**MANE 3332.04**

# Lecture 17, March 31

**Agenda**

- Midterm not graded: still contacting students who missed exam

- Continue working on Technical Report One Assignment

- Chapter Six

- Attendance

- Questions?

## Grades

25% mid-term

○  mid-term  →  $-100(.25) = \underline{-25}$

↓

could make a $\underline{75}$

**Schedule**

| Monday Lecture | Wednesday Lecture |
|---|---|
| 3/31: Chapter 6 | 4/2: Chapter 5 |
| 4/7: Chapter 7 & 8 | 4/9: Chapter 8, Case 1 |
| 4/14: Chapter 8: Case 2 | 4/16: Chapter 8: Case 3 |
| 4/21: Chapter 9, case 1 | 4/23: Chapter 9, Case 2 |
| 4/28: Chpater 9, Case 3 | 4/30: Chapter 11 |
| 5/5: Chapter 11 | 5/7: Review |

**12 classroom sessions plus Final Exam**

**Handouts**

- [Chapter 6 Slides](#)
- Chapter 6 Slides marked

# Data Analysis

1) location of Data $\rightarrow$ mean, median, mode

2) Variability/Spread $\rightarrow$ variance/standard dev.

3) Shape of data



uniform          normal          exponential

# Numerical Summaries

- Called Descriptive Statistics in Chapter 6
  - Descriptive statistics help us understand the location or central tendency of data and the scatter or variability in data
  - Included in all statistical software packages, R does a good job calculating descriptive statistics

**Central Tendency**

- Ostle, et. al. (1996) define central tendency as "the tendency of sample data to cluster about a particular numerical value"
- Population mean

Greek letter → $\mu$

$N$ - population size

$$\mu = \frac{1}{N}\sum_{i=1}^{N} x_i$$

hat notation

- Sample mean

$n$ is now lower-case

$\bar{x}$

$$\bar{x} = \hat{\mu} = \frac{1}{n}\sum_{i=1}^{n} x_i$$

$\hat{\mu}$ - est

estimator of $\mu$

- Sample median - middle value → $\tilde{x}$
- Sample mode - most commonly occuring number(s)

$$Range = X_{max} - X_{min}$$

**Measures of Variability**
- There are several statistics that measure the variability or spread present in data
- Population variance

$$\sigma^2 = \frac{\sum_{i=1}^{N}(x_i - \mu)^2}{N}$$

Calculator: $\xi\ \sigma_N$ or $\sigma_n$

- Sample variance

$$s^2 = \hat{\sigma}^2 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}$$

Calculator: $\sigma_{n-1}$

- Shortcut (Computational) Formula

$$s^2 = \frac{\sum_{i=1}^{n} x_i^2 - \frac{(\sum_{i=1}^{n} x_i)^2}{n}}{n-1}$$

- Standard deviation is often used because it is measured in the original units

$$\sigma = \sqrt{\sigma^2}; \ s = \sqrt{s^2}$$

all columns

midterm has 3 columns

**R Function Su...**

- R code

  `summary...`

- Output is fr...

```
26 ▾ ```{r}
27    summary(midterm)
28 ▾ ```
28:4    # Import Dataset ▾                                      R Markdown ▾

Console    Terminal ×    Render ×    Background Jobs ×

R  R 4.3.1 · /Volumes/SAMSUNG T7/wfscsBackup/Teaching2/AY_2023_2024/MANE3332_spring2024/PartTwo/l...

> knitr::opts_chunk$set(echo = TRUE)
> library(readxl)
> midterm <- read_excel("/Volumes/NO NAME/midterm.xlsx")
> View(midterm)
> summary(midterm)
 Participation       QuizAverage       MidtermExam
 Min.   :  2.941   Min.   :  0.00   Min.   :28.00
 1st Qu.: 67.500   1st Qu.: 60.00   1st Qu.:59.75
 Median : 87.941   Median : 80.00   Median :65.00
 Mean   : 77.096   Mean   : 74.42   Mean   :66.07
 3rd Qu.: 95.147   3rd Qu.: 92.00   3rd Qu.:74.75
 Max.   :100.000   Max.   :100.00   Max.   :92.00
                                     NA's   :5
>
```

Descriptive Statistics

$x = Q_1$

**R Function Su**

- R code
  `summary`

- Output is fr



Descriptive Statistics

**R Funct...**

- Sum...
- Desc...
- Desc...
- R Co...
  - lib...
  - des...
- Psych package output from Spring 2024



Describe() Output

*(handwritten annotations: "imported library", "mean absolute deviation", "forecasting", and symbols $\bar{x}$, $s$, $\tilde{x}$)*

```r
34  ```{r}
35  library(psych)
36  describe(midterm)
37
```

Description: df [3 × 13]

| | vars <dbl> | n <dbl> | mean <dbl> | sd <dbl> | median <dbl> | trimmed <dbl> | mad <dbl> |
|---|---|---|---|---|---|---|---|
| Participation | 1 | 33 | 77.10 | 25.65 | 87.94 | 81.35 | 16.13 |
| QuizAverage | 2 | 33 | 74.42 | 23.49 | 80.00 | 77.48 | 23.72 |
| MidtermExam | 3 | 28 | 66.07 | 13.73 | 65.00 | 66.62 | 13.34 |

3 rows | 1–8 of 13 columns

37:4   # Import Datset ⬍                                   R Markdown ⬍

$$\sigma^2 \sim (x - N)^2 \rightarrow 2^{nd} \text{ order}$$

$$3^{rd} \text{ order} \quad 4^{th} \text{ order}$$

Description: df [3 × 13]

Max

| | median | trimmed | mad | min | m... | range | skew | kurtosis | se |
|---|---|---|---|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 87.94 | 81.35 | 16.13 | 2.94 | 100 | 97.06 | −1.31 | 0.79 | 4.47 |
| | 80.00 | 77.48 | 23.72 | 0.00 | 100 | 100.00 | −1.24 | 1.28 | 4.09 |
| | 65.00 | 66.62 | 13.34 | 28.00 | 92 | 64.00 | −0.46 | 0.39 | 2.59 |

3 rows | 6–14 of 13 columns

Describe Output    Standard error

Chapters 8 & 9    $\dfrac{s}{\sqrt{n}}$

# Calculating Quantiles

$$x_p = x_{(i)} + d[x_{(i+1)} - x_{(i)}]$$

$$p(n+1) = i + d$$

reference for calculating quantiles

1) Sort data

2) find $p = .5^-$ (median)

$X_1 = 6, X_2 = 4, \ldots, X_8 = 6 \rightarrow$ order in data

$X_{(1)} = 3, X_{(2)} = 4, X_{(3)} = 4, \ldots, X_{(8)} = 7$

**Qua**

8 Observations from binomial distribution with

6, 4, 5, 7, 3, 5, 4, 6 $\rightarrow$ 3, 4, 4, 5, 5, 6, 6, 7

sorted

$X_{(4)} \quad X_{(5)}$

Quantile Example

$$p = .5, \quad n = 8$$

$i \rightarrow$ integer

$P(n+1) = \boxed{i + d}$

$d \rightarrow$ decimal. remainder

$$.5(8+1) = \underline{4.5}$$

$\rightarrow i = 4, \quad d = .5$

$$X_{.5} = X_{(4)} + .5[X_{(5)} - X_{(4)}]$$

$$= 5 + .5[5 - 5]$$

$\rightarrow$ what am I doing?

$$= \underline{5}$$

linear interpolation

② add to note card for final

**Exploratory Data (Graphical) Analysis**

- Exploratory data analysis (EDA) is the use of graphical procedures to analyze data.

- John Tukey was a pioneer in this field and invented several of the procedures

- Tools include stem-and-leaf diagrams, box plots, time series plots and digidot plots ⟶ *dot plots*

**Stem and Leaf Diagram**

- Excellent tool that maintains data integrity
- The stem is the leading digit or digits
- The leaf is the remaining digit
- Make sure to include units
- R Code

```
stem(midterm$MidtermExam)
```

*look at graph and get data values*

**Stem a**

The decimal point is 1 digit(s) to the right of the |

units

- R ou

```
stem   leaf
 2 | 8
 3 |
 4 | 4
 5 | 11566
 6 | 13334446679
 7 | 2247
 8 | 00147
 9 | 2
```

Stem and Leaf Plot of Midterm Exam Scores

Stem - 7
leaf - 4
units → 7 (10)
value 7 (10) + 4 = 74

**Histogram**

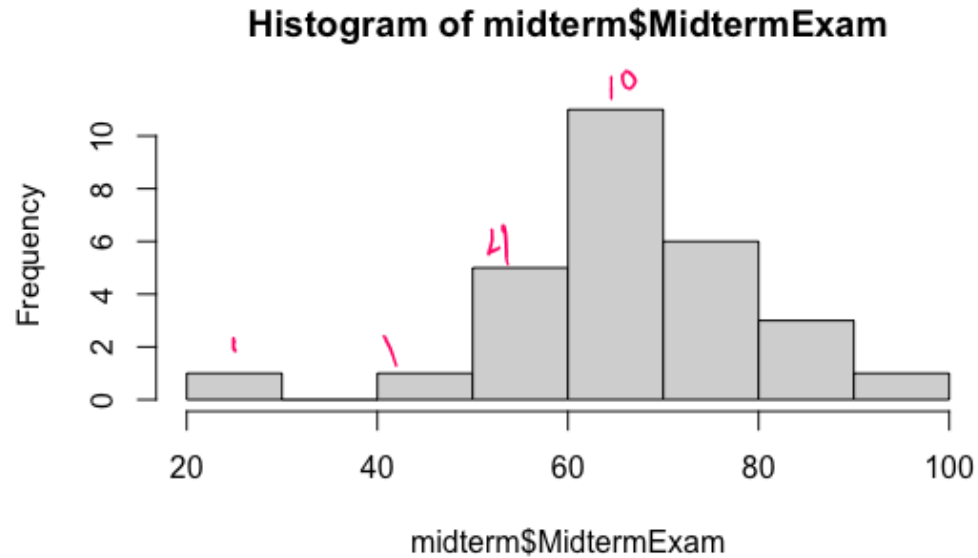- A histogram is a barchart displaying the frequency distribution information

- There are three types of histograms: frequency, relative frequency and cumulative relative frequency

- R code

```
hist(midterm$MidtermExam)
```

*Handwritten annotations:*

? bins

↗ counts

% of observations

↳ total %

→ needs lots of data to get shape

# Histogram E

- R output



Histogram of Midterm Exam Scores

estimators to be robust → not overly influenced by extreme values

**Boxplot**

Inter Quartile Range

$$IQR = Q_3 - Q_1$$

median

Whisker extends to smallest data point within 1.5 interquartile ranges from first quartile

$Q_1$    $Q_2$    $Q_3$

First quartile    Second quartile    Third quartile

Whisker extends to largest data point within 1.5 interquartile ranges from third quartile

Outliers

Figure 2-11
Description of a box plot.

Outliers    Extreme outlier

|← 1.5 IQR →|← 1.5 IQR →|← IQR →|← 1.5 IQR →|← 1.5 IQR →|

**Box Plot 1**

← famous Statistician Statistician

Box plot with explanation

Box & whiskers

Q1: which plant has highest Quality index? Plant 1 has highest QI

**Box Plot 2**

Q2: what's difference in QI between plants 2 & 3? by location, same median, so **no** difference
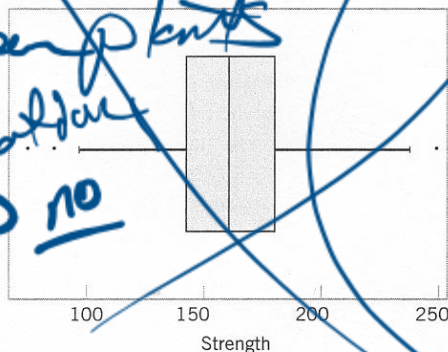


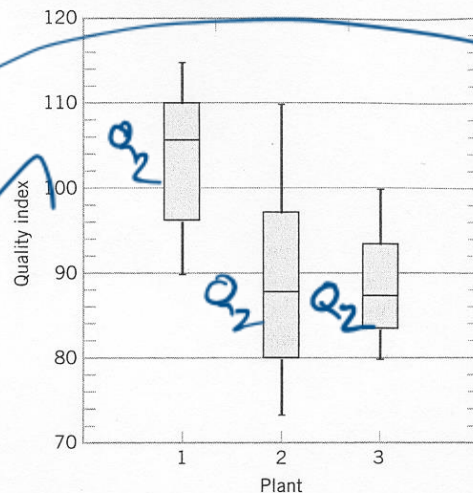Figure 2-12 Box plot for compressive strength data in Table 2-2.



Figure 2-13 Comparative box plots of a quality index at three plants.
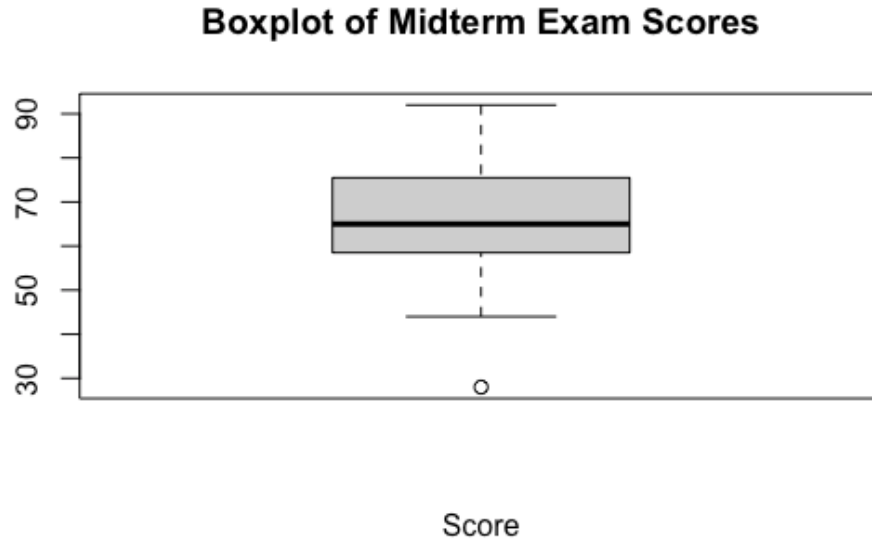
examples of boxplots

Q3: which plant has most Variability?
Q2 → rectangle is largest & whiskers are largest

# Box Plot 3

- R code for I

```
boxplot                              :'Score',m
ain='Box                             es')
```

- R Box Plot



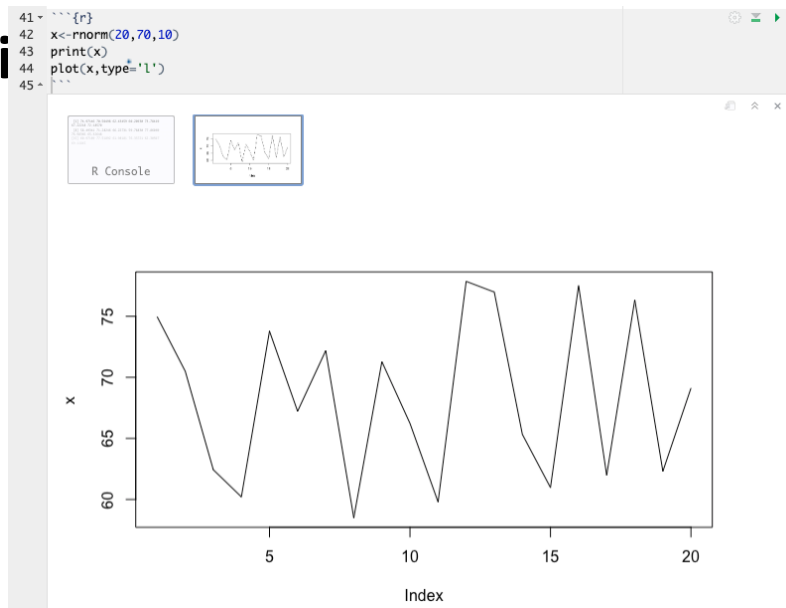Boxplot of Midterm Exam Scores

**Time Series Plot**

- A **time series plot** is a graph in which the vertical axis denotes the observed value of the variable (say $x$) and the horizontal axis denotes time
- Excellent tool for detecting:
  - trends,
  - cycles,
  - other non-random patterns

# Time Series Plot i

```{r}
x<-rnorm(20,70,10)
print(x)
plot(x,type='l')
```

R Console



Time Series Plot

**Probability Plotting**

- **Probability plotting** is a graphical method of determining whether sample data conform to a hypothesized distribution

- Used for validating assumptions

- Alternative to hypothesis testing

**Construction**

_ignore this ancient history_

1. Sort the data from smallest to largest, .

2. $x_{(1)}, x_{(2)}, \ldots, x_{(n)}$

3. Calculate the observed cumulative frequency $(j - 0.5)/n$

_replaced by computer graphics_

For the normal distribution find $z_j$ that satisfies

$$\frac{j - 0.5}{n} = P(Z \le z_j) = \Phi(z_j)$$

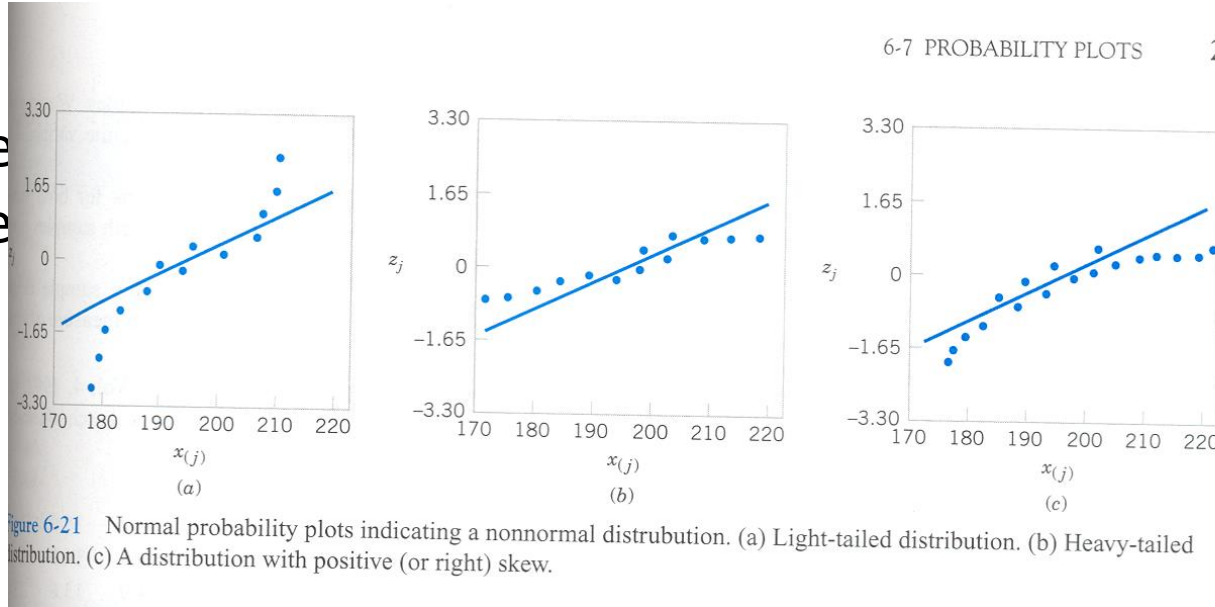3. Plot $z_j$ versus $x_{(j)}$ on special graph paper

weakness: Subjectivity

**Usage**

- If the ... tion is corre...



6-7 PROBABILITY PLOTS 215

Figure 6-21 Normal probability plots indicating a nonnormal distribution. (a) Light-tailed distribution. (b) Heavy-tailed distribution. (c) A distribution with positive (or right) skew.
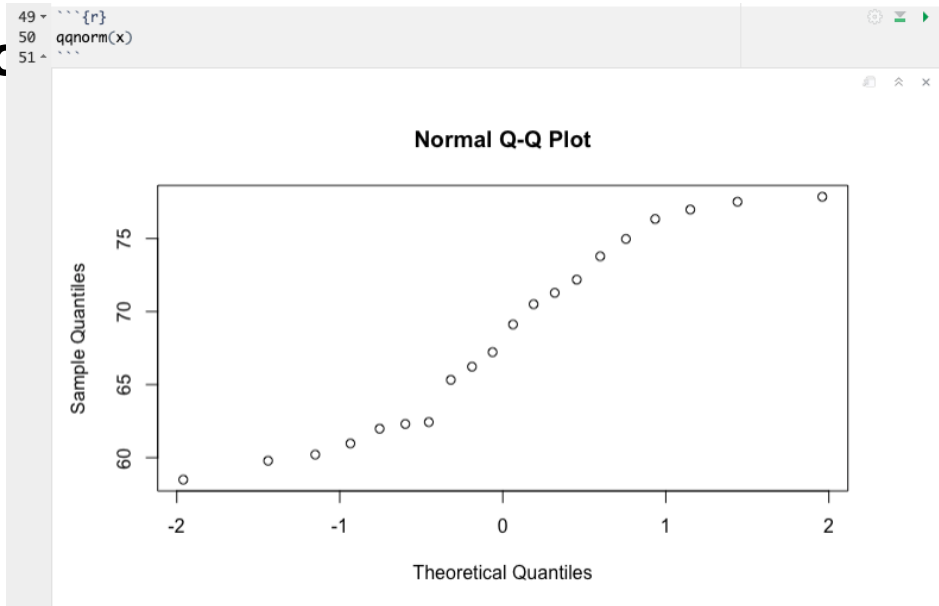
normal probability plots from textbook, figure 6.21 on page 215
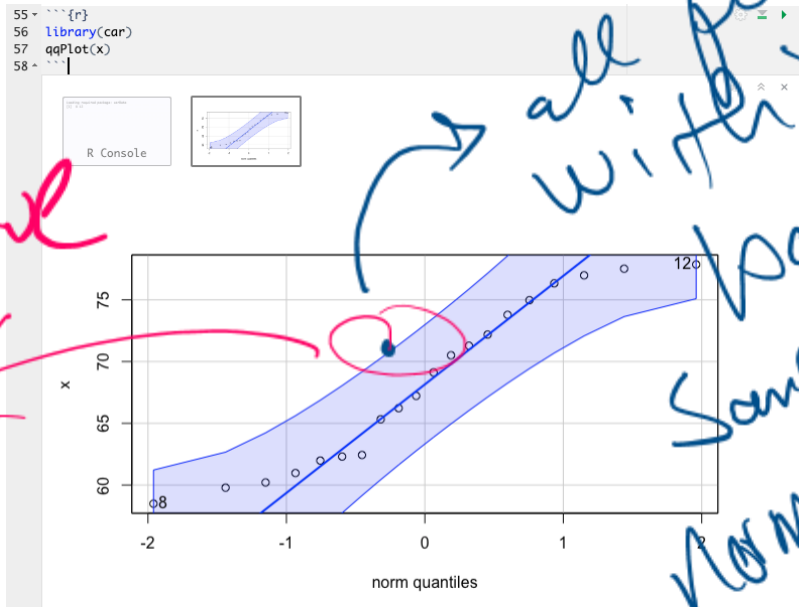
# Probability Plot



Normal Probability Plot

**Probability Plot Example 2**

- Difficulty from example one is how close to straight is "good enough"

- Add confidence bands to normal probability plot
  - Requires package car to be added to R
  - If all points are within the band, we are 95% confident that the sample is from a normal distribution. However if one or more points are not within band, the data is not from a normal distribution

```r
{r}
library(car)
qqPlot(x)
```
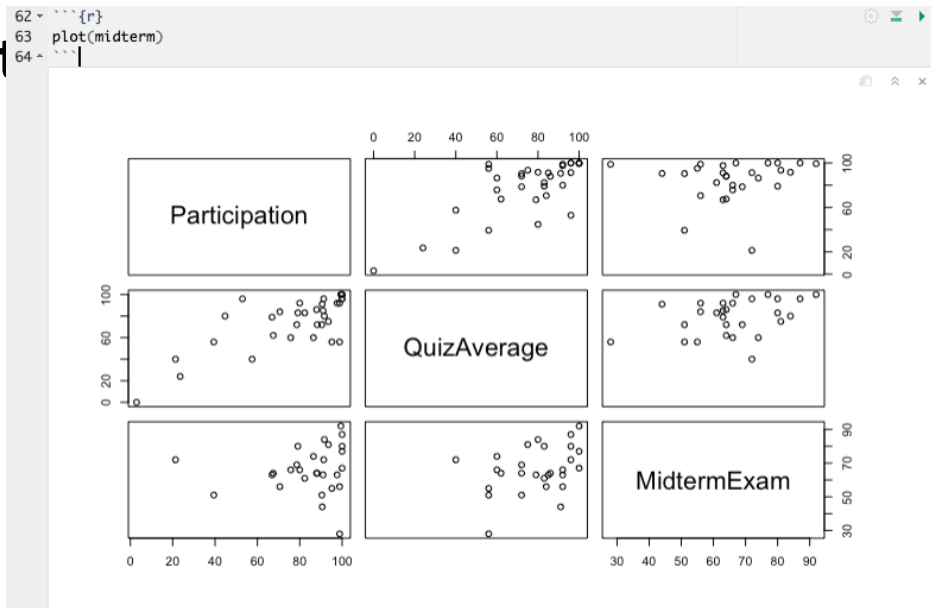
R Console

this would be Subjective

→ all points within the bands so sample is from Normal



QQ Plot with band

# Multivariate Data

not cover need Chapter 5

**Matrix of Scatterplots**

```r
62    ```{r}
63    plot(midterm)
64    ```
```



Scatter Plots

# Covariance in R

```{r}
midterm_NA <- na.omit(midterm)
print(cov(midterm_NA))
```

```
                Participation QuizAverage MidtermExam
Participation      340.16778     193.7847    28.75699
QuizAverage        193.78474     269.0899    81.17460
MidtermExam         28.75699      81.1746   188.43915
```

Covariance Matrix

# Correlation

```{r}
print(cor(midterm_NA))
```

```
              Participation QuizAverage MidtermExam
Participation     1.0000000   0.6405076   0.1135825
QuizAverage       0.6405076   1.0000000   0.3604839
MidtermExam       0.1135825   0.3604839   1.0000000
```

Correlation Matrix