Linear Regression

- ◆ Linear regression predicts a real value output based on the input value.
 - Simple linear regression only has one input/variable and we try to fit a line to best describe the output and input relationship
 - Multiple linear regression has multiple variables as input and we try to fit a hyperplane instead.
- ♦ A example of simple linear regression is salary prediction, with the salary as output, and years of working experience as input variable.
- ◆ The algorithm can be used for simple linear regression or multiple linear regression, and in this notebook, it will be implemented for simple linear regression using sklearn linear regression algorithm

Import libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn.linear_model import LinearRegression
```

Read data

```
In [2]:
    df = pd.read_csv('Salary_Data.csv')
    # shuffle data
    df = df.sample(frac=1, random_state=42).reset_index(drop=True)

In [3]:
    df.head()
```

```
      Out[3]:
      YearsExperience
      Salary

      0
      9.6
      112635.0

      1
      4.9
      67938.0

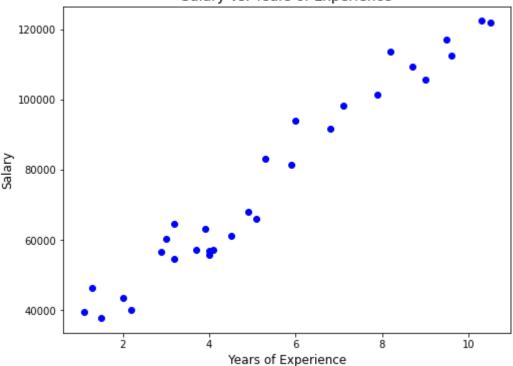
      2
      8.2
      113812.0

      3
      5.3
      83088.0

      4
      3.2
      64445.0
```

```
In [4]: # visualize data
fig = plt.figure(figsize=(8,6))
plt.scatter(df.iloc[:,0], df.iloc[:,1], c='b')
plt.title('Salary vs. Years of Experience', fontsize=14)
plt.xlabel('Years of Experience', fontsize=12)
plt.ylabel('Salary', fontsize=12)
plt.show()
```

Salary vs. Years of Experience



Split data to train and test dataset

```
In [5]: x = np.array(df.iloc[:,0:1])
y = np.array(df.iloc[:,1])

In [6]: train_frac = 0.7
    train_size = int(df.shape[0] * train_frac)
    x_train = x[0:train_size,:]
    y_train = y[0:train_size]
    x_test = x[train_size:,:]
    y_test = y[train_size:]

In [7]: print ('training size is {}\ntest size is {}\'.format(x_train.shape[0], x_test.shape[0])
    training size is 21
    test size is 9

Initialize the model
```

train_accuracy = lr.score(x_train, y_train)*100

Predict the test data output

lr = LinearRegression()
lr.fit(x_train, y_train)

In [8]:

```
y_pred_train = lr.predict(x_train)
y_pred_test = lr.predict(x_test)
```

Visulize fitted results

```
In [11]:
    fig = plt.figure(figsize=(8,6))
    plt.scatter(df.iloc[:,0], df.iloc[:,1], c='b')
    plt.plot(x_train, y_pred_train, c='r', lw=2)
    plt.title('', fontsize=14)
    plt.xlabel('Years of Experience', fontsize=12)
    plt.ylabel('Salary', fontsize=12)
    plt.legend(['fitted line using sklearn LinearRegression'], fontsize=12)
    plt.show()
```

