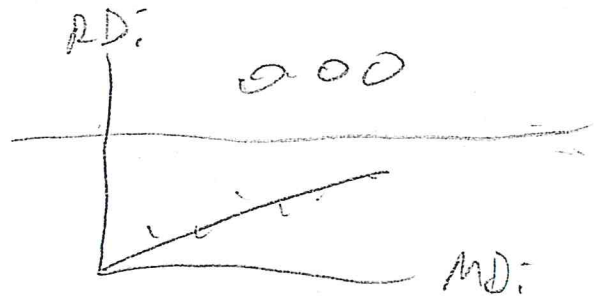


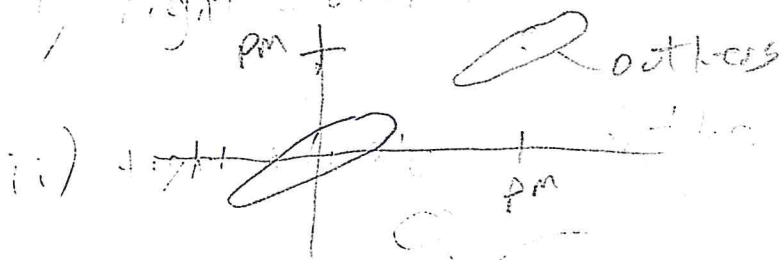
64) To compare outlier resistance of MB, FCH, RFCH, RMVN, DBK, MBA, FMCD etcetera, generate M outlier data sets and count how many times the RD_i for all of the outliers \rightarrow largest RD_i of the clean data. Then the outliers can be separated from the clean data with a horizontal line in the DD plot.



Let the data set have n cases and let $\gamma =$ outlier percentage. Counts near M with γ near but less than 0.5 are good.

65) Clean cases: $\underline{X} \sim N(\underline{0}, \text{diag}(1, \dots, p))$.

a) mean shift outliers $\underline{X} \sim N(\underline{\mu}, \text{diag}(1, \dots, p))$



b) $\underline{X} \sim N(\underline{0}, \text{diag}(1, \dots, 0, pm))$, $0.001 I_p$ = near point mass



→ ... then ...

are like $\underline{x} \sim N_p(\underline{\mu}, \frac{p \times p}{J_A})^T, dI_p$

type 4) outliers

P	γ	Λ	PM	MBA	FCH	R FCH	RMVN	DGK	FMU	MB
10	.1	100	40	73	92	92	92	100	95	100
10	.25	100	25	0	99	99	90	0	0	99
60	.4	200	150	0	100	100	100	0	0	100

MBA does not use location criterion and selects DGK. FCH uses location criterion and usually selects MB

see AWS: midsim 6

1) know Central Limit Theorem (CLT):

Let (Y_1, \dots, Y_n) be iid with $E(Y_i) = \mu$ and $V(Y_i) = \sigma^2$,

Then $\sqrt{n}(\bar{Y}_n - \mu) \xrightarrow{D} N(0, \sigma^2)$.

2) If $W_n \xrightarrow{D} X$, then X is the asymptotic distribution or limiting distribution of W_n .

X does not depend on n . The approximate distribution of \bar{Y} is $\bar{Y} \sim AN(\mu, \frac{\sigma^2}{n}) \Rightarrow$

$\bar{Y} \approx N(\mu, \frac{\sigma^2}{n})$, which does depend on n .

3) Let $\{Z_n, n=1, 2, \dots\}$ be a sequence of RVs with CDFs F_n and let X be a RV with CDF F .

Then Z_n converges in dist to X , written

$Z_n \xrightarrow{D} X$, if $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ at each continuity point of F .

4) X_n converges in probability to X , written

$X_n \xrightarrow{P} X$, if $\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0 \quad \forall \epsilon > 0$

$\Leftrightarrow \lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1 \quad \forall \epsilon > 0$.

5) A sequence of estimators T_n is a consistent estimator if



$\lim_{n \rightarrow \infty} U_{\theta}(\underline{T}_n) = 0$ and $\lim_{n \rightarrow \infty} E_{\theta}(\underline{T}_n) = \underline{T}(\theta)$ for all $\theta \in \Theta$ (293)
 both $\forall \theta \in \Theta$, then \underline{T}_n is a consistent estimator of $\underline{T}(\theta)$.

Multivariate limit theorems

\Rightarrow Let \underline{X}_n have (joint) CDF $F_n(\underline{x})$ and \underline{X} have (joint) CDF $F(\underline{x})$.

a) $\underline{X}_n \xrightarrow{D} \underline{X}$ if $F_n(\underline{x}) \rightarrow F(\underline{x})$ at all continuity points \underline{x} of F .

b) $\underline{X}_n \xrightarrow{P} \underline{X}$ if $\forall \epsilon > 0, P(\|\underline{X}_n - \underline{X}\| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$.

\Rightarrow Multivariate CLT (MULT): If $\underline{X}_1, \dots, \underline{X}_n$ are iid $k \times 1$ random vectors with $E(\underline{X}) = \underline{\mu}$ and $\text{cov}(\underline{X}) = \underline{\Sigma}_X$, then $\sqrt{n}(\underline{\bar{X}} - \underline{\mu}) \xrightarrow{D} N_k(\underline{\mu}, \underline{\Sigma}_X)$.

a) If estimator $\underline{g}(\underline{T}_n) \xrightarrow{P} \underline{g}(\theta) \forall \theta \in \Theta$, then $\underline{g}(\underline{T}_n)$ is a consistent estimator of $\underline{g}(\theta)$.

b) If $\sqrt{n}(\underline{g}(\underline{T}_n) - \underline{g}(\theta)) \xrightarrow{D} \underline{X}$, then $\underline{g}(\underline{T}_n) \xrightarrow{P} \underline{g}(\theta)$.

c) If $\underline{X}_n \xrightarrow{P} \underline{X}$ then $\underline{X}_n \xrightarrow{D} \underline{X}$.

d) $\underline{X}_n \xrightarrow{P} \underline{g}(\theta)$ if $\underline{X}_n \xrightarrow{D} \underline{g}(\theta)$.

10) CONTINUOUS MAPPING THEOREM.

Let g be a continuous function $g: \mathbb{R}^k \rightarrow \mathbb{R}^d$.

If $\underline{x}_n \xrightarrow{D} \underline{x}$, then $g(\underline{x}_n) \xrightarrow{D} g(\underline{x})$.

11) Know for exam 2

a) If $\sqrt{n}(\underline{T}_n - \underline{\mu}) \xrightarrow{D} N_p(\underline{0}, \Sigma)$

and A is a $q \times p$ constant matrix, then

$$A \sqrt{n}(\underline{T}_n - \underline{\mu}) \xrightarrow{D} N_q(A\underline{0}, A\Sigma A^T)$$

b) Let $\Sigma > 0$ be positive definite.

If $\sqrt{n}(\underline{T}_n - \underline{\mu}) \xrightarrow{D} N_p(\underline{0}, \Sigma)$ and \underline{d}_n is a consistent estimator of Σ , then

$$D_{\underline{T}_n}(\underline{\mu}, \underline{d}_n) = n(\underline{T}_n - \underline{\mu})^T \underline{d}_n^{-1} (\underline{T}_n - \underline{\mu}) \xrightarrow{D} \chi_p^2$$

c) If $\Sigma > 0$, $\underline{T}_n \xrightarrow{D} \underline{\mu}$ and $\underline{d}_n \xrightarrow{D} \Sigma$, then

$$\begin{aligned} D_{\underline{x}}^2(\underline{T}_n, \underline{d}_n) &= (\underline{x} - \underline{T}_n)^T \underline{d}_n^{-1} (\underline{x}_n - \underline{T}_n) \xrightarrow{D} D_{\underline{x}}^2(\underline{\mu}, \Sigma) \\ &= (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}). \quad \text{So} \end{aligned}$$

$$D_{(\underline{T}_n, \underline{d}_n)}^2 \xrightarrow{P} D_{(\underline{\mu}, \Sigma)}^2 \text{ where } P(D_{\underline{x}}^2(\underline{\mu}, \Sigma) \leq D_{\underline{x}, 1-\alpha}^2) = 1-\alpha$$

↑
often $\chi_{p, 1-\alpha}^2$

Remark: Let $\underline{z}_n = \sqrt{n}(\underline{T}_n - \underline{\mu}) \rightarrow N_p(\underline{0}, \Sigma) = \underline{z} \quad \left. \vphantom{\underline{z}} \right\} 305$

Then i) $A \underline{z} \sim N_K(A\underline{0}, A \Sigma A^T)$

and ii) $A \underline{z}_n \xrightarrow{D} N_K(A\underline{0}, A \Sigma A^T) \stackrel{D}{=} A \underline{z}$.
can't depend on n

So i) and ii) behave similarly:

ex) see HW 5 11.33

ex) Suppose $\underline{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix}$ where

$x_{i1} \sim \text{pois}(\lambda_1) \perp x_{i2} \sim \text{pois}(\lambda_2)$,

If $w \sim \text{pois}(\lambda)$, $E(w) = V(w) = \lambda$.

Find the limiting dist of $\sqrt{n}(\bar{\underline{x}} - \underline{c})$
for appropriate vector \underline{c} .

soln) $\sqrt{n}(\bar{\underline{x}} - E(\underline{x})) \xrightarrow{D} N_p(\underline{0}, \Sigma_x)$

$E(\underline{x}) = \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix} = \underline{c}$, $\Sigma_x = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$ since

$x_{i1} \perp x_{i2} \Rightarrow \text{cov}(x_{i1}, x_{i2}) = 0$. So

$\sqrt{n}(\bar{\underline{x}} - \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix}) \xrightarrow{D} N_2\left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}\right]$ by MCLT

1) A large sample $100(1-\delta)\%$ prediction region for \underline{x} is a set A_n such that

$P(\underline{x} \in A_n)$ is eventually bounded below by $1-\delta$ as $n \rightarrow \infty$.

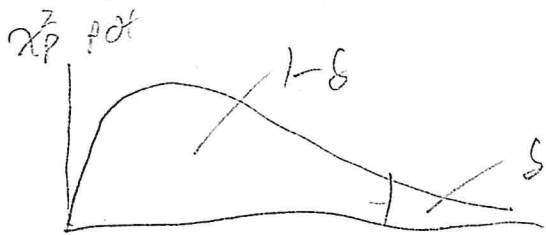
2) An asymptotically optimal prediction region has volume that converges to the $100(1-\delta)\%$ highest density region, and we often want $P(\underline{x} \in A_n) \rightarrow 1-\delta$ as $n \rightarrow \infty$.

3) The classical MVN prediction region is

$$\{ \underline{z} : (\underline{z} - \bar{x})^T S^{-1} (\underline{z} - \bar{x}) \leq \chi_{p, 1-\delta}^2 \}$$

$$= \{ \underline{z} : D_{\underline{z}}^2(\bar{x}, S) \leq \chi_{p, 1-\delta}^2 \} \quad \text{where}$$

$$P(X \leq \chi_{p, 1-\delta}^2) = 1-\delta \quad \text{if } X \sim \chi_p^2$$



This region is $\chi_{p, 1-\delta}^2$ asymptotically optimal if

$\underline{x}_1, \dots, \underline{x}_n, \underline{x}_{n+1}$ are iid $N_p(\mu, \Sigma)$

training data: (\bar{x}, S)

$$\cup \text{Let } 0\alpha = \min\left(\frac{1-\alpha}{2}, \frac{1-\delta}{2}\right), \delta \geq 0.1$$

$$\min\left(\frac{1-\delta}{2}, \frac{1-\delta}{2} + \frac{10\delta P}{n}\right) \delta \leq 1.$$

If $\delta_n < \frac{1-\delta}{2} + 0.001$, use $\delta_n = \frac{1-\delta}{2}$.

Let $D_{(n)}^2$ be the $100\delta_n$ th percentile of $D_i^2(\bar{x}, S)$. The $100(1-\delta)\%$ large sample nonparametric prediction region

$$\text{is } \left\{ \underline{z} : D_{\underline{z}}^2(\bar{x}, S) \leq D_{(n)}^2 \right\} =$$

$$\left\{ \underline{z} : (\underline{z} - \bar{x})^T S^{-1} (\underline{z} - \bar{x}) \leq D_{(n)}^2 \right\}.$$

5) This region replaces $\chi_{p, 1-\delta}^2$ by $D_{(n)}^2$ and is a large sample prediction region if x_1, \dots, x_n are iid with nonsingular covariance matrix Σ_x .

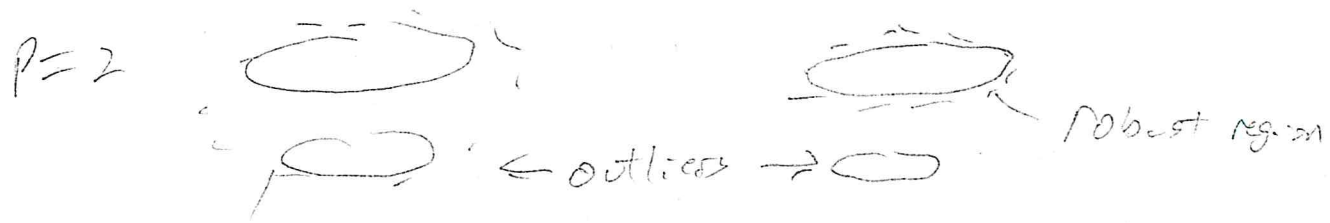
6) The semiparametric prediction region

$$\text{is } \left\{ \underline{z} : D_{\underline{z}}^2(T_{RMVN}, C_{RMVN}) \leq D_{(n)}^2(T_{RMVN}, C_{RMVN}) \right\}$$

→ Both 4) and 6) are asymptotically

optimal on a large class of EC distributions with nonsingular covariance matrices.

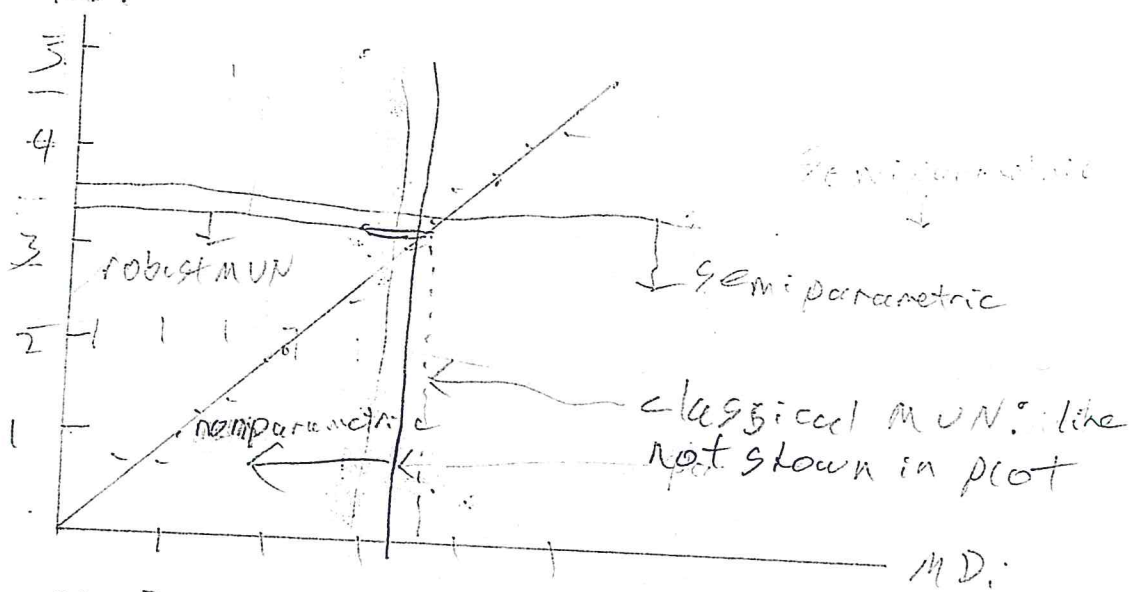
8) The semi-parametric prediction region can tolerate about $100(1-\alpha)^{1/p}$ outliers.



classical 95% region

RD:

ex)



- $p=5, n=82$ Nonparametric $MD_i \leq D_{(.95, n)} \approx 3.09$
- semi-parametric $RD_i \leq 3.44$
- robust MUN $RD_i \leq 3.33$
- classical $MD_i \leq 3.33$

4.2 Bootstrap confidence regions

9) Non parametric bootstrap:

Let data be Z_1, \dots, Z_n , often iid.

Draw a sample of size n Z_1^*, \dots, Z_n^*

with replacement from Z_1, \dots, Z_n

COMPUTE $T_1^* = T_1^* = T_1^* = \dots = T_n^*$, and repeat
B times to get a bootstrap sample.

T_1^*, \dots, T_B^* .

10) ch 2 did this for random variables.
see quiz 2 and Hw 2.

11) Let (\bar{T}^*, S^*) be the sample mean and
covariance matrix computed from T_1^*, \dots, T_B^* .

$$\text{Let } g_B = \begin{cases} \text{mean} \left(1 - \frac{\delta}{2}, 1 - \delta + 10 \frac{\delta \delta}{B} \right), & \delta \leq 0.1 \\ \text{mean} \left(1 - \delta + 0.05, 1 - \delta + \frac{\delta}{B} \right) & \delta > 0.1 \end{cases}$$

If $t_\delta \leq 0.999$ and $g_B \leq 1 - \delta + 0.001$, set $g_B = t_\delta$.

12) A large sample $100(1-\delta)\%$ confidence region
for θ is a set A_n such that $P(\theta \in A_n)$
is eventually bounded below by $1-\delta$ as $n \rightarrow \infty$.

13) often want $P(\theta \in A_n) \rightarrow 1-\delta$ as $n \rightarrow \infty$.

14) If the actual coverage $100(t_{\delta n}) < 100(1-\delta)$,
the region is liberal. If $100(t_{\delta n}) > 100(1-\delta)$,
the region is conservative. A region that
is 5% conservative is "much better" than a
region that is 5% liberal.

15) Let $D_{(1-\delta)}$ be the $100(1-\delta)$ th sample quantile
of the D_i .

		Sample G.
--	--	----------------------