

37) The squared Euclidean distance

$$\text{of } \underline{x} \text{ from } T \text{ is } (\underline{x} - T)^T (\underline{x} - T) = \underset{\sim}{D}_x^2(T, \mathbb{I}_p)$$

38) Outlier detection if $p > n$, but more than $\frac{n}{2}$ cases in the bulk of the data.

a) Use Euclidean distances from the coordinatewise median $D_i(\text{MED}(W), \mathbb{I}_p)$.

b) Let MED_j be the coordinatewise median computed from cases \underline{x}_i with

$$D_i^2 \leq D_i^2(\text{MED}_{j-1}, \mathbb{I}_p) \text{ where } \text{MED}_0 = \text{MED}(W)$$

often use $j=0$ or $j=9$. Let $D_i = D_i(\text{MED}_j, \mathbb{I}_p)$.

$$\text{Let } w_i = \begin{cases} 1 & \text{if } D_i \leq \text{MED}(D_1, \dots, D_n) + 5\text{MAD}(D_1, \dots, D_n) \\ 0 & \text{else} \end{cases}$$

The covmb2 set B consists of the $m \geq \frac{n}{2}$ cases with weight $w_i = 1$. The covmb2 estimator $(T)_c$ is the

sample mean and sample covariance matrix of the m cases: 22.5

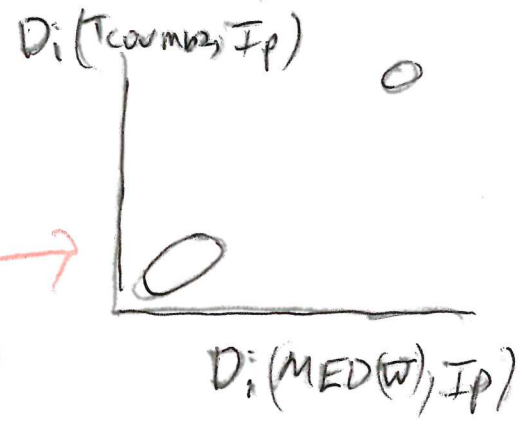
$$T = \frac{\sum_{i=1}^n w_i \underline{x}_i}{\sum_{i=1}^n w_i} \quad , \quad Q = \frac{\sum_{i=1}^n w_i (\underline{x}_i - T)(\underline{x}_i - T)^T}{\sum_{i=1}^n w_i - 1}$$

lots of w 's

39] Let $\underline{w}_i = (y_i, \underline{x}_i^T)^T$ and let the continuous predictors from \underline{x}_i be \underline{u}_i (the predictors that take on many values, so not gender). Apply the regression (or classification) method to the m cases \underline{w}_i corresponding to the covmb2 set $B = \{i_1, \dots, i_m\}$ applied to $\underline{u}_1, \dots, \underline{u}_n$. MLR, GLMs, GAMs, LDA, QDA, KNN etc.

40] The function DD plots plots $D_i(\text{MED}(w), I_p)$ vs $D_i(T_{\text{covmb2}}, I_p)$.

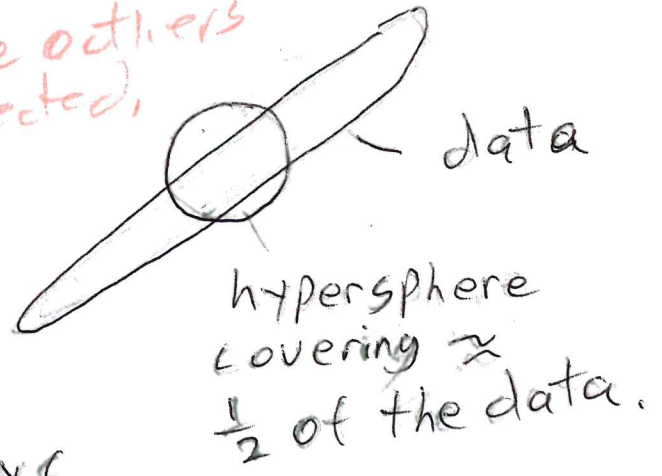
The plotted points tend to cluster about the identity line with outliers in the upper right corner of the plot with a gap between the bulk of the data and the outliers.



see HW 3

41) To detect outliers in one group and the bulk of the data in another group, the distance of the outliers from the bulk of the data increases roughly with \sqrt{p} .

otherwise the outliers can't be detected.



outliers in the hypersphere can sometimes be detected with j steps, but not always.

42) For $\begin{pmatrix} y \\ z \end{pmatrix} \sim N_{p+1}(\underline{\mu}, \Sigma)$,

$$1.4826 \text{ MAD}(z_j) \xrightarrow{p} \sqrt{V(z_j)} = \text{SD}(z_j)$$

where $z_j = y$ or $z_j = x_j$ (j th variable).

Since $\text{COV}(X, Y) \stackrel{\leftarrow \text{HW1}}{=} \frac{V(X+Y) - V(X-Y)}{4}$, 23.5
 $\hat{\text{COV}}(X, Y)$.

$$\approx \frac{(1.4826)^2}{4} \left([\text{MAD}(X+Y)]^2 - [\text{MAD}(X-Y)]^2 \right)$$

$$\approx \text{COV}(X_i, Y)$$

So $\hat{\underline{m}}_{\sim M} = \begin{pmatrix} \hat{\text{COV}}(X_1, Y) \\ \vdots \\ \hat{\text{COV}}(X_p, Y) \end{pmatrix}$ is a

robust estimator of $\underline{\beta}_{XY}$.

see HW2
C)

43) If the predictors $(\underline{x}_1, \dots, \underline{x}_p)$ are continuous, get the covmbz set B applied on $(\underline{x}_1, \dots, \underline{x}_p, Y)$ ← assuming all continuous and let

$\hat{\underline{m}}_{\sim MB}$ be the sample covariance vector $\hat{\underline{\beta}}_{\sim MB}(MB)$ applied to the cases in B .

(unlike 39), apply covmbz to \underline{x} and Y)

44) ^{know} compute pop OLS $\underline{\Sigma}_x^{-1} \underline{\Sigma}_{xy}$ HD 24

and pop OPLS $\lambda \underline{\Sigma}_{xy}$ given

$\underline{\Sigma}_x$ and either $\underline{\Sigma}_{xy}$ or $\underline{\beta} = \underline{\beta}_{OLS}$.

use $\underline{\Sigma}_{xy} = \underline{\Sigma}_x \underline{\beta}$.

OPLS paper near bottom of course webpage

ex) $p=4, \underline{x}_i \sim N_p \left[\underline{0}, \underbrace{\text{diag}(1, 2, 3, 4)}_{\underline{\Sigma}_x} \right]$

$\tilde{x}_i = y_i = x_{i1} + x_{i2} + e_i, i=1, \dots, n, \text{ ind. cases}$
not needed

$e_i \sim N(0, 1) \perp \underline{x}_i$
not needed

So $\alpha = 0, \underline{\beta} = (1, 1, 0, 0)^T, (\alpha, \underline{\beta}) = (\alpha_{OLS}, \underline{\beta}_{OLS})$.

$$\underline{\Sigma}_{xy} = \underline{\Sigma}_x \underline{\beta}_{OLS} = \begin{pmatrix} 1 & & & \\ & 2 & & \\ & & 3 & \\ & & & 4 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 0 \\ 0 \end{pmatrix}$$

$\lambda_{OPLS} =$

$$\lambda = \frac{\underline{\Sigma}_{xy}^T \underline{\Sigma}_{xy}}{\underline{\Sigma}_x^T \underline{\Sigma}_x \underline{\Sigma}_{xy}} = \frac{5}{9}, \quad S = 1^2 + 2^2 + 0^2 + 0^2$$

$$(1, 2, 0, 0) \begin{pmatrix} 1 \\ 4 \\ 0 \\ 0 \end{pmatrix} = 1(1) + 2(4) = 9$$

$$\text{So } \beta_{OLS} = \lambda \underline{\underline{1}}_{xy} = \left(\frac{5}{9}, \frac{10}{9}, 0, 0 \right)^T \quad (24.5)$$

$$\neq \beta_{OLS} = (1, 1, 0, 0)^T$$

know for E1, Hw4, Q4

45) Given $y | \underline{\underline{m}}^T \underline{\underline{x}} \sim N(\mu_w, \sigma_w^2)$ as in 23), be able to compute

$$\sigma_w^2 = \sigma_{\underline{\underline{m}}}^2 = \sigma_y^2 - \frac{(\underline{\underline{1}}_{xy}^T \underline{\underline{m}})^2}{\underline{\underline{m}}^T \underline{\underline{1}}_x \underline{\underline{m}}}$$

if $\underline{\underline{m}} = \underline{\underline{1}}_{xy}$ or $\underline{\underline{m}} = \beta_{OLS}$.

ex) In last ex, if $\underline{\underline{m}} = \underline{\underline{1}}_{xy}$, then

$$\sigma_w^2 = \sigma_y^2 - \frac{(\underline{\underline{1}}_{xy}^T \underline{\underline{1}}_{xy})^2}{\underline{\underline{1}}_{xy}^T \underline{\underline{1}}_x \underline{\underline{1}}_{xy}} =$$

$$\sigma_y^2 - \frac{5^2}{9} = \sigma_y^2 - \frac{25}{9}$$

If $\hat{\beta} = \hat{\beta}_{OLS}$,

then $\sigma_w^2 = V(Y|\underline{x}) = \sigma_y^2 - \underline{x}_{xy}^T \underline{x}_x^{-1} \underline{x}_{xy} =$
 $\text{diag}(1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}) \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}$

$$\left(\sigma_y^2 - \frac{(\underline{x}_{xy}^T \underline{x}_x^{-1} \underline{x}_{xy})^2}{\underline{x}_{xy}^T \underline{x}_x^{-1} \underline{x}_x \underline{x}_x^{-1} \underline{x}_{xy}} \right) = \sigma_y^2 - (1, 2, 0, 0) \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

$$= \sigma_y^2 - 3 = \sigma_y^2 - \frac{27}{9}$$

§ 3.10

46] The marginal maximum likelihood estimator (MMLE) is simple:

regress Y on X_i to get $\hat{\beta}_i$,
marginal regression estimator

then $\hat{\beta}_{MMLE} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^T$

47] For MLR, the MMLE is sometimes

called the marginal least squares estimator, and $\hat{\beta}_i$ comes from the simple linear regression (SLR) of y_i on x_i . (29.5)

48] Pop SLR is OLS pop reg

$$\text{of } Y \text{ on } X_i \therefore \beta_i = \Sigma_{X_i}^{-1} \Sigma_{X_i, Y}$$

$$= \frac{\text{COV}(X_i, Y)}{V(X_i)}$$

$$\hat{\beta}_i = \frac{\widehat{\text{COV}}(X_i, Y)}{\widehat{V}(X_i)}$$

diag matrix with diagonal of Σ_X

$$\text{Hence } \underline{\beta}_{\text{MMLE}} = \left[\text{diag}(\underline{\Sigma}_X) \right]^{-1} \underline{\Sigma}_{X, Y}$$

$$\text{and } \hat{\underline{\beta}}_{\text{MMLE}} = \left[\text{diag}(\hat{\underline{\Sigma}}_X) \right]^{-1} \hat{\underline{\Sigma}}_{X, Y}$$

49] If the \underline{w}_i are vectors of standardized predictors,

such that the sample variances

$$= 1, \text{ then } \hat{\beta}_{\text{MMLE}} = \hat{\beta}_{\text{MMLE}}(\underline{w}, \underline{y})$$

$$= \hat{\beta}_{\underline{w}, \underline{y}} = \underline{I}^{-1} \hat{\beta}_{\underline{w}, \underline{y}} = \hat{\beta}_{\text{OPLS}}(\underline{w}, \underline{y}).$$

50] Let $V = \text{diag}(\underline{I}_X) = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$.

Let $D = \text{diag}(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_p})$.

Let $\text{diag}(\hat{\beta}_X) = \text{diag}(s_1^2, \dots, s_p^2)$

and $\hat{D} = \text{diag}(\frac{1}{s_1}, \dots, \frac{1}{s_p})$.

Let $\underline{w}_i = \hat{D} \underline{x}_i$ and $\underline{v}_i = D \underline{x}_i$.

Then \underline{I}_V is the pop correlation matrix of the \underline{x}_i . For iid cases,

$$\hat{\beta}_{\text{MMLE}} = V^{-1} \hat{\beta}_{X,Y} = V^{-1} \hat{\beta}_X \hat{\beta}_X^{-1} \hat{\beta}_{X,Y} =$$

$V^{-1} \hat{\beta}_X \hat{\beta}_{\text{OLS}}$. Then $\hat{\beta}_{\text{MMLE}} = \hat{\beta}_{\text{OLS}}$

if $\beta_{OLS} = \underline{0}$, if $(\underline{V}^{-1} - \underline{K}_X^{-1}) \underline{K}_{xy} = \underline{0}$,

26.5

or if β_{OLS} is an eigenvector of $\underline{V}^{-1} \underline{K}_X$ with eigenvalue 1.

5) MML theory: For MLR with iid cases, let \underline{w}_i be the standardized predictors and assume $\hat{\underline{K}}_{wy} \xrightarrow{P} \underline{K}_{wy}$ and $\hat{\underline{K}}_w \xrightarrow{P} \underline{K}_w$ where the $\hat{\underline{K}}_w$ are non singular for large enough n and \underline{K}_w is non singular.

$$a) \hat{\beta}_{MML}(\underline{w}, \gamma) = \hat{\underline{K}}_{wy} = \hat{\underline{m}}_{OPLS}(\underline{w}, \gamma)$$

$$\xrightarrow{P} \underline{K}_{wy} = \underline{m}_{OPLS}(\underline{w}, \gamma) =$$

$$\underline{K}_w \underline{K}_w^{-1} \underline{K}_{wy} = \underline{K}_w \beta_{OLS}(\underline{w}, \gamma)$$

b) Let $\beta_{OLS} = \beta_{OLS}(\underline{w}, \gamma)$. Then

$$\underline{\beta}_{MML} = \beta_{MML}(\underline{w}, \gamma) = \underline{K}_w \beta_{OLS} = \beta_{OLS}$$

3) variable selection is

27.5

a search for a subset of predictors that can be deleted

{ without important loss of information if $n \geq 10p$
so that model I is good for prediction if $n \leq 5p$
↑
remaining predictors, in the model

4) Often there are J candidate submodels

I_1, \dots, I_J . Let I_{\min} be the model that minimizes criterion $C(I)$

$I \in I_1, \dots, I_J$. $AIC(I)$, $BIC(I)$, $EBIC(I)$ are common. $C_p(I)$ for MLR.

5) Forward selection; models have $sp = \alpha + \beta^T x$
(take $\alpha = 0$ if there is no constant), predictors x_1, \dots, x_p .

I_1 uses x_1^* where x_1^* minimizes criterion $C(x_i)$, $i = 1, \dots, p$. p models fit

I_2 uses (x_1^*, x_2^*) where x_2^* minimizes

Criterion $C(x_1^*, x_j^*) : p-1$ models fit

I_J uses x_1^*, \dots, x_J^*

where x_j^* minimizes $C(x_1^*, x_2^*, \dots, x_{j-1}^*, x_j^*) : p-j+1$ models fit

Often $J = \min(n-r, p)$

where $r \in \{0, 1, 2, 3, 4, 5\}$.

(could use $\min(n, p)$
eg $\alpha = 0.1$)

I_{\min} has the smallest value of $C(x_1^*, \dots, x_J^*)$
 $J = 1, \dots, J$.

6] Forward selection can be too slow if n and p are large!

fit $p + p-1 + p-2 + \dots + p-J+1$ models

step $\begin{matrix} 1 & 2 & 3 & & J \end{matrix}$
 $= \sum_{i=1}^J (p-i+1) = J(p+1) - \sum_{i=1}^J i$

$$= J(p+1) - \frac{J(J+1)}{2} \approx np - \frac{n^2}{2} \approx n(p - \frac{n}{2}) \text{ if } p \gg n.$$

more details later

(28.5)

7) Lasso variable selection uses a grid of λ values to get J submodels I_1, \dots, I_J . Often $J=100$, but larger J may be useful in high dimensions.

8) MMLE variable selection: Standardize the predictors. Take the J variables x_1^*, \dots, x_J^* that have the largest $|\hat{\beta}_i|$. This method can be fast in high dimensions.

~~These variables are not standardized.~~

9) Using 8) with $r \approx \min(n, p)$ and then forward selection or lasso variable selection on the $x_{i_1}^*, \dots, x_{i_r}^*$ selected by 8) to get x_1^*, \dots, x_a^* can reduce the multicollinearity that can result from 8).
eg. lots of highly correlated predictors.

PhD topic

10) Let $\hat{\beta}_I$ be $a \times 1$. Form the

$p \times 1$ vector $\hat{\beta}_{I,0}$ by adding 0's corresponding to the omitted variables.

ex) $p=4$, $\hat{\beta}_{I_{\min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$. Then

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

11) As a statistic, $\hat{\beta}_{vs} = \hat{\beta}_{I_k,0}$

with prob's $\pi_{kn} = P(I_{\min} = I_k)$, $k=1, \dots, J$.

ex) $Y = \text{height}$, $x_1 = \text{height at right shoulder}$, $x_2 = \text{height at left shoulder}$, x_3, \dots, x_p . For MMLE, if $\hat{\beta}_1 = \frac{6}{5}$,

then $\hat{\beta}_2 \approx \frac{6}{5}$. For OPLS, if $\hat{\beta}_1 = \hat{\lambda} \frac{6}{5}$

then $\hat{\beta}_2 = \hat{\lambda} \frac{6}{5}$. For forward selection and lasso vs, probably only one of x_1 and x_2

is included in I_{\min} ,

MMLE variable selection will likely include x_1 and x_2 which have correlation ≈ 1 .

29.5

$p = 4$

ex) Model I_j	x_1	x_2	x_3	x_4	$\hat{\beta}_{I_j, 0}$
I_1			*		$(0, 0, \hat{\beta}_3, 0)^T$
I_2			*	*	$(0, 0, \hat{\beta}_3, \hat{\beta}_4)^T$
I_3	*		*	*	$(\hat{\beta}_1, 0, \hat{\beta}_3, \hat{\beta}_4)^T$
I_4	*	*	*	*	$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4)^T$

Full model

Model	I_1	I_2	I_3	I_4
$C(I)$	33.21	17.43	3.41	5.00

$$I_{\min} = I_3 = \{1, 3, 4\}$$

$$\hat{\beta}_{I_{\min}} = (\hat{\beta}_1, \hat{\beta}_3, \hat{\beta}_4)^T$$

$$\hat{\beta}_{I_{\min}, 0} = (\hat{\beta}_1, 0, \hat{\beta}_3, \hat{\beta}_4)^T$$

see HW4

($\hat{\beta}_i$ depends on model I_j)