



Universidad  
del Cauca



# GRAPHS AND DATA DESCRIPTION

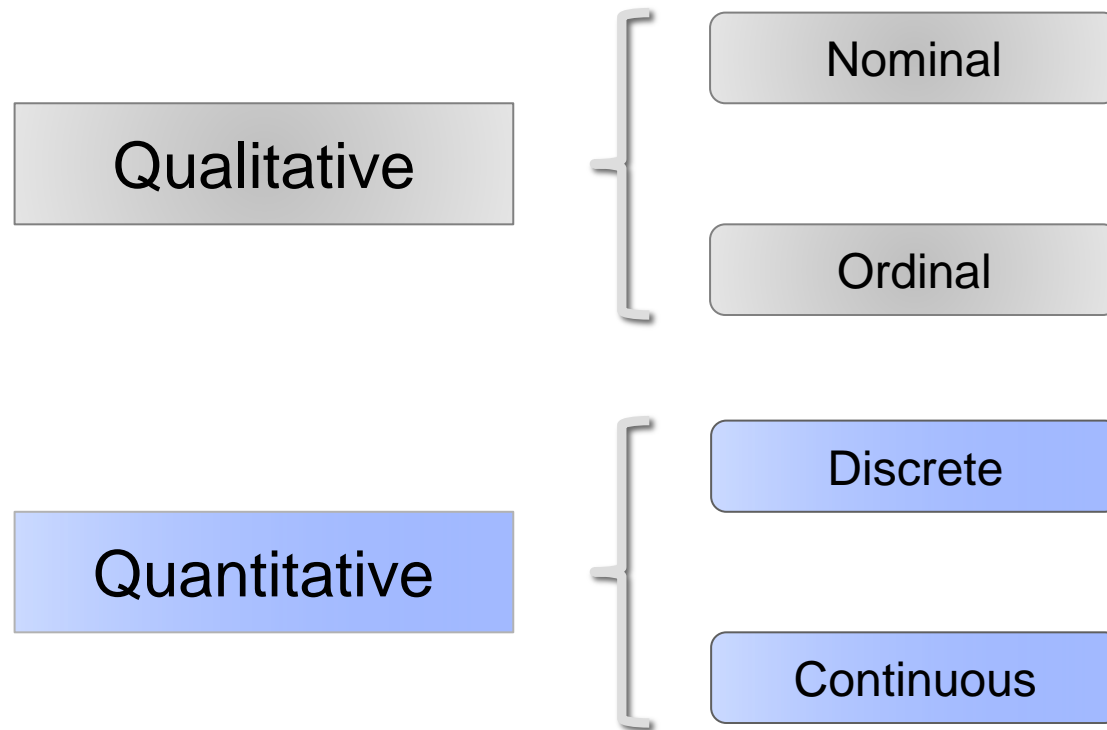
---

***Jose Andres Calvache MD MSc PhD***

*Department of Anesthesiology, Universidad del Cauca, Colombia*

*Department of Anesthesiology, Easmus University Medical Centre, Rotterdam, The Netherlands*

# Variable



# Qualitative variables (n, %)

- Absolute frequencies (n) and proportions (%).
- Sample ( $X_1, X_2, \dots, X_n$ ).
- Each category  $j(1, 2, \dots, k)$ .

$$p_j = \frac{n_j \times 100\%}{n} \quad (j = 1, \dots, k)$$

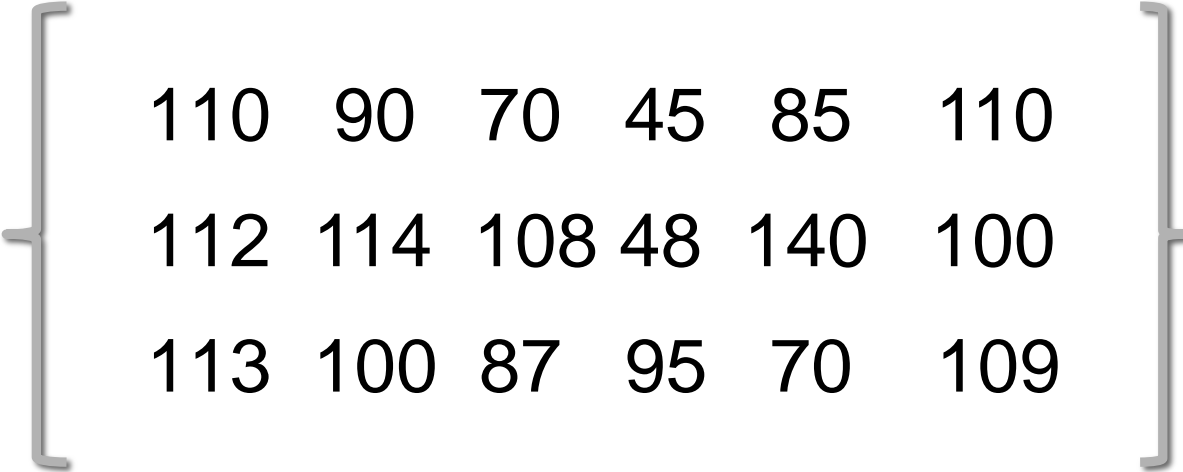
$$\frac{10 \times 100\%}{50}$$



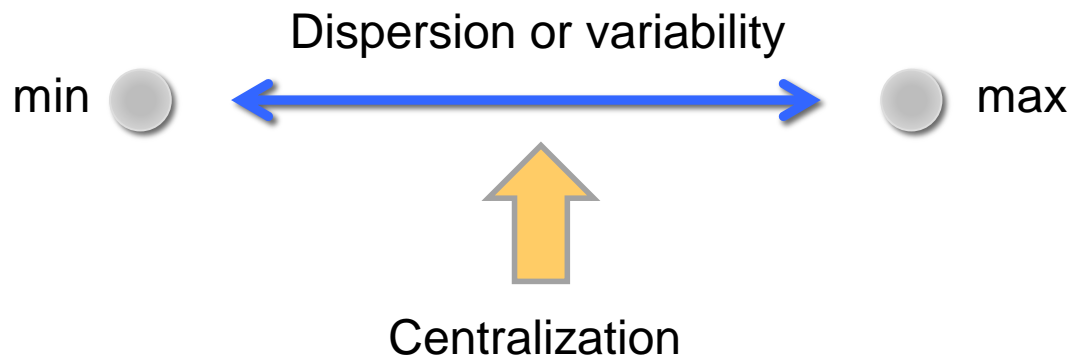
Pain level	Frequency (n)	%
None	10	20
Mild	15	30
Moderate	15	30
Severe	10	20

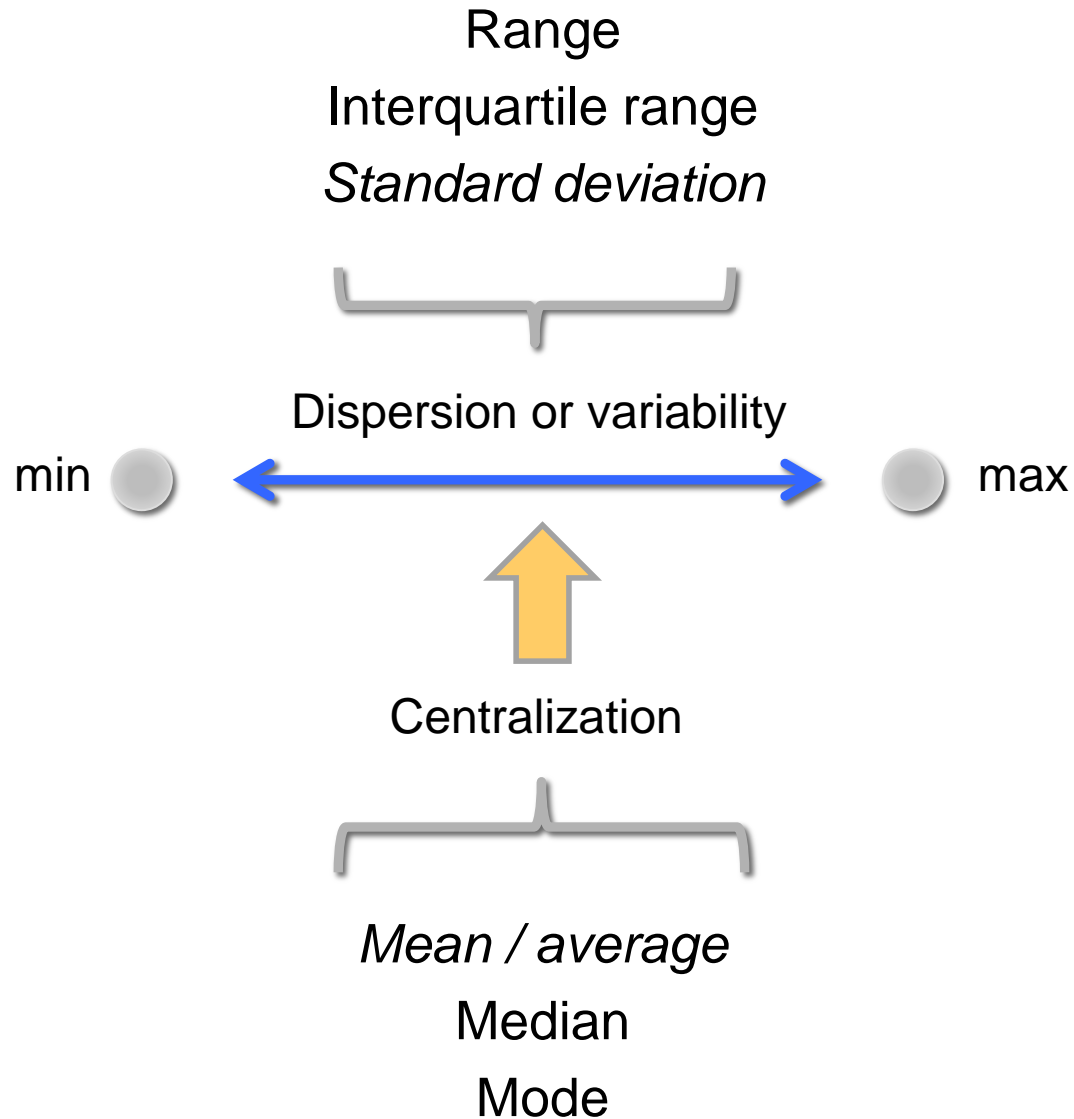
# Quantitative variables

- Glycemic values of 20 medical students (mg/dL)



110	90	70	45	85	110
112	114	108	48	140	100
113	100	87	95	70	109





# Mean ( $\bar{x}$ )

- Arithmetic mean or average
- Sample ( $X_1, X_2, \dots, X_n$ )

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

- Centralization statistic
- Sensitive to extreme values of data
- Gravity center of data

# Standard deviation (*SD*)

- Lets consider a sample ( $X_1, X_2, \dots, X_n$ ) with mean ( $\bar{x}$ )

$$\begin{aligned} SD &= \sqrt{\frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1}} = s \\ &= \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}} \end{aligned}$$

- Measure of variability :  $s \geq 0$
- Sensitive to extreme values
- Num: Sum of squares
- $s^2$ : variance

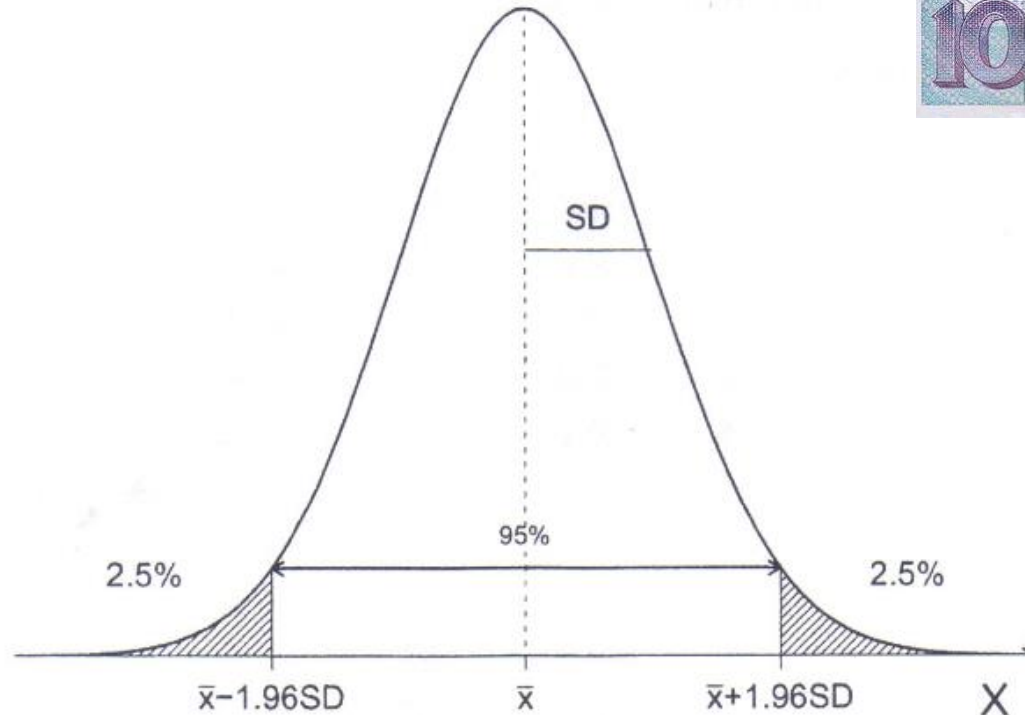
$$\longrightarrow \sum (x - \bar{x})^2$$

## Characteristics of the 355 participating doctors

Characteristics	Mean $\pm$ SD	Frequency	%
Age (years)	43.8 $\pm$ 10.		
Sex			
Male		295	83
Female		60	17
Certification			
Yes		326	93
No		26	7
Professional experience (years)	18.0 $\pm$ 10.5		
Size of practice (patients per week)			
$\leq$ 30		27	8
31–60		70	20
61–90		84	24
> 90		171	48

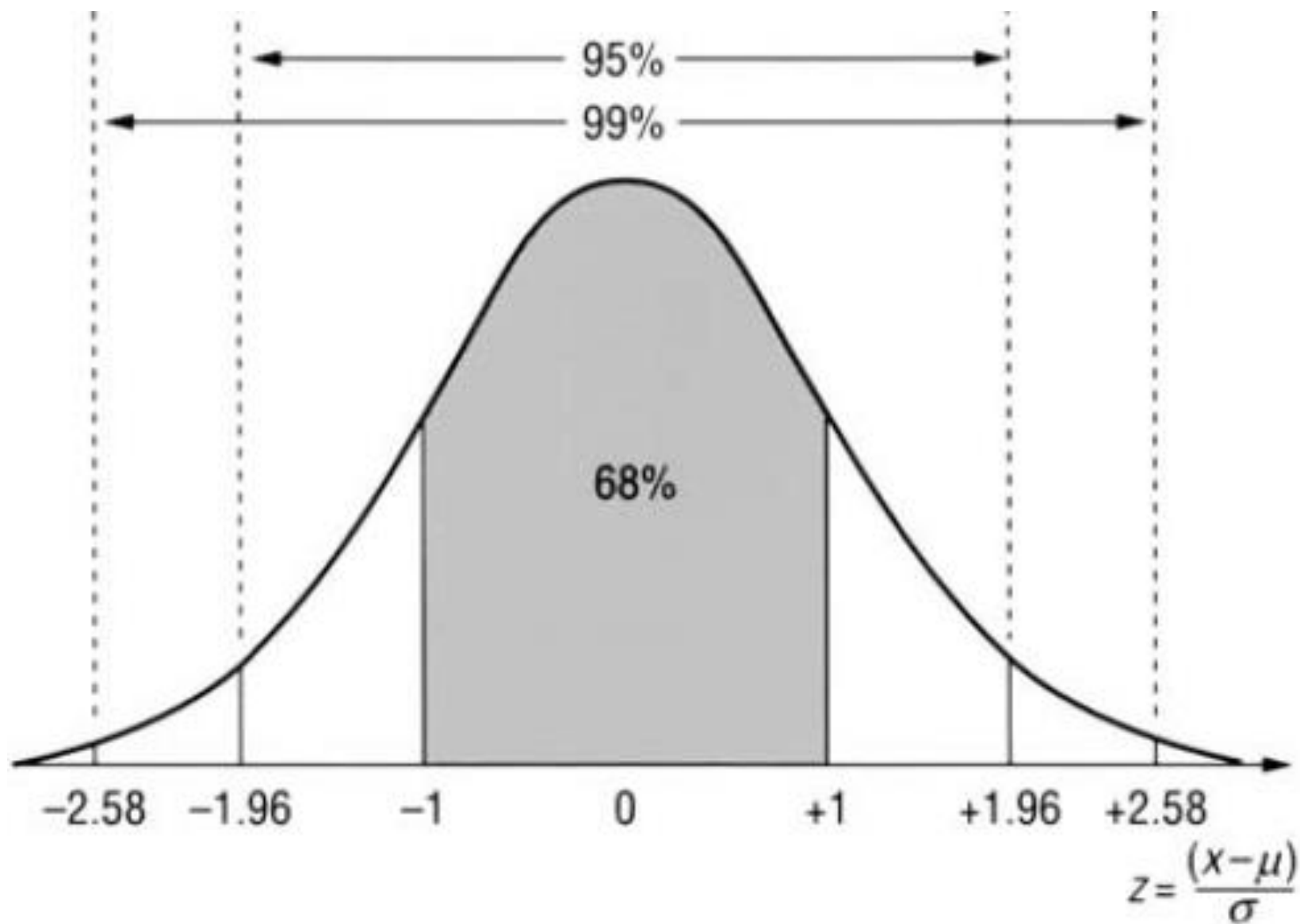


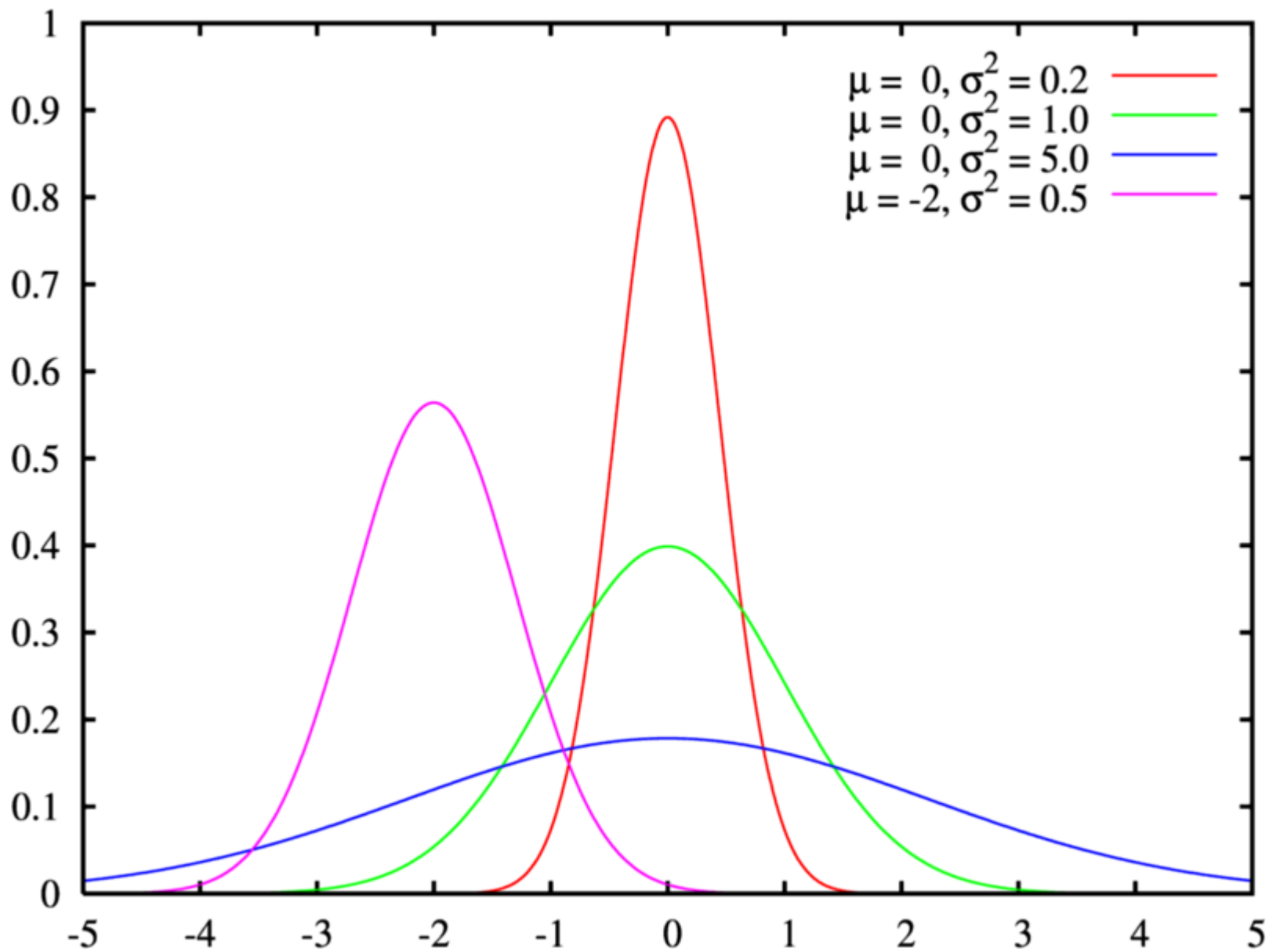
# Normal distribution

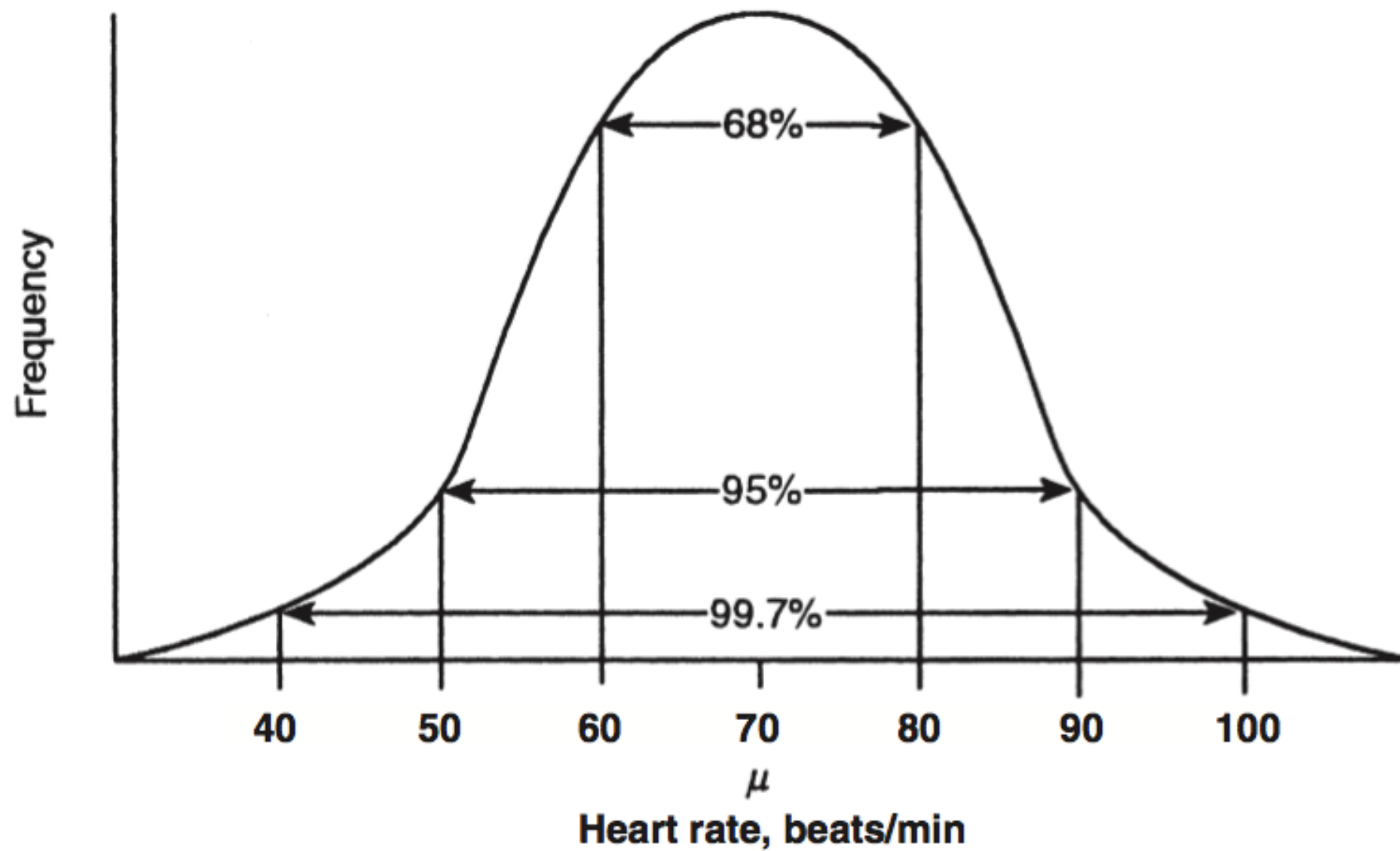


1777 – 1855  
Carl Friedrich Gauss

- Two parameters: mean and SD
- Symmetry around its mean







# SD and normal distribution

- Mean  $\pm 1$  SD contains 68% of the data
- Mean  $\pm 2$  SD contains 95% of the data
- Mean  $\pm 3$  SD contains 99,7% of the data

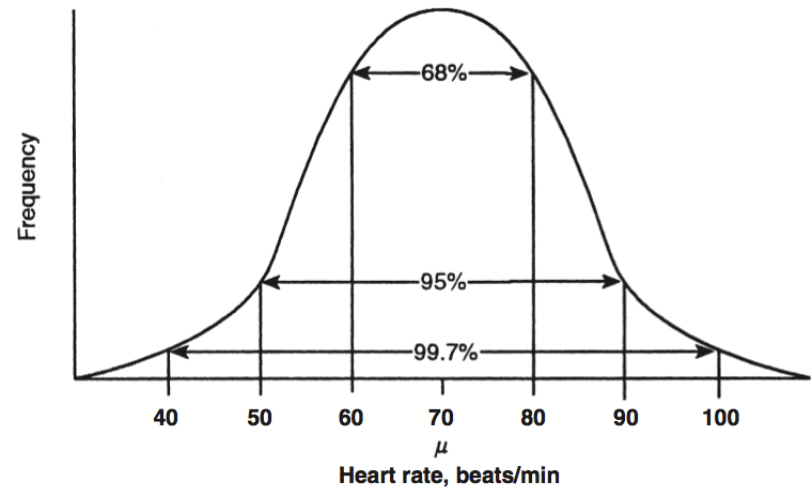
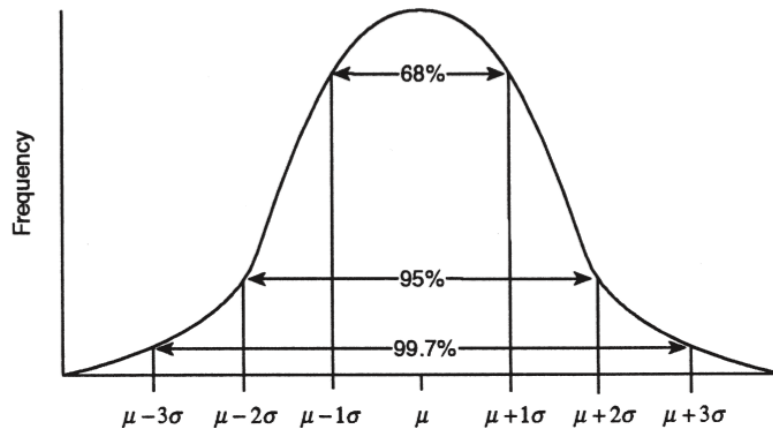
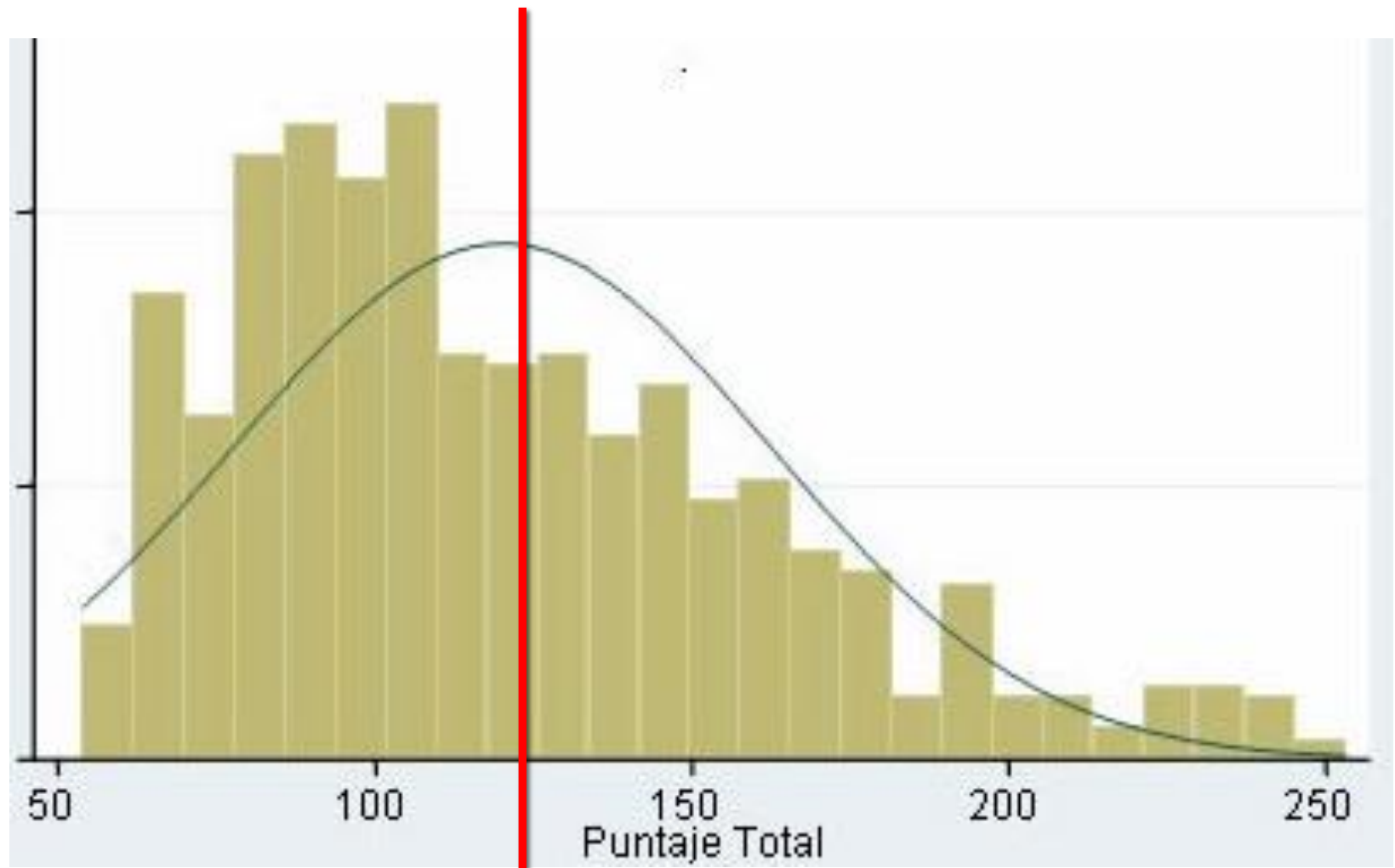


Table 1

Demographics and peri-operative data.

	Placebo <i>n</i> = 34	Morphine 100 µg <i>n</i> = 34
Age (year)	51 ± 8 (38–69)	51 ± 8 (38–66)
Weight (kg)	71 ± 12	70 ± 9
Length (cm)	166 ± 6	166 ± 7
Diagnosis ( <i>n</i> ) myoma uteri	24	23
sarcoma/carcinoma uteri	10	11
Duration of surgery (min)	95 ± 27	102 ± 29
Range	(42–163)	(60–170)
Blood loss during surgery (ml)	307 ± 246	314 ± 221
Range	(0–1050)	(0–800)

Values are expressed as mean, SD (range), and number of patients (*n*).

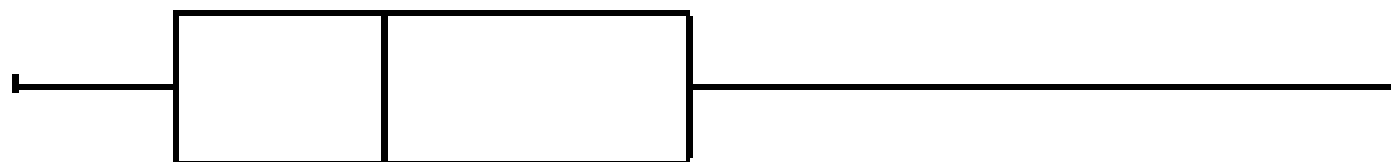
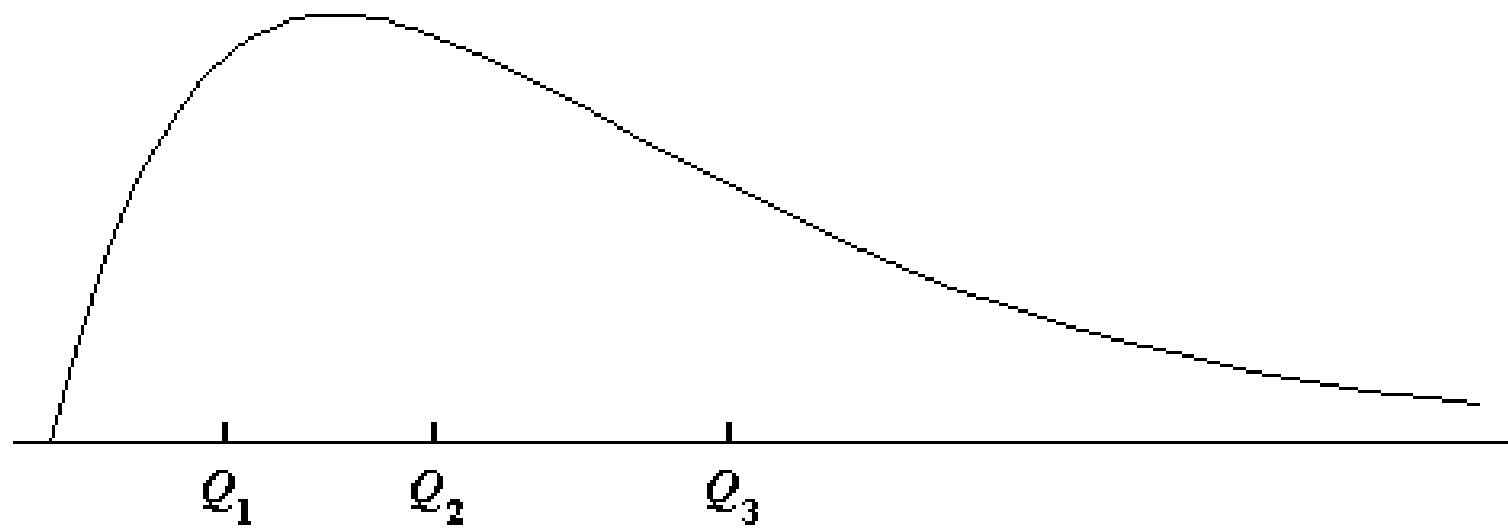


$\bar{x}$

# Median and quartiles

- Asymmetric distribution
- Non parametric distribution
- Median: centralization
- Quartiles: variability
- Interquartile range (IQR): variability
- 10 percentile
- $P_{10}$ : left 10% of data below this boundary and 90% above.





# Median

- Sample  $(X_1, X_2, \dots, X_n)$
- Reorder data from lowest to highest

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

- $X_{(1)}$ : minimum; y  $X_{(n)}$ : maximum. Between them «range»

$$x_{(\frac{n+1}{2})} \quad \Bigg| \quad \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2}$$

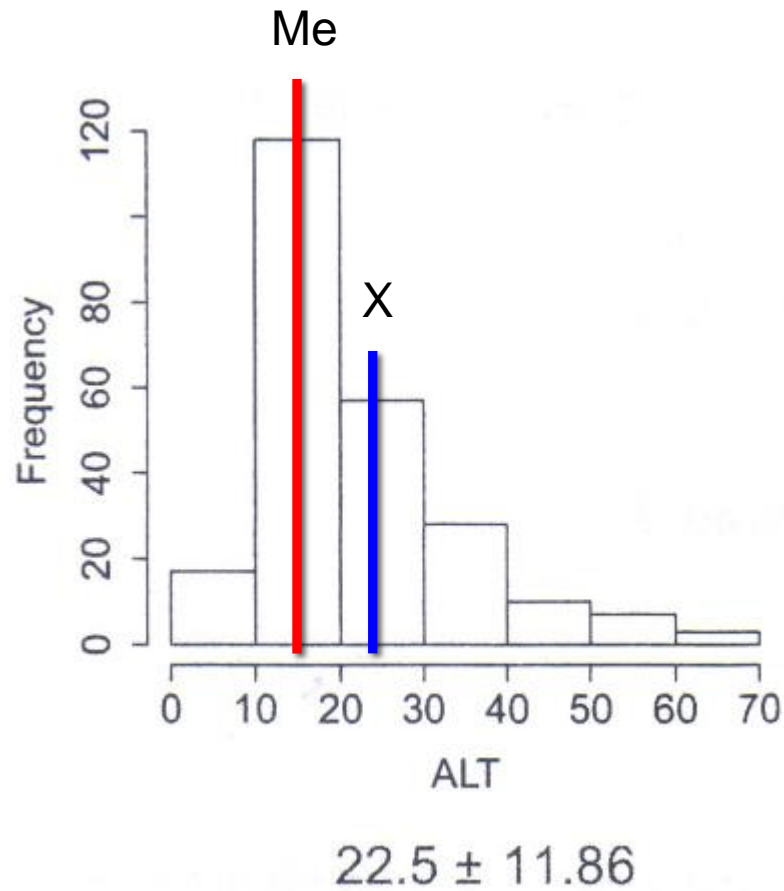
- It is not affected by the extreme data. Robust
- Useful in asymmetric distributions
- $P_{(50)} = \text{Median} = \text{Quartile 2 (Q2)}$

## Data from 240 medical students (ALT)

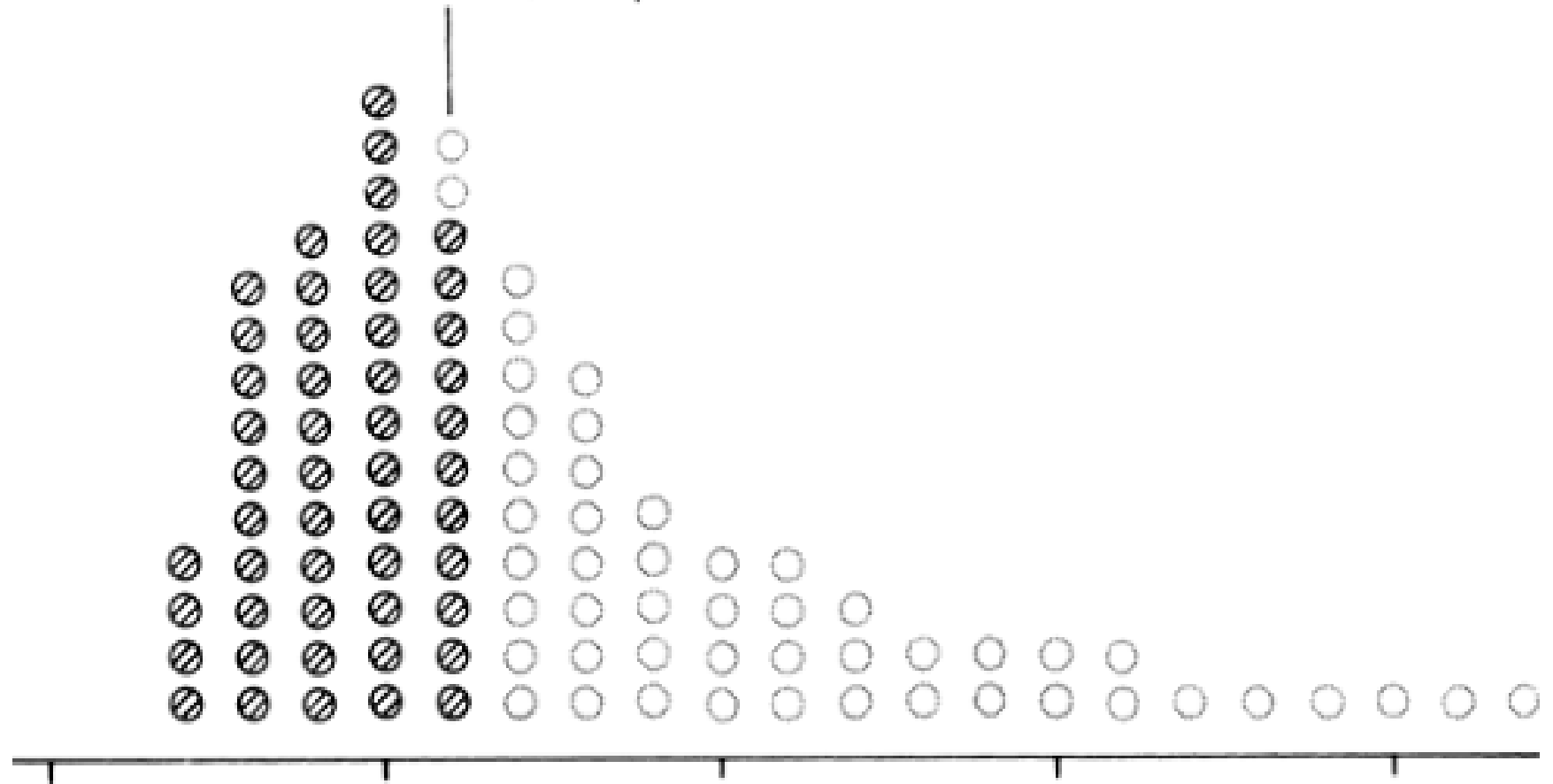
5	6	6	6	7	8	8	8	8	8	9	9	9	10	10	10	10	11	11	11
11	11	11	11	11	11	11	11	12	12	12	12	12	12	12	12	12	12	12	12
12	13	13	13	13	13	13	13	13	13	13	13	13	13	14	14	14	14	14	14
14	14	14	14	14	14	14	15	15	15	15	15	15	15	15	15	15	16	16	16
16	16	16	16	16	16	16	16	17	17	17	17	17	17	17	17	17	18	18	18
18	18	18	18	18	18	18	19	19	19	19	19	19	19	19	19	19	19	19	19
<u>20</u>	20	20	20	20	20	20	20	20	20	20	20	20	20	20	21	21	21	21	<u>21</u>
21	21	21	21	21	21	22	22	22	22	22	22	22	22	23	23	23	23	23	24
24	24	25	25	25	25	25	25	25	25	25	25	25	26	26	26	26	26	27	28
28	28	28	28	28	29	29	30	30	30	30	30	31	31	31	31	31	32	33	34
34	35	35	36	36	36	36	36	36	37	37	37	38	38	39	39	39	40	40	40
41	42	45	45	46	47	47	48	48	49	51	51	51	53	54	55	55	62	65	69

$$\text{Median} = \frac{x_{(120)} + x_{(121)}}{2} = \frac{19 + 20}{2} = 19.5 \text{ IU/L}$$

## Data from 240 medical students (ALT)



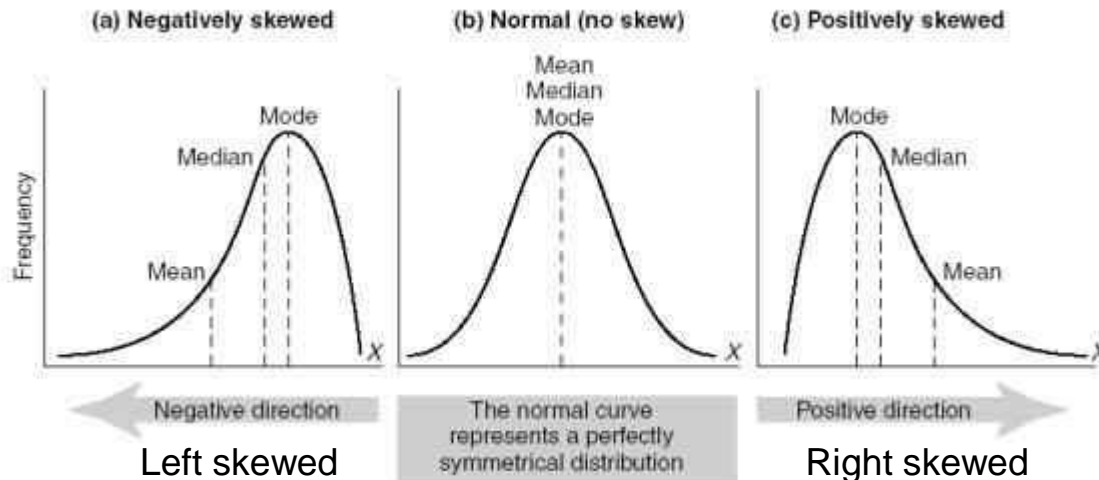
Median (50th percentile)



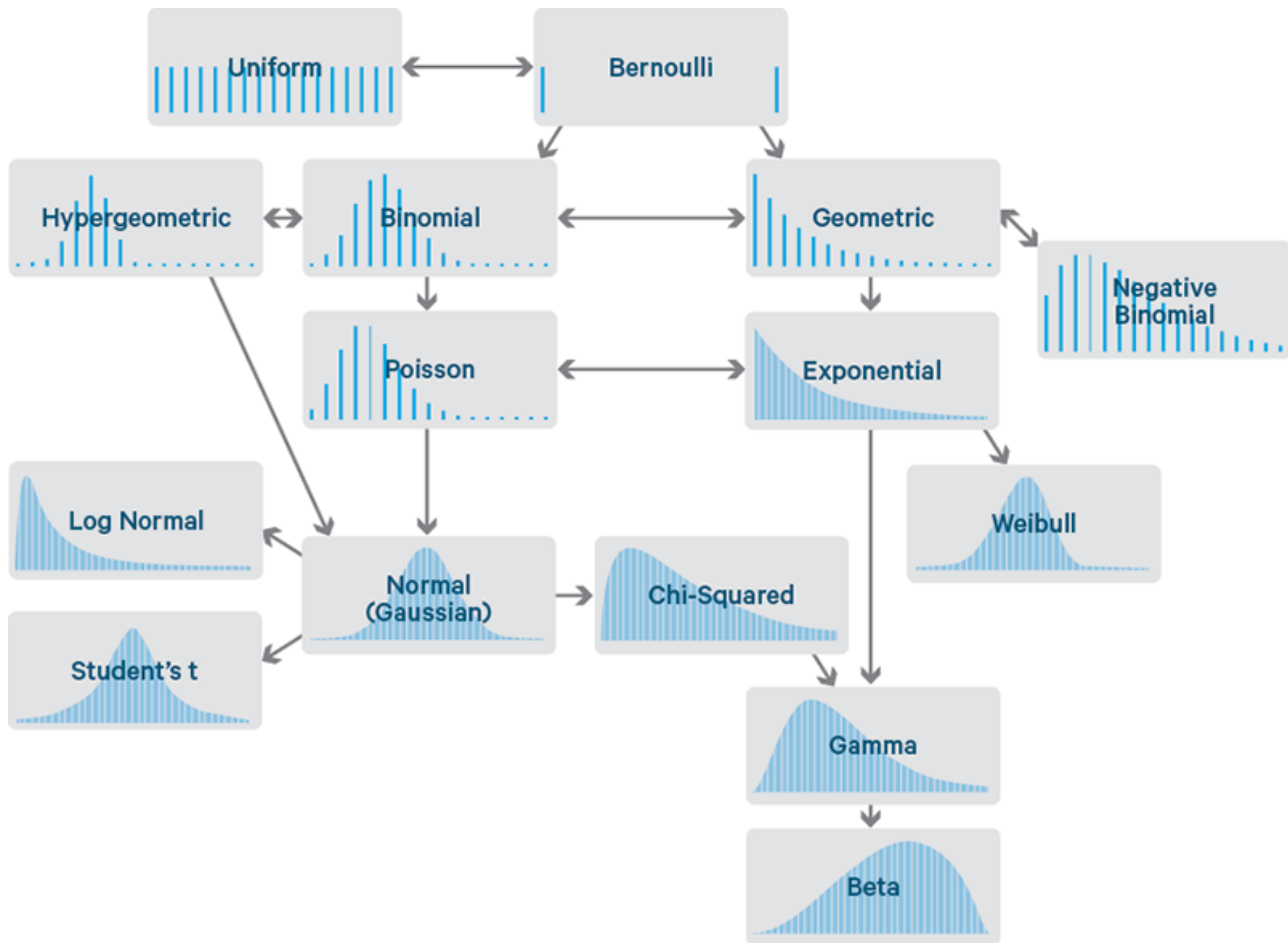
# Median and IQR

- Useful in asymmetric distributions
- Robust parameters
- IQR contains 50% of data
- Components of the boxplot

# Distribution of the data



- $\text{Mean} > \text{Median}$  : Right skewed distribution
- $\text{Mean} < \text{Median}$  : Left skewed distribution





# Graphs in statistics

- Fisher (1890-1962)
  - Type of graphs depends on type of data to summarize
  - Very useful to check and present the data
  - First step during analysis
- 
1. Qualitative: Bar graph, Pie graph
  2. Ordinal: Bar graph, (order)
  3. Discrete: Bar graph
  4. Continuous: Histogram, density plot, frequency plot, box plot, scatter plot.

# Graphs in statistics

## **Categorical Variables**

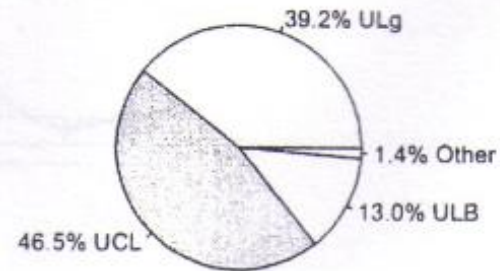
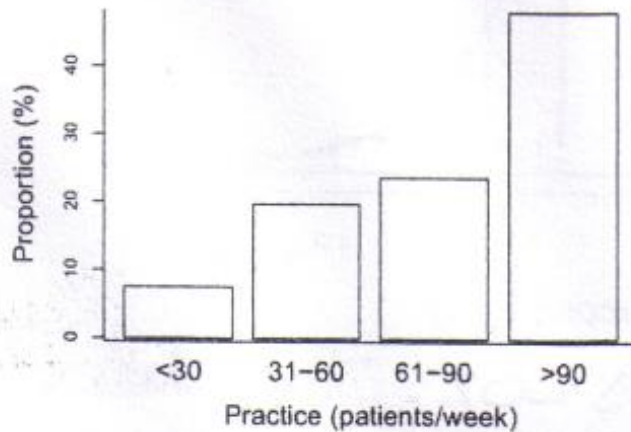
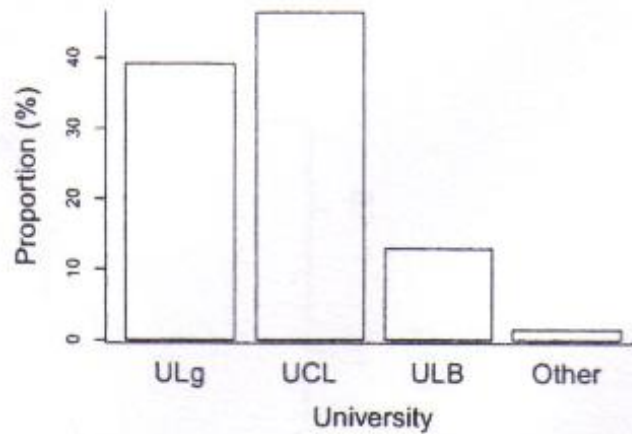
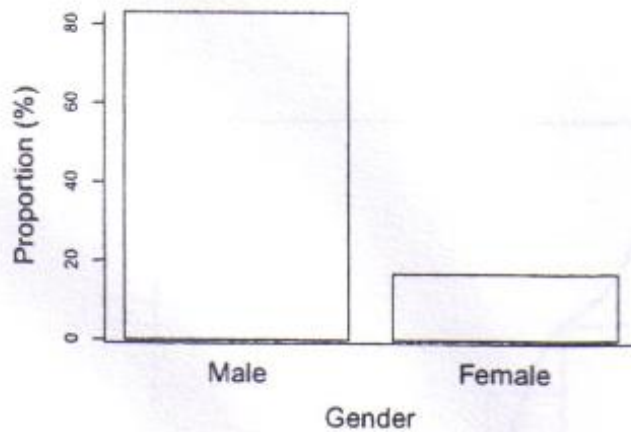
- Frequency distribution
- Bar chart
- Pie chart
- Pareto diagram

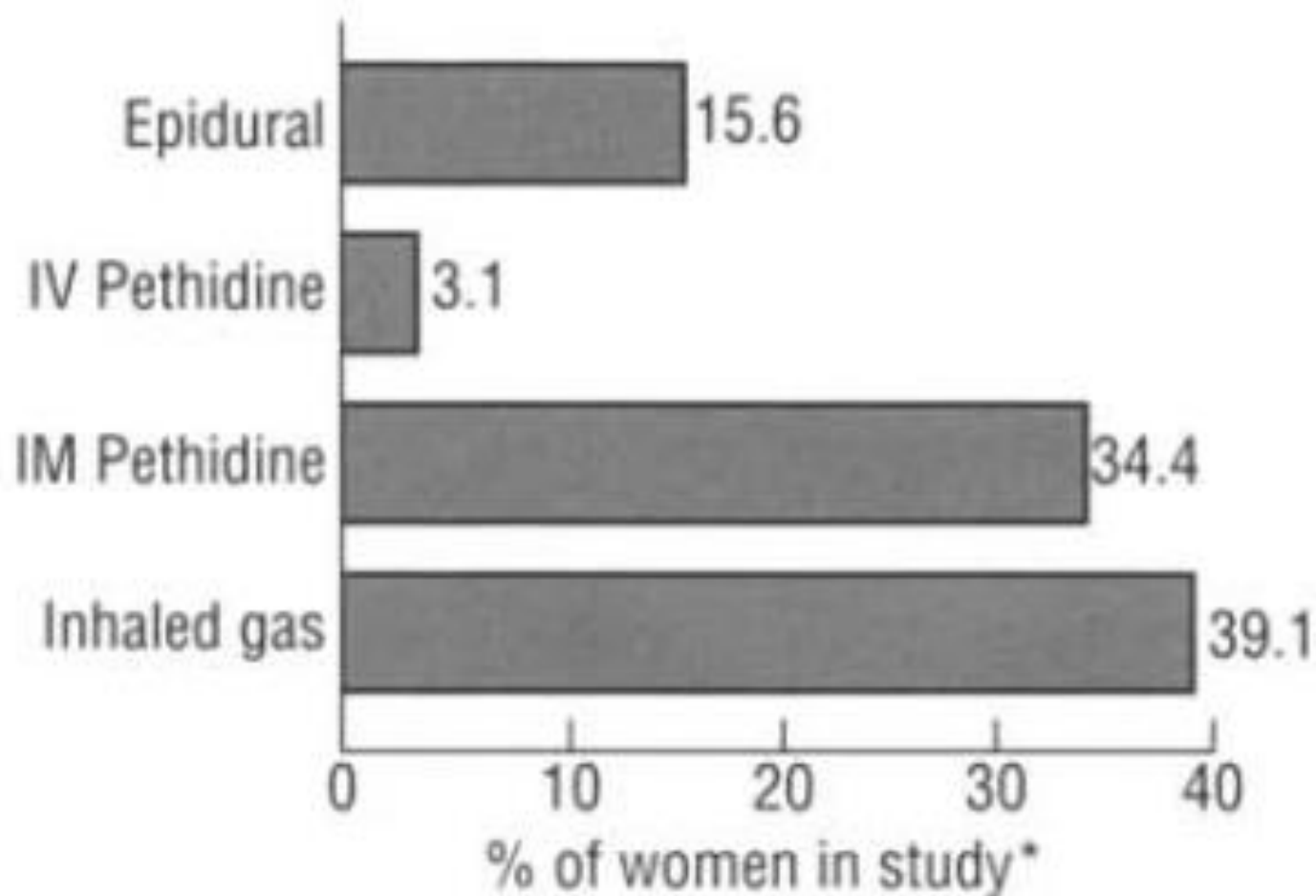
## **Numerical Variables**

- Line chart
- Frequency distribution
- Histogram and ogive
- Scatter plot

# Qualitative

N = 355 GPs





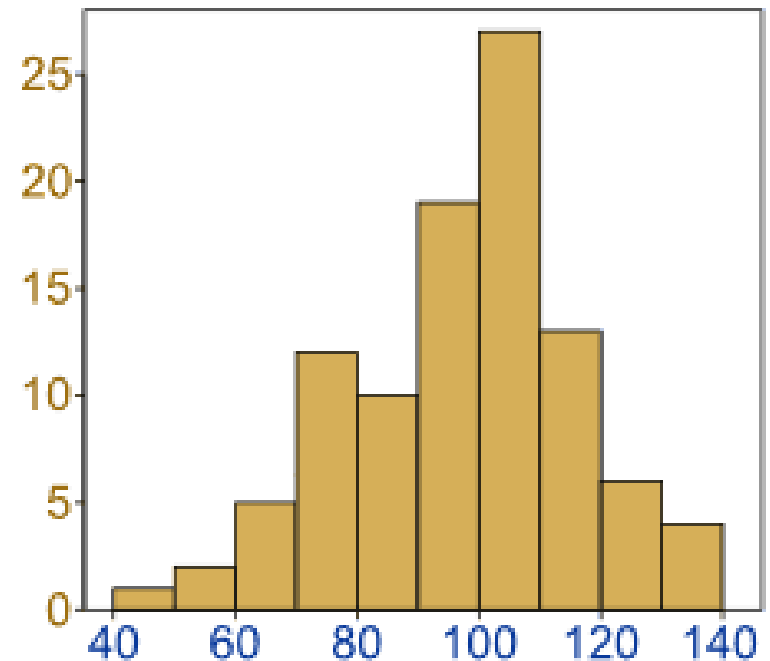
\*Based on 48 women with pregnancies

# Quantitative data

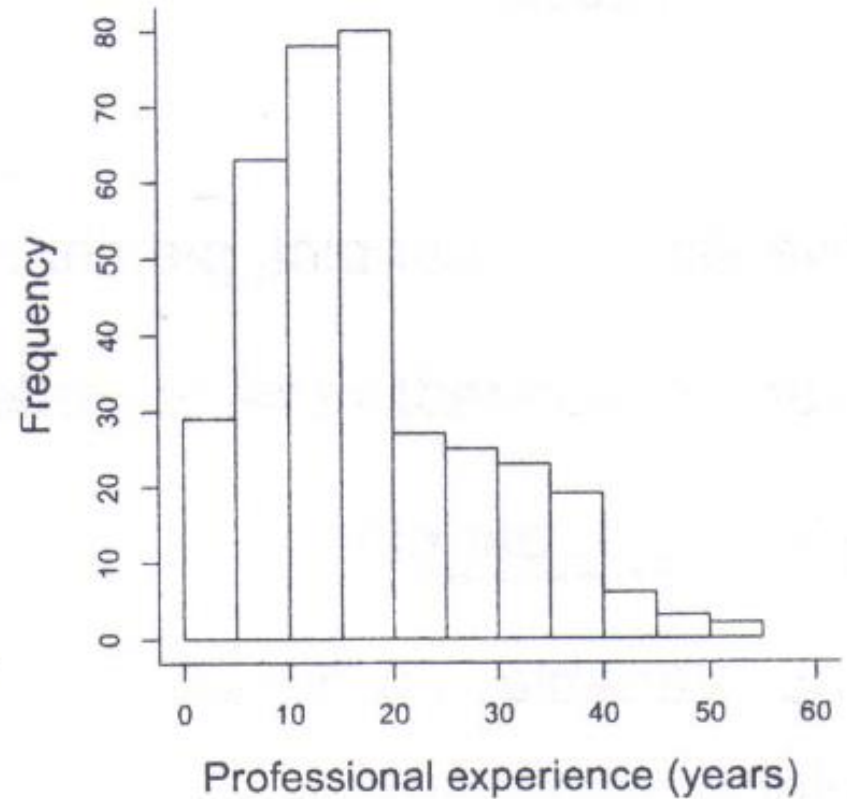
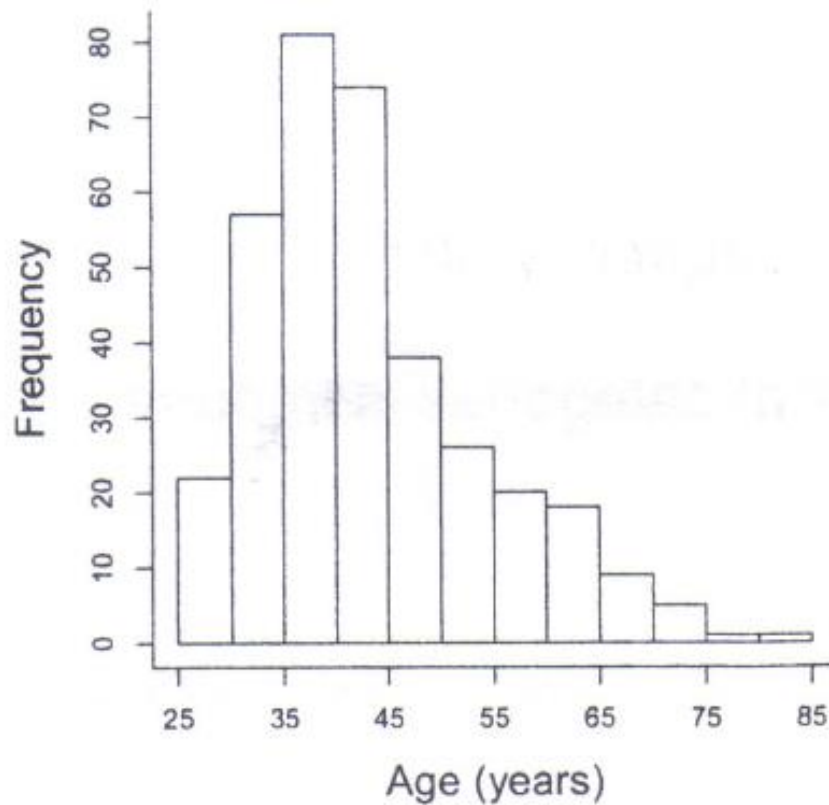
- Discrete / continuous
- Absolute or relative frequency
- Histogram

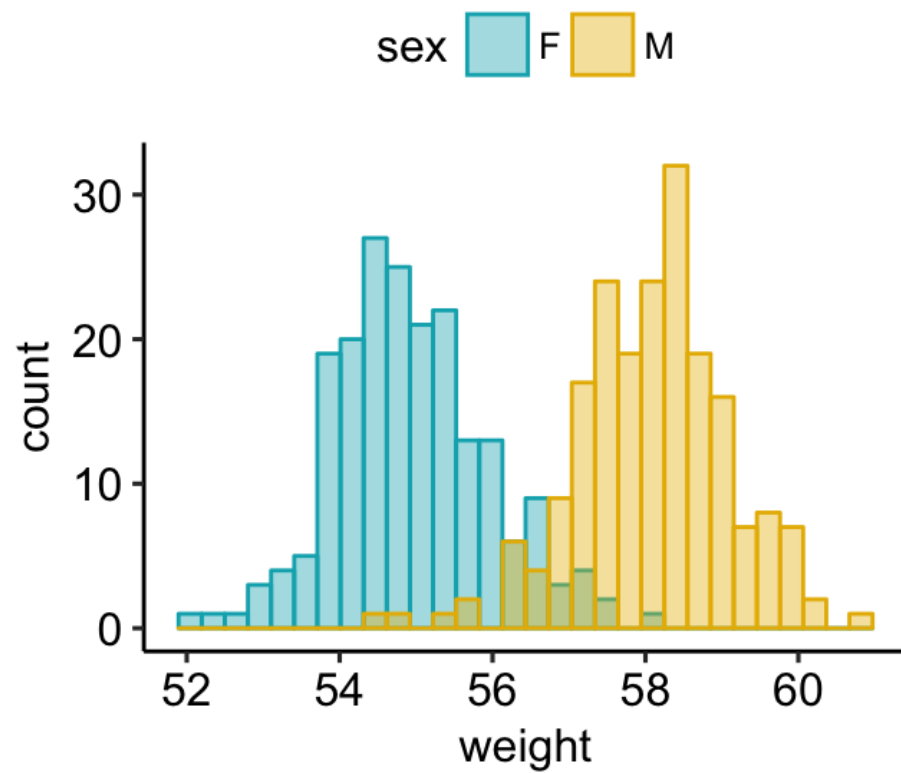
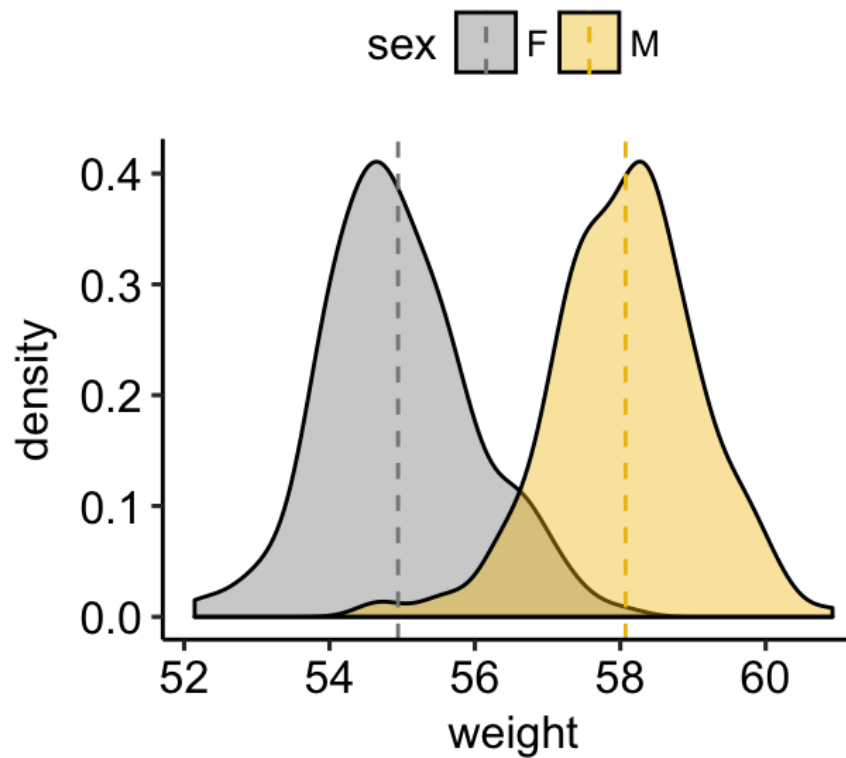
110	90	70	45	85	110
112	114	108	48	140	100
113	100	87	95	70	109

n = 20



## Histograms - N = 355 GPs

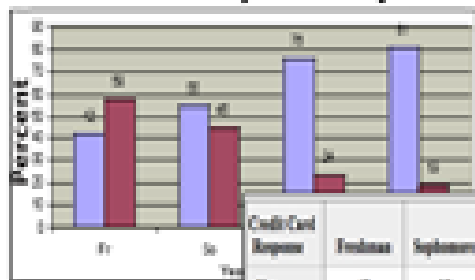




# Relationships between Variables

## Two Categorical Variables

Bar Graph of Percents by Groups

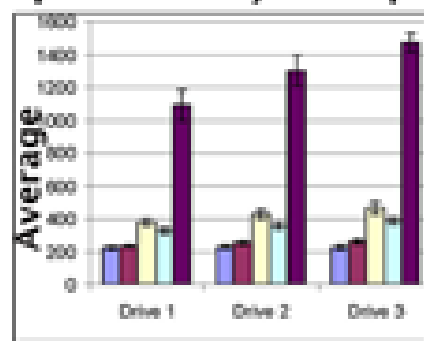


Table

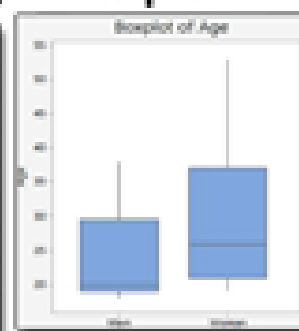
Credit Card Response	Ford Focus	Saturn	Jaguar	Tesla	Total
Yes	40	50	75	80	245
No	60	50	25	20	155
Total	100	100	100	100	400

## One Measurement, One Categorical Variable

Bar Graph of Proportions by Groups

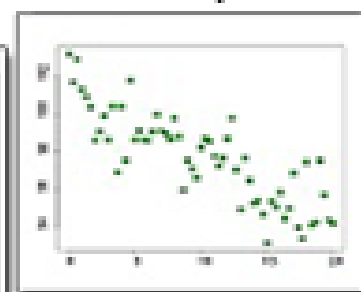


Side by side Boxplot

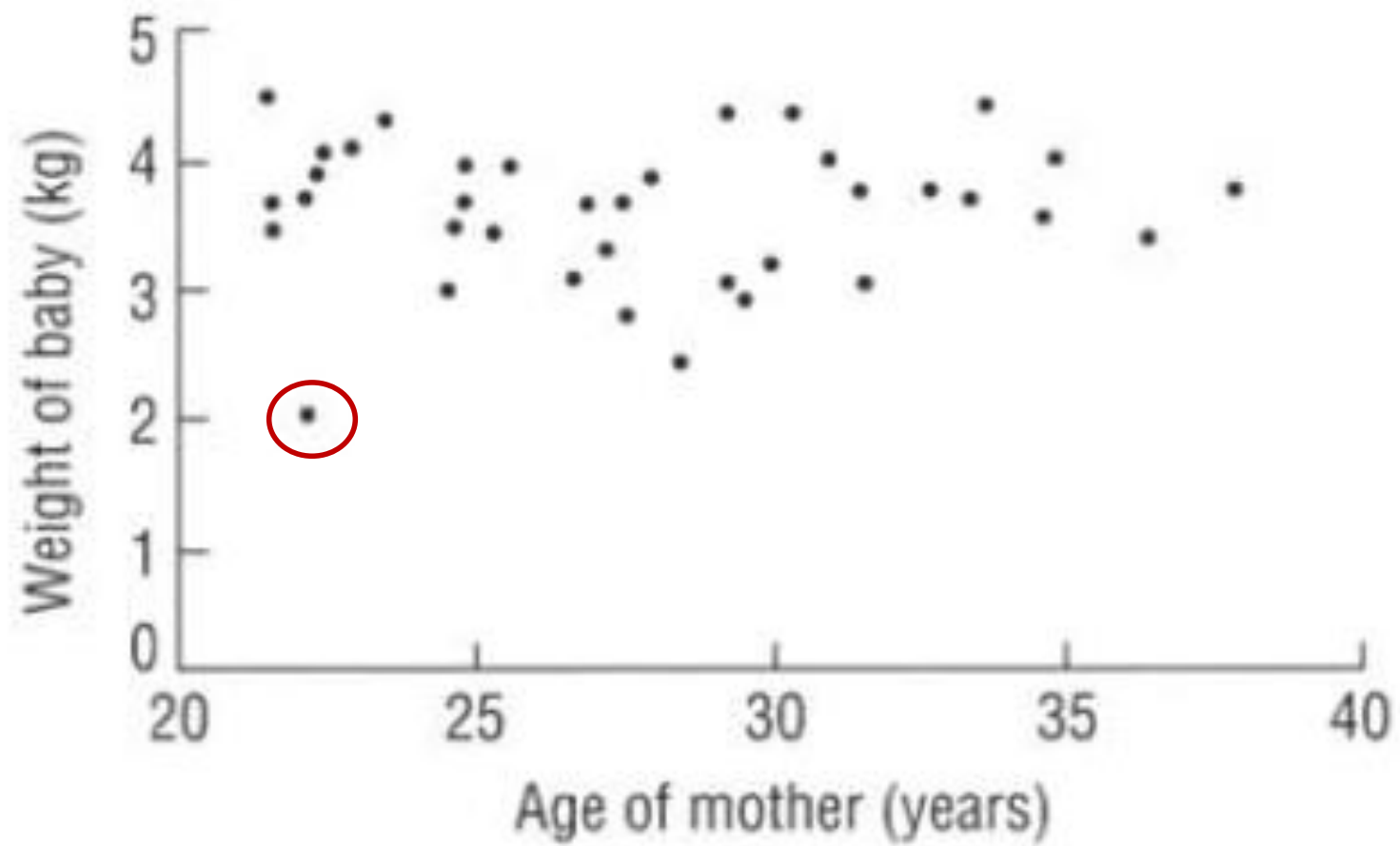


## Two Measurement Variables

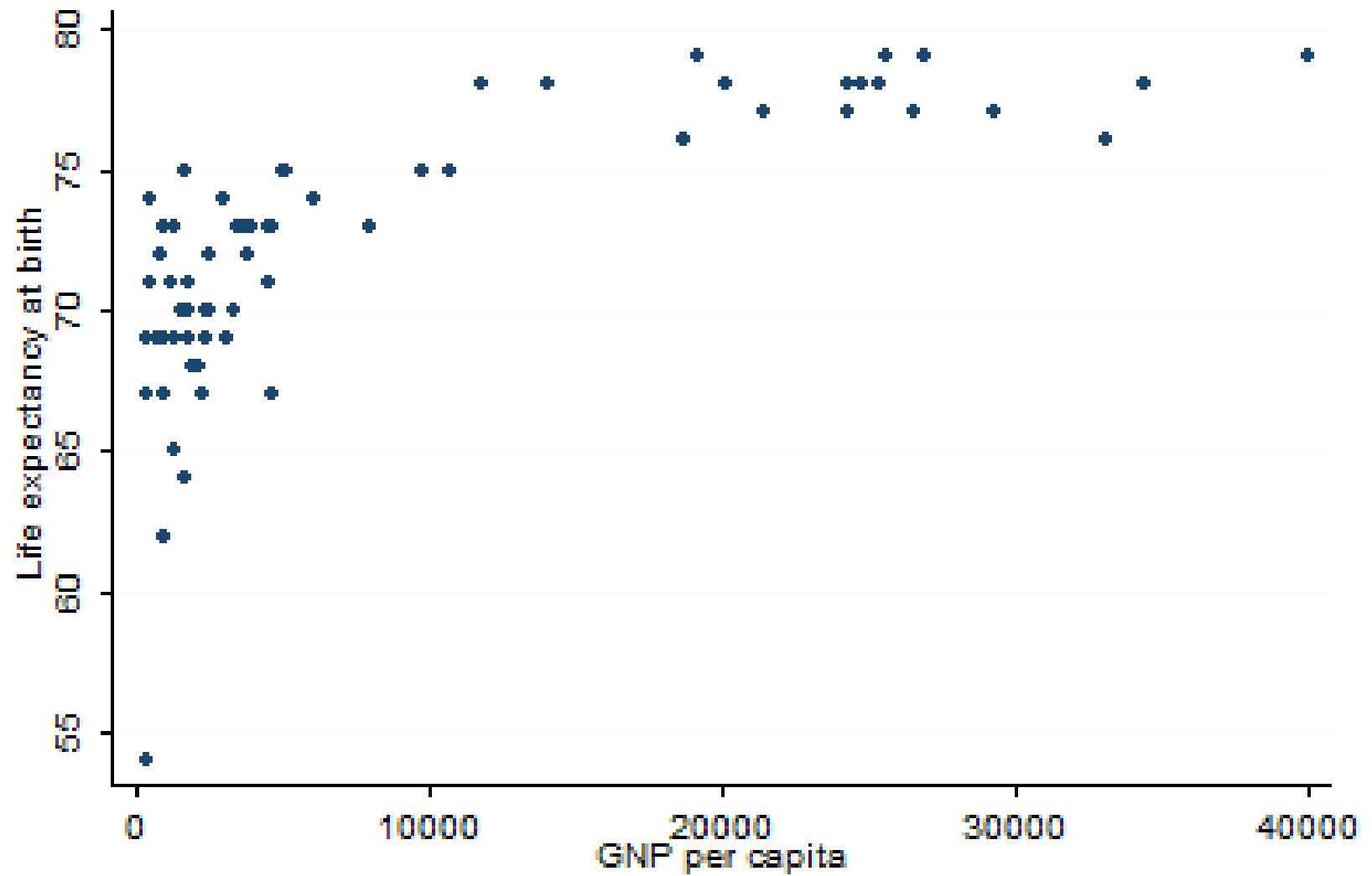
Scatterplot

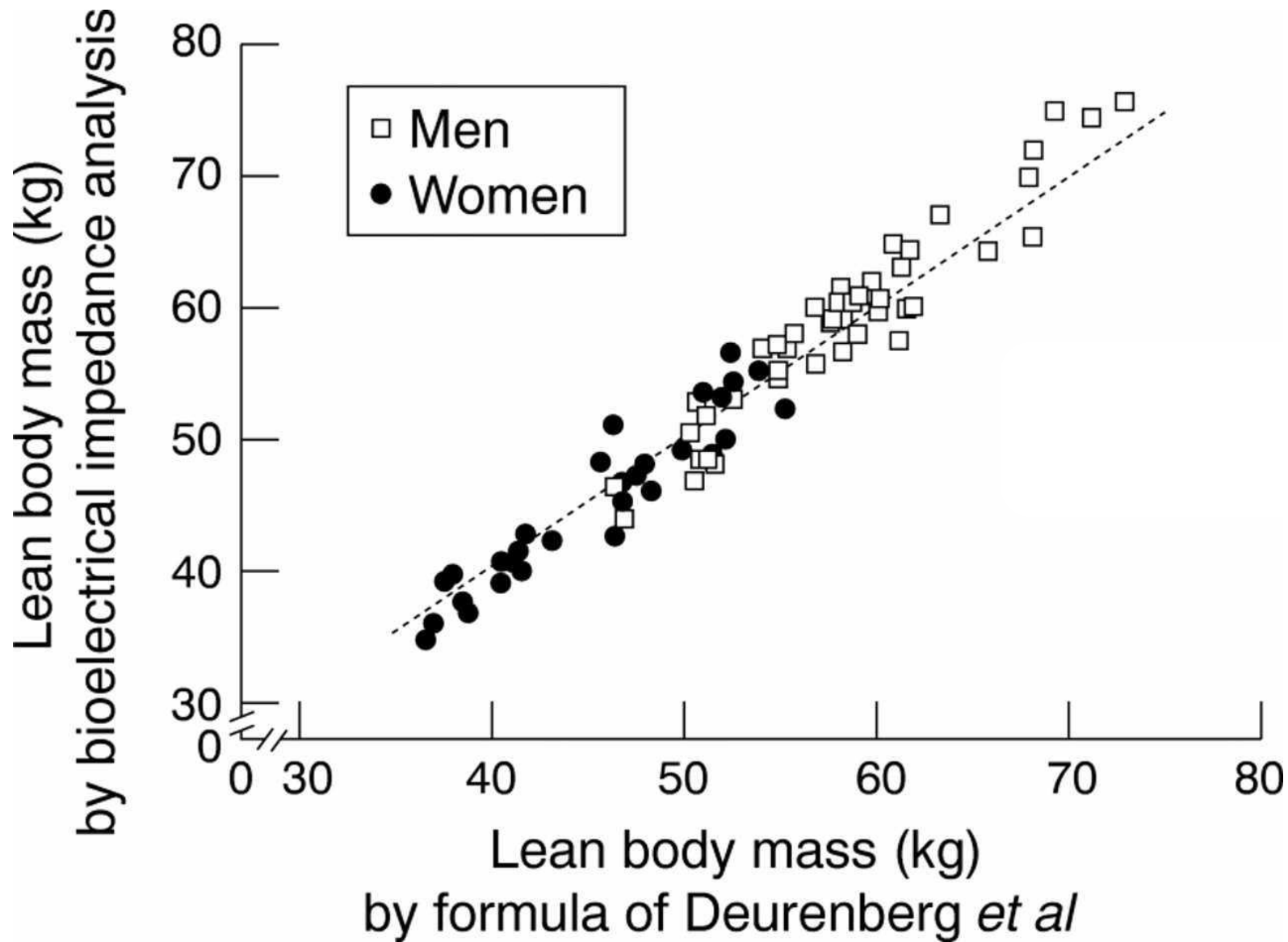




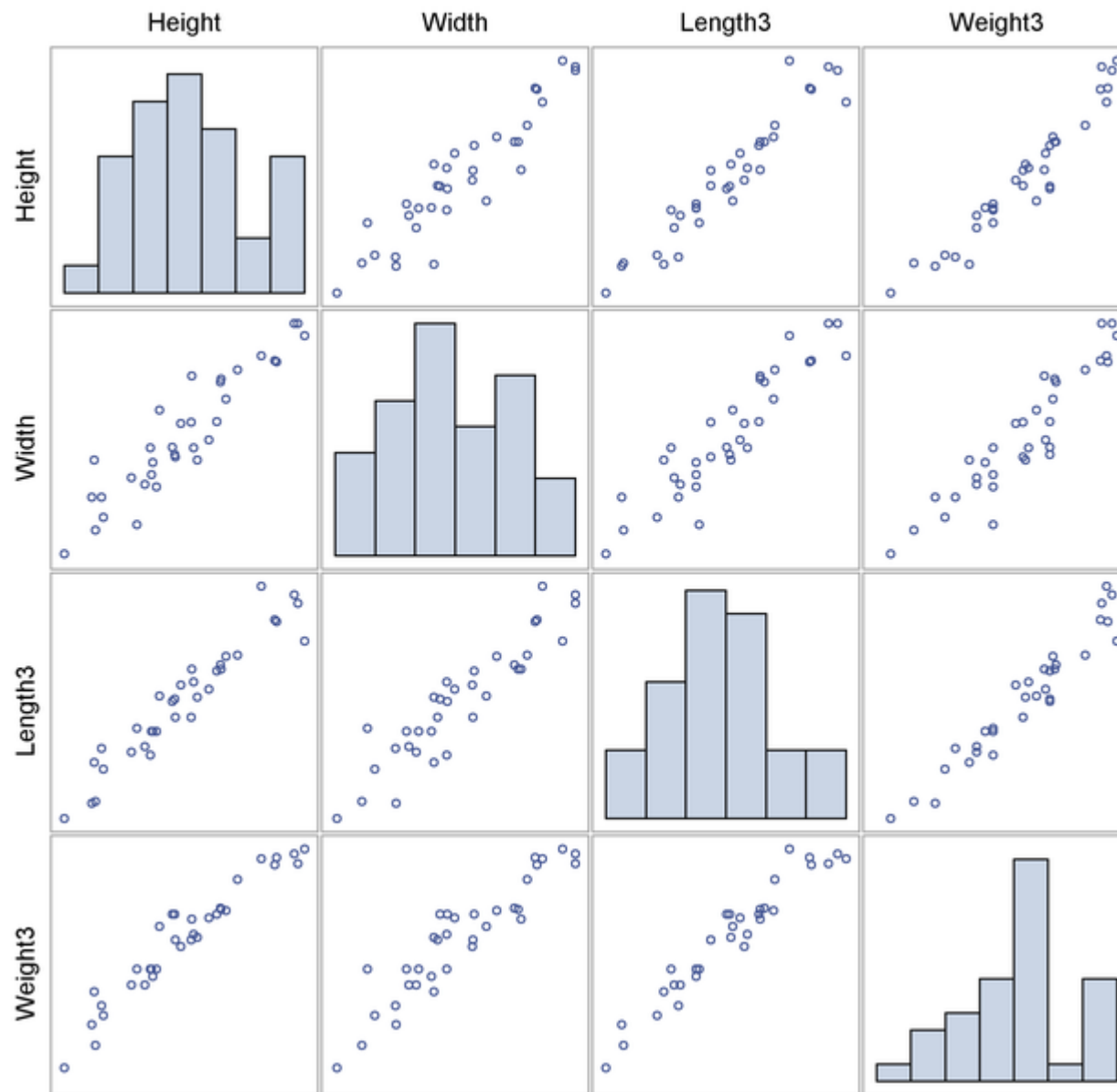


## Life Expectancy and GNP

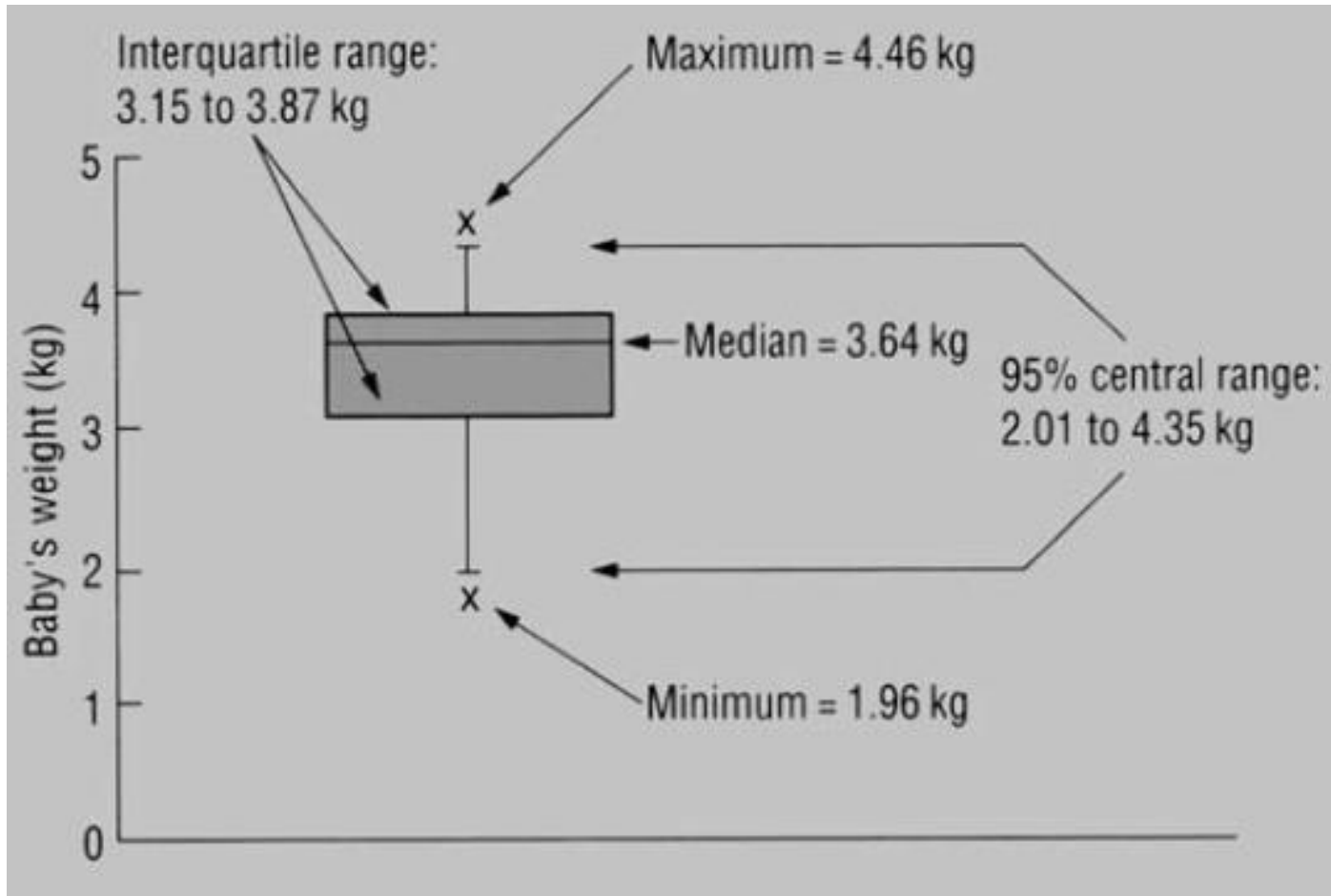




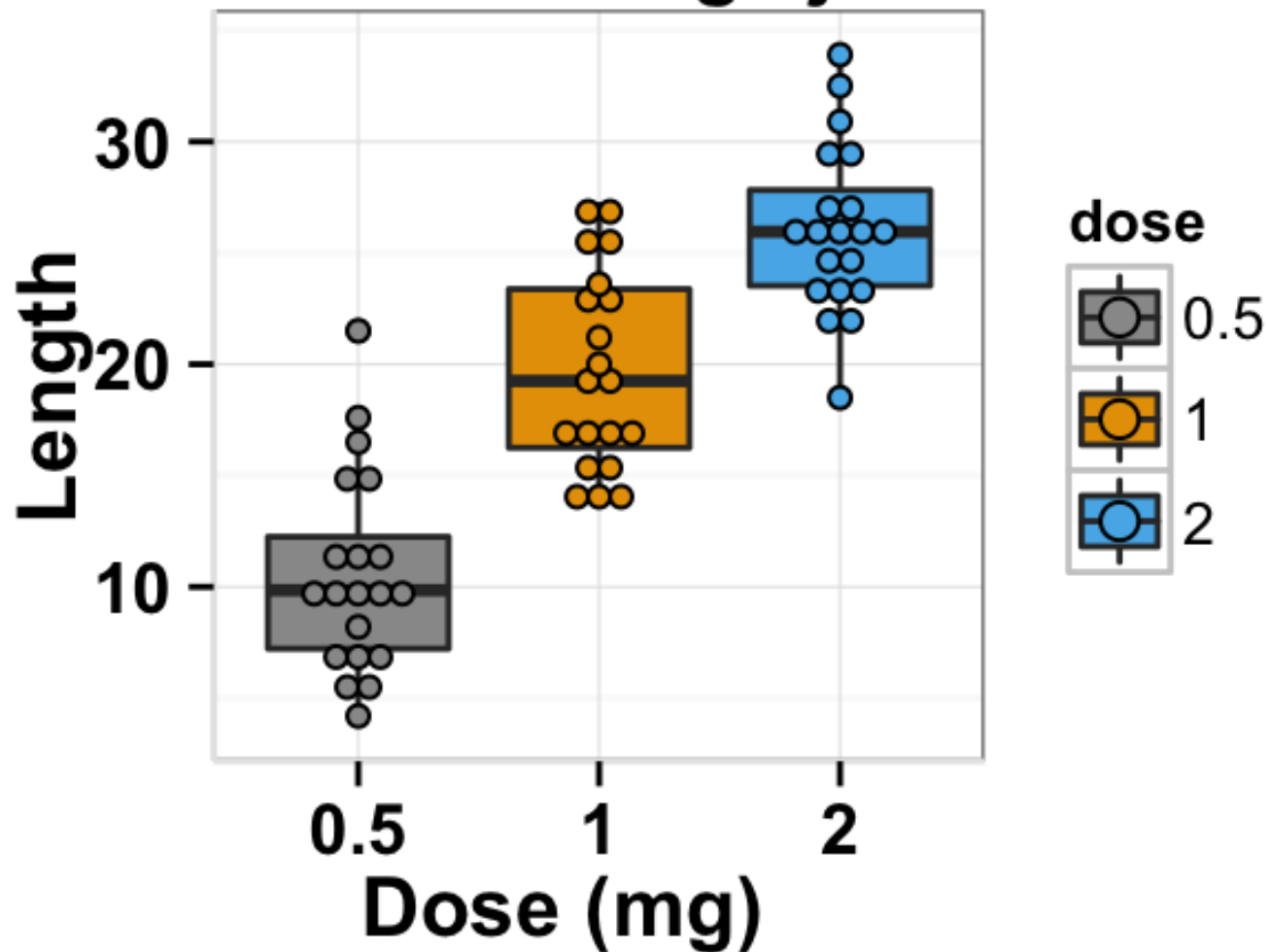
Scatter Plot Matrix



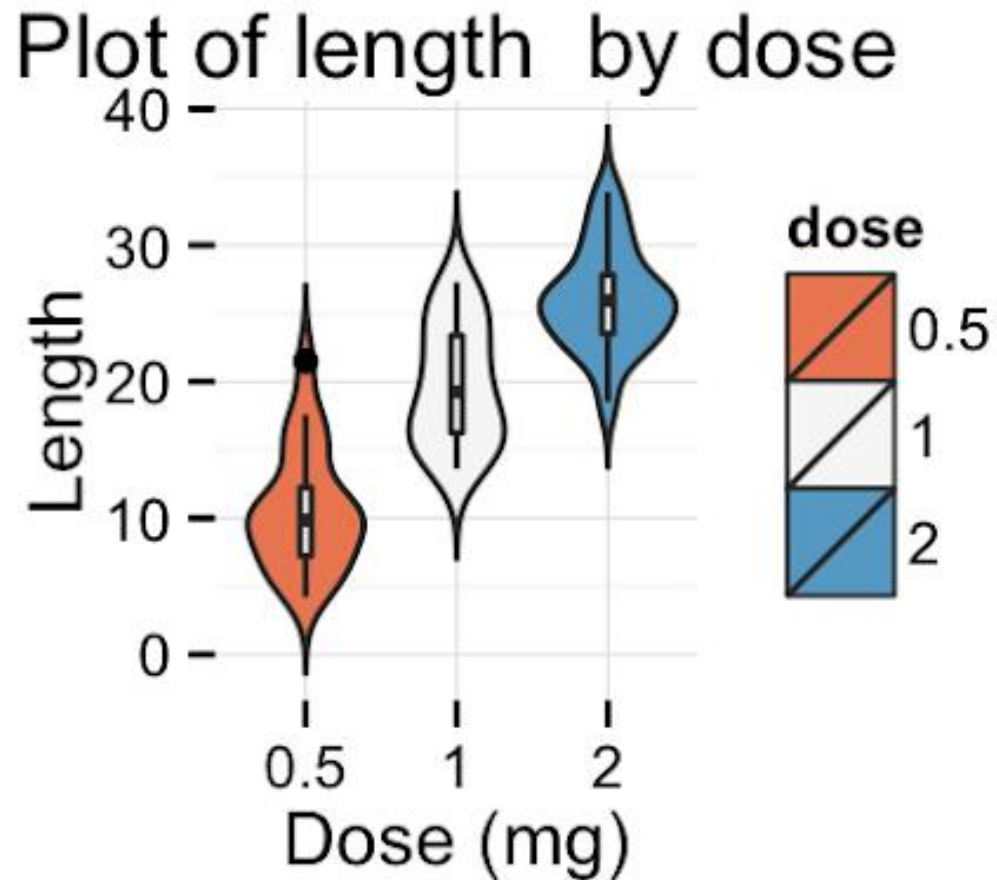
# Boxplot



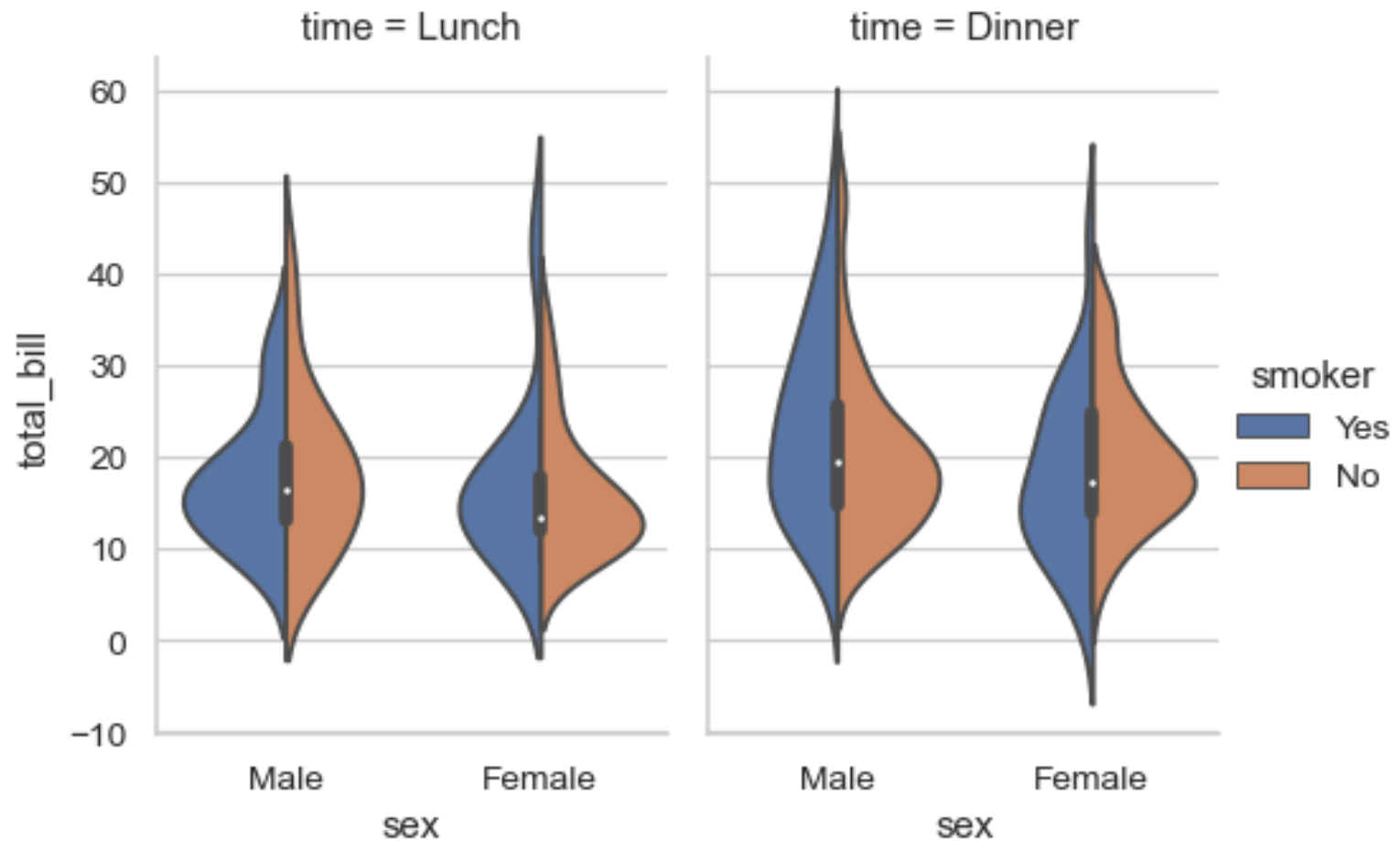
## Plot of teeth length according to vitamin C/Orange juice dose



# Violin plot 1

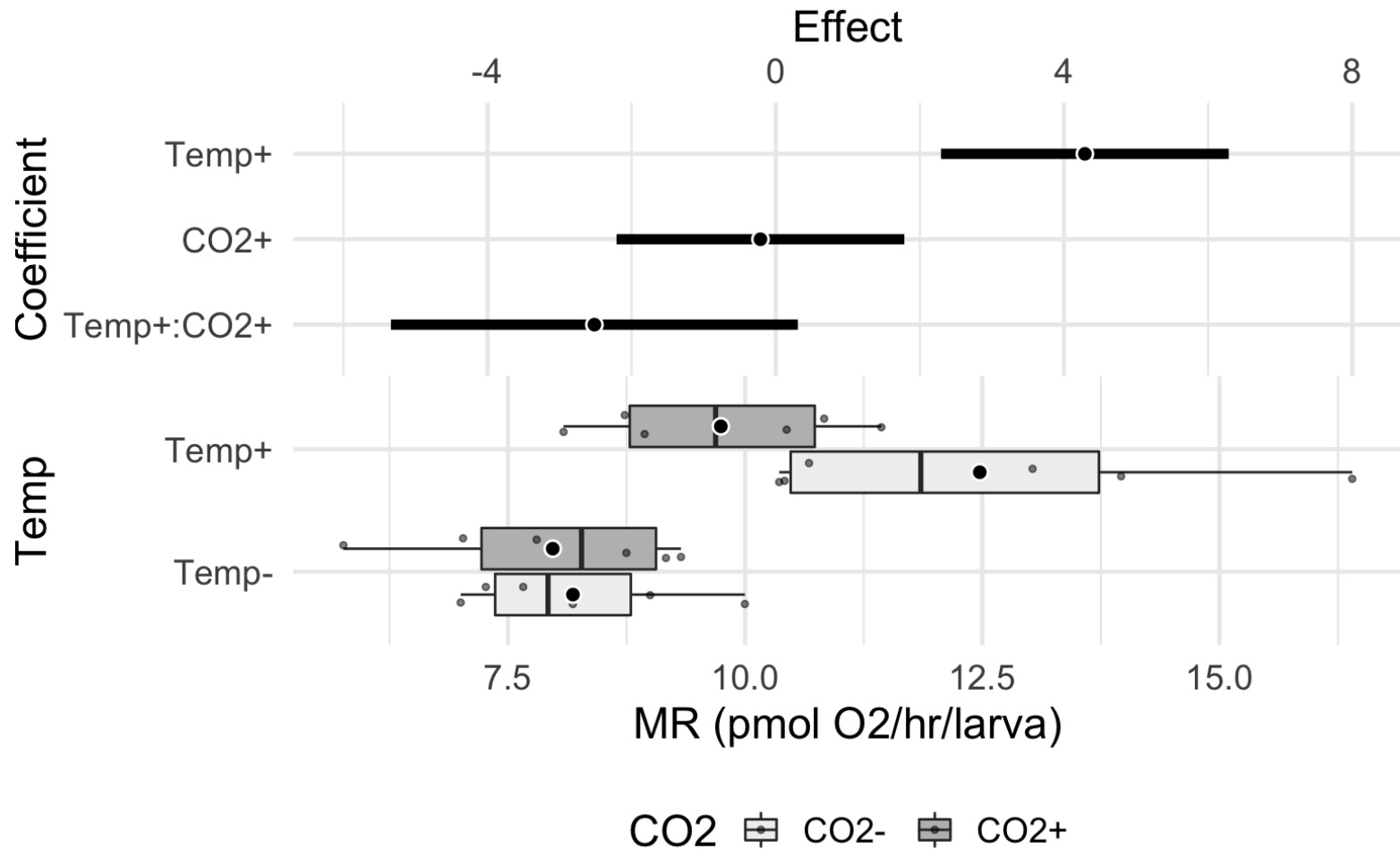


# Violin plot 2

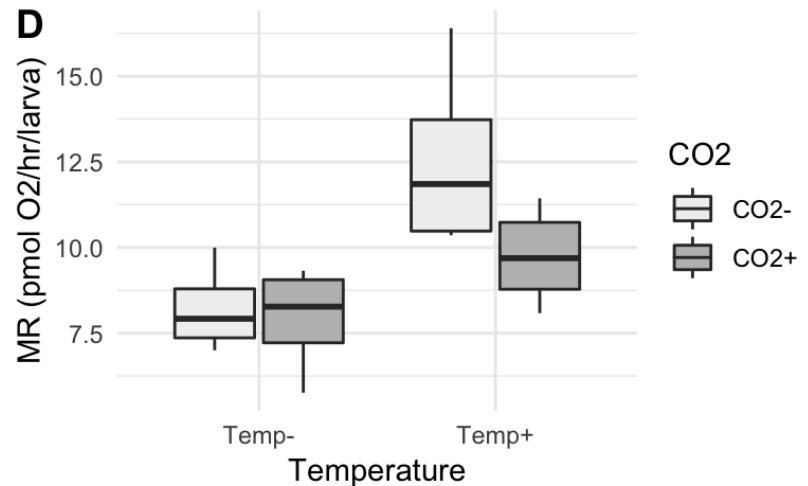
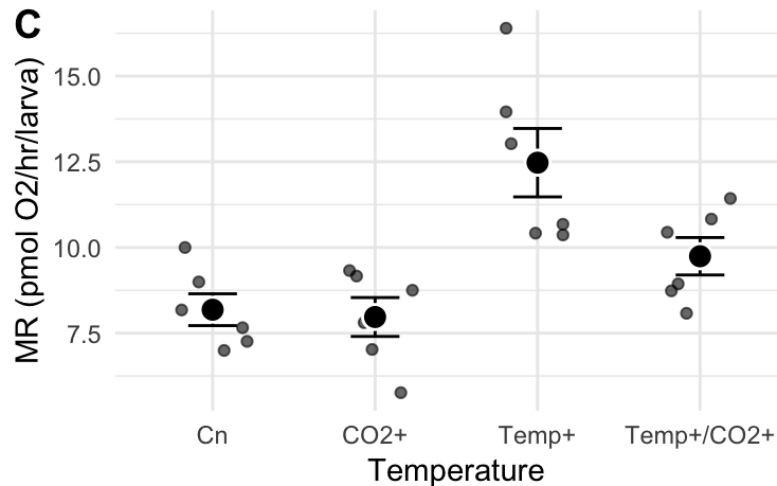
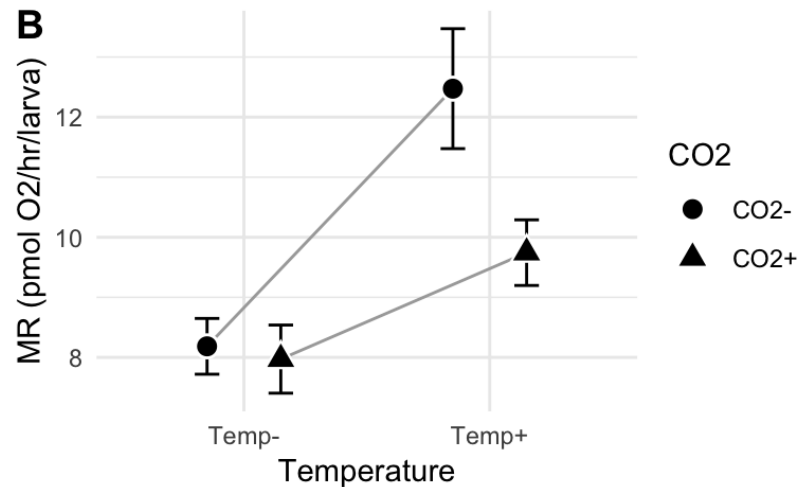
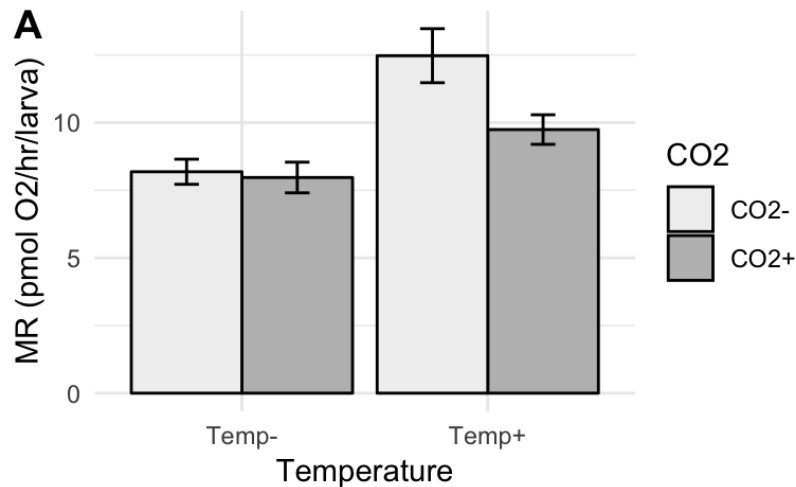




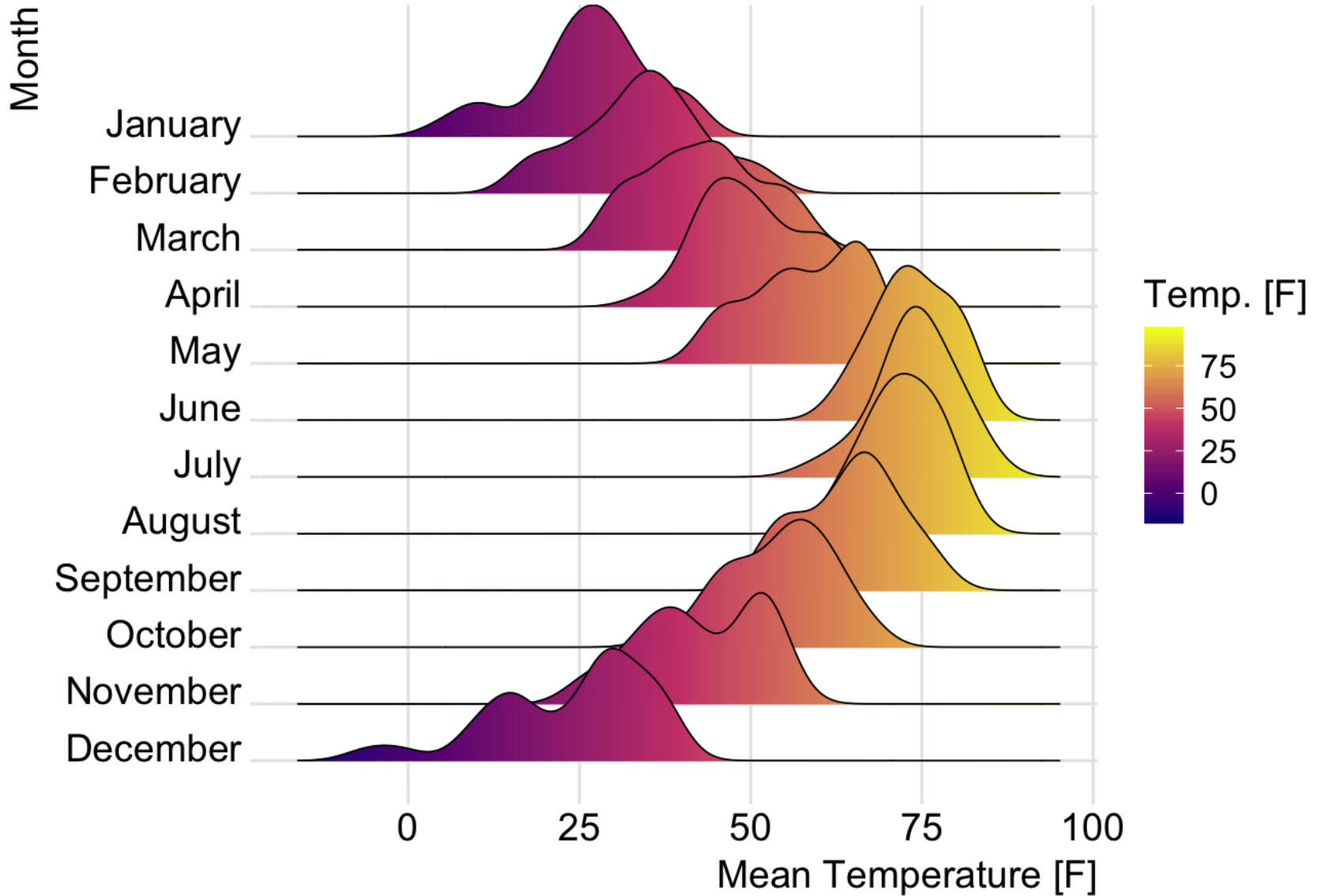
# Harrell Plot

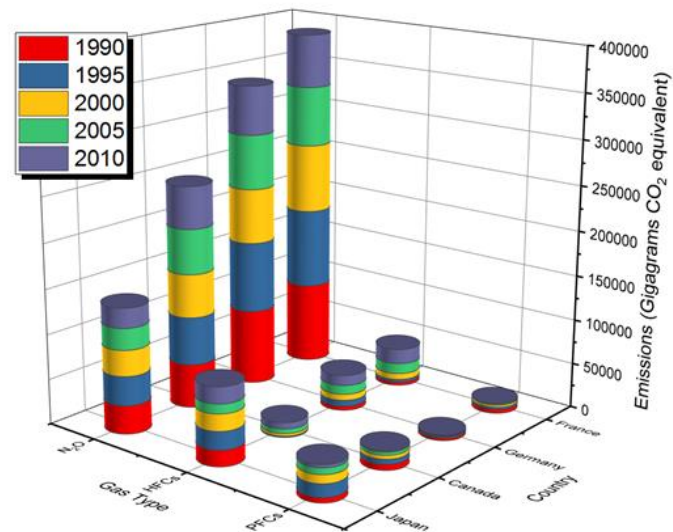
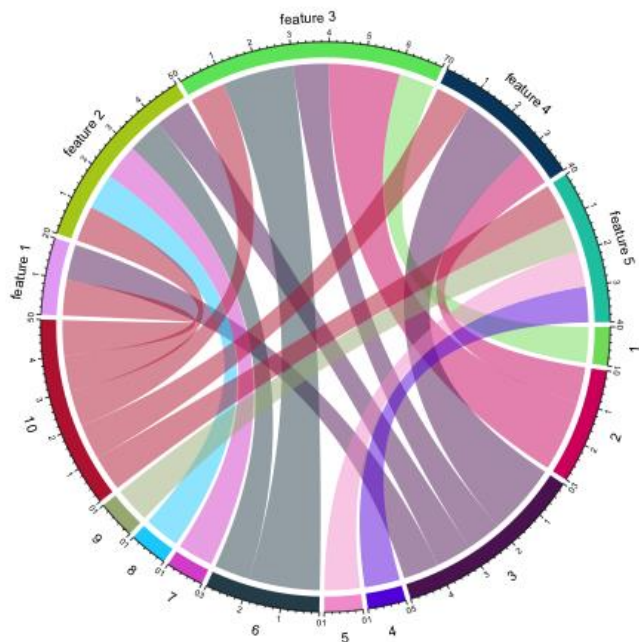


# Harrell plots vs. common alternatives



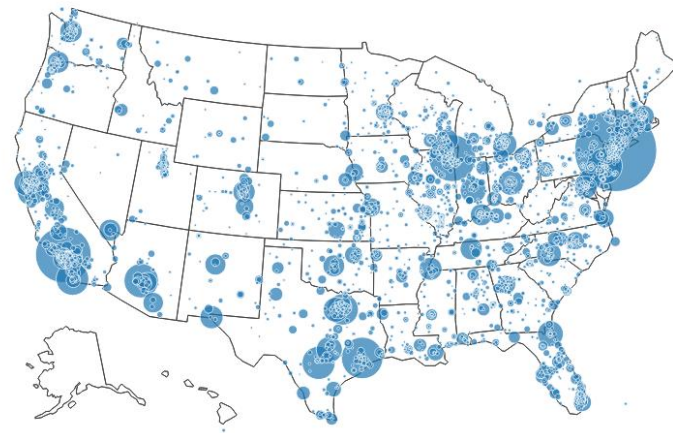
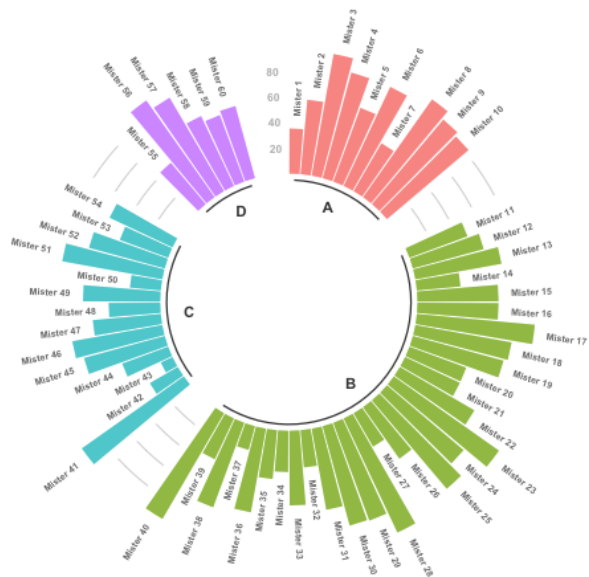
## Temperatures in Lincoln NE





Source - <http://data.un.org>

2014 U.S. City Populations



# Tufte's Views on Graphical Excellence

“Excellence in statistical graphics consists of complex ideas communicated with clarity, precision, and efficiency.

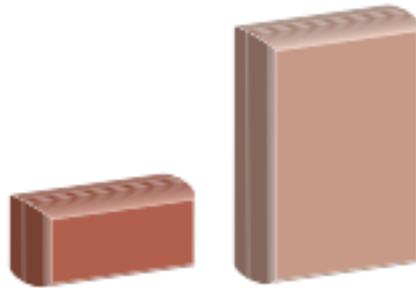
- Graphical displays should
  - Show the data
  - Induce the viewer to think about the substance rather than about methodology, graphic design, the technology of graphic production, or something else
  - Avoid distorting what the data have to say
  - Present many numbers in a small space
  - Make large data sets coherent
  - Encourage the eye to compare different pieces of data
  - Reveal the data at several levels of detail, from a broad overview to the fine structure
  - Serve a reasonably clear purpose: description, exploration, tabulation, or decoration
  - Be closely integrated with the statistical and verbal descriptions of a data set.

# Variable

## Categorical

### Nominal

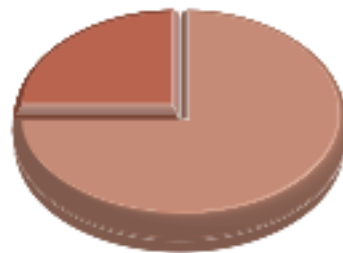
Ex. gender



Descriptives:  
absolute and  
relative frequencies,  
mode

### Ordinal

Ex. educational  
level, NYHA  
classification

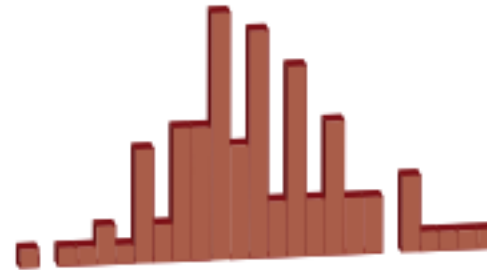


Descriptives:  
absolute and  
relative frequencies,  
median,  
interquartile range,  
minimum,  
maximum,  
percentile

## Numerical

### Continuous

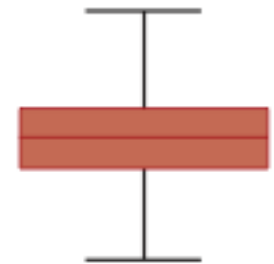
Ex. weight, systolic  
blood pressure



Descriptives:  
mean,  
median,  
SD,  
minimum,  
maximum,  
percentile

### Discrete

Ex. number of visits



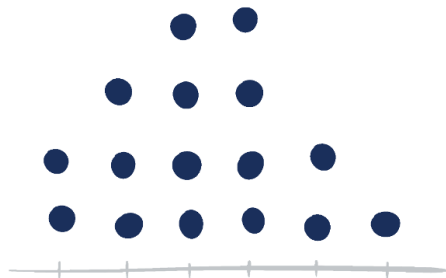
# Banning bar graphs



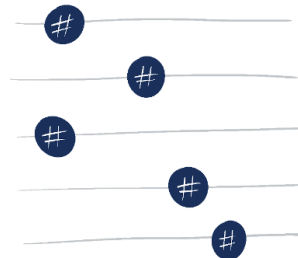
Frank Harrell  
@f2harrell

Banning bar graphs - a great idea. High ink:information ratio, poor perception, hard-to-read labels, optical distortion (humans add part of error bar to main value), hard to show 2-sided uncertainty intervals. Replace w/Bill Cleveland dot charts.

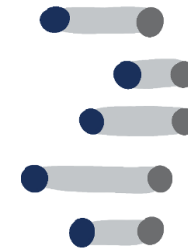
Don't let variability, outliers & skewness hide: ban bar graphs



DOT PLOT



CLEVELAND DOT PLOT

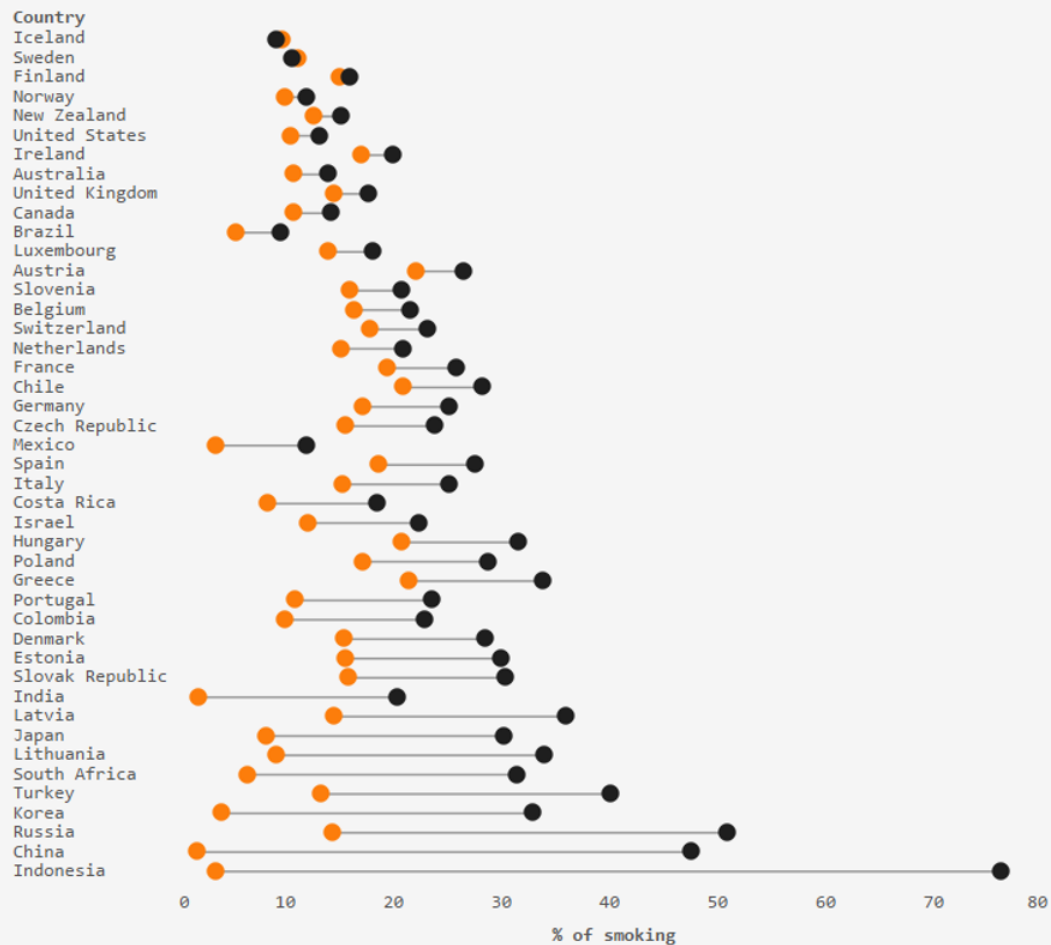


CONNECTED DOT PLOT

# In Iceland and Sweden there are more smoking women than men aged 15+

Sort by  
Difference between smoking males and females aged 15+

Sort Order  
Descending





# Descriptive statistics / Describing data

- To look slowly the data
- To discover trends
- To show relationships
- To what ?
  - To find the best approach to show it
  - To increase the comprehension of the data
  - To give utility to the data