

LS 87

$$\frac{\lambda_1^T \Sigma_{xy}}{\lambda_1^T \Sigma_{xx} \lambda_1^T} \frac{\lambda_1}{P_{X1}} = \lambda_1 (\lambda_1^T \Sigma_{xx} \lambda_1)^{-1} \lambda_1^T \Sigma_{xy} = \lambda_1^{-1} \Sigma_{xy}$$

for  $\theta = 1, \dots, P$ ,  
 $\lambda_1$  scalar

PLS theory shows  $\lambda_1 = \Sigma_{xy}$ .

$$\text{Hence } \lambda_1 \Sigma_{xy} = \lambda_1^{-1} \Sigma_{xy} = P.$$

eigenvalue  $\uparrow$  eigenvector.

17) There are P+1 PLS estimators

$$(\hat{\alpha}_{\theta, \text{PLS}}, \hat{\beta}_{\theta, \text{PLS}}) \text{ for } \theta = 1, \dots, P \text{ where}$$

$$(\hat{\alpha}_{P, \text{PLS}}, \hat{\beta}_{P, \text{PLS}}) = (\hat{\alpha}_{\text{OLS}}, \hat{\beta}_{\text{OLS}})$$

and  $(\hat{\alpha}_{\theta^*, \text{PLS}}, \hat{\beta}_{\theta^*, \text{PLS}})$  where  $\theta^*$  is

found from a model selection technique,  
such as k-fold cross validation.

Interest is in the one component PLS  
estimator  $\hat{\beta}_{1, \text{PLS}} = \lambda_1 \hat{\Sigma}_{xy} = \lambda_1 \hat{\alpha}_1$ .

$$\hat{\lambda} = \frac{\hat{\Sigma}_{XY}^T \hat{\Sigma}_{XY}}{\hat{\Sigma}_{XY}^T \hat{\Sigma}_X \hat{\Sigma}_{XY}}, \text{ Compute}$$

87.5

$$\hat{\eta} = \hat{\Sigma}_{XY} = \widehat{\text{cov}(X, Y)}. \text{ Then compute}$$

$$w_i = \hat{\eta}^T X_i \quad \text{for } i=1, \dots, n.$$

Fit the working simple linear regression

$$\text{Model } Y_i = \alpha + \lambda w_i + \varepsilon_i \quad i=1, \dots, n$$

with OLS to get  $\hat{\alpha}_{IPLS}$  and  $\hat{\beta}_{IPLS}$ .

$$18) \text{ Let } \hat{\Sigma}_W = \hat{\Sigma}_n. \text{ Let } \hat{\beta}_{OPLS} = \hat{\beta}_{IPLS}.$$

19) Th: Assume  $\hat{\beta}_{IPLS}$  is <sup>gaussian</sup>

$$\sqrt{n} \left[ \left( \begin{matrix} \hat{\lambda} \\ \hat{\eta} \end{matrix} \right) - \left( \begin{matrix} \lambda \\ \eta \end{matrix} \right) \right] \xrightarrow{D} N_{p+1} \left[ \left( \begin{matrix} 0 \\ 0 \end{matrix} \right), \left( \begin{matrix} \hat{\Sigma}_{XX} & \hat{\Sigma}_{XY} \\ \hat{\Sigma}_{YX} & \hat{\Sigma}_{YY} \end{matrix} \right) \right]$$

$$\approx N_{p+1}(\underline{0}, \underline{\Sigma}).$$

$$a) \sqrt{n} (\hat{\underline{\beta}} - \underline{\beta}) \xrightarrow{D} N_p(0, \Sigma_{\beta}) = N_p(0, \Sigma_w)$$

by <sup>arguing</sup> LS theory for  $\hat{\underline{\beta}} = \hat{\Sigma}_{xy}$  given in ex  
on notes ??.

See Exam 3 review.

$$b) \sqrt{n} (\hat{\underline{\alpha}} - \underline{\alpha}) = \sqrt{n} (\hat{\underline{\beta}}_{OLS} - \underline{\beta})$$

$$\xrightarrow{D} N_p(0, D \Sigma_{\beta} D^T) \text{ where by}$$

HW 7 # 1:

$$D = \begin{pmatrix} \frac{\partial \hat{\alpha}_1}{\partial \alpha_1} & \frac{\partial \hat{\alpha}_1}{\partial \alpha_2} & \dots & \frac{\partial \hat{\alpha}_1}{\partial \alpha_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \hat{\alpha}_p}{\partial \alpha_1} & \frac{\partial \hat{\alpha}_p}{\partial \alpha_2} & \dots & \frac{\partial \hat{\alpha}_p}{\partial \alpha_p} \end{pmatrix} = (n \ 2I_p)$$

c) Let  $A$  be a  $K \times P$  constant matrix with full rank  $K$ : ( $1 \leq K \leq P$ ) where  $\underline{AB} = A\underline{B} = \underline{0}$ .

Then  $\underline{\sigma^2}(\underline{AB}_{OLS} - \underline{0}) \rightarrow N_K(\underline{0}, \underline{\sigma^2} \underline{A} \underline{\mathbb{I}_m} \underline{A}^\top)$

See HW7 1b),

d) Estimator:  $\hat{\sigma^2} \underline{A} \underline{\mathbb{I}_m} \underline{A}^\top$

where  $\underline{\mathbb{I}_m} = \underline{\mathbb{I}_w} = \underline{\mathbb{I}_z} = \frac{1}{n-1} \sum_{i=1}^n \underbrace{x_i(y_i - \bar{y})}_{z_i}$

e)  $H_0 \underline{AB} = \underline{0}$  iff  $\underline{H_0 A \mathbb{I}} = \underline{0}$ .

slightly harder test  
if  $n \gg p$

easier test

203 For high dimensional data and big data, often  $\frac{n}{p}$  is not large, e.g.  $p \gg n$ .

with  $\hat{\Sigma} = \widehat{\text{cov}}(\underline{x}, \underline{y})$ , can still do a lot of testing because

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{D} N_p(0, \hat{\Sigma}_w) \text{ and } \hat{\Sigma}_w$$

is not an inverse matrix

(MLEs and OLS estimate  $\hat{\beta}$  with an inverse matrix like  $(\underline{x}^T \underline{x})^{-1}$  or

$$\hat{I}_1^{-1}(\hat{\beta}).)$$

If  $n < sp$  we can't get a good estimator of  $\text{cov}(\hat{\beta}) = \text{cov}(\hat{\Sigma}_{xy}) = \hat{\Sigma}_w$ , but we can get good nonsingular estimators of  $\text{cov}(\hat{\Sigma}_{oy}) = \text{cov}(\hat{m}_{11}, \dots, \hat{m}_{tk})^T$

with  $\underline{v}_i = (x_{i1}, \dots, x_{ik})^T$  where  $k \geq 10t$ :

use the sample covariance matrix of the vectors  $\underline{v}_i(\underline{y}_i - \bar{y})$ . Hence we can test

hypotheses like  $H_0: \beta_i = 0$  or  $H_0: \beta_i - \beta_j = 0$ .

213 Variable Selection is useful

in low and high dimensions.

Let  $\underline{Y} \sim \mathcal{N}(\underline{B}^T \underline{x})$  eg  $\underline{Y} = \underline{x}^T \underline{B} + \epsilon$

or  $\underline{Y} \sim \text{Pois}[\exp(\underline{x}^T \underline{B})]$ . Assume a constant  $B_0$ ,  $x_{0i} \equiv 1$  is always in the model.

A model for variable selection is

$$\underline{x}^T \underline{B} = \underline{x}_S^T \underline{B}_S + \underline{x}_E^T \underline{B}_E = \underline{x}_S^T \underline{B}_S \quad (*)$$

where  $\underline{x} = (\underline{x}_S^T, \underline{x}_E^T)^T$ ,  $\underline{x}_S$  is  $as \times 1$

$\underline{x}_E$  is  $(p-as) \times 1$ . Given  $\underline{x}_S$  is in the model,  $\underline{B}_E = \underline{0}$  and  $E$  denotes the subset of terms in the model that can be eliminated given  $S$  is in the model. Often take  $S$  to be the unique subset of important variables that should be in the model.

Since  $S$  is unknown, candidate LS 90  
 subsets are examined. Let  $\underline{x}_I$  be  
 a candidate  $\alpha \times 1$  vector indexed  
 by  $I$  and let  $\underline{x}_0$  be the vector  
 of predictors out of the model.

$$\text{Then } \underline{x}^T \underline{B} = \underline{x}_I^T \underline{B}_I + \underline{x}_0^T \underline{B}_0.$$

If  $S \subseteq I$  and (\*) holds, then

$$\underline{x}^T \underline{B} = \underline{x}_S^T \underline{B} = \underline{x}_S^T \underline{B}_S + \underline{x}_{I/S}^T \underline{B}_{I/S} + \underline{x}_0^T \underline{\Omega}$$

slash

$= \underline{x}_I^T \underline{B}$  where  $\underline{x}_{I/S}$  denotes  
 the predictors in  $I$  that are not in  $S$ .

$I$  is a submodel ( $Y \perp\!\!\!\perp \underline{x}_I | \underline{x}_I^T \underline{B}$ ), and  
 the full model is a submodel.

22) Forward selection  $I_1$  uses  $\underline{x}_I^* \equiv I = \underline{x}_I$   
 the constant  $I_2$  uses  $\underline{x}_1^*, \underline{x}_2^*$  where  
 $\underline{x}_2^*$  minimizes criterion  $C(\underline{x}_1^* \underline{x}_2^*), j \in \{1, \dots, p\}$ .

90.5

$I_K$  uses  $\underline{x}_{I_K}^* = \underline{x}_{K-1}^*, \underline{x}_K^*$  where  
in model

$\underline{x}_K^*$  minimizes  $C(\underline{x}_{I_K}^*, \dots, \underline{x}_{K-1}^*, \underline{x}_j^*)$

$j \in \{1, \dots, p\} / \{\underline{x}_1^*, \dots, \underline{x}_{K-1}^*\}$ . ( $C = AIC$  or  $BIC$ .)  
Often

Forward Selection has  $r$  models

$I_1, \dots, I_r$ , Often  $r = p$  if  $n \geq 10p$ .

ex)  $P=4$   $x_i$  corresponds to  $B_i$ ,

is always in the model. Suppose

$S = \{1, 2\} = I_2$ . There are  $3 = 2^{P-1} = 8$

possible subsets of  $\{1, \dots, p\}$  that always

contain  $1: \frac{1}{1} \frac{1}{2}, \dots, \frac{1}{p}$ . There

are  $2^{P-|S|}$  subsets that contain  $S$ ,

$\hat{\underline{x}}_{I_7} = (\hat{B}_1, \hat{B}_3, \hat{B}_4)^T$  is obtained by

regressing  $Y$  on  $\underline{x}_{I_7} = (x_1, x_3, x_4)^T$ .

Let  $I_{min}$  correspond to the predictors

selected by the variable selection method such as forward selection or Lasso (later). If  $\hat{\beta}_I$  is a xl, form the pxl vector  $\hat{\beta}_{I,0}$  by adding 0's corresponding to the omitted variables eg  $p=4$   $\hat{\beta}_{I_{\min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$ . Then

$$\hat{\beta}_{Vs} = \hat{\beta}_{I_{\min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$$

As a statistic  $\hat{\beta}_{Vs} = \hat{\beta}_{I_K,0}$  with prob's  $\pi_{K,n} = P(I_{\min} = I_K)$  for  $k=1, \dots, J$

$$\text{eg } J = 2^{p-1}$$

23) A random vector  $\underline{v}$  has a mixture distribution of random vectors  $\underline{v}_j$  with probabilities  $\pi_j$  if  $\underline{v}$  equals the randomly selected vector  $\underline{v}_j$  with prob's  $\pi_j$  for  $j=1, \dots, J$ . (The selection process must not change the dist of the  $\underline{v}_j$ .)

$F_U(t) = \sum_{j=1}^J \pi_j F_{U_j}(t)$  where

$0 \leq \pi_j \leq 1$ ,  $\sum_{j=1}^J \pi_j = 1$ ,  $J \geq 2$  and

$F_{U_j}(t)$  is the cdf of  $U_j$ .

Suppose both  $E(h(U))$  and  $E(h(U_j))$

exist. Then  $E[h(U)] = \sum_{j=1}^J \pi_j E[h(U_j)]$

$$E[U] = \sum_{j=1}^J \pi_j E(U_j)$$

$$\text{cov}(U) = \sum_{j=1}^J \pi_j \text{cov}(U_j) + \sum_{j=1}^J \pi_j E[U_j] E[U_j]^T - E[U] E[U]^T.$$

If  $E(U_j) = \theta$  for  $j=1, \dots, J$ , then

$$E[U] = \theta \text{ and } \text{cov}(U) = \sum_{j=1}^J \pi_j \text{cov}(U_j),$$

24) Law of Total Probability:

Let  $A_1, \dots, A_J$  form a partition of the sample space  $\Omega$  (The  $A_i$  are disjoint,  $P(A_i) > 0$  and  $\bigcup_{i=1}^J A_i = \Omega$ ). Then

$$P(B) = \sum_{k=1}^J P(B \cap A_k) = \sum_{k=1}^J P(B|A_k) P(A_k)$$

LS 92

Variant: if  $P(A_k) = 0$  define

$$P(B|A_k) P(A_k) = 0.$$

25) Let  $w = k$  if  $\hat{B}_{vs} = \hat{B}_{I_{k,0}}$  where

$$P(w=k) = \pi_{kn} \text{ for } k=1, \dots, J.$$

Let  $\hat{B}_{I_{k,0}}^C$  be a random vector from the conditional dist  $\hat{B}_{I_{k,0}} | w=k$ .

Let  $w_n = \sqrt{n}(\hat{B}_{vs} - \hat{B})$  and

$$\underline{w}_{kn} = \sqrt{n}(\hat{B}_{I_{k,0}} - \hat{B}) | w=k \stackrel{D}{=} \sqrt{n}(\hat{B}_{I_{k,0}}^C - \hat{B}).$$

Let col  $F_{\underline{z}}(\underline{t}) = P(z_1 \leq t_1, \dots, z_p \leq t_p)$ .

$$\text{Then } F_{w_n}(\underline{t}) = \sum_{k=1}^J F_{\underline{w}_{kn}}(\underline{t}) \pi_{kn}.$$

Hence  $\hat{B}_{vs}$  has a mixture dist of the  $\hat{B}_{I_{k,0}}^C$  with prob's  $\pi_{kn}$  and

925

$w_n$  has a mixture dist of the  $w_n$  with  
probs  $\pi_{kn}$ .

$$\begin{aligned}
 & \text{Proof} \quad F_{w_n}(t) = P[\overline{\sqrt{n}}(\hat{B}_{vs} - B) \leq t] = \\
 & \sum_{k=1}^J P\left(\overline{\sqrt{n}}(\hat{B}_{vs} - B) \leq t \mid \hat{B}_{vs} = \hat{B}_{I_k}\right) P(\hat{B}_{vs} = \hat{B}_{I_k}) \\
 & = \sum_{k=1}^J P\left[\overline{\sqrt{n}}(\hat{B}_{I_k} - B) \leq t \mid \hat{B}_{vs} = \hat{B}_{I_k}\right] \pi_{kn} \\
 & = \sum_{k=1}^J P\left[\overline{\sqrt{n}}(\hat{B}_{I_k}^c - B) \leq t\right] \pi_{kn} \\
 & = \sum_{k=1}^J F_{w_{kn}}(t) \pi_{kn}. \quad \square
 \end{aligned}$$

law of tot prob

263 Assume that if  $S \subseteq I_j$  where

$\dim(I_j) = a_j$ , then  $\sqrt{n}(\hat{B}_{I_j} - B_{I_j}) \xrightarrow{D} N_{a_j}(0, V_j)$ .

Then  $\sqrt{n}(\hat{B}_{I_{j0}} - B) \xrightarrow{D} N_p(0, V_{j0})$

where  $V_{j0}$  adds rows and columns

of zeroes corresponding to the tree (LS 93)  
 $x_j$  not in  $I_j$ .  $V_{j0}$  is singular unless  
 $I_j$  is the full model.

Let  $\hat{B}_{\text{MIX}}$  have a mixture dist of  
the  $\hat{B}_{I_{k0}}$  with prob's  $\pi_{kn}$  but the  
 $I_k$  are randomly selected.

e.g. in simulation, generate a data set  
perform VS and pick  $I_{k0}$  and  $B_{VS^+}$ , generate  
another data set and use  $I_{k0}$  to get  $\hat{B}_{\text{MIX}}$ .

27) This Assume  $P(S \subseteq I_{\min}) \rightarrow 1$  as  $n \rightarrow \infty$ .

Let  $\hat{B}_{\text{MIX}} = \hat{B}_{I_{k0}}$  with prob's  $\pi_{kn}$  where  
 $\pi_{kn} \rightarrow \pi_k$  as  $n \rightarrow \infty$ . Denote the  
positive  $\pi_k$  by  $\pi_j$ . Assume  
 $U_{jn} = \sqrt{n}(\hat{B}_{I_{j0}} - B) \xrightarrow{D} U_j \sim N_p(0, V_{j0})$ .

a)  $\underline{v}_n = \sqrt{n} (\hat{B}_{\text{mix}} - B) \xrightarrow{D} \underline{v}$  where

the cdf  $F_{\underline{v}}(\underline{z}) = \sum_j \pi_j F_{v_j}(\underline{z})$ .

mixture dist of the  $v_j$

b) Let  $A$  be a full rank  $g \times p$  matrix with  $1 \leq g \leq p$ . Then

$$\underline{v}_n = A \underline{v}_n = \sqrt{n} (A \hat{B}_{\text{mix}} - AB)$$

$\xrightarrow{D} A \underline{v} = \underline{\Sigma}$  where  $\underline{\Sigma}$  has a mixture dist of the  $\underline{v}_j =$

$$A v_j \sim N_g(0, A V_{j0} A^T).$$

c)  $\hat{B}_{\text{VS}}$  is a  $\sqrt{n}$  consistent estimator of  $B$ ;  $\sqrt{n} (\hat{B}_{\text{VS}} - B) = O_p(1)$ .

d) If  $\Pi_d = 1$ , then  $\sqrt{n} (\hat{B}_{\text{SEL}} - B) \xrightarrow{D} N_p(0, V_{d0})$

where  $\underline{S}_{k,n}$  is MIX or VS.

LS94

c) Let  $\underline{B}_{k,n} = \hat{B}_{I_{k,0}}$  with probs  $\pi_{k,n}$ .

Assume  $\underline{w}_{k,n} = \sqrt{n}(\hat{B}_{I_{k,0}}^c - B) \xrightarrow{D} w_j$ .

Then  $\underline{w}_n = \sqrt{n}(\underline{B}_{k,n} - B) \xrightarrow{D} \underline{w}$

where  $F_{\underline{w}}(t) = \sum_{j=1}^J \pi_j F_{w_j}(t)$ .

Thus  $\underline{w}$  is a mixture dist of  
the  $w_j$  with probs  $\pi_j$ .

Proof a) Since  $\underline{v}_n$  has a mixture dist  
of the  $v_{k,n}$  with probs  $\pi_{k,n}$  the  
cdf of  $\underline{v}_n$  is  $\sum_k \pi_{k,n} F_{v_{k,n}}(t) \rightarrow$

$F_{\underline{v}}(t) = \sum_j \pi_j F_{v_j}(t)$  at continuity

points of  $F_{v_j}(t)$  as  $n \rightarrow \infty$ .

- b) Since  $\underline{y}_n \not\rightarrow \underline{y}$ ,  $A\underline{y}_n \not\rightarrow A\underline{y}$ .  
c) Selecting from a finite # of  $\sqrt{n}$  consistent estimators (even on a set that goes to 1 in prob), results in a  $\sqrt{n}$  consistent estimator by Pratt. See HW7 #4.

d) If  $\Pi d = 1$ , there is no selection bias asymptotically.

e) Proof almost the same as a)  $\square$

28) The assumption  $P(S \subseteq I_{\min}) \rightarrow 1$  as  $n \rightarrow \infty$  is important. The assumption has been proved for AIC and BIC for MLR, GLMS, AR(P) time series. A necessary condition for estimators (see lasso and elastic

net to be consistent is LS 95  
 that  $P(S \subseteq I_{\min}) \rightarrow 1$ , so OLS  
 or OPLS after elastic net or lasso  
 satisfies the assumption for MLR  
 under mild conditions.

29] Let  $\underline{X} = (1, \underline{U}^T)^T$  and  $\underline{Y} = \underline{X}\underline{B} + \underline{\epsilon}$   
 MLR. Let  $W$  be the matrix of  
 standardized  $\underline{U}$  so that

$R_U = \frac{W^T W}{n}$  is the sample  
 correlation matrix of the  $\underline{U}$ . Let  
 $\underline{Z} = \underline{Y} - \bar{\underline{Y}}$ . Then regression  
 through the origin is used for  
 $\underline{Z} = W \underline{\beta} + \underline{\epsilon}$ .

$\hat{\underline{Y}} = \bar{\underline{Y}} + \hat{\underline{Z}}$ .  $W$  does not contain a  
 column of 1s.  $\hat{\underline{\beta}}$  can be obtained from  $\hat{\underline{\beta}}$   
 and  $\bar{\underline{Y}}$ .

This method is used so that people who use the same units of measurement (eg could use feet or meters etc) get the same answer from the software. 95.5

30) Assume  $R_U = \frac{W^T W}{n} \xrightarrow{P} V^{-1}$ .

Then  $V^{-1} = S_U = \text{pop corr matrix of the non-trivial predictors } U_i$  if the  $U_i$  are iid. Let  $H = W (W^T W)^{-1} W^T = (h_{ij})$ .

OLS CLT: Assume  $\max_{i=1, \dots, n} h_{ii} \xrightarrow{P} 0$

as  $n \rightarrow \infty$ . Then

$$\sqrt{n} (\hat{\beta}_{OLS} - \beta) \xrightarrow{D} N_{p+1}(\mathbf{0}, \sigma^2 V),$$

↑  
(no longer  $\perp \!\!\! \perp$   $X, Y$ )

31)  $\underline{Y} = X \beta + \underline{\epsilon}$ ,  $\underline{Z} = W \underline{m} + \underline{\epsilon}$