

$$P(\text{rank } D_{\mathcal{X}}^2 \leq \nu) = \frac{\nu}{1+n_{\nu}}$$

HD 41

	rank						
	1	2	...	ν	...	n_{ν}	$n_{\nu+1}$
prob	$\frac{1}{n_{\nu+1}}$	$\frac{1}{n_{\nu+1}}$...	$\frac{1}{n_{\nu+1}}$...	$\frac{1}{n_{\nu+1}}$	$\frac{1}{n_{\nu+1}}$

discrete uniform

Back to regression

47) Suppose data splitting is used resulting in model $I_H = I$, β_I a x 1,

$n_{\nu} \geq J a$, $J = 10, 20$ or higher. (sparse I)

If the model is good (check with response plot, etc)

$$\sqrt{n_{\nu}} (\hat{\beta}_I - \beta_I) \xrightarrow{D} N_a(0, V_I)$$

and do not need $S \subseteq I$ or

$$\beta_{I_0} = \beta = \beta_F$$

with iid cases) $\beta_{I,OLS} = \left(\sum_{\mathcal{X}_I}^{-1} \sum_{\mathcal{X}_I} \right)$

$\left(\sum_{\mathcal{X}_I, OLS} = \sum_{\mathcal{X}_I} Y \right)$ and

$\beta_{I, OPLS} = \lambda_I \beta_{x_I y}$ with

$$\lambda_I = \frac{\sum_{x_I y}^T \sum_{x_I y}}{\sum_{x_I y}^T \sum_{x_I} \sum_{x_I y}}$$

48) In low dimensions, want to estimate $\beta = \beta_F$ that uses all p predictors, and find predictors that are not needed. In high dimensions, can't estimate β_F well, but with iid cases, want to greatly outperform the null model that uses iid y with no predictors.

nontrivial predictors don't count the constant

49) we will say a fitted or POP model is sparse if a of the predictors are active (β_i , $\hat{\beta}_i$) so a of the predictors have nonzero β_i or $\hat{\beta}_i$, and $n \geq J a$ with $J \geq 10$.

otherwise the model is nonsparse.

A high dimensional pop reg full model is dense or abundant if $n \leq 5p$ and nearly all p predictors are active with non-zero β_i .

50} what $\hat{\beta}_F$, $\hat{\beta}_I$, or $\hat{\beta}_{I,0}$ estimate

low dimensional	data splitting variable selection HD but sparse I	HD error
general $\beta_F(x, y) = \beta_{I,0}(x_I, y)$	$\beta_I(x_I, y)$	$\beta_F = \beta(x, y) = \beta(x_I, y) \sim I, 0$
lassos $\beta_F(x, y) = \beta_{I,0}(x_I, y)$	$\beta_I(x_I, y)$	$\beta_F = \beta(x, y) = \beta_{I,0}(x_I, y)$
OLS $\beta_{I,0}(x_I, y) = \beta_{OLS_F}(x, y)$	$\beta_I(x_I, y)$ or $\sum_{-I} x_{-I} z_{x_I y}$ <small>iid cases</small>	$\beta_F = \beta_{OLS} = \beta_F$ <small>β estimator</small>
OPLS $\beta_{OPLS} = \lambda \sum x y$	$\beta_{I,OPLS} = \lambda I \sum x_I y$	$\beta_F = \beta_{OPLS} = \beta_{OLS}$
MMLE $\beta_{MMLE} = \sum u y$ standardized data	$\beta_{I,MMLE} = \sum u_I y$	$\beta_F = \beta_{MMLE} = \beta_{OLS}$

In low dimensions, $\hat{\beta}_F \xrightarrow{P} \beta_F$
 and $\hat{\beta}_{I,0} \xrightarrow{P} \beta_F$ often hold (OLS, GLMs, PH).
 A common error is to assume $\hat{\beta}_{I,0} \xrightarrow{P} \beta_F$ in high dimensions.

$$1] \text{ Let } \underline{y} = \underline{X}^T \underline{\beta} + e \quad (\text{MLR1}).$$

OLS minimizes the residual sum of squares

$$Q_{\text{OLS}}(\underline{\beta}) = \sum (y_i - x_i^T \underline{\beta})^2$$

$$= (\underline{y} - \underline{X}\underline{\beta})^T (\underline{y} - \underline{X}\underline{\beta}) = \text{RSS}(\underline{\beta}).$$

2] One ridge regression estimator $\hat{\underline{\beta}}_R$

minimizes the criterion

$$Q_R(\underline{\beta}) = \frac{1}{\alpha} \underbrace{(\underline{y} - \underline{X}\underline{\beta})^T (\underline{y} - \underline{X}\underline{\beta})}_{\text{RSS}(\underline{\beta})} + \frac{\lambda \ln}{\alpha} \underbrace{\underline{\beta}^T \underline{\beta}}_{\sum_{i=1}^p \beta_i^2} = \lambda \|\underline{\beta}\|_2^2$$

with $\alpha = 1, 2, n, 2n$ common.

If $\lambda \ln = 0$, $\hat{\underline{\beta}}_R = \hat{\underline{\beta}}_{\text{OLS}}$.

It can be shown that

$$\hat{\underline{\beta}}_R = (\underline{X}^T \underline{X} + \lambda \ln \underline{I}_p)^{-1} \underline{X}^T \underline{y} = (\underline{X}^T \underline{X} + \lambda \ln \underline{I}_p)^{-1} \underline{X}^T \underline{X} \underbrace{(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}}_{\hat{\underline{\beta}}_{\text{OLS}}}$$

HW6

if $(X^T X)^{-1}$ exists.

HD 43

3) $X^T X$ is symmetric and square.

Hence $X^T X \geq 0$ (positive semidefinite)

with eigenvalues $\psi_1 \geq \psi_2 \geq \dots \geq \psi_p \geq 0$.

If $\psi_p = 0$, then $(X^T X)^{-1}$ does not exist.

Let (ψ, \underline{g}) be an eigenvalue eigenvector pair of $X^T X$. Then $(X^T X + \lambda_{in} I_p) \underline{g}$

$$= X^T X \underline{g} + \lambda_{in} \underline{g} = \psi \underline{g} + \lambda_{in} \underline{g} = (\psi + \lambda_{in}) \underline{g}.$$

Hence $(\underbrace{\psi + \lambda_{in}}_{> 0 \text{ if } \lambda_{in} > 0}, \underline{g})$ is an eigenvalue eigenvector

pair of $X^T X + \lambda_{in} I_p > 0$ if $\lambda_{in} > 0$.
positive definite

Hence $(X^T X + \lambda_{in} I_p)^{-1}$ exists if $\lambda_{in} > 0$

$\forall \lambda_{in} > 0$, even if $X^T X$ is singular or illconditioned.

4) The score equations for $QR(\beta)$ give

Hwb $\nabla QR(\beta) \stackrel{set}{=} 0$

like
normal
eq's for
OLS

$$-X^T (y - X\hat{\beta}_R) + \lambda \hat{\beta}_R = \underline{0}$$

claim: $\hat{\beta}_R = X^T \tilde{z}$ for some \tilde{z} .

plug in ↓

$$\text{proof: } -X^T (y - X X^T \tilde{z}) + \lambda X^T \tilde{z} = \underline{0}$$

$$\text{or } -X^T y + X^T X X^T \tilde{z} + \lambda X^T \tilde{z} = \underline{0}$$

$p \times n$

$$\text{or } X^T y = X^T (X X^T + \lambda I_n) \tilde{z} \quad \text{and}$$

$$\tilde{z} = (X X^T + \lambda I_n)^{-1} y \quad \text{works. } \square$$

$$\text{Thus } \hat{\beta}_R = \underbrace{X^T (X X^T + \lambda I_n)^{-1}}_{n \times n \text{ matrix}} y$$

$$\text{and } \hat{\beta}_R = \underbrace{(X^T X + \lambda I_p)^{-1}}_{p \times p \text{ matrix}} X^T y$$

If $n \gg p$, use the 2nd formula.

If $p \gg n$ use the 1st formula.

5] In AD statistics, sometimes the method

can be computed using either a $p \times p$ matrix B or an $n \times n$ matrix A as in 4). PCA and PCR give more examples.

$$6) \text{ By 2) } \hat{\beta}_R = (X^T X + \lambda_n I_p)^{-1} X^T Y \hat{\beta}_{OLS} \\ = A_n \hat{\beta}_{OLS} = B_n \hat{\beta}_{OLS} = \\ \hat{\beta}_{OLS} - \frac{\lambda_n}{n} n (X^T X + \lambda_n I_p)^{-1} \hat{\beta}_{OLS}$$

$$\text{where } B_n = I_p - \lambda_n (X^T X + \lambda_n I_p)^{-1}$$

See 11.5; show $A_n - B_n = 0$.

$$7) \text{ Suppose } \frac{X^T X}{n} \xrightarrow{P} V^{-1} \text{ so}$$

$$\sqrt{n} (\hat{\beta}_{OLS} - \beta) \xrightarrow{D} N_p(0, \sigma^2 V)$$

$$\text{If } \frac{\lambda_n}{n} \rightarrow 0, \text{ then } \frac{X^T X + \lambda_n I_p}{n} \xrightarrow{P} V^{-1} \text{ and}$$

$$n (X^T X + \lambda_n I_p)^{-1} \xrightarrow{P} V.$$

$$\text{Then } A_n = A_{n\lambda} = \left(\frac{X^T X + \lambda_{in} I_p}{n} \right)^{-1} \frac{X^T X}{n} P_\gamma$$

$$V V^{-1} = I_p.$$

8) RR CLT: Assume \mathcal{T} holds with p fixed.
for model $Y = X\beta + e$.

a) If $\frac{\hat{\lambda}_{in}}{\sqrt{n}} \xrightarrow{P} 0$, then $\sqrt{n}(\hat{\beta}_R - \beta) \xrightarrow{D} N_p(0, \sigma^2 V)$
 \parallel
 $\hat{\beta}_{OLS}$

b) If $\frac{\hat{\lambda}_{in}}{\sqrt{n}} \xrightarrow{P} \gamma \geq 0$, then

$$\sqrt{n}(\hat{\beta}_R - \beta) \xrightarrow{D} N_p(-\gamma V \beta, \sigma^2 V)$$

Proof] If $\frac{\hat{\lambda}_{in}}{\sqrt{n}} \xrightarrow{P} \gamma \geq 0$, then by 6],

$$\hat{\beta}_R = \left[I_p - \hat{\lambda}_{in} (X^T X + \hat{\lambda}_{in} I_p)^{-1} \right] \hat{\beta}_{OLS}$$

$$\text{Hence } \sqrt{n}(\hat{\beta}_R - \beta) = \sqrt{n} \left[\underbrace{\hat{\beta}_R - \hat{\beta}_{OLS}}_0 + \underbrace{\hat{\beta}_{OLS} - \beta}_{\parallel \hat{\beta}_{OLS}}$$

$$= \sqrt{n}(\hat{\beta}_{OLS} - \beta) - \sqrt{n} \frac{\hat{\lambda}_{in}}{n} (X^T X + \hat{\lambda}_{in} I_p)^{-1} \hat{\beta}_{OLS}$$

$$\xrightarrow{D} N_p(0, \sigma^2 V) - \gamma V \beta \sim N_p(-\gamma V \beta, \sigma^2 V).$$



The estimators could be poor for β

HD 45

9] Suppose we have J estimators (of β) such as $\hat{\beta}_R(\lambda_1), \dots, \hat{\beta}_R(\lambda_J)$

where ridge regression is computed on a grid of J λ values

$0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_J$ where $\lambda_i = \lambda_{ini}$.

Then k -fold CV (cross validation) is often used to select $\hat{\lambda}_{in} = \lambda_i^*$.

Randomly divide the (training) data into k groups or folds of

approx equal size $n_j \approx \frac{n}{k}$ for $j = 1, \dots, k$. Leave out the 1st fold, fit the method to the $k-1$ remaining folds, then compute some criterion for the 1st fold. Repeat for folds 2, 3, ..., k .

like data splitting but repetition causes dependence

10] For MLR, compute $\hat{y}_i(j)$ for each y_i in fold j left out.

45.5

$$\text{Then } MSE_j = \frac{1}{n_j} \sum_{i=1}^{n_j} (Y_i - \hat{Y}_i(j))^2$$

and the overall criterion is

$$CV_{(k)} = \frac{1}{k} \sum_{j=1}^k MSE_j.$$

If each $n_j = \frac{n}{k}$, then $CV_{(k)} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i(j))^2$.

Pick the model I that minimizes

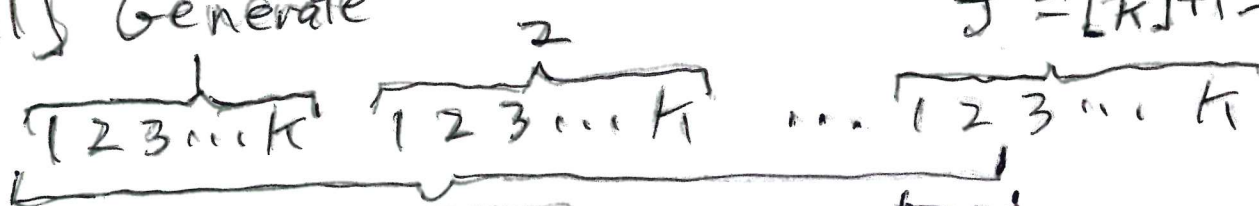
$CV_{(k)}(I)$. Usually $k=10$ or 5 .

If $a_1 \leq a_2 \leq \dots \leq a_J$,

then I_j corresponds to B_{a_j} .

Hence for MLR, k -fold CV picks a model that is good for prediction of future Y_i . (Multitude of models)

ii) Generate



get vector $\begin{matrix} \boxed{k \text{ folds}} \\ n \times 1 \end{matrix}$

fold 1 and 2 get 1 extra case