

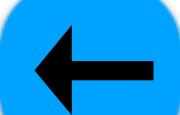
Confidence Intervals



Points → Intervals

Distribution or population

Estimate parameters



Point estimates

$\mu \approx 3.14$

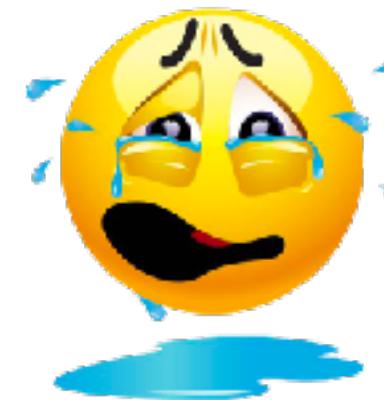
$p \approx 0.48$



Precise



Certainly wrong



No confidence



Confidence intervals



Precision



Confidence



With 95% confidence (probability) $\mu \in (3.1, 3.18)$

Back to Normal

CLT

Averages normally distributed

Intuition

Almost everything

$X_1, \dots, X_n \perp, \sim$ any distribution with mean μ , and stdv σ

$$\overline{X}^n \stackrel{\text{def}}{=} \frac{X_1 + \dots + X_n}{n}$$

Sample mean

$$Z_n \stackrel{\text{def}}{=} \frac{(X_1 + \dots + X_n) - n\mu}{\sigma\sqrt{n}}$$

Typically ≥ 30

CLT

For sufficiently large n

Roughly

$$Z_n \stackrel{\text{distr}}{\sim} \mathcal{N}(0, 1)$$

Standard Normal Variable

Predicting Standard Normal

Standard Normal Variable

Predict value of Z

Point prediction

Highest probability

Unbiased

Precise

Wrong

Interval

$$-a \leq Z \leq a$$

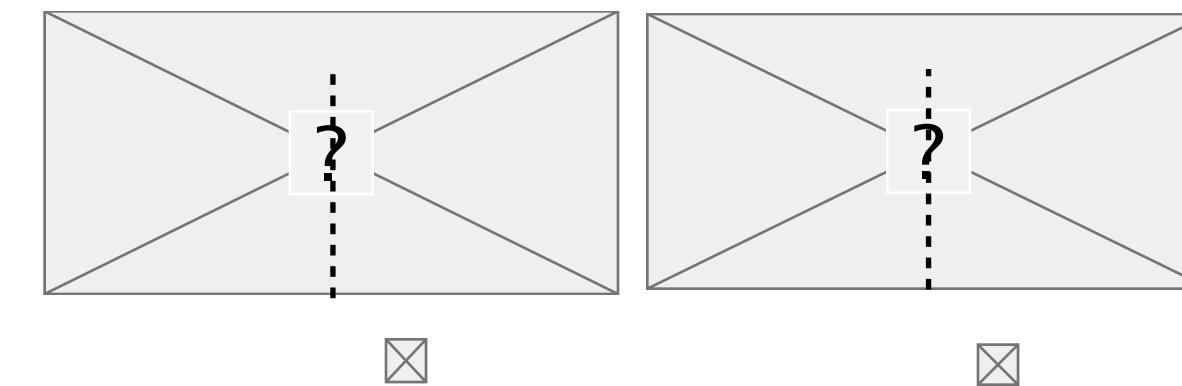
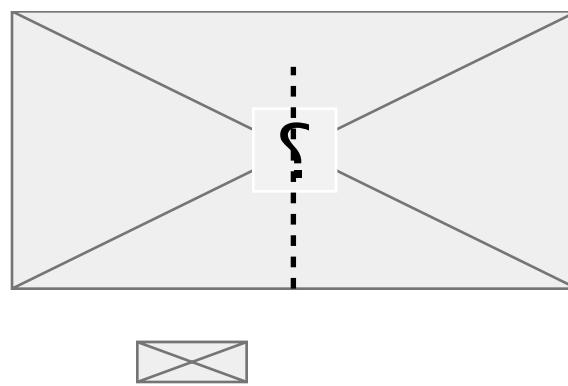
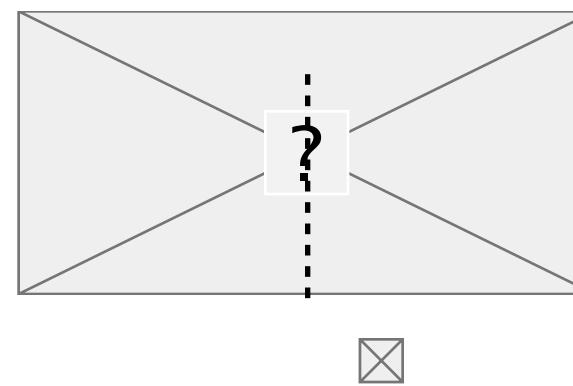
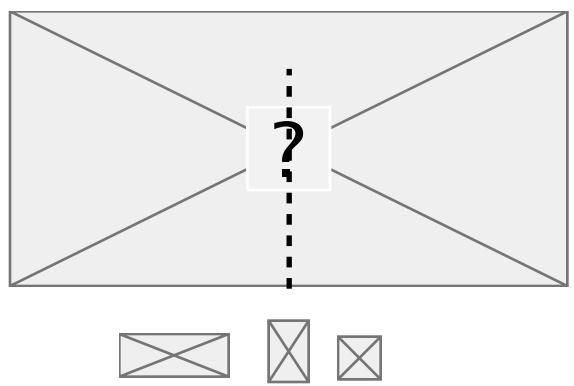
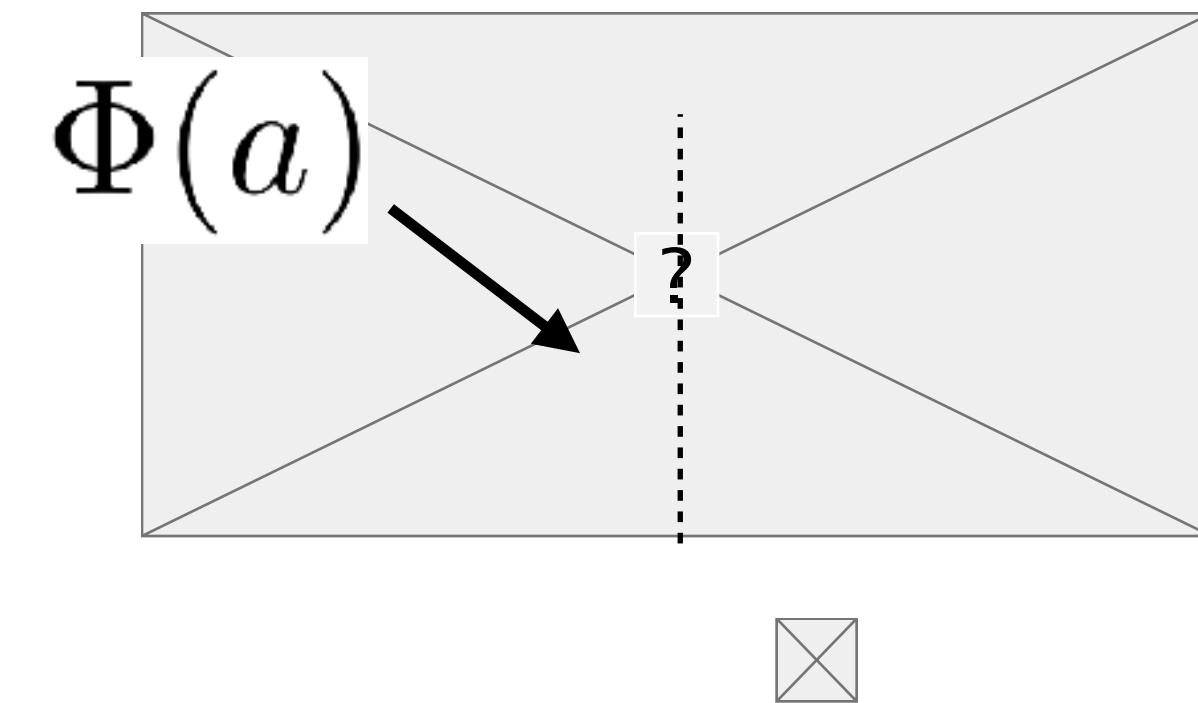
$$P(-a \leq Z \leq a) > 0$$

P = ?

Interval Probability

$$Z \sim \mathcal{N}(0, 1)$$

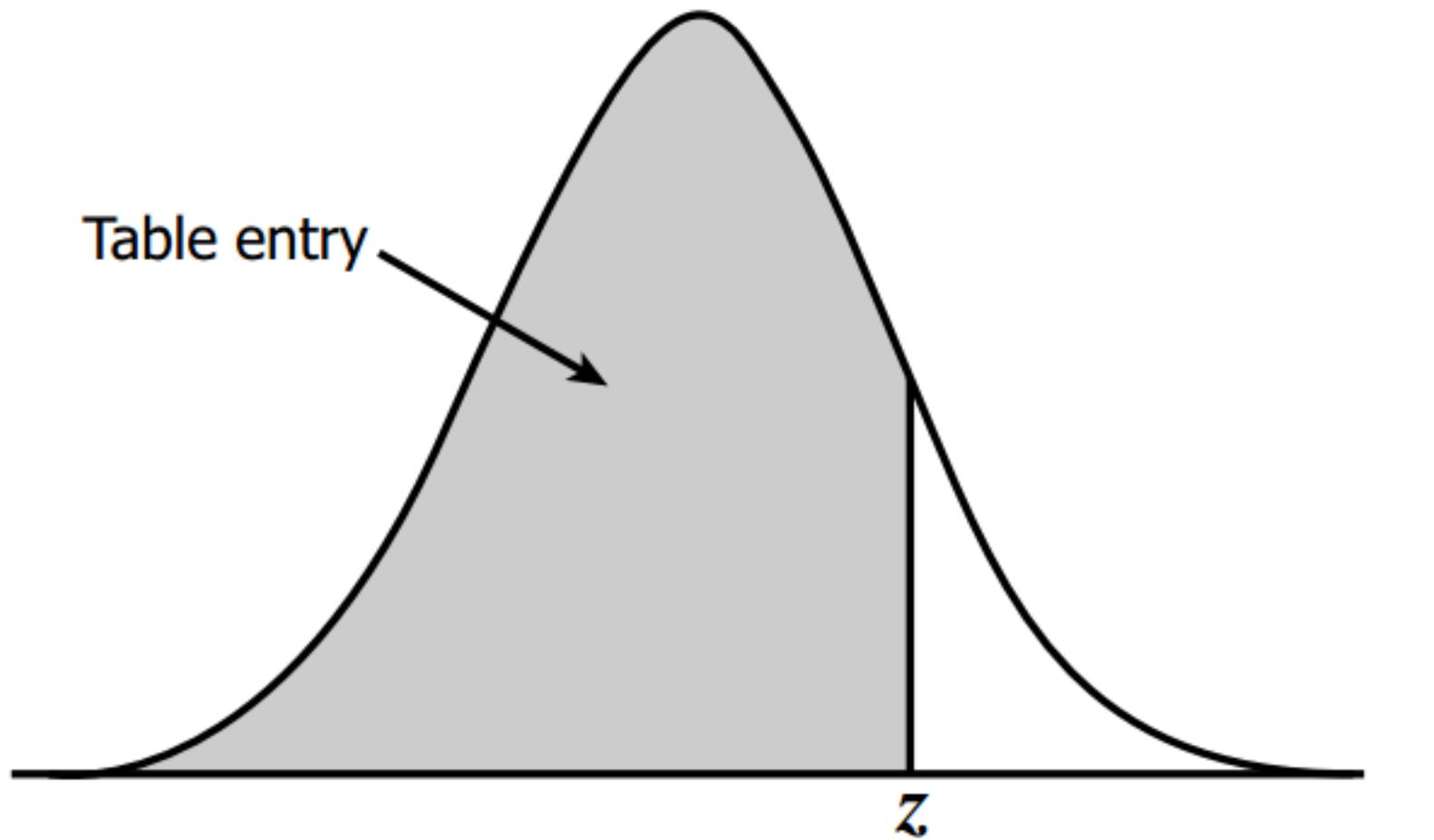
$$\Phi(a) = F(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-t^2/2} dt$$



Calculating $\Phi(a)$

No known formula

Use Z / Standard Normal Table



$$\Phi(1) = 0.8413$$

Program

Python

In `scipy.stats`

`norm.cdf(x)`

cumulative distribution function

$\Phi(1)$

```
from scipy.stats import norm  
norm.cdf(1)  
0.841344746069
```

$\Phi(2)$

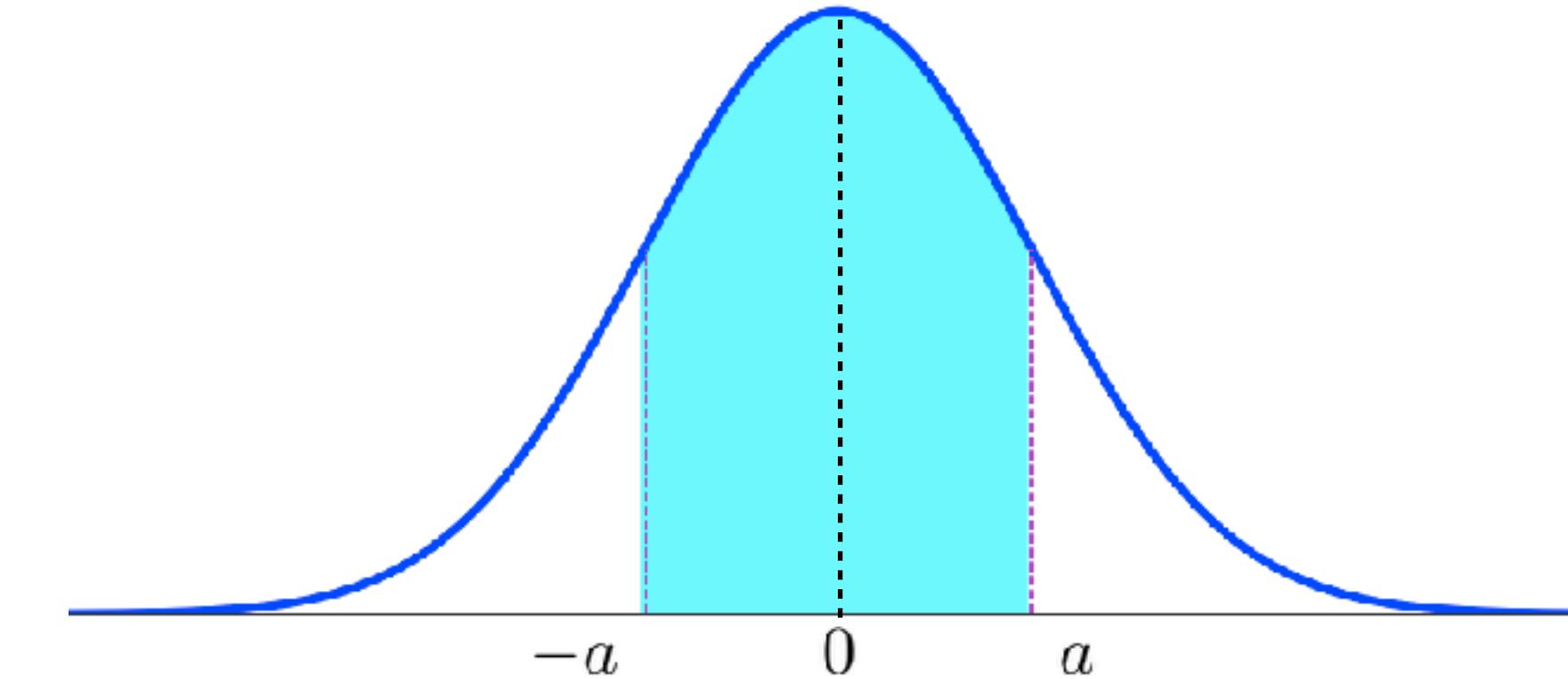
```
norm.cdf(2)  
0.977249868052
```

$\Phi(3)$

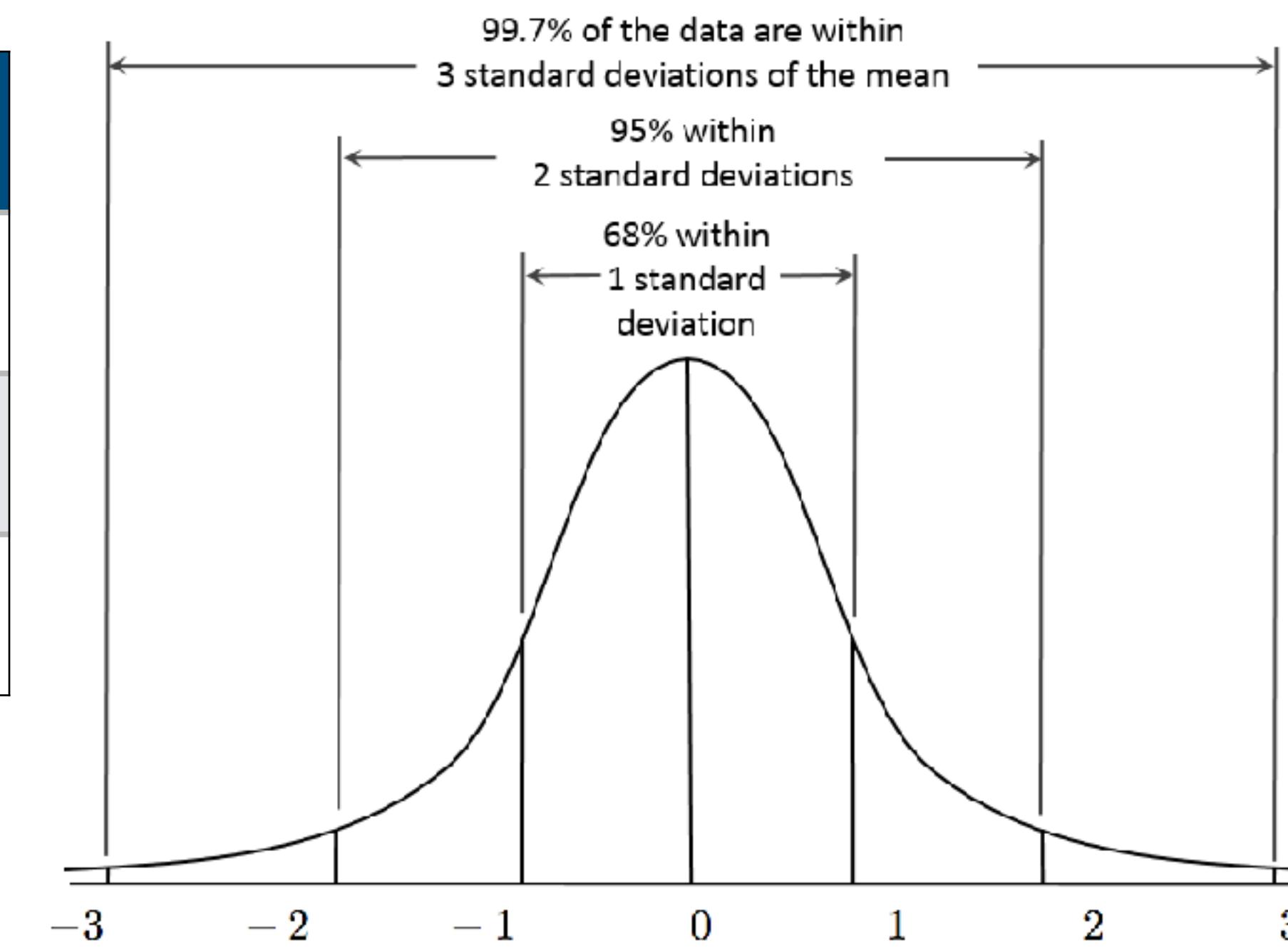
```
norm.cdf(3)  
0.998650101968
```

68 - 95 - 99.7 Rule

$$P(-a \leq Z \leq a) = 2\Phi(a) - 1$$



a	$P(-a \leq Z \leq a)$
1	$2 \cdot 0.8413 - 1 = 0.682$
2	$2 \cdot 0.9772 - 1 = 0.9544$
3	$2 \cdot 0.9987 - 1 = 0.9974$



Interval → Probability

Typically

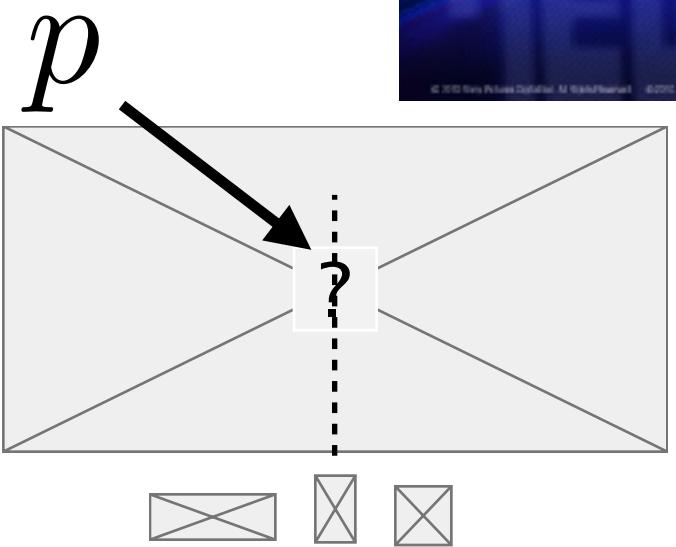
Given desired probability p

Find

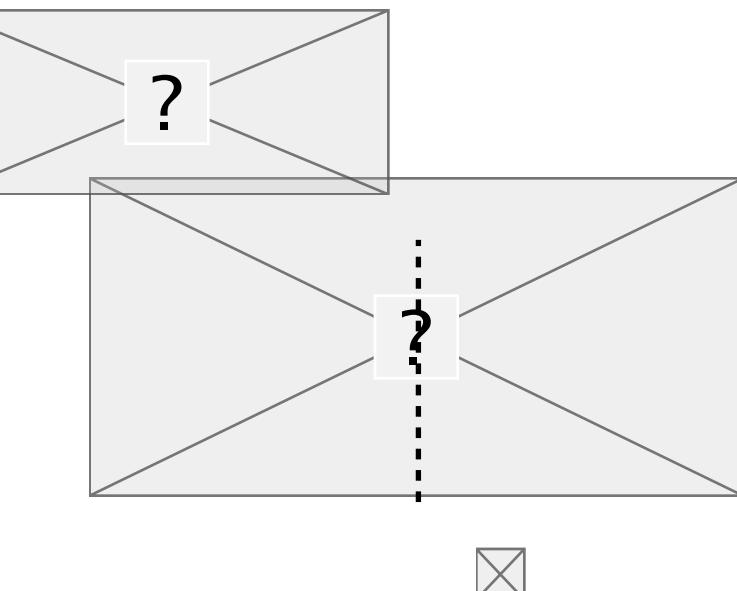
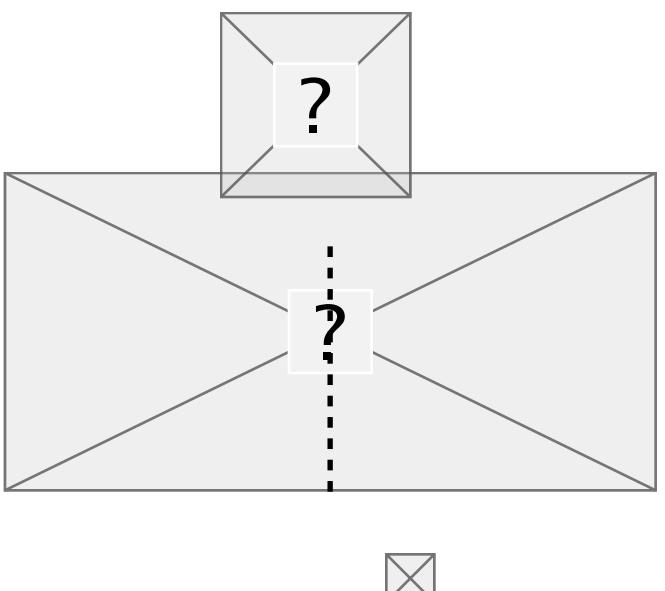
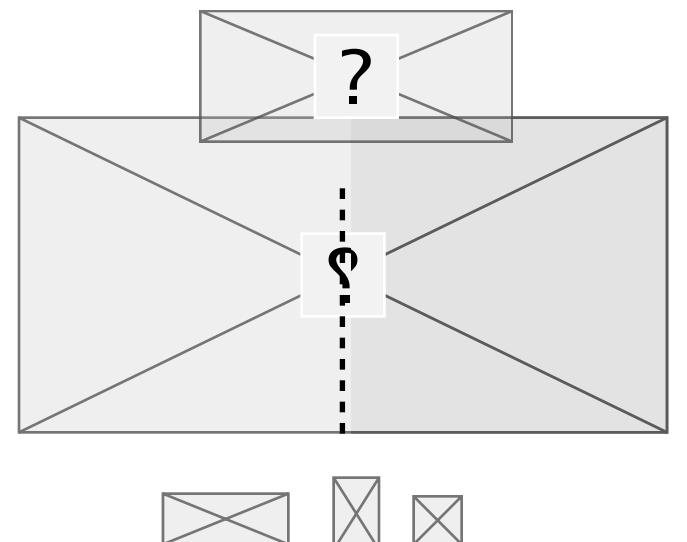
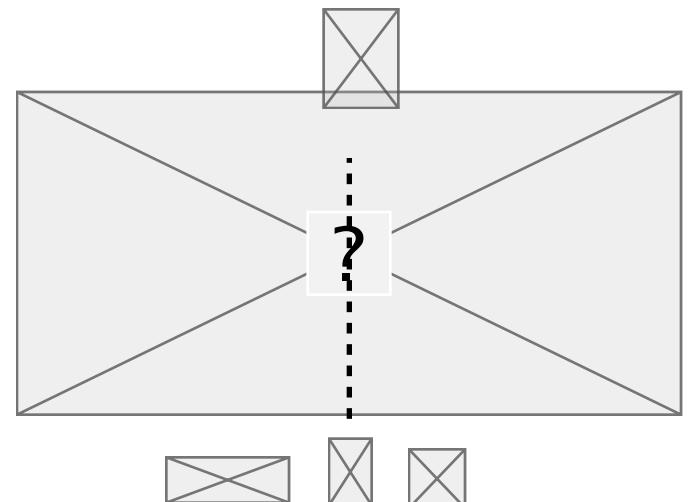
a

s.t.

$$P(-a \leq Z \leq a) = p$$



Saw



Python

norm.ppf(p)

percent point function

converts percentile to a point

$\Phi^{-1}(0.95)$

```
from scipy.stats import norm  
norm.ppf(0.95)  
1.64485362695
```

$\Phi^{-1}(0.975)$

```
norm.ppf(0.975)  
1.95996398454
```

$\Phi^{-1}(0.99)$

```
norm.ppf(0.99)  
2.32634787404
```

Common Values

p=95%

`norm.ppf(0.975)`
1.95996398454

$$a = \Phi^{-1} \left(\frac{1 + p}{2} \right) = \Phi^{-1}(0.975) \approx 1.96$$

$$P(-1.96 \leq Z \leq 1.96) \approx 0.95$$

68 - 95 - 99.7

$$P(-2 \leq Z \leq 2) \approx 0.95$$

p	$\frac{1 + p}{2}$	$\Phi^{-1} \left(\frac{1 + p}{2} \right)$
90	0.95	1.645
95	0.975	1.960
98	0.99	2.056

General Normal Distributions

$$X \sim \mathcal{N}_{\mu, \sigma^2} \quad \mathcal{N}(\mu, \sigma^2)$$

$$Z \stackrel{\text{def}}{=} \frac{X - \mu}{\sigma} \sim \mathcal{N}_{0,1}$$

Standard Normal

$$P(\mu - a\sigma \leq X \leq \mu + a\sigma) = P(-a\sigma \leq X - \mu \leq a\sigma)$$

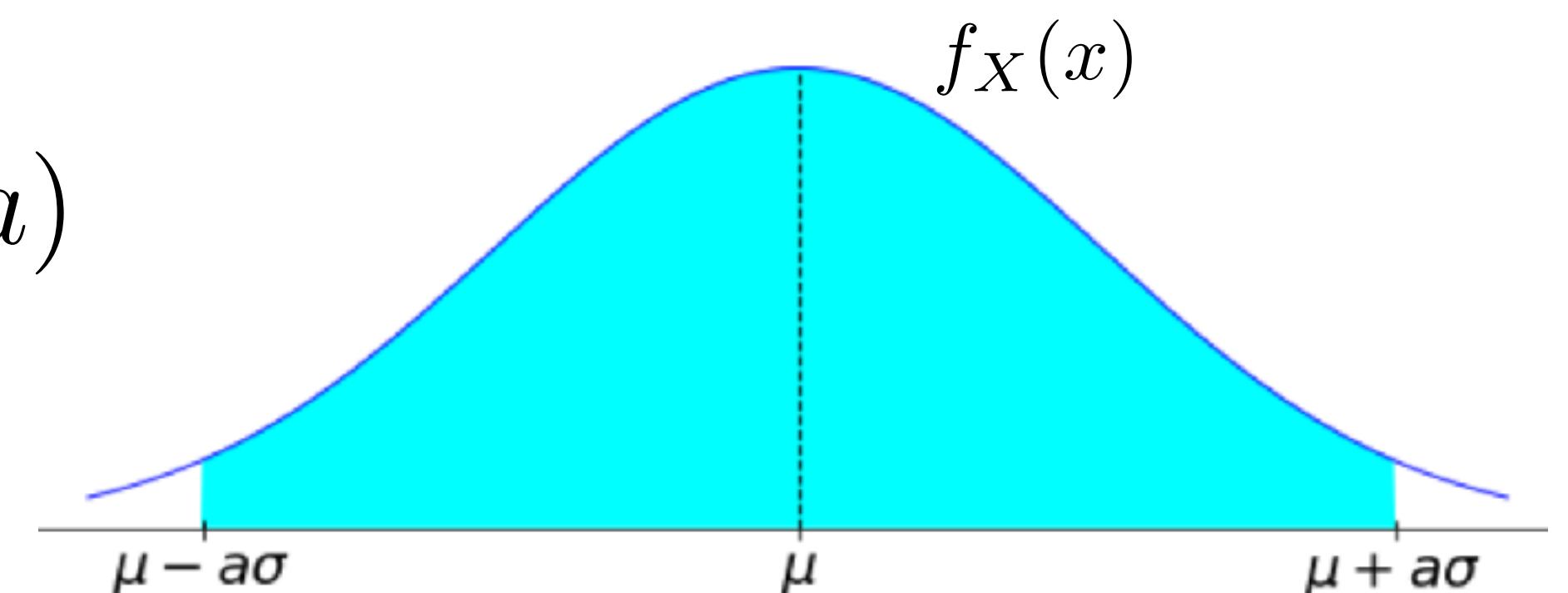
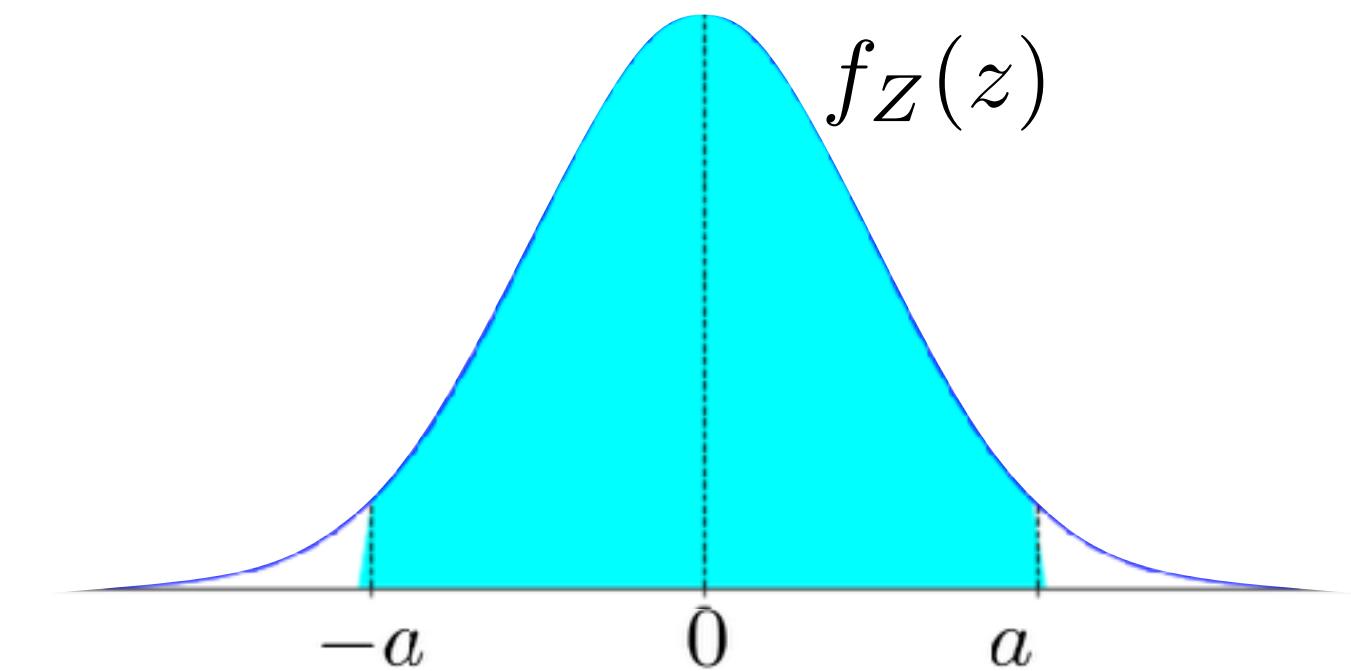
X within “a” std
from its mean

$$= P(-a \leq \frac{X - \mu}{\sigma} \leq a)$$

$$= P(-a \leq Z \leq a)$$

Z within “a” std
from its mean

X, Z
normal



Example

$$X \sim N(1, 4)$$

$$p = 0.95$$

$$\mu = 1$$

$$\sigma = 2$$

Z within 1.96 std
from its mean

X within 1.96 std
from its mean

$$\begin{aligned} 0.95 &\approx P(-1.96 \leq Z \leq 1.96) = P(\mu - 1.96\sigma \leq X \leq \mu + 1.96\sigma) \\ &= P(-2.92 \leq X \leq 4.92) \end{aligned}$$

Confidence Intervals

Any parameter

Simplest and by far most common

mean μ

proportion p

Given a sample X_1, \dots, X_n

Find an interval containing μ

First

σ known

Next lecture

σ unknown

Sample-Mean Distribution

$$\times \frac{\sigma}{\sqrt{n}}$$

$$\frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}} \stackrel{d}{\sim} \mathcal{N}(0, 1)$$

$$+ \mu$$

$$\frac{X_1 + \dots + X_n - n\mu}{n} \stackrel{d}{\sim} \mathcal{N}\left(0, \frac{\sigma^2}{n}\right)$$

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} \stackrel{d}{\sim} \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Roughly normal

$$\bar{X}$$

Centered at sample mean

Standard deviation

$$V(\bar{X}) = \frac{\sigma^2}{n} \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

Sampling
Distribution of the
sample mean

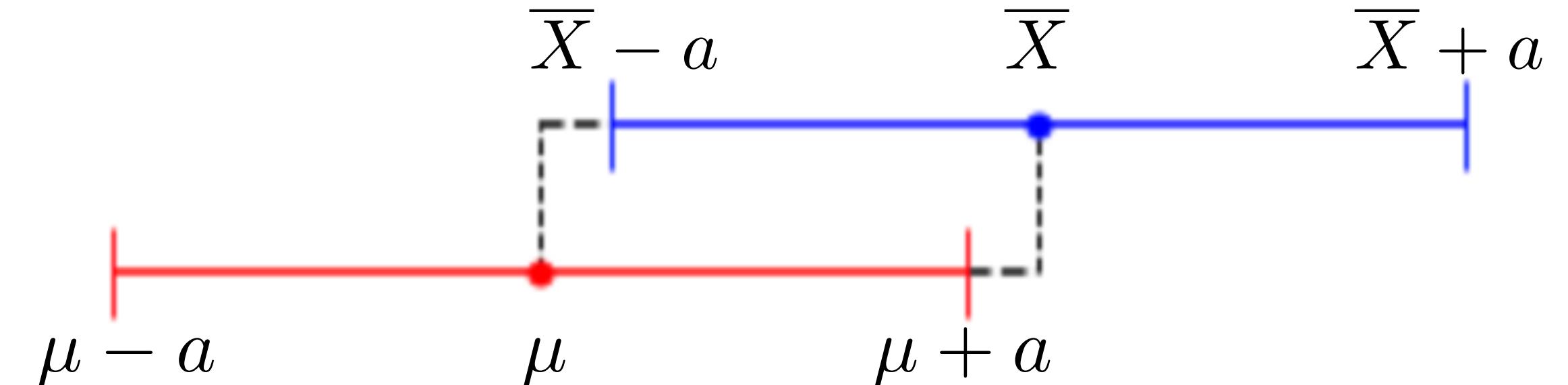
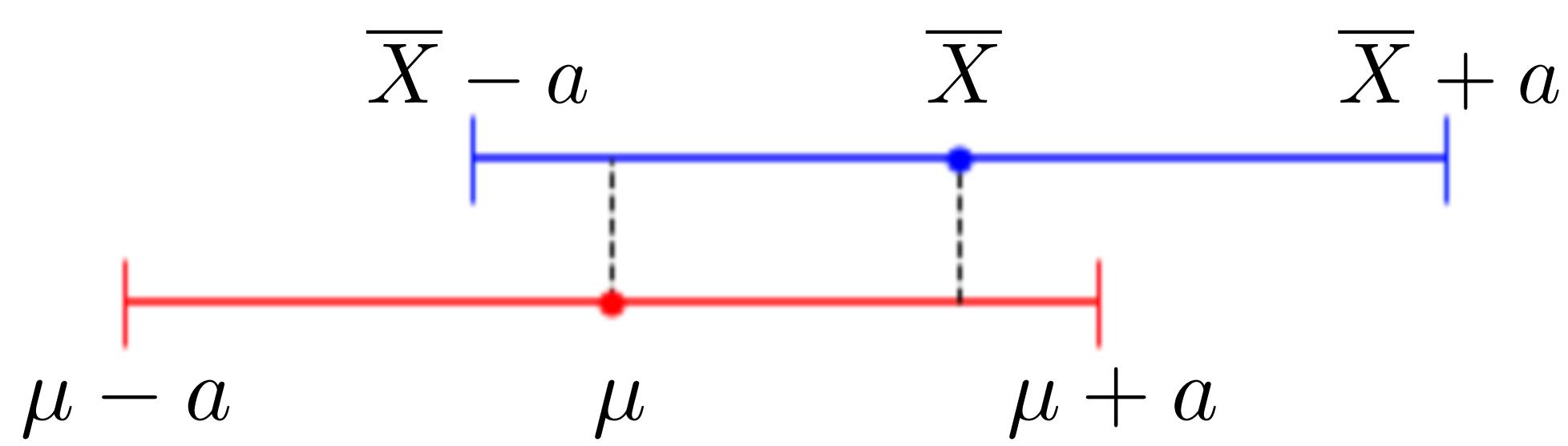
Standard
error

Proximity is Reciprocal

← \bar{X} near μ

→ μ near \bar{X}

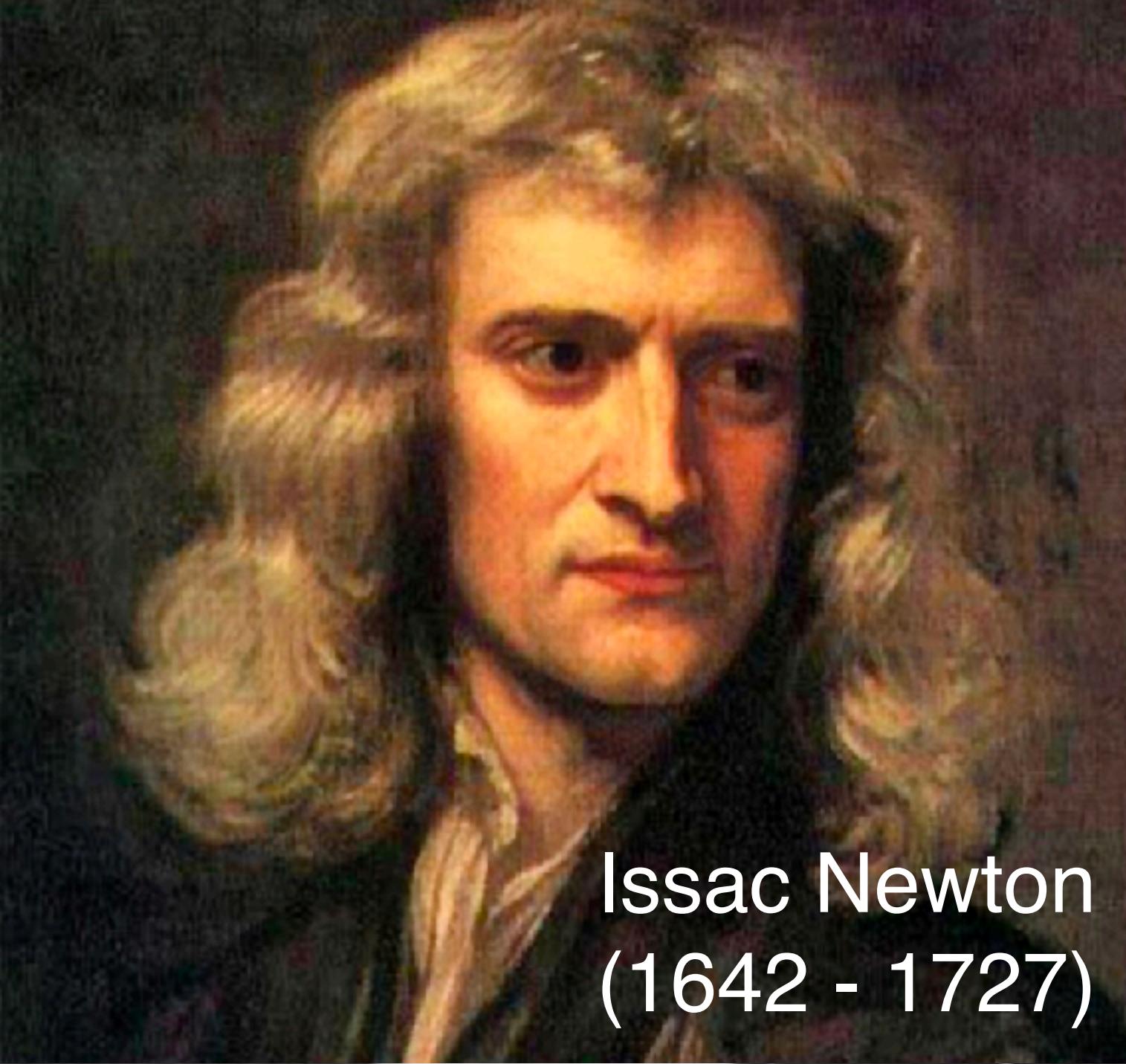
$$\bar{X} \in (\mu - a, \mu + a) \quad |\bar{X} - \mu| < a \quad \mu \in (\bar{X} - a, \bar{X} + a)$$



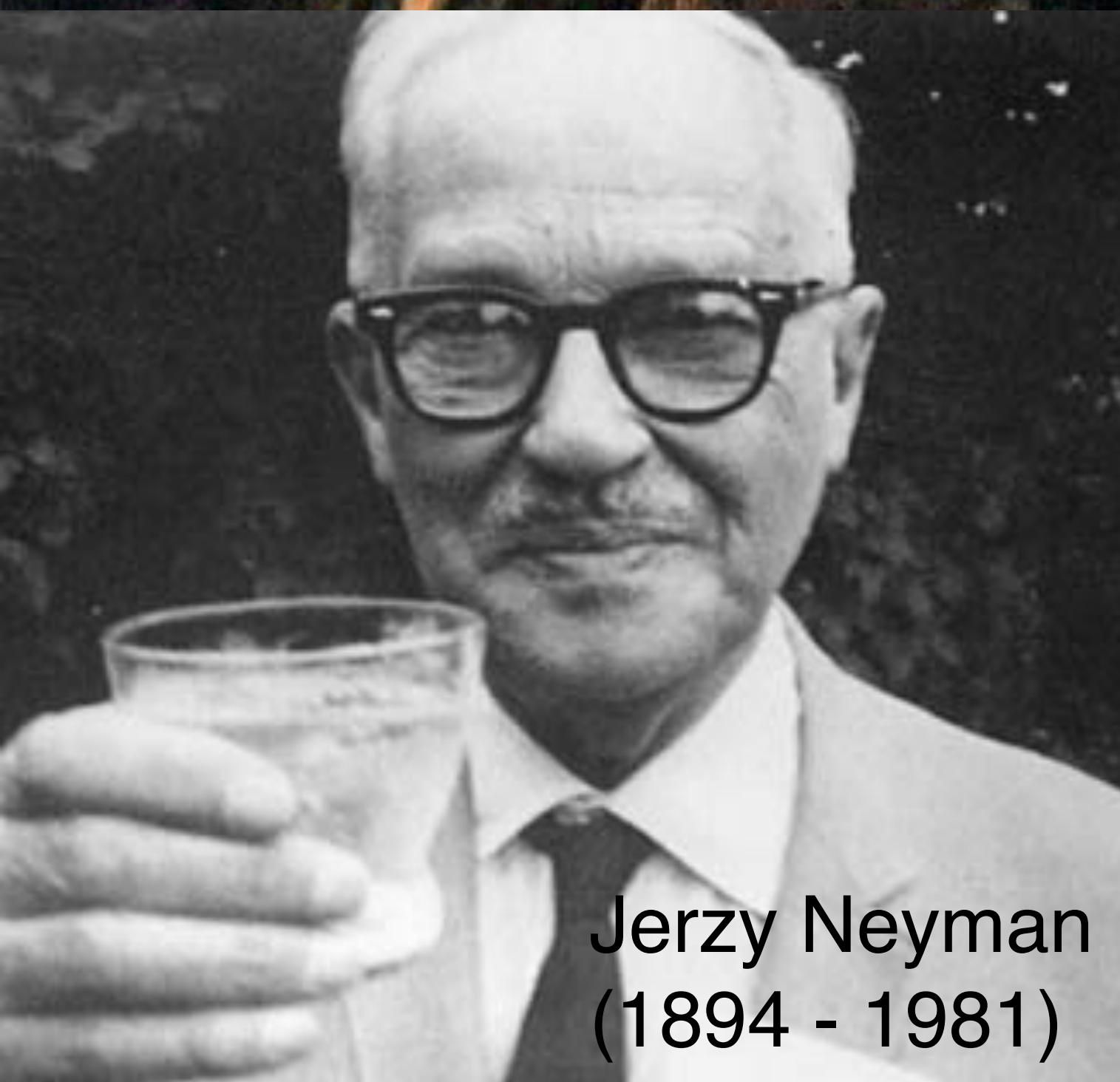
$$\bar{X} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$$

With high probability \bar{X} near μ

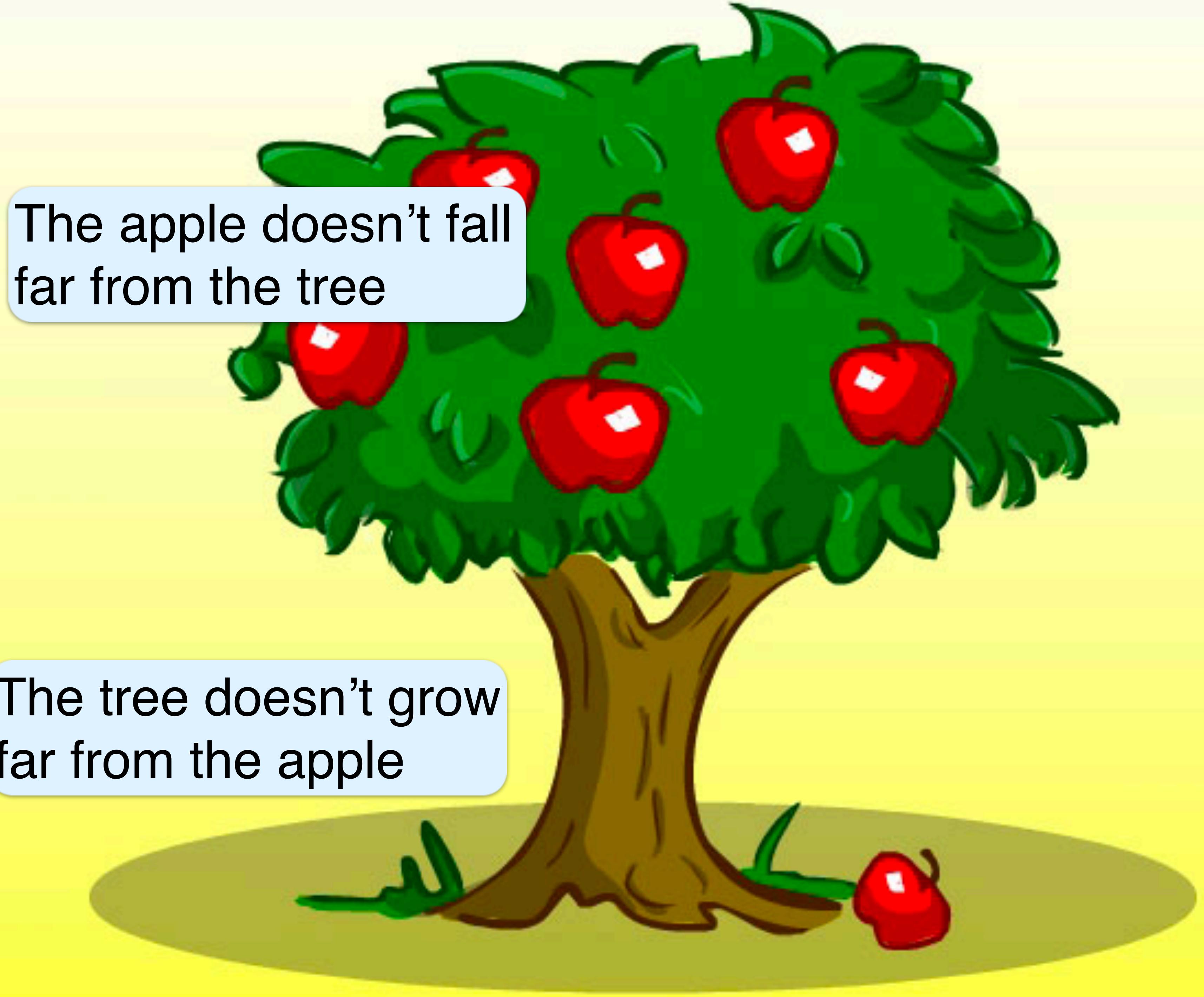
With high probability μ near \bar{X}



Issac Newton
(1642 - 1727)



Jerzy Neyman
(1894 - 1981)



The apple doesn't fall
far from the tree

The tree doesn't grow
far from the apple

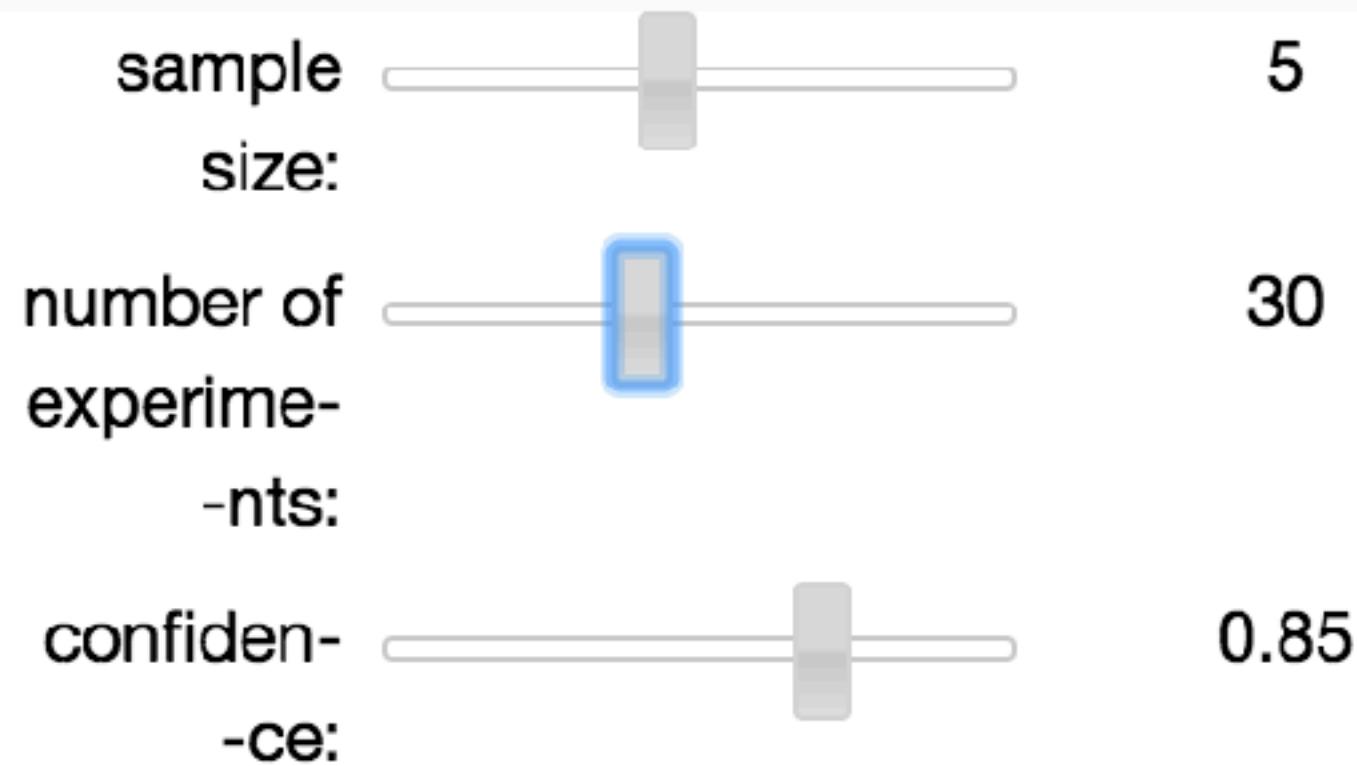
Confidence Interval

With probability p

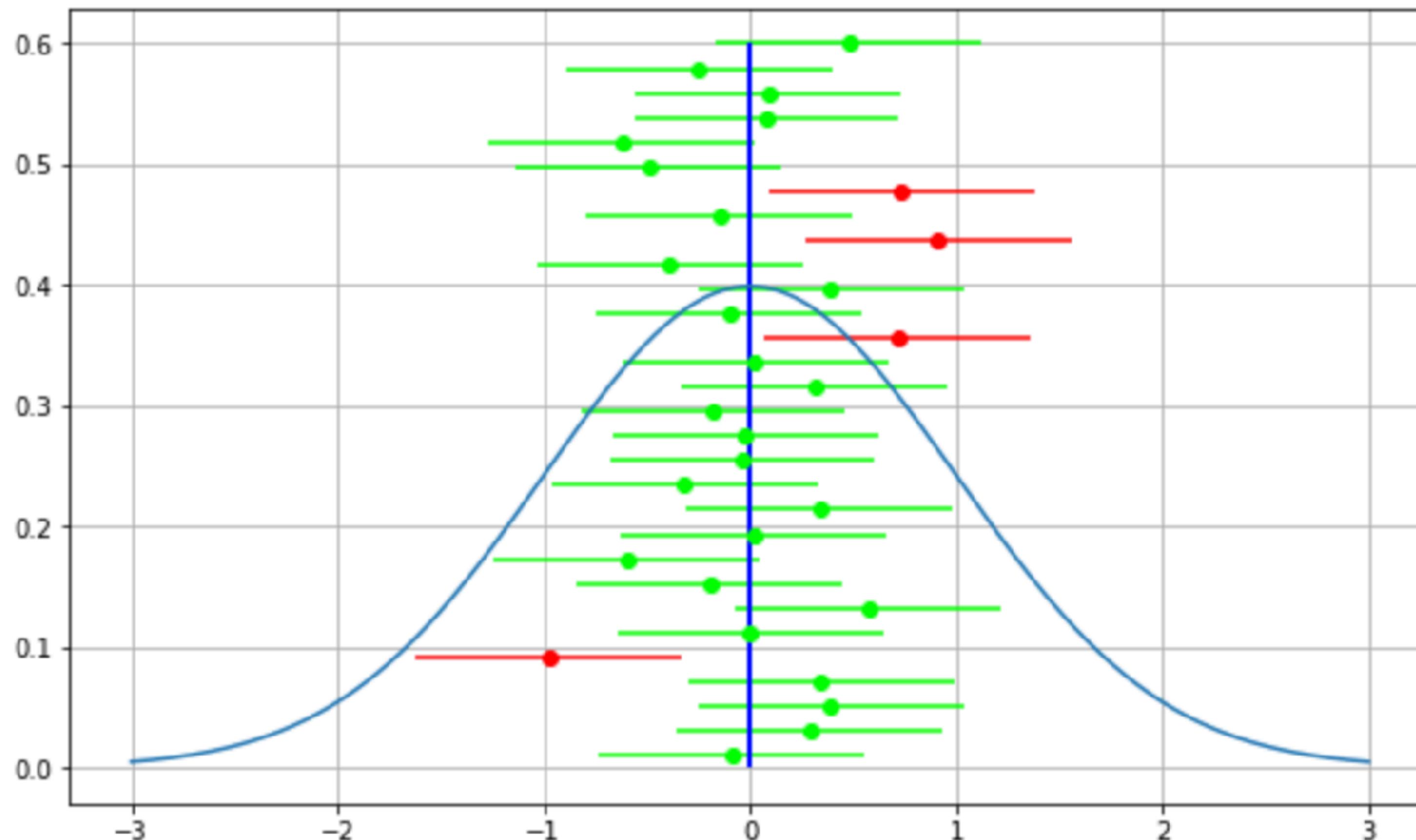
$$\bar{X} \in (\mu - z_p \sigma_{\bar{X}}, \mu + z_p \sigma_{\bar{X}})$$

$$|\bar{X} - \mu| < z_p \sigma_{\bar{X}}$$

$$\mu \in \left(\bar{X} - z_p \frac{\sigma}{\sqrt{n}}, \bar{X} + z_p \frac{\sigma}{\sqrt{n}} \right)$$



Confidence level = 85.00%, 0 falls in 86.67% of the intervals



Daily Tweets

tweets of a random Tweeter user is a random variable with $\sigma=2$

In a sample of 121 users the sample mean was 3.7

Find the 95% confidence interval for the distribution mean

$$z_p = \Phi^{-1}\left(\frac{1 + p}{2}\right) = \Phi^{-1}(0.975) = 1.96$$

95% confidence interval for mean

$$(\bar{X} - z_p \sigma_{\bar{X}}, \bar{X} + z_p \sigma_{\bar{X}}) = (\bar{X} - z_p \frac{\sigma}{\sqrt{n}}, \bar{X} + z_p \frac{\sigma}{\sqrt{n}})$$

Margin of error

$$= (3.344, 4.056)$$

Heart Rate per Minute

Adult heart rate has standard deviation $\sigma=7.5$ beats per minute

Estimate average heart rate within margin of error < 2

With confidence level 90%

$$z_p = \Phi^{-1}\left(\frac{1+p}{2}\right) = \Phi^{-1}(0.95) = 1.645$$

Sample size

$$z_p \sigma_{\bar{X}} = z_p \frac{\sigma}{\sqrt{n}} = 2 \quad n = \left(z_p \frac{\sigma}{2}\right)^2 = 38.05$$

Confidence Intervals

