

# Confidence Intervals Unknown $\sigma$



Estimate  $\sigma$

Student's t-distribution

Step-by-step instructions

Example

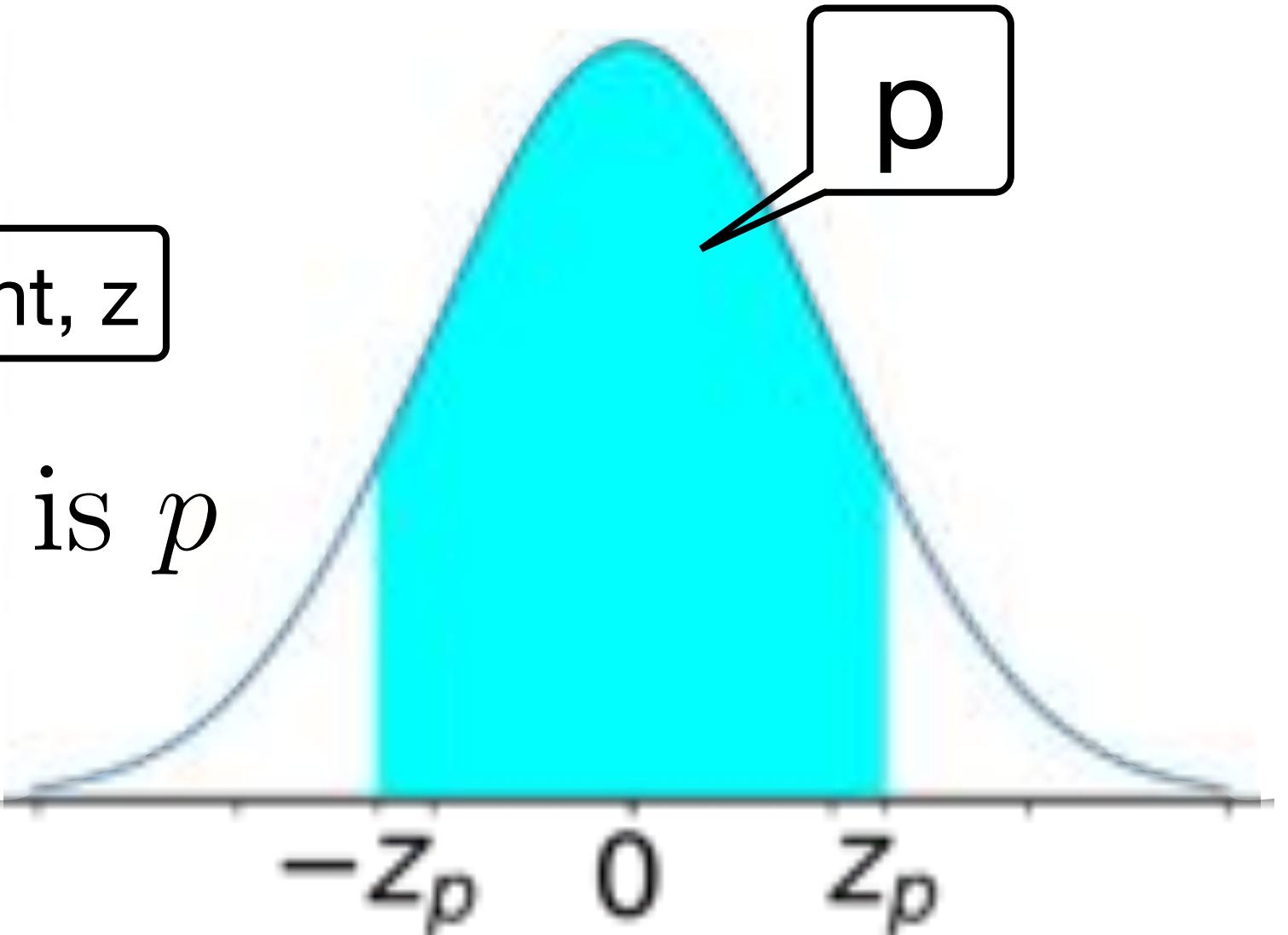
# Confidence Intervals - Known $\sigma$

Standard normal distribution

$$\mathcal{N}_{0,1}$$

$0 \leq p \leq 1$      $z_p$  : point s.t. area between  $-z_p$  and  $z_p$  is  $p$

$$Z \sim \mathcal{N}_{0,1} \quad p = P(|Z| \leq z_p) = P(-z_p \leq Z \leq z_p)$$



$$= \Phi(z_p) - \Phi(-z_p) = \Phi(z_p) - (1 - \Phi(z_p)) = 2\Phi(z_p) - 1$$

CDF of standard normal

$$\Phi(z_p) = \frac{1+p}{2}$$

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

$$z_{0.9} = \Phi^{-1}\left(\frac{1+0.9}{2}\right) = \Phi^{-1}(0.95) = 1.645$$

$$P(|Z| \leq 1.645) = 0.9$$

# Sample Mean $\approx$ Normal

$X_1, X_2, \dots, X_n$

i.i.d.

known  $\sigma$

unknown  $\mu$

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}$$

Sample mean

$$\mu_{\bar{X}} = \mu$$

Unbiased

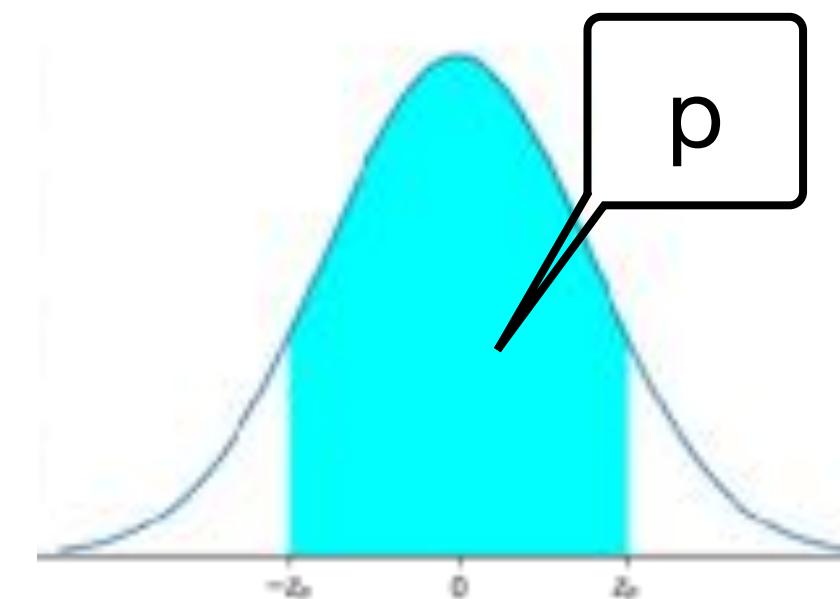
$$V(\bar{X}) = \frac{\sigma^2}{n} \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

$$\frac{\bar{X} - \mu}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \underset{\text{CLT}}{\sim} N_{0,1}$$

mean 0  
std 1

# Confidence Interval

Standard normal



$$Z \sim \mathcal{N}_{0,1} \quad P(|Z| \leq z_p) = p$$

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}_{0,1}$$

$$P\left(\left|\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}\right| \leq z_p\right) \approx p$$

$$P\left(\left|\bar{X} - \mu\right| \leq z_p \frac{\sigma}{\sqrt{n}}\right) \approx p$$

With probability  $\approx p$

$$\left|\bar{X} - \mu\right| \leq z_p \frac{\sigma}{\sqrt{n}}$$

Margin of error

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

# Confidence → Interval

Given

confidence p

samples  $X_1, \dots, X_n$

Determine

Critical value (z)  $z_p = \Phi^{-1} \left( \frac{1+p}{2} \right)$



Sample mean  $\bar{X} = \frac{X_1 + \dots + X_n}{n}$

Margin of error  $z_p \frac{\sigma}{\sqrt{n}}$   $\sigma$  known

Confidence interval  $\left[ \bar{X} - z_p \frac{\sigma}{\sqrt{n}}, \bar{X} + z_p \frac{\sigma}{\sqrt{n}} \right]$

Problem?

$\sigma$  almost never known

# Unknown $\sigma$

$$X_1, X_2, \dots, X_n \perp \mathcal{N}_{\mu, \sigma}$$

Neither  $\sigma$  nor  $\mu$  known

$\mu = ?$

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}$$

Sample mean

$$\mu_{\bar{X}} = \mu$$

Unbiased

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

$$\frac{\bar{X} - \mu}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}_{0,1}$$

Standard normal

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Sample variance

Bessel corrected

$$\mu_{S^2} = \sigma^2$$

Unbiased

$$\frac{\bar{X} - \mu}{S/\sqrt{n}}$$

Almost

standard

$$S \approx \sigma$$

normal

$S - r.v.$

Student's t-distribution

n-1 degrees of freedom

# Student's t-distribution

$$T_\nu = \frac{\bar{X} - \mu}{S/\sqrt{\nu+1}}$$

Student's t-distribution,  $\nu$  degrees of freedom

PDF

$$T_\nu \sim f_\nu(t) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\cdot\Gamma(\frac{\nu}{2})} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

Gamma function

Only dependence on t

Symmetric around 0

See a bit more



t object

in `scipy.stats` module



Degrees of freedom

probability density function

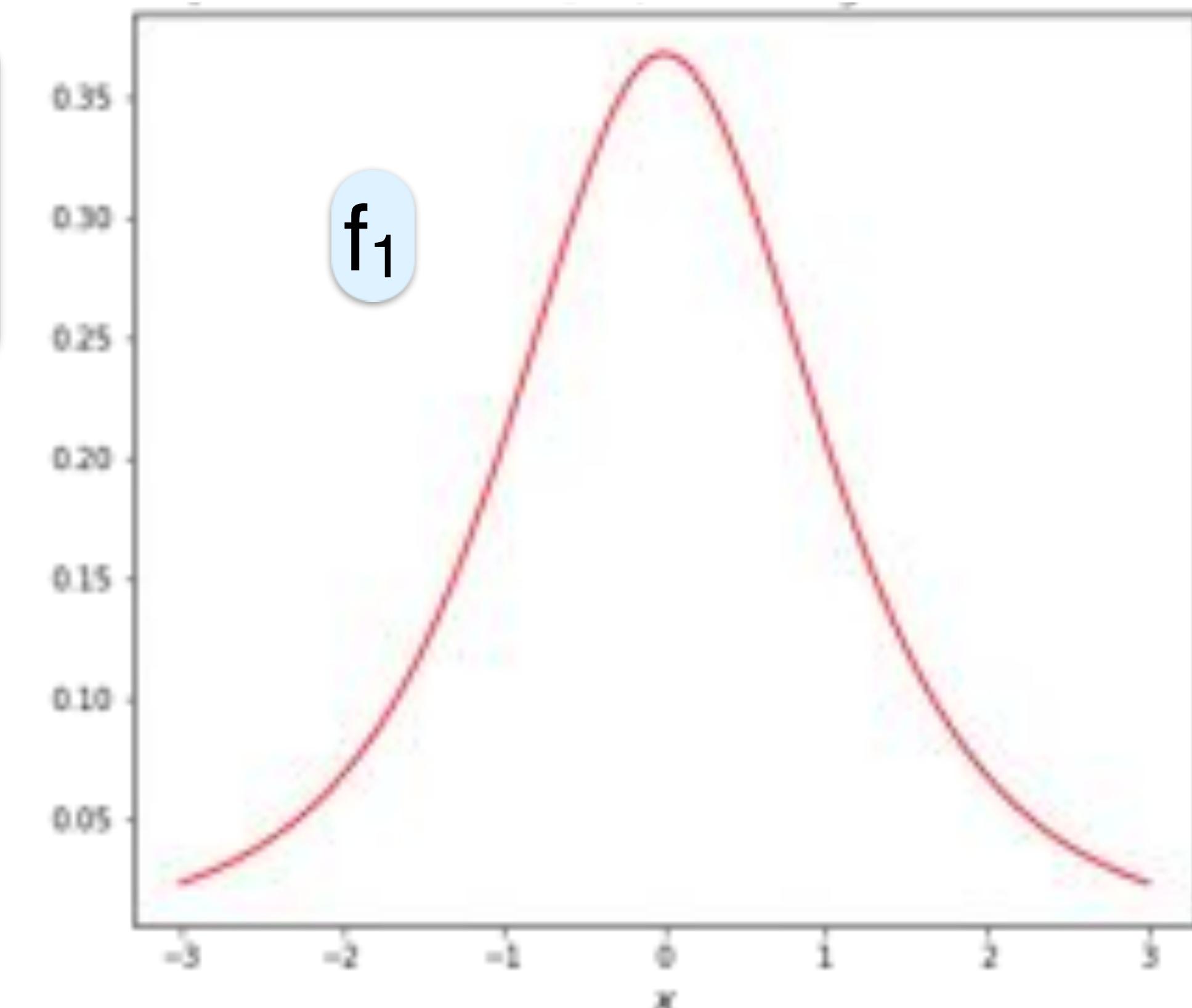
`t.pdf(x, ν)`

$f_3(1)$

```
from scipy.stats import t  
t.pdf(1, 3)  
0.20674833578317203
```

Bell shaped

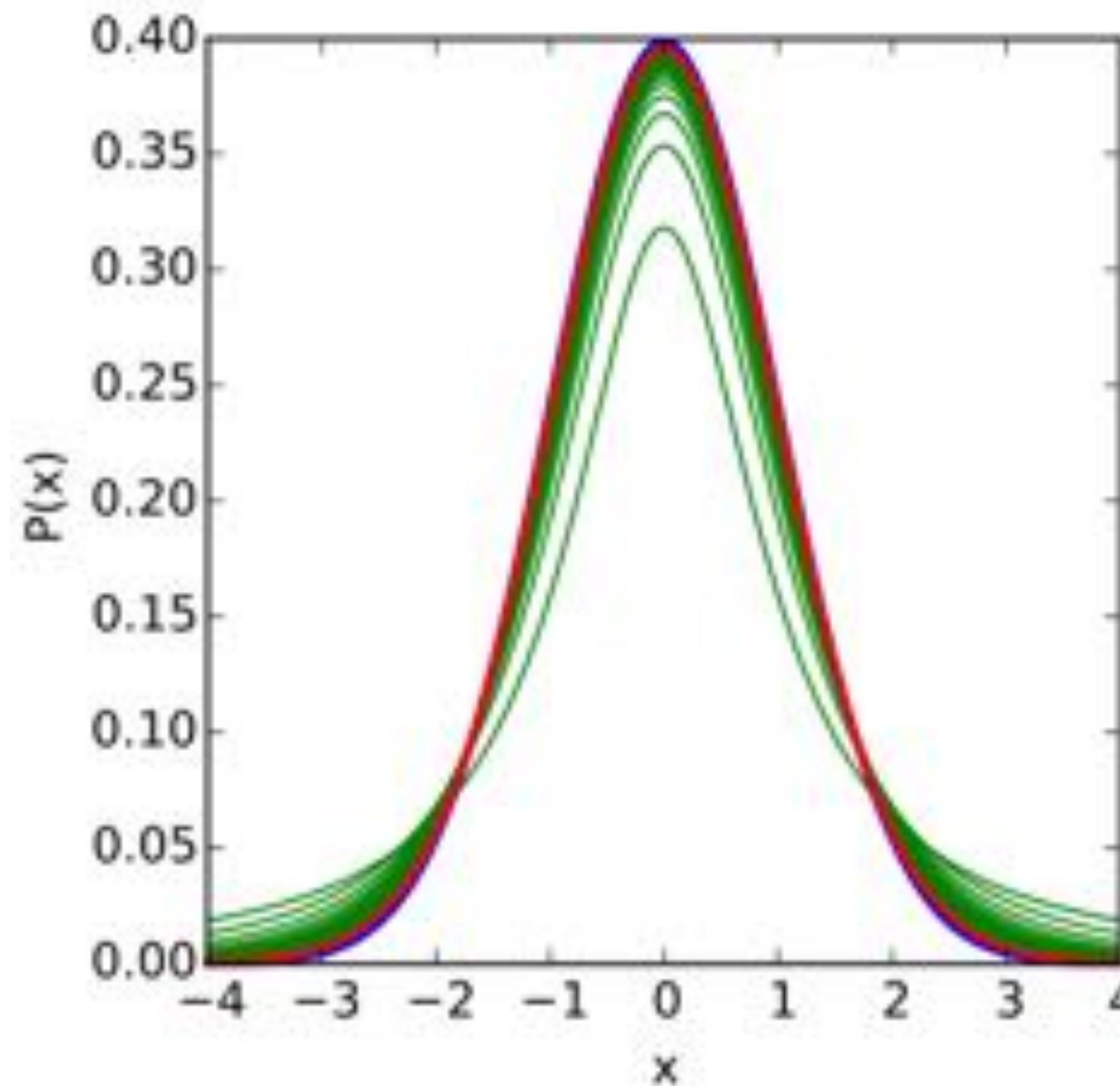
Similar to Gaussian



# Dependence on $\nu$

As  $\nu$  increases

$$f_\nu(t) \rightarrow \phi(t)$$



— standard normal distribution

— t-distribution  $\nu=30$

Logical

$S \rightarrow$

constant

$\sigma$

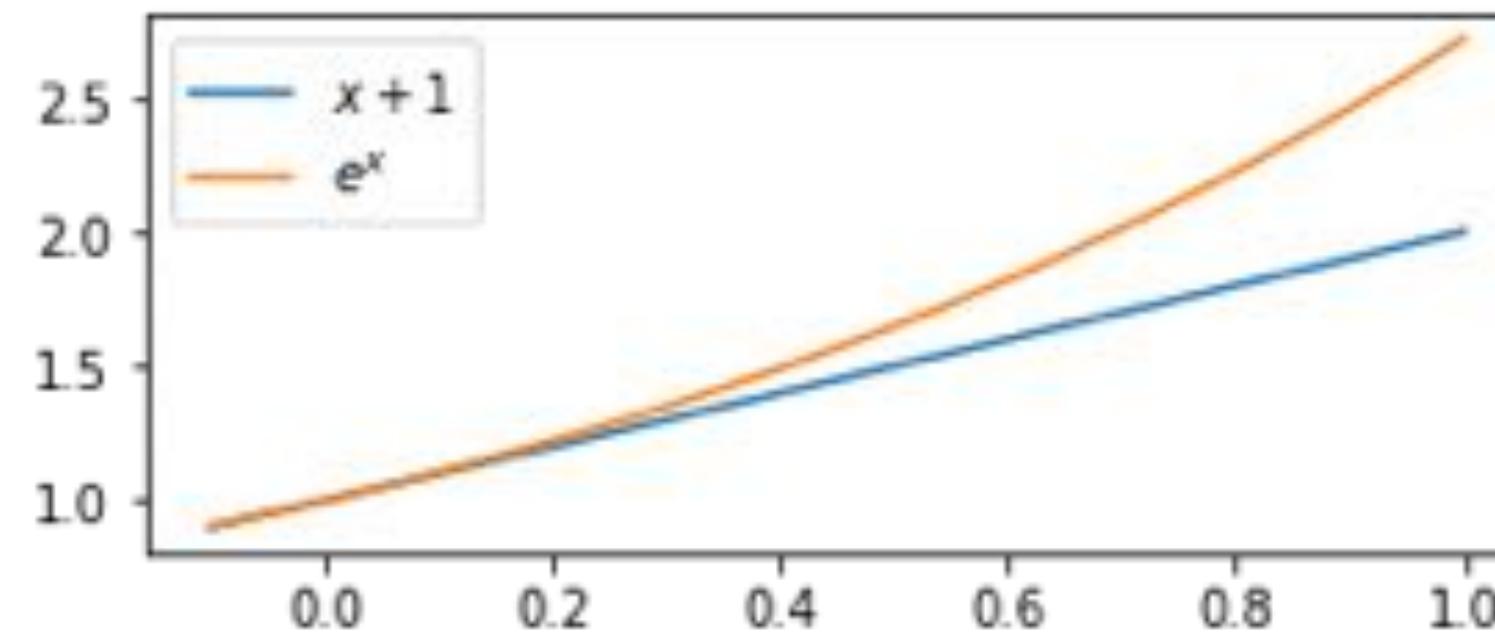
Examples in notebook

# Analytical Argument

$$f_\nu(t) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\cdot\Gamma(\frac{\nu}{2})} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

For small x

$$1 + x \sim e^x$$



As  $\nu \rightarrow \infty$   $\left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}} \approx \left(e^{\frac{t^2}{\nu}}\right)^{-\frac{\nu+1}{2}} \approx e^{-\frac{t^2}{\nu} \frac{\nu+1}{2}} \approx e^{-\frac{t^2}{2}}$

Standard Normal

Distribution

Constant same

$$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \approx \frac{1}{\sqrt{2\pi}}$$

# William Sealy Gosset

Guinness Brewery

World's largest

Billion pints / year

~1.75 Billion bottles

Quality

Consistency

Statisticians

Trade secret

Publish

Pseudonym



1876-1937

VOLUME VI

MARCH, 1908

No. 1

## BIOMETRIKA.

### THE PROBABLE ERROR OF A MEAN.

By STUDENT.

*Introduction.*

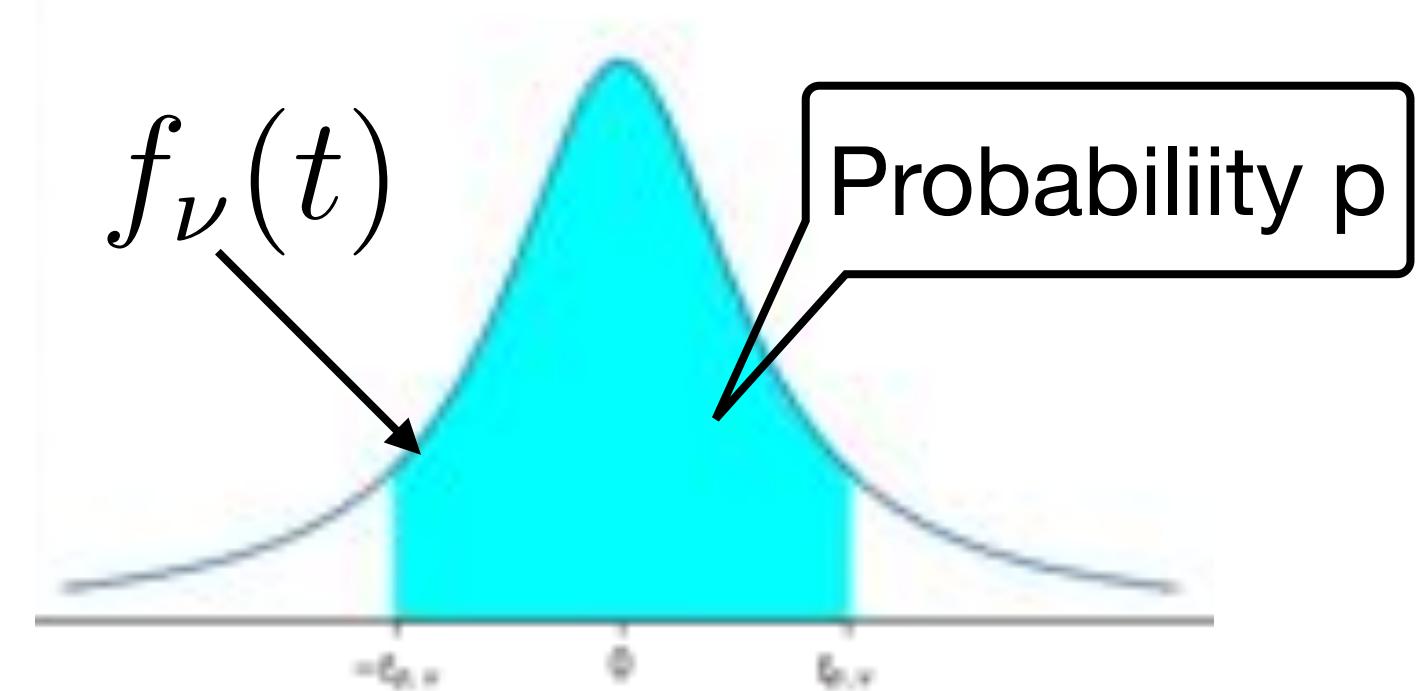
ANY experiment may be regarded as forming an individual of a "population" of experiments which might be performed under the same conditions. A series of experiments is a sample drawn from this population.

# Déjà vu

$T_\nu$  Student's t-distribution,  $\nu$  degrees of freedom

Critical value,  $t$

$$t_{p,\nu} : P(|T_\nu| \leq t_{p,\nu}) = p$$

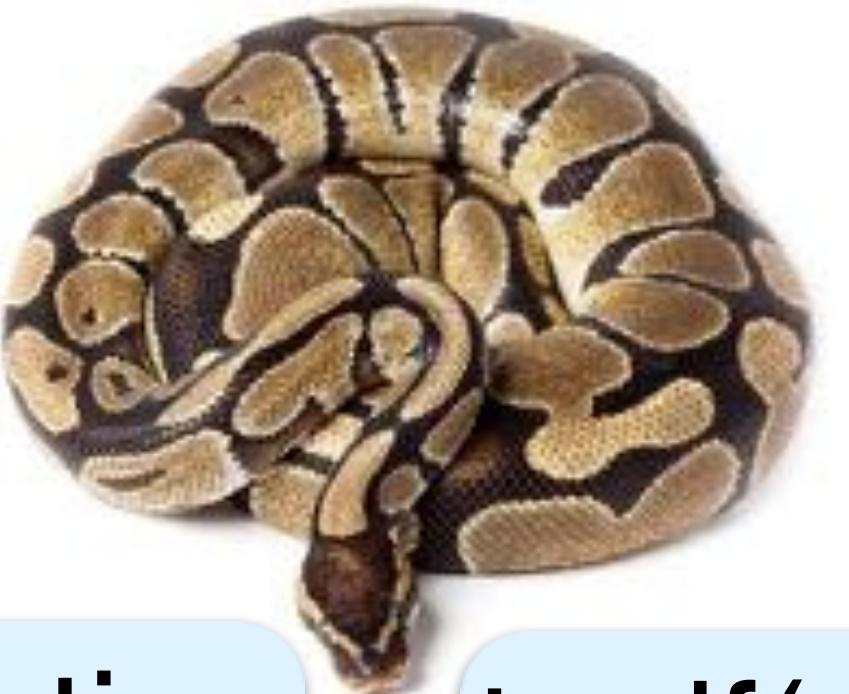


$$p = P(|T_\nu| \leq t_{p,\nu}) = P(-t_{p,\nu} \leq T_\nu \leq t_{p,\nu}) = 2F(t_{p,\nu}) - 1$$

$$F(t_{p,\nu}) = \frac{1+p}{2}$$

$$t_{p,\nu} = F_\nu^{-1} \left( \frac{1+p}{2} \right)$$

CDF of  $T_\nu$



Degrees of freedom

Cumulative distribution function

`t.cdf(x, ν)`

$F_3(1)$

```
from scipy.stats import t  
t.cdf(1, 3)  
0.80449889052211476
```

Inverse cdf

percent point function

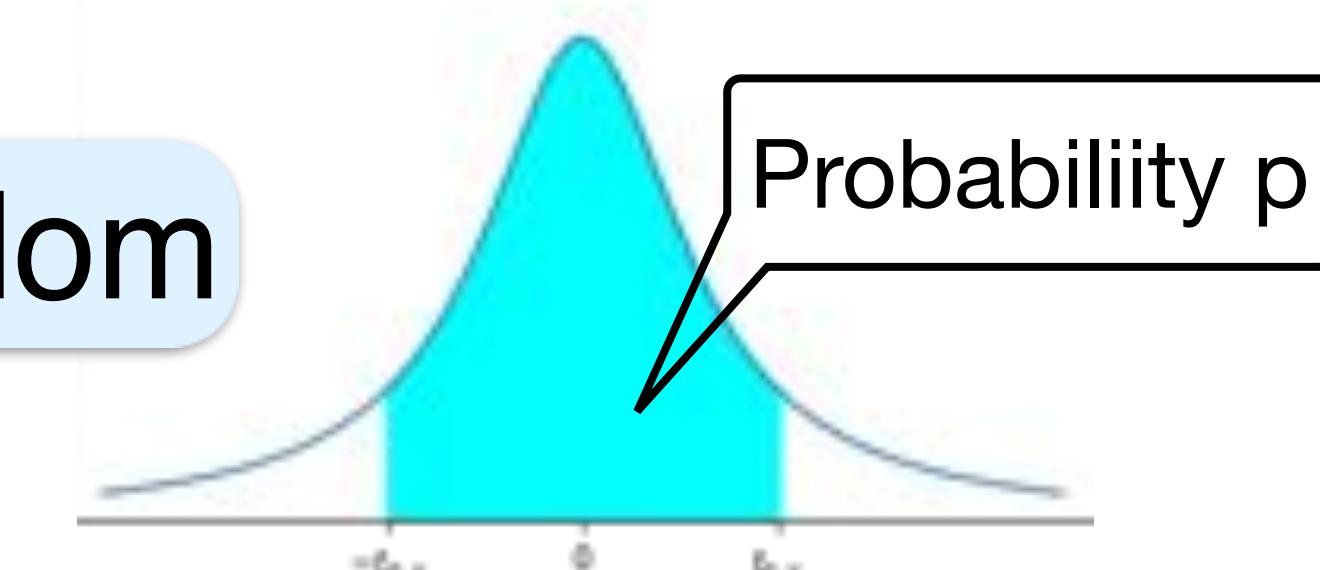
`t.ppf(x, ν)`

$F_3^{-1}(0.95)$

```
t.ppf(0.95, 3)  
2.3533634348018264
```

# Confidence Interval

t-distribution,  $\nu$  degrees of freedom



t-statistic

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim f_{n-1}(t) \quad P\left(\left|\frac{\bar{X} - \mu}{S/\sqrt{n}}\right| \leq t_{p,n-1}\right) = p$$

$$P\left(|\bar{X} - \mu| \leq t_{p,n-1} \frac{S}{\sqrt{n}}\right) = p$$

With probability p

$$|\bar{X} - \mu| \leq t_{p,n-1} \frac{S}{\sqrt{n}}$$

margin of error

$$\mu \in \left[\bar{X} - t_{p,n-1} \frac{S}{\sqrt{n}}, \bar{X} + t_{p,n-1} \frac{S}{\sqrt{n}}\right]$$

# Confidence → Interval

Given confidence p samples  $X_1, \dots, X_n$

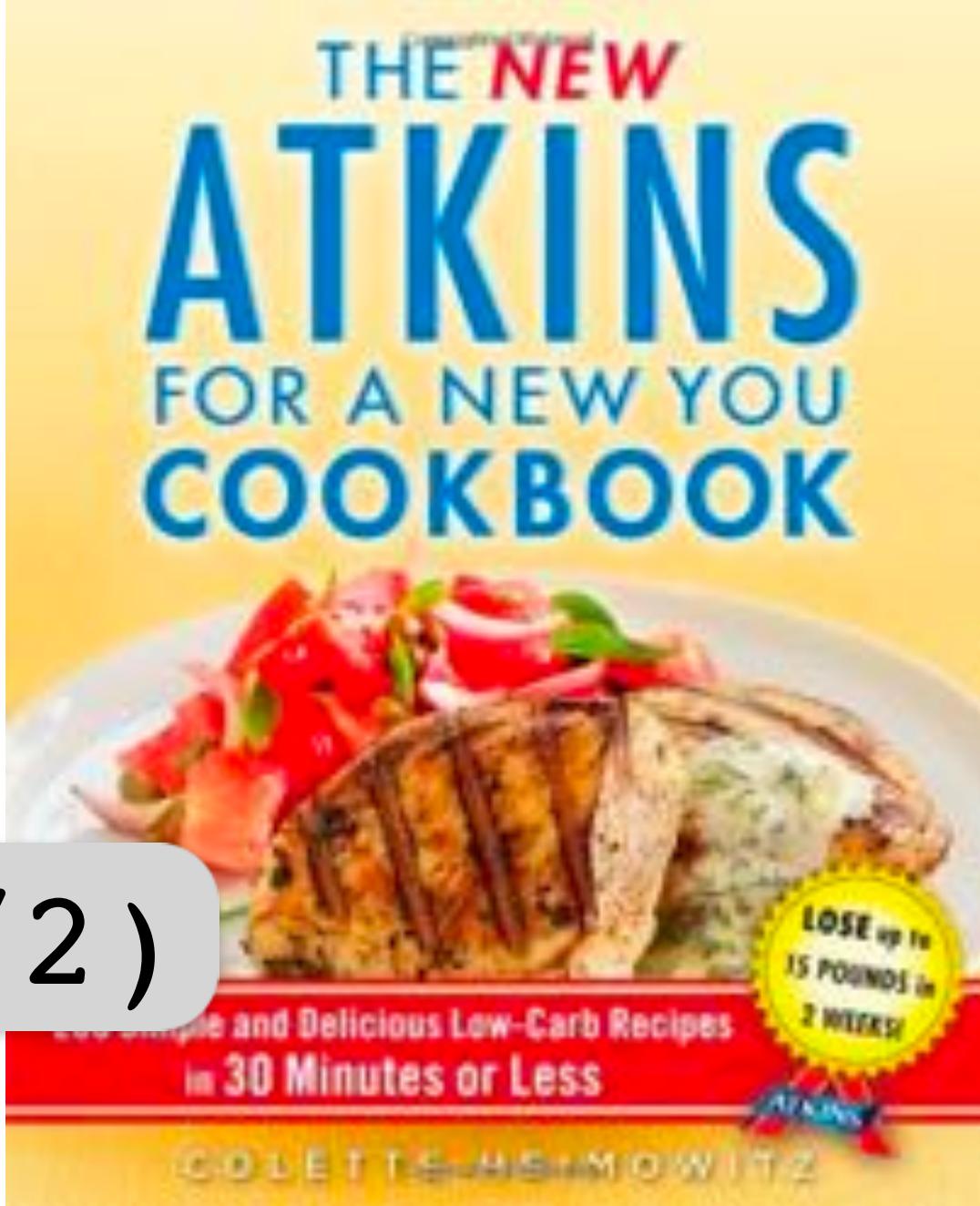
Determine critical t  $t_{p,n-1} = F_{n-1}^{-1}\left(\frac{1+p}{2}\right)$   $t \cdot \text{ppf}((1+p)/2)$

Sample mean  $\bar{X} = \frac{X_1 + \dots + X_n}{n}$

Sample variance  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

Margin of error  $t_{p,n-1} \frac{S}{\sqrt{n}}$   $\sigma$  unnecessary

Confidence interval  $\left[ \bar{X} - t_{p,n-1} \frac{S}{\sqrt{n}}, \bar{X} + t_{p,n-1} \frac{S}{\sqrt{n}} \right]$



# Mature African Elephant Trunk Length



# Mature African Elephant Trunk Length

8 measurements

5.62, 6.07, 6.64, 5.91, 6.30, 6.55, 6.19, 5.48 feet

Find

95% confidence interval for distribution mean

Critical t  $t_{p,n-1} = F_{n-1}^{-1}\left(\frac{1+p}{2}\right) = F_7^{-1}(0.975) \approx 2.3646$  `t.ppf(0.975, 7)`

2.3646

Sample mean

$$\bar{X} = 6.095$$

Sample variance

$$S^2 \approx 0.1705 \quad S = 0.4130$$

Margin of error

$$t_{p,n-1} \frac{S}{\sqrt{n}} \approx 0.3453$$

Confidence interval

$$\left(\bar{X} - t_{p,n-1} \frac{S}{\sqrt{n}}, \bar{X} + t_{p,n-1} \frac{S}{\sqrt{n}}\right) \\ \approx (5.7497, 6.4403)$$

# Observations

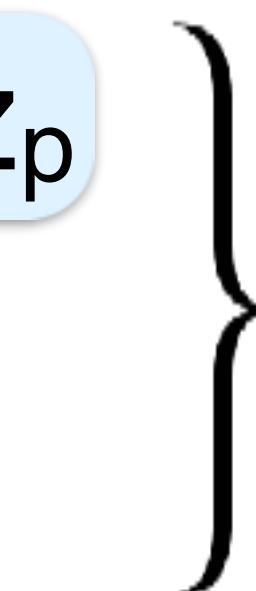
n large

$$f_{n-1}(t) \rightarrow \phi(t)$$

$$t_{p,n-1} \rightarrow z_p$$

$$S \rightarrow$$

$$\sigma$$



Can use z-based techniques

n small

t-distribution more accurate

Yields larger margin of error than known  $\sigma$

Assumed  $X_i \sim N$ , best when this roughly holds

# Confidence Intervals Unknown $\sigma$



Estimate  $\sigma$

Student's t-distribution

Step-by-step instructions

Example

# Confidence Intervals



# Points → Intervals

Distribution or population

Estimate parameters



Point estimates

$\mu \approx 3.14$

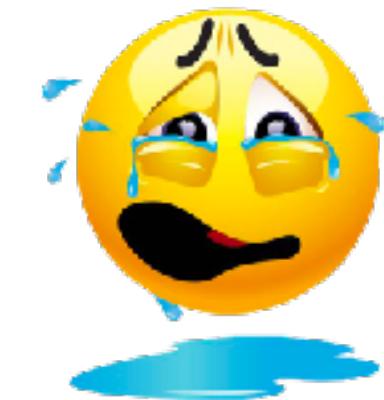
$p \approx 0.48$



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  $\mu$ , and stdv  $\sigma$

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

Sample mean

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

Typically  $\geq 30$

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

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

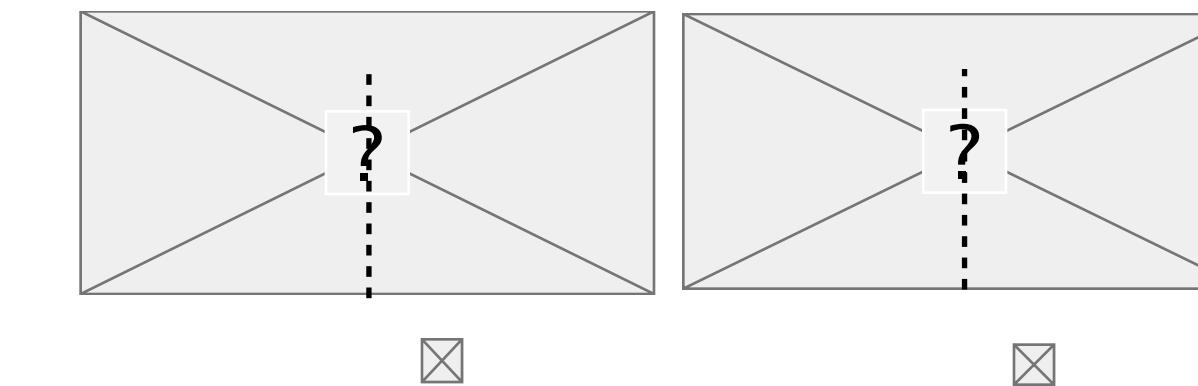
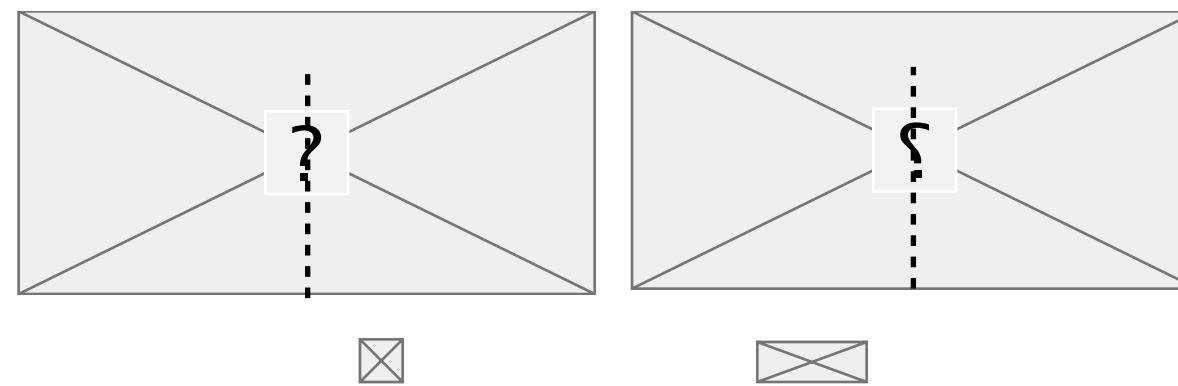
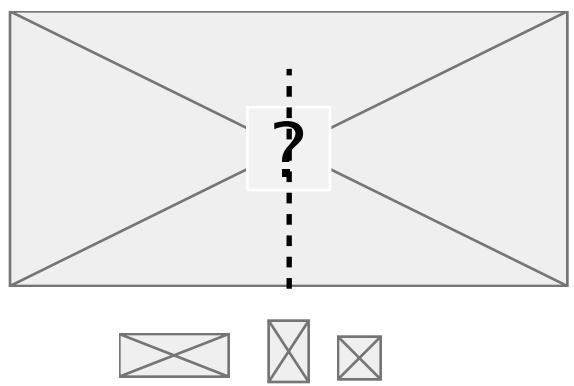
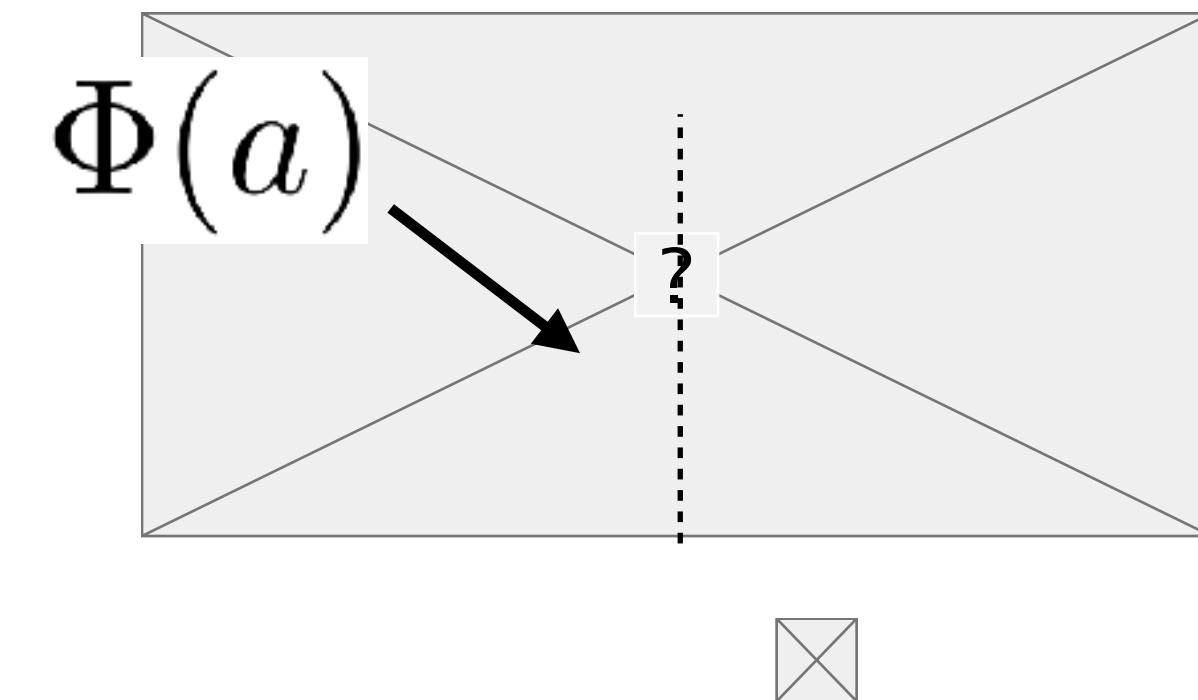
$$P(-a \leq Z \leq a) > 0$$

P = ?

# Interval Probability

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

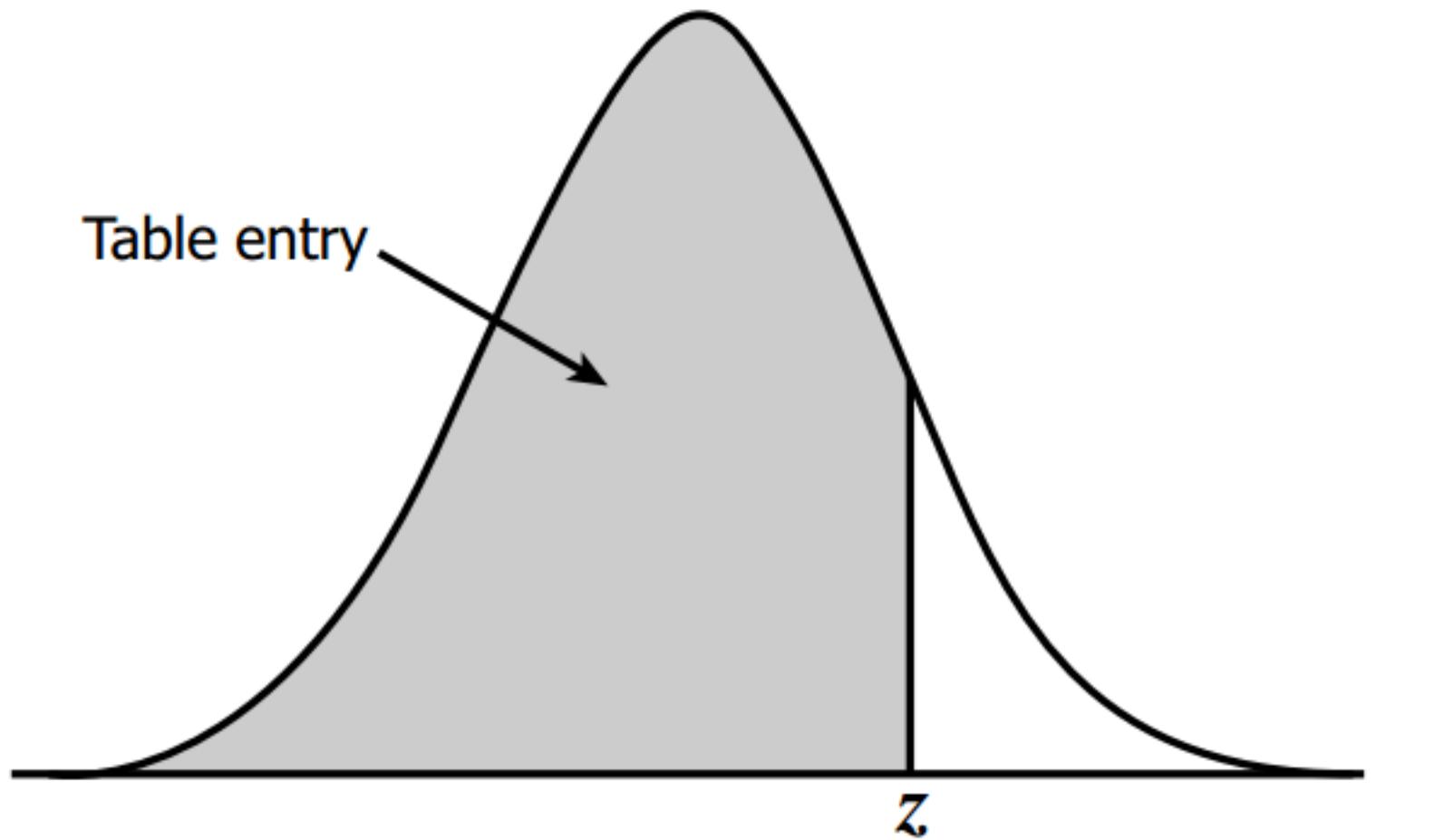
$$\Phi(a) = F(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-t^2/2} dt$$



# 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
```

$\Phi(2)$

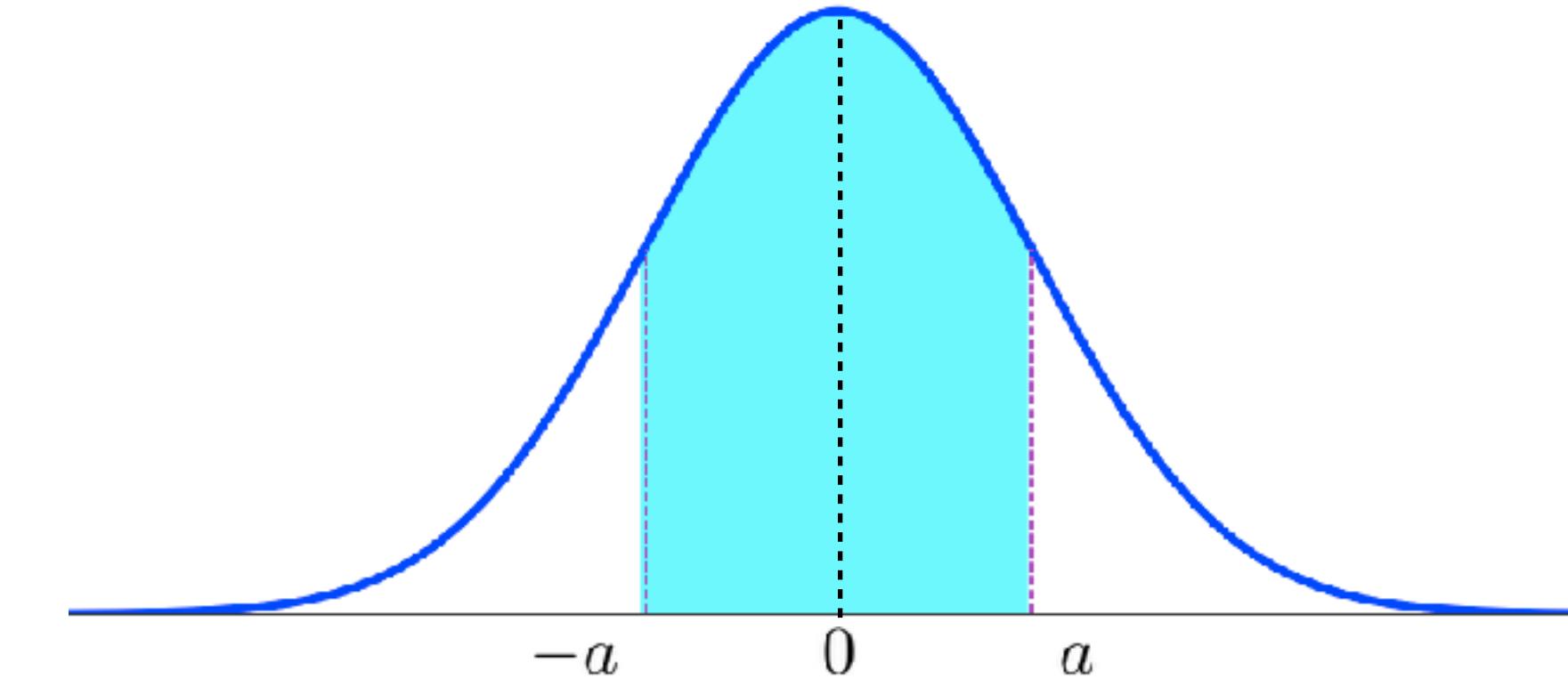
```
norm.cdf(2)  
0.977249868052
```

$\Phi(3)$

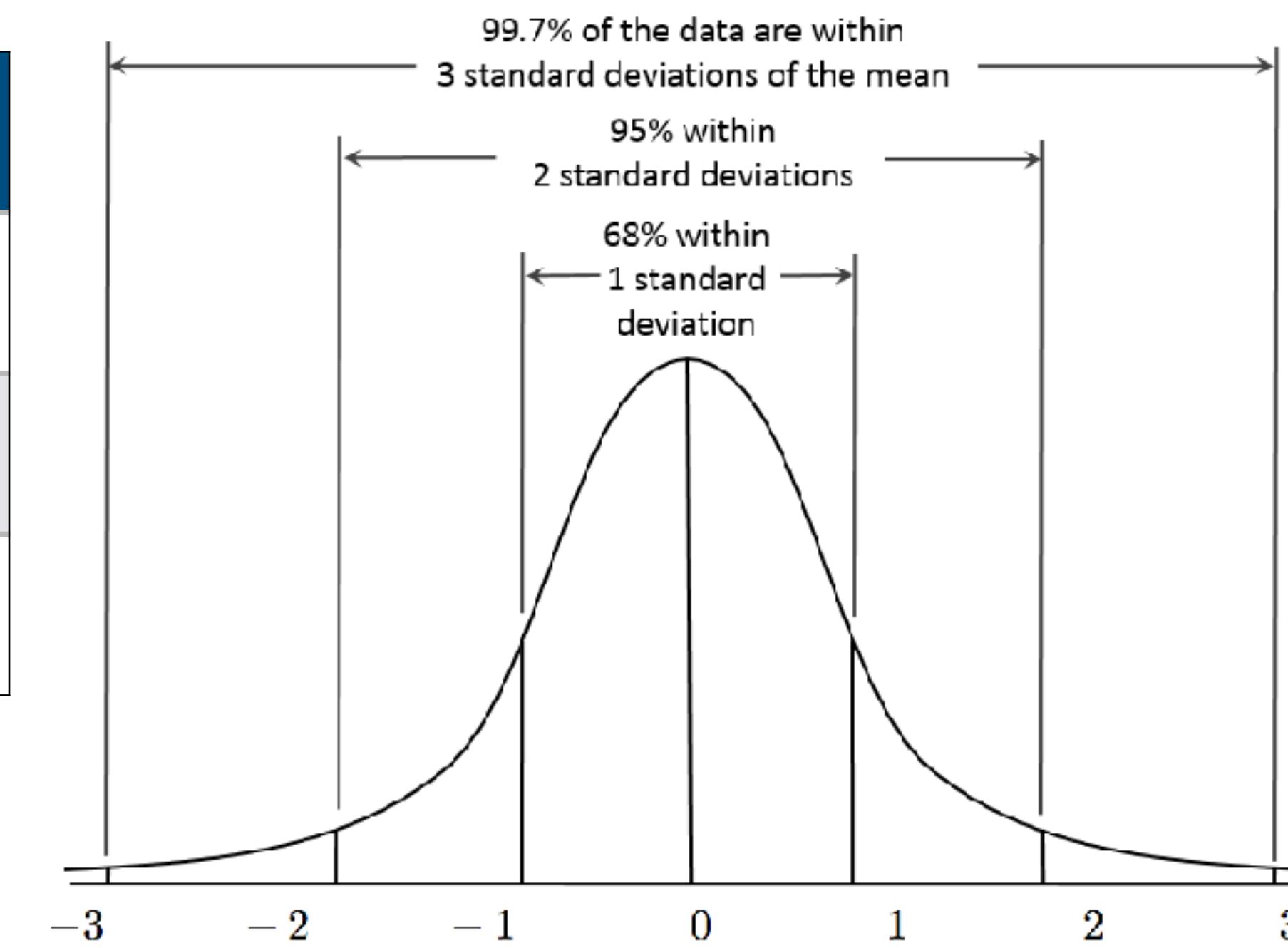
```
norm.cdf(3)  
0.998650101968
```

# 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

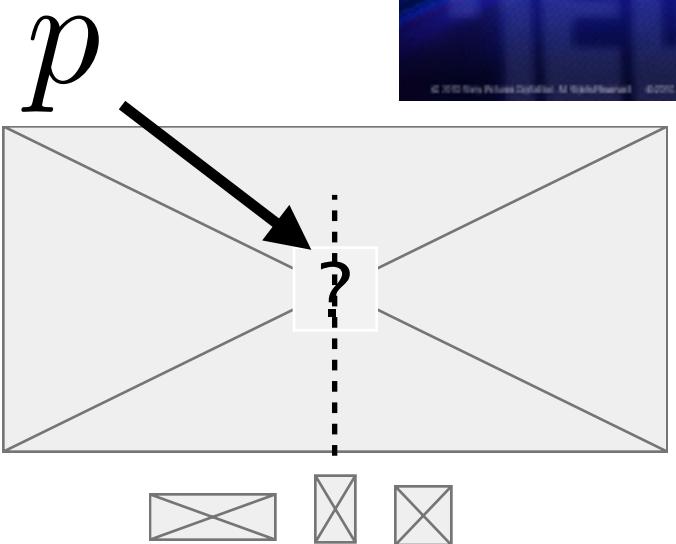
Given desired probability  $p$

Find

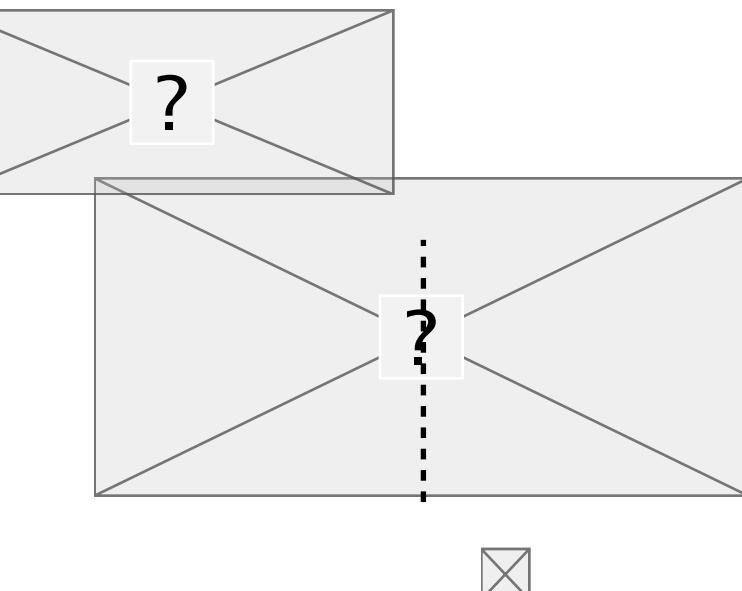
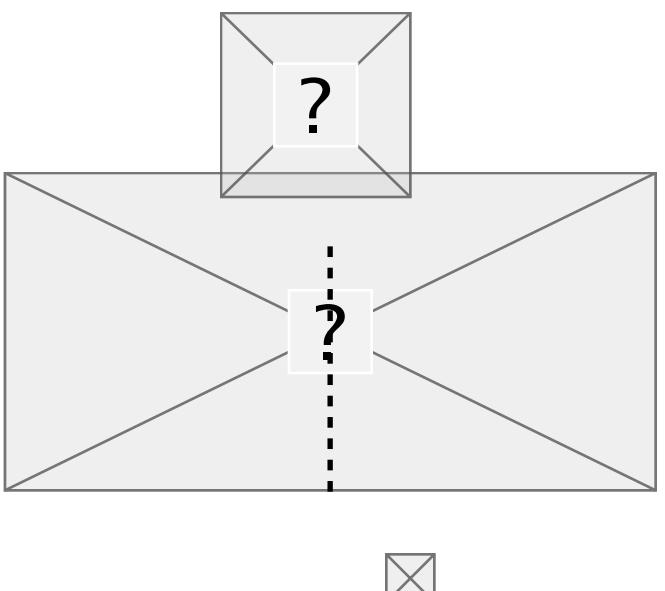
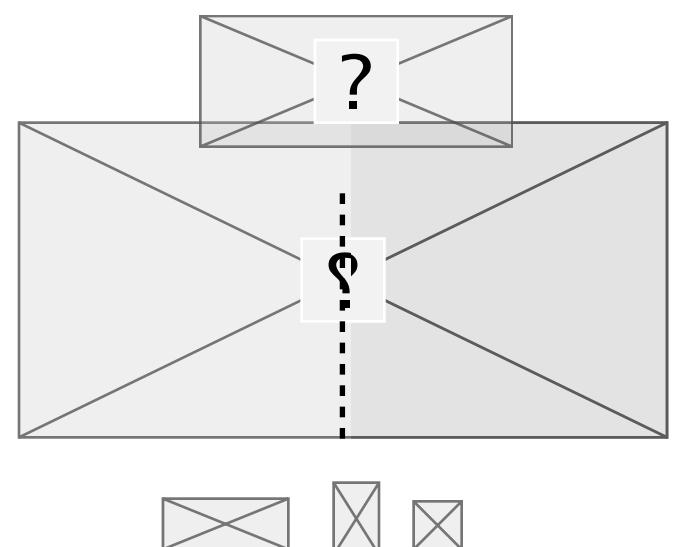
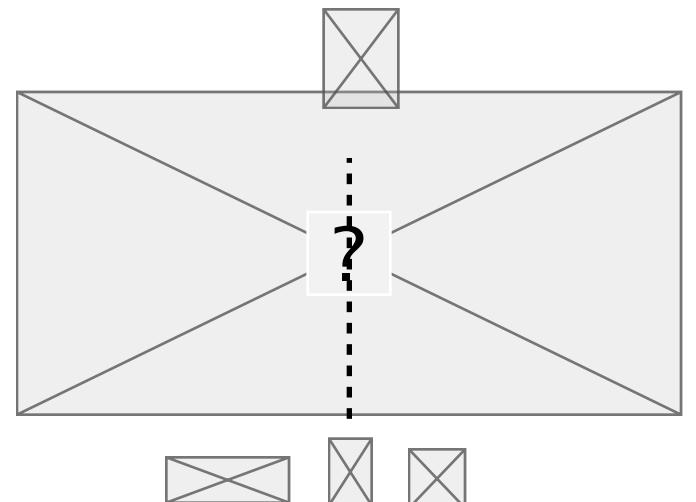
$a$

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

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

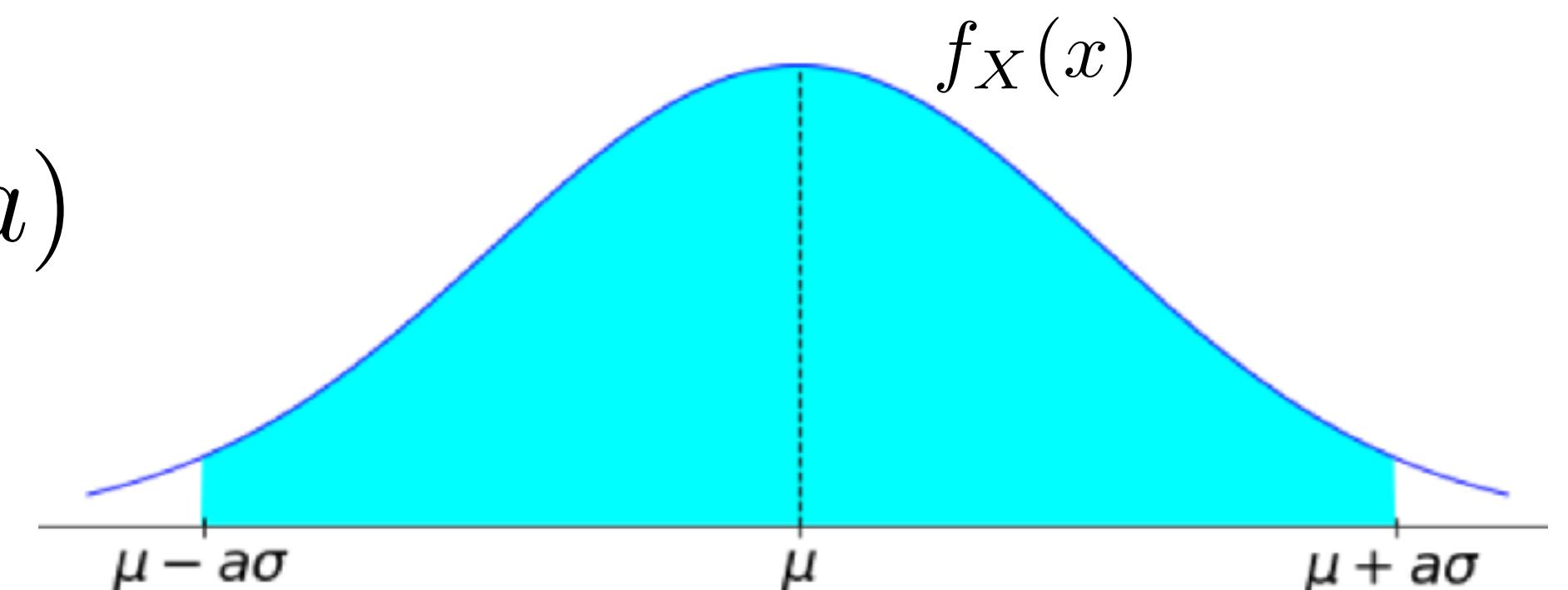
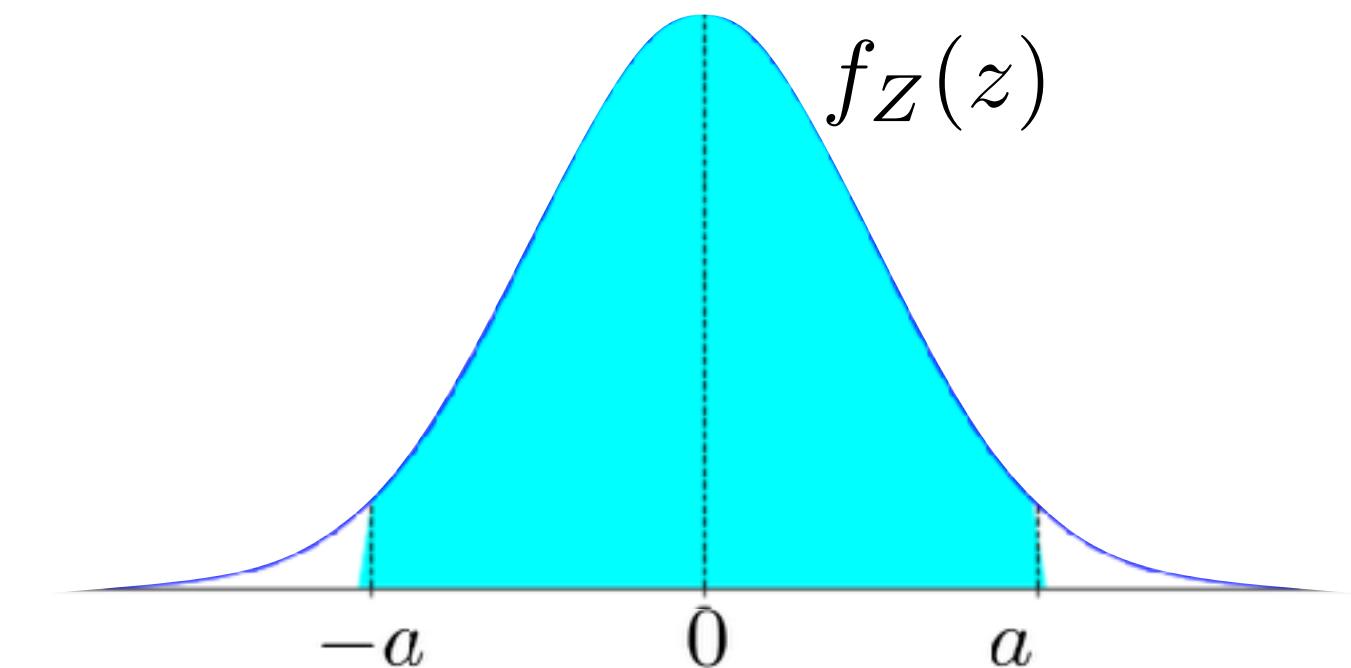
X within “a” std  
from its mean

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

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

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  $\mu$

proportion  $p$

Given a sample  $X_1, \dots, X_n$

Find an interval containing  $\mu$

First

$\sigma$  known

Next lecture

$\sigma$  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

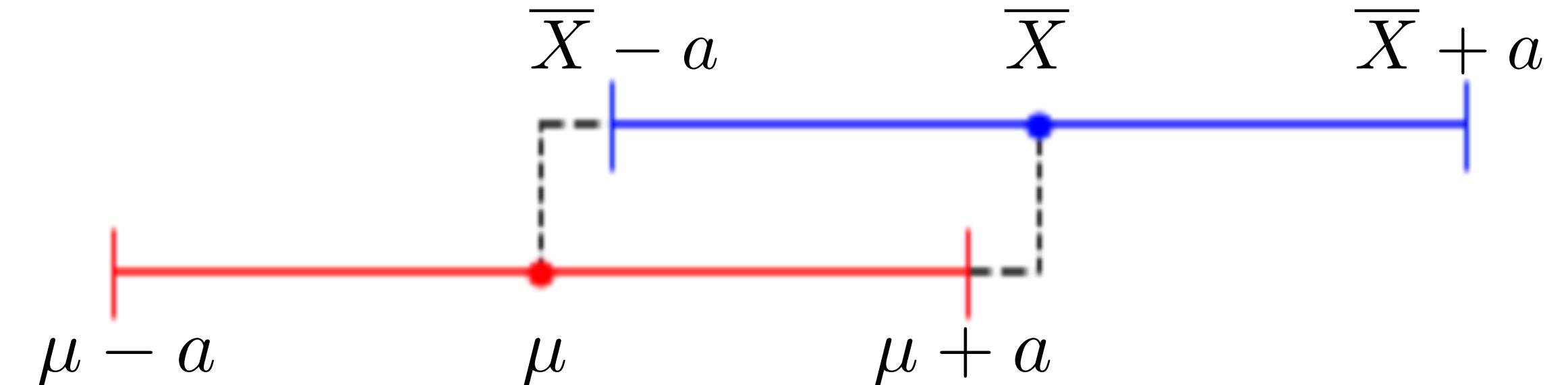
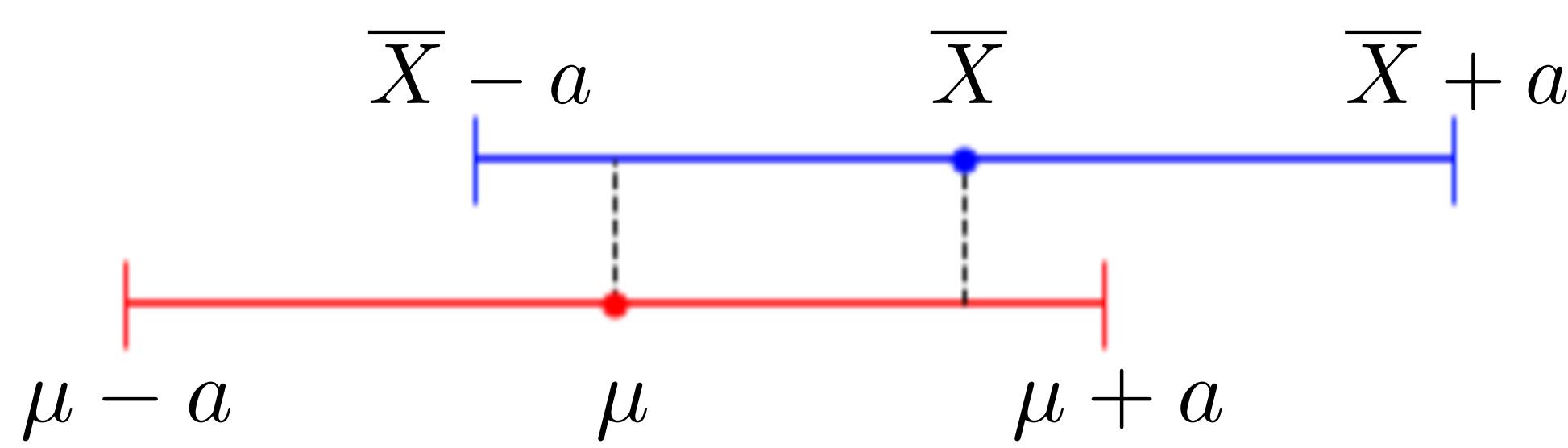
Standard  
error

# Proximity is Reciprocal

←  $\bar{X}$  near  $\mu$

→  $\mu$  near  $\bar{X}$

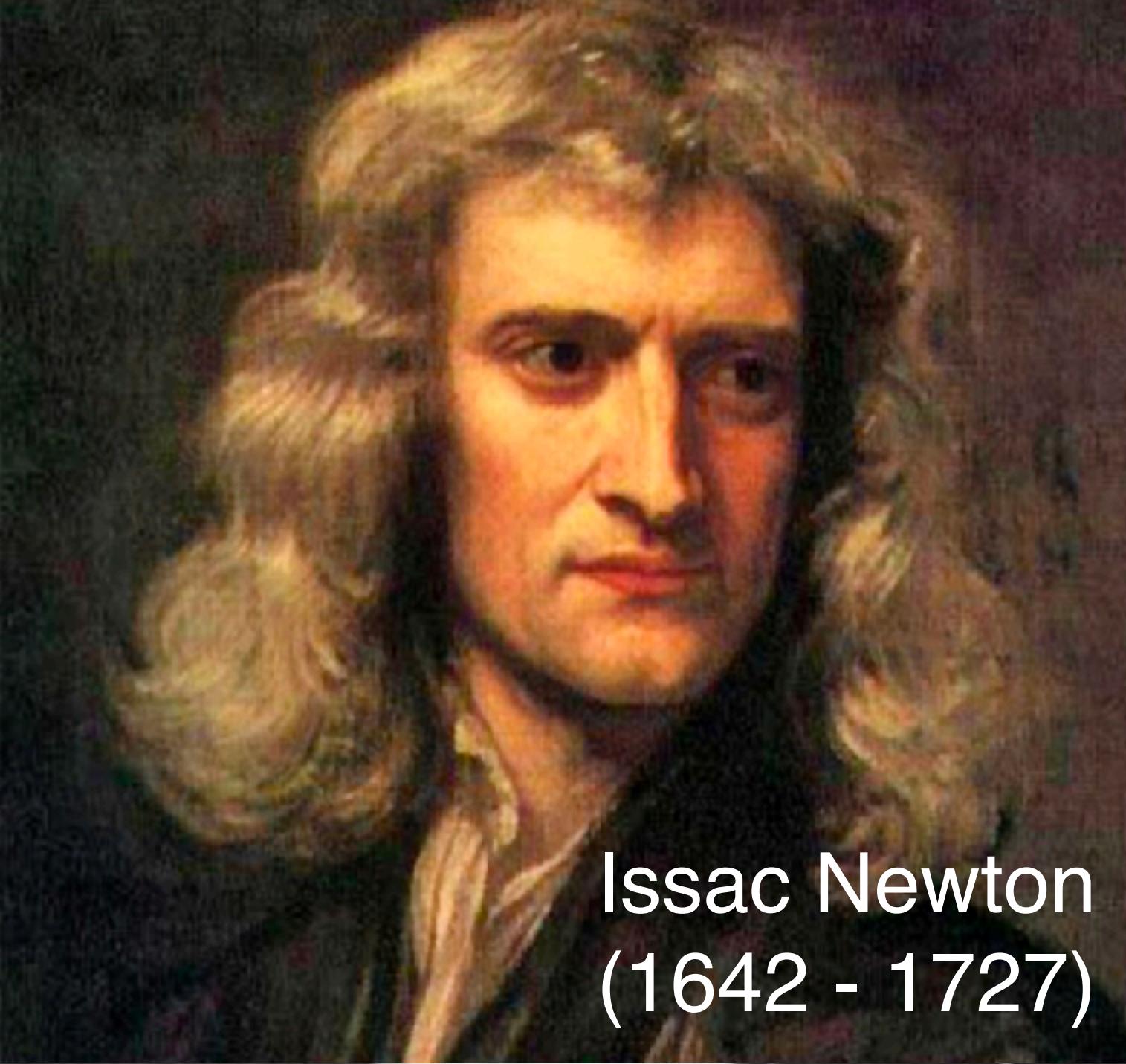
$$\bar{X} \in (\mu - a, \mu + a) \quad |\bar{X} - \mu| < a \quad \mu \in (\bar{X} - a, \bar{X} + a)$$



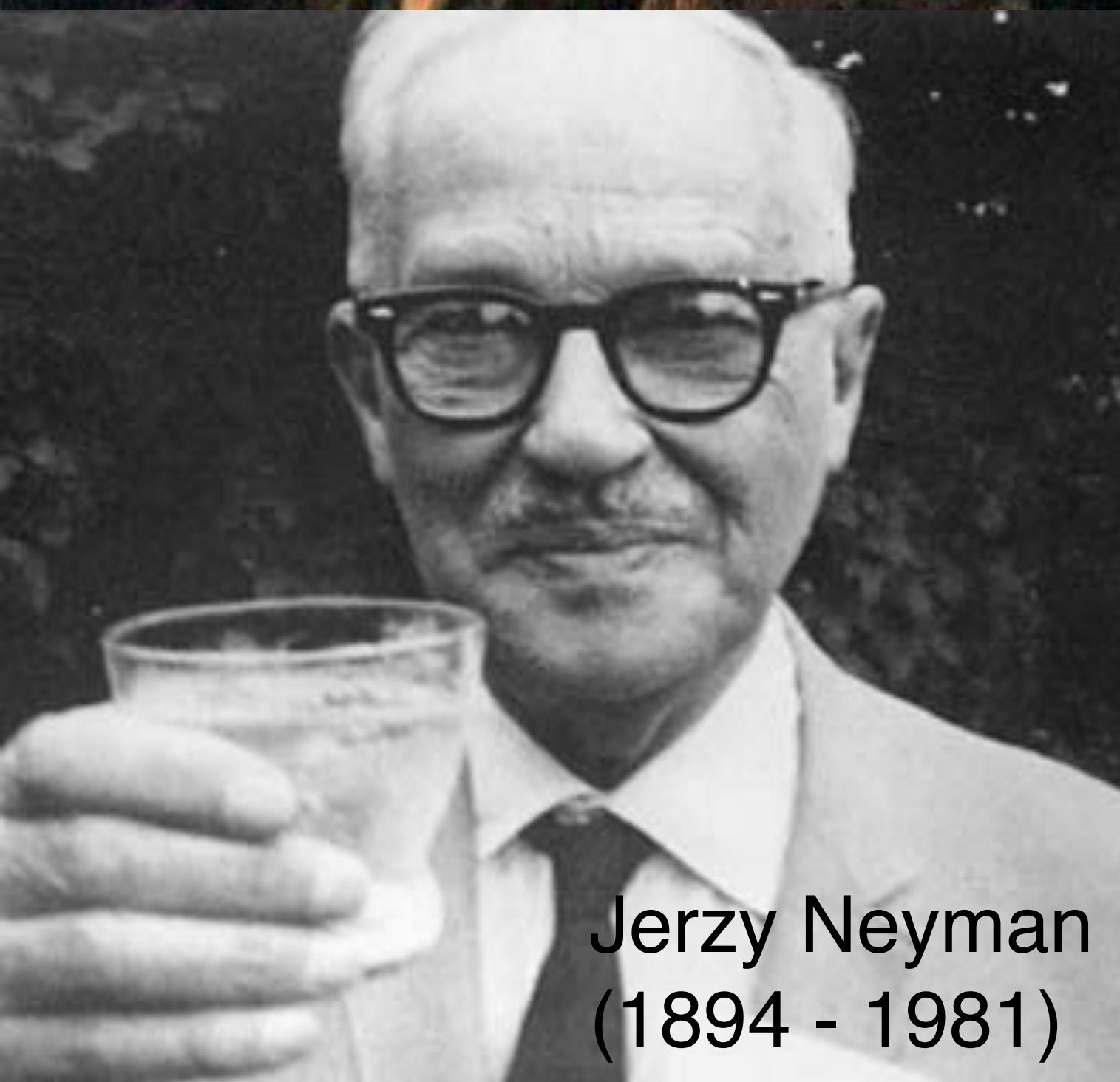
$$\bar{X} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$$

With high probability  $\bar{X}$  near  $\mu$

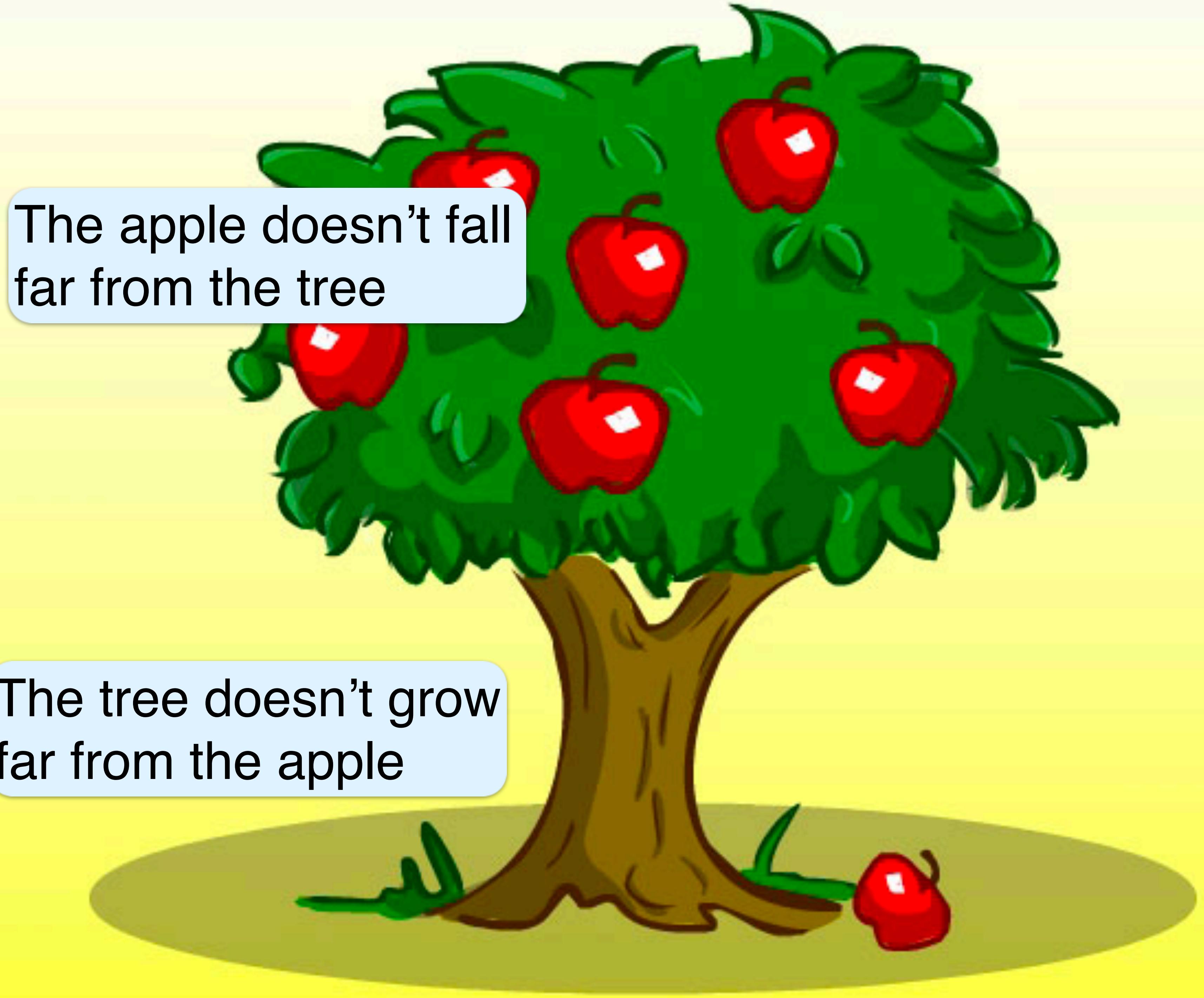
With high probability  $\mu$  near  $\bar{X}$



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

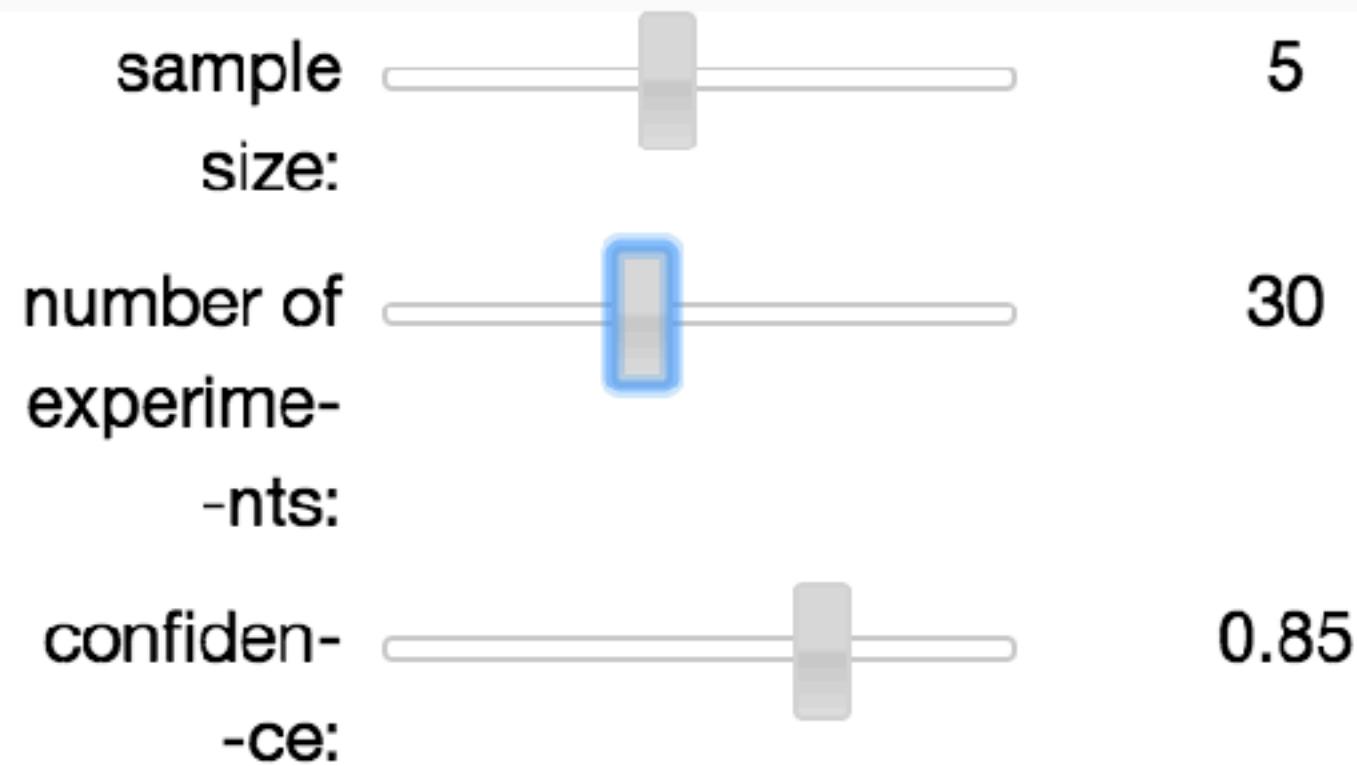
# Confidence Interval

With probability p

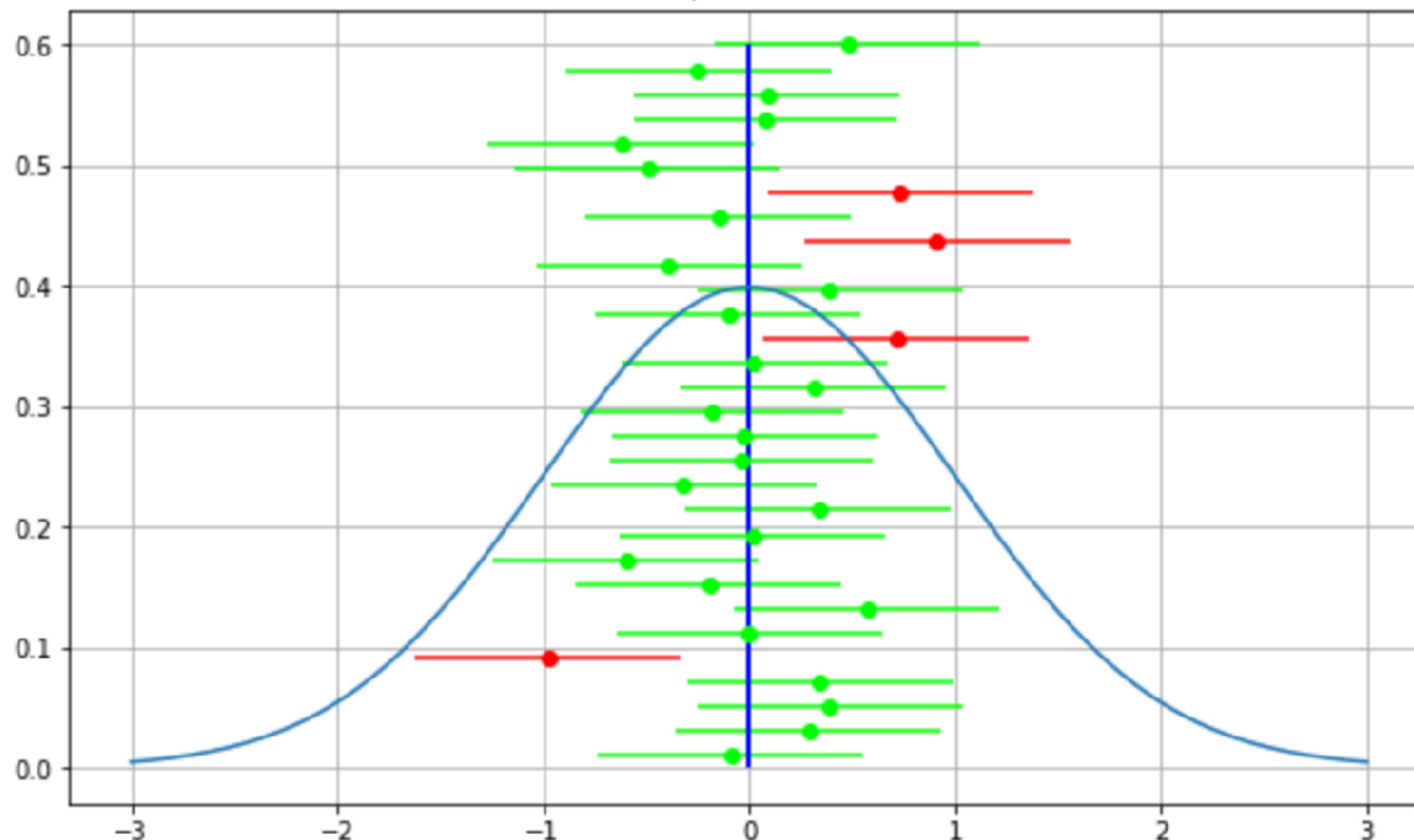
$$\bar{X} \in (\mu - z_p \sigma_{\bar{X}}, \mu + z_p \sigma_{\bar{X}})$$

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

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



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

