# Experiments and Evaluation

### Part 3: Performing and Evaluating Experiments

Manfred Jaeger

Aalborg University

Review

**Hypothesis about**

**Test**

Accept/Reject

Interpretation,
Documentation,...

**Algorithms**

**Experiment**

**Users**

**Hardware**

**Observations**

**Data**

**Implementations**

Warnings and Disclaimers

A more fundamental concern is that statistical significance is too often confused with *practical significance*—whether an observed effect is large enough to matter. The tiniest of effects can be statistically significant if enough points are sampled to make the likely random error even tinier. That does not mean the effect is important, or even worth reporting in a formal presentation of experimental results. A practically significant effect should not be claimed unless it is statistically significant, because it may be a random accident, but a demonstration of statistical significance does not prove practical importance. Analysts must take care not to let formal statistical tests dominate their investigations or substitute for their judgment and insight about the problems and algorithms addressed.

[R. L. Rardin and R. Uzsoy: Experimental Evaluation of Heuristic Optimization Algorithms: A Tutorial]

- Go to a bank and get 1000 5Kr. coins

▶ Go to a bank and get 1000 5Kr. coins



▶ Toss each coin 100 times, and record the number of times it comes up *heads* or *tails*

- Go to a bank and get 1000 5Kr. coins



- Toss each coin 100 times, and record the number of times it comes up *heads* or *tails*
- Use the Binomial test of significance at level 5% to test for each coin the hypothesis that

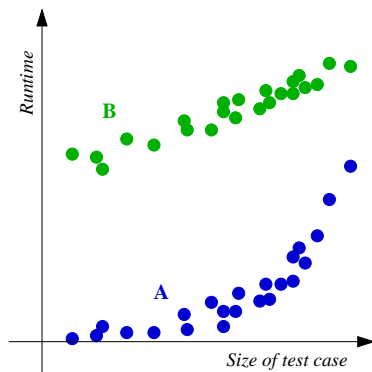$$P(heads) = P(tails) = 1/2$$

▶ Go to a bank and get 1000 5Kr. coins



▶ Toss each coin 100 times, and record the number of times it comes up *heads* or *tails*
▶ Use the Binomial test of significance at level 5% to test for each coin the hypothesis that

$$P(heads) = P(tails) = 1/2$$

▶ If all coins are fair, then the hypothesis will be rejected in approximately 5% = 50 of the experiments

- Go to a bank and get 1000 5Kr. coins



- Toss each coin 100 times, and record the number of times it comes up *heads* or *tails*
- Use the Binomial test of significance at level 5% to test for each coin the hypothesis that

$$P(heads) = P(tails) = 1/2$$

- If all coins are fair, then the hypothesis will be rejected in approximately 5% = 50 of the experiments
- Bring these 50 coins back to the bank ...

- Go to a bank and get 1000 5Kr. coins



- Toss each coin 100 times, and record the number of times it comes up *heads* or *tails*
- Use the Binomial test of significance at level 5% to test for each coin the hypothesis that

$$P(heads) = P(tails) = 1/2$$

- If all coins are fair, then the hypothesis will be rejected in approximately 5% = 50 of the experiments
- Bring these 50 coins back to the bank ...

**Multiple Testing:** when a large number of tests are performed, then the outcomes will contain some *false rejections* of the null Hypothesis. Any given (true) hypothesis can be (eventually) rejected by testing it often enough (if tests are independent).
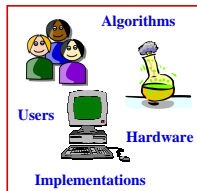
**Complexity analysis: B** (linear) is better than **A** (quadratic?).

**Testing:** Average runtime of **A** is less than of **B** (statistically significant).

► Which statement is more relevant depends on whether test cases are representative for size of future cases
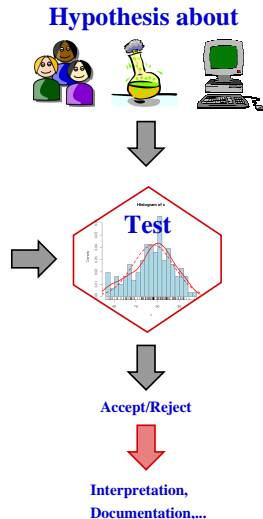
"Experimental Design"

Hypothesis about

Algorithms

Experiment

Users

Hardware

Observations

Implementations

Data

Test

Accept/Reject

Interpretation,
Documentation,...

Admission statistics from Berkeley University, 1973:

| Applicant Gender | Admitted |
|------------------|----------|
| Male             | 44%      |
| Female           | 35%      |

Is there gender discrimination against female applicants?

Admission statistics from Berkeley University, 1973:

| Applicant Gender | Admitted |
| --- | --- |
| Male | 44% |
| Female | 35% |

Is there gender discrimination against female applicants?

By departments:

| Department | Male admitted | Female admitted |
| --- | --- | --- |
| A | 62% | 82% |
| B | 63% | 68% |
| C | 37% | 34% |
| D | 33% | 35% |
| E | 28% | 24% |
| F | 6% | 7% |

▶ There seems to be gender discrimination at the university as a whole, but not at any single department ...

Adding distribution of applicants to the picture (hypothetical numbers):

| | | | % of applicants applying to department | |
|---|---|---|---|---|
| Department | Male admitted | Female admitted | male | female |
| A | 62% | 82% | 23% | 10% |
| B | 63% | 68% | 20% | 10% |
| C | 37% | 34% | 22% | 15% |
| D | 33% | 35% | 20% | 25% |
| E | 28% | 24% | 10% | 20% |
| F | 6% | 7% | 5% | 20% |

☞ Male applicants have a higher admission rate overall, but a lower one in (almost) all departments, because female students apply for the programs that are harder to get into.

When collecting data for evaluating the differences between setups A and B, make sure that test cases do not differ systematically with regard to factors other than defined by A and B (**confounding factors**).

► Run algorithms A and B on the same hardware
► Randomly assign test cases or test users to setups A and B
► Paired samples (within-subject design): apply both setups to each test case/user
  ► Easy when test cases are computational problems
  ► For test users: randomize the order in which setups A,B are shown to the users
► **Stratified Sampling:** Include in the sample used for statistical analysis equal numbers of test cases for every configuration of the confounding factor(s).

Analyzing data that is not collected in **controlled experiments**: try to identify all factors that could potentially have influenced the observations.

Test Cases: Data

For some types of computational problems there exist repositories of benchmark problems:

Satisfiability problems – SATLIB:

```
http://www.satlib.org/
```

Machine Learning data sets – UCI ML Repository:

```
http://archive.ics.uci.edu/ml/
```

Optimization:

```
http://plato.la.asu.edu/bench.html
```

Plus:

- ▶ Availability
- ▶ Reproducibility
- ▶ Historical results available

Minus:

- ▶ Older and less interesting datasets
- ▶ Development can be driven by optimization on benchmarks
- ▶ "Best published" results may be difficult to improve upon (and may be too good to be true)
- ▶ Invites multiple testing
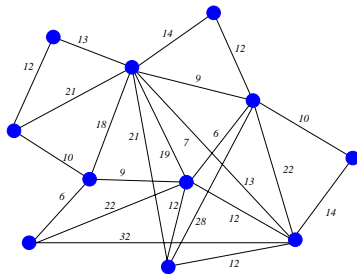
Create a generator for (random) data cases.

Plus:

- ► Easy to generate in large quantities
- ► Some data characteristics can be precisely tuned (size/complexity of problem instances)
- ► Can be shared with others (no proprietary or privacy concerns)

Minus:

- ► Difficult to get the structure of "real" data in random data (TSP example)
- ► Easy to manipulate to obtain desired outcomes

Input case for Traveling Salesman Problem:



Creating random instances for the Traveling Salesman Problem (cf. Rardin,Uzsoy):

- ▶ 1st approach: create distance matrix with random entries. Problems: distances will not satisfy triangle inequality, and most tours will have about equal length
- ▶ 2nd approach: randomly distribute points in the plane, and use the Euclidean distance between points. Problems: how to distribute points, so that resulting graph resembles typical TSP problem?

The data used as experimental test cases should be representative for the inputs of the intended application domain.

▶ For general algorithmic problems (e.g. TSP): test cases should cover wide spectrum of possible applications

▶ For concrete application scenarios (e.g. design of web shop homepage): test cases must be closely representative for use cases (best: live data collection in deployed system)

Temptation: use test cases that make your setup look good!

The data used as experimental test cases should be representative for the inputs of the intended application domain.

► For general algorithmic problems (e.g. TSP): test cases should cover wide spectrum of possible applications

► For concrete application scenarios (e.g. design of web shop homepage): test cases must be closely representative for use cases (best: live data collection in deployed system)

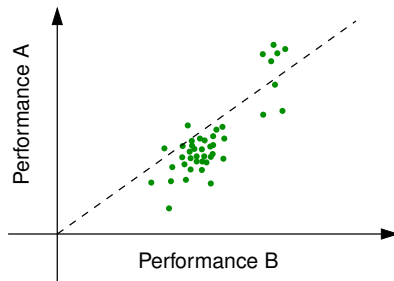Temptation: use test cases that make your setup look good!

**Example:**

► A first experiments yields these results

The data used as experimental test cases should be representative for the inputs of the intended application domain.

- ► For general algorithmic problems (e.g. TSP): test cases should cover wide spectrum of possible applications
- ► For concrete application scenarios (e.g. design of web shop homepage): test cases must be closely representative for use cases (best: live data collection in deployed system)

Temptation: use test cases that make your setup look good!

**Example:**

- ► A first experiments yields these results
- ► Temptation: conduct another experiment with test cases resembling those in subset G.
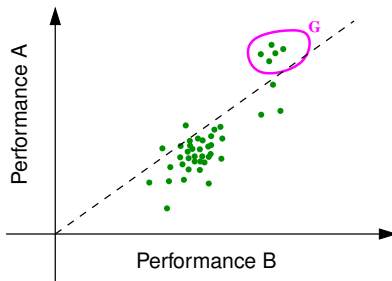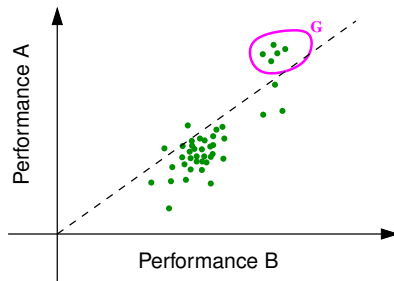
The data used as experimental test cases should be representative for the inputs of the intended application domain.

► For general algorithmic problems (e.g. TSP): test cases should cover wide spectrum of possible applications
► For concrete application scenarios (e.g. design of web shop homepage): test cases must be closely representative for use cases (best: live data collection in deployed system)

Temptation: use test cases that make your setup look good!

**Example:**

► A first experiments yields these results
► Temptation: conduct another experiment with test cases resembling those in subset G.
► This need not be unethical, if all results are reported (and/or special nature of cases in G is acknowledged)

Experimental Procedures

**Experiments: Evaluation vs. Development**

A setup can be controlled by parameters:

- ▶ Parameters controlling time/space/accuracy tradeoff of optimization algorithm
- ▶ Parameters controlling layout of web shop homepage
- ▶ . . .
- ▶ Parameters representing more high-level design choices

Experiments often interleave:

- ▶ Performance evaluation of setup
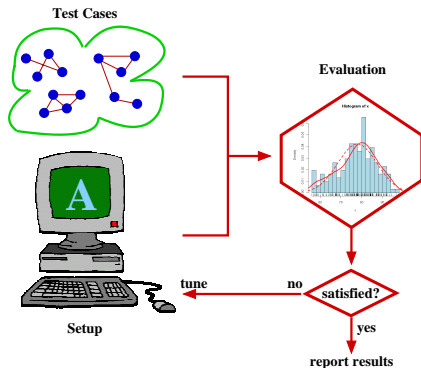- ▶ Tuning of parameters to optimize performance

A setup may contain a large number of tunable parameters:

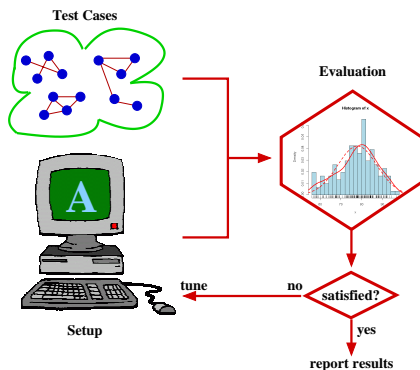$$A \in \mathbb{R}, \ B \in \mathbb{Z}, \ \gamma \in [0, 1], \ t \in \{0, 1\}, \ \ldots$$

Usually impossible to try all combinations of parameter values! Simple procedure:

- Using some "best guess" default setting for other parameters, find the best value of the first parameter
- Using best value of first parameter, default settings for others, find best value of second parameter
- ...

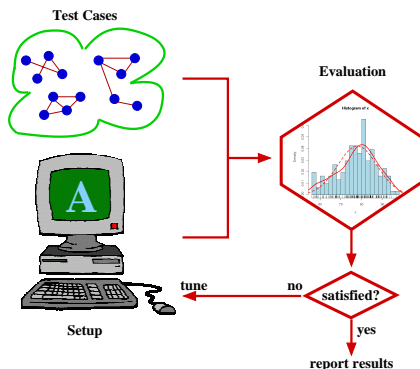Not guaranteed to find best setting, but good heuristic.

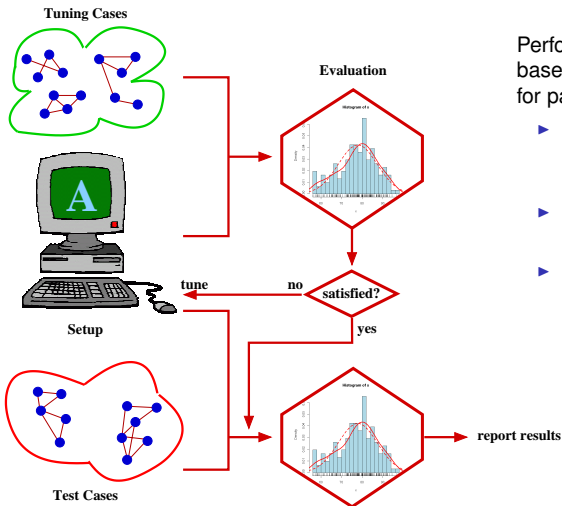Any problems with this process?

Any problems with this process?

**Problem 1: Overfitting.** The setup is optimized to perform best on the test cases. Performance on future cases may be worse than expected.

Any problems with this process?

**Problem 1: Overfitting.** The setup is optimized to perform best on the test cases. Performance on future cases may be worse than expected.

**Problem 2: Multiple Testing.** If evaluation contains statistical tests, then spurious rejections of null hypothesis may occur (mitigated by the fact that tests for different parameter settings are not independent – not necessarily 5% false rejections!).

**Tuning Cases**

**Evaluation**

**Setup**

**tune**     **no**    **satisfied?**

**yes**

**Test Cases**

**report results**

Perform final performance evaluation based on test cases that were not used for parameter tuning

► Avoids overfitting: test cases in final evaluation representative of future cases

► Final significance analysis based on a single test

► . . . but, in reality: almost impossible to completely eliminate the "tuning on the test cases" fallacy

Make sure:

- ▶ 'Our" system **A** and competitor **B** tested on the same hardware
- ▶ No other processes running on machine during testing

Even then, one often finds:

- ▶ "Our" **A** finely tuned
- ▶ Competitor **B** run in default configuration
- ▶ Expertise lacking to get best performance out of **B**
- ▶ ⤳ comparison biased in favor of our system

Documentation

One often reads (here: in Machine Learning papers):

*. . . we conducted experiments on 10 datasets from the UCI repository . . .*

Suspicion: experiments were conducted on more datasets, but these 10 gave the best results.

Ideally:

▶ All experiments and their results are reported, including negative or inconclusive ones

Second best option:

▶ Admit that the selected test cases are well suited for proposed setup A. Explain why.

Other researchers should be able to independently reproduce experimental results:

- ▶ In publication: precisely describe experimental procedure
    - ▶ Test cases
    - ▶ Preprocessing (data cleaning, outlier removal, . . . )
    - ▶ Tune/Test approach
    - ▶ Hardware platform

- ▶ Make all relevant information publicly available (supplementary material published on the web)
    - ▶ Code for your implementation/algorithm
    - ▶ Test data (if not already publicly available)
    - ▶ Questionnaires used in user studies
    - ▶ . . .