## Design Matrix

$$A = \begin{bmatrix} 1 & X_1 & X_1^2 & \cdots & X_1^{m-1} \\ & X_2 & X_2^2 & \cdots & X_2^{m-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ X_1 & X_n & X_n^2 & & X_n^{m-1} \end{bmatrix}_{n \times m}$$

## Weight Vector

$$W = \begin{bmatrix} W_0 \\ W_1 \\ \vdots \\ W_{m-1} \end{bmatrix}_{m \times 1}$$

## target Vector

$$b = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

⓪ LSE is We hope

Find a W to make

$$Aw = b$$

But A is not invertible ( perhaps not )

→ We only find $\tilde{W}$

Such $Aw = b$ to $\min \| Aw - b \|^2$

$$\|Aw - b\|^2 = (Aw - b)^T (Aw - b)$$

$$= W^T A^T A W - 2 W^T A^T b + b^T b$$

☆ $W^T A^T b$ is "Scalar" not a Vetocr

So, $(W^T A^T b)^T = b^T A W = W^T A^T b$

$$\min_{W_{LSE}} \|Aw - b\|^2 \Rightarrow \frac{\partial}{\partial W}(W^T A^T A W) - \frac{\partial}{\partial W}(2^T A^T b) + \frac{\partial}{\partial W}(b^T b)$$

$$\Rightarrow 2 A^T A W - 2 A^T b + 0$$

$$\Rightarrow A^T A W = A^T b$$

$$\Rightarrow W_{LSE} = (A^T A)^{-1} A^T b \quad \#$$

---

But $W = (A^T A)^{-1} A^T b \Rightarrow$ We only know $\det(A^T A)^{-1} \geq 0$

Can't 100% make sure full rank (invertable)

We must add regularized term "$L_1$ norm" or "$L_2$ norm"

To 100% make sure full rank (invertable)

We usually call rLSE,

We choose **L2 norm** = $\lambda \|W\|^2 \Rightarrow \frac{d}{dW} L_2 = 2\lambda W$

$$\min_{W} \|Aw-b\|^2 \Rightarrow 2A^TAW - 2A^Tb + 2\lambda W = 0$$

$$\Rightarrow (A^TA + \lambda I)W = A^Tb \Rightarrow W_{rLSE} = (A^TA + \lambda I)^{-1}A^Tb \#$$

While $\lambda > 0$, $A^TA + \lambda I$ 100% $\det > 0$, inertable #

Both LSE / rLSE We all need to calculate

Inverse Matirx, I will choose LU

Because only We ① slove I each row $e_i$

To slove $LUx_i = e_i$ ② Assemble togather We

Let $A^{-1}$ For example : $A = LU \Rightarrow A^{-1} = U^{-1}L^{-1}$

$A^{-1} = [X_1, X_2 \cdots X_n]$, $X_i$ is a column Vector

By $AA^{-1} = I \Rightarrow Ax_i = e_i$, $i = 1 \cdots n$ ①

$\Rightarrow$ Solve $y_i$

$\Rightarrow Ax_i = e_i \Rightarrow LUx_i = e_i$, let $Ux_i = y_i$ ② Slove $x_i$

Summary ① Foward substitution $Ly_i = e_i$, Slove $y_i$
       ② Backward substitution $Ux_i = y_i$, slove $x_i$

$A^{-1} = [X_1, X_2 \cdots X_n] \rightarrow$ We can get $A^T$ #

---

## ⑥ Steepest decent

The Update fomula of steepest decent

$$W_{t+1} = W_t - \eta \nabla f(W_t)$$

$W_t$ is parameter of $t$ steps
$\eta$ is learning rate (Steps size)

$\nabla f(W_t)$ is $W_t$'s gradient

Lipschitz contiuns with $L>0$ ( L is Lipschitz constant)
This make sure the change will not so fierce

From $f(W_t)$ We expansion of $f$ ( $f$ is Loss function)

by Tayler

$$f(W_{t+1}) \leq f(W_t) + \nabla f(W_t)^T (W_{t+1} - W_t) + \frac{L}{2} \| W_{t+1} - W_t \|^2$$

$$= f(W_t) - \eta \| \nabla f(W_t) \|^2 + \frac{L\eta^2}{2} \| \nabla f(W_t) \|^2$$

$$= f(W_t) - \eta(1 - \frac{L\eta}{2}) \| \nabla f(W_t) \|^2$$

$\| W_{t+1} - W_t \|$ must small enough, which implies

$\eta$ also must small enough

While $\eta \leq \frac{L}{2}$, this is monotonically increasing

( unless $\nabla f(W_t) = 0$ )

$\Rightarrow$ Best learning rate : $\eta = \frac{1}{L}$

$\Rightarrow f(W_{t+1}) \leq f(W_t) - \frac{1}{2L} \| \nabla f(W_t) \|^2$  this is discription
                                                                        of convergance
                                                                        properties

Regularized term in L1 norm
the gradient of it can be written

as sign funtion $Sign(W_i) = \begin{cases} 1, & \text{if } W_i > 0 \\ -1, & \text{if } W_i < 0 \\ 0, & \text{if } W_i = 0 \end{cases}$

The gradient in total is $2A^T(A\vec{W}-b) + \lambda \cdot sign(\vec{W})$

---

# 6 Newton's method

Basic type:
Based Talyers expesion

$$f(X) \approx f(X_0) + \nabla f(X_0)^T(X-X_0) + \frac{1}{2}(X-X_0)^T H(X_0)(X-X_0)$$

$H(X_0) = \nabla^2 f(X_0)$ is Hessian Matrix

# Newton's method optimization problem :

Same We wanna slove $\min_{w} f(w)$

We must find $\nabla f(w) = 0$

$\nabla f(w_t) + \underline{H(t)} (w - w_t) = 0$

Slove step : $w - w_t = -H(t)^{-1} \nabla f(w_t)$

Update fomula : $w_{t+1} = w_t - H(t)^{-1} \nabla f(w_t)$

# For LSE problem

We already knew $f(w) = \|Aw - b\|^2$

$\nabla f(w) = 2A^{\top}(Aw - b)$

$\nabla^2 f(w) \overset{\triangle}{=} H(w) = \underline{2A^{\top}A}$ ( it not related with w ! )

Apply Update fomula $w_{t+1} = w_t - (2A^{\top}A)^{-1} \cdot 2A^{\top}(Aw_t - b)$

$\Rightarrow w_{t+1} = w_t - (A^{\top}A)^{-1} A^{\top}(Aw_t - b)$

# Special propeties of Newton's method

For any quadratic func (ie. LSE)

One step is enough!

$$W_1 = W_0 - (A^TA)^{-1} A^T(AW_0 - b)$$
$$= W_0 - (A^TA)^{-1} A^T AW_0 + (A^TA)^{-1}A^Tb$$
$$= W_0 - W_0 + (A^TA)^{-1}A^Tb$$
$$= (A^TA)^{-1}A^T b$$

Does it look familer ~

Because this is optimal Solution of LSE, as shown in previous We calcuted #