



# Linear Regression

# Regression

- Classification model output value is **discrete**.
- Regression model output value is **continuous**.

gender

house price

crop yield

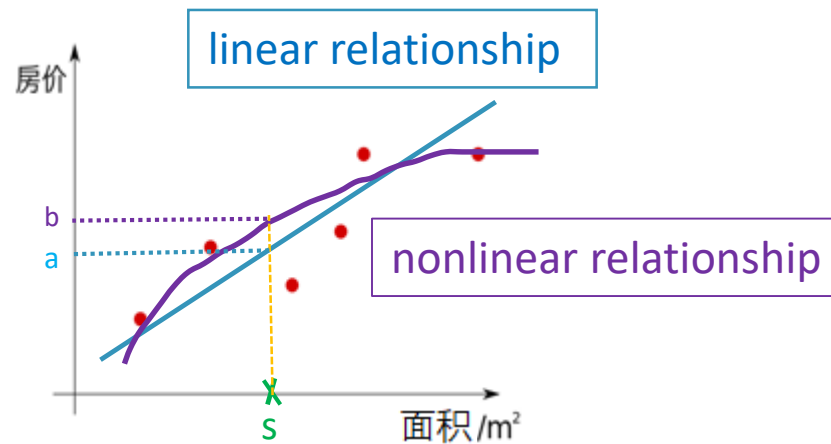
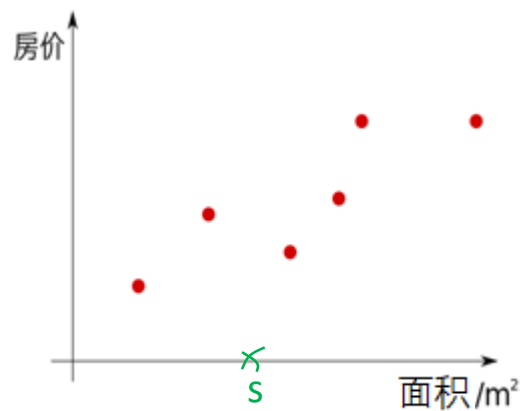
height

movie type

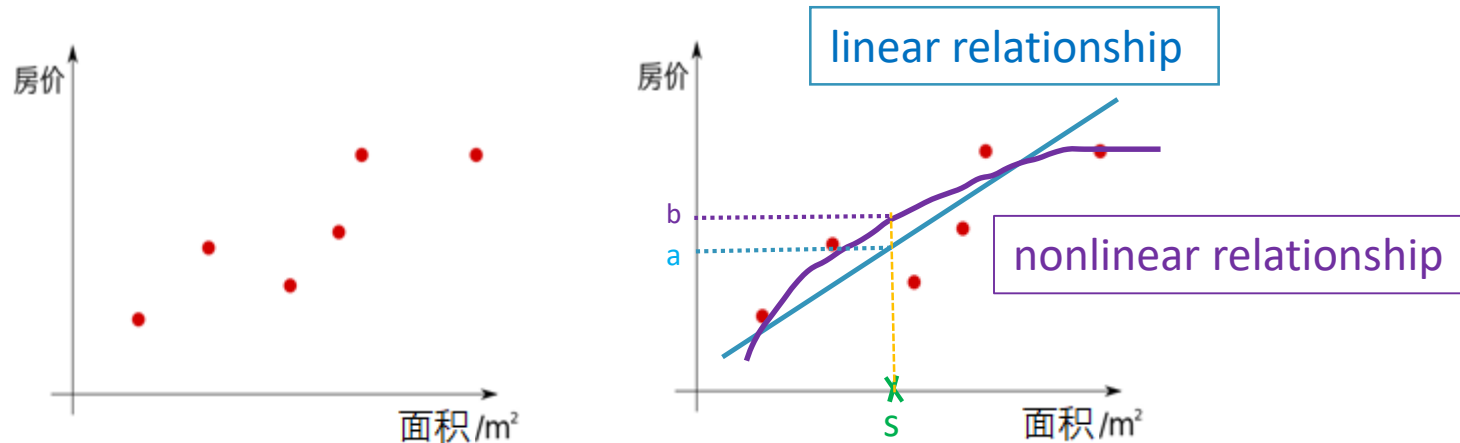
watermelon:  
good or bad

regression	classification
house price	gender: male or female
crop yield	movie type
height	good or bad watermelon
...	...

# Regression



# Regression



- Regression **estimates the relationship** between input values and output value, and establishes a mathematical model in order to accurately **predict** the output value of **new sample**.
- Regression is a **supervised** learning question.

# Regression

## Regression

the number of input values

single regression(单元回归)

multiple regression(多元回归)

## Regression

the relationship of input values  
and output value

linear regression

nonlinear regression

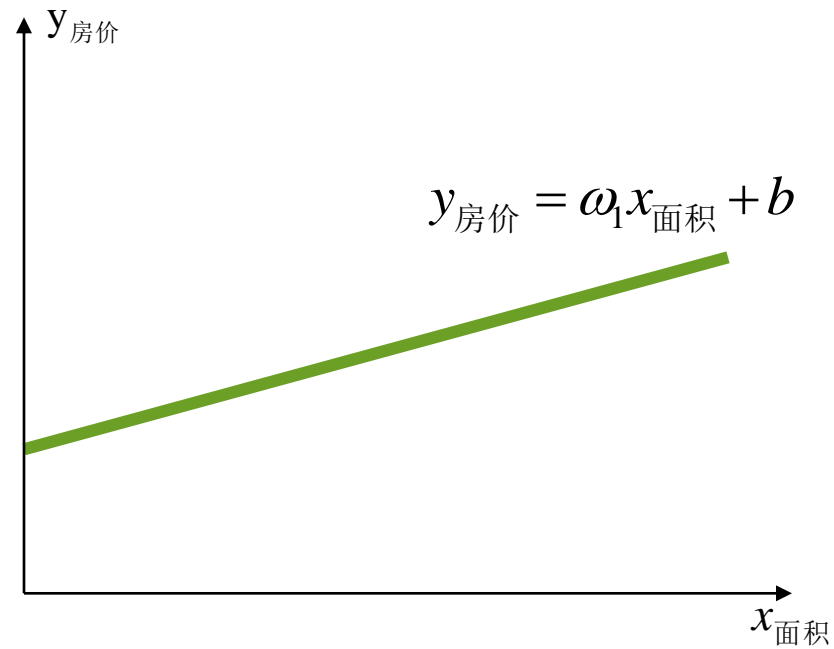
# Outline

- linear regression model
  - linear model
  - least square method --> basic linear regression
  - ridge regression(岭回归)
  - lasso regression

# Outline

- linear regression model
  - linear model
  - least square method --> basic linear regression
  - ridge regression(岭回归)
  - lasso regression

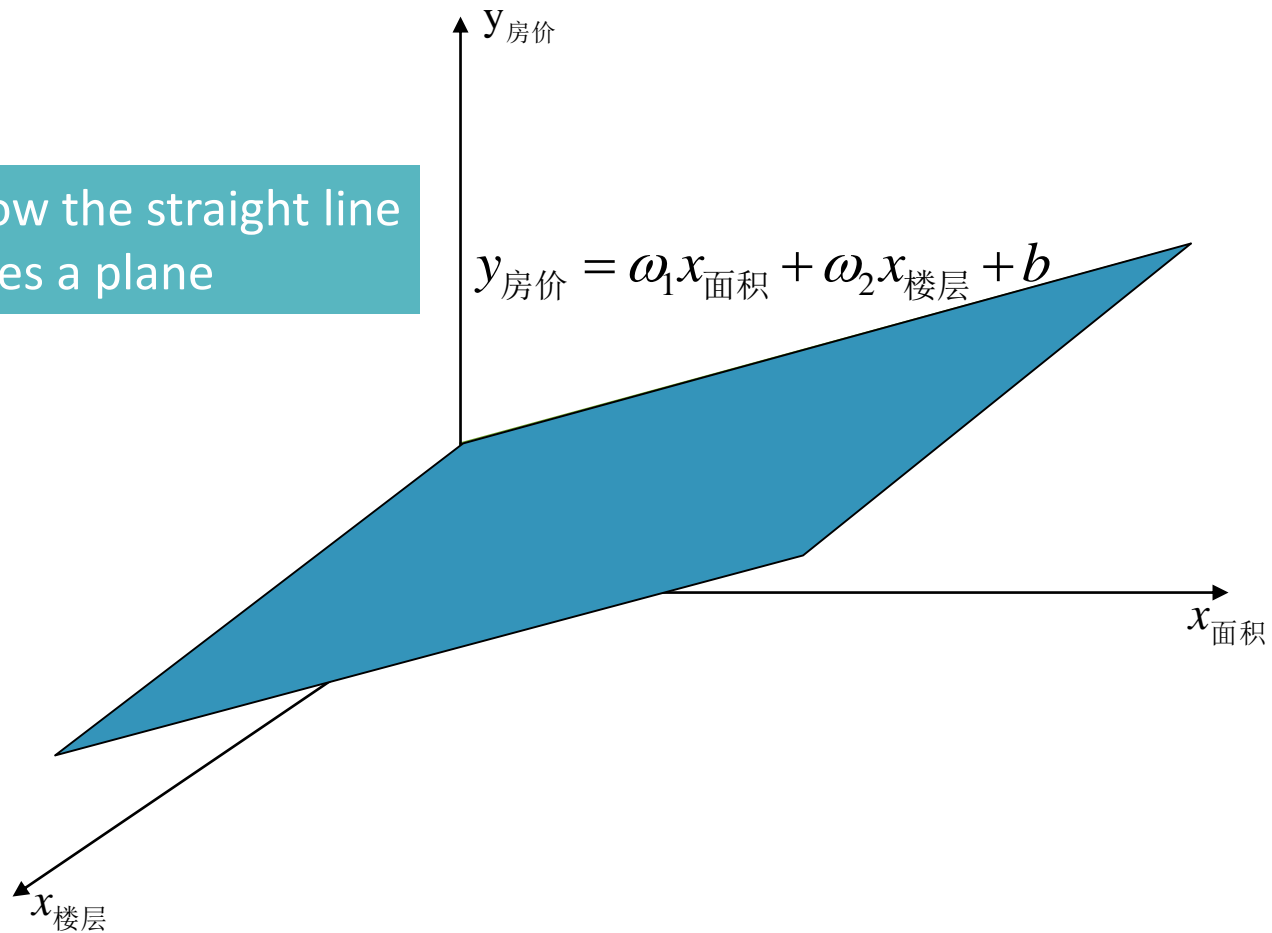
# linear regression model





# linear regression model

Note how the straight line becomes a plane



# linear regression model

two-dimension:  $y_{\text{房价}} = \omega_1 x_{\text{面积}} + b$

three-dimension:  $y_{\text{房价}} = \omega_1 x_{\text{面积}} + \omega_2 x_{\text{楼层}} + b$

four-dimension:  $y_{\text{房价}} = \omega_1 x_{\text{面积}} + \omega_2 x_{\text{楼层}} + \omega_3 x_{\text{地段}} + b$

# linear regression model

general model:  $f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$

$x_1, x_2, \dots, x_d$ : the feature of sample

$w_1, w_2, \dots, w_d$ : **weight**, represent the importance of corresponding feature

$$f_{\text{房价}}(x) = 0.2 \cdot x_{\text{楼层}} + 0.5 \cdot x_{\text{地段}} + 0.3 \cdot x_{\text{面积}} + 1$$

importance:  $x_{\text{地段}} > x_{\text{面积}} > x_{\text{楼层}}$



# linear model's vector formal

general model:  $f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$

vector formal:  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$

$$\mathbf{w} = (w_1; w_2; \dots; w_d)$$

$$\mathbf{x} = (x_1; x_2; \dots; x_d)$$

# linear model's vector formal

training set:

S1: 楼层=6, 地段=5, 面积=120。房价=100

S2: 楼层=3, 地段=1, 面积=100。房价=200

S3: 楼层=16, 地段=8, 面积=200。房价=120

$$\begin{cases} 6\omega_1 + 5\omega_2 + 120\omega_3 = 100 \\ 3\omega_1 + \omega_2 + 100\omega_3 = 200 \\ 16\omega_1 + 8\omega_2 + 200\omega_3 = 120 \end{cases} \quad \longrightarrow \quad \begin{aligned} \omega_1 &= a \\ \omega_2 &= b \\ \omega_3 &= c \end{aligned}$$

$$f_{\text{房价}}(x) = a \cdot x_{\text{楼层}} + b \cdot x_{\text{地段}} + c \cdot x_{\text{面积}}$$

- Solve weight vector  $\mathbf{w}$  according to training set.

# Which model is better?

$$f_{\text{房价}}(x) = 0.2 \cdot x_{\text{楼层}} + 0.5 \cdot x_{\text{地段}} + 0.3 \cdot x_{\text{面积}} + 1$$

$$g_{\text{房价}}(x) = 0.1 \cdot x_{\text{楼层}} + 0.4 \cdot x_{\text{地段}} + 0.5 \cdot x_{\text{面积}} + 3$$

a known data: 楼层=6, 地段=5, 面积=120, 房价=100

model  $f$  predictive value:  $0.2 \cdot 6 + 0.5 \cdot 5 + 0.3 \cdot 120 + 1 = 40.7$

model  $g$  predictive value:  $0.1 \cdot 6 + 0.4 \cdot 5 + 0.5 \cdot 120 + 3 = 65.6$

model  $f$  error:  $100 - 40.7 = 59.3$

model  $g$  error:  $100 - 65.6 = 34.4$

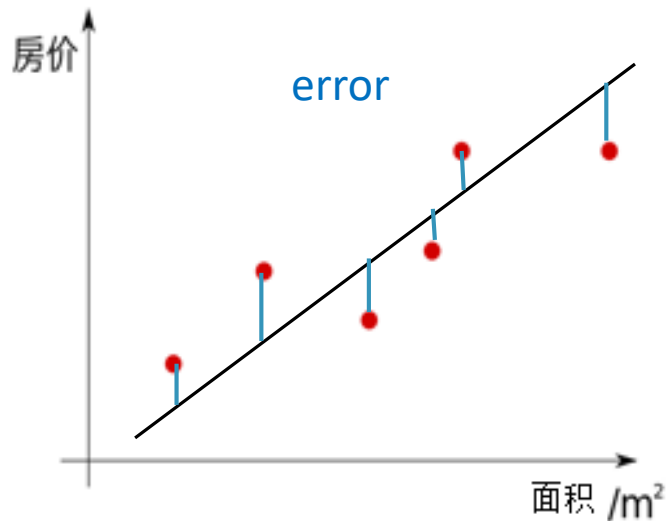
$$34.4 < 59.3$$

model  $g$  is better

- the model which can **minimize error** between model predictive value and true value.

# least square method

- In regression task, the **least square method** is often used measuring error between model predictive value and true value.



$$\text{error} = \sum_{i=1}^m (f(x_i) - y_i)^2$$

predictive value

true  
value

**least square method:**

$$\min \sum_{i=1}^m (f(x_i) - y_i)^2$$

# basic linear regression

- solve  $w$  and  $b$

single linear regression:

model:  $f(x_i) = wx_i + b$

loss function: 
$$E_{(w,b)} = \sum_{i=1}^m (f(x_i) - y_i)^2$$
$$= \sum_{i=1}^m (y_i - wx_i - b)^2$$

$$\min E_{(w,b)}$$



# basic linear regression

- solve  $w$  and  $b$


single linear regression:

$$E_{(w,b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$$

$$\frac{\partial E_{(w,b)}}{\partial w} = 2 \left( w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right) = 0$$

$$\frac{\partial E_{(w,b)}}{\partial b} = 2 \left( mb - \sum_{i=1}^m (y_i - wx_i) \right) = 0$$

new sample


$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left( \sum_{i=1}^m x_i \right)^2}$$
$$b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i) \quad \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

model:  $f(x_i) = wx_i + b$

dataset:  $D = \{(x_i, y_i)\}_{i=1}^m$

# basic linear regression

- solve  $w$  and  $b$  multiple linear regression:

$$\text{model: } f(\mathbf{x}_i) = w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} + b$$

$$\text{vector model: } f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$$



$* m$

$$f(\mathbf{x}) = \mathbf{X}\hat{\mathbf{w}}$$

$\hat{\mathbf{w}} = (\mathbf{w}; b)$   $(d+1)$  column vector

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix} \quad m \times (d+1) \text{ matrix}$$

$\mathbf{y} = (y_1; y_2; \dots; y_m)$   $m$  column vector, true values of samples

# basic linear regression

- solve  $\hat{\mathbf{w}}$

multiple linear regression:

model:  $f(\mathbf{x}) = \mathbf{X}\hat{\mathbf{w}}$

least square method

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$

$$\hat{\mathbf{w}} = (\mathbf{w}; b)$$

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix} \quad (1 \quad 2 \quad 3 \quad 4) \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} = 1^2 + 2^2 + 3^2 + 4^2$$

$$\mathbf{y} = (y_1; y_2; \dots; y_m)$$

# basic linear regression

• solve  $\hat{\mathbf{w}}$

multiple linear regression:

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$

向量求导公式:

$$\frac{\partial(\mathbf{X}\theta)}{\partial\theta} = \mathbf{X}^T$$

$$\frac{\partial(\theta^T \mathbf{X})}{\partial\theta} = \mathbf{X}$$

$$\frac{\partial(\theta^T \mathbf{X}\theta)}{\partial\theta} = (\mathbf{X}^T + \mathbf{X})\theta$$

$$\begin{aligned} &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}} - \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{y} + \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} \\ &\quad - 2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} \end{aligned}$$

令  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$ , 对  $\hat{\mathbf{w}}$  求导得到

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2 \mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) = 0$$



analytical solution:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$   
解析解

# basic linear regression

- multiple linear regression:

dataset:  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$

1. substitute

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

2. get

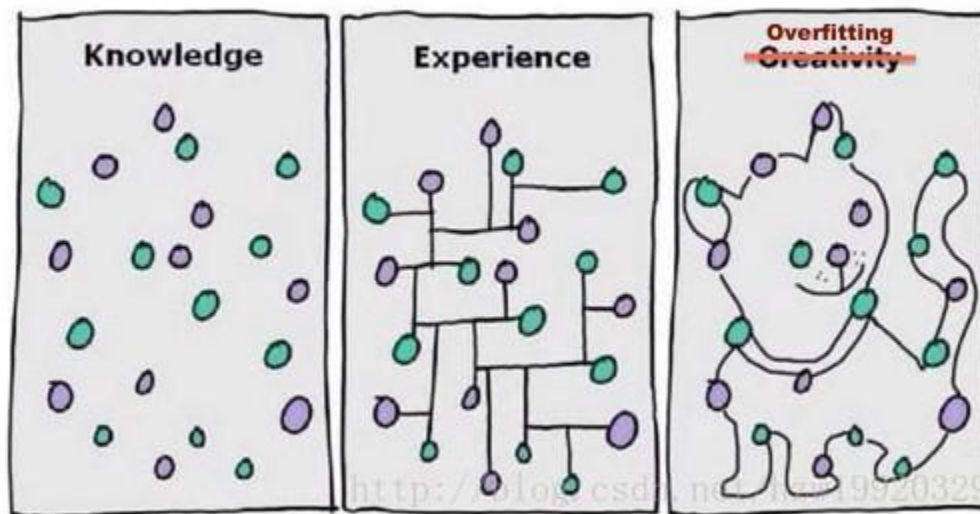
model:  $f(\mathbf{x}) = \mathbf{X}\hat{\mathbf{w}}$

new sample

# least square method's drawback

loss function: 
$$E_{(w,b)} = \sum_{i=1}^m (f(x_i) - y_i)^2$$

(1)最小二乘法的损失函数有可能会造成模型过拟合。



# least square method's drawback

loss function:  $E_{(w,b)} = \sum_{i=1}^m (f(x_i) - y_i)^2$

analytical solution:

$$\hat{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- (2)解析解中  $\mathbf{X}^T \mathbf{X}$  有可能不是满秩矩阵, 那样  $(\mathbf{X}^T \mathbf{X})^{-1}$  不存在,
- 样本的特征数远远超过样本数
- 这样,  $\hat{w}$  会有多个解。

# least square method's application

1. 已有样本中的值比较准确
2. 样本的特征数小于样本数





# Outline

- linear regression model
  - linear model
  - least square method --> basic linear regression
  - ridge regression(岭回归)
  - lasso regression

# ridge regression & lasso regression

**model:**  $f(\mathbf{x}_i) = w_1x_{i1} + w_2x_{i2} + \dots + w_dx_{id} + b$

  
vector model:  $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$

  
 $f(\mathbf{x}) = \mathbf{X}\hat{\mathbf{w}}$

**loss function:**

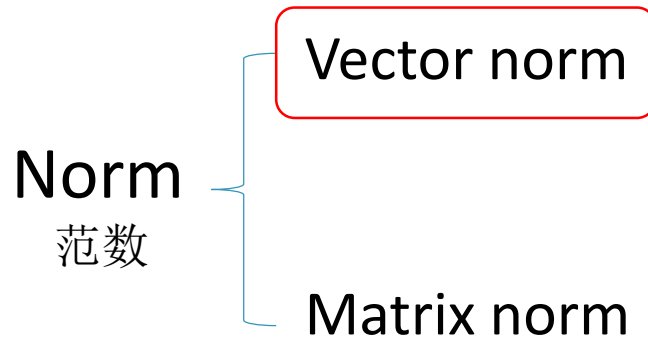
basic linear regression:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$

ridge regression:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_2^2 \quad \lambda > 0$

lasso regression:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_1 \quad \lambda > 0$

ridge 回归和lasso回归可以防止模型过拟合，也可以解决样本的特征数远超过样本数的问题。

# Norm



- to measure the size of vector and matrix
- Norm is a value.

# Vector norm

$$\mathbf{w} \in \mathbf{R}^d$$

$L_0$ -norm: the number of non-zero elements in the vector

$$L_1\text{-norm: } \|\mathbf{w}\|_1 = |w_1| + |w_2| + \dots + |w_d| = \sum_{i=1}^d |w_i|$$

$$L_2\text{-norm: } \|\mathbf{w}\|_2 = \sqrt{w_1^2 + w_2^2 + \dots + w_d^2} = \left( \sum_{i=1}^d w_i^2 \right)^{\frac{1}{2}}$$

$$L_\infty\text{-norm: } \|\mathbf{w}\|_\infty = \max\{|w_1|, |w_2|, \dots, |w_d|\} = \max_{1 \leq i \leq d} \{|w_i|\}$$



$$\mathbf{w} = [1 \quad 2 \quad -2]$$

$$\|\mathbf{w}\|_0 = 3 \quad \|\mathbf{w}\|_1 = 5$$

$$\|\mathbf{w}\|_2 = 3 \quad \|\mathbf{w}\|_\infty = 2$$

# Outline

- linear regression model
  - linear model
  - least square method --> basic linear regression
  - ridge regression(岭回归)
  - lasso regression

# ridge regression

$$\text{model: } f(\mathbf{x}_i) = w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} + b$$

$$\begin{array}{c} \updownarrow \\ f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b \end{array}$$

**loss function:**  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_2^2 \quad \lambda > 0$

向量求导公式:

$$\frac{\partial(\mathbf{X}\theta)}{\partial\theta} = \mathbf{X}^T$$

$$\frac{\partial(\theta^T \mathbf{X})}{\partial\theta} = \mathbf{X}$$

$$\frac{\partial(\theta^T \mathbf{X}\theta)}{\partial\theta} = (\mathbf{X}^T + \mathbf{X})\theta$$

$$= (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \hat{\mathbf{w}}^T \hat{\mathbf{w}}$$

$$= (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \hat{\mathbf{w}}^T \mathbf{I} \hat{\mathbf{w}}$$

$\mathbf{I}$ : identity matrix

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2 \mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) + \lambda \mathbf{I} \hat{\mathbf{w}} = 0$$



analytical solution:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$   
解析解

# ridge regression

## basic linear regression:

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$

analytical solution:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

## ridge regression:

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_2^2 \quad \lambda > 0$

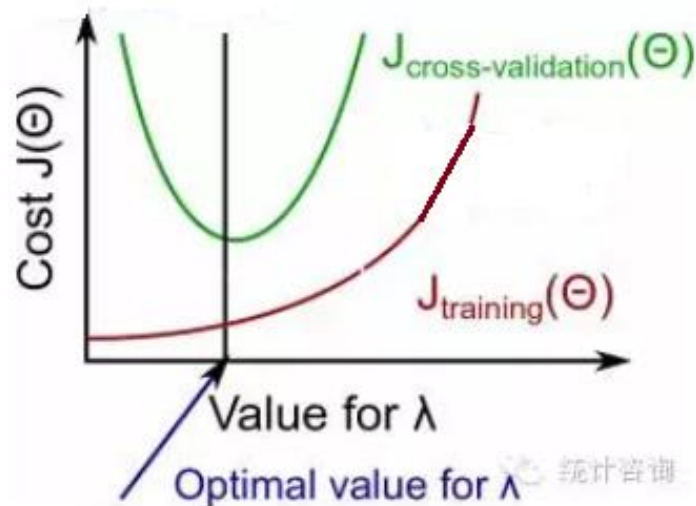
analytical solution:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$  is full rank matrix,  $\hat{\mathbf{w}}$  can get analytical solution.

可以解决特征数大于样本数的问题

# ridge regression

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_2^2 \quad \lambda > 0$






# Outline


- linear regression model
  - linear model
  - least square method --> basic linear regression
  - ridge regression(岭回归)
  - lasso regression

# lasso regression

- Least Absolute Shrinkage and Selection Operator,最小绝对收敛算子

**model:**  $f(\mathbf{x}_i) = w_1x_{i1} + w_2x_{i2} + \dots + w_dx_{id} + b$

  
vector model:  $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$

  
 $f(\mathbf{x}) = \mathbf{X}\hat{\mathbf{w}}$

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_1 \quad \lambda > 0$

$$\|\mathbf{w}\|_1 = |w_1| + |w_2| + \dots + |w_d| = \sum_{i=1}^d |w_i|$$

# lasso regression

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_1 \quad \lambda > 0$

函数在某点可导条件:

- 1) 函数在该点连续
- 2) 函数在该点左右两侧导数都存在并且相等。

$$y = |\omega|, \omega \in R$$

$$y'(0^-) = -1, y'(0^+) = 1$$

# lasso regression

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_1 \quad \lambda > 0$

- **Solving** loss function

Due to  $\|\mathbf{w}\|_1$  having not derivative(导数), we use **proximal gradient descent method**(近端梯度下降法) to minimize the loss function of lasso regression.

**analytical solution:**

$$w_{k+1}^i = \begin{cases} z^i - \lambda/L, & \lambda/L < z^i; \\ 0, & |z^i| \leq \lambda/L; \\ z^i + \lambda/L, & z^i < -\lambda/L, \end{cases} \quad L > 0 \text{ 的常数}$$
$$\mathbf{z} = \mathbf{w}_k - 1/L \nabla f(\mathbf{w}_k)$$

可以解决特征数大于样本数的问题

# lasso regression

Lasso regression is easier to get **sparse solution**  $\hat{w}$  .

sparse solution(稀疏解):

- $w$  contains **less non-zero** elements.
- It can use **feature selection**.  
 $w_d = 0$  means the feature  $x_{id}$  is **not important** to the task.

# lasso regression

model:  $f(\mathbf{x}_i) = w_1x_{i1} + w_2x_{i2} + \dots + w_dx_{id} + b$

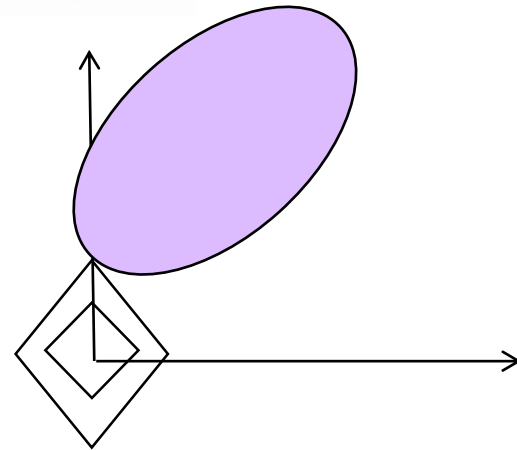
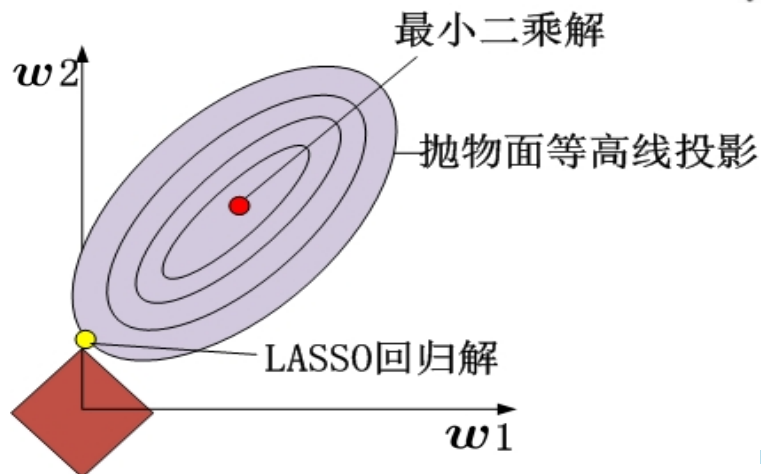


$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$$

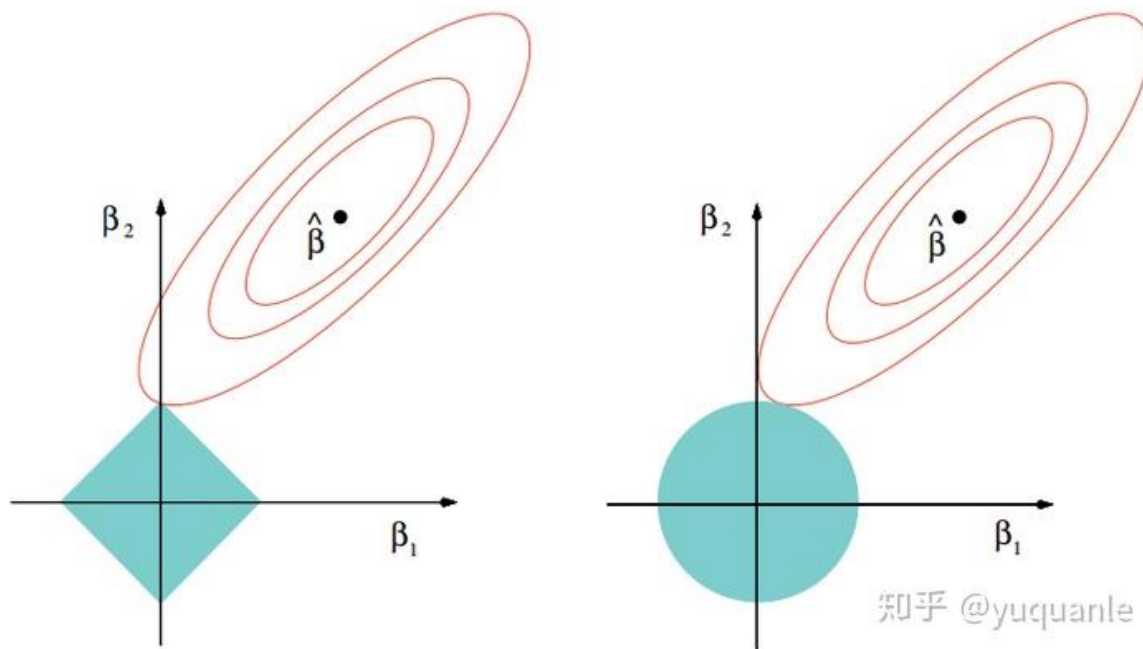
loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_1 \quad \lambda > 0$

$$\|\mathbf{w}\|_1 = |w_1| + |w_2| + \dots + |w_d|$$

basic linear regression:  $E_{(w,b)} = \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2$



# lasso vs ridge regression



Lasso和岭回归的区别很好理解，在优化过程中，最优解为函数等值线与约束空间的交集，正则项可以看作是约束空间。可以看出二范的约束空间是一个球形，一范的约束空间是一个方形，这也就是二范会得到很多参数接近0的值，而一范会尽可能非零参数最少。

# lasso vs ridge regression

- 1) **Prevent** model overfitting.
- 2) solve the problem that the number of features is larger than the number of samples.

ridge regression

lasso regression

It is easier to get **sparse solution**  $w$ .



# function

matlab:

```
b = ridge(y,X,k); %ridge regression
```

```
B = lasso(X,Y); %lasso regression
```

# Summary

- linear regression model

- linear model

model:  $f(\mathbf{x}_i) = w_1x_{i1} + w_2x_{i2} + \dots + w_dx_{id} + b$

vector model:  $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$

$f(\mathbf{x}) = \mathbf{X}\hat{\mathbf{w}}$

- least square method --> basic linear regression

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$

analytical solution:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

# Summary

- linear regression model
  - ridge regression(岭回归)

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_2^2 \quad \lambda > 0$

analytical solution:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

- lasso regression

loss function:  $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_1 \quad \lambda > 0$

Lasso regression is easier to get **sparse solution**  $\hat{\mathbf{w}}$  .

# Reference

- [1]周志华. 机器学习[M]. 北京：清华大学出版社，2016
- [2]李航. 统计学习方法[M]. 北京：清华大学出版社，2012
- [3]于剑. 机器学习：从公理到算法[M]. 北京：清华大学出版社，2017



Thank you !