

# VISUAL DATA ANALYSIS



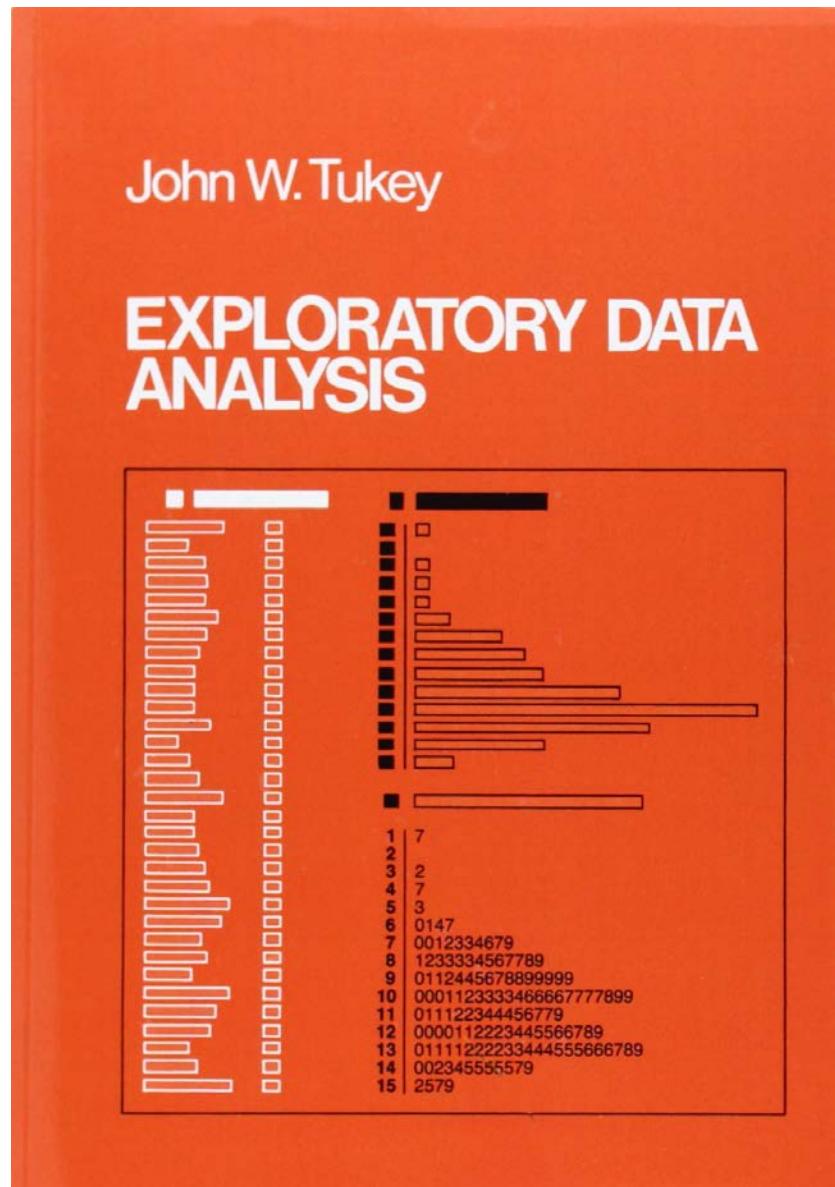
UNIVERSITY OF  
**CALGARY**

# **“EXPLORATORY DATA ANALYSIS”**



**JOHN TUKEY**

(IN CONTRAST TO “CONFIRMATORY” DATA ANALYSIS)



Based on insights developed at Bell Labs in the 60's

Introduced a number of novel techniques for **visualizing** and **summarizing** data:

- 5-number summary
- Box plots
- Stem and leaf diagrams

# EXPLORATORY ANALYSIS IS ABOUT UNDERSTANDING DATA AND CHECKING ASSUMPTIONS

- IS THE DATA CORRECT?
- DOES IT MATCH OUR PREVIOUS EXPECTATIONS?
- IS THERE A RELATIONSHIP?
  - A CORRELATION?
  - A TREND?
  - ETC.?



## E.D.A. CIRCA ~1970

- Mostly done by hand  
(computation is expensive and inaccessible)
- Simple statistical summaries and charts



# TUKEY'S 5-NUMBER SUMMARY

The sample minimum (smallest observation)

The lower quartile

The median (middle value)

The upper quartile

The sample maximum (largest observation)



# STEM-AND-LEAF PLOTS

## Volcano heights:

900 feet  
1957 feet  
823 feet  
2620 feet  
19300 feet  
730 feet  
1753 feet  
603 feet  
2930 feet  
12400 feet  
650 feet  
3663 feet

0 | 9 = 900 feet

Stem-and-leaf displays:  
heights of 218 volcanoes, unit 100 feet.

19 | 3 = 19,300 feet

0   98766562
1   97719630
2   69987766544422211009850
3   876655412099551426
4   9998844331929433361107
5   97666666554422210097731
6   898665441077761065
7   98855431100652108073
8   653322122937
9   377655421000493
10   0984433165212
11   4963201631
12   45421164
13   47830
14   00
15   676
16   52
17   92
18   5
19   39730

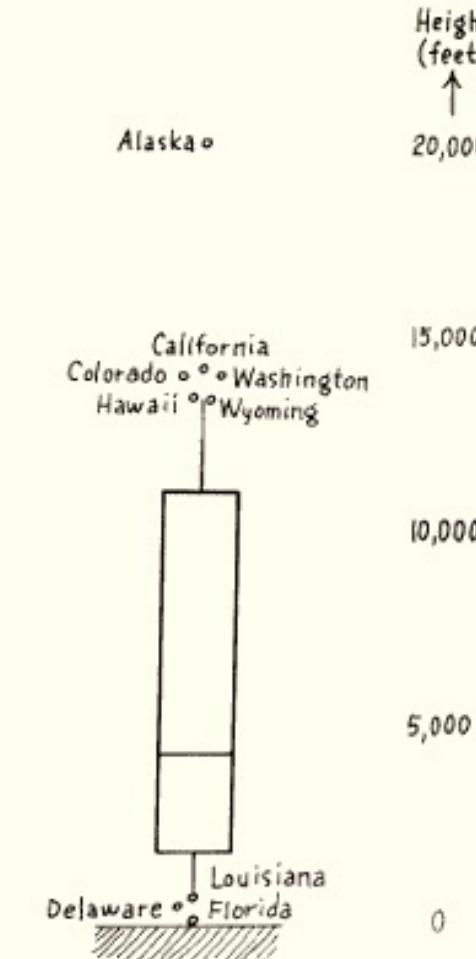


# BOX PLOTS

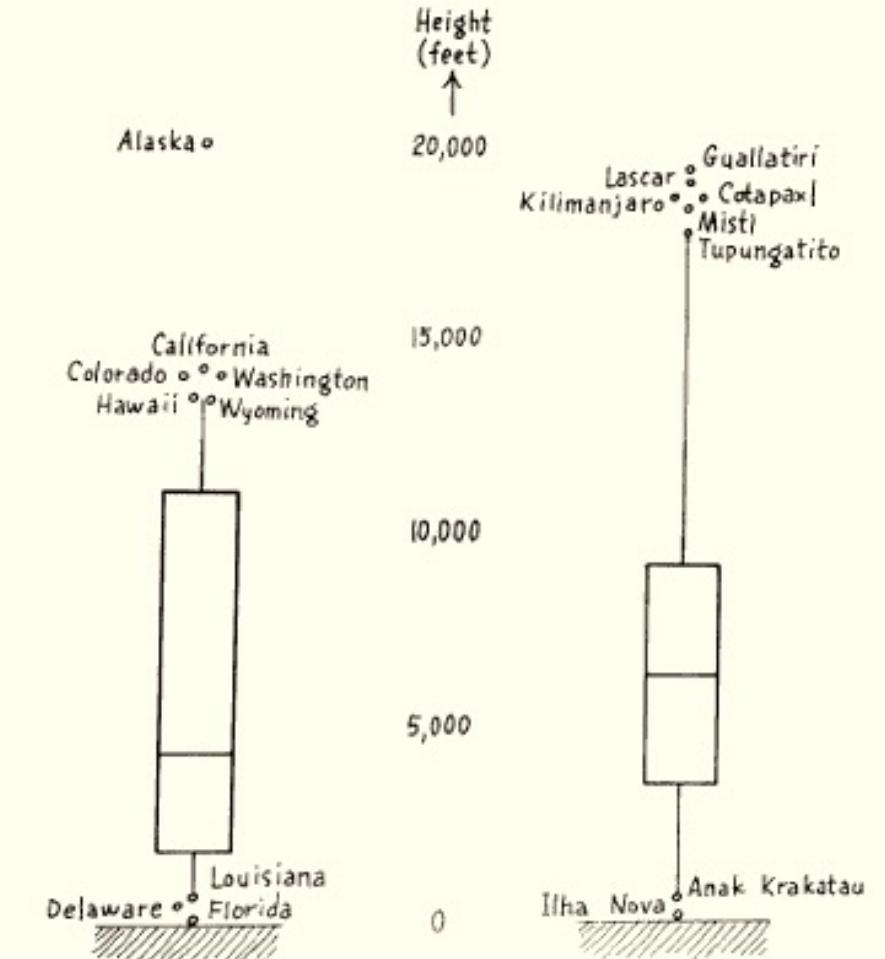
exhibit 6 of chapter 2: various heights

Box-and-whisker plots with end values identified

A) HEIGHTS of 50 STATES



B) HEIGHTS of 219 VOLCANOS



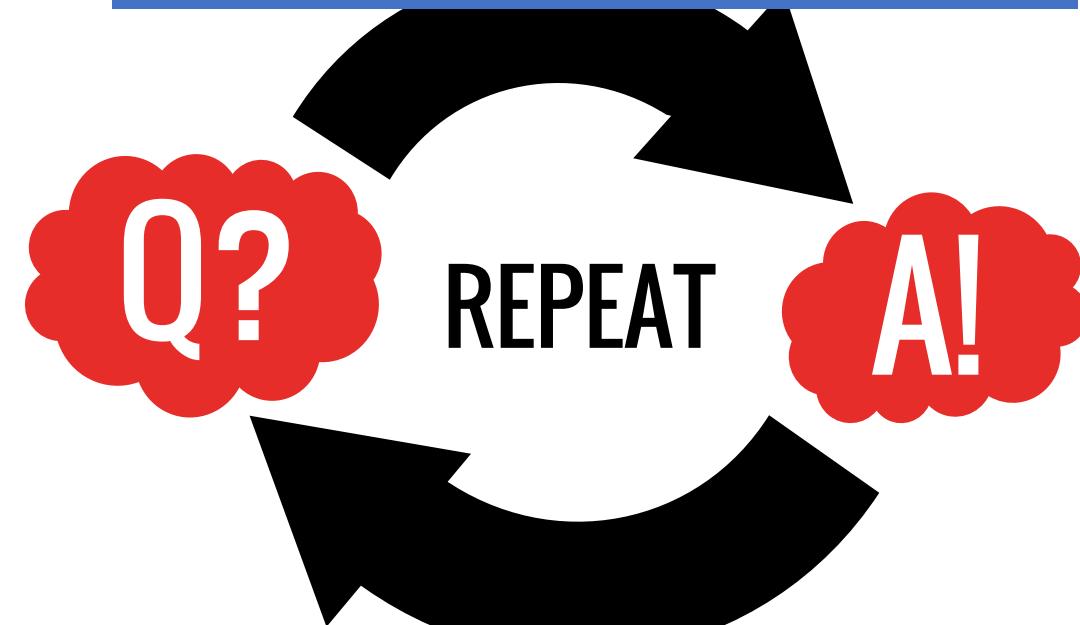
# EXPLORATORY ANALYSIS IS ABOUT UNDERSTANDING DATA AND CHECKING ASSUMPTIONS

- IS THE DATA CORRECT?
- DOES IT MATCH OUR PREVIOUS EXPECTATIONS?
- IS THERE A RELATIONSHIP?  
A CORRELATION?  
A TREND?  
ETC.?

BUT, HOW SHOULD WE GO  
ABOUT DOING THIS?

# ANALYSIS IS A CYCLE

GATHERING DATA,  
APPLYING STATISTICAL TOOLS, AND  
CONSTRUCTING GRAPHICS TO  
ADDRESS QUESTIONS



INSPECT “ANSWERS” AND  
ASSESS NEW QUESTIONS

# START SIMPLE

IT'S EASY TO GET SIDETRACKED TRYING TO DO  
COMPLICATED ANALYSES AND MISS THE BASIC STUFF



# SOME FIRST STEPS TO START WITH

1. Plot the raw data
2. Plot simple statistics
3. Examine distributions
4. Look at attributes together

DON'T TRY TO CREATE A WHOLE  
NEW CHART ALL AT ONCE!  
CHECK YOUR LOGIC AT EVERY STEP.

**LOOKING AT DATA WITH  
“THE PAINTER’S EYE”**



J. BERTIN

**EMBRACING  
“SLOW DATA”**



STEPHEN FEW

# PLOT THE RAW DATA

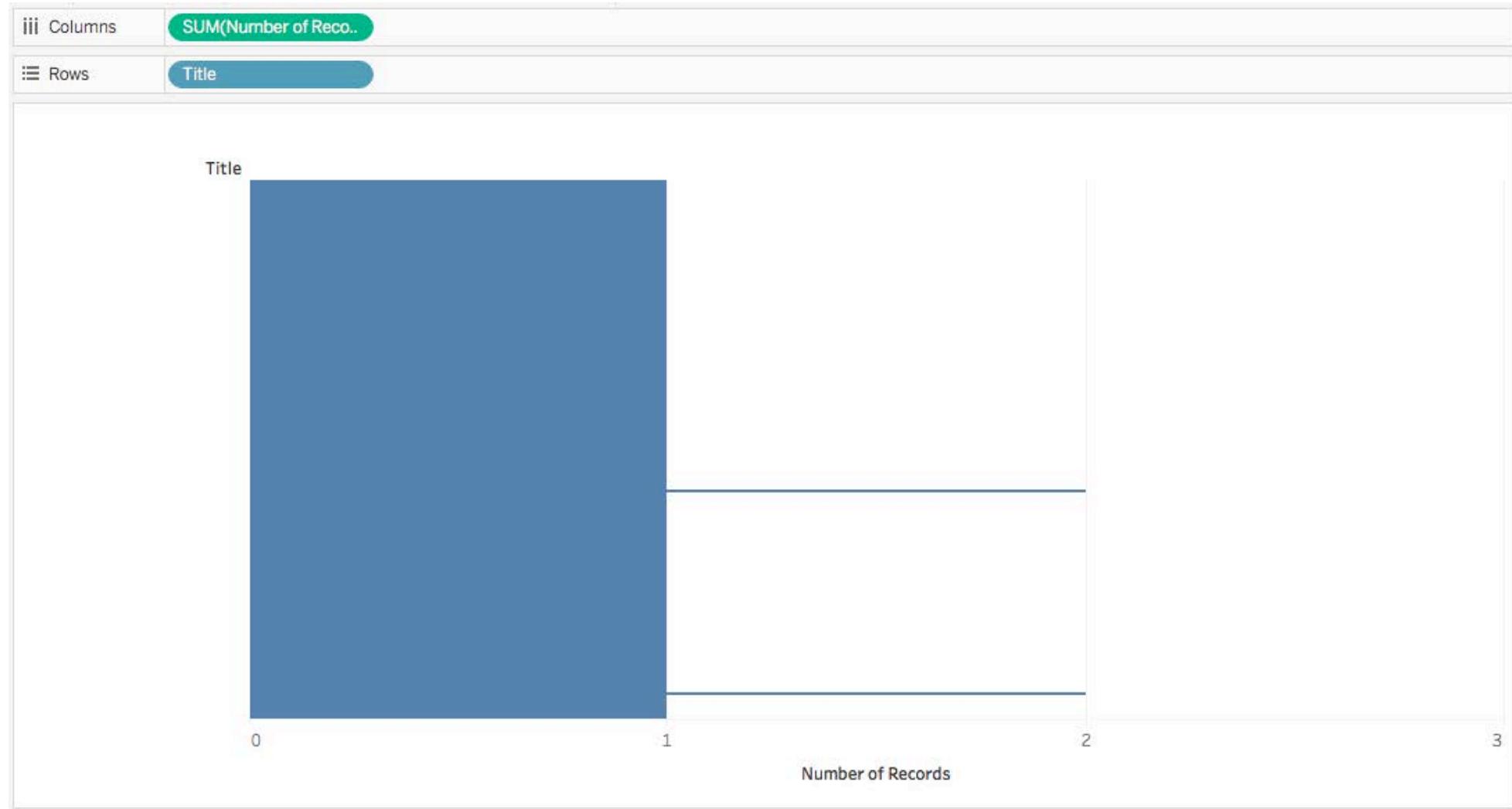
ARE THE FIELDS CORRECT?

# movies.csv Movie Id	Abc movies.csv Title	E Genres	# movies.csv User Id	# ratings.csv movielId (ratings.c... Rating	# ratings.csv Timestamp	# Calculation Year
3	Grumpier Old Men (1995)	Comedy Romance	2	1	5.00000	859,046,895 1995.00
4	Waiting to Exhale (1995)	Comedy Drama Rom...	16	2	3.00000	849,188,326 1995.00
5	Father of the Bride Part II (1995)	Comedy	2	3	2.00000	859,046,959 1995.00
6	Heat (1995)	Action Crime Thriller	80	4	3.50000	1,253,152,402 1995.00
7	Sabrina (1995)	Comedy Romance	2	5	3.00000	859,046,959 1995.00
8	Tom and Huck (1995)	Adventure Children	9	6	4.00000	842,686,600 1995.00
7	Sudden Death (1995)	Action	3	7	3.00000	841,484,087 1995.00
8	GoldenEye (1995)	Action Adventure Thriller	156	8	4.00000	1,322,062,970 1995.00
9	American President, The (1995)	Comedy Drama Rom...	16	9	4.00000	841,483,689 1995.00
10	Dracula: Dead and Lo... (1995)	Comedy Horror	7	10	4.00000	840,548,213 1995.00

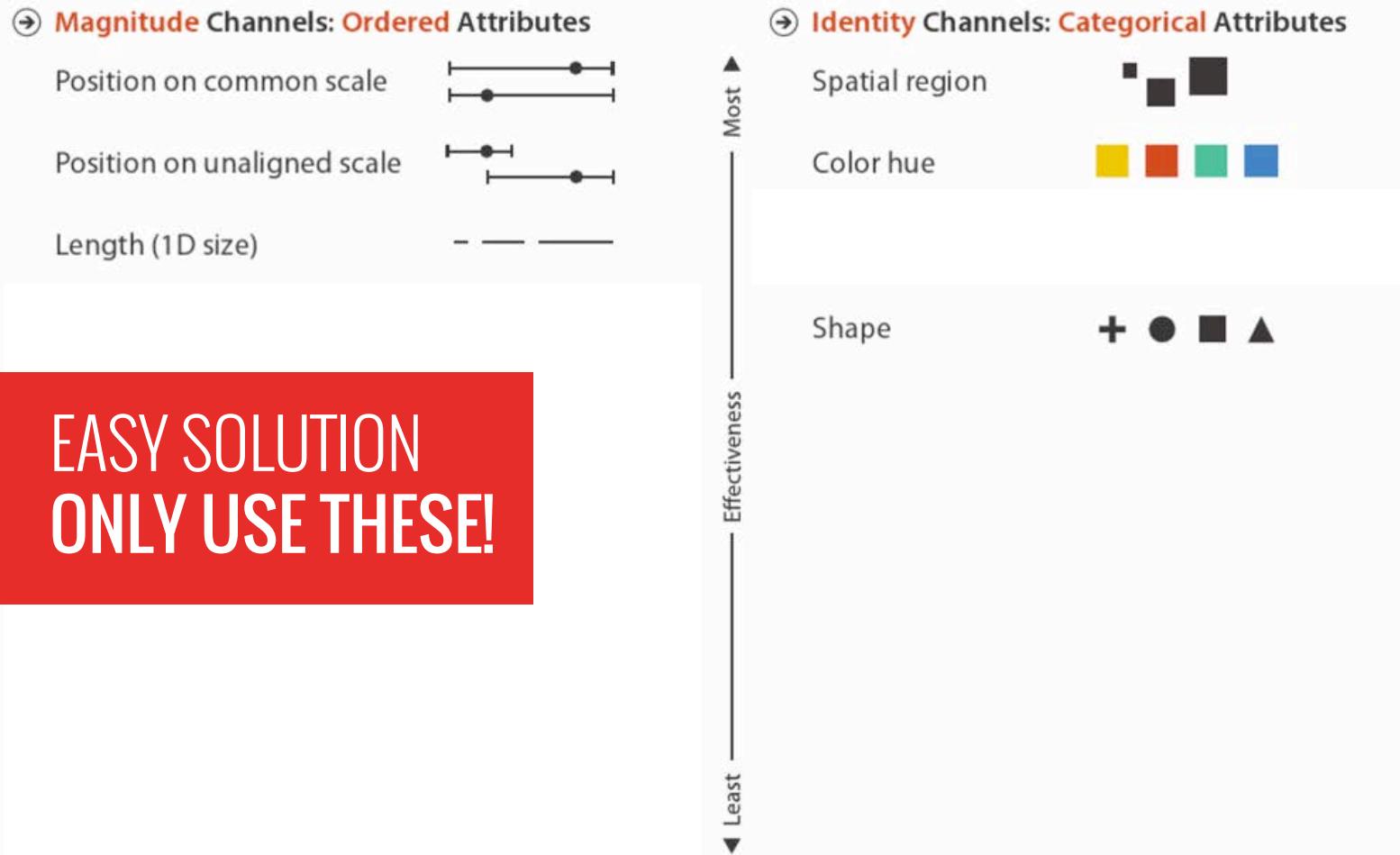
WHAT ABOUT THE DATA TYPES?

WHAT ABOUT THE VALUES?

# USE THE SIMPLEST REPRESENTATION YOU CAN TO EVALUATE ALL OF THE DATA



# CHOOSE REPRESENTATIONS THAT MAKE IT EASY TO COMPARE DIFFERENCES AND SEE PATTERNS



TAMARA MUNZNER

# DEFAULT TO SIMPLE AND EFFECTIVE CHART TYPES

*the BAR*



*the LINE*



*the SCATTER*



**+ COLOUR & SHAPE  
TO SHOW CATEGORIES**

# SOME FIRST STEPS TO START WITH

1. Plot the raw data
2. Plot simple statistics
3. Examine distributions
4. Look at attributes together

# CHECK SIMPLE STATISTICS

Measures

# Rating

CNT(Rating)

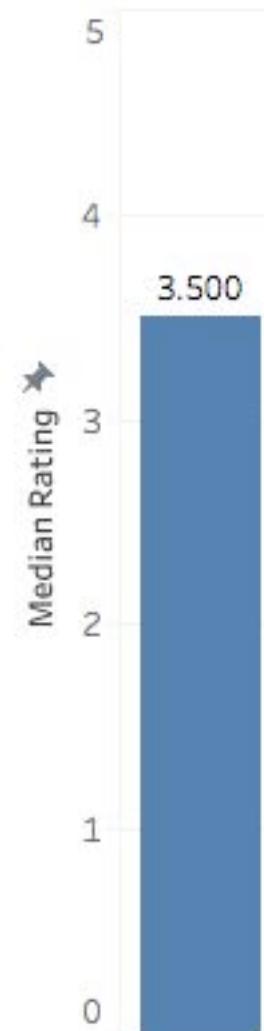
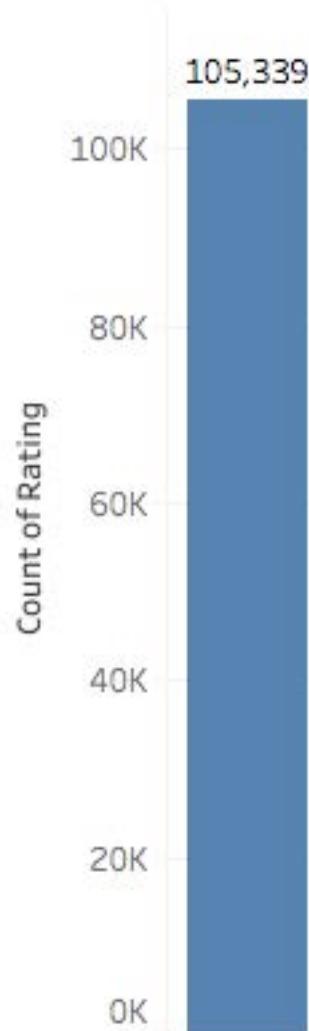
Avg(Rating)

MEDIAN(Rating)

MIN(Rating)

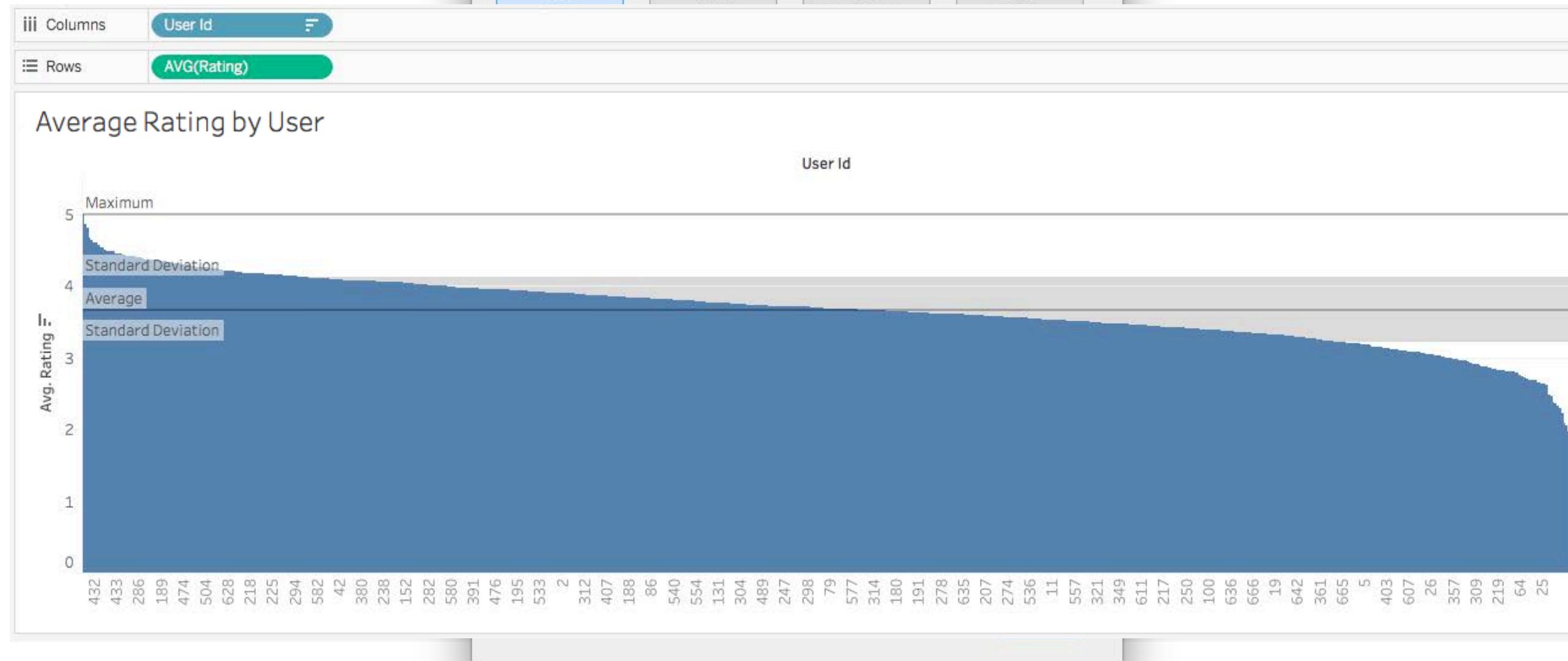
MAX(Rating)

STDEV(Rating)



WHAT'S ONE MORE EASY THING WE SHOULD DO?

CHECK SIMP





JACQUES BERTIN

***LA GRAPHIQUE ET LE  
TRAITEMENT GRAPHIQUE  
DE L'INFORMATION (1977)***

# MATRIX REORGANIZATION [BERTIN 1977]



# LA MATRICE ORDONNABLE

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1															
COLLÈGE															
COOPÉRATIVE AGRIC.															
GARE															
ÉCOLE CLASSE UNIQUE															
VÉTÉRINAIRE															
PAS DE MÉDECIN															
PAS D'ADDUCTION D'EAU															
GENDARMERIE															
REMENBREMENT															

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1															
11															
12															
13															
14															
15															
16															
17															
18															
19															

2

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1															
3															
8															
2															
5															
9															
4															
6															
7															

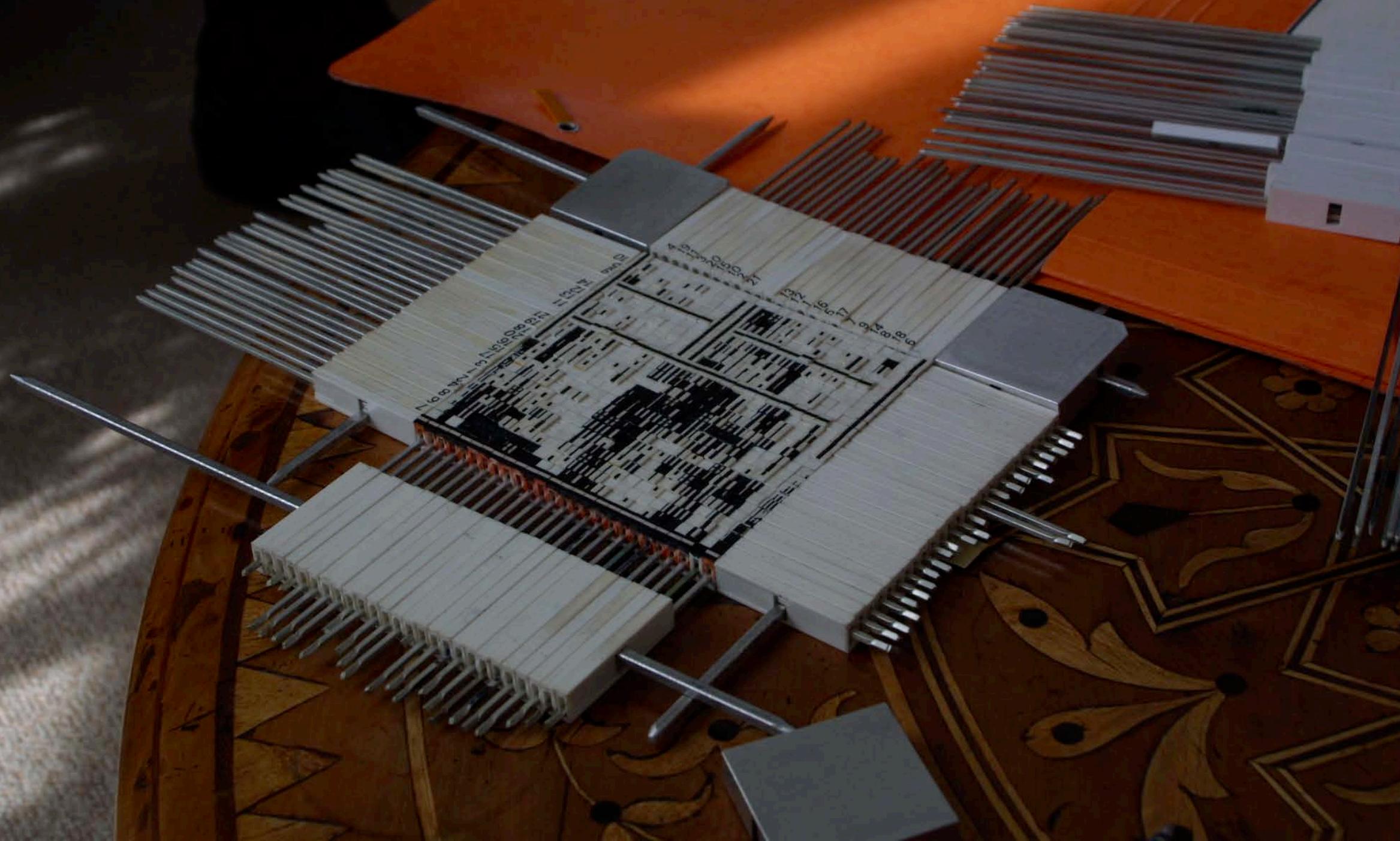
3

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1															
3															
8															
2															
5															
9															
4															
6															
7															

4

N	J	P	M	I	F	E	A	B	O	L	G	D	C	H	K
1															
3															
8															
2															
5															
9															
4															
6															
7															

5



# SOME FIRST STEPS TO START WITH

1. Plot the raw data
2. Plot simple statistics
3. Examine distributions
4. Look at attributes together

# UNDERSTANDING DISTRIBUTIONS

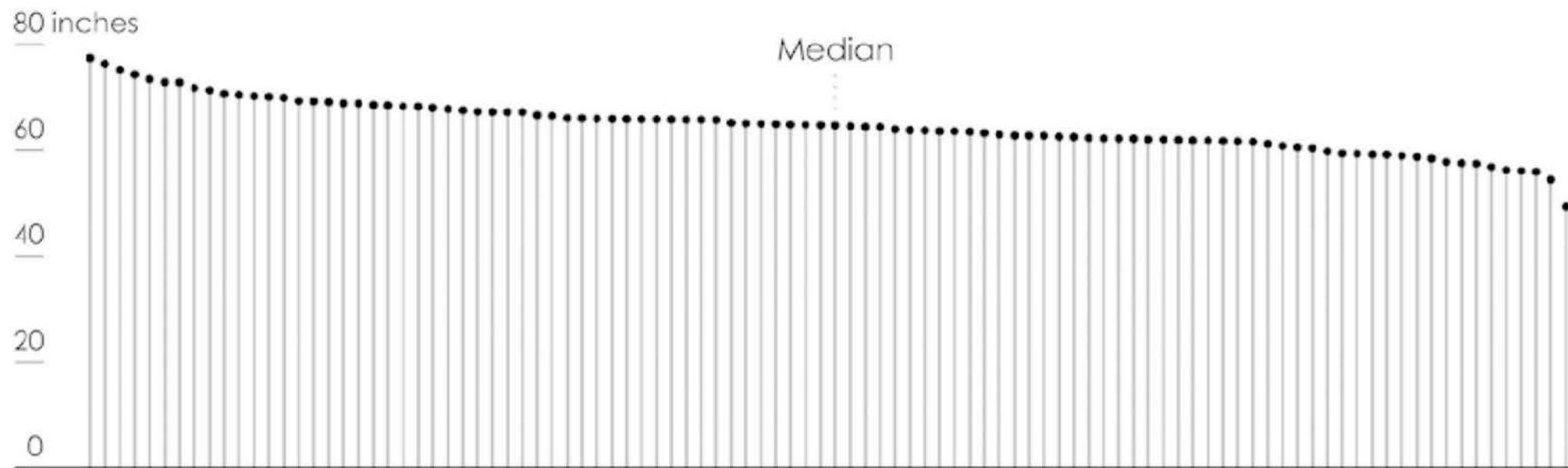
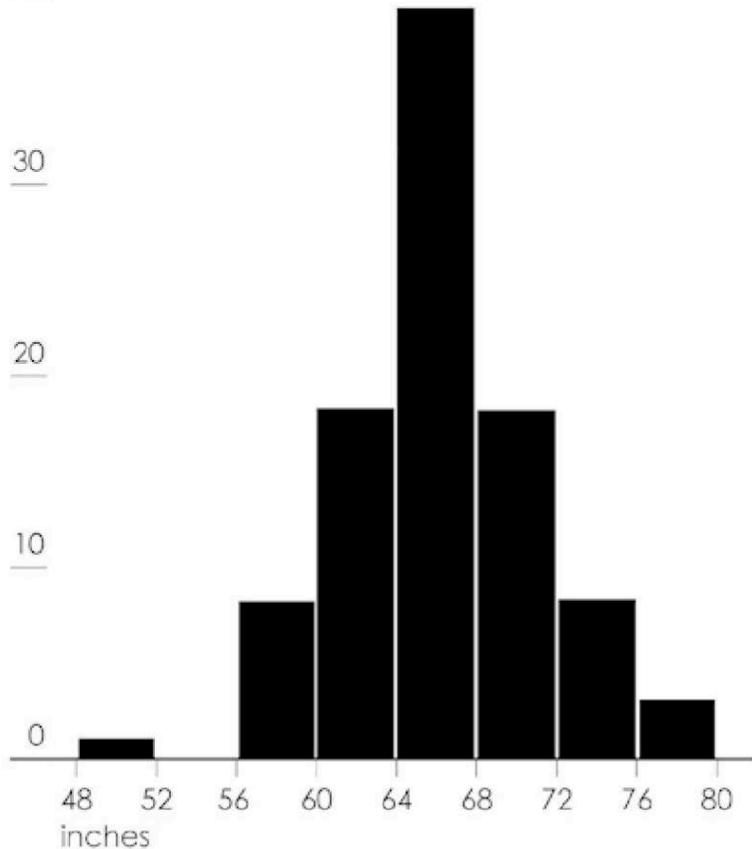
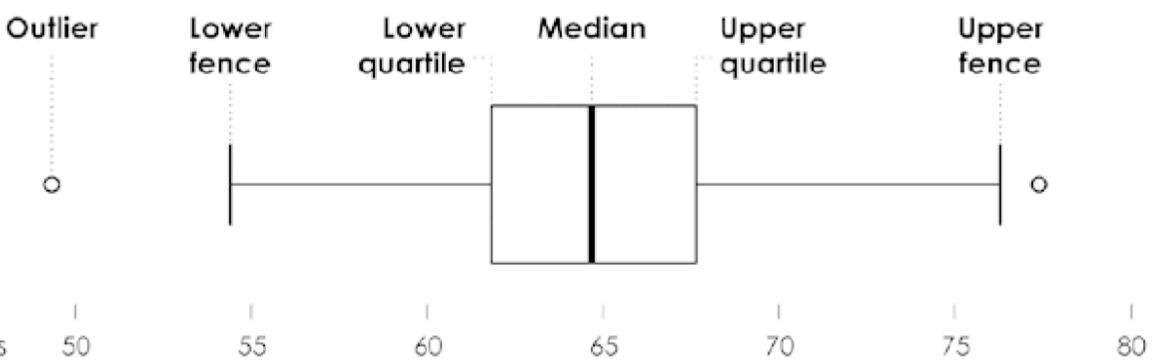


FIGURE 4-52 Heights of imaginary people, sorted from shortest to tallest

40 people



## HISTOGRAMS



## BOX PLOTS

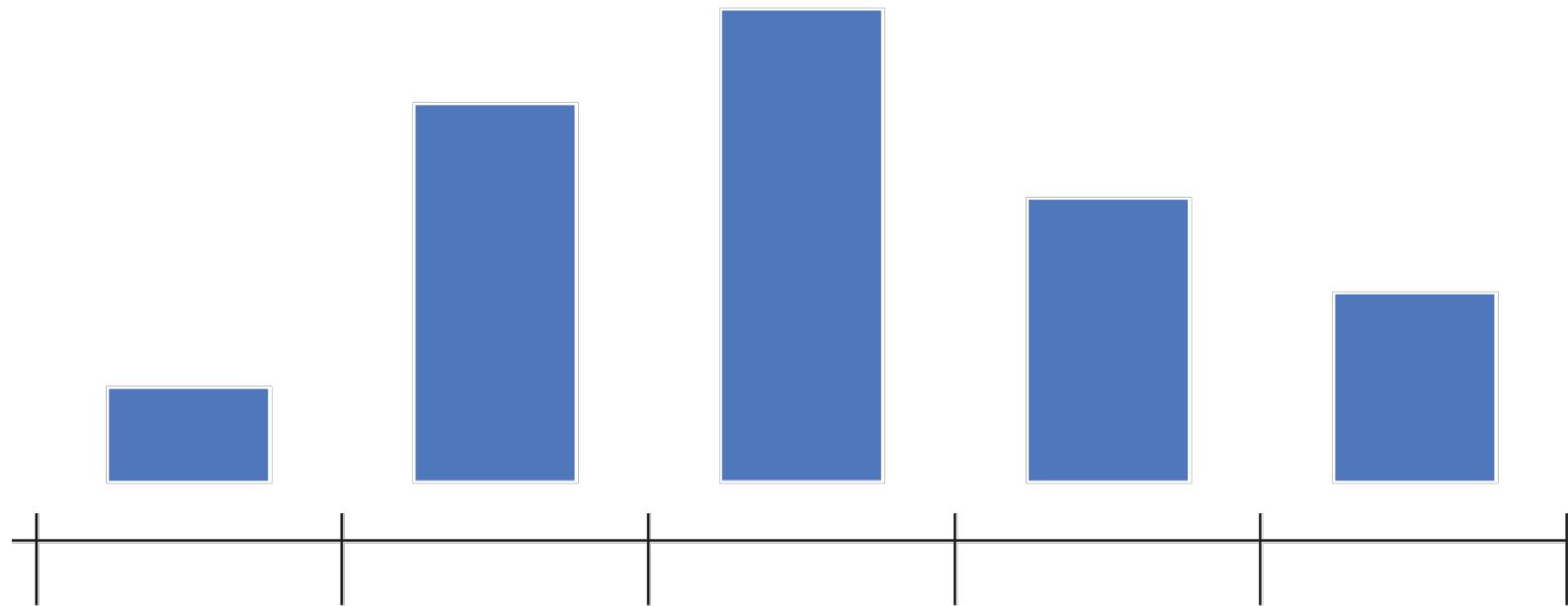
# HISTOGRAMS

---

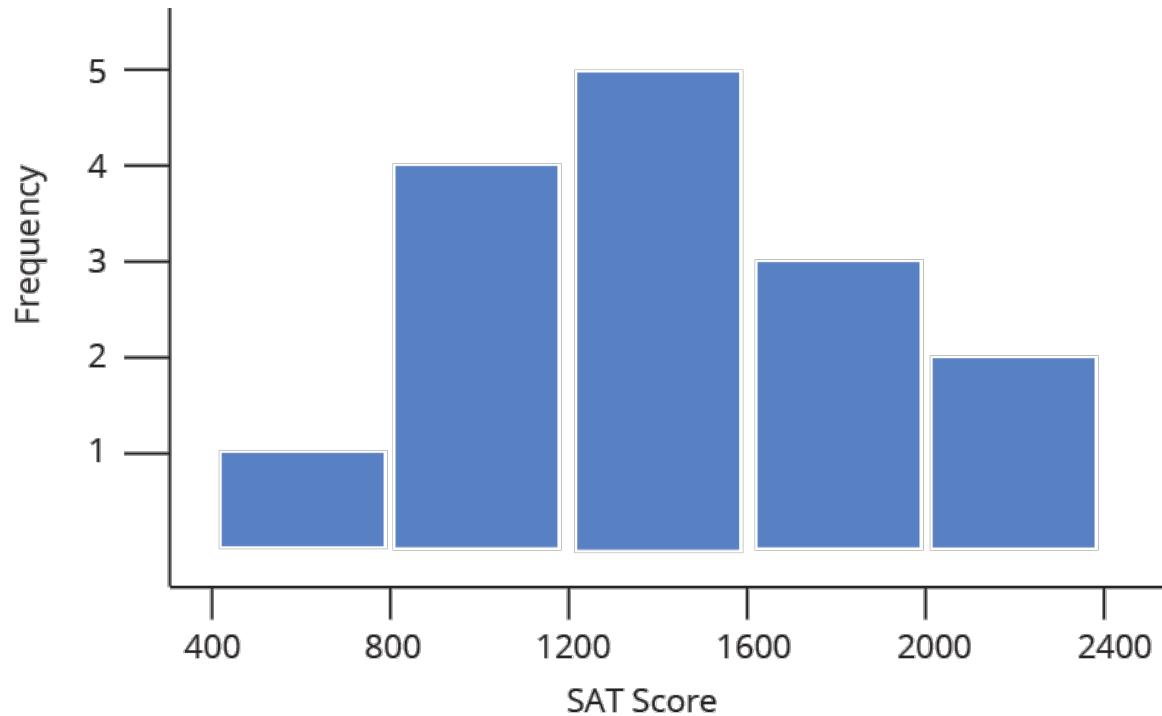
# HISTOGRAMS



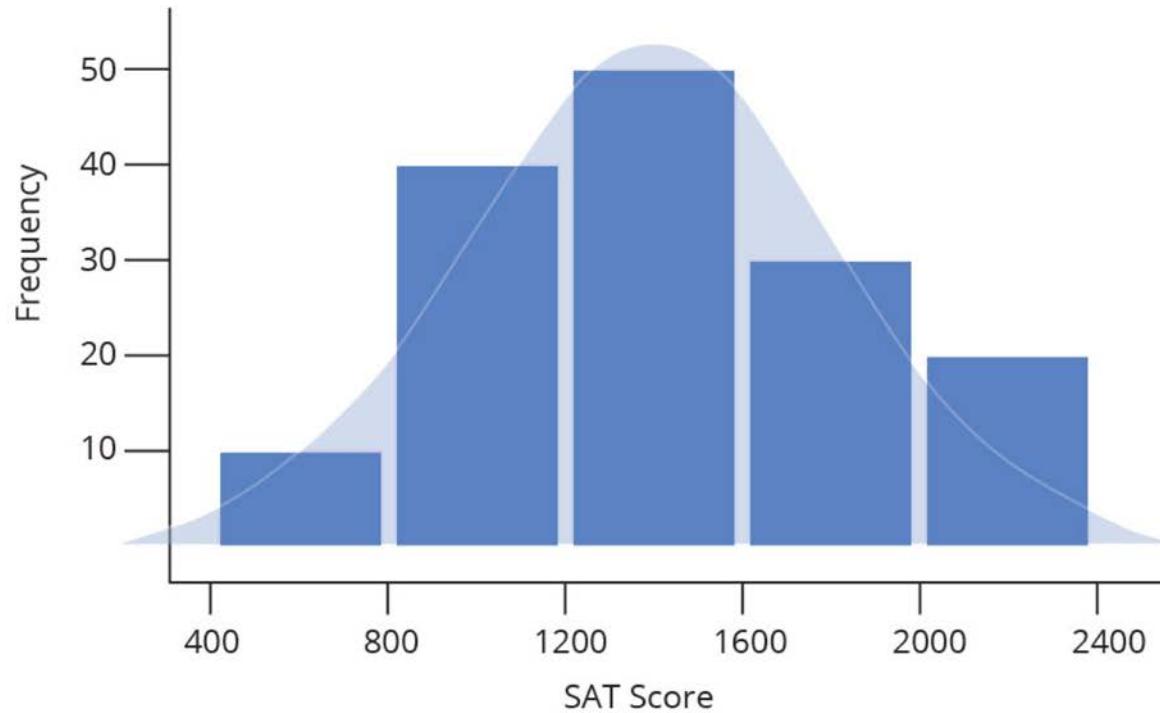
# HISTOGRAMS



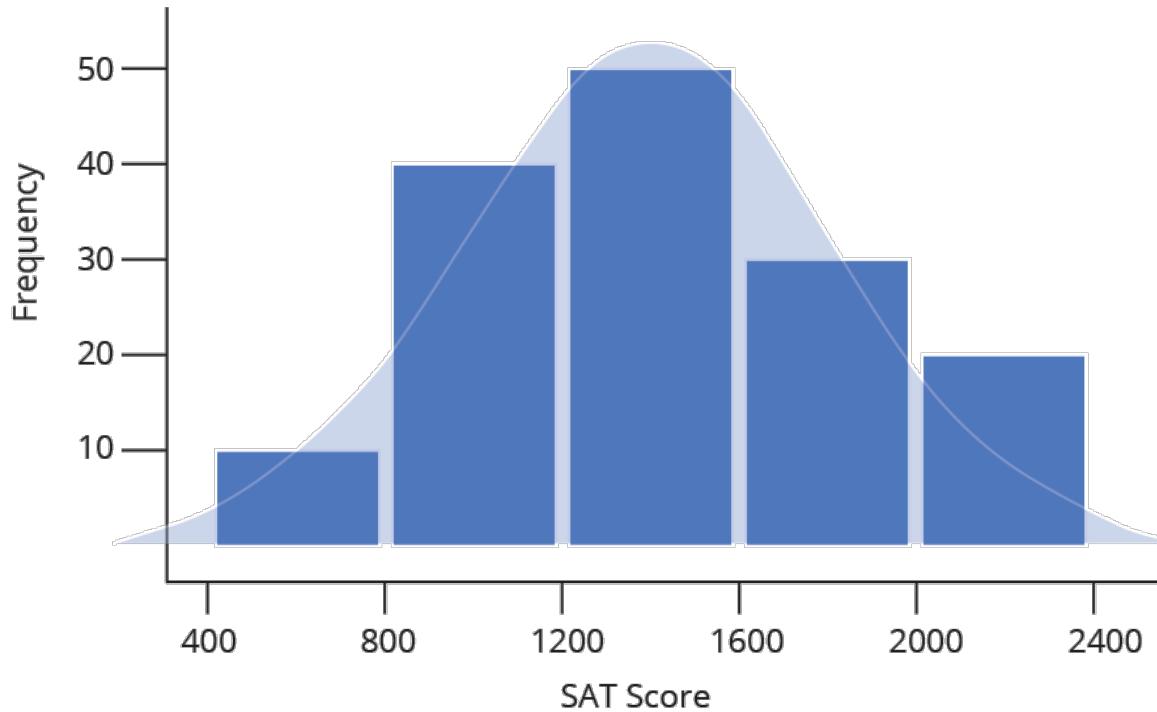
# HISTOGRAMS



# HISTOGRAMS

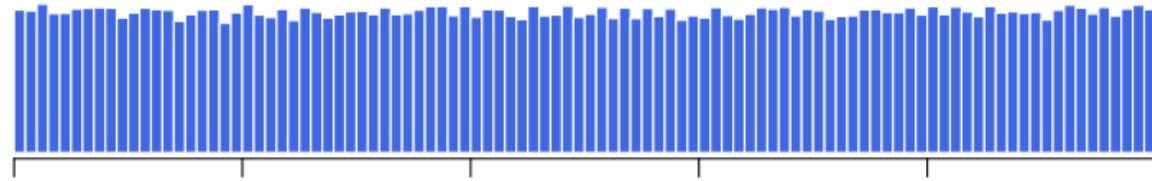


# HISTOGRAMS



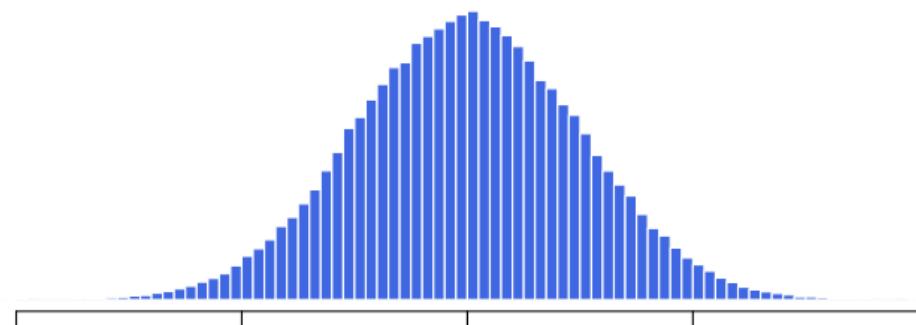
# SOME SIMPLE DISTRIBUTIONS

## UNIFORM



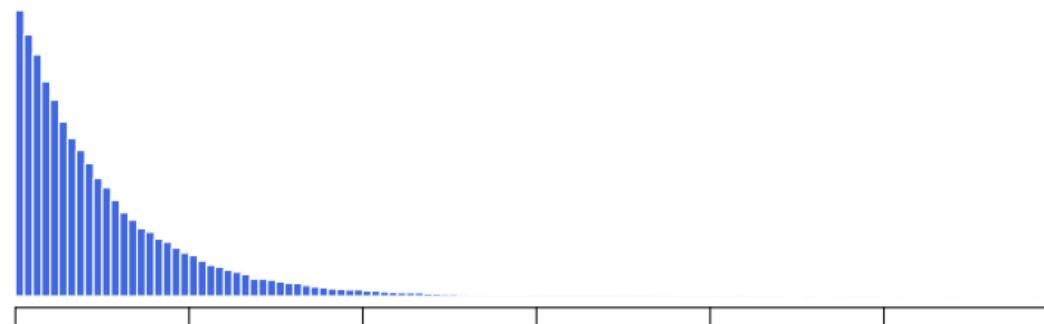
## NORMAL

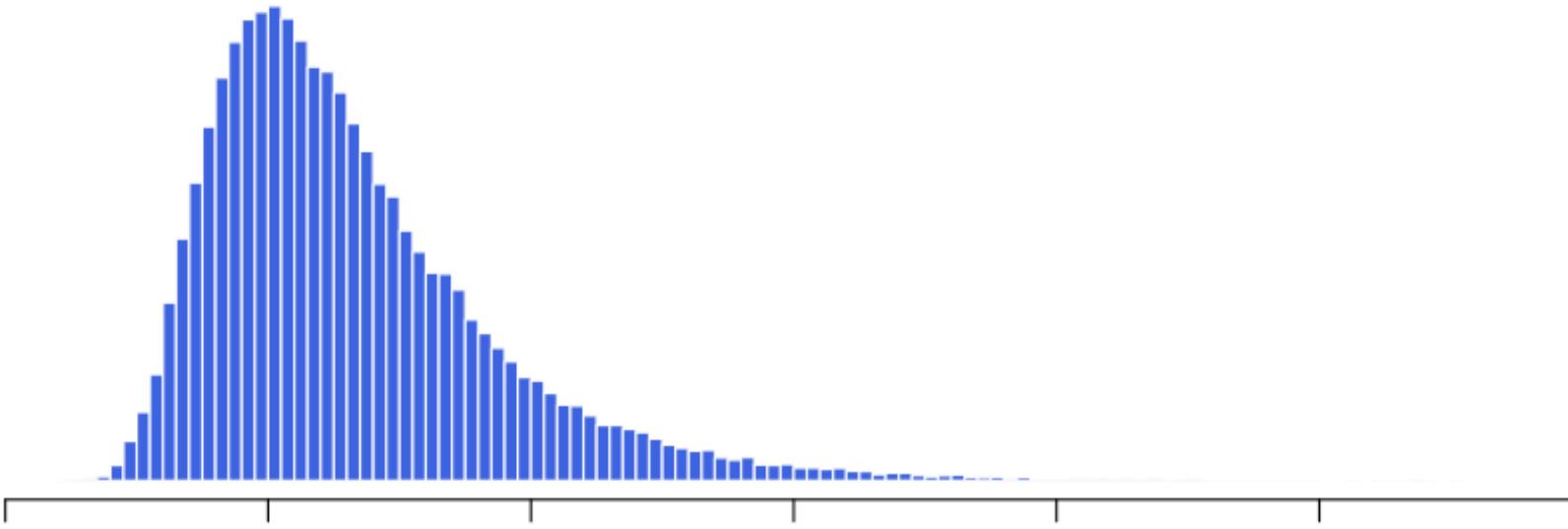
REPEATED MEASURES, VARIATION  
IN POPULATIONS, ETC.



## EXPONENTIAL

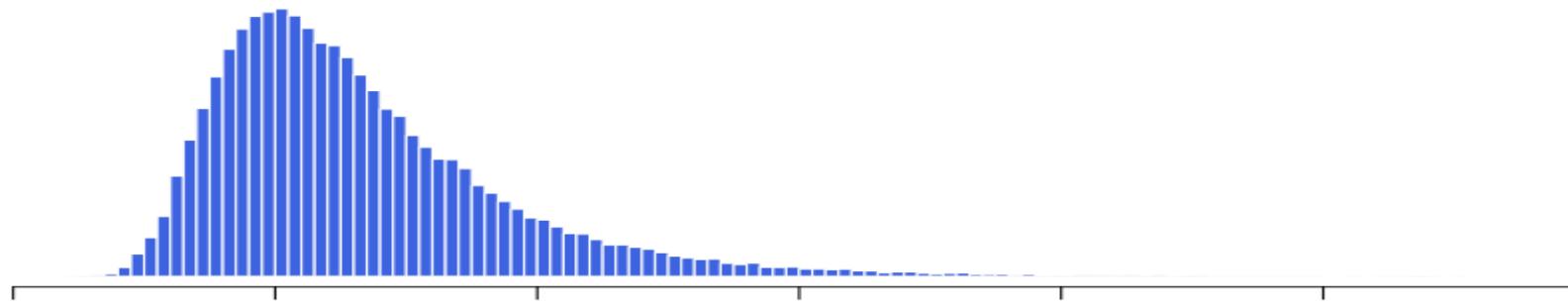
DURATIONS BETWEEN EVENTS, ETC.



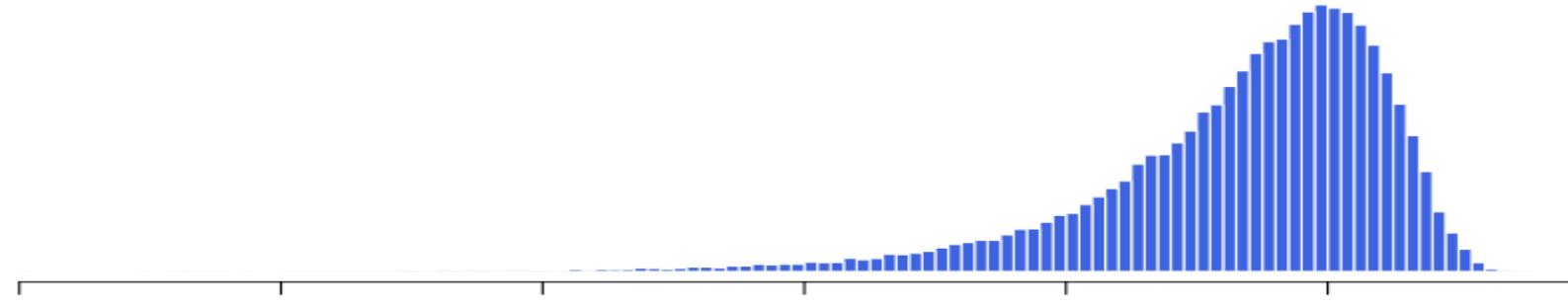


# LOG-NORMAL

MANY BIOLOGICAL AND SOCIAL PROCESSES  
(TISSUE GROWTH, CITY SIZE, # OF FOLLOWERS)

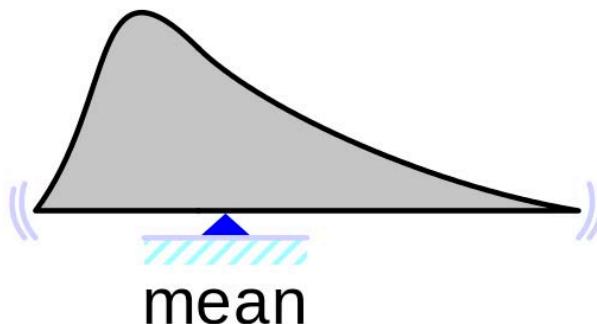
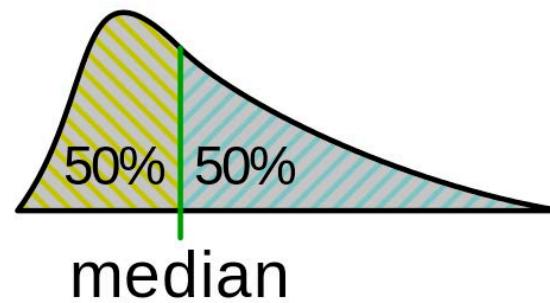
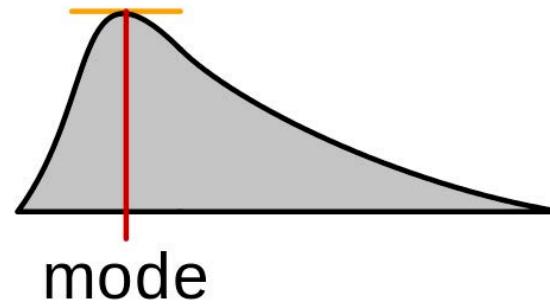


SKEWED LEFT

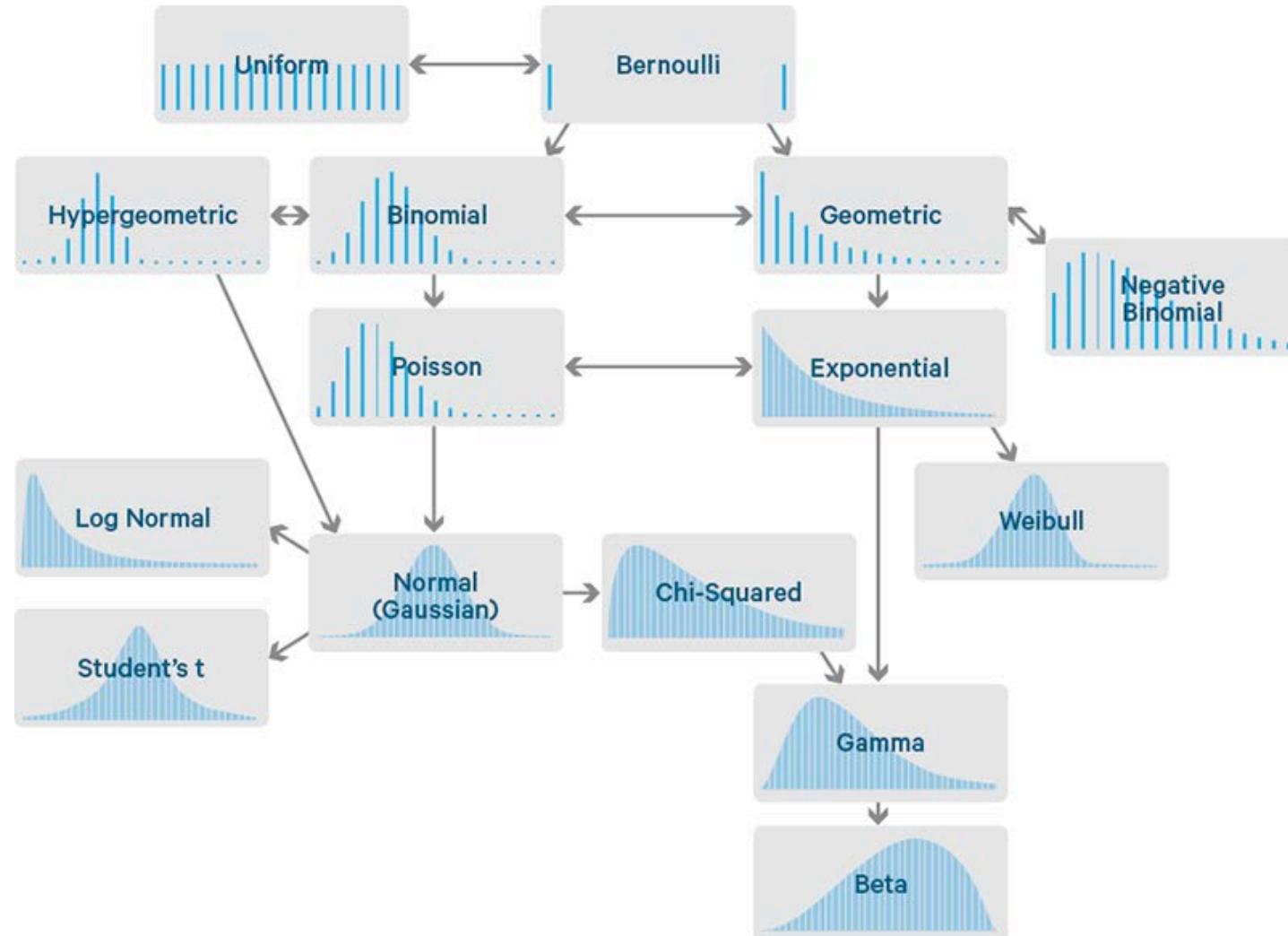


SKEWED RIGHT

# DISTRIBUTIONS AND SOME COMMON MEASURES OF CENTRAL TENDENCY

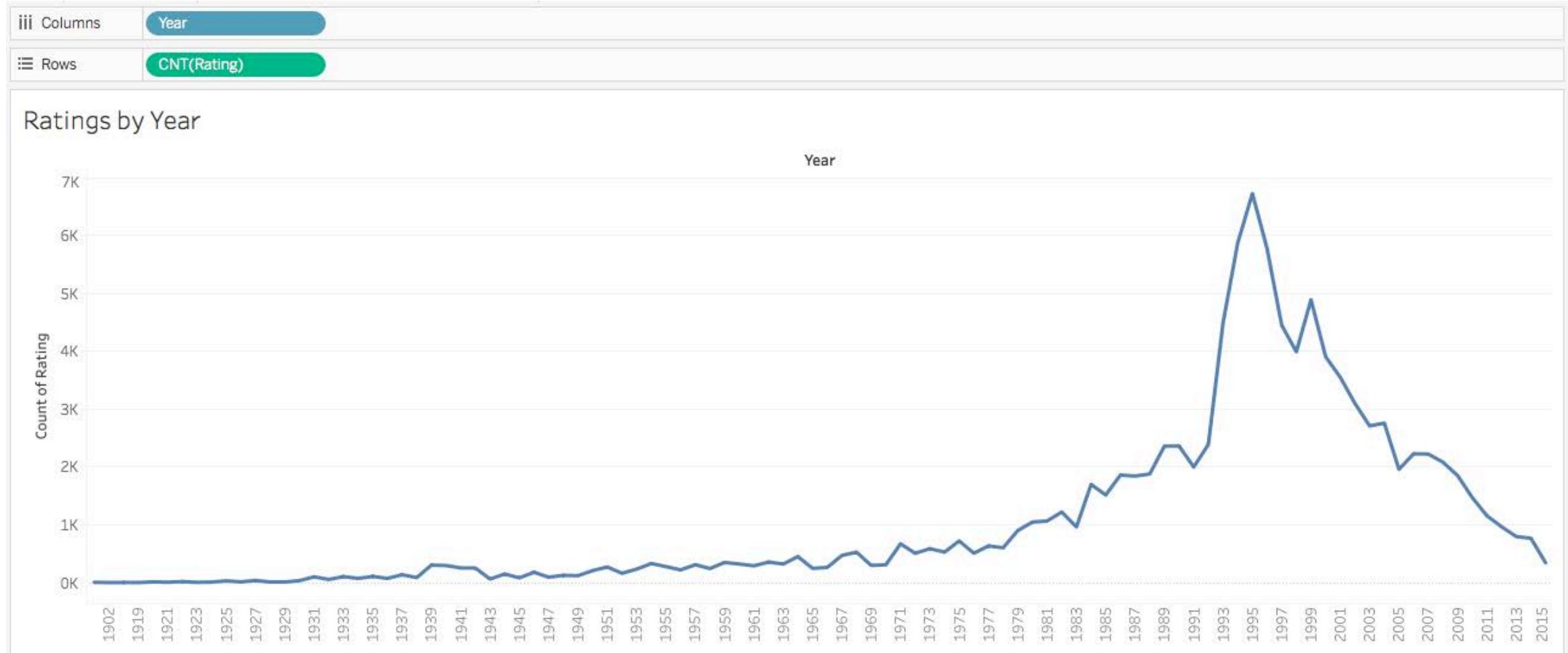


# COMMON DISTRIBUTIONS



CLOUDERA

# DISTRIBUTIONS OVER TIME



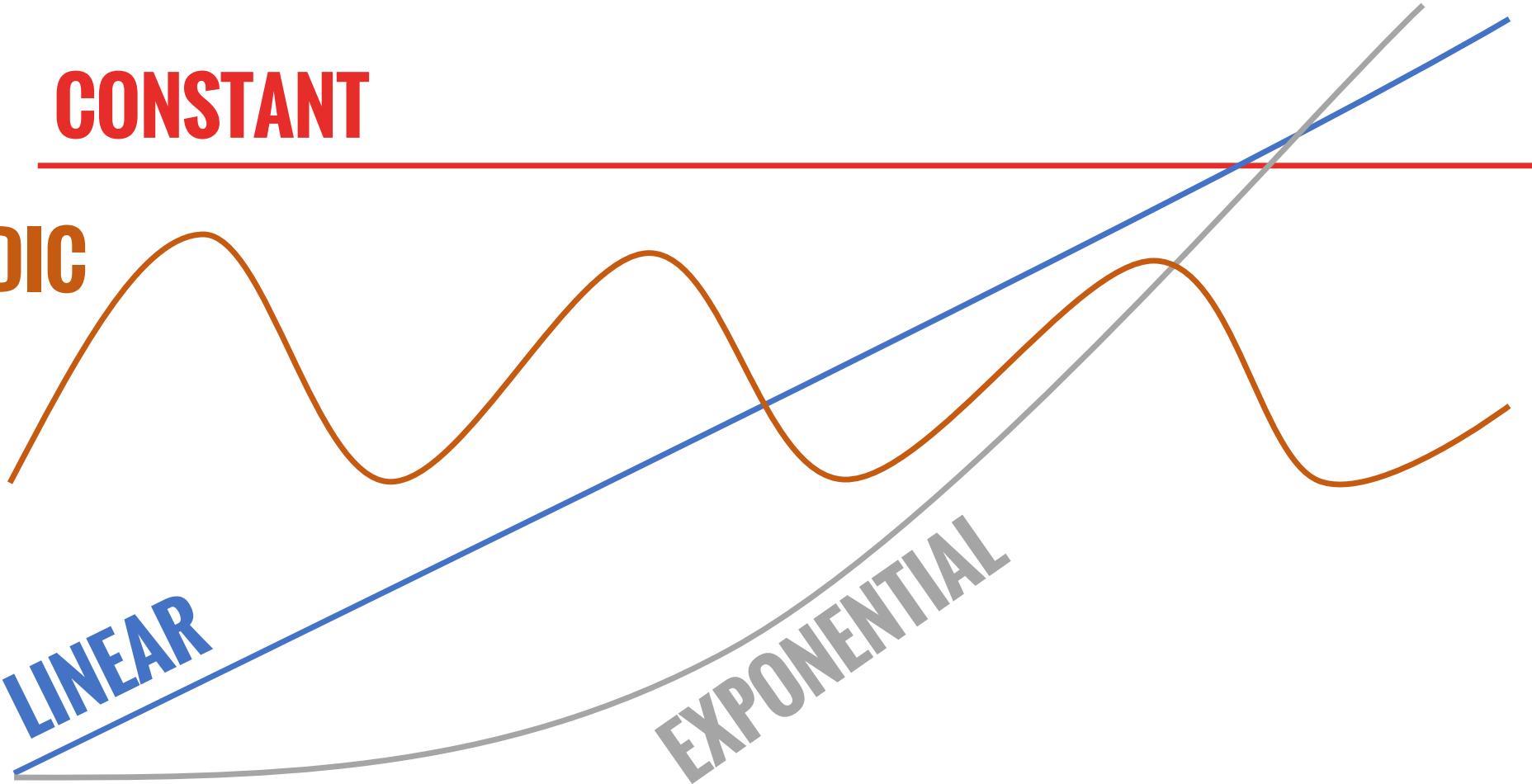
# IDENTIFYING TRENDS

**CONSTANT**

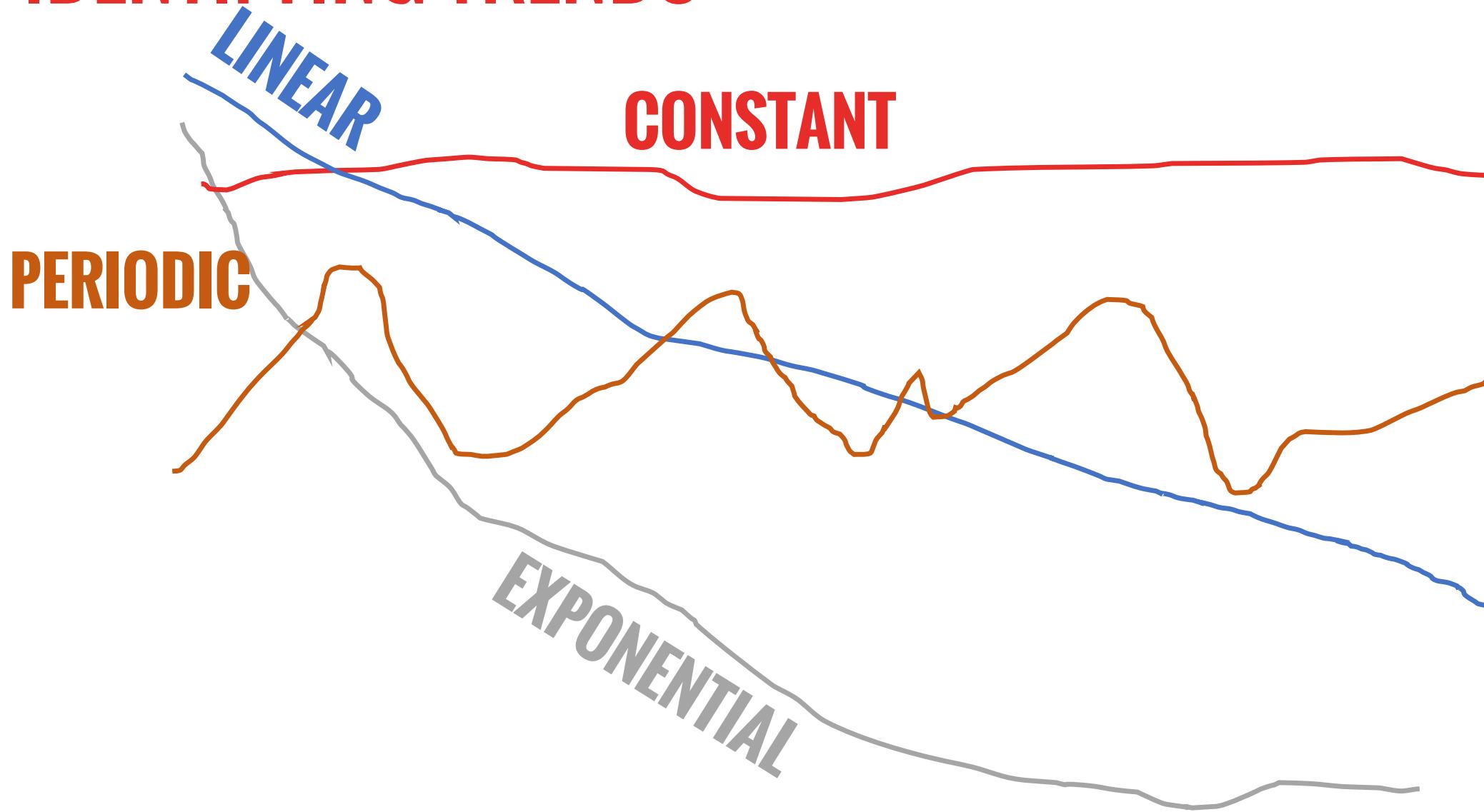
**PERIODIC**

**LINEAR**

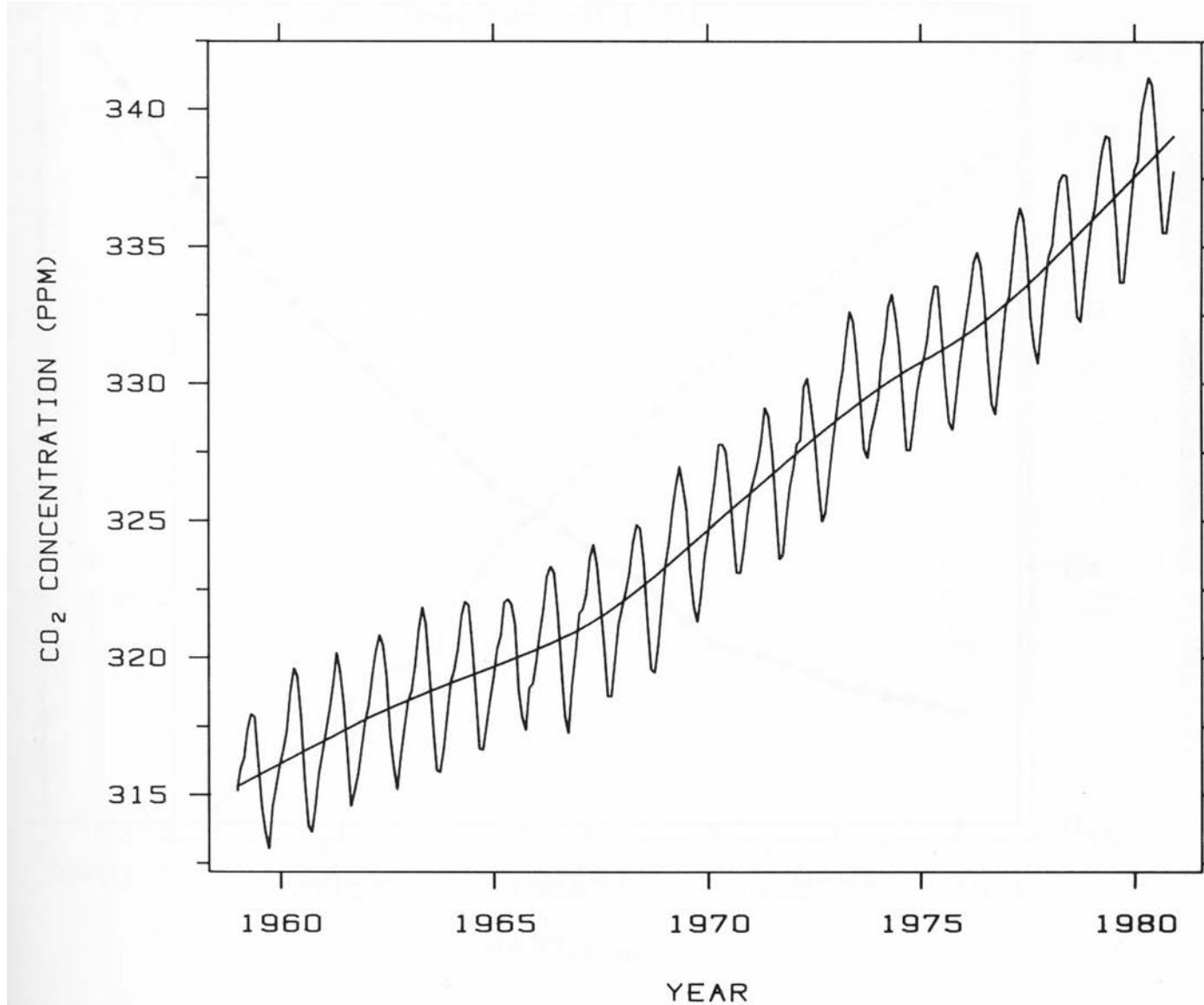
**EXPONENTIAL**



# IDENTIFYING TRENDS

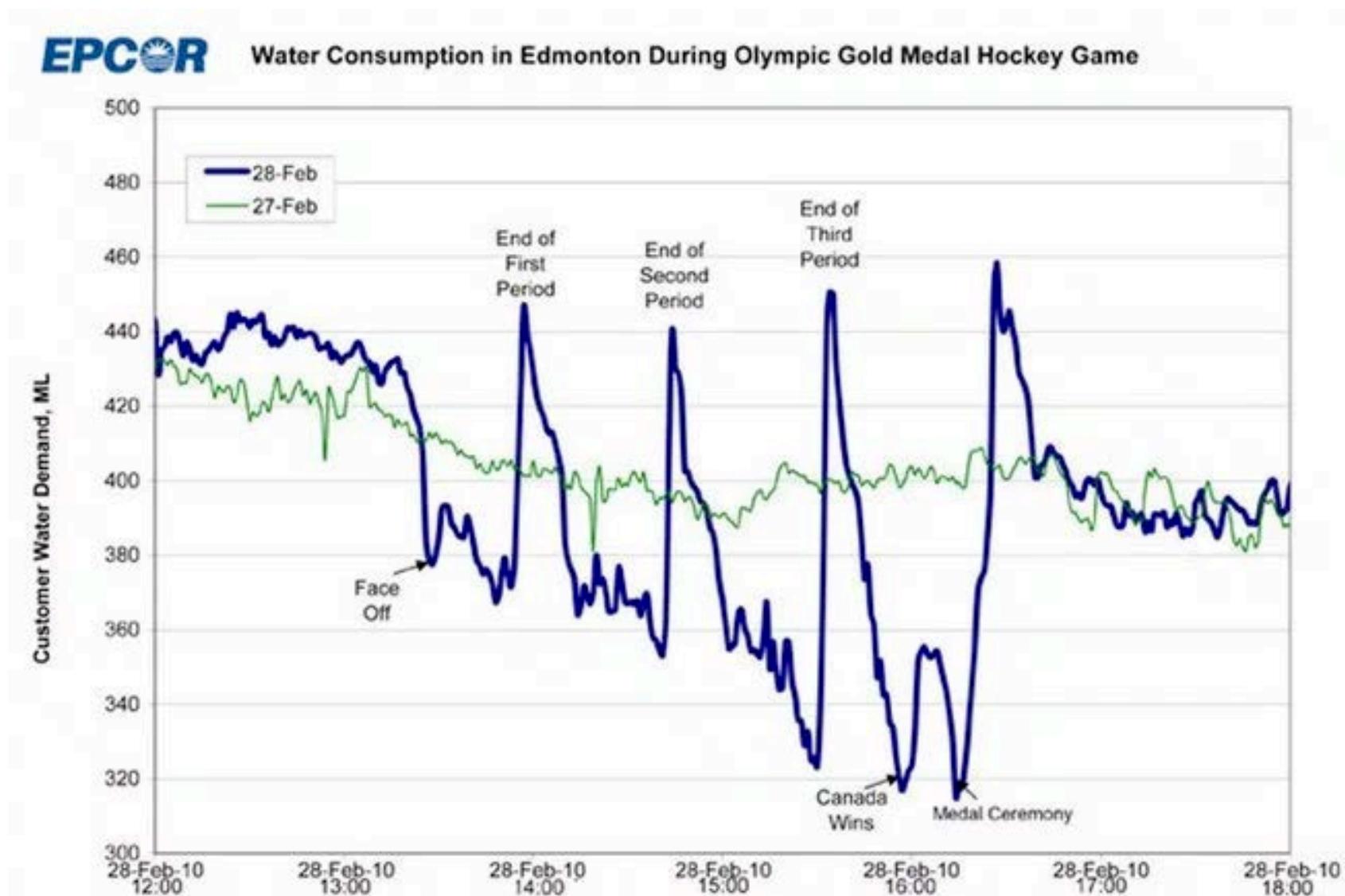


# COMBINATIONS



**YEARLY CO<sub>2</sub> CONCENTRATIONS**  
CLEVELAND 85

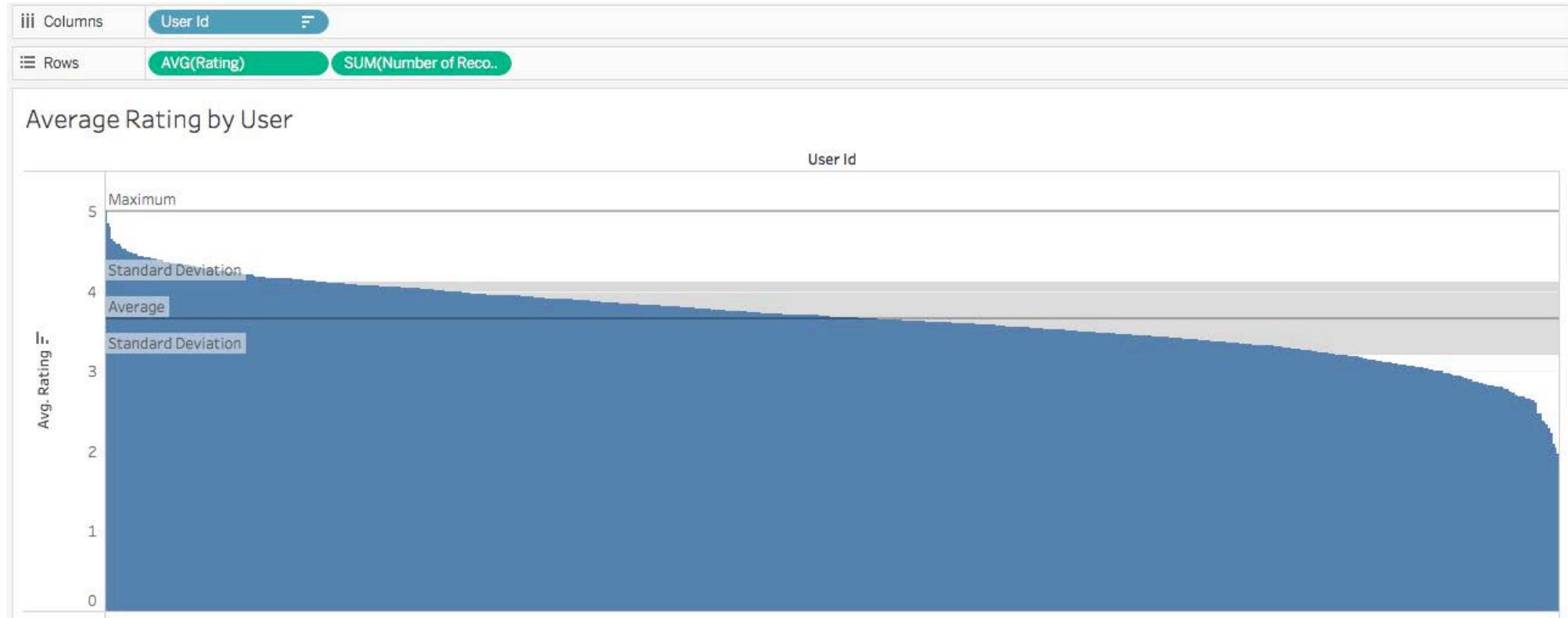
# COMPARISON AGAINST A KNOWN BASELINE



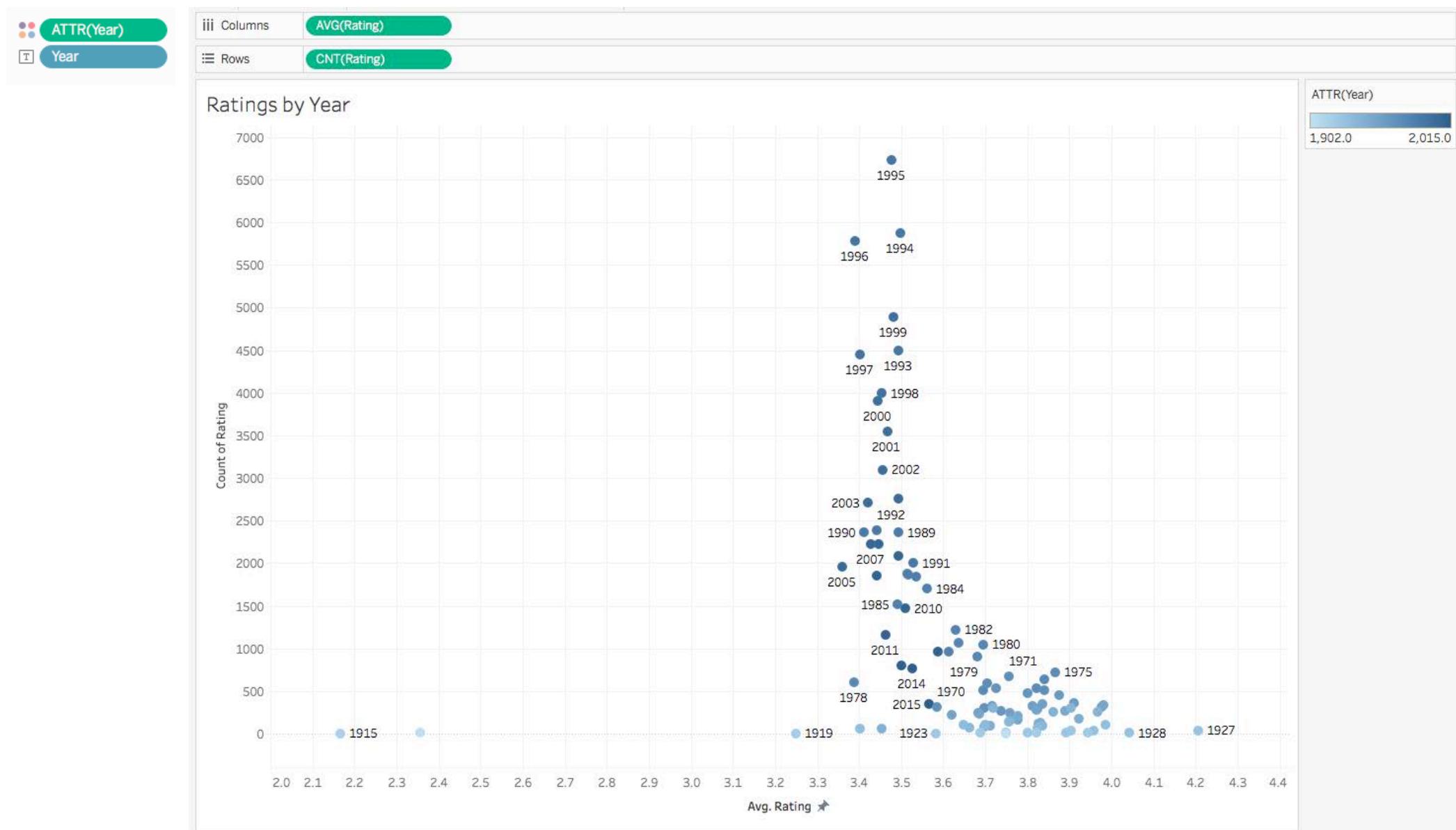
# SOME FIRST STEPS TO START WITH

1. Plot the raw data
2. Plot simple statistics
3. Examine distributions
4. Look at attributes together

# COMPARE MULTIPLE PLOTS



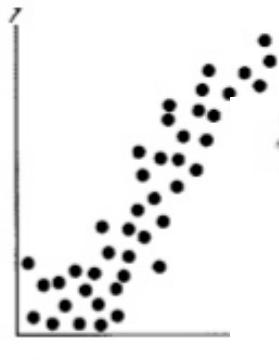
# COMPARING TWO MEASURES



# IDENTIFYING CORRELATIONS



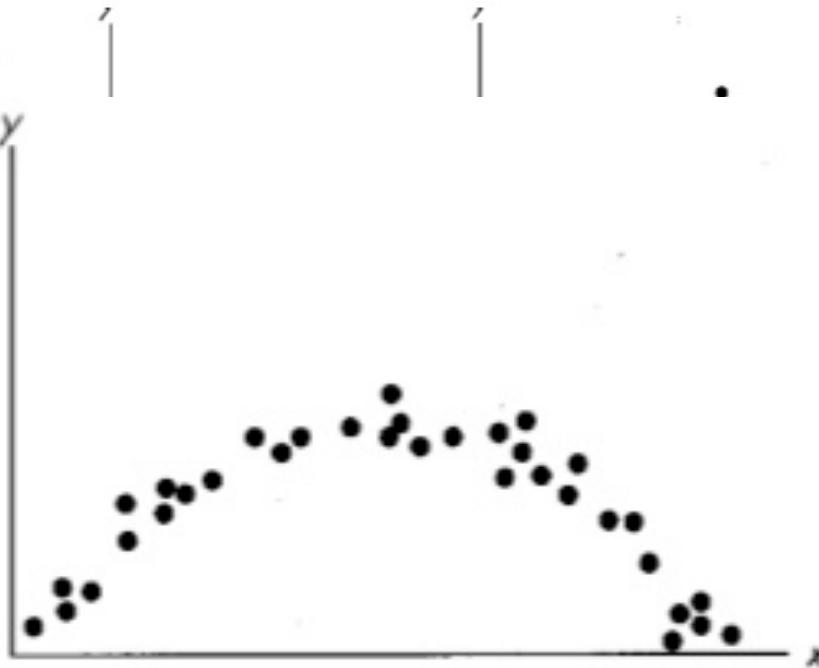
(g) No correlation  
between  $x$  and  $y$



(a) Positive corre  
between  $x$  and  $y$



(d) Negative cor  
between  $x$  and  $y$



(h) Nonlinear correlation  
between  $x$  and  $y$

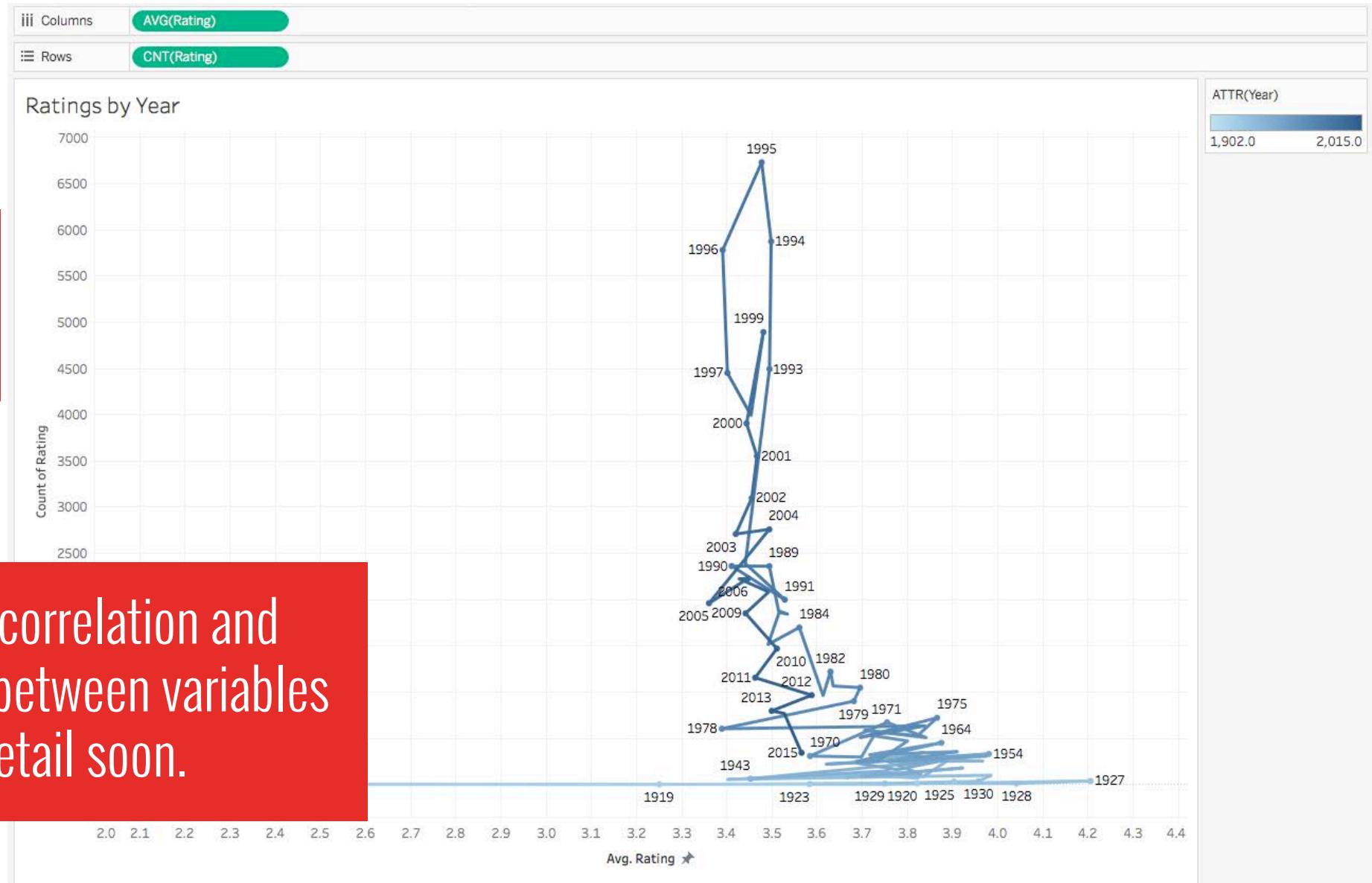
correlation between  
 $x$  and  $y$

correlation between  
 $x$  and  $y$

# COMPARING TWO MEASURES

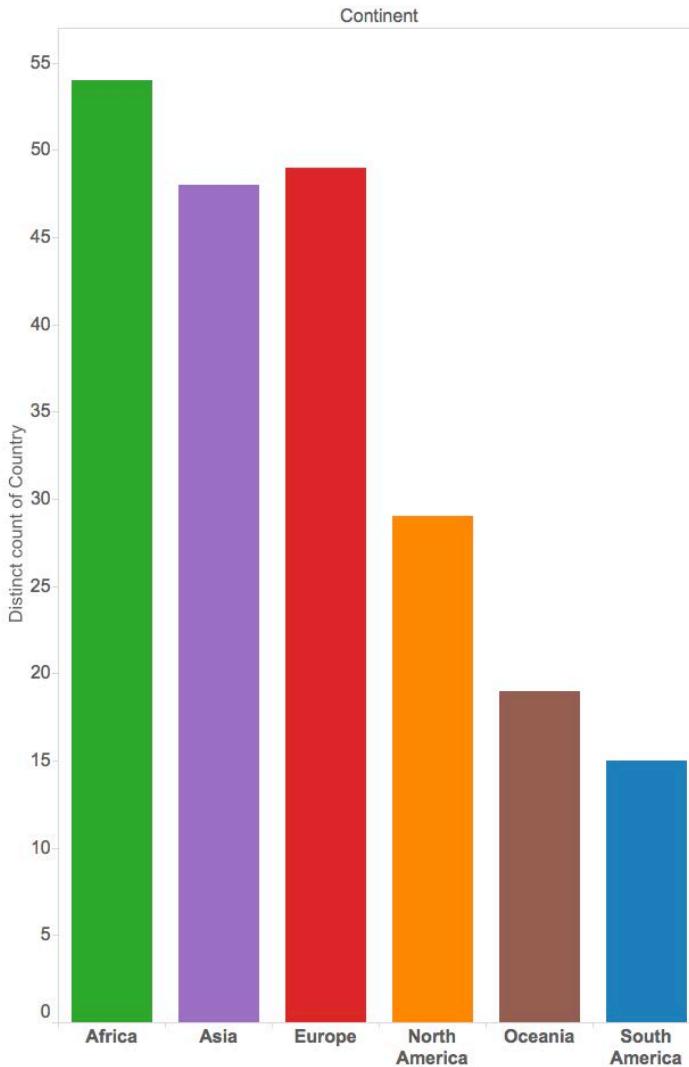
ANY  
HYPOTHESES?

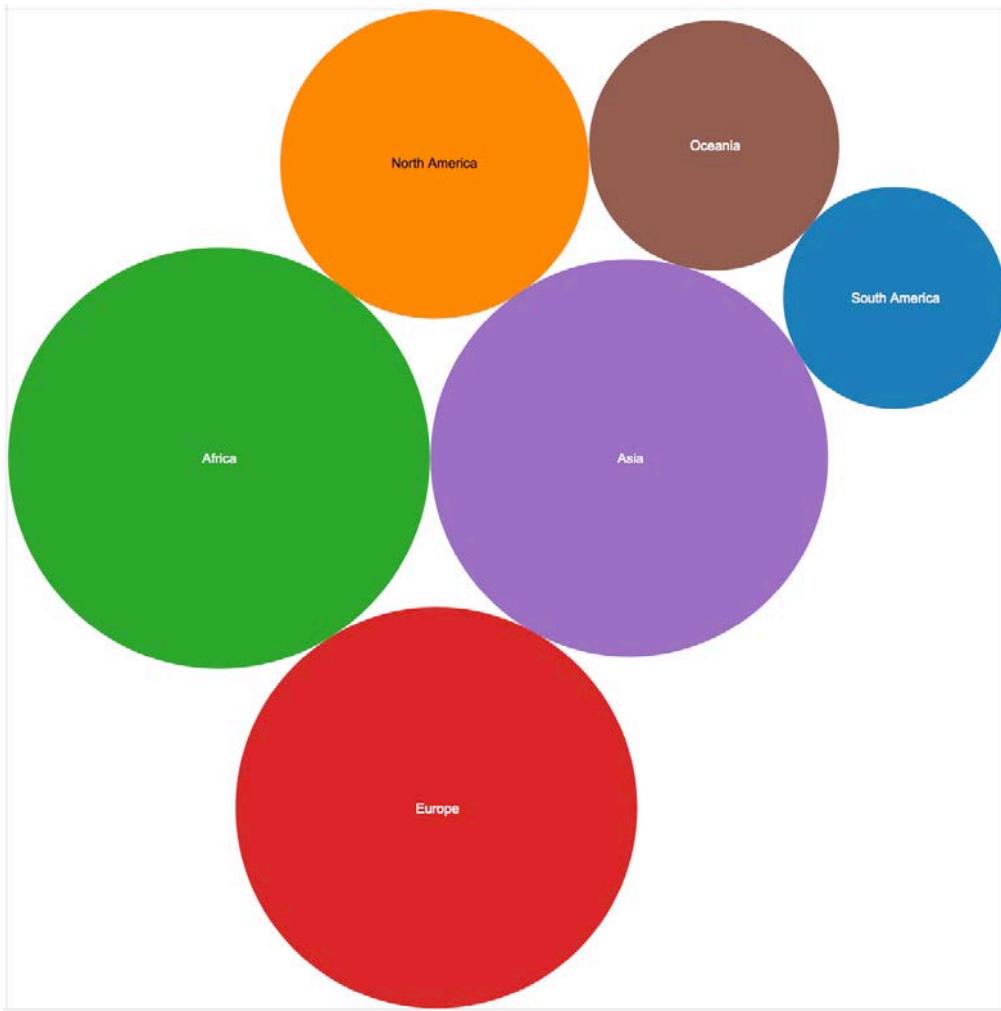
We will revisit correlation and  
Relationships between variables  
In a lot more detail soon.



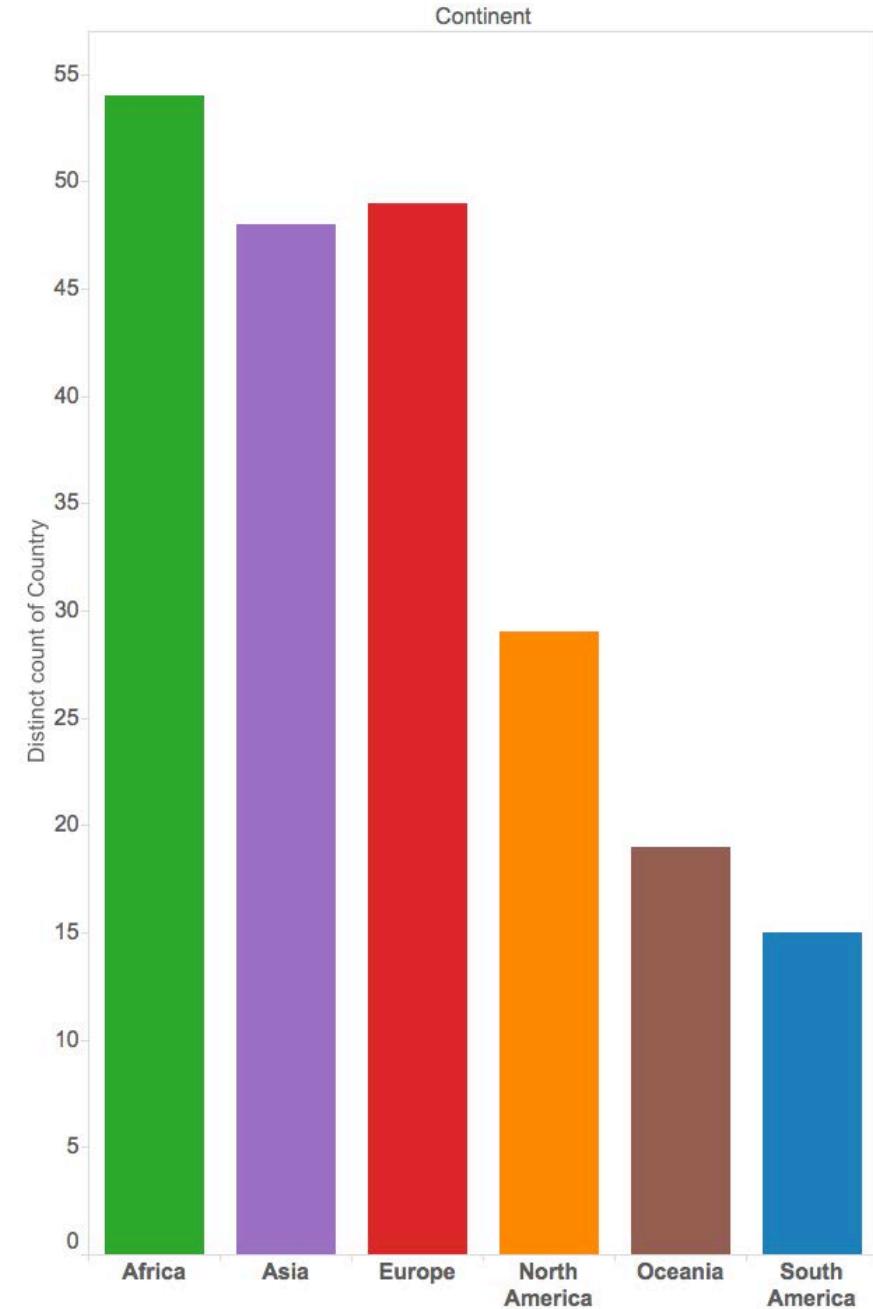
# SOME COMMON VISUAL MAPPING MISTAKES

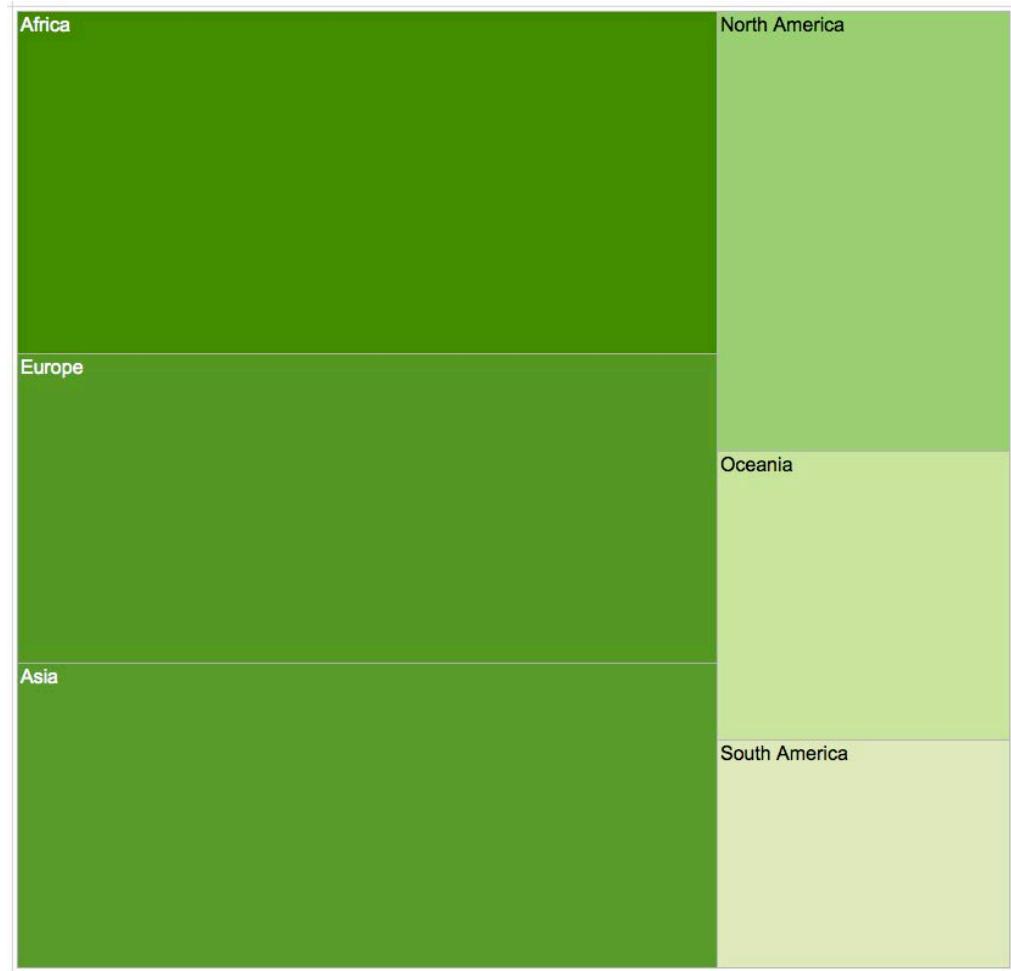
# COMPARING QUANTITIES



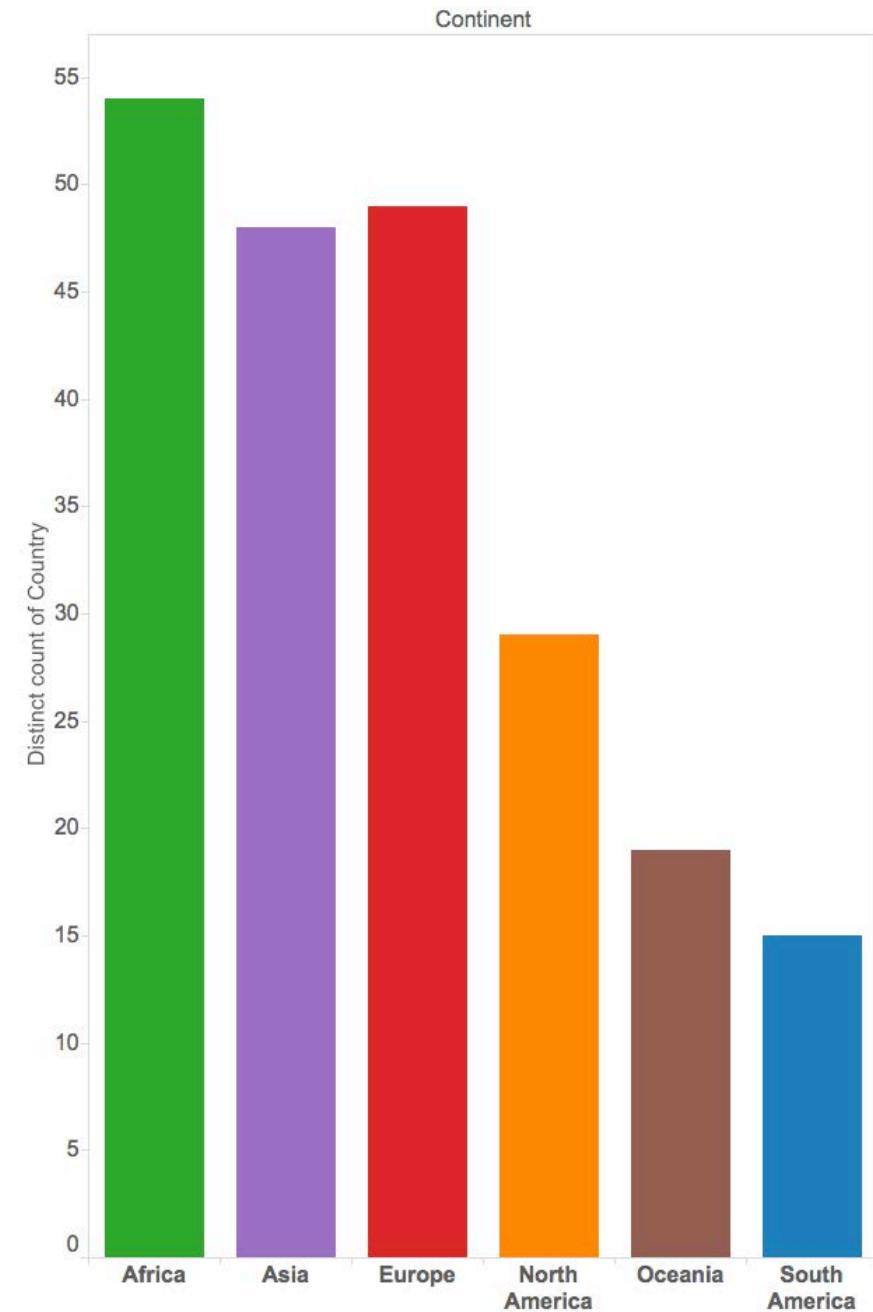


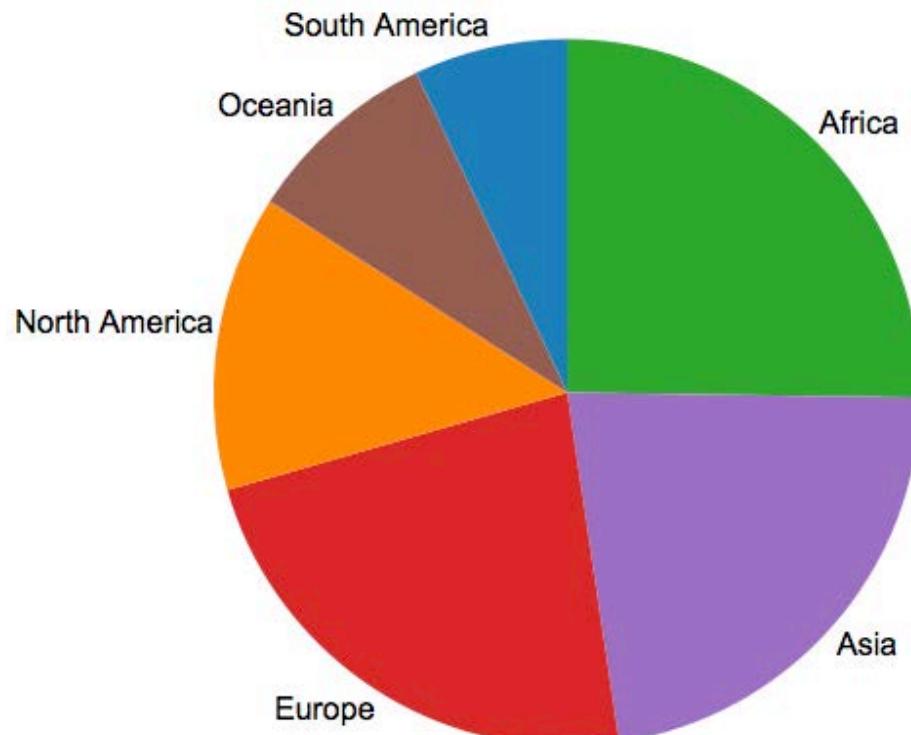
VS



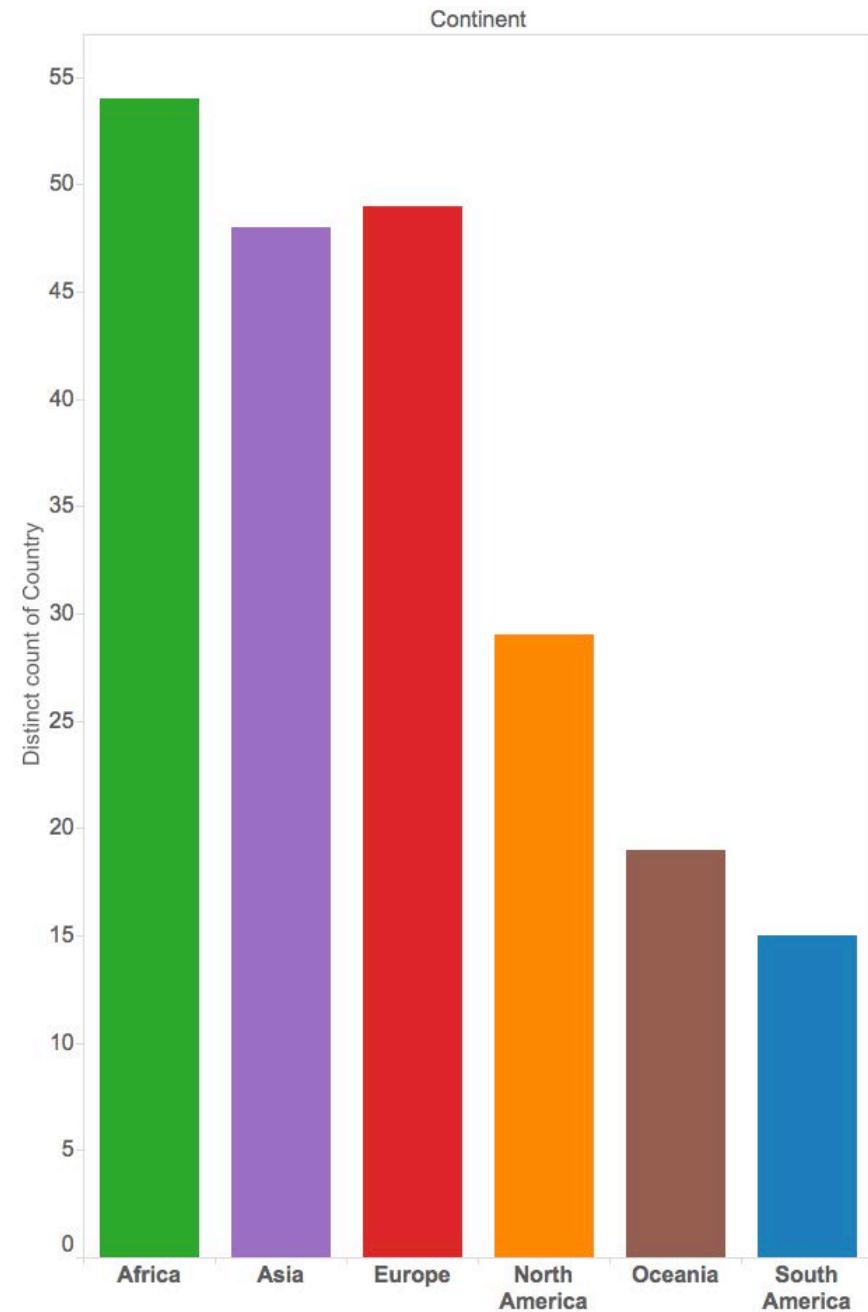


VS

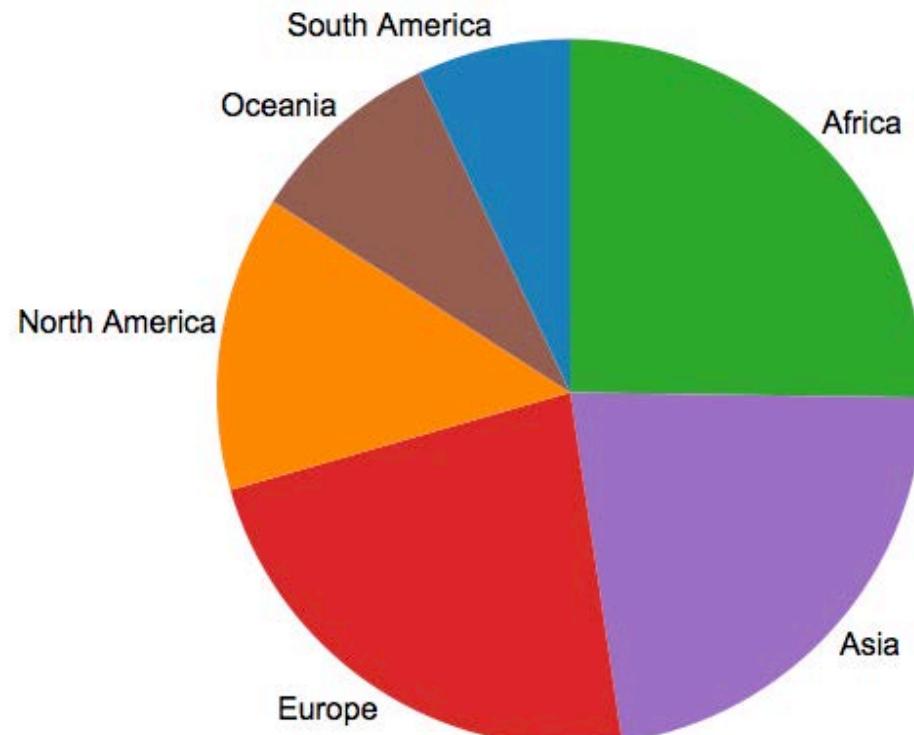




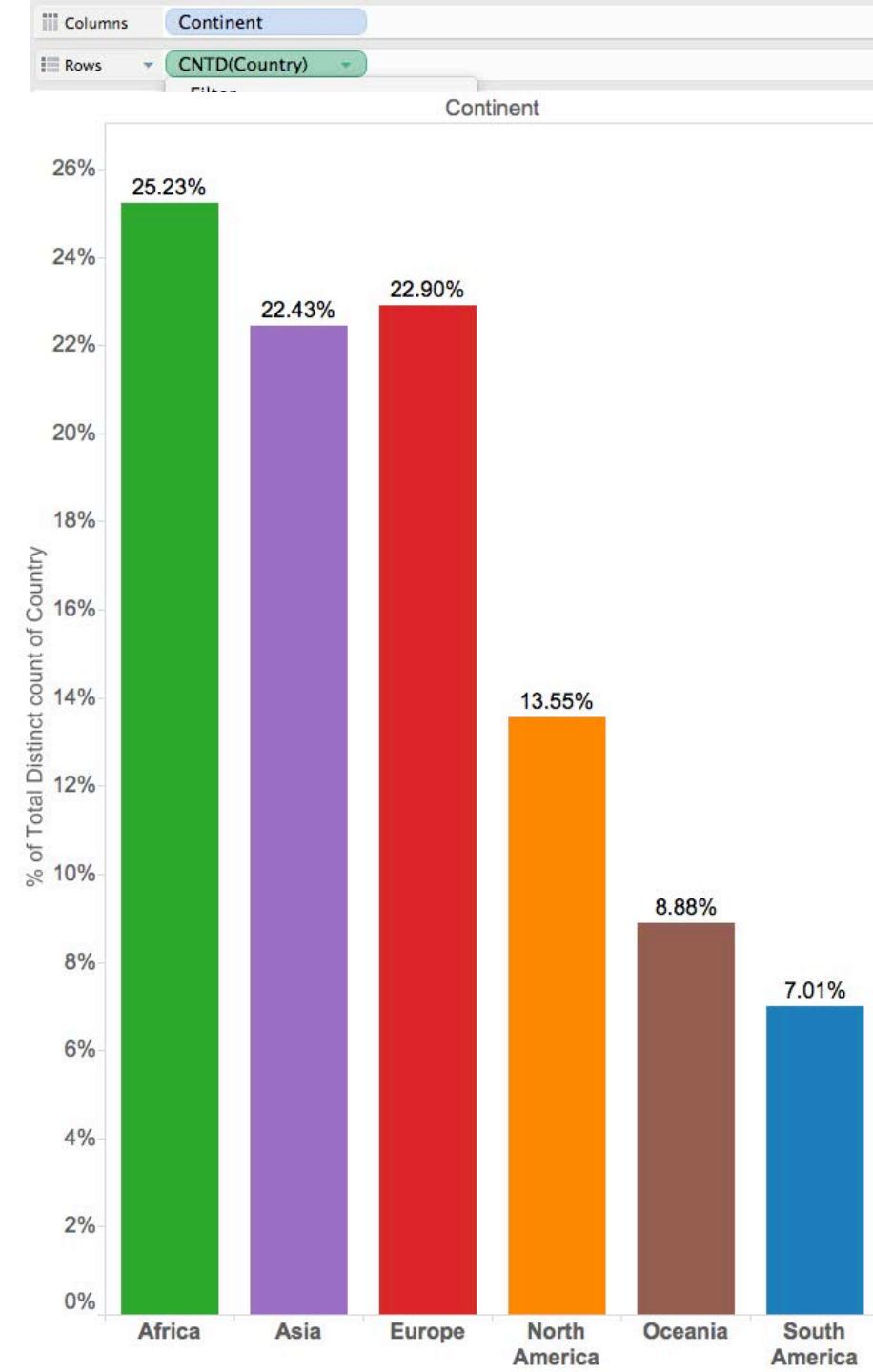
VS

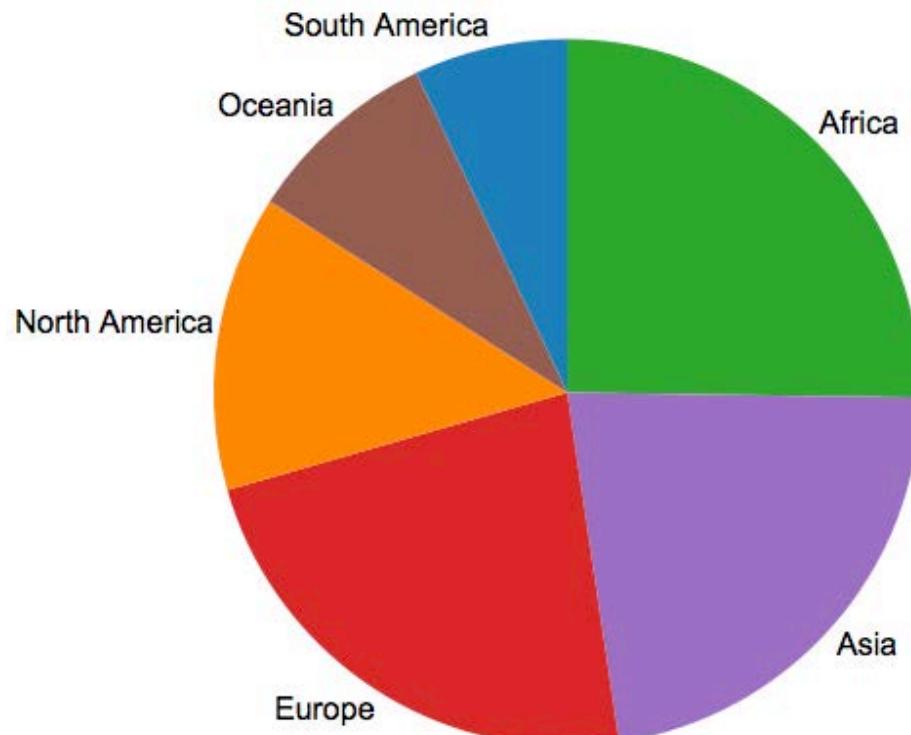


**PARTS-OF-A-WHOLE  
RELATIONSHIPS**

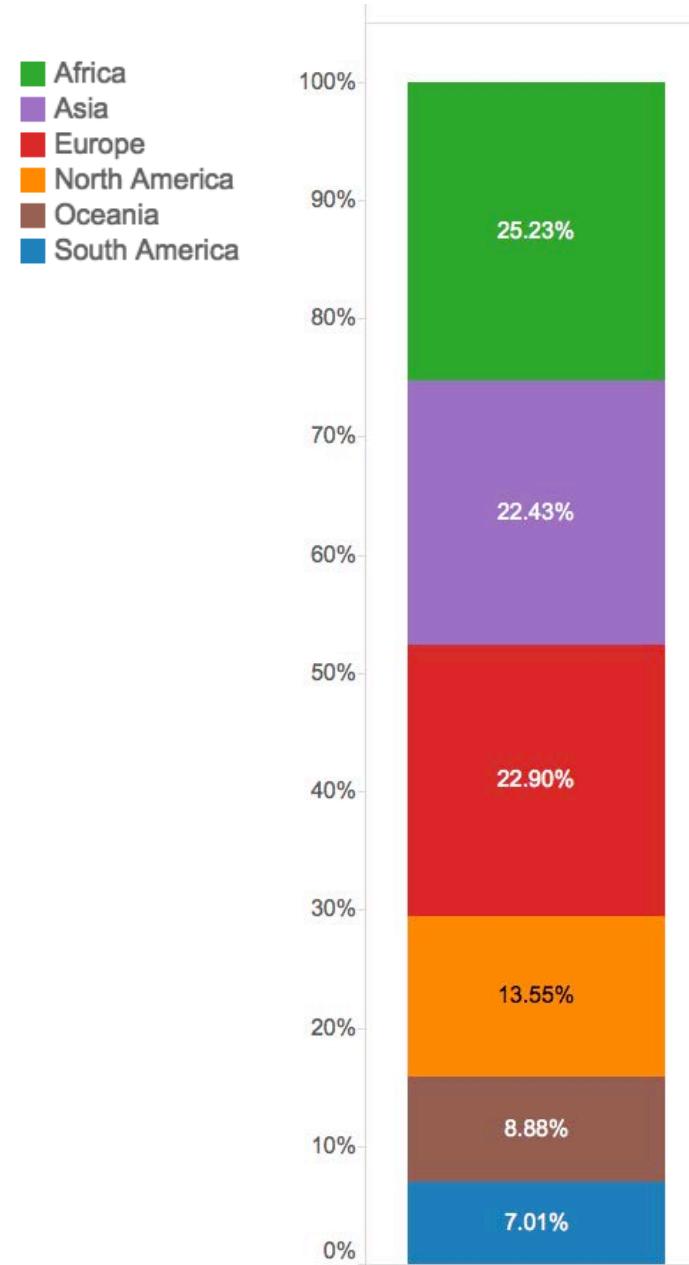


VS

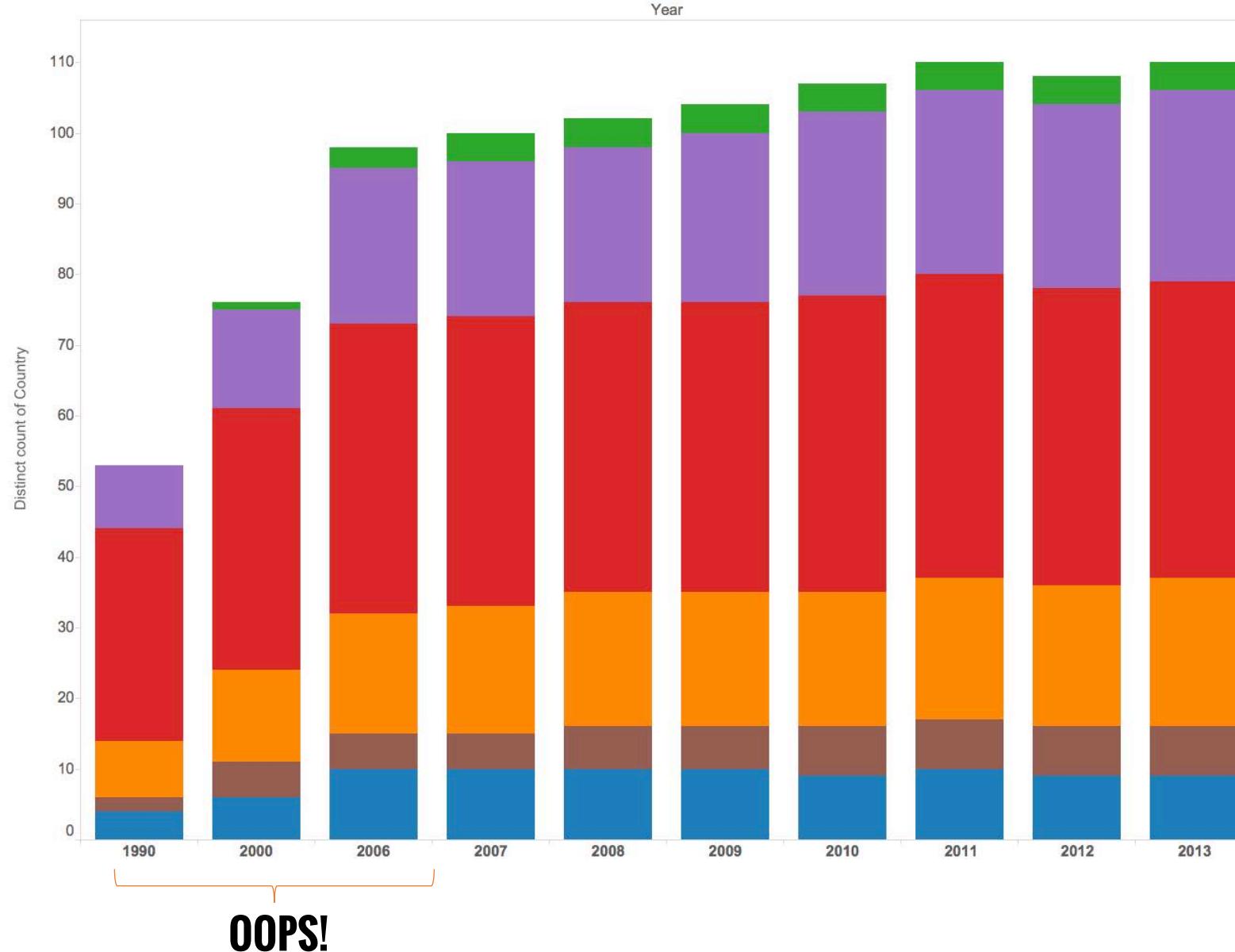




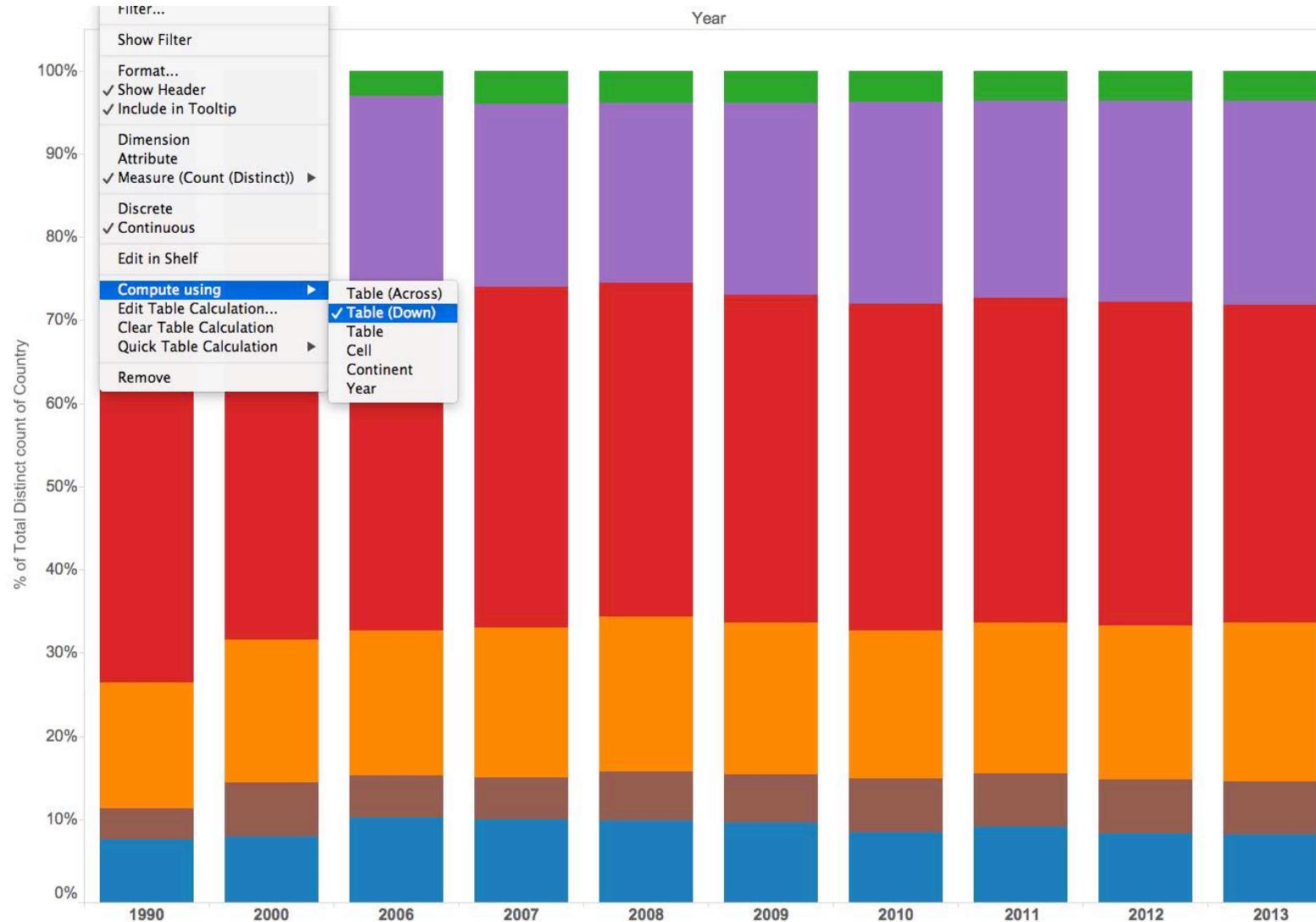
VS



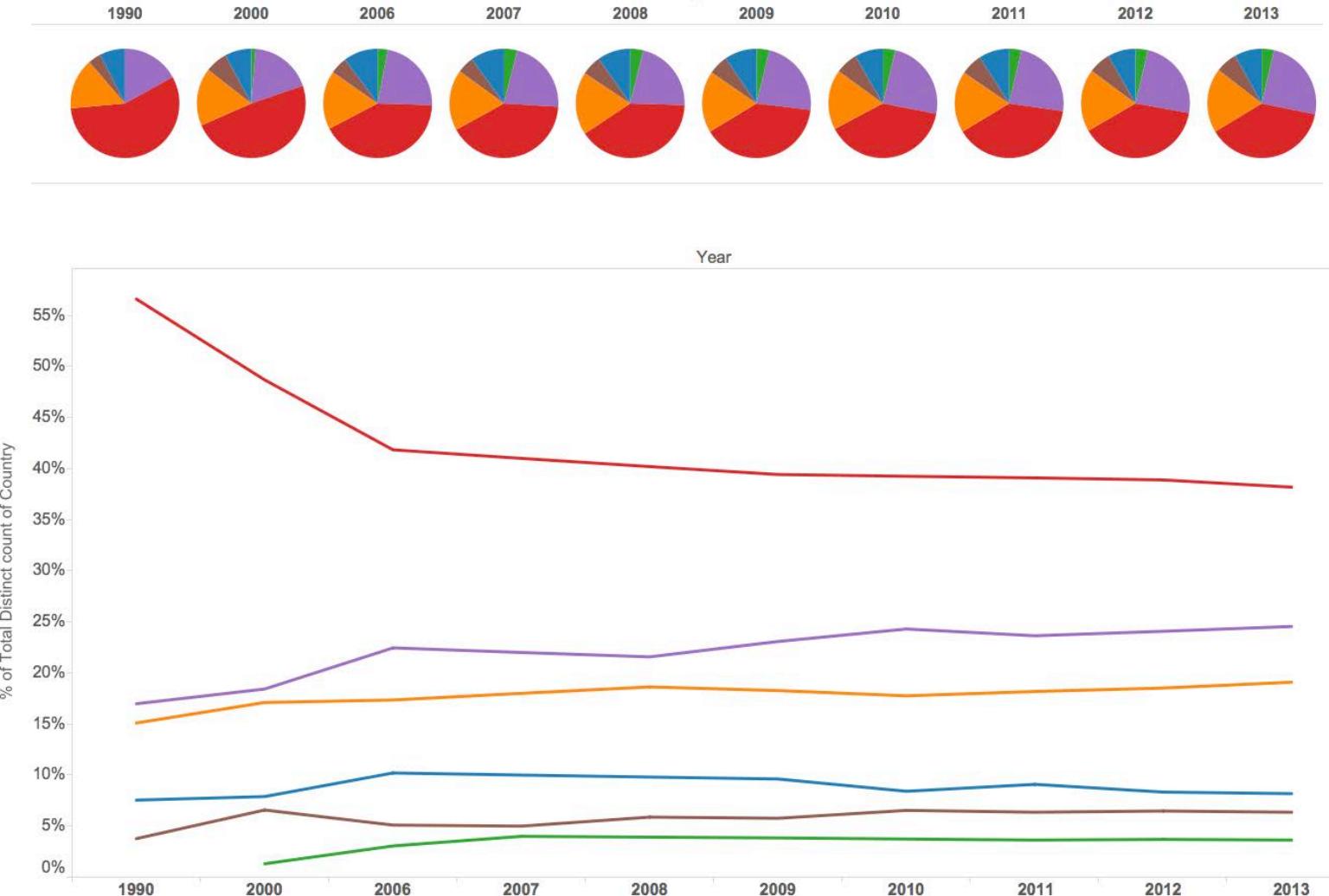
# # OF COUNTRIES W/ LIFE EXPECTANCY > 70 YEARS



# # OF COUNTRIES W/ LIFE EXPECTANCY > 70 YEARS

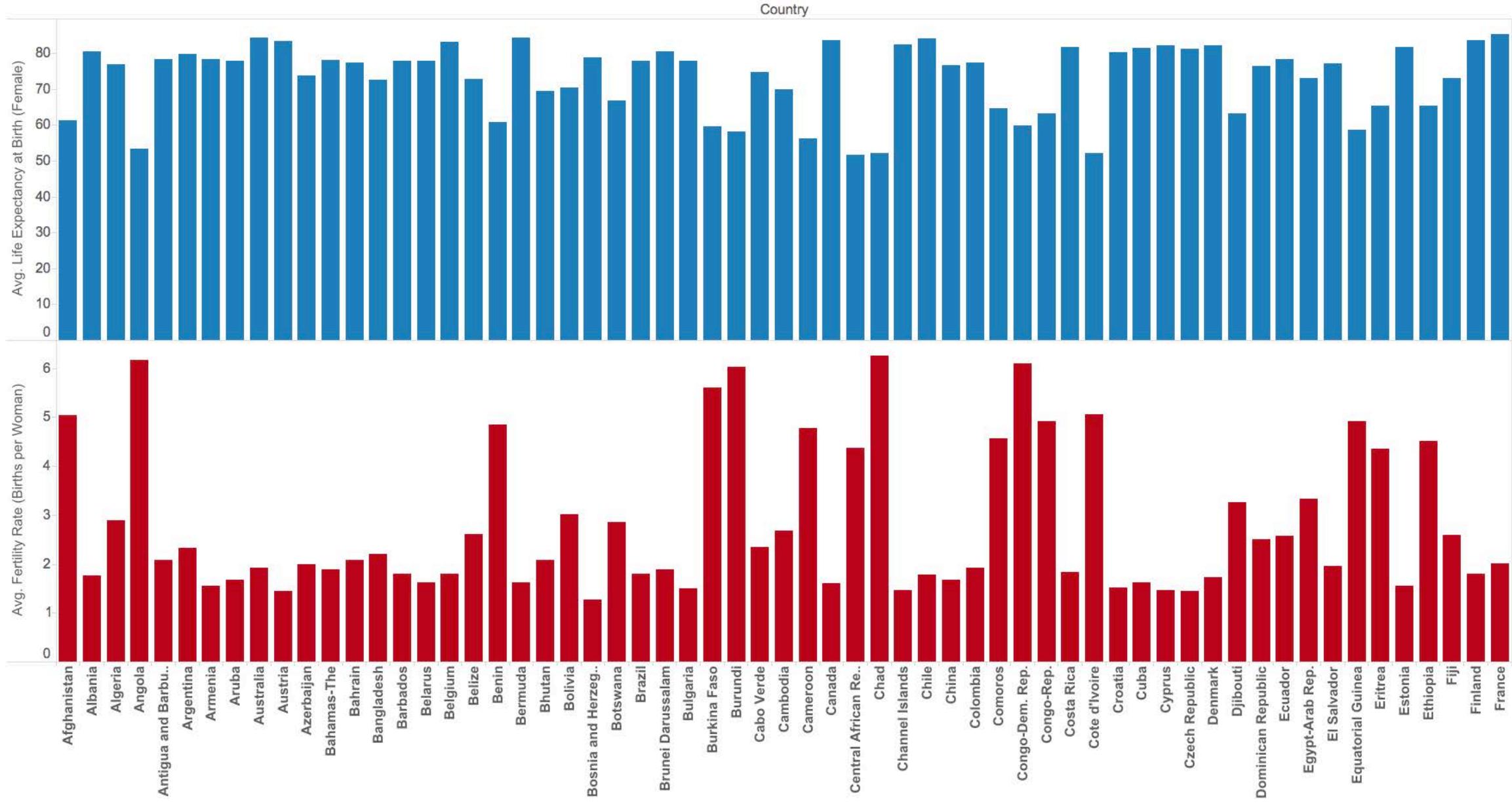


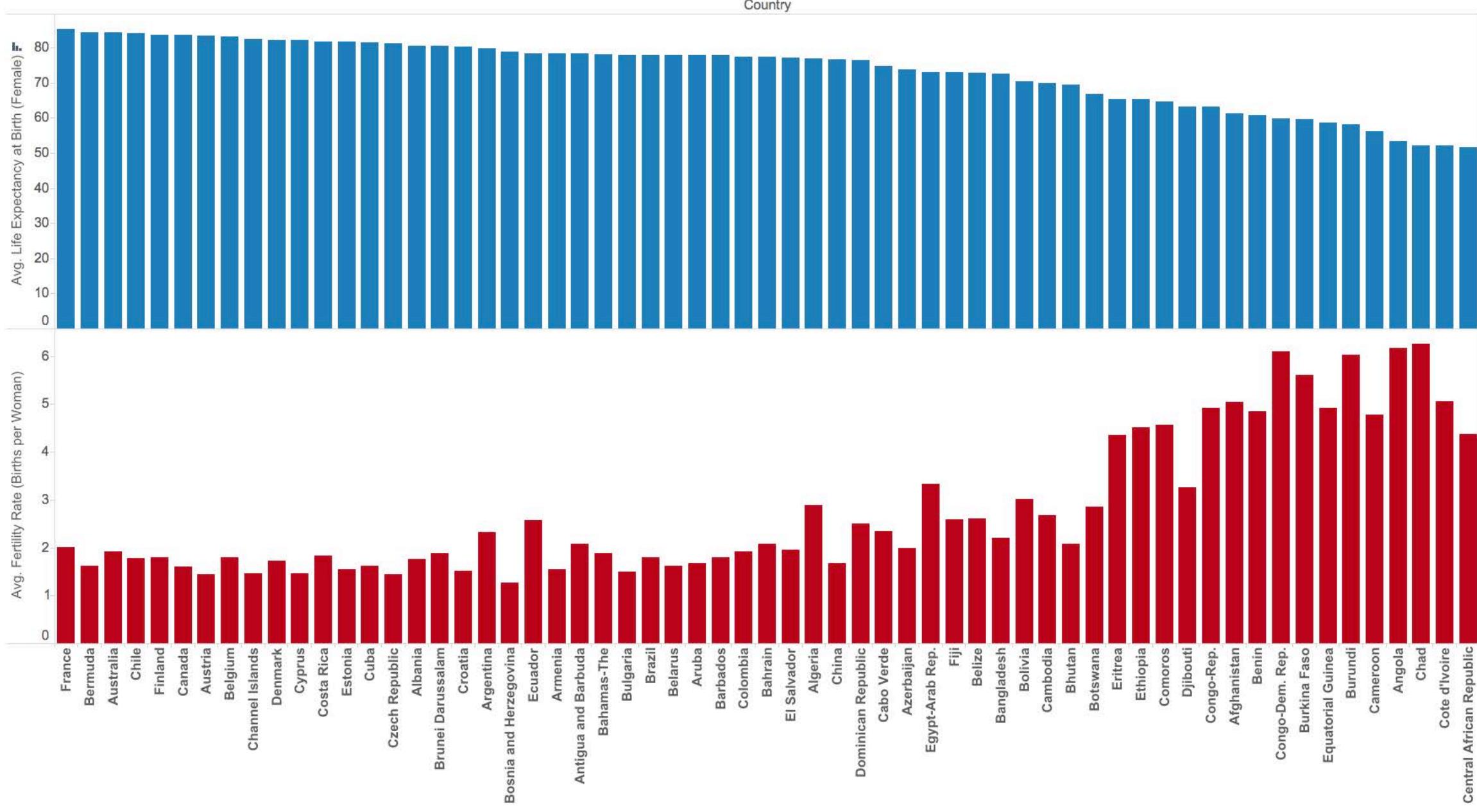
# # OF COUNTRIES W/ LIFE EXPECTANCY > 70 YEARS

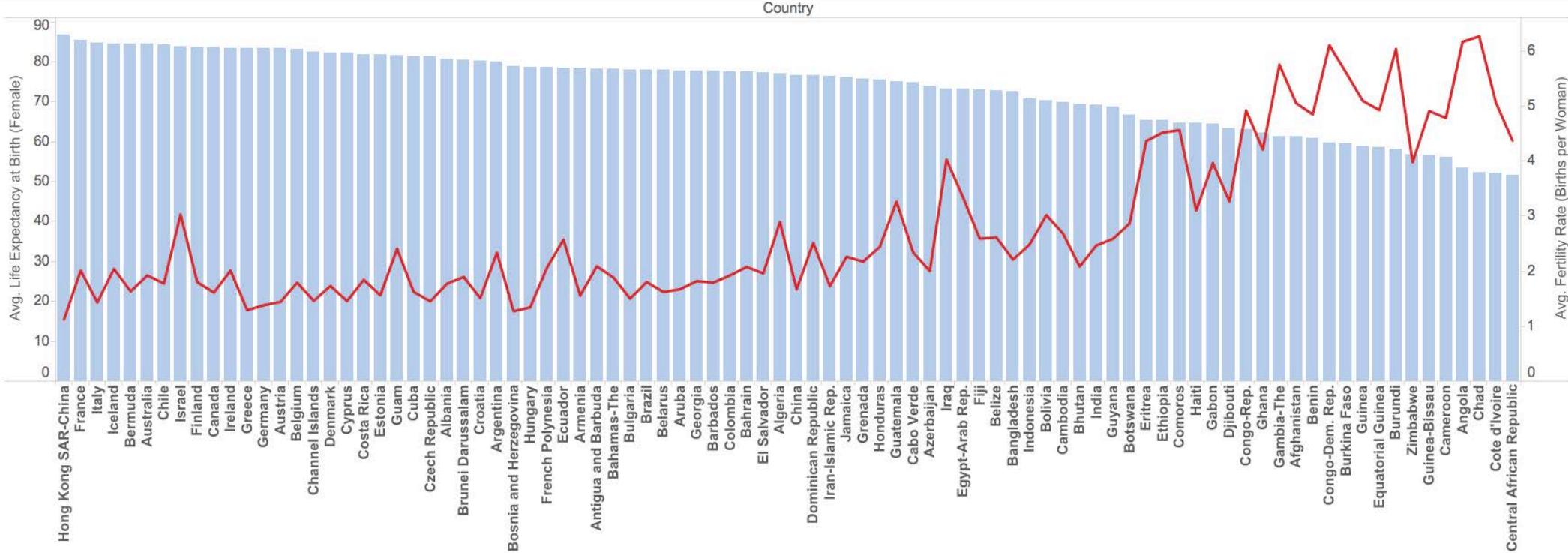


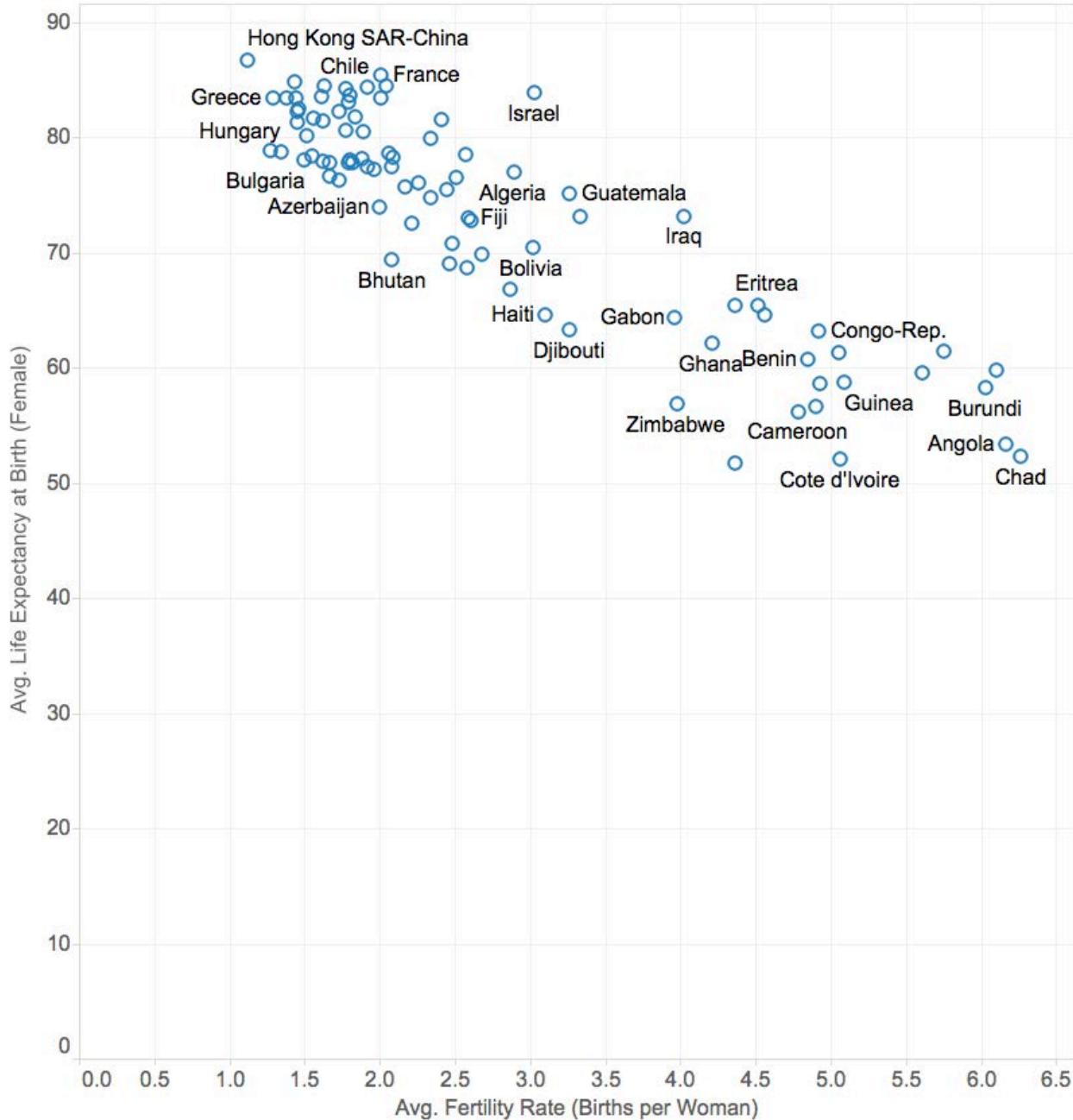
# SHOWING RELATIONSHIPS BETWEEN MEASURES

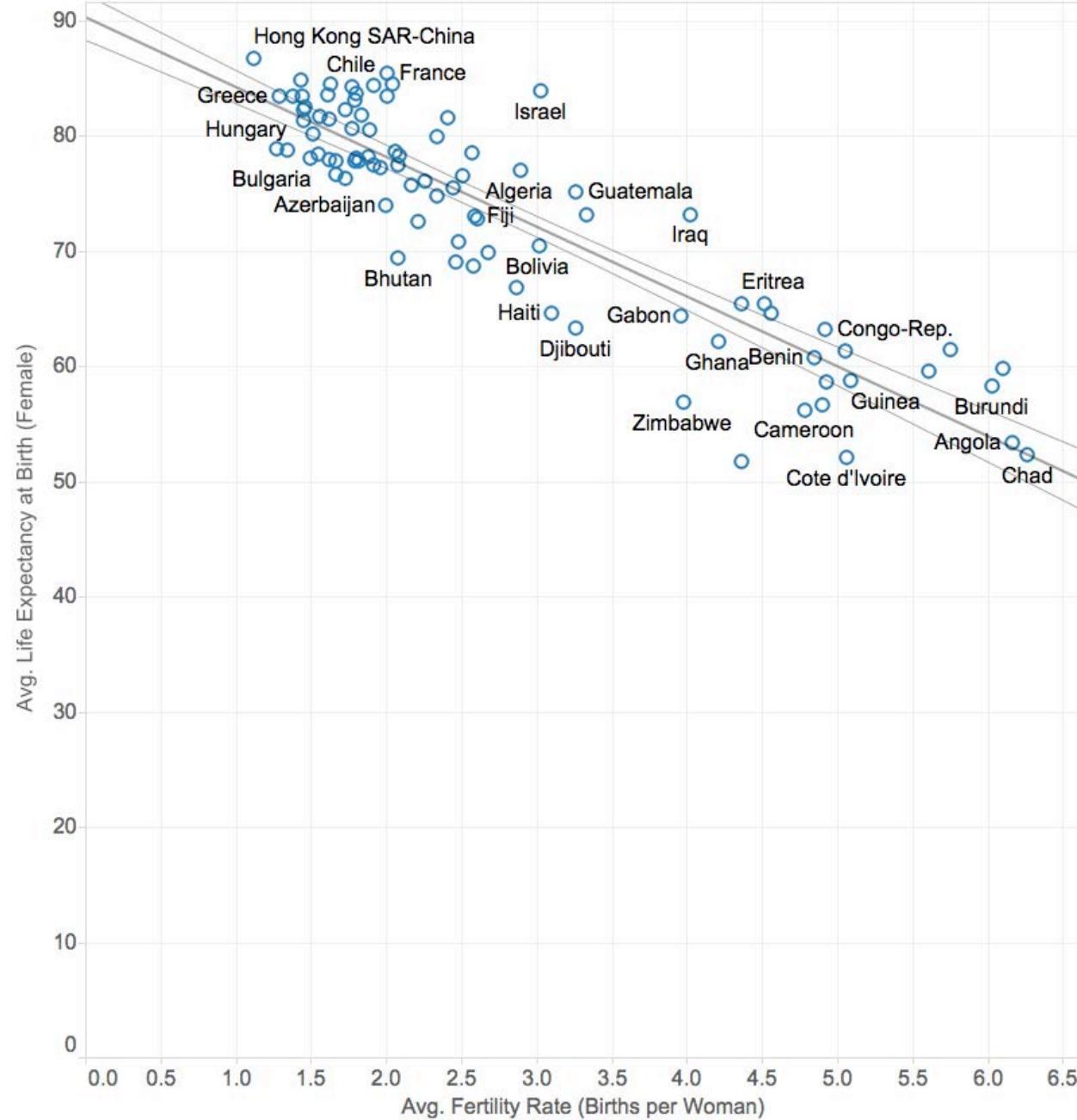
(CORRELATIONS)

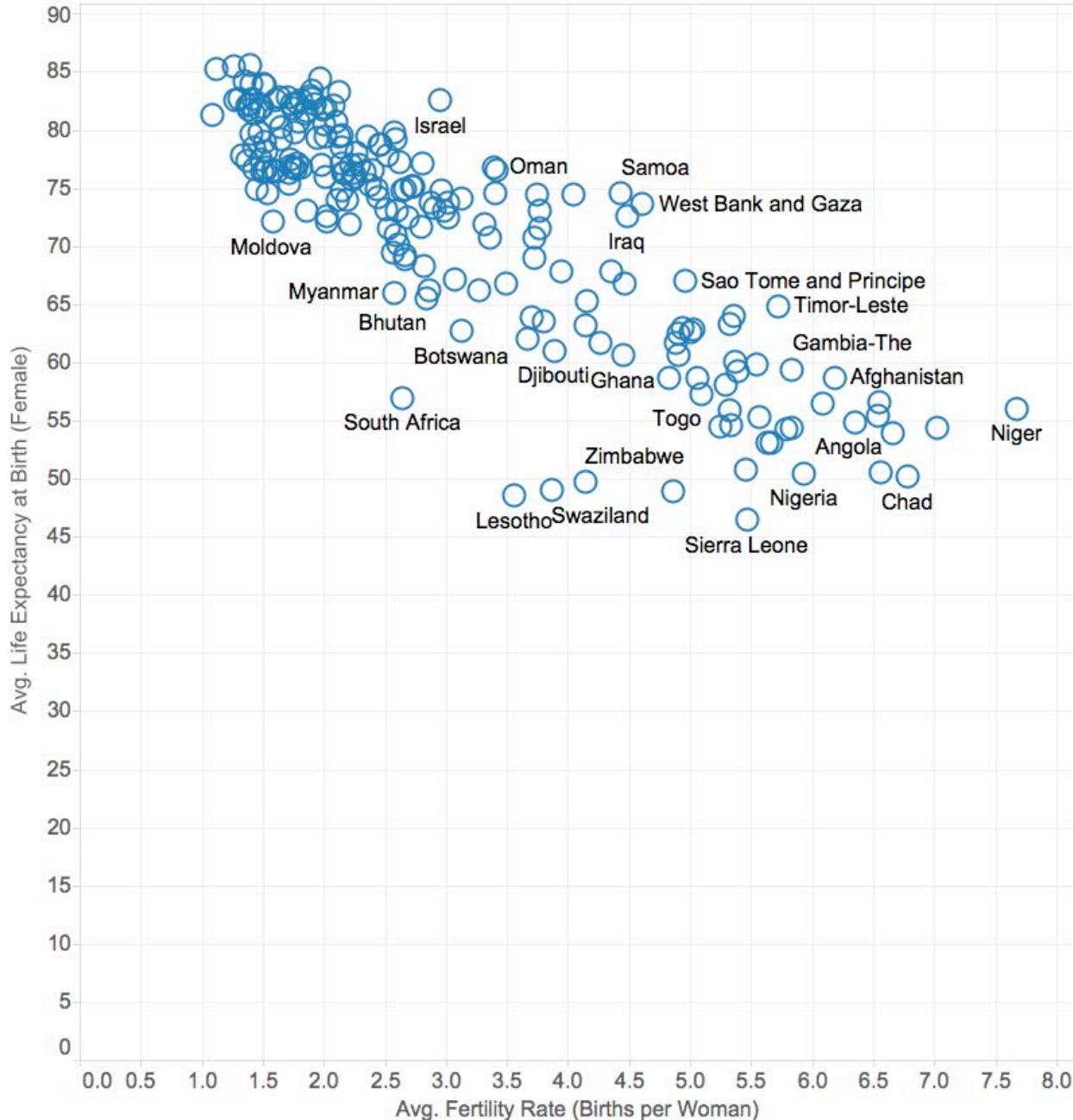


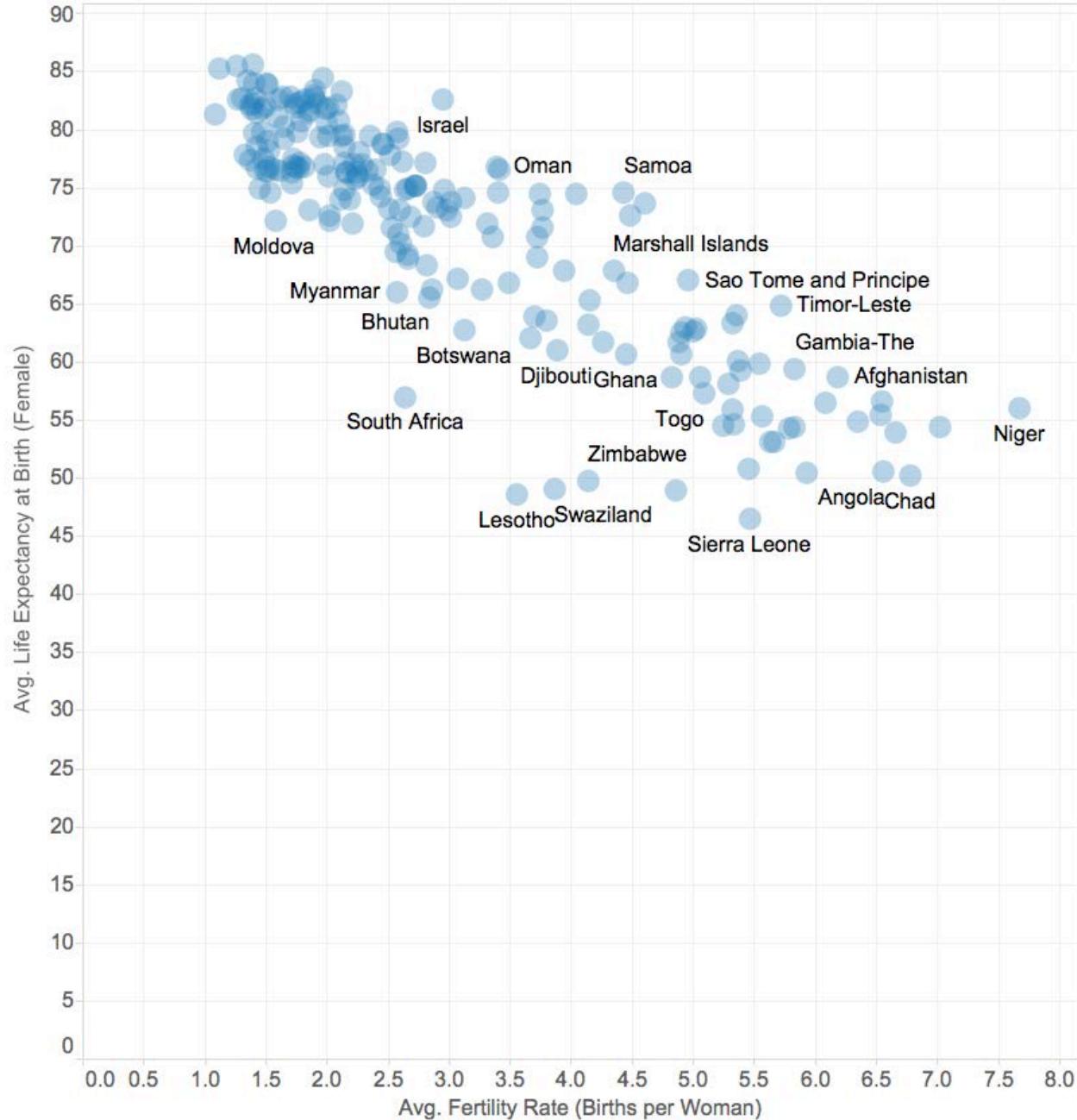


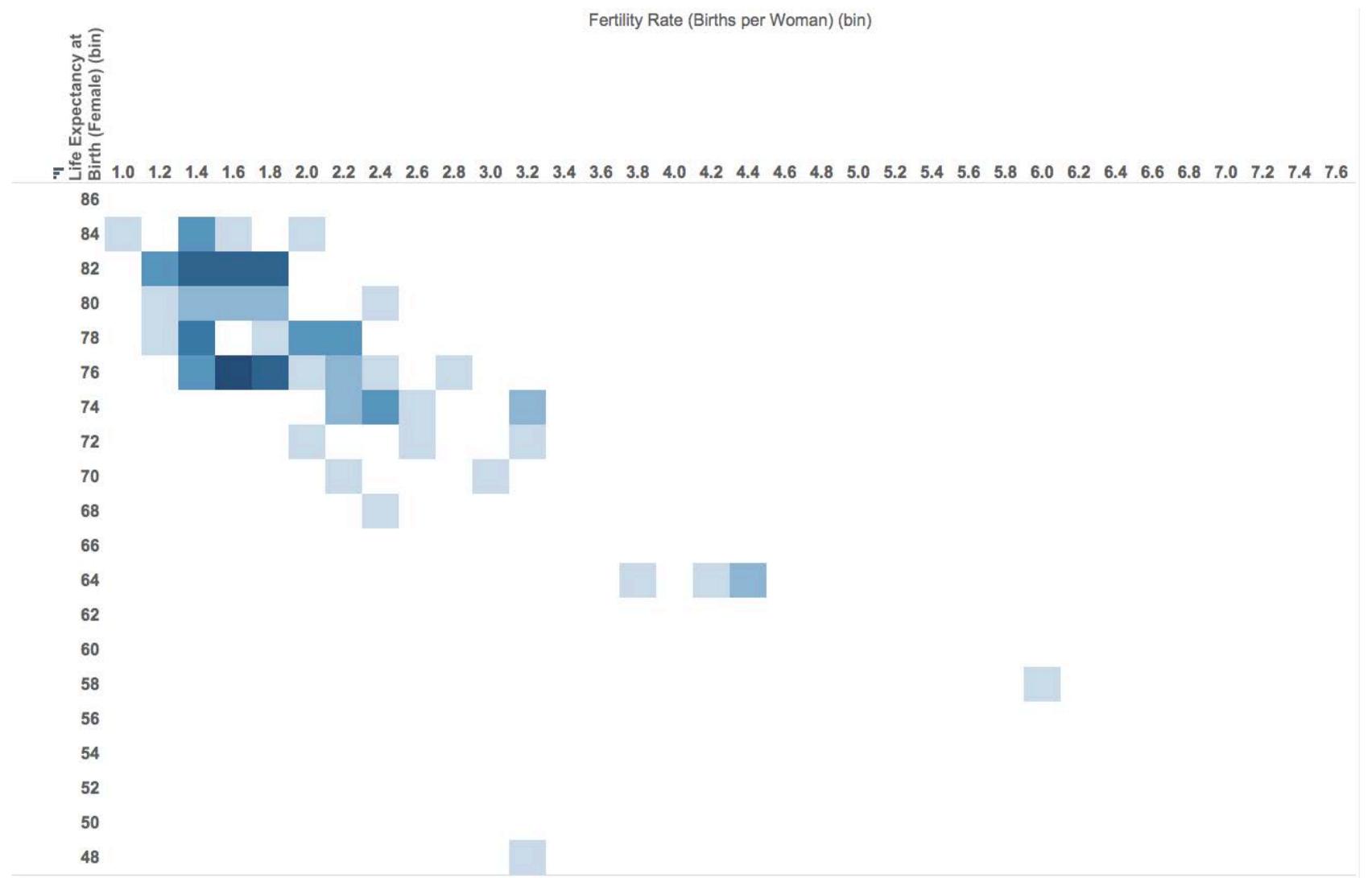












# QUESTIONS?