

A PICTURE IS WORTH

THOUSAND WORDS

S U A L S A

THE GREATEST VALUE OF

A PICTURE IS WHEN IT

FORCES US TO NOTICE WHAT

WE NEVER EXPECTED TO

SEE

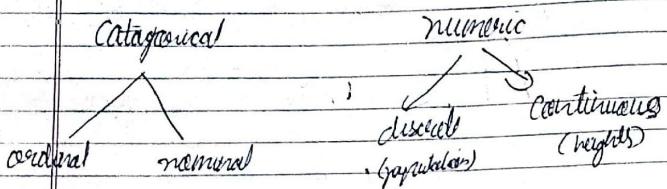
JOHN W. TUKEY

(FATHER OF EDA)

## SECTION-1

### INTRO TO DATA VISUALIZATION AND DISTRIBUTION

$$x + \text{SD} \rightarrow \text{std. dev}$$



- names (heights). # to put out column names
- \* discrete numeric data can be considered as ordinal.
- unique (?) # to find and unique values in column
- tabb (x) # nos of times a variable is occurring in table  
variable  $\rightarrow x$   
nos of rows  $\rightarrow n$

# new hypothetical use case :- Describe heights to ET  
Collectiveensus (Alice)

\* To find proportions of male & female  
prop. table (table (heights & sex))  
♀ female      male  
0.227      0.773

\* cumulative distribution function (CDF)  
 $F(a) = P(x \leq a)$   
proportion/probability

$\Rightarrow$  frequency table  $\Rightarrow$  for categorical data  
 $\Rightarrow$  CDF  $\Rightarrow$  for numeric data  
# ECDF  $\rightarrow$  Empirical CDF

\* Histograms are preferred over CDF plots

$\Rightarrow$  CDF is essential for calculating probability related to continuous data.

\* to represent large data on whole page  
using CDF is not a good idea, we can use CDF to obtain summary, such as probability that a student is b/w 69.5 & 69.5 inches

• 67.5

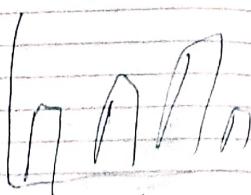
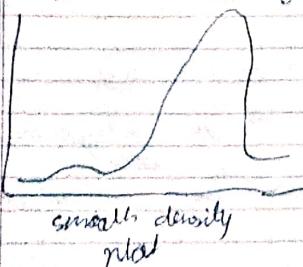
`a <- seq(min(my-data), max(my-data), length = 100)`

`cdf.function <- function(x) {  
 mean(my-data) <= x  
}`

`cdf.values <- sapply(a, cdf.function)`

`plot(a, cdf.values)`

### Smooth density plot



⇒ degree of smoothness should be taken care of.

### Normal distribution

(gaussian distribution)

(bell curve)

- for any interval  $a \& b$ , the proportion of values in the interval, ~~is~~ using the formula

$$P_a(b) = \int_a^b \frac{1}{\sqrt{2\pi}s} e^{-\frac{1}{2} \left(\frac{x-m}{s}\right)^2} dx$$

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

$m \rightarrow$  mean       $s \rightarrow$  std deviation  
 $\pi, e \rightarrow$  constants       $a, b \rightarrow$  intervals

- we have built-in functions in R to compute this formula
- If distribution is defined by just 2 parameters, i.e. mean & std deviation.

\* a low std deviation means, values are closer to mean of the set, while a high std deviation indicates that the values are spread out over a wider range.

\* `order <- heights $ sex = "Male"`  
`x <- heights $ height[order]`

`average <- mean(x)`  
`SD <- sd(x)`  
~~average~~  $\approx$  average,  $SD = SD$

~~average~~ 69.31       $SD$  3.61

### standard units (z-score)

- for data which is approximately normal, it is common to think in terms of standard units. No standard units of value tells us how many standard units away from the average this value is.

\* example

`z <- scale(x)`

⇒ To see how many nos are within 2 std deviation from the average.

$$\text{mean}(z) < z \\ 0.95$$

~ 95% (This is what a normal distribution predicts)

⇒ Z-scores are useful to quickly evaluate whether an observation is average or extreme. Z-scores near 0 are average, Z-scores below 2 are above-2 are significantly above or below the mean, and Z-scores above 3 or below -3 are extremely rare.

$$\Rightarrow \mu \pm \sigma \Rightarrow 68.3\% \text{ have } |z| \leq 1$$

$$\Rightarrow \mu \pm 2\sigma \Rightarrow 95.4\% \text{ have } |z| \leq 2$$

$$\Rightarrow \mu \pm 3\sigma \Rightarrow 99.7\% \text{ have } |z| \leq 3$$

# ~~2~~ cumulative distributions for the normal distribution

pnorm(1) # function for

$$F(a) = \text{pnorm}(a, \text{mean}, \text{sd}) \quad \# \text{ distribution}$$

$$\Rightarrow 1 - \text{pnorm}(70.5, \text{mean}(x), \text{sd}(x))$$

# how many students have height greater than

⇒ for intervals that include exactly one round number

• using data

$$\text{mean}(z \leq 68.5) - \text{mean}(z \leq 67.5) \\ [1] 0.114532$$

• using approx method

$$\text{pnorm}(68.5, \text{mean}(x), \text{sd}(x)) - \text{pnorm}(67.5, \text{mean}(x), \text{sd}(x)) \\ [1] 0.1031677$$

\* These data which don't include integers (a property where all integers) is not good for using approx (pnorm)

$$\text{mean}(z \leq 70.9) - \text{mean}(z \leq 70.1) \\ [1] 0.07716749$$

$$\text{pnorm}(70.9, \text{mean}(x), \text{sd}(x)) - \text{pnorm}(70.1, \text{mean}(x), \text{sd}(x)) \\ [1] 0.08359562$$

⇒ This situation is known as discretization

# make QQ-plot with scaled values

observed-quantiles <- quantile(z, p)

theoretical-quantiles <- qnorm(p)  
plot(theoretical-quantiles, observed-quantiles)

abline(c(0, 1))

→ for normal distributions median & mean are same

\* Exploratory data analysis (in brief)

distribution of small kid height

data visualization helps us to find pattern in our data.

→ mad() ⇒ mean absolute deviation

\* ggplot2

→ can be loaded by tidyverse & dplyr

→ works with data tables where rows are observations & columns are variables

→ grammar of graphics

\* ggplot2 consists of 3 main components

→ Data: The dataset being summarized

→ geometry: the type of plot (scatterplot, barplot, etc)

→ aesthetic mapping: variables mapped to visual cues such as position, size, color & shape

→ There are additional components

• Scale

• Labels, title, legend

• Theme / style

→ p <- ggplot(data = murders) +  
geom\_point(p)

Layers  
= nppsymbol

data : > ggplot() + layer1 + layer2 + ... + layerN  
usually  
geometry

→ geom - name of geometry and quick example

→ aes # for aesthetic mappings

→ ?geom\_point for jupyter explanation  
of population

→ scatterplot

murders : > ggplot() +  
geom\_point(aes(x = population / 10^6, y = total))

→ we can add layers to previously defined objects

p <- ggplot(data = murders)

p + geom\_point(aes(population / 10^6, total))

geom - label  $\#$  adds a small rectangle  
geom - text  $\#$  simply adds a text

p + geom\_label(aes(population / 10<sup>6</sup>, total, label = abbr))  
 $\downarrow$   
abbr  
per state

### \* tinkering

$\Rightarrow$  call size = 3 (~~biggest~~ biggest scatter plot now)  
outside aes.

$\Rightarrow$  nudge - x to move the label a little bit  
geom scatter plot (also out from aes)

$\Rightarrow$  global aes  $\#$  to give x, y only once  
in geometry, not everywhere

p <- murders %>% ggplot(aes(population / 10<sup>6</sup>,  
total, label = abbr))

$\Rightarrow$  now simplified code will be :-

p + geom\_point(size = 3) + geom\_text()

$\Rightarrow$  local ~~points~~ to outside global ones

p + geom\_point(size = 3) + geom\_text(aes(G = 1,  
y = 800, label = "Hello there!"))

### Scales, labels and colors

do do the log<sub>10</sub> transformation we can use

p + scale\_x\_continuous(trans = "log10") +  
scale\_y\_continuous(trans = "log10")

more efficient

p + scale\_x\_log10() + scale\_y\_log10()

$\Rightarrow$  to add labels

p + xlab("Population in millions(log scale)") +  
ylab("Total number of murders(log scale)") +  
ggtitle("US gun numbers as US 2010")

$\Rightarrow$  avg murders rate

sr <- murders %>% summarise(rate = sum(total)/  
sum(population) \* 10<sup>6</sup>) %>% pull(rate)

$\Rightarrow$  do change colours according to expenses

p + geom\_point(aes(col = region), size = 3) +  
geom\_text()

$\Rightarrow$  too adding a st line

p + geom\_abline(intercept = log(810(9)))

# slope is about -1

⇒ to call line before prints & changing line  
to draw lines & do change colors  
to darkgray

$P \leftarrow p + geom\_abline (intercept = \log(10(3)), slope = "darkgray") + geom_raster (aes("col = my))$

# we have called line first, then prints (get it?)



⇒ to calculate no ~~the~~ legend

$PL \leftarrow p + scale_color_discrete (name = "Region")$

⇒ add ~~on~~ packages (for more functionality in ggplot)

⇒ library (gridExtra) # to change flow

$p + theme_economist()$

⇒ to make geometry as such that tables don't overlap each other

library (gridExtra)

~~p + gridExtra::grid.arrange()~~

$p + geom_text_repel()$

other examples

see section 2 → customizing plots → other examples

\* display

⇒ summarize is group-by  
⇒ datplaceholder  
⇒ arrange

⇒ summarizing ~~between~~ dataframes

$S \leftarrow heights \%>\%> filters (scr = "Male") \%>\%> summarizing (average = mean (height), standard deviation = sd (height))$

⇒ datplaceholder functions numeric

US\_murder\_rate  $\leftarrow$  murders \%>\%>  
+ summarizing (rate = sum (total) / sum (population)  
\* 100 000) \%>\%> \$ rate

→ this plot is datplaceholder

⇒ heights \%>\%> group\_by (scr) \%>\%> summarizing  
(average = mean (height), standard deviation = sd (height))

scr	average	standard deviation
Female	64.9	3.76
Male	69.3	3.61

⇒ arrange has ordering

murders %>% arrange(murder\_rate) %>% head()

# ascending order

murders %>% arrange(desc(murder\_rate)) %>% head()

# descending order

⇒ nested sorting (for no use)

murders %>% arrange(explan, murder\_rate) %>% head()

sorted by regions, thus  
lowest murder  
rate

⇒ murders %>% top\_n(10, murder\_rate)

# top 10 nos of exes, but not ordered

⇒ murders %>% arrange(desc(murder\_rate)) %>% top\_n(10)

\* groupby

disclaimer  
library(gapminder)  
data(gapminder)  
head(gapminder)

⇒ country year, w/ modality, by expectancy, fertility,  
population, gdp, continent, region

• gapminder %>% filter(year == 2015) is country  
%>% c("Sri Lanka", "Turkey") %>% select  
(country, region, modality)

⇒ scatterplot of life expectancy versus fertility.

ds - thor - set()

filter(gapminder, year == 1962) %>%

ggplot(aes(fertility, life\_expectancy)) + geom\_point()

color = continent

\* faceting # side by side plots

facet\_grid() # purchase per sf

gapminder %>% filter(gapminder\$year %in% 1962, 2012) %>% + ggplot(aes(fertility,

life\_expectancy, color = continent)) + geom\_point()

+ facet\_grid(continent ~ years)

↓ rows ↓ columns

\* facet-grid (. ~ year)

⇒ facet-wrap() # automatically wraps by  
 scales of labels, so that wrap  
 displays has wrapped dimensions

• changes in necessary code

> years <- c(1962, 1980, 1990, 2000, 2010)

> continents <- c("Europe", "Asia")

> facet-wrap (~ year)

> filter(year %in% years & continent %in% continents)

\* time series plots (time at x-axis)

⇒ single time series

gapminder %>% filter(country == "United States") %>%  
ggplot(aes(x = year, fertility)) + geom\_line()

⇒ multiple time series

countries <- c("South Korea", "Germany")

gapminder %>% filter(country %in% countries) %>%  
ggplot(aes(x = year, fertility, col = country)) +  
geom\_line()

⇒ adding facet labels

labels <- doto.frame(Tcountry = countries, x = c(1975, 1985), y =

gapminder %>% filter(country %in% countries) %>%

ggplot(aes(x = year, life\_expectancy, col = country)) + geom\_line()

+ geom\_text(data = labels, aes(x, y, label = country), size = 5)

\* T transformations (useful for better understanding  
of distribution)

gapminder <- gapminder %>% mutate(dollars\_per\_day = gdp/population / 365)

⇒ log transformations

• using base 2 for example means that every  
time a value doubles.. log transformation  
increases by one



# two modes of income  
(bimodal distribution)

⇒ in normal distribution, mode = average

• use log base 2, log base 10 for EDA

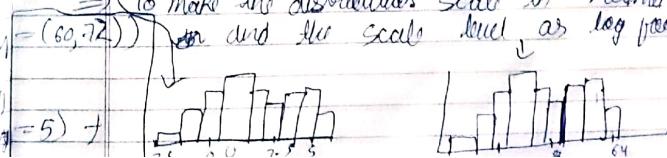
past\_year <- 1970

gapminder %>%

filter(year == past\_year & !is.na(gdp)) %>%

ggplot(aes(log2(dollars\_per\_day))) +  
geom\_histogram(binwidth = 1, color = "black")

⇒ To make the distribution scale in normal form  
on and off scales need as log form



gapminder %>%

```
filter(year == past_year & !is.na(gdp)) %>%  
ggplot(aes(dollars_per_day)) +  
geom_histogram(binwidth=1, color = "black") +  
scale_x_continuous(trans = "log2")
```

( )  
# for proper normal scaling

## \* Stratify and Boxplot

P <- gapminder %>%

```
filter(year == past_year & !is.na(gdp)) %>%  
mutate(region = gapminder$region, dollars_per_day, fill = median) %>% # reorder according to median values  
ggplot(aes(region, dollars_per_day, fill = continent)) + # colors by continent  
* geom_boxplot() +  
theme(axis.text.x = element_text(angle = 90, hjust = -1)) + # rotates the labels by 90°  
xlab("") # to ignore the x axis label  
+ scale_y_continuous(trans = "log2") # to properly show the graph
```

## quantiles

`quantile(data, q)`

### Percentiles

`p <- seq(0.01, 0.99, 0.01)`  
`quantile(data, p)`

### Quartiles

divide dataset into 4 parts  
25% → 1st quartile  
50% → median  
75% → 3rd quartile

→ `summary()` returns minimum, quartiles & maximum of a vector.

### summary (heights \$height)

`p <- seq(0.01, 0.99, 0.01)`  
`percentiles <- quantile(heights $height, p)`

→ 25% & 75% percentiles match 1st & 3rd quartile to verify it

`percentiles[names(Percentiles)] == "25%"`  
`percentiles[names(Percentiles)] == "75%"`

→ `qnorm()` function gives the theoretical value of quantile with probability  $p$  of observing a value equal to or less than that quantile value given a normal distribution with mean  $\rightarrow \mu$  std.dev  $\rightarrow \sigma$

`qnorm`: `qnorm(p, mu, sigma)`  
`qnorm(p)`

\* quantiles are defined as such that  $p$  is the probability of random observation less than or equal to the quantile

\* Relations of `pnorm()` & `qnorm()`

`pnorm(-1.96) ≈ 0.025`

→ (result of `pnorm()` is quantile)  
`qnorm(0.025) ≈ -1.96`

→ `qnorm()` & `pnorm()` are inverse functions

`pnorm(qnorm(0.025)) = 0.025`

### # Quantile - Q-Q plots

`male <- heights $height[sex == "Male"]`

`dc <- heights $height [index]`

`z <- scale(dc)`

# proportion of data below 69.5  
`mean(z <= 69.5)`

# calculating observed & theoretical quantiles  
`p <- seq(0.05, 0.95, 0.05)`

`observed-quantiles <- quantile(x, p)`  
`theoretical-quantiles <- qnorm(p, mean = mean(x),  
sd = sd(x))`

# make Q-Q plot  
`plot(theoretical-quantiles, observed-quantiles)`