

# AMS 572 Data Analysis I

## Summarizing and Exploring Data

Pei-Fen Kuan

Applied Math and Stats, Stony Brook University

# Content

- ▶ Types of variables
- ▶ Measures of location
- ▶ Measures of spread, shape
- ▶ Data displays

# Types of variables

- ▶ A *variable* is a quantity that may vary from object to object
- ▶ A *sample* or *data set* is a collection of values of one or more variables.
- ▶ Quantitative variable intrinsically numerical  
e.g. *age, height, weight*
- ▶ Qualitative (categorical) - intrinsically nonnumerical  
e.g. *race, gender*

# Types of variables

- ▶ Qualitative (categorical) - intrinsically nonnumerical
  - ▶ Binary, dichotomous  
e.g., *female/male*
  - ▶ Ordinal - natural ordering  
e.g., *attitude (strongly agree, agree, neutral, disagree, strongly disagree)*
  - ▶ Nominal - no natural ordering  
e.g., *race*
- ▶ In recording qualitative data, numerical values may be assigned

# Measures of Location

- ▶ (Arithmetic) Mean
- ▶ Percentiles
- ▶ Median
- ▶ Mode
- ▶ Geometric mean

# Arithmetic mean

- Data:

$$x_1, x_2, \dots, x_n$$

- Mean:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Example: Duration of hospital stay in days:

$$x_1 = 5, x_2 = 10, x_3 = 6, x_4 = 11$$

Mean:

$$\bar{x} = \frac{1}{4}(5+10+6+11) = \frac{32}{4} = 8$$

# Properties of Mean

- ▶ Let  $c$  be any constant
- ▶ If

$$y_i = x_i + c \text{ for } i = 1, 2, 3, \dots, n,$$

then

$$\bar{y} = \bar{x} + c$$

- ▶ If

$$y_i = cx_i \text{ for } i = 1, 2, 3, \dots, n,$$

then

$$\bar{y} = c\bar{x}$$

# Properties of Mean - Example

- ▶ A sample of birth weights in a hospital found

$$\bar{y} = 3166.9 \text{ grams}$$

- ▶ 1 oz = 28.35 g
- ▶ Therefore the mean in ozs. is

$$\bar{x} = \frac{3166.9}{28.35} = 111.7 \text{ oz}$$



# Order Statistics

dependence

- ▶ Data:  $x_1, x_2, \dots, x_n$
- ▶ Order data from smallest to largest

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

- ▶  $x_{(1)}, x_{(2)}, \dots, x_{(n)}$  are *order statistics*
- ▶ Note

$$x_{(1)} = \min\{x_1, x_2, \dots, x_n\}$$

$$x_{(n)} = \max\{x_1, x_2, \dots, x_n\}$$

- ▶ Example: Duration of hospital stay in days:

$$x_1 = 5, x_2 = 10, x_3 = 6, x_4 = 11$$

Order statistics:

$$x_{(1)} = 5, x_{(2)} = 6, x_{(3)} = 10, x_{(4)} = 11$$

# Percentiles

- $y_1, y_2, \dots, y_n$
- ▶ Intuitive definition: the  $x$  percentile is such that  $x\%$  of the observations are less than that value
  - ▶ Also known as sample *quantile*
  - ▶ The  $(p \times 100)^{th}$  percentile of a sample

$$\hat{\zeta}_p = \begin{cases} y_{(np+p)} & \text{if } np+p \text{ is an integer} \\ \{y_{(\lfloor np+p \rfloor)} + y_{(\lceil np+p \rceil)}\}/2 & \text{otherwise} \end{cases}$$

for  $0 < p < 1$

*floor*  $\lfloor \rfloor$       *ceiling*  $\lceil \rceil$

## Percentiles: General form

- ▶ General form (Hyndman and Fan, *Am Stat* 1996)

$$\hat{\zeta}_p = (1 - \gamma)y_{(j)} + \gamma y_{(j+1)}$$

where  $j = \lfloor pn + k \rfloor$  for some  $k \in \mathbb{R}$  and  $0 \leq \gamma \leq 1$ .

- ▶ Your textbook

$$\hat{\zeta}_p = \begin{cases} y_{(np+p)} & \text{if } np + p \text{ is an integer} \\ y_{(m)} + [np + p - m](y_{(m+1)} - y_{(m)}) & \text{otherwise} \end{cases}$$

where  $m = \lfloor np + p \rfloor$

## Example

- Suppose  $n = 278$  and we want the 75th percentile

$$np = 278 \times .75 = 208.5 (\text{or } np + p = 209.25)$$

such that

$$\hat{\zeta}_{.75} = x_{(209)}$$

- R

```
> x <- 1:278
> quantile(x,.75,type=2)
75%
209
```

1st quantile : 25%  
3rd quartile : 75%

# Median

- ▶ The sample median is the 50th percentile

$$\hat{\zeta}_{.5} = \begin{cases} y_{(\frac{n+1}{2})} & \text{if } n \text{ is odd} \\ \{y_{(n/2)} + y_{(n/2+1)}\}/2 & \text{if } n \text{ is even} \end{cases}$$

for  $0 < p < 1$

- ▶ Example: Duration of hospital stay in days:

$$x_1 = 5, x_2 = 10, x_3 = 6, x_4 = 11$$

Median:

$$\hat{\zeta}_{.5} = \frac{x_{(2)} + x_{(3)}}{2} = 8$$

# Mode

- ▶ The mode is the most frequently occurring value in the data set
- ▶ E.g., if

$$x_1 = 5, x_2 = 11, x_3 = 6, x_4 = 11$$

then mode is 11

# Geometric Mean

- ▶ Data:  $x_1, x_2, \dots, x_n$
- ▶ The geometric mean of  $x$  is

$$\bar{x}_g = (x_1 x_2 \cdots x_n)^{1/n}$$

- ▶ Let  $y_i = \log(x_i)$  for  $i = 1, 2, \dots, n$ . Then

$$\bar{y} = \frac{\log(x_1) + \cdots + \log(x_n)}{n} = \log(\pi x_i)^{\frac{1}{n}} \quad \bar{x}_g = \exp(\bar{y})$$

- ▶ Eg, suppose  $x_1 = 10$  and  $x_2 = 0.1$ . Then

$$\bar{x}_g = (10 \times 0.1)^{\frac{1}{2}}$$

*trimmed mean*

## *Arithmetic*

- ▶ Mean is most often used measure
- ▶ Median is better if there are influential observations (more robust to extreme values)
- ▶ Mode rarely used (exception: nominal data)



## Example

- Duration of hospital stay in days:

$$x_1 = 5, x_2 = 10, x_3 = 6, x_4 = 11$$

$$m = \hat{\xi}_{0.5} = \bar{x} = 8, \quad \bar{x}_g = 7.6$$

Alter the last observation:

$$x_1 = 5, x_2 = 10, x_3 = 6, x_4 = 50$$

$$m = \hat{\xi}_{0.5} = 8, \quad \bar{x} = 17.7, \quad \bar{x}_g = 11.1$$

# Measures of Spread, Shape

- ▶ Range
- ▶ Variance and standard deviation
- ▶ Interquartile range
- ▶ Skewness, Kurtosis

# Range

- ▶ Range:

$$r_a = x_{(n)} - x_{(1)}$$

- ▶ Easy to calculate
- ▶ Sensitive to unusual observations (outliers)
- ▶ Usually, the larger  $n$  is, the larger  $r_a$

# Sample Variance and Standard Deviation

- ▶ Want to measure deviation from mean
- ▶ Sample variance

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

- ▶ Sample standard deviation

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

# Sample Variance and Standard Deviation

- ▶ An alternative form of the sample variance is

$$s_1^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

MLE

- ▶ We have shown that  $s^2$  is unbiased for population variance  $\sigma^2$ , however

$$E(s_1^2) = \sigma^2 - \frac{\sigma^2}{n}$$

# Sample Standard Deviation

$s_y$ : sample sd of " $y_i$ "  
 $s_x$ : sample sd of " $x_i$ "

- ▶ The units of  $s$  are the same as the units of  $x_i$
- ▶ If  $s$  is large, the data are spread over a wide range
- ▶ If  $c$  is a constant and

$$y_i = x_i + c,$$

then  $s_y = s_x$

- ▶ If

$$y_i = cx_i$$

then

$$s_y = c s_x$$

## Some approximations

qnorm(0.75)

- ▶ The interval  $\bar{x} \pm s$  will contain approx 68% of the observations
- ▶ The interval  $\bar{x} \pm 2s$  will contain approx 95% of the observations
- ▶ Approx  $s$  by

$$s \approx \frac{\hat{\zeta}_{.75} - \hat{\zeta}_{.25}}{1.35}$$

- ▶ Note

$\hat{\zeta}_{.75} - \hat{\zeta}_{.25} \rightarrow$  25th percentile  
is called *interquartile range*  $\downarrow$   
75th percentile

# Symmetry and Skewness

- ▶ Informally, define *symmetry* to indicate having a uniform or even distribution about the mean
- ▶ If a distribution is symmetric,

$$\text{mean} = \text{median}$$

- ▶ Data sets that are not symmetric are said to be *skewed*
- ▶ *Skewness* is a measurement of the degree to which a data set is skewed



## Skewness

In R, library(e1071)  
skewness( )

- Define  $r$ th sample moment about the mean

$$m_r = \frac{\sum_i (y_i - \bar{y})^r}{n} \text{ for } r = 1, 2, 3, \dots$$

- Text definition of sample skewness:

$$a_3 = \frac{\sum_i (y_i - \bar{y})^3 / n}{\{\sum_i (y_i - \bar{y})^2 / (n - 1)\}^{3/2}}$$

- $a_3 > 0$  indicates skewness to the right



mean > mode -

# Kurtosis

- ▶ *Kurtosis* is a measure of the flatness or peakedness of a distribution; degree of archedness; thickness of tails
- ▶ Text definition of *sample* kurtosis:

$$a_4 = \frac{\sum_i (y_i - \bar{y})^4 / n}{\{\sum_i (y_i - \bar{y})^2 / (n - 1)\}^2}$$

- ▶  $a_4 > 3$  indicates the distribution has heavier tails than the normal distribution.

eg. T distribution, Cauchy

# Data display

- ▶ Simplest form is a line listing
- ▶ A *frequency table* gives the frequency of observations within a set of ordered intervals
- ▶ Intervals should be mutually exclusive and exhaustive
- ▶ 8 to 10 intervals is usually sufficient
- ▶ With the exception of the end intervals, the length of the intervals should be constant

## Frequency Table - Example

Blood Pressure	Pop1	Pop2	Pop3
< 106	218	4	23
106-114 <del>5</del>	272	23	132
116-124 <del>5</del>	337	49	290
126-134 <del>5</del>	362	33	347
136-144 <del>5</del>	302	41	346
146-154 <del>5</del>	261	38	202
156-164 <del>5</del>	166	23	109
> 164 <del>5</del>	314	52	112
Total	2232	263	1561

# Frequency Tables

- ▶ Table on previous slide example of *empirical frequency distribution*
- ▶ Difficult to compare blood pressure distributions due to different sample sizes
- ▶ Divide by sample size to get *empirical relative frequency distribution*

## ERFD - Example

Blood Pressure	Pop1	Pop2	Pop3
< 106	0.098	0.015	0.015
106-114	0.122	0.087	0.085
116-124	0.151	0.186	0.186
126-134	0.162	0.125	0.222
136-144	0.135	0.156	0.222
146-154	0.117	0.144	0.129
156-164	0.074	0.087	0.070
> 164	0.141	0.198	0.072
Total	2232	263	1561

# Graphs

- ▶ Histogram
- ▶ Stem and leaf plot
- ▶ Box plot
- ▶ Trellis/conditional plots

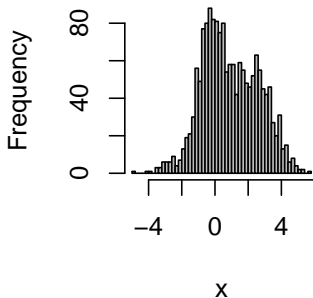
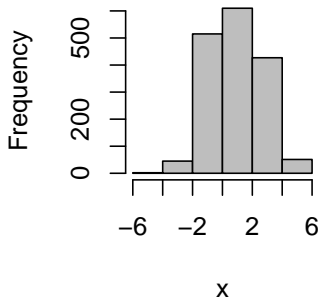
# Histogram

- ▶ Data are divided into intervals as in a frequency table
- ▶ A histogram is a bar graph with the area of each bar equal to the relative frequency in the interval.
- ▶ Can compare histograms from samples of different size
- ▶ Intervals need not be the same width
- ▶ Beware effect of choice of interval width



## *random norm*

```
> x <- c(rnorm(50,-2.5),rnorm(1000),rnorm(600,2.7))  
> hist(x,breaks=5,col="gray",xlab="x",main="")  
> hist(x,breaks=50,col="gray",xlab="x",main="")
```



# Stem and Leaf Plot: Example

```
> stem(mileage)
```

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

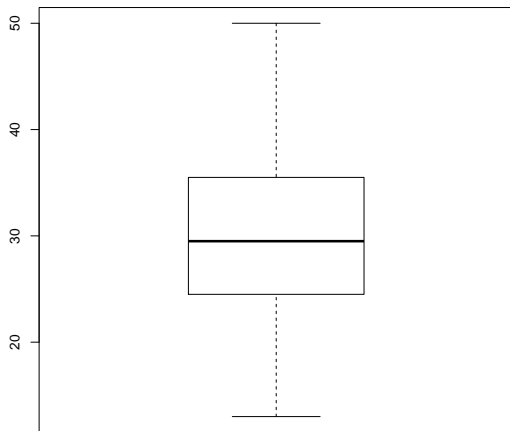
```
1 | 37
2 | 0134557889
3 | 011125678
4 | 00
5 | 0
```

# Box plot

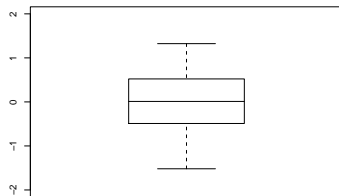
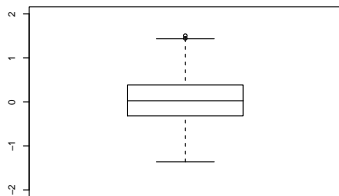
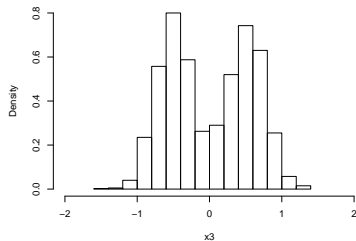
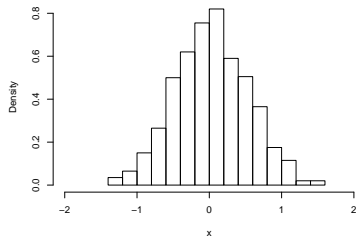
- ▶ The top of the box is the 75th percentile ( $\hat{\zeta}_{.75}$ ); the bottom is the 25th percentile ( $\hat{\zeta}_{.25}$ )
- ▶ A line through the box is drawn at the median
- ▶ The lines extending out of the box (*whiskers*) may extend to
  - ▶ the 90th and 10th percentiles
  - ▶ the largest and smallest values
  - ▶ largest observation  $\leq \hat{\zeta}_{.75} + 1.5 \times \text{IQR}$ ;  
smallest observation  $\geq \hat{\zeta}_{.25} - 1.5 \times \text{IQR}$
- ▶ Data beyond whiskers may be plotted individually

# Boxplot Example

```
> boxplot(mileage)
```

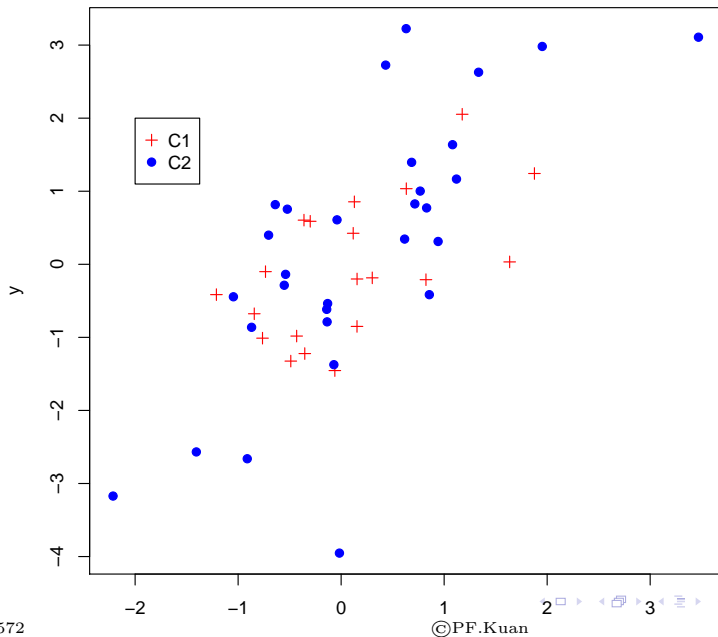


# Box plot and Histogram Example

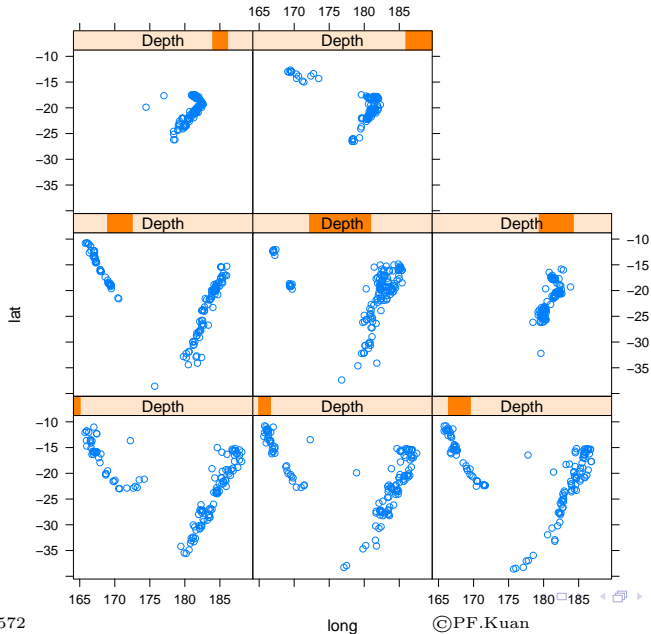


# Multivariate plots

- ▶ Describe relationships/associations between more than one variable
- ▶ Scatterplots
  - ▶ Simple for two variables
  - ▶ Add color, symbols for  $> 2$  variables



# Trellis plots





# Tables or graphs?

- ▶ Tables best suited for looking up specific information
- ▶ Graphs better for perceiving trends, making comparisons and predictions