

Statistics and data analysis

Zohar Yakhini

IDC, Herzeliya

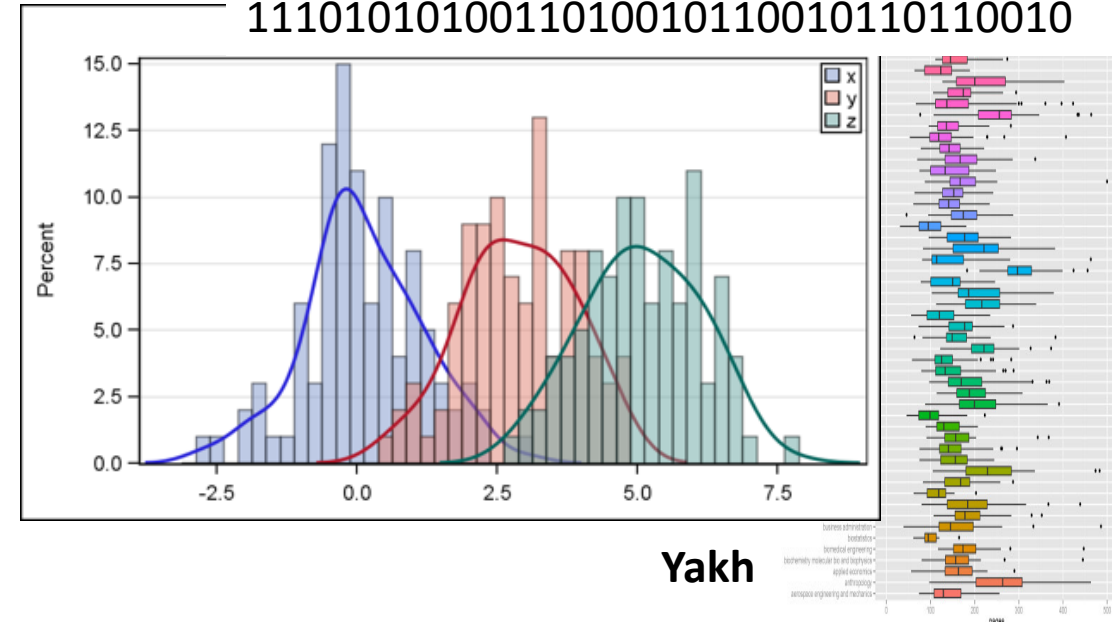


0001010110111100101010010100111000101011110010101010010100
101101001010010000111100101001000101001010100111101100101101001
0101001010101010100010101000010101110010101001010011100010
10111110010101010100101000101010010100100100001111001010010001010010
1010100111110110010110100101001010010101001010100100110010101
01010101001011110100111001101010111001101010101010011110010
11100101010010001010010101001111011001011001010101010101010
01010101001001100101010101010101010101011100101011010101001010
01000011110010101001000101010101010111101100101100101010101010
01010101001010101010011001010101010111101101011100110101101111
1110110010110010101010101010101010100001010111100101010101
01001101001110001010111100101010101010101001010100000111100
1010100100010101010101011111011001011000101010101010101010
10101010011001010101010101010111001110010101010111010011010
111001
110010101001000101010101010111101100101101001011100101010010100
110100111000101011110010
10010001010010101010111110110010101010101010101010101010101010
0100110010
10100111100101110010
01010010
110101010010100100001111001111101100



Intro Class

0010011101010100101010100100100010
1010100010101111101011010011001001
1110101010011010010110010110110010



Yakh

Data analysis

Analysis of data is the process of inspecting, cleaning, transforming, and modeling data with the goal of discovering useful information, suggesting conclusions, and supporting decision-making and inference of principles.

Analysis of data is a process of applying mathematical tools to extract useful and faithful information from observed data

Statistics is a scientific domain that investigates and characterizes tools that can be used in the process of analyzing data and in communicating and interpreting analysis results



Data analysis:

In analyzing observed data from the month of September we found that Covid19 positive testing rate in Haifa was 5.3% and in Jerusalem was 7.1%.

Statistics:

Do these results represent a significant difference?
Could they represent some random effect?

Further data analysis:

Are the differences related to any demographic parameters?
Get data about such potential parameters from more locations and compute correlations.

Statistics:

Are the observed correlations significant?
Can they be the result of some random effect?
Are there confounding factors involved?



Statistics and data analysis in the age of computers

The story of statistics is changing since:

- More efficient algorithms and computers make deeper and more elaborate calculations possible and practical
- The scope of data is changing in many respects, most notably volume

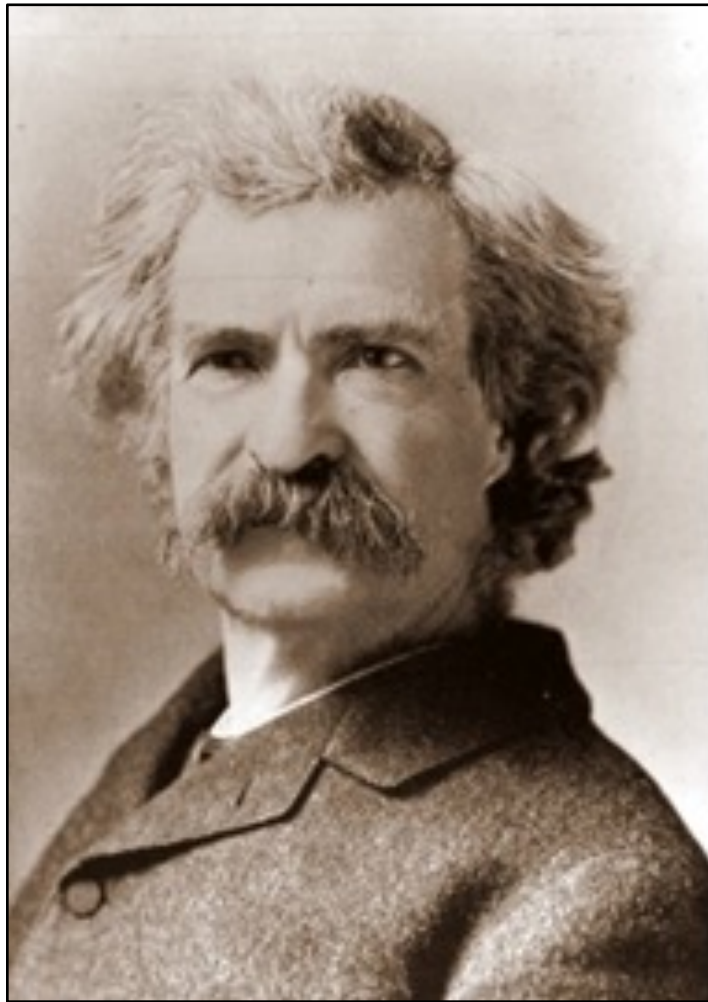


We will see many examples, including intro examples in this class ...

Statistics - defined



Statistics



**There are three kinds
of lies: lies, damned
lies, and statistics.**

**–Mark Twain
Chapters from My
Autobiography**



**The science of effectively drawing conclusions from data
OR
The science of effectively and convincingly lying**

Statistics is a common bond supporting all other sciences as well as all quantitative social and business investigations. It provides standards of empirical proof and a language for communicating results in these domains.

The process of statistical investigation includes:

- Designing experiments to maximize information
- Using models to describe observations and assess their significance
- Efficiently and effectively answering questions of interest
- Verifying the validity of the process
- Snooping around for more to be learned ...



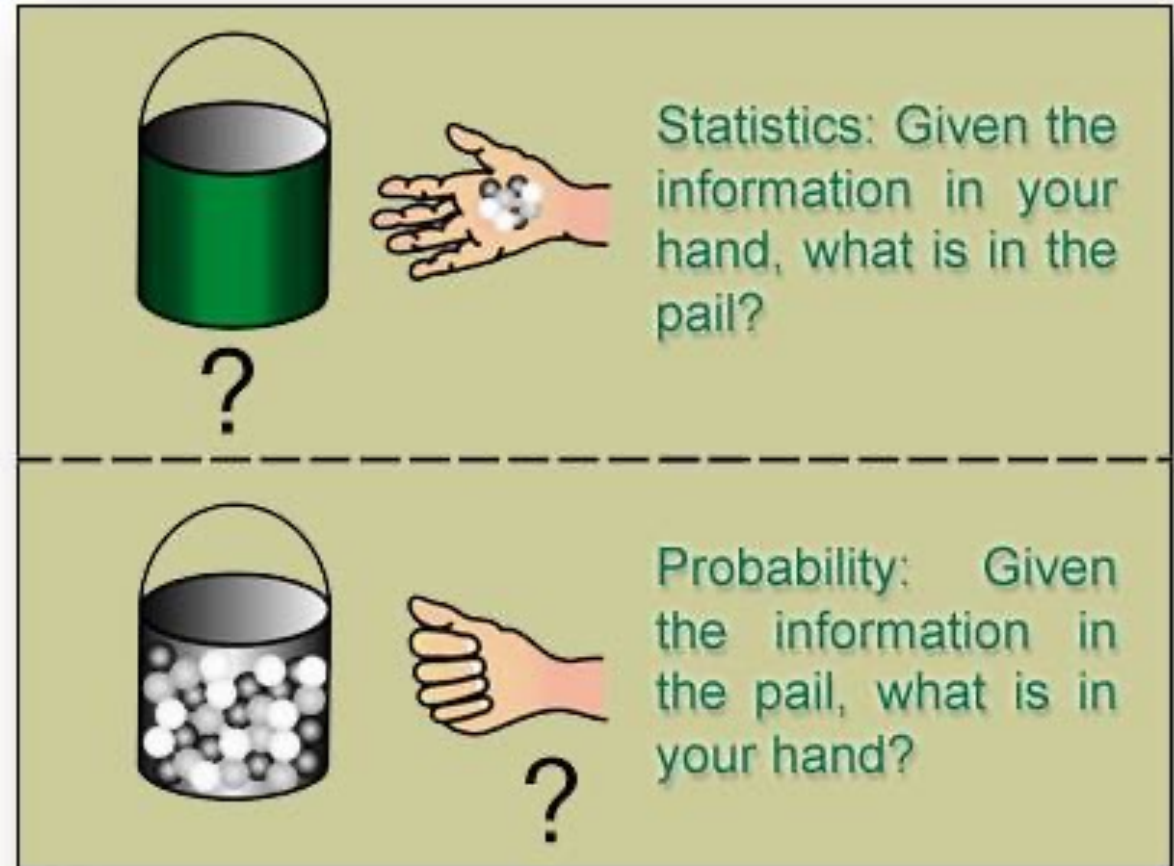
Adapted from F Ramsey and D Schafer, Oregon State Univ

Yakhini TASHPA

Probability theory and statistics

Statistics – given observations, what can we say about the underlying mechanism/system that gave rise to these observations?

Probability – assuming a model - what is the expected behavior of observations from the model?



In the age of computers, there are two separate aspects of the statistical enterprise:

1. Algorithmic developments aimed to draw conclusions from data.

For example:

Efficiently computing correlations for millions of quantitative records associated with users of a system or with a population of members of a healthcare service provider

OR

Using decision trees or random forests to make predictions about quantities of interest

2. Inference and assessment methodology.

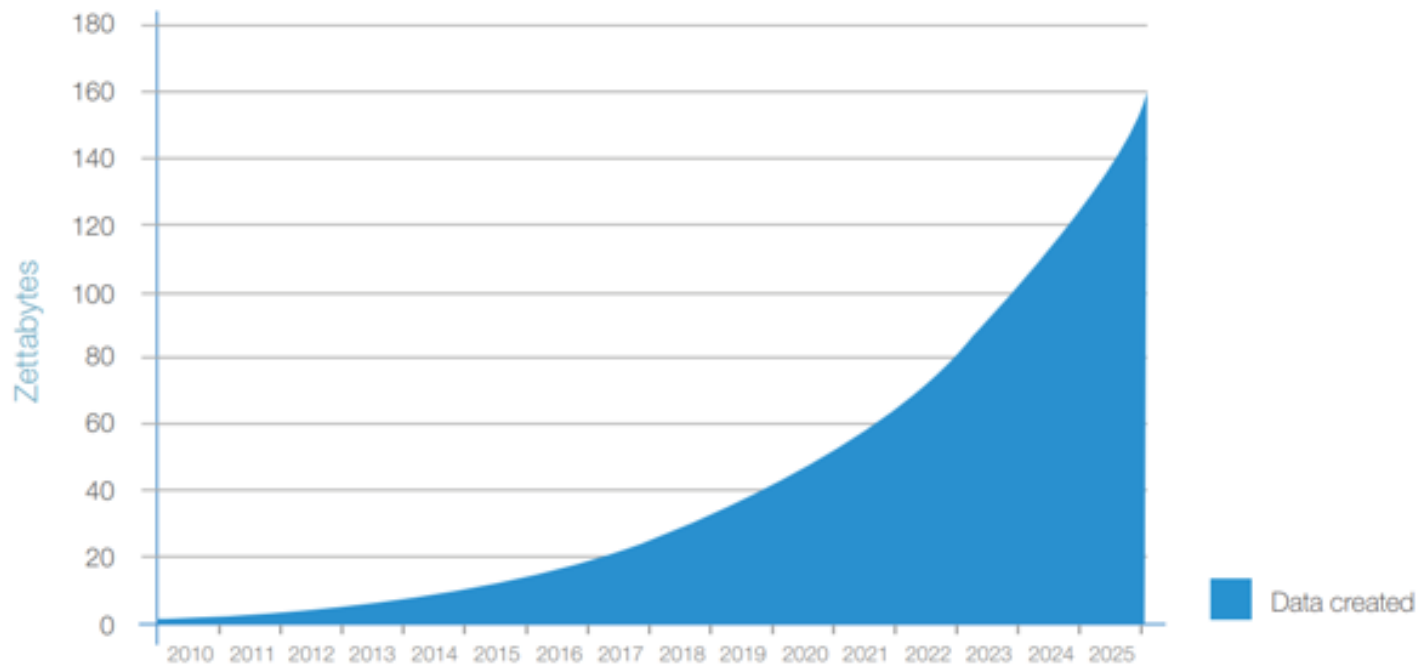
Aimed to test the validity of the results of the algorithm and to provide arguments to support the conclusions.



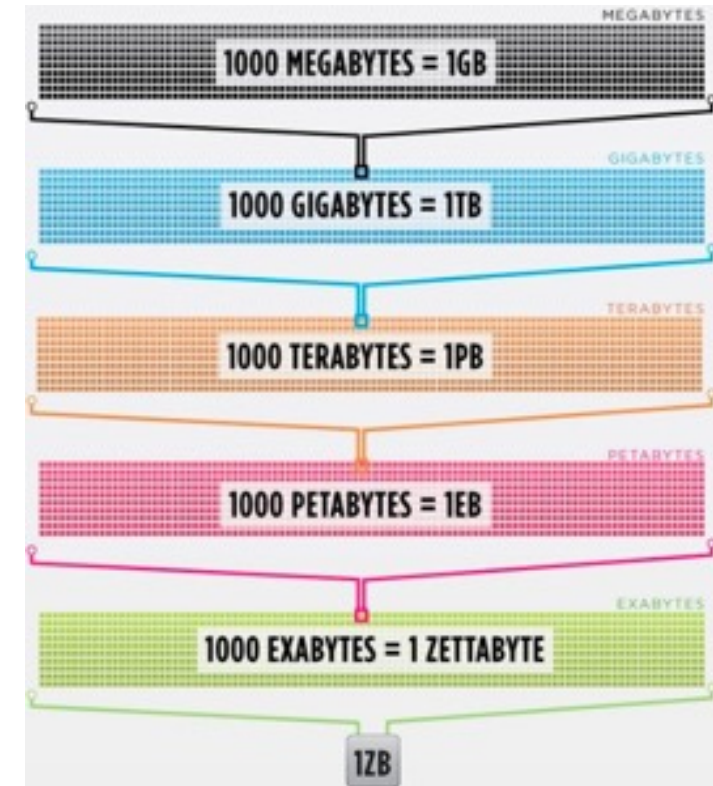
The science of effectively drawing conclusions from data
OR
The science of effectively and convincingly lying

DATA

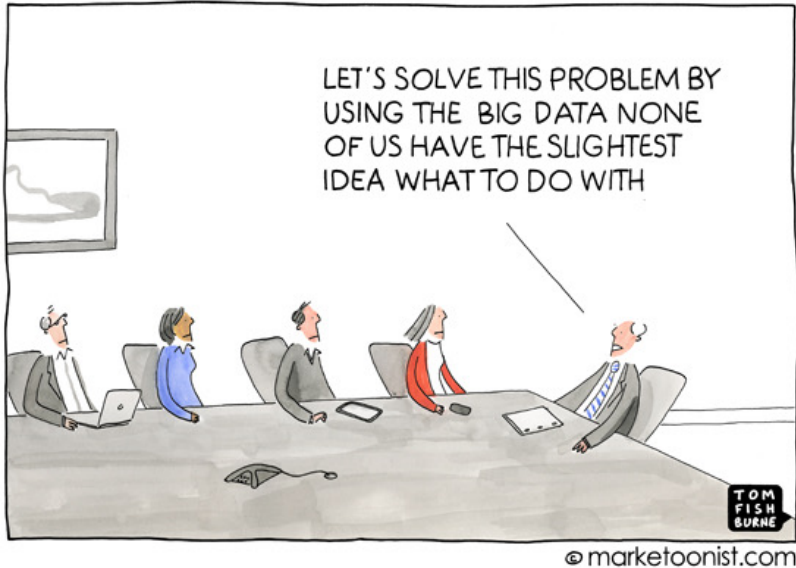
- Data is eating the world: 163 Zettabytes will be created in 2025



Source: IDC's Data Age 2025 study, sponsored by Seagate, April 2017



BIG DATA

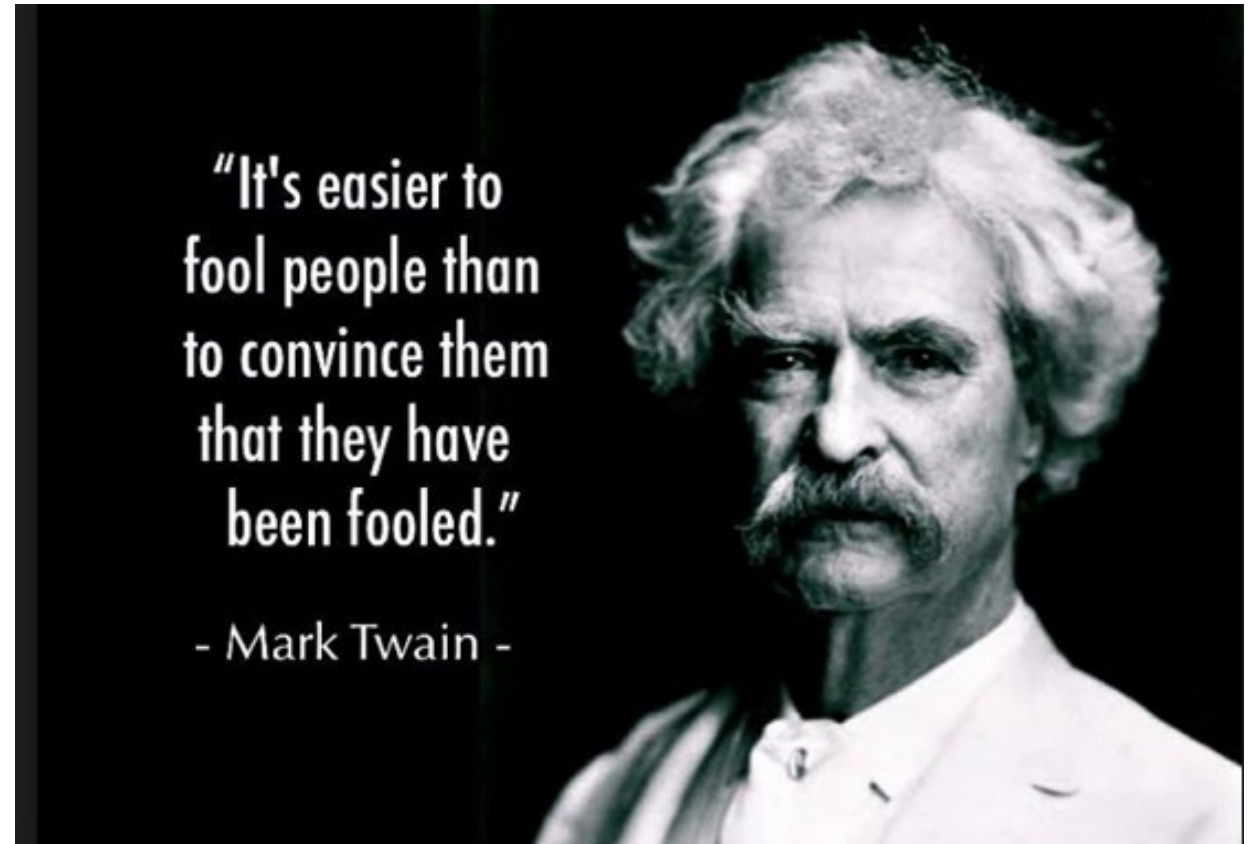


Yes - its a buzz word ...

But still – the utilization of data collected in many disciplines, from physics and molecular biology to environmental and social science, as well as by business oriented organizations, involves many steps and each one of them needs to be done using efficient and effective methods to get to desirable results:

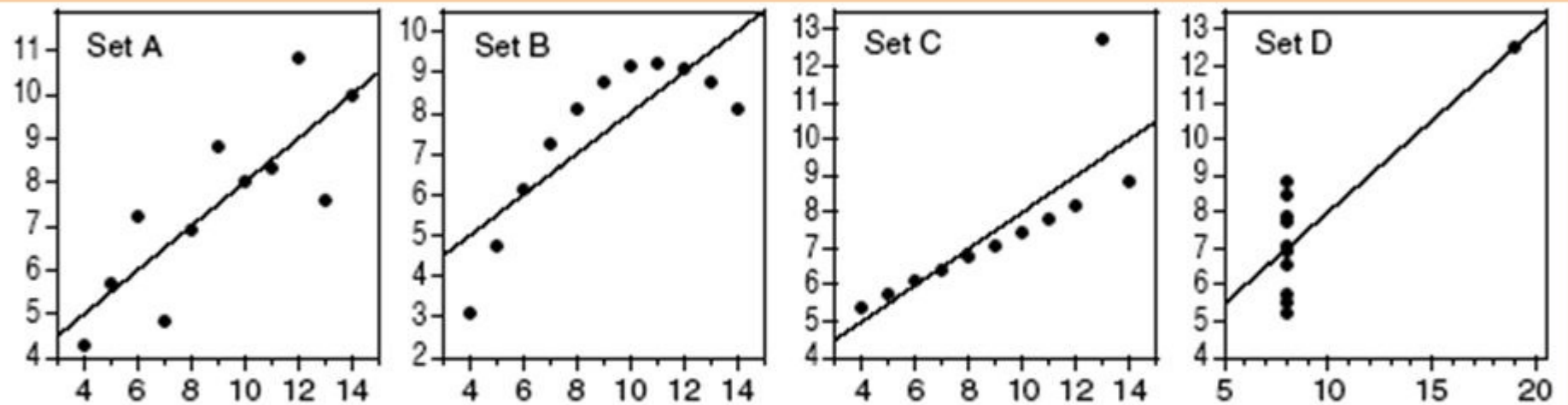
- Data collection and acquisition; Experiment design
 - + Garbage in garbage out
 - + Confounding factors, Representation
- Data storage, processing
 - + Accessibility and interaction
 - + Data cleaning
- **Statistical data analysis and efficient accurate inference**
- Conclusions
 - + **Visualization**
 - + **Reporting**
- Feedback
 - + Into data collection and other steps above

- Can we make statistical arguments simple and convincing?
- Visualization and presentation
- Clear and simple statements
- Rigorous and accurate methodologies
- The data scientist should know, to the greatest possible extent, what is standing behind any stated conclusion.

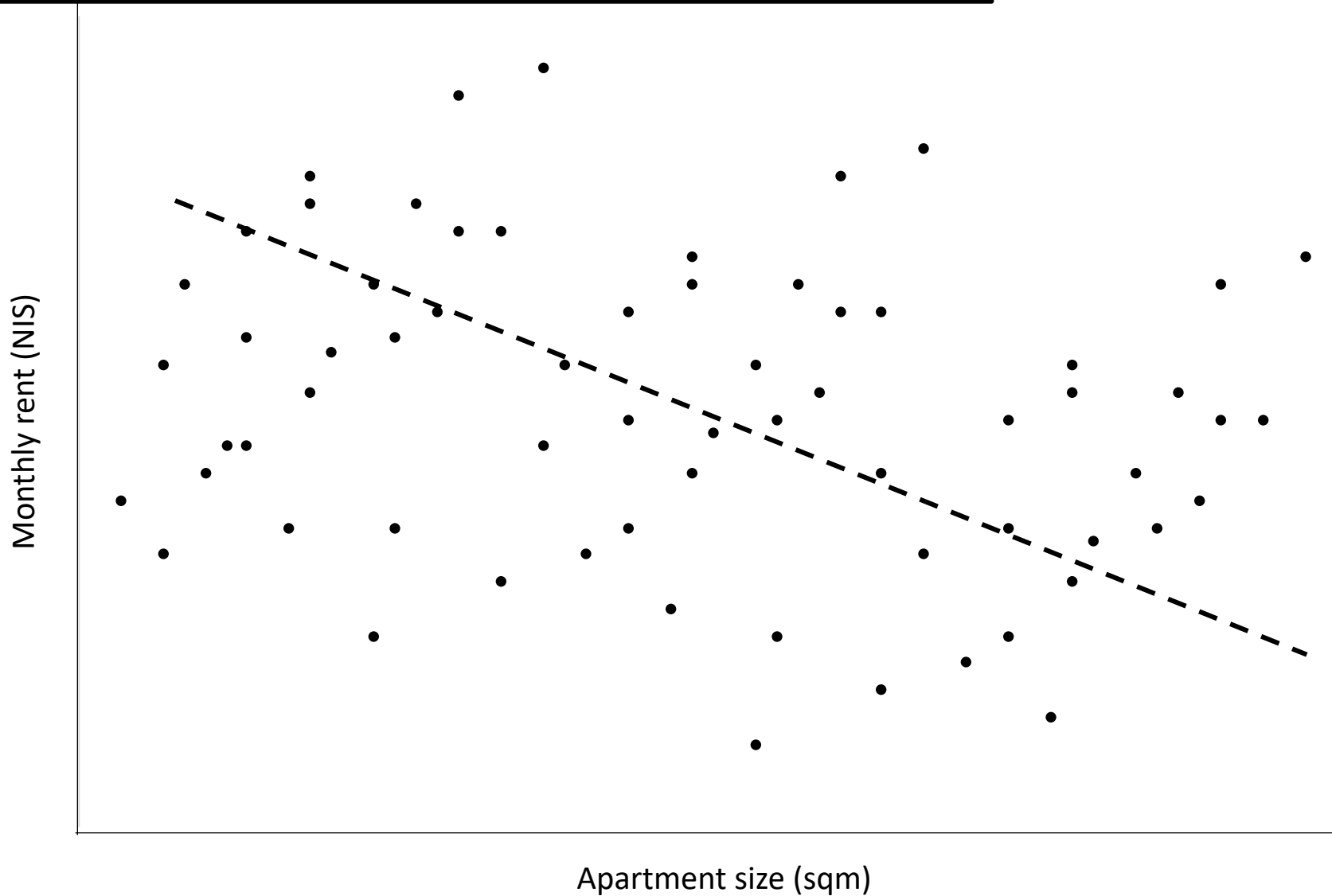


Effective inference – Example 1

Applying the same formal statistical tool on all of these yields the same result – high correlation between the variables.
But clearly – the actual conclusion should be different for each one of them



Example 2 –correlations??



Basketball

The Randomistan basketball association conducted a test in its top basketball leagues.

They recorded the number of basketball players who managed to score 5/5 free shots from the line.

The data is segmented by height and by gender.

	Less than 1.70m	1.70-1.90	Taller than 1.90
Women	4/6	4/6	8/9
Men	1/2	1/2	23/27



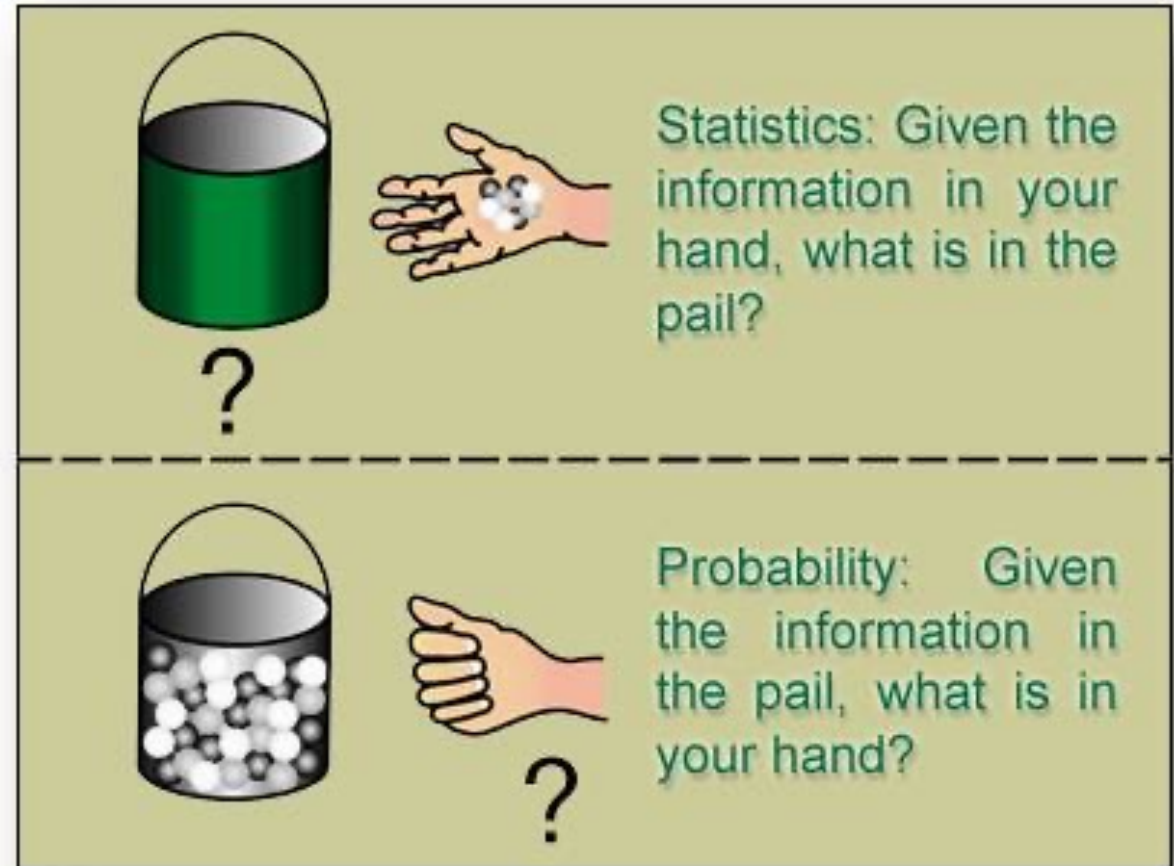
Who scores better from the line in Randomistan – men or women?



Probability theory and statistics

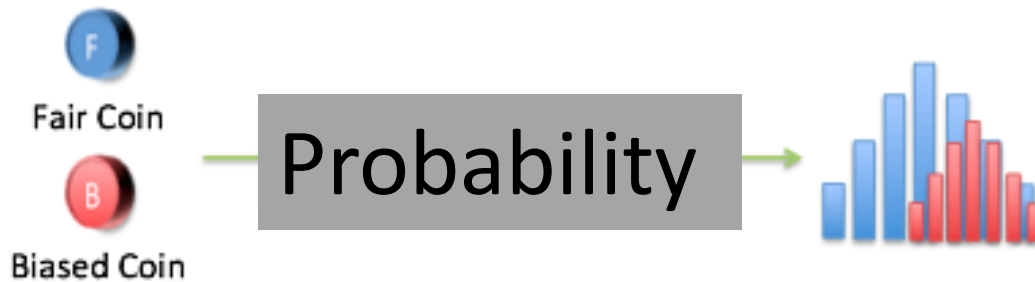
Statistics – given observations, what can we say about the underlying mechanism/system that gave rise to these observations?

Probability – assuming a model - what is the expected behavior of observations from the model?



Inference

Probability & Statistics



Probability

Given a model – determine the probability of occurrence of various data related events (including functions)



Statistics

Given observed data – Infer a model mechanism that could generate it; Quantify the inference

Course structure and formalities

- Time – Thu 1830-21hrs (H) OR Fri 845-1115 (E), with a 15 mins break
- 13 weeks
- 6 recitations: recitation will be once every two weeks (on average) – exact dates and times in Moodle (**first** recitation on Fri 11/11 and Wed 16/11)
- Prof Zohar Yakhini, office hours: C123 CS Bdg, Please email to coordinate.
Ben Galili, office hours: C314 CS Bdg, Please email to coordinate.
- zohar.yakhini@gmail.com , ben.galili@post.runi.ac.il
- Shuli Finely, Saar Buchnik.
Course Assistants will publish their office hours
- There will be 4 HW assignments that include theoretical aspects as well as practical data analysis mini-projects
- HW schedule will evolve as we go.
- HWAs will be due 2-3 weeks after assigned; Must be in pairs.
- Each HWA will be 18pts from the 100 grade points.
- The exam will be 28pts (MUST pass to pass class)

Topics to be covered

1. Introduction and review of probability theory
2. Important distributions
3. Statistical independence and what it means; Marginals and copulas
4. Data presentation and visualization
5. The binomial distribution and CLT
6. Foundations of statistical inference: confidence intervals, p-values, hypothesis testing
7. Parameter estimation and the EM algorithm
8. Correlation measures and how to use and misuse them
9. The hypergeometric distribution, inc examples and approximation schemes
10. Ranked lists, Wilcoxon rank sum, mHG
11. Multiple testing, Bonferroni and FDR corrections
12. Survival analysis, KM curves and the Mantel-Haenzel test
13. Entropy and information, KL and distances between distributions