

42. v) MD Plot is

a plot of classical versus robust Mahalanobis distances (MD_i vs RD_i).

$$46) MD_i = \sqrt{(x_i - \bar{x})^T S^{-1} (x_i - \bar{x})} = D_i(\bar{x}, S)$$

47) Let (T_A, C_A) be a robust estimator,

$$RD_i = \frac{\sqrt{\chi_{p, 0.5}^2}}{\text{MED}(D_i(A))} D_i(A) \quad \text{Where } D_i(A) = D_i(T_A, C_A)$$

and $\chi_{p, 0.5}^2$ is the median of the χ_p^2 distribution.

Equivalently, $RD_i = D_i(T_R, C_R)$ where

$$T_R = T_A \quad \text{and} \quad C_R = \frac{C_A}{T^2} \quad \text{where}$$

$$T = \frac{\sqrt{\chi_{p, 0.5}^2}}{\text{MED}(D_i(A))}$$

← sample median of $D_1(A), \dots, D_n(A)$

48) Suppose the data x_1, \dots, x_n are iid from a $EC_p(\underline{\mu}, \Sigma, \sigma)$ distribution with 2nd moment. Then (\bar{x}, S) is a consistent estimator for $(\underline{\mu}, C_2 \Sigma) = (\underline{\mu}, \text{cov}(x))$. Suppose that the robust estimator (T_A, C_A) is a consistent estimator for $(\underline{\mu}, C_2 \Sigma)$ for some constant c_A .

Then $MD_i \leftarrow \frac{1}{\sigma_A} (x_i - \mu)' \Sigma^{-1} (x_i - \mu)$ (24.9)

$= \frac{\sigma_A}{\sigma_A} \frac{1}{\sigma_X} (x_i - \mu)' \Sigma^{-1} (x_i - \mu)$

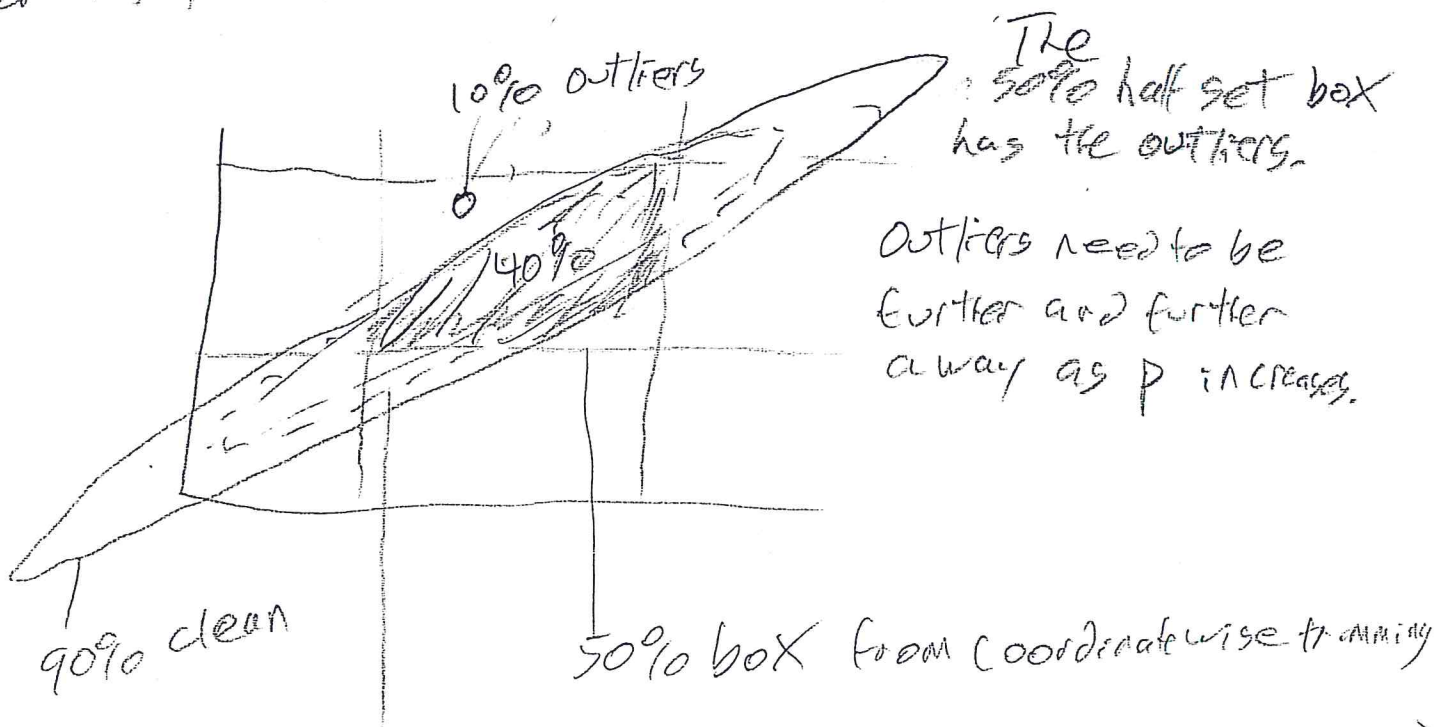
$= \frac{\sigma_A}{\sigma_X} (x_i - \mu)' (\sigma_A \Sigma^{-1}) (x_i - \mu) \approx \frac{\sigma_A}{\sigma_X} D_i^2(A)$

Problem: $\frac{\sigma_A}{\sigma_X}$ depends on the EL distribution and

σ_A, σ_X and $\frac{\sigma_A}{\sigma_X}$ are usually unknown.

GLP ↓

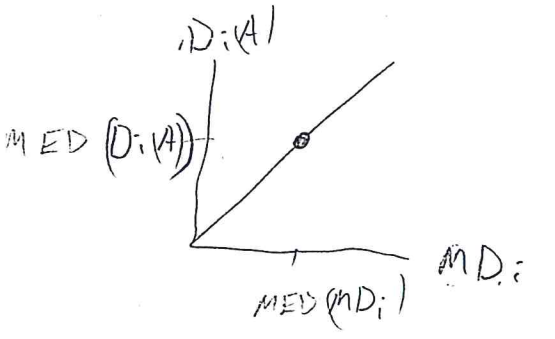
(Coordinatewise trimming! $p=1$ 25% trimmed mean uses half set. As $p \uparrow$ need to decrease amount of trimming to get a half set. Also the method is poor for highly correlated data where $\Sigma \neq \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$.)



50% hyperellipsoid does not contain the outliers,

↑ GLP

eg) P 423 A plot of MD_i vs $D_i(A)$ will have plotted points that cluster tightly about the line through $(0,0)^T$ and $(MED(MD_i), MED(D_i(A)))^T$. (The correlation goes to 1 as $n \rightarrow \infty$.)



Hence $D_i(A) \approx \frac{MED(D_i(A))}{MED(MD_i)} MD_i$.

30) If $\underline{x}_1, \dots, \underline{x}_n$ are iid $N_p(\underline{\mu}, \underline{\Sigma})$, then

$MED(MD_i) \approx \sqrt{\chi^2_{p, 0.5}}$. Hence

sample median pop median

$D_i(A) \approx \frac{MED(D_i(A))}{\sqrt{\chi^2_{p, 0.5}}} MD_i$

or $MD_i \approx \frac{\sqrt{\chi^2_{p, 0.5}}}{MED(D_i(A))} D_i(A) = \gamma D_i(A) = R D_i$.

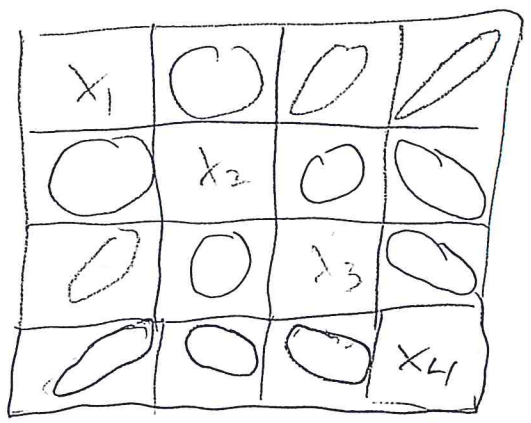
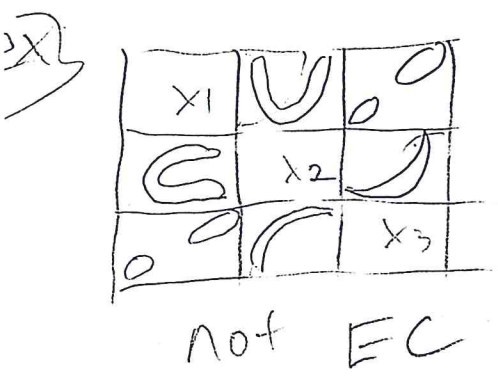
31) $R D_i = \gamma D_i(A) = \gamma D_i(\underline{T}_A, C_A) = \gamma \sqrt{(\underline{x}_i - \underline{T}_R)^T C_A^{-1} (\underline{x}_i - \underline{T}_R)} = \gamma \sqrt{(\underline{x}_i - \underline{T}_R)^T \left(\frac{C_A}{\gamma^2}\right)^{-1} (\underline{x}_i - \underline{T}_R)} = \sqrt{(\underline{x}_i - \underline{T}_R)^T C_R^{-1} (\underline{x}_i - \underline{T}_R)} = D_i(\underline{T}_R, C_R)$.

$\gamma^2 C_A^{-1} = \left(\frac{C_A}{\gamma^2}\right)^{-1}$

then the plotted points in the DD plot will cluster tightly about the identity line as $n \rightarrow \infty$.
 If the data are EC with nonbingular $\text{cov}(x)$, but not MVN, then the data will cluster tightly about some line through the origin, usually with slope > 1 .

Know If x_1, \dots, x_n are iid EC

and $\underline{x} = (x_1, \dots, x_p)^T$, then there should be no strong nonlinearities in a scatterplot matrix of x_1, \dots, x_p .



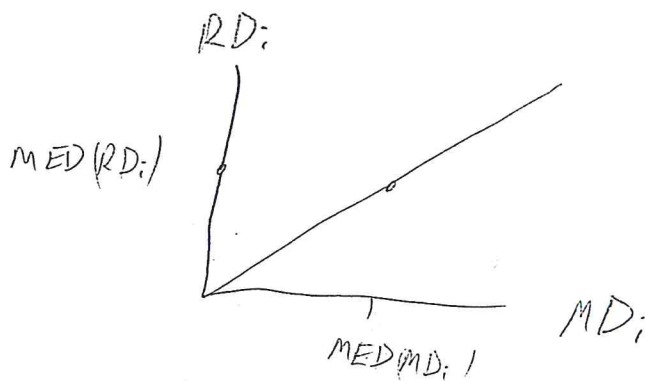
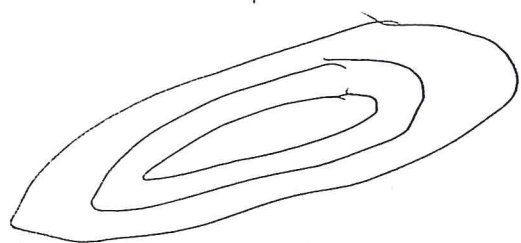
ex) \nwarrow MVN \swarrow EC not MVN EC (see Fig 11.1)

For EC data, the ellipsoids

$$E_{RH} = \{ \underline{x} \mid (\underline{x} - \underline{\tau}_R)^T C_R^{-1} (\underline{x} - \underline{\tau}_R) \leq R D_{(H)}^2 \} \text{ and}$$

$$E_{mH} = \{ \underline{x} \mid (\underline{x} - \bar{\underline{x}})^T S^{-1} (\underline{x} - \bar{\underline{x}}) \leq M D_{(H)}^2 \} \text{ approximate}$$

the highest (by 75% if $h \approx .95n$) density regions. Hence the 2 ellipsoids are approximately concentric.



39) In the DD plot, points below $RD_{(h)}$ correspond to points in E_{Rh} while points to the left of $MD_{(h)}$ correspond to points in E_{Mh} .

3b) know p 128 If (T_A, C_A) is a robust estimator, the DD plot is useful for detecting multivariate outliers.

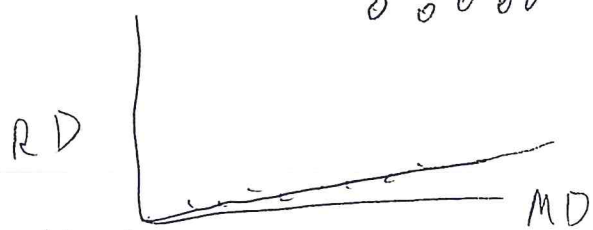
ex) ^{p 128} Buxton data $x_1 = \text{head length}$ $x_2 = \text{nasal height}$

$x_3 = \text{bigonal breadth}$ $x_4 = \text{cephalic index}$

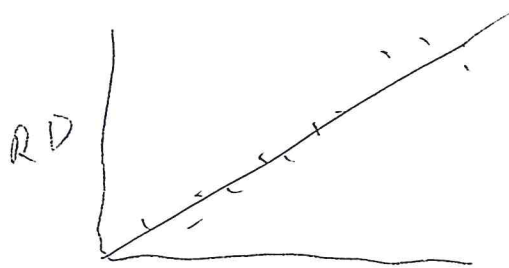
$$\underline{x}_i = (x_{1i}, x_{2i}, x_{3i}, x_{4i})^T$$

Over 5 feet long,

0 0 0 0 0



DD plot

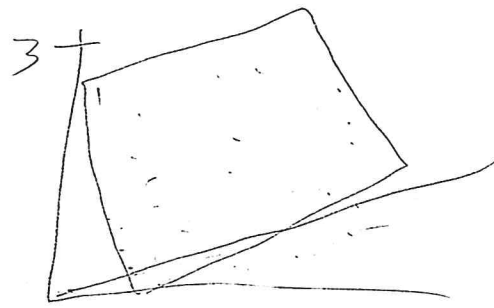
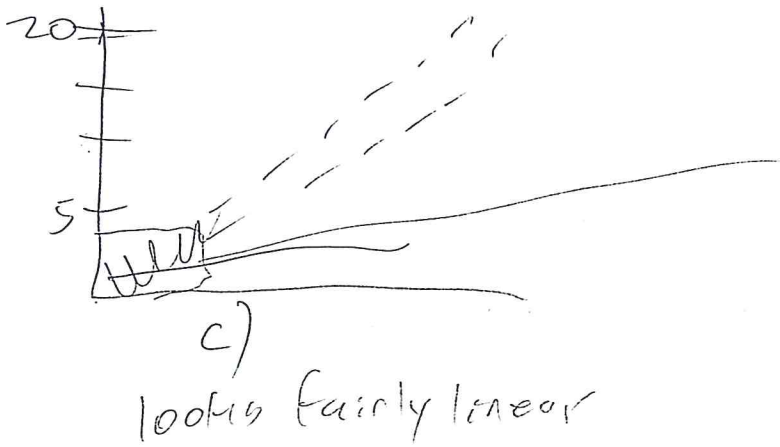


DD plot with outliers deleted

Fig 12.2

Southwest corner of the plot. This can be useful to see whether the data is 26.9 tightly correlated about some line.

ex Fig 3.4 c) d) on p126



Weighted DD plot shows data is not tightly clustered about any line. The data dist. is not EC.

§ B.10 58) - Outlier detection if $p > 1$:

More than $\frac{1}{2}$ cases are in the bulk of the data.

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

b) Let MED_j be the coordinatewise median computed from the cases X_i , with

$D_i^2 \leq \text{MED}(D_i^2(\text{MED}_{j-1}, IP))$ where $\text{MED}_0 = \text{MED}(W)$. Often use $j = 0$ or 9 .

Let $D_i = D_i(\text{MED}_j, IP)$. Let

58327)

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

The covmb2 set B consists of the $m \geq \frac{n}{2}$ cases with weight $w_i = 1$. The covmb2 estimator (T, C) has $T = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$

$$C = \frac{\sum_{i=1}^n w_i (x_i - T)(x_i - T)^T}{\sum_{i=1}^n w_i - 1}$$

which is the

sample mean and covariance matrix computed from the cases in B.

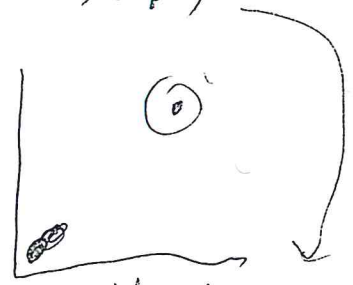
This is like concentration using hyperspheres:

The dispersion matrix IP does not change.

59] Can make a DD plot of $D_i(\text{MED}(w), IP)$

$$US \ D_i(\text{MED}_j, IP) = D_i(T_{\text{covmb2}}, IP)$$

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



- 60] RMVN! RMVN Set U
- RFCH! RFCH Set V
- COVMB2! COVMB2 Set B

the outliers. To detect the outliers, the distance of the outliers from the bulk of the data increases roughly with \sqrt{p} .

contain $\geq \frac{n}{2}$ cases and (\bar{x}_A, s_A)

27.9

$= (\bar{x}_A, s_A)$ is computed from those cases.

$$(T, C) = (\bar{x}_A, s_A) ; d_A > 0, n > 0,$$

61) The two main methods to handle

outliers are a) apply the classical statistical method to the cleaned data

b) plug in robust estimators for classical estimators: eg. (T, C) for (\bar{x}, s) ,

62) The multivariate method applied to sets U and V is often the classical method applied to the cleaned cases in the sets and the plug in method using (\bar{x}_A, s_A) instead of (\bar{x}, s) .

63) Robust regression! Let $\underline{w}_i = (1, \underline{x}_i^T)^T$ 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 method to the m cases \underline{w}_i corresponding to the set $D = \{i_1, \dots, i_m\}$ applied

to $\underline{u}_1, \dots, \underline{u}_m$: MLR, GLMs, GAMs, LDA, QDA, F, NN etc,