# Experiments and Evaluation

## Part 1: Statistical Testing

Manfred Jaeger

Aalborg University

Introduction

**Schedule**

- ▶ Three lectures with exercises
- ▶ Exercises in two groups (morning/afternoon)
    - ▶ 10:15-12:00: SW
    - ▶ 12:30-14:15: DAT

**Hand-in Exercises**

- ▶ Started in exercise sessions following lectures (in small groups)
- ▶ **Individual Solutions** finished in last exercise session and extra time
    - ▶ Individual solution means: you must use your own words to describe your solution.
- ▶ Solutions for both exercises in one pdf document uploaded by Friday, March 29, 24:00.
- ▶ Put your name on the solution sheet!
- ▶ Pass/fail evaluation of solutions

**Survey Article** [mostly relevant for second/third part]

*R. L. Rardin and R. Uzsoy:*
*Experimental Evaluation of Heuristic Optimization Algorithms: A Tutorial. Journal of*
*Heuristics, 7 (2001)*

**Book** [further reading – when needed]

*A. B. Downey: Think Stats – Probability and Statistics for Programmers. Green Tea*
*Press, 2011. Available online:* http://greenteapress.com/thinkstats/

**Online Book**

http://onlinestatbook.com/index.html

**Wikipedia**

► Statistical hypothesis testing
► Student's t-test
► Wilcoxon signed-rank test

Claims:

- My algorithm returns a correct (optimal) solution on more than 90% of its inputs
- My program runs on average in less than 10s
- My algorithm/implementation is better than your algorithm/implementation
- The users of my web-site are happier than the users of your web-site
- . . .

How do we determine the validity of such claims, based on experimental data?

Often: Final section in a scientific article.

Goal for these lectures:

- ▶ Understanding the basic principles of (statistical) empirical evaluations
- ▶ Being able to perform basic tests on your own data
- ▶ Being able to identify possible strengths and weaknesses of a presented evaluation

**Program**

- ▶ First lecture: the more technical/analytical core
- ▶ Second lecture: technical/analytical core, experimental process
- ▶ Third lecture: experimental process (data collection)

Empirical Evaluations are based on experimental or observational data:

**Measurements of performance on test cases**



**My System**

| My System | | | | | |
|---|---|---|---|---|---|
| 1 | 0.55 | 6 | 1.2 | 11 | 0.52 |
| 2 | 0.01 | 7 | 0.91 | 12 | 0.88 |
| 3 | 0.87 | 8 | 0.02 | 13 | 0.45 |
| 4 | 1.3 | 9 | 1.01 | 14 | 0.03 |
| 5 | 0.54 | 10 | 0.76 | 15 | 0.65 |

☞ One sample tests

**Measurements of performance on test cases**

| A | | | | | |
|---|---|---|---|---|---|
| 1 | 0.55 | 6 | 1.2 | 11 | 0.52 |
| 2 | 0.01 | 7 | 0.91 | 12 | 0.88 |
| 3 | 0.87 | 8 | 0.02 | 13 | 0.45 |
| 4 | 1.3 | 9 | 1.01 | 14 | 0.03 |
| 5 | 0.54 | 10 | 0.76 | 15 | 0.65 |

**Setup A**

| B | | | | | |
|---|---|---|---|---|---|
| 1 | 0.88 | 6 | 0.67 | 11 | 0.52 |
| 2 | 0.21 | 7 | 0.63 | 12 | 0.73 |
| 3 | 0.54 | 8 | 0.87 | 13 | 1.05 |
| 4 | 1.22 | 9 | 0.41 | 14 | 0.87 |
| 5 | 1.54 | 10 | 1.76 | 15 | 1.56 |

**Setup B**

☞ Two sample tests

**Setups** can be ... different algorithms, different implementations of the same algorithm, different hardware platforms, different user interfaces, . . .

**Test cases** can be ... multiple runs on different inputs, interactions by different users with a web site, . . .

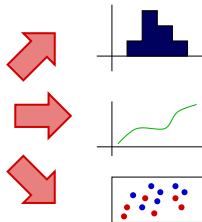**Measurements** can be ... time and/or space consumption, user satisfaction, quality of solution, . . .
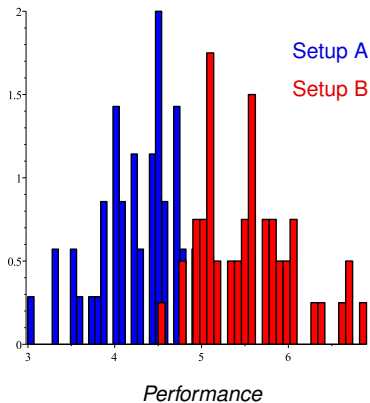
Descriptive Statistics

**Descriptive Statistics:** present/summarize most important aspects of the given data using suitable

- quantitative summarizations (means, extreme values ...)
- visualization tools

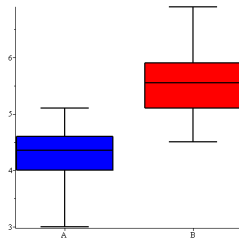*Performance*

- ▶ Quite detailed data summary
- ▶ Can be difficult to compare 2 or more setups

- ▶ $Q_1$: First quartile: 25% of data lies below this point, 75% above
- ▶ $Q_3$: Third quartile: 75% of data lies below this point, 25% above
- ▶ Median = $Q_2$: 50% of data lies below this point, 50% above
- ▶ Other features can be added to the Box Plot

- ▶ For paired measurements
- ▶ Can reveal existence of input types with different performance characteristics for A,B
- ▶ Also used for plotting two measurements, e.g. *time* and *space* for a single algorithm

y's same magnitude as $x's$, no correlation



y's same magnitude as $x's$, correlated



y's larger than $x's$, no correlation



y's larger than $x's$, correlated

Hypothesis Testing

**Inferential Statistics**

Making inferences from the data about the underlying populations/processes ... that produced the data. Types of inferential statistics:

- ▶ Testing (our topic)
- ▶ Estimation

**Hypothesis Testing**

- ▶ At some point, a decision has to be made:
  - ▶ Do we treat patients with the new drug, or do we keep the old?
  - ▶ Is our system good enough to be sold to customers?
  - ▶ Which algorithm should we use?
- ▶ For this, the information content of the experimental data has to be reduced to a simple binary decision:

    Statistical decision: *accepting* or *rejecting* a hypothesis
    ⇓
    Real-life (e.g. business) decision: deploy system, . . .

  ☞ The hypotheses that we want to test refer to *the underlying processes or populations*, not the known data.

Hypothesis relates to a **quantitative performance measure**:

- ▶ Runtime
- ▶ Memory usage
- ▶ User satisfaction
- ▶ Average \$ amount spend by visitors of web site
- ▶ ...

**One sample**: hypotheses make a statement about the performance:

*My Algorithm/Implementation/Setup/Design satisfies performance condition X*

**Two sample**: hypotheses make a comparison:

*[Algorithm/Implementation/Setup/Design] A*
*performs better than*
*[Algorithm/Implementation/Setup/Design] B*

One Sample Tests

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1      | 1       | 0       |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1      | 1       | 0       |
| 5      | 3       | 2       |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1      | 1       | 0       |
| 5      | 3       | 2       |
| 10     | 5       | 5       |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1      | 1       | 0       |
| 5      | 3       | 2       |
| 10     | 5       | 5       |
| 20     | 12      | 8       |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1      | 1       | 0       |
| 5      | 3       | 2       |
| 10     | 5       | 5       |
| 20     | 12      | 8       |
| 50     | 26      | 24      |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1 | 1 | 0 |
| 5 | 3 | 2 |
| 10 | 5 | 5 |
| 20 | 12 | 8 |
| 50 | 26 | 24 |
| 100 | 52 | 48 |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1 | 1 | 0 |
| 5 | 3 | 2 |
| 10 | 5 | 5 |
| 20 | 12 | 8 |
| 50 | 26 | 24 |
| 100 | 52 | 48 |
| 500 | 261 | 239 |

Is this coin fair?

| Trials | # Heads | # Tails |
|--------|---------|---------|
| 1      | 1       | 0       |
| 5      | 3       | 2       |
| 10     | 5       | 5       |
| 20     | 12      | 8       |
| 50     | 26      | 24      |
| 100    | 52      | 48      |
| 500    | 261     | 239     |
| 1000   | 557     | 443     |

If the coin lands *Heads* with probability $p$, then the probability of observing in $N$ tosses exactly $k$ heads is

$$B(N,p)(k) = \binom{N}{k} p^k (1-p)^{N-k}$$

Plotted for $N = 20$ and $p = 0.5$:



From these numbers:

$$P(\#H \geq 12) \approx 0.12 + 0.073 + 0.0369 + 0.0147 + 0.0046 = 0.249$$

We can read off sums

$$P(\#H \geq k) = P(\#H = k) + P(\#H = k+1) + \ldots + P(\#H = N)$$

directly from the *Cumulative Distribution Function (CDF)*:



$$CDF(k) = \sum_{j \leq k} P(\#H = j); \qquad P(\#H \geq k) = 1 - CDF(k-1)$$

If the coin is fair, then

- the probability of seeing in 20 tosses a number of heads greater or equal the observed $\#H = 12$ is $\approx 0.25$
- the probability of seeing in 20 tosses a number of heads that deviates from the expected number of 10 heads by the observed difference 12-10 =2 is

$$P(\#H \leq 8) + P(\#H \geq 12) \approx 2 \cdot 0.25 = 0.5$$

☞ The observed experimental outcome is not very unlikely under the fairness hypothesis

☞ We should not **reject** the fairness hypothesis ($p = 0.5$) on the basis of the data

Looking at $N = 500$ and $\#H = 261$:



$P(\#H \geq 261) \approx 0.18 \qquad P(\#H \geq 261 \text{ or } \#H \leq 239) \approx 0.36$

☞ Also no strong indication that coin is not fair

Looking at $N = 1000$ and $\#H = 557$:



$$P(\#H \geq 557) \approx 0.00015 \qquad P(\#H \geq 557 \text{ or } \#H \leq 443) \approx 0.0003$$

☞ The observed data is extremely unlikely if the fairness hypothesis was true

☞ We should reject the fairness hypothesis

**p-value (one sided)**: Probability under the fairness hypothesis of observing at least as many heads as in the sample

**p-value (two sided)**: Probability under the fairness hypothesis of observing at least as extreme a deviation from the expected number of heads as in the sample

| Trials | # Heads | # Tails | *p-value* (two-sided) |
|--------|---------|---------|-----------------------|
| 1      | 1       | 0       |                       |
| 5      | 3       | 2       |                       |
| 10     | 5       | 5       | 1.0                   |
| 20     | 12      | 8       | 0.5034                |
| 50     | 26      | 24      | 0.8877                |
| 100    | 52      | 48      | 0.7643                |
| 500    | 261     | 239     | 0.3476                |
| 1000   | 557     | 443     | 0.000347              |

☞ Reject the hypothesis when the p-value is sufficiently small

▶ Fix a level of significance $\alpha$. Typical:

$$\alpha = 0.05 \qquad \text{or} \qquad \alpha = 0.01$$

▶ *Reject* the hypothesis if the p-value obtained from the sample is $\leq \alpha$

☞ The significance level should be set *before* the experiment or data analysis begins. *Not:*

*O.k. — I get a p-value of 0.013; I would really like to reject the hypothesis, so let's set $\alpha = 0.05$*

**One-sided or Two-sided?**

One-sided: "upper/lower bound" hypothesis, such as: $p \leq 0.5$

Two-sided: "point" hypothesis, such as: $p = 0.5$

Observed runtimes of implementation on test cases:

| Trials | Times (in ms) |
|---|---|
| 1 | 4.4 |
| 5 | 6.41, 3.36, 4.71, 4.07, 6.41 |
| 100 | 4.17, 7.05, 4.38, 5.44, 5.69, 5.45, 5.76, . . . |
| 1000 | 4.76, 4.73, 6.42, 3.30, 6.44, 3.78, 7.62, 4.99, . . . |

Hypothesis: the average runtime is greater or equal 7ms (e.g.: customer defined performance requirement).

Binomial: each sample generated by a 0/1 valued Bernoulli variable characterized by parameter $p$

Assume that runtimes follow a **Gaussian distribution** with some *mean $\mu$* and *standard deviation $\sigma$*



N(7,1) density function

Plot for $\mu = 7.0, \sigma = 1.0$

Binomial: just count the number of heads in a sample

Compute from the raw data the relevant *test statistic*. Here from a sample of size *N*

$$x_1, x_2, \ldots, x_N$$

compute:

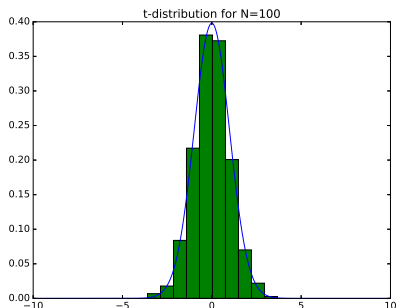**Sample Average** : $\bar{x} := \frac{1}{N} \sum_{i=1}^{N} x_i$

**Sample Variance**: $\bar{s} := \frac{1}{N} \sum_{i=1}^{N} (\bar{x} - x_i)^2$

**t-Statistic**: $t := \sqrt{\frac{N}{\bar{s}}} (\bar{x} - 7.0)$ (contains the hypothesized mean value of 7.0!)

| Trials | $\bar{x}$ | $\bar{s}$ | $t$ |
|--------|------|------|---------|
| 1      | 4.4  | 0.0  | $-\infty$ |
| 5      | 4.99 | 1.51 | -3.644  |
| 100    | 5.44 | 0.75 | -17.903 |
| 1000   | 5.54 | 1.0  | -45.93  |

Binomial: if every sample point comes from a Bernoulli distribution with parameter $p$ then #successes has a Binomial distribution with parameters $N, p$

If every sample point comes from a normal distribution with mean $\mu$, then the t-statistic has a **t-distribution with N-1 degrees of freedom**



t-distribution for N=100
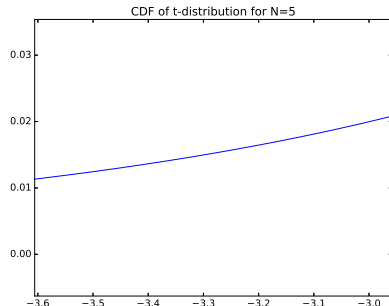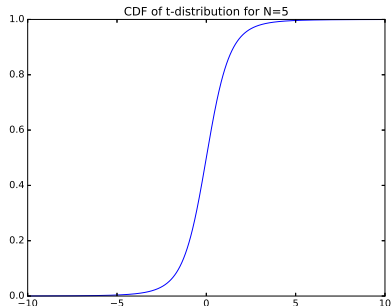
Blue line: t-distribution with 99 degrees of freedom

Green histogram: values of t-statistics computed for 1000 datasets of size 100 (datapoints sampled from normal distribution with mean 7.0).

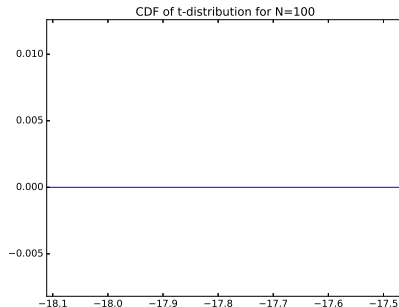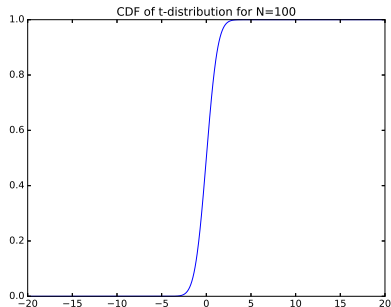Binomal: evaluating the CDF of the binomial distribution under the (null) hypothesis at the observed # successes

Evaluate the CDF of t-distribution at the observed value of the t-statistic (computed for the (null) hypothesis).

Our data at $N = 5$, $t = -3.644$



CDF of t-distribution for N=5

CDF of t-distribution for N=5

☞ p-value (one-sided) is $\approx 0.011$.

Our data at $N = 100$, $t = -17.903$



☞ p-value (one-sided) is $\approx 0.0$.

|  | **Binomial** | **Normal** |
|---|---|---|
| **Data** | $\{0, 1\}$ | Real |
|  | ⇓ | ⇓ |
| **Statistics** | #Successes | t-statistic |
|  | ⇓ | ⇓ |
| **Distribution** | Binomial | t-distribution |

⇓: place where we injected the null hypothesis