



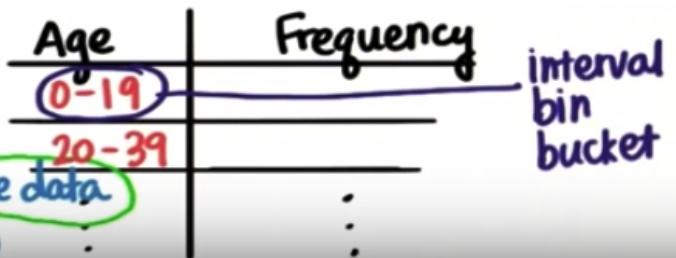
How many rows would you need in your table?

- o 50 (one for each student)
- o 66 (one for each age, 10-75)
- o 8-10 (easy to understand)
- It depends on how you group the data
- 2 (over 50 yrs, under 50 yrs)

Sample of student ages

15	19	18	14	13
27	16	65	15	31
22	15	24	22	51
24	20	45	22	33
24	27	18	66	15
18	39	10	30	13
19	28	53	28	65
30	20	21	20	18
20	23	18	41	52
75	19	63	14	18

n=50



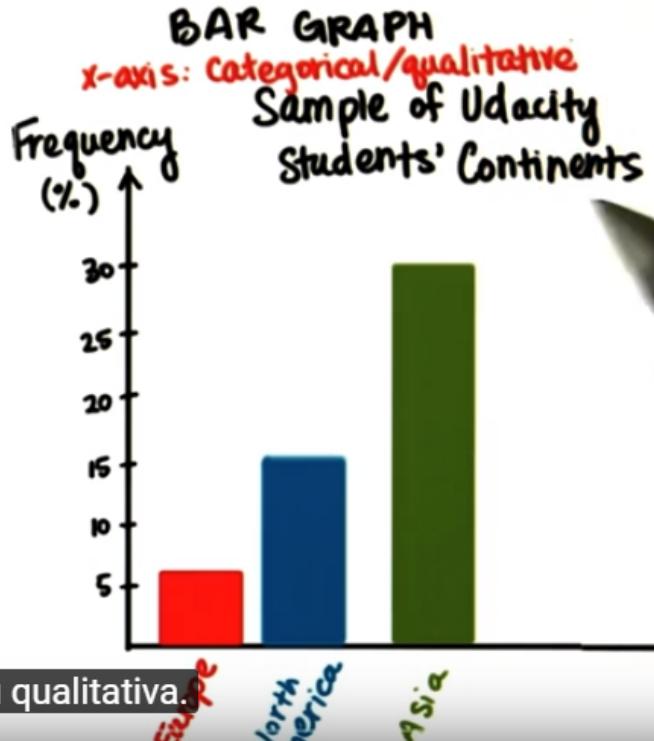
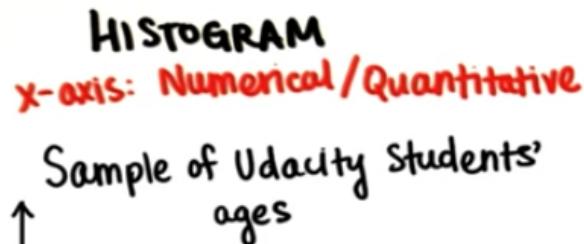
Sobre os eixos do grafico, a frequencia/variável dependente(resultado sempre fica no axes Y e a variável no axes X

Histograma

- Frequencia: Variável eixo Y
- Origem: Interseção dos eixos, coordenadas cartesianas (0,0)

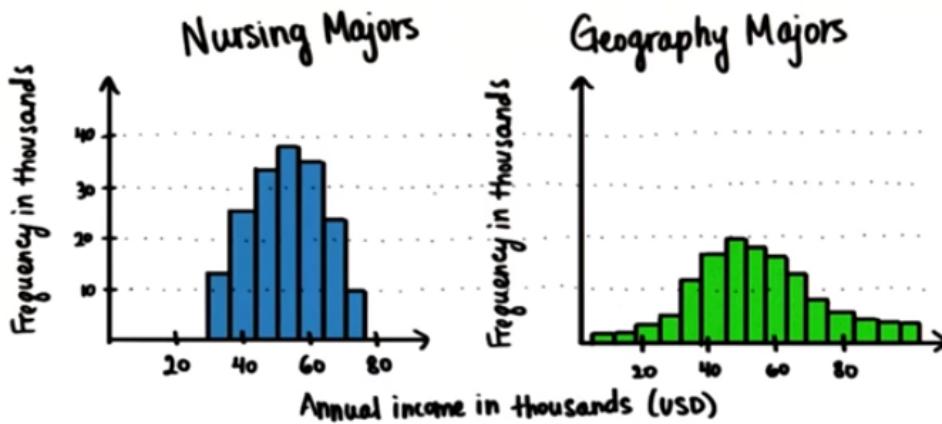
Em um histograma a variável do eixo X é numérica e quantitativa

Em um grafico de barras o eixo X é categórico e qualitativo



2005 United States
Income Distribution (Bottom 98%)
Each  equals 500,000 households





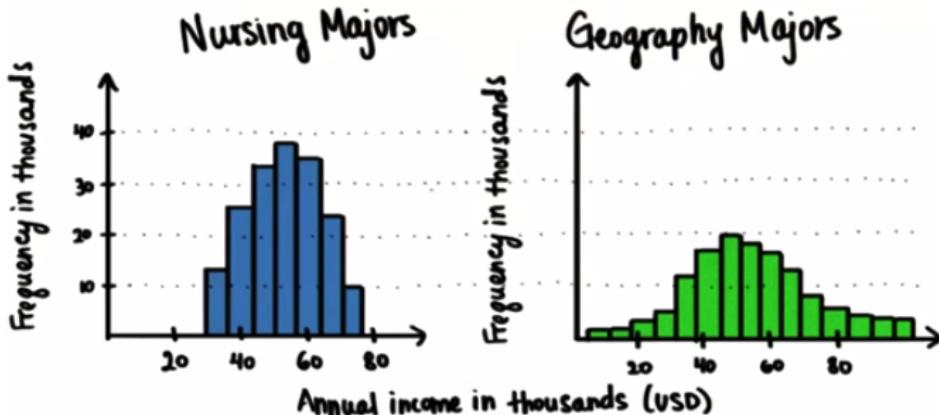
How would you choose one number (or a small range of numbers) that accurately represents the typical salary of nursing or geography majors?

The value at which frequency is highest

O valor onde a frequência é maior

O Average

↳ Biggest value on x-axis



How would you choose one number (or a small range of numbers) that accurately represents the typical salary of nursing or geography majors?

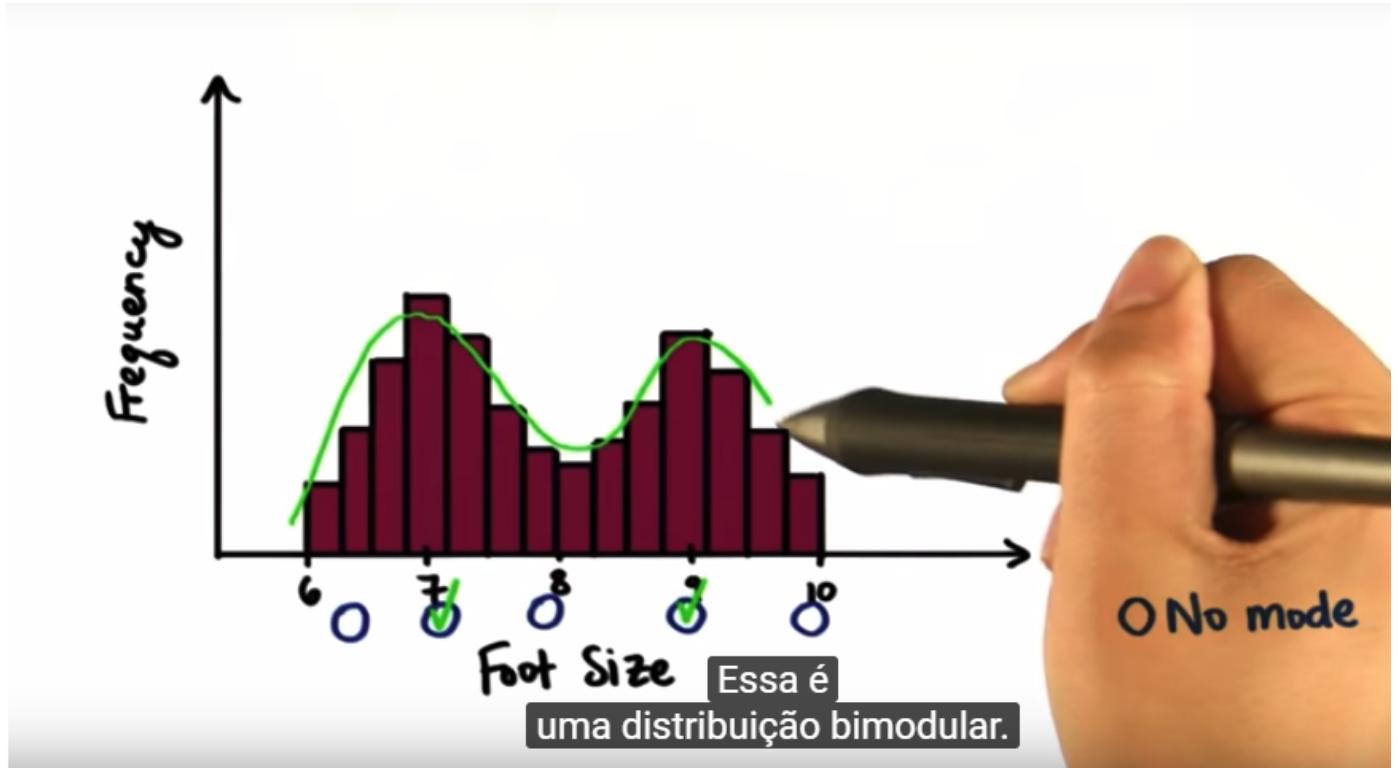
The value at which frequency is highest

O valor no meio da distribuição

O Average

↳ Biggest value on x-axis

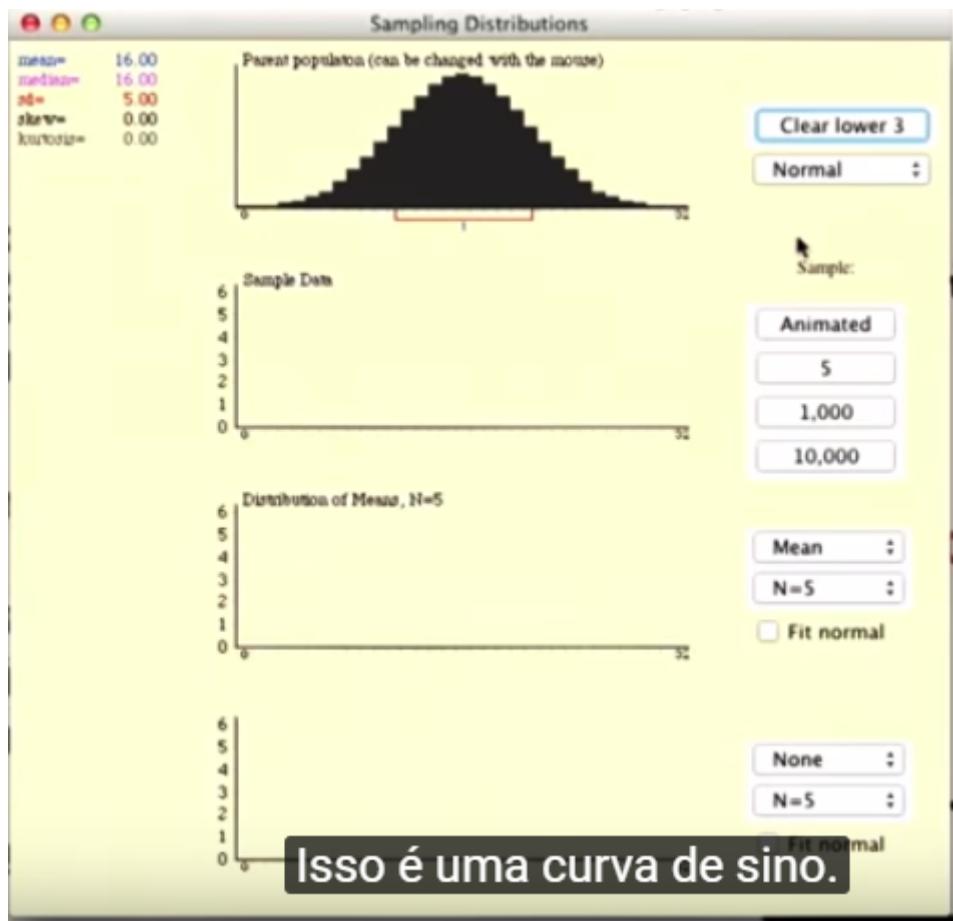
Quando a distribuição possui varias MODAS, neste caso chamados de **DISTRIBUIÇÃO UNIFORME**



BIMODULAR

Answer Remember, the mode occurs on the X-axis, so you are looking for whatever value has the highest frequency.

The numbers 7,000 and 1,000 are the actual frequencies. The mode, itself, is "Plain."



Nursing	Geography
\$58,350	\$48,670
\$63,120	\$57,320
\$44,640	\$38,150
\$56,380	\$41,290
\$72,250	\$53,160

$$\bar{x} = \frac{\sum x}{n}$$

Sum (Salary of geography majors)
(number of geography majors)

Σ n



interessante? Agora, lembre-se: isso é uma amostra, e você viu na Lição 1 que a

Nursing	Geography
\$58,350	\$48,670
\$63,120	\$57,320
\$44,640	\$38,150
\$56,380	\$41,290
\$72,250	\$53,160

$$\bar{x} = \frac{\sum x}{n}$$

Sum (Salary of geography majors)
(number of geography majors)

Σ n

$$\mu = \frac{\sum x}{N}$$



de toda a população, enquanto esse é o número da amostra. Agora, isso

Nursing	Geography
\$58,350	\$48,670 = x_1
\$63,120	\$57,320 = x_2
\$44,640	\$38,150 = x_3
\$56,380	\$41,290 = x_4
\$72,250	\$53,160 = x_5

$$\bar{x} = \frac{\sum x_i}{n}$$

Sum (Salary of geography majors)
(number of geography majors)

$$\mu = \frac{\Sigma x}{N}$$

$$= \frac{x_1 + x_2 + \dots + x_n}{n}$$

abcde seu raciocínio matemático em todas as aulas. Os símbolos são como o alfabeto.

Nursing	Geography
\$58,350	\$48,670
\$63,120	\$57,320
\$44,640	\$38,150
\$56,380	\$41,290
\$72,250	\$53,160
	\$500,000

median = \$58,350

mean = \$47,718
median = \$48,670

Mean = \$123,098
Median = \$50,915
↑ ROBUST

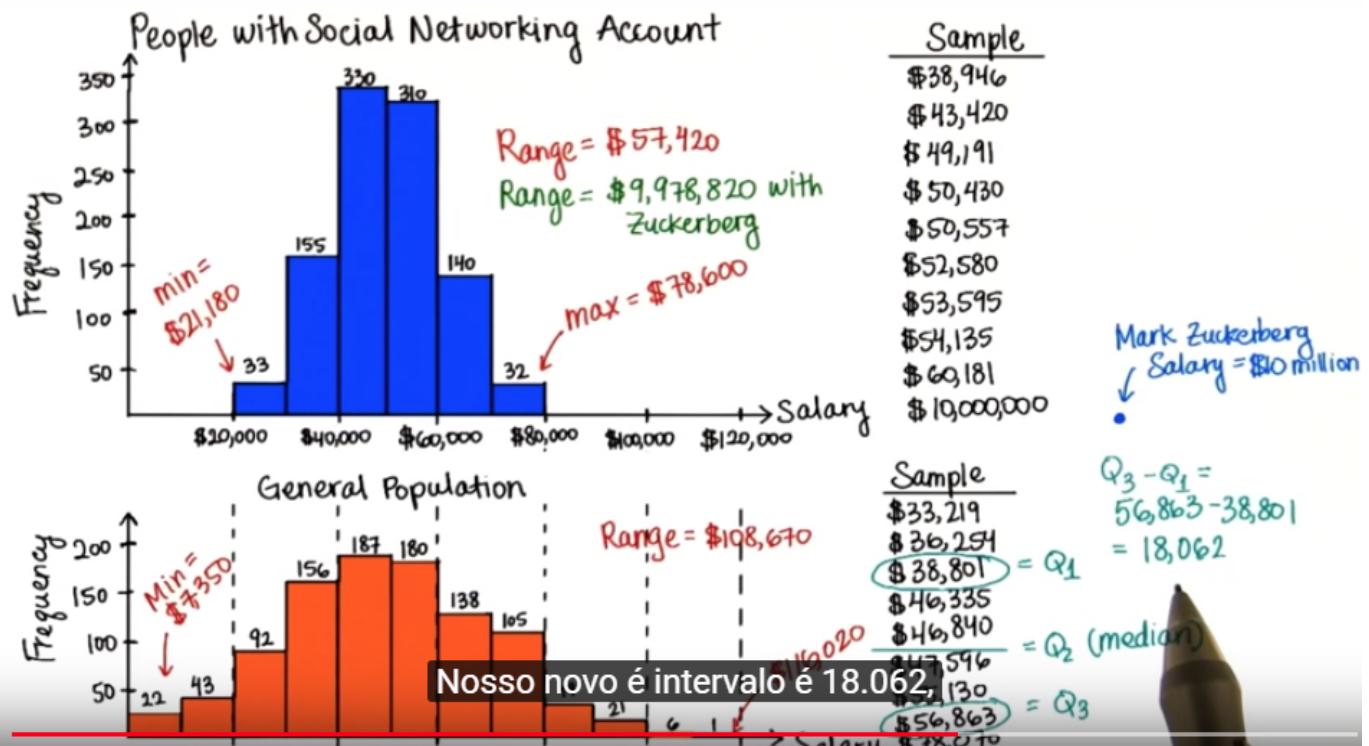
Na linguagem diária, robusto significa forte e resistente, o que faz sentido para sua

MEAN MEDIAN MODE
(Measures of center)

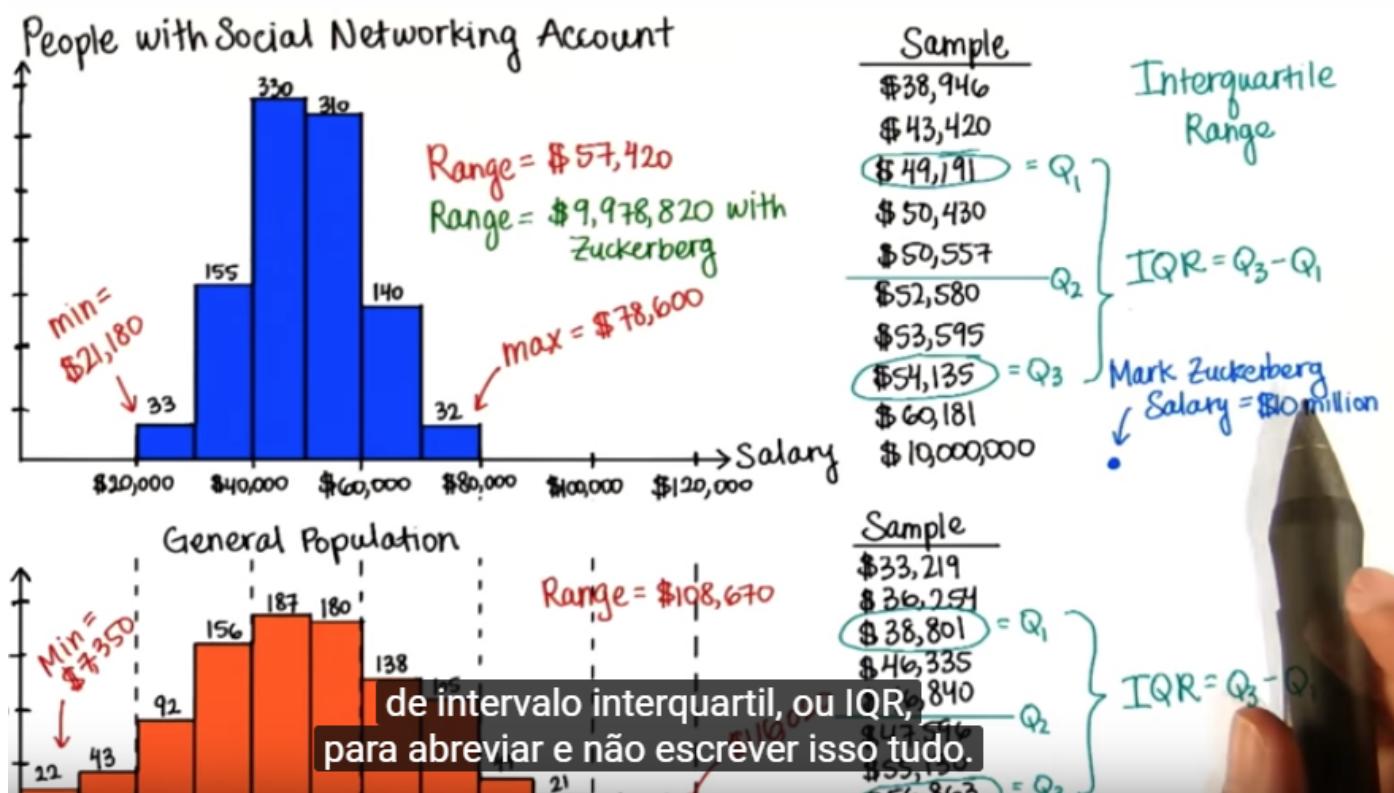
Como visto em aula, a média pode não descrever o centro devido a um outlier, às vezes a moda também não descreve o centro e a mediana não abrange todos os pontos de dados.

Muitas vezes a mediana pode ser mais adequada para trabalhar com gráficos altamente enviesados

Nesta etapa é calculado os quartis, para isso se divide os valores da amostra em 2 partes e depois cada valor mediano de cada conjunto são os Q1 = Primeiro quartil e Q3 = Terceiro quartil



Para calcular a nova distribuição faça $Q_3 - Q_1$, assim vc terá o que chamamos de intervalo interquartil



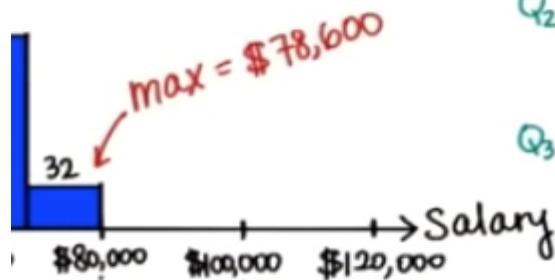
$$\text{Outlier} < Q_1 - 1.5 \text{ (IQR)}$$

$$> Q_3 + 1.5 \text{ (IQR)}$$

What values do you think are outliers for this dataset?

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

Working Account



Sample

\$38,946

\$43,420

Q_1 \$49,191

\$50,430

\$50,557

Q_2 \$52,580

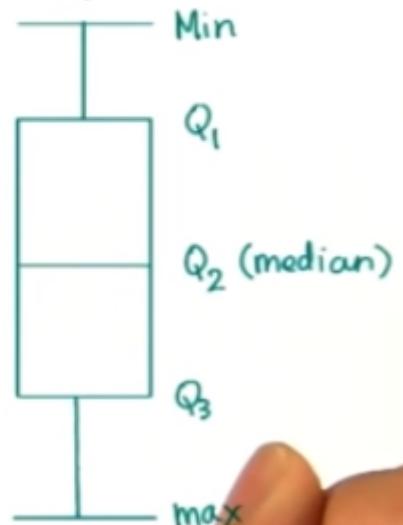
\$53,595

Q_3 \$54,135

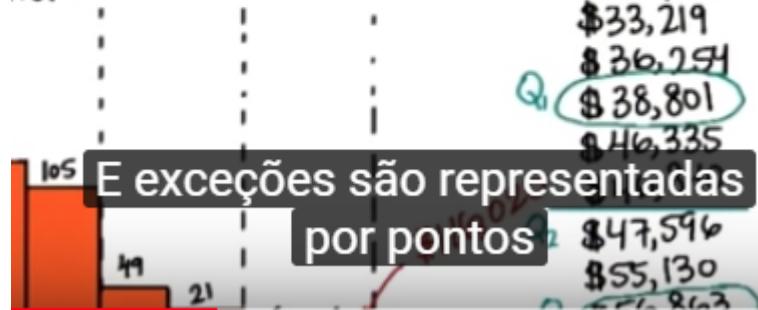
\$60,181

\$100,000,000

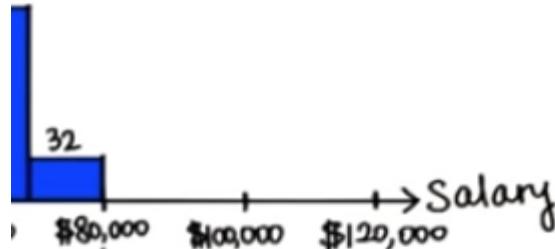
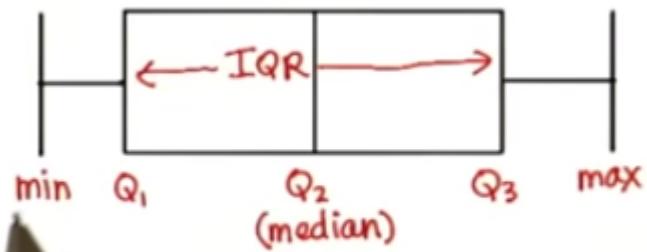
Boxplots



Median



Working Account



non



Sample

\$33,219
\$36,254
\$38,801
\$46,335
\$46,840
\$47,596
\$55,130
\$56,863
\$78,070
\$88,830

Deviation from Mean ($x_i - \bar{x}$)

-19,574.80
-16,539.80
-13,992.80
-6,458.80
-5953.80
-5197.80
2336.20
4069.20
25,276.20
36,036.20

$$\bar{x} = \frac{\sum x}{10} = \$52,793.80$$

ao invés de fazer o oposto,
 x barra menos x_i .

Sample	Deviation from Mean ($x_i - \bar{x}$)	Absolute Deviations $ x_i - \bar{x} $
\$33,219	-19,574.80	19,574.80
\$36,254	-16,539.80	16,539.80
\$38,801	-13,992.80	13,992.80
\$46,335	-6,458.80	6,458.80
\$46,840	-5953.80	5953.80
\$47,596	-5197.80	5197.80
\$55,130	2336.20	2336.20
\$56,863	4069.20	4069.20
\$78,070	25,276.20	25,276.20
\$88,830	36,036.20	36,036.20

$$\bar{x} = \frac{\sum x}{10} = \$52,793.80$$

A pode descrever o que você fez para achar a média do desvio padrão.

$\frac{|\sum x_i|}{n}$

$\frac{\sum |x_i - \bar{x}|}{n}$

$\sum \left(\frac{|x_i - \bar{x}|}{n} \right)$

$\frac{\sum (|\bar{x} - x_i|)}{n}$

$\frac{|x_i - \bar{x}|}{n}$

Sample	Deviation from Mean ($x_i - \bar{x}$)	Squared Deviations
\$33,219	-19,574.80	383172795
\$36,254	-16,539.80	273,564,984
\$38,801	-13,992.80	195,798,452
\$46,335	-6,458.80	41,716,097
\$46,840	-5953.80	354,477,34
\$47,596	-5197.80	27,017,124
\$55,130	2336.20	5,457,830
\$56,863	4069.20	16,558,389
\$78,070	25,276.20	638,886,286
\$88,830	36,036.20	1,298,607,710

$$\bar{x} = \frac{\sum x}{10} = \$52,793.80$$

$$\text{Avg. dev.} = \frac{\sum (x_i - \bar{x})}{10} = 0$$

SS (sum of squares)

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

A variancia é a soma dos desvios ao quadrado dividido pelo total de itens (média)

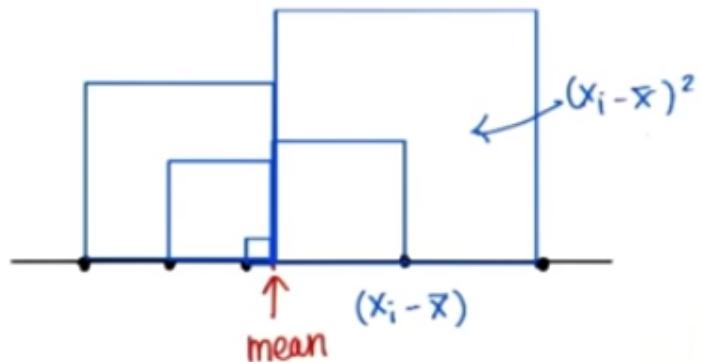
Sample	from Mean ($x_i - \bar{x}$)	Deviations
\$33,219	-19,574.80	383172795
\$36,254	-16,539.80	273,564,984
\$38,801	-13,992.80	195,798,452
\$46,335	-6,458.80	41,716,097
\$46,840	-5953.80	35,447,734
\$47,596	-5197.80	27,017,124
\$55,130	2336.20	5457,830
\$56,863	4069.20	16,558,389
\$78,070	25,276.20	638886,286
\$88,830	36,036.20	1298607,710
$\bar{x} = \frac{\sum x}{10} =$	Avg. dev. = $\frac{\sum (x_i - \bar{x})}{n}$	Avg. squared dev = $\frac{\sum (x_i - \bar{x})^2}{n}$
		Se chama variância.

SS (sum of squares)
 $= \sum (x_i - \bar{x})^2$

VARIANCE

Deviation from Mean ($x_i - \bar{x}$)	Squared Deviations
-19,574.80	383172795
-16,539.80	273,564,984
-13,992.80	195,798,452
-6,458.80	41,716,097
-5953.80	35,447,734
-5197.80	27,017,124
2336.20	5457,830
4069.20	16,558,389
25,276.20	638886,286
36,036.20	1298607,710

Avg. dev = é cada área dev /
 $\sum (x_i - \bar{x})^2 = \text{é } X_i \text{ menos } X\text{-linha ao quadrado.}$



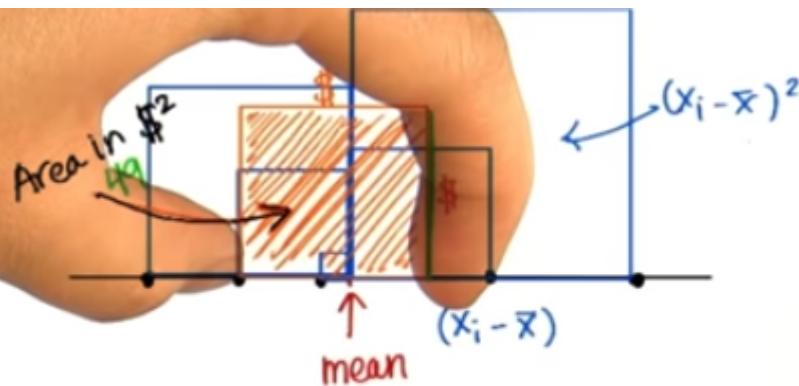
$$SS = \sum (\text{area of each square}) =$$

VARIANCE

	Squared Deviations
80	383172795
80	273,564,984
80	195,798,452
80	41,716,097
80	35447734
30	27,017,124
20	5457830
20	16558,389
20	638886286
20	1298607710

Avg. squared dev

Isso se chama desvio padrão.



How do we convert the dimensions of the average square back to one dimension (\$) from two (\$^2)?

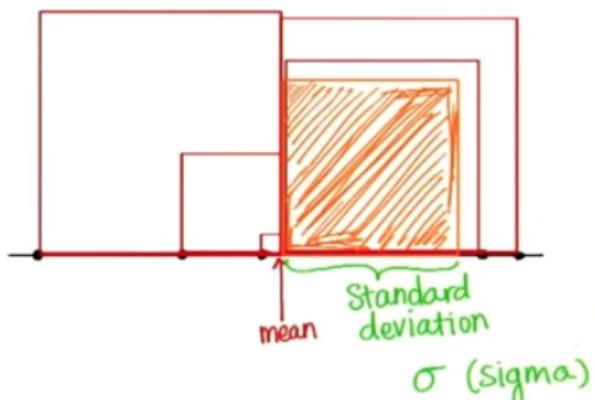
Take the square root

- o Divide by 2
- o Subtract one \$

VARIANCE

Desvio Padrão = Variância ao quadrado

Sample	Deviation from Mean $(x_i - \bar{x})$	Squared Deviations
\$33,219	-19,574.80	383172795
\$36,254	-16,539.80	273,564,984
\$38,801	-13,992.80	195,798,452
\$46,335	-6,458.80	41,716,097
\$46,840	-5953.80	35447734
\$47,596	-5197.80	27,017,124
\$55,130	2336.20	5457830
\$56,863	4069.20	16558,389
\$78,070	25,276.20	638886286
\$88,830	36,036.20	1298607710

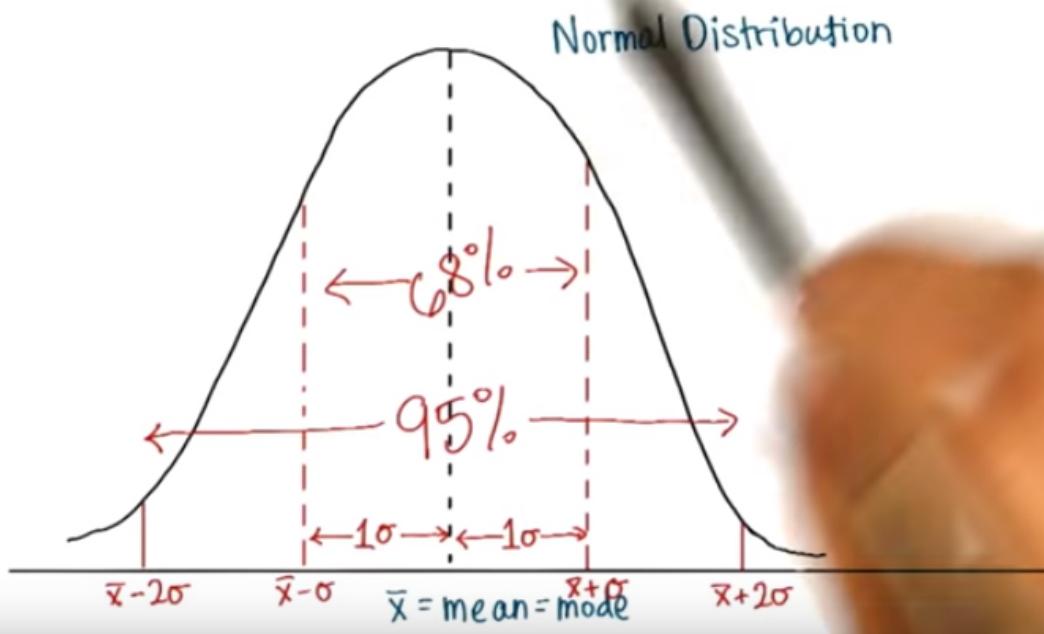


VARIANCE

σ (sigma)

$\bar{x} = \frac{\sum x}{n}$ Avg. dev. Avg. squared dev
comum de dispersão, e o símbolo é o sigma grego minúsculo. Lembre-se de que

What's so great about the standard deviation anyway?



18

20

18

23

22

Sample
 $n = 9$

21

Population

$$\mu = 18.97$$

$$\sigma = 5.99$$

Bessel's Correction

$$\text{Standard deviation} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

onde em vez de dividir por n ,
dividimos por n menos 1.

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

Você viu que representamos este desvio padrão corrigido com o s minúsculo.
Quando

Sample standard deviation

$$S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} \approx \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

$$(n=5) \quad \bar{x} = 3$$

5

4

2

1

1

4

0

9

7

16

squared
deviations

$$\sum (x_i - \bar{x})^2 = 34 \text{ (sum of squares)}$$

$$\text{standard deviation of sample} = \sqrt{\frac{34}{5}}$$

$$\text{standard deviation of population} \approx \sqrt{\frac{34}{4}}$$

Quando o desvio padrão precisa estar aproximado ao valor da população faça o n-1 para aumentar a distribuição, mas quando não houver esta necessidade faça a formula padrão

Sample standard deviation

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} \approx \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

(n=5) $\bar{x} = 3$

5

4

2

1

1

4

0

9

7

16

} squared deviations

$$\sum (x_i - \bar{x})^2 = 34 \text{ (sum of squares)}$$

$$\text{Standard deviation of sample} = \sqrt{\frac{34}{5}}$$

Standard deviation

Sample standard dev.
(s)

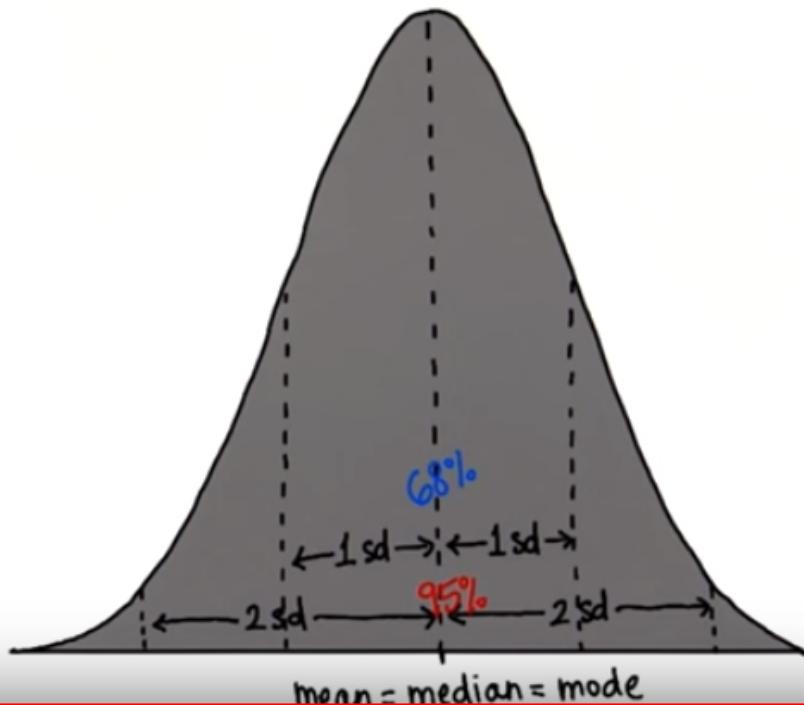
$$\approx \sqrt{\frac{34}{4}}$$

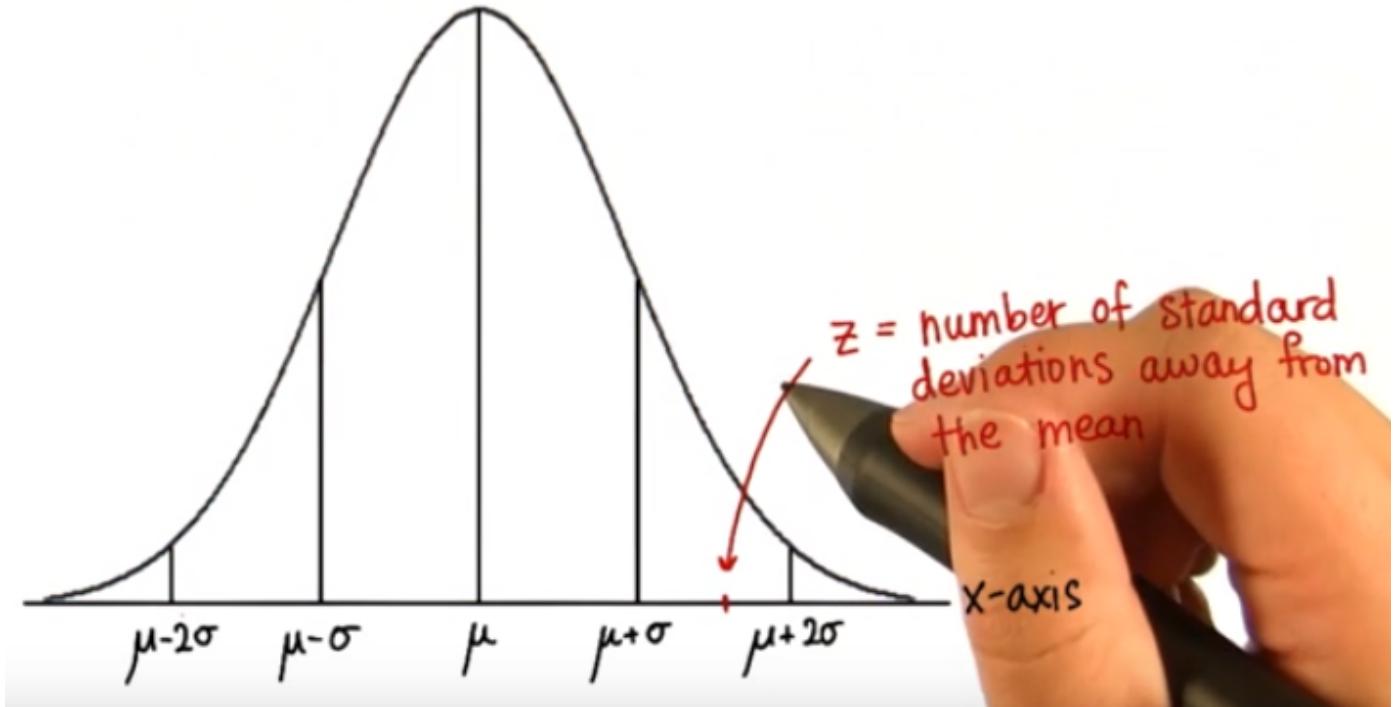
s minúsculo. Então, lembre-se: se você receber uma amostra e precisar

Modelo para calcular um valor % dentro da amostra

Área sempre será 1

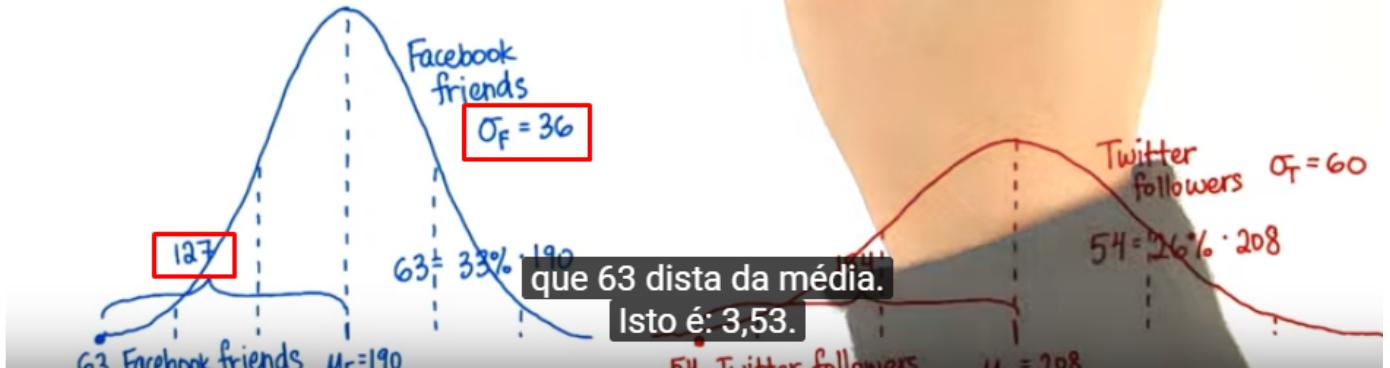
Normal Distribution





How many standard deviations is Katie's number of Facebook friends from the mean number of Facebook friends?

$$\frac{127}{36} = 3.53$$



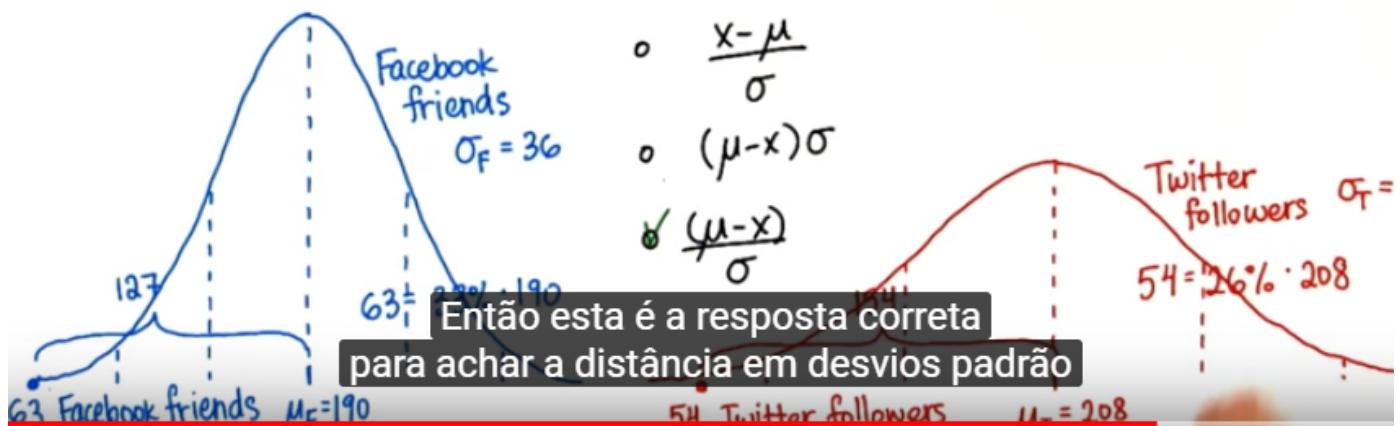
If Andy only uses Twitter and I only use Facebook,
can we necessarily say that Andy is more unpopular than me?

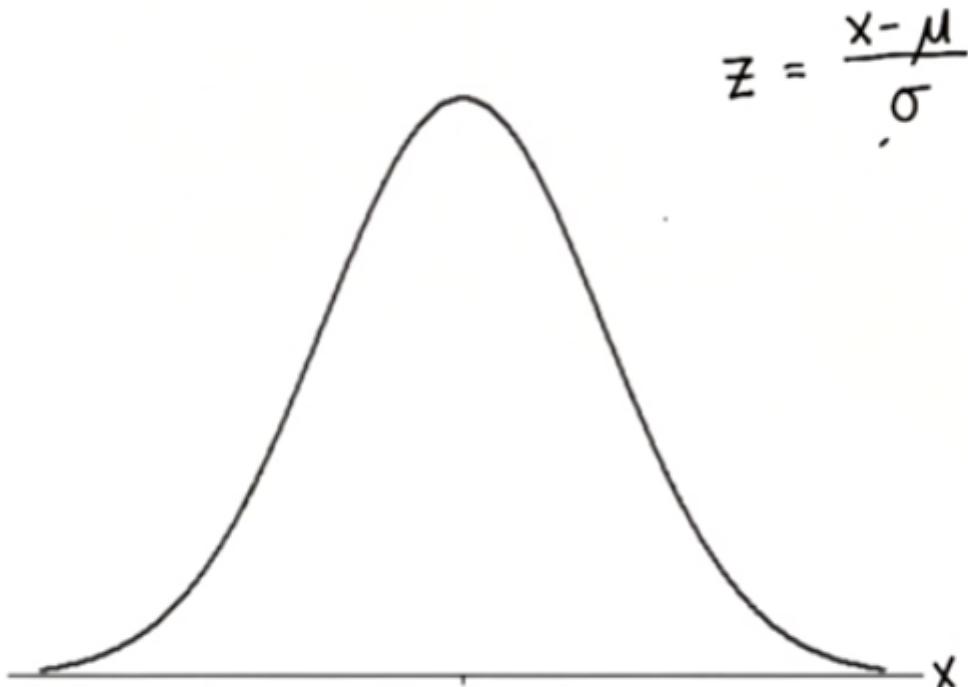
- Yes Why?
- No

Isso se chama padronizar as distribuições.

What formula describes what you did to find the number of standard deviations each value (let's call it x) is from the mean?

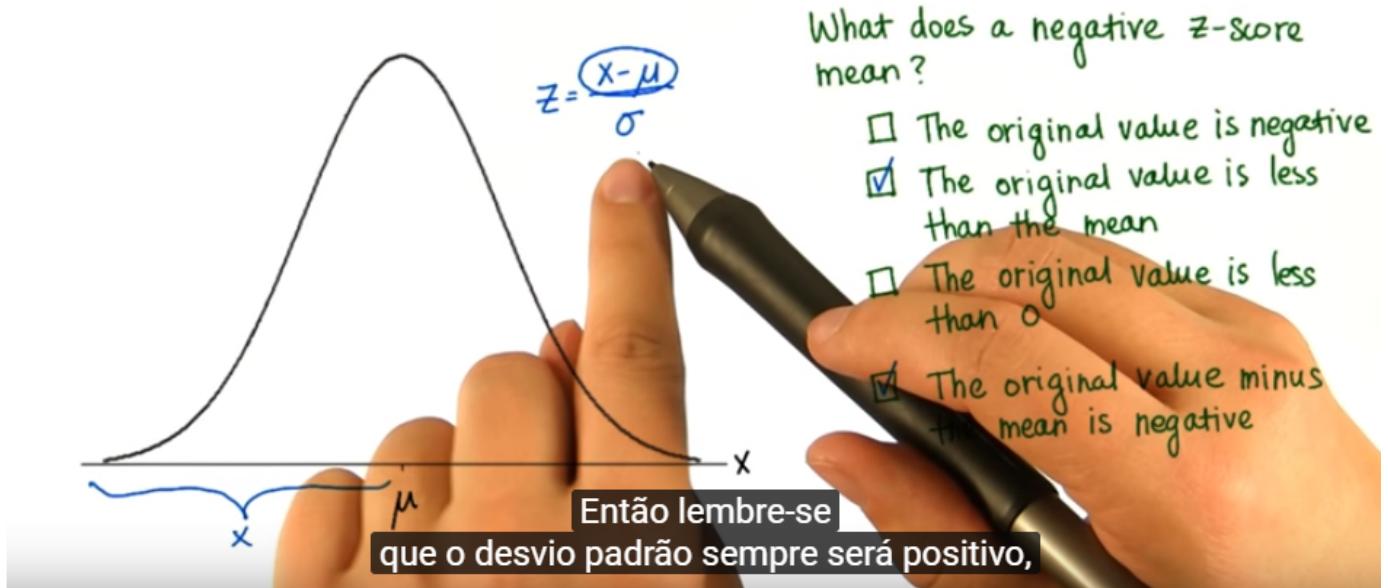
- $\mu - x$
- $x - \mu$
- $(x - \mu)\sigma$
- $\frac{x - \mu}{\sigma}$
- $\frac{(\mu - x)}{\sigma}$



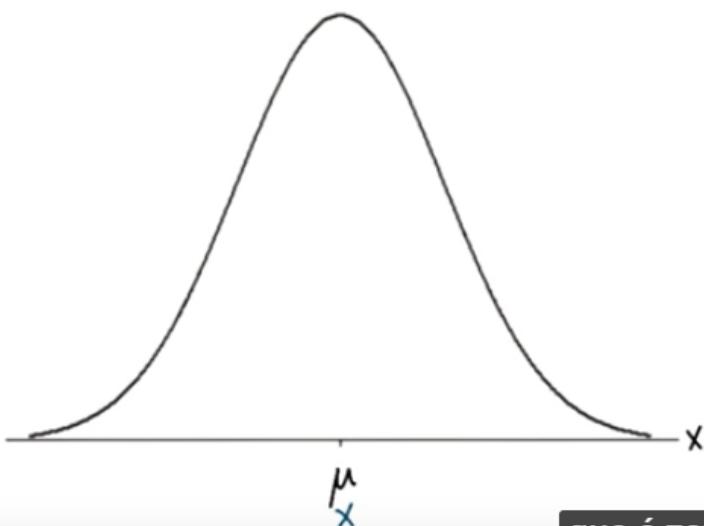


O Escore Z é, basicamente,
o número de desvios padrão

Sempre será positivo porque é elevado ao quadrado



It's all from formula of Z. $z = \frac{x - \mu}{\sigma}$. when $x = \mu$ we get $Z=0$, so mean of Standardized distribution is 0. You can find in literature the phrase: "Distribution with mean = 0 and standard deviation = 1" So it's Standardized distribution



If we standardize a distribution by converting every value to a z-score, what will be the new mean of this standardized distribution?

$$\mu = 0$$

$$z = \frac{x - \mu}{\sigma}$$

$$z = \frac{\mu - \mu}{\sigma} = \frac{0}{\sigma} = 0$$



que é zero.

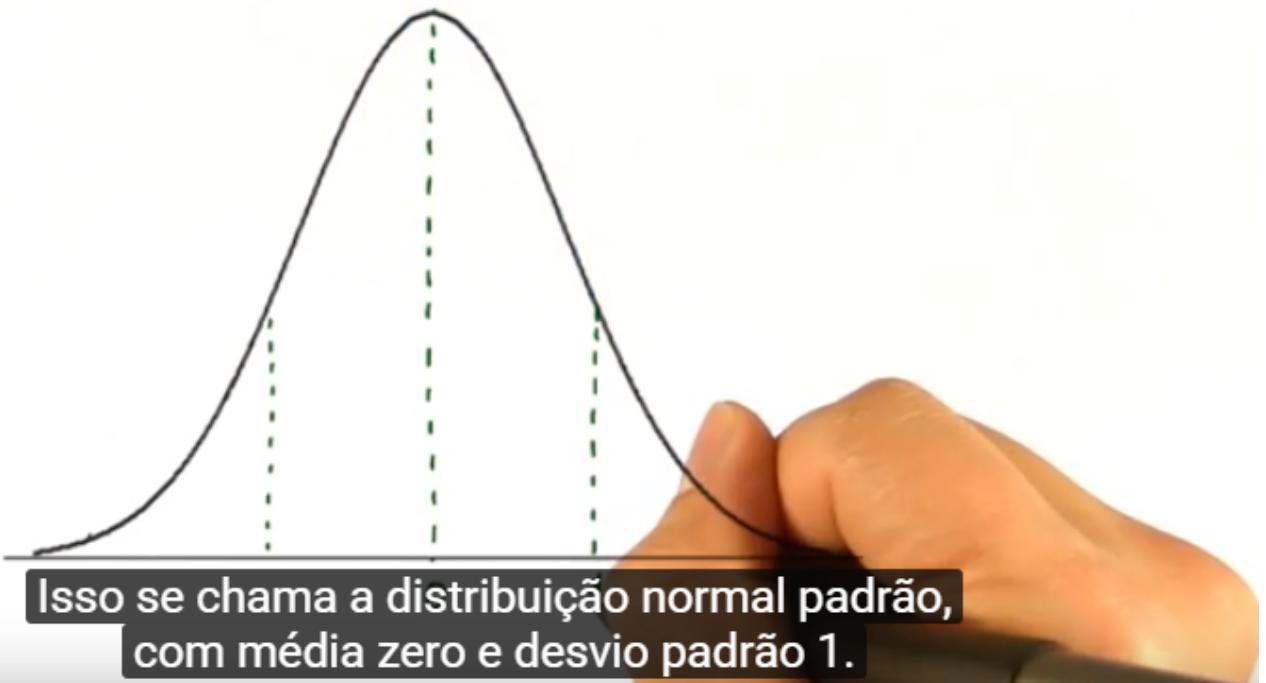
If we standardize a distribution by converting every value to a z-score, what will be the new standard deviation of this standardized distribution?

$$z = \frac{x - \mu}{\sigma}$$

$$\frac{\sigma - 0}{\sigma} = \frac{\sigma}{\sigma} = 1$$

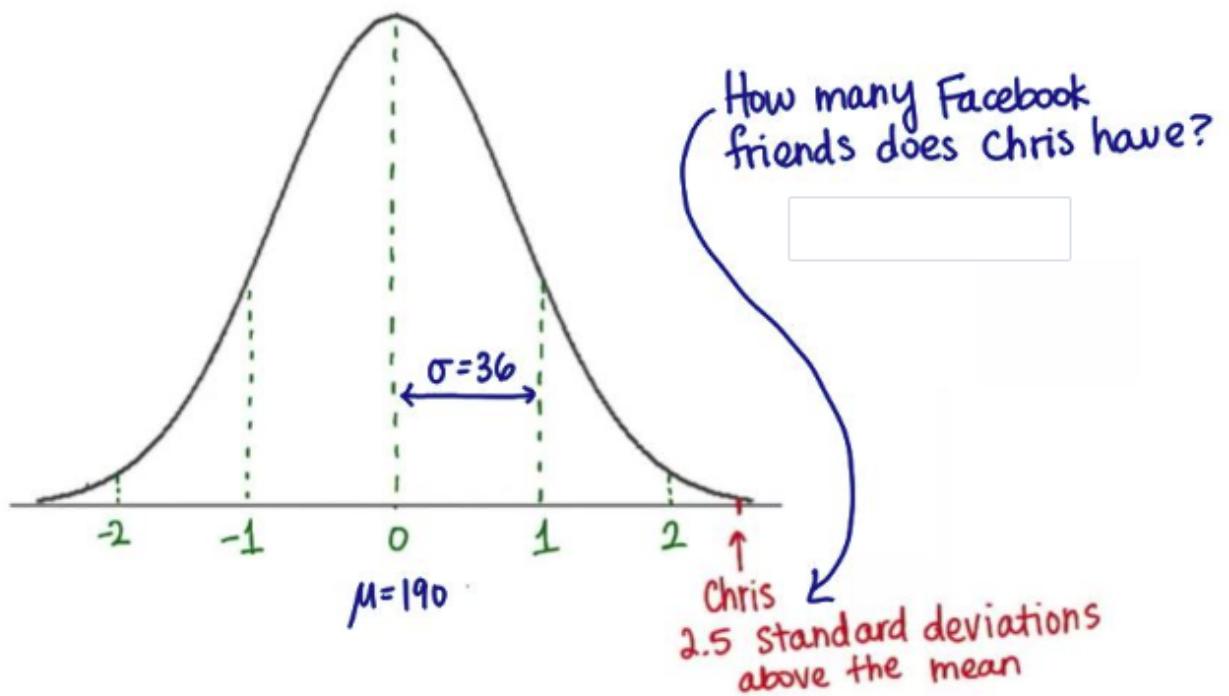
Então, a pontuação Z de qualquer valor que estiver a um desvio padrão da média

Standard Normal Distribution



Standard Normal Distribution

$$2.5 * 36 = 90$$



$$\text{SOMA } 190 + 90 = 280$$

How many Facebook friends does Chris have if the mean number of friends is 190, the standard deviation is 36, and the number of friends Chris has is 2.5 standard deviations above the mean?

$$\begin{aligned}2.5 \text{ std.dev.} &= 36 \times 2.5 \\&= 90\end{aligned}$$

$$190 + 90 = 280$$

Chris
2.5 standard deviations

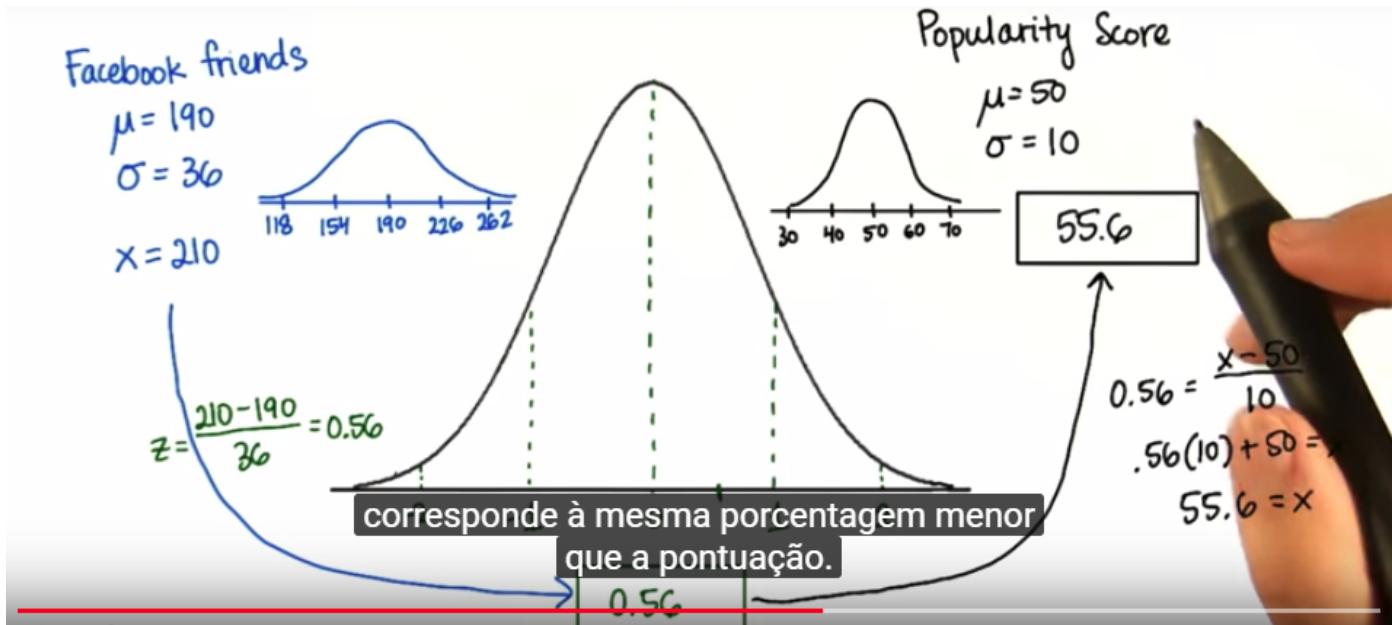
$$\begin{aligned}z &= \frac{x - \mu}{\sigma} \\2.5 &= \frac{x - 190}{36}\end{aligned}$$

Se resolver algebraicamente multiplicando em cruz e somando 190, Chris

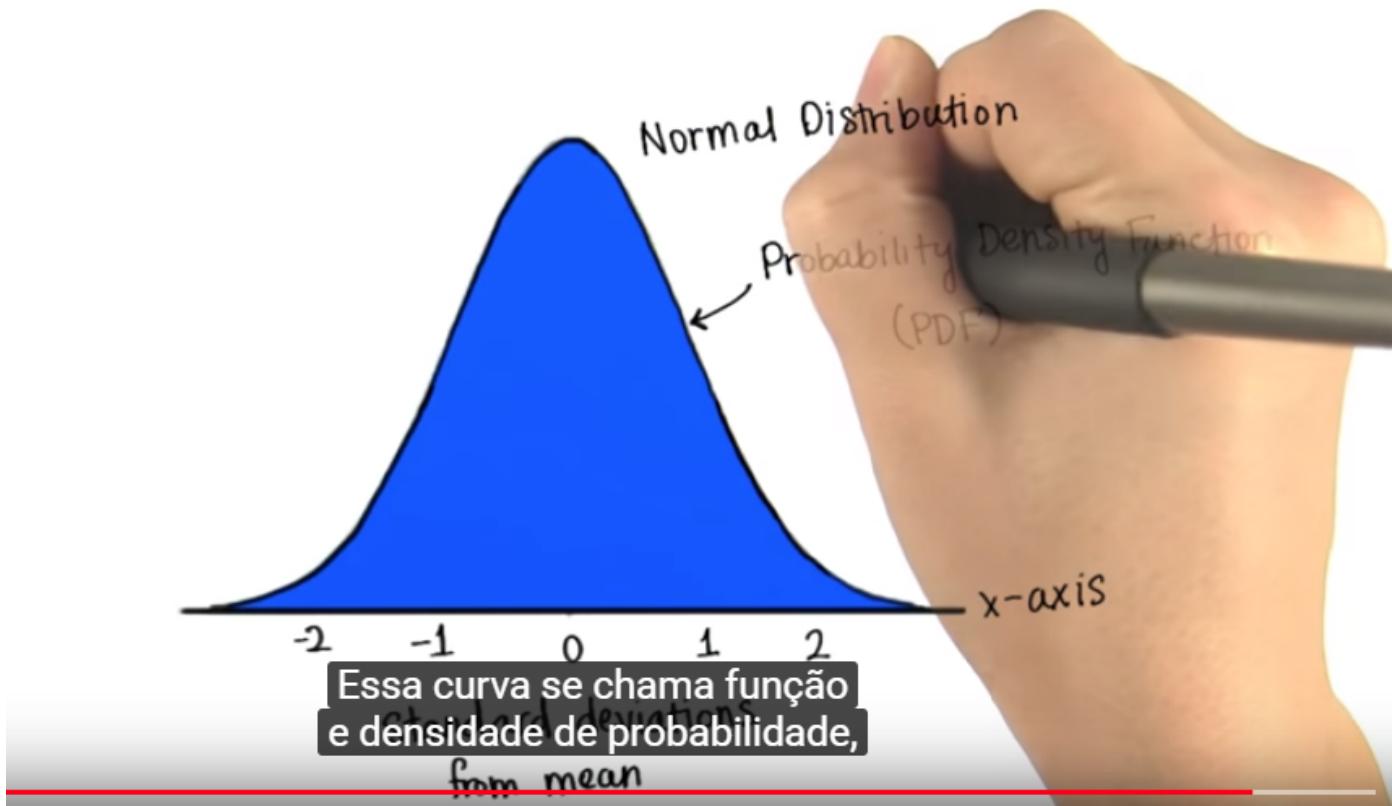
How many Facebook friends does Chris have if the mean number of friends is 190, the standard deviation is 36, and the number of friends Chris has is 2.5 standard deviations above the mean?

$$\begin{aligned}2.5 \text{ std.dev.} &= 36 \times 2.5 \\&= 90\end{aligned}$$

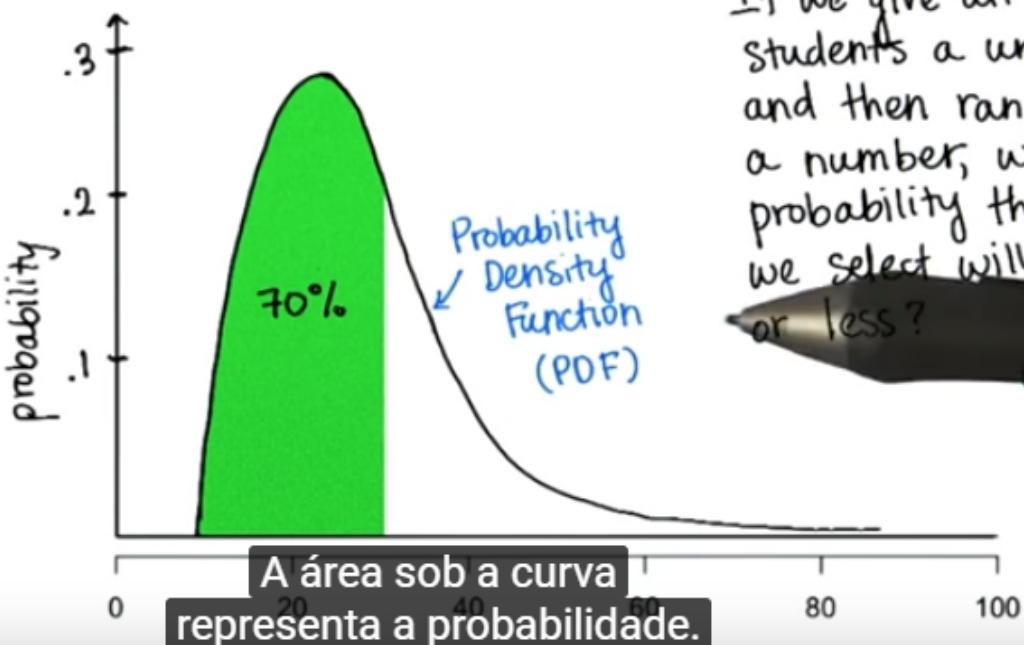
$$190 + 90 = 280$$



PDF



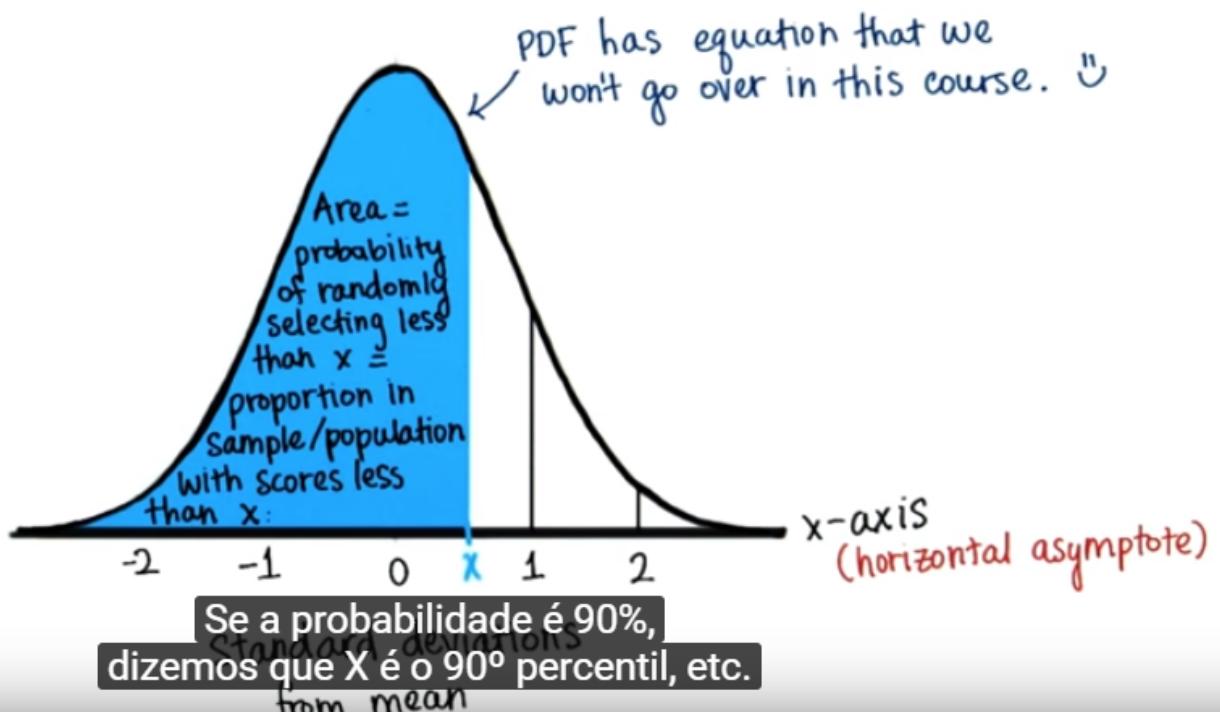
Udacity Students' Ages

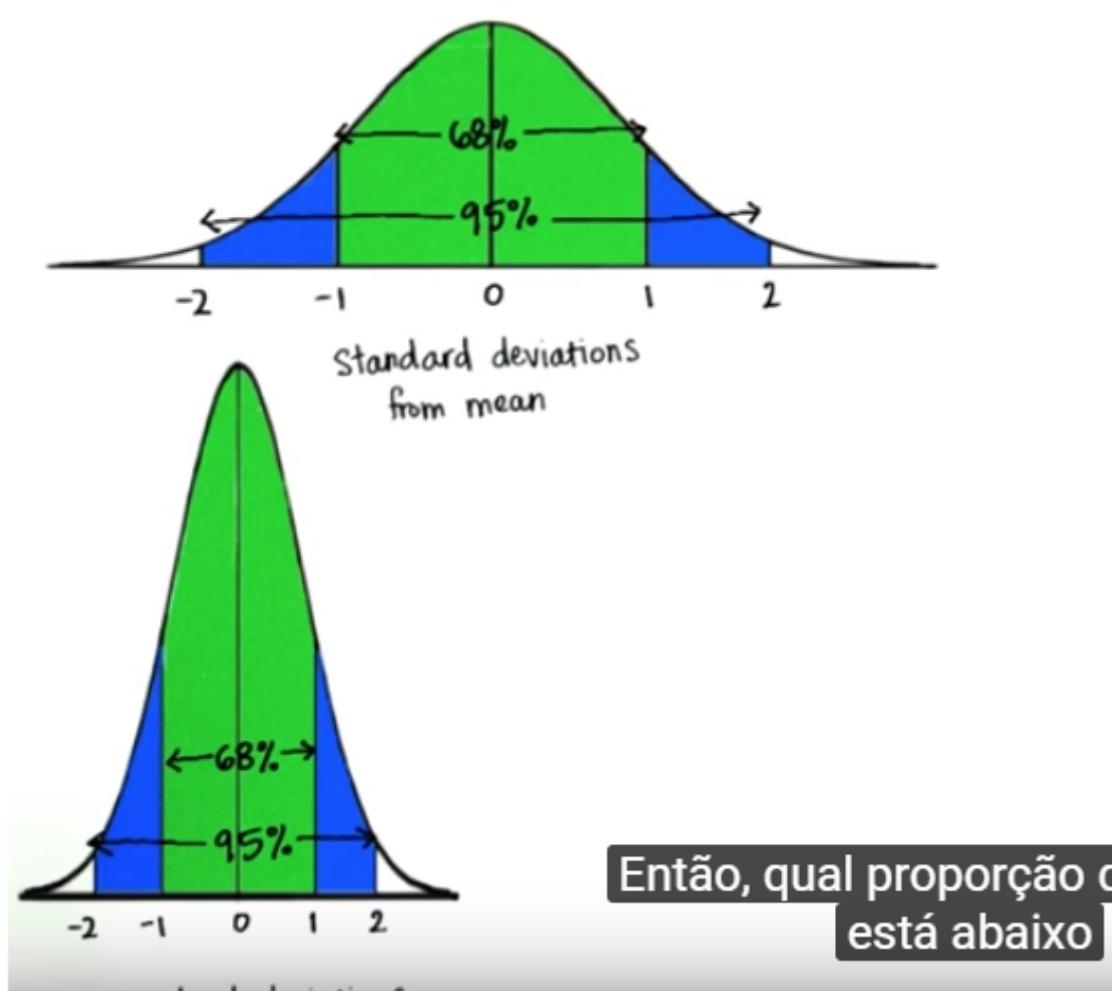
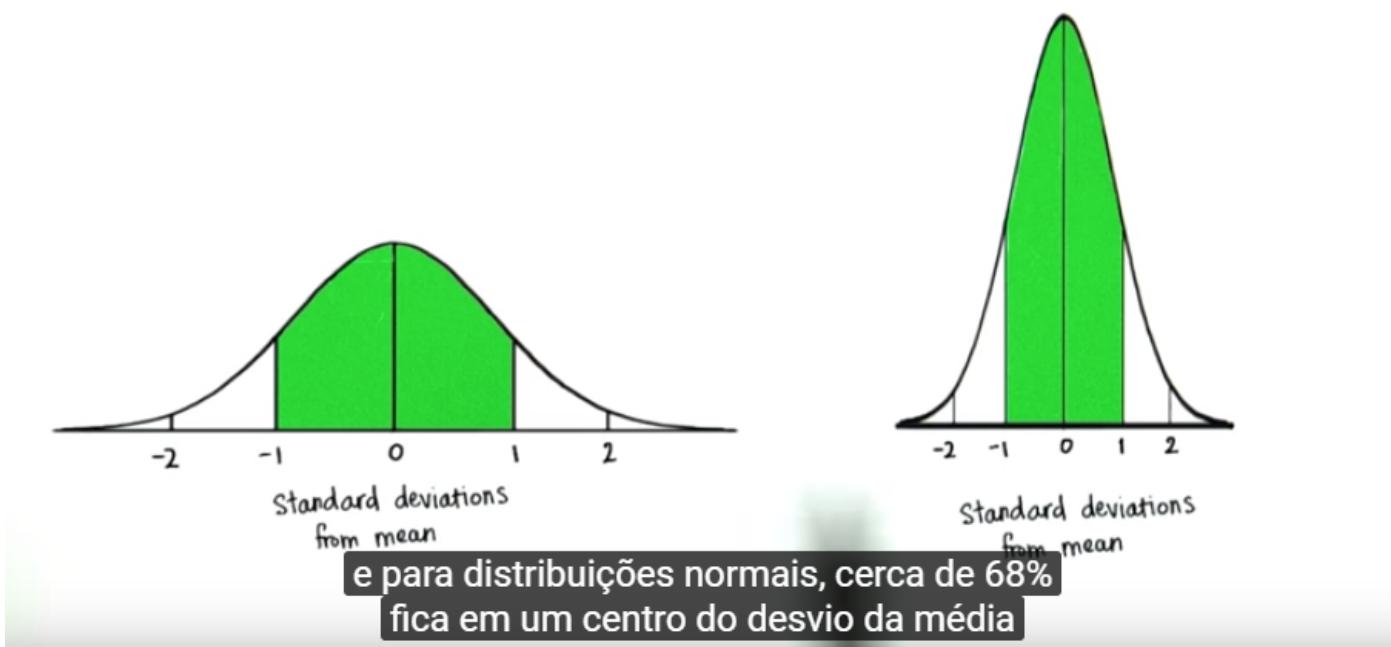


PDF regras

- As extremidades nunca tocam o eixo X (Consideramos isso, pois não podemos ter 100% de algo)
- Se existe um certo valor (x) a área sob a curva do negativo até X

Normal Distribution

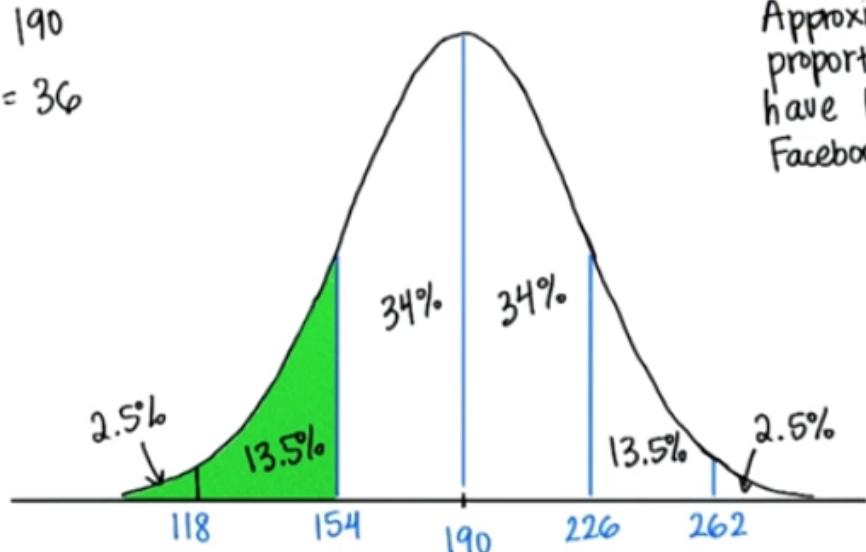




Distribution of Facebook friends

$$\mu = 190$$

$$\sigma = 36$$



Approximately what proportion of people have less than 154 Facebook friends?

Então é uma proporção de 0,16.

In []: