

Mathematical Statistics II (MAT323)

Ch6. Point Estimation

Jungsoon Choi

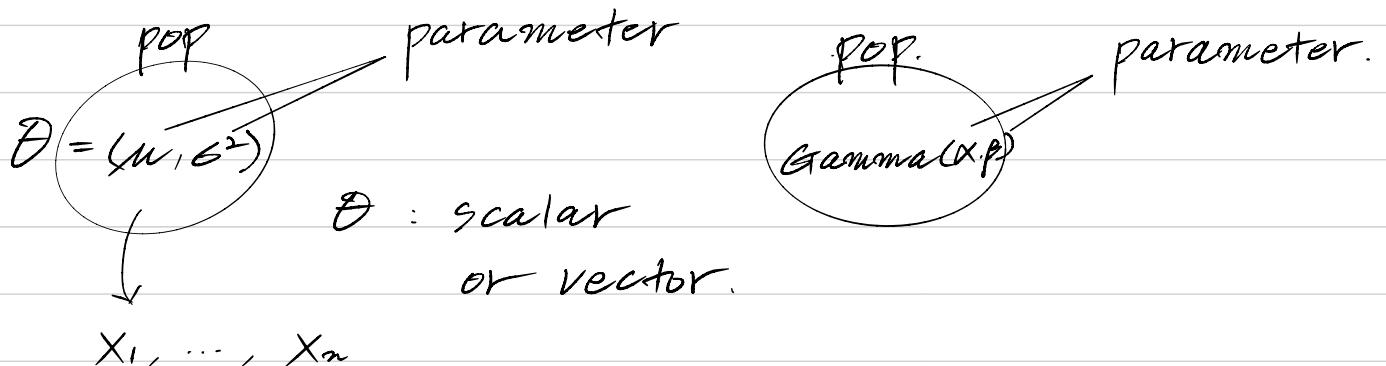
jungsoonchoi@hanyang.ac.kr

Table of Contents

- Order Statistics
- Q-Q plot

Ch6.3 Order Statistics

parameter : (unknown) interest in pop.



$$\begin{cases} \bar{x}_n & x_n = 2.5 \\ & \text{point estimation} \\ & C.I \text{ estimation.} \\ & (3.4, 3.6) \end{cases}$$



Testing.

Order statistics

Definition

Let X_1, \dots, X_n denote a random sample from a continuous distribution of $f(x)$. A random sample Y_1, \dots, Y_n is called the **order statistics** of X_1, \dots, X_n if

$$Y_1 < Y_2 < \cdots < Y_n$$

$$X_{(1)} < X_{(2)} < \cdots < X_{(n)}$$

Y_i or $X_{(i)}$ is called the i th order statistic of the random sample X_1, \dots, X_n .

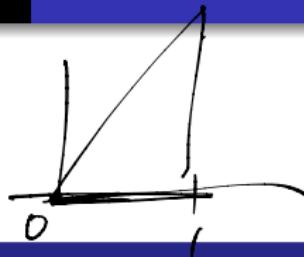
POP $\rightarrow X_1, \dots, X_n$ random sample.

2 가지 statistic

$$X_{(1)} = Y_1 < X_{(2)} = Y_2 < \dots < X_{(n)} = Y_n$$
$$\min_{i=1, \dots, n} \{x_i\}$$
$$\max_{i=1, \dots, n} \{x_i\}$$

$X_{(1)}, \dots, X_{(n)}$: order statistic

want to find distⁿ of order statistic



Example 6.3-2

$Y_1 < Y_2 < \dots < Y_5$ are order statistics of X_1, \dots, X_5 with a pdf $f_x(x) = 2x, 0 < x < 1$.

- Find the probability of $P(Y_4 \leq 0.5)$.
- Find the cdf of Y_4 , $F_{Y_4}(y) = P(Y_4 \leq y), 0 < y < 1$.
- Find the pdf of Y_4 , $f_{Y_4}(y)$.

$$f_X(x) = 2x \quad 0 < x < 1.$$

$$1) P(Y_4 \leq 0.5).$$

$$2) F_{Y_4}(y). \quad 0 < y < 1.$$

$$3) f_{Y_4}(y).$$

X_1, \dots, X_5

$Y_1 < \dots < Y_5. \quad Y_4 \leq 0.5 \Leftrightarrow P(X \leq 0.5)$, at least 4.

$$1) P_s = P(X \leq 0.5) = \int_0^{0.5} 2x \, dx = [x^2]_0^{0.5} = \left(\frac{1}{4}\right).$$

$$Z \sim \text{Bin}(5, \frac{1}{4}).$$

$$\begin{aligned} P(Y_4 \leq 0.5) &= P(Z=4) + P(Z=5) \\ &= \binom{5}{4} \left(\frac{1}{4}\right)^4 \left(\frac{3}{4}\right) + \left(\frac{1}{4}\right)^5 \end{aligned}$$

$$2) Z \sim \text{Bin}(5, y^2).$$

$$F_{Y_4}(y) = \binom{5}{4} (y^2)^4 (1-y^2) + (y^2)^5$$

$$3) f_{Y_4}(y) = \cancel{\frac{20}{5} (y^2)^3 (2y)(1-y^2)} + \cancel{\frac{5}{5} (y^2)^4 (2y)} \\ + 5 \cdot (y^2)^4 (-2y)$$

$$= \binom{5}{3} (F_X(y))^3 \cdot f_X(y) \cdot \binom{2}{1} (1 - F_X(y))$$

$$= 20 (y^2)^3 \cdot 2y \cdot (1 - y^2).$$

$$X_i \sim f_X(x) = 2x, 0 < x < 1 \quad Y_1 < Y_2 < Y_3 < Y_4 < Y_5 \\ < \frac{1}{2}$$

$\{Y_4 \leq \frac{1}{2}\}$ happens

\Leftrightarrow {at least 4 r.v.s among $X_1, \dots, X_5 \leq \frac{1}{2}$ }

Assume $\{X_i \leq \frac{1}{2}\}$ is a "success event"

prob of success \downarrow

Let Z be the number of success events $\sim \text{Bin}(5, p_s)$

$$p_s = P(X_i \leq \frac{1}{2}) = F_{X_i}(\frac{1}{2}) = \int_0^{\frac{1}{2}} 2x \, dx = [x^2]_0^{\frac{1}{2}} = \frac{1}{4} \\ \therefore Z \sim \text{Bin}(5, \frac{1}{4}).$$

$$1) P(Y_4 \leq \frac{1}{2}) = P(\text{at least 4 r.v.s among } X_1, \dots, X_5 \leq \frac{1}{2})$$

$$= P(4 \text{ r.v.s among } X_i \leq \frac{1}{2}) + P(5 \text{ r.v.s among } X_i \leq \frac{1}{2})$$

$$= P(Z=4) + P(Z=5)$$

$$= \left(\frac{5}{4}\right)\left(\frac{1}{4}\right)^4\left(\frac{3}{4}\right)^1 + \left(\frac{1}{4}\right)^5$$

2) consider $\{X_i \leq y\}$ as a "success event"

$$Z \sim \text{Bin}(5, P(X_i \leq y)) = F_{X_i}(y) = \int_0^y 2x \, dx = y^2$$

$$\underline{P(Y_4 \leq y)} = P(Z=4) + P(Z=5)$$

$$= F_{Y_4}(y)$$

$$= \left(\frac{5}{4}\right)(y^2)^4(1-y^2) + (y^2)^5, 0 < y < 1$$

$$ii) f_{Y_4}(y) = \left(\frac{5}{4}\right) 4(y^2)^3(2y)(1-y^2) + \left(\frac{5}{4}\right)(y^2)^4(-2y) + 5(y^2)^4(2y)$$

$$= 20(y^2)^3(2y)(1-y^2) \quad y^2 = F_X(y)$$

Since $F_X(y) = y^2$

$$\frac{f_{Y_4}(y)}{\text{pdf}} = \frac{\frac{5!}{3!2!} \times \frac{2!}{2!} \times [F_X(y)]^3 [1-F_X(y)]' \times f_X(y)}{5C_5} \quad 0 < y < 1$$

*(commonalities)
x thus same // pdf.*

$P(Y_4 = y)$

Next job : generalize the cdf & pdf of order statistic



cdf of order statistics

Suppose that $Y_1 < Y_2 < \dots < Y_n$ are order statistics of continuous r.v X_1, \dots, X_n with a pdf $f_X(x)$ and a cdf $F_X(x)$, $a < x < b$. Since $P(X_i \leq y) = F_X(y)$,

$$F_{Y_r}(y) = P(Y_r \leq y) = \sum_{k=r}^n \binom{n}{k} [F_X(y)]^k [1 - F_X(y)]^{n-k}$$

$$= \sum_{k=r}^{n-1} \binom{n}{k} [F_X(y)]^k [1 - F_X(y)]^{n-k} + [F_X(y)]^n$$



pdf of order statistics

$$\begin{aligned}
 f_{Y_r}(y) &= \sum_{k=r}^{n-1} \binom{n}{k} (k) [F_X(y)]^{k-1} f_X(y) [1 - F_X(y)]^{n-k} \\
 &\quad + \sum_{k=r}^{n-1} \binom{n}{k} [F_X(y)]^k (n-k) [1 - F_X(y)]^{n-k-1} [-f_X(y)] \\
 &\quad + n [F_X(y)]^{n-1} f_X(y)
 \end{aligned}$$

With some calculations,

$$f_{Y_r}(y) = \frac{n!}{(r-1)!(n-r)!} [F_X(y)]^{r-1} [1 - F_X(y)]^{n-r} f_X(y)$$

$$a < y < b$$

$$X_i \sim f_x(x), F_x(x), i=1, \dots, n$$

$$F_{Y_r}(y) = P(Y_r \leq y) = P(\text{at least } r \text{ rvs among } X_i \leq y) \\ = P(Z=r) + P(Z=r+1) + \dots + P(Z=n).$$

where $Z \sim \text{Bin}(n, P(X_i \leq y) = F_x(y))$

$$= \sum_{k=r}^{n-1} \binom{n}{k} [F_x(y)]^k [1 - F_x(y)]^{n-k} + [F_x(y)]^n$$

$\underbrace{\dots}_{\text{pdf}} \sum_{k=r}^{n-1} \binom{n}{k} [F_x(y)]^k [1 - F_x(y)]^{n-k}$

pdf \Rightarrow

$$f_{Y_r}(y) = \frac{\partial F_{Y_r}(y)}{\partial y} = \sum_{k=r}^{n-1} \binom{n}{k} k [F_x(y)]^{k-1} f_x(y) [1 - F_x(y)]^{n-k}$$

$$+ \sum_{k=r}^{n-1} \binom{n}{k} (n-k) [F_x(y)]^k [1 - F_x(y)]^{n-(k+1)} (-f_x(y)) + n [F_x(y)]^{n-1} f_x(y).$$

$$\text{Since } \binom{n}{k} k = \frac{n!}{(k-1)! (n-k)!}, \quad \& \quad \binom{n}{k} (n-k) = \frac{n!}{k! (n-(k+1))!}$$

$$f_{Y_r}(y) = \frac{n!}{(r-1)! (n-r)!} [F_x(y)]^{r-1} [1 - F_x(y)]^{n-r} f_x(y)$$

$$+ \underbrace{\sum_{k=r}^{n-1} \binom{n}{k} \frac{n!}{(k-1)! (n-k)!} [F_x(y)]^{k-1} [1 - F_x(y)]^{n-k} f_x(y)}$$

$$- \sum_{k=r}^{n-2} \frac{n!}{k! (n-(k+1))!} [F_x(y)]^k [1 - F_x(y)]^{n-(k+1)} f_x(y)$$

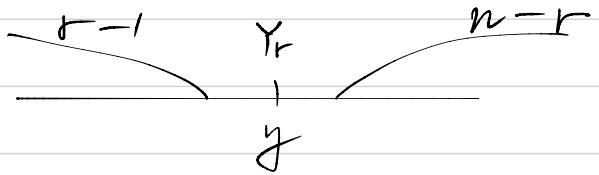
$$- \frac{n!}{(n-1)!} [F_x(y)]^{n-1} f_x(y) + n [F_x(y)]^{n-1} f_x(y) = 0.$$

$$= \frac{n!}{(r-1)! (n-r)!} [F_x(y)]^{r-1} [1 - F_x(y)]^{n-r} f_x(y)$$

$$+ \sum_{k=r}^{n-2} \frac{n!}{k'! (n-(k'+1))!} [F_x(y)]^{k'} [1 - F_x(y)]^{n-(k'+1)} f_x(y)$$

$$= 0 \quad - \sum_{k=r}^{n-2} \frac{n!}{k! (n-(k+1))!} [F_x(y)]^k [1 - F_x(y)]^{n-(k+1)} f_x(y).$$

$$= \frac{n!}{(r-1)! (n-r)!} [F_x(y)]^{r-1} [1 - F_x(y)]^{n-r} f_x(y). \quad a < y < b.$$



$$\begin{array}{ccc} r-1 & r.v & < y \\ n < & n-r & r.v > y \\ & / & r.v = y \end{array}$$

같은 이유로 개념적으로 이해.

$$f_{Y_T=y} \neq P(Y_T = y).$$

$$P(X > y)$$

$$P(X = y).$$

$$\begin{aligned}
 & P(X < y) = P(X > y) + P(X = y). \\
 & = \binom{n}{r-1} [F_X(y)]^{r-1} \binom{n-r+1}{n-r} [1 - F_X(y)]^{n-r} \cdot f_X(y). \\
 & = \frac{n!}{(r-1)! (n-r+1)!} \frac{(n-r+1)!}{(n-r)!} [F_X(y)]^{r-1} [1 - F_X(y)]^{n-r} f_X(y).
 \end{aligned}$$

$$f_{Y_r}(y) = \frac{n!}{(r-1)!(n-r)!} [F_X(y)]^{r-1} [1-F_X(y)]^{n-r} \cdot f_X(y).$$



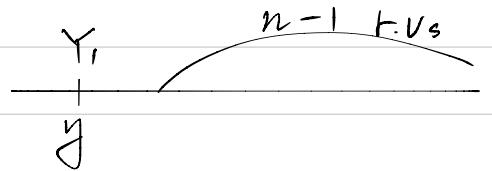
pdfs of Y_1 and Y_n

$$f_{Y_1}(y) = n [1 - F_X(y)]^{n-1} f_X(y), \quad a < y < b$$

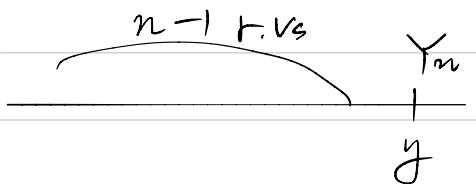
$$f_{Y_n}(y) = n [F_X(y)]^{n-1} f_X(y), \quad a < y < b$$

X_1, \dots, X_n $Y_1 < Y_2 < \dots < Y_n$
min{ X_i } max{ X_i } $f_{Y_1}(y) \neq P(Y_1 = y)$.

$$= \binom{n}{n-1} [1 - F_X(y)]^{n-1} f_X(y).$$

 $P(X > y).$  $f_{Y_n}(y) \neq P(Y_n = y)$.

$$= \binom{n}{n-1} [F_X(y)]^{n-1} f_X(y).$$



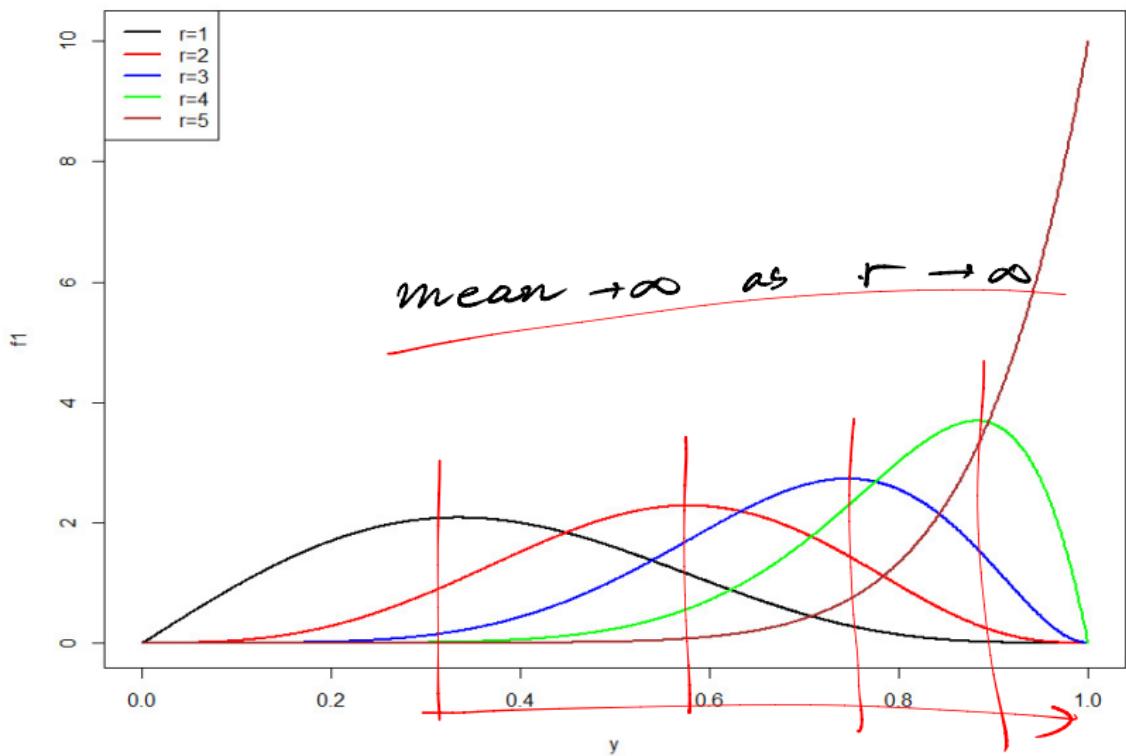
Example 6.3-3

In Example 6.3-2, find the pdfs of each Y_1, \dots, Y_5 .

$$f_{Y_1} = 5 \cdot (1 - F_x(y))^4 \cdot f_x(y) = 5 \cdot (1 - y^2)^4 \cdot 2y$$

$$f_{Y_3} = \frac{5!}{2!3!} [F_x(y)]^2 [1 - F_x(y)]^2 f_x(y) = 10 \cdot (y^2)^2 (1 - y^2)^2 y$$

$$f_{Y_5} = 5 \cdot (F_x(y))^4 \cdot f_x(y) = 5 \cdot (y^2)^4 \cdot 2y$$

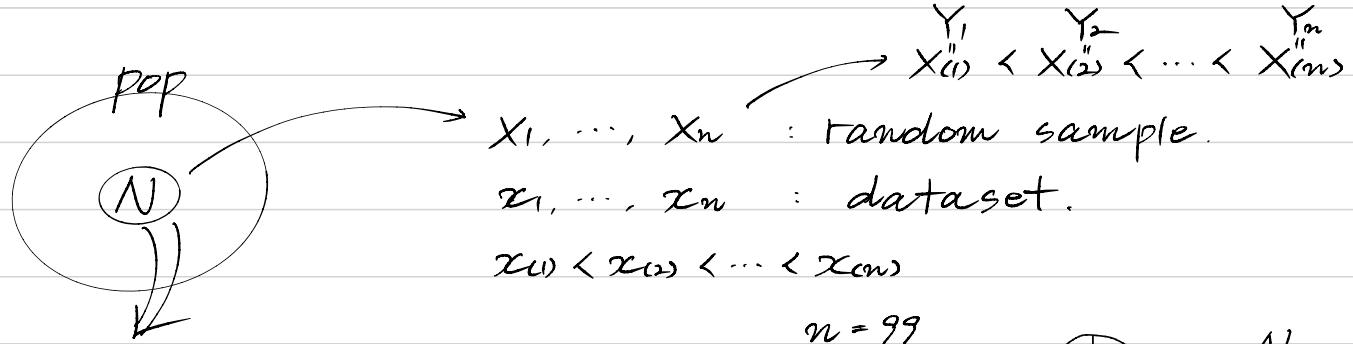
pdfs of order statistics of X , $f(x)=2x$ 

[Review] π_p , (100p)th percentile

The (100p)th percentile is a number π_p such that the area under $f_X(x)$ to the left of π_p is p .

$$p = \int_{-\infty}^{\pi_p} f_X(x)dx = F_X(\pi_p).$$

- median=50th percentile=the second quartile, $m = \pi_{0.5}$
- the first quartile=25th percentile, $q_1 = \pi_{0.25}$
- the third quartile=75th percentile, $q_3 = \pi_{0.75}$



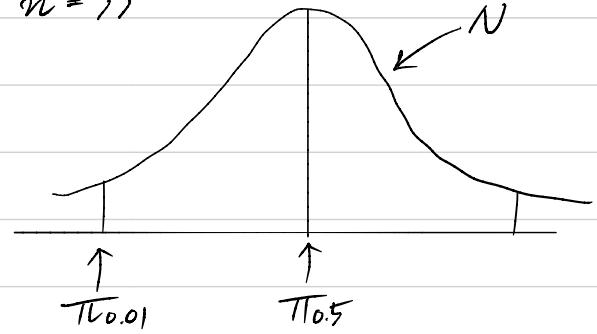
$T_{P} = (100p)^{th}$ percentile

$T_{0.01} = 1^{st}$ percentile.

:

$T_{10} = 100^{th}$ "

$n = 99$



Sample percentiles (using the data)

from dataset.

Sample Percentiles ($0 < p < 1$)

The $100p^{\text{th}}$ sample percentile ($= Y_r$) is defined as

$$\begin{cases} (n+1)p^{\text{th}} \text{ order statistic} & \text{if } r = (n+1)p \text{ is integer} \\ \text{weighted average of } Y_r \text{ and } Y_{r+1} & \text{if } (n+1)p \text{ is not integer} \end{cases}$$

$r = [(n+1)p]$: the greatest integer

W 99

Notation

$Q1=25\text{th percentile}= \text{First quartile}$

$Q2=50\text{th percentile}= \text{Second quartile}= \text{Median}$

$Q3=75\text{th percentile}= \text{Third quartile}$

$$n = 99 \quad p = 0.01$$

$$(n+1)p = 100 \times 0.01 = 1$$

$(100p)^{th}$ sample percentile

= 1st sample percentile. = 1st order statistic.
2nd " = 2nd " = $x_{(2)}$.

$$99^{th} \quad " \quad = x_{(99)}$$

$$\text{e.g. } n = 10. \quad x_1, \dots, x_{10}$$

$$\Rightarrow x_{(1)} < \dots < x_{(10)}$$

$$p = 0.5$$

$x_{(5)} = (100p)^{th}$ sample percentile.

$$r = 11 \times 0.5 = 5.5 = 5.$$

$$\Rightarrow \frac{y_5 + y_6}{2}$$

$$(12+1)0.25 = \frac{13}{4} = 3.25$$

$$\frac{3}{4}Y_3 + \frac{1}{4}Y_4$$



Example 6.3-4

Let X be the weight of soap in a “1000-gram” bottle. A random sample of $n = 12$ observations of X yielded the following weights, which have been ordered:

1013	1019	1021	1024	1026	1028
1033	1035	1039	1040	1043	1047

36

$$8 \quad 9 \quad \frac{39}{4}$$

$$3 - 0.15 \quad \frac{3}{4}$$

$$Y_3 + (Y_4 - Y_3) \cdot 0.15$$

Find the first, second, and third quartiles.

$$25\text{th percentile} = (12+1)0.25 = \frac{13}{4} = 3.25$$

$$Y_3 + (Y_4 - Y_3) \cdot 0.25$$

$$= 0.75Y_3 + 0.25Y_4$$

Q-Q plot:

To check if the assumed population distribution is reasonable based on the observed samples

Q-Q plot

Q-Q plot

A Q-Q plot is a probability plot which is a graphical method to compare two probability distributions by plotting their quantiles against each other. If the points in the Q-Q plot approximately lie on the straight line, then it is said that two distributions are similar.

sample percentile vs theoretical percentile.
dataset. compare.

To check if the observed samples follow the normal distribution, a Q-Q plot can be used.

Suppose that $x_i \sim N(\mu, \sigma^2)$, $i = 1, \dots, n$.

x-axis: theoretical quantile from the normal distribution

y-axis: data quantile

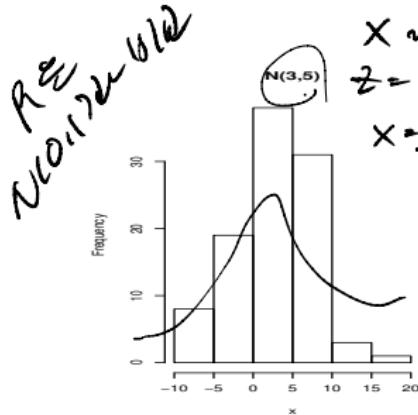
$$\text{plot of } \left(z_i = \Phi^{-1} \left(\frac{i}{n+1} \right), x_{(i)} \right), \quad i = 1, \dots, n$$

where $x_{(i)}$ is the i th ordered sample.

Instead of $(i/(n+1))$, $(i - 0.5)/n$ or $(i - 3/8)/(n + 1/4)$ can be used.

If $z_i \sim N(0, 1)$ and $x_i \sim N(\mu, \sigma^2)$, $x_i = \mu + \sigma z_i$.

To check if the observed samples follow the normal distribution

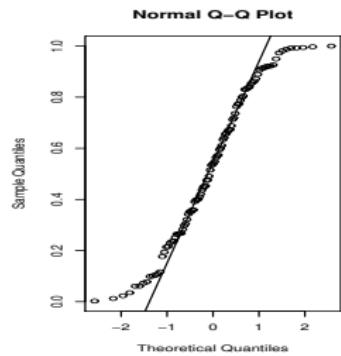
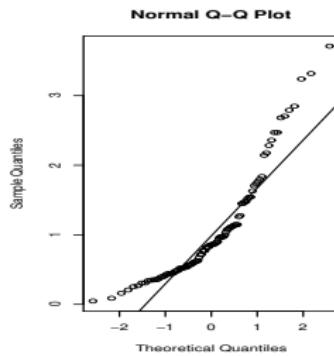
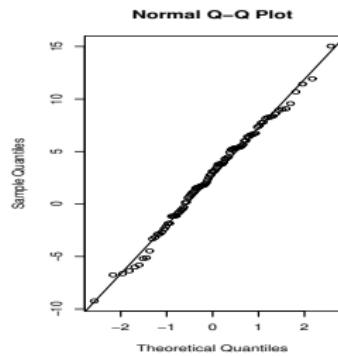
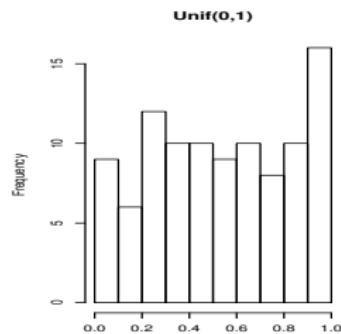
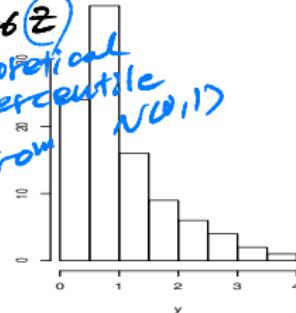


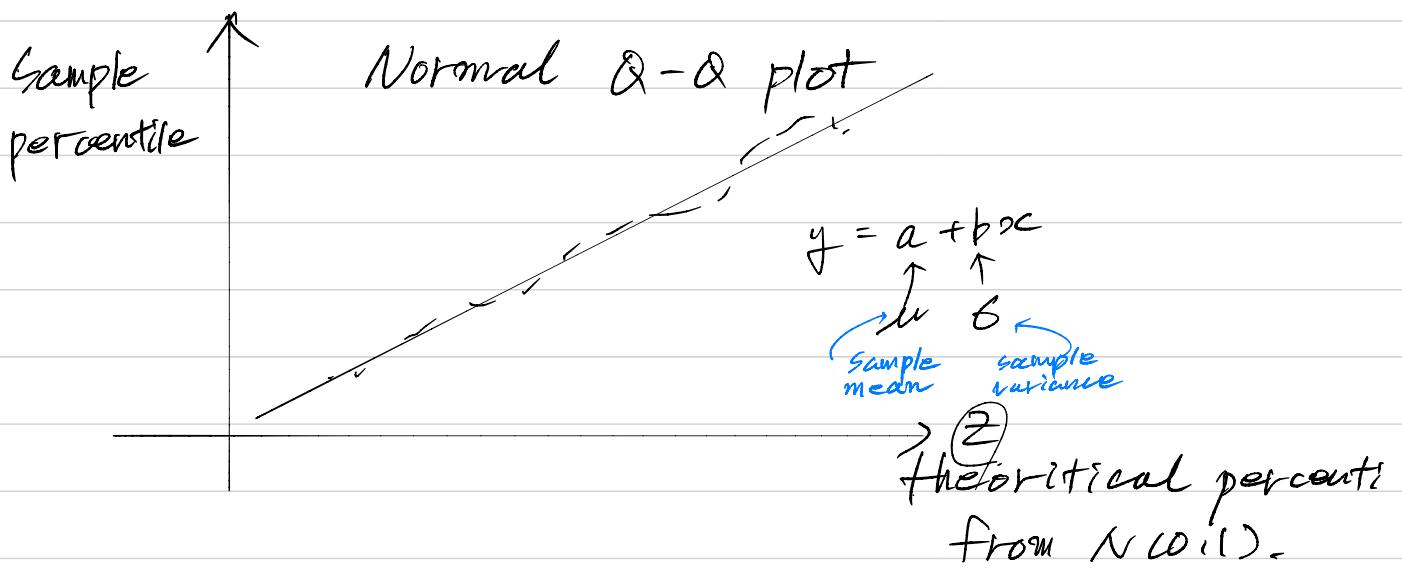
$$X \sim N(\mu=6^2)$$

$$z = \frac{X-\mu}{\sigma} \sim N(0,1)$$

$$X = \mu + \sigma z$$

theoretical percentile
from $N(0,1)$





Example 6.3-5

Here are the diameters (in mm) of 30 randomly selected grains of soil:

1.24	1.36	1.28	1.31	1.35	1.20	1.39	1.35	1.41	1.31
1.28	1.26	1.37	1.49	1.32	1.40	1.33	1.28	1.25	1.39
1.38	1.34	1.40	1.27	1.33	1.36	1.43	1.33	1.29	1.34

For these data, $\bar{x} = 1.33$ and $s^2 = 0.0040$. May we assume that these are observations of a random variable X that is $N(1.33, 0.0040)$?

$$P = \Phi(C\bar{X}^P) \sim \Phi(\bar{X})$$

\bar{X}

To draw a Q-Q plot, we obtain the following table:

k	y (data sorted)	$p = k/31$	z_{1-p}
1	1.20	0.0323	-1.85
2	1.24	0.0645	-1.52
3	1.25	0.0968	-1.30
	...		
29	1.43	0.9355	1.52
30	1.49	0.9677	1.85

maximum
of sample
but not maximum of pop

