

```
import sqlalchemy as sa
import pandas as pd
```

```
pip install pymysql
```

```
Collecting pymysqlNote: you may need to restart the kernel to use updated packages.

Obtaining dependency information for pymysql from https://files.pythonhosted.org/packages/e5/30/20467e39523d0cfc2b6227902d3687a1636436
Downloading PyMySQL-1.1.0-py3-none-any.whl.metadata (4.4 kB)
Downloading PyMySQL-1.1.0-py3-none-any.whl (44 kB)
----- 0.0/44.8 kB ? eta -:--:--
----- 20.5/44.8 kB 330.3 kB/s eta 0:00:01
----- 41.0/44.8 kB 393.8 kB/s eta 0:00:01
----- 44.8/44.8 kB 276.2 kB/s eta 0:00:00
Installing collected packages: pymysql
Successfully installed pymysql-1.1.0
```

```
### create_engine('mysql+pymysql://username:password@hostname:port number/db name')
engine=sa.create_engine("mysql+pymysql://root:Arc111han555@localhost:3306/db_dev")
engine

Engine(mysql+pymysql://root:***@localhost:3306/db_dev)
```

Importing Libraries:

```
df=pd.read_sql_table('medical_insurance',engine)
df
```

	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
...
1333	50	male	30.970	3	no	northwest	10600.54830
1334	18	female	31.920	0	no	northeast	2205.98080
1335	18	female	36.850	0	no	southeast	1629.83350
1336	21	female	25.800	0	no	southwest	2007.94500
1337	61	female	29.070	0	yes	northwest	29141.36030

1338 rows × 7 columns

```
import pandas as pd
import numpy as np

from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder

from sklearn.metrics import mean_squared_error,mean_absolute_error, r2_score

import matplotlib.pyplot as plt
import seaborn as sns

import pickle
import json

# Loading Train dataset:
train_data = df
```

```
# Shape of dataset:
train_data.shape

(1338, 7)

# Cheacking for NaN Values (Missing Values):
train_data.isnull().sum()

age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64
```

```
# Insights of dataset:
train_data.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
```

```
# Description of dataset:
train_data.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

Double-click (or enter) to edit

```
encoder = LabelEncoder()
labels = encoder.fit_transform(train_data.sex)

train_data['sex'] = labels

train_data.head()
```

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

```
labels = encoder.fit_transform(train_data.region)
```

```
train_data['region'] = labels
train_data.head()
```

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

```
labels = encoder.fit_transform(train_data.smoker)
```

```
train_data['smoker'] = labels
train_data.head()
```

	age	sex	bmi	children	smoker	region	charges
0	19	0	27.900	0	1	3	16884.92400
1	18	1	33.770	1	0	2	1725.55230
2	28	1	33.000	3	0	2	4449.46200
3	33	1	22.705	0	0	1	21984.47061
4	32	1	28.880	0	0	1	3866.85520

```
train_data.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   int32
2    bmi         1338 non-null   float64
3    children    1338 non-null   int64
4    smoker      1338 non-null   int32
5    region      1338 non-null   int32
6    charges     1338 non-null   float64
dtypes: float64(2), int32(3), int64(2)
memory usage: 57.6 KB
```

▼ Train Test split

```
df = train_data.select_dtypes(exclude=object)
x = train_data.drop('charges',axis = 1)
y = train_data['charges']
x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.2,random_state=34)
```

```
x_train
```

	age	sex	bmi	children	smoker	region
414	19	0	35.150	0	0	1
1279	25	0	26.790	2	0	1
647	40	0	23.370	3	0	0
764	45	0	25.175	2	0	0

```
x_train.columns
```

```
Index(['age', 'sex', 'bmi', 'children', 'smoker', 'region'], dtype='object')
```

```
df
```

	age	sex	bmi	children	smoker	region	charges
0	19	0	27.900	0	1	3	16884.92400
1	18	1	33.770	1	0	2	1725.55230
2	28	1	33.000	3	0	2	4449.46200
3	33	1	22.705	0	0	1	21984.47061
4	32	1	28.880	0	0	1	3866.85520
...
1333	50	1	30.970	3	0	1	10600.54830
1334	18	0	31.920	0	0	0	2205.98080
1335	18	0	36.850	0	0	2	1629.83350
1336	21	0	25.800	0	0	3	2007.94500
1337	61	0	29.070	0	1	1	29141.36030

```
1338 rows × 7 columns
```

▼ Model Training

```
model = LinearRegression()
model.fit(x_train, y_train)
```

```
▼ LinearRegression
LinearRegression()
```

```
# Testing Data Evaluation
y_pred = model.predict(x_test)
```

```
mse = mean_squared_error(y_test, y_pred)
print("MSE :",mse)
```

```
rmse = np.sqrt(mse)
print("RMSE :",rmse)
```

```
mae = mean_absolute_error(y_test, y_pred)
print("MAE :",mae)
```

```
r2 = r2_score(y_test, y_pred)
print('R-Squared :',r2)
```

```
MSE : 41271154.57832547
RMSE : 6424.26295992976
MAE : 4410.013263731576
R-Squared : 0.7461578203319277
```

```
# Training Data Evaluation
```

```
y_pred_train = model.predict(x_train)
mse = mean_squared_error(y_train, y_pred_train)
print("MSE :",mse)
```

```
rmse = np.sqrt(mse)
print("RMSE :",rmse)

mae = mean_absolute_error(y_train, y_pred_train)
print("MAE :",mae)

r2 = r2_score(y_train, y_pred_train)
print('R-Squared :',r2)

MSE : 35365859.39407277
RMSE : 5946.92016039166
MAE : 4094.433690405064
R-Squared : 0.7514552383513079

filename = 'medical_insurance_cost_predictor.pkl'
pickle.dump(model, open(filename,'wb'))
```