

### 3 Variable Selection

3.1

$S \rightarrow$  fixed integer  $\rightarrow$  most relevant variables.

$X_{\cdot j} \rightarrow$  all the entries of  $\text{any } j^{\text{th}}$  column

Correlation  $c_j = \text{Cor}(X_{\cdot j}, Y)$

Variables such that  $\Rightarrow$  ① Correlations are largest in absolute value;  
 $|S| \neq S \& \text{② } |c_j| > |c_i| \text{ for any } i \in S, j \notin S$ .

$\Rightarrow$  Correlation  $\rightarrow$  is a statistical measure that expresses till what extent two variables are linearly related or degree to which a pair of variables are linearly related.

But selecting a variable simply because its correlation with target is higher does not justify that they are actually good variables for our model.

For example, variables like house price as Y can be dependent on features such as house area, location, etc. but does not have any direct relation with the car parked in the garage. However when we calculate the correlation it may show otherwise, so it doesn't appear to be a good procedure for variable selection.

~~There also~~ A spurious relationship occurs in such variables that shows association but not causally related, may occur due to either coincidence or presence of a  $3^{\text{rd}}$  unseen factor.

It is important that we consider advanced techniques for variable selection like Lasso, as it has model selection consistency for highly correlated data.

3.2 Let  $u \rightarrow p$ -dimensional vector  $\notin \mathbb{R}^{(0,0)}$   
 $\underset{U \in \mathbb{R}^p}{\text{min}} \|u - v\|_2^2 + \lambda \|u\|_2^2 \quad u \in \mathbb{R}^p$

Q2 This function is similar to the ridge regression loss function.

Solution using multivariate calculus is as follows:-

$$U^* = (I_p + \lambda J_p)^{-1} U \quad I_p = \text{Identity matrix}$$

$\text{Since } (I_p + \lambda I_p)^{-1}$  is always invertible &  $\lambda \in [0, \infty)$  which shows that the solution is a non-zero positive & solution always exist.

If  $U \neq 0$  then  $(I_p + \lambda I_p)^{-1} U \neq 0$ , our  $U \neq 0$ .  
So for any  $j \in \mathbb{R}^p$   $U_j \neq 0$ .

4.2

We use beta\_refit instead of beta\_Lasso.  
As Lasso is basically a technique for variable selection and does not have a closed form. Our goal using beta\_Lasso is to only select variables which would truly impact the output & push the other variables (or suppress) to achieve 0 by pushing their  $\beta$ 's or co-efficients to 0.

On the other hand beta\_refit is using ridge regression formula  $\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I_p)^{-1} X^T Y$ . Hence

we have a closed form & we can use these coefficients to actually estimate  $Y$  & get the test error using mean squared error.

5.1

It is an incorrect assumption that we will always find  $X\beta = Y$  as the error term defined for this expression, i.e.  $\|X\beta - Y\|_2^2$  is not actually always 0, a lot of times it is tending to 0 & fit a curve & goal to minimize this error but only after mistakes of data has achieved the error as 0 but in normal data scenarios we try to minimize the term as much as possible so that  $\|X\beta - Y\|_2^2 \rightarrow 0$  but achieving 0 may be after fetching assumptions even though the parameters  $p, n$  as usually there are instances of unaccounted noise in the data.

Hence by summary,

A  $p$ -dimensional vector  $\beta$ ,  $\beta \in \mathbb{R}^p$  such that

(1)  $x_p \neq 1$  &  $p < n$  does not necessarily have a solution to satisfy these set of equations perfectly.

(2)  $X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{pmatrix} \in \mathbb{R}^{3 \times 2}$  matrix

$$X^T X = X_{1.} X_{1.}^T + X_{2.} X_{2.}^T + X_{3.} X_{3.}^T$$

$$(AB)^T = B^T A^T \therefore (X_{1.} X_{1.}^T)^T = (X_{1.}^T X_{1.})^T$$

Similarly  $X_{2.} X_{2.}^T = (X_{2.}^T X_{2.})^T$   
&  $X_{3.} X_{3.}^T = (X_{3.}^T X_{3.})^T$

$$\text{Now } (X_{1.} X_{1.}^T)^T = \left[ \begin{pmatrix} x_{11} \\ x_{12} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} \end{pmatrix} \right]^T = \begin{bmatrix} x_{11}^2 & x_{11} x_{12} \\ x_{11} x_{12} & x_{12}^2 \end{bmatrix} = \begin{bmatrix} x_{11}^2 & x_{11} x_{12} \\ x_{11} x_{12} & x_{12}^2 \end{bmatrix}$$

$$(X_{2.} X_{2.}^T)^T = \begin{pmatrix} x_{21} \\ x_{22} \end{pmatrix} \begin{pmatrix} x_{21} & x_{22} \end{pmatrix}^T = \begin{pmatrix} x_{21}^2 & x_{21} x_{22} \\ x_{21} x_{22} & x_{22}^2 \end{pmatrix} = \begin{pmatrix} x_{21}^2 & x_{21} x_{22} \\ x_{21} x_{22} & x_{22}^2 \end{pmatrix}$$

&  $(X_{3.} X_{3.}^T)^T = \begin{pmatrix} x_{31} \\ x_{32} \end{pmatrix} \begin{pmatrix} x_{31} & x_{32} \end{pmatrix} = \begin{pmatrix} x_{31}^2 & x_{31} x_{32} \\ x_{31} x_{32} & x_{32}^2 \end{pmatrix} \rightarrow (3)$

Addn (1) & (3)  $\Rightarrow \begin{pmatrix} x_{11}^2 + x_{21}^2 + x_{31}^2 & x_{11} x_{12} + x_{21} x_{22} + x_{31} x_{32} \\ x_{11} x_{12} + x_{21} x_{22} + x_{31} x_{32} & x_{12}^2 + x_{22}^2 + x_{32}^2 \end{pmatrix} \rightarrow (4)$

$$X^T X = \begin{pmatrix} x_{11} & x_{21} & x_{31} \\ x_{12} & x_{22} & x_{32} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{pmatrix} = \begin{pmatrix} x_{11}^2 + x_{21}^2 + x_{31}^2 & x_{11} x_{12} + x_{21} x_{22} + x_{31} x_{32} \\ x_{11} x_{12} + x_{21} x_{22} + x_{31} x_{32} & x_{12}^2 + x_{22}^2 + x_{32}^2 \end{pmatrix}$$

From (4) & (5)  $X^T X = (X_{1.} X_{1.}^T)^T + (X_{2.} X_{2.}^T)^T + (X_{3.} X_{3.}^T)^T$   
 $= X_{1.} X_{1.}^T + X_{2.} X_{2.}^T + X_{3.} X_{3.}^T$  Hence proved.

(Q3)

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} \quad B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix}$$

$$AB = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} \end{pmatrix}$$

$$\begin{aligned} BA &= \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} \\ &= \begin{pmatrix} a_{11}b_{11} + a_{21}b_{12} & a_{12}b_{11} + a_{22}b_{12} & a_{13}b_{11} + a_{23}b_{12} \\ a_{11}b_{21} + a_{21}b_{22} & a_{12}b_{21} + a_{22}b_{22} & a_{13}b_{21} + a_{23}b_{22} \\ a_{11}b_{31} + a_{21}b_{32} & a_{12}b_{31} + a_{22}b_{32} & a_{13}b_{31} + a_{23}b_{32} \end{pmatrix} \end{aligned}$$

$$t_r(AB) = a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} + a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32}$$

$$t_r(BA) = a_{11}b_{11} + a_{21}b_{12} + a_{12}b_{21} + a_{22}b_{22} + a_{31}b_{31} + a_{23}b_{32}$$

Therefore,  $t_r(AB) = -t_r(BA)$

~~5.4 Let  $\hat{p} = x^T(xx^T)^{-1}y$~~

To prove:  $\hat{p} \rightarrow \text{solution for } \min_{f \in \mathbb{R}^p} \|f\|_2$

such that  $X\hat{p} = Y$ ,  $X \in \mathbb{R}^{n \times p}$  where  $p \geq n$ .  
 $\& Xx^T \rightarrow \text{invertible}$

now,  $\alpha \in \mathbb{R}^p$  where  $X\alpha = 0$ , we need  $X\hat{p}\alpha = 0$

$$\begin{aligned} X^T \hat{p} &= X^T x^T (xx^T)^{-1} y \\ &= (X\alpha)^T (xx^T)^{-1} y \\ &= 0^T (xx^T)^{-1} y \quad \because X\alpha = 0 \text{ (given)} \end{aligned}$$

PAGE No.	
DATE	/ /

$\hat{\beta} \Rightarrow$  orthogonal to any vector that satisfies  $x\alpha = 0$   
 & if can be expressed as linear combination of  
 - the rows of  $x$

$X X^T \rightarrow$  invertible so  $\hat{\beta}$  can be expressed as a linear  
 combination of rows of  $x$

$\hat{\beta}$  is the least square estimate to  $x\beta = Y$  i.e. the  
 solution of  $\min_{\beta \in \mathbb{R}^n} \|\beta\|_2^2$  such that  $x\beta = Y$

= ~~homework~~