

# Test Design, Effect Size and Test Power

**Matthias RAUTERBERG**

Eindhoven University of Technology

**2017**

# Causal Relationships

- ONLY if Independent Variables (IV) are specified and implemented UPFRONT.
- IV's are mainly Nominal Scaled.
- Any other measurements are either a Dependent Variable (DV) or a Co-Variate (CV).
- DV and CV can be of any scale type.

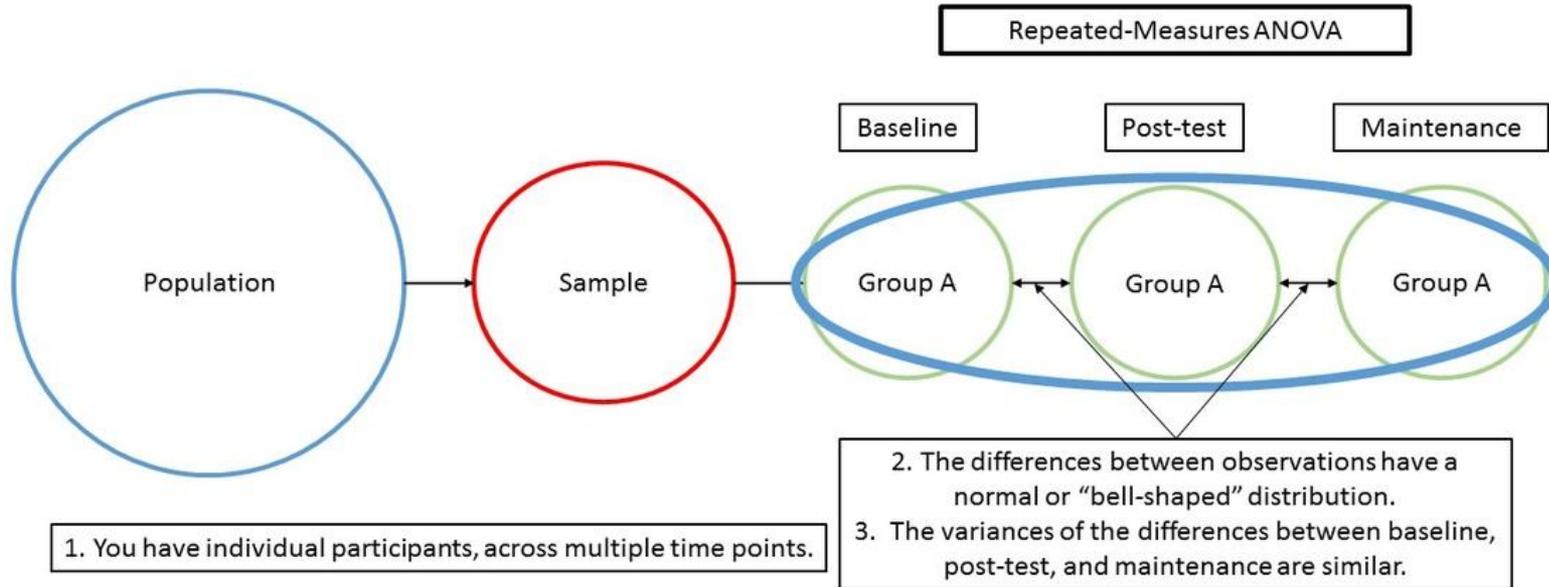
# Test Design

- **Between-subject design:**  
measurements are taken from different samples.
- **Within-subject design:**  
measurements are taken from same sample.

# How to deal with Learning effects?

- If the main concept under investigation (e.g., a human being) is adaptive and can learn, any measurement might change it's internal state.
- This has to be addressed for within-subject test design with repeated measurements.

# Repeated Measurements



# Counter Balancing

Assuming we have 4 test conditions: A, B, C, D  
Counterbalancing means to test ALL permutations:  
[ABCD] [BACD] [BCAD] [BCDA] [...]

Total number of possible sequences is  $4! = 24$

In case of a 2x2 factorial test design, it's called LATIN SQUARE.

We can easily calculate a factorial from the previous one:

n	n!		
1	<b>1</b>	1	1
2	$2 \times \mathbf{1}$	$= 2 \times \mathbf{1!}$	$= 2$
3	$3 \times \mathbf{2} \times \mathbf{1}$	$= 3 \times \mathbf{2!}$	$= 6$
4	$4 \times \mathbf{3} \times \mathbf{2} \times \mathbf{1}$	$= 4 \times \mathbf{3!}$	$= 24$
5	$5 \times \mathbf{4} \times \mathbf{3} \times \mathbf{2} \times \mathbf{1}$	$= 5 \times \mathbf{4!}$	$= 120$
6	etc	etc	

# Treatment vs Control Group

- How to find an appropriate control test condition?
  - take best solution available (state of art)
  - avoid unfair comparison



# Design of Experiments

- In the design of experiments, treatments are applied to experimental units in the **treatment group(s)**.  
In *comparative* experiments, members of the complementary group, the **control group**, receive either *no* treatment or a *standard* treatment.
- For the conclusions drawn from the results of an experiment to have validity, it is essential that the items or test subjects are assigned to *treatment and control groups* be representative of the same population.  
In some experiments, such as many in agriculture or psychology, this can be achieved by randomly assigning items from a common population to one of the treatment and control groups. In studies of twins involving just one treatment group and a control group, it is statistically efficient to do this random assignment separately for each pair of twins, so that one is in the treatment group and one in the control group.
- In some medical studies, where it may be unethical not to treat patients who present with symptoms, controls may be given a standard treatment, rather than no treatment at all. Another alternative is to select controls from a wider population, provided that this population is well-defined and that those presenting with symptoms at the clinic are representative of those in the wider population.

# Definition of Test Power

- Power (of a test) is the probability of detecting an effect, given that the effect is really there.
- Or likewise, the probability of rejecting the null hypothesis when it is in fact false.
- An example;
  - Power of 0.8 = if we performed a study 1000 times, we would see a statistically significant difference 80% of the time.

# Outcomes of a statistical test

**H<sub>0</sub> True**  
(no correlation)

**H<sub>1</sub> True**  
(correlation)

**Do not reject H<sub>0</sub>**  
(not stat. sig.)

**Reject H<sub>0</sub>**  
(stat. sig.)

<b>Correct decision</b>	<b>Type II (beta error)</b>
<b>Type I (alpha error)</b>	<b>Correct decision</b>

# What you need to know

- Core elements
  - Effect size
  - Power
  - Sample size
  - Significance
- These elements are all inter-related such that;
  - If you know three you can estimate the fourth
  - Manipulating one influences the others

# Introduction to sample size and power calculations

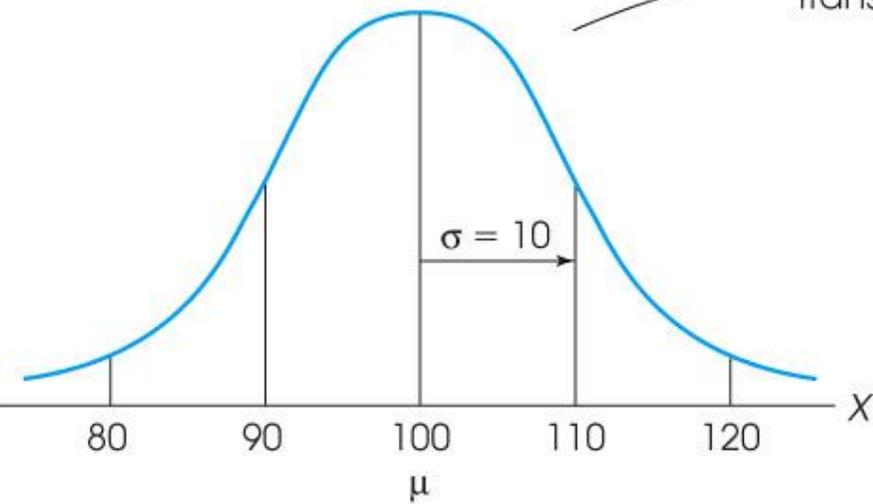
How much chance do we have to reject the null hypothesis when the alternative is in fact true?

(what's the probability of detecting a real effect?)

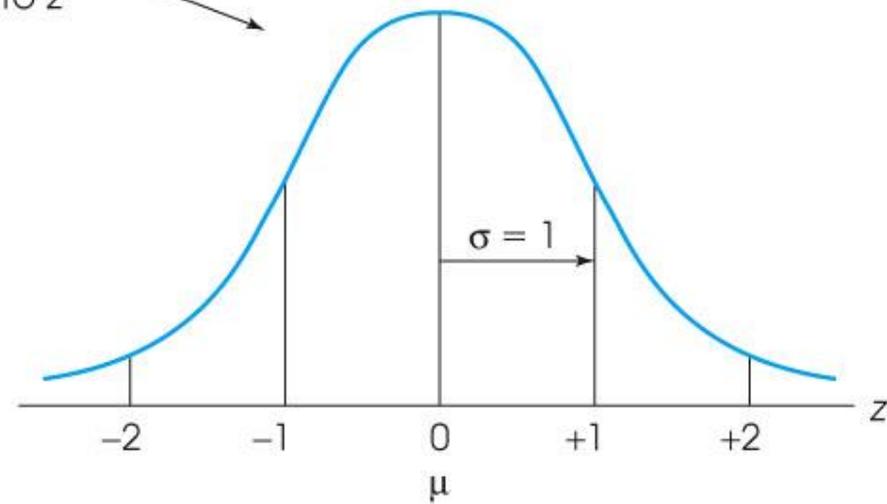
Can we quantify how much power we have for given sample and effect sizes?

# The z-transformation of variables

Population of scores  
(X values)



Population of z-scores  
(z values)



Transform X to z

## Transforming back and forth between X and z

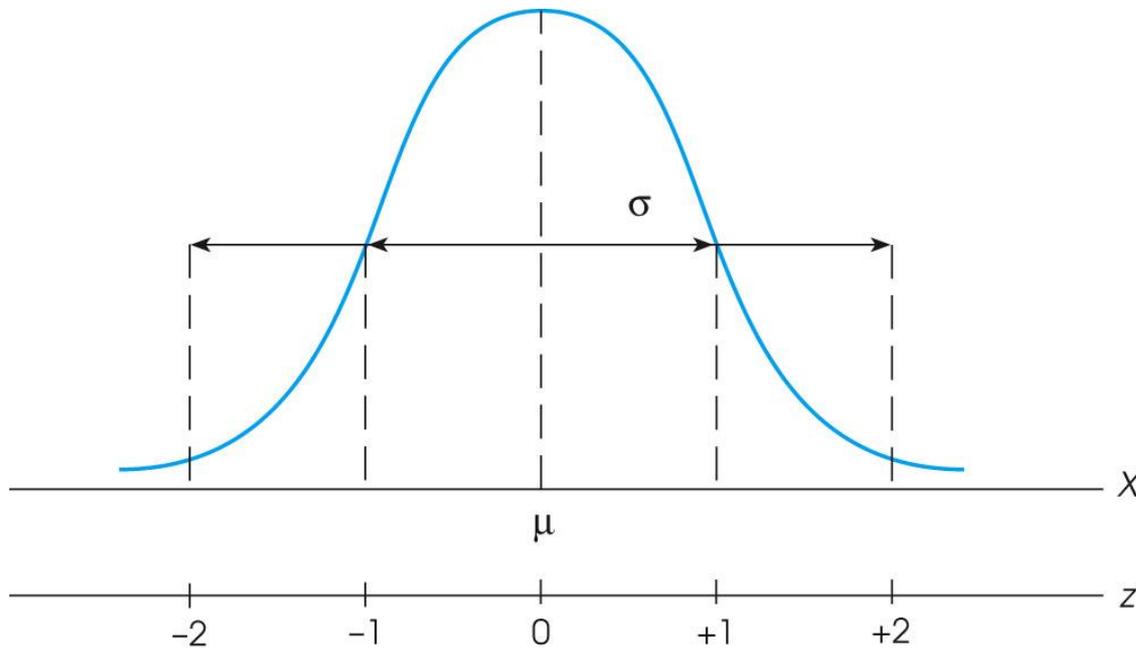
- The basic z-score definition is usually sufficient to complete most z-score transformations. However, the definition can be written in mathematical notation to create a formula for computing the z-score for any value of X.

- 

$$z = \frac{X - \mu}{\sigma} \quad \text{for population}$$

$$z = \frac{X - M}{s} \quad \text{for sample}$$

# The z-distribution

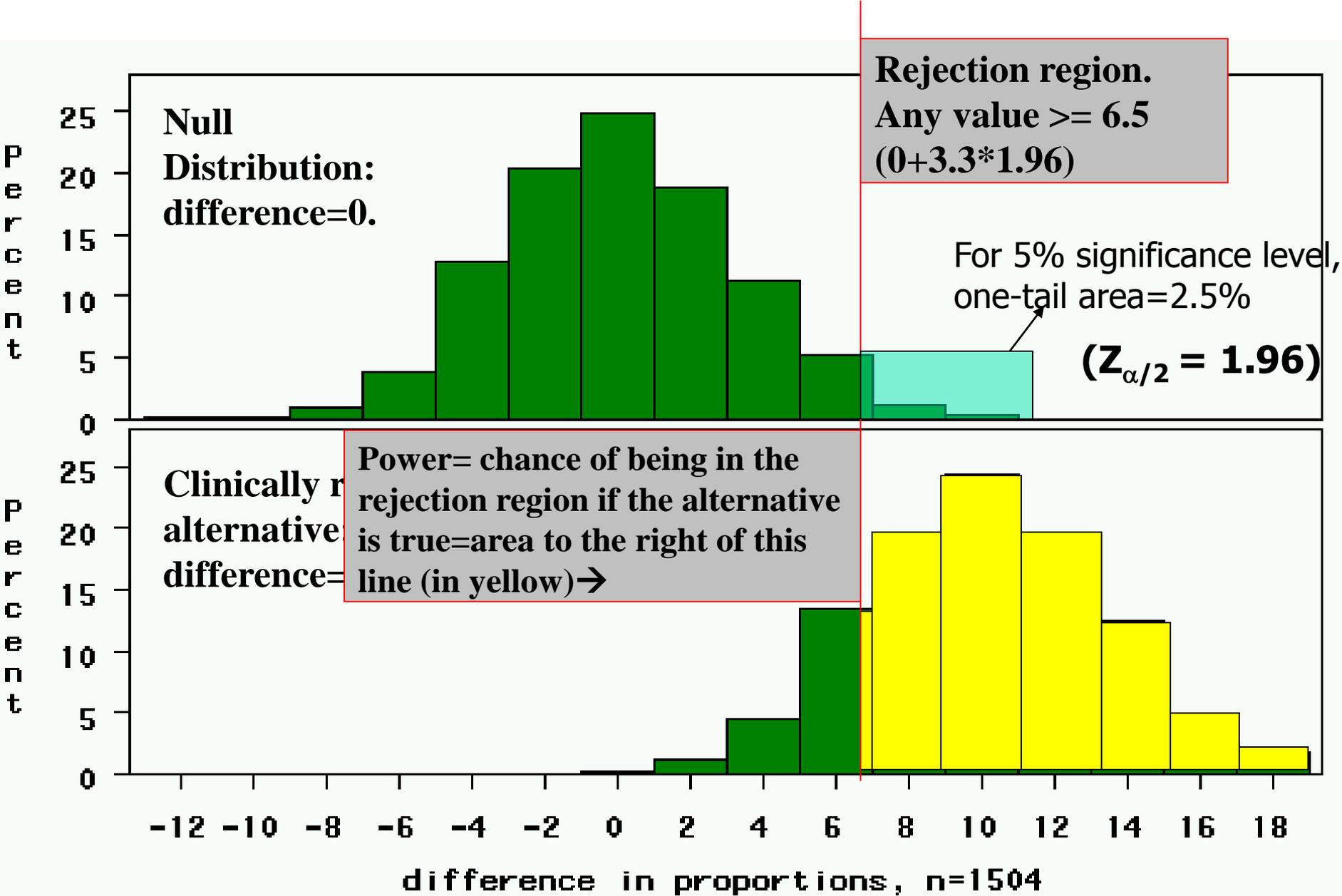


$Z=1.96$  for  $p=0.05$  (two sided)

# Deciding on levels of $\alpha$ and $\beta$

- **Power** (sensitivity) [ $1-\beta$ ]
  - Probability of finding a true effect when one does exist
  - Type 2 error [ $\beta$ ]: incorrectly accepting the null hypothesis (false negative)
  - Aim to minimise the risk of failing to detect a real effect
  - Typical values for power are 80%, 90% and 95%
- **Significance** (p-value) [ $\alpha$ ]
  - Probability that an effect occurred by chance alone
  - Type 1 error [ $\alpha$ ]: incorrectly rejecting the null hypothesis (false positive)
  - Aim to minimise the risk of detecting a non-real/spurious effect
  - Typical values are 0.05, 0.01
- **Reducing the risk of type 1 error  $\leftrightarrow$  increased risk of type 2 error (i.e. reduced power)**

study 1: 263 cases, 1241 controls



# study 1: 263 cases, 1241 controls

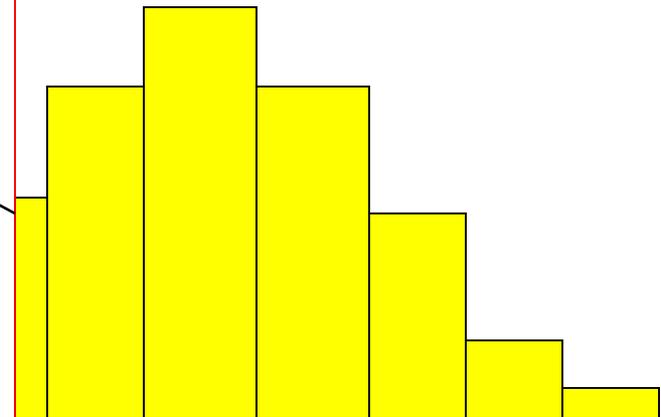
**Rejection region.  
Any value  $\geq 6.5$   
( $0+3.3*1.96$ )**

**Power= chance of being in the  
rejection region if the alternative  
is true=area to the right of this  
line (in yellow)**

Power here:

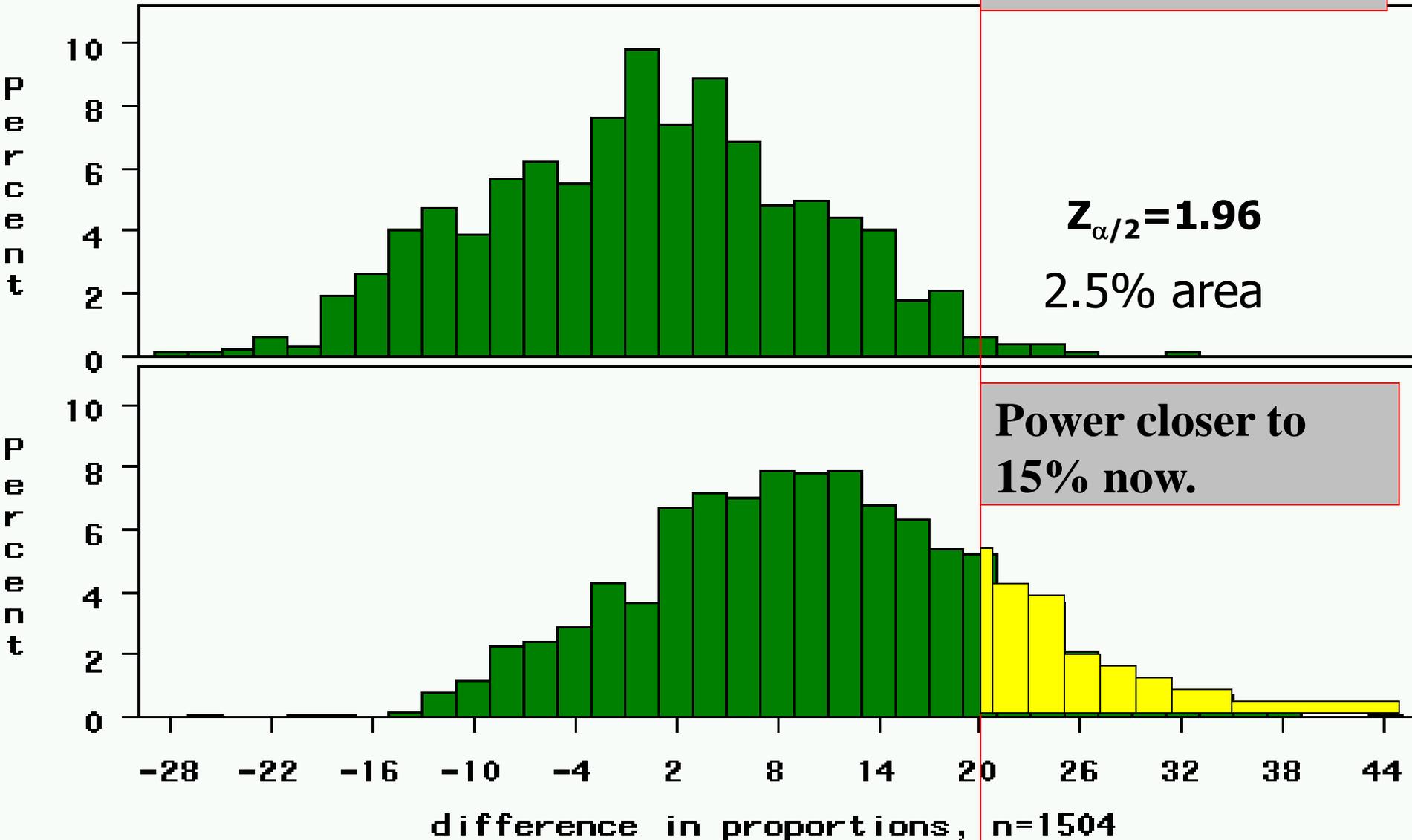
$$P(Z > \frac{6.5 - 10}{3.3}) =$$

$$P(Z > -1.06) = 85\%$$

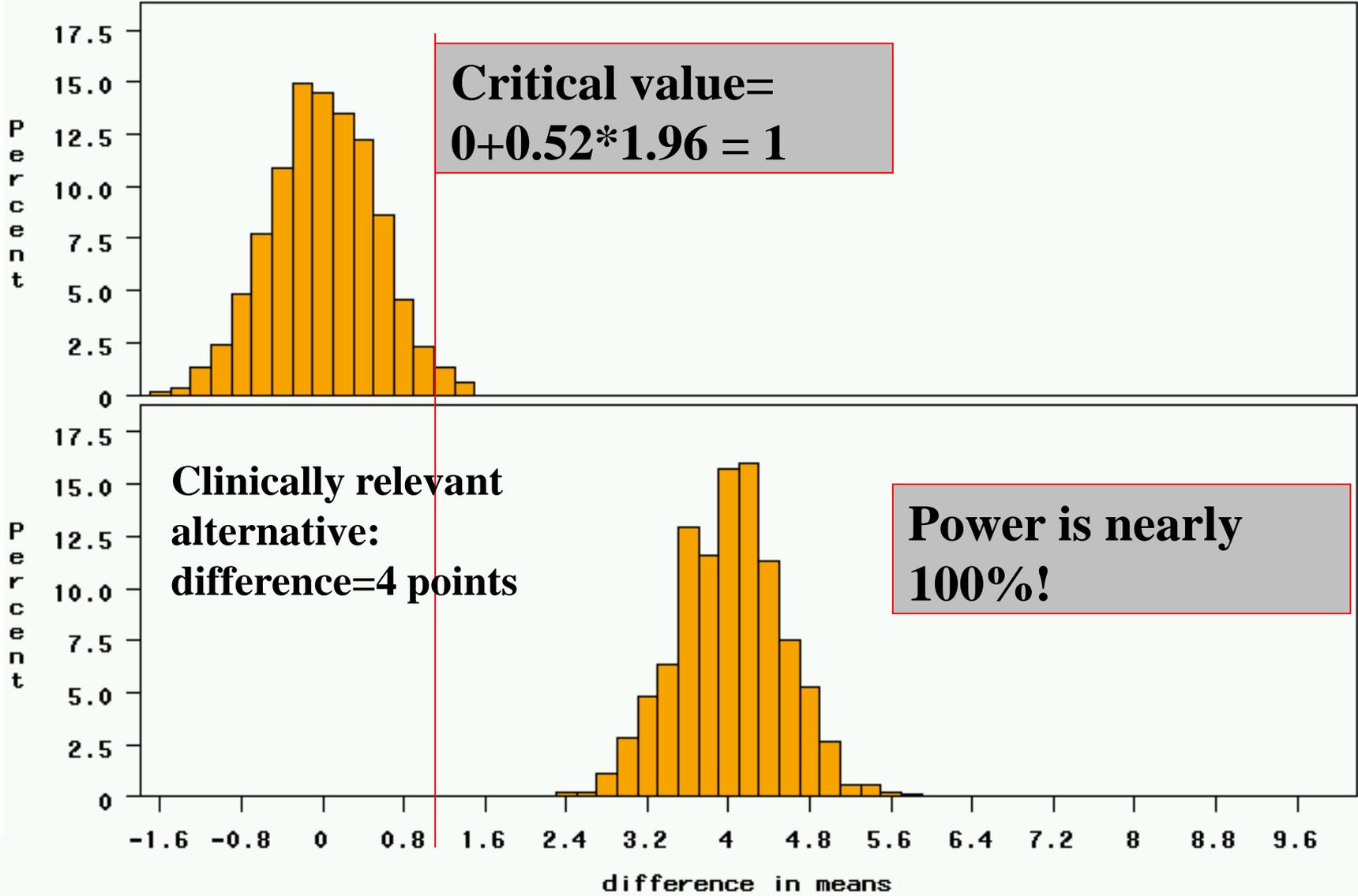


study 1: 50 cases, 50 controls

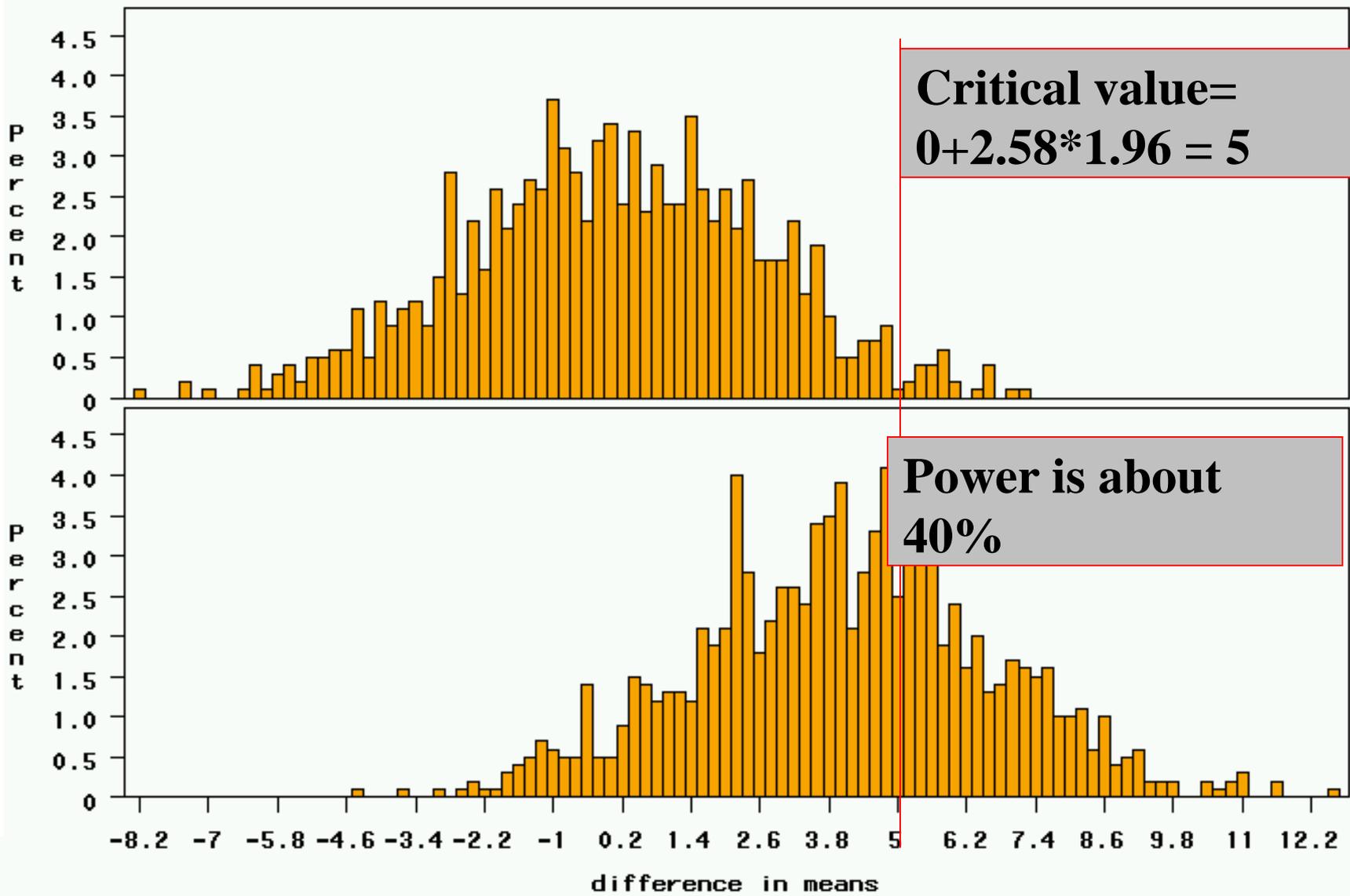
**Critical value=  
 $0+10*1.96=20$**



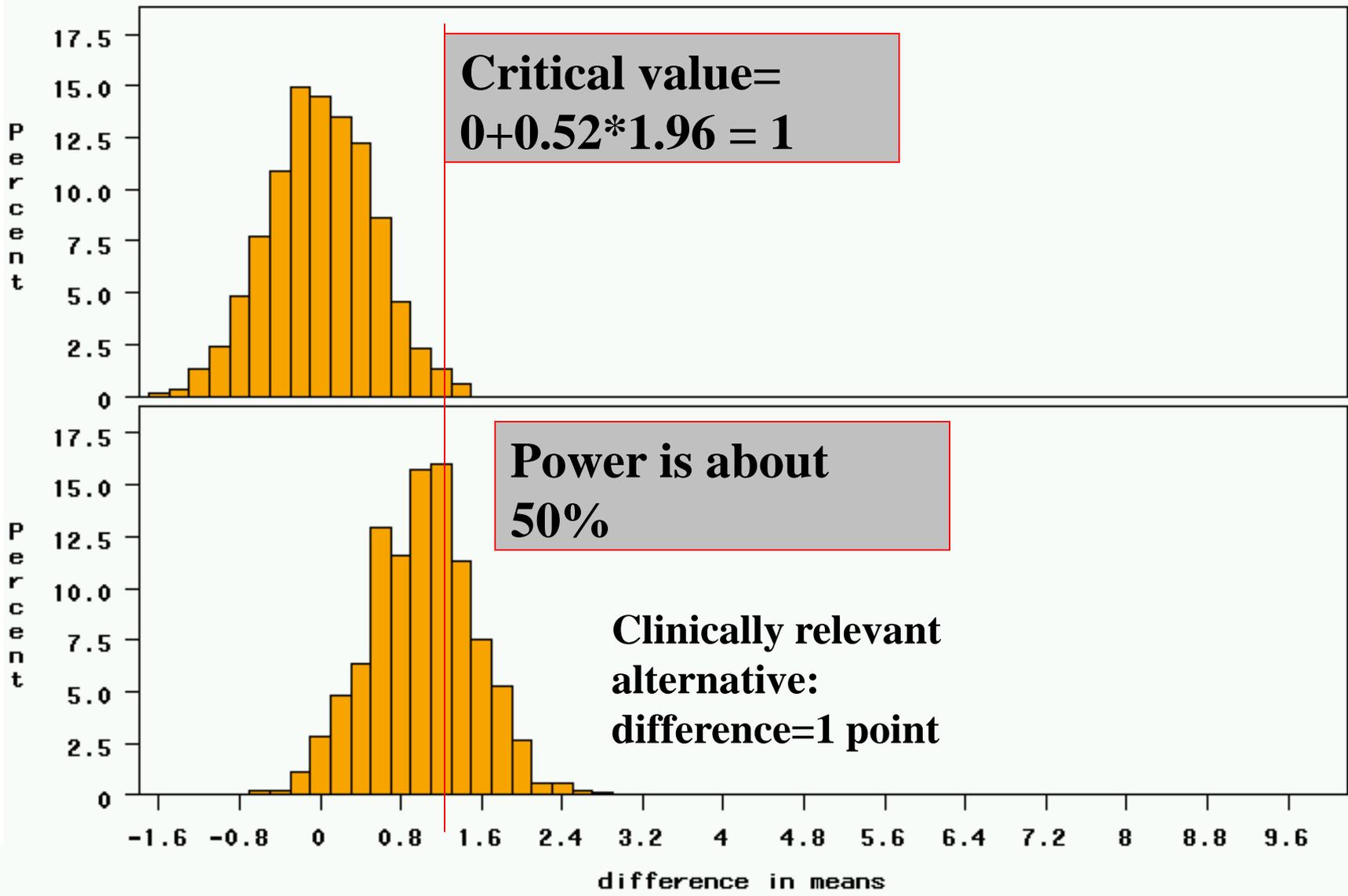
Study 2: 18 treated, 72 controls, STD DEV = 2



# Study 2: 18 treated, 72 controls, STD DEV=10



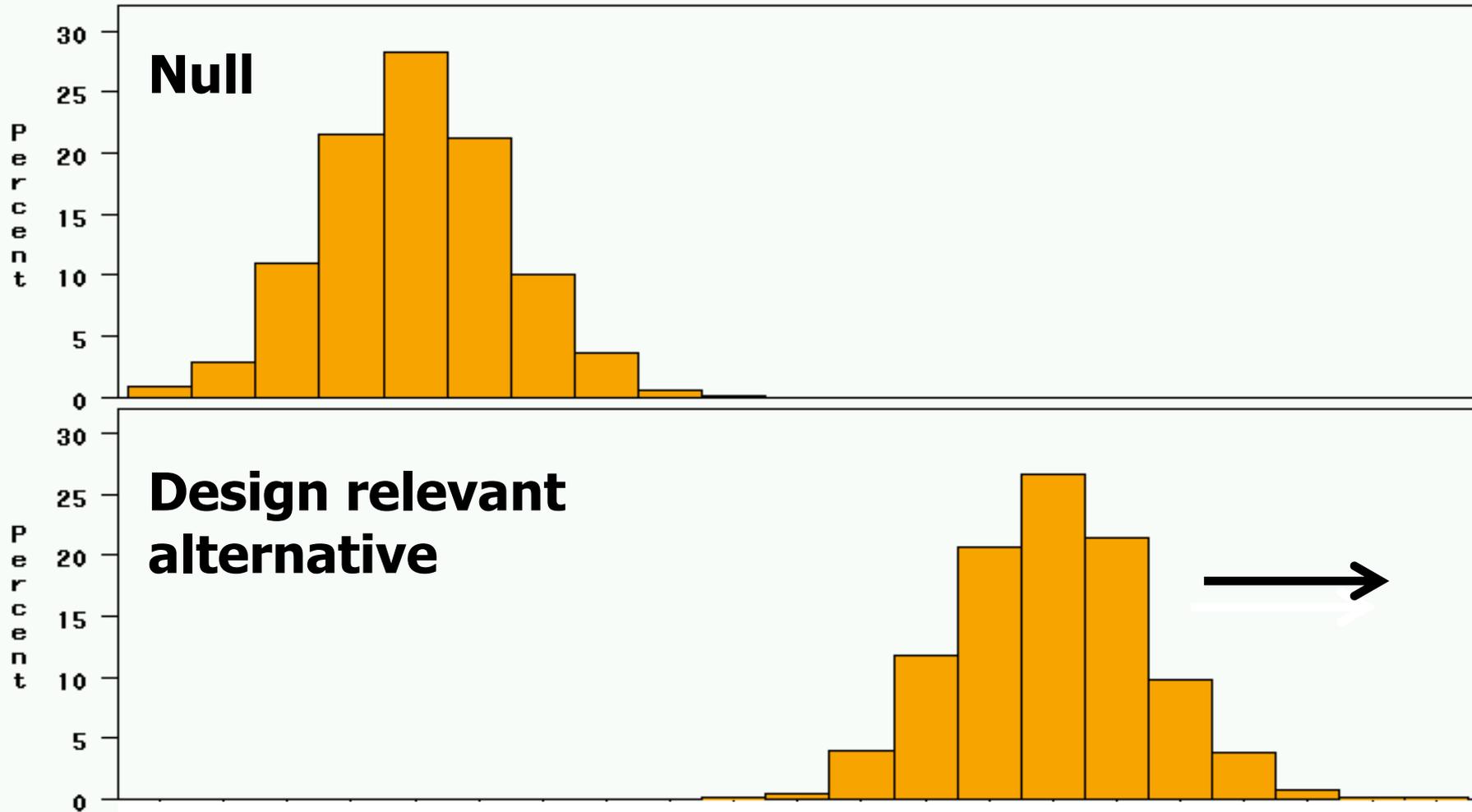
# Study 2: 18 treated, 72 controls, effect size=1.0



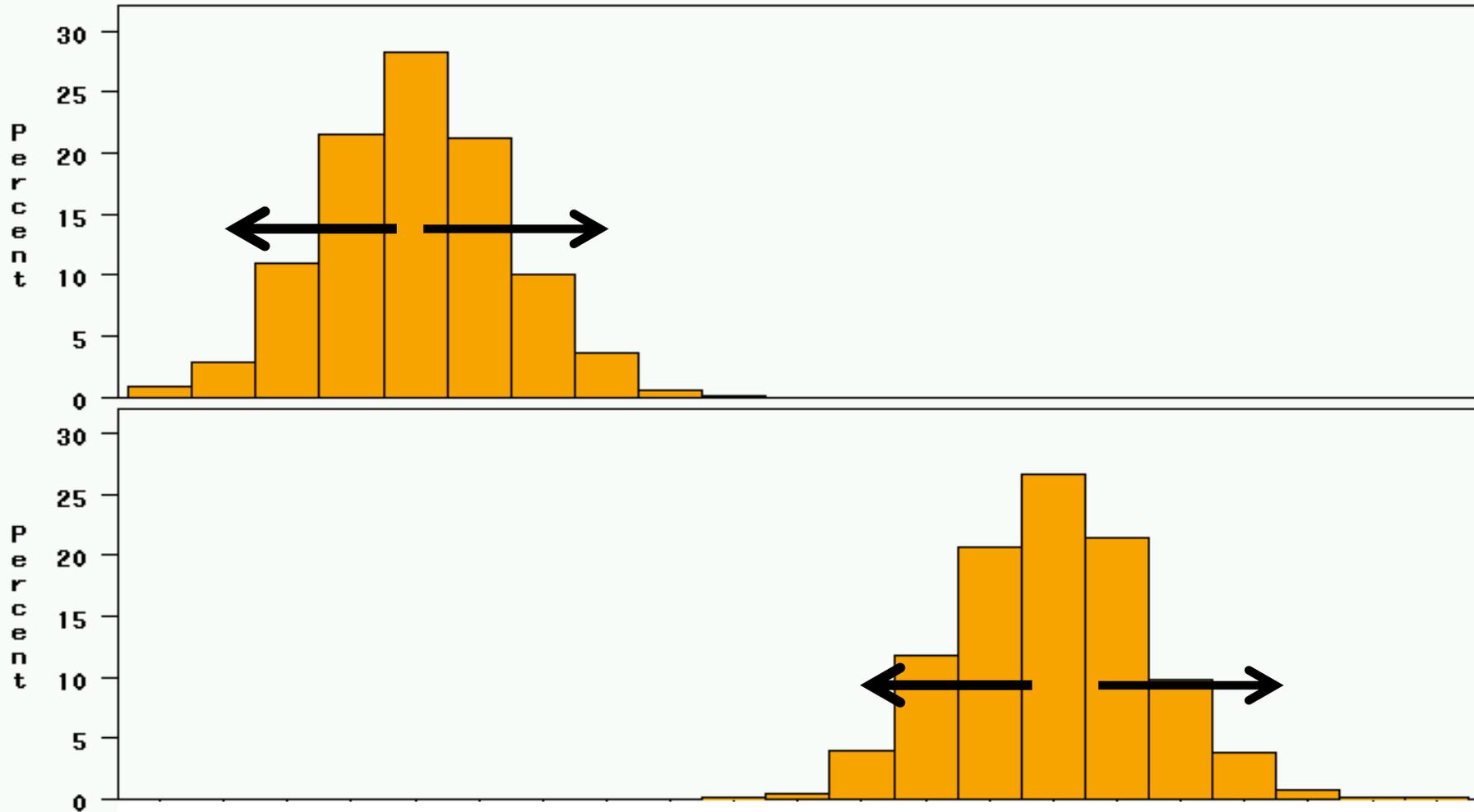
# Factors Affecting Power

1. Size of the effect ↑
2. Standard deviation of the variable ↓
3. Bigger sample size ↑
4. Significance level desired ↓

# 1. Bigger difference from the null mean

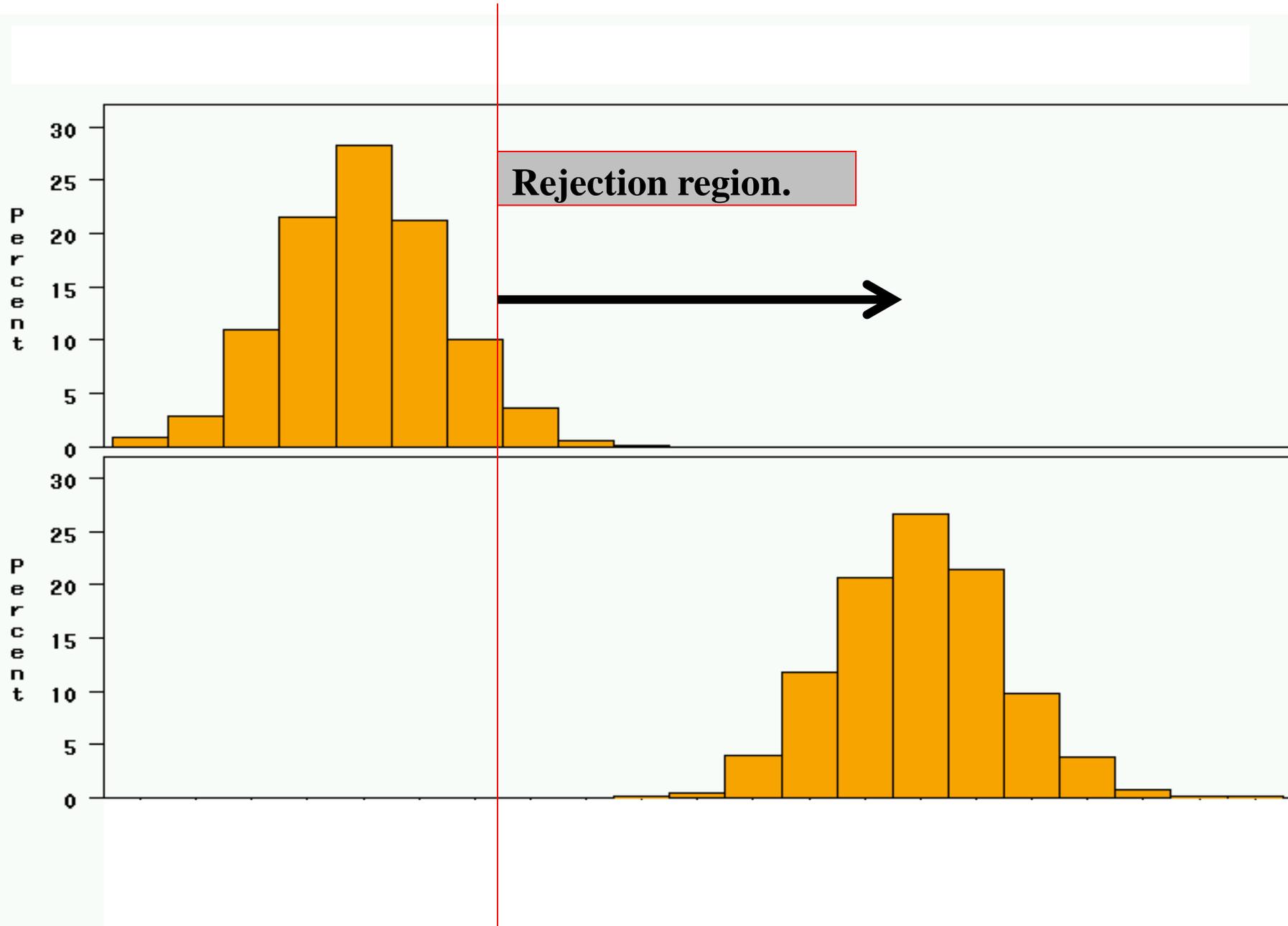


## 2. Bigger standard deviation





# 4. Higher significance level



# Effect Size

- This is an indication of the size of the treatment effect, its meaningfulness.
- With a large effect size, it will be easy to detect differences and statistical power will be high.
- But, if the treatment effect is small, it will be difficult to detect differences and power will be low.

# Effect Size

- Numerous authors have indicated the need to estimate the magnitude of differences between groups as well as to report the significance of the effects
- One way to describe the strength of a treatment effect, or meaningfulness of the findings, is the computation of “effect size” (ES)

$$ES = \frac{M_1 - M_2}{SD}$$

**Note: SD represents the standard deviation of the control group or the pooled standard deviation if there is no control group**

# A note on Effect Size

- There are many ways to define and calculate effect size
  - Difference in means
  - Variance explained
  - Odds ratio
- Standardised vs. unstandardised measures
  - If possible use unstandardized measures
    - Raw difference between group means
    - Raw regression coefficients
- Use standardised effect sizes as a last resort
  - Standardised difference (d): difference in means/SD
  - Pearson's correlation coefficient (r)

**Cohen's recommendations**

Effect	d	r
Small	≥0.2	≥0.1
Medium	≥0.5	≥0.3
Large	≥0.8	≥0.5

# Sample size calculations

- Based on these elements, you can write a formal mathematical equation that relates power, sample size, effect size, standard deviation, and significance level...



The higher the amplifying factor (=sample size)  
the smaller the structures to be seen!

# Formula for difference in *Proportions*

**Sample size** in each group  
(assumes equal sized groups)

Represents the **desired power** (typically .84 for 80% power).

$$n = \frac{2(\bar{p})(1 - \bar{p})(Z_{\beta} + Z_{\alpha/2})^2}{(p_1 - p_2)^2}$$

A measure of **variability**  
(similar to standard deviation)

difference in proportions

Represents the desired **level of statistical significance** (typically 1.96).

# Formula for difference in *Means*

**Sample size** in each group  
(assumes equal sized groups)

Represents the **desired power** (typically .84 for 80% power).

$$n = \frac{2\sigma^2 (Z_\beta + Z_{\alpha/2})^2}{\text{difference}^2}$$

**Standard deviation** of the outcome variable

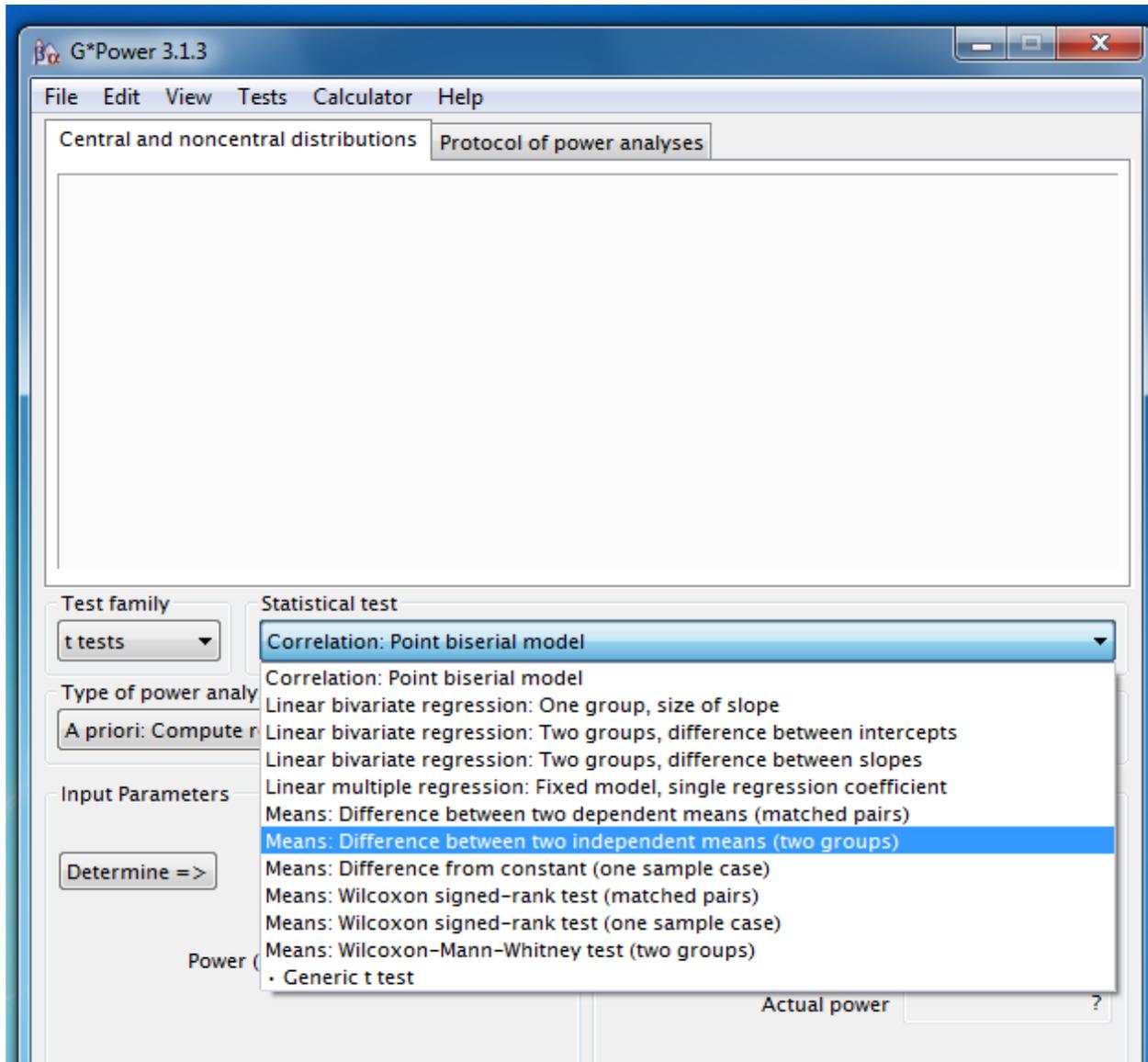
difference in means

Represents the desired **level of statistical significance** (typically 1.96).

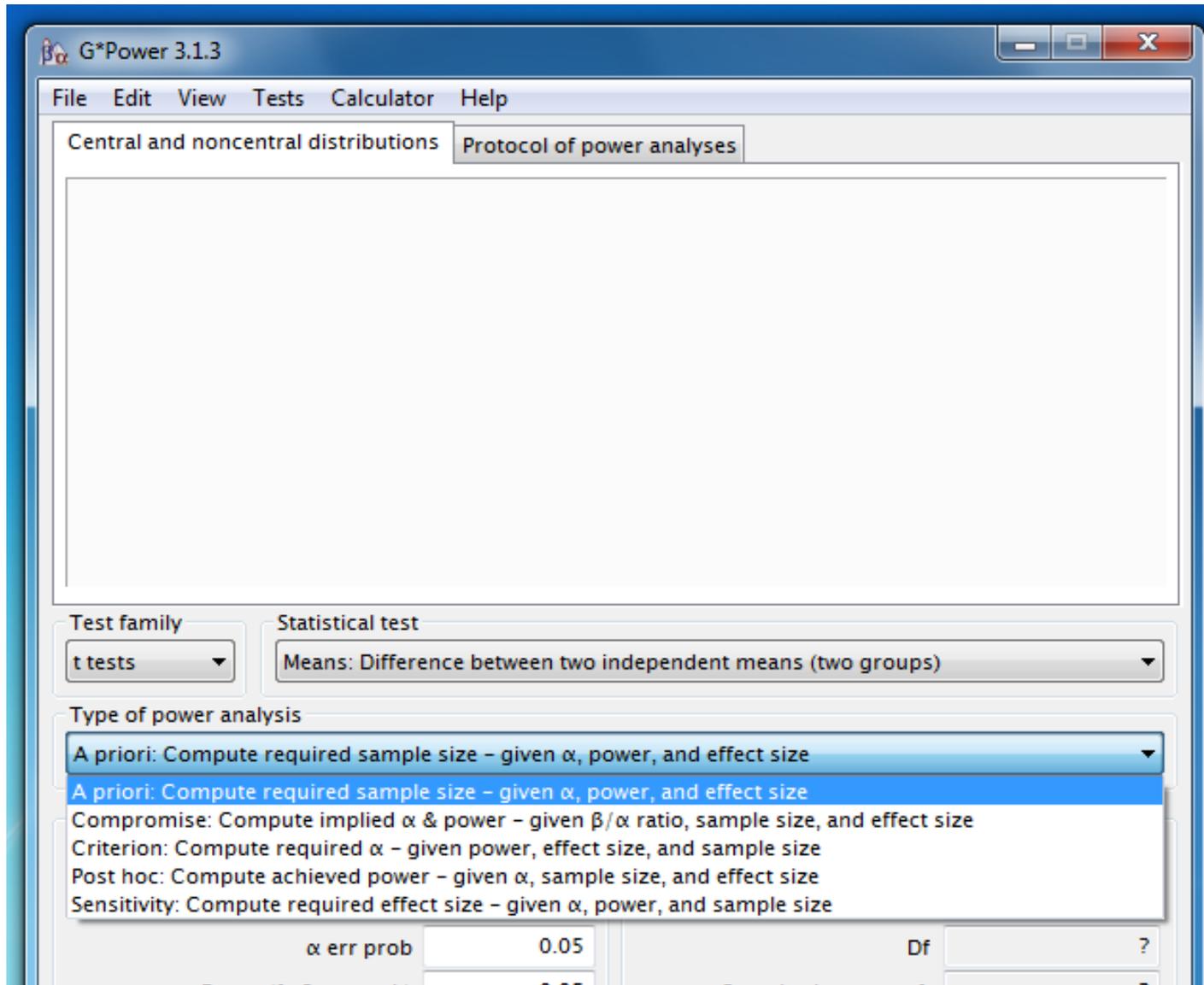
# Question

- How many test subjects would we need to sample if we want to have 80% power to detect an average increase in MCAT biology score of 1 point, if the average change without instruction (just due to chance) is plus or minus 3 points (=standard deviation of change)?
- Freely available software
  - G\*Power [www.psych.uni-duesseldorf.de/abteilungen/aap/gpower3/](http://www.psych.uni-duesseldorf.de/abteilungen/aap/gpower3/)
  - Quanto <http://hydra.usc.edu/gxe/>

# 1. Select the statistical test



## 2. Select the type of power analysis



# 3. Input the data characteristics to determine the effect size

The screenshot shows the G\*Power 3.1.3 software interface. The main window is titled "G\*Power 3.1.3" and has a menu bar with "File", "Edit", "View", "Tests", "Calculator", and "Help". The "Tests" menu is open, showing "Central and noncentral distributions" and "Protocol of power analyses".

The "Test family" is set to "t tests" and the "Statistical test" is "Means: Difference between two independent means (two groups)". The "Type of power analysis" is "A priori: Compute required sample size - given  $\alpha$ , power, and effect size".

**Input Parameters:**

- Tail(s): One
- Effect size d: 1.2000000
- $\alpha$  err prob: 0.05
- Power ( $1 - \beta$  err prob): 0.95
- Allocation ratio N2/N1: 1

**Output Parameters:**

- Noncentrality parameter  $\delta$ : ?
- Critical t: ?
- Df: ?
- Sample size group 1: ?
- Sample size group 2: ?
- Total sample size: ?
- Actual power: ?

On the right side, there are two radio buttons for "n1 != n2" and "n1 = n2". The "n1 = n2" option is selected. Below these are input fields for "Mean group 1", "Mean group 2", "SD  $\sigma$  within each group", "Mean group 1", "Mean group 2", "SD  $\sigma$  group 1", and "SD  $\sigma$  group 2".

At the bottom right, there is a "Calculate" button, an "Effect size d" field with the value 1.2, a "Calculate and transfer to main window" button, and a "Close" button.

At the bottom center, there is an "X-Y plot for a range of values" button and a "Calculate" button.

# 4. Input power parameters

The screenshot displays the G\*Power 3.1.3 software interface. The window title is "G\*Power 3.1.3" and the menu bar includes "File", "Edit", "View", "Tests", "Calculator", and "Help". The main area is divided into two tabs: "Central and noncentral distributions" and "Protocol of power analyses".

**Test family:** t tests  
**Statistical test:** Means: Difference between two independent means (two groups)

**Type of power analysis:** A priori: Compute required sample size – given  $\alpha$ , power, and effect size

**Input Parameters:**

- Tail(s): Two
- Determine =>
- Effect size d: 1.2000000
- $\alpha$  err prob: 0.05
- Power (1- $\beta$  err prob): 0.90
- Allocation ratio N2/N1: 1

**Output Parameters:**

- Noncentrality parameter  $\delta$ : ?
- Critical t: ?
- Df: ?
- Sample size group 1: ?
- Sample size group 2: ?
- Total sample size: ?
- Actual power: ?

**Group Parameters:**

- n1 != n2
  - Mean group 1: 0
  - Mean group 2: 1
  - SD  $\sigma$  within each group: 0.5
- n1 = n2
  - Mean group 1: 0.95
  - Mean group 2: 0.83
  - SD  $\sigma$  group 1: 0.10
  - SD  $\sigma$  group 2: 0.10

**Buttons:** Calculate, Calculate and transfer to main window, Close

**Footer:** X-Y plot for a range of values, Calculate

# 5. Draw plot for a range of values

