

# 机器学习与人工智能

# Machine Learning and Artificial Intelligence

## Lecture 2

## k-Nearest Neighbors & Regression

Yingjie Zhang (张颖婕)

Peking University

[yingjiezhang@gsm.pku.edu.cn](mailto:yingjiezhang@gsm.pku.edu.cn)

2024 Fall

# ML Model Types

Criterion

Whether or not they are trained with human supervision



## Supervised Learning

Fraud detection  
Prediction of stock markets



## Unsupervised Learning

Customer segmentation  
Recommendation

## Semi-supervised Learning

Photo-hosting service  
Speech analysis  
Web-content classification

## Reinforcement Learning

Robotics  
Go games  
Self-driving cars

# Data Division

- **Training dataset**: the sample of data used to fit the model
- **Validation Dataset**: the sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyperparameters.
- **Test Dataset**: a set of examples used only to assess the performance (i.e. generalization) of a fully specified classifier

# Supervised Learning

- k-Nearest Neighbors
- Regressions
- Decision Trees
- Support Vector Machine
- Naïve Bayes
- ...

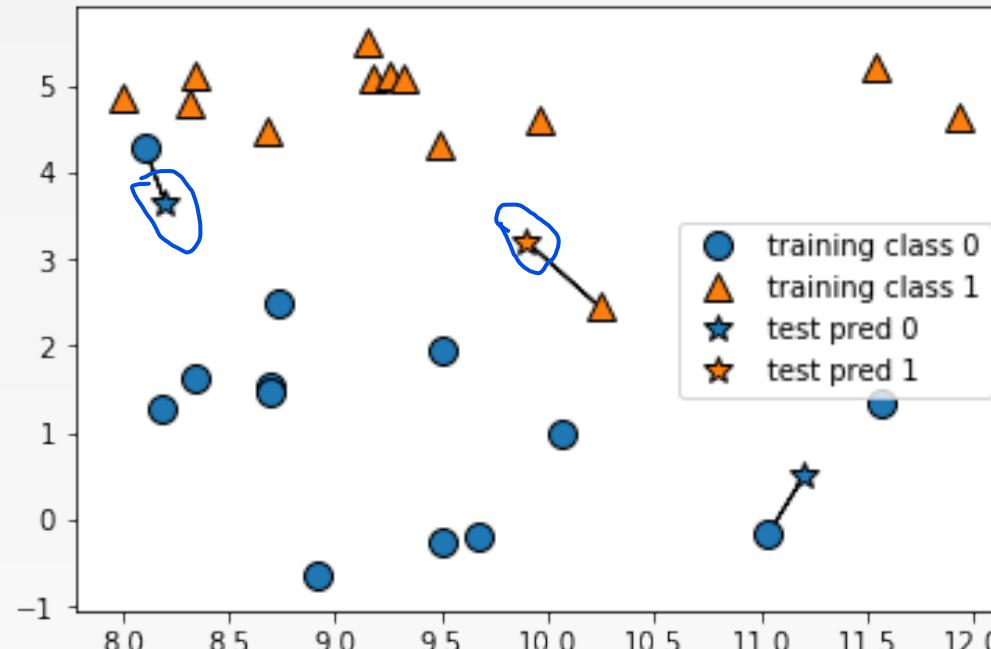
前四种为 diciminate , 判别式的模型 ,  $x-y$

后面为 generated model 模拟数据的来源 , 同时做出判断

# k-Nearest Neighbors

# An Illustrative Example

K : 在预测时的近邻 neighbor  
two features



预测两个五角星的颜色 , 根据最近的data point 颜色预测

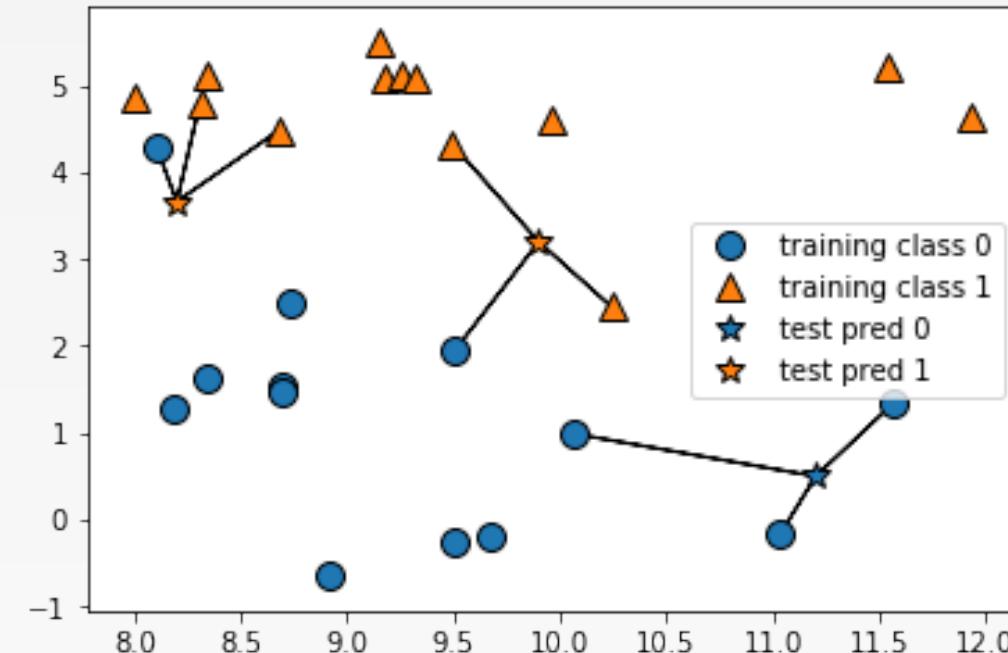
**K = 1**

( X\_train,Y\_train ) -KNN->X\_test=>{y^hat}\_test

K如何影响算大效果

- 1.K=N 很稳定 , 无异质性 , 所有的预测都相同
- 2.K=1 永远取最近的点 可能依赖的信息是outlier , 原始数据错误导致预测错误 , outlier的影响很大 受噪音、波动影响大,忽視整个pattern

K太大 丧失异质性 ; 太小 , 受噪音影响大  
underfitting overfitting



综合k=3 个 data point 的颜色 , 预测 五角星的颜色

**K = 3**

# KNN: Classification Algorithm

- Training: Store all the examples ( $X_{train}, Y_{train}$ )



如果数据点有很多features  
MD计算量更小  
看model中哪个更好

- Prediction:  $X^{(new)}$

- Let  $X^{(1)}, \dots, X^{(k)}$  be the  $k$  most similar examples to  $X^{(new)}$

- Use certain method to determine  $y^{(new)}$  based on  $(y^{(1)}, \dots, y^{(k)})$

category binary  
如何处理?  
jaccard diatance =  $|A \cap B| / |A \cup B|$

## Keys

- A distance metric

X、Y 数值范围不同，如何处理？outlier 不要太影响结果 python boxplot给出数据分布  
1. 标准化  $x - \mu/\sigma$   
2. 归一化  $\min\max x - \min(x)/\max(x)$

惯用 ED Euclidean distance  $d(X^{(j)}, X^{(k)}) = \sqrt{\sum_i (X_i^{(j)} - X_i^{(k)})^2}$

Manhattan distance  $d(X^{(j)}, X^{(k)}) = \sum_i |X_i^{(j)} - X_i^{(k)}|$

cosine diatance



夹角sita

- Value of “K”

如何选到最好的K

Cross validation: larger k? smaller k?

调参：先set K 为某个值，再进行其他操作

K太大 丧失异质性；太小，受噪音影响大  
underfitting overfitting

- Aggregation of the classes of neighbor points

【WiYi, Wi=1/di

Majority vote

# k-NN: Pros and Cons

- **Advantages:**

- Very simple and intuitive
- The cost of the learning process is zero
- No assumption about the characteristics/distributions
- Works on both classification and regression tasks

- **Drawbacks:**

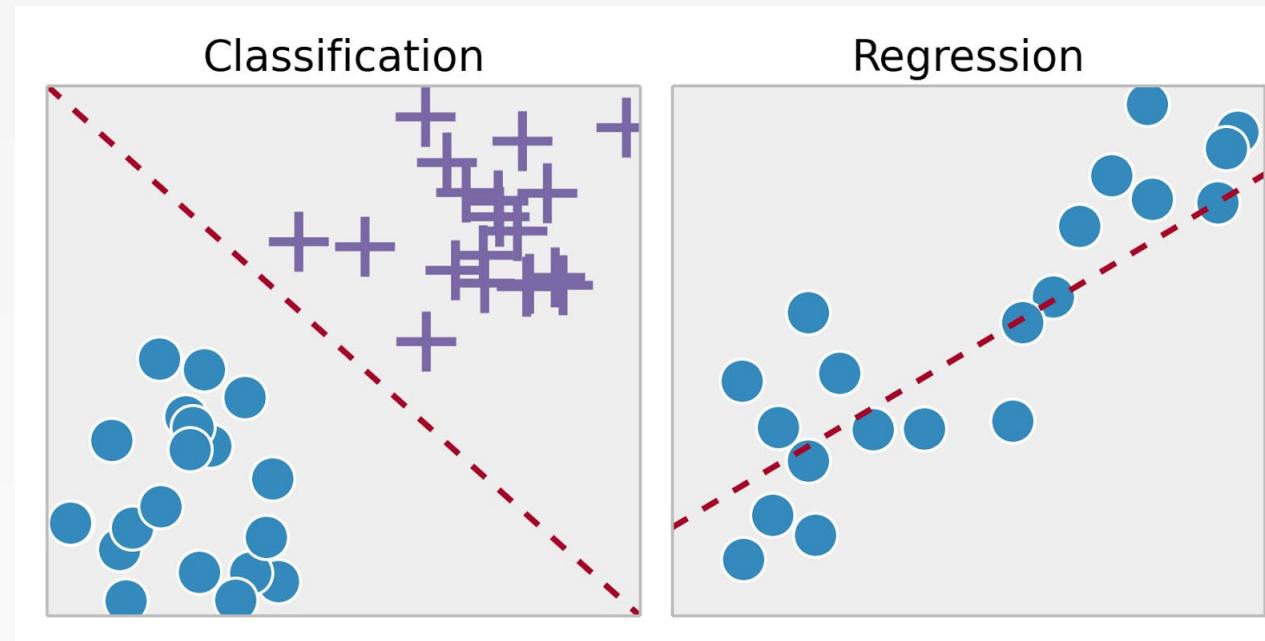
- Computationally expensive when the dataset is very large
  - Need to calculate the compare distance from new example to all other examples
- Sensitive to outliers

计算成本高，无论是行、列增加，都要重新计算，不可复用；  
行、列过多都要遍历，难计算

# Python Review (k-NN)

# Classification vs. Regression

- **Classification:** the goal is to predict a *class label*, which is a choice from a predefined list of possibilities
- **Regression:** the goal is to predict a continuous number, or a *floating-point number (real number)* in programming (math) terms



# Regression

- Given the value of an input  $X$ , the output  $Y$  belongs to the set of real value R.
- Evaluation: predict output accurately
- Examples:
  - Predict housing price
  - Forecast precipitation

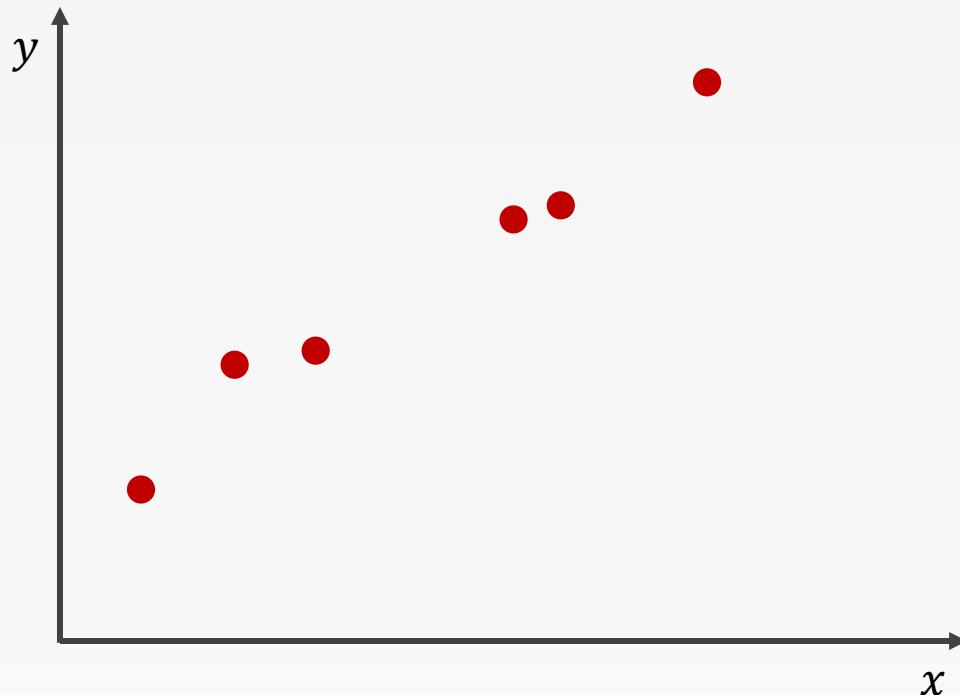
Reg :

1.KNN 前面是分类问题，但也可以做回归  
2.lineal regression : definition、Optimazation OLS、侧重GD梯度下降、  
regularization (讨论overfitting Ridge、LASSQ)  
引申：feature engine、polynomial、logistic reg

# Regression

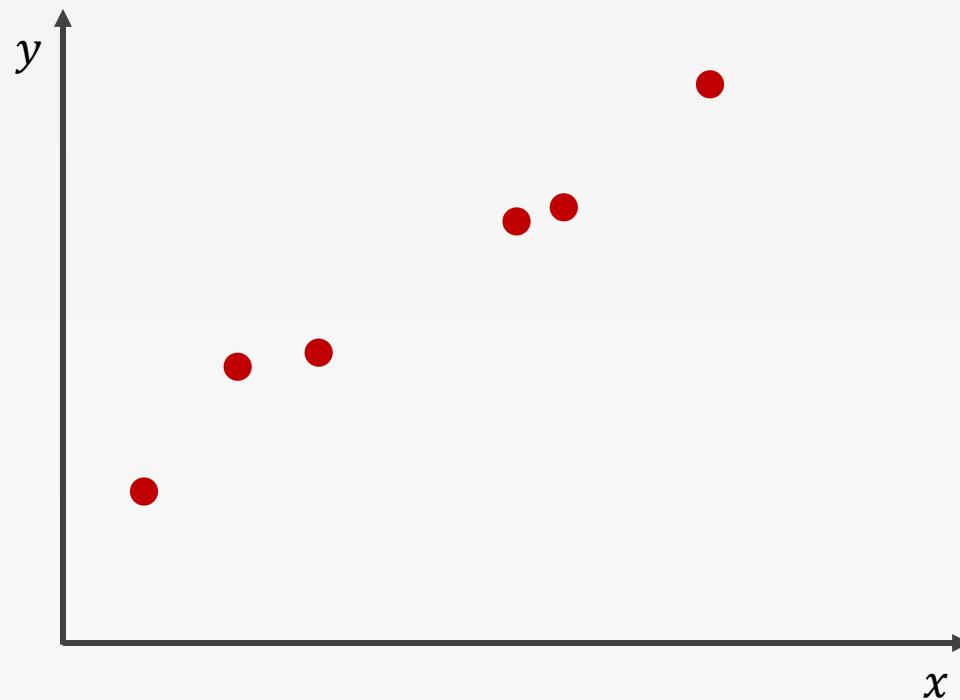
**Example:** Dataset with only one feature  $x$  and one scalar output  $y$

**Q:** What is the function that best fits these points?



# k-NN Regression

**Example:** Dataset with only one feature  $x$  and one scalar output  $y$



KNN regression:  
K个近邻  
(y<sub>1</sub>, y<sub>2</sub>, ..... ) → y<sup>hat</sup>

$$k = 1$$

- *Train:* store all  $(x, y)$  pairs
- *Predict:* pick the nearest  $x$  in the training data and return its  $y$

$k = 2$  Nearest Neighbor Distance Weighted Regression

- *Train:* store all  $(x, y)$  pairs
- *Predict:* pick the nearest two instances  $x^{(n1)}$  and  $x^{(n2)}$  in training data and return the weighted average of their  $y$  values

# Linear Regression

# Linear Regression

- Linear relationship: outcome (dependent) variable is a linear combination of predictor (independent) variables.

$$Y = b + w_1X_1 + w_2X_2 + \cdots + w_nX_n$$

Outcome

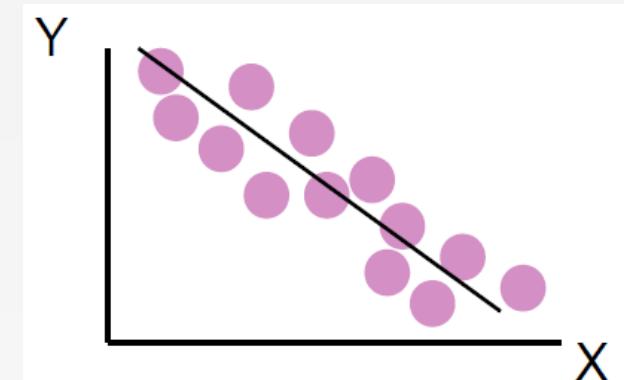
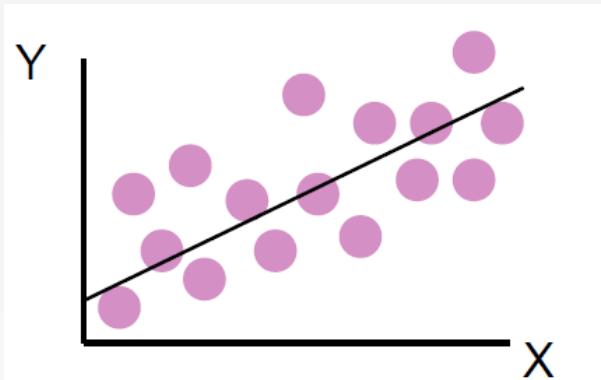
Intercept

Coefficients

Predictors

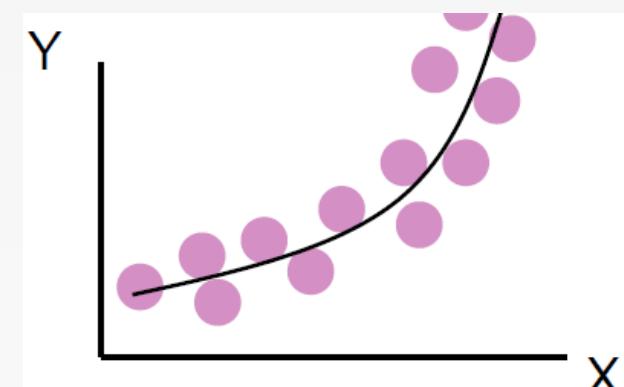
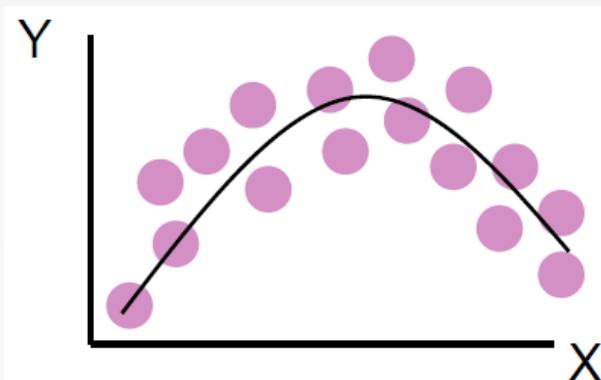
# Linear Regression?

$$Y^{(i)} = b + wX^{(i)}$$



$$Y^{(i)} = b + w_1 X^{(i)} + w_2 (X^{(i)})^2$$

ML中只要求参数线性



A linear model means linear in parameters (not X)!

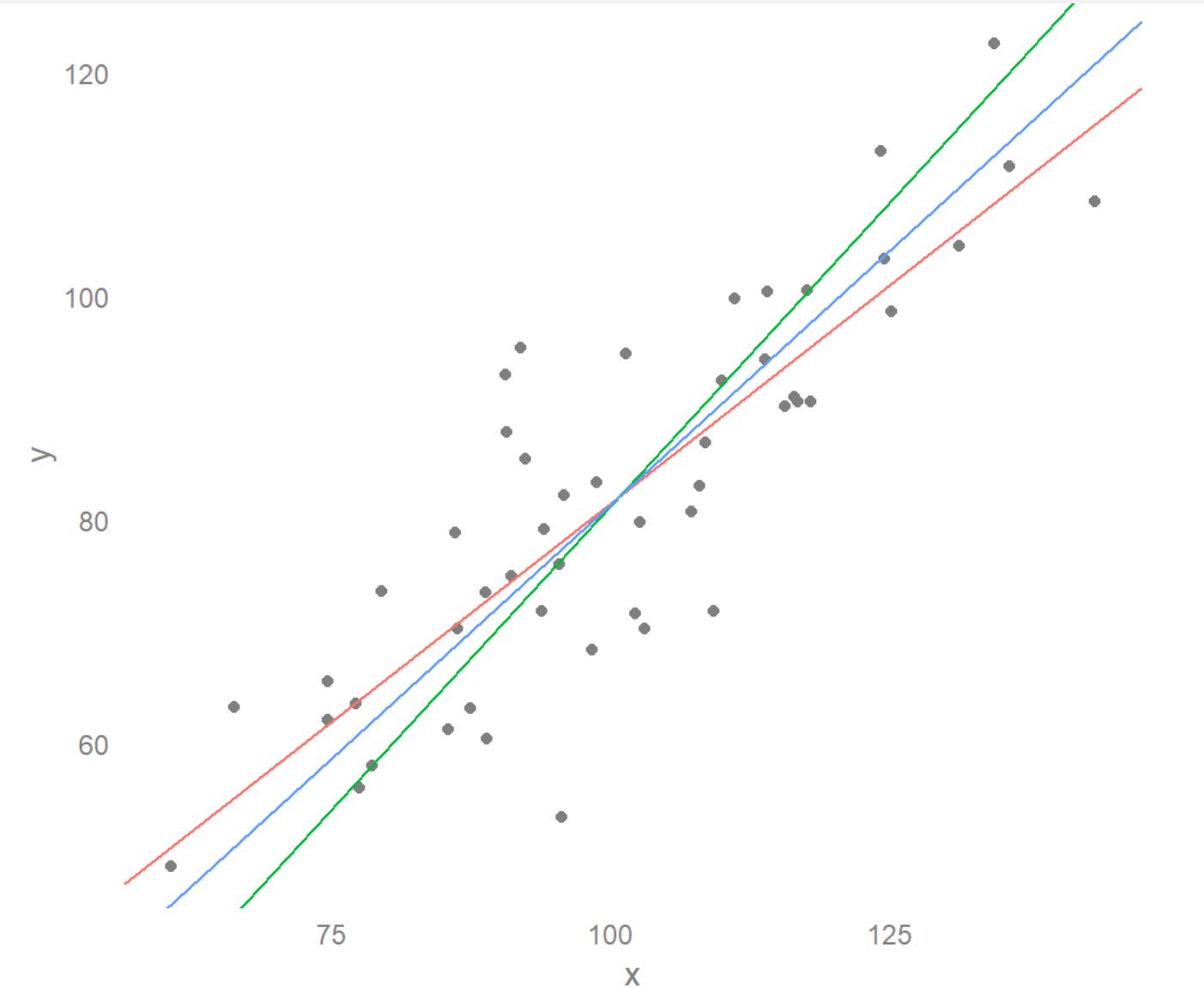
# Linear Regression?

- $y = \sum_j \omega_j f_j(x)$  ✓
- $y = \sum_j \omega_j x_j$  ✓
- $y = \sum_j e^{\textcircled{w}_j} x_j$  ✗
- $y = \sum_j w_j \sin(j^2 x_j)$  ✓

a linear combination of predictor variable.

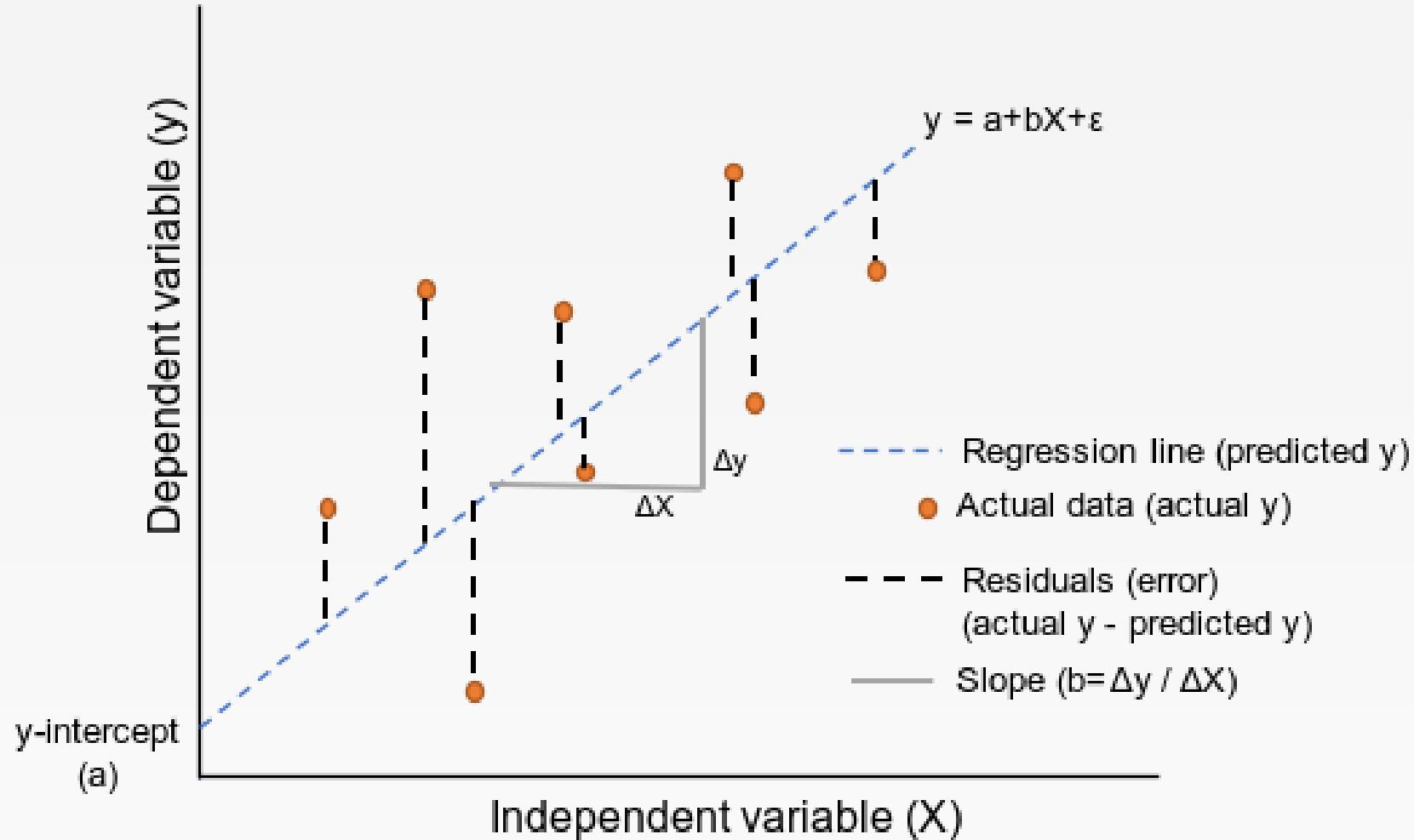
# Training?

to find the parameters.



# Residuals

Residuals  $e = \text{observed } (y) - \text{predicted } (y)$



# Train a Linear Model

- Goal: minimize the Error

$$\text{residual } e = \underset{\text{observed}}{y} - \underset{\text{predicted}}{\hat{y}}.$$

- Potential ways:

- Sum (mean) of absolute errors  $|e_1| + |e_2| + |e_3| + \dots$
- Sum (mean) of squared errors  $e_1^2 + e_2^2 + e_3^2 + \dots$

# Function Approximation

- Objective function: mean squared error (MSE)

$$\begin{aligned} J(\theta) &= \frac{1}{N} \sum_{i=1}^N e_i^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2 \end{aligned}$$

$$x' = [1, x_1, x_2, \dots, x_M]^T$$

$$\theta = [b, w_1, \dots, w_M]^T$$

$$(1, x_1, x_2, \dots, x_M) \cdot \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_M \end{pmatrix} = x^T \theta$$

- Solve the unconstrained optimization problem

- Closed form OLS
- Gradient descent
  - Stochastic gradient descent

$$\min_{\theta} J(\theta) \Rightarrow \hat{\theta}$$

$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta)$$

- Test time

- Given a new  $x$ , make a prediction  $\hat{y} = \hat{\theta}^T x$

# Closed-form Solution (OLS)

最小二乘法

**Criteria** Minimize the sum of squared errors

$$\min \sum_n (y^{(i)} - \hat{y}^{(i)})^2$$

**Solution**  $y_i = b + w_1 x_i + \varepsilon_i$

$w_1$ : slope for the estimated regression equation  $w_1 = \frac{\sum_n (x^{(i)} - \bar{x})(y^{(i)} - \bar{y})}{\sum_n (x^{(i)} - \bar{x})^2}$

$b$ : intercept for the estimated regression equation  $b = \bar{y} - w_1 \bar{x}$

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

$(x, y) \xrightarrow{\text{Train}} \hat{\beta} = (\quad)$   
 $\hat{y}_{\text{new}}$

TPE

$$D = (\tilde{x}^{\text{train}}, y^{\text{train}}) = \{(\tilde{x}^{(1)}, y), \dots, (\tilde{x}^{(m)}, y)\} \quad \xrightarrow{\text{training size}}$$

$$\text{eg } (\tilde{x}^{(1)}, y^{(1)}) = (x_1^{(m)}, \dots, x_D, y) \quad \xrightarrow{\text{feature space}}$$

$y \in \mathbb{R}$  or  $y \in \{-1, 1\}$

$$\text{MLOPs: } C^*: \tilde{x} \rightarrow y \quad y = C^*(\tilde{x})$$

Hypothesis set  $H$ . find  $h^* \in H$  best approximate  $C^*$ ; loss function  $\hat{y} = h^*(\tilde{x})$ ,  $y = C^*(\tilde{x})$

**Gradient Descent** 梯度下降

$$f(\hat{y} - y)$$



# Linear Regression Solutions

- Closed-form solution:

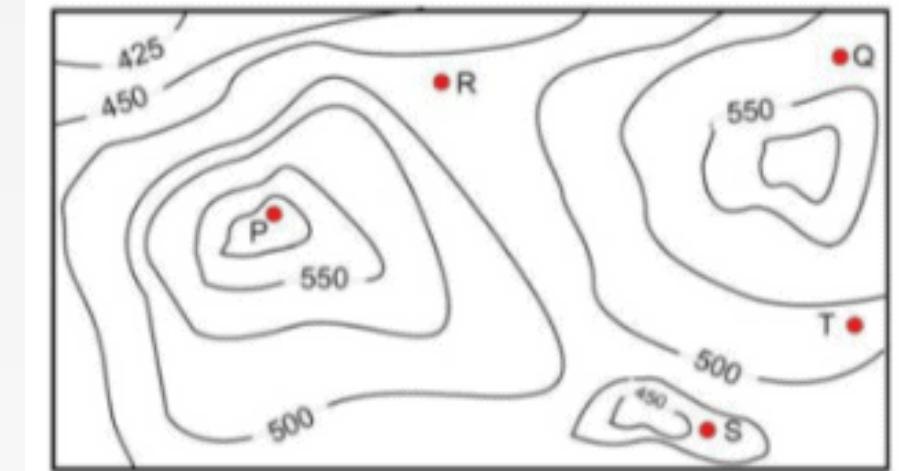
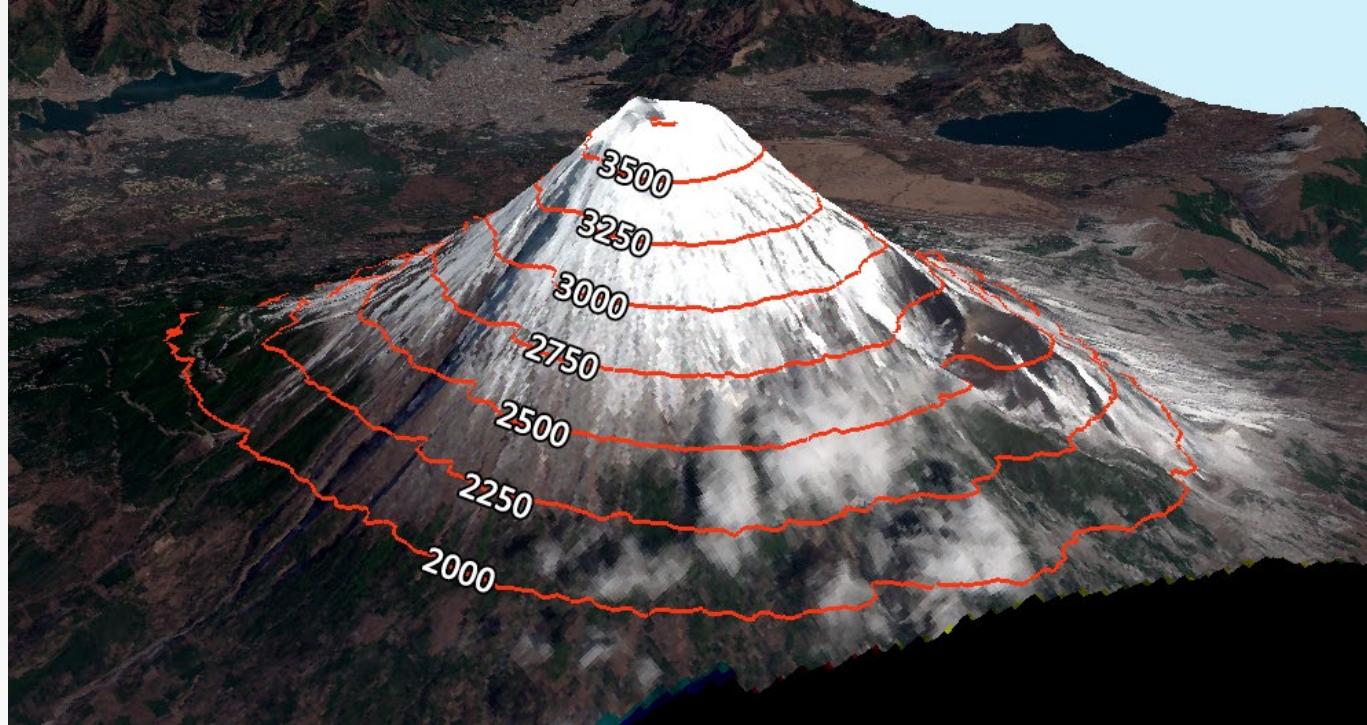
- Computational complexity
- Stability

OLS 算得慢; 不稳定, 不一定有解

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

- Gradient Descent for Linear Regression

# Contour Plots



1. Each level curve labeled with value
2. Value label indicates the value of the function for all points lying on that level curve

# Optimization by Random Guessing

随机

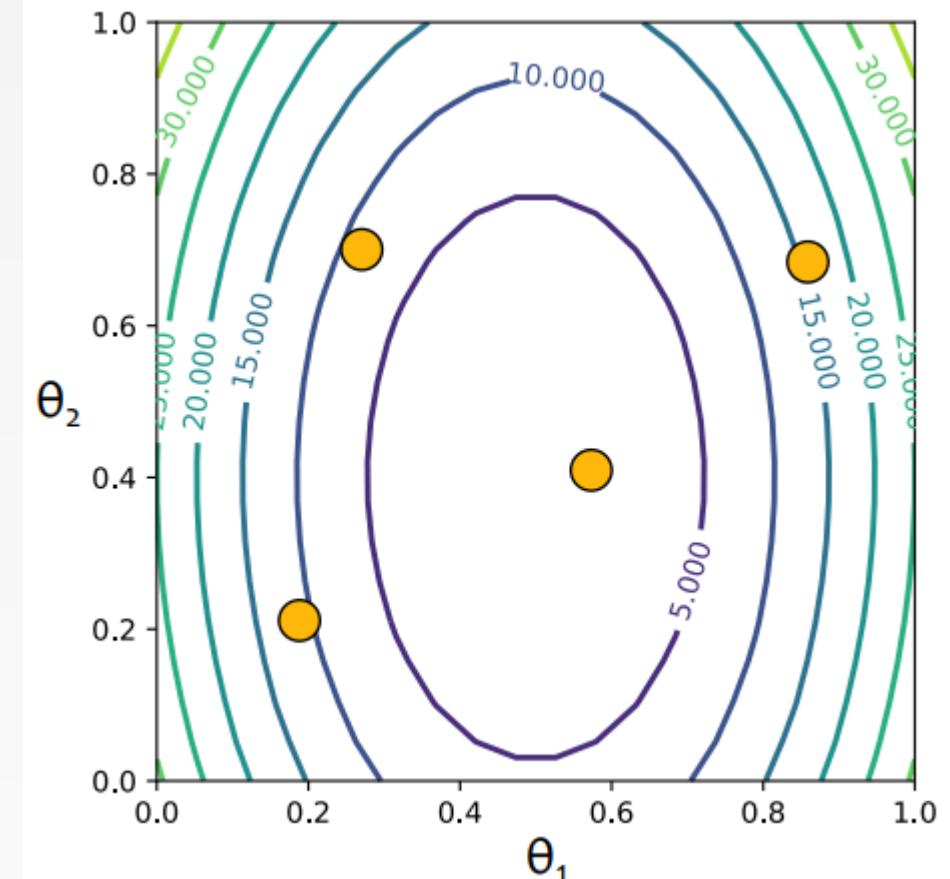
## Random guessing:

1. Pick a random  $\theta$
2. Evaluate  $J(\theta)$
3. Repeat steps 1 and 2 many times
4. Return  $\theta$  that gives smallest  $J(\theta)$

## Linear Regression:

1. Objective function: MSE
2. contour plot: each line labeled with MSE – lower means a better fit
3. **minimum** corresponds to parameters  $(w, b) = (\theta_1, \theta_2)$  that **best fit** some training dataset

$$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2) = (10(\theta_1 - 0.5))^2 + (6(\theta_2 - 0.4))^2$$





最陡的路，速度最快  $\Rightarrow$  梯度

$$\theta^{(j)} \rightarrow \theta^{(j+1)} \rightarrow \dots \theta^*$$

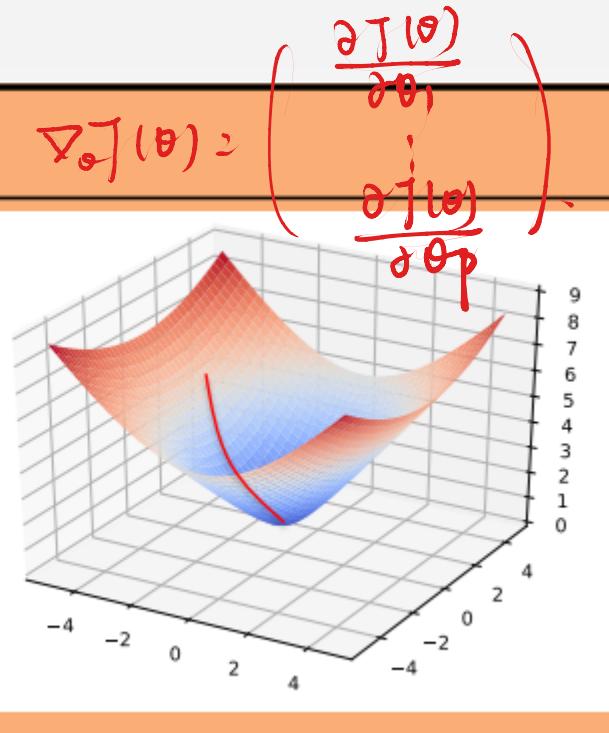
# Algorithm

## Algorithm 1 Gradient Descent

```

1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)} = 0$  or random
3:   while not converged do
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$  调参

```



**Convergence Criteria** (one example):  $\|\nabla_{\theta} J(\theta)\|_2 \leq \epsilon$   $|\theta^{(i+1)} - \theta^{(i)}| \leq \varepsilon; |J(\theta)|$

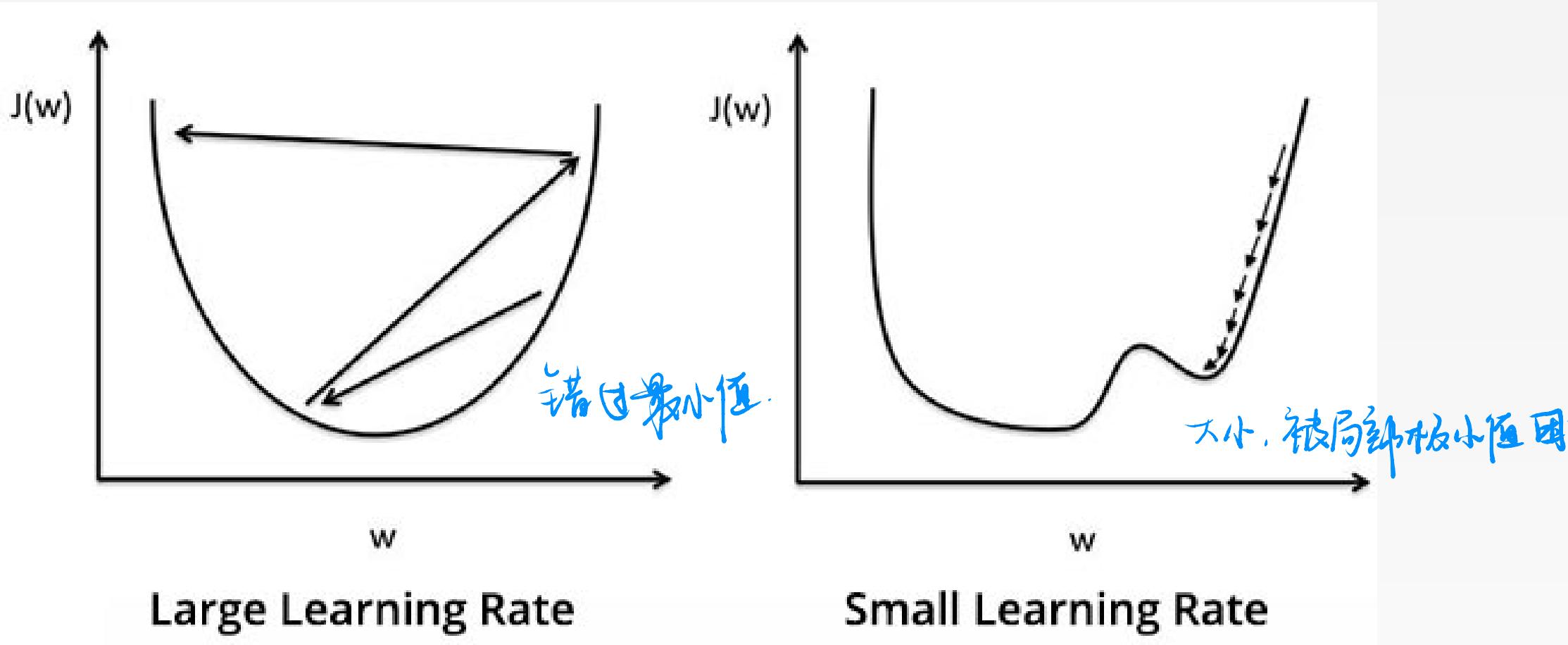
$$\min J(\theta)$$

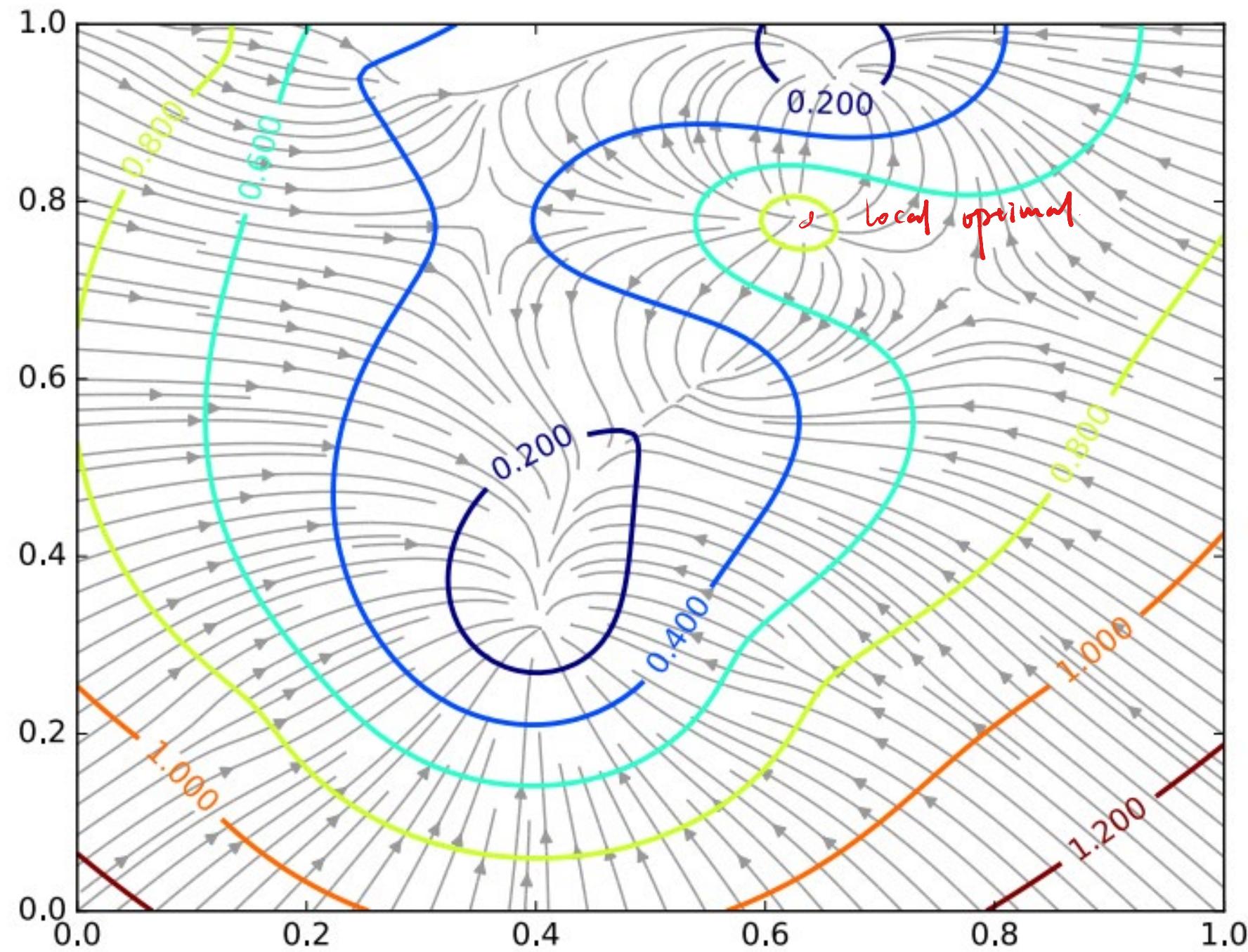
①  $\vec{\theta}^{(0)}$   $\rightarrow$  initial step size  
 ②  $\vec{\theta}^{(1)} = \vec{\theta}^{(0)} - \gamma \nabla J(\vec{\theta}^{(0)})$   
 repeat

# Learning Rate

① stop criteria

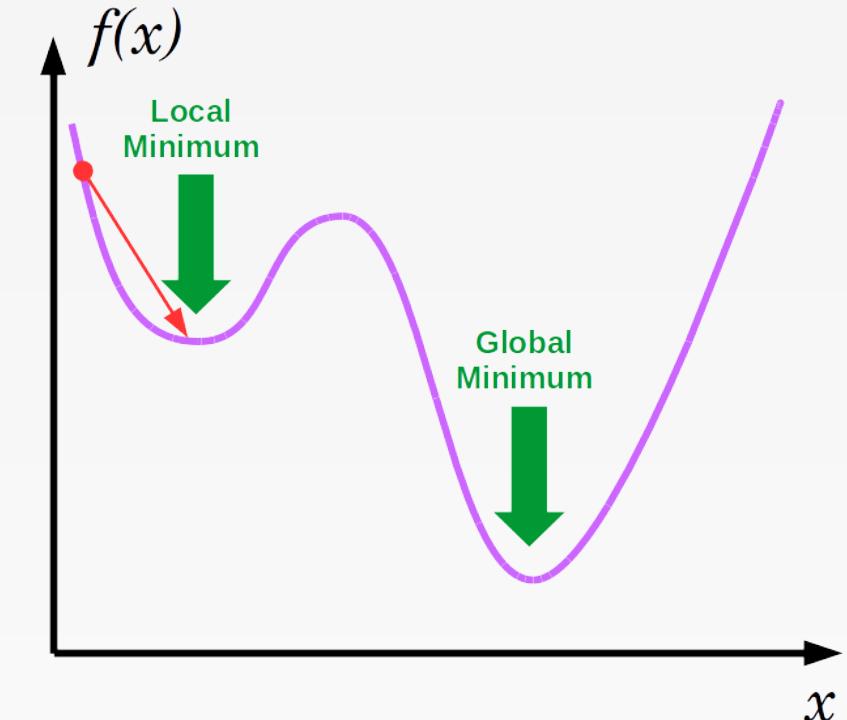
$$\text{eg } \delta^t = \frac{r^o}{t} \text{ 步长随时间减小}$$





# Pros and Cons

- Advantages:
  - Simple and often quite effective on ML tasks
  - Often very scalable  
*可扩展性好。*
- Drawbacks
  - Might find a local minimum  
*假局部*
  - Only applies to smooth function (differentiable)



# GD for Linear Regression

$$J(\theta) = \sum_{i=1}^N (\hat{y}^i - y^i)^2$$

## Algorithm 1 GD for Linear Regression

```

1: procedure GDLR( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$                                  $\triangleright$  Initialize parameters
3:   while not converged do  $\nabla J(\theta)$ ?
4:      $\mathbf{g} \leftarrow \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$        $\triangleright$  Compute gradient
5:      $\theta \leftarrow \theta - \gamma \mathbf{g}$                                  $\triangleright$  Update parameters
6:   return  $\theta$ 

```

# Stochastic Gradient Descent (SGD)

## Algorithm 2 Stochastic Gradient Descent (SGD)

```

1: procedure SGD ( $D, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \sim \text{Uniform}(\{1, 2, 3, \dots, N\})$ 
5:        $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 

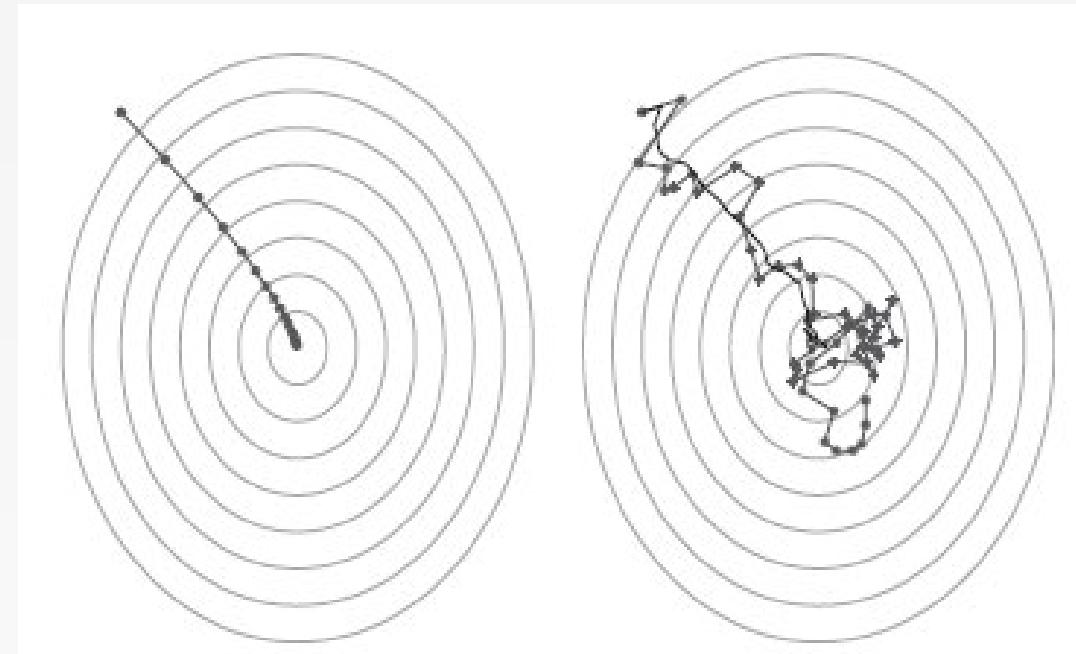
```

minbatch GD  
 }  
 SGD       $\sum_{i=1}^n$   
 算一个 i 然后 GD.  
 训练量显著下降，但精度下降

为什么 SGD 可行？  
 (期望)

# Stochastic Gradient Descent

- Just picks a random instance (or sampling) in the training set to compute the gradient



SGD v

# Mini-Batch SGD

- Gradient Descent:
  - Compute true gradient exactly from all N examples
- Stochastic Gradient Descent (SGD):
  - Approximate true gradient by the gradient of one randomly chosen example
- Mini-Batch SGD:
  - Approximate true gradient by the average gradient of K randomly chosen examples