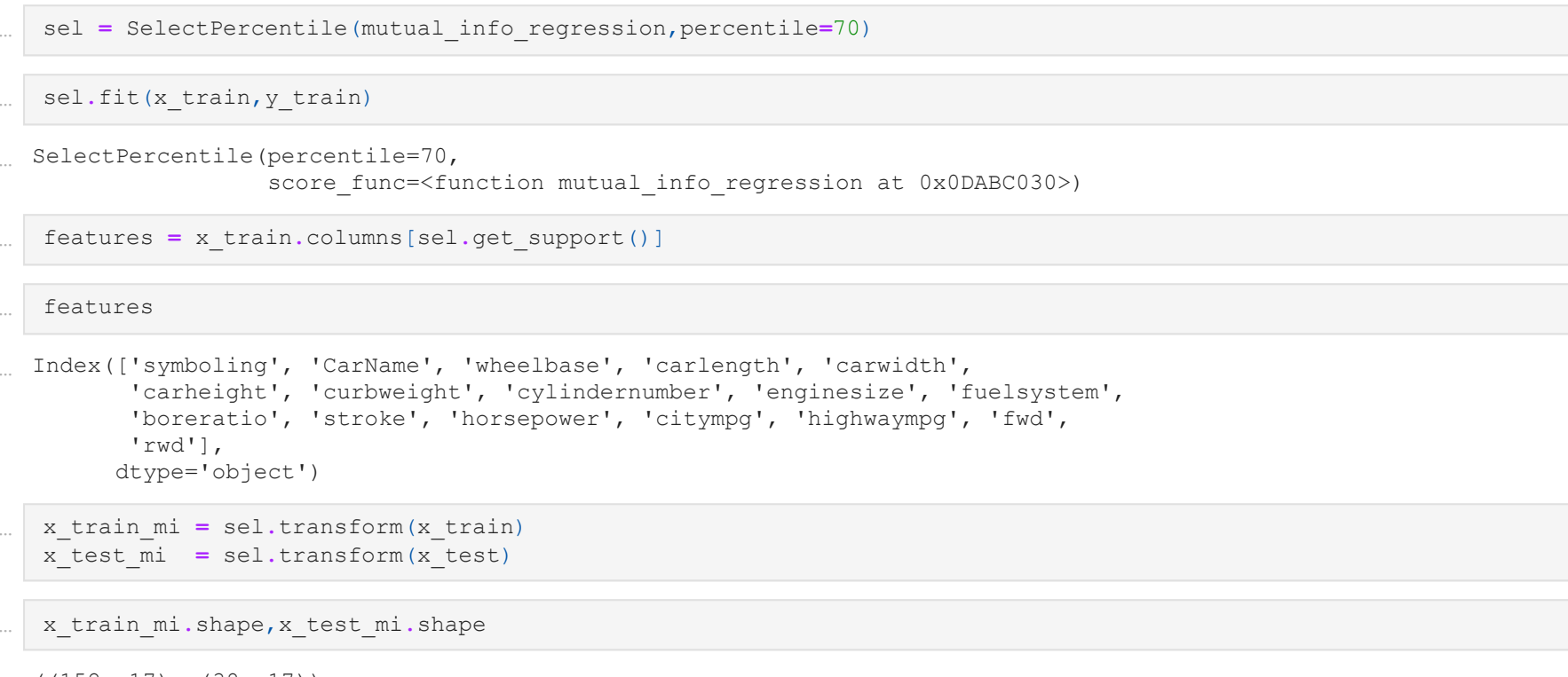
In [129]: mi = mutual_info_regression(x_train,y_train,n_neighbors=3)
mi

Out[129]: array([0.23099133, 0.20735529, 0.01525175, 0.08124114, 0.
0.01140548, 0.
0.31365507, 0.4377656 , 0.22362044, 0.5913376 , 0.17722742,
0.32568874, 0.21814842, 0.
0.58536724, 0.6733307 , 0.21442722, 0.205485 ])
```

```
In [130]: mi = pd.Series(mi,index=x_train.columns)

In [131]: mi.sort_values(ascending=False,inplace=True)

In [134]: plt.figure(figsize=(20,4))
mi.plot(kind='bar')
plt.show()
```



If we see the higher information gain contain by those features which are important for us.

```
In [178]: sel = SelectPercentile(mutual_info_regression,percentile=70)

In [179]: sel.fit(x_train,y_train)

Out[179]: SelectPercentile(percentile=70,
score_func=<function mutual_info_regression at 0x0DABC030>)

In [180]: features = x_train.columns[sel.get_support()]

In [181]: features

Out[181]: Index(['symboling', 'CarName', 'wheelbase', 'carlength', 'carwidth', 'curbweight', 'horsepower', 'engine', 'cylindernumber', 'enginesize', 'fuelsystem', 'boreratio', 'stroke', 'horsepower', 'citympg', 'highwaympg', 'fwd', 'rwd'],
dtype='object')

In [182]: x_train_mi = sel.transform(x_train)
x_test_mi = sel.predict(x_test)

In [183]: x_train_mi.shape,x_test_mi.shape

Out[183]: ((152, 17), (38, 17))

In [184]: from sklearn.linear_model import LinearRegression

In [185]: def LinearModel(x_train,x_test,y_train,y_test):
model = LinearRegression()
model.fit(x_train,y_train)
y_pred = model.predict(x_test_mi)
print('Accuracy :',r2_score(y_test,y_pred))

In [186]: LinearModel(x_train_mi,x_test_mi,y_train,y_test)

Accuracy : 0.81125822578604
Wall time: 0 ns
```

Mutual Infomation Classification Problem.

```
In [187]: data = pd.read_csv('dataset01_eurusd4h.csv')

In [188]: data.head()

Out[188]:      rs1  rs2  rs3  rs4  rs5  rs6  stoch1  stoch2  stoch3  stoch4  ...  WPR6  close1  close2  close3  close4  close5  close6  hour d
1  2890  3107  4001  4051  3995  4198  13.53  29.27  46.80  43.52  ... -56.09  1.3132  1.3132  1.3132  1.3132  1.3132  0
2  2890  3107  4001  4051  3995  4198  13.53  29.27  46.80  43.52  ... -56.93  1.3125  1.3125  1.3125  1.3125  1.3125  4
3  2841  2739  2890  3107  4001  4051  4.27  3.93  13.53  29.27  ... -66.07  1.3127  1.3127  1.3127  1.3127  1.3127  8
4  3448  2841  2739  2890  3107  4001  22.99  4.27  3.93  13.53  ... -74.34  1.3138  1.3138  1.3138  1.3138  1.3138  12
5  3385  3448  2841  2739  2890  3107  14.48  12.49  4.27  3.93  ... -92.30  1.3134  1.3134  1.3134  1.3134  1.3134  16

5 rows × 99 columns

In [189]: data.tail()
```

	rs1	rs2	rs3	rs4	rs5	rs6	stoch1	stoch2	stoch3	stoch4	...	WPR6	close1	close2	close3	close4	close5	close6	hour
4474	6240	6140	6984	7707	8166	7633	26.52	38.85	60.16	74.48	...	-1.88	1.1974	1.1974	1.1974	1.1974	1.1974	1.1974	4
4475	5995	6240	6140	6984	7707	8166	19.21	26.52	38.85	60.16	...	-10.44	1.1964	1.1964	1.1964	1.1964	1.1964	1.1964	8
4476	5654	5995	6240	6140	6984	7707	17.90	19.21	26.52	38.85	...	-15.83	1.1950	1.1950	1.1950	1.1950	1.1950	1.1950	12
4477	4839	5654	5995	6240	6140	6984	15.75	17.90	19.21	26.52	...	-24.95	1.1912	1.1912	1.1912	1.1912	1.1912	1.1912	16
4478	4579	4839	5654	5995	6240	6140	12.88	15.75	17.90	19.21	...	-49.51	1.1898	1.1898	1.1898	1.1898	1.1898	1.1898	20

5 rows × 99 columns

```
In [189]: data.columns

Out[189]: Index(['rs1', 'rs2', 'rs3', 'rs4', 'rs5', 'rs6', 'stoch1', 'stoch2', 'stoch3', 'stoch4', 'WPR6', 'close1', 'close2', 'close3', 'close4', 'close5', 'close6', 'hour'],
dtype='object')

In [191]: len(data.columns)

Out[191]: 99

In [195]: data.dtypes.head(20)

Out[195]: rs1          float64
rs2          float64
rs3          float64
rs4          float64
rs5          float64
rs6          float64
stoch1       float64
stoch2       float64
stoch3       float64
stoch4       float64
stoch5       float64
stoch6       float64
ema20Slope1 float64
ema20Slope2 float64
ema20Slope3 float64
ema20Slope4 float64
ema20Slope5 float64
ema20Slope6 float64
ema50Slope1 float64
ema50Slope2 float64
dtype: object

In [198]: data.isnull().sum()
0

Out[198]: tipo          0
ema100Slope1  0
ema100Slope3  0
ema100Slope4  0
ema100Slope5  0
ema100Slope6  0
ema20Slope1  0
ema20Slope2  0
ema20Slope3  0
ema20Slope4  0
dtype: int64

In [200]: data.nunique()

Out[200]: rs1          2918
rs2          2928
rs3          2929
rs4          2930
rs5          2926
close5       1536
close6       1536
hour         24
dayOfWeek    6
tipo         2
Length: 99, dtype: int64

In [201]: x = data.drop(columns='tipo')
y = data['tipo']

In [202]: x.shape,y.shape

Out[202]: ((4479, 98), (4479,))

In [203]: x_train,x_test,y_train,y_test = train_test_split(x,y,train_size=0.80,stratify=y,random_state=42)

In [204]: x_train.shape,x_test.shape

Out[204]: ((3583, 98), (896, 98))

In [205]: from sklearn.feature_selection import VarianceThreshold

In [206]: vt = VarianceThreshold(threshold=0.01)
x_train_unique = vt.fit_transform(x_train)
x_test_unique = vt.transform(x_test)

In [214]: x_train_unique = pd.DataFrame(x_train_unique)
x_test_unique = pd.DataFrame(x_test_unique)

In [208]: x_train_T = x_train_unique.T

In [209]: x_train_T = pd.DataFrame(x_train_T)

In [210]: x_train_T.duplicated().sum()

Out[210]: 0

In [237]: y_train.value_counts()

Out[237]: 0    1817
1     1766
Name: tipo, dtype: int64

In [248]: scaler = StandardScaler()
x_train_unique = scaler.fit_transform(x_train_unique)
x_test_unique = scaler.transform(x_test_unique)

In [249]: x_train_unique = pd.DataFrame(x_train_unique)
x_test_unique = pd.DataFrame(x_test_unique)

There is no need find the duplicated.

In [250]: mi =mutual_info_classif(x_train_unique,y_train,n_neighbors=3,random_state=42)

In [251]: mi

Out[251]: array([0.00787631, 0.
0.
0.00874687, 0.
0.
0.0141222 , 0.00506768, 0.01428357,
0.
0.00821254, 0.00745152, 0.
0.00243273,
0.00053892, 0.00826799, 0.
0.01421783, 0.01113077,
0.00346849, 0.
0.
0.
0.
0.00484005,
0.01190763, 0.
0.01351362, 0.00897282, 0.00416849,
0.0155984 , 0.
0.
0.0047066 ])
```

```
In [252]: mi = pd.Series(mi,index=x_train_unique.columns)

In [253]: mi.sort_values(ascending=False,inplace=True)

In [254]: plt.figure(figsize=(20,4))
mi.plot(kind='bar')
plt.show()
```



```
In [265]: sel = SelectKBest(mutual_info_classif,k=30)

In [266]: sel.fit(x_train_unique,y_train)

In [267]: SelectKBest(k=30, score_func=<function mutual_info_classif at 0x0DABC078>)

In [268]: features = x_train_unique.columns[sel.get_support()]

In [269]: mi

Out[269]: 40 0.015510
14 0.014284
23 0.014218
12 0.014112
37 0.013514
30 0.011908
24 0.011131
3 0.010968
31 0.009497
7 0.008747
21 0.008268
16 0.008213
0 0.007876
43 0.004707
39 0.004168
25 0.003468
19 0.002433
20 0.000539
4 0.000000
5 0.000000
6 0.000000
2 0.000000
8 0.000000
9 0.000000
31 0.000000
41 0.000000
21 0.000000
36 0.000000
29 0.000000
18 0.000000
22 0.000000
32 0.000000
26 0.000000
27 0.000000
15 0.000000
dtype: float64

In [269]: features

Out[269]: Int64Index([ 0, 6, 7, 12, 13, 14, 15, 19, 20, 21, 23, 24, 25, 27, 28, 29, 30,
31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43],
dtype='int64')

In [270]: x_train_sel = sel.transform(x_train_unique)
x_test_sel = sel.transform(x_test_unique)

In [271]: x_train_sel.shape,x_test_sel.shape

Out[271]: ((3583, 30), (896, 30))

In [281]: from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression

In [279]: def SVCClassifier(x_train,x_test,y_train,y_test):
clf = SVC(C=0.01, kernel='sigmoid',degree=1)
clf.fit(x_train,y_train)
y_pred = clf.predict(x_test)
print('Accuracy :',accuracy_score(y_test,y_pred))

In [280]: SVCClassifier(x_train_sel,x_test_sel,y_train,y_test)

Accuracy :- 0.5066964285714286
Wall time: 4.26 s

In [282]: def Logistics(x_train,x_test,y_train,y_test):
clf = LogisticRegression(random_state=42,max_iter=200,solver='lbfgs',verbose=1)
y_pred = clf.predict(x_test)
print('Accuracy :',accuracy_score(y_test,y_pred))

In [400]: Logistics(x_train_sel,x_test_sel,y_train,y_test)

Accuracy:- 0.5022321428571429
Wall time: 92 ms

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.0s finished
```