

LAB4 tutorial for Machine Learning

Linear Regression

The document description are designed by Jla Yanhong in 2022. Sept. 22th

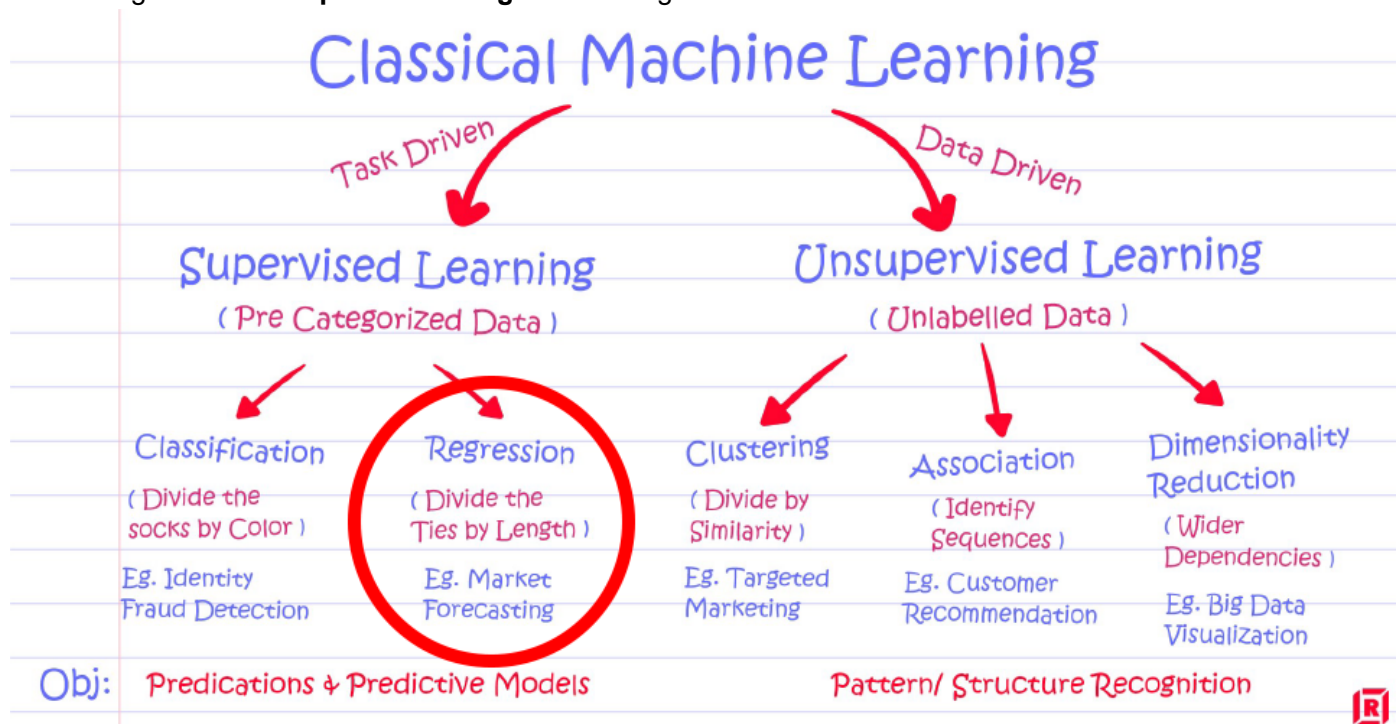
Objective

- Master the linear regression algorithm.
- Understanding Gradient Descent
- Polynomial Regression
- Learn how to evaluate regression models
- Complete the LAB assignment and submit it to BB.

1 Linear Regression

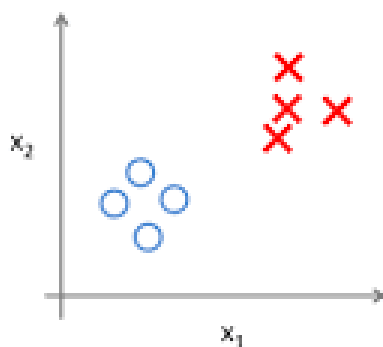
1.1 Conceptual Overview

Linear Regression is “supervised” “regression” algorithm.

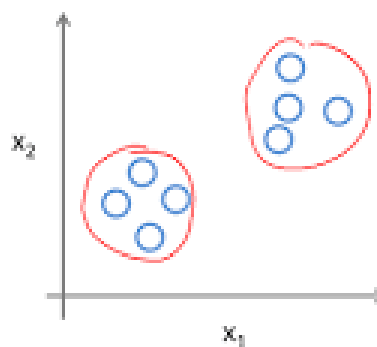


Supervised meaning we use labeled data to train the model.

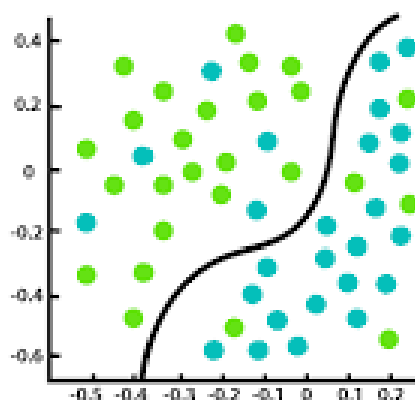
Supervised Learning



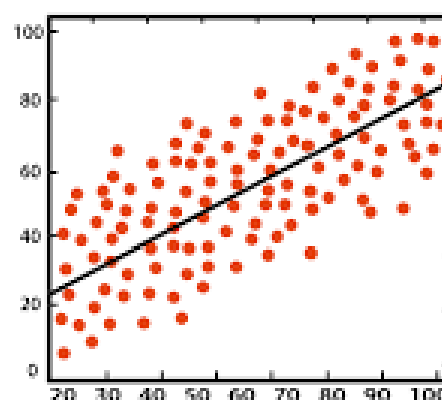
Unsupervised Learning



Regression meaning we predict a numerical value, instead of a “class”.



Classification



Regression

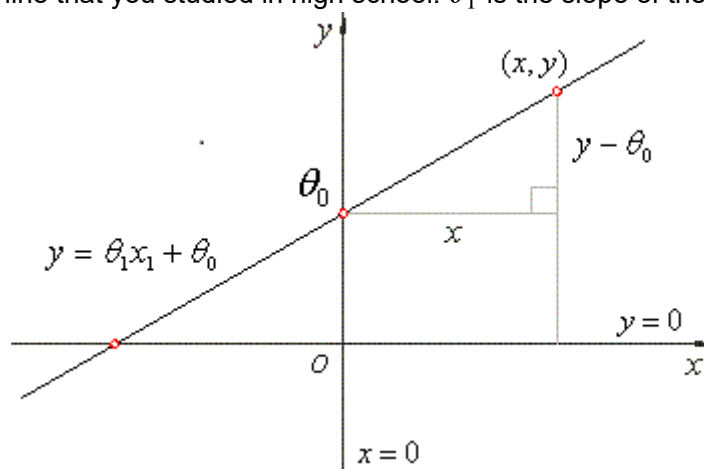
1.2 Linear Regression

In statistics, linear regression is a linear approach to modelling the relationship between a dependent variable and one or more independent variables.

We will define a linear relationship between these two variables as follows:

$$y = \theta_0 + \theta_1 x_1$$

This is the equation for a line that you studied in high school. θ_1 is the slope of the line and θ_0 is the y intercept.



When there are multiple independent variables, the linear relationship becomes as follows:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Simple Linear Regression

$$y = \theta_0 + \theta_1 x_1$$

Multiple Linear Regression

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Let \vec{x} be the independent variable and \mathbf{y} be the dependent variable. The general form of the model's prediction:

$$\hat{y} = h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \theta^T \cdot x$$

θ is the model's parameter vector, containing the bias term θ_0 and the feature weights θ_1 to θ_n . (n = number of features).

$\theta^T \cdot x$ is the dot product of the vectors θ and x , which is, of course, equal to $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$.

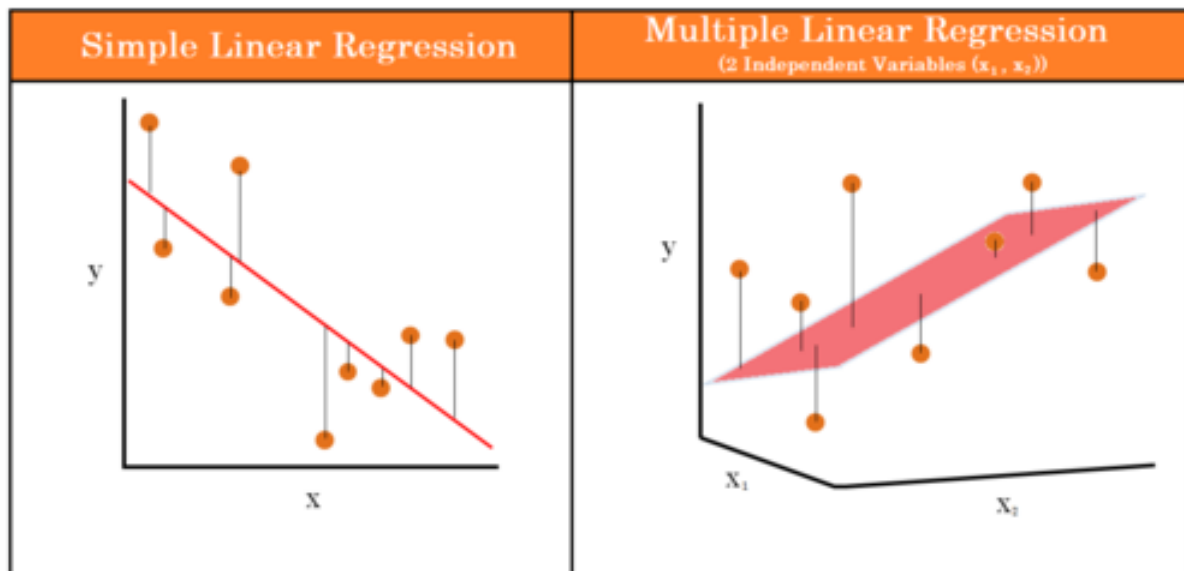
There is an error between the real y and the predicted \hat{y} :

$$y = \hat{y} + \varepsilon = h_{\theta}(x) + \varepsilon = \theta^T \cdot x + \varepsilon$$

Our challenge today is to determine the value of θ .

For a single variable, such that the line corresponding to those values is the best fitting line .

For multiple variables, such as 2 variables, the surface corresponding to these values is the surface of best fit.



1.3 Cost Function

Cost Function evaluates the model's predictions and tells us how accurate are the model's predictions. The lower the value of the cost function, better accurate the predictions of the model. There are many cost functions to choose, but we will use the Mean Squared Error (MSE) cost function.

The MSE function calculates the average of the squared difference between the prediction and the actual value (y).

$$J(\theta) = MSE = \frac{1}{2m} \sum_{j=1}^m (h_{\theta}(x^j) - y^j)^2 = \frac{1}{2m} \sum_{j=1}^m (\theta_0 + \theta_1 x_1^j + \theta_2 x_2^j + \dots + \theta_n x_n^j - y^j)^2$$

x_i^j : The i_{th} feature of the j^{th} sample

y_i^j : the i_{th} label of the j^{th} sample

(x^j, y^j) : the j^{th} sample

m : total number of instances in your dataset

1.4 Least-squares estimation

Now that we have determined the cost function, the only thing left to do is minimize it. This is done by finding the partial derivative of $J(\Theta)$, equating it to 0 and then finding an expression for Θ .

The loss function can be written as:

$$\begin{aligned} J(\theta) &= \frac{1}{2m} \sum_{j=1}^m (h_{\theta}(x^j) - y^j)^2 \\ &= \frac{1}{2m} \sum_{j=1}^m (\theta^T \cdot x^j - y^j)^2 \\ &= \frac{1}{2m} (X\theta - y)^T (X\theta - y) \end{aligned}$$

As the loss is convex the optimum solution lies at gradient zero. The gradient of the loss function is :

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta} &= \frac{\partial \frac{1}{2m} (X\theta - y)^T (X\theta - y)}{\partial \theta} \\ &= \frac{\partial \frac{1}{2m} (\theta^T X^T - y^T)(X\theta - y)}{\partial \theta} \\ &= \frac{\partial \frac{1}{2m} (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y)}{\partial \theta} \\ &= \frac{1}{2m} (2X^T X\theta - X^T y - (y^T X)^T) \\ &= \frac{1}{m} X^T X\theta - X^T y \end{aligned}$$

When X is a matrix of full rank, setting the gradient to zero produces the optimum parameter:

$$\theta^* = (X^T X)^{-1} X^T y$$

$\begin{aligned} X = \end{aligned}$

$$\begin{aligned} \left\{ \begin{array}{cccccccccccccccc} 1 & x_1^1 & x_2^1 & \dots & x_n^1 & 1 & x_1^2 & x_2^2 & \dots & x_n^2 & \dots & \dots & 1 & x_1^m & x_2^m & \dots & x_n^m \end{array} \right\} \\ , y = \end{aligned}$$

$$\left\{ y_0 \ y_1 \ y_2 \ \dots \ y_m \right\}$$

$\theta =$

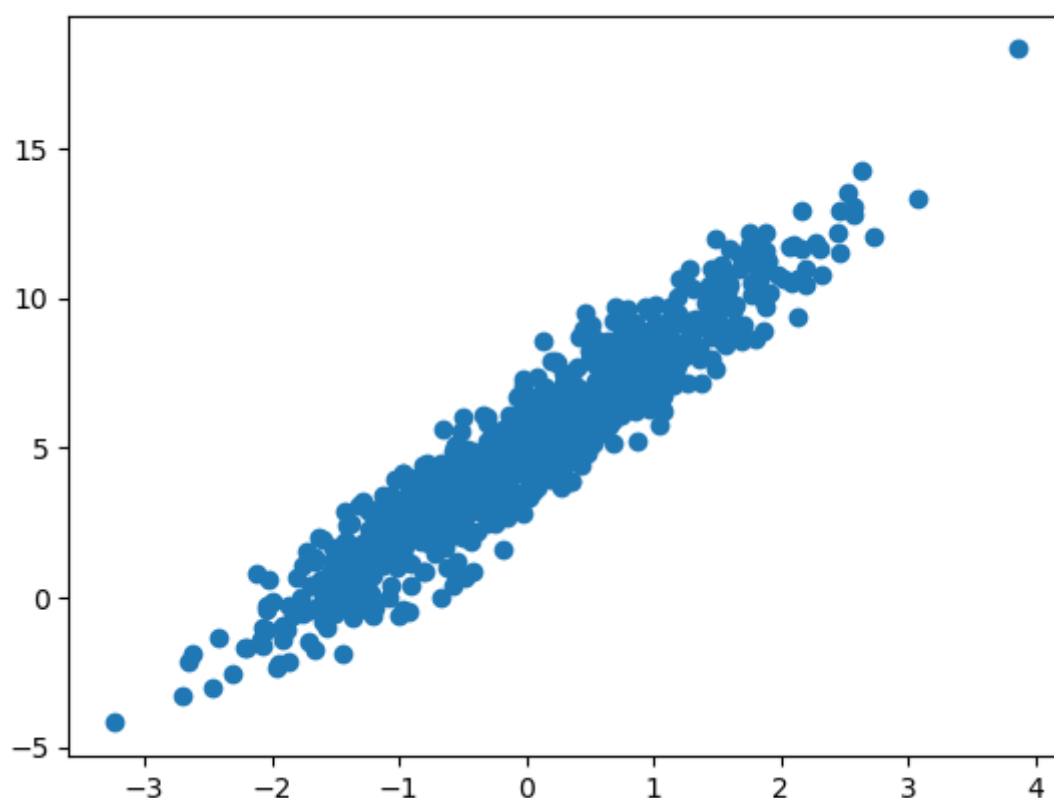
$$\left\{ \theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n \right\}$$

$\end{aligned}$

Hands-on Coding

In [105]:

```
1 import numpy as np
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import r2_score
4
5 import matplotlib.pyplot as plt
6 np.random.seed(42)
7 X = np.random.randn(1000, 1)
8
9 # y = 5 + 3 * X + Gaussian noise -> because in real-world it is very unlikely to get data that
10 y = 5 + 3 * X + np.random.randn(1000, 1)
11 plt.scatter(X, y)
12 plt.show()
13
```



In [106]:

```

1 X_b = np.c_[np.ones((1000, 1)), X] # Adding the bias term which is equal to 1
2
3 # Dividing the data into train and test sets
4 X_train, X_test, y_train, y_test = train_test_split(X_b, y, test_size=0.2, random_state=42)
5
6 theta_optimize = np.linalg.inv(X_train.T.dot(X_train)).dot(X_train.T).dot(y_train)
7
8 #Output
9 #obtained feature weights w0 = 5.04, w1 = 2.94
10 print(theta_optimize)
11
12 # Predicting new data with the obtained feature weights
13 y_pred = X_test.dot(theta_optimize)
14 r2_score(y_test, y_pred)
15

```

```

[[5.08703256]
 [2.93815133]]

```

Out[106]:

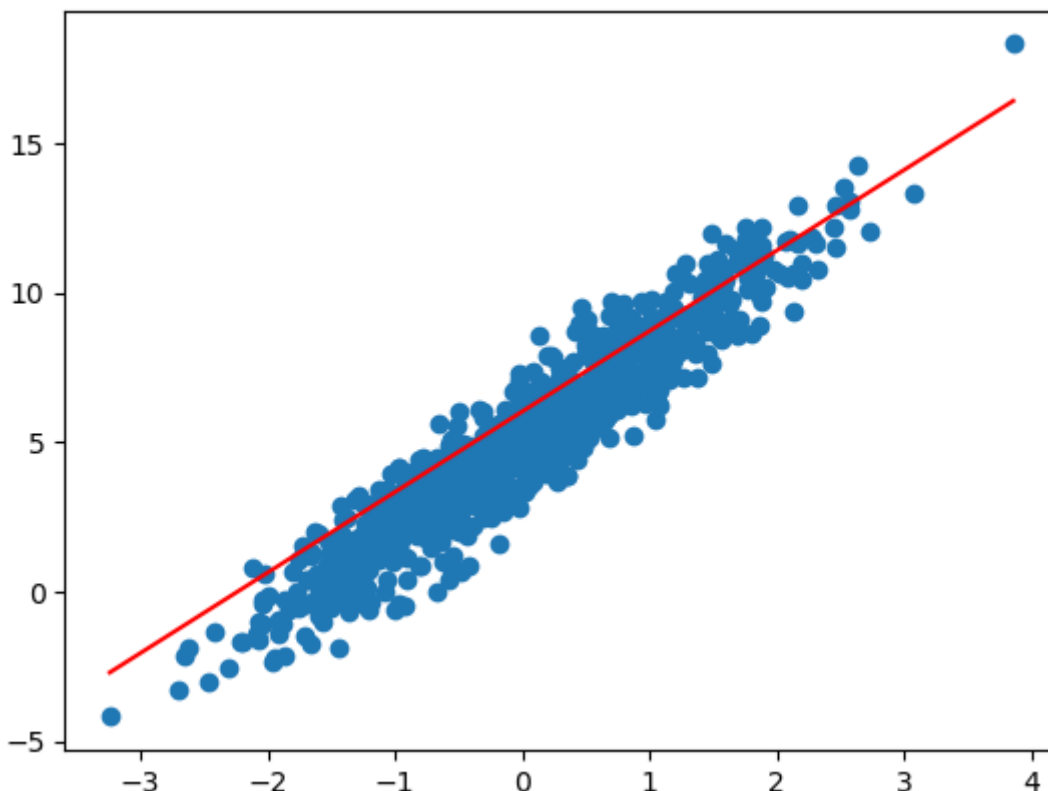
0.8961012486926588

In [107]:

```

1 plt.scatter(X, y)
2 plt.plot([min(X), max(X)], [min(y_pred), max(y_pred)], color='red') # regression line
3 plt.show()

```



You have seen it has predicted the feature weights very close to the actual values ($y = 5 + 3 * X + \text{Gaussian noise}$), but due to the noise in the data it is unable to predict the exact values, but the predictions were close enough.

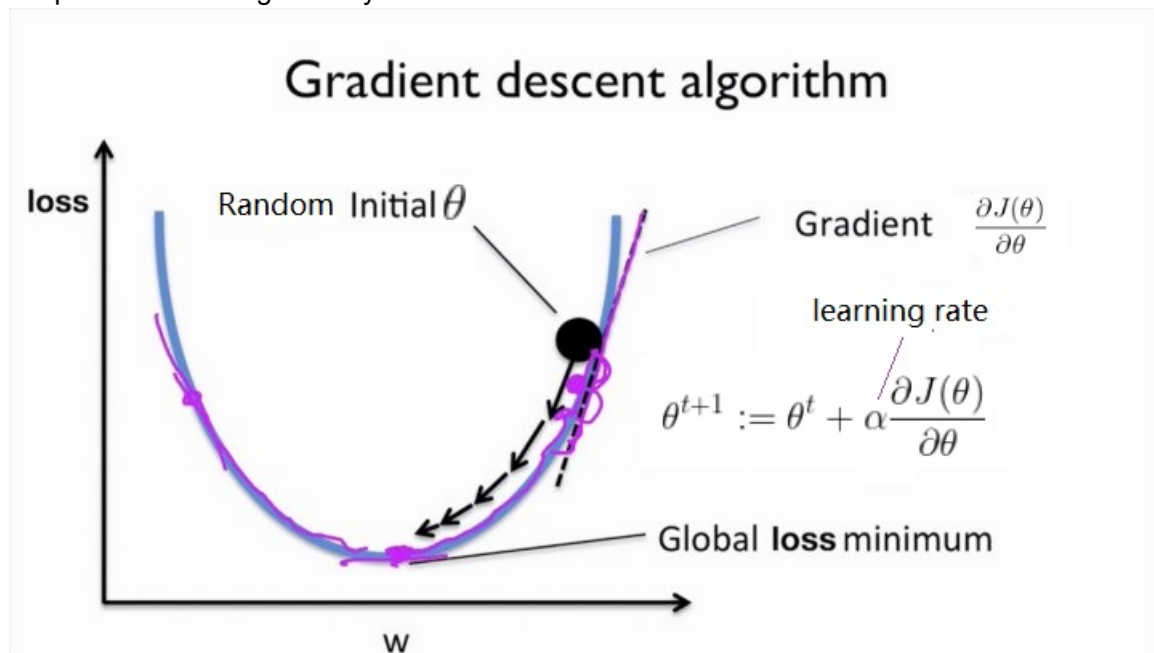
Disadvantages

- It is computationally expensive if you have a large number of features.
- If there are any redundant features in your data, then the matrix inversion in the normal equation is not possible. In that case, the inverse can be replaced by the pseudo inverse.

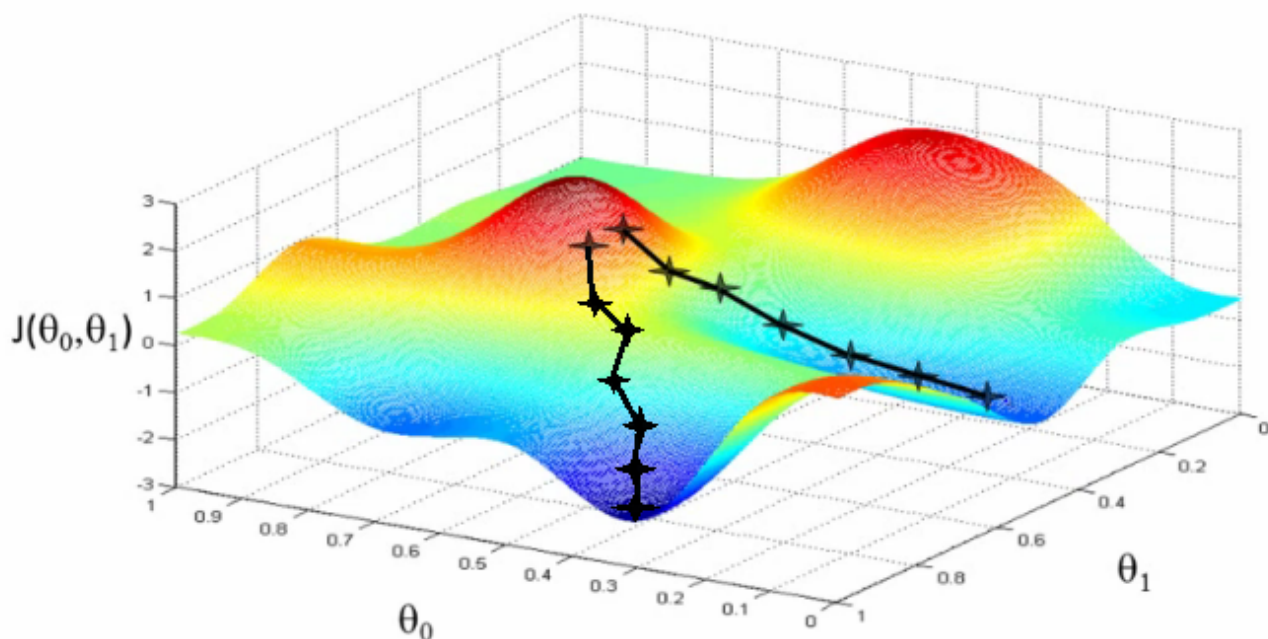
1.5 Gradient Descent

In this section you can learn how the gradient descent algorithm works and implement it from scratch in python.

Gradient Descent minimizes the cost function by iteratively moving in the direction of steepest descent, updating the parameters along the way.



In a real world example, it is similar to find out a best direction to take a step downhill.



We take a step towards the direction to get down. From the each step, you look out the direction again to get down faster and downhill quickly.

To find the best minimum, repeat steps to apply various values for θ . In other words, repeat steps until convergence.

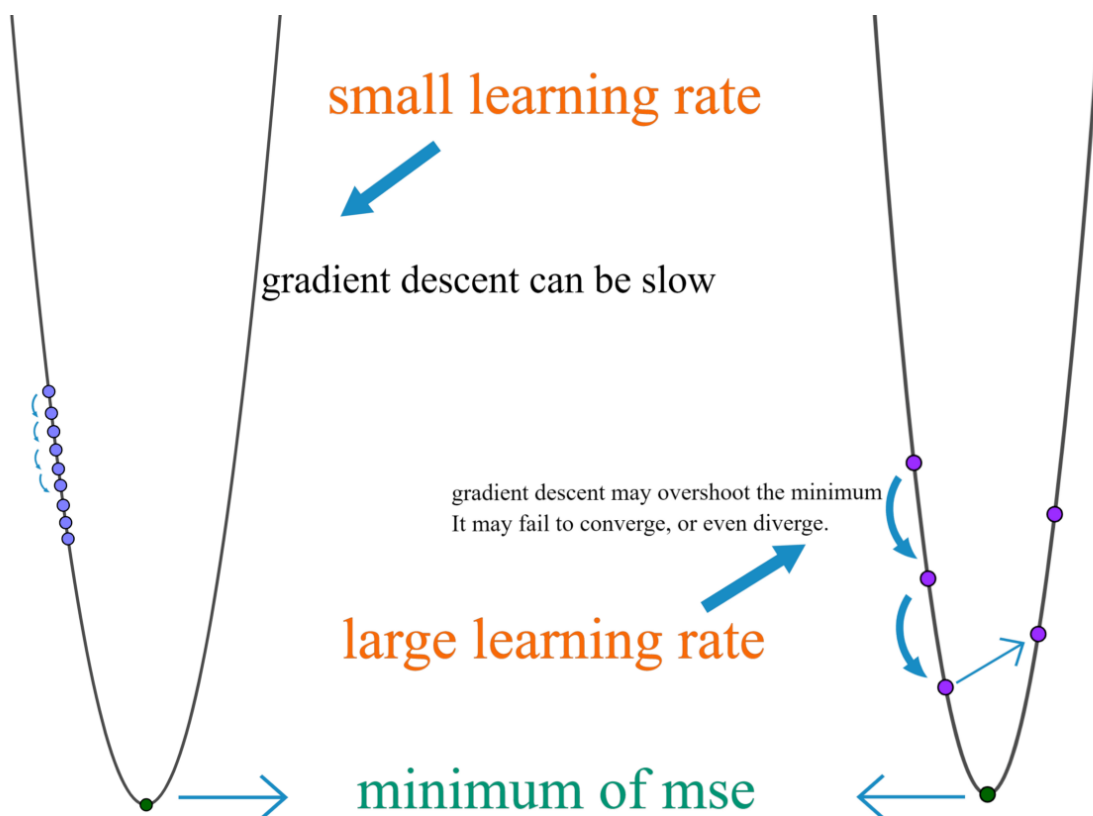
repeat until convergence {

$$\theta^{t+1} := \theta^t - \alpha \frac{1}{m} \frac{\partial J(\theta)}{\partial \theta}$$

}

The choice of correct learning rate is very important as it ensures that Gradient Descent converges in a reasonable time.

- If we choose α to be very small, Gradient Descent will take small steps to reach local minima and will take a longer time to reach minima.
- If we choose **α to be very large**, Gradient Descent can overshoot the minimum. It may fail to converge or even diverge.



Hands-on Coding

implement Linear Regression from scratch

In [108]:

```
1 import numpy as np
2 from sklearn.metrics import r2_score
3
4 np.random.seed(42)
5
6 learning_rate = 0.1
7 iterations = 50
8 m = 100 # total number of samples
9 theta = np.random.randn(2,1) # random initialization
10 for iteration in range(iterations):
11     gradient = 2/m * X_train.T.dot(X_train.dot(theta) - y_train)
12     theta = theta - learning_rate * gradient
13
14 # Output
15 print(theta)
16
17
18 # array([[5.08703256],
19 #         [2.93815133]])
20
21 # Predicting new values with gradient descent
22 y_pred = X_test.dot(theta)
23 r2_score(y_test, y_pred)
```

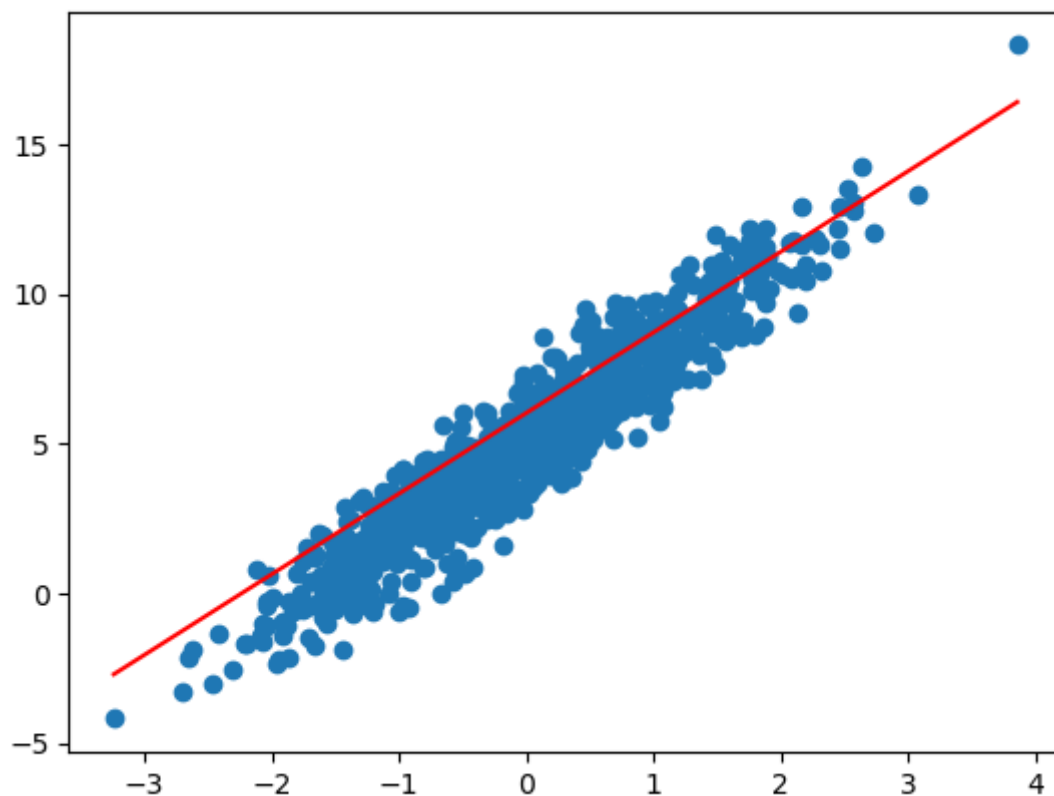
```
[[5.08703256]
 [2.93815133]]
```

Out[108]:

0.8961012486937926

In [109]:

```
1 plt.scatter(X, y)
2 plt.plot([min(X), max(X)], [min(y_pred), max(y_pred)], color='red') # regression line
3 plt.show()
```



Implementing Linear Regression in Scikit-Learn

In [110]:

```
1 from sklearn.linear_model import LinearRegression
2
3 X_train = X_train[:, 1] # The sklearn model will automatically add the bias term so we donot ha
4 X_test = X_test[:, 1]
5
6
7 linear_regression = LinearRegression()
8 linear_regression.fit(X_train.reshape(-1, 1), y_train)
9 print(linear_regression.intercept_)
10 # Output
11 #array([5.08703256])
12
13 print(linear_regression.coef_)
14 # Ouput
15 #array([[2.93815133]])
16
17 # Predicting new values and calculating the r2_score
18 linear_regression.score(X_test.reshape(-1, 1), y_test)
19
20 # output
21 #0.8961012486926588
```

```
[5.08703256]
[[2.93815133]]
```

Out[110]:

0.8961012486926588

As you can see the values we got from the normal equation, gradient descent, sklearn are nearly the same.

Disadvantages

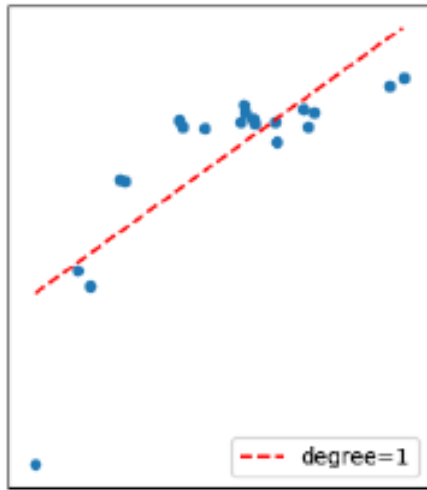
- There are possibilities for the gradient descent to stuck in a local minimum if you use another cost function that is not of a convex shape.
- You should find the appropriate value for the learning rate.

Well if you have read this far and everything makes sense pat yourself on the back!. You have learned all the underlying concepts of linear regression. A

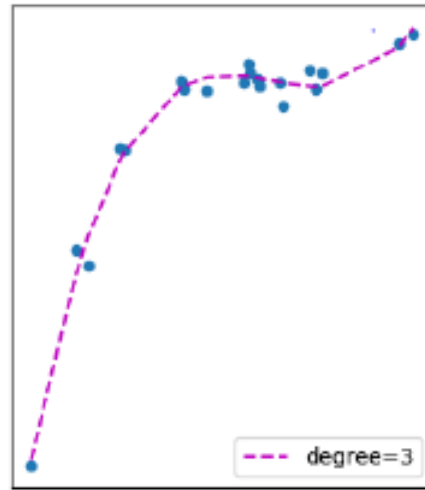
1.6 Polynomial Regression

(https://en.wikipedia.org/wiki/Polynomial_regression)

We use polynomial regression when the relationship between the independent and dependent variables is nonlinear.



Simple Linear
Regression



Polynomial
Regression

This is accomplished by "exponentiating" our variable by taking it to powers greater than 1.

Simple Linear Regression

$$y = \theta_0 + \theta_1 x_1$$

Multiple Linear Regression

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Polynomial Linear
Regression

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \dots + \theta_n x_1^n$$

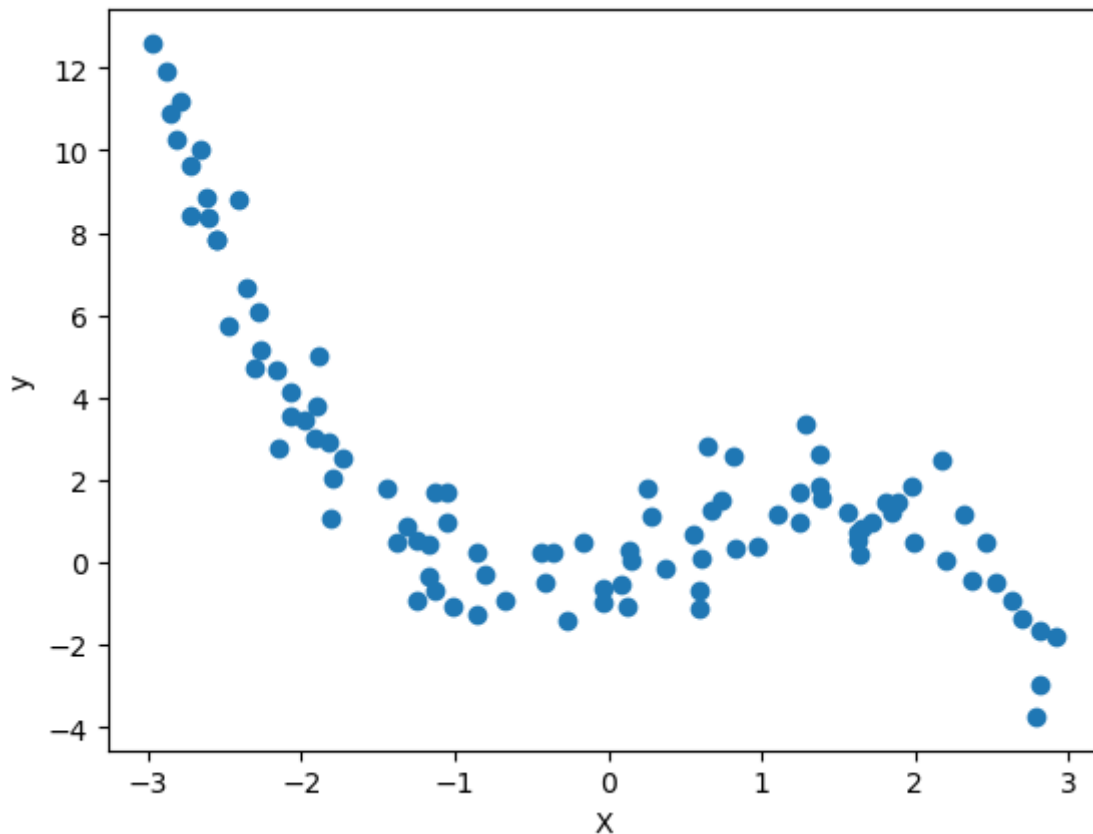
Hands-on Coding

In [111]:

```
1 #step 1: Import the required libraries
2 import numpy as np
3 from scipy import stats
4 import matplotlib.pyplot as plt
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import r2_score, mean_squared_error
8 from sklearn.model_selection import train_test_split
```

In [112]:

```
1 #step2: Load the data set. Here, We're generating data
2 sample_cnt= 100
3 #X = np.linspace(start = -8, stop = 8, num = sample_cnt)
4 X = np.random.uniform(-3, 3, size=100)
5
6 # curve using polynomial
7  $\theta_0, \theta_1, \theta_2, \theta_3 = 0.2, 1, 0.5, -0.4$ 
8  $y = \theta_0 + \theta_1 * X + \theta_2 * (X**2) + \theta_3 * (X**3) + \text{np.random.normal}(0, 1, \text{size} = \text{sample\_cnt})$ 
9
10
11
12 plt.scatter(X, y)
13 plt.xlabel('X')
14 plt.ylabel('y')
15 plt.savefig('regu-0.png', dpi=200)
16 plt.show()
17
```



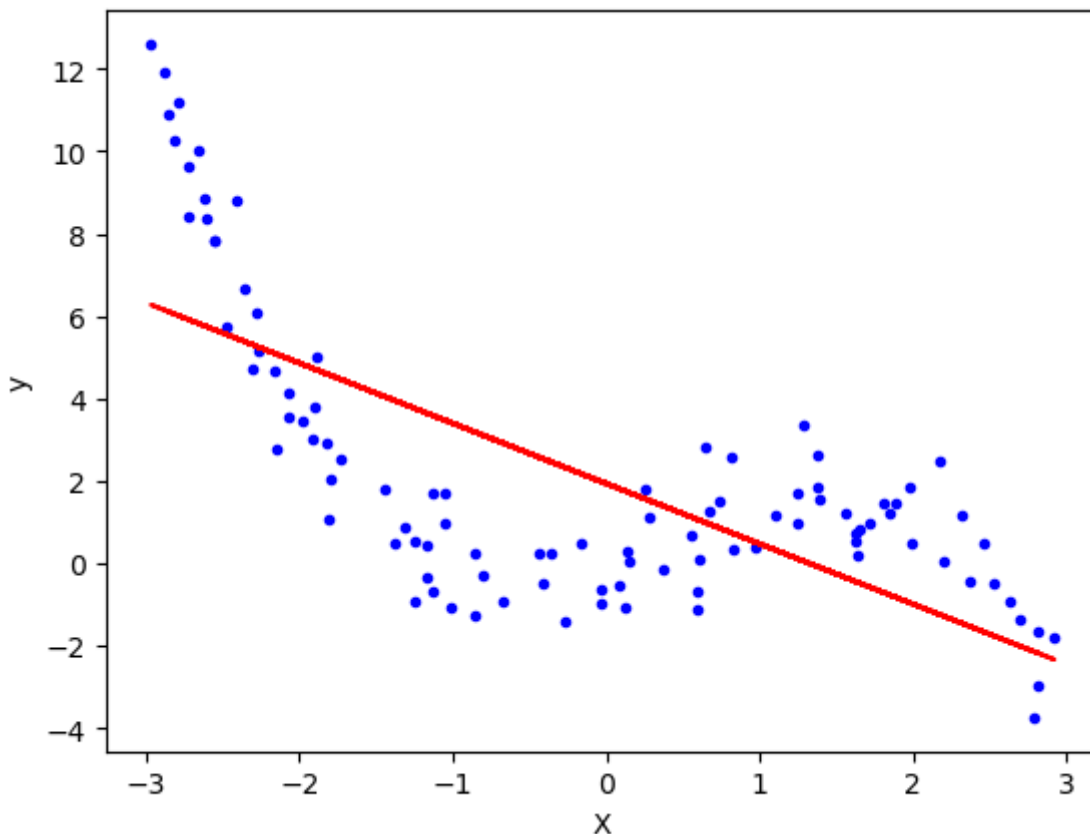
In [113]:

```

1 #split data
2 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8, random_state=1)
3
4 # step 3: using a linear regression model to predict
5 lin_reg = LinearRegression()
6 lin_reg.fit(X_train.reshape(-1,1), y_train.reshape(-1,1))
7
8 y_pred = lin_reg.predict(X_train.reshape(-1,1))
9
10 a = lin_reg.coef_[0][0]
11 b = lin_reg.intercept_[0]
12
13 plt.plot(X, y, 'b.')
14 # plt.plot([min(X_train), max(X_train)], [min(y_pred), max(y_pred)], color='red') # regression
15 plt.plot(X_train, y_pred, 'r-')
16 plt.xlabel('X')
17 plt.ylabel('y')
18 plt.savefig('regu-1.png', dpi=200)
19 print("The equation: y = {}x+{}".format(a,b))
20 #Calculate the error and evaluate the model
21 print(lin_reg.intercept_, lin_reg.coef_)
22 print(r2_score(y_pred, y_train))
23 plt.show()
24

```

The equation: y = -1.4623971689071653x+1.9344273366540403
 [1.93442734] [[-1.46239717]]
 0.0640863835493316



In [114]:

```

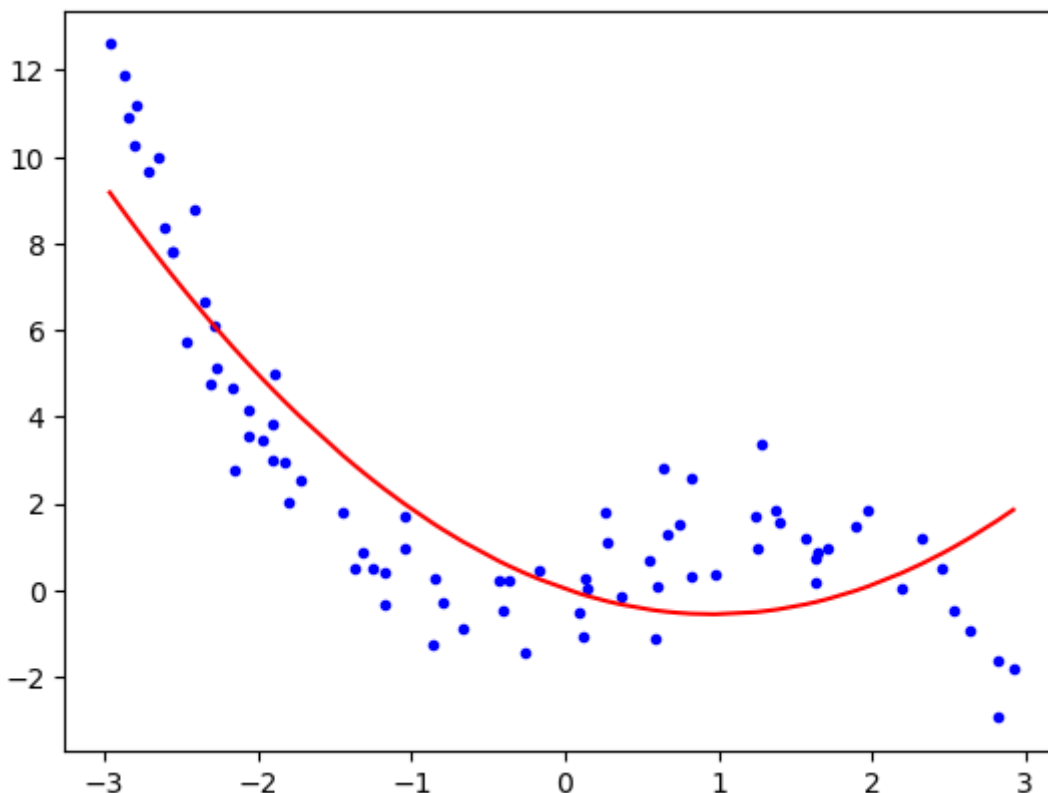
1 # Use a polynomial with degree of 2
2 ploy = PolynomialFeatures(degree=2, include_bias=False)
3 X_train2 = ploy.fit_transform(X_train.reshape(-1,1))
4
5 lin_reg = LinearRegression()
6 lin_reg.fit(X_train2, y_train.reshape(-1,1))
7 print(lin_reg.intercept_, lin_reg.coef_) # [ 2.60996757] [[-0.12759678  0.9144504 ]]
8 a = lin_reg.coef_[0][0]
9 b = lin_reg.coef_[0][1]
10 c = lin_reg.intercept_[0]
11
12
13
14 y_pred = np.dot(X_train2, lin_reg.coef_.T) + lin_reg.intercept_
15
16 plt.plot(X_train, y_train, 'b.')
17 plt.plot(np.sort(X_train), y_pred[np.argsort(X_train)], color='r')
18
19 plt.savefig('regu-2.png', dpi=200)
20 plt.show()
21 print("The equation: y = {:.3f}x^2+{:.3f}x+{:.3f}".format(a,b,c))
22 print(r2_score(y_pred, y_train))
23 print(r2_score(np.dot(ploy.transform(X_test.reshape(-1,1)), lin_reg.coef_.T) + lin_reg.intercept_

```

```

[0.0253991] [[-1.21577111  0.62996561]]

```



The equation: $y = -1.216x^2 + 0.630x + 0.025$
 0.6395847705797559
 0.1295465340821086

As the data extending to polynomial features, the value would be extremely large or small because of the power operation. That will influence the use of gradient descent which runs in background when we call `fit()`. So a normalization or standardization is necessary. See `StandardScaler` in preprocessing.

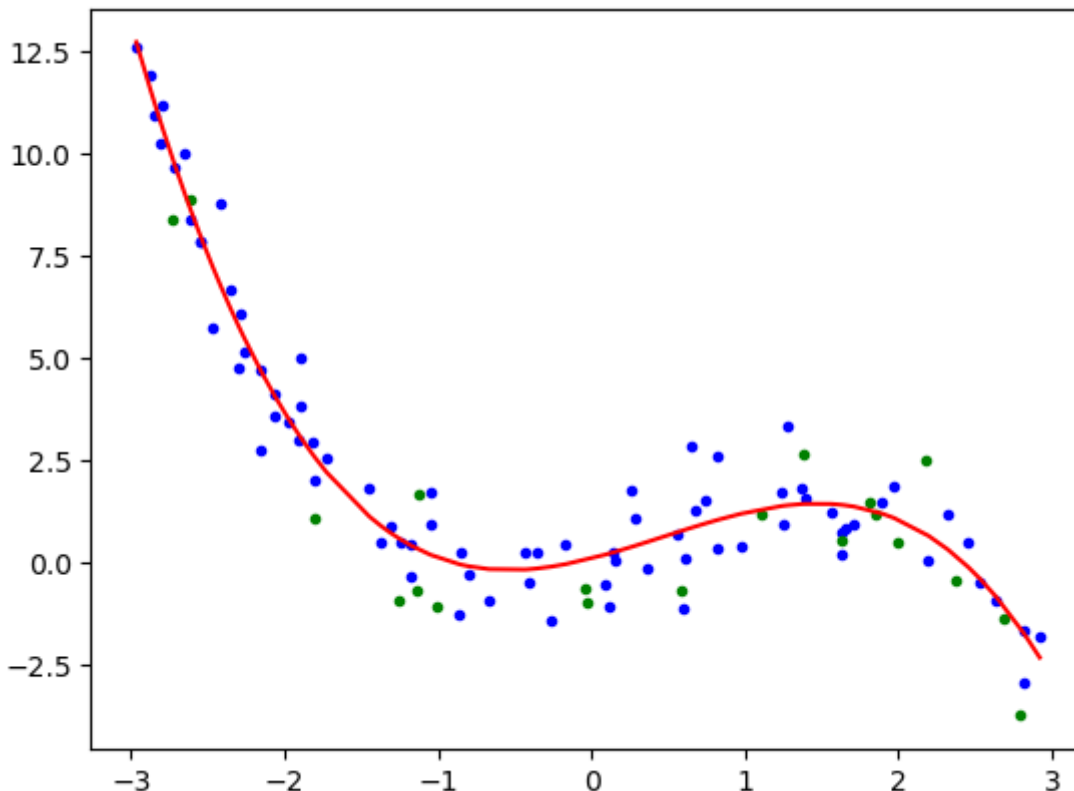
`Pipeline` can help us assemble several preprocessing functions and the learning process together.

In [115]:

```

1 # Use a polynomial with degree of 3
2 from sklearn.pipeline import Pipeline
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.preprocessing import PolynomialFeatures
5
6 poly_reg = Pipeline([
7     ('poly', PolynomialFeatures(degree=3)),
8     ('std_scaler', StandardScaler()),
9     ('lin_reg', LinearRegression())
10 ])
11
12 poly_reg.fit(X_train.reshape(-1,1), y_train.reshape(-1,1))
13
14
15 y_pred = poly_reg.predict(X_train.reshape(-1,1))
16
17
18 plt.plot(X_train, y_train, 'b.')
19 plt.plot(X_test, y_test, 'g.')
20 plt.plot(np.sort(X_train), y_pred[np.argsort(X_train)], color='r')
21
22 plt.savefig('regu-3.png', dpi=200)
23 plt.show()
24 params = poly_reg.named_steps['lin_reg']
25 print("The equation:", poly_reg.named_steps['lin_reg'].coef_, poly_reg.named_steps['lin_reg'].int
26 print("The equation: y = {:.3f}x^3+{:.3f}x^2+{:.3f}x+{:.3f}".format(params.coef_[0][0], params.c
27 print(r2_score(y_pred, y_train))
28 print(r2_score(poly_reg.predict(X_test.reshape(-1,1)), y_test))
29 #print(mean_squared_error(y_pred, y))

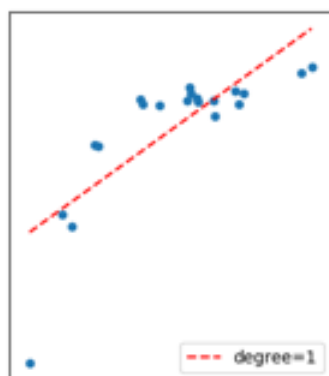
```



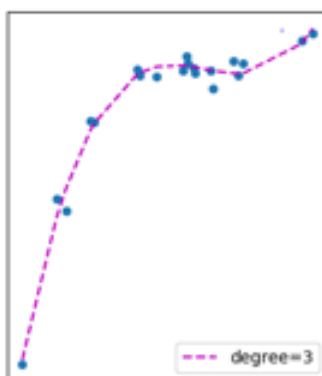
The equation: $[[0. \quad 1.65603834 \quad 1.52914098 \quad -4.15373307]] \quad [2.49147196]$
The equation: $y = 0.000x^3 + 1.656x^2 + 1.529x - 4.154$

0. 9413635998834047
0. 8303023309883463

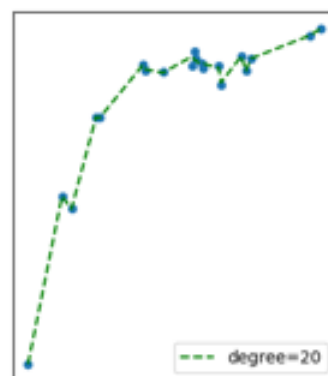
In practice, we need to select the "degree" of the polynomial.



... x

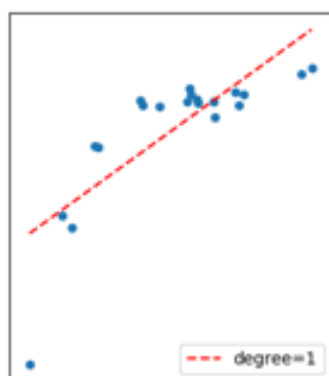


... x^3

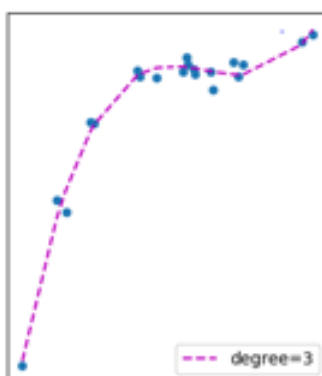


... x^{20}

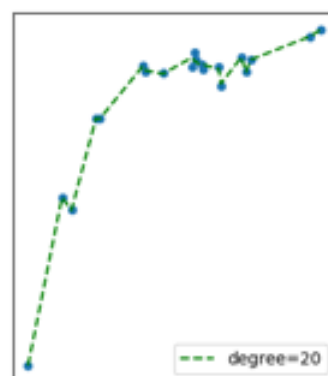
Our selection should seek to fit the current data well, and generalize to new data.



Under fit



Good fit



Over fit

This challenge is known as the "bias variance tradeoff".

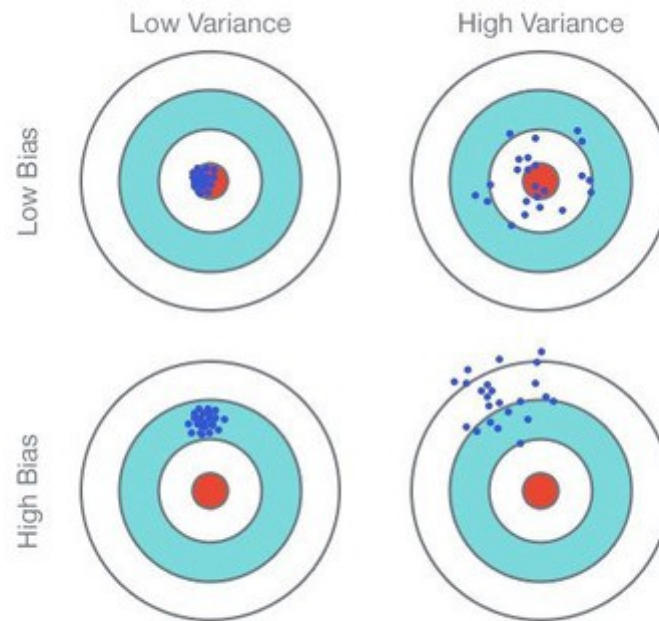
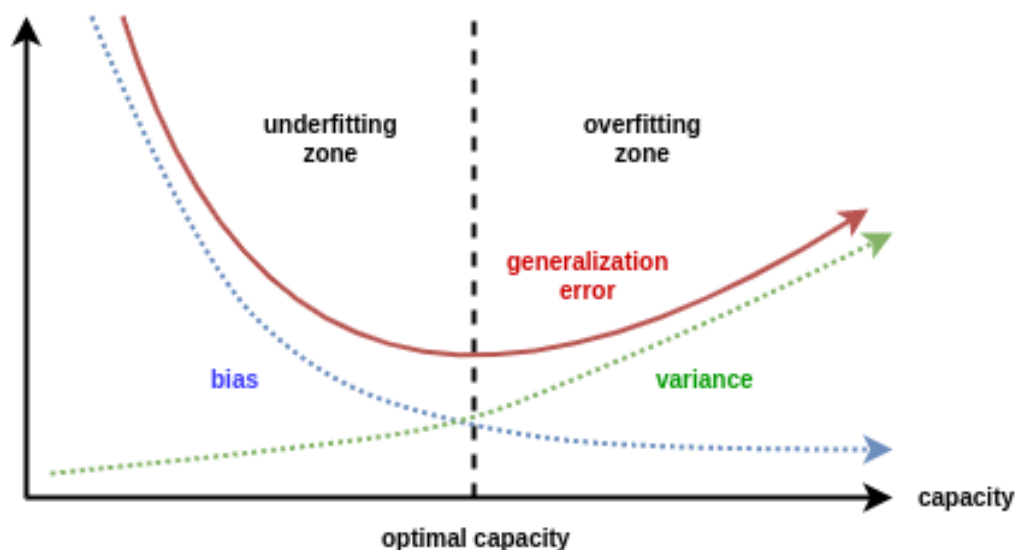


Fig. 1: Graphical Illustration of bias-variance trade-off, Source: Scott Fortmann-Roe., Understanding Bias-Variance Trade-off

Bias usually caused by underfitting, Variance caused by overfitting. What can we do to solve this problem?

- 1. Choose a better polynomial degree. As we increase the degree, the bias decreases, but the variance increases. We want to stop where these two factors are minimized.



Note: polynomial regression is very sensitive to outliers, and we must take care when selecting the degree to avoid overfitting.

- 2. Regularization .Regularization in scikit learn is [RidgeRegression](https://en.wikipedia.org/wiki/Ridge_regression#:~:text=Ridge%20regression%20is%20a%20method%20,which%20is%20in%20the%20scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html) (https://en.wikipedia.org/wiki/Ridge_regression#:~:text=Ridge%20regression%20is%20a%20method%20,which%20is%20in%20the%20scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html), which is in [linear_model](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html) (https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html). Use it if you need regularization in your model.

Ridge regression is also known as L2 regularization. It is a regularized version of linear regression to find a better fitting line. It adds L2 penalty terms in the cost function and thereby reducing coefficients lower towards zero and minimizing their impact on the training data. It is useful to avoid over-fitting of

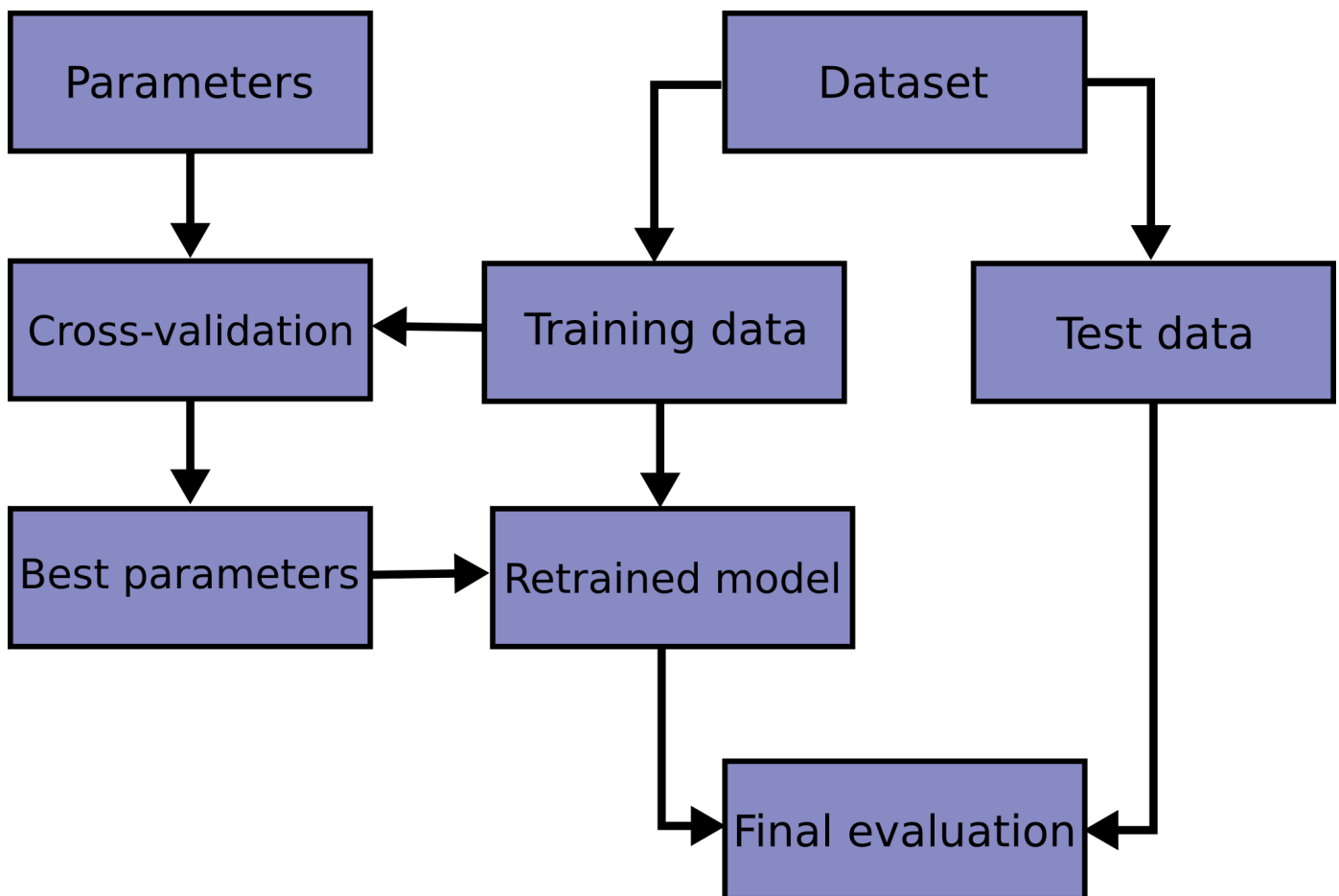
the data in a model. The alpha parameter of a model is used to calculate cost function to manage bias-variance tradeoff and reduce an error in prediction.

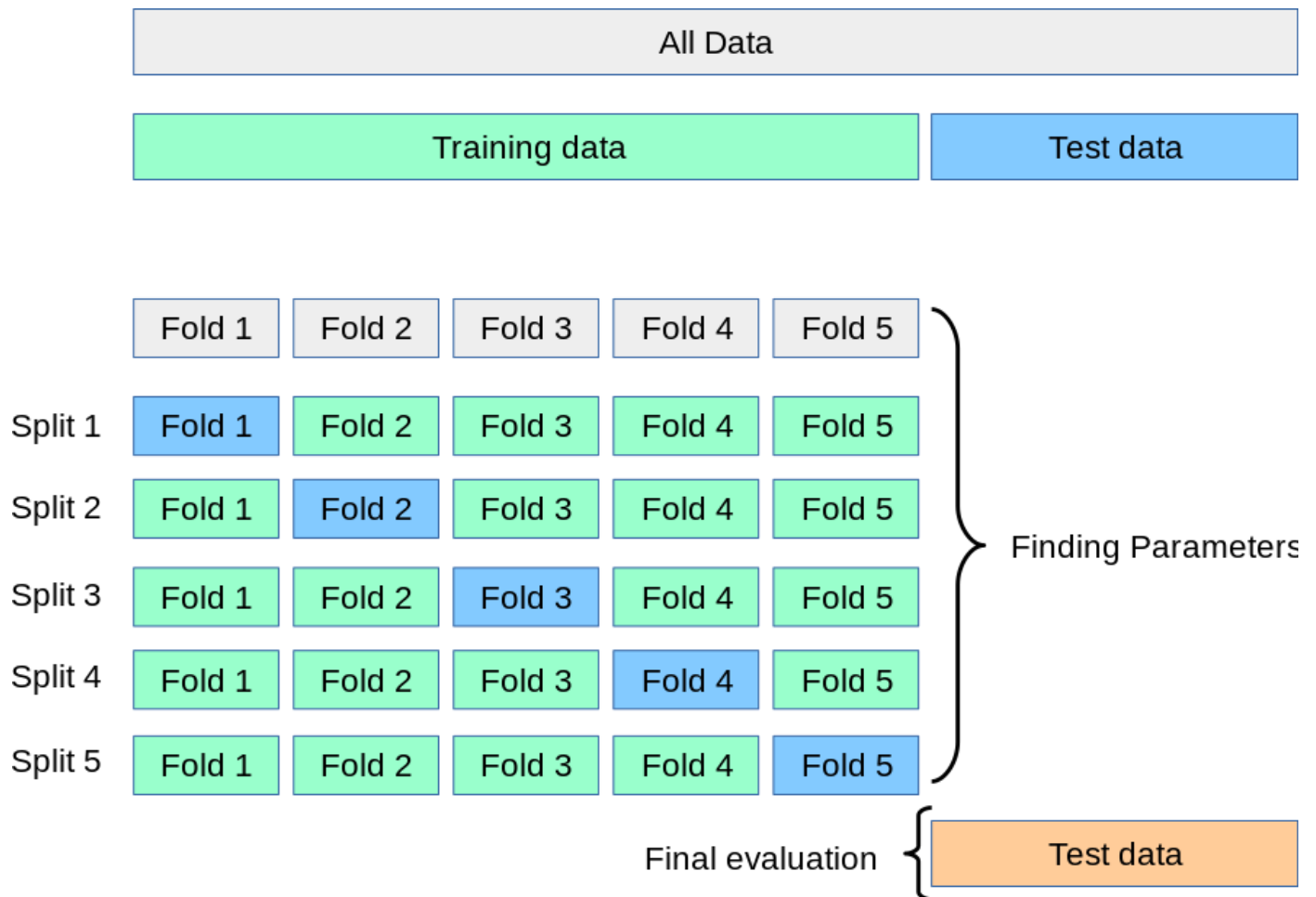
$$J(\theta) = \frac{1}{2m} \sum_{j=1}^m (h_{\theta}(x^j) - y^j)^2 + \alpha \sum_{i=1}^n \theta_i$$
$$= \frac{1}{2m} (X\theta - y)^T (X\theta - y) + \alpha \theta^T \theta$$

- 3. Use training set, validation set and test set or [cross validation \(https://scikit-learn.org/stable/modules/cross_validation.html\)](https://scikit-learn.org/stable/modules/cross_validation.html) to acquire a better model.

2 [Cross-validation \(https://scikit-learn.org/stable/modules/cross_validation.html\)](https://scikit-learn.org/stable/modules/cross_validation.html)

We often use cross validation to select parameters or to select models.





For more detail:

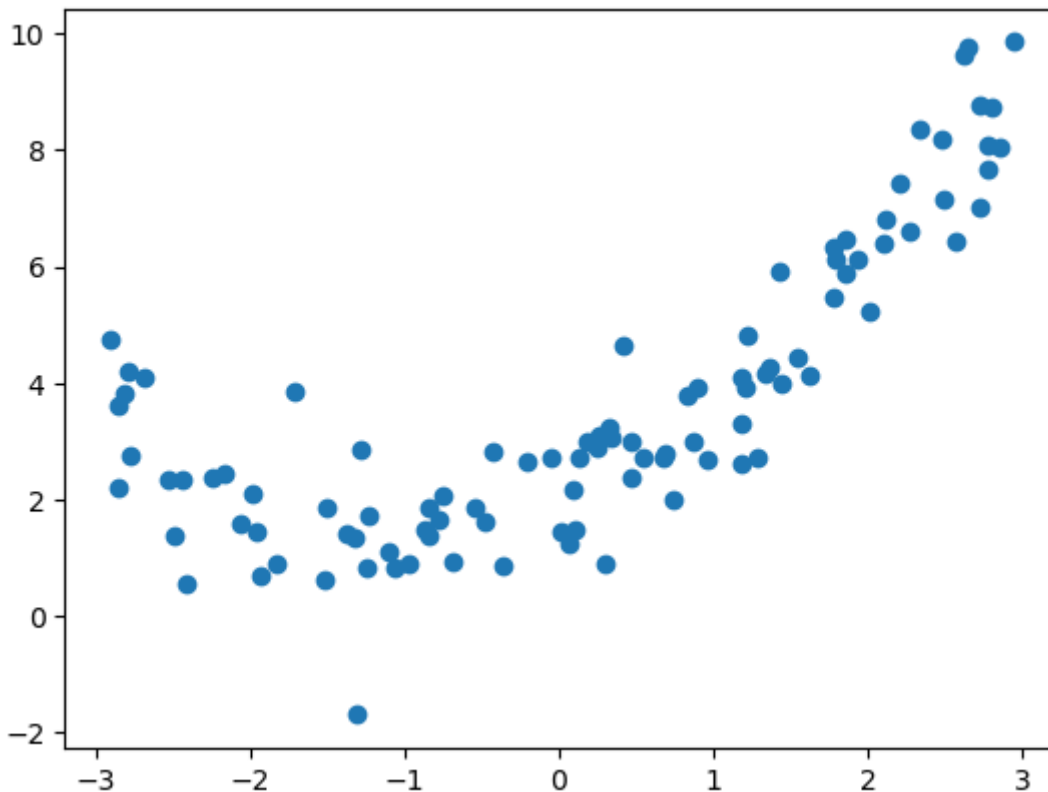
- https://scikit-learn.org/stable/modules/cross_validation.html (https://scikit-learn.org/stable/modules/cross_validation.html)
- [https://en.wikipedia.org/wiki/Cross-validation_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics)) ([https://en.wikipedia.org/wiki/Cross-validation_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics)))
- https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)

In [117]:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import Lasso, Ridge
6 from sklearn.model_selection import GridSearchCV
7 from sklearn.pipeline import Pipeline
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.preprocessing import PolynomialFeatures
10 from sklearn.linear_model import LinearRegression
11 sample_cnt= 100
12 X = np.random.uniform(-3,3,size=100).reshape(-1, 1)
13
14
15 # curve using polynomial
16  $\theta_0, \theta_1, \theta_2, \theta_3 = 2, 1, 0.5, 0$ 
17  $y = \theta_0 + \theta_1 * X + \theta_2 * (X**2) + \theta_3 * (X**3)$ 
18 y += np.random.normal(0,1,size = sample_cnt).reshape(-1, 1)
19 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1)
20
21 plt.scatter(X, y)
22 plt.show()

```



In [118]:

```
1
2
3 model = Pipeline([
4     ('poly', PolynomialFeatures()),
5     ('std_scaler', StandardScaler()),
6     ('ridge', Ridge())
7 ])
8 params = { 'poly__degree': [1, 2, 3, 4],
9 'ridge__alpha': list(x / 10 for x in range(0, 101)),
10            'ridge__max_iter': [5, 10, 100, 1000]}
11
12
13 grid_search = GridSearchCV(model, param_grid=params, cv=5)
14 grid_search.fit(X_train, y_train)
15 print('The highest score of cross validation: {:.3f}\n'.format(grid_search.best_score_))
16 print('The best parameters: {}\n'.format( grid_search.best_params_))
17 print('The highest score: {:.3f}\n'.format(grid_search.score(X_test, y_test)))
18
19 # y_hat = grid_search.predict(np.array(X_test))
20 # t = np.arange(len(X_test))
21 # plt.plot(t, y_test, 'r-', linewidth=2, label='Test')
22 # plt.plot(t, y_hat, 'g-', linewidth=2, label='Predict')
23 # plt.legend(loc='upper right')
24 # plt.grid()
25 # plt.show()
```

The highest score of cross validation:0.837

The best parameters: {'poly__degree': 2, 'ridge__alpha': 2.5, 'ridge__max_iter': 5}

The highest score:0.897

3 Evaluation for Regression Models

Metrics commonly used to evaluate regression models are:

- Mean Absolute Error
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- R-Squared

1) Mean Absolute Error

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

```
def MAE(y, y_pre):
    return np.mean(np.abs(y - y_pre))
```

Advantages

- MAE is not sensitive to outliers. Use MAE when you do not want outliers to play a big role in error calculated.

Disadvantages

- MAE is not differentiable globally. This is not convenient when we use it as a loss function, due to the gradient optimization method.

2) Mean Squared Error (MSE)

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

```
def MSE(y, y_pre):
    return np.mean((y - y_pre) ** 2)
```

Advantages

- Graph of MSE is differentiable which means it can be easily used as a loss function.
- MSE can be decomposed into variance and bias squared. This helps us understand the effect of variance or bias in data to the overall error.

$$MSE(\hat{y}) = Var(\hat{y}) + (Bias(\hat{y}))^2$$

** Disadvantages **

- The value calculated MSE has a different unit than the target variable since it is squared. (Ex. meter → meter²)
- If there exists outliers in the data, then they are going to result in a larger error. Therefore, MSE is not robust to outliers (this can also be an advantage if you are looking to penalize outliers).

3) Root Mean Squared Error (RMSE)

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

```
def RMSE(y, y_pre):
    return np.sqrt(np.mean((y - y_pre) ** 2))
```

Advantages

- The error calculated has the same unit as the target variables making the interpretation relatively easier.

Disadvantages

- Just like MSE, RMSE is also susceptible to outliers.

4) R-Squared

$$R^2 = 1 - \frac{\sum_i (\hat{y}^{(i)} - y^{(i)})^2}{\sum_i (\bar{y} - y^{(i)})^2}$$

$$R^2 = 1 - \frac{(\sum_i (\hat{y}_i - y_i)^2) / m}{(\sum_i (\hat{y}_i - y_i)^2) / m}$$

$$= 1 - \frac{MSE(y, y)}{Var(y)}$$

```
def R2(y, y_pre):
    u = np.sum((y - y_pre) ** 2)
    v = np.sum((y - np.mean(y)) ** 2)
    return 1 - (u / v)
```

Value closer to 1 is better.

4 LAB Assignment

Now it's time to implement linear regression techniques in practice. In this lab, you will use linear regression to fit a house price model. You will use some real-world data as the test set to evaluate your model.

4.1 Before Assignment

4.1.1 Load dataset & Import the required libraries

Datasets: scikit-learn provides a number of datasets which can be directly loaded by using a function. First we load a dataset as an example.

In [119]:

```
1 import warnings
2 from sklearn import datasets
3 boston = datasets.load_boston()
4 print(boston.DESCR)
```

.. _boston_dataset:

Boston house prices dataset

****Data Set Characteristics:****

:Number of Instances: 506

:Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.

:Attribute Information (in order):

- CRIM per capita crime rate by town
- ZN proportion of residential land zoned for lots over 25,000 sq.ft.
- INDUS proportion of non-retail business acres per town
- CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
- NOX nitric oxides concentration (parts per 10 million)
- RM average number of rooms per dwelling
- AGE proportion of owner-occupied units built prior to 1940
- DIS weighted distances to five Boston employment centres
- RAD index of accessibility to radial highways
- TAX full-value property-tax rate per \$10,000
- PTRATIO pupil-teacher ratio by town
- B $1000(B_k - 0.63)^2$ where B_k is the proportion of black people by town
- LSTAT % lower status of the population
- MEDV Median value of owner-occupied homes in \$1000's

:Missing Attribute Values: None

:Creator: Harrison, D. and Rubinfeld, D.L.

This is a copy of UCI ML housing dataset.

<https://archive.ics.uci.edu/ml/machine-learning-databases/housing/> (<https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>)

This dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.

The Boston house-price data of Harrison, D. and Rubinfeld, D.L. 'Hedonic prices and the demand for clean air', J. Environ. Economics & Management, vol. 5, 81-102, 1978. Used in Belsley, Kuh & Welsch, 'Regression diagnostics ...', Wiley, 1980. N.B. Various transformations are used in the table on pages 244-261 of the latter.

The Boston house-price data has been used in many machine learning papers that address regression problems.

.. topic:: References

- Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and Sources of Collinearity', Wiley, 1980. 244-261.
- Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Conference of Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.

d:\miniconda3\envs\cs405\lib\site-packages\sklearn\utils\deprecation.py:87: FutureWarning: Function load_boston is deprecated; `load_boston` is deprecated in 1.0 and will be removed in 1.2.

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::

```
import pandas as pd
import numpy as np

data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)
data = np.hstack([raw_df.values[::2, :], raw_df.values[1::2, :2]])
target = raw_df.values[1::2, 2]
```

Alternative datasets include the California housing dataset (i.e. :func:`~sklearn.datasets.fetch_california_housing`) and the Ames housing dataset. You can load the datasets as follows::

```
from sklearn.datasets import fetch_california_housing
housing = fetch_california_housing()
```

for the California housing dataset and::

```
from sklearn.datasets import fetch_openml
housing = fetch_openml(name="house_prices", as_frame=True)
```

for the Ames housing dataset.
warnings.warn(msg, category=FutureWarning)

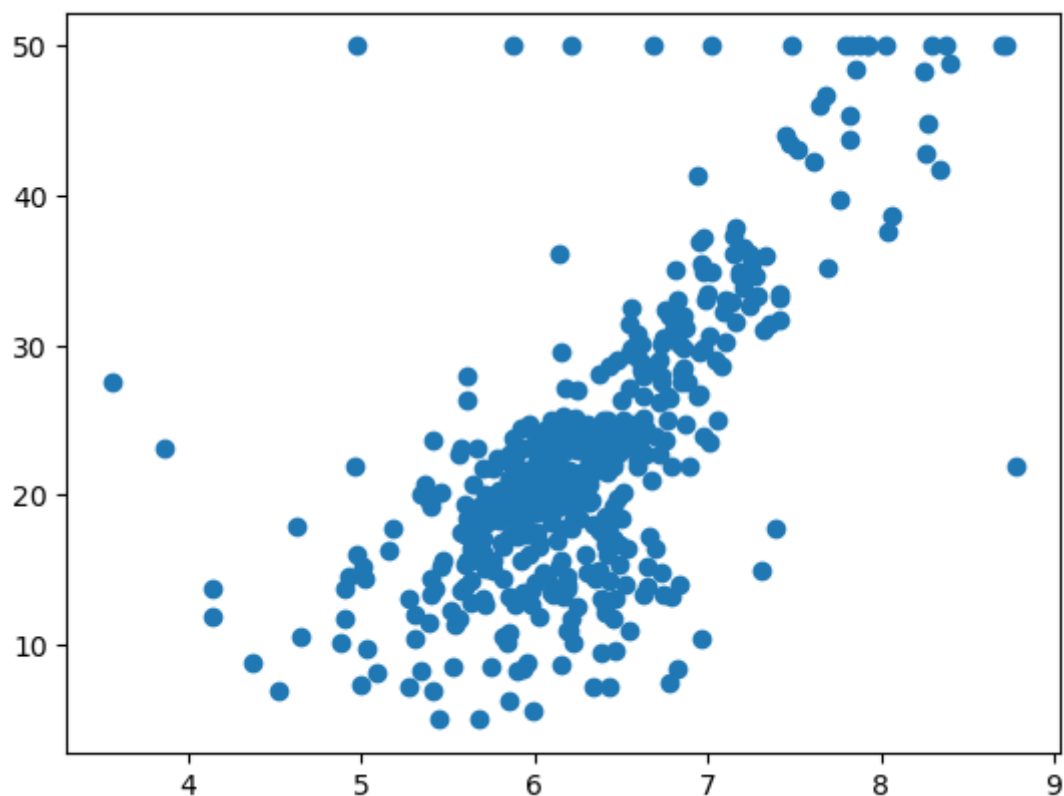
See [sklearn website \(https://scikit-learn.org/stable/modules/classes.html#module-sklearn.datasets\)](https://scikit-learn.org/stable/modules/classes.html#module-sklearn.datasets) for details. To do this you have to import right packages and modules.

4.1.2 Preprocessing data

This is a small dataset containing 506 samples and 13 attributes. We need to use proper visualization methods to have an intuitive understanding. We choose the sixth attribute and draw a scattering plot to see the distribution of samples. We use *matplotlib* for data visualization.

In [120]:

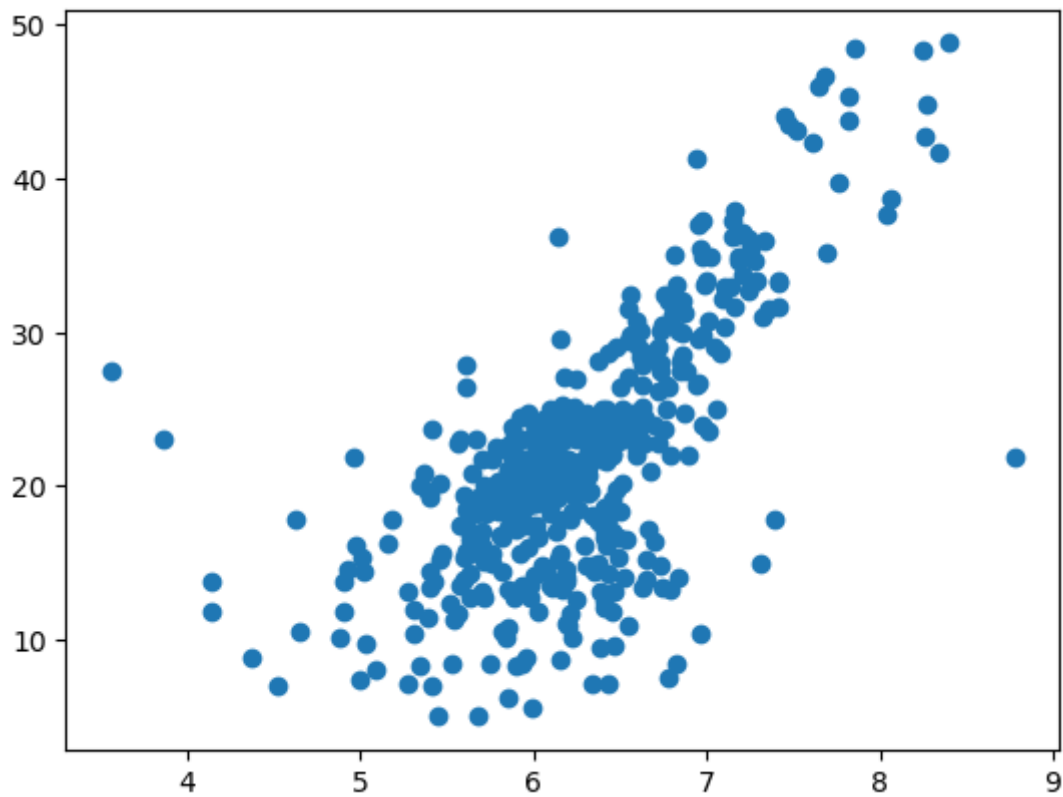
```
1 # Use one feature for visualization
2 x = boston.data[:,5]
3
4 # Get the target vector
5 y = boston.target
6
7 # Scattering plot of price vs. room number
8 from matplotlib import pyplot as plt
9 plt.scatter(x, y)
10 plt.show()
```



It can be seen that the samples have some exceptional distributions at the top of the plot. They may be outliers owing to some practical operation during the data input (e.g., convert any price larger than 50 into 50). However, these data are harmful to the model training, and should be removed.

In [121]:

```
1 x = x[y<50.0]
2 y = y[y<50.0]
3
4 plt.scatter(x, y)
5 plt.show()
```



Now it can be seen that the data is nearly linear, although just in one dimension. Now we use X to denote all attributes

In [122]:

```
1 X = boston.data
2 y = boston.target
3
4 X = X[y<50.0]
5 y = y[y<50.0]
6
7 X.shape
```

Out[122]:

(490, 13)

4.1.3 Split data

Now we divide the whole dataset into a training set and a test set using the the scikit-learn `model_selection` module.

In [123]:

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y)
3 y_train.shape
```

Out[123]:

(367,)

Usually we also use a validation set. When we use the test set for evaluation, the model will not be changed after the evaluation. However, sometime we want to optimize our model by changing its parameters according to prediction results. The solution is to split a validation set from the training set for adjusting our model. When we believe that the model is good enough, then we evaluate our model on the test set. A more rigorous and costly way is cross validation. With that method, the training set is divided into several pieces in the same size and take every piece as a validation set in turn.

4.1.4 Training

1) Linear Regression

Now we try to implement a simple linear regression model because the dataset seems linear.

In [124]:

```
1 from sklearn.linear_model import LinearRegression
2 lin_reg = LinearRegression()
3 lin_reg.fit(X_train, y_train)
```

Out[124]:

LinearRegression()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

The model has been trained just by using a few lines of codes. Now let's make a prediction for testing

In [125]:

```
1 # Make a prediction
2 y_0_hat = lin_reg.predict(X_test[0].reshape(1, -1))
3 y_0_hat
```

Out[125]:

```
array([17.99145573])
```

In [126]:

```
1 y_test[0]
```

Out[126]:

```
14.1
```

Notice that in scikit-learn, the standard interface for machine learning is

1) instantiate a learner with super parameters or none; 2) use `fit()` method and feed the learner with training data; 3) use `predict()` for prediction.

Moreover, the data preprocessing algorithms also have the same interface, they just use `transform()` instead of `predict()`.

Below are the trained parameters.

In [127]:

```
1 lin_reg.coef_
```

Out[127]:

```
array([-1.24143014e-01,  3.44659702e-02, -1.08808701e-01,  1.81305890e-01,
        -1.07954247e+01,  3.58033630e+00, -1.70292872e-02, -1.19026724e+00,
         3.19775631e-01, -1.60081950e-02, -7.83309917e-01,  7.67810088e-03,
        -3.39052027e-01])
```

In [128]:

```
1 lin_reg.intercept_
```

Out[128]:

```
32.00270292551771
```

Use the evaluation method to see if it is a good model. The `score()` method uses R-square.

In [129]:

```
1 lin_reg.score(X_test, y_test)
```

Out[129]:

```
0.7595299927902152
```

2) Polynomial Regression

If you have understood the concept of linear regression, you can easily implement polynomial regression.

4.1.5 Evaluation model

Checking the results on test set.

4.2 LAB Assignment

Please use the real world dataset, **California housing price**, for model training and evaluate the model's prediction performance. You can use simple linear regression, polynomial regression or more complicated base functions such as Gaussian function or use regularization methods. Make sure at least **20% data for testing** and choose one evaluation method you think good. **Please do not just train your model and say that is good enough, you need to analyze the bias and variance.** For that end, validation or cross validation is needed. Compare the score in the training set and the validation set. If they are both good enough, then use the model on the test set.

Your test set can only be used for final evaluation!

In [130]:

```
1 ##### Write Your Code Here #####
2
3 #####
```

4.3 Questions

1) Describe another real-world application where the regression method can be applied 2) What are the strengths of the linear/polynomial regression methods; when do they perform well? 3) What are the weaknesses of the linear/polynomial regression methods; when do they perform poorly? 4) What makes the linear regression method a good candidate for the regression problem, if you have enough knowledge about the data?

Please complete lab4 Assignment , and submit the result to bb as required (Lab02_Assignment_Template.ipynb)

Well done! 🙌 You have made it.