

▼ Data Mining - Lab 06 - Feature Scaling and Regression

Nguyễn Đức Vũ Duy - 18110004

```
1 #Import libraries
2 # Supress Warnings
3
4 import warnings
5 warnings.filterwarnings('ignore')
6
7 # Import the numpy and pandas package
8 import numpy as np
9 import pandas as pd
10
11 # Data Visualisation
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14
15 # Import library for label encoder
16 from sklearn.preprocessing import LabelEncoder
17
18 # Import library for scaling
19 from sklearn.preprocessing import MinMaxScaler
20
21 # Import library for splitting training and testing data
22 from sklearn.model_selection import train_test_split
23
24 # Import library of RMSE
25 from sklearn.metrics import mean_squared_error
26
27 # Import libraries for model Building
28 from sklearn.linear_model import LinearRegression
29 from sklearn.linear_model import Ridge, Lasso
30 from sklearn.tree import DecisionTreeRegressor
31 from sklearn.neighbors import KNeighborsRegressor
32 from sklearn.ensemble import RandomForestRegressor
```

Sau khi làm lại 2 bài trên thì áp dụng kết hợp các bước ở trên để thực hiện giải bài toán Regression tốt nhất cho dữ liệu sau "insurance.csv". Cần dự đoán cột charges.

Trong bài tập "Regression Methods (Customer Churn Analysis)" ta có các bước :
Data Integration

Data Cleansing

EDA

Encoding

Feature Scaling

Feature Selection (Feature Importance)

Model Building

Model Evaluation

▼ Data Integration

```
1 #Import dataset
2 path='https://raw.githubusercontent.com/duynguyenhcmus/Repository/main/Hock
3 df=pd.read_csv(path)
4 df.head()
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520

```
1 #Print data information
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  -
0   age         1338 non-null   int64
1   sex         1338 non-null   object
2   bmi         1338 non-null   float64
3   children    1338 non-null   int64
4   smoker      1338 non-null   object
5   region      1338 non-null   object
6   charges     1338 non-null   float64
```

```
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
```

```
1 #Print data describe
2 df.describe()
```

	age	bmi	children	charges
count	1338.000000	1338.000000	1338.000000	1338.000000
mean	39.207025	30.663397	1.094918	13270.422265
std	14.049960	6.098187	1.205493	12110.011237
min	18.000000	15.960000	0.000000	1121.873900
25%	27.000000	26.296250	0.000000	4740.287150
50%	39.000000	30.400000	1.000000	9382.033000
75%	51.000000	34.693750	2.000000	16639.912515
max	64.000000	53.130000	5.000000	63770.428010

▼ Data cleaning

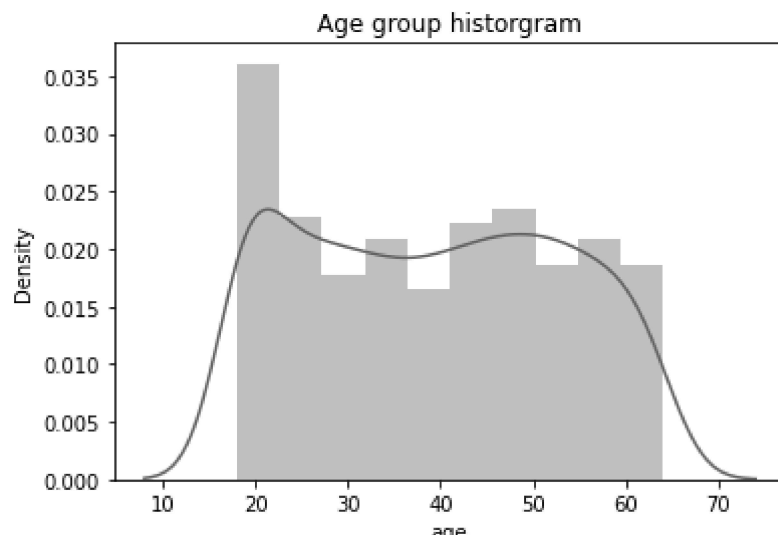
```
1 #Check for null values
2 df.isnull().sum()*100/df.shape[0]
3 #No null values
```

```
age          0.0
sex          0.0
bmi          0.0
children     0.0
smoker       0.0
region       0.0
charges      0.0
dtype: float64
```

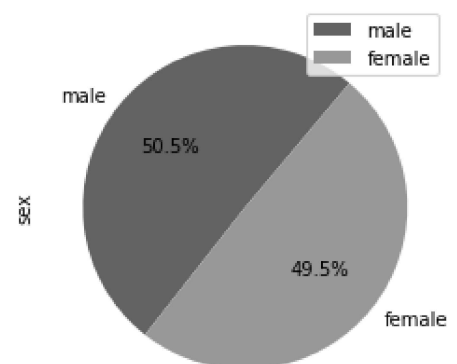
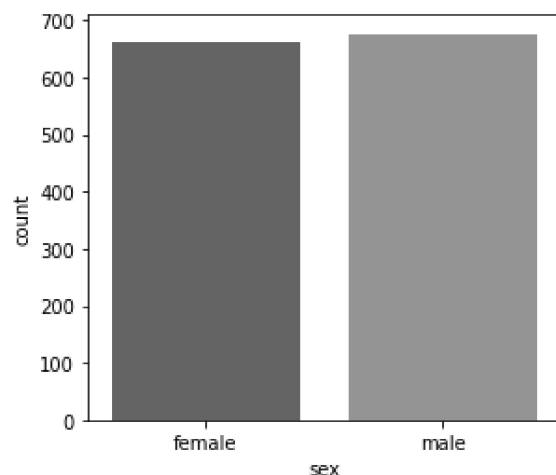
▼ EDA

```
1 #Age features
2 fig,ax=plt.subplots()
3 sns.distplot(a=df.age, bins=10,ax=ax)
4 plt.title('Age group histogram')
5 #Not a normal distribution with a majority of people is at the age range of
```

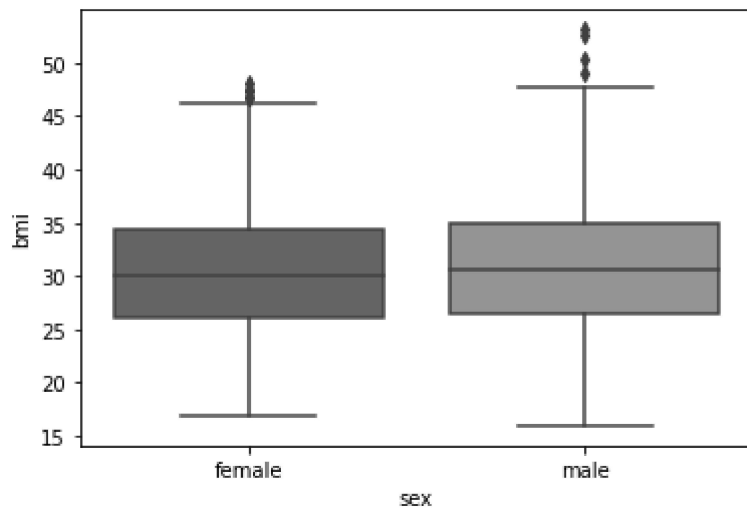
```
Text(0.5, 1.0, 'Age group histogram')
```



```
1 #Sex features
2 fig, axs = plt.subplots(1,2, figsize = (10,4))
3 plt1 = sns.countplot(df['sex'], ax = axs[0])
4
5 pie_churn = pd.DataFrame(df['sex'].value_counts())
6 pie_churn.plot.pie( subplots=True, labels = pie_churn.index.values, autopct=
7 # Unsquish the pie.
8 plt.gca().set_aspect('equal')
9
10 plt.show()
11 #It's balance between 2 sex group
```

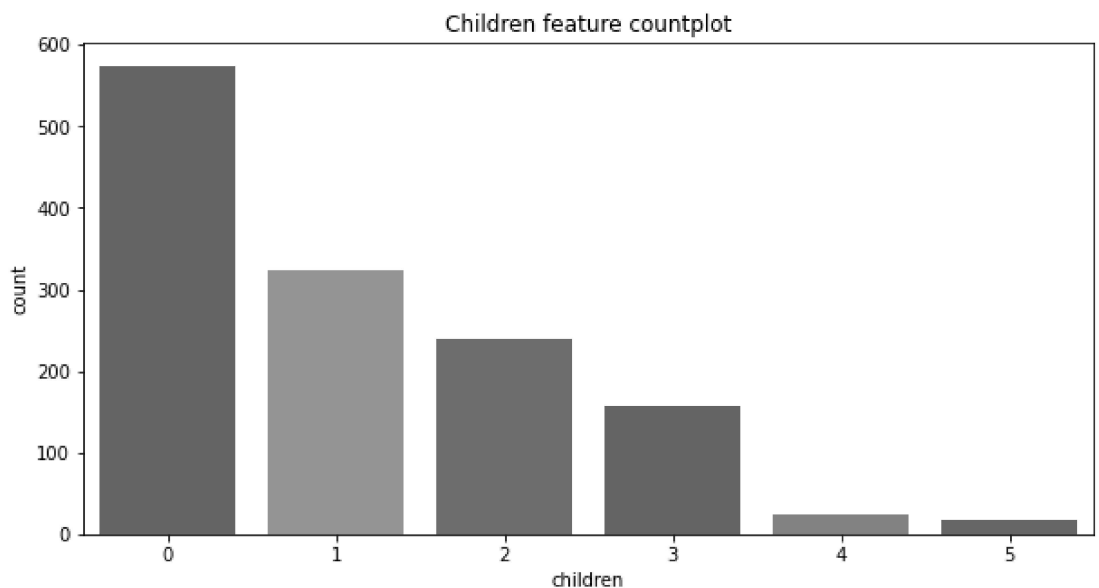


```
1 #bmi features
2 sns.boxplot(x = 'sex', y = 'bmi', data = df)
3 plt.show()
4 #Male tend to get more outlier than woman. BMI score between 2 sexes is bal
```



```
1 #children feature
2 fig, axs = plt.subplots(figsize = (10,5))
3 sns.countplot(df['children'], ax = axs)
4 plt.title('Children feature countplot')
5 #It's not balance with children at group 4 and 5
```

```
Text(0.5, 1.0, 'Children feature countplot')
```

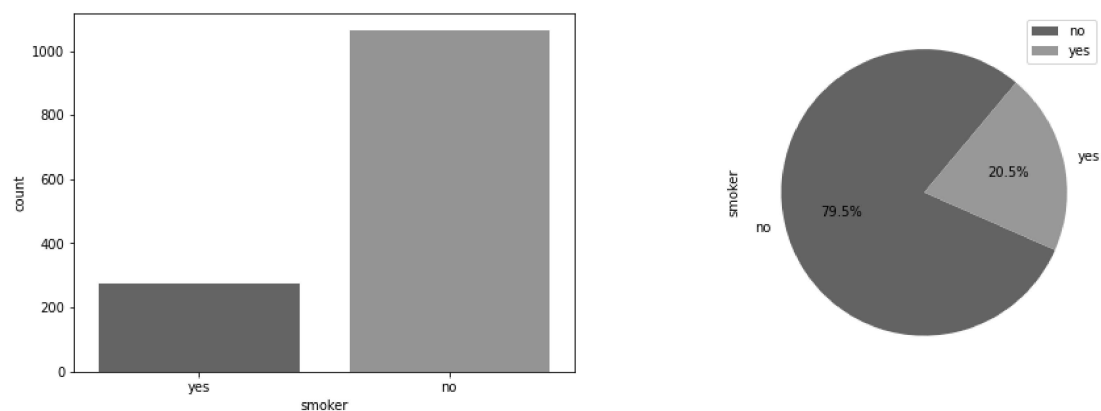


```
1 #smoking feature
2 fig, axs = plt.subplots(1,2, figsize = (15,5))
3 plt1 = sns.countplot(df['smoker'], ax = axs[0])
4
5 pie_churn = pd.DataFrame(df['smoker'].value_counts())
6 pie_churn.plot.pie( subplots=True, labels = pie_churn.index.values, autopct=
7 # Unsquish the pie.
8 plt.gca().set_aspect('equal')
```

9

```
10 plt.show()
```

```
11 #A majority of people is non-smoker
```



```
1 #region feature
```

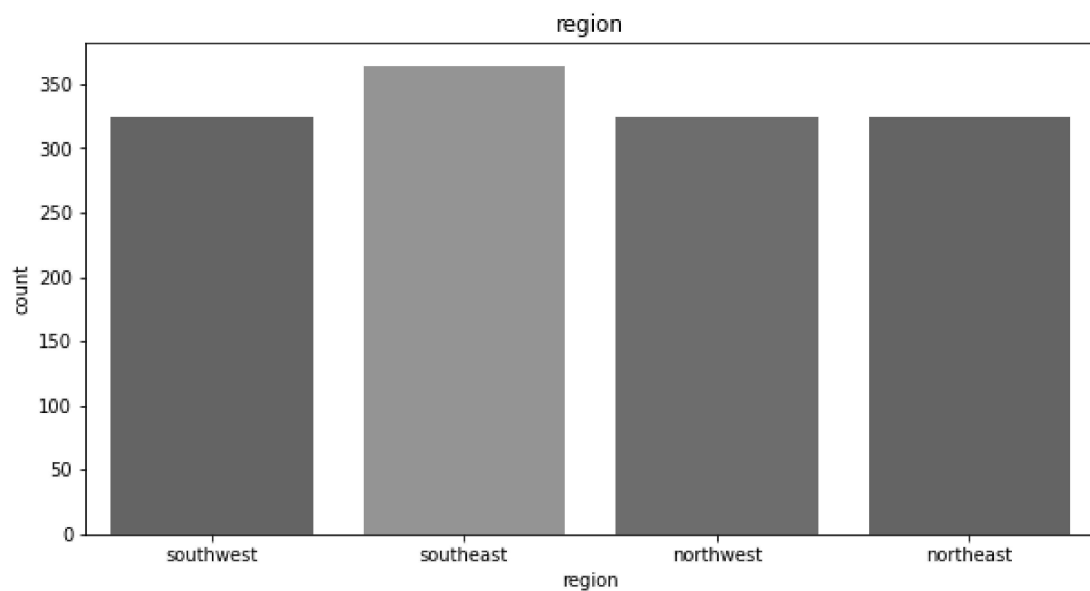
```
2 fig, axs = plt.subplots(figsize = (10,5))
```

```
3 sns.countplot(df['region'], ax = axs)
```

```
4 plt.title('region')
```

```
5 #It's quite balance between the number of different groups
```

```
Text(0.5, 1.0, 'region')
```

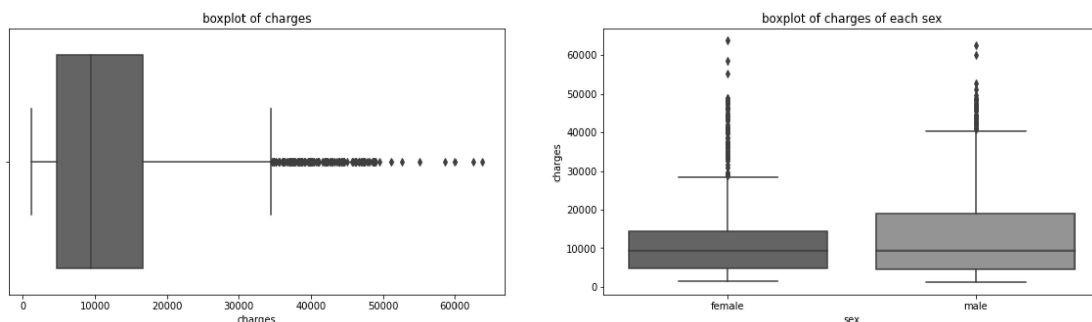


```
1 #Charge as target
```

```

2 fig, axs = plt.subplots(1,2, figsize = (20,5))
3 plt1 = sns.boxplot(df['charges'], ax = axs[0])
4 plt1.set_title('boxplot of charges')
5 plt2 = sns.boxplot(data=df, x='sex',y='charges',ax=axs[1])
6 plt2.set_title('boxplot of charges of each sex')
7
8 plt.show()
9 #Both sex and each sex have a lot of outliers

```



```

1 #Removing outliers
2 Q1 = df.quantile(0.25)
3 Q3 = df.quantile(0.75)
4 IQR = Q3 - Q1
5 print('> Shape of data before handling outlier values: ', df.shape)
6 df = df[~((df<(Q1-1.5*IQR))|(df>(Q3+1.5*IQR))).any(axis=1)]
7 print('> Shape of data after handling outlier values: ', df.shape)

> Shape of data before handling outlier values: (1338, 7)
> Shape of data after handling outlier values: (1193, 7)

```

▼ Encoding

```

1 #Using label encoding from sklearn
2 label = LabelEncoder()
3 data_columns = df.dtypes.pipe(lambda X: X[X=='object']).index
4 for col in data_columns:
5     df[col] = label.fit_transform(df[col])
6

```

```

7 df.info()
8 #It's all numeric data now

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1193 entries, 0 to 1337
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  -
0    age         1193 non-null   int64
1    sex         1193 non-null   int64
2    bmi         1193 non-null   float64
3    children    1193 non-null   int64
4    smoker      1193 non-null   int64
5    region      1193 non-null   int64
6    charges     1193 non-null   float64
dtypes: float64(2), int64(5)
memory usage: 74.6 KB

```

▼ Feature scaling

```

1 #Train test split for feature scaling
2 X=df.iloc[:, :-1].copy()
3 y=df.iloc[:, -1].copy()
4 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,rand

```

```

1 # copy of datasets
2 X_train_norm = X_train.copy()
3 X_test_norm = X_test.copy()
4
5 # numerical features
6 num_cols = ['age', 'bmi']
7 # apply standardization on numerical features
8 for i in num_cols:
9     # fit scaler on training data
10    norm = MinMaxScaler().fit(X_train_norm[[i]])
11    # transform the training data column
12    X_train_norm[i] = norm.transform(X_train_norm[[i]])
13    # transform the testing data column
14    X_test_norm[i] = norm.transform(X_test_norm[[i]])

```

▼ Feature Selection

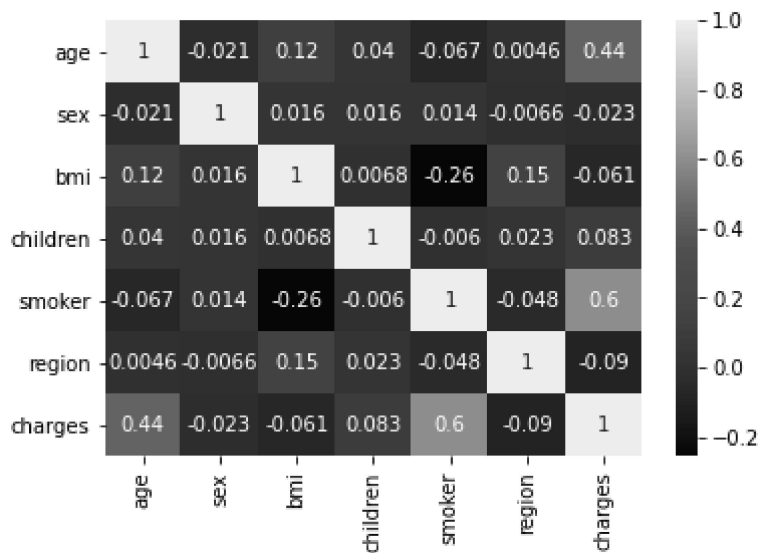
```

1 # Let's see the correlation matrix
2 plt.figure()
3 sns.heatmap(df.corr(),annot = True)

```



```
4 plt.show()
5 #The correlation between features is quite low. So there is no implication
6 #However, sex, bmi, children and region have very low correlation with the
```



```
1 #Remove uncorrelated features
2 df=df.drop(['sex','bmi','children','region'],axis=1)
3 df.head()
```

	age	smoker	charges
0	19	1	16884.92400
1	18	0	1725.55230
2	28	0	4449.46200
3	33	0	21984.47061
4	32	0	3866.85520

Hãy sử dụng các thuật toán Regression : Linear, Ridge, Lasso, Decision Tree, K Neighbour, Random Forest Regression và so sánh RMSE.

▼ Model building

```
1 #Build Model
2 rmse=[]
3
```

```

4 #Linear Regression
5 linear_regression=LinearRegression().fit(X_train, y_train)
6 y_linear_regression=linear_regression.predict(X_test)
7 rmse.append(mean_squared_error(y_linear_regression,y_test,squared=False))
8
9 #Ridge Regression
10 ridge_regression=Ridge(alpha=1).fit(X_train,y_train)
11 y_ridge_regression=ridge_regression.predict(X_test)
12 rmse.append(mean_squared_error(y_ridge_regression,y_test,squared=False))
13
14 #Lasso Regression
15 lasso_regression=Lasso(alpha=1).fit(X_train,y_train)
16 y_lasso_regression=lasso_regression.predict(X_test)
17 rmse.append(mean_squared_error(y_lasso_regression,y_test,squared=False))
18
19 #Decision Tree Regression
20 decision_tree=DecisionTreeRegressor(max_depth=5, random_state=0).fit(X_train,y_train)
21 y_decision_tree=decision_tree.predict(X_test)
22 rmse.append(mean_squared_error(y_decision_tree,y_test,squared=False))
23
24 #K Neighbors Regression
25 k_neighbors=KNeighborsRegressor(n_neighbors=10).fit(X_train,y_train)
26 y_k_neighbors=k_neighbors.predict(X_test)
27 rmse.append(mean_squared_error(y_k_neighbors,y_test,squared=False))
28
29 #Random forest regression
30 random_forest=RandomForestRegressor(n_estimators=50).fit(X_train,y_train)
31 y_random_forest=random_forest.predict(X_test)
32 rmse.append(mean_squared_error(y_random_forest,y_test,squared=False))
33

```

▼ Model evaluation

```

1 # visualizing the result
2 df_dt = pd.DataFrame({'RMSE':rmse},index=['Linear regression','Ridge Regression','Lasso Regression','Decision Tree Regression','K Neighbors Regression','Random Forest Regression'])
3 df_dt
4
5

```

	RMSE
Linear regression	4626.282460
Ridge Regression	4624.517125
Lasso Regression	4625.986412
Decision Tree Regression	4542.606272
K Neighbors Regression	5998.767660
Random Forest Regression	4879.058564