



# Machine Learning

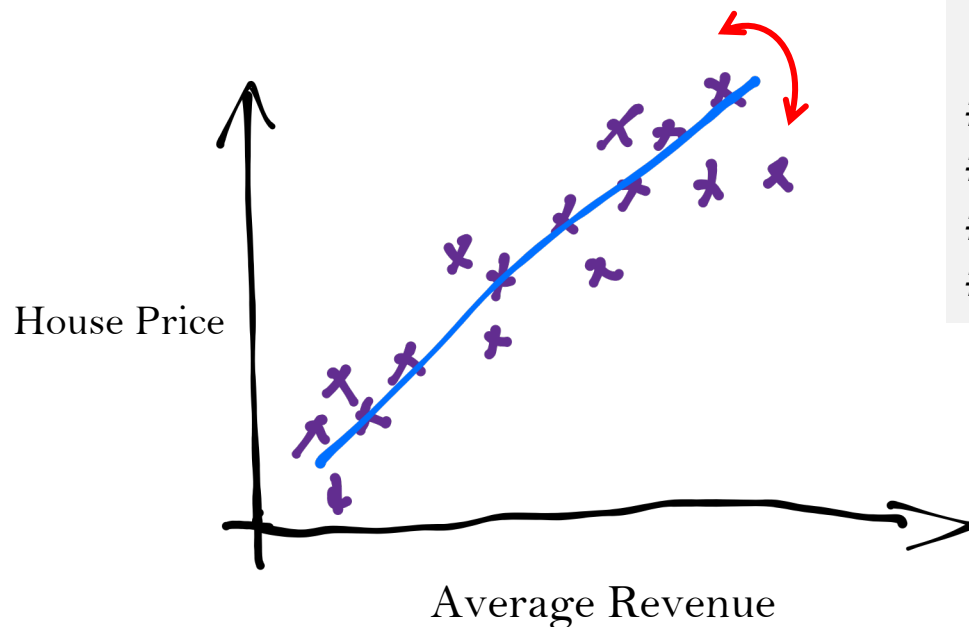
## Chapter 2: Linear Regression

Fall 2022

Instructor: Xiaodong Gu



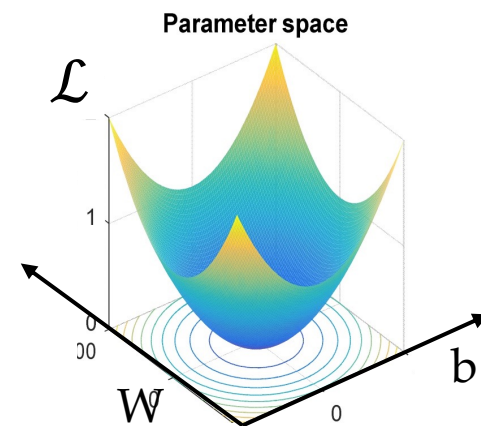
# Recall: Key Elements of Machine Learning



## Elements:

- #1 Data (Experience)
- #2 Model (Hypothesis)
- #3 Loss Function (Objective)
- #4 Optimization (Improve)

$$\theta^* = \operatorname{argmin}_{\theta} \mathcal{L}(\theta|\mathcal{D})$$



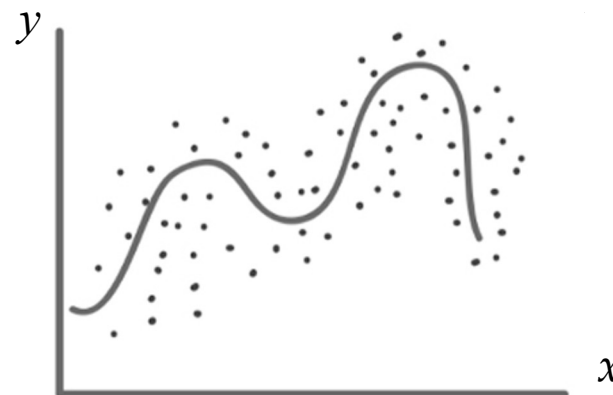


# Regression in Machine Learning

## Machine Learning

- **Supervised Learning**
  - Regression (✓)
  - Classification
  - ...
- Unsupervised Learning
- Reinforcement Learning

| Advertisement | Sales  |
|---------------|--------|
| \$90          | \$1000 |
| \$120         | \$1300 |
| \$150         | \$1800 |
| \$100         | \$1200 |
| \$130         | \$1380 |
| \$200         | ??     |

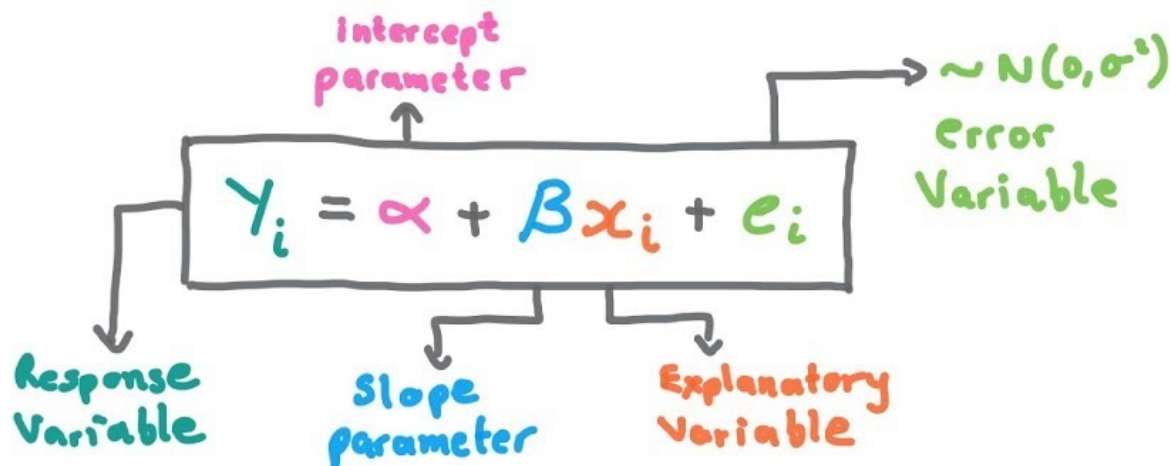


**Regression:** predicts real-valued labels



# Regression Model

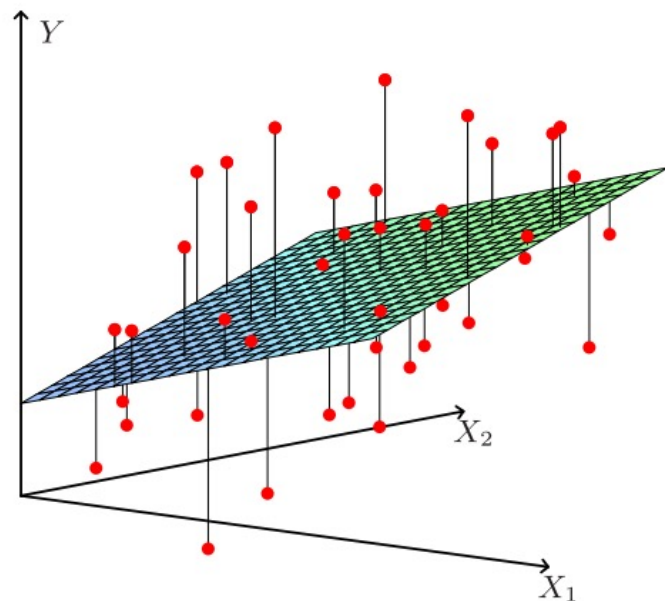
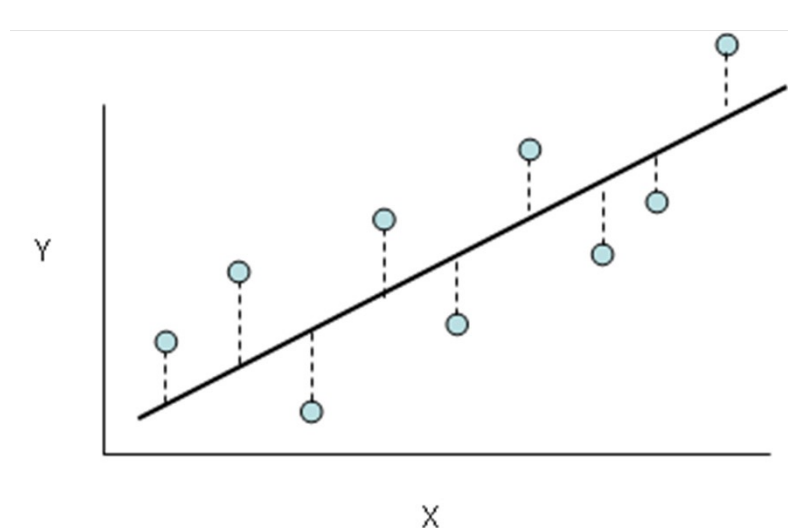
- A regression model provides a function that describes the relationship between one or more **independent variables** and a response, **dependent**, or target variable.



# Linear Regression

A linear function for regression

$$y = f(x) = \mathbf{w}^T \mathbf{x} + w_0$$

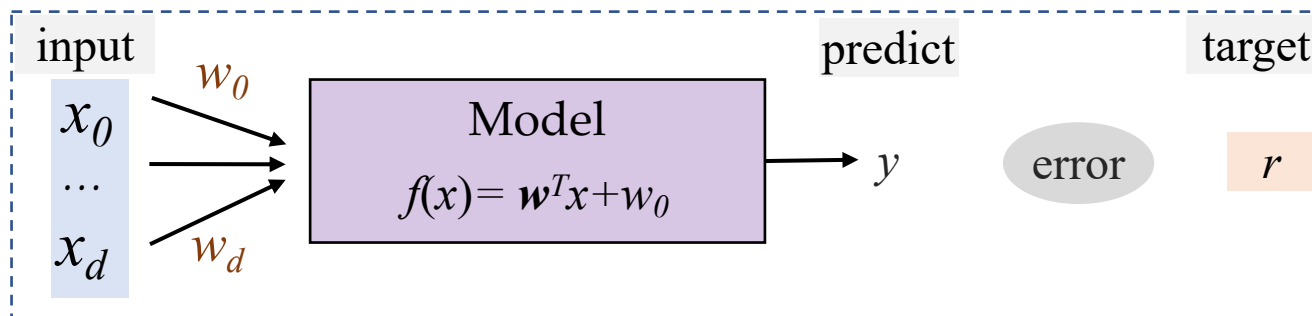


Linear model for regression is a  $(d+1)$ -dimensional hyperplane



# Model Architecture

A simple linear function.



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



# Loss Function

---

- 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)**:

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

这里的2据老师说只是为了之后求导好看一点。。。



# Optimization

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

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

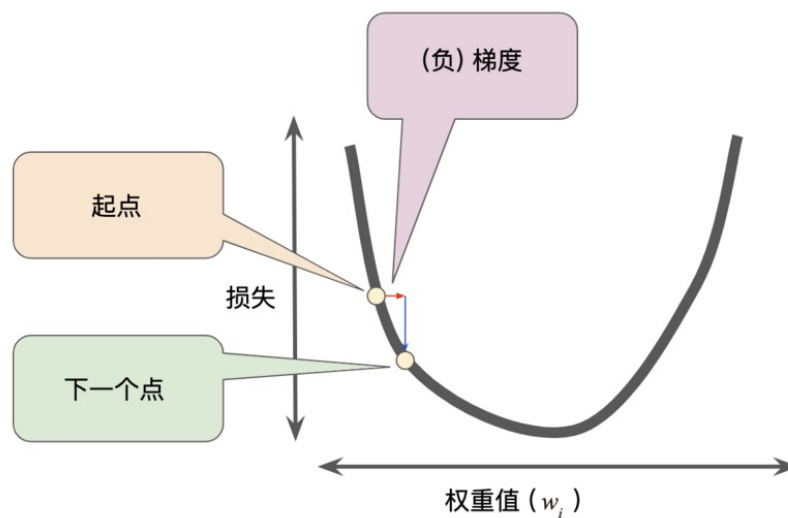
- Goal:

$$\min_w L(w)$$

- Iteration:

$$w_{t+1} = w_t - \eta_t \frac{\partial L}{\partial w}$$

What is  $\frac{\partial L}{\partial w}$ ? 梯度







# Optimization – Gradient Descend

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

What is  $\frac{\partial L}{\partial w}$  ?

For each  $w_j$  ( $j=1, \dots, d$ ): 此处将 $w$ 作为自变量， $y$ 作为因变量，则将 $y=w^*x+w_0$ 两边求导就可以得到下面式子。

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

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



# The Algorithm

## Gradient Descend for Liner Regression

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

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

$w_j \leftarrow \text{rand}(-0.01, 0.01)$

**repeat**

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

$\Delta w_j \leftarrow 0$

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

$y \leftarrow 0$

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

$y \leftarrow y + w_j x_j^{(l)}$

$\Delta w_j \leftarrow \Delta w_j + (r^{(l)} - y) x_j^{(l)}$

$\Delta w_j = \Delta w_j / N$

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

$w_j \leftarrow w_j + \eta \Delta w_j$

**until** convergence



# The Matrix Form

---

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

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



# The Matrix Form

---

- Gradient

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

- Solution

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



# The Matrix Form

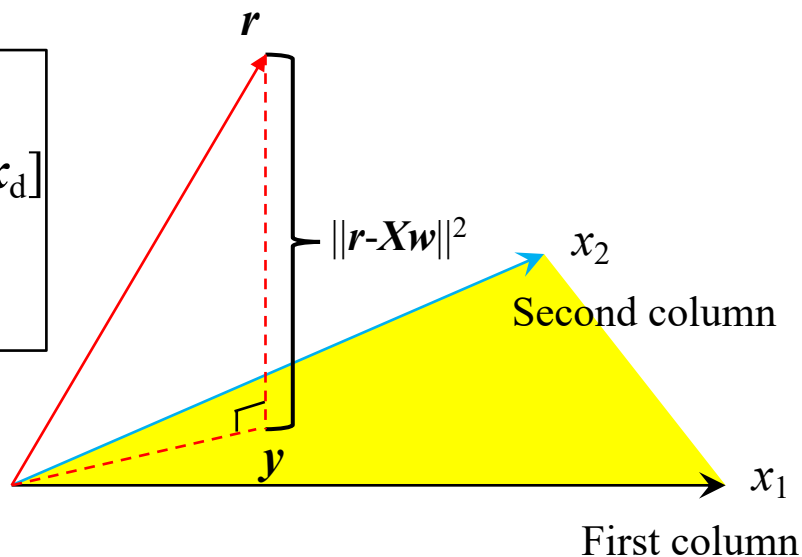
- Then the predicted values are

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

## 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

几何意义：偏差最小时即为x平面与偏差垂直的时候，此时y正好为r在x方向的分量。





# Matrix Form with Regularization

奇异矩阵，行列式为0

**Problem:**  $X^T X$  might be singular

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

**Solution:** Regularization

奇异矩阵-->正则化

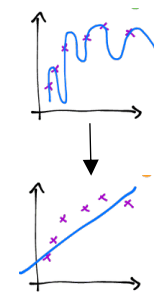
$$L(\mathbf{w}) = \frac{1}{2} (\mathbf{r} - \mathbf{y})^T (\mathbf{r} - \mathbf{y}) = \frac{1}{2} (\mathbf{r} - X\mathbf{w})^T (\mathbf{r} - X\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

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

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

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

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



# Python Tutorial

---



<https://colab.research.google.com/github/cs231n/cs231n.github.io/blob/master/python-colab.ipynb>

<https://www.kaggle.com/code/sudhirn17/linear-regression-tutorial/data?select=insurance.csv>

# What's Next?



## Classifications

Find a decision boundary that **maximizes the margin** between two classes.

### Machine Learning

- **Supervised Learning**
  - Regression
  - Classification(✓)
  - ...
- Unsupervised Learning
- Reinforcement Learning

