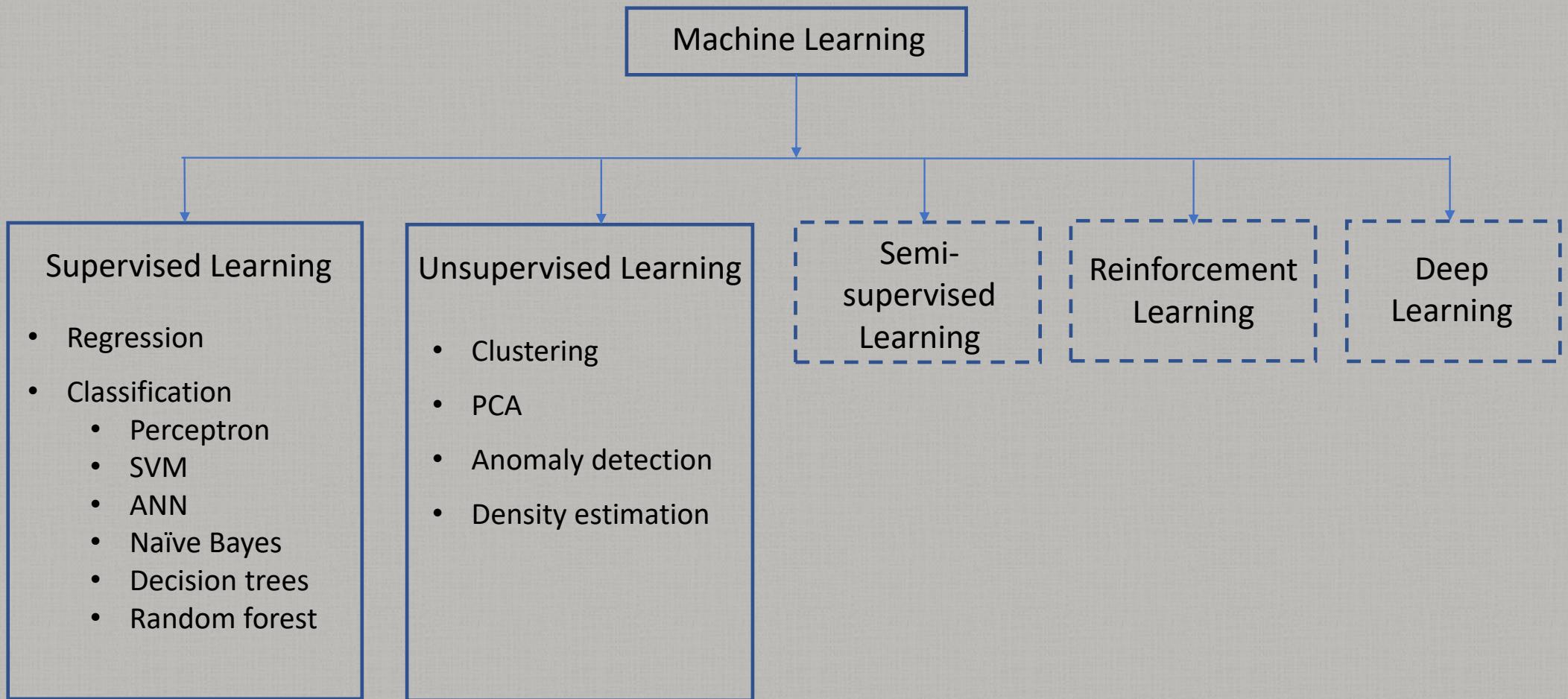


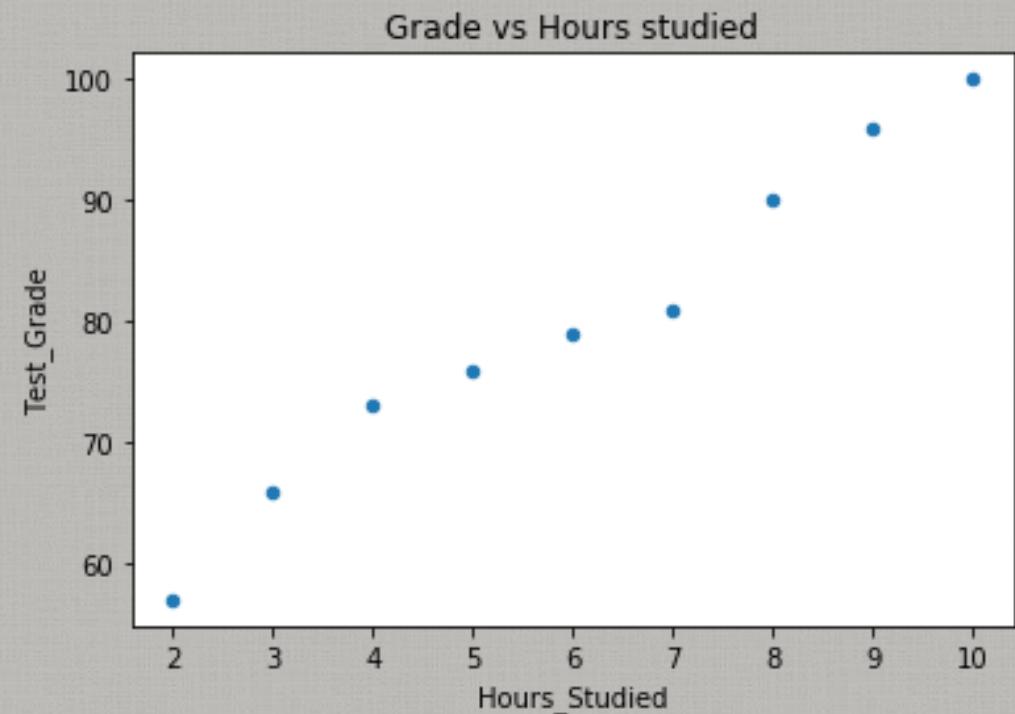
# Types of Machine Learning Problems



# Regression Problems

- **Regression** is a supervised learning problem with real valued target values  $y = f(x)$
- Find the hypothesis  $\hat{y} = h(x)$  to predict the target value of unseen instance  $x$

input	output
Hours_Studied	Test_Grade
2	57
3	66
4	73
5	76
6	79
7	81
8	90
9	96
10	99



# House Price Prediction

output

*features*

price	lotsize	bedrooms	bathrms	stories	driveway	recroom	fullbase	gashw	airco	garagepl	prefarea
42000	5850	3	1	two	yes	no	yes	no	no	1	no
38500	4000	2	1	one	yes	no	no	no	no	0	no
49500	3060	3	1	one	yes	no	no	no	no	0	no
60500	6650	3	1	two	yes	yes	no	no	no	0	no
61000	6360	2	1	one	yes	no	no	no	no	0	no
66000	4160	3	1	one	yes	yes	yes	no	yes	0	no
66000	3880	3	2	two	yes	no	yes	no	no	2	no
69000	4160	3	1	three	yes	no	no	no	no	0	no
83800	4800	3	1	one	yes	yes	yes	no	no	0	no

- Notice some features in the data set are categorical which need to be turned into numerical values. (<https://scikit-learn.org/stable/modules/preprocessing.html#encoding-categorical-features>)

# Regression Models

最小二乘法

- **Ordinary Least Squares (OLS)** method, the simplest linear regression problem
- Expanded linear least squares using nonlinear **feature mapping**
- **Polynomial regression** 多项式回归  
岭回归 收缩 正则化 特征映射
- **Ridge regression** (Shrinkage method for **regularization**)
- **Maximum Likelihood Estimate** 极大似然估计
- **Maximum a posterior (MAP) method** 最大后验估计
- The relations among these methods

## The *regression problem*:

- Given the training data set  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ , with  $x_i \in \mathbb{R}^d$  the input (feature) vector, and  $y_i \in \mathbb{R}$  the output (target value)
- Governed by unknown mapping  $y = f(x)$
- Hypothesis set  $h_w(x)$  characterized by the weight vector  $w \in \mathbb{R}^d$
- Want to find  $\hat{y} = h_{w^*}(x)$  that best approximate  $y = f(x)$
- We concern about the prediction accuracy, i.e., how accurate the estimated model can make prediction on unseen data

$\nearrow (w_0, w_1)$

$$y = a \times + b = (1, x) \begin{pmatrix} a \\ b \end{pmatrix} = \vec{x}^T \vec{w}$$

**Linear hypothesis model in one-dimensional case:**  $\vec{x} = [1 \ x]^T$   $\vec{w} = [w_0 \ w_1]^T$

- When  $d = 1$ , i.e.,  $x = x_1$ . The linear hypothesis model for this case should be

$$h_w(x) = w_0 + w_1 x_1 \Rightarrow (1 \ x) \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$$

Where  $w_0$  is the **interception** and  $w_1$  is the **slope**.

- The above expression can be written in the vector form as

$$h_w(x) = [1 \ x_1] \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \vec{x}^T \vec{w}$$

Where  $\vec{x} = [1 \ x_1]^T$ ,  $\vec{w} = [w_0 \ w_1]^T$  are the **expanded feature vector** and **weight vector**.

## ~~Linear hypothesis model in multi-dimensional case:~~

- In general, the linear hypothesis model in high dimensional feature space is,

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_dx_d$$

- let's expand the original feature vector and weight vector as

$$\mathbf{x} = [1 \ x_1 \ x_2 \ \dots \ x_d]^T, \mathbf{w} = [w_0 \ w_1 \ w_2 \ \dots \ w_d]^T$$

- We then have,

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = [1 \ x_1 \ x_2 \ \dots \ x_d] \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \sum_{j=0}^d w_j x_j$$

- Since  $h_{\mathbf{w}}(\mathbf{x})$  is linear with respect to  $\mathbf{w}$  and  $\mathbf{x}$ , the problem is called a ***linear regression problem***.



最小二乘法

**Ordinary Least Squares (OLS) method**

- The linear hypothesis model is used:

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$$

- For each input  $\mathbf{x}_i$  in the data set, the predicted target value is:

$$\hat{y}_i = h_{\mathbf{w}}(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}, \quad i = 1, 2, \dots, N.$$

And we want that:

$$\hat{y}_i \approx y_i, \quad i = 1, 2, \dots, N$$

- We want to determine the value  $\mathbf{w}^*$  of the weight vector so that the model prediction is a close approximation of the target value in the data set

- Let's represent the predictions over the data set as:

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ 1 & x_{21} & \cdots & x_{2d} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & x_{N1} & \cdots & x_{Nd} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} = \mathbf{X}\mathbf{w}$$

Where,

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ 1 & x_{21} & \cdots & x_{2d} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & x_{N1} & \cdots & x_{Nd} \end{bmatrix}$$

*1st data point*

*d - number of features*

The matrix  $\mathbf{X} \in \mathcal{R}^{N \times (d+1)}$  has the input data point  $x_i$  as its  $i^{th}$  row.  $\mathbf{X}$  is sometimes called the ***design matrix***.

Let's specify the **loss (cost) function**:

- We want to find the value of  $w$  that minimizes the **sum of squared error (SSE)**:

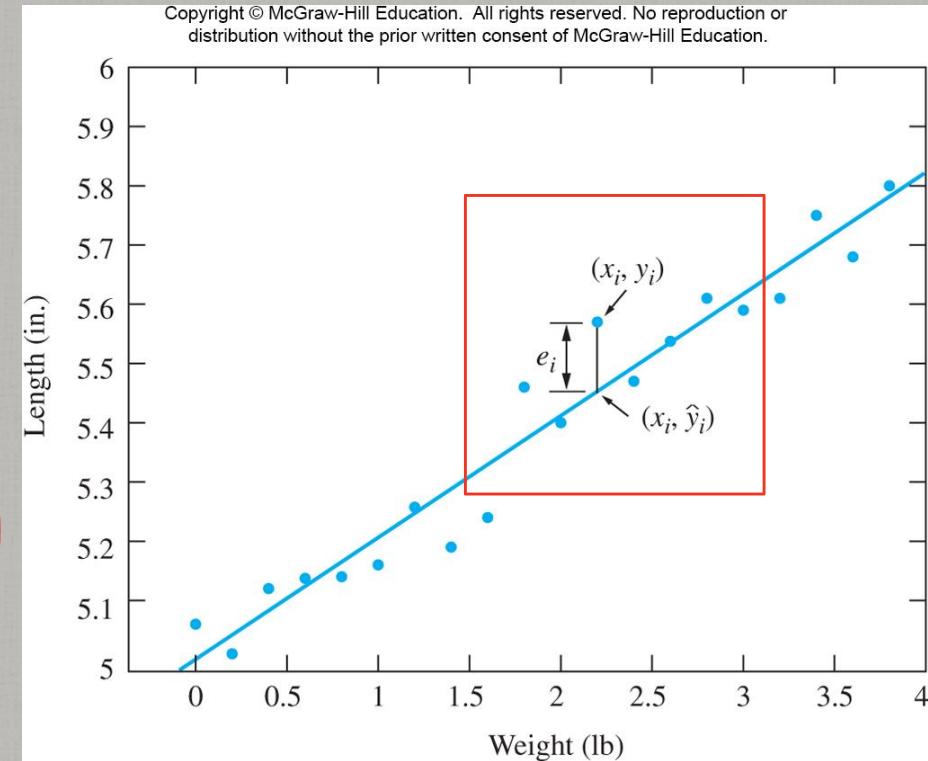
$$L(w) = \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \sum_{i=1}^N (h_w(x_i) - f(x_i))^2$$

$$= \sum_{i=1}^N (x_i^T w - y_i)^2 = \|Xw - y\|_2^2$$

$$\begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_N \end{pmatrix}$$

Where,

- $y = [y_1, y_2, \dots, y_N]^T$  are the **target values** in the training data set.
- $e_i = \hat{y}_i - y_i$  is called the **residual** of the  $i^{th}$  instance.   
 *target value*   
 *predicted value*   
 *残差*



- Now, we have the following ***optimization*** problem: (**Least Squares**)

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \{L(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2\}$$

使  $L(\mathbf{w})$  取得最小值时  $\mathbf{w}$  的值

This problem can be solved to have a closed-form solution!

- Let's solve this optimization problem:

- Let's find the gradient of the loss function  $L(\mathbf{w})$ :

$$\begin{aligned}
 L(\mathbf{w}) &= \|X\mathbf{w} - \mathbf{y}\|_2^2 = (X\mathbf{w} - \mathbf{y})^T(X\mathbf{w} - \mathbf{y}) \\
 &= (X\mathbf{w})^T X\mathbf{w} - (X\mathbf{w})^T \mathbf{y} - \mathbf{y}^T X\mathbf{w} + \mathbf{y}^T \mathbf{y} \\
 &= \mathbf{w}^T X^T X\mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}
 \end{aligned}$$

- Using the following results from matrix calculus:

表示对 $\mathbf{x}$ 求梯度，对向量 $\mathbf{x}$ 中的每一个分量进行偏导

$$\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{b}) = \mathbf{b}$$

$$\nabla_{\mathbf{x}}(\mathbf{x}^T A \mathbf{x}) = 2A\mathbf{x}$$

- Then we have the gradient of  $L(\mathbf{w})$  to be:

$$\begin{aligned}
 \nabla_{\mathbf{w}} L(\mathbf{w}) &= \nabla_{\mathbf{w}}(\mathbf{w}^T X^T X\mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \\
 &= \nabla_{\mathbf{w}}(\mathbf{w}^T X^T X\mathbf{w}) - 2\nabla_{\mathbf{w}}(\mathbf{w}^T X^T \mathbf{y}) + \nabla_{\mathbf{w}}(\mathbf{y}^T \mathbf{y}) \\
 &= 2X^T X\mathbf{w} - 2X^T \mathbf{y}
 \end{aligned}$$

$L(w)$  取小

- Setting the gradient to  $\mathbf{0}$ , we have,

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

滿秩

- if  $\mathbf{X}^T \mathbf{X}$  is *full rank*, then,


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

解析解 / 闭合解

(closed-form solution!)

analytical solution

- Is  $\mathbf{w}^*$  a global minima?

Since the loss function is convex,  $\mathbf{w}^*$  is the global minima.

- To show this, it is sufficed to compute the Hessian of  $L$  and show it is positive semi-definite:

HL

$$\nabla^2 L(\mathbf{w}) = 2\mathbf{X}^T \mathbf{X}$$

$$\forall \mathbf{w}, \mathbf{w}^T (2\mathbf{X}^T \mathbf{X}) \mathbf{w} = 2(\mathbf{X}\mathbf{w})^T \mathbf{X}\mathbf{w} = 2\|\mathbf{X}\mathbf{w}\|_2^2 \geq 0$$

- The OLS problem has a ***closed-form*** solution:

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

- Notice that  $\mathbf{X}$  and  $\mathbf{y}$  can be constructed using the given data set!
- Can you recognize some potential problems using this closed-form solution?
  - Computing  $(\mathbf{X}^T \mathbf{X})^{-1}$  can be very expensive when the data set is large and has many features.  
    滿秩
  - We get a problem when  $\mathbf{X}^T \mathbf{X}$  is not full rank or close to not full rank
- Do we have other method to find  $\mathbf{w}^*$ ?

Example

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

input	output
Hours_Studied	Test_Grade
2	57
3	66
4	73
5	76
6	79
7	81
8	90
9	96
10	99

$$X = \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \\ 1 & 6 \\ 1 & 7 \\ 1 & 8 \\ 1 & 9 \\ 1 & 10 \end{bmatrix}, \quad y = \begin{bmatrix} 57 \\ 66 \\ 73 \\ 76 \\ 79 \\ 81 \\ 90 \\ 96 \\ 99 \end{bmatrix}$$

$$\mathbf{w}^* = \begin{bmatrix} 50 \\ 4.95 \end{bmatrix}$$

- Discuss on the linear models

➤ ***Advantages:***

- Simple, light-weighted, good ***generalization*** (especially when the training data set is small)

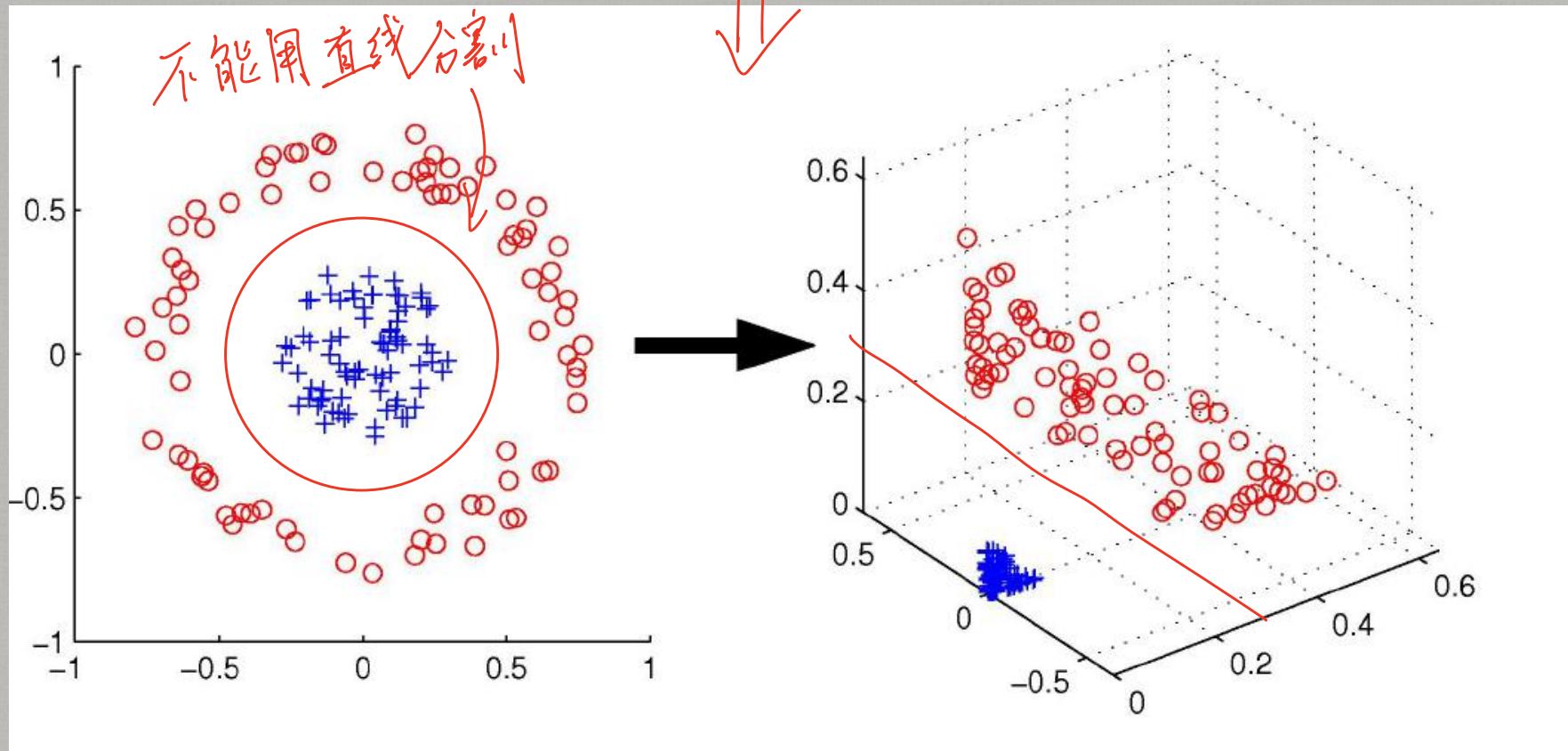
➤ ***Disadvantages:***

- Simple, lack of ***expressive power***, large ***bias*** for sophisticated mapping

表现力

非线性特征映射

## Use Nonlinear Feature Mapping



# Non-linear Feature Map

- The true input-output relationship  $y = f(x)$  maybe **nonlinear**. It is useful to consider nonlinear model as well. This can be achieved by **augmenting the data with nonlinear feature mapping**.  
*in it*
- We devise some function  $\phi: \mathcal{R}^l \rightarrow \mathcal{R}^{d+1}$ , (usually  $d > l$ ) called a **feature map**, that maps each raw data point  $x_i \in \mathcal{R}^l$  into a vector of new features  $\phi(x_i) \in \mathcal{R}^{d+1}$ :

$$\phi(x_i)^T = [\phi_0(x_i) \quad \phi_1(x_i) \quad \cdots \quad \phi_d(x_i)]$$



## Non-linear Feature Map

$$\begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \\ \vdots \\ \boldsymbol{x}_N \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \dots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nl} \end{bmatrix} \Rightarrow \begin{bmatrix} \phi_0(\boldsymbol{x}_1) & \phi_1(\boldsymbol{x}_1) & \dots & \phi_d(\boldsymbol{x}_1) \\ \phi_0(\boldsymbol{x}_2) & \phi_1(\boldsymbol{x}_2) & \dots & \phi_d(\boldsymbol{x}_2) \\ \vdots & \vdots & \dots & \vdots \\ \phi_0(\boldsymbol{x}_N) & \phi_1(\boldsymbol{x}_N) & \dots & \phi_d(\boldsymbol{x}_N) \end{bmatrix}$$

*$\phi(x)$*

## Expanded Linear Least Squares

- The hypothesis function then becomes:

$$h_{\mathbf{w}}(\mathbf{x}) = \sum_{j=0}^d w_j \phi_j(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}$$

(linear combination of nonlinear mapping functions)

- Note that this model is still *linear* with respect to the mapped features, but it is nonlinear with respect to the original data if  $\phi_j$ s are nonlinear. This model is ***linear with respect to the weights!***

- The component functions  $\phi_j$  are sometimes called **basis functions**. (polynomial, piecewise linear, Gaussian, sigmoid, etc.)
- Example **basis functions**:

$$\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s^2}} \quad (\text{Gaussian})$$

$$\phi_j(x) = \frac{1}{1 + e^{-\frac{x-\mu_j}{s}}} \quad (\text{sigmoid})$$

基函数

分段线性

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

- we can then use least squares to estimate the weights  $\mathbf{w}$  to minimize  $\|\Phi\mathbf{w} - \mathbf{y}\|^2$ , where,  $\Phi \in \mathcal{R}^{N \times (d+1)}$  is the ***design matrix*** as follows:

$X \Rightarrow \phi$

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_d(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_d(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_d(x_N) \end{bmatrix}$$

- The solution of the ***expanded least square problem*** is:

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



## Polynomial Feature Map (Transformation)

多项式特征映射

- **Polynomial feature mapping** for *one dimensional* input:

$$x \Rightarrow \phi(x) = [1 \quad x \quad x^2 \quad x^3 \quad \dots \quad x^m]^T$$

- The hypothesis models set becomes:

$$h_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_mx^m$$

$$= [1 \quad x \quad x^2 \quad x^3 \quad \dots \quad x^m] \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$$

$$= \phi(x)^T w$$

## 单变量

- The **univariate** polynomial model with degree m:

$$y = w_0 + w_1 x \quad (m=1)$$

$$y = w_0 + w_1 x + w_2 x^2 \quad (m=2)$$

$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3 \quad (m=3)$$

⋮

$$y = w_0 + w_1 x + w_2 x^2 + \cdots + w_m x^m$$

## 双变量

- Bivariate polynomial model** (two-dimensional input, m = 2):

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_{11} x_1^2 + w_{22} x_2^2 + w_{12} x_1 x_2$$

# Polynomial Features

- A big reason that we care about polynomial features is that ***any smooth function can be approximated arbitrarily closely by some polynomial.***
- For this reason, polynomial are said to be a ***universal approximator.***

通用近似器

- For simplicity, we use a 1D input as an example to find the solution.
- The data set:  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  governed by the unknown relation:  $y = f(x)$
- The hypothesis model:

$$\hat{y} = h_{\mathbf{w}}(x) = w_0 + w_1x + w_2x^2 + \cdots + w_mx^m = \boldsymbol{\phi}(x)^T \mathbf{w}$$

Where,  $\mathbf{w} = [w_0 \quad w_1 \quad w_2 \quad \dots \quad w_m]^T$  is the weight vector, and

$\boldsymbol{\phi}(x)^T = [1 \quad x \quad x^2 \quad x^3 \quad \dots \quad x^m]^T$  is the mapped feature vector.

- This model is still linear with respect to the weight vector  $\mathbf{w}$

- From the linear hypothesis model, we have,

$$\hat{y}_1 = w_0 + w_1 x_1 + w_2 x_1^2 + \cdots + w_m x_1^m$$

$$\hat{y}_2 = w_0 + w_1 x_2 + w_2 x_2^2 + \cdots + w_m x_2^m$$

⋮

$$\hat{y}_N = w_0 + w_1 x_N + w_2 x_N^2 + \cdots + w_m x_N^m$$

In matrix form,

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \cdots & x_N^m \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$$

i.e.,

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$$

- To minimize the *sum of squared error*:

$$L(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

- We have,

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

- Note that  $m < N$ , i.e., the number of data is greater than the number of features.

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \cdots & x_N^m \end{bmatrix}$$

is the *design matrix*.

- Construction of the ***design matrix X***:

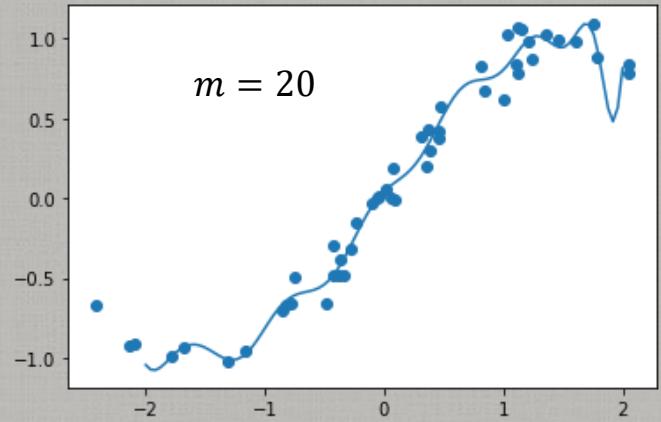
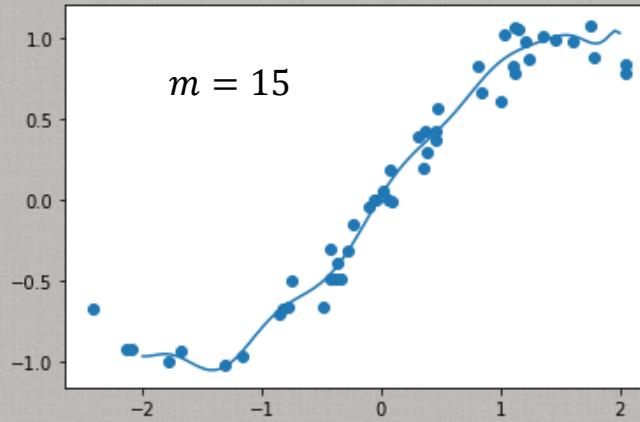
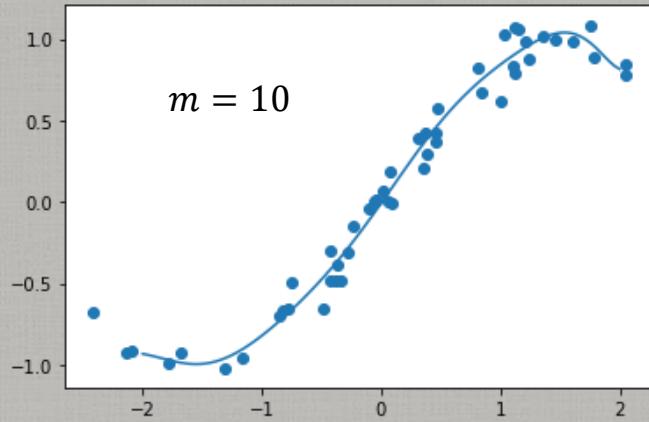
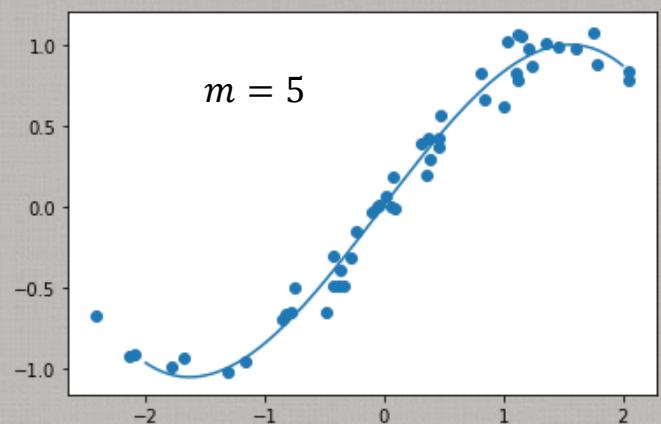
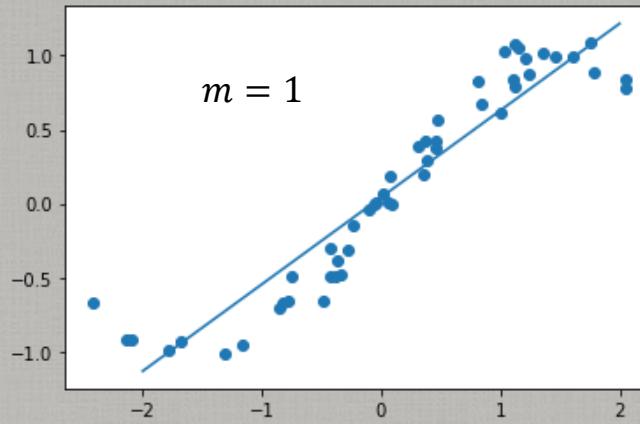
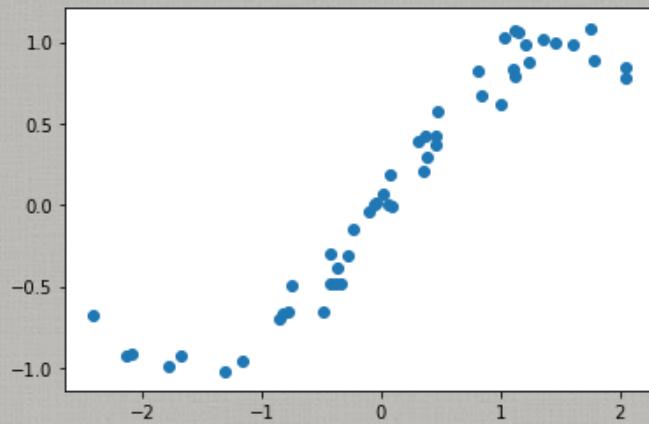
$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \cdots & x_N^m \end{bmatrix}$$

X	Y
Hours_Studied	Test_Grade
2	57
3	66
4	73
5	76
6	79
7	81
8	90
9	96
10	100

When  $m = 2$

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \cdots & x_N^2 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \\ 1 & 5 & 25 \\ 1 & 6 & 36 \\ 1 & 7 & 49 \\ 1 & 8 & 64 \\ 1 & 9 & 81 \\ 1 & 10 & 100 \end{bmatrix}$$

# Polynomial Regression Example



# Regularization

- **Numerical instability** with the closed-form solution formula:

数值不稳定性

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

共线性

*Numerical instability* arise when the features of the data are close to **collinear**, i.e., some features are linear combinations of others, causing the input matrix  $\mathbf{X}$  to lose its rank or have singular values that very close to 0

This will cause  $\mathbf{X}^T \mathbf{X}$  to be not invertible or it is invertible but  $(\mathbf{X}^T \mathbf{X})^{-1}$  will be huge, and the variance of the estimated coefficients  $\mathbf{w}^*$  will be enormous.

过拟合

- **Overfitting happens when complex hypothesis models are used (some components of the weight vector becomes extremely large!)**
- **Bad generalization** (bad prediction when applying the estimated model to unseen examples)

# Ridge Regression (a shrinkage method)

岭回归

模型复杂度通常由参数决定。当参数  $w$  过大时，过度拟合训练集的噪声点和异常点。

- There is a very simple solution to these issues: *penalize the large entries of  $w$  to prevent them from becoming too large.*
- We can do this by adding a penalty term constraining the norm of  $w$ . -- regularization
- For a fixed small scalar  $\lambda > 0$ , we now have,


$$L(w) = \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

超参数 (训练过程中需要预先设定的参数)

- $\lambda$  is called a *hyper parameter* that measures the sensitivity to the values in  $w$ .
- $\lambda$  is a value that we choose through *validation*. 验证

# Ridge Regression

- L2 regularization:

$$\begin{aligned} L(\mathbf{w}) &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} + \lambda \mathbf{w}^T \mathbf{w} \end{aligned}$$

$(w_1^2 + \dots + w_n^2)$

- Take the gradient of  $L(\mathbf{w})$  and let it be 0, we have,

$$\begin{aligned} \nabla_{\mathbf{w}} L(\mathbf{w}) &= 0 \\ \Rightarrow 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} &= 0 \\ \Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\ \Rightarrow \mathbf{w}_{ridge} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

- This value is guaranteed to achieve the unique global minimum, because the objective function is convex.

## Ridge Regression ( $L_2$ regularization)

- Let's find the **Hessian** of  $L(w)$ :

$$HL(w)$$

$$\nabla^2 L(w) = 2X^T X + 2\lambda I$$

$$\forall w \neq 0, w^T (2X^T X + 2\lambda I) w = (Xw)^T Xw + \lambda w^T w = \|Xw\|^2 + \lambda \|w\|^2 > 0$$

i.e.,  $\nabla^2 L(w)$  is **positive definite** (PD). The minimum is unique!

- Now with our slight tweak, the matrix  $X^T X + \lambda I$  has become full rank and invertible, meaning the numerical instability is solved.

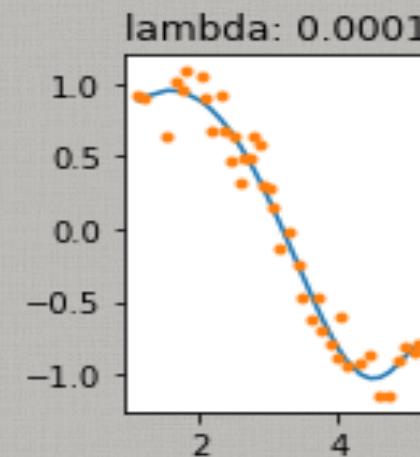
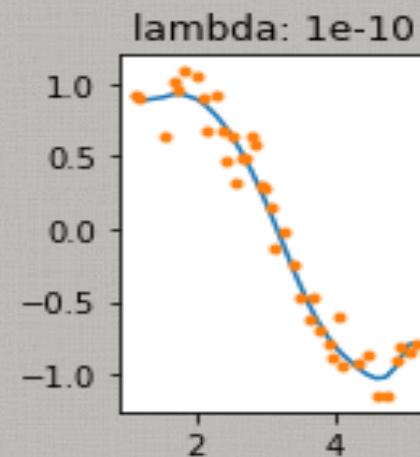
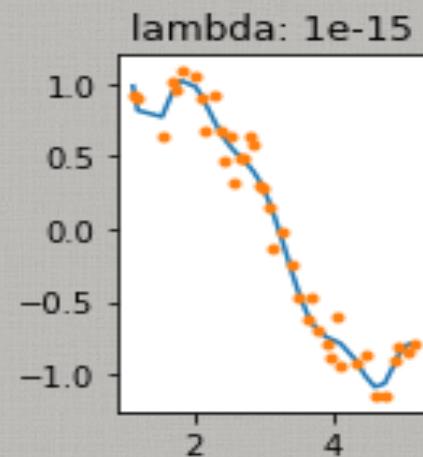


## Lasso Regression (L1 regularization)

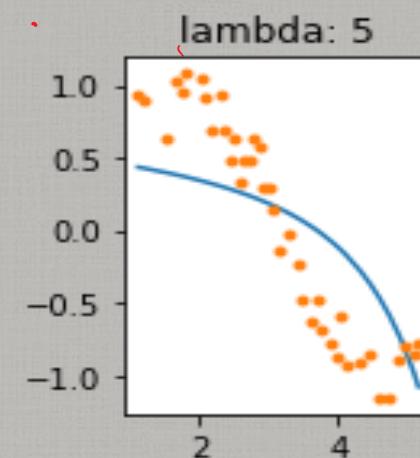
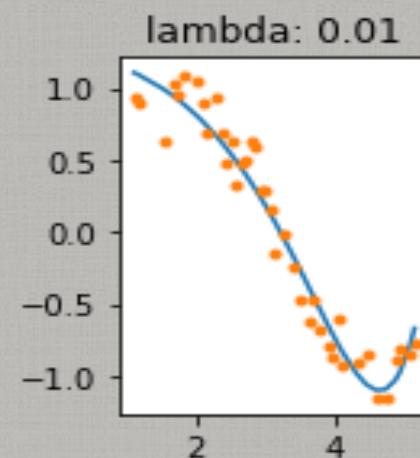
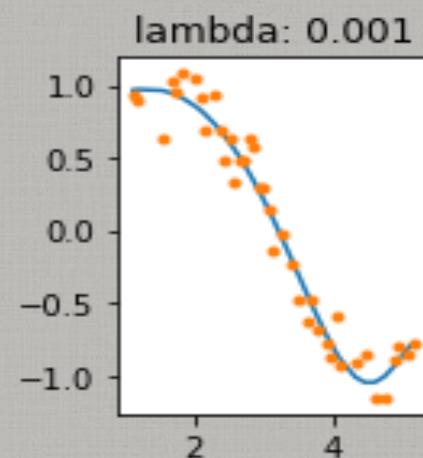
$$\mathbf{w}_{Lasso} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=0}^d |w_j| \right\}$$


Regularization technique can be applied to overcome overfitting!!!

$$L(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

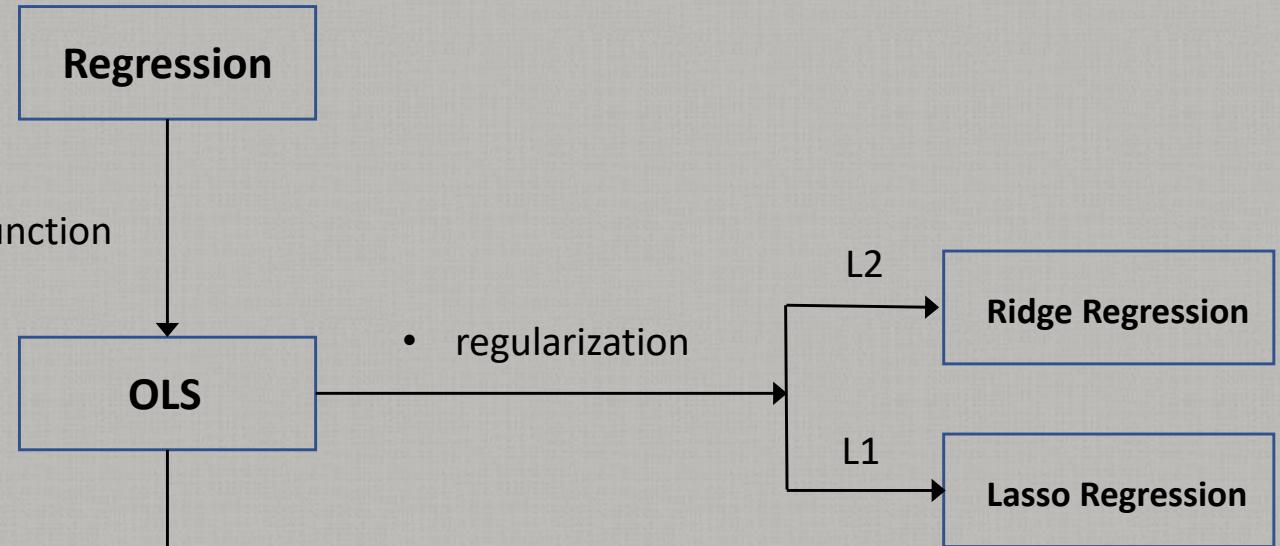


$\lambda \uparrow$



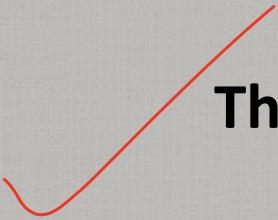
18

- Linear hypothesis model
- sum of squares error loss function
- Nonlinear feature mapping



Polynomial Regression

**The models developed so far do not involve any assumptions on distributions of data**



## The Maximum Likelihood Estimate (MLE) Model

mean 均值      方差 Variance  
 $N(\mu, \sigma^2)$

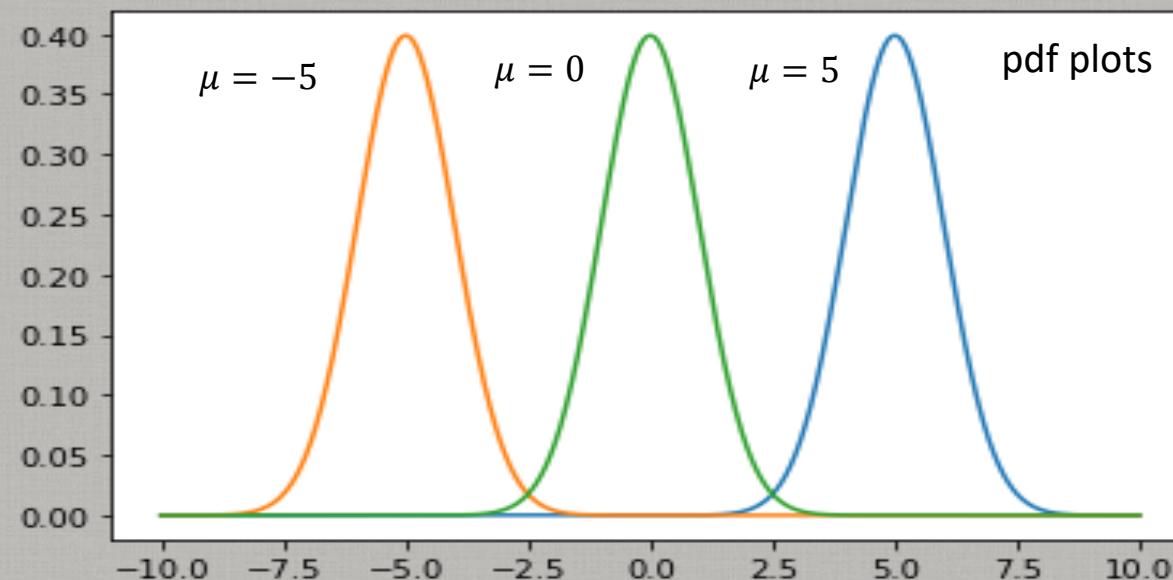
## Normal/Gaussian distribution

- Let random variable  $X \sim N(\mu, 1)$ . To estimate the value of  $\mu$ , a random sample of size 5 was drawn from  $X$  as follows:

$$\mu = 0.02$$

0.1	-0.9	0.8	0.2	-0.1
-----	------	-----	-----	------

- Let's also assume that  $\mu$  can only take three possible values,  $\mu = -5, 0, 5$ . Based on the sampled data, which is the most likely value of  $\mu$ ?



48

- Let random variable  $X \sim N(\mu, 1)$ . To estimate the value of  $\mu$ , a random sample of size 5 was drawn from  $X$  as follows:

0.1	-0.9	0.8	0.2	-0.1
-----	------	-----	-----	------

观测到数据  $X$  的联合分布概率，作为参数  $\theta$  的函数

- Let's define the following likelihood function:

$L(\theta | X) = P(X|\theta)$  在参数  $\theta$  下，数  
据  $X$  的联合分布

$$L(\mu) = p((x_1 = 0.1) \cap (x_2 = -0.9) \cap (x_3 = 0.8) \cap (x_4 = 0.2) \cap (x_5 = -0.1))$$

$$= p(x_1 = 0.1|\mu)p(x_2 = -0.9|\mu)p(x_3 = 0.8|\mu)p(x_4 = 0.2|\mu)p(x_5 = -0.1|\mu)$$

$$= \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(0.1-\mu)^2}{2}} \right) \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(-0.9-\mu)^2}{2}} \right) \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(0.8-\mu)^2}{2}} \right) \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(0.2-\mu)^2}{2}} \right) \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(-0.1-\mu)^2}{2}} \right)$$

- Maximizing ***likelihood function*** is equivalent to maximizing the ***log likelihood***, which is:

$$\begin{aligned}\ln(L(\mu)) &= \ln\left(\frac{1}{\sqrt{2\pi}}e^{-\frac{(0.1-\mu)^2}{2}}\right) + \ln\left(\frac{1}{\sqrt{2\pi}}e^{-\frac{(-0.9-\mu)^2}{2}}\right) + \ln\left(\frac{1}{\sqrt{2\pi}}e^{-\frac{(0.8-\mu)^2}{2}}\right) + \ln\left(\frac{1}{\sqrt{2\pi}}e^{-\frac{(0.2-\mu)^2}{2}}\right) + \ln\left(\frac{1}{\sqrt{2\pi}}e^{-\frac{(-0.1-\mu)^2}{2}}\right) \\ &= -\frac{5}{2}\ln(2\pi) - \frac{(0.1-\mu)^2}{2} - \frac{(-0.9-\mu)^2}{2} - \frac{(0.8-\mu)^2}{2} - \frac{(0.2-\mu)^2}{2} - \frac{(-0.1-\mu)^2}{2}\end{aligned}$$

*with respect to*

- Take the **derivative** of the ***log likelihood function*** w.r.t  $\mu$  and let it equal to 0:

$$(0.1 - \mu) + (-0.9 - \mu) + (0.8 - \mu) + (0.2 - \mu) + (-0.1 - \mu) = 0$$

$$\frac{\partial \ln(L(\mu))}{\partial \mu} = 0 \quad \Rightarrow -0.1 - 5\mu = 0 \quad \Rightarrow \hat{\mu} = -0.02$$



## Maximum Likelihood Estimate (MLE) – Probability Models

- In supervised learning, we assume that there exists a true underlying model mapping inputs to outputs:  $f(x)$
- The only information we have about the true model is via a data set:

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$$

Where  $x_i \in \mathcal{R}^d$  is the input and  $y_i \in \mathcal{R}$  is the observation of a random variable  $Y_i$ , i.e.,

$$Y_i = h_w(x_i) + Z_i$$

- We assume that  $x_i$  is a fixed value, while  $Z_i$  is a **random variable**.
- In most contexts, we assume that  $Z_i \sim N(0, \sigma^2), i = 1, 2, \dots, N$ . i.e.,  $Z_i$  are **independent identically distributed (i.i.d.) zero mean Gaussians**. Then,

$$Y_i \sim N(h_w(x_i), \sigma^2)$$

- **Maximum Likelihood Estimate (MLE)**, the goal is to find the hypothesis model that maximizes the *likelihood function* of the data, i.e.,

$$\widehat{\mathbf{w}}_{MLE} = \underset{\mathbf{w}}{\operatorname{argmax}} L(\mathbf{w}; \mathcal{D})$$

Where, the *likelihood function* is:

$$L(\mathbf{w}, \mathcal{D}) = p(Y_1 = y_1, \dots, Y_N = y_n | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}) = p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w})$$

$$= p(y_1 | \mathbf{x}_1, \mathbf{w}) p(y_2 | \mathbf{x}_2, \mathbf{w}) \cdots p(y_N | \mathbf{x}_N, \mathbf{w}) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \mathbf{w})$$

- Maximizing *likelihood function* is equivalent to maximizing the *log likelihood*, which is:

$$l(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \ln(L(\mathbf{w}; \mathcal{D})) = \sum_{i=1}^N \ln(p(y_i | \mathbf{x}_i, \mathbf{w}))$$

- When  $Y_i \sim N(h_w(x_i), \sigma^2)$ , we have, the **probability density function (pdf)** of  $Y_i$  is:

$$p(Y_i = y_i | x_i, w) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - h_w(x_i))^2}{2\sigma^2}}$$

- Then, the **nature log likelihood** becomes:

$$\max_l l(w; X, y) = -\sum_{i=1}^N \frac{(y_i - h_w(x_i))^2}{2\sigma^2} - N \ln(\sqrt{2\pi}\sigma)$$

- Maximizing  $l(w; X, y)$  is equivalent to the following:

$$\hat{w}_{MLE} = \underset{w}{\operatorname{argmin}} \left[ \sum_{i=1}^N \frac{(y_i - h_w(x_i))^2}{2\sigma^2} + N \ln(\sqrt{2\pi}\sigma) \right] = \underset{w}{\operatorname{argmin}} \left[ \sum_{i=1}^N (y_i - h_w(x_i))^2 \right]$$

target value  
predict value  
residual

- In the linear regression problem, our hypothesis model has the form:

$$h_{\mathbf{w}}(x_i) = \mathbf{x}_i^T \mathbf{w}$$

Then, the problem becomes:

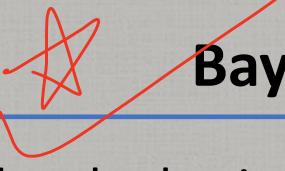
$$\hat{\mathbf{w}}_{MLE} = \underset{\mathbf{w}}{\operatorname{argmin}} \left[ \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \right] = \underset{\mathbf{w}}{\operatorname{argmin}} \{ \|X\mathbf{w} - \mathbf{y}\|_2^2 \}$$

This is just the **OLS** problem!

- This means that **MLE** is a probabilistic justification for why **sum of squared error (SSE)** is a good metric for evaluating a regression model!!!

- **Question:** what are the differences between **MLE** model and **OLS** model of regression?

- OLS is a special case of MLE when Gaussian distribution is assumed for  $Y_i$  and linear hypothesis model is used
- Other assumption can also be used for MLE!



## Bayes' Theorem

- **Bayesian reasoning** provides the basis for learning algorithms that directly manipulate probabilities.
- **Bayes' Theorem:**

Given a **hypothesis**  $h$  and a piece of **evidence**  $E$ , *Bayes Theorem* states that

$$P[h|E] = \frac{P[E|h]P[h]}{P[E]}$$

Where,

- $P[h|E]$  is the **posterior probability**,
- $P[h]$  is the **prior probability**,
- $P[E|h]$  is called **likelihood**.
- **Bayes' theorem** describes how to update the **probabilities of hypothesis when given an evidence**.

- ***Bayes' Theorem:***

Given multiple ***hypothesis***  $h_i, i = 1, 2, \dots M$  that are mutual exclusive and exhaustive, and a piece of **evidence**  $E$ , *Bayes Theorem* states that

$$P[h_i|E] = \frac{P[E|h_i]P[h_i]}{\sum_{k=1}^M P[E|h_k]P[h_k]}$$

- ***Bayes' theorem*** describes how to update the ***probabilities of hypothesis when given an evidence.***

# An example on Bayes' Theorem applied on diagnosis problem

Consider a simple medical diagnosis problem as follows:

- Two alternative hypothesis:
  - $h$ : the patient has a particular form of disease;
  - $h^c$ : the patient does not have the disease;
- We have prior knowledge that over the entire population of people only 0.8% have this disease
- The available data is from a particular lab test with two possible outcomes: *positive* or *negative*
- Furthermore, the lab test is only an imperfect indicator of the disease
  - The test returns a correct positive result in only 99% of the cases where the disease is present
  - The test returns a correct negative result in only 99.5% of the cases in which the disease is not present

- Define the following events:
  - $h$ : a person has the disease;  $h^c$ : a person does not have the disease;
  - $\oplus$ : the lab test result is positive;  $\ominus$ : the lab test result is negative
- Then, the following probabilities are given:

$$P(h) = 0.008; \quad P(h^c) = 0.992 \text{ (**prior**)}$$

$$P(\oplus | h) = 0.99 \Rightarrow P(\ominus | h) = 0.01 \text{ (**false negative rate**)} \\$$

$$P(\ominus | h^c) = 0.995 \Rightarrow P(\oplus | h^c) = 0.005 \text{ (**false positive rate**)} \\$$

- Suppose we now have a patient for whom the lab test returns a positive result. Should we diagnose the patient as having the disease or not?
- Using the ***Bayes' Theorem***, we can find the posterior probability of  $h$  and  $h^c$ :

$$P(h|\oplus) = \frac{P(\oplus|h)P(h)}{P(\oplus)} = \frac{P(\oplus|h)P(h)}{P(\oplus|h)P(h) + P(\oplus|h^c)P(h^c)}$$

$$= \frac{0.99 \times 0.008}{0.99 \times 0.008 + 0.005 \times 0.992} = 0.6149$$

$$P(h^c|\oplus) = \frac{P(\oplus|h^c)P(h^c)}{P(\oplus)} = \frac{P(\oplus|h^c)P(h^c)}{P(\oplus|h)P(h) + P(\oplus|h^c)P(h^c)}$$

$$= \frac{0.003 \times 0.992}{0.99 \times 0.008 + 0.005 \times 0.992} = 0.3851$$

- Since  $P(h|\oplus) > P(h^c|\oplus)$ , we have  $h^* = h$ , i.e., we should diagnose the patient as ***having the disease***.

- Alternatively, ratio of the posterior probabilities can be calculated as:

$$\frac{P(h|\oplus)}{P(h^c|\oplus)} = \frac{P(\oplus|h)P(h)}{P(\oplus|h^c)P(h^c)} = \frac{0.99 \times 0.008}{0.005 \times 0.992} = 1.6$$

- What is your conclusion based on this ratio of posterior probabilities?

## Maximum a Posterior (MAP) method – A Probability model

- *In many learning scenarios, the learner considers some set of candidate hypothesis  $H$  and is interested in finding the most probable hypothesis  $h \in H$  given the observed data  $\mathcal{D}$ .*
- Any such maximally probable hypothesis is called a ***maximum a posterior (MAP) hypothesis.***

- Under this setting, applying *Bayes' Theorem* gives:

$$\begin{aligned} h_{MAP} &= \underset{h \in H}{argmax} P[h|\mathcal{D}] \\ &= \underset{h \in H}{argmax} \left\{ \frac{P[\mathcal{D}|h]P[h]}{P[\mathcal{D}]} \right\} = \underset{h \in H}{argmax} \{P[\mathcal{D}|h]P[h]\} \end{aligned}$$

- Notice that in the last step, we dropped  $P[\mathcal{D}]$  because it is a constant independent of  $h$ .

- When the hypothesis models are parameterized by the weight vector  $w$ , the problem becomes:

$$\begin{aligned}
 \hat{w}_{MAP} &= \underset{w}{\operatorname{argmax}}\{P[h_w | \mathcal{D}]\} = \underset{w}{\operatorname{argmax}}\{P[\mathcal{D} | h_w]P[h_w]\} \\
 &= \underset{w}{\operatorname{argmax}}\{\ln P[\mathcal{D} | h_w] + \ln P[h_w]\} \\
 &= \underset{w}{\operatorname{argmin}}\{-\ln P[\mathcal{D} | h_w] - \ln P[h_w]\}
 \end{aligned}$$

Where  $P[\mathcal{D} | h_w]$  is the ***likelihood of data***.  $P[h_w]$  is the ***prior*** over the hypothesis model

- Assuming that every hypothesis in  $H$  is ***equally likely***, i.e., every  $h \in H$  has ***equal prior probability***, i.e.,

$$P[h_w] = \text{constant}$$

The ***MAP hypothesis*** becomes the ***maximum likelihood (ML)*** hypothesis:

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P[\mathcal{D}|h] = h_{ML}$$

- Hence, we have,

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ - \sum_{i=1}^N \ln[p(y_i | \mathbf{x}_i, \mathbf{w})] - \ln[p(\mathbf{w})] \right\}$$

- Let's assume that  $Y_i \sim N(h_{\mathbf{w}}(\mathbf{x}_i), \sigma^2)$ , and for the prior  $p(\mathbf{w})$ , we assume that the components  $w_j$  are **i.i.d Gaussian**, i.e.,  $w_j \sim N(w_{j0}, \sigma_h^2)$ . Then, we have,

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{\sum_{i=1}^N (y_i - h_{\mathbf{w}}(\mathbf{x}_i))^2}{2\sigma^2} + \frac{\sum_{j=1}^d (w_j - w_{j0})^2}{2\sigma_h^2} \right\}$$

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - h_{\mathbf{w}}(\mathbf{x}_i))^2 + \frac{\sigma^2}{\sigma_h^2} \left( \sum_{j=1}^d (w_j - w_{j0})^2 \right) \right\}$$

- In the linear regression problem, we have  $h_{\mathbf{w}}(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}$ . When  $w_{j0} = 0$ , the problem becomes:

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \frac{\sigma^2}{\sigma_h^2} \left( \sum_{j=1}^d w_j^2 \right) \right\}$$

- Let  $\lambda = \frac{\sigma^2}{\sigma_h^2}$ , the above becomes:

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \right\}$$

This is just the *Ridge Regression* model!

- MAP is a probabilistic justification for adding the penalized term in Ridge Regression!***
- MAP model is robust to overfitting!***

