

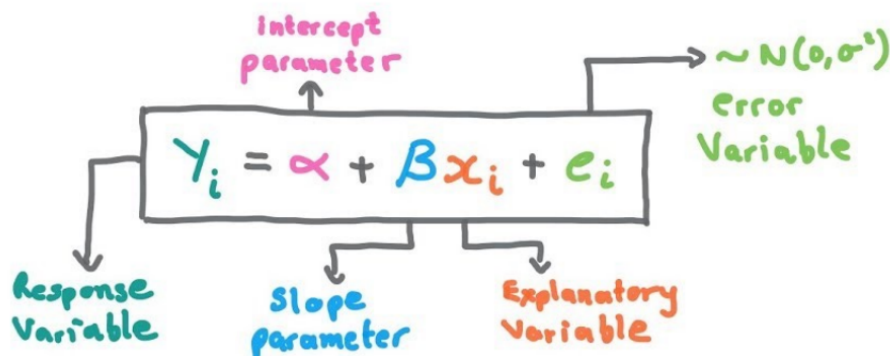


# ch2: Linear Regression

## Regression Model Definition(P4)

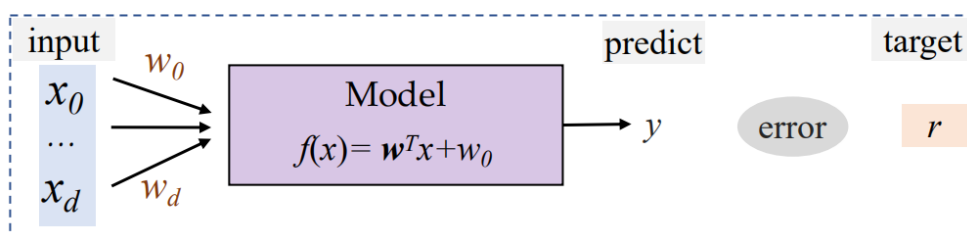
- A regression model provides a function that describes the relationship between one or more **independent variables** and a response, **dependent**, or target variable.

回归模型提供了一个函数，用来描述一个或多个自变量与响应、因变量或目标变量之间的关系



## Model Architecture(P6)

A simple linear function.



- **Train:**
  - estimate the parameters  $w$  and  $w_0$  from data
- **Test:**
  - calculate  $f(x) = w^T x + w_0$ .

## Loss Function(P7)

- For a given input  $x$ , the model outputs a real value  $y$ . Let  $r \in \mathbb{R}$  be target value, the square error is :

$$l(\mathbf{w}, w_0 | x, r) = (r - y)^2$$

- Given:  $D = \{(x^{(1)}, r^{(1)}), \dots, (x^{(N)}, r^{(N)})\}$ , the loss over the dataset is defined as the **mean square error (MSE)**:

$$L(\mathbf{w}, w_0 | D) = \frac{1}{2N} \sum_{\ell=1}^N (r^{(\ell)} - y^{(\ell)})^2$$

### Optimization(P8-9)

Given:  $D = \{(x^{(1)}, r^{(1)}), \dots, (x^{(N)}, r^{(N)})\}$

minimize the loss function using **gradient descend**:

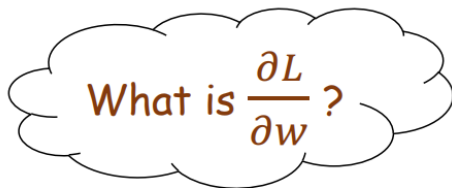
- Goal:

$$\min_{\mathbf{w}} L(\mathbf{w})$$

- Iteration:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \frac{\partial L}{\partial \mathbf{w}}$$

$$L(\mathbf{w}, w_0 | D) = -1/2N \sum_{\ell=1}^N (r^{(\ell)} - y^{(\ell)})^2$$



For each  $\mathbf{w}_j$  ( $j=1, \dots, d$ ):

$$\frac{\partial L}{\partial w_j} = -\frac{1}{N} \sum_{\ell} \underbrace{(r^{(\ell)} - y^{(\ell)})}_{\text{Chain rule}} \frac{\partial y^{(\ell)}}{\partial w_j} = -\frac{1}{N} \sum_{\ell} (r^{(\ell)} - y^{(\ell)}) x^{(\ell)}$$

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \frac{1}{N} \sum_{\ell=1}^N (r^{(\ell)} - y^{(\ell)}) \mathbf{x}^{(\ell)}$$

- Gradient Descend for Liner Regression(P9)(没细看，应该不会考吧)

### Gradient Descend for Liner Regression

**Input:**  $D = \{(\mathbf{x}^{(l)}, r^{(l)})\} \ (l=1:N)$

**for**  $j = 0, \dots, d$

$\mathbf{w}_j \leftarrow \text{rand}(-0.01, 0.01)$

**repeat**

**for**  $j = 0, \dots, d$

$\Delta \mathbf{w}_j \leftarrow 0$

**for**  $l = 1, \dots, N$

$y \leftarrow 0$

**for**  $j = 0, \dots, d$

$y \leftarrow y + \mathbf{w}_j x_j^{(l)}$

$\Delta \mathbf{w}_j \leftarrow \Delta \mathbf{w}_j + (r^{(l)} - y) x_j^{(l)}$

$\Delta \mathbf{w}_j = \Delta \mathbf{w}_j / N$

**for**  $j = 0, \dots, d$

$\mathbf{w}_j \leftarrow \mathbf{w}_j + \eta \Delta \mathbf{w}_j$

**until** convergence

### The Matrix Form (P11)

$$X = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(N)} \end{bmatrix} = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_0^{(N)} & x_1^{(N)} & x_2^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^{(1)} \\ r^{(2)} \\ \vdots \\ r^{(N)} \end{bmatrix}$$

- Prediction:  $\mathbf{y} = X\mathbf{w} = \begin{bmatrix} \mathbf{x}^{(1)}\mathbf{w} \\ \mathbf{x}^{(2)}\mathbf{w} \\ \vdots \\ \mathbf{x}^{(N)}\mathbf{w} \end{bmatrix}$

- Objective:  $L(\mathbf{w}) = \frac{1}{2} (\mathbf{r} - \mathbf{y})^T (\mathbf{r} - \mathbf{y}) = \frac{1}{2} (\mathbf{r} - X\mathbf{w})^T (\mathbf{r} - X\mathbf{w})$

- Gradient

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -X^T(\mathbf{r} - X\mathbf{w})$$

- Solution

loss最小的时候

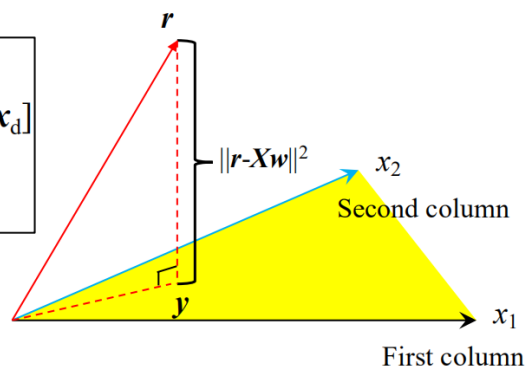
$$\begin{aligned} \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0 &\Rightarrow X^T(\mathbf{r} - X\mathbf{w}) = 0 \\ &\Rightarrow X^T\mathbf{r} = X^TX\mathbf{w} \\ &\Rightarrow \mathbf{w}^* = (X^TX)^{-1}X^T\mathbf{r} \end{aligned}$$

- Then the predicted values are

$$\begin{aligned} \mathbf{y} &= X(X^TX)^{-1}X^T\mathbf{r} \\ &= \mathbf{H}\mathbf{r} \end{aligned}$$

#### Geometrical Explanation

- The column vectors  $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$  form a subspace of  $\mathbb{R}^n$ .
- $\mathbf{H}$  is a least square projection



出现奇异值无法直接计算的时候：

When some column vectors are not independent (e.g.,  $\mathbf{x}_2=3\mathbf{x}_1$ ), then  $\mathbf{X}^T\mathbf{X}$  is singular, thus  $\mathbf{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{r}$  cannot be directly calculated.

**Solution:** Regularization

$$L(\mathbf{w}) = \frac{1}{2} (\mathbf{r} - \mathbf{y})^T (\mathbf{r} - \mathbf{y}) = \frac{1}{2} (\mathbf{r} - \mathbf{X}\mathbf{w})^T (\mathbf{r} - \mathbf{X}\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

New gradient:  $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -\mathbf{X}^T(\mathbf{r} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}$

New optimal solution:  $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0 \Rightarrow -\mathbf{X}^T(\mathbf{r} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w} = 0$

$$\Rightarrow \mathbf{X}^T \mathbf{r} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w}$$

$$\Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{r}$$

