

## 1. Artificial Intelligence. (wiki).

Intelligence demonstrated by machine.

vs. natural intelligence displayed by humans or animals.

Goals. ①. perception 感知. use inputs from sensors

(e.g., speech recognition. facial recognition) 计算视觉.

②. Reasoning deduction early researchers. 逻辑推理 (let machine think through some complex puzzles).  
inductive systematic search③. Control / motion / manipulation robots. (interact w/ physical world)  
how to move efficiently. (触觉. 视觉. 处理 (拾起. 抬臂)).

④. planning. 智能规划. 导航. 自动驾驶. 车间调度

⑤. learning. computer algorithms that improve automatically through experience  
static. dynamic. overtime.⑥. Natural language process (communication)  
allow machines to read and understand human language.

## 2. Data type. 统计学意义.

Numerical. 具有实际测量意义. (身高. 体重. 面积) Quantitative data.

离散: 不可数 (有限. / 无限) → 描述统计直到100次人头朝上的次数. (100, 10)  
连续: 不可数. 只能用区间表示.

Categorical: 被描述对象的性质. (可用 numerical data 描述. e.g., 1 = F, 0 = M).  
无数学意义  
qualitative data

nominal: 无排序的定性  
ordinal: 有排序的

Boolean - Binary classification - 二分类.

Categorical - multiclass -

ordinal - ordinal -

(train multi binary classification)

$$P(Y=1) = 1 - P(T > 1)$$

$$P(Y=2) = P(T > 1) - P(T > 2)$$

movie rating

real - regression

ordering - ranking.

....



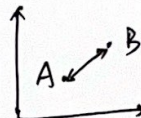
### 3. K-Nearest Neighbor

(1). Example. loan application.

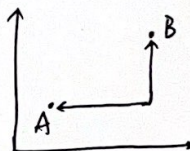
User ID.	$x_1^{(u)}$ age.	$x_2^{(u)}$ loan amount.	$x_3^{(u)}$ income.	$x_4^{(u)}$ purchase Freq.	$y$ default.
1.	25 $x_1^{(1)}$	40K $x_2^{(1)}$	35K	39.	N
2.	35 $x_1^{(2)}$	60K $x_2^{(2)}$	50K	18.	N
3.	23	95K $x_2^{(3)}$	200K	100.	Y.
4.	40.	62K.	170K.	28.	Y.
5.	45	80K	40K	12.	N.

(2). Distance.

(a). Type. (i) Euclidean distance. ( $\sqrt{x^2 + y^2}$ )



(ii) Manhattan distance



# features  $J$ . prefer (ii).

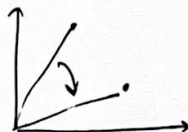
(iii). Hamming distance. (2 binary strings)

- XOR  $a \oplus b = (\neg a \wedge b) \vee (a \wedge \neg b)$   
- count #1.

$$d(11011001, 10011101) = 2$$

$$\oplus = 01000100.$$

(iv). cosine distance



mainly used in collaborative filtering based recommendation systems.

$$D(x, y) = \cos(\theta) = \frac{x \cdot y}{\|x\| \|y\|}$$

(b). special case: categorical data

(i) transform to use Hamming.

(ii) Cosine.

(iii) Python categorical-similarity-measures library

(c). Normalization



身高 + 性别  $\Rightarrow$  Gender

A [179, 42, M]. B [178, 43, M]. C [165, 36, F]. D [177, 42, M]. E [160, 35, F]

Test (167, 43)

$$AF = \sqrt{145} \quad BF = \sqrt{121} \quad CF = \sqrt{53} \quad DF = \sqrt{161} \quad EF = \sqrt{103}$$

$K=3 \Rightarrow C, D, E \Rightarrow$  Female. Not realistic

$$M_i = \max_{j \in \text{train}} x_{j,i} - \min_{j \in \text{train}} x_{j,i}$$

$$d(x_j, x_k) = \sqrt{\sum_i \left( \frac{x_{j,i}}{M_i} - \frac{x_{k,i}}{M_i} \right)^2}$$

sklearn.preprocessing pkg.

(a) standard scalar.  $z = (x - \mu) / s$ .  $\rightarrow$  并不会变成 normal dist. 对于特征分布

(b) min max scale  $z'_i = (z_i - z_i^{\min}) / (z_i^{\max} - z_i^{\min})$  (normalize)

$\rightarrow$  归一化  $\rightarrow [0, 1]$  or  $[-1, 1]$ .  $\rightarrow$  or mean 和极值系数大.

标准归  $\rightarrow$  mean = 0. std. dev = 1 ~~normal dist.~~

对数据范围有严格要求  $\rightarrow$  归一化 (image processing, a typical NN requires 0-1 scale)  
数据不稳定, 存在极端值  $\rightarrow$  标准归  $\rightarrow$  image 0-255 RGB

distance  $\rightarrow$  标准归  
PCA

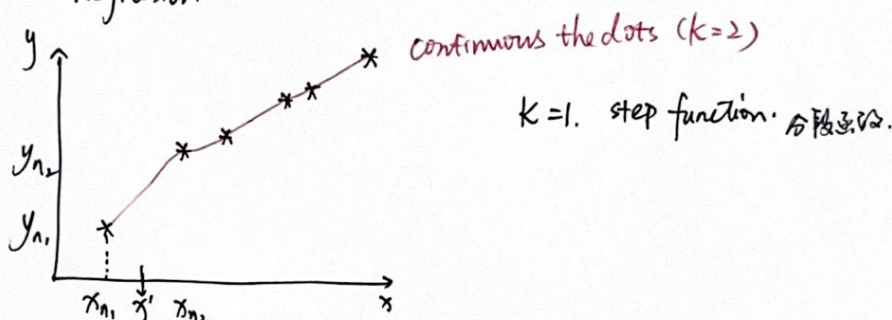
when in doubt

When in doubt, just standardize the data, it shouldn't hurt

## 1. Regression.

Example: stock price prediction (if predict price  $\uparrow$  or  $\downarrow \Rightarrow$  classification).

## 2. KNN. Regression

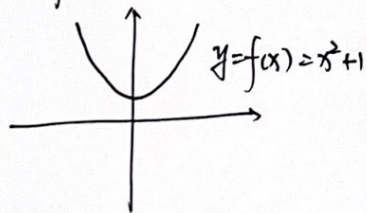


3. Definition  $x \rightarrow$  input. features. independent variables 自变量  
 $y \rightarrow$  output. value. dependent variables. 因变量

## 4. Key ideas of linear regression:

1. find a linear function (parameters  $\vec{w}, b$ )
2. minimize residuals for a training dataset (sum of square)

## 5. Optimization.



$$v^* = \min_x f(x) = 1$$

$$x^* = \operatorname{argmin}_x f(x) = 0$$

6. Data  $D = (x^{(i)}, y^{(i)})_{i=1}^N$

Assume  $x^{(i)} \sim p^*(\cdot)$   
 $y^{(i)} = h^*(x^{(i)})$  } unknown.

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

$$\mathcal{H} = \{h_{\theta} : h_{\theta}(x) = \theta^T x, \theta \in \mathbb{R}^M\}.$$



closed-form solution. 数值解.

① Solve  $\nabla J(\vec{\theta}) = 0$  for  $\vec{\theta}$

MSE: 均方误差

② Test min/max using second derivative.

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \quad X = \begin{bmatrix} x_1^{(1)} & \dots & x_m^{(1)} \\ \vdots & & \vdots \\ x_1^{(N)} & \dots & x_m^{(N)} \end{bmatrix} \begin{matrix} \leftarrow \vec{x}^{(1)} \\ \\ \leftarrow \vec{x}^{(N)} \end{matrix}$$

$$J(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2 \quad \leftarrow \text{ordinary Least Square (OLS)} \quad \text{最小二乘法}$$

$$= \frac{1}{N} \cdot \frac{1}{2} (\vec{X}\vec{\theta} - \vec{y})^T (\vec{X}\vec{\theta} - \vec{y}).$$

② Gradient  $\nabla J(\vec{\theta}) = \vec{X}^T \vec{X} \vec{\theta} - \vec{X}^T \vec{y} = 0.$

③  $\vec{\theta}^{MLE} = (\vec{X}^T \vec{X})^{-1} (\vec{X}^T \vec{y}) = \arg\min_{\vec{\theta}} J(\vec{\theta})$

→ 满秩矩阵时可逆. (当变量数 > 样本数, 即列数多于行数, 不满秩. 需正则化 Regularization)

computational complexity

$$\underbrace{\begin{pmatrix} \vec{X}^T \vec{X} \end{pmatrix}^{-1}}_{M \times M} \underbrace{(\vec{X}^T \vec{y})}_{M \times 1}$$

$$\begin{aligned} \vec{X}^T \vec{X} &: O(M^2 N) \\ (\quad)^{-1} &: O(M^{2.5}) \sim M^3. \\ \vec{X}^T \vec{y} &: O(MN) \\ \frac{(\quad)^{-1} (\quad)}{O(M^2 N + M^{2.5})} & \end{aligned}$$

Linear in # of examples  $N$

Polynomial in # of features  $M$ .

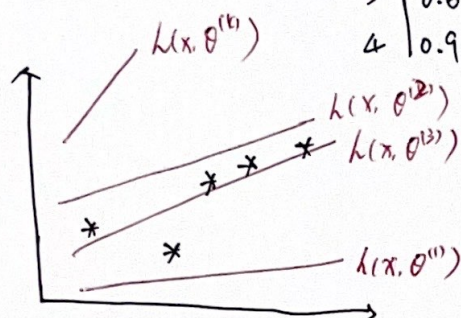


## 7. Gradient Descent.

(1). Random Guessing.

$k$	$\theta_1$	$\theta_2$	$J(\theta_1, \theta_2)$
1	0.2	0.2	10.4
2	0.3	0.7	7.2
3	0.6	0.4	1.0
4	0.9	0.7	19.2

when  $\# \theta \uparrow$ , you cannot directly look at/plot contour plots.



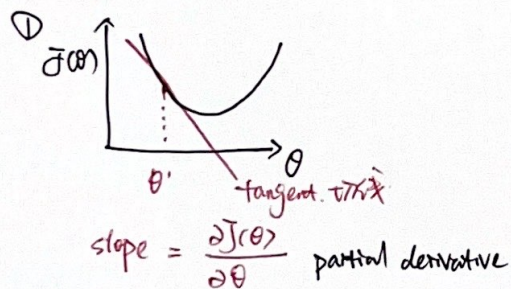
(2). unconstrained optimization

Given function  $J(\vec{\theta})$   $J: \mathbb{R}^m \rightarrow \mathbb{R}$ .

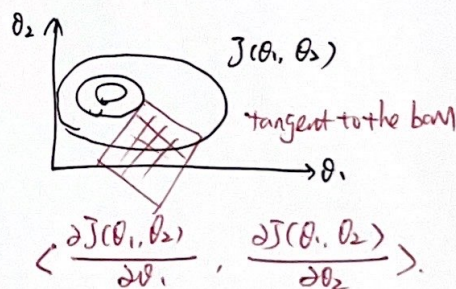
Goal.  $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta)$ .

$\theta$   $\rightarrow$  objective function  
 $\rightarrow$  parameters.

Derivative.



② Tangent plane



vector

Gradient

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J(\vec{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\vec{\theta})}{\partial \theta_m} \end{bmatrix} \quad \text{first order partial derivatives}$$



## (3) Algorithm

① choose an initial point  $\vec{\theta}$ ② Repeat. a, compute gradient  $\vec{g} = \nabla J(\vec{\theta})$ b, choose a step size  $\gamma > 0$  (a real value)c, update  $\vec{\theta} \leftarrow \vec{\theta} - \gamma \vec{g}$  (taking steps on the contour plots)③ return  $\vec{\theta}$  when stopping criterion is met.

## Remarks.

① starting points. a,  $\vec{\theta} = 0$  b  $\vec{\theta}$  randomly② stopping.  $\|\nabla J(\vec{\theta})\|_2 < \varepsilon$   $\varepsilon = 10^{-8}$ 

$$\|x\|_2 = \sqrt{\sum_{k=1}^n |x_k|^2}$$

 $L_2$  norm. 2-norm③ step size. a, fixed value  $\gamma = 0.1$ 

b, exact line search 线搜索

c, backtracking line search 求得一个包含理想的步长的区间. 2分法

d, schedule  $\gamma_t = \gamma_0 / (t-1)\gamma_0 + 1$  (gradually shrink the step)

## (4) Gradient for Linear Regression.

$$\text{MSE } J(\theta) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta) \quad \text{where } J^{(i)}(\theta) = \frac{1}{2} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})^2$$

→ doesn't affect argmin. just for mathematical convenience

$$\frac{\partial J^{(i)}(\theta)}{\partial \theta_j} = \frac{1}{2} \cdot 2 (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \frac{\partial}{\partial \theta_j} (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})$$

$$= (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \frac{\partial}{\partial \theta_j} (y^{(i)} - \sum_{m=1}^M \theta_m x_m^{(i)})$$

$$= (y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) x_j^{(i)}$$

$$\nabla J^{(i)}(\theta) = \begin{bmatrix} \frac{\partial J^{(i)}(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J^{(i)}(\theta)}{\partial \theta_M} \end{bmatrix} = \underbrace{-(y^{(i)} - \vec{\theta}^T \vec{x}^{(i)})}_{\text{scalar}} \underbrace{\vec{x}^{(i)}}_{\text{vector}}$$

$$\nabla J(\theta) = \nabla \left( \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta) \right) = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\theta) = \frac{1}{N} \sum_{i=1}^N -(y^{(i)} - \vec{\theta}^T \vec{x}^{(i)}) \vec{x}^{(i)}$$

if error, you would put a larger weight on the specific data point



(t) why SGD works.

Expectation of Gradients.

$$I \sim \text{Uniform}(\{1, 2, \dots, N\})$$

$$\mathbb{E}_I[\nabla J_I(\vec{\theta})] = \sum_{i=1}^N P(I=i) \nabla J_i(\vec{\theta})$$

$$= \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta}) = \nabla J(\vec{\theta}).$$

Mini-batch  
批梯度下降

8. MLE.

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} P(D|\theta) \quad \left. \begin{array}{l} \log \text{ is monotonic} \end{array} \right\}$$

$$= \underset{\theta}{\operatorname{argmax}} \log P(D|\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \ell(\theta) \quad \text{where } \ell(\theta) \triangleq \log P(D|\theta)$$

log-likelihood.

9. Logistic Regression.

Odds ratio 赔率 从数据看, 一个特定事件比另一个发生的  
可能性大小.

(1) Binary

$$P(y|\vec{x}) = \begin{cases} \sigma(\vec{\theta}^T \vec{x}) & \text{if } y=1 \\ 1 - \sigma(\vec{\theta}^T \vec{x}) & \text{if } y=0 \end{cases}$$

$$\sigma(u) = \frac{1}{1 + \exp(-u)},$$

$$y \sim \text{Bernoulli}(\sigma)$$

objective

$$\ell(\vec{\theta}) = \sum_{i=1}^N \log P(y^{(i)} | x^{(i)}, \vec{\theta}) \quad \leftarrow P(y^{(i)} | x^{(i)}).$$

$$J(\vec{\theta}) = -\frac{1}{N} \ell(\vec{\theta}) = -\frac{1}{N} \sum_{i=1}^N \log P(y^{(i)} | x^{(i)}, \vec{\theta}).$$

$J^{(i)}(\vec{\theta})$  for SGD.

(2) Multinomial  $y \in \{P, S, W\}$

$$P(y=P|\vec{x}) = \exp(\vec{\theta}_P \cdot \vec{x}) / Z(\vec{x}, \vec{\theta}) \quad Z(\vec{x}, \vec{\theta}) = \exp(\vec{\theta}_P \cdot \vec{x}) + \exp(\vec{\theta}_S \cdot \vec{x}) + \exp(\vec{\theta}_W \cdot \vec{x})$$



## Derivatives (Logistic Regression)

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial \theta_m} &= \frac{\partial}{\partial \theta_m} \left( -\log P(y^{(i)} | x^{(i)}, \theta) \right) \\ &= \begin{cases} \frac{\partial}{\partial \theta_m} -\log [\sigma(\theta^T x^{(i)})] & \text{if } y^{(i)} = 1 \\ \frac{\partial}{\partial \theta_m} -\log [1 - \sigma(\theta^T x^{(i)})] & \text{if } y^{(i)} = 0 \end{cases} \end{aligned}$$

= ...

$$= - \underbrace{\left( \underbrace{y^{(i)}}_{\text{truth}} - \underbrace{\sigma(\theta^T \vec{x}^{(i)})}_{\text{prob of } y=1} \right)}_{\text{scalar}} \underbrace{x_m^{(i)}}_{m^{\text{th}} \text{ feature}}$$