



SKKU Biostats and Big data

CHOONG-WAN WOO | COCOAN lab | <http://cocoanlab.github.io>



Lecture 03

Data visualization

Quiz!

<https://forms.gle/rUbz1gbXAowQLsfy8>

A statistician's manifesto

(From T. Hastie, via J. McAuliffe, [via Jordan Boyd-Graber](#))

- Understand the ideas behind the statistical methods, so you know how to use them, when to use them, when not to use them.
- Complicated methods build on simple methods. Understand simple methods first.
- The results of a method are of little use without an assessment of how well or poorly it is doing.

Let's go over the survey results

'Biostat and Big Data' course survey (Spring 2022)

This survey isn't meant to test your knowledge. Rather, it'll help us adjust the class just for you, and make it fun! (이 설문조사는 여러분을 평가하기 위한 것이 아닙니다. 수업이 여러분에게 더 도움이 되고, 즐겁게 만들기 위한 것이니 솔직하고 편하게 작성해주세요!)

 choongwan.woo@gmail.com (not shared) [Switch account](#) 

* Required

Name *

Your answer

Student ID (e.g., 2019123456) *

Your answer

Do you have any experience with data analysis? (데이터 분석을 해본 경험이 있나요?) *

Yes

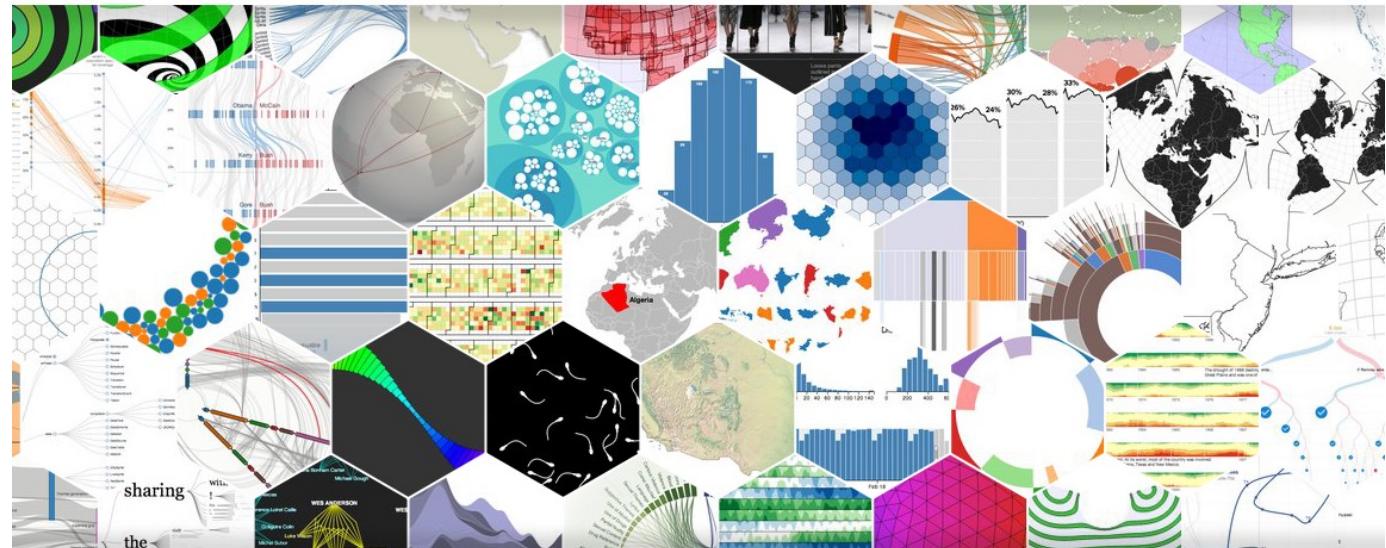
No



Three Rules of Data Analysis

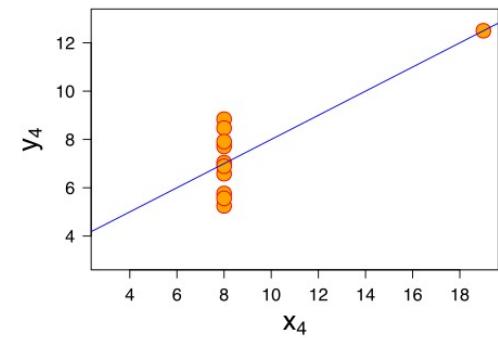
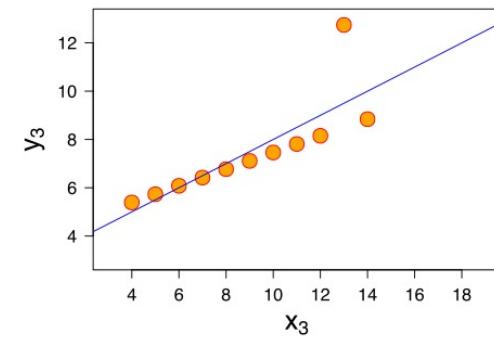
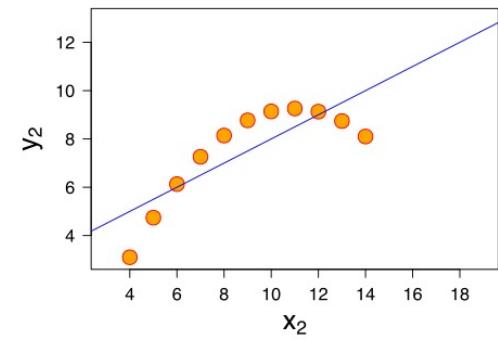
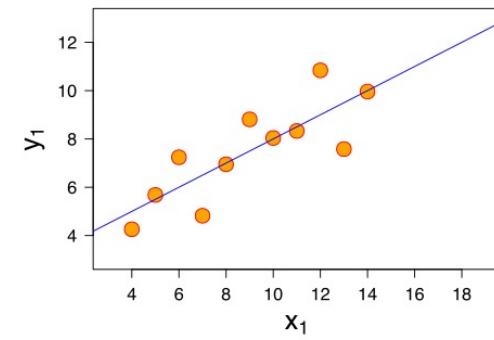
- Make a picture
- Make a picture
- Make a picture

- To think, show, and tell...
- And to make it cool!
- E.g., <https://d3js.org>



Why visualization matters

- We're lazy and don't like to read
- Some information isn't easy to describe verbally
- Statistical summaries can be misleading
- Aesthetically-pleasing visuals are engaging
- Anscombe's quartet:
- https://en.wikipedia.org/wiki/Anscombe%27s_quartet



Slide credit: Tal Yarkoni

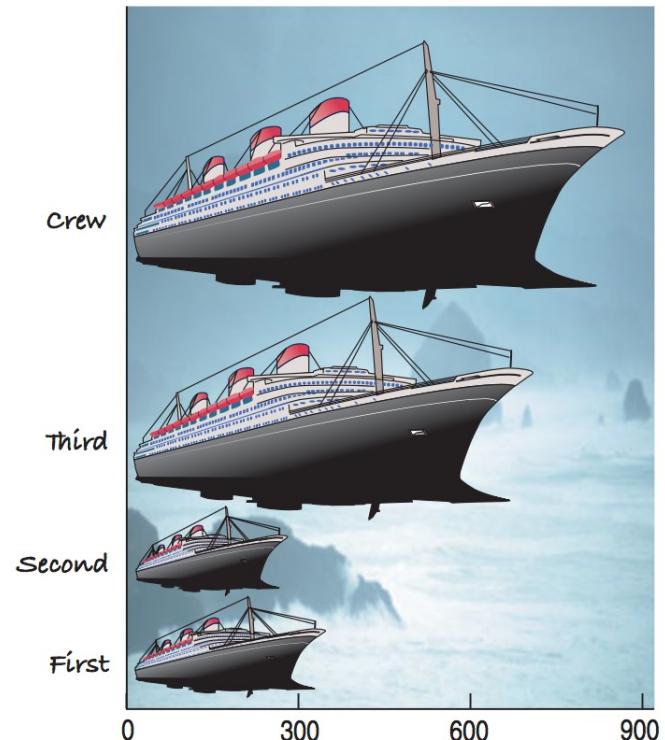
CHOONG-WAN WOO | COCOAN lab | <http://cocoanlab.github.io>

Textbook's example data

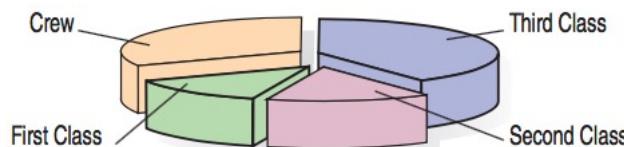
- Frequency/contingency table of “ticket class” and “Survival” for the *Titanic* passengers

Class	Count
First	325
Second	285
Third	706
Crew	885

Survival	Class					Total
	First	Second	Third	Crew		
Alive	203	118	178	212	711	
Dead	122	167	528	673	1490	
Total	325	285	706	885	2201	



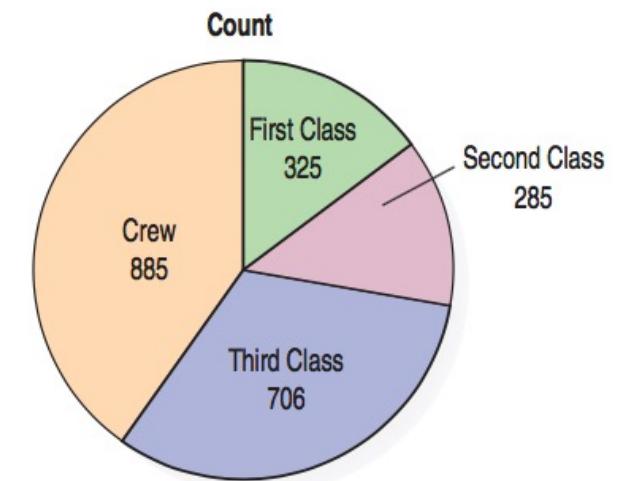
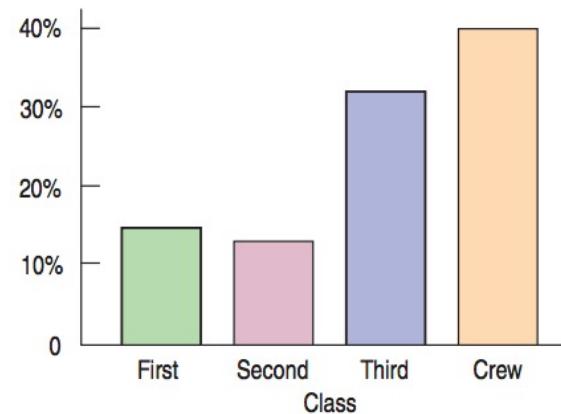
- What's wrong with this plot?
- “Area principle”: the area occupied by a part of the graph should correspond to the magnitude of the value it represents.



Bar Charts, Pie Charts



Relative frequency bar chart



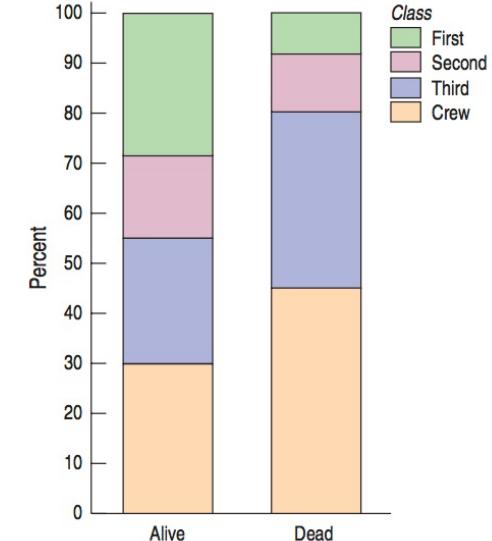
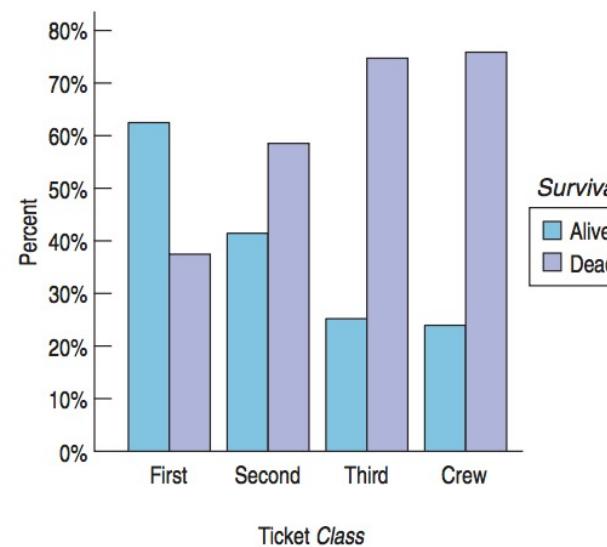
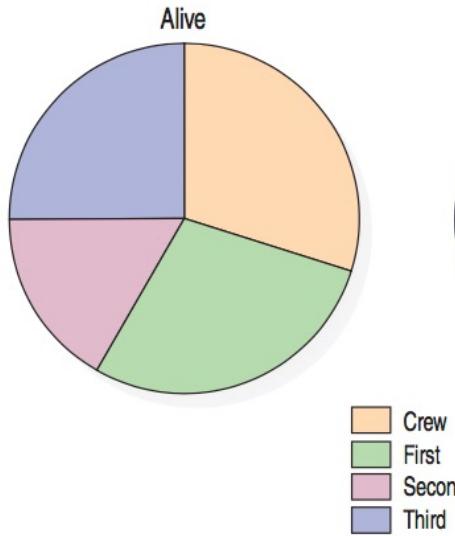
Marginal distribution, conditional distribution

Contingency table

		Class				Marginal distribution of "Survival"	
		First	Second	Third	Crew		
Survival	Alive	203	118	178	212	711	Conditional distribution (satisfy a condition on another variable)
	Dead	122	167	528	673	1490	
	Total	325	285	706	885	2201	

Marginal distribution of "Class"

Charts for conditional distribution



- The variables are **independent**: “when the distribution of one variable tells us nothing about the distribution of the other variable”

Simpson's Paradox

- Two pilots, Moe and Jill
- Who is the better pilot?

Proportion of on-time flights

Pilot	Time of Day		
	Day	Night	Overall
Moe			100 out of 120 83%
Jill			94 out of 120 78%

Simpson's Paradox

- Two pilots, Moe and Jill
- Who is the better pilot?
- Now who is the better pilot?

Proportion of on-time flights

		Time of Day		
		Day	Night	Overall
Pilot	Moe	90 out of 100 90%	10 out of 20 50%	100 out of 120 83%
	Jill	19 out of 20 95%	75 out of 100 75%	94 out of 120 78%

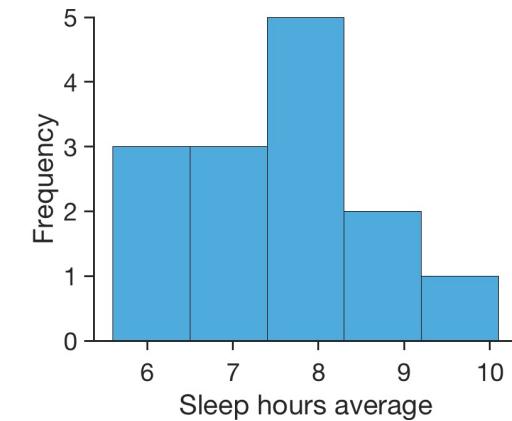
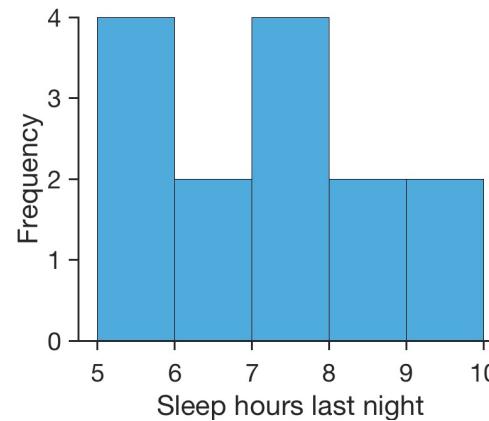
- Don't average over things *unfairly*
- *Different* numbers in *different* categories for Day and Night
- *Unfair* to average across categories

Quiz!

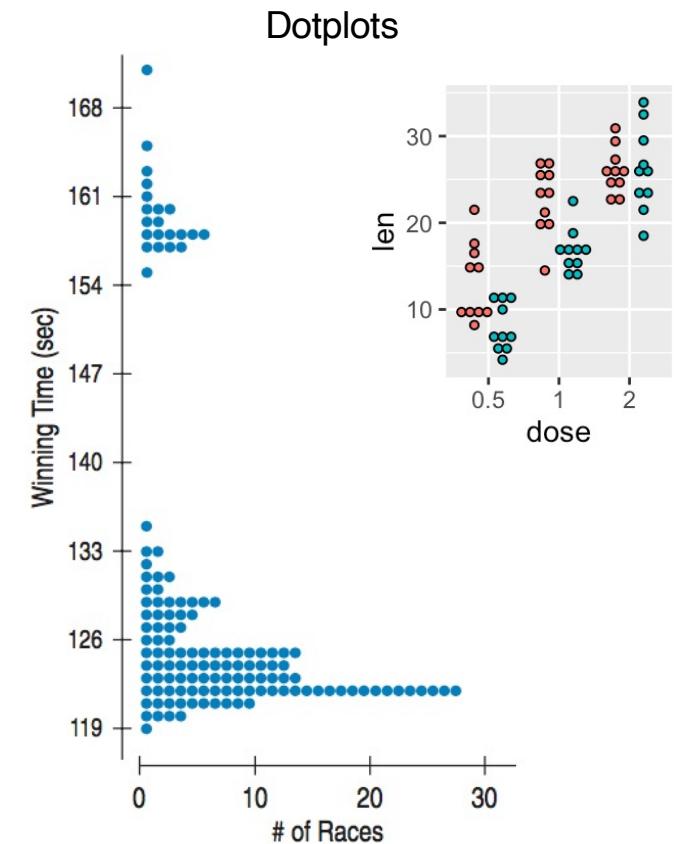
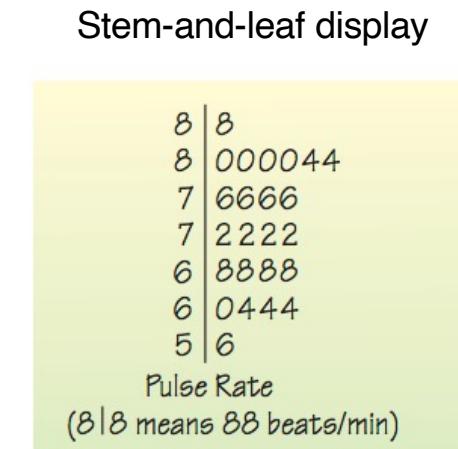
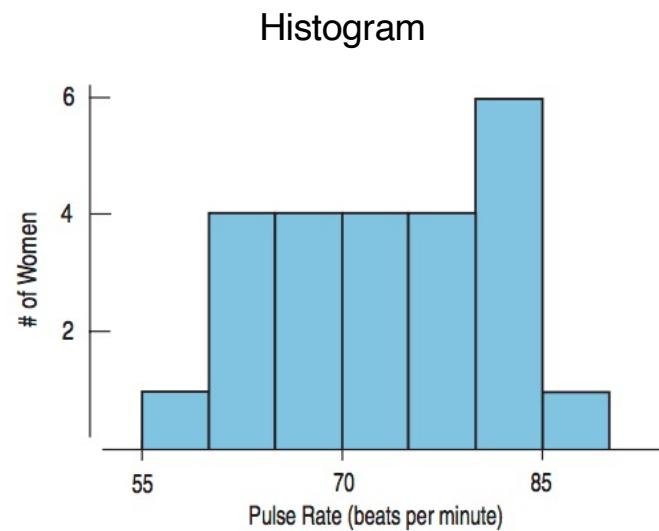
<https://forms.gle/rAdXTtcgRoA7KjGJ6>

Histograms for quantitative data

	sleep_hours_last_night	sleep_hours_average
Indiv 1	7	6
Indiv 2	5.5	6.5
Indiv 3	9	7
Indiv 4	5	7
Indiv 5	5	6
Indiv 6	7	8
Indiv 7	8	7
Indiv 8	7	8
Indiv 9	8	10
Indiv 10	10	9
Indiv 11	6	7.5
Indiv 12	7	9
Indiv 13	6	8
Indiv 14	5	8



Stem/leaf, dot plot



Shape

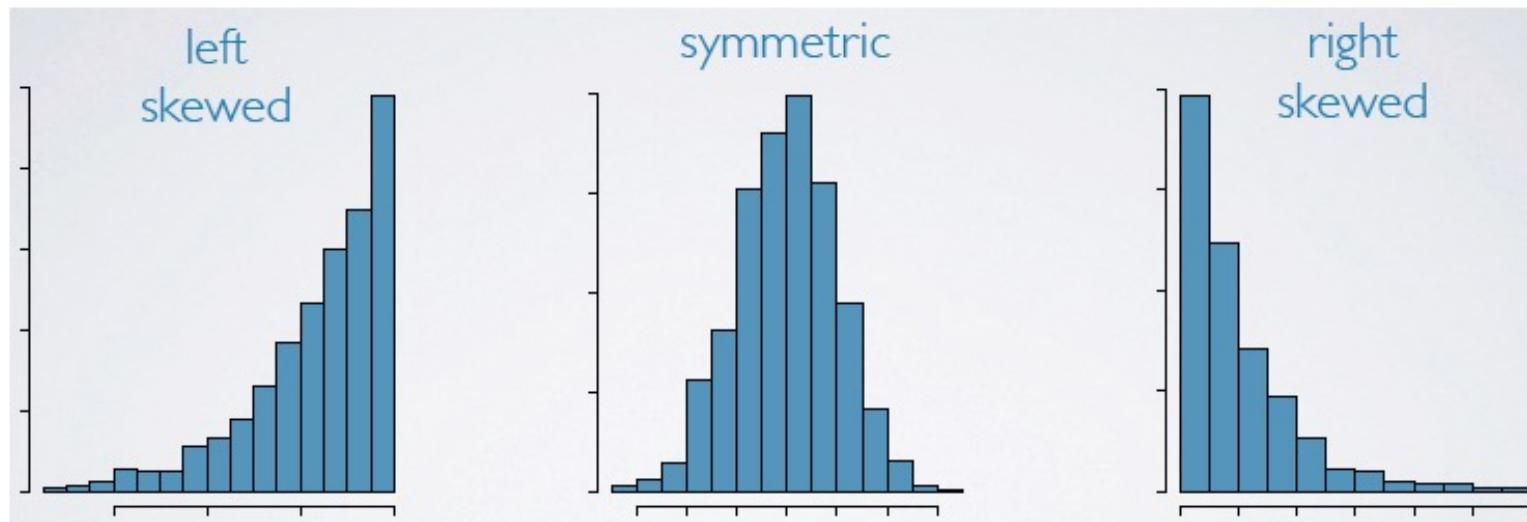
- Modes?



<http://researchhubs.com/post/ai/data-analysis-and-statistical-inference/visualizing-numerical-data.html>

Shape

- Modes?
- Symmetric?



<http://researchhubs.com/post/ai/data-analysis-and-statistical-inference/visualizing-numerical-data.html>

Center

- Median

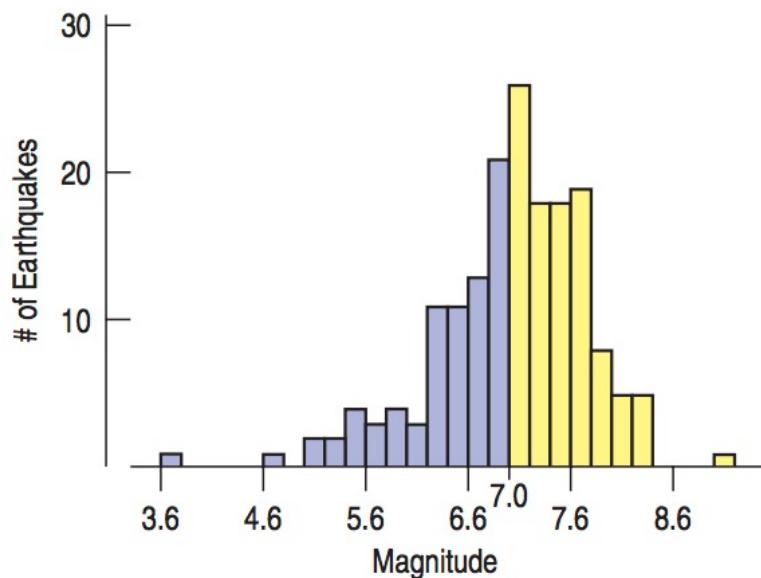


FIGURE 4.10

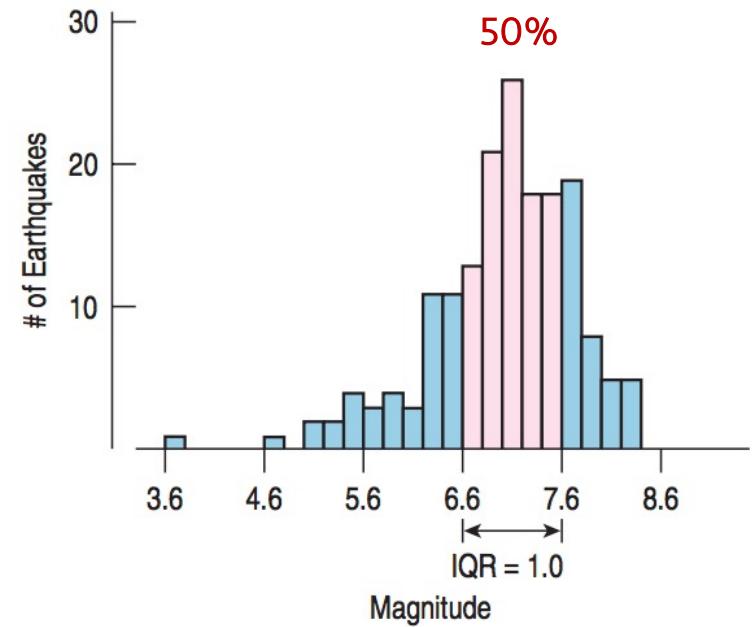
Tsunami-causing earthquakes (1981–2005).

The median splits the histogram into two halves of equal area.

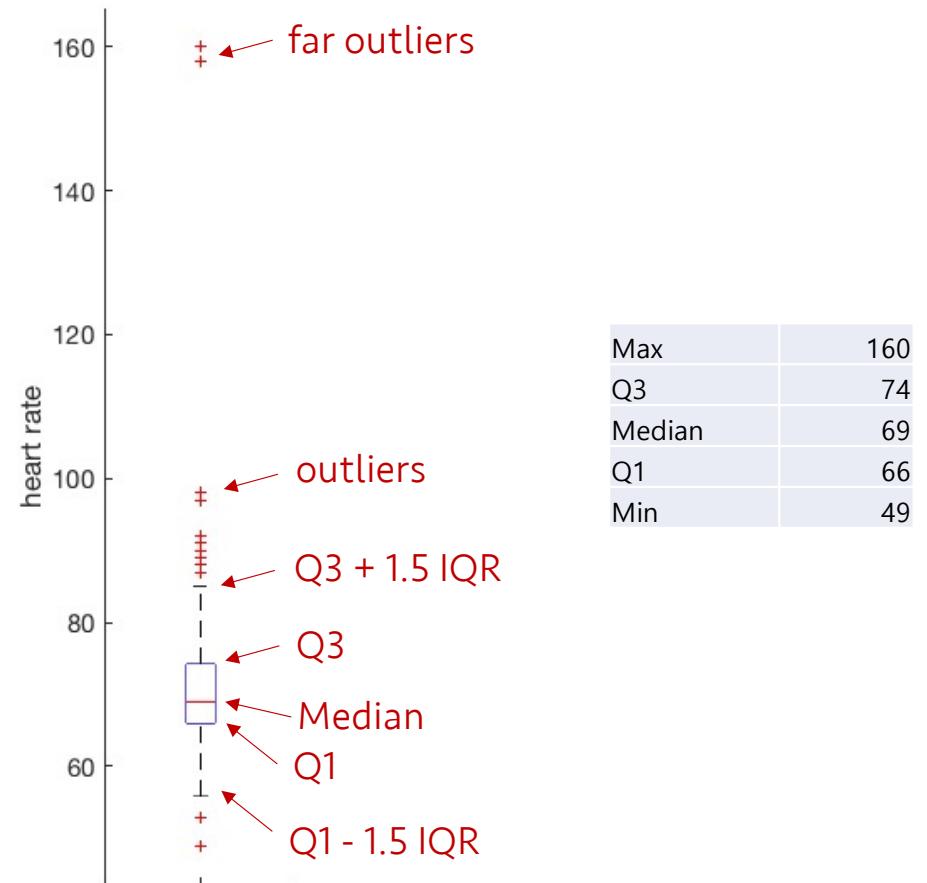
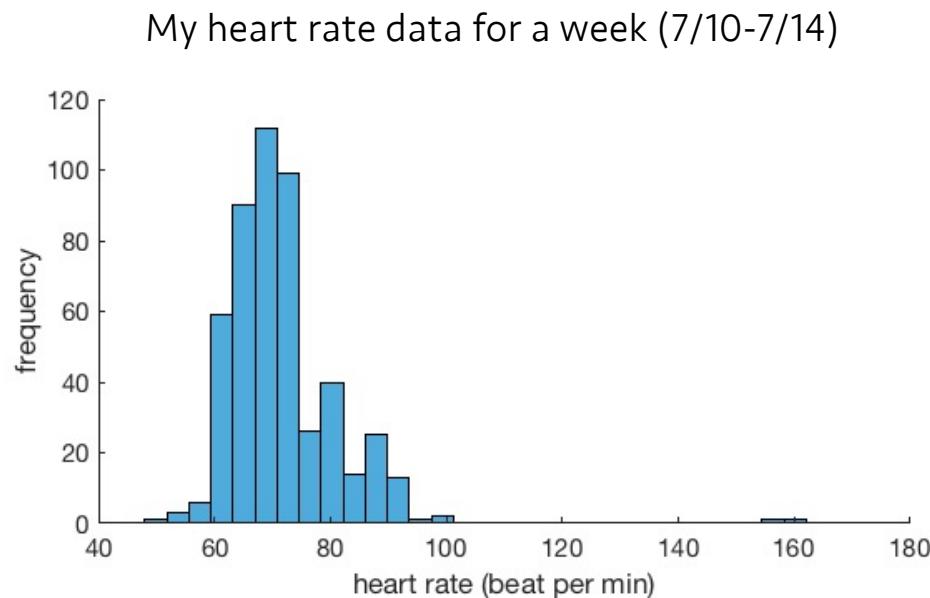
- 176 earthquakes
- Median: $(176+1)/2 = 88.5^{\text{th}}$ value in the sorted data
- ".5" = average of the two values (88^{th} and 89^{th})
- If there was 221 earthquakes
- Median: $(221+1)/2 = 111^{\text{th}}$ value in the sorted data

Spread

- Range
 - Range = max - min
- Interquartile range
 - Interquartile range (IQR) = upper quartile - lower quartile

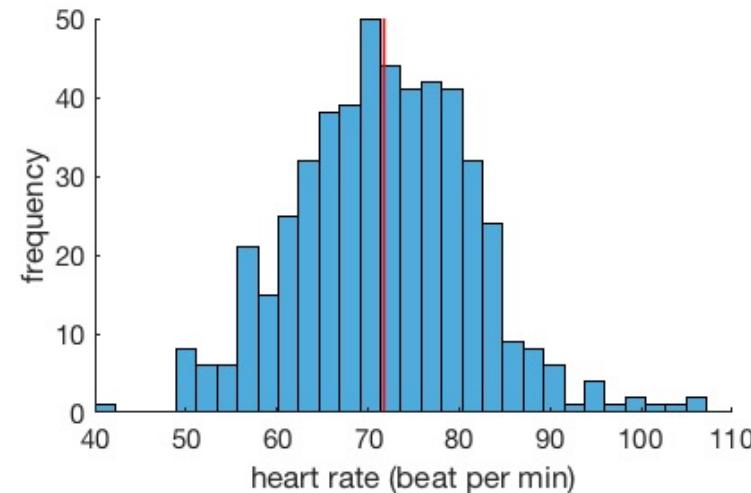
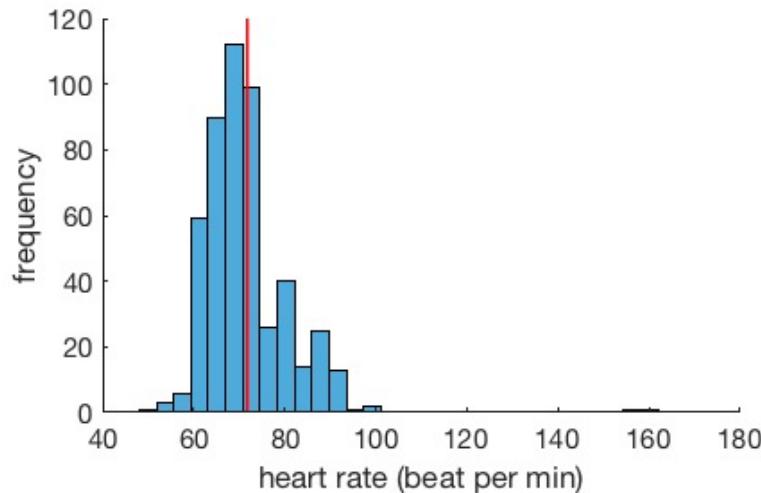


Boxplots and 5-Number Summaries



Center of Symmetric Distribution: Mean

$$\bar{y} = \frac{\text{Total}}{n} = \frac{\sum y}{n}$$



- If the histogram is symmetric and there are no outliers, the mean will be preferable.
- However, if the histogram is skewed or has outliers, the median might be better.

Quiz!

<https://forms.gle/jxEvNUTLTpZpfPnv9>

c.f.

$$\text{Percentile} = 100 \times \frac{i-0.5}{N}$$

i = rank after sorting the values in an ascending order

N = the number of values

Quiz!

With the following data, which of the following is WRONG summary statistics? *

Data: 2, 26, 29, 21, 23, 23, 12, 20, 6, 22

- range = 27
- Q3=12
- median = 21.5
- IQR=11
- mean = 18.4

c.f.

$$\text{Percentile} = 100 \times \frac{i-0.5}{N}$$

i = rank after sorting the values in an ascending order

N = the number of values

Spread of Symmetric Distribution: Standard Deviation

Variance

$$s^2 = \frac{\sum (y - \bar{y})^2}{n - 1}$$

Standard Deviation

$$s = \sqrt{\frac{\sum (y - \bar{y})^2}{n - 1}}$$

	Original Values	Deviations	Squared Deviations
Mean = 17	14	$14 - 17 = -3$	$(-3)^2 = 9$
	13	$13 - 17 = -4$	$(-4)^2 = 16$
	20	$20 - 17 = 3$	9
	22	$22 - 17 = 5$	25
	18	$18 - 17 = 1$	1
	19	$19 - 17 = 2$	4
	13	$13 - 17 = -4$	16

Add up the squared deviations: $9 + 16 + 9 + 25 + 1 + 4 + 16 = 80$.

Now divide by $n - 1$:

Finally, take the square root:

$$\frac{80}{6} = 13.33.$$

$$s = \sqrt{13.33} = 3.65$$

Quiz!

<https://forms.gle/7UnA2spepxwb4brw6>

Textbook for the class?

바이오통계와 빅데이터

(부제: 수학없는 통계)

우충완 (성균관대학교 글로벌바이오메디컬공학과)

목차

- 1... 서론
- 4... 1장. 통계학이란?
- 9... 2장. 데이터
- 14... 3장. 데이터 시각화
- 27... 4장. 분포의 비교와 정상 모델
- 36... 5장. 스캐터 도표(scatterplot)와 상관값(correlation)
- 49... 6장. 선형 회귀 (linear regression)
- 65... 7장. 샘플링(sampling)
- 71... 8장. 실험 디자인(Designing experiments)
- 77... 9장. 확률과 확률 모델
- 88... 10장. 무선 변수(random variables)
- 95... 11장. 확률 모델(probability models)
- 107... 12장. 샘플링 분포와 중심극한정리
- 123... 13장. 비율의 신뢰 구간
- 131... 14장. 비율에 대한 가설 검정 및 p값 논쟁
- 142... 15장. 평균에 대한 가설 검정, t-test
- 159... 16장. 가설검정과 신뢰구간에 대한 추가적인 고려사항들
- 173... 17장. 두 집단을 비교하기
- 185... 18장. 대응 표본 t 검정 (Paired t-test)
- 197... 19장. 개수를 비교하기 (카이제곱 검정)
- 208... 20장. 회귀에 대한 주론
- 218... 21장. 분산 분석 (아노바, ANOVA)

235... 22장. 다료인 분산분석(Multifactor ANOVA)

251... 23장. 다중회귀(multiple regression)

261... 24장. 다중회귀에 대한 추가적인 사항들

271... 부록 1. 통계 소프트웨어와 프로그래밍 언어

290... 부록 2. 가설 검정 개발의 역사

- 3 -

Key Points

Chapter 3: Displaying categorical data

- Bar chart for categorical data
- Pie chart for proportions of whole
- Faithful reporting and the area principle
- Contingency tables
- Simpson's paradox

Chapter 4: Displaying quantitative data

- Histograms, Stem-leaf, dot plots
- Shape (mode, symmetrical)
- Center (median, mean)
- Spread (range, IQR, variance, standard deviation)
- Box plots