LINEAR LEAST SQUARES        (LLS)

- FINDING SOL. FOR LINEAR REGRESSION
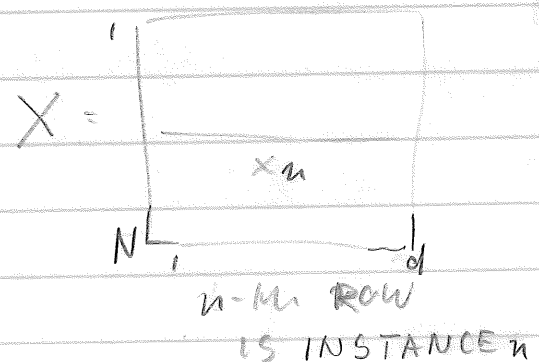  (BATCH SETTING)
  BY SOLVING THE QUADRADIC EQUATIONS

$$\min_{w} \left( \sum_{u=1}^{N} \underbrace{(\bar{x}_t \cdot \bar{w}}_{d\times 1 \ d\times 1} - \underset{1\times 1}{y_t})^2 \right)$$

$$\underset{T\times n}{\|} X \underset{n\times 1}{\bar{w}} - \underset{T\times 1}{\bar{y}} \|^2$$

$$X = \begin{bmatrix} & & \\ & x_u & \\ & & \end{bmatrix}$$

$$N \underset{\substack{n\text{-th ROW}\\ \text{IS INSTANCE } u}}{\Big\lfloor \quad\quad\quad \sim d}$$

$$= (X\bar{w} - \bar{y})^T (X\bar{w} - \bar{y})$$

$$= (\bar{w}^T X^T - \bar{y})^T (X\bar{w} - \bar{y})$$

$$= \bar{w}^T X^T X \bar{w} - \bar{w}^T X \bar{y} - \bar{y}^T X \bar{w} + \bar{y}^T \bar{y}$$

$$= \bar{w}^T X^T X \bar{w} - 2\bar{y}^T X \bar{w} + \bar{y}^T \bar{y} \qquad\qquad \text{QUADRADIC}$$

$$\nabla_w \|X\bar{w} - \bar{y}\|^2 = 2 X^T X \bar{w} - 2 X^T \bar{y} = 0 \qquad (*)$$

$$X^T X \bar{w} = X^T \bar{y} \qquad\qquad \text{NORMAL EQUATIONS}$$

$$Xw = y \quad \text{MIGHT NOT HAVE SOLUTION}$$

BUT NORMAL EQATIONS ALWAYS HAVE
SOLUTION  (Why?)

$$X^T X \bar{w}^x = X^T y$$

$$\bar{w}^x = (X^T X)^{-1} X^T y$$
$$\underset{d \times 1}{} \quad \underset{d \times N \, N \times d}{} \quad \underset{d \times N \, N \times 1}{}$$

PROBLEM: $X^T X$ MIGHT NOT HAVE FULL RANK

$$r(X^T X) = r(X)$$

FIX 1: REGULARIZE

$$\min_{w} \left( \lambda \|w\|^2 + \|X\bar{w} - \bar{y}\|^2 \right) \quad \overset{\text{CALLED RIDGE}}{\underset{\text{REGRESSION}}{\leftarrow}}$$

$$\nabla_w \left( \lambda \|w\|^2 + \|X\bar{w} - \bar{y}\|^2 \right)$$

$$= 2\lambda\bar{w} + 2X^T X\bar{w} - 2X^T \bar{y}$$

$$= 2\left( (\lambda I + X^T X)\bar{w} - X^T \bar{y} \right)$$

$$= 0$$

$$\bar{w}^* = \underbrace{(\lambda I + X^T X)^{-1}}_{\substack{\text{ALWAYS} \\ \text{FULL RANK}}} X^T y$$

FIX 2: PSEUDO INVERSE

DIGRESSION

— MATRIX DECOMPOSITION:

A SYMMETRIC IF $A = A^T$

U ORTHOGONAL IF $UU^T = U^T U = I$ } A, U SQUARE

$\forall$ SYM. $A$ : $\exists$ ORTH. $U$ & REAL $\sigma$ S.t.
$M \times M$    MATR.     DIAG MATR.

$$A = U \sigma U^T$$

$$= \boxed{U} \begin{pmatrix} 0 & 0 \\ & \sigma \end{pmatrix} \boxed{U}^T$$

THE $m$ COLUMNS $u_i$ OF $U$
ARE THE EIGENVECTORS OF $A$
AND THE DIAGONAL ELEMENTS OF
$\sigma$ THE EIGENVALS

$u_i$'s ARE ORTHOGONAL :

$$u_i \cdot u_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & = \end{cases}$$

$$U \sigma U^T$$

$$= U \sum_i \sigma_i e_i e_i^T U^T$$

$$= \sum_i \sigma_i \underset{M \times M}{U} \underset{M \times 1}{e_i} \underset{1 \times M}{e_i^T} \underset{M \times M}{U^T}$$

$$= \sum_i \sigma_i u_i u_i^T$$

$u_j$ IS EIGENVEC $j$

$$\left( \underbrace{\sum_i \sigma_i u_i u_i^T}_{A} \right) u_j = \sum_i \sigma_i u_i (u_i^T u_j)$$

$$= \sigma_j u_j \underbrace{u_j^T u_j}_{1}$$

$$= \sigma_j u_j$$

$\uparrow$

EIGENVALUE

WHAT IF $A$ non-square ?

$$\boxed{A}_{M \times N} = \boxed{U}_{M \times M} \quad \boxed{\begin{smallmatrix} \sigma & & \\ 0 & & 0 \\ & & 0 \end{smallmatrix}}_{M \times N} \quad \boxed{V^T}_{N \times N}$$

$UU^T = I_{M \times M}$

NON-NEG. REAL
DIAGONAL
MATRIX

$VV^T = I_{N \times N}$

- DIAGONAL ENTRIES OF $\sigma$ ARE NON-NEGATIVE
- CALLED SINGULAR VALUES

$$= \sum_{i=1}^{\min(M,N)} \sigma_i \, u_i \, v_i^T$$

MATLAB: $[U, D] = eig(A)$

$[U, D, V] = svd(A)$

RETURN TO FIX 2 :

$$\underbrace{X^T}_{d\times N} \underbrace{X}_{N\times d} \underbrace{\vec{w}^*}_{d\times 1} = \underbrace{X^T}_{d\times N} \underbrace{\vec{y}}_{N\times 1}$$

SOLUTION:  $\vec{w}^* = \underset{\uparrow}{X^+} \vec{y}$

PSEUDO INVERSE OF $X$

If $X = USV^T$
$X^+ := VS^+U^T$

$\uparrow$ INVERSE OF NON-ZERO ELMT'S OF DIAG
LEAVE 0'S AS 0'S

$$\begin{pmatrix} 6 & & \\ & 2 & \\ & & 0 \end{pmatrix}^+ = \begin{pmatrix} \frac{1}{6} & & \\ & \frac{1}{2} & \\ & & 0 \end{pmatrix}$$

CHECK NORMAL EQUATIONS !

$$\underbrace{VS^+U^T}_{X^T} \underbrace{USV^T}_{X} \underbrace{VS^+U^T\vec{y}}_{\vec{w}^*}$$

$$= V \underbrace{S^T}_{d\times d} \underbrace{S}_{d\times N} \underbrace{S^+}_{N\times d} \underbrace{U^T}_{d\times N} \underbrace{}_{N\times N} \underbrace{\vec{y}}_{N\times 1}$$

$$= \underbrace{V \quad S^T \quad U^T \quad \vec{y}}_{\underset{X^T}{d\times N}}$$

CONNECTION BETWEEN

FIX 1 & FIX 2

$$\lim_{\lambda = 0} (X^T X + \lambda I)^{-1} X^T = X^+$$
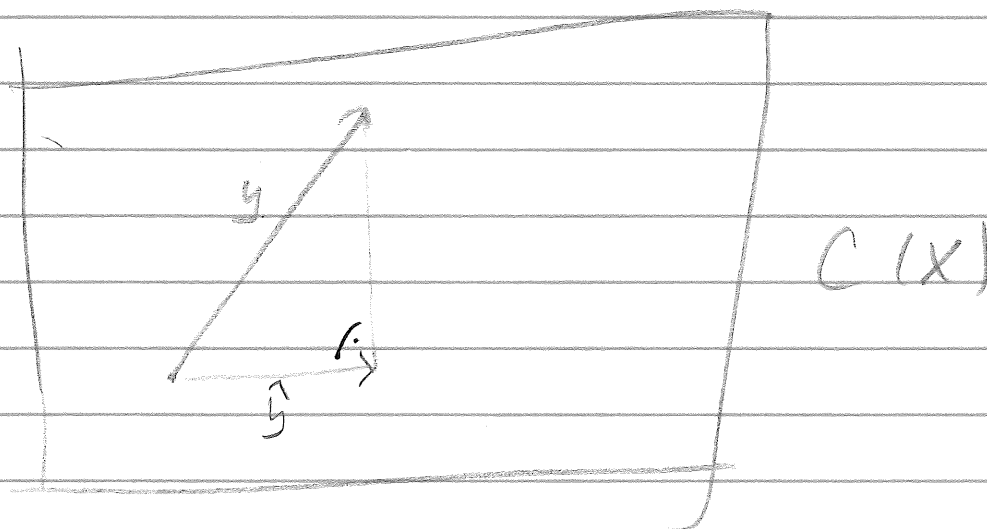
# GEOMETRIC VIEW

$$\min_{w} \quad \| \underbrace{Xw}_{\hat{y}} - y \|_2^2$$

$\hat{y} \in$ COLUMN SPACE OF $X$

$$C(X)$$

OPT $\hat{y}$ IS CLOSEST VECTOR TO $y$
THAT LIES IN $C(X)$



$\hat{y}$ IS PROJECTION OF $y$ ONTO $C(X)$

$$\hat{y} = Py$$
$$= \underbrace{X^+ X^T}_{\text{PROJ. MATRIX}} y$$

$$X = U\sigma V'$$
$$X^+ X^T = U\sigma^+ V' V \sigma U'$$
$$= \tilde{U}\tilde{U}^T$$

ALL COL. WITH
NON-ZERO EIGEN
VAL

$$= \sum_{i\,:\,\lambda_i \neq 0} u_i u_i^T$$

BREGMAN PROJ WRT $\| \|_2^2$

Recall GD update for minimizing batch loss

$$L(w) = \sum_u L_{y_u}(w \cdot x_u)$$

$$w_{q+1} = w_q - \eta \sum_u L'_{y_u}(w_q \cdot x_u) x_u$$

Motivation:

$$w_{q+1} = \arg\min_w \left( \frac{1}{2} \| w - w_q \|^2 + \eta \sum_u L_{y_u}(w \cdot x_u) \right)$$

$$\nabla '' \Big|_{w = w_{q+1}} = 0$$

$$w_{q+1} - w_q + \eta \sum_u L'_{y_u}(w_{q+1} \cdot x_u) x_u$$

$$w_{q+1} = w_q - \eta \sum_u L'_{y_u}(w_{q+1} \cdot x_u) x_u$$

IMPLICIT UPDATE

$$\approx w_q - \eta \sum_u L'_{y_u}(w_q \cdot x_u) x_u$$

APPROXIMATION

EXPLICIT UPDATE

IMPLICIT USES GRADIENT OF LOSS AT CURRENT ITERATION
EXPLICIT         "                              LAST         "

PRECISE WAY TO MOTIVATE EXPLICIT UPDATE:

$$w_{q+1} = \underset{w}{\arg\inf} \left( \|w - w_q\|^2 + \eta \left( L(w_q) + (w - w_q)^T \Delta L(w_q) \right) \right)$$

INERTIA TERM     FIRST ORDER APPROX.
NEEDED BECAUSE    OF BATCH LOSS
LOSS LINEAR       AT $w_q$

BATCH GD. US STOCHASTIC GD BASED ON
SINGLE EXAMPLE

TIME    1 ITERATION $\approx$ 1 PASS

↑
converges faster
because it uses
"more recent" gradients

IDEA: WHY NOT USE 2ND ORDER APPROXIM.
OF LOSS

$$L(w) \approx L(w_q) + (w-w_q)^T \nabla L(w_q)$$
$$\phantom{L(w) \approx L(w_q) + } \underset{1 \times d}{} \phantom{(w-w_q)^T} \underset{d \times 1}{}$$

$$+ \tfrac{1}{2}(w-w_q)^T \nabla^2 L(w_q)(w-w_q)$$
$$\phantom{+ \tfrac{1}{2}(w-w_q)^T} \underset{d \times d}{}$$

quadratic approx. of $L(w)$ at $w_q$

$$\hat{L}_q(w)$$
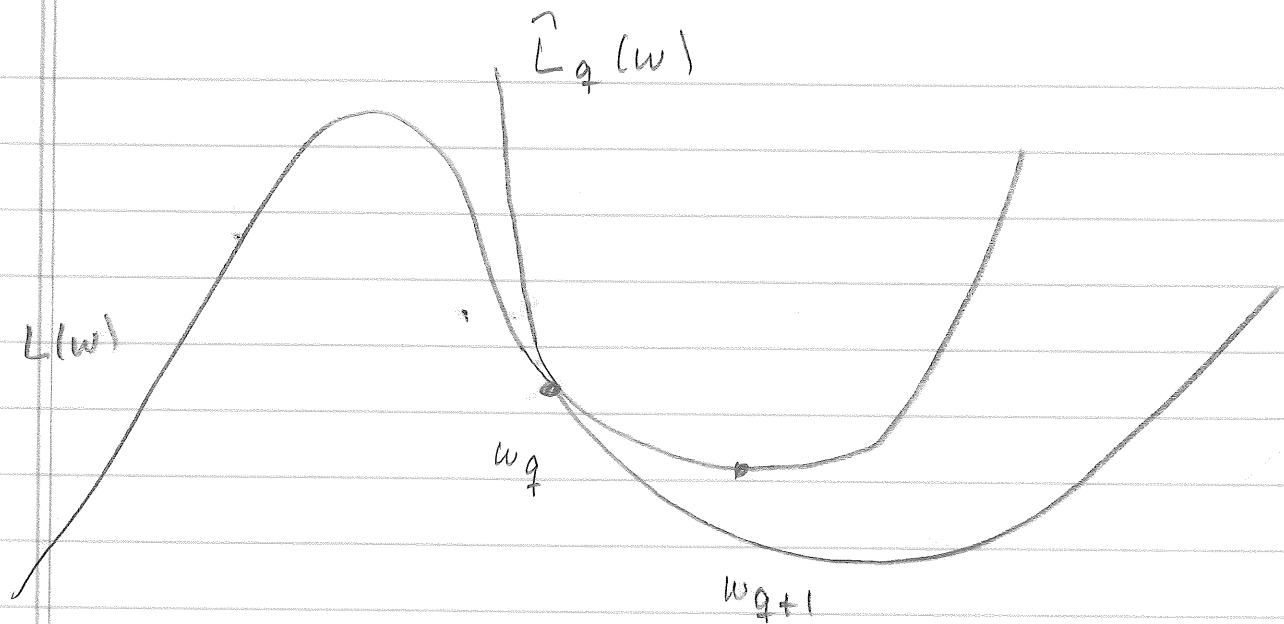
$$w_{q+1} = \underset{w}{\arg\min}\ \hat{L}_q(w) \qquad\qquad (1)$$

WITH INERTIA TERM:
$$w_{q+1} = \underset{w}{\arg\min}\ \|w-w_q\|^2 + y\,\hat{L}_q(w) \qquad (2)$$

(1) SIMPLER
  FOCUS ON (1) FIRST

$\hat{L}_q(w)$

$L(w)$

$w_q$

$w_{q+1}$

$$\underset{d \times 1}{\nabla \hat{L}_q(w)} = \underset{d \times 1}{\nabla L(w_q)} + \underset{d \times d}{\nabla^2 L(w_q)} \underset{d \times 1}{(w - w_q)}$$

$$\nabla \hat{L}(w) \Big|_{w = w_{q+1}} = 0$$

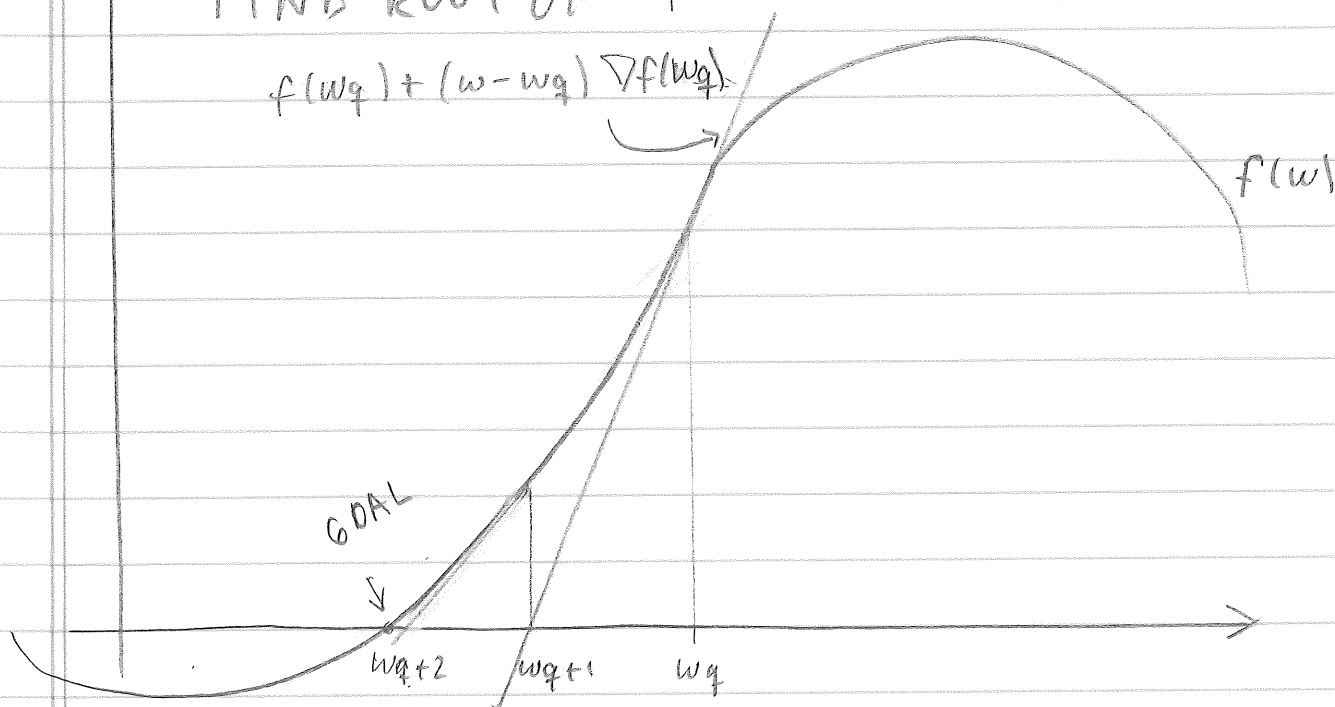$$\nabla L(w_q) + \nabla^2 L(w_q)(w_{q+1} - w_q) = 0$$

$$w_{q+1} - w_q = -\left(\nabla^2 L(w_q)\right)^{-1} \nabla L(w_q)$$

$$w_{q+1} = w_q - \left(\nabla^2 L(w_q)\right)^{-1} \nabla L(w_q)$$

SECOND VIEW OF NEWTON

FIND ROOT OF $f(w) := \nabla L(w)!$

$f(w_q) + (w - w_q)\nabla f(w_q)$

$f(w)$

GOAL

$w_{q+2}$   $w_{q+1}$   $w_q$

$$f(w_q) + (w - w_q)\nabla f(w_q)\Big|_{w = w_{q+1}} = 0$$

$$w_{q+1} - w_q = -(\nabla f(w_q))^{-1} f(w_q)$$

$$w_{q+1} = w_q - (\nabla f(w_q))^{-1} f(w_q)$$

# NEWTON W. REGULARIZATION

$$w_{q+1} = \inf_{w} \left( \frac{\lambda}{2} \| w - w_q \|^2 + L(w_q) + (w - w_q)^\top \nabla L(w_q) \right.$$
$$\left. + \frac{1}{2} (w - w_q)^\top \nabla^2 L(w_q) (w - w_q) \right)$$

$$\nabla'' = \lambda (w - w_q) + \nabla L(w_q) + \nabla^2 L(w_q)(w - w_q)$$

$$\nabla'' |_{w = w_{q+1}} = 0$$

$$\nabla L(w_q) + \left( \nabla^2 L(w_q) + \lambda I \right) (w_{q+1} - w_q) = 0$$

$$w_{q+1} = w_q - \left( \nabla^2 L(w_q) + \lambda I \right)^{-1} \nabla L(w_q)$$

LLS   = 1 STEP OF NEWTON

RIDGE REGRESSION
        = 1 STEP OF REG. NEWTON

# NEWTON-RAPSON FOR LOGISTIC REGRESSION

BATCH LOSS $\qquad L(w) = \sum_n L_{y_n}(w \cdot x_n)$

WHERE $\quad L_y(a) = \ln(1 + e^a) - ya$

$$\nabla L(w) = X^T(\hat{y} - y) = \sum_n x_n (\hat{y}_n - y_n)$$

X HAS EXAMPLE AS ROWS
$X^T$ " " COLS
$\begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$

$$\nabla_w \times \left( \sigma(\overbrace{w \cdot x}^{a}) - y \right)$$

$$= X x' \; \sigma'(a)$$

$$= X x' \; \sigma(a)(1 - \sigma(a))$$

$$\sigma(a) = \frac{e^a}{1 + e^a} \qquad \sigma'(a) = \frac{e^a(1 + e^a) - e^{2a}}{(1 + e^a)^2}$$

$$= \sigma(a)(1 - \sigma(a))$$

$$\nabla \nabla L(w)$$

$$= \nabla \; X^T (\hat{y} - y)$$

$$= \nabla \sum_n x_n (\hat{y}_n - y_n)$$

$$= \sum_n x_n x_n^T \; \sigma(w x_n)(1 - \sigma(w x_n))$$

$$= X^T R X$$

$$\uparrow$$

$$\left( \hat{y}_n (1 - \hat{y}_n) \right)$$

LINEAR REGR.

$$\nabla^2 L(w) = X^T X$$

R depends on weight vector

NEWTON RAPSON

$$w_{q+1} = w_q - (X^T R_q X)^{-1} X^T (\hat{y}_q - y)$$

$$= (X^T R_q X)^{-1} (X^T R_q X w_q - X^T (\hat{y}_q - y))$$

$$= (X^T R_q X)^{-1} X^T R_q z_q$$

$$\text{WHERE} \quad z_q = X w_q - R_q^{-1}(\hat{y}_q - y)$$

LLS $\quad w^* = (X^T X)^{-1} X y$

> CALLED ITERATED REWEIGHTED
  LEAST SQUARES

MANY NUMERICAL ISSUES

EXPLICIT →GD    USES OLD GRADIENTS
   NEWTON    "          "      + HESSIANS

IMPLICIT GD    USES CURRENT GRADIENTS

Conjecture: IMPLICIT GD BEATS NEWTON