

Lecture 8: Least Squares, Subspaces, & Classification

①

Last time model $\hat{y} = X \underline{w}$ $\xrightarrow{\hspace{2cm}}$ $X \in \mathbb{R}^{n \times p}$

$$\hat{y} = w_1 \underline{x}_1 + w_2 \underline{x}_2 + \dots + w_p \underline{x}_p$$

\uparrow 2nd col of X

$\hat{y} \in \text{Span}\{\underline{x}_1, \dots, \underline{x}_p\}$ ← this a subspace.

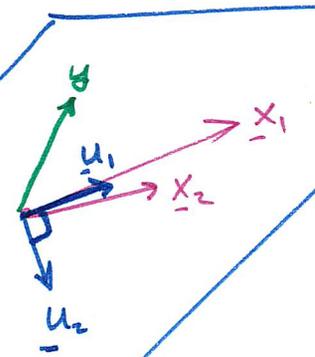
Find an orthonormal basis for subspace (via Gram-Schmidt orthogonalization); let U be matrix w/ basis vectors as columns: $U = [\underline{u}_1, \underline{u}_2, \dots, \underline{u}_r]$. Then $U^T U = I$.

each col of U , u_i is $\in \text{Span}\{\underline{x}_1, \dots, \underline{x}_p\}$

Let $\underline{u}_i = \alpha_{i1} \underline{x}_1 + \alpha_{i2} \underline{x}_2 + \alpha_{i3} \underline{x}_3 + \dots + \alpha_{ip} \underline{x}_p$

or $\underline{U} = \underline{X} \underline{A}$ What i^{th} col of \underline{A} is $\underline{a}_i = \begin{bmatrix} \alpha_{i1} \\ \alpha_{i2} \\ \vdots \\ \alpha_{ip} \end{bmatrix} \iff \underline{u}_i = \underline{X} \underline{a}_i$

$n \times r$ $n \times p$ $p \times r$



$$\hat{\underline{y}} = w_1 \underline{x}_1 + w_2 \underline{x}_2 + \dots + w_p \underline{x}_p$$

$$= v_1 \underline{u}_1 + v_2 \underline{u}_2 + \dots + v_r \underline{u}_r$$

$$= v_1 \underline{X} \underline{a}_1 + v_2 \underline{X} \underline{a}_2 + \dots + v_r \underline{X} \underline{a}_r$$

$$= v_1 (\alpha_{11} \underline{x}_1 + \alpha_{12} \underline{x}_2 + \dots + \alpha_{1p} \underline{x}_p)$$

$$+ v_2 (\alpha_{21} \underline{x}_1 + \alpha_{22} \underline{x}_2 + \dots + \alpha_{2p} \underline{x}_p)$$

$$\vdots$$

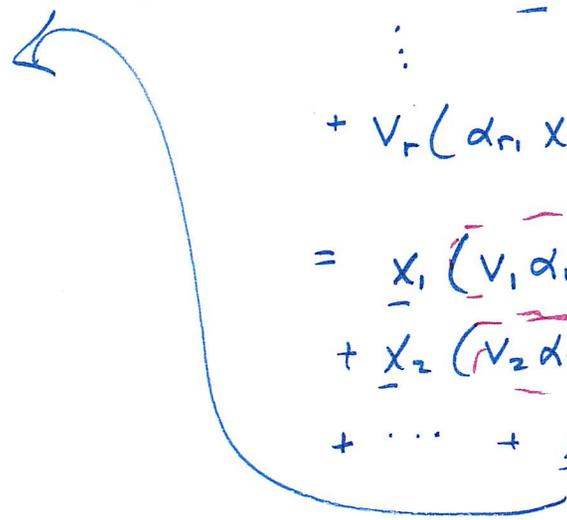
$$+ v_r (\alpha_{r1} \underline{x}_1 + \alpha_{r2} \underline{x}_2 + \dots + \alpha_{rp} \underline{x}_p)$$

$$= \underline{x}_1 (v_1 \alpha_{11} + v_2 \alpha_{21} + \dots + v_r \alpha_{r1})$$

$$+ \underline{x}_2 (v_1 \alpha_{12} + v_2 \alpha_{22} + \dots + v_r \alpha_{r2})$$

$$+ \dots + \underline{x}_p (v_1 \alpha_{1p} + v_2 \alpha_{2p} + \dots + v_r \alpha_{rp})$$

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



Given new sample $x_{new} \in \mathbb{R}^p$, how to predict label y_{new} ?

(A) $\hat{y}_{new} = \langle \underline{x}_{new}, \hat{\underline{w}} \rangle$

(B) ~~is~~ what if using orthonormal basis U ?

recall $\underset{n \times p}{U} = \underset{n \times p}{X} \underset{p \times r}{A} \implies \underline{u}_{new}^T = \underline{x}_{new}^T A$
or $\underline{u}_{new} = A \underline{x}_{new}$

$\implies \hat{y}_{new} = \langle \underline{u}_{new}, \hat{\underline{v}} \rangle$

So what is A ? if cols of X are linearly independent ($r = p$)

Recall $\underline{u}_i = X \underline{a}_i$. To get \underline{a}_i , use LS: $\underline{a}_i = (X^T X)^{-1} X^T \underline{u}_i$

Prediction using orthonormal basis: Given $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$

① compute basis U from X (Gram Schmidt)

② $\hat{\underline{v}} = \underset{\underline{v}}{\operatorname{argmin}} \|y - U \underline{v}\|_2^2 = U^T y$

③ out of sample prediction: $x_{new} \rightarrow u_{new}$; $\hat{y}_{new} = \langle \underline{u}_{new}, \hat{\underline{v}} \rangle$

Theorem: Let $X \in \mathbb{R}^{n \times p}$, $n \geq p$, X is full rank
(p cols are lin. indep)

(4)

and let $y \in \mathbb{R}^n$
let u_1, \dots, u_p be orthonormal basis vectors

for $\text{span}\{x_1, \dots, x_p\} = \text{span}\{u_1, \dots, u_p\}$

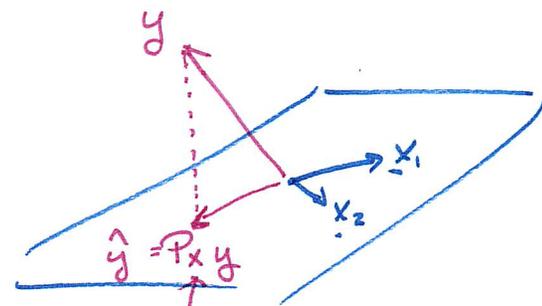
Then $\hat{y} = X \hat{w}$ where $\hat{w} = \underset{w}{\text{argmin}} \|y - Xw\|_2^2$

is given by $\hat{y} = UU^T y$ where $U = [u_1, \dots, u_p]$

note: because U is an orthogonal matrix, $U^T U = I \neq U U^T$

Proof:

$$\hat{y} = X \hat{w} = \underbrace{X (X^T X)^{-1} X^T}_{\text{"projection matrix"} P_X} y$$



projection of y onto
subspace spanned by
cols of X

recall $\text{span}\{\text{cols of } X\} = \text{span}\{u_1, \dots, u_p\}$

$$\Rightarrow P_X y = P_U y$$

$$P_x = P_v = U \underbrace{(U^T U)^{-1}}_{=I} U^T$$

$$= UU^T$$

$$\hat{y} = P_x y = P_v y = UU^T y$$

Least Squares & Classification

Setup: n training samples $(\underline{x}_i, y_i) \in \mathbb{R}^p \times \{-1, +1\}$ for $i=1, \dots, n$

$$y \approx Xw \quad X = \begin{bmatrix} - \underline{x}_1^T - \\ - \underline{x}_2^T - \\ \vdots \\ - \underline{x}_n^T - \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \{-1, +1\}^n$$

\uparrow $\underline{x}_i \in \mathbb{R}^p$ \uparrow $y_i \in \{-1, +1\}$

Assume $n \geq p$, X full rank

$$\hat{\underline{w}} = \arg \min_{\underline{w}} \|y - X\underline{w}\|_2^2 = (X^T X)^{-1} X^T y$$

\hat{y} = predicted labels on training sample

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

Predictions

given new sample $\underline{x}_{\text{new}} \in \mathbb{R}^p$, want to predict y_{new} .

$$\textcircled{A} \quad \hat{y}_{\text{new}} = \underline{x}_{\text{new}}^T \hat{\underline{w}} \in \mathbb{R}$$

$$\textcircled{B} \quad \hat{y}_{\text{new}} = \text{sign}(\underline{x}_{\text{new}}^T \hat{\underline{w}}) \in \{+1, -1\}$$

} "linear classifier"

Classification rule: predict $+1$ if $\underline{x}_{\text{new}}^T \underline{w} > 0$
 -1 if $\underline{x}_{\text{new}}^T \underline{w} < 0$