# Using the normal distribution for inference
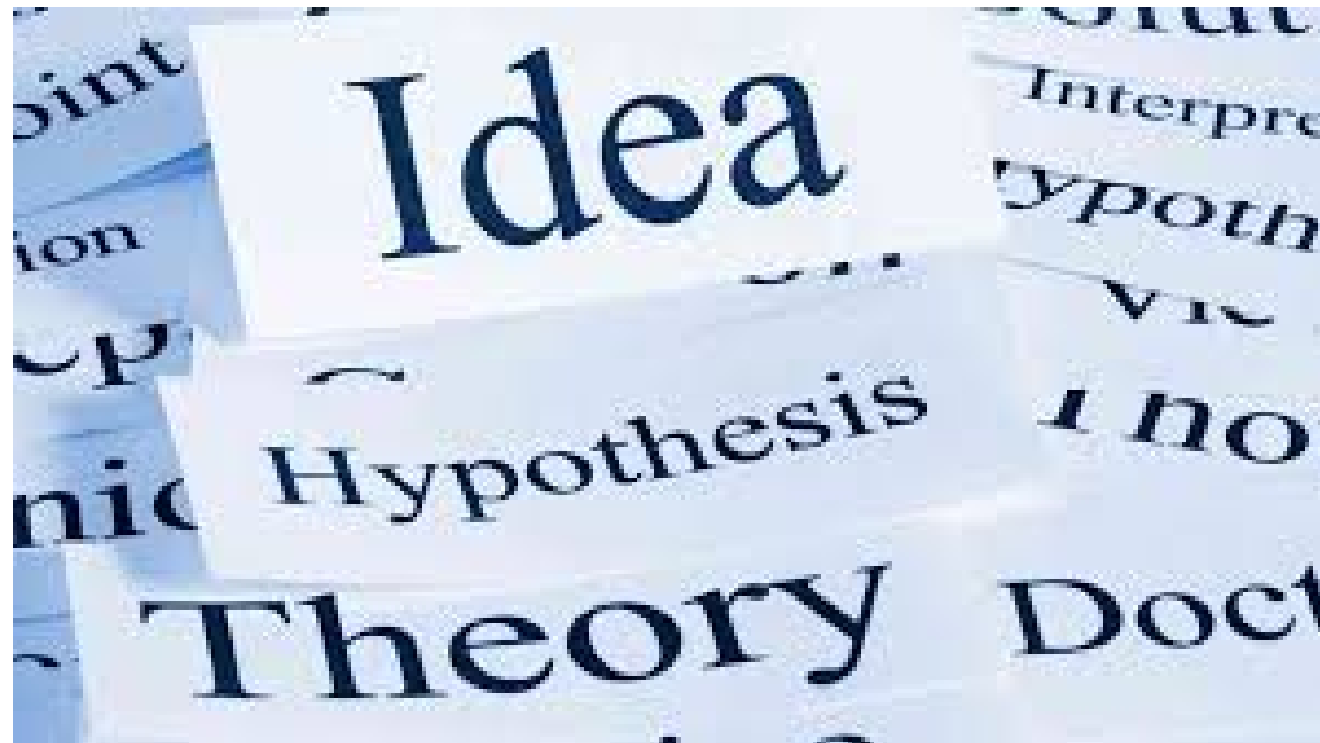
# Overview

Quick review of theories/concepts in hypothesis testing

Review and continuation of the normal distribution

Using the normal distribution for inference
- Hypothesis tests
- Confidence intervals

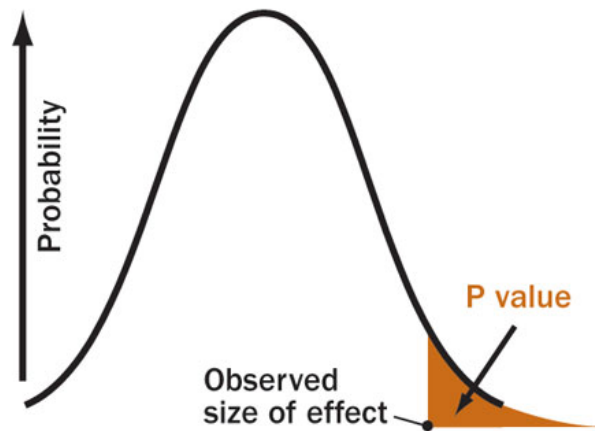# Quick review of theories of hypothesis tests

# Two theories of hypothesis testing

1. **Significance testing** of Ronald Fisher
   - p-value as strength of evidence against the null hypothesis

2. **Hypothesis testing** of Jezy Neyman and Egon Pearson
   - Make a formal decision of whether to reject $H_0$ (if p-value < predefined $\alpha$ value)
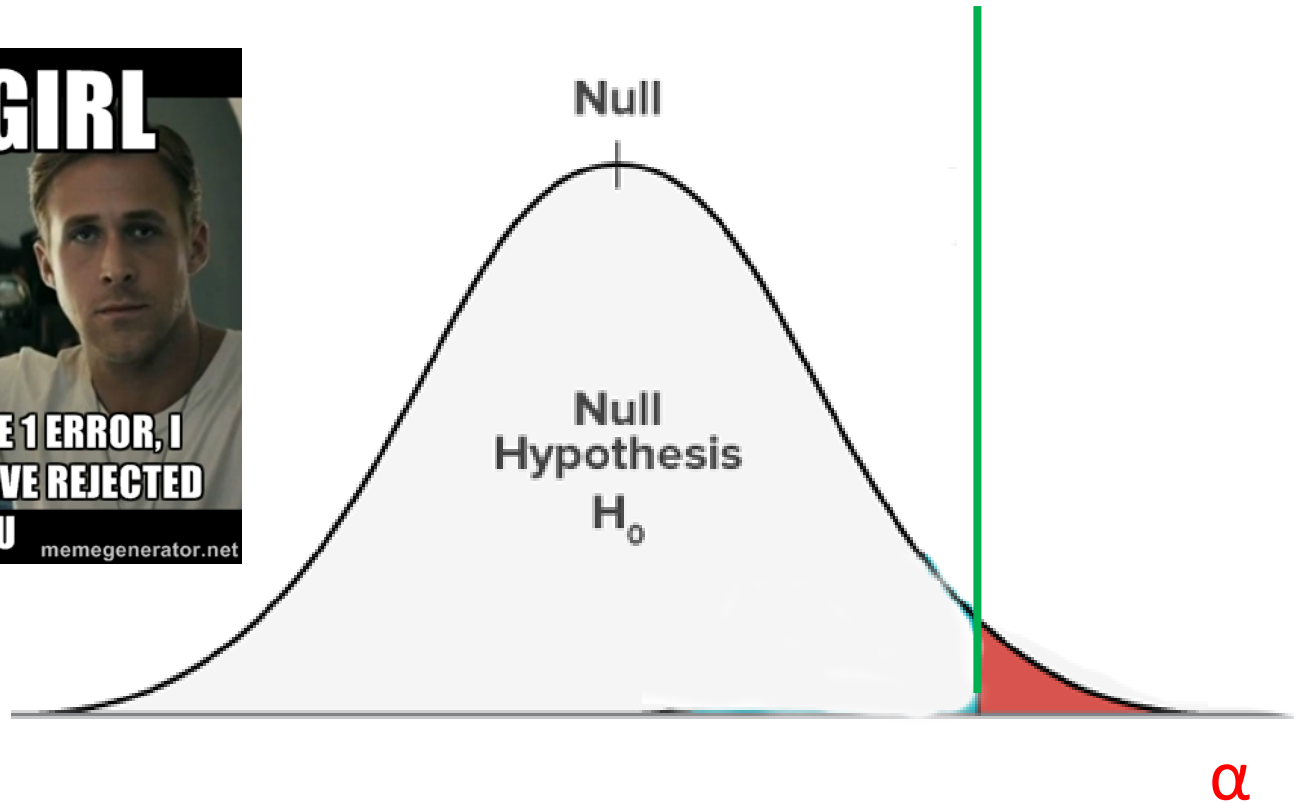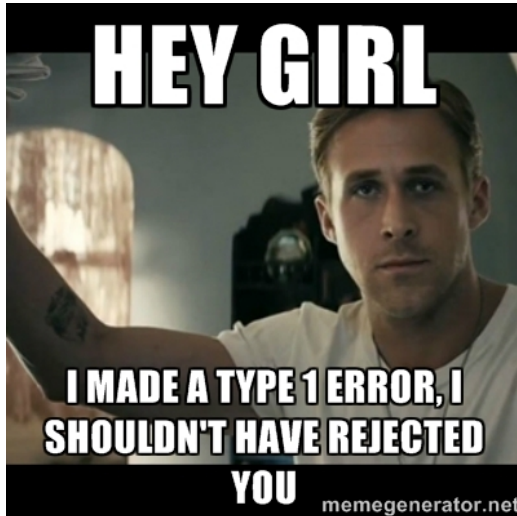


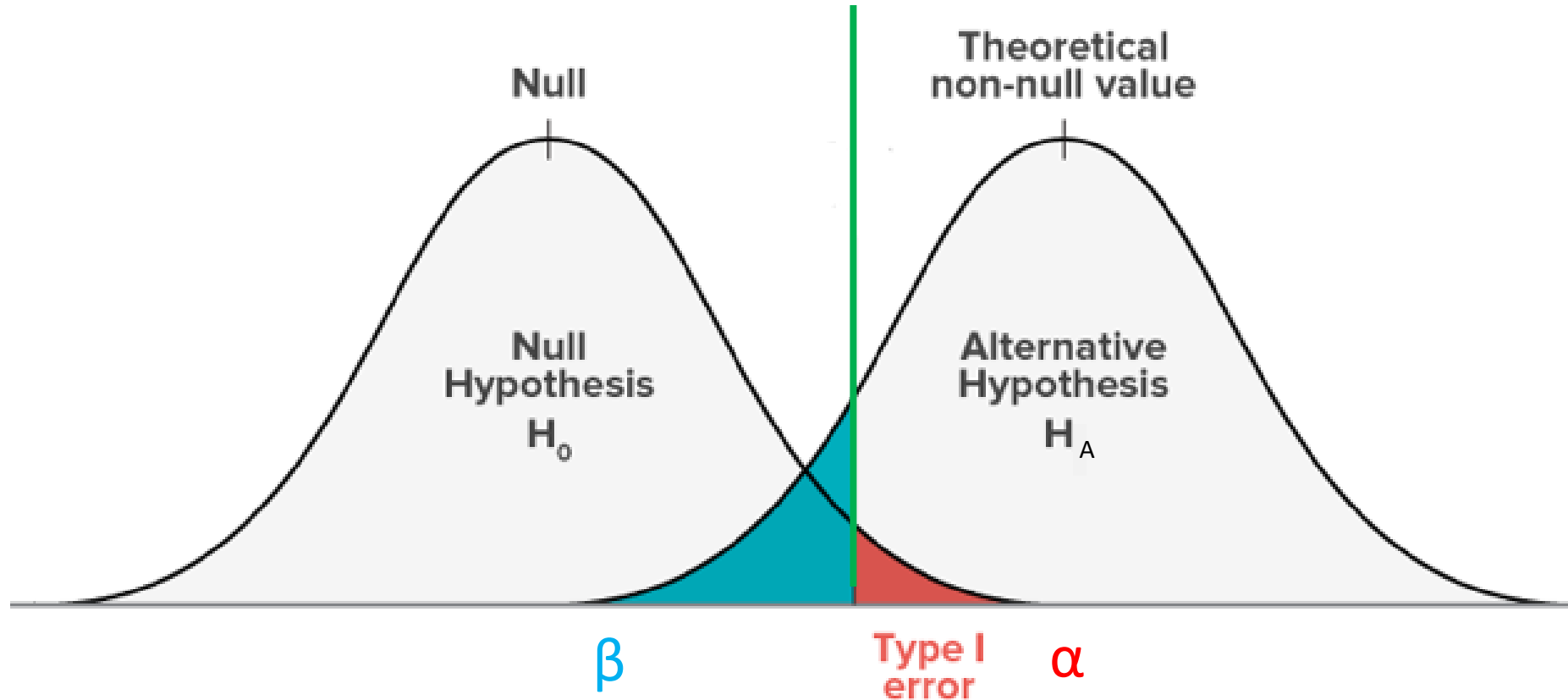**Significance testing**     **Hypothesis testing**

# Neyman-Pearson Frequentist logic



If Neyman-Pearson null hypothesis testing paradigm was followed perfectly, then only ~5% of all published research findings should be wrong (for $\alpha = 0.05$)
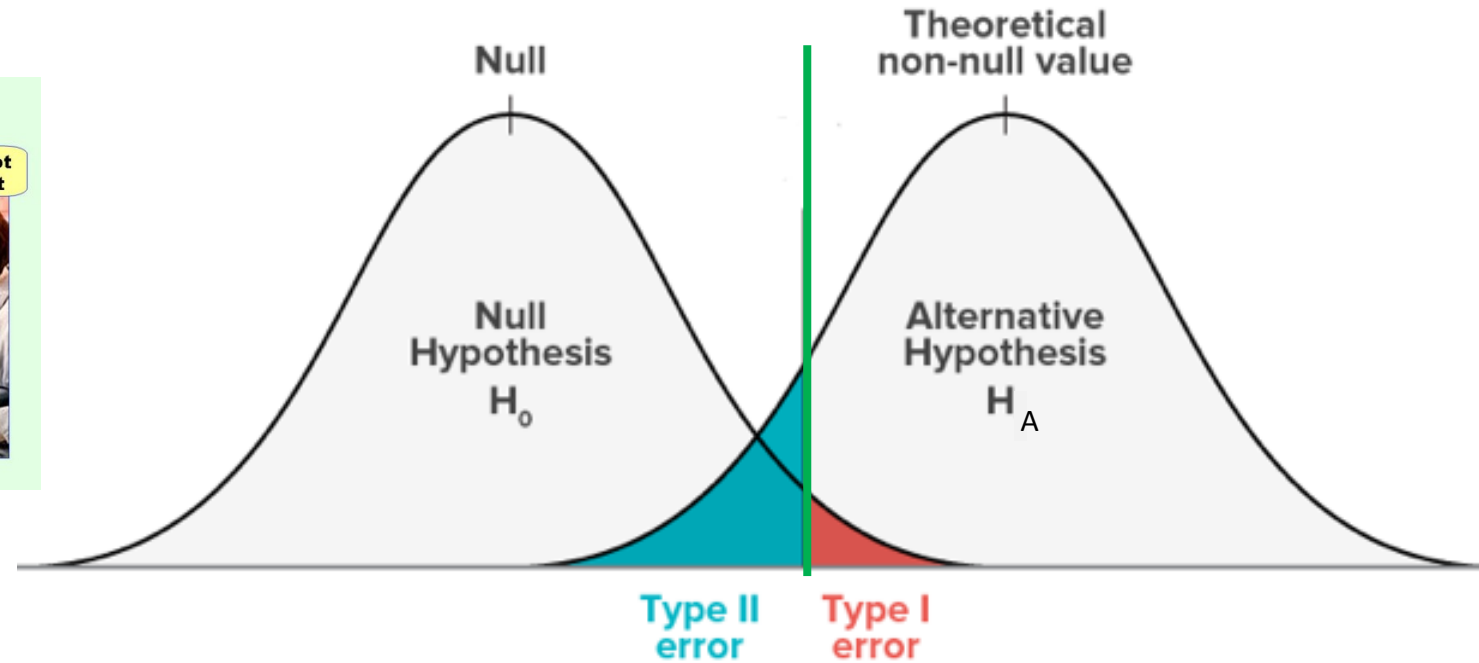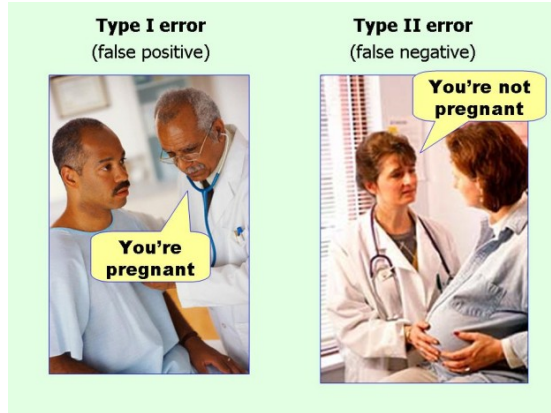
**Type I error**: incorrectly rejecting the null hypothesis when it is true

# Neyman-Pearson Frequentist logic



**Type 2 error**: incorrectly rejecting failing to reject $H_0$ when it is false

# Type I and Type II Errors



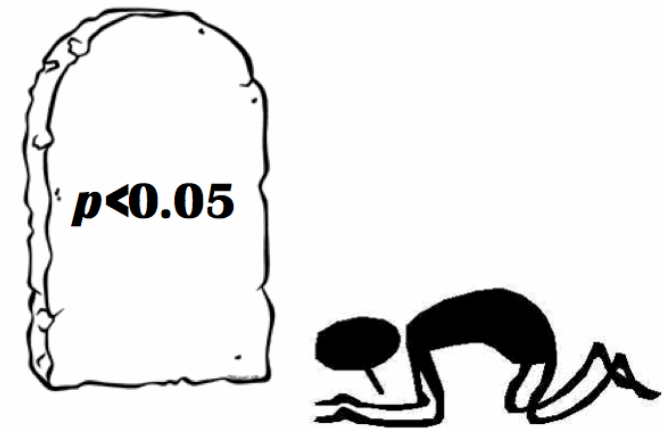| | Reject $H_0$ | Do not reject $H_0$ |
|---|---|---|
| $H_0$ is true | Type I error ($\alpha$) (false positive) | No error |

# Problems with the NP hypothesis tests

Problem 1: we are interested in the results of a specific experiment, not whether we are right most of the time

- E.g., 95% of these statements are true:
  - Calcium is good for your heart, Paul is psychic, Buzz and Doris can communicate, …
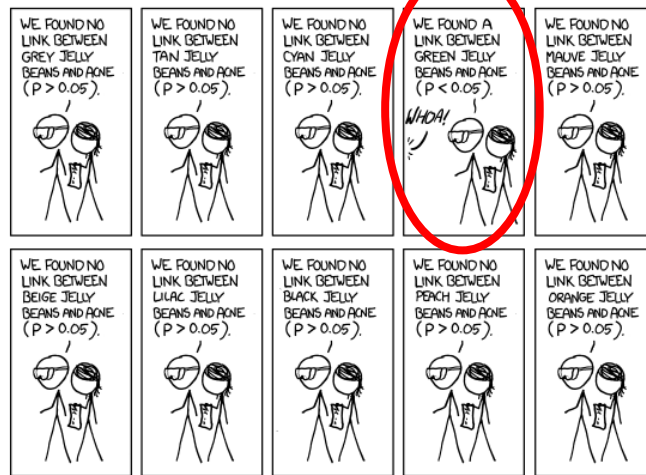
Problem 2:  Arbitrary thresholds for alpha levels
- P-value = 0.051, we don't reject $H_0$?

*p*<0.05
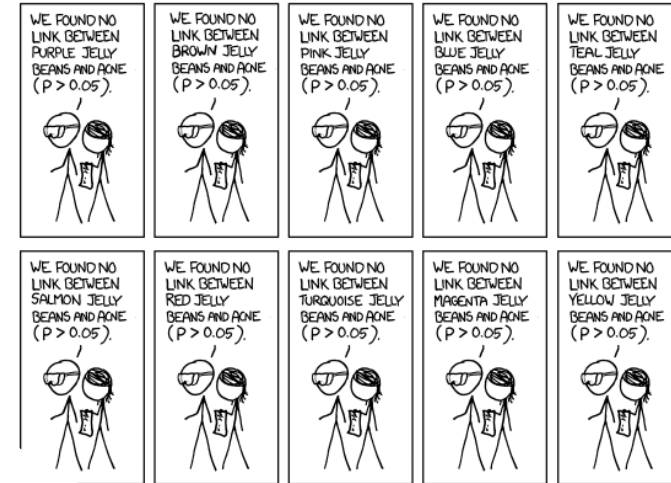
Problem 3: running many tests can give rise to a high number of type 1 errors

# Multiple hypothesis tests

# Replication crisis

**Essay**

## Why Most Published Research Findings Are False

John P. A. Ioannidis

### The file drawer effect

...and this is where we put the non-significant results.

somee cards
user card

American Statistical Association's 'Statement on p-values

# Some thoughts…

Better to have hypothesis tests than none at all. Just need to think carefully and use your judgment.

Report effect size in most cases – i.e., confidence intervals

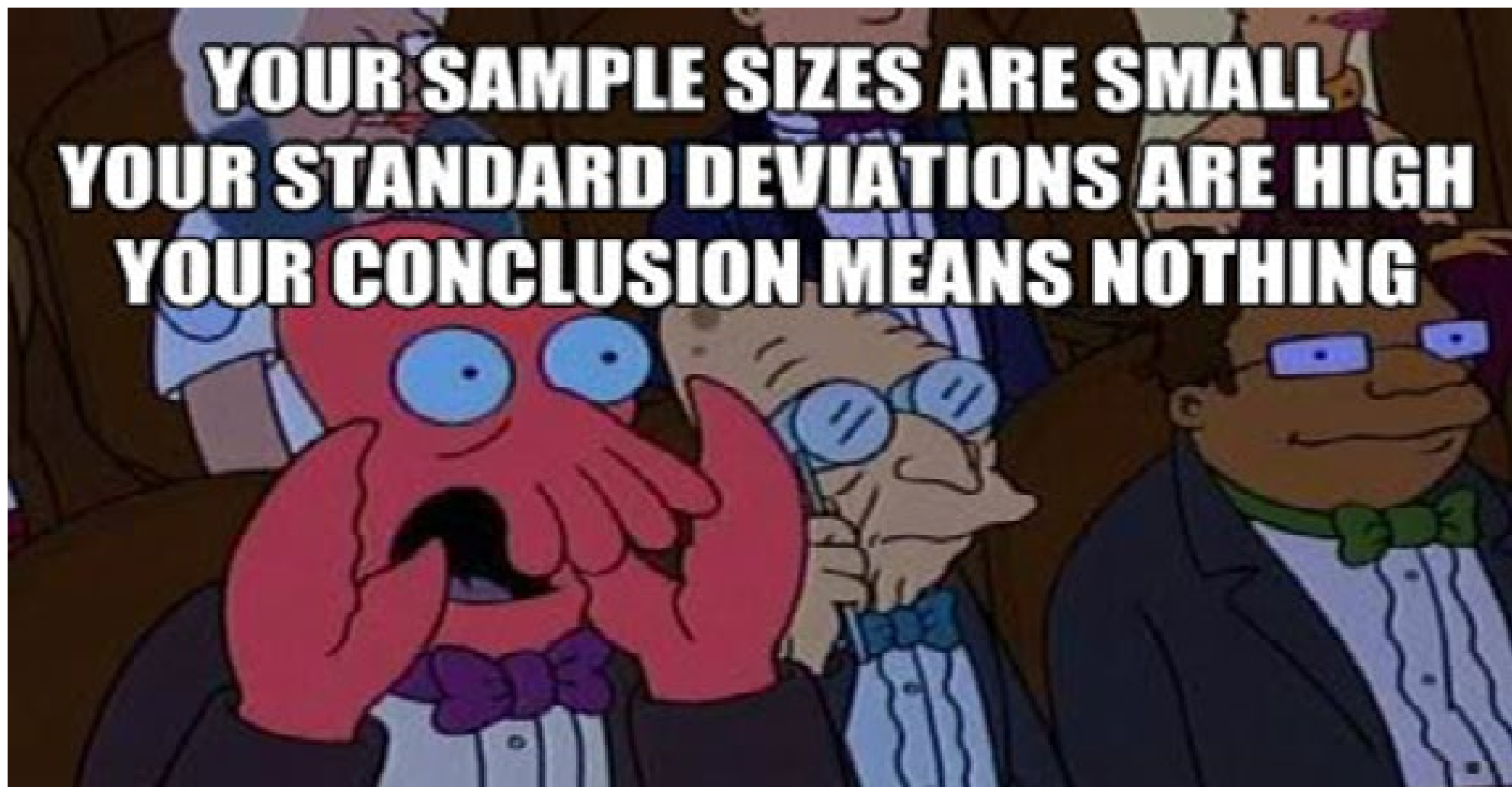Report the p-values rather than accept/reject $H_0$
- i.e., report   p = 0.023  not   p < 0.05
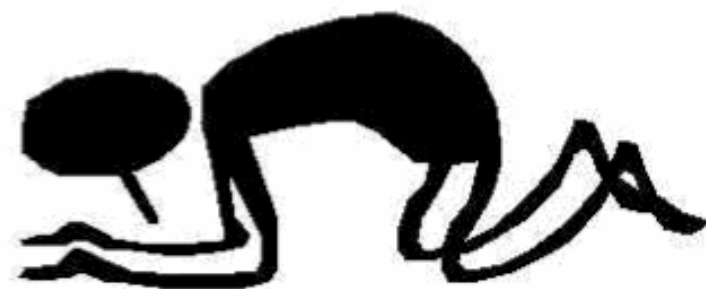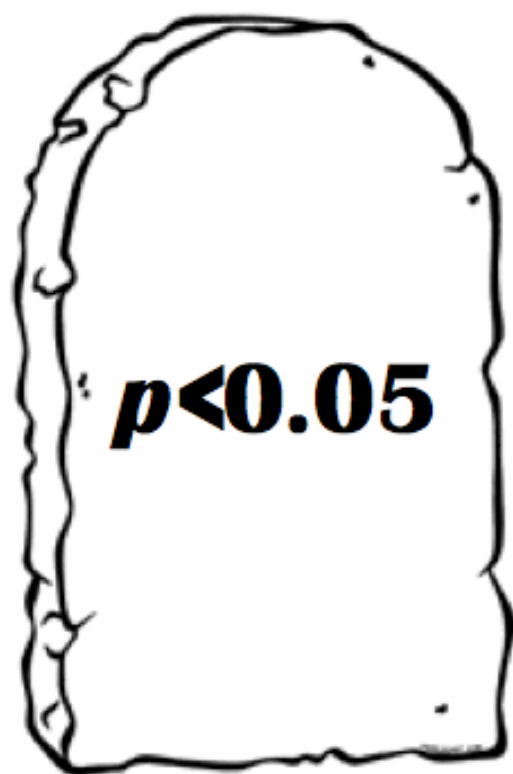
Replicate findings (perhaps in different contexts) to make sure you get the same results

Be a good/honest scientists and try to get at the Truth!

# Inference using parametric probability distributions



Total Area = 1

# Inference using parametric probability distributions

We can use mathematical functions called **probability distributions** to do inference

- e.g. instead of running computer simulations to create null distributions we can just mathematical probability distributions

A **density curve** is a mathematical function $f(x)$ that has two important properties:

1. The total area under the curve $f(x)$ is equal to 1

2. The curve is always $\geq 0$

# Density Curves

The <u>area under the curve</u> in an interval [a, b] models the probability that a random number X will be in the interval

Pr(a < X < b)  is the area under the curve from a to b

# The Normal Density Curve

Normal distributions are a family of bell-shaped curves with two parameters

- The mean: μ
- The standard deviation: σ



Notation: X ~ N(μ, σ)



**Changing μ**



**Changing σ**

# Graphing Normal Curves

Plotting IQ scores

```
x <- 40:150

                              μ        σ

y <- dnorm(x, 100, 15)

plot(x, y, type = "l")
```

# Finding normal probabilities of a normal curve

To get the probability (area) from a normal distribution we can use the pnorm function

```
pnorm(x, mu, sigma)
```

Pr(X < 9;  11, 3)  $\mu$   $\sigma$

```
pnorm(9, 11, 3)
```



P(X<9)

# Calculate the probability a random person you meet has an IQ less than 88

`pnorm(88, 100, 15)`



Normal area Pr(X ≤ x) app                    Normal area Pr(a < X < b) app

# Probability practice questions

1. What is probability a randomly chosen person will have an IQ greater than 96?

    pnorm(96, 100, 15, lower.tail = FALSE)

- Answer:  0.605

2. What is the probability a randomly chosen person will have an IQ between 88 and 96?

    pnorm(96, 100, 15) - pnorm(88, 100, 15)

- Answer:  0.183

# Calculating quantiles

To find quantiles of the normal distribution we can use the quantile function:

```
qnorm(quantile, mu, sigma)
```

What are the IQ scores (interval) that demark the middle 50% of the IQ range?

What about the middle 95%?

Normal quantile app

# Middle 50% of IQ scores



```
qnorm(c(.25, .75), 100, 15)
```

Middle 50%:  89.9  to  110.1                    Middle 95%:  70.6 to 129.3

# Summary of R functions

Plot the actually density curve
- dnorm(x_vec, mu, sigma)

Get the probability that we would get a random value less than x
- pnorm(x_vec, mu, sigma)

Get the quantile value for a given proportion of the distribution
- qnorm(area, mu, sigma)

Note: pnorm and qnorm are inverses of each other
- y = pnorm(x, mu, sigma)
- qnorm(y, mu, sigma)          # the output value here is x

# The Standard Normal distribution and the Central Limit Theorem

# Standard Normal N(0, 1)

Since all normal distributions have the same shape, it is convenient to convert them to a standard scale with:

$$\mu = 0, \quad \sigma = 1$$

This is called the **standard normal** distribution:

$$Z \sim N(0, 1)$$

# Converting to the standard normal distribution

We can use a z-score transformation to any normally distributed random variable $X \sim N(\mu, \sigma)$ to the standard normal distribution $Z \sim N(0, 1)$:

$$Z = (X - \mu)/\sigma$$

To convert from $Z \sim N(0, 1)$ to any $X \sim N(\mu, \sigma)$, we reverse the standardization with:

$$X = \mu + Z \cdot \sigma$$

# Converting to the standard normal distribution

1. What is the Z-score of someone who has an IQ score of 112?

   $Z = (X - \mu)/\sigma$

2. What if someone has an Z-score of 2.2, what is their IQ score?

   $X = \mu + Z \cdot \sigma$

**Answer 1**: $Z = (112 - 100)/15 = .8$

**Answer 2**: $IQ = 100 + 2.2 * 15 = 133$

# Central limit theorem

For random samples with a sufficiently large sample size (n), the distribution of sample statistics for a <span style="color:red">mean ($\bar{x}$)</span> or a <span style="color:red">proportion ($\hat{p}$)</span> is:

- normally distributed
- centered at the value of the population parameter

Stated again: the sampling distribution for means or proportions will be a normal distribution

- so we don't need to do resampling to get a bootstrap or null distribution!

# Central limit theorem

## proportion (p̂)



## mean (x̄)



Proportion sampling distribution app

Sampling/Bootstrap distribution app

# Summary of standard normal and CLT

For large n, the sampling distributions of $\bar{x}$ and $\hat{p}$ are normal

We can convert any normal distribution $N(\mu, \sigma)$, into a standard normal distribution $N(0, 1)$

We are now (almost) ready to run hypothesis tests and compute confidence intervals for $\bar{x}$ and $\hat{p}$ using normal distributions

# Hypothesis tests using a normal distribution

# Hypothesis tests based on a Normal Distribution

When the null distribution is normal, it is often convenient to use a standard normal test statistic using:

$$z = \frac{Sample\ Statistic\ -\ Null\ Parameter}{SE}$$

The p-value for the test is the probability a standard normal value is beyond this standardized test statistic



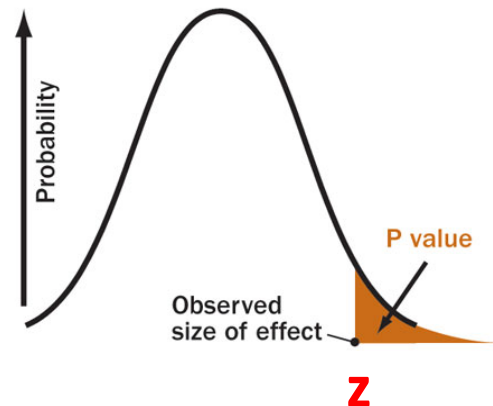$Pr(\ Z \geq z_{obs}\ ;\ \mu = 0,\ \sigma = 1)$

pnorm(z, 0, 1, lower.tail = FALSE)

# Hypothesis tests based on a Normal Distribution

To repeat what was on the last slide:  we can transform our obs_stat to a z-statistic that comes from a standard normal distribution N(0, 1) using:

$$z = \frac{stat_{obs} - param_0}{SE}$$

The p-value is then the probability of obtaining a value from a standard normal distribution beyond this z statistic

> pnorm(z, 0, 1)              if     $H_A$:  $\mu$ < $param_0$

> 1 - pnorm(z, 0, 1)          if     $H_A$:  $\mu$ > $param_0$

> 2 * (1 - pnorm(abs(z), 0, 1))    if     $H_A$:  $\mu$ ≠ $param_0$

# Do greater than 40% of Americans go without using cash in a typical week?

A survey of 1,000 Americans reported that 43% said they went an entire week without using cash, with a SE = 0.016

Assuming the distribution of the statistic is normal, calculate whether the proportion of all Americans going a week without using cash is greater than 40%

**1. Start by stating H$_0$ and H$_A$**

$H_0: \pi = .4$

$H_A: \pi > .4$

# Do greater than 40% of Americans go without using cash in a typical week?

A survey of 1,000 Americans reported that 43% said they went an entire week without using cash, with a SE = 0.016

Assuming the distribution of the statistic is normal, calculate whether the proportion of all Americans going a week without using cash is greater than 40%

## 2. Can you compute the z statistic?

$$z = \frac{Sample\ Statistic\ -\ Null\ Parameter}{SE}$$

# Do greater than 40% of Americans go without using cash in a typical week?

A survey of 1,000 Americans reported that 43% said they went an entire week without using cash, with a SE = 0.016

Assuming the distribution of the statistic is normal, calculate whether the proportion of all Americans going a week without using cash is greater than 40%

**2. Can you compute the z statistic?**

$$z = (.43 - .4)/.016 = 1.875$$

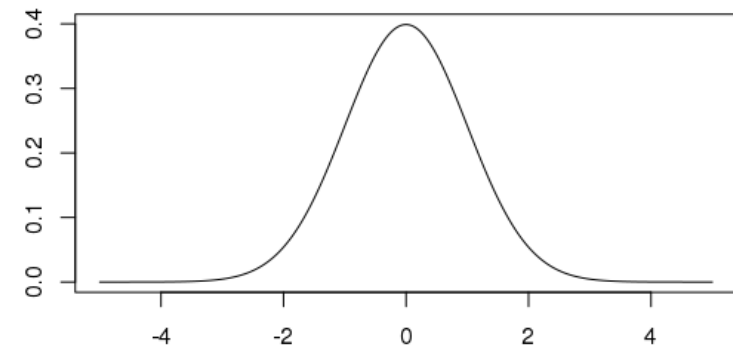# Do greater than 40% of Americans go without using cash in a typical week?

**Steps: 3-4**. What is the probability one would get a z-statistic as larger or larger than 1.875 from a standard normal distribution?

> pnorm(1.875, 0, 1, lower.tail = FALSE)

> 1 – pnorm(1.875, 0, 1)

p-value = .0304

Normal area app  Pr(X ≤ x)

Standard normal null distribution



**Step 5?**
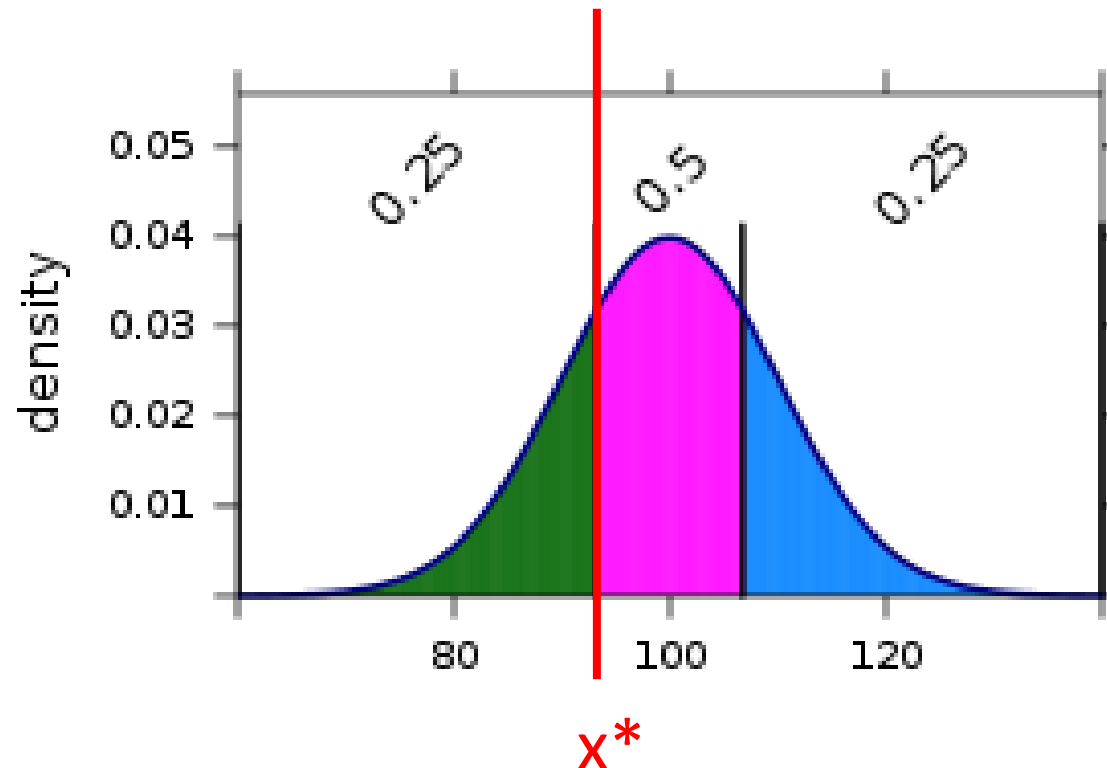
# Confidence intervals using a normal distribution

# Finding quantile values

We can find the quantile value from a normal distribution with:

qnorm(q, mu, sigma)

The 'q' in qnorm stands for quantile

What is the max and min that q can be?
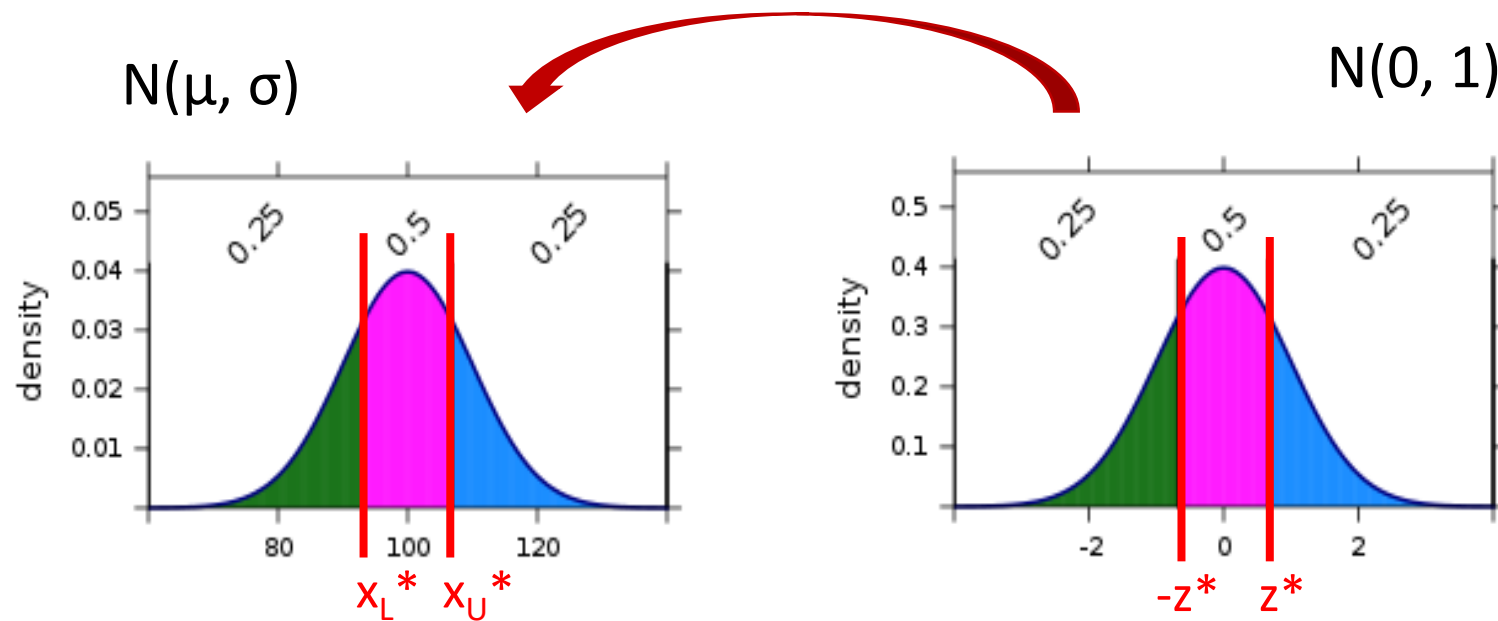
# Standard Normal N(0, 1)

It is often convenient to find quantiles on the standard normal distribution $Z \sim N(0, 1)$ and then to transform them to an arbitrary normal distribution $X \sim N(\mu, \sigma)$, using :

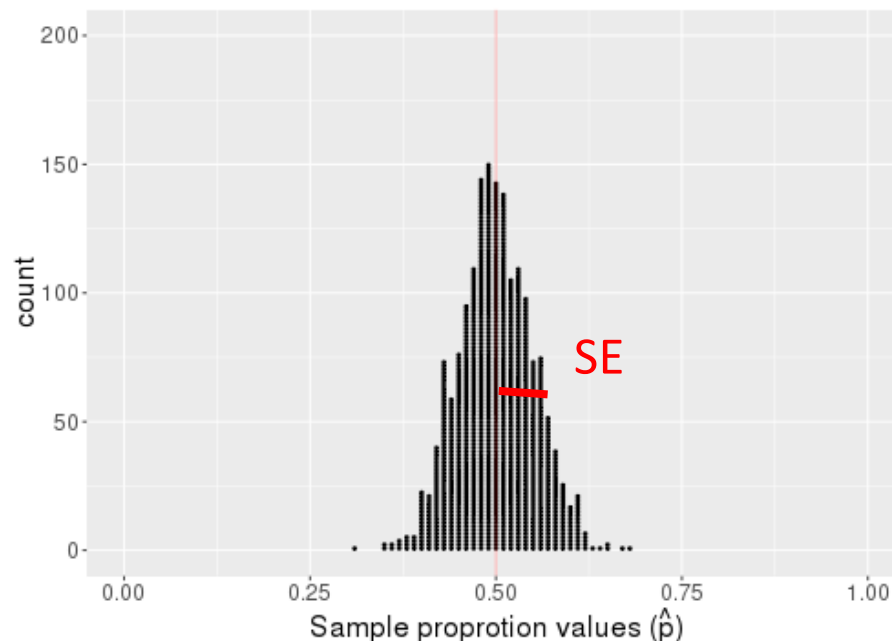$$X = \mu + Z \cdot \sigma$$

# Central limit theorem

Questions:
1. What is the standard deviation of these sampling distributions called?
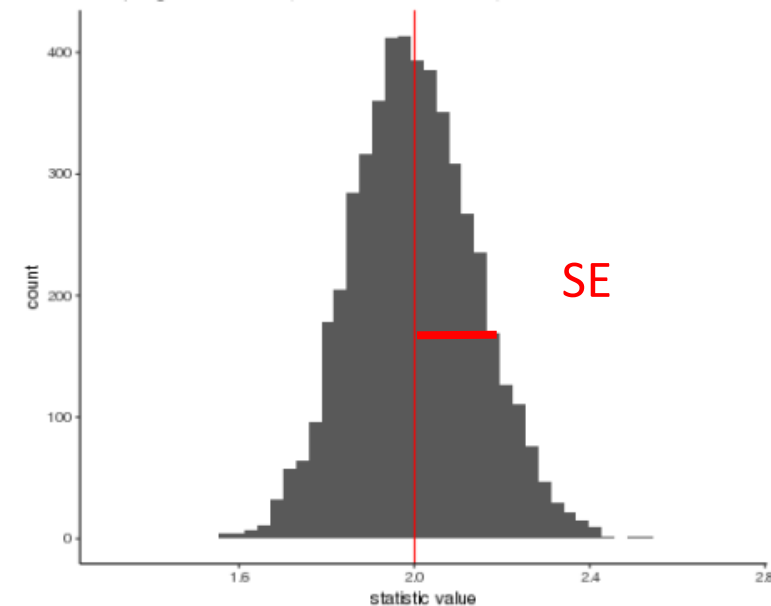2. Suppose we have a $\hat{p}$ or $\bar{x}$ and know the SE, how can we create a 95% CI?

For a proportion $\pi$:   $CI_{95} = \hat{p} \pm 2 \cdot SE$        For a mean $\mu$:   $CI_{95} = \bar{x} \pm 2 \cdot SE$



proportion ($\hat{p}$)



mean ($\bar{x}$)

# Confidence intervals based on a Normal Distribution

If the distribution for a statistic is normal with a standard error SE, we can find a confidence interval for the parameter using:
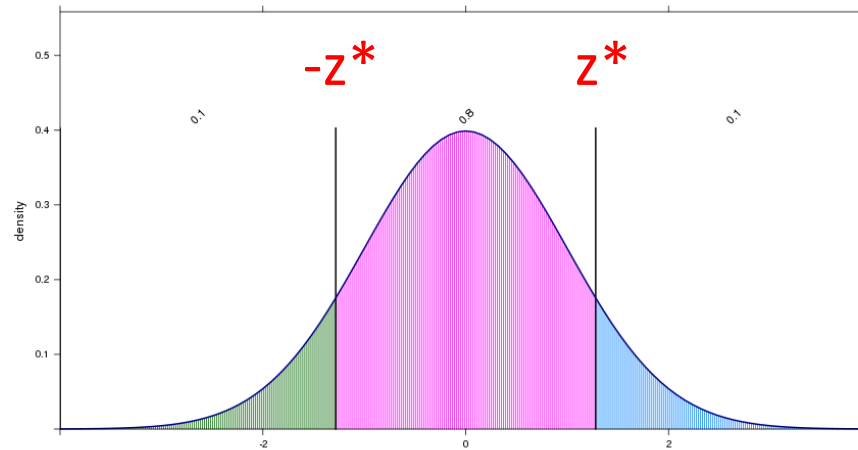
sample statistic  ±  z* ×  SE

where z* is chosen so that the area between –z* and + z* in the standard normal distribution is the desired confidence level

- i.e., z* is chosen such that say 95% of the distribution is between ± z*

# Confidence intervals based on a Normal Distribution

Suppose we are interested in 80% confidence intervals for $\mu$

We calculate the $\pm z_{80}$ that has 80% of the data on N(0, 1)



Let's assume we know the SE but don't know $\mu$. If we have an observed statistic from:

$x_{obs} \sim$ N($\mu$, SE)

We can create an interval that will capture $\mu$ 80% of the time using:

$x_{obs} \pm z_{80} \cdot$ SE

# Normal percentiles for common confidence levels

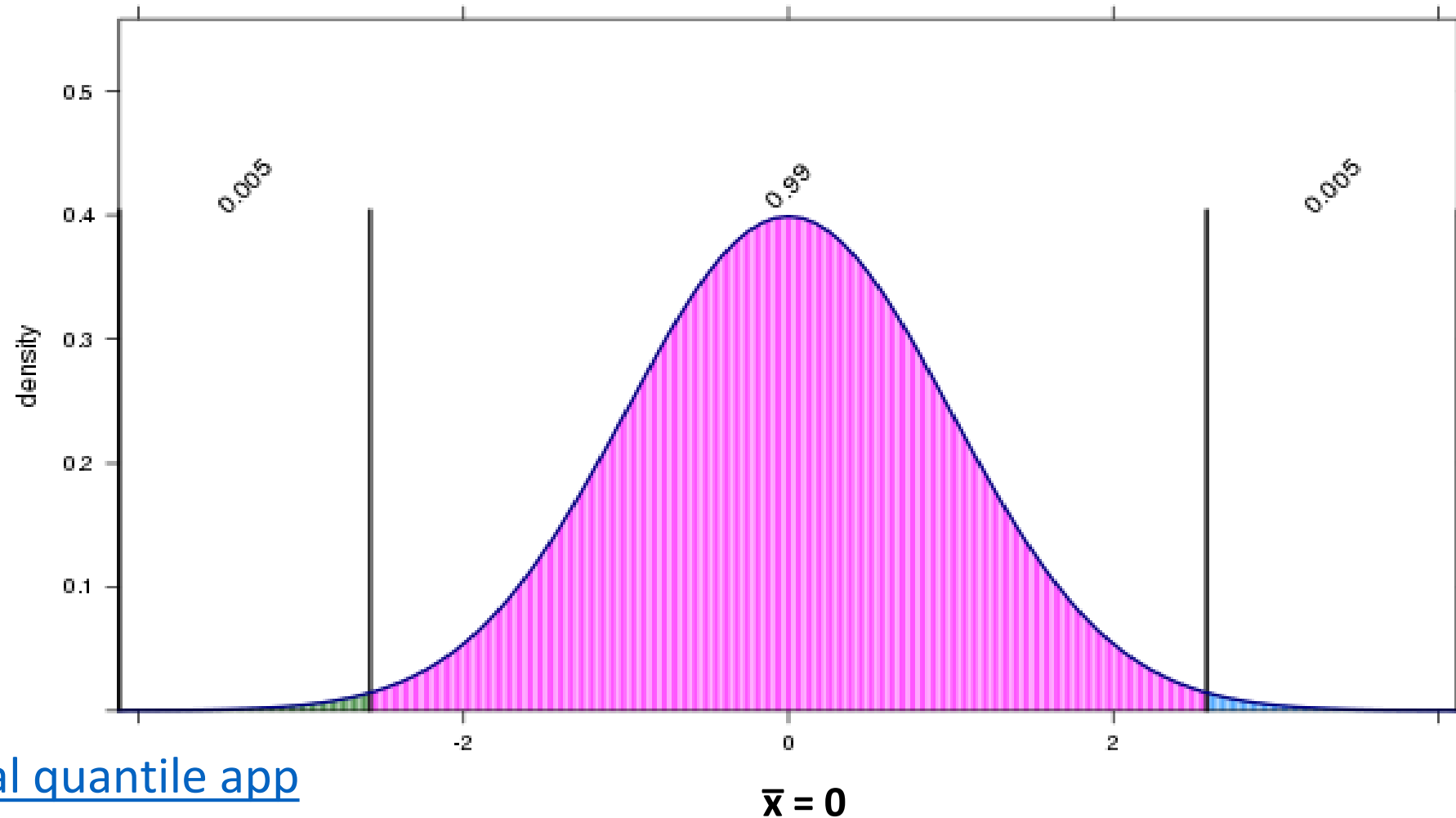| Confidence level | 80% | 90% | 95% | 98% | 99% |
|---|---|---|---|---|---|
| Z* | 1.282 | 1.645 | 1.960 | 2.326 | 2.576 |

z_stars <- qnorm(c(.90, .95, .975, .99, .995), 0, 1)

Normal quantile app

# .99 quantile values



$\bar{x} = 0$,     SE = 1

Quantile values:    [-2.576   2.576]

Normal quantile app

# What is the most preferred seat?

A survey of 1,000 air travelers found that 60% prefer a window seat, with a bootstrap standard error of SE = 0.015

Use the normal distribution to compute a 90%, 95% and 99% CIs for the proportion of people who prefer a window seat

sample statistic $\pm$ z* $\times$ SE

| Confidence level | 80% | 90% | 95% | 98% | 99% |
|---|---|---|---|---|---|
| Z* | 1.282 | 1.645 | 1.960 | 2.326 | 2.576 |

# What is the most preferred seat?

A survey of 1,000 air travelers found that 60% prefer a window seat, with a bootstrap standard error of SE = 0.015.

90% CI =  .6 ± 1.645 × .015    =    [.575   .625]

95% CI =  .6 ± 1.96 × .015     =    [.571   .629]

99% CI =  .6 ± 2.576 × .015    =    [.569   .638]

Sample statistics  ± z* ×  SE

| Confidence level | 80% | 90% | 95% | 98% | 99% |
|---|---|---|---|---|---|
| Z* | 1.282 | 1.645 | 1.960 | 2.326 | 2.576 |