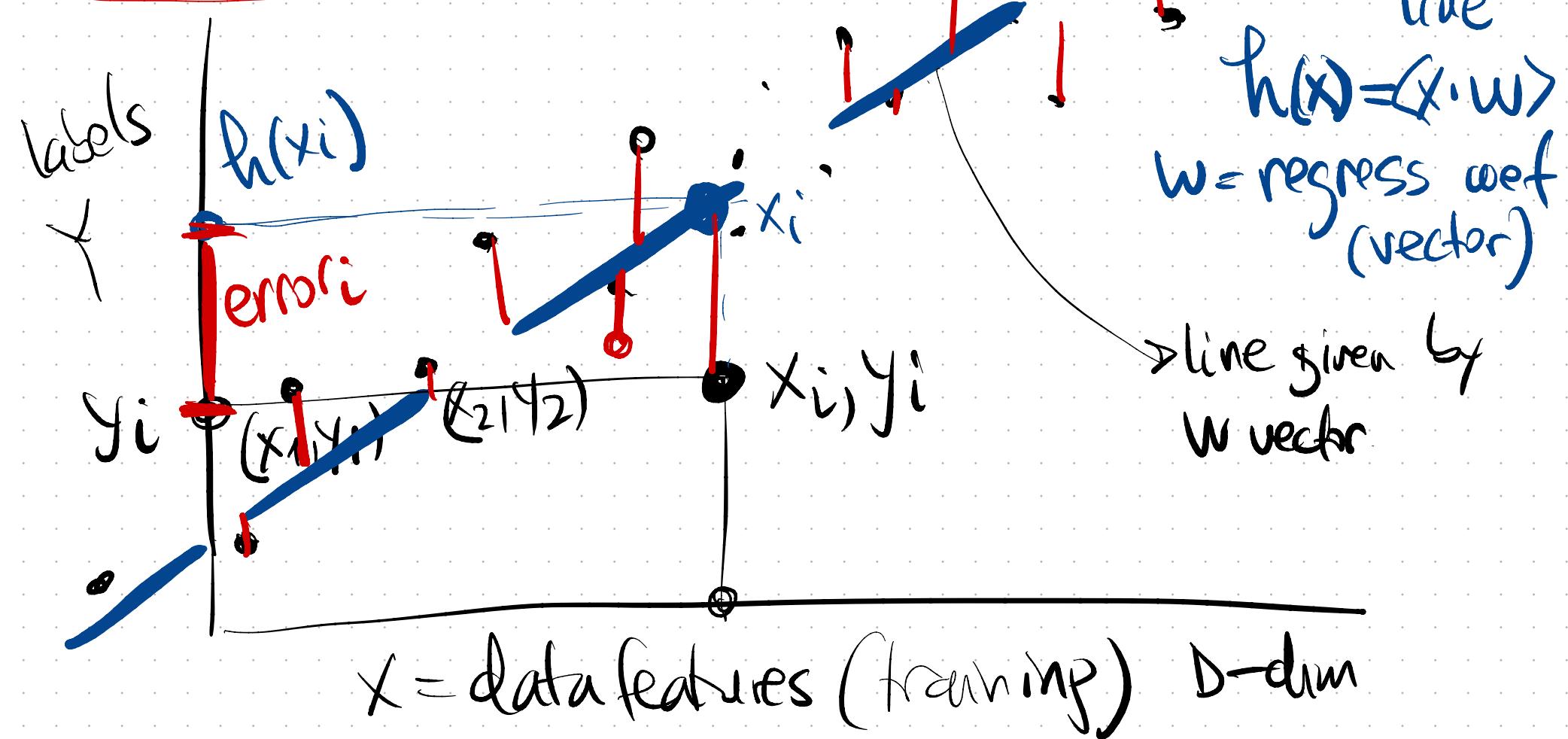


Lecture 5/14/25

- Linear regression
- HW1 (demo brief) 7:20 - Caleb
- HW1 math pb5: convex hulls



Regression line (classifier) by coef $\mathbf{w} = (w^1, w^2, \dots, w^D)$

datapoint x_i

$$h(x_i) = \langle x_i \cdot \mathbf{w} \rangle = \sum_{d=1}^D x_i^d \cdot w^d$$

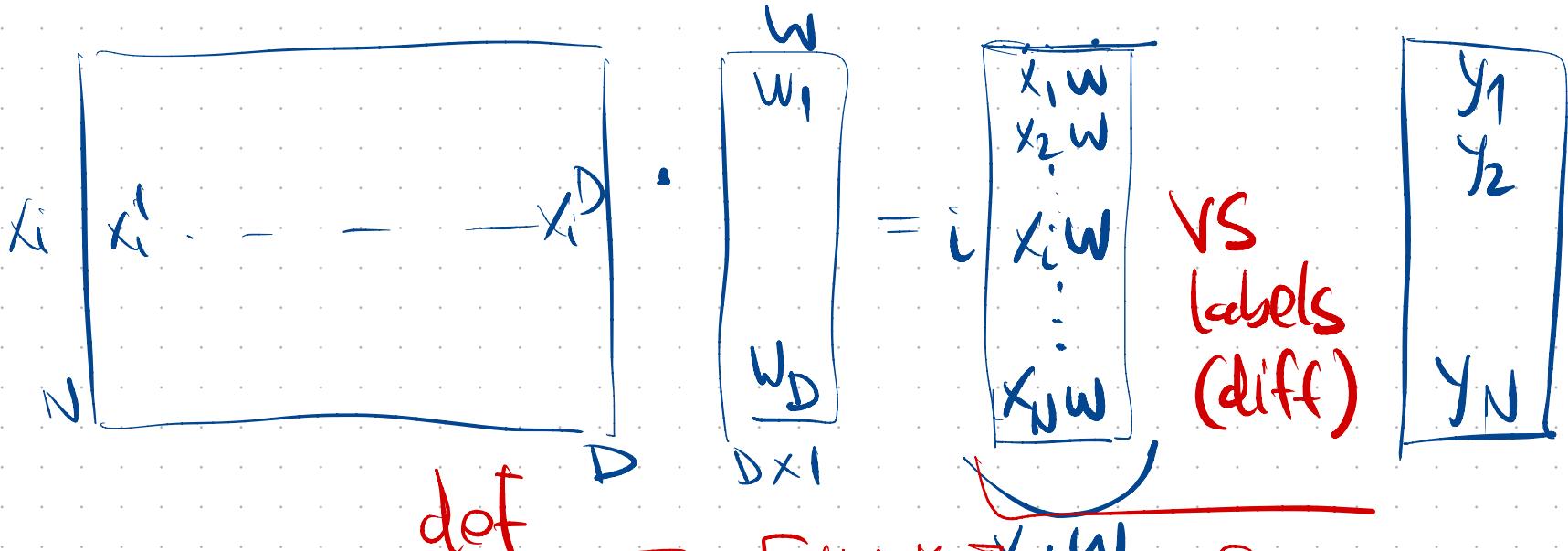
want $h(x_i) \approx y_i$
predict label

sq error_i = $(h(x_i) - y_i)^2$
 $= (\langle x_i \cdot \mathbf{w} \rangle - y_i)^2$

$$x_i = (x_i^1, x_i^2, \dots, x_i^D)$$

All datapoints \rightarrow matrix $X_{N \times D}$

	1	2	D	
1	x_1^1	x_1^2	x_1^D	y_1
2	x_2^1	x_2^2	x_2^D	y_2
N	x_N^1	x_N^2	x_N^D	y_N



$$\underbrace{X \cdot W - Y = E}_{N \times D \quad D \times 1 \quad N \times 1} ; \quad \underbrace{\begin{bmatrix} x_1 \cdot w - y_1 \\ x_2 \cdot w - y_2 \\ \vdots \\ x_N \cdot w - y_N \end{bmatrix}}_{E^T} ; \quad E = \begin{bmatrix} x_1 \cdot w - y_1 & x_2 \cdot w - y_2 & \cdots & x_N \cdot w - y_N \end{bmatrix}$$

want sq error: $J(w) = \frac{1}{2} \sum_{i=1}^N (h(x_i) - y_i)^2 = \frac{1}{2} E^T \cdot E = \frac{1}{2} (xw - y)^T (xw - y)$

$$\langle E^T \cdot E \rangle = \sum_{i=1}^N e_i \cdot e_i = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (x_i \cdot w - y_i)^2 \quad E^T \quad E$$

want: to find w vector $w = (w^1, w^2, \dots, w^D)$ to minimize the sq error $\frac{1}{2} (xw - y)^T \cdot (xw - y)$

\rightarrow CONVEX: $\frac{\partial J}{\partial w} > 0$ positive

intuition $J(w) = \frac{1}{2}(xw - y)^T(xw - y)$ convex in w

$\Rightarrow \frac{\partial J}{\partial w} = 0$ when $J(w)$ is minimum.

$\text{dim } (1 \times D)(D \times D)(N \times 1)$
same
transposed

\Leftrightarrow reg term

convex

$J(w)$

$$\frac{\partial J}{\partial w} = \frac{1}{2} \frac{\partial}{\partial w} (xw - y)^T(xw - y)$$

$$= \frac{1}{2} \frac{\partial}{\partial w} (w^T x^T x w - w^T x^T y - y^T x w)$$

$$+ \frac{1}{2} \frac{\partial}{\partial w} w^T w \cdot \lambda \frac{\partial}{\partial w} =$$

$$+ \frac{1}{2} \frac{\partial}{\partial w} w^T w \cdot \lambda \downarrow \text{pretend 1-dim}$$

$$+ \lambda I \cdot w$$

$$? \quad \frac{\partial}{\partial w} (y^T y) = 0$$

$$= \frac{1}{2} \frac{\partial}{\partial w} (w^T x^T x w - 2 w^T x^T y - y^T y)$$

$$= \frac{1}{2} \frac{\partial}{\partial w} w^T x^T x w - \frac{1}{2} \frac{\partial}{\partial w} w^T x^T y$$

$$= \frac{1}{2} \frac{\partial}{\partial w} w^T x^T x w - \frac{1}{2} \frac{\partial}{\partial w} w^T x^T y$$

$$= \frac{1}{2} 2 x^T x w - x^T y + \lambda I \cdot w$$

$$= x^T x w - x^T y + \lambda I w$$

want $= 0$ \rightarrow sq $D \times D$, pos definite (pos eigen val)

$$x^T x w = x^T y$$

invert-left by $(x^T x)^{-1}$

$$(x^T x + \lambda I) w = x^T y$$

$$w = (x^T x + \lambda I)^{-1} \cdot x^T y$$

$$w = (x^T x)^{-1} x^T y$$

NORMAL EQ (closed form)

Best (Th) w for min square error of reg. line

1-dim w differentiable

$$\frac{\partial}{\partial w} (\underline{w^T \cdot C \cdot w}) = \frac{\partial}{\partial w} (w^2 \cdot c) = \underline{2wc}$$

$\frac{\partial}{\partial w} = ? \Rightarrow$ let math people worry about this!
 \downarrow vector
 \rightarrow If interested, talk about at OHL (not repurc
 for course)

Naive way: we want $X \cdot w \approx y$

why not (if $X \approx$ sq matrix)
 make it square?

$N \times D$ $D \times 1$

$N \times 1$

$$w = X^{-1}y ?$$

unstable,

w not min sgn.

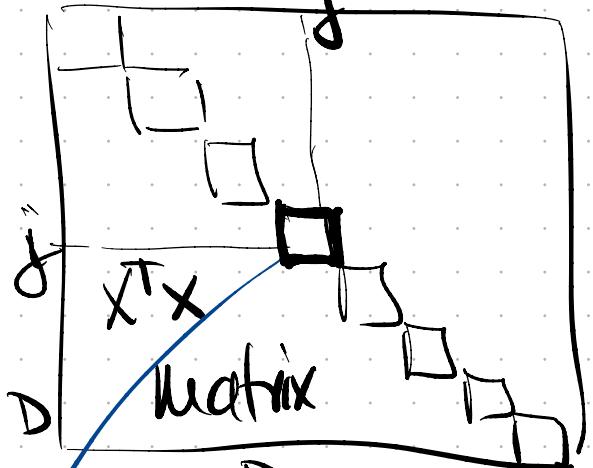
X^{-1} terrible
 to invert

Correct way

$X^T X$ stable matrix, pos. def. \rightarrow

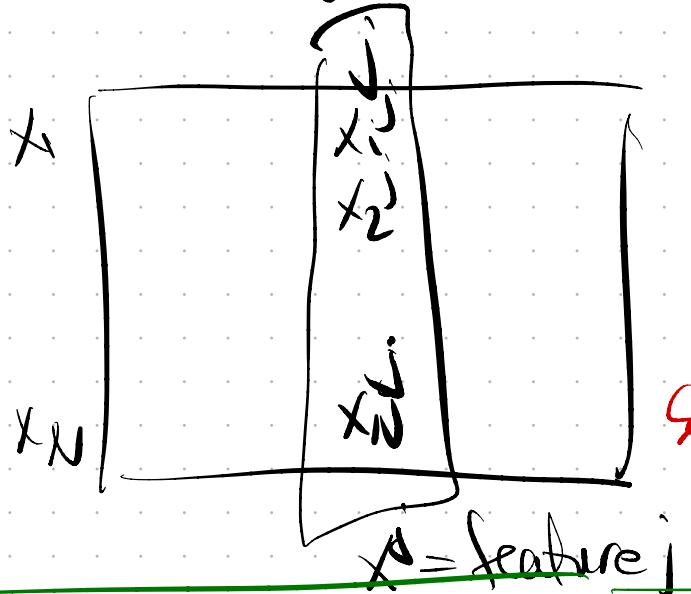
$D \times D$ \Rightarrow invert this matrix

$x = \text{centered}$ (all column mean=0) $\Rightarrow X^T X = N \cdot \text{COVAR}(x)$

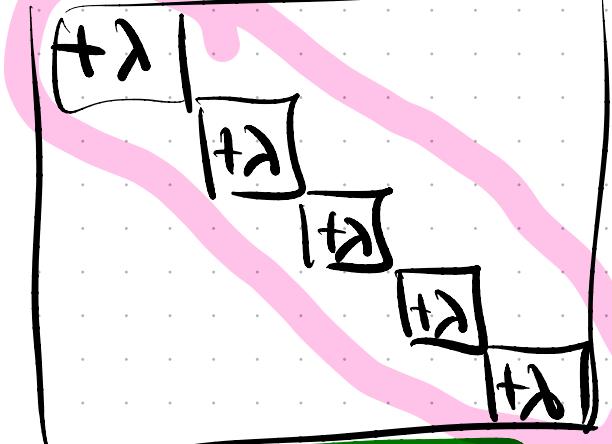


add stability
 (for inversion)
 $X^T X^{-1}$

diagonal vals (pos j)
 Column j , column j
 $m \times m$
 $\text{var}(j \text{ feature}) = \text{var}(X^j)$



λ = hyperparameter of L_2 reg



Increase
 wak
 diagonal

$$X^T X + \lambda I$$

more stable
 $\neq X^T X$

identity

$$\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 0 \end{bmatrix}$$

this op $X^T X + \lambda I$ actually
 corresponds to regularized L2 error

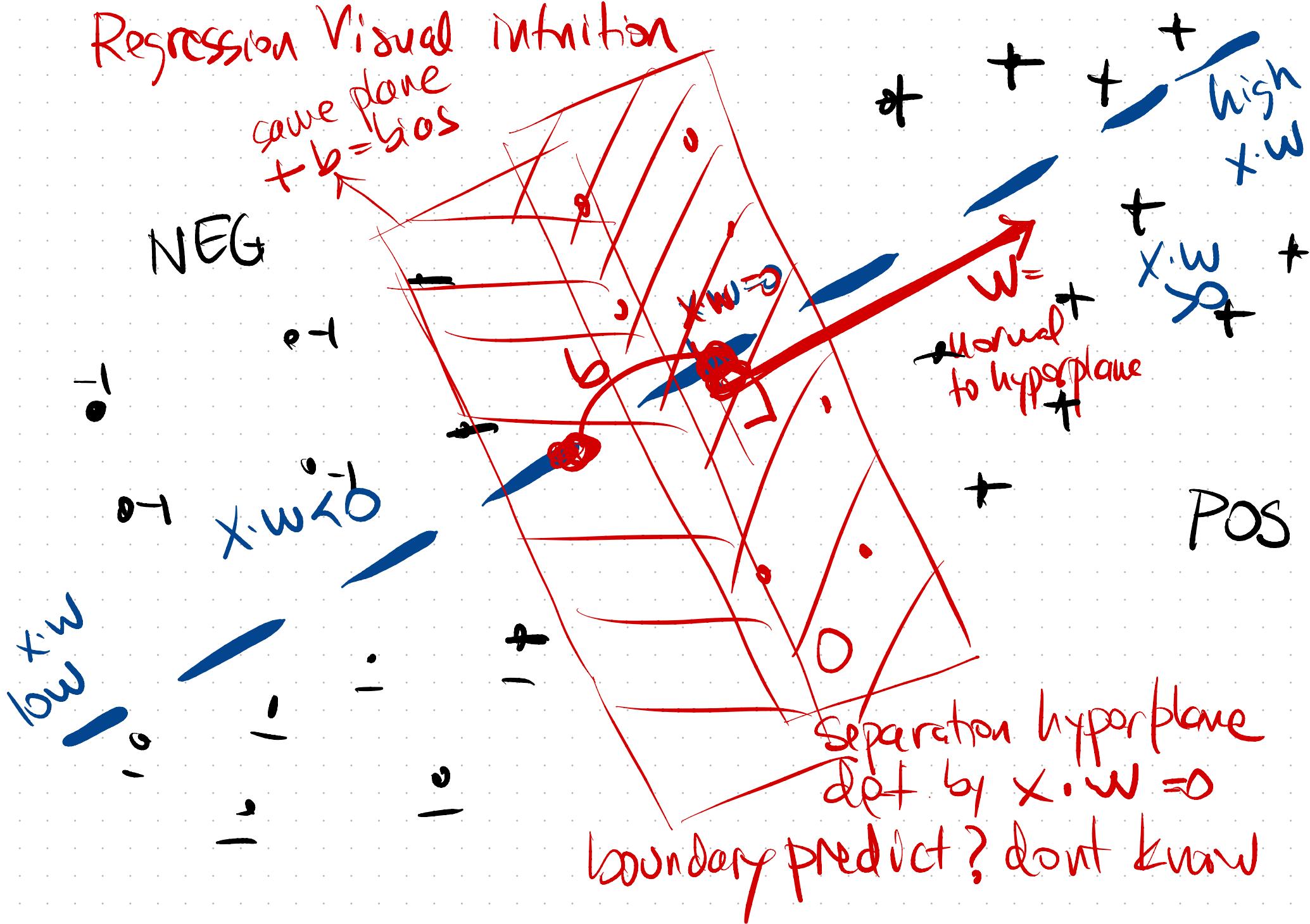
$$J(w) = \frac{1}{2} (xw - y)^T (xw - y) + \frac{\lambda}{2} \sum_{d=1}^D (w_d)^2$$

error

$$= \frac{1}{2} (xw - y)^T (xw - y) + \frac{\lambda}{2} w^T w$$

prevents large w values
 $\Leftrightarrow L_2$ regularization

Regression Visual Intuition



Hw 1 PB 3: 1-dim w = a bias = b

Reg $f_i(x_i) = x_i \cdot a + b \stackrel{\text{want}}{\approx} y_i \quad \forall i=1:n$

1-dim "1" red

x_1	1
x_2	1
x_N	1

$$\begin{bmatrix} a & b \end{bmatrix}$$

=

$$\begin{array}{|c|} \hline x_1 \cdot a + b \\ x_2 \cdot a + b \\ \vdots \\ x_N \cdot a + b \\ \hline \end{array}$$

vs

$$\begin{array}{|c|} \hline y_1 \\ y_2 \\ \vdots \\ y_N \\ \hline \end{array}$$