

IMD0105 - Special Issues in Information Technology VI

Introduction to Statistics III

Natal-RN
April 2017

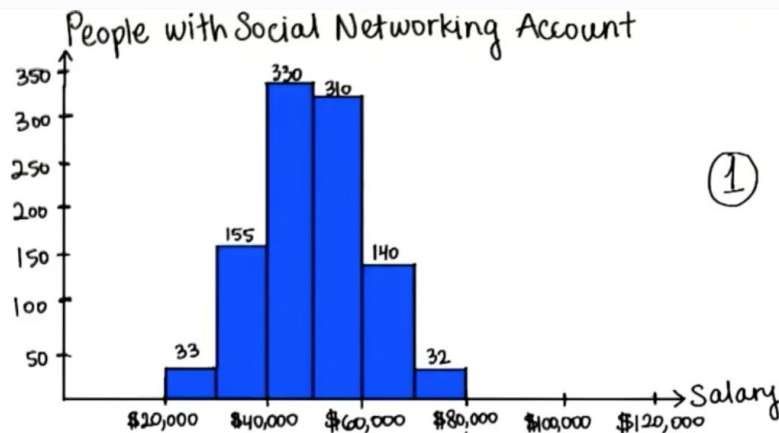


Agenda

- Quantify frequency
- Quatile
- Boxplot
- Measure variability (mean, variance, standard deviation)
- Z-Score
- Z-Table
- Facebook friends

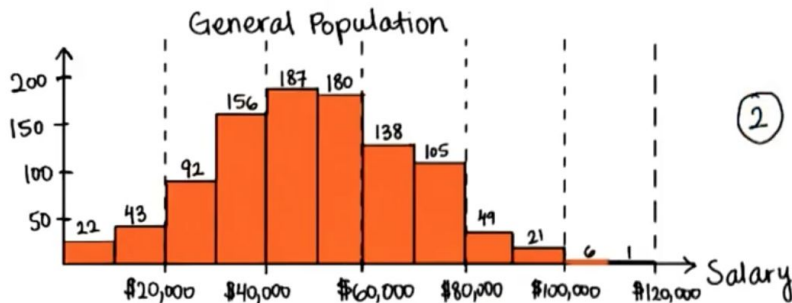
Previously on last class (...)

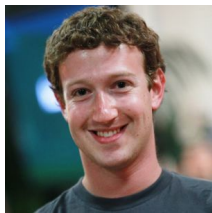
Social Networkers' Salaries



Should you get a networking account?

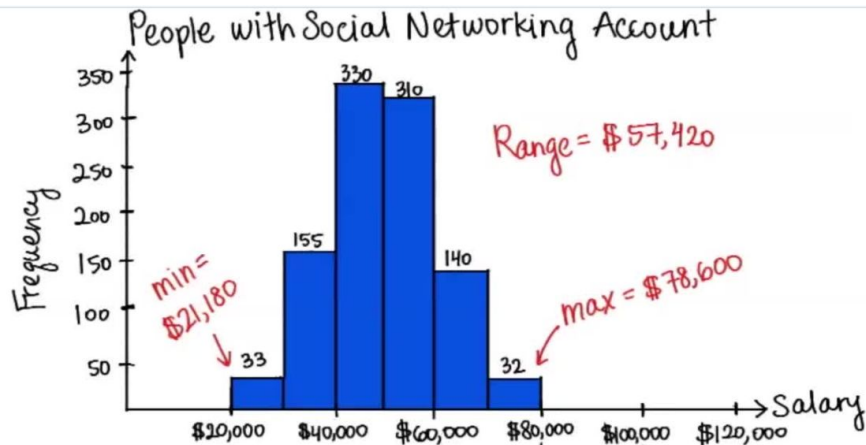
What's the difference between these two distributions?





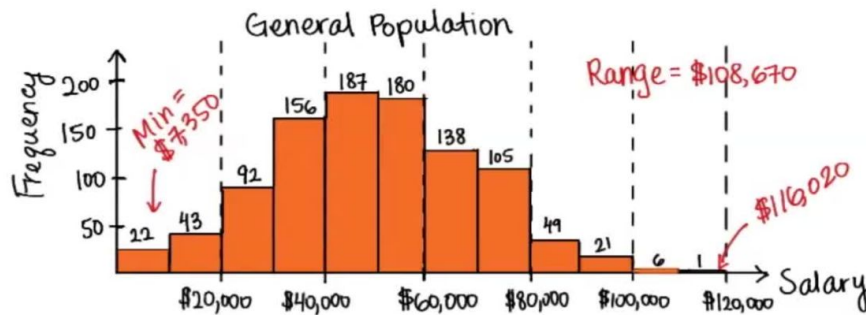
\$10M Salary

Quantify Spread



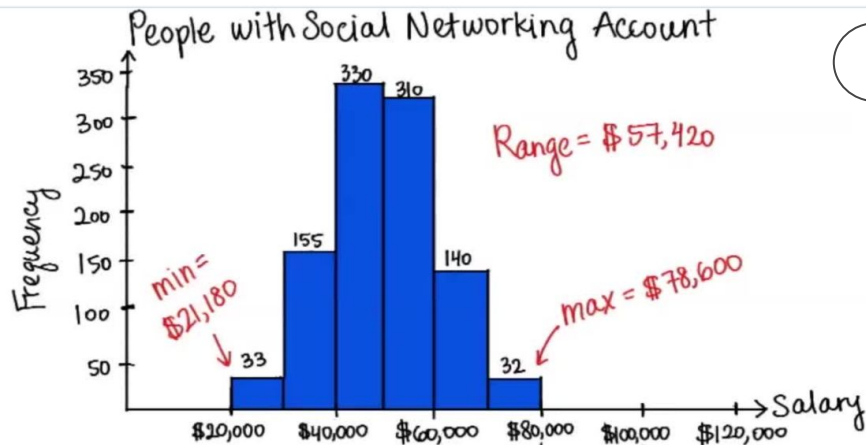
Range changes when we add new data to the dataset?

- 1) Always
- 2) Sometimes
- 3) Never



Quantify Spread (Quartile Q1 - Q3)

Chop of the tails



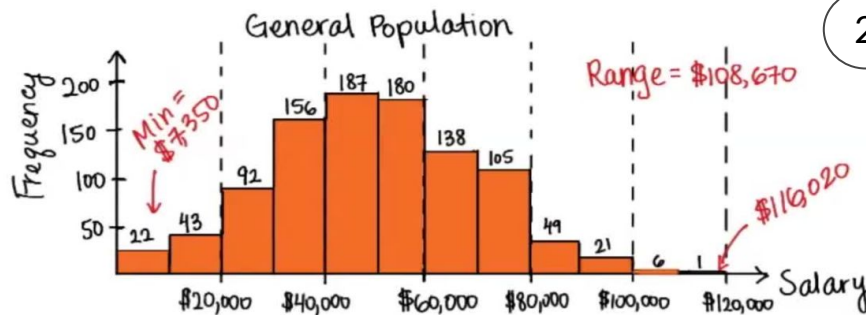
1

Sample 1

38,946
 43,420
 49,191 Q1
 50,430
 50,557 Q2 Median
 52,580
 53,595
 54,135 Q3
 60,181
 10,000,000

Sample 2

33,219
 36,254
 38,801 Q1
 46,335
 46,840 Q2 Median
 47,596
 55,130
 56,863 Q3
 78,070
 88,830



2

Q3 - Q1 = Interquartile range (IQR)

Sample 1

38,946

43,420

49,191 Q1

50,430

50,557 Q2 Mediana

52,580

53,595

54,135 Q3

60,181

10,000,000

$$Q3 - Q1 = IQR$$

$$4944$$

IQR

Quiz: True or False



Room: 4PSX6ZWHG

Sample 2

33,219

36,254

38,801 Q1

46,335

46,840 Q2 Mediana

47,596

55,130

56,863 Q3

78,070

88,830

$$Q3 - Q1 = IQR$$

$$18,062$$

What is an outlier?

Sample 1

38,946

43,420

49,191 Q1

50,430

50,557 Q2 Mediana

52,580

53,595

54,135 Q3

60,181

10,000,000

Q3 - Q1 = IQR
4944

What values do you think are outliers for this dataset?

- \$60,000
- \$80,000
- \$100,000
- \$200,000

Definition:

Outlier < $Q1 - 1.5 \times IQR$

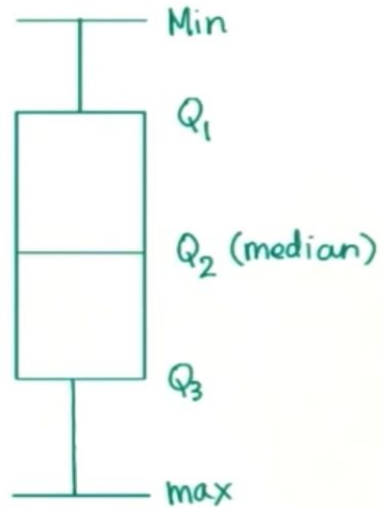
Outlier > $Q3 + 1.5 \times IQR$

41,775

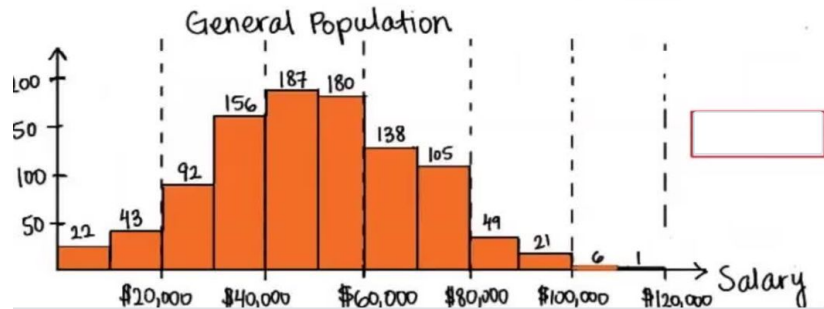
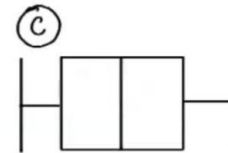
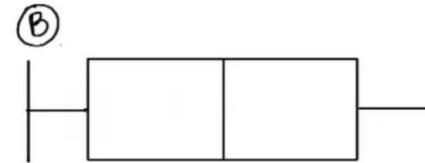
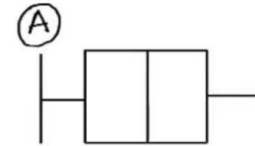
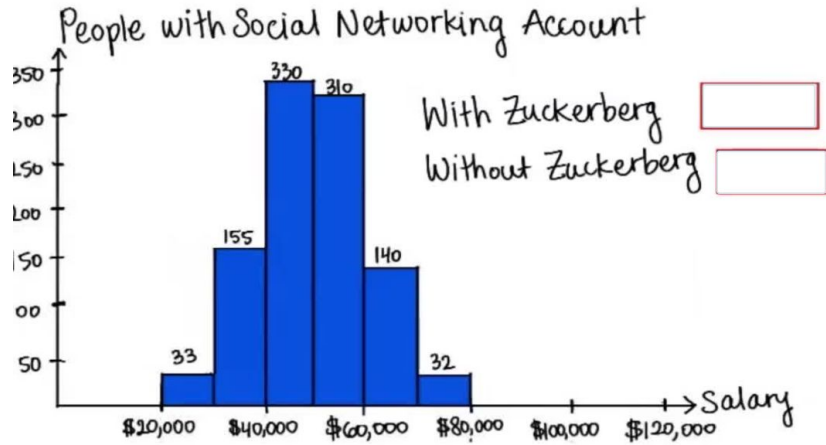
61,551

Match Boxplots

Boxplots

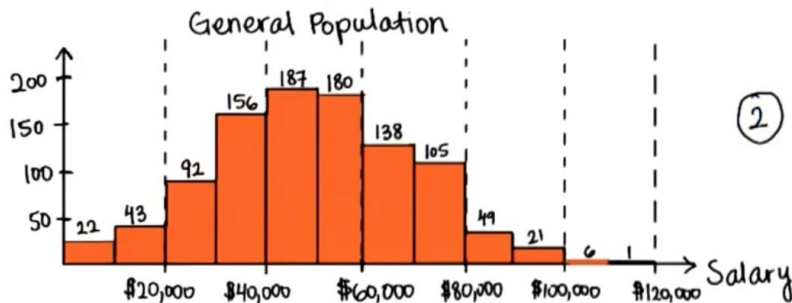
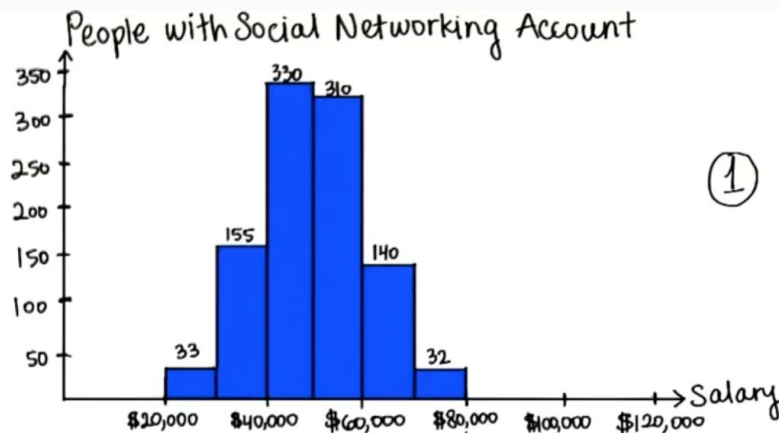


Quiz: Match Boxplots



Room:
4PSX6ZWHG

Problem with IQR



Will the mean always be between Q_1 and Q_3 ?

- Yes
- No

Coding



<https://goo.gl/ib1F2s>

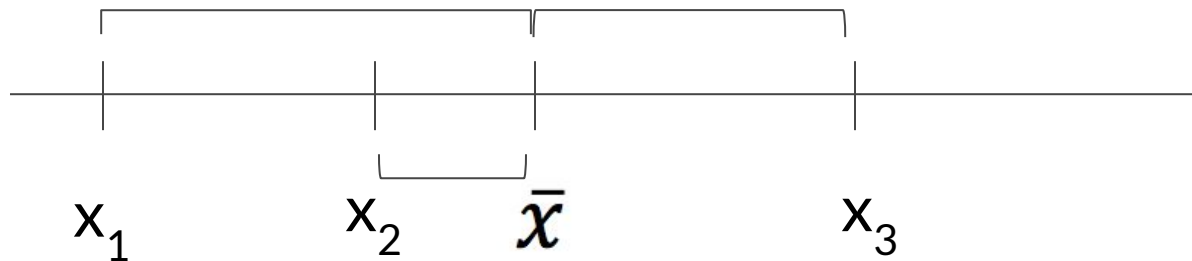
Measure Variability

We need one number that describe the spread of sample that takes all the data into account.

~~Range~~

~~IQR~~

Measure Variability (idea)



Variance

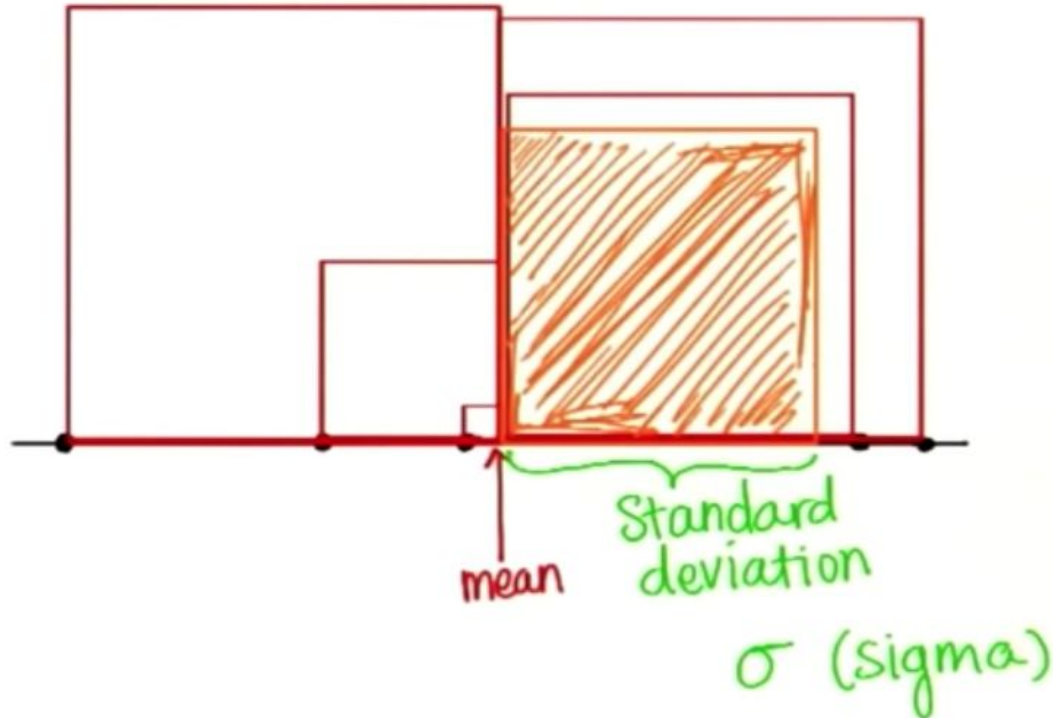
$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

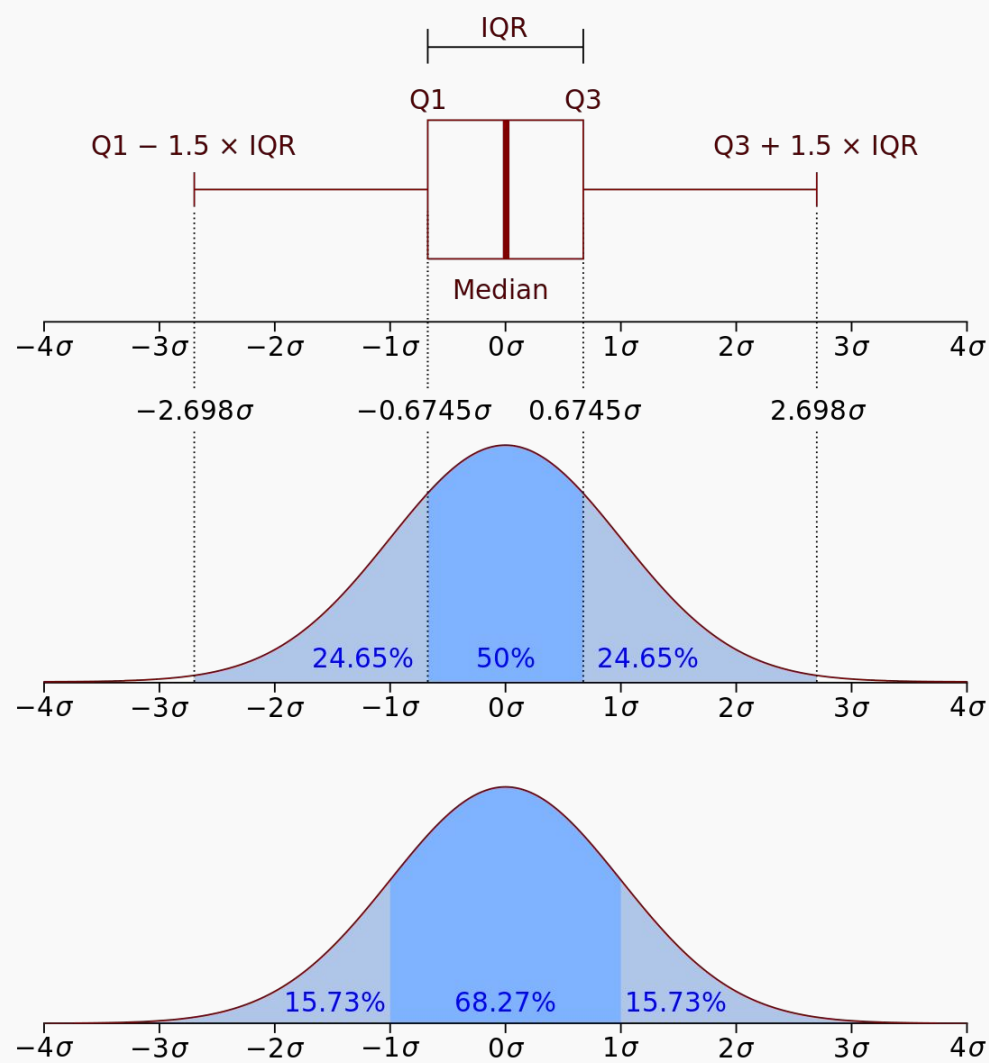
Standard
Deviation

σ

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

Measure Variability (idea)

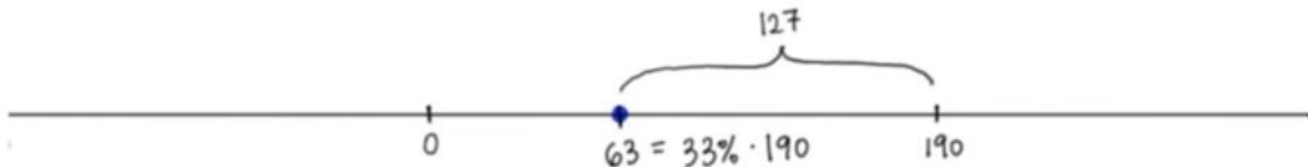




Quiz: Who is the most popular?

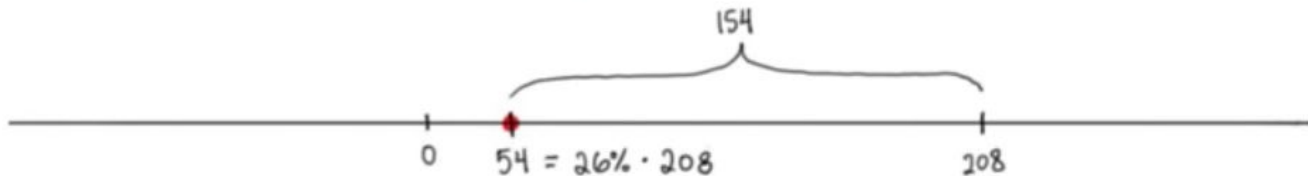
Javanildo

Facebook friends

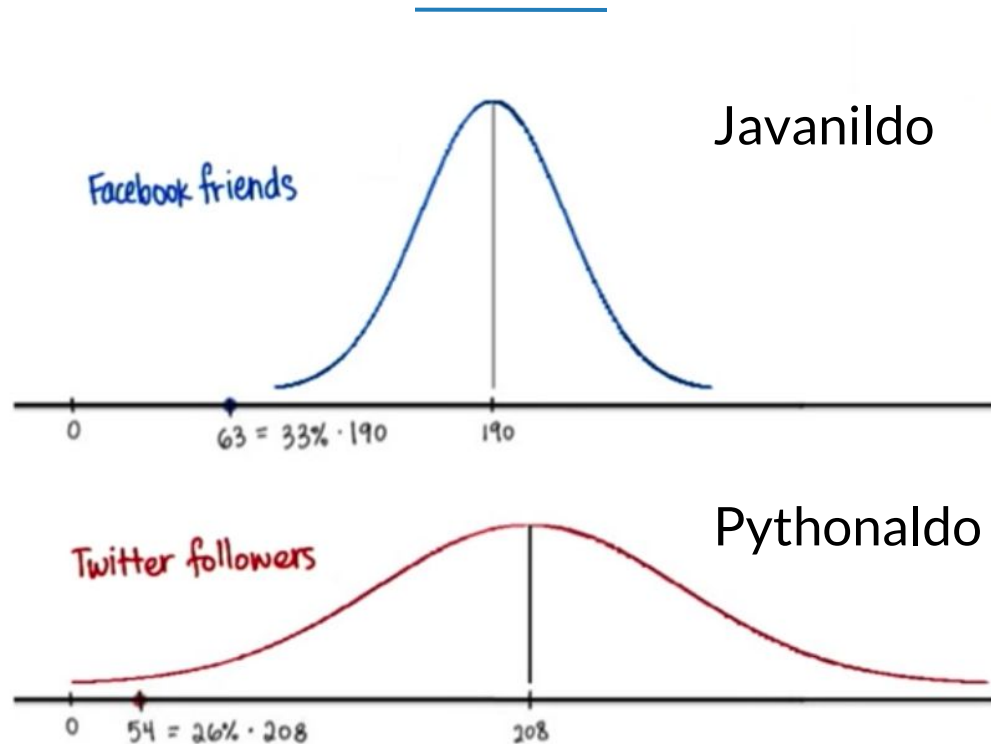


Pythonaldo

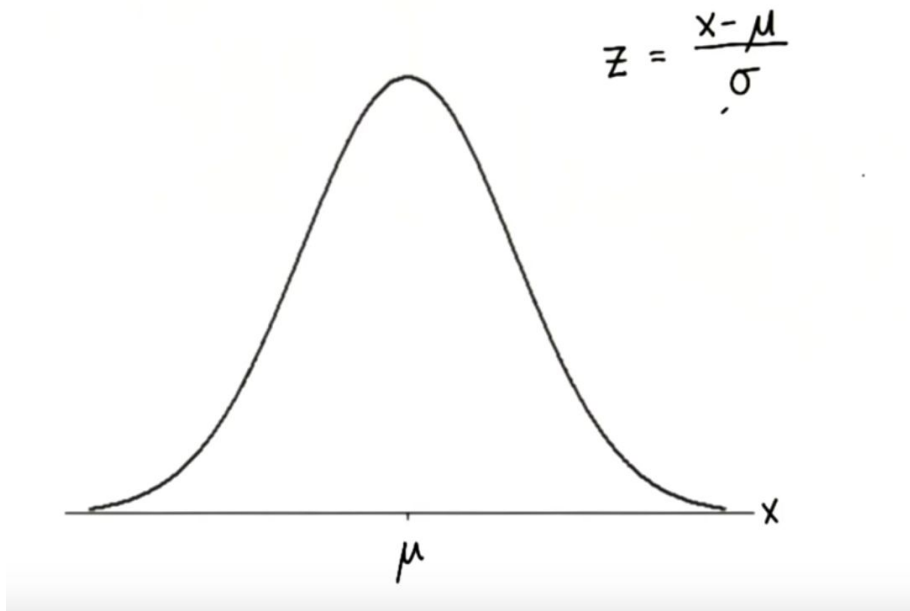
Twitter followers



Quiz: Who is the most popular?



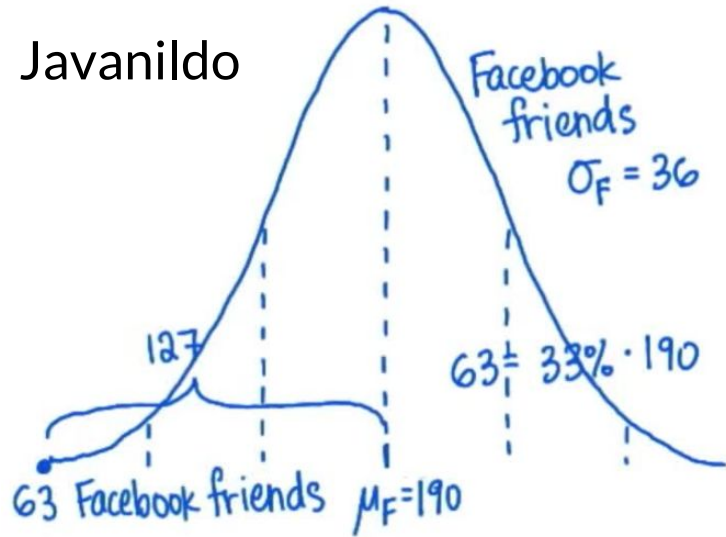
Z-Score



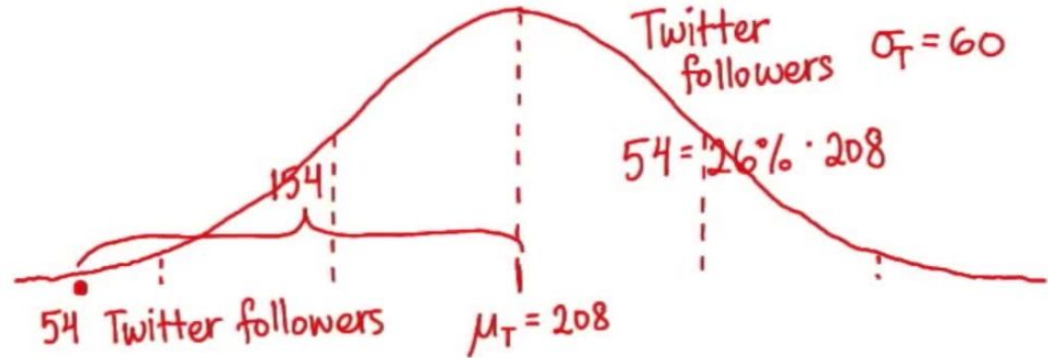
Quiz: Who is the most popular?

How many standard deviation is the number of friends from the mean?

Javanildo



Pythonaldo

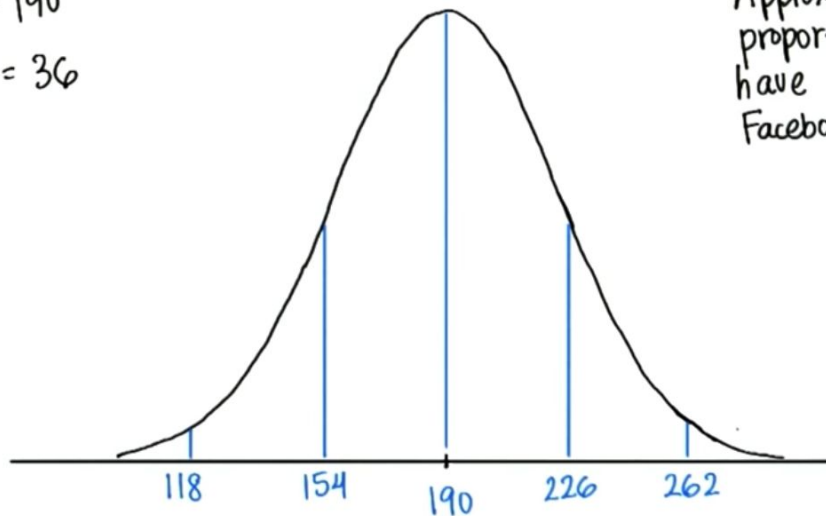


Quiz: Facebook friends (Table Z)

Distribution of Facebook friends

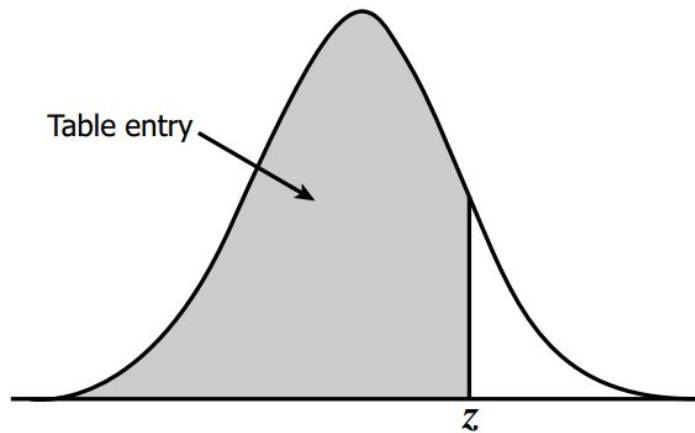
$$\mu = 190$$
$$\sigma = 36$$

Approximately what proportion of people have less than 240 Facebook friends?



Number of Facebook friends

Z-Table



<http://www.z-table.com/>

Facebook example:

$$\mu = 190$$

$$\sigma = 36$$

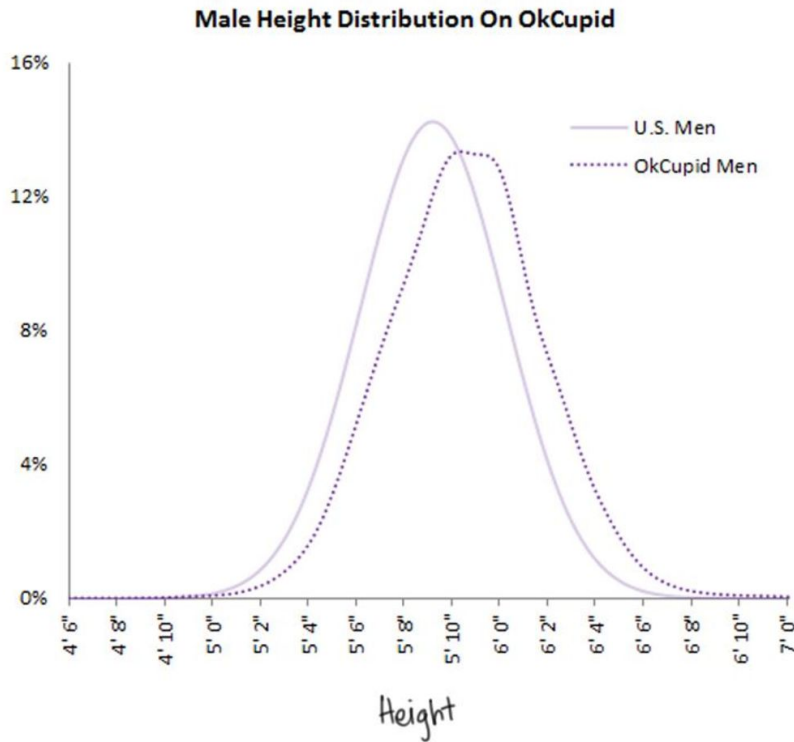
$$X_i = 240$$

Coding



<https://goo.gl/9opGD6>

Curious Quiz: OkCupid



The OkCupid blog shared this graph. What is going on here?