

Math 583 Exam 1 is on Wednesday, Sept. 27 and covers homeworks 1-4 and quizzes 1-4. You are allowed 9 sheets of notes and a calculator. Any needed tables will be provided. CHECK FORMULAS: YOU ARE RESPONSIBLE FOR ANY ERRORS ON THIS HANDOUT!

1) For classical regression and multivariate analysis, we often want  $n \geq Jp$ , with  $J \geq 5$ . Often much larger  $J$  is needed. High dimensional methods have  $n \leq 5p$ .

2) *Regression* investigates how the response variable  $Y$  changes with the value of a  $p \times 1$  vector  $\mathbf{x}$  of predictors. The *estimated sufficient predictor ESP*  $= \hat{\alpha} + \mathbf{x}^T \hat{\boldsymbol{\beta}}$ .

3) A *response plot* is a plot of ESP vs  $Y$  and a *residual plot* is a plot of ESP vs.  $r$ .

4) A plot of  $w$  vs.  $z$  puts  $w$  on the horizontal axis and  $z$  on the vertical axis.

5) A *model for variable selection* is  $\mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_S^T \boldsymbol{\beta}_S + \mathbf{x}_E^T \boldsymbol{\beta}_E = \mathbf{x}_S^T \boldsymbol{\beta}_S$  where  $\mathbf{x} = (\mathbf{x}_S^T, \mathbf{x}_E^T)^T$ ,  $\mathbf{x}_S$  is an  $a_S \times 1$  vector, and  $\mathbf{x}_E$  is a  $(p - a_S) \times 1$  vector. Let  $\mathbf{x}_I$  be the vector of  $a$  terms from a candidate subset indexed by  $I$ , and let  $\mathbf{x}_O$  be the vector of the remaining predictors (out of the candidate submodel). If  $S \subseteq I$ , then  $\mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_S^T \boldsymbol{\beta}_S = \mathbf{x}_S^T \boldsymbol{\beta}_S + \mathbf{x}_{I/S}^T \boldsymbol{\beta}_{(I/S)} + \mathbf{x}_O^T \mathbf{0} = \mathbf{x}_I^T \boldsymbol{\beta}_I$  where  $\mathbf{x}_{I/S}$  denotes the predictors in  $I$  that are not in  $S$ . Since this is true regardless of the values of the predictors,  $\boldsymbol{\beta}_O = \mathbf{0}$  if  $S \subseteq I$ . Note that  $\boldsymbol{\beta}_E = \mathbf{0}$ . Let  $k_S = a_S - 1 =$  the number of population active nontrivial predictors. Then  $k = a - 1$  is the number of active predictors in the candidate submodel  $I$ .

6) For multiple linear regression (MLR), model MLR 1) is

$$Y_i = \beta_1 + x_{i,2}\beta_2 + \cdots + x_{i,p}\beta_p + e_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i$$

for  $i = 1, \dots, n$ . Here  $n$  is the sample size and the random variable  $e_i$  is the  $i$ th error. Assume that the  $e_i$  are iid with expected value  $E(e_i) = 0$  and variance  $V(e_i) = \sigma^2$ . In matrix notation, these  $n$  equations become  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  where  $\mathbf{Y}$  is an  $n \times 1$  vector of dependent variables,  $\mathbf{X}$  is an  $n \times p$  matrix of predictors,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of unknown coefficients, and  $\mathbf{e}$  is an  $n \times 1$  vector of unknown errors.

Model MLR 2) is

$$Y_i = \alpha + x_{i,1}\beta_1 + \cdots + x_{i,p}\beta_p + e_i = \alpha + \mathbf{x}_i^T \boldsymbol{\beta} + e_i$$

for  $i = 1, \dots, n$ . For this model, we may use  $\boldsymbol{\phi} = (\alpha, \boldsymbol{\beta}^T)^T$  with  $\mathbf{Y} = \mathbf{X}\boldsymbol{\phi} + \mathbf{e}$ .

7) For model MLR 1)  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ , the ordinary least squares (OLS) estimator is  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , which exists if  $n > p$  and  $\mathbf{X}$  has full rank  $p$ .

8) For model MLR 2)  $Y_i = \alpha + \mathbf{x}_i^T \boldsymbol{\beta} + e_i$ , let  $\mathbf{x}_i^T = (x_{i1}, \dots, x_{ip})$ , let  $\alpha$  be the intercept, and let the slopes vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ . Let the population covariance matrices

$$\text{Cov}(\mathbf{x}) = E[(\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^T] = \boldsymbol{\Sigma}_{\mathbf{x}}, \quad \text{and}$$

$$\text{Cov}(\mathbf{x}, Y) = E[(\mathbf{x} - E(\mathbf{x}))(Y - E(Y))] = \boldsymbol{\Sigma}_{\mathbf{x}Y} = \boldsymbol{\Sigma}_{Y\mathbf{x}}.$$

If the cases  $(\mathbf{x}_i, Y_i)$  are iid from some population where  $\boldsymbol{\Sigma}_{\mathbf{x}Y}$  exists and  $\boldsymbol{\Sigma}_{\mathbf{x}}$  is nonsingular, then the population coefficients from an OLS regression of  $Y$  on  $\mathbf{x}$  (even if a linear model does not hold) are

$$\alpha = \alpha_{OLS} = E(Y) - \boldsymbol{\beta}^T E(\mathbf{x}) \quad \text{and} \quad \boldsymbol{\beta} = \boldsymbol{\beta}_{OLS} = \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \boldsymbol{\Sigma}_{\mathbf{x}Y}.$$

9) The **sample mean**  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = (\bar{X}_1, \dots, \bar{X}_p)^T$  where  $\bar{X}_i$  is the sample mean of the data in column  $i$  corresponding to variable  $X_i$ . Let the sample covariance matrices be

$$\hat{\Sigma}_{\mathbf{x}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \text{ and } \hat{\Sigma}_{\mathbf{xY}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(Y_i - \bar{Y}).$$

Then  $\Sigma_{\mathbf{x}} = \mathbf{S} = (S_{ij})$ . That is, the  $ij$  entry of  $\mathbf{S}$  is the sample covariance  $S_{ij}$ . Let the method of moments estimators be  $\tilde{\Sigma}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$  and

$$\tilde{\Sigma}_{\mathbf{xY}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(Y_i - \bar{Y}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i Y_i - \bar{\mathbf{x}} \bar{Y}.$$

10) For model MLR 2)  $Y_i = \alpha + \mathbf{x}_i^T \boldsymbol{\beta} + e_i$ , a) If  $\hat{\Sigma}_{\mathbf{x}}^{-1}$  exists, then  $\hat{\alpha} = \bar{Y} - \hat{\boldsymbol{\beta}}^T \bar{\mathbf{x}}$  and

$$\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{OLS} = \tilde{\Sigma}_{\mathbf{x}}^{-1} \tilde{\Sigma}_{\mathbf{xY}} = \hat{\Sigma}_{\mathbf{x}}^{-1} \hat{\Sigma}_{\mathbf{xY}}.$$

b) Suppose that  $(Y_i, \mathbf{x}_i^T)^T$  are iid random vectors such that  $\sigma_Y^2$ ,  $\Sigma_{\mathbf{x}}^{-1}$ , and  $\Sigma_{\mathbf{xY}}$  exist. Then  $\hat{\alpha} \xrightarrow{P} \alpha$  and

$$\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\beta} \text{ as } n \rightarrow \infty$$

where  $\alpha$  and  $\boldsymbol{\beta}$  are given by 8).

11) If  $\mathbf{X}$  and  $\mathbf{Y}$  are  $p \times 1$  random vectors,  $\mathbf{a}$  a conformable constant vector, and  $\mathbf{A}$  and  $\mathbf{B}$  are conformable constant matrices, then

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}), \quad E(\mathbf{a} + \mathbf{Y}) = \mathbf{a} + E(\mathbf{Y}), \quad \& \quad E(\mathbf{A}\mathbf{X}\mathbf{B}) = \mathbf{A}E(\mathbf{X})\mathbf{B}.$$

Also

$$\text{Cov}(\mathbf{a} + \mathbf{A}\mathbf{X}) = \text{Cov}(\mathbf{A}\mathbf{X}) = \mathbf{A}\text{Cov}(\mathbf{X})\mathbf{A}^T.$$

Note that  $E(\mathbf{A}\mathbf{Y}) = \mathbf{A}E(\mathbf{Y})$  and  $\text{Cov}(\mathbf{A}\mathbf{Y}) = \mathbf{A}\text{Cov}(\mathbf{Y})\mathbf{A}^T$ .

12) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $E(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ .

13) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and if  $\mathbf{A}$  is a  $q \times p$  matrix, then  $\mathbf{A}\mathbf{X} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . If  $\mathbf{a}$  is a  $p \times 1$  vector of constants, then  $\mathbf{X} + \mathbf{a} \sim N_p(\boldsymbol{\mu} + \mathbf{a}, \boldsymbol{\Sigma})$ .

14) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and if  $\mathbf{A}$  is a  $q \times p$  matrix, then  $\mathbf{A}\mathbf{X} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . If  $\mathbf{a}$  is a  $p \times 1$  vector of constants, then  $\mathbf{a} + \mathbf{X} \sim N_p(\mathbf{a} + \boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

Suppose

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \text{and} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

15)  $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ .

16) Given a MVN distribution, be able to find the MVN distribution of subsets, pairs of independent random variables and the correlation

$$\rho(X_i, X_j) = \frac{\sigma_{i,j}}{\sqrt{\sigma_{ii}\sigma_{jj}}} = \text{Cov}(X_i, X_j) / \sqrt{V(X_i)V(X_j)}.$$

17) If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the conditional distribution of  $\mathbf{X}_1$  given that  $\mathbf{X}_2 = \mathbf{x}_2$  is multivariate normal with mean  $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$  and covariance  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ . That is,

$$\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

Often  $X_1 = Y$  and  $X_2 = X$ . Then find  $E(Y|X)$ ,  $V(Y|X)$  and  $\rho(Y, X)$ .

18) Know that if  $Y_1, \dots, Y_n$  are iid with  $E(Y) = \mu$  and  $V(Y) = \sigma^2$ , then  $E(\bar{Y}) = \mu$  and  $V(\bar{Y}) = \sigma^2/n$ . Know  $E(S^2) = \sigma^2$ .

19) Let  $\mathbf{X}_n$  be a sequence of random vectors with joint cdfs  $F_n(\mathbf{x})$  and let  $\mathbf{X}$  be a random vector with joint cdf  $F(\mathbf{x})$ .

a)  $\mathbf{X}_n$  converges in distribution to  $\mathbf{X}$ , written  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ , if  $F_n(\mathbf{x}) \rightarrow F(\mathbf{x})$  as  $n \rightarrow \infty$  for all points  $\mathbf{x}$  at which  $F(\mathbf{x})$  is continuous. The distribution of  $\mathbf{X}$  is the limiting distribution or asymptotic distribution of  $\mathbf{X}_n$ .

b)  $\mathbf{X}_n$  converges in probability to  $\mathbf{X}$ , written  $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$ , if for every  $\epsilon > 0$ ,  $P(\|\mathbf{X}_n - \mathbf{X}\| > \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$ .

20) Multivariate Central Limit Theorem (MCLT): If  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are iid  $k \times 1$  random vectors with  $E(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ , then

$$\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$$

where the sample mean

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i.$$

21) Suppose  $\sqrt{n}(T_n - \boldsymbol{\mu}) \xrightarrow{D} N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ . Let  $\mathbf{A}$  be a  $q \times p$  constant matrix. Then  $\mathbf{A}\sqrt{n}(T_n - \boldsymbol{\mu}) = \sqrt{n}(\mathbf{A}T_n - \mathbf{A}\boldsymbol{\mu}) \xrightarrow{D} N_q(\mathbf{A}\boldsymbol{\theta}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

22) Suppose  $\mathbf{A}$  is a conformable constant matrix and  $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ . Then  $\mathbf{A}\mathbf{X}_n \xrightarrow{D} \mathbf{A}\mathbf{X}$ .

23) The behavior of convergence in distribution to a MVN distribution in B) is much like the behavior of the MVN distributions in A). The results in B) can be proven using the multivariate delta method. Let  $\mathbf{A}$  be a  $q \times k$  constant matrix,  $b$  a constant,  $\mathbf{a}$  a  $k \times 1$  constant vector, and  $\mathbf{d}$  a  $q \times 1$  constant vector. Note that  $\mathbf{a} + b\mathbf{X}_n = \mathbf{a} + \mathbf{A}\mathbf{X}_n$  with  $\mathbf{A} = b\mathbf{I}$ . Thus i) and ii) follow from iii).

A) Suppose  $\mathbf{X} \sim N_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then

i)  $\mathbf{A}\mathbf{X} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

ii)  $\mathbf{a} + b\mathbf{X} \sim N_k(\mathbf{a} + b\boldsymbol{\mu}, b^2\boldsymbol{\Sigma})$ .

iii)  $\mathbf{A}\mathbf{X} + \mathbf{d} \sim N_q(\mathbf{A}\boldsymbol{\mu} + \mathbf{d}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

(Find the mean and covariance matrix of the left hand side and plug in those values for the right hand side. **Be careful with the dimension  $k$  or  $q$ .**)

B) Suppose  $\mathbf{X}_n \xrightarrow{D} N_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then

i)  $\mathbf{A}\mathbf{X}_n \xrightarrow{D} N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

ii)  $\mathbf{a} + b\mathbf{X}_n \xrightarrow{D} N_k(\mathbf{a} + b\boldsymbol{\mu}, b^2\boldsymbol{\Sigma})$ .

iii)  $\mathbf{A}\mathbf{X}_n + \mathbf{d} \xrightarrow{D} N_q(\mathbf{A}\boldsymbol{\mu} + \mathbf{d}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .

24)  $\Sigma_{\mathbf{x}, \mathbf{y}} = \text{Cov}(\mathbf{x}, \mathbf{y}) = E[(\mathbf{x} - E(\mathbf{x}))(\mathbf{y} - E(\mathbf{y}))^T]$   
 $\text{Cov}(\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}) = \mathbf{A}\text{Cov}(\mathbf{x}, \mathbf{y})\mathbf{B}^T$   
 $\Sigma_{\mathbf{A}\mathbf{x}, \mathbf{Y}} = \mathbf{A}\Sigma_{\mathbf{x}, \mathbf{Y}}$  (take  $\mathbf{B} = \mathbf{I}_1 = 1$ )  
 $\Sigma_{\mathbf{x}, \mathbf{B}\mathbf{y}} = \Sigma_{\mathbf{x}, \mathbf{y}}\mathbf{B}^T$  (take  $\mathbf{A} = \mathbf{I}$ )  
If  $\mathbf{w}_i = \mathbf{A}\mathbf{x}_i$  for  $i = 1, \dots, n$ , then  
 $\bar{\mathbf{w}} = \mathbf{A}\bar{\mathbf{x}}$ ,  
 $\hat{\Sigma}_{\mathbf{w}} = \mathbf{A}\hat{\Sigma}_{\mathbf{x}}\mathbf{A}^T$ ,  
 $\tilde{\Sigma}_{\mathbf{w}} = \mathbf{A}\tilde{\Sigma}_{\mathbf{x}}\mathbf{A}^T$ ,  
 $\hat{\Sigma}_{\mathbf{w}, \mathbf{Y}} = \mathbf{A}\hat{\Sigma}_{\mathbf{x}, \mathbf{Y}}$ ,  
 $\tilde{\Sigma}_{\mathbf{w}, \mathbf{Y}} = \mathbf{A}\tilde{\Sigma}_{\mathbf{x}, \mathbf{Y}}$ ,  
 $\text{Cov}(\sum_{i=1}^n \mathbf{x}_i, \sum_{j=1}^m \mathbf{z}_j) = \sum_{i=1}^n \sum_{j=1}^m \text{Cov}(\mathbf{x}_i, \mathbf{z}_j)$ ,

25) All subsets of a MVN distribution are MVN. Suppose

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \text{and } \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

If  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$ , and if  $\mathbf{X}_1$  is a  $q \times 1$  vector, then  $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \Sigma_{11})$ .

26) Let  $\mathbf{x}_n = (x_{1n}, \dots, x_{pn})^T$  and  $\mathbf{x} = (x_1, \dots, x_p)^T$  be random vectors. Then  $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$  implies  $x_{in} \xrightarrow{D} x_i$  for  $i = 1, \dots, p$ . Hence all subsets of  $\mathbf{x}_n$  converge in distribution to the corresponding subsets of  $\mathbf{x}$ : if  $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$ , then

$$\begin{pmatrix} x_{i_1, n} \\ \vdots \\ x_{i_k, n} \end{pmatrix} \xrightarrow{D} \begin{pmatrix} x_{i_1} \\ \vdots \\ x_{i_k} \end{pmatrix}.$$

Typically marginal convergence in distribution  $x_{in} \xrightarrow{D} x_i$  for  $i = 1, \dots, p$  **does not imply**  $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$ .

27) The OLS regression of  $Y$  on  $\mathbf{w} = \mathbf{A}\mathbf{x}$ , where  $\mathbf{A}$  is a  $k \times p$  constant matrix with full rank  $k$ , is  $\hat{\boldsymbol{\beta}}_{OLS}(\mathbf{w}, Y) = \hat{\Sigma}_{\mathbf{w}}^{-1} \hat{\Sigma}_{\mathbf{w}, \mathbf{Y}} = (\mathbf{A}\hat{\Sigma}_{\mathbf{x}}\mathbf{A}^T)^{-1} \mathbf{A}\hat{\Sigma}_{\mathbf{x}, \mathbf{Y}}$ , provided the inverse matrices exist.

28) Under the conditions of 27), if the cases  $(\mathbf{x}_i, Y_i)$  are iid, then the population OLS regression of  $Y$  on  $\mathbf{w}$  is  $\boldsymbol{\beta}_{OLS}(\mathbf{w}, Y) = \Sigma_{\mathbf{w}}^{-1} \Sigma_{\mathbf{w}, \mathbf{Y}} = (\mathbf{A}\Sigma_{\mathbf{x}}\mathbf{A}^T)^{-1} \mathbf{A}\Sigma_{\mathbf{x}, \mathbf{Y}}$ , provided the inverse matrices exist.

29) **OLS CLTs.** Consider the MLR model and assume that the zero mean errors are iid with  $E(e_i) = 0$  and  $\text{VAR}(e_i) = \sigma^2$ . If the  $\mathbf{x}_i$  are random vectors, assume that the cases  $(\mathbf{x}_i, Y_i)$  are independent, and that the  $\mathbf{e}_i$  and  $\mathbf{x}_i$  are independent. Also assume that  $\max_i(h_1, \dots, h_n) \rightarrow 0$  and

$$\frac{\mathbf{X}^T \mathbf{X}}{n} \rightarrow \mathbf{V}^{-1}$$

as  $n \rightarrow \infty$  where the convergence is in probability if the  $\mathbf{x}_i$  are random vectors (instead of nonstochastic constant vectors).

a) For MLR model 1)  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ , the OLS estimator  $\hat{\boldsymbol{\beta}}$  satisfies

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

b) For MLR model 2)  $\mathbf{Y} = \mathbf{X}\boldsymbol{\phi} + \mathbf{e}$ , the OLS estimator  $\hat{\boldsymbol{\phi}}$  satisfies

$$\sqrt{n}(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi}) \xrightarrow{D} N_{p+1}(\mathbf{0}, \sigma^2 \mathbf{V}).$$

c) Suppose the cases  $(\mathbf{x}_i, Y_i)$  are iid from some population and MLR model 2)  $Y_i = \alpha + \mathbf{x}_i^T \boldsymbol{\beta} + e_i$  holds. Assume that  $\boldsymbol{\Sigma}_{\mathbf{x}}^{-1}$  and  $\boldsymbol{\Sigma}_{\mathbf{x}, Y}$  exist. Then b) holds and

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}_{\mathbf{x}}^{-1})$$

where  $\boldsymbol{\beta} = \boldsymbol{\beta}_{OLS} = \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \boldsymbol{\Sigma}_{\mathbf{x}, Y}$ .

30) If the  $e_i$  are iid and  $p$  is fixed, then under mild condition  $\hat{\boldsymbol{\beta}}_{OLS} \xrightarrow{P} \boldsymbol{\beta}_{OLS}$ . Need iid cases for  $\boldsymbol{\beta}_{OLS} = \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \boldsymbol{\Sigma}_{\mathbf{x}, Y}$ .

31) The one component partial least squares (OPLS) estimator is easy to compute in high and low dimensions. Compute  $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y} = \hat{\boldsymbol{\eta}}_{OPLS}$ , then compute  $W_i = \hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}^T \mathbf{x}_i$  for  $i = 1, \dots, n$ . Then do the OLS regression of  $Y$  on the  $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}^T \mathbf{x}_i$  using the working model  $Y = \alpha + \lambda W + \epsilon$  to get  $\hat{\alpha}$  and  $\hat{\lambda}$ . Then  $\hat{\alpha}_{OPLS} = \hat{\alpha}$  and

$$\hat{\boldsymbol{\beta}}_{OPLS} = \hat{\lambda} \hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y} \quad \text{where} \quad \hat{\lambda} = \frac{\hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}^T \hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}}{\hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}^T \hat{\boldsymbol{\Sigma}}_{\mathbf{x}} \hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}}.$$

Under iid cases,

$$\boldsymbol{\beta}_{OPLS} = \lambda \boldsymbol{\Sigma}_{\mathbf{x}Y} \quad \text{where} \quad \lambda = \frac{\boldsymbol{\Sigma}_{\mathbf{x}Y}^T \boldsymbol{\Sigma}_{\mathbf{x}Y}}{\boldsymbol{\Sigma}_{\mathbf{x}Y}^T \boldsymbol{\Sigma}_{\mathbf{x}} \boldsymbol{\Sigma}_{\mathbf{x}Y}}$$

for  $\boldsymbol{\Sigma}_{\mathbf{x}Y} \neq \mathbf{0}$ . If  $\boldsymbol{\Sigma}_{\mathbf{x}Y} = \mathbf{0}$ , then  $\boldsymbol{\beta}_{OPLS} = \mathbf{0}$ . The OPLS MLR model is  $Y = Y | \boldsymbol{\beta}_{OPLS}^T \mathbf{x} = \alpha_{OPLS} + \boldsymbol{\beta}_{OPLS}^T \mathbf{x} + e$ .

32) CLT for  $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}Y}$ : Assume the cases  $(\mathbf{x}_i^T, Y_i)^T$  are iid. Assume  $E(x_{ij}^k Y_i^m)$  exist for  $j = 1, \dots, p$  and  $k, m = 0, 1, 2$ . Let  $\boldsymbol{\mu}_{\mathbf{x}} = E(\mathbf{x})$  and  $\mu_Y = E(Y)$ . Let  $\mathbf{w}_i = (\mathbf{x}_i - \boldsymbol{\mu}_{\mathbf{x}})(Y_i - \mu_Y)$  with sample mean  $\bar{\mathbf{w}}_n$ . Let  $\boldsymbol{\eta} = \boldsymbol{\Sigma}_{\mathbf{x}, Y}$ . Then a)

$$\sqrt{n}(\bar{\mathbf{w}}_n - \boldsymbol{\eta}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{w}}), \quad \sqrt{n}(\hat{\boldsymbol{\eta}}_n - \boldsymbol{\eta}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{w}}),$$

$$\text{and} \quad \sqrt{n}(\tilde{\boldsymbol{\eta}}_n - \boldsymbol{\eta}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{w}}).$$

b) Let  $\mathbf{z}_i = \mathbf{x}_i(Y_i - \bar{Y}_n)$  and  $\mathbf{v}_i = (\mathbf{x}_i - \bar{\mathbf{x}}_n)(Y_i - \bar{Y}_n)$ . Then  $\hat{\boldsymbol{\Sigma}}_{\mathbf{w}} = \hat{\boldsymbol{\Sigma}}_{\mathbf{z}} = \hat{\boldsymbol{\Sigma}}_{\mathbf{v}}$ . Hence  $\tilde{\boldsymbol{\Sigma}}_{\mathbf{w}} = \tilde{\boldsymbol{\Sigma}}_{\mathbf{z}} = \tilde{\boldsymbol{\Sigma}}_{\mathbf{v}}$ .

c) Let  $\mathbf{A}$  be a  $k \times p$  full rank constant matrix with  $k \leq p$ , assume  $H_0 : \mathbf{A}\boldsymbol{\beta}_{OPLS} = \mathbf{0}$  is true, and assume  $\hat{\lambda} \xrightarrow{P} \lambda \neq 0$ . Then

$$\sqrt{n}\mathbf{A}(\hat{\boldsymbol{\beta}}_{OPLS} - \boldsymbol{\beta}_{OPLS}) \xrightarrow{D} N_k(\mathbf{0}, \lambda^2 \mathbf{A}\boldsymbol{\Sigma}_{\mathbf{w}}\mathbf{A}^T).$$

33) **Slutsky's Theorem:** Let  $\mathbf{x}_n = (X_{1n}, \dots, X_{kn})^T$  be a sequence of  $k \times 1$  random vectors, let  $\mathbf{y}_n$  be a sequence of  $k \times 1$  random vectors, and let  $\mathbf{x} = (X_1, \dots, X_k)^T$  be a  $k \times 1$  random vector. Let  $\mathbf{W}_n$  be a sequence of  $k \times k$  nonsingular random matrices, and let  $\mathbf{C}$  be a  $k \times k$  constant nonsingular matrix. If  $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$  and  $\mathbf{y}_n \xrightarrow{P} \mathbf{d}$  for some constant  $k \times 1$  vector  $\mathbf{c}$ , then i)  $\mathbf{x}_n + \mathbf{y}_n \xrightarrow{D} \mathbf{x} + \mathbf{d}$  and

ii)  $\mathbf{y}_n^T \mathbf{x}_n \xrightarrow{D} \mathbf{d}^T \mathbf{x}$ .

c) If  $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$  and  $\mathbf{W}_n \xrightarrow{P} \mathbf{C}$ , then  $\mathbf{W}_n \mathbf{x}_n \xrightarrow{D} \mathbf{C}\mathbf{x}$ ,  $\mathbf{x}_n^T \mathbf{W}_n \xrightarrow{D} \mathbf{x}^T \mathbf{C}$ ,  $\mathbf{W}_n^{-1} \mathbf{x}_n \xrightarrow{D} \mathbf{C}^{-1} \mathbf{x}$ , and  $\mathbf{x}_n^T \mathbf{W}_n^{-1} \xrightarrow{D} \mathbf{x}^T \mathbf{C}^{-1}$ .

34) Let  $\mathbf{u} = (x_{i_1}, \dots, x_{i_k})^T$  where  $n \geq Jk$  with  $J \geq 5$ . Sometimes much larger  $J$  will be needed. Then apply 32ab) with  $\mathbf{u}$  in place of  $\mathbf{x}$  to do a hypothesis test of the form  $H_0 : \mathbf{A}\boldsymbol{\beta}_{OPLS} = \mathbf{0}$  or  $H_0 : \mathbf{A}\boldsymbol{\Sigma}\mathbf{x},Y = \mathbf{0}$  where  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{x},Y = \mathbf{B}\boldsymbol{\Sigma}\mathbf{u},Y$ . In particular, if  $\boldsymbol{\beta}_{OPLS} = (\beta_1, \dots, \beta_p)^T$ , use  $\mathbf{A} = (0, \dots, 0, 1, 0, \dots, 0)$  with a 1 in the  $i$ th position to test  $H_0 : \beta_i = 0$  versus  $H_A : \beta_i \neq 0$ . If  $i_1 < i_2 < \dots < i_k$ , then the  $j$ th row of  $\mathbf{A}$  has a 1 in the  $i_j$  position and all other row entries equal to 0 to test  $H_0 : (\beta_{i_1}, \dots, \beta_{i_k})^T = \mathbf{0}$ . Note that  $\mathbf{A} = \mathbf{I}_p$  can be used to test  $H_0 : \boldsymbol{\beta}_{OPLS} = \mathbf{0}$  in low dimensions ( $n \geq Jp$ ), but not in high dimensions. Similarly,  $\mathbf{A} = [\mathbf{0} \ \mathbf{I}_k]$  can be used to test whether the last  $k$  elements  $(\beta_{p-k+1}, \dots, \beta_p)^T = \mathbf{0}$ . To test  $H_0 : \beta_3 - \beta_4 = 0$ , use  $\mathbf{A} = (0, 0, 1, -1, 0, \dots, 0)$ . Using a test based on 32c) may not work well in high dimensions because  $\hat{\lambda}$  may not be close to  $\lambda$ .

35) If the data  $Y_1, \dots, Y_n$  are arranged in ascending order from smallest to largest and written as  $Y_{(1)} \leq \dots \leq Y_{(n)}$ , then  $Y_{(i)}$  is the  $i$ th order statistic and the  $Y_{(i)}$ 's are called the *order statistics*. The *sample median*

$$\begin{aligned} \text{MED}(n) &= Y_{((n+1)/2)} \quad \text{if } n \text{ is odd,} \\ \text{MED}(n) &= \frac{Y_{(n/2)} + Y_{((n/2)+1)}}{2} \quad \text{if } n \text{ is even.} \end{aligned}$$

The notation  $\text{MED}(n) = \text{MED}(Y_1, \dots, Y_n)$  will also be used. The *sample median absolute deviation* is  $\text{MAD}(n) = \text{MED}(|Y_i - \text{MED}(n)|, i = 1, \dots, n)$ .

36) Suppose the multivariate data has been collected into an  $n \times p$  matrix

$$\mathbf{W} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}.$$

The *coordinatewise median*  $\text{MED}(\mathbf{W}) = (\text{MED}(X_1), \dots, \text{MED}(X_p))^T$  where  $\text{MED}(X_i)$  is the sample median of the data in column  $i$  corresponding to variable  $X_i$ .

37) Let  $(T, \mathbf{C}) = (T(\mathbf{W}), \mathbf{C}(\mathbf{W}))$  be an estimator of multivariate location and dispersion. The  $i$ th *Mahalanobis distance*  $D_i = \sqrt{D_i^2}$  where the  $i$ th *squared Mahalanobis distance* is  $D_i^2 = D_i^2(T(\mathbf{W}), \mathbf{C}(\mathbf{W})) = (\mathbf{x}_i - T(\mathbf{W}))^T \mathbf{C}^{-1}(\mathbf{W})(\mathbf{x}_i - T(\mathbf{W}))$ .

38) The squared Euclidean distances of the  $\mathbf{x}_i$  from the coordinatewise median is  $D_i^2 = D_i^2(\text{MED}(\mathbf{W}), \mathbf{I}_p)$ . Concentration type steps compute the weighted median  $\text{MED}_j$ : the coordinatewise median computed from the cases  $\mathbf{x}_i$  with  $D_i^2 \leq \text{MED}(D_i^2(\text{MED}_{j-1}, \mathbf{I}_p))$  where  $\text{MED}_0 = \text{MED}(\mathbf{W})$ . Often used  $j = 0$  (no concentration type steps) or  $j = 9$ . Let  $D_i = D_i(\text{MED}_j, \mathbf{I}_p)$ . Let  $W_i = 1$  if  $D_i \leq \text{MED}(D_1, \dots, D_n) + k\text{MAD}(D_1, \dots, D_n)$  where  $k \geq 0$  and  $k = 5$  is the default choice. Let  $W_i = 0$ , otherwise.

39) Let the *covmb2 set*  $B$  of at least  $n/2$  cases correspond to the cases with weight  $W_i = 1$ . Then the *covmb2 estimator*  $(T, \mathbf{C})$  is the sample mean and sample covariance matrix applied to the cases in set  $B$ . Hence

$$T = \frac{\sum_{i=1}^n W_i \mathbf{x}_i}{\sum_{i=1}^n W_i} \quad \text{and} \quad \mathbf{C} = \frac{\sum_{i=1}^n W_i (\mathbf{x}_i - T)(\mathbf{x}_i - T)^T}{\sum_{i=1}^n W_i - 1}.$$

The function `ddplot5` plots the Euclidean distances from the coordinatewise median versus the Euclidean distances from the `covmb2` location estimator. Typically the plotted points in this DD plot cluster about the identity line, and outliers appear in the upper right corner of the plot with a gap between the bulk of the data and the outliers.

Table 1: 40) OPLS MLR Results

General	$\beta_{OLS} = \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x},Y} = \lambda \Sigma_{\mathbf{x},Y} = \beta_{OPLS}$
$\beta_{OLS} = \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x},Y} = \frac{1}{\lambda} [Cov(\mathbf{x})]^{-1} \beta_{OPLS}$	$\beta_{OLS}$ is an eigenvector of $\Sigma_{\mathbf{x}}$
$\beta_{OPLS} = \lambda \Sigma_{\mathbf{x},Y} = \lambda Cov(\mathbf{x}) \beta_{OLS}$	$\beta_{OPLS}$ is an eigenvector of $\Sigma_{\mathbf{x}}$
$\Sigma_{\mathbf{x},Y} = Cov(\mathbf{x}) \beta_{OLS} = \Sigma_{\mathbf{x}} \beta_{OLS}$	$\Sigma_{\mathbf{x},Y}$ is an eigenvector of $\Sigma_{\mathbf{x}}$

41) **Multitude of MLR models:** Suppose

$$\begin{pmatrix} Y \\ \mathbf{x} \end{pmatrix} \sim N_{p+1} \left( \begin{pmatrix} \mu_Y \\ \boldsymbol{\mu}_x \end{pmatrix}, \begin{pmatrix} \Sigma_Y & \Sigma_{Y\mathbf{x}} \\ \Sigma_{\mathbf{x}Y} & \Sigma_{\mathbf{x}} \end{pmatrix} \right).$$

Let  $\mathbf{A}$  be a constant  $k \times p$  matrix with full rank  $k$ . Let  $\mathbf{w} = \mathbf{A}\mathbf{x}$ . Assume  $\Sigma_{\mathbf{x}}$  is nonsingular. Then

$$\begin{pmatrix} Y \\ \mathbf{w} \end{pmatrix} \sim N_{k+1} \left( \begin{pmatrix} \mu_Y \\ \boldsymbol{\mu}_w \end{pmatrix}, \begin{pmatrix} \Sigma_Y & \Sigma_{Y\mathbf{w}} \\ \Sigma_{\mathbf{w}Y} & \Sigma_{\mathbf{w}} \end{pmatrix} \right) \sim N_{k+1} \left( \begin{pmatrix} \mu_Y \\ \mathbf{A}\boldsymbol{\mu}_x \end{pmatrix}, \begin{pmatrix} \Sigma_Y & \Sigma_{Y\mathbf{x}}\mathbf{A}^T \\ \mathbf{A}\Sigma_{\mathbf{x}Y} & \mathbf{A}\Sigma_{\mathbf{x}}\mathbf{A}^T \end{pmatrix} \right).$$

If  $\mathbf{A} = \boldsymbol{\eta}^T$  and  $\mathbf{w} = \boldsymbol{\eta}^T \mathbf{x}$  is a random variable, then

$$\begin{pmatrix} Y \\ \boldsymbol{\eta}^T \mathbf{x} \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_Y \\ \boldsymbol{\eta}^T \boldsymbol{\mu}_x \end{pmatrix}, \begin{pmatrix} \Sigma_Y & \Sigma_{\mathbf{x}Y}^T \boldsymbol{\eta} \\ \boldsymbol{\eta}^T \Sigma_{\mathbf{x}Y} & \boldsymbol{\eta}^T \Sigma_{\mathbf{x}} \boldsymbol{\eta} \end{pmatrix} \right).$$

Using 17),  $Y | \boldsymbol{\eta}^T \mathbf{x} \sim N(\alpha \boldsymbol{\eta} + \beta \boldsymbol{\eta}^T \boldsymbol{\mu}_x, \sigma_{\boldsymbol{\eta}}^2)$  where  $\alpha \boldsymbol{\eta} = \mu_Y - \beta \boldsymbol{\eta}^T \boldsymbol{\mu}_x$ ,  $\beta \boldsymbol{\eta} = \lambda \boldsymbol{\eta}$ ,

$$\sigma_{\boldsymbol{\eta}}^2 = \Sigma_Y - \beta \boldsymbol{\eta}^T \Sigma_{\mathbf{x}Y} = \Sigma_Y - \lambda \boldsymbol{\eta}^T \Sigma_{\mathbf{x}Y} = \Sigma_Y - \frac{(\Sigma_{\mathbf{x}Y}^T \boldsymbol{\eta})^2}{\boldsymbol{\eta}^T \Sigma_{\mathbf{x}} \boldsymbol{\eta}},$$

and

$$\lambda = \frac{\Sigma_{\mathbf{x}Y}^T \boldsymbol{\eta}}{\boldsymbol{\eta}^T \Sigma_{\mathbf{x}} \boldsymbol{\eta}}.$$

Hence  $Y | \boldsymbol{\eta}^T \mathbf{x}$  follows a population OLS regression model for the regression of  $Y$  on  $\boldsymbol{\eta}^T \mathbf{x}$ :  $Y | \boldsymbol{\eta}^T \mathbf{x} = \alpha \boldsymbol{\eta} + \mathbf{x}^T (\lambda \boldsymbol{\eta}) + e$  where  $e \sim N(0, \sigma_{\boldsymbol{\eta}}^2)$ . Using  $\boldsymbol{\eta} = \Sigma_{\mathbf{x}Y}$  corresponds to OPLS. Using  $\boldsymbol{\eta} = \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x}Y} = \beta_{OLS}$  gives  $Y | \beta_{OLS}^T \mathbf{x} \sim Y | \mathbf{x} \sim N(E(Y | \mathbf{x}), V(Y | \mathbf{x}))$  or  $Y | \mathbf{x} = \alpha_{OLS} + \mathbf{x}^T \beta_{OLS} + e$  where  $e \sim N(0, V(Y | \mathbf{x}))$  and  $V(Y | \mathbf{x}) = \Sigma_Y - \Sigma_{\mathbf{x}Y}^T \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x}Y}$ .

Using 17),  $Y | \mathbf{w}$  also follows a population OLS regression model for the regression of  $Y$  on  $\mathbf{w}$ :  $Y | \mathbf{w} = \alpha \mathbf{w}_{,OLS} + \mathbf{w}^T \beta_{\mathbf{w},OLS} + e$  where  $e \sim N(0, \sigma_{Y|\mathbf{w}}^2)$ . Here

$$\beta_{\mathbf{w},OLS} = \Sigma_{\mathbf{w}}^{-1} \Sigma_{\mathbf{w}Y} = (\mathbf{A} \Sigma_{\mathbf{x}} \mathbf{A}^T)^{-1} \mathbf{A} \Sigma_{\mathbf{x}Y},$$

and  $\alpha \mathbf{w}_{,OLS} = E(Y) - \beta_{\mathbf{w},OLS}^T E(\mathbf{w})$ .

42) Referring to 40) and 41), be able to compute  $\lambda_{OPLS}$  from 31),  $\beta_{OPLS} = \lambda \Sigma_{\mathbf{x}Y}$ ,  $\beta_{OLS} = \Sigma_{\mathbf{x}}^{-1} \Sigma_{\mathbf{x}Y}$ , and  $\Sigma_{\mathbf{x}Y} = \Sigma_{\mathbf{x}} \beta_{OLS}$  if the cases are iid, the MLR model holds, and  $\Sigma_{\mathbf{x}}$  is a diagonal matrix. Be able to recognize  $\beta = \beta_{OLS}$  from the MLR model.

43) Let  $\|\mathbf{a}\|_2 = \sqrt{\mathbf{a}^T \mathbf{a}}$  be the Euclidean norm.

44) The marginal maximum likelihood estimator (MMLE or marginal least squares estimator) computes the marginal regression of  $Y$  on  $x_i$  resulting in the estimator  $(\hat{\alpha}_{i,M}, \hat{\beta}_{i,M})$  for  $i = 1, \dots, p$ . Then  $\hat{\boldsymbol{\beta}}_{MMLE} = (\hat{\beta}_{1,M}, \dots, \hat{\beta}_{p,M})^T$ .

45) For multiple linear regression, the marginal estimators are the simple linear regression (SLR) estimators, and  $(\hat{\alpha}_{i,M}, \hat{\beta}_{i,M}) = (\hat{\alpha}_{i,SLR}, \hat{\beta}_{i,SLR})$ . Hence

$$\hat{\boldsymbol{\beta}}_{MMLE} = [\text{diag}(\hat{\boldsymbol{\Sigma}}\mathbf{x})]^{-1} \hat{\boldsymbol{\Sigma}}\mathbf{x}_Y.$$

If the  $\mathbf{w}_i$  are the predictors standardized to have unit sample variances, then

$$\hat{\boldsymbol{\beta}}_{MMLE} = \hat{\boldsymbol{\beta}}_{MMLE}(\mathbf{w}, Y) = \hat{\boldsymbol{\Sigma}}\mathbf{w}_Y = \mathbf{I}^{-1} \hat{\boldsymbol{\Sigma}}\mathbf{w}_Y = \hat{\boldsymbol{\eta}}_{OPLS}(\mathbf{w}, Y)$$

where  $(\mathbf{w}, Y)$  denotes that  $Y$  was regressed on  $\mathbf{w}$ , and  $\mathbf{I}$  is the  $p \times p$  identity matrix.

46) The MMLE is interesting since if each predictor satisfies a marginal model, then the marginal model theory can be used to find a confidence interval for  $\beta_i$  for  $i = 1, \dots, p$  where  $\beta_i$  is the  $i$ th component of  $\boldsymbol{\beta}_{MMLE}$ . For MLR, let  $\mathbf{V} = \text{diag}(\boldsymbol{\Sigma}\mathbf{x}) = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$ . For iid cases,  $\boldsymbol{\beta}_{MMLE} = \mathbf{V}^{-1}\boldsymbol{\Sigma}\mathbf{x}_Y = \mathbf{V}^{-1}\boldsymbol{\Sigma}\mathbf{x}\boldsymbol{\beta}_{OLS}$ .

For standardized predictors, let  $s_j$  and  $\sigma_j$  be the sample and population standard deviations of  $x_j$ . Let  $\mathbf{w}_i = \hat{\mathbf{D}}\mathbf{x}_i = \text{diag}(1/s_1, \dots, 1/s_p)\mathbf{x}_i$  and  $\mathbf{u}_i = \mathbf{D}\mathbf{x}_i = \text{diag}(1/\sigma_1, \dots, 1/\sigma_p)\mathbf{x}_i$ . Note that  $\sqrt{n}(\hat{\boldsymbol{\Sigma}}\mathbf{w}_Y - \boldsymbol{\Sigma}\mathbf{u}_Y) = \sqrt{n}(\hat{\boldsymbol{\Sigma}}\mathbf{w}_Y - \hat{\boldsymbol{\Sigma}}\mathbf{u}_Y) + \sqrt{n}(\hat{\boldsymbol{\Sigma}}\mathbf{u}_Y - \boldsymbol{\Sigma}\mathbf{u}_Y) = O_P(1) + \sqrt{n}(\hat{\boldsymbol{\Sigma}}\mathbf{u}_Y - \boldsymbol{\Sigma}\mathbf{u}_Y)$  under mild regularity conditions for iid cases. Hence  $\hat{\boldsymbol{\Sigma}}\mathbf{w}_Y$  is a  $\sqrt{n}$  consistent estimator of  $\boldsymbol{\Sigma}\mathbf{u}_Y$ . Note that  $\boldsymbol{\Sigma}\mathbf{u}$  is the correlation matrix of  $\mathbf{x}$ .

47) Theorem: Consider the MMLE for MLR. Suppose the cases  $(Y_i, \mathbf{x}_i^T)^T$  are iid from some distribution. Let  $\mathbf{w}_i$  be the standardized predictors and assume  $\hat{\boldsymbol{\Sigma}}\mathbf{w}_Y \xrightarrow{P} \boldsymbol{\Sigma}\mathbf{u}_Y$  and  $\hat{\boldsymbol{\Sigma}}\mathbf{w} \xrightarrow{P} \boldsymbol{\Sigma}\mathbf{u}$  where the  $\hat{\boldsymbol{\Sigma}}\mathbf{w}$  are nonsingular for large enough  $n$  and  $\boldsymbol{\Sigma}\mathbf{u}$  is nonsingular.

$$\begin{aligned} a) \hat{\boldsymbol{\beta}}_{MMLE} &= \hat{\boldsymbol{\beta}}_{MMLE}(\mathbf{w}, Y) = \hat{\boldsymbol{\Sigma}}\mathbf{w}_Y = \hat{\boldsymbol{\eta}}_{OPLS}(\mathbf{w}, Y) \xrightarrow{P} \boldsymbol{\Sigma}\mathbf{u}_Y = \\ &\boldsymbol{\eta}_{OPLS}(\mathbf{u}, Y) = \boldsymbol{\beta}_{MMLE} = \boldsymbol{\Sigma}\mathbf{u}[\boldsymbol{\Sigma}\mathbf{u}]^{-1}\boldsymbol{\Sigma}\mathbf{u}_Y = \boldsymbol{\Sigma}\mathbf{u}\boldsymbol{\beta}_{OLS}(\mathbf{u}, Y). \end{aligned}$$

b) Let  $\boldsymbol{\beta}_{OLS} = \boldsymbol{\beta}_{OLS}(\mathbf{u}, Y)$ . Then  $\boldsymbol{\beta}_{MMLE} = \boldsymbol{\Sigma}\mathbf{u}\boldsymbol{\beta}_{OLS} = \boldsymbol{\beta}_{OLS}$  if  $\boldsymbol{\beta}_{OLS} = \mathbf{0}$  or if  $\boldsymbol{\beta}_{OLS}$  is an eigenvector of  $\boldsymbol{\Sigma}\mathbf{u}$  with eigenvalue = 1.

## Ch. 2

48) Refer to point 5) for a model for variable selection:  $\mathbf{x}^T\boldsymbol{\beta} = \mathbf{x}_S^T\boldsymbol{\beta}_S + \mathbf{x}_E^T\boldsymbol{\beta}_E = \mathbf{x}_S^T\boldsymbol{\beta}_S$ . If  $S \subseteq I$ , then  $\mathbf{x}^T\boldsymbol{\beta} = \mathbf{x}_S^T\boldsymbol{\beta}_S = \mathbf{x}_I^T\boldsymbol{\beta}_I + \mathbf{x}_O^T\mathbf{0} = \mathbf{x}_I^T\boldsymbol{\beta}_I$ . To clarify notation, suppose  $p = 3$ , a constant  $\alpha$  is always in the model, and  $\boldsymbol{\beta} = (\beta_1, 0, 0)^T$ . Then the  $J = 2^p = 8$  possible subsets of  $\{1, 2, \dots, p\}$  are  $I_1 = \emptyset$ ,  $S = I_2 = \{1\}$ ,  $I_3 = \{2\}$ ,  $I_4 = \{3\}$ ,  $I_5 = \{1, 2\}$ ,  $I_6 = \{1, 3\}$ ,  $I_7 = \{2, 3\}$ , and  $I_8 = \{1, 2, 3\}$ . There are  $2^{p-a_S} = 4$  subsets  $I_2, I_5, I_6$ , and  $I_8$  such that  $S \subseteq I_j$ . Let  $\hat{\boldsymbol{\beta}}_{I_7} = (\hat{\beta}_2, \hat{\beta}_3)^T$  and  $\mathbf{x}_{I_7} = (x_2, x_3)^T$ .

Let  $I_{min}$  correspond to the set of predictors selected by a variable selection method such as forward selection or lasso variable selection. If  $\hat{\boldsymbol{\beta}}_I$  is  $a \times 1$ , use zero padding to form the  $p \times 1$  vector  $\hat{\boldsymbol{\beta}}_{I,0}$  from  $\hat{\boldsymbol{\beta}}_I$  by adding 0s corresponding to the omitted variables. For example, if  $p = 4$  and  $\hat{\boldsymbol{\beta}}_{I_{min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$ , then the observed variable selection estimator  $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$ . As a statistic,  $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}$  with probabilities  $\pi_{kn} = P(I_{min} = I_k)$  for  $k = 1, \dots, J$  where there are  $J$  subsets, e.g.  $J = 2^p$ .



$I_j$	model	$x_1$	$x_2$	$x_3$	$x_4$	$\hat{\beta} = \hat{\beta}_{I_j}$
49) $I_1$	1		*			$(0, \hat{\beta}_2, 0, 0)^T$
$I_2$	2		*	*		$(0, \hat{\beta}_2, \hat{\beta}_3, 0)^T$
$I_3$	3	*	*	*		$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, 0)^T$
$I_4$	4	*	*	*	*	$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4)^T = \hat{\beta}_{OLS}$

Model  $I_{min}$  is the model, among  $p$  candidates, that minimizes  $C_p$ , AIC, BIC if  $n \geq 10$ , or EBIC if  $n < 10p$ . Model  $I_j$  contains  $j$  predictors,  $x_1^*, x_2^*, \dots, x_j^*$  if forward selection is used.

50) Variable selection is a search for a subset of predictors that can be deleted without important loss of information if  $n \geq 10p$  and such that model  $I$  (containing the remaining predictors that were not deleted) is good for prediction if  $n < 10p$ . Note that the “100%” shorth CI for a  $\beta_i$  that is a component of  $\beta_O$  is  $[0,0]$ .

51) Underfitting occurs if  $S \not\subseteq I$  so that  $\mathbf{x}_I$  is missing important predictors. Underfitting will occur if  $\mathbf{x}_I$  is  $k \times 1$  with  $d = k < a_S$ . Overfitting occurs if  $S \subset I$  with  $S \neq I$  or if  $n < 5k$ .

52) In 49) sometimes TRUE = \* and FALSE = blank. The  $x_i$  may be replaced by the variable name or letters like a b c d.

$I_j$	model	$x_2$	$x_3$	$x_4$	$x_5$
$I_2$	1	FALSE	TRUE	FALSE	FALSE
$I_3$	2	FALSE	TRUE	TRUE	FALSE
$I_4$	3	TRUE	TRUE	TRUE	FALSE
$I_5$	4	TRUE	TRUE	TRUE	TRUE

53) The out\$cp line gives  $C_p(I_2), C_p(I_3), \dots, C_p(I_p) = p$  and  $I_{min}$  is the  $I_j$  with the smallest  $C_p$ .

54) Typical bootstrap output for forward selection, lasso, and elastic net is shown below. The SE column is usually omitted except possibly for forward selection. The term “coef” might be replaced by “Estimate.” This column gives  $\hat{\beta}_{I,0}$  where  $I = I_{min}$  for forward selection,  $I = L$  for lasso, and  $I = EN$  for elastic net. Note that the SE entry is omitted if  $\hat{\beta}_i = 0$  so variable  $x_i$  was omitted by the variable selection method. In the output below,  $\hat{\beta}_2 = \hat{\beta}_3 = 0$ . The SE column corresponds to the OLS SE obtained by acting as if the OLS full model contains a constant and the variables not omitted by the variable selection method. The OLS SE is incorrect unless the variables were selected before looking at the data for forward selection.

Label	Estimate or coef	SE	shorth 95% CI for $\beta_i$
Constant=intercept= $x_1$	$\hat{\beta}_1$	$SE(\hat{\beta}_1)$	$[\hat{L}_1, \hat{U}_1]$
$x_2$	$\hat{\beta}_2$	$SE(\hat{\beta}_2)$	$[\hat{L}_2, \hat{U}_2]$
$x_3$	0		$[\hat{L}_3, \hat{U}_3]$
$x_4$	0		$[\hat{L}_4, \hat{U}_4]$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_p$	$\hat{\beta}_p$	$SE(\hat{\beta}_p)$	$[\hat{L}_p, \hat{U}_p]$

55) Forward selection generates a sequence of  $J$  models. Assume all models with a constant  $\alpha$  contain a constant. To form  $I_1$ , consider all models  $I$  with one predictor. Compute criterion  $C(x_i)$  and let  $x_1^*$  minimize the criterion. Then  $I_1$  corresponds to  $x_1^*$  and  $p$  models are fit. To form  $I_2$ , compute  $C(x_1^*, x_j)$  for the  $p - 1$  models where  $x_j \neq x_1^*$ . Continue in this manner. The last model fit is  $I_J$  using  $C(x_1^*, \dots, x_{j-1}^*, x_j)$  for the  $p - J + 1$  models where  $x_j$  is not one of the  $J - 1$  variables already selected.  $I_{min}$  is the model corresponding to the smallest criterion. Often  $J = \min(n - r, p)$  for some integer  $r \geq 0$ .

The same regression method, eg OLS or a GLM, is used to compute each fitted model  $\hat{\beta}_I$ . For MLR, typically  $C(I) = RSS(I) + p(I)$  where  $SSE(I) = RSS(I) = \mathbf{r}^T(I)\mathbf{r}(I) = \sum_{i=1}^n r_i^2(I) = \sum_{i=1}^n (Y_i - \hat{Y}_i(I))^2$ .

56) Forward selection can be slow if  $n$  and  $p$  are large with the number of fitted models close to  $n(p - n/2)$  if  $p \gg n$ .

57) Lasso variable selection uses a grid of  $J$  models depending on a parameter  $\lambda_i$  for  $i = 1, \dots, J$  to form  $I_1, \dots, I_J$ . The `glmnet` default appears to be 100, but researchers often use much larger  $J$  for high dimensional data.

58) MMLE variable selection: Let the  $\mathbf{x}_i$  be the predictors. Let the  $\mathbf{w}_i$  be the standardized predictors such that the sample variance of each predictor is 1. Find the  $J$  variables  $x_1^*, \dots, x_j^*$  corresponding to the largest  $|\hat{\beta}_i|$ . Note that these variables are not standardized.

59) A two stage variable selection method is to use 58) to find  $r$  variables, then use forward selection or lasso variable selection on the  $r$  variables to get model  $I_{min}$  with variables  $x_1^*, \dots, x_a^*$ .

**Math 583 Exam 1 is on Wednesday, Sept. 27 and covers homeworks 1-4 and quizzes 1-4. You are allowed 9 sheets of notes and a calculator.**