

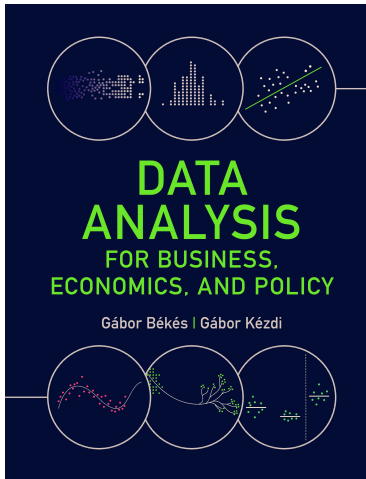
# 06 Testing hypotheses

**Gabor Bekes**

Data Analysis 1: Exploration

2019

# Slideshow for the Békés-Kézdi Data Analysis textbook



- ▶ Cambridge University Press, 2021 January
- ▶ Available in paperback, hardcover and e-book
- ▶ Slideshow be used and modified for educational purposes only
- ▶ [gabors-data-analysis.com](https://gabors-data-analysis.com)
  - ▶ Download all data and code
  - ▶ Additional material, links to references

## Motivation

- *The internet allowed the emergence of specialized online retailers while brick-and-mortar shops also sell goods on the main street. How to measure price inflation in the age of these options? To help answer this, we can collect and compare online and offline prices of the same products and test if they are the same.*

## The logic of hypothesis testing

- ▶ A hypothesis is a statement about a general pattern, of which we are not sure if true or not.
- ▶ **Hypothesis testing** = analyze our data to make a decision on the hypothesis
- ▶ Reject the hypothesis if there is enough evidence against it.
- ▶ Don't reject it if there isn't enough evidence against it.
- ▶ We may not have enough evidence against a hypothesis
  - ▶ if the hypothesis is true
  - ▶ or it is not true only the evidence is weak
- ▶ Important asymmetry here: rejecting a hypothesis is a more conclusive decision than not rejecting it.

## The logic of hypothesis testing: inference

- ▶ Testing a hypothesis: making inference with a focus on a specific statement.
- ▶ Can answer questions about the population, or general pattern, represented by our data.
- ▶ It is an inference: have to assess external validity

## The logic of hypothesis testing: the setup

- ▶ Define the *the statistic we want to test*,  $s$  (e.g. mean).
- ▶ We are interested in the true value of  $s$ ,  $s_{true}$ .
- ▶ The value the statistic in our data is its estimated value, denoted by a hat on top  $\hat{s}$ .

## The logic of hypothesis testing: $H_0$ vs $H_A$

- ▶ Formally stating the question as two competing hypotheses of which only one can be true: a **null** hypothesis  $H_0$  and an **alternative** hypothesis  $H_A$ .
- ▶ Formulated in terms of the unknown true value of the statistic.
- ▶ The null specifies some value/ range; the alternative specifies other possible values.
- ▶ Together, the null and the alternative cover all the possibilities we are interested in
- ▶ One example is null:  $s$  is zero, alternative:  $s$  is not zero.

$$H_0 : s_{true} = 0$$

$$H_A : s_{true} \neq 0$$

## The logic of hypothesis testing: $H_0$ vs $H_A$

- ▶ Our case study research question: Do the online and offline prices of the same products differ or are they the same?
- ▶ We have the price difference as our statistic and  $H_0 : s_{true} = 0$
- ▶ Logic: testing a hypothesis = see if there is enough evidence in our data to reject the null.
- ▶ The null is protected: it has to be hard to reject it otherwise the conclusions of hypothesis testing would not be strong.



## The logic of hypothesis testing: The criminal court example

- ▶ Logic of testing like a criminal court procedure.
  - ▶ Decide if the accused is guilty or innocent of a certain crime.
  - ▶ Assumption of innocence: accused judged guilty only if enough evidence against innocence
  - ▶ Even though the accused in court because of suspicion of guilt.
  
- ▶ To translate this procedure to the language of hypothesis testing,
  - ▶  $H_0$  is that the person is innocent
  - ▶  $H_A$  is that the person is guilty.

## The logic of hypothesis testing: $H_0$ vs $H_A$

- ▶ **Two-sided** alternative: The case when we test if  $H_A : s_{true} \neq 0$  - allows for  $s_{true}$  to be either greater than zero or less than zero. Not interested if the difference is positive or negative.
- ▶ **One-sided** alternative: interested if a statistic is positive or not.
- ▶ Different setup: the hypothesis we are interested in is in the alternative set.

$$H_0 : s_{true} \leq 0$$

$$H_A : s_{true} > 0$$

## Case Study - Comparing online and offline prices: Testing hypotheses

- ▶ Question: Do the online and offline prices of the same products differ?
- ▶ this data includes 10 to 50 products in each retail store included in the survey (the largest retailers in the U.S. that sell their products both online and offline).
- ▶ The products were selected by the data collectors in offline stores, and they were matched to the same products the same stores sold online.
- ▶ Let define our statistic as the difference in average prices.

## Case Study - Comparing online and offline prices: Testing hypotheses

- Descriptive statistics of the difference
- Each product  $i$  has both an online and an offline price in the data,  $p_{i,online}$  and  $p_{i,offline}$ ,  $pdiff$  is their difference:

$$pdiff_i = p_{i,online} - p_{i,offline} \quad (1)$$

The statistic with  $n$  observations (products) in the data, is:

$$s = \overline{pdiff} = \frac{1}{n} \sum_{i=1}^n (p_{i,online} - p_{i,offline}) \quad (2)$$

## Case Study - Comparing online and offline prices: Testing hypotheses

- The average of the price differences is equal to the difference of the average prices
- s statistic also measures the difference between the average of online prices and the average of offline prices among products with both kinds of price

$$\frac{1}{n} \sum_{i=1}^n (p_{i,online} - p_{i,offline}) = \frac{1}{n} \sum_{i=1}^n p_{i,online} - \frac{1}{n} \sum_{i=1}^n p_{i,offline}$$

## Case Study - Comparing online and offline prices: Testing hypotheses

### Descriptive statistics of the difference

- ▶ The mean difference is USD -0.05: online prices are, on average, 5 cents lower in this dataset.
- ▶ Spread around this average: Std: USD 10
- ▶ Extreme values matter: Range: -380 — USD +415.
- ▶ Of the 6439 products, 64% have the same online and offline price, for 87%, the difference within  $\pm 1$  dollars.

## Case Study - Comparing online and offline prices: the setup

### External validity

- ▶ The products in the data may not represent all products sold at these stores.
  - ▶ Could be a bias. **Example?**
- ▶ Strictly: The general pattern of the statistic represented by this dataset is average online-offline price differences in large retail store chains for the kind of products that data collectors would select with a high likelihood.
- ▶ More broadly: price differences among *all* products in the U.S. sold both online and offline by the same retailers.
  - ▶ Need an assumption. **What would it be?**

## Case Study - Comparing online and offline prices: the setup

Do average prices differ in the general pattern represented by the data?

$$H_0 : s_{true} = \bar{p}_{online\ true} - \bar{p}_{offline\ true} = 0 \quad (3)$$

$$H_A : s_{true} = \bar{p}_{online\ true} - \bar{p}_{offline\ true} \neq 0 \quad (4)$$



## The logic of hypothesis testing

- ▶ The **t-test** is the testing procedure based on the **t-statistic**
- ▶ We compare the estimated value of the statistic  $\hat{s}$  (our best guess of  $s$ ) to zero.
- ▶ Evidence to reject the null = difference between  $\hat{s}$  and zero.
- ▶ Reject if large: large difference means it is unlikely to be zero.
- ▶ Not reject the null if the estimate is not very far, i.e., when there is not enough evidence against it.

# T-test

- ▶ The **test statistic** is a statistic that measures the distance of the estimated value from what the true value would be if  $H_0$  was true.
- ▶ Uses estimated value of  $s$  ( $\hat{s}$ ) and the standard error of estimate ( $SE(\hat{s})$ ).
- ▶ Consider  $H_0 : s_{true} = 0, H_A : s_{true} \neq 0$ . The t-statistic for this hypotheses is:

$$t = \frac{\hat{s}}{SE(\hat{s})} \tag{5}$$

- ▶ The test statistic summarizes all the information needed to make the decision.
- ▶ When hypotheses are about value of one coefficient the test statistic = t-statistic

# T-test

When  $\hat{s}$  is the average of a variable  $x$ , the t-statistic is simply

$$t = \frac{\bar{x}}{SE(\bar{x})} \quad (6)$$

When  $\hat{s}$  is the average of a variable  $x$  minus a number, the t-statistic is

$$t = \frac{\bar{x} - number}{SE(\bar{x})} \quad (7)$$

When  $\hat{s}$  is the difference between two averages, say,  $\bar{x}_A$  and  $\bar{x}_B$ , the t-statistic is

$$t = \frac{\bar{x}_A - \bar{x}_B}{SE(\bar{x}_A - \bar{x}_B)} \quad (8)$$

# T-test

- ▶ If  $\hat{s} > 0$  = the t-statistic is positive; if  $\hat{s} < 0$  = the t-statistic is negative.
- ▶ With a two-sided alternative ( $H_A : s_{true} \neq 0$ ) it is the magnitude not the sign of the t-statistic that matters.
- ▶ If  $\hat{s}$  were zero the t-statistic would be zero. Never exactly zero. If the null is correct and thus  $s_{true}$  is zero we expect the  $\hat{s}$  estimate to be close to zero.
- ▶ Conversely, if the null is incorrect and thus  $s_{true}$  is not zero we expect the  $\hat{s}$  estimate to be far from zero.

# T-test

- ▶ The magnitude of the t-statistic - the distance of  $\hat{s}$  from what it should be if the null was true.
- ▶ We once again standardize, use  $SE(\bar{x})$
- ▶ May use  $SE(\bar{x}) = \sqrt{\frac{1}{n}} Std[x]$ .
- ▶ When  $\hat{s}$  is the difference of two averages, SE formula more complicated (R, Stata has it).
- ▶ Sometimes no appropriate SE formula for statistic interested in → What do we do?

## Making a decision

- ▶ In hypothesis testing the decision is based on a clear rule specified in advance.
- ▶ A decision rule makes the decision straightforward + transparent
- ▶ Helps avoid personal bias: put more weight on the evidence that supports our prejudices.
- ▶ Clear decision rules are designed to minimize the room for such temptations.

## Making a decision

- ▶ The decision rule = comparing the test statistic to a pre-defined **critical value**.
- ▶ Is test statistic is large enough to reject the null.
- ▶ Null rejected if the test statistic is larger than the critical value
- ▶ Critical value - between being too strict or too lenient.

## Making a decision

- ▶ When we make the decision, we may be right or wrong, don't know
- ▶ Need to think about it
- ▶ We can be right in our decision in two ways:
  - ▶ we reject the null when it is not true,
  - ▶ or we do not reject the null when it is true.
- ▶ We can be wrong in our decision in two ways, too:
  - ▶ we reject the null even though it is true,
  - ▶ or we do not reject the null even though it is not true.



## Making a decision

- ▶ We can be right in our decision in two ways:
  - ▶ we reject the null when it is not true,
  - ▶ or we do not reject the null when it is true.
- ▶ We can be wrong in our decision in two ways, too:
  - ▶ we reject the null even though it is true,
  - ▶ or we do not reject the null even though it is not true.

Retailer ID:	44	45	46	47	48	49	50	51
Diff	3.74	-1.2	-0.43	0.05	0.42	2.41	0.61	0.28
p-value	0.04	0.22	0.00	0.10	0.04	0.20	0.10	0.06
Retailer ID:	53	54	56	57	58	59	60	62
Diff	-0.97	-0.03	-0.49	0.93	-0.17	-0.53	-0.14	1.36
p-value	0.01	0.80	0.04	0.00	0.00	0.00	0.70	0.12

## Making a decision

- ▶ We say that our decision is a *false positive* if we reject the null when it is true.
  - ▶ “positive” because we take the active decision to reject the protected null.
  - ▶ medical: person has the condition that they were tested against
  - ▶ False positive = type-I error;
- ▶ Our decision is a *false negative* if we do not reject the null even though we should.
  - ▶ “negative” because we do not take the active decision
  - ▶ medical: result is “negative” = not have the condition
  - ▶ False negative = type-II error.

## Making a decision

- ▶ False positives and false negatives: both wrong, but not equally.
- ▶ Testing procedure protects the null: reject it only if evidence is strong
- ▶ The background assumption - wrongly rejecting the null (a false positive) is a bigger mistake than wrongly accepting it (a false negative).
- ▶ Decision rule (critical value) is chosen in a way that makes false positives rare.

## Making a decision

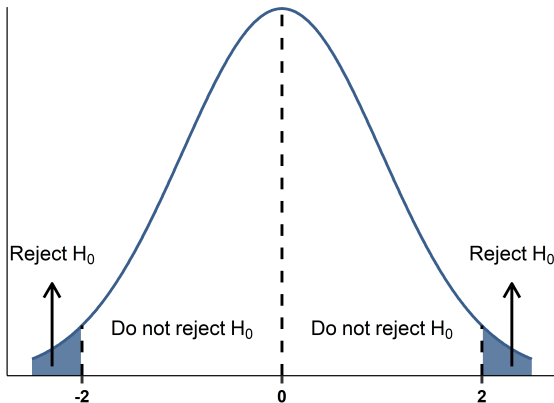
- ▶ A commonly applied critical value for a t-statistic is  $\pm 2$ :
  - ▶ reject the null if the t-statistic is smaller than  $-2$  or larger than  $+2$ ;
  - ▶ don't reject the null if the t-statistic is between  $-2$  and  $+2$ .
- ▶  $\text{Prob}(\text{t-statistic} < -2)$  or  $\text{Prob}(\text{t-statistic} > 2)$  are both appr 2.5%
- ▶ If the null is true: Probability t-statistic is below  $-2$  or above  $+2$  is 5%
- ▶ With  $\pm 2$  critical value - 5% is the probability of false positives - we have 5% as the probability that we would reject the null if it was true (False positive).
- ▶ If we make the critical values  $-2.6$  and  $+2.6$  the chance of the false positive is 1%.

## Critical values and generalization

- ▶ Why?
- ▶ We can calculate the likelihood of a false positive because we know what the sampling distribution of the test statistic would be if the null were true.
- ▶ The sampling distribution of a statistic is its distribution across repeated samples of the same size from the same population.
- ▶ Average: approximately normal, its mean is equal to the true mean, and its standard deviation is called the standard error.
- ▶ The t-statistic has the average in its numerator, (distribution is also approximately normal), its standard deviation is one because the denominator is the SE of  $\hat{s}$ .
- ▶ How the sampling distribution would look if the null hypothesis were true.
- ▶ Distribution of the t-statistic would be standard normal  $N(0, 1)$

## Sampling distribution of the test statistic when the null is true

- Distribution of the t-statistic would be standard normal  $N(0, 1)$
- Prob t-statistic  $< -2$  ( $> 2$ ) is approximately 2.5%. Prob t-statistic is  $< -2$  or  $> +2$  is 5% if the null is true.
- 5% = prob of false positives if we apply the critical values of  $\pm 2$  (=prob we reject the null if it was true)



## Critical values and generalization

- ▶ Can set other critical values that correspond to different probabilities of a false positive.
- ▶ That choice of 5% means that we tolerate a 5% chance for being wrong when rejecting the null
- ▶ Data analysts avoid biases when testing hypotheses: use the same critical value regardless of the data and hypothesis they are testing.

## Making a decision

- ▶ Fixing the chance of false positives affects the chance of false negatives at the same time.
- ▶ A false negative arises when the t-statistic is within the critical values and we don't reject the null even though the null is not true.
- ▶ Making a false negative call is more likely when it is harder to make a decision
- ▶ In what situations?



## Making a decision

- ▶ Fixing the chance of false positives affects the chance of false negatives at the same time.
- ▶ A false negative arises when the t-statistic is within the critical values and we don't reject the null even though the null is not true.
- ▶ Making a false negative call is more likely when it is harder to make a decision
  - ▶ Sample is small
  - ▶ The difference between true value and null is small

## Making a decision: size and power of the test

- ▶ **size** of the test: the probability of a false positive
- ▶ **level of significance**: The maximum probability of false positives we tolerate
- ▶ When we fix the level of significance at 5% and end up rejecting the null we say that the statistic we tested is significant at 5%
- ▶ **power** of the test: the probability of avoiding a false negative
- ▶ We usually fix the level of significance at 5% and hope for a high power (ie high probability of avoiding a false negative)
- ▶ High power is more likely when (i) the sample is large and (ii) and the further away the true value is from what's in a null.

# The p-value

- ▶ The p-value makes testing easier - captures info for reject/accept calls.
  - ▶ Instead of calculating test statistics and specify critical values, we can make an informed decision based on the p-value only.
- ▶ **p-value** is the smallest significance level at which we can reject  $H_0$  given the value of the test statistic in the sample.
- ▶ The p-value tells us the largest probability of a false positive.
- ▶ The p-value depends on
  1. the test statistic,
  2. the critical value
  3. the sampling distribution of the test statistic

# The p-value

- ▶ If the p-value is 0.05 the maximum probability that we make a false positive decision is 5%.
- ▶ If we are willing to take that chance, we should reject the null; if we aren't, we shouldn't.
- ▶ If the p-value is, say, 0.001 there is at most a 0.1% chance of being wrong if we were to reject the null.
- ▶ We can never be certain!  $p$  is never zero.
- ▶ For a reject/accept decision, one should pick a level of significance before the test
- ▶ What we can accept depends on the setting: what is the cost of a false positive.

## Case Study - Comparing online and offline prices: Testing hypotheses

- ▶ Let's fix the level of significance at 5%.
- ▶ Doing so we tolerate a 5% chance for a false positive.
- ▶ Allow a 5% chance to be wrong if we reject the null hypothesis of zero average price difference.
- ▶ A 5% level of significance translates to  $\pm 2$  bound for the t-statistic.
- ▶ The value of the statistic in the dataset is -0.054. Its standard error is 0.124.
- ▶ Thus the t-statistic is 0.44. This is well within  $\pm 2$ .
- ▶ Don't reject the null hypothesis of zero difference.
- ▶ (we do **not** say we proved it's zero. We showed we cannot tell it apart from zero. )

## Case Study - Comparing online and offline prices: Testing hypotheses

- ▶ Conclude that the average price difference is not different from zero in the general pattern represented by the data.
- ▶ Large dataset, good power. What we see in t-statistic is not because of very small sample size
- ▶ It is still possible that prices are indeed different, just the difference is very small. A few cent difference would not matter economically ...

## Case Study - Comparing online and offline prices: Testing hypotheses

- ▶ The p-value of the test is 0.66.
- ▶ That means that the smallest level of significance at which we can reject the null is 66%.
- ▶ The chance that we would make a mistake if we rejected the null is at most 66%.
- ▶ So we don't reject the null

## Case Study - Management quality - China vs India

- ▶ Where is management quality higher, in India or in China?
- ▶ This is testing difference of means across two sub-samples. (Also called Welch's t-test)
- ▶ To get the t-value we divide the difference between means with the SE of this difference.

Retailer ID:	44	45	46	47	48	49	50	51
Diff	3.74	-1.2	-0.43	0.05	0.42	2.41	0.61	0.28
p-value	0.04	0.22	0.00	0.10	0.04	0.20	0.10	0.06
Retailer ID:	53	54	56	57	58	59	60	62
Diff	-0.97	-0.03	-0.49	0.93	-0.17	-0.53	-0.14	1.36
p-value	0.01	0.80	0.04	0.00	0.00	0.00	0.70	0.12

Note: Source: Management quality is an average score of 18 variables. *wms-management-survey data*.



## Also good to know

Welch's t-test,  $t$  is:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\text{Var}[x_1]}{n_1} + \frac{\text{Var}[x_2]}{n_2}}}$$

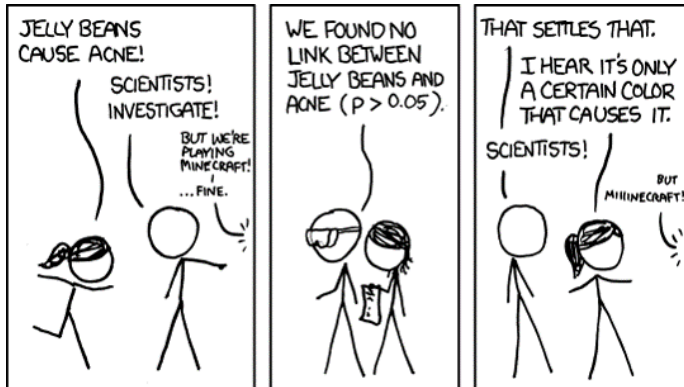
## Multiple testing

- ▶ Medical dataset: data on 400 patients
- ▶ A particular heart disease binary variable and 100 feature of life style (sport, eating, health background, socio-economic factors)
- ▶ Look for a pattern – is the heart disease equally likely for poor vs rich, take vitamins vs not, etc.
- ▶ You test one-by-one
- ▶ You find that for half a dozen factors, there is a difference
- ▶ Any special issue?

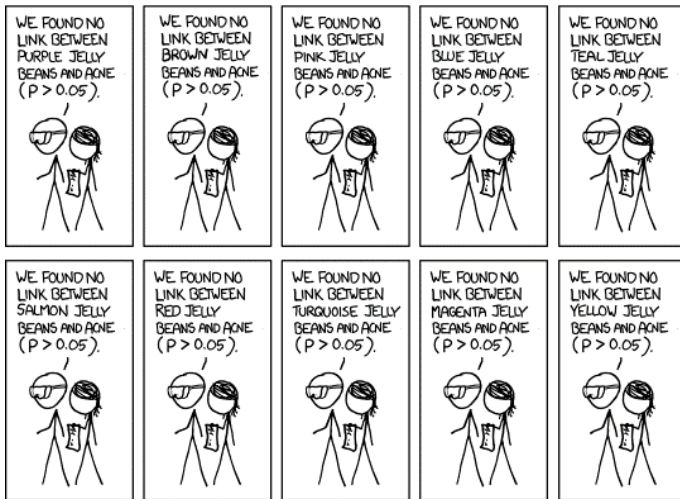
## Multiple testing

- ▶ The pre-set level of significance / p-value are defined for a single test
- ▶ In many cases, you will consider doing many many tests.
  - ▶ Different measures (mean, median, range, etc)
  - ▶ Different products, retailers, countries
  - ▶ Different measures of management quality
- ▶ For multiple tests, you cannot use the same approach as for a single one.

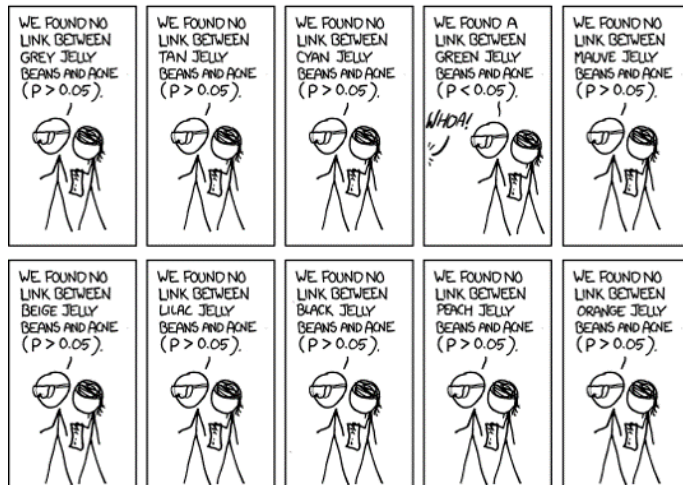
## Multiple testing - a serious example



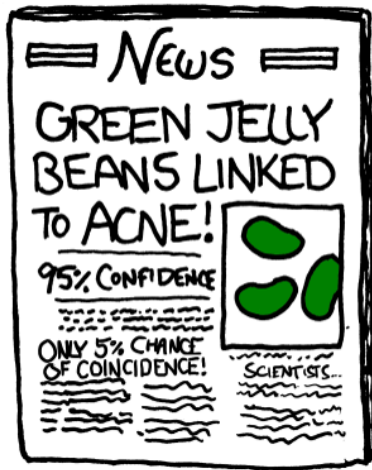
## Multiple testing - a serious example



## Multiple testing - a serious example



## Multiple testing - a serious example



## Multiple testing

- ▶ Consider a situation in which we test 100 hypotheses.
- ▶ Assume that all of those 100 null hypotheses are true.
- ▶ Set significance - we accept 5% chance to be wrong when rejecting the null. That means that we tolerate if we are wrong 5 out of 100 times.
- ▶ We can expect the null to be rejected 5 times when we test our 100 null hypotheses, all of which are true.
- ▶ In practice that would appear in 5 out of the 100 tests
- ▶ We could pick those five null hypotheses and say there is enough evidence to reject.
- ▶ But that is wrong: we started out assuming that all 100 nulls are true.
- ▶ Simply by chance, we will see cases when we would reject the null, but we should not



## Multiple testing

- ▶ There are various ways to deal with probabilities of false positives when testing multiple hypotheses.
- ▶ Often complicated.
- ▶ Solution 1: If you have a few dozens of cases, just use a strict criteria (such as 0.1-0.5% instead than 1-5%) for rejecting null hypotheses.
- ▶ A very strict such adjustment is the Bonferroni correction that suggests dividing the single hypothesis value by the number of hypotheses.
  - ▶ For example, if you have 20 hypotheses and aim for a  $p=.05$
  - ▶ reject the null only if you get a  $p=0.05/20=0.0025$
  - ▶ It is typically two strict

## Testing when data is very big

- ▶ With very large datasets some aspects of statistical inference lose their relevance.
- ▶ When the data has millions of observations generalizing to the general pattern does not add much.
- ▶ That is true for testing hypotheses, too.
- ▶ If, for example, two averages calculated from millions of observations are different to a meaningful extent they are almost surely different in the general pattern represented by the dataset.
- ▶ So: if you have millions of observations, just look at meaningful difference - do not worry about hypotheses testing (unless you care about very very small differences)

## Testing when data is very big

- ▶ Having few or many observations don't affect external validity in any way
- ▶ Recall step 2 of hypothesis testing: defining the general pattern represented by our data and comparing it to the general pattern we are interested in.
- ▶ So do not let the hype fool you. Just because the dataset is very large, it does not have to be representative and it does not necessarily have high external validity.