

CS189 HW2

Hanze Yao

3033093286

1. Linear Regression: Projection and Pseudoinverse

(a) Proof: $\|y - u\|^2 = \|y - P_X(y) + P_X(y) - u\|^2$

$$= \|y - P_X(y)\|^2 + \|P_X(y) - u\|^2$$

$$\geq \|y - P_X(y)\|^2$$

$$\therefore P_X(y) = \underset{u \in \text{range}(X)}{\text{argmin}} \|y - u\|^2$$

(b) First, Proof: P is a rank- d orthogonal projection matrix \Rightarrow

there exists
a $U \in \mathbb{R}^{n \times d}$
 $P = U U^T$
 $U^T U = I$

$\therefore P$ is a rank- d orthogonal projection matrix

$$\text{rank}(P) = d \quad P = P^T \quad P^2 = P$$

$$\therefore P = V \Sigma V^T$$

$$\therefore P^T P = (V \Sigma V^T)^T (V \Sigma V^T) = V \Sigma V^T V \Sigma V^T = V \Sigma \Sigma^T V^T$$

$$= P = V \Sigma V^T$$

$\therefore \Sigma^2 = \Sigma \quad \therefore \Sigma$ is a identity matrix

$$\therefore P^T P = V V^T$$

$$\therefore U = V$$

$$\therefore V^T V = I \quad \therefore U^T U = I$$

Second proof: there exists a $U \in \mathbb{R}^{n \times d}$ $\Rightarrow P$ is a rank- d orthogonal projection matrix

$$\therefore U \in \mathbb{R}^{n \times d} \quad U^T U \in \mathbb{R}^{d \times d} = I \quad \therefore \text{rank}(U) = d$$

$$\therefore P = U U^T \quad \therefore P^T = (U U^T)^T = (U^T)^T U = U U^T = P$$

$$P^2 = U U^T U U^T = U U^T = P$$

$$\therefore \text{rank}(U) = d \quad \therefore \text{rank}(U^T) = d \quad \therefore \text{rank}(U U^T) = \text{rank}(P) = d$$

$\therefore P$ is a rank- d orthogonal projection matrix

(c) Proof: $\text{tr}(P) = \text{tr}(U U^T) = \text{tr}(U^T U) = \text{tr}(I) = d$

$$\therefore U \in \mathbb{R}^{n \times d} \quad U^T U \in \mathbb{R}^{d \times d} \quad \therefore \text{tr}(I) = d = \text{tr}(P)$$

(d) Proof: $\therefore X \in \mathbb{R}^{n \times d} \quad \text{rank}(X) = d \quad n > d \quad \therefore \text{rank}(X^T X) = d \quad \text{rank}(X^T X) = d$

$$\text{rank}((X^T X)^+) = d \quad \therefore \text{rank}(X(X^T X)^+ X^T) = d$$

$$\therefore (X(X^T X)^+ X^T)^T = X(X^T X)^+ X^T$$

$\therefore (X^T X)^T = X^T X \quad \therefore X^T X$ is symmetric $\therefore (X^T X)^+$ is also symmetric since $X^T X$ has full rank (d).

$$\therefore (X(X^T X)^+ X^T)^T = X(X^T X)^+ X^T$$

$$\therefore (X(X^T X)^+ X^T)^2 = X(X^T X)^+ X^T X(X^T X)^+ X^T = X(X^T X)^+ X^T$$

$\therefore X(X^T X)^T X^T$ is a rank- d orthogonal projection matrix

$$\therefore UU^T = X(X^T X)^T X^T, U^T U = I$$

$$\therefore UU^T U = X(X^T X)^T X^T U$$

$$U = X(X^T X)^T X^T U$$

$$U = X(X^T X)^T X^T$$

(e)

$$(a) \therefore X = \sum_{i=1}^{\min(n,d)} \sigma_i U_i V_i^T$$

$$= \sigma_1 U_1 V_1^T + \sigma_2 U_2 V_2^T + \dots + \sigma_d U_d V_d^T \quad (\text{given } n > d)$$

$\therefore V_1^T, V_2^T, \dots, V_d^T$ are all orthonormal $R^{d \times 1}$ row vectors

$\therefore \sigma_i > 0$
 $\therefore \sigma_i U_i \neq 0$
 \therefore Each row of X is a linear combination of $V_1^T, V_2^T, \dots, V_d^T$

$\therefore \{V_i : \sigma_i > 0\}$ are an orthonormal basis for the row space of X .

$$(b) X^T = \sum_{i=1}^{\min(n,d)} \sigma_i V_i U_i^T$$

For the same reason as above

$\{U_i : \sigma_i > 0\}$ are an orthonormal basis for the row space of X^T

$\therefore \{U_i : \sigma_i > 0\}$ are an orthonormal basis for the column space of X .

$$(f) X^+ = \sum_{i: \sigma_i > 0} \sigma_i^{-1} V_i U_i^T = V \Sigma^+ U^T$$

$$\tilde{P}_X(y) = X X^+ y = U \Sigma V^T V \Sigma^+ U^T y$$

If $\text{rank}(X) = d$, then $V \in R^{d \times d}$ which is full rank.

but $U \in R^{n \times n}$ which is not full rank

$$\therefore \tilde{P}_X(y) = U U^T y$$

If $\text{rank}(X) = d$ then there are multiple solutions for P_X

$\tilde{P}_X(y)$ will be one of them.

If $\text{rank}(X) = d = n$ then there is only one solution for P_X ,
 and $\tilde{P}_X(y) = P_X$

2. The Least Norm Solution

(a) Using Lagrange Multiplier Method

Lagrangian function should be

$$L(w, \lambda) = w^T w + \lambda (X^T X w - X^T y)$$

$$\therefore \frac{\partial}{\partial w} L(w, \lambda) = 2w^T + \lambda X^T X = 0$$

↳ Because the constraint

is w must minimize

$$\|Xw - y\|_2^2 \therefore \nabla_w \|Xw - y\|_2^2 = 0$$

which gives $X^T X w - X^T y = 0$

$$w^T = -\frac{\lambda}{2} X^T X$$

$$w = -X^T \left[X \left(-\frac{\lambda}{2} \right)^T \right]$$

$\therefore \hat{w}_{LN}$ is on the row space of X .

Next, proof of uniqueness.

Assume w_1 and w_2 are two solutions for \hat{w}_{LN}

Then $\|w_1 + \lambda(w_2 - w_1)\|_2^2$ has two minimum values

at $\lambda = 1$ and $\lambda = 0$.

$$\|w_1 + \lambda(w_2 - w_1)\|_2^2$$

$$= \|w_1\|_2^2 + 2\lambda \|w_1\|_2 \|w_2 - w_1\|_2 + \lambda^2 \|w_2 - w_1\|_2^2$$

$$\therefore w_1 \neq w_2 \quad \|w_2 - w_1\|_2^2 \neq 0 \text{ and } > 0$$

$\therefore \|w_1 + \lambda(w_2 - w_1)\|_2^2$ has single minimum value

which contradicts my assumption

\hat{w}_{LN} is unique.

$$(b) X = \sum_{i: \sigma_i > 0} \sigma_i u_i v_i^T \quad X^T = \sum_{i: \sigma_i > 0} \sigma_i v_i u_i^T$$

$$\therefore \tilde{w} = \sum_{i: \sigma_i > 0} \frac{1}{\sigma_i} v_i (u_i^T y) = \sum_{i: \sigma_i > 0} (\sigma_i v_i u_i^T) \times \left(\frac{1}{\sigma_i^2} y \right)$$

$$= X^T (\Sigma^{-2} y)$$

$\therefore \tilde{w}$ is in the row space of X

Then check the optimality.

$$X = U \Sigma V^T \quad \tilde{w} = V \Sigma^{-1} U^T y$$

check $X^T X \tilde{w} = X^T y$

$$(U \Sigma V^T)^T U \Sigma V^T V \Sigma^{-1} U^T y = V \Sigma U^T U \Sigma V^T V \Sigma^{-1} U^T y$$

$$= V \Sigma U^T y = X^T y$$

$$\therefore \tilde{w} = \hat{w}_{LN}$$

(C)

$$(1) X^+ X = V \Sigma^+ U^T U \Sigma V^T = V V^T$$

$\{V_i : \sigma_i > 0\}$ is an orthonormal basis for the row space of X .

$\therefore X^+ X$ is the orthogonal projection matrix onto the row space of X .

$$(2) X = \sum_i \sigma_i U_i V_i^T \quad X^T = \sum_i \sigma_i V_i U_i^T$$

$$X^T X = \sum_i \sigma_i^2 V_i V_i^T$$

$$(X^T X)^+ = \sum_i \frac{1}{\sigma_i^2} V_i V_i^T$$

$$(X^T X)^+ X^T = \sum_i \frac{1}{\sigma_i} V_i V_i^T (\sigma_i V_i U_i^T)$$

$$= \sum_i \frac{1}{\sigma_i} V_i U_i^T = X^+$$

(3) Need to show that $X^+ X w = X^+ y$ when w satisfies the least square optimality conditions.

w satisfies the least square optimality

$$\therefore (X^T X) w = X^T y$$

$$w = (X^T X)^+ X^T y$$

$$\therefore X^+ X w = X^+ X (X^T X)^+ X^T y = X^+ X X^+ y = X^+ y$$

\rightarrow according to C (2)

$$\therefore P_X(w) = X^+ X w = X^+ y$$

(4) \therefore According to Z(1), $X^+ y$ is a orthogonal projection onto the row space of X .

According to Z(3), $X^+ y$ satisfies the least square optimality

$\therefore X^+ y = \hat{w}_{LS}$ is consistent with previous part of the problem

$$V^T V = I \\ V^{-1} = V^T$$

3. The Ridge Regression Estimator

$$(a) \hat{W}_\lambda = (X^T X + \lambda I)^{-1} X^T Y$$

$\therefore \lambda I$ is full rank $\therefore X^T X + \lambda I$ is also full rank

$\therefore (X^T X + \lambda I)$ is invertible $\therefore \hat{W}_\lambda$ is unique.

$$X = USV^T \quad X^T = V^T U^T S$$

$$X^T X = V^T U^T U S V^T = V^T S^2 V^T$$

$$\begin{aligned} \therefore \hat{W}_\lambda &= (V^T S^2 V^T + \lambda I)^{-1} V^T U^T Y \\ &= (V^T) (V^T S^2 V^T + V^T \lambda I V)^{-1} (V^T)^{-1} V^T U^T Y \\ &= (V^T) (S^2 + \lambda I)^{-1} V V^T U^T Y \\ &= \frac{S}{S^2 + \lambda I} V U^T Y \\ &= \sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + \lambda} V_i \langle U_i, Y \rangle \end{aligned}$$

$$\begin{aligned} (b) \|\hat{W}_\lambda\|^2 &= \hat{W}_\lambda^T \hat{W}_\lambda = \sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + \lambda} V_i^T \langle U_i, Y \rangle \times \underbrace{\frac{\sigma_i}{\sigma_i^2 + \lambda} V_i \langle U_i, Y \rangle}_{\text{is a scalar}} \\ &= \sum_{i=1}^d \left(\frac{\sigma_i}{\sigma_i^2 + \lambda} \right)^2 V_i^T V_i \langle U_i, Y \rangle^2 \\ &= \sum_{i=1}^d \left(\frac{\sigma_i}{\sigma_i^2 + \lambda} \right)^2 \langle U_i, Y \rangle^2 \end{aligned}$$

$$(c) \therefore \hat{W}_{\lambda \rightarrow \infty} = 0 \quad \therefore \sum_{i: \sigma_i > 0} \sigma_i^{-1} \langle U_i, Y \rangle V_i = 0$$

$\therefore \sigma_i > 0$, $\langle U_i, Y \rangle$ is a scalar

$$\therefore \langle U_i, Y \rangle = 0$$

$$\therefore \hat{W}_\lambda = \sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + \lambda} V_i \underbrace{\langle U_i, Y \rangle}_0 = 0 \text{ for all } \lambda$$

(d) If $\hat{W}_{\lambda \rightarrow \infty} \neq 0$, then $\langle U_i, Y \rangle \neq 0$, then $\hat{W}_\lambda \neq 0$

\therefore as λ increases and $\sigma_i > 0$, $\frac{\sigma_i}{\sigma_i^2 + \lambda}$ strictly decreases

$$\therefore \lambda \in (0, \infty) \text{ and } \sigma_i > 0 \quad \therefore \frac{\sigma_i}{\sigma_i^2 + \lambda} > 0$$

\therefore as λ increases, $\left(\frac{\sigma_i}{\sigma_i^2 + \lambda} \right)^2$ decreases $\Rightarrow \|\hat{W}_\lambda\|^2$ decreases

$$\because \hat{W}_\lambda \neq 0 \quad \therefore \|\hat{W}_\lambda\|^2 > 0$$

\therefore map $\lambda \mapsto \|\hat{W}_\lambda\|^2$ is strictly decreasing and positive on $(0, \infty)$

(e) As $\lambda \rightarrow 0$

$$\hat{W}_\lambda = \sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + \lambda} V_i \langle U_i, y \rangle$$

$$= \sum_{i=1}^d \frac{1}{\sigma_i} V_i \langle U_i, y \rangle$$

$$= \sum_{i=1}^d \frac{1}{\sigma_i} V_i U_i^T y = X^T y = \hat{W}_{\text{LS}}$$

$$\therefore \lim_{\lambda \rightarrow 0} \hat{W}_\lambda = \hat{W}_{\text{LS}}$$

(f) Because of part (d), as λ increases, $\|\hat{W}_\lambda\|^2$ will decrease. Therefore, giving λ a reasonable value can control the value of $\|W\|^2$, which controls the complexity of the model.

4. Patrick vs Alvin

(a) For Patrick's problem,

x_1, x_2, \dots, x_n are inputs, which is X

y_1, y_2, \dots, y_n are observations with noise, which is

$$Y = f(X) + N$$

Now we want to fit a degree d polynomial to this data so this is a linear regression problem with feature matrix as

$$\begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ 1 & x_2 & x_2^2 & \dots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{pmatrix}$$

(b) Seen in the code

(c) The average training error decrease suddenly as the increase of d at first, then the average training error starts to decrease slowly as the d increases. (There is a slight increase from 18 to 19)

This is because as the increase of d , the model becomes more complex and thus fits more closely to the data (with noise) and thus the average training error will decrease.

If I try to fit a polynomial of degree n with a standard matrix inversion method, in the plot, the training error will suddenly increase from degree 18 to degree 19.

(d) For fresh error, after degree 4, the error starts to increase as degree d increases.

I think this is because, as d increases, Patrick starts to fit the noise in the data. As a result, overfit problem will apply, the fresh error (which is like test error) will increase.

(e) With these 2 plots, I think using degree 4 polynomial may be a good choice.

(f) Not enough space. To next page.

According to the first part the best degree is 4 (has the smallest validation error at $\lambda = 0.1$)

(f) According to the output, the best degree is 4, best lambda is 0.15

With $\lambda = 0.1$, select the degree with the lowest Average validation error. Model has degree 4

(g) The degree of p is 4, the best lambda is 0.15

For the degree the best λ and λ are 4 and 0.15