

Lecture 9 - LS, Classification, Ortho Bases, Regularization

LS & Classification

Given training $\{\underline{x}_i, y_i\}_{i=1}^n$, $\underline{x}_i \in \mathbb{R}^P$ (e.g. $\underline{x}_i = \begin{bmatrix} \text{left mouth to brow} \\ \text{right " " } \\ \text{mouth width} \end{bmatrix}$)

$y_i \in \mathbb{R}$ (degree of happiness)

or

$y_i \in \{-1, 1\}$ (happy or sad)

Linear prediction

$$\begin{aligned}\hat{y}_i &= w_1 x_{i1} + w_2 x_{i2} + \cdots + w_p x_{ip} \\ &= \langle \underline{w}, \underline{x}_i \rangle = \underline{w}^\top \underline{x}_i\end{aligned}$$

(2)

A) Build a "data matrix" or "feature matrix" and label vector

$$X = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times p}, \quad \underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

linear model

$$\hat{\underline{y}} = \underline{X} \underline{w}$$

B) solve least squares opt:

$$\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{argmin}} \|\underline{y} - \underline{X} \underline{w}\|_2^2 \rightarrow \sum_{i=1}^n (y_i - \underline{x}_i^T \underline{w})^2$$

if cols of X are linearly independent $\Rightarrow X^T X$ is positive definite

$$\Rightarrow \underline{a}^T X^T X \underline{a} > 0 \quad \forall \underline{a} \neq 0 \quad \Rightarrow X^T X \text{ is invertible}$$

e.g. $X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \rightarrow \star \underline{a}^T X^T \underline{X} \underline{a} \text{ w/ } \underline{a} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \Rightarrow \underline{X} \underline{a} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow \underline{a}^T \underline{X}^T \underline{X} \underline{a} = 0$
 $\Rightarrow X^T X \text{ is NOT pos. def.}$

(3)

if $X^T X$ is pos. def, then \exists unique LS solution

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

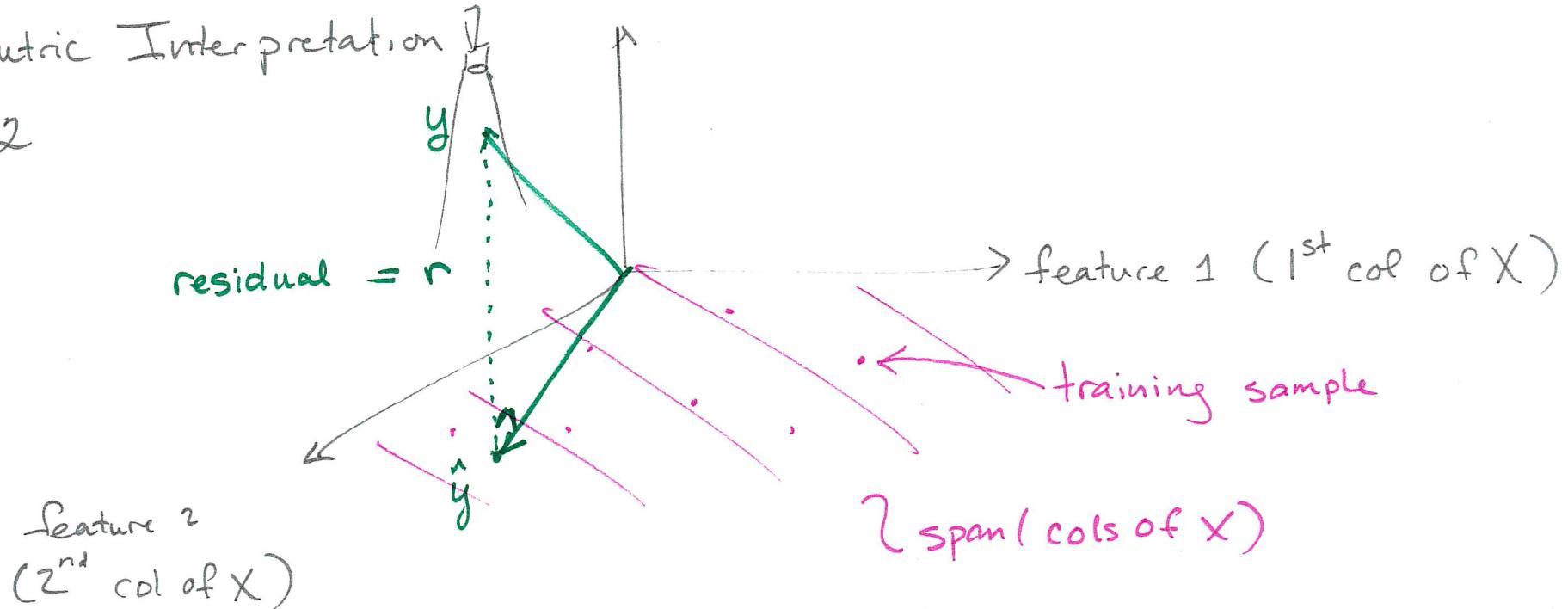
c) Validate w/ test / hold out data

$$\{\tilde{x}_i, \tilde{y}_i\}_{i=1}^{n_T}$$

$$\sum_{i=1}^{n_T} (\tilde{y}_i - \tilde{x}_i^T \hat{w})^2 = \text{error on test data}$$

Geometric Interpretation

$$P = 2$$



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

$$\hat{\underline{y}} = \underline{X} \hat{\underline{w}}$$

$$= \underline{X} \underbrace{(\underline{X}^T \underline{X})^{-1}}_{\text{projection matrix } P_x} \underline{X}^T \underline{y}$$

projection matrix P_x

Binary classification

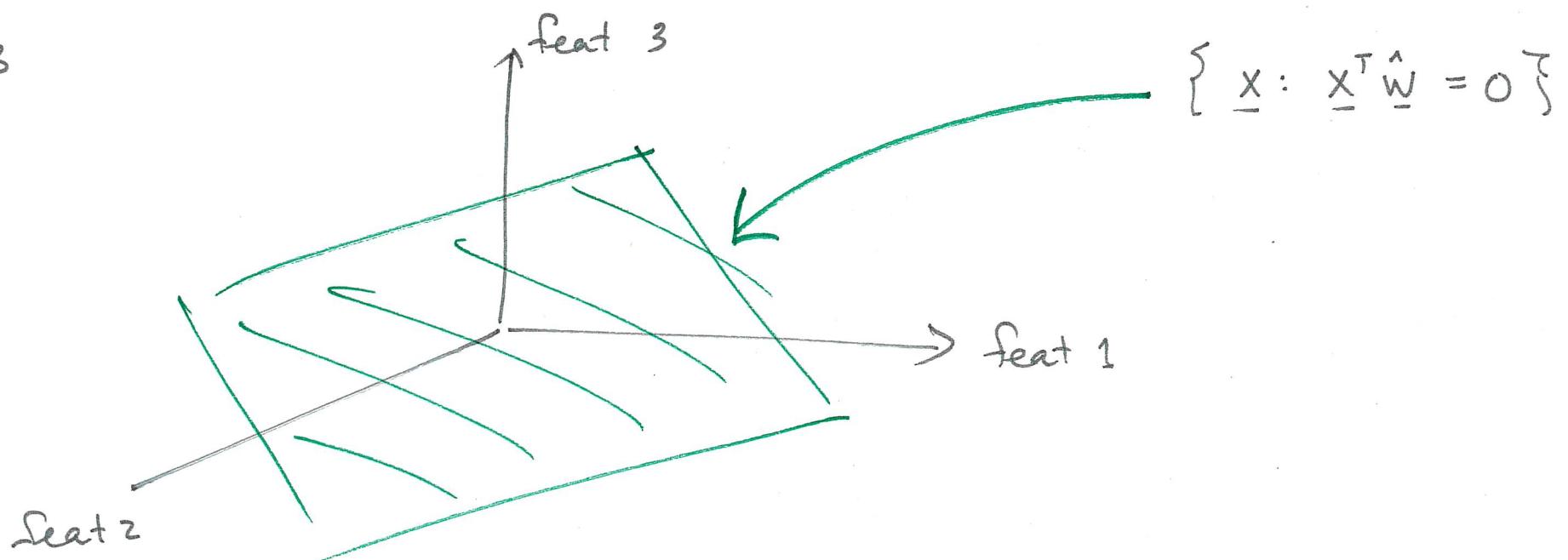
$$y_i = \begin{cases} +1 & \text{class 1} \\ -1 & \text{class 2} \end{cases}$$

Classification rule : $\hat{y}_i = \begin{cases} +1 & \text{if } \hat{y}_i \text{ is closer to } +1 \text{ than } -1 \\ -1 & \text{otherwise} \end{cases}$

$$= \text{sign}(\underline{x}_i^T \hat{\underline{w}})$$

Decision Boundary in Feature Space

$p=3$



What sort of set is $B = \{\underline{x}: \underline{x}^T \hat{w} = 0\}$?

$$\textcircled{1} \quad \underline{x} = 0 \Rightarrow \underline{x} \in B \quad (\underline{0} \in B)$$

$$\textcircled{2} \quad \text{if } \underline{x}_1 \text{ and } \underline{x}_2 \in B, \text{ then } \underline{x}_1 + \underline{x}_2 \in B$$

$$\underline{x}_1^T w = 0 = \underline{x}_2^T w \quad (\underline{x}_1 + \underline{x}_2)^T w = \underline{x}_1^T w + \underline{x}_2^T w = 0 + 0 = 0$$

$$\textcircled{3} \quad \text{if } \underline{x} \in B, \alpha \in \mathbb{R}, \quad \underline{\alpha x} \in B$$

$$\underline{x}^T \hat{w} = 0 \Rightarrow \underline{\alpha x}^T \hat{w} = 0$$

$\boxed{B \text{ is a subspace}}$

(6)

What if we include an offset?

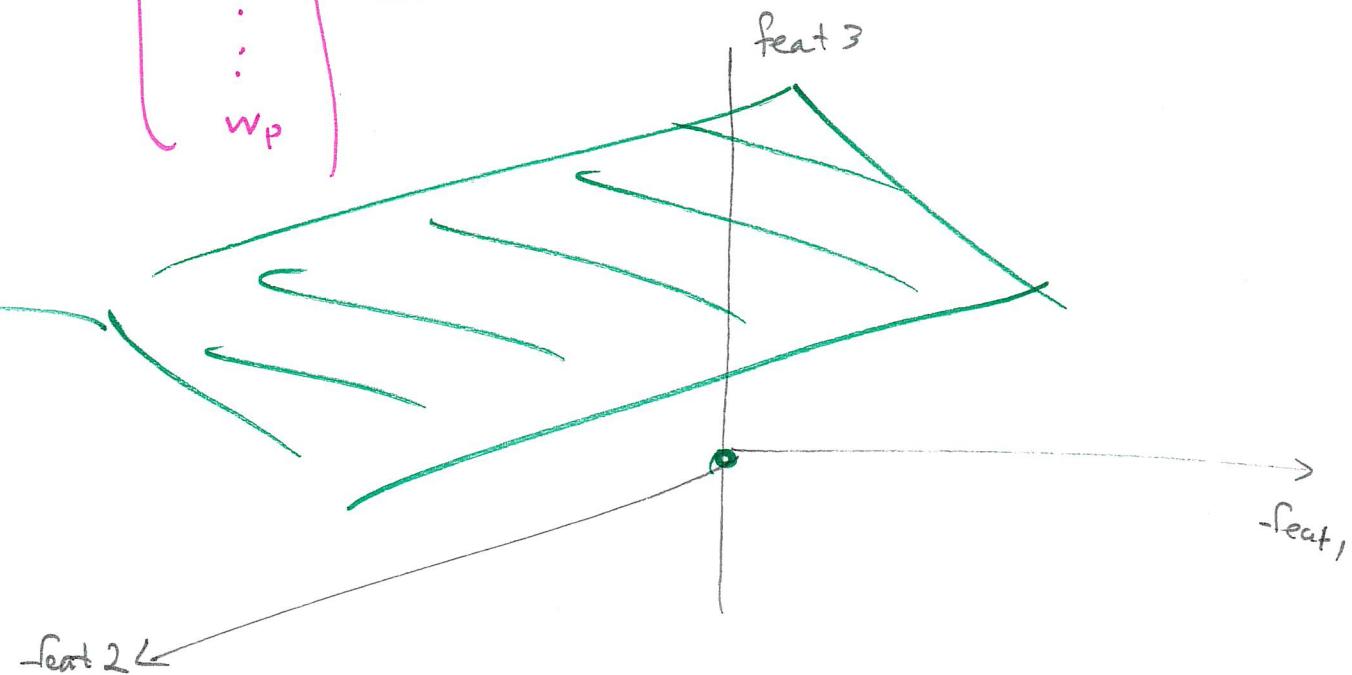
$$\hat{y}_i = w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip} + w_0$$

$$= \underline{w^T x_i} + \underline{w_0} = \underline{\underline{w}}^T \underline{x_i}$$

$$\underline{x_i} = \begin{bmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}$$

$$\underline{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix} \in \mathbb{R}^{P+1}$$

Shifted version
of \mathbf{B}



(7)

$$X = \begin{bmatrix} 98 & 102 & 300 \\ 102 & 98 & 300 \\ 96 & 104 & 100 \\ 104 & 96 & 100 \end{bmatrix}, \quad \underline{y} = \begin{bmatrix} 350 \\ 300 \\ 250 \\ 250 \end{bmatrix}$$

↑ ↑ ↑
 left eye to right eye mouth
 lip dist. to lip dist. width

cols of X are linearly independent $\Rightarrow X^T X$ invertible

$$\hat{\underline{w}} = (X^T X)^{-1} X^T \underline{y} = \begin{bmatrix} 1 \\ 1 \\ \frac{1}{2} \end{bmatrix}$$

$\Rightarrow U$ is an orthonormal basis matrix $\Rightarrow U^T U = I$

$$U = \begin{bmatrix} \frac{1}{2} & -\frac{1}{\sqrt{10}} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{\sqrt{10}} & \frac{1}{2} \\ \frac{1}{2} & -\frac{2}{\sqrt{10}} & -\frac{1}{2} \\ \frac{1}{2} & \frac{2}{\sqrt{10}} & -\frac{1}{2} \end{bmatrix}, \quad \begin{aligned} \hat{\underline{v}} &= (U^T U)^{-1} U^T \underline{y} \\ &= U^T \underline{y} = \begin{bmatrix} 600 \\ 0 \\ 100 \end{bmatrix} \end{aligned}$$

$$U = XA, \quad A = \begin{bmatrix} \frac{1}{400} & \frac{1}{4\sqrt{10}} & -\frac{1}{200} \\ \frac{1}{400} & -\frac{1}{4\sqrt{10}} & -\frac{1}{200} \\ 0 & 0 & \frac{1}{200} \end{bmatrix}$$

$$\underline{u}_1 = \frac{x_1 + x_2}{400} \approx \text{avg of left + right eye to lip distance}$$

$$\underline{u}_2 = \frac{x_1 - x_2}{4\sqrt{10}} \approx \text{diff between left \& right}$$

$$\underline{u}_3 \approx \text{mouth width} - \text{avg of L + R eye to lip dist}$$

$$= \frac{x_3 - (x_1 + x_2)}{200}$$

(7)

w/ X , get weight vect. $X \hat{w}$ good predictions

suggests that left and right features equally important

w/ U , clear that avg. of L+R important

diff between L+R unimportant.

but u_3 less interpretable.

| Q: how should we represent data / features ? |

Regularization

$$\text{LS: } \hat{\underline{w}} = \underset{\underline{w}}{\operatorname{arg\,min}} \|y - X\underline{w}\|_2^2$$

what if $X = \boxed{\quad}$ (short & fat)
 $n < p$

fewer equations than unknowns.

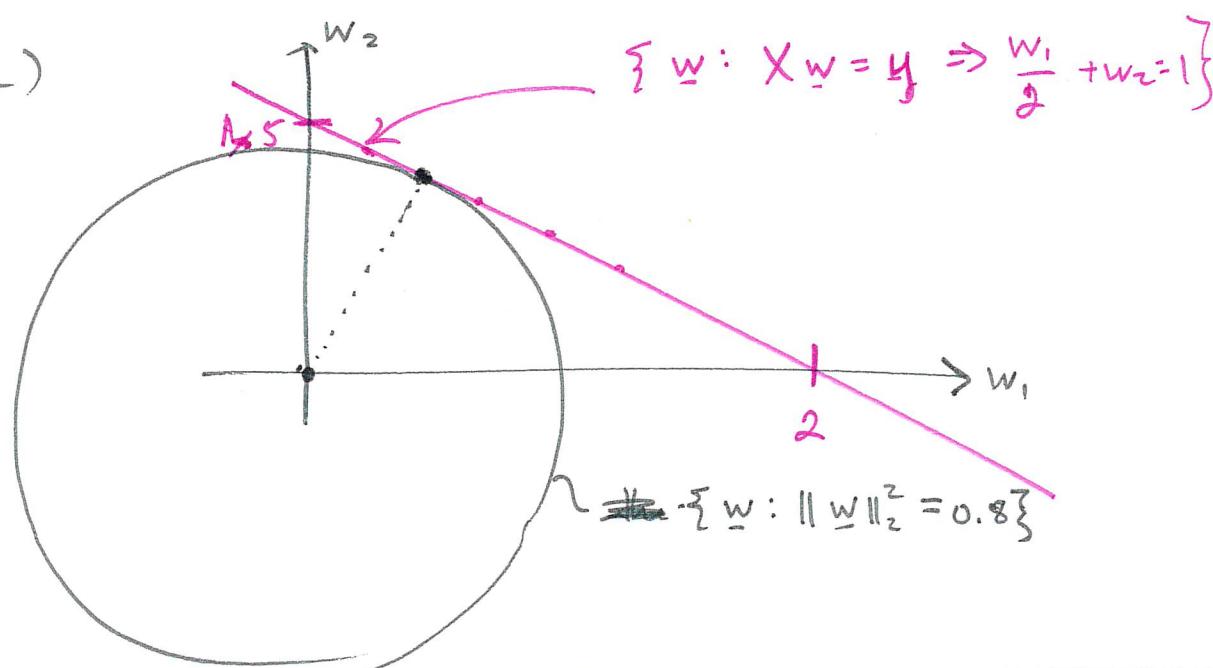
$X^T X$ is not invertible

one approach: find \underline{w} that has small norm $\|\underline{w}\|_2^2$

$$X = \begin{bmatrix} 1/2 & 1 \end{bmatrix} \quad (n=1, p=2)$$

$$y = 1$$

$$\frac{w_1}{2} + w_2 = 1$$



(11)

minimize $\|\underline{w}\|_2^2$ subject to $\|\underline{y} - \underline{X}\underline{w}\|_2^2$

alt: find $\hat{\underline{w}}$ so that $\|\underline{y} - \underline{X}\underline{w}\|_2^2$ is small and $\|\underline{w}\|_2^2$ small

$$\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{argmin}} \|\underline{y} - \underline{X}\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2 \quad \text{for } \lambda > 0$$

λ tells us relative importance of two terms.

$\lambda \|\underline{w}\|_2^2$ = "regularizer"

~~Tikhonov~~ Tikhonov regularization

Ridge regression

$$\|y - \underline{x}_w\|_2^2 + \lambda \|\underline{w}\|_2^2 = \left\| \begin{bmatrix} y - \underline{x}_w \\ \underline{w} \end{bmatrix} \right\|_2^2$$

$$= \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} \underline{x}_w \\ \cancel{\lambda \underline{w}} \end{bmatrix} \right\|_2^2 = \|\tilde{y} - \tilde{\underline{x}}_w\|_2^2$$

$\underbrace{\tilde{y}}$ $\begin{bmatrix} \underline{x} \\ \cancel{\lambda \underline{w}} \end{bmatrix}$
 $\tilde{\underline{x}}$

$$\hat{w} = \underset{w}{\arg \min} \quad \|\tilde{y} - \tilde{\underline{x}}_w\|_2^2 = \underbrace{(\tilde{\underline{x}}^\top \tilde{\underline{x}})}^{-1} \tilde{\underline{x}} \tilde{y}$$

always
 invertible!

(15)

$$\left(\tilde{X}^T \tilde{X} \right)^{-1} \tilde{X}^T \tilde{y} =$$

$$= \left(X^T X + \underline{\lambda I} \right)^{-1} X^T y$$

$$\tilde{X}^T \tilde{X} = \left[X^T \quad \sqrt{\lambda} I \right] \begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix}$$

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