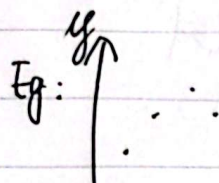


Regression 回归

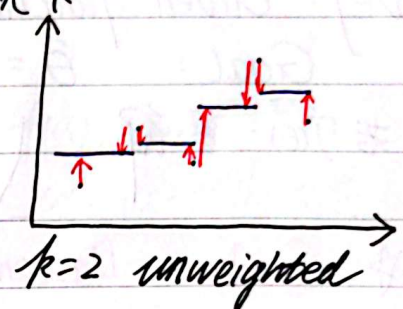
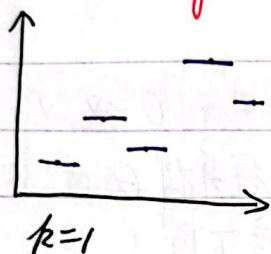
Goal: Given $\{(\vec{x}, y)\}$ dataset, learn a function:
 $y' = h(\vec{x})$, 应与 y 接近

Eg:  找函数拟合各点曲线

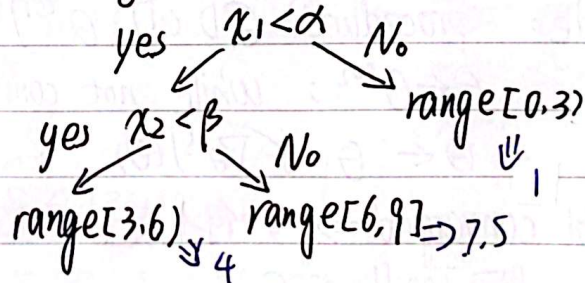
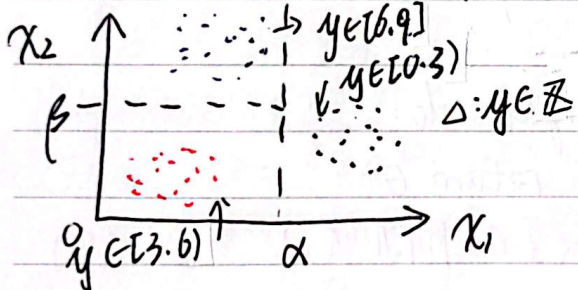
Method ①: k -NN 回归

$k=1$: pick the nearest x , return its y

or: $k=2$: weighted / unweighted 效果见下



Method ②: Decision Tree Regression



划分 x_i ; 然后 $\{y_i\}$ 聚为一类, 看 mean (Decision Tree 分支到最后时)

Method ③: Linear Functions, Residuals & MSE

Δ Linear functions \neq Linear decision Boundaries

Def: Regression is predicting real-valued outputs

$$D = \{(\vec{x}^{(i)}, y^{(i)})\}_{i=1}^n, \quad \vec{x}^{(i)} \in \mathbb{R}^M, y^{(i)} \in \mathbb{R}$$

Eg: $y = w^T x + b$

Eg: $y = \text{sign}(w^T x + b)$



目标

Key idea: Find the linear function h (with w, b) that minimizes the squares of residual for training set.

Residual: $e_i = |y^{(i)} - h(w^T x^{(i)} + b)|$ 不是点向直线
作垂线

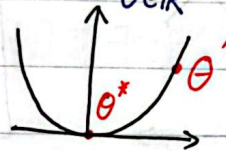
$$\text{MSE: Loss} = \frac{1}{n} \sum_{i=1}^n e_i^2$$

有MSE有什么用? The big picture: Optimization for ML

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

Goal: $\hat{\theta} = \underset{\theta \in \mathbb{R}^m}{\operatorname{argmin}} J(\theta)$

形象理解:



在 θ' 如何一步步到 θ^* 呢?

Naive one: 随机采样 θ

★ Gradient Descent: $\nabla J(\vec{\theta}) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_m} \end{bmatrix}$ 梯度下降!

$\vec{\theta} \leftarrow \vec{\theta} - \delta t \nabla J(\vec{\theta})$, where δt is step size, $\in \mathbb{R}$

Pseudocode: procedure GD($J, \theta^{(0)}$):

$\theta \leftarrow \theta^{(0)}$; while not converged do:

$\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$; return θ

那么 not converged \Rightarrow 如何判定 θ 已收敛至 θ^* 附近呢?

$$\Delta: \|\nabla_{\theta} J(\theta)\|_2 \leq \epsilon$$

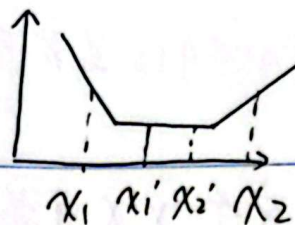
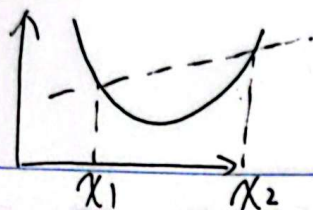
而 MSE Loss: $J^{(i)}(\theta) = (y^{(i)} - \theta^T x^{(i)})^2$, $J(\theta) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta)$

$$\text{则 } \frac{d}{d\theta_k} J^{(i)}(\theta) = 2(\theta^T x^{(i)} - y^{(i)}) \frac{d}{d\theta_k} (\theta^T x^{(i)} - y^{(i)})$$

$$= 2(\theta^T x^{(i)} - y^{(i)}) \frac{d}{d\theta_k} \left(\sum_{j=1}^K \theta_j x_j^{(i)} - y^{(i)} \right) = 2(\theta^T x^{(i)} - y^{(i)}) x_k^{(i)}$$

$$\therefore \frac{d}{d\theta_k} J(\theta) = \sum_{i=1}^N (\theta^T x^{(i)} - y^{(i)}) x_k^{(i)}, \quad \nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \vdots \\ \frac{d}{d\theta_m} J(\theta) \end{bmatrix} = \sum_{i=1}^N (\theta^T x^{(i)} - y^{(i)}) x^{(i)}$$





No.

Date

Optimization for ML:

Convexity: A function $f: \mathbb{R}^D \rightarrow \mathbb{R}$, convex if:
 $\forall x^{(1)} \in \mathbb{R}^D, x^{(2)} \in \mathbb{R}^D, c \in [0, 1]$, 有:

$$f(cx^{(1)} + (1-c)x^{(2)}) \leq cf(x^{(1)}) + (1-c)f(x^{(2)})$$

上述为 not strictly convex^{*}; 下述为 strictly^{*}: 如 fig 1

$c \in (0, 1)$, 有: $f(cx^{(1)} + (1-c)x^{(2)}) < cf(x^{(1)}) + (1-c)f(x^{(2)})$

^{*}: 如 fig 2, 取 x_1', x_2' , 上面这行实现不了!

对于凸函数来说: global min: $f(x^*) \leq f(x), \forall x \in \mathbb{R}^D$

local min: $\exists \varepsilon$ s.t. $f(x^*) \leq f(x)$ 且 s.t. $\forall x, \|x - x^*\|_2 < \varepsilon$

Closed form Optimization:

key: 找到 θ^* , s.t. $\nabla J(\theta^*) = 0$, 其中 J 凸

若无法确保 J 凸, 则找 $\{\theta_i\}$, s.t. $\nabla J(\theta_i) = 0$ for $i \in \text{range}$

并确认它们是极大 or 小值

Given $D = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$, 令: $X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_D^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times D+1}$

而 $y = [y^{(1)}, \dots, y^{(N)}] \in \mathbb{R}^N$ 是目标:

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

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

$$\nabla_{\theta} J(\theta)^* = \frac{1}{2N} (2X^T X \theta - X^T y) = 0 \quad \therefore X^T X \hat{\theta} = X^T y$$

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

^{*}: 矩阵求导!!

^{*}: $\frac{1}{2}$ 无伤大雅; $J(\theta)$ 含义为: 最小 Mean Square Error!!

Δ : 原: $w^T x + b$, 令 $x = \begin{bmatrix} 1 \\ x \end{bmatrix}$, $\theta = \begin{bmatrix} b \\ w \end{bmatrix}$, 则 $\Rightarrow \theta^T x$

KOKUYO



扫描全能王 创建

No. *: 直觉上: $N=2$, 那么超平面显然有很多种划分

Date

$\hat{\theta} = (X^T X)^{-1} X^T y$, 这个 $X^T X$ 可逆么? 一般认为 $N \gg D+1$ 时, 几乎总是 full-rank 的*. 但是关键在于计算 $X^T X$ 逆时间复杂度多少? 算 $X^T X: O(ND^2)$; 算它的逆: $O(D^3)$.

那么, 有没有不那么精确但复杂度能够接受的近似求法呢?

GD & SGD:

Gradient Descent:

$\theta \leftarrow \theta^{(0)}$
while not converged do

$\theta \leftarrow \theta - \gamma \nabla J(\theta)$

return θ

N 项, $\theta \in \mathbb{R}^M$

而 Stochastic Gradient Descent:

$\theta \leftarrow \theta^{(0)}$

while not converged do

$i \sim \text{uniform}(\{1, 2, \dots, N\})$

$\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$

return θ

$$\Delta: J(\theta) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta)$$

实质: 原来用全部样本,

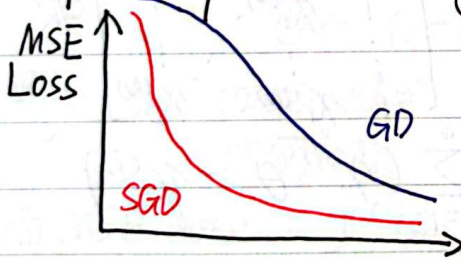
现在一步用一个

why SGD 能 work? 因为 $E[\nabla_{\theta} J^{(i)}(\theta)] = \sum_{i=1}^N p(x^{(i)}, y^{(i)}) \nabla_{\theta} J^{(i)}(\theta)$
 $= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} J^{(i)}(\theta) = \nabla_{\theta} J(\theta)$



黑: GD
蓝: SGD

Def: An epoch is a single passing through training data



GD: One update per epoch

SGD: N update per epoch

step to convergence

Computation per step

Theoretical Comparison:

GD

$O(\log 1/\epsilon)$

$O(NM)$

SGD

$O(1/\epsilon)$

$O(M)$

理论上 SGD 收敛更慢, 但实践中更快

那么在上-面的 Linear Regression 情景中, $J(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$

则 SGD:

$\theta \leftarrow \theta^{(0)}$
while not converged do:

为了保证 N 次可过一个 epoch

故是 range $(1, N)$ shuffle

以模拟 Uniform.

GD: $\theta \leftarrow \theta - \gamma \frac{1}{N} \sum_{i=1}^N (\theta^T x^{(i)} - y^{(i)}) x^{(i)}$
Campus $x^{(i)}$
return θ



扫描全能王 创建