

L3

2. Recap: Linear Regression

Today:

① Gradient descent

② Application: Linear Regression (again)

Fit a linear model:

$$y = w_0 + w_1 x$$

- loss function: $MSE(w) = \frac{1}{n} \sum [y_i - (w_0 + w_1 x_i)]^2$

$$w_0 = \bar{y} - w_1 \bar{x}$$

$$w_1 = \frac{\sum xy}{\sum x^2}$$

Univariate 单变量
Case

Multivariate:

$$X = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{pmatrix} \approx y$$

Linear model: $y = w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)}$

\uparrow intercept \uparrow coefficients

Assumption: Data & label are zero mean. ignore w_0 . $y = \sum_{i=1}^d w_i x_i = \langle w, x \rangle$

Step 1

$$y = \langle w, x \rangle$$

Step 2

$$MSE = \frac{1}{2} \sum_{i=1}^n (y_i - \langle w, x \rangle)^2$$

d 变量, n 次测量 $n \times d$ matrix

Step 3 Minimize MSE with respect to W

$$X = \begin{bmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_n- \end{bmatrix}$$

$$\begin{bmatrix} \langle w_1, x_1 \rangle \\ \langle w_2, x_1 \rangle \\ \vdots \\ \langle w_n, x_n \rangle \end{bmatrix} XW = \begin{bmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_n- \end{bmatrix} \begin{bmatrix} w \\ \vdots \end{bmatrix}$$

本集是 $\frac{1}{n}$,

n 为样本数目是常数, 设为 $\frac{1}{2}$

方便后续有系数

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$MSE = \left(\frac{1}{2} \right) \|y - XW\|_2^2$$

矩阵不能直接除

$$W^* = \left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \left[\sum_{i=1}^n y_i x_i \right]$$

$n \gg d$ $X^T X$ 可逆

intuitive:
we need

$n \gg d$ equations

to solve d variable

$$W^* = (X^T X)^{-1} X^T y$$

\uparrow \uparrow \uparrow \uparrow
 $d \times 1$ $d \times d$ $d \times n$ $n \times 1$

Multivariate

Case

Bad! \uparrow

inversion takes $d(d+1)$ time

这样求 MSE 最小的系数
太慢

so today:

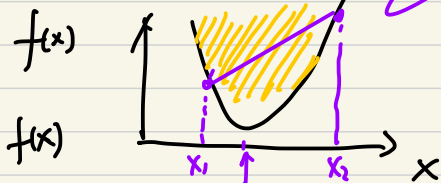
Goal: Improve this running time sufficiently using

Gradient descent

梯度下降 \Rightarrow 找最优值
(极值)

First:

convexity



Convex function (凸函数)

$\alpha x_1 + (1-\alpha)x_2$ (if $\alpha=0.5 \Rightarrow$ mid point)

Words: straight line

$(x_1, f(x_1))$ & $(x_2, f(x_2))$ should be

above $f(x)$

Integer Redm 变量空间

Math: 定义:

映射到/维

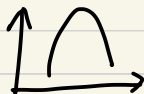
$f: \mathbb{R}^d \rightarrow \mathbb{R}$ for any $x, y \in \mathbb{R}^d$

$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$

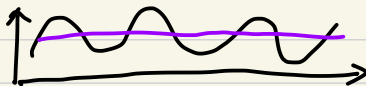
all α in $[0,1]$

Opposite of convex:

Concave (凹)

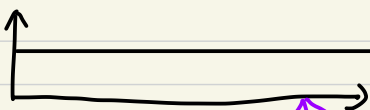


$\sin x$



Neither concave nor convex

Constant



Both concave and convex

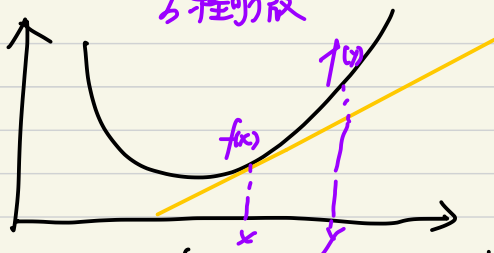
Straight



property

Assume: function has differentiability/smoothness

方程可微



Words:

tangent of $f(x)$ at any point lies
"below" the curve

Math:

Scalar

$$\text{For any } x, y, \quad f(y) \geq f(x) + \frac{df(x)}{dx}(y-x)$$

Vector

∇f

$$\begin{bmatrix} \frac{\partial f}{\partial x^{(1)}} \\ \frac{\partial f}{\partial x^{(2)}} \\ \vdots \\ \frac{\partial f}{\partial x^{(n)}} \end{bmatrix} = \nabla f(x)$$

$$\text{For all } x, y: f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle$$

x, y 横坐标

Observation:

If $\nabla f(x) = 0$, then $f(y) \geq f(x)$ for all y

$\Rightarrow x$ is a global minimum of f

此时 x 为最小值点

Gradient descent

寻找最小MSE对应系数
(MSE最小值而又是最小)

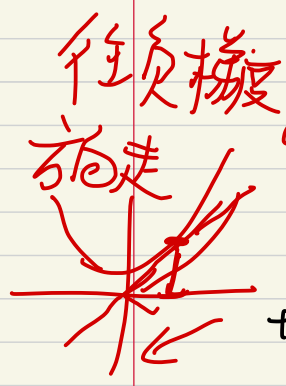


goal: find x^* that minimize $f(x)$

Approach: find α , $f(w + \alpha) < f(w)$ 往更小的方向找最小值

- $w \leftarrow w + \alpha$
- Repeat

$\nabla f(x)$ points in the direction of maximum change in f (该点的梯度方向)



$w_k - \alpha \nabla f(w_k)$ defines how much I want to move on that direction

"step size"
"learning rate"

for $k = 0, 1, 2, \dots, T$
number of epochs

terminate after T epochs 纪元更新次数
or when ∇f becomes small

$$(AB)^T = B^T A^T$$

$$\frac{\partial Ax}{\partial x} = A^T$$

Application of GD: Linear Regression

(MSE) $L(w) = \frac{1}{2} \|y - Xw\|_2^2$ 对 w 求导, 矩阵求导就是梯度

GD: - initialize w .
- iterate

$$\nabla L(w) = -X^T(y - Xw) = X^T(Xw - y)$$

往梯度下降方向 负梯度方向

$$w_{k+1} \leftarrow w_k - \alpha_k X^T(Xw_k - y) \text{ until } T \text{ or } \nabla \text{ became small}$$

迭代更新 w stop size

Running time : $O(nd) \cdot T$

- Set error parameter ϵ , terminate when $\|\nabla f(x)\|_2 \leq \epsilon$

- can prove that this happens after $T = \log \left[\frac{\|w_0\|_2}{\epsilon} \right]$ Fix x 可以知道 T (更新步数) 停止的步数
T 为运行次数时, 上面 \log 老和 ϵ (停止的步数)

Where $\rho = 1 - \frac{l}{L}$ (ρ is a number between $0 \sim 1$)
 $\rho = \frac{L-l}{L+l}$ l & L are smallest and largest eigen value of $X^T X$

Eigen value : $A v = \lambda v$
特征值 \uparrow \uparrow eigen vector eigen value