



合肥學院
HEFEI UNIVERSITY



Evolutionary Computation

Comparing Optimization Algorithms

Thomas Weise · 汤卫思

tweise@hfu.edu.cn · <http://iao.hfu.edu.cn/5>

Institute of Applied Optimization (IAO)
School of Artificial Intelligence and Big Data
Hefei University
Hefei, Anhui, China

应用优化研究所
人工智能与大数据学院
合肥学院
中国安徽省合肥市

Outline

1. Introduction
2. Views on Performance and Time
3. Statistical Measures
4. Statistical Comparisons
5. Testing is Not Enough
6. Other Stuff
7. Summary



Introduction



Introduction

- There are many optimization algorithms.

Introduction

- There are many optimization algorithms.
- For solving an optimization problem, we want to use the algorithm most suitable for it.

Introduction

- There are many optimization algorithms.
- For solving an optimization problem, we want to use the algorithm most suitable for it.
- What does this mean?

Introduction

- There are many optimization algorithms.
- For solving an optimization problem, we want to use the algorithm most suitable for it.
- What does this mean?
- And how do we find this algorithm?

Introduction

- There are many optimization algorithms.
- For solving an optimization problem, we want to use the algorithm most suitable for it.
- What does this mean?
- And how do we find this algorithm?
- Hopefully this lesson will answer these questions.

Introduction

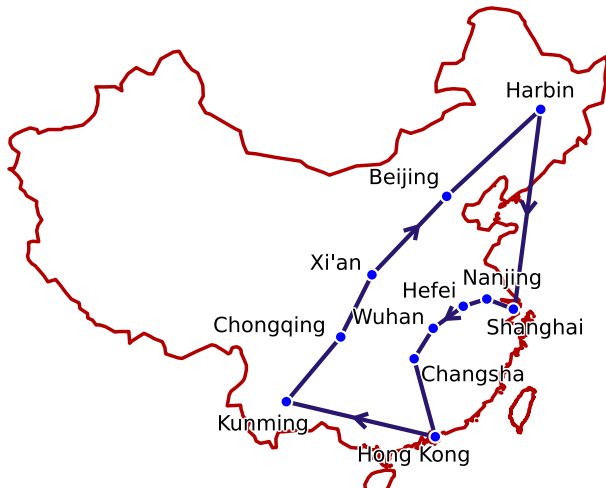
- There are many optimization algorithms.
- For solving an optimization problem, we want to use the algorithm most suitable for it.
- What does this mean?
- And how do we find this algorithm?
- Hopefully this lesson will answer these questions.
- As a complement to this lesson, I suggest the report *“Benchmarking in Optimization: Best Practice and Open Issues”*³ on Arxiv.

Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.

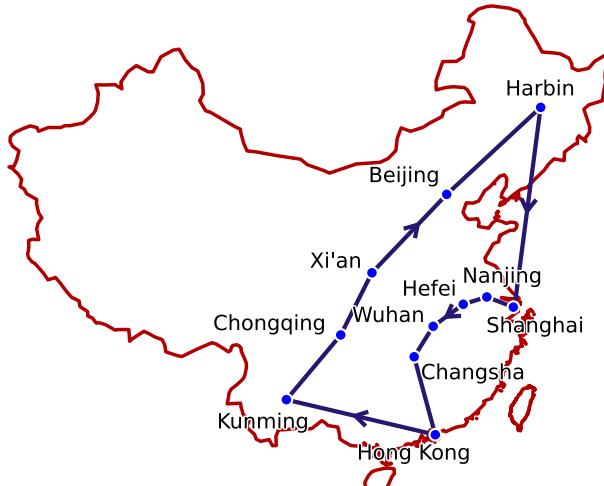
Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).



Exact vs. Heuristic Algorithms

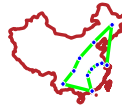
- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Clearly, there is (at least) one shortest tour.



Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Clearly, there is (at least) one shortest tour.

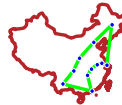
getting the optimal solution
for a TSP



Exact vs. Heuristic Algorithms

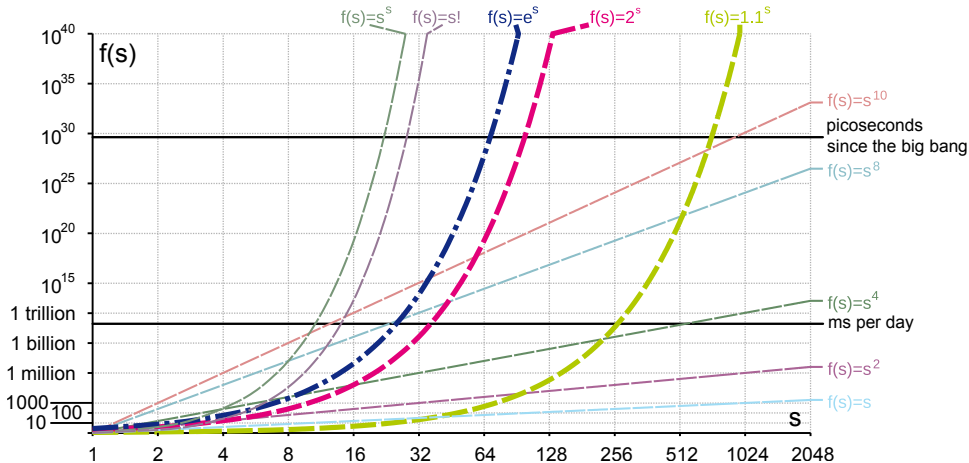
- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Clearly, there is (at least) one shortest tour.
 - Theory proofs that the time to find this tour may grow exponentially with the number of cities we want to visit in the worst case.⁴⁻⁸

getting the optimal solution
for a TSP



Exact vs. Heuristic Algorithms

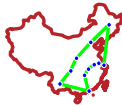
- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).



Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Finding the best tour (what exact algorithms do) may take too long.

getting the optimal solution
for a TSP may take too long

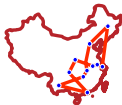


consumed runtime:

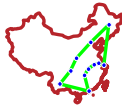
very much / too (?) long

Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Finding the best tour (what exact algorithms do) may take too long.
 - But we can find just **some** tour very quickly.



getting the optimal solution
for a TSP may take too long



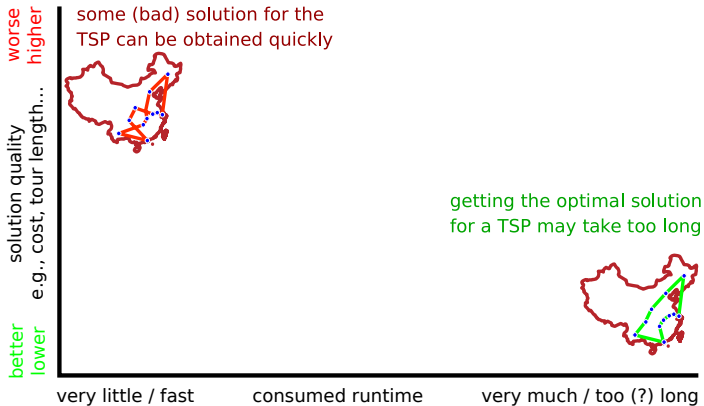
very little / fast

consumed runtime

very much / too (?) long

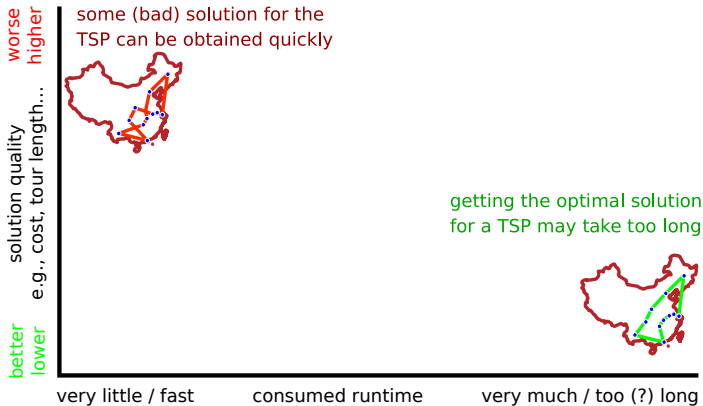
Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - But we can find just **some** tour very quickly.
 - Of course the quality of that tour will be lower.



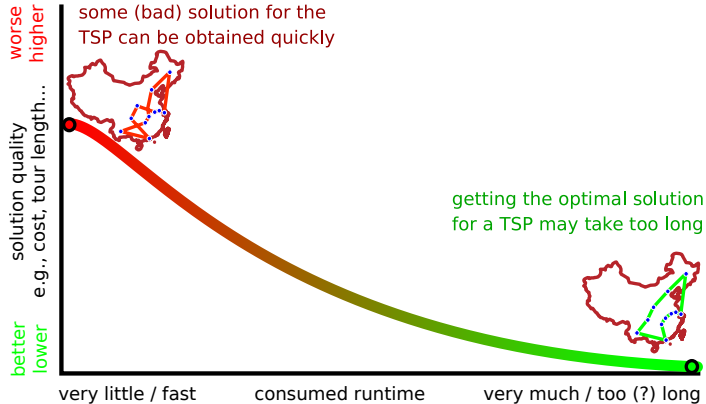
Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Of course the quality of that tour will be lower: the tour will be longer than the best one.



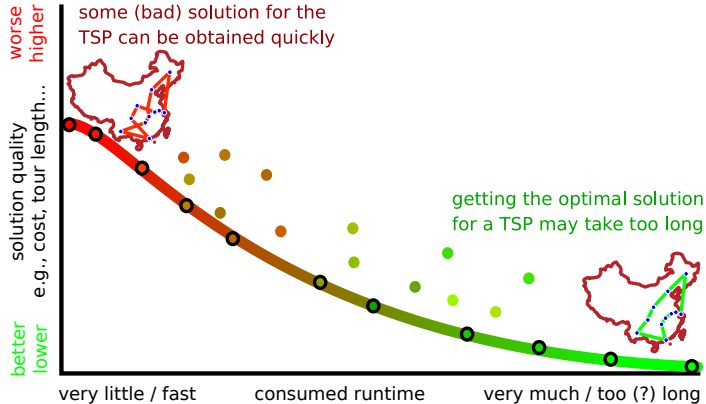
Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Of course the quality of that tour will be lower.
 - Is there something inbetween?



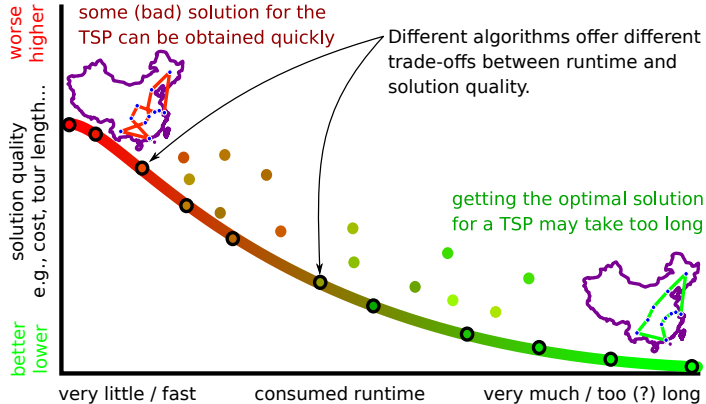
Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Is there something inbetween?
 - (Meta-)Heuristic optimization algorithms try to find solutions which are as good as possible as fast as possible.



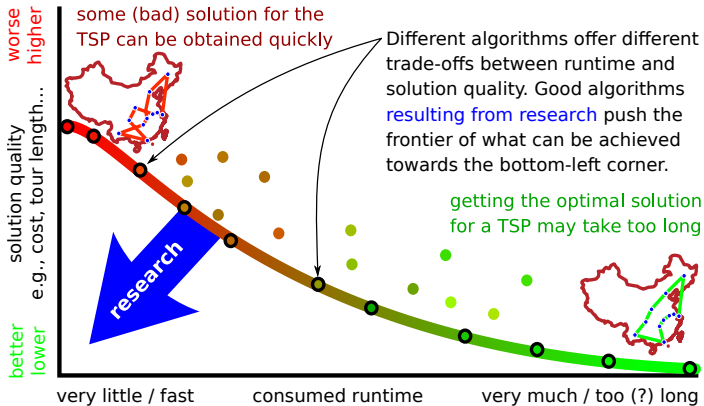
Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - (Meta-)Heuristic optimization algorithms try to find solutions which are as good as possible as fast as possible.
 - **Optimization often means to make a trade-off between solution quality and runtime.**



Exact vs. Heuristic Algorithms

- In optimization, there exist **exact** and **heuristic** algorithms.
- Let's look at the classical Traveling Salesperson Problem (TSP).
 - Optimization often means to make a trade-off between solution quality and runtime.



Views on Performance and Time



Views on Performance

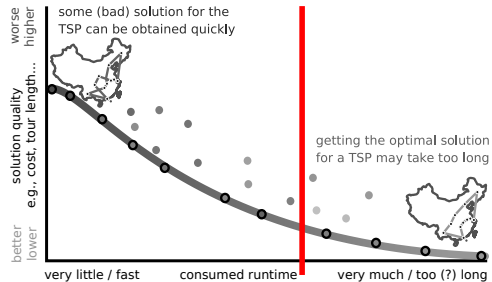
- Runtime and solution quality in optimization are intertwined and should never be considered separately.

Views on Performance

- Runtime and solution quality in optimization are intertwined and should never be considered separately.
- Views⁹⁻¹²

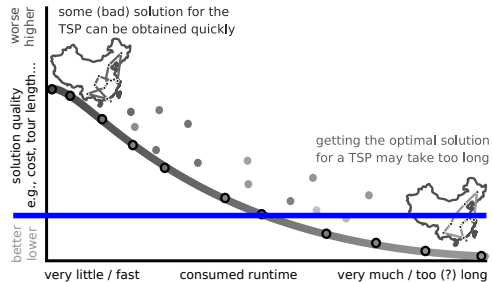
Views on Performance

- Runtime and solution quality in optimization are intertwined and should never be considered separately.
- Views⁹⁻¹²:
 1. Solution quality reached after a certain runtime



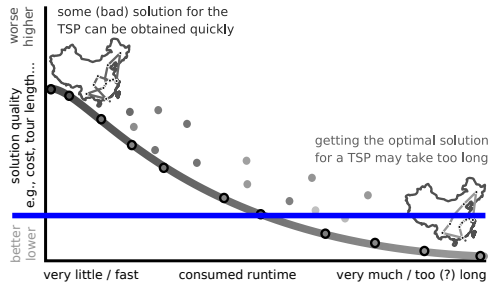
Views on Performance

- Runtime and solution quality in optimization are intertwined and should never be considered separately.
- Two view⁹⁻¹²:
 1. Solution quality reached after a certain runtime
 2. Runtime to reach a certain solution quality



Views on Performance

- Runtime and solution quality in optimization are intertwined and should never be considered separately.
- Two view⁹⁻¹²:
 1. Solution quality reached after a certain **runtime**
 2. **Runtime** to reach a certain solution quality



What is Runtime?

- What actually is runtime?

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm in ms.

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the **algorithm implementation** in ms.

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- Advantages

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- Advantages:
 - Results in many works reported in this format

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- Advantages:
 - Results in many works reported in this format
 - A quantity that makes physical sense

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- Advantages:
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- Advantages:
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation
- Disadvantages

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- **Advantages:**
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation
- **Disadvantages:**
 - Strongly machine dependent and inherently incomparable over different machines

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- **Advantages:**
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation
- **Disadvantages:**
 - Strongly machine dependent and inherently incomparable over different machines
 - Measurements are only valuable for a few years

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- **Advantages:**
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation
- **Disadvantages:**
 - Strongly machine dependent and inherently incomparable over different machines
 - Measurements are only valuable for a few years
 - Can be biased by “outside effects,” e.g., OS, scheduling, other processes, I/O, swapping, ...

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- **Advantages:**
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation
- **Disadvantages:**
 - Strongly machine dependent and inherently incomparable over different machines
 - Measurements are only valuable for a few years
 - Can be biased by “outside effects,” e.g., OS, scheduling, other processes, I/O, swapping, ...
- Hardware, software, OS, programming language, etc. all have nothing to do with the **optimization algorithm** itself and are relevant only in a specific application...

Clock Time as Absolute Runtime

Measure the (absolute) consumed runtime of the algorithm implementation in ms.

- **Advantages:**
 - Results in many works reported in this format
 - A quantity that makes physical sense
 - Includes all “hidden complexities” of an algorithm implementation
- **Disadvantages:**
 - Strongly machine dependent and inherently incomparable over different machines
 - Measurements are only valuable for a few years
 - Can be biased by “outside effects,” e.g., OS, scheduling, other processes, I/O, swapping, ...
- Hardware, software, OS, programming language, etc. all have nothing to do with the **optimization algorithm** itself and are relevant only in a specific application...
- ...for **research** they may be less interesting, while for a **specific application** they do matter.

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- Advantages

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- Advantages:
 - Results in many works reported in this format (or FEs can be deduced)

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- Advantages:
 - Results in many works reported in this format (or FEs can be deduced)
 - Machine-independent, theory-related measure

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- Advantages:
 - Results in many works reported in this format (or FEs can be deduced)
 - Machine-independent, theory-related measure
 - Cannot be influenced by “outside effects”

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**
 - Results in many works reported in this format (or FEs can be deduced)
 - **Machine-independent, theory-related measure**
 - Cannot be influenced by “outside effects”
 - In many optimization problems, computing the objective value is the most time consuming task

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- Advantages:
 - Results in many works reported in this format (or FEs can be deduced)
 - Machine-independent, theory-related measure
 - Cannot be influenced by “outside effects”
 - In many optimization problems, computing the objective value is the most time consuming task
- Disadvantages

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**
 - Results in many works reported in this format (or FEs can be deduced)
 - **Machine-independent, theory-related measure**
 - Cannot be influenced by “outside effects”
 - In many optimization problems, computing the objective value is the most time consuming task
- **Disadvantages:**
 - No clear relationship to real runtime

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**
 - Results in many works reported in this format (or FEs can be deduced)
 - **Machine-independent, theory-related measure**
 - Cannot be influenced by “outside effects”
 - In many optimization problems, computing the objective value is the most time consuming task
- **Disadvantages:**
 - No clear relationship to real runtime
 - Does not contain “hidden complexities” of algorithm

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**
 - Results in many works reported in this format (or FEs can be deduced)
 - **Machine-independent, theory-related measure**
 - Cannot be influenced by “outside effects”
 - In many optimization problems, computing the objective value is the most time consuming task
- **Disadvantages:**
 - No clear relationship to real runtime
 - Does not contain “hidden complexities” of algorithm
 - 1 FE: very different costs in different situations!¹³

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**

- Results in many works reported in this format (or FEs can be deduced)
- **Machine-independent, theory-related measure**
- Cannot be influenced by “outside effects”
- In many optimization problems, computing the objective value is the most time consuming task

- **Disadvantages:**

- No clear relationship to real runtime
- Does not contain “hidden complexities” of algorithm
- 1 FE: very different costs in different situations!¹³
 - When applying a local search that swaps two cities in each move to the Traveling Salesperson Problem (TSP), one FE can be done in $\mathcal{O}(1)$.

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**

- Results in many works reported in this format (or FEs can be deduced)
- **Machine-independent, theory-related measure**
- Cannot be influenced by “outside effects”
- In many optimization problems, computing the objective value is the most time consuming task

- **Disadvantages:**

- No clear relationship to real runtime
- Does not contain “hidden complexities” of algorithm
- 1 FE: very different costs in different situations!¹³
 - When applying a local search that swaps two cities in each move to the Traveling Salesperson Problem (TSP), one FE can be done in $\mathcal{O}(1)$.
 - When applying Ant Colony Optimization instead, each FE takes $\mathcal{O}(n^2)$.

Function Evaluations: FEs

Measure (count) the number of fully constructed and tested candidate solutions.

- **Advantages:**

- Results in many works reported in this format (or FEs can be deduced)
- **Machine-independent, theory-related measure**
- Cannot be influenced by “outside effects”
- In many optimization problems, computing the objective value is the most time consuming task

- **Disadvantages:**

- No clear relationship to real runtime
- Does not contain “hidden complexities” of algorithm
- 1 FE: very different costs in different situations!¹³
 - When applying a local search that swaps two cities in each move to the Traveling Salesperson Problem (TSP), one FE can be done in $\mathcal{O}(1)$.
 - When applying Ant Colony Optimization instead, each FE takes $\mathcal{O}(n^2)$.
- Relevant for comparing algorithms, but not so much for the practical application or comparing implementations.

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs.

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.
- Often, the mapping between generations and FEs is not clear

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.
- Often, the mapping between generations and FEs is not clear, for example
 - Do you evaluate offspring solutions that are identical to their parents?

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.
- Often, the mapping between generations and FEs is not clear, for example
 - Do you evaluate offspring solutions that are identical to their parents?
 - Is a local search involved that refines some or all solutions in the population?

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.
- Often, the mapping between generations and FEs is not clear, for example
 - Do you evaluate offspring solutions that are identical to their parents?
 - Is a local search involved that refines some or all solutions in the population?
 - In a $(\mu + \lambda)$ -EA, is the first population of size $\mu + \lambda$, λ , or μ ?

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.
- Often, the mapping between generations and FEs is not clear, for example
 - Do you evaluate offspring solutions that are identical to their parents?
 - Is a local search involved that refines some or all solutions in the population?
 - In a $(\mu + \lambda)$ -EA, is the first population of size $\mu + \lambda$, λ , or μ ?
 - What if the population size changes adaptively?

Do not count generations

- Do not use the number of *generations* in your EA as time measure! Instead count the FEs, because:
- The “number of generations” are not really comparable for different population sizes or with algorithms that do not use populations.
- Often, the mapping between generations and FEs is not clear, for example
 - Do you evaluate offspring solutions that are identical to their parents?
 - Is a local search involved that refines some or all solutions in the population?
 - In a $(\mu + \lambda)$ -EA, is the first population of size $\mu + \lambda$, λ , or μ ?
 - What if the population size changes adaptively?
- I suggest to prefer FEs over generations if you want to count algorithm steps.

Runtime

- I suggest to always measure both the consumed FEs and the runtime in milliseconds.

Runtime

- I suggest to always measure both the consumed FEs and the runtime in milliseconds.
- Anyway, with what we have learned, we can rewrite the two views by choosing a time measure^{9 11}

Runtime

- I suggest to always measure both the consumed FEs and the runtime in milliseconds.
- Anyway, with what we have learned, we can rewrite the two views by choosing a time measure^{9 11}, e.g.:
 1. Solution quality reached after a certain **number of FEs**

Runtime

- I suggest to always measure both the consumed FEs and the runtime in milliseconds.
- Anyway, with what we have learned, we can rewrite the two views by choosing a time measure^{9 11}, e.g.:
 1. Solution quality reached after a certain **number of FEs**
 2. **Milliseconds** needed to reach a certain solution quality

Solution Quality

- Common measure of solution quality: Objective function value of best solution discovered.

Solution Quality

- Common measure of solution quality: Objective function value of best solution discovered.
- Rewrite the two views^{9 11}

Solution Quality

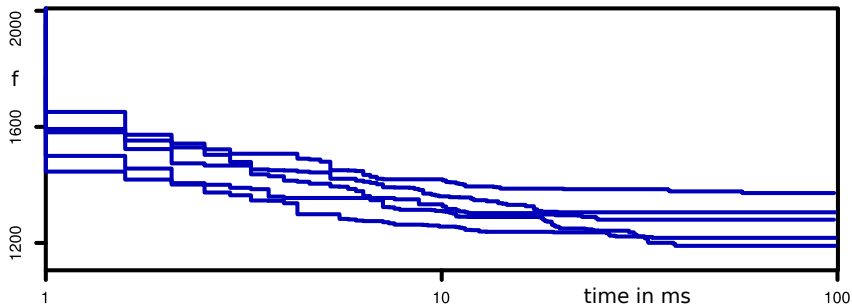
- Common measure of solution quality: Objective function value of best solution discovered.
- Rewrite the two views^{9 11}:
 1. **Best objective function value** reached after a certain number of milliseconds

Solution Quality

- Common measure of solution quality: Objective function value of best solution discovered.
- Rewrite the two views^{9 11}:
 1. **Best objective function value** reached after a certain number of milliseconds
 2. Number FEs needed to reach a certain **objective function value**

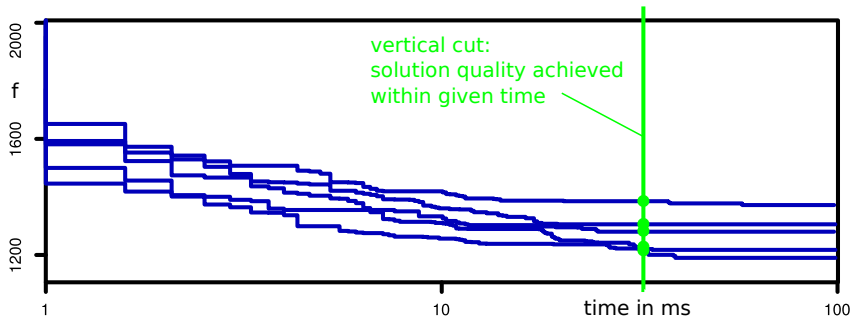
Views on Performance

- Which one is the “better” view on performance?



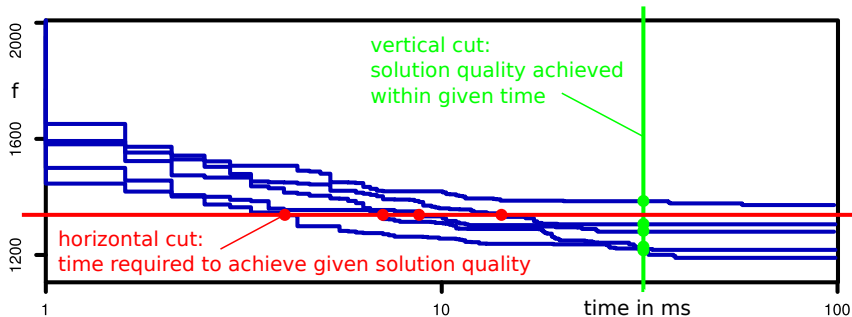
Views on Performance

- Which one is the “better” view on performance?
 1. Best objective function value reached after a certain number of FEs



Views on Performance

- Which one is the “better” view on performance?
 1. Best objective function value reached after a certain number of FEs
 2. **Number FEs needed to reach a certain objective function value**



Views on Performance

- Which one is the “better” view on performance?
 1. Best objective function value reached after a certain number of FEs
 2. Number FEs needed to reach a certain objective function value
- This question is still debated in research. . .

Which view is better?

- Number FEs needed to reach a certain objective function value
- Preferred by, e.g., the BBOB/COCO benchmark suite⁹

Which view is better?

- Number FEs needed to reach a certain objective function value
- Preferred by, e.g., the BBOB/COCO benchmark suite⁹:
 - Measures the time needed to reach a target function value allows meaningful statements such as “Algorithm A is two/ten/hundred times faster than Algorithm B in solving this problem.”

Which view is better?

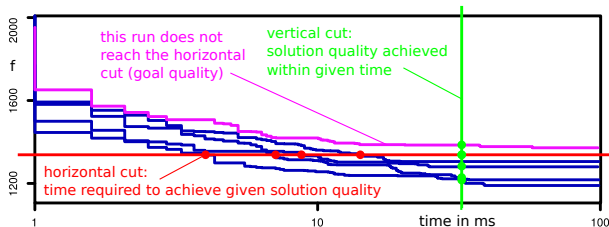
- Number FEs needed to reach a certain objective function value
- Preferred by, e.g., the BBOB/COCO benchmark suite⁹:
 - Measures the time needed to reach a target function value allows meaningful statements such as “Algorithm A is two/ten/hundred times faster than Algorithm B in solving this problem.”
 - However, there is no interpretable meaning to the fact that Algorithm A reaches a function value that is two/ten/hundred times smaller than the one reached by Algorithm B .

Which view is better?

- Number FEs needed to reach a certain objective function value
- Preferred by, e.g., the BBOB/COCO benchmark suite⁹:
 - Measures the time needed to reach a target function value allows meaningful statements such as “Algorithm A is two/ten/hundred times faster than Algorithm B in solving this problem.”
 - However, there is no interpretable meaning to the fact that Algorithm A reaches a function value that is two/ten/hundred times smaller than the one reached by Algorithm B .
 - “Benchmarking Theory Perspective”

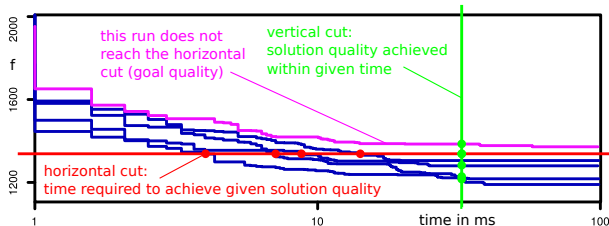
Which view is better?

- Number FEs needed to reach a certain objective function value
- Preferred by, e.g., the BBOB/COCO benchmark suite⁹:
 - Measures the time needed to reach a target function value allows meaningful statements such as “Algorithm A is two/ten/hundred times faster than Algorithm B in solving this problem.”
 - However, there is no interpretable meaning to the fact that Algorithm A reaches a function value that is two/ten/hundred times smaller than the one reached by Algorithm B .
 - “Benchmarking Theory Perspective”
- Sometimes problematic: What if one run does not reach the goal quality?



Which view is better?

- Number FEs needed to reach a certain objective function value
- Preferred by, e.g., the BBOB/COCO benchmark suite⁹:
 - Measures the time needed to reach a target function value allows meaningful statements such as “Algorithm A is two/ten/hundred times faster than Algorithm B in solving this problem.”
 - However, there is no interpretable meaning to the fact that Algorithm A reaches a function value that is two/ten/hundred times smaller than the one reached by Algorithm B .
 - “Benchmarking Theory Perspective”
- Sometimes problematic: What if one run does not reach the goal quality?
- Then, alternative measures need to be computed, such as the ERT^{14 15} or PAR2 and PAR10^{16 17}.



Which view is better?

- Best objective function value reached after a certain number of FEs

Which view is better?

- Best objective function value reached after a certain number of FEs
- Preferred by many benchmark suites such as¹⁸.

Which view is better?

- Best objective function value reached after a certain number of FEs
- Preferred by many benchmark suites such as¹⁸.
- Practice Perspective: Best results achievable with given time budget wins.

Which view is better?

- Best objective function value reached after a certain number of FEs
- Preferred by many benchmark suites such as¹⁸.
- Practice Perspective: Best results achievable with given time budget wins.
- This perspective maybe less suitable for scientific benchmarking, but surely is useful in practice.

Which view is better?

- Best objective function value reached after a certain number of FEs
- Preferred by many benchmark suites such as¹⁸.
- Practice Perspective: Best results achievable with given time budget wins.
- This perspective maybe less suitable for scientific benchmarking, but surely is useful in practice.
- “How good is the tour for the TSP that we can find in 5 minutes with our algorithm?”

Which view is better?

- Best objective function value reached after a certain number of FEs
- Preferred by many benchmark suites such as¹⁸.
- Practice Perspective: Best results achievable with given time budget wins.
- This perspective maybe less suitable for scientific benchmarking, but surely is useful in practice.
- “How good is the tour for the TSP that we can find in 5 minutes with our algorithm?”
- Always well-defined, because vertical cuts can always be reached.

Views on Performance

- No official consensus on which view is “better.”

Views on Performance

- No official consensus on which view is “better.”
- This also strongly depends on the situation.

Views on Performance

- No official consensus on which view is “better.”
- This also strongly depends on the situation.
- If we can actually always solve the problem to a “natural” goal quality (e.g., to optimality), then we should prefer the horizontal cut (time-to-target) method.

Views on Performance

- No official consensus on which view is “better.”
- This also strongly depends on the situation.
- If we can actually always solve the problem to a “natural” goal quality (e.g., to optimality), then we should prefer the horizontal cut (time-to-target) method.
- If we have clear application requirements specifying a fixed budget, then we should prefer the fixed-budget approach.

Views on Performance

- No official consensus on which view is “better.”
- This also strongly depends on the situation.
- If we can actually always solve the problem to a “natural” goal quality (e.g., to optimality), then we should prefer the horizontal cut (time-to-target) method.
- If we have clear application requirements specifying a fixed budget, then we should prefer the fixed-budget approach.
- Otherwise, the best approach may be: Evaluate algorithm according to both methods.¹¹⁻¹³

Views on Performance

- No official consensus on which view is “better.”
- This also strongly depends on the situation.
- If we can actually always solve the problem to a “natural” goal quality (e.g., to optimality), then we should prefer the horizontal cut (time-to-target) method.
- If we have clear application requirements specifying a fixed budget, then we should prefer the fixed-budget approach.
- Otherwise, the best approach may be: Evaluate algorithm according to both methods.^{11–13}
- Maybe cast a net of several horizontal and vertical cuts, to get a better picture. . .

Determining Target Values

- How to determine the right maximum FEs or target function values?

Determining Target Values

- How to determine the right maximum FEs or target function values?
 1. from the constraints of a practical application

Determining Target Values

- How to determine the right maximum FEs or target function values?
 1. from the constraints of a practical application
 2. from studies in literature regarding similar or the same problem

Determining Target Values

- How to determine the right maximum FEs or target function values?
 1. from the constraints of a practical application
 2. from studies in literature regarding similar or the same problem
 3. from simple or well-known algorithms

Determining Target Values

- How to determine the right maximum FEs or target function values?
 1. from the constraints of a practical application
 2. from studies in literature regarding similar or the same problem
 3. from simple or well-known algorithms
 4. from experience

Determining Target Values

- How to determine the right maximum FEs or target function values?
 1. from the constraints of a practical application
 2. from studies in literature regarding similar or the same problem
 3. from simple or well-known algorithms
 4. from experience
 5. from prior, small-scale experiments

Determining Target Values

- How to determine the right maximum FEs or target function values?
 1. from the constraints of a practical application
 2. from studies in literature regarding similar or the same problem
 3. from simple or well-known algorithms
 4. from experience
 5. from prior, small-scale experiments
 6. based on known results or well-accepted bounds

Statistical Measures



Problem Instances and Randomized Algorithms

- For each **optimization problem** (like the TSP) there are several **instances** (e.g., different sets of cities that need to be visited).

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.
- Special situation: Randomized Algorithms

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.
- Special situation: Randomized Algorithms:
 - Performance values cannot be given absolute!

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.
- Special situation: Randomized Algorithms:
 - Performance values cannot be given absolute!
 - 1 run = 1 application of an optimization algorithm to a problem, runs are independent from all prior runs.

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.
- Special situation: Randomized Algorithms:
 - Performance values cannot be given absolute!
 - 1 run = 1 application of an optimization algorithm to a problem, runs are independent from all prior runs.
 - Results can be different for each run!

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.
- Special situation: Randomized Algorithms:
 - Performance values cannot be given absolute!
 - 1 run = 1 application of an optimization algorithm to a problem, runs are independent from all prior runs.
 - Results can be different for each run!
 - Executing a randomized algorithm one time does not give reliable information.

Problem Instances and Randomized Algorithms

- For each optimization problem (like the TSP) there are several instances (e.g., different sets of cities that need to be visited).
 - Some instances will be easy, some will be hard.
 - We always must use multiple different problem instances to get reliable results.
 - Performance indicators need to be computed for each instance and also summarized over several instances.
- Special situation: Randomized Algorithms:
 - Performance values cannot be given absolute!
 - 1 run = 1 application of an optimization algorithm to a problem, runs are independent from all prior runs.
 - Results can be different for each run!
 - Executing a randomized algorithm one time does not give reliable information.
 - Statistical evaluation over sets of runs necessary.

Important Distinction

- Crucial Difference: **distribution** and **sample**

Important Distinction

- Crucial Difference: distribution and sample
- A **sample** is what we *measure*.

Important Distinction

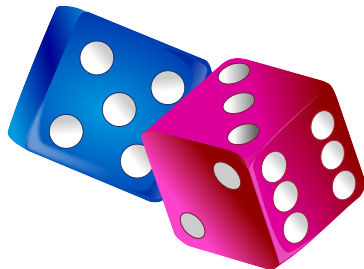
- Crucial Difference: distribution and sample
- A sample is what we *measure*.
- A **distribution** is the asymptotic result of the ideal process.

Important Distinction

- Crucial Difference: distribution and sample
- A sample is what we *measure*.
- A distribution is the asymptotic result of the ideal process.
- Statistical parameters of the distribution can be **estimated** from a sample.

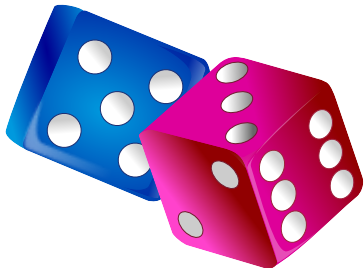
Important Distinction

- Crucial Difference: distribution and sample
- A sample is what we *measure*.
- A distribution is the asymptotic result of the ideal process.
- Statistical parameters of the distribution can be estimated from a sample.
- Example: Dice Throw



Important Distinction

- Crucial Difference: distribution and sample
- A sample is what we *measure*.
- A distribution is the asymptotic result of the ideal process.
- Statistical parameters of the distribution can be estimated from a sample.
- Example: Dice Throw
- How likely is it to roll a 1, 2, 3, 4, 5, or 6?



Important Distinction

# throws	number	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000



based on <http://www.freestockphotos.biz/stockphoto/16223>

Important Distinction

# throws	number	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000



Important Distinction

# throws	number	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000



Important Distinction

# throws	number	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000



Important Distinction

# throws	number	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000
5	3	0.2000	0.0000	0.2000	0.4000	0.2000	0.0000



Important Distinction

# throws	number	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000
5	3	0.2000	0.0000	0.2000	0.4000	0.2000	0.0000
6	3	0.1667	0.0000	0.3333	0.3333	0.1667	0.0000



Important Distinction

# throws	number	f(1)	f(2)	f(3)	f(4)	f(5)	f(6)
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000
5	3	0.2000	0.0000	0.2000	0.4000	0.2000	0.0000
6	3	0.1667	0.0000	0.3333	0.3333	0.1667	0.0000
7	2	0.1429	0.1429	0.2857	0.2857	0.1429	0.0000
8	1	0.2500	0.1250	0.2500	0.2500	0.1250	0.0000
9	4	0.2222	0.1111	0.2222	0.3333	0.1111	0.0000
10	2	0.2000	0.2000	0.2000	0.3000	0.1000	0.0000



Important Distinction

# throws	number	f(1)	f(2)	f(3)	f(4)	f(5)	f(6)
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000
5	3	0.2000	0.0000	0.2000	0.4000	0.2000	0.0000
6	3	0.1667	0.0000	0.3333	0.3333	0.1667	0.0000
7	2	0.1429	0.1429	0.2857	0.2857	0.1429	0.0000
8	1	0.2500	0.1250	0.2500	0.2500	0.1250	0.0000
9	4	0.2222	0.1111	0.2222	0.3333	0.1111	0.0000
10	2	0.2000	0.2000	0.2000	0.3000	0.1000	0.0000
11	6	0.1818	0.1818	0.1818	0.2727	0.0909	0.0909
12	3	0.1667	0.1667	0.2500	0.2500	0.0833	0.0833



Important Distinction

# throws	number	f(1)	f(2)	f(3)	f(4)	f(5)	f(6)
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000
5	3	0.2000	0.0000	0.2000	0.4000	0.2000	0.0000
6	3	0.1667	0.0000	0.3333	0.3333	0.1667	0.0000
7	2	0.1429	0.1429	0.2857	0.2857	0.1429	0.0000
8	1	0.2500	0.1250	0.2500	0.2500	0.1250	0.0000
9	4	0.2222	0.1111	0.2222	0.3333	0.1111	0.0000
10	2	0.2000	0.2000	0.2000	0.3000	0.1000	0.0000
11	6	0.1818	0.1818	0.1818	0.2727	0.0909	0.0909
12	3	0.1667	0.1667	0.2500	0.2500	0.0833	0.0833
100	...	0.1900	0.2100	0.1500	0.1600	0.1200	0.1700
1'000	...	0.1700	0.1670	0.1620	0.1670	0.1570	0.1770
10'000	...	0.1682	0.1699	0.1680	0.1661	0.1655	0.1623
100'000	...	0.1671	0.1649	0.1664	0.1676	0.1668	0.1672
1'000'000	...	0.1673	0.1663	0.1662	0.1673	0.1666	0.1664



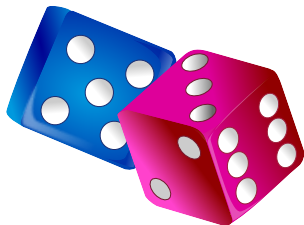
Important Distinction

# throws	number	f(1)	f(2)	f(3)	f(4)	f(5)	f(6)
1	5	0.0000	0.0000	0.0000	0.0000	1.0000	0.0000
2	4	0.0000	0.0000	0.0000	0.5000	0.5000	0.0000
3	1	0.3333	0.0000	0.0000	0.3333	0.3333	0.0000
4	4	0.2500	0.0000	0.0000	0.5000	0.2500	0.0000
5	3	0.2000	0.0000	0.2000	0.4000	0.2000	0.0000
6	3	0.1667	0.0000	0.3333	0.3333	0.1667	0.0000
7	2	0.1429	0.1429	0.2857	0.2857	0.1429	0.0000
8	1	0.2500	0.1250	0.2500	0.2500	0.1250	0.0000
9	4	0.2222	0.1111	0.2222	0.3333	0.1111	0.0000
10	2	0.2000	0.2000	0.2000	0.3000	0.1000	0.0000
11	6	0.1818	0.1818	0.1818	0.2727	0.0909	0.0909
12	3	0.1667	0.1667	0.2500	0.2500	0.0833	0.0833
100	...	0.1900	0.2100	0.1500	0.1600	0.1200	0.1700
1'000	...	0.1700	0.1670	0.1620	0.1670	0.1570	0.1770
10'000	...	0.1682	0.1699	0.1680	0.1661	0.1655	0.1623
100'000	...	0.1671	0.1649	0.1664	0.1676	0.1668	0.1672
1'000'000	...	0.1673	0.1663	0.1662	0.1673	0.1666	0.1664
10'000'000	...	0.1667	0.1667	0.1666	0.1668	0.1667	0.1665
100'000'000	...	0.1667	0.1666	0.1666	0.1667	0.1667	0.1667
1'000'000'000	...	0.1667	0.1667	0.1667	0.1667	0.1667	0.1667



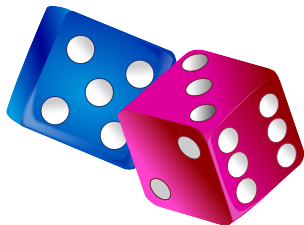
Important Distinction

- Crucial Difference: distribution and sample
- A sample is what we *measure*.
- A distribution is the asymptotic result of the ideal process.
- Statistical parameters of the distribution can be estimated from a sample.
- Example: Dice Throw
- How likely is it to roll a 1, 2, 3, 4, 5, or 6?
- All statistically determined parameters are just estimates based on measurements.



Important Distinction

- Crucial Difference: distribution and sample
- A sample is what we *measure*.
- A distribution is the asymptotic result of the ideal process.
- Statistical parameters of the distribution can be estimated from a sample.
- Example: Dice Throw
- How likely is it to roll a 1, 2, 3, 4, 5, or 6?
- All statistically determined parameters are just estimates based on measurements.
- The parameters of a random process cannot be measured directly, but only be **estimated** from multiple measures.



Measures of the Average

- Assume that we have obtained a sample $A = (a_0, a_1, \dots, a_{n-1})$ of n observations from an experiment.

Measures of the Average

- Assume that we have obtained a sample $A = (a_0, a_1, \dots, a_{n-1})$ of n observations from an experiment, e.g., we have measured the qualities a_i of the best discovered solutions of $n = 101$ independent runs of an optimization algorithm.

Measures of the Average

- Assume that we have obtained a sample $A = (a_0, a_1, \dots, a_{n-1})$ of n observations from an experiment, e.g., we have measured the qualities a_i of the best discovered solutions of $n = 101$ independent runs of an optimization algorithm.
- We usually want to reduce this set of numbers to a single value which can give us an impression of what the “average outcome” (or result quality is).

Measures of the Average

- Assume that we have obtained a sample $A = (a_0, a_1, \dots, a_{n-1})$ of n observations from an experiment, e.g., we have measured the qualities a_i of the best discovered solutions of $n = 101$ independent runs of an optimization algorithm.
- We usually want to reduce this set of numbers to a single value which can give us an impression of what the “average outcome” (or result quality is).
- Three of the most common options for doing so, for estimating the “center” of a distribution, are to either compute the **arithmetic mean**, the **median**, or the **geometric mean**.

Arithmetic Mean

Definition (Arithmetic Mean)

The arithmetic mean $\text{mean}(A)$ is an **estimate** of the expected value of a distribution.

Arithmetic Mean

Definition (Arithmetic Mean)

The arithmetic mean $\text{mean}(A)$ is an **estimate** of the expected value of a distribution. Its is computed on data sample $A = (a_0, a_1, \dots, a_{n-1})$ as the sum of all n elements a_i in the sample data A divided by the total number n of values.

Arithmetic Mean

Definition (Arithmetic Mean)

The arithmetic mean $\text{mean}(A)$ is an **estimate** of the expected value of a distribution. Its is computed on data sample $A = (a_0, a_1, \dots, a_{n-1})$ as the sum of all n elements a_i in the sample data A divided by the total number n of values.

$$\text{mean}(A) = \frac{1}{n} \sum_{i=0}^{n-1} a_i \quad (1)$$

Median

Definition (Median)

The median $\text{med}(A)$ is the value separating the bigger half from the lower half of a data sample or distribution.

Median

Definition (Median)

The median $\text{med}(A)$ is the value separating the bigger half from the lower half of a data sample or distribution. Its estimate is the value right in the middle of a *sorted* data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \ \forall i \in 1 \dots (n-1)$.

Median

Definition (Median)

The median $\text{med}(A)$ is the value separating the bigger half from the lower half of a data sample or distribution. Its estimate is the value right in the middle of a *sorted* data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \ \forall i \in 1 \dots (n-1)$.

$$\text{med}(A) = \begin{cases} a_{\frac{n-1}{2}} & \text{if } n \text{ is odd} \\ \frac{1}{2} (a_{\frac{n}{2}-1} + a_{\frac{n}{2}}) & \text{otherwise} \end{cases} \quad \text{if } a_{i-1} \leq a_i \ \forall i \in 1 \dots (n-1) \quad (2)$$

Outliers

- Sometimes the data contains outliers^{19 20}.

Outliers

- Sometimes the data contains outliers^{19 20}, i.e., observations which are much different from the other measurements.

Outliers

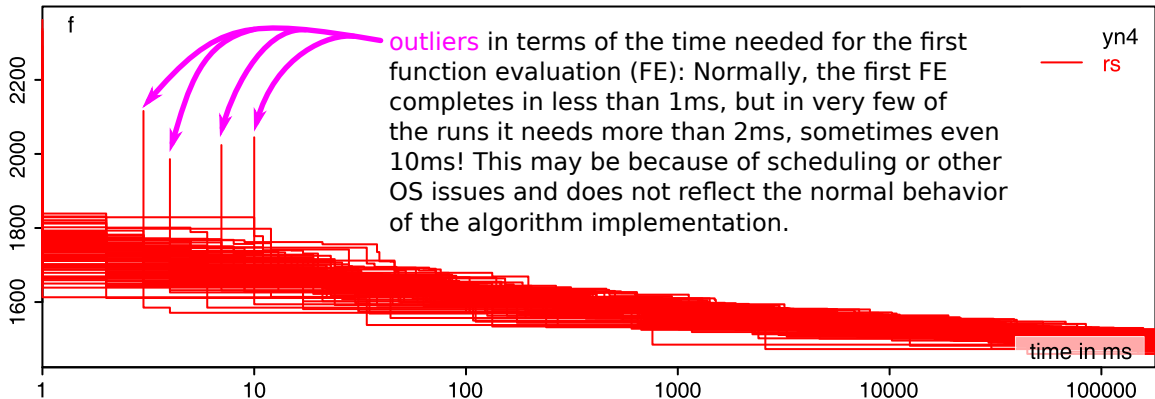
- Sometimes the data contains outliers^{19 20}, i.e., observations which are much different from the other measurements.
- They may represent measurement errors or observations which have been disturbed by unusual effects.

Outliers

- Sometimes the data contains outliers^{19 20}, i.e., observations which are much different from the other measurements.
- They may represent measurement errors or observations which have been disturbed by unusual effects.
- For example, maybe the operating system was updating itself during a run of one of our algorithms and, thus, took away some of the computation budget.

Outliers

- For example, maybe the operating system was updating itself during a run of one of our algorithms and, thus, took away some of the computation budget.
- In my experiments here, there are sometimes outliers in the time that it takes to create and evaluate the first candidate solution.



Outliers

- For example, maybe the operating system was updating itself during a run of one of our algorithms and, thus, took away some of the computation budget.
- In my experiments here, there are sometimes outliers in the time that it takes to create and evaluate the first candidate solution.
- But outliers are actually important. So I say this right now. I will also say it again later. But I am afraid that you may tune out during the following example. So remember: Outliers are important. Anyway...

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$

$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$

$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$

- We find that

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

- We find that

- $\text{mean}(A) = \frac{1}{19} \sum_{i=0}^{18} a_i = \frac{133}{19} = 7$

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

- We find that
 - $\text{mean}(A) = \frac{1}{19} \sum_{i=0}^{18} a_i = \frac{133}{19} = 7$ and
 - $\text{mean}(B) = \frac{1}{19} \sum_{i=0}^{18} b_i = \frac{10'127}{19} = 553$

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

- We find that
 - $\text{mean}(A) = \frac{1}{19} \sum_{i=0}^{18} a_i = \frac{133}{19} = 7$ and
 - $\text{mean}(B) = \frac{1}{19} \sum_{i=0}^{18} b_i = \frac{10'127}{19} = 553$, while
 - $\text{med}(A) = a_9 = 6$

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

- We find that
 - $\text{mean}(A) = \frac{1}{19} \sum_{i=0}^{18} a_i = \frac{133}{19} = 7$ and
 - $\text{mean}(B) = \frac{1}{19} \sum_{i=0}^{18} b_i = \frac{10'127}{19} = 553$, while
 - $\text{med}(A) = a_9 = 6$ and
 - $\text{med}(B) = b_9 = 6$.

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

- We find that
 - $\text{mean}(A) = \frac{1}{19} \sum_{i=0}^{18} a_i = \frac{133}{19} = 7$ and
 - $\text{mean}(B) = \frac{1}{19} \sum_{i=0}^{18} b_i = \frac{10'127}{19} = 553$, while
 - $\text{med}(A) = a_9 = 6$ and
 - $\text{med}(B) = b_9 = 6$.
- The median is not affected by the outliers.

Example for Data Samples w/o Outlier

- Two sets of data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

- We find that
 - $\text{mean}(A) = \frac{1}{19} \sum_{i=0}^{18} a_i = \frac{133}{19} = 7$ and
 - $\text{mean}(B) = \frac{1}{19} \sum_{i=0}^{18} b_i = \frac{10'127}{19} = 553$, while
 - $\text{med}(A) = a_9 = 6$ and
 - $\text{med}(B) = b_9 = 6$.
- The median is not affected by the outliers.
- $\text{mean}(B) = 553$ is a value completely different from anything that actually occurs in B . . .
... it gives us a completely wrong impression.

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing.

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results.

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Imagine that:

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Imagine that: Your algorithm can actually solve the TSP or MaxSat problem in polynomial time on 90% of all instances...

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Imagine that: Your algorithm can actually solve the TSP or MaxSat problem in polynomial time on 90% of all instances... ...but on 10%, it needs exponential time.

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Imagine that: Your algorithm can actually solve the TSP or MaxSat problem in polynomial time on 90% of all instances... ..but on 10%, it needs exponential time. If you just look at the median runtime, you may think you discovered something awesome.

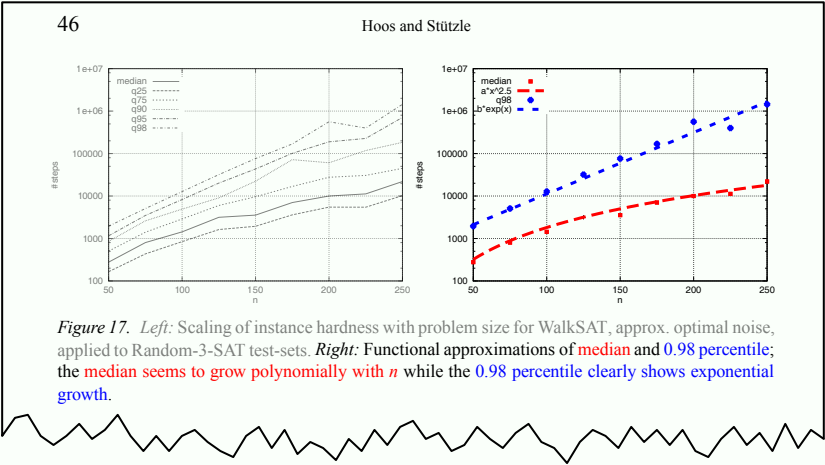
Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Imagine that: Your algorithm can actually solve the TSP or MaxSat problem in polynomial time on 90% of all instances... ..but on 10%, it needs exponential time. If you just look at the median runtime, you may think you discovered something awesome. Actually, this is quite common...

Outliers can be important!

- If you
1. Th
co
2. (U
ob
efl
• Instead
• bu



(Taken from the paper “Local Search Algorithms for SAT: An Empirical Evaluation” by Hoos and Stützle, coloring added manually²¹.)

ing. This
y, or the
tside”

in
tial
ing

Outliers can be important!

- If you think about it, where could outliers in **our** experiments come from?
 1. The operating systems scheduling or other strange effects could mess with our timing. This could cause worse results. But usually this is already it.
 2. (Unless your objective function is noisy, e.g., if you measure some physical quantity, or the objective function involves randomized simulations, there are hardly any other “outside” effects that could mess up our results!)
- Instead, most likely there could be
 - **bugs in our code!**
 - Bugs in our code are *the* most important number one reason for outliers!
 - Yes, also in your code! (Btw: Please use unit tests.)
 - Or: **bad (but rare) worst-case behaviors of our algorithm!**

Imagine that: Your algorithm can actually solve the TSP or MaxSat problem in polynomial time on 90% of all instances... ..but on 10%, it needs exponential time. If you just look at the median runtime, you may think you discovered something awesome. Actually, this is quite common...
- Thus, we may actually **want** that outliers influence our statistics...

Geometric Mean

OK, arithmetic mean, median . . . but what about the geometric mean?

Geometric Mean

OK, arithmetic mean, median ... but what about the geometric mean?

Definition (Geometric Mean)

The geometric mean $\text{geom}(A)$ is the n^{th} root of the product of n **positive** values.

Geometric Mean

OK, arithmetic mean, median ... but what about the geometric mean?

Definition (Geometric Mean)

The geometric mean $\text{geom}(A)$ is the n^{th} root of the product of n **positive** values.

$$\text{geom}(A) = \sqrt[n]{\prod_{i=0}^{n-1} a_i} \quad (3)$$

(4)

Geometric Mean

OK, arithmetic mean, median ... but what about the geometric mean?

Definition (Geometric Mean)

The geometric mean $\text{geom}(A)$ is the n^{th} root of the product of n **positive** values.

$$\text{geom}(A) = \sqrt[n]{\prod_{i=0}^{n-1} a_i} \quad (3)$$

$$\text{geom}(A) = \exp \left(\frac{1}{n} \sum_{i=0}^{n-1} \log a_i \right) \quad (4)$$

Normalized Data

- Often, our data is somehow **normalized**.

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1			
I_2			
I_3			

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s		
I_2			
I_3			

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s		
I_2	20 s		
I_3			

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s		
I_2	20 s		
I_3	40 s		

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s	20 s	
I_2	20 s		
I_3	40 s		

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s	20 s	
I_2	20 s	40 s	
I_3	40 s		

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s	20 s	
I_2	20 s	40 s	
I_3	40 s	10 s	

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	
I_3	40 s	10 s	

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say we solve the problem instances I_1 to I_3 with the different algorithms A_1 to A_3 .
- We measure the required runtimes as follows:

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s

Normalized Data

- Often, our data is somehow **normalized**.
- We measure the required runtimes as follows:
- The arithmetic mean values are the same.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s

Normalized Data

- Often, our data is somehow **normalized**.
- The arithmetic mean and the median values are the same.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- The arithmetic mean, the median, and the geometric mean values are the same.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- The arithmetic mean, the median, and the geometric mean values are the same.
- We can conclude that the three algorithms offer the same performance in average over these benchmark instances.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- We can conclude that the three algorithms offer the same performance in average over these benchmark instances.
- But often the measured numbers “look messier” and are harder to compare at first glance.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- But often the measured numbers “look messier” and are harder to compare at first glance.
- So often we want to normalize them by picking one algorithm as “standard” and dividing them by its measurements.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- But often the measured numbers “look messier” and are harder to compare at first glance.
- So often we want to normalize them by picking one algorithm as “standard” and dividing them by its measurements.
- Let's say A_1 was a well-known heuristic, maybe we even took its results from a paper, and we want to use it as baseline for comparison and normalize our data by it.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- So often we want to normalize them by picking one algorithm as “standard” and dividing them by its measurements.
- Let's say A_1 was a well-known heuristic, maybe we even took its results from a paper, and we want to use it as baseline for comparison and normalize our data by it.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- Let's say A_1 was a well-known heuristic, maybe we even took its results from a paper, and we want to use it as baseline for comparison and normalize our data by it.
- OK, so we get this table with normalized values, which allow us to make sense of the data at first glance.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values, which allow us to make sense of the data at first glance.
- If we now compute the arithmetic mean

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values, which allow us to make sense of the data at first glance.
- If we now compute the arithmetic mean, then **A_1 is best**

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values, which allow us to make sense of the data at first glance.
- If we now compute the arithmetic mean, then **A_1 is best** and **A_3 looks worst**.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values, which allow us to make sense of the data at first glance.
- If we now compute the arithmetic mean, then **A_1 is best** and **A_3 looks worst**.
- According to the median

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67
med:	1.00	2.00	0.50

Normalized Data

- Often, our data is somehow **normalized**.
- If we now compute the arithmetic mean, then A_1 is best and A_3 looks worst.
- According to the median, A_3 is best

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67
med:	1.00	2.00	0.50

Normalized Data

- Often, our data is somehow **normalized**.
- If we now compute the arithmetic mean, then A_1 is best and A_3 looks worst.
- According to the median, A_3 is best and A_2 is worst!

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67
med:	1.00	2.00	0.50

Normalized Data

- Often, our data is somehow **normalized**.
- If we now compute the arithmetic mean, then A_1 is best and A_3 looks worst.
- According to the median, A_3 is best and A_2 is worst!
- Only the geometric mean still indicates that the algorithms perform the same. . .

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67
med:	1.00	2.00	0.50
geom:	1.00	1.00	1.00

Normalized Data

- Often, our data is somehow **normalized**.
- Hm.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	1.00	2.00	4.00
I_2	1.00	2.00	0.50
I_3	1.00	0.25	0.50
<hr/>			
mean:	1.00	1.42	1.67
med:	1.00	2.00	0.50
geom:	1.00	1.00	1.00

Normalized Data

- Often, our data is somehow **normalized**.
- Hm. OK, then let's normalize using the results of A_2 instead.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

Normalized Data

- Often, our data is somehow **normalized**.
- Hm. OK, then let's normalize using the results of A_2 instead.
- OK, so we get this table with normalized values.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values.
- If we now compute the arithmetic mean

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values.
- If we now compute the arithmetic mean, then A_2 is best

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values.
- If we now compute the arithmetic mean, then **A_2 is best** and **A_1 looks worst**.

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42

Normalized Data

- Often, our data is somehow **normalized**.
- OK, so we get this table with normalized values.
- If we now compute the arithmetic mean, then A_2 is best and A_1 looks worst.
- According to the median

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42
med:	0.50	1.00	2.00

Normalized Data

- Often, our data is somehow **normalized**.
- If we now compute the arithmetic mean, then A_2 is best and A_1 looks worst.
- According to the median, A_1 is best

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42
med:	0.50	1.00	2.00

Normalized Data

- Often, our data is somehow **normalized**.
- If we now compute the arithmetic mean, then A_2 is best and A_1 looks worst.
- According to the median, A_1 is best and A_3 is worst!

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42
med:	0.50	1.00	2.00

Normalized Data

- Often, our data is somehow **normalized**.
- If we now compute the arithmetic mean, then A_2 is best and A_1 looks worst.
- According to the median, A_1 is best and A_3 is worst!
- Only the geometric mean still indicates that the algorithms perform the same. . .

	A_1	A_2	A_3
I_1	10 s	20 s	40 s
I_2	20 s	40 s	10 s
I_3	40 s	10 s	20 s
<hr/>			
mean:	23.33 s	23.33 s	23.33 s
med:	20.00 s	20.00 s	20.00 s
geom:	20.00 s	20.00 s	20.00 s

	A_1	A_2	A_3
I_1	0.50	1.00	2.00
I_2	0.50	1.00	0.25
I_3	4.00	1.00	2.00
<hr/>			
mean:	1.67	1.00	1.42
med:	0.50	1.00	2.00
geom:	1.00	1.00	1.00

Normalized Data

- Often, our data is somehow **normalized**.
- Only the geometric mean still indicates that the algorithms perform the same. . .
- The geometric mean is the only meaningful average if we have **normalized** data!²²

Normalized Data

- Often, our data is somehow **normalized**.
- The geometric mean is the only meaningful average if we have **normalized** data!²²
- And we very often have normalized data.

Normalized Data

- Often, our data is somehow **normalized**.
- The geometric mean is the only meaningful average if we have **normalized** data!²²
- And we very often have normalized data.
- For example, at least half of the papers on the Job Shop Scheduling Problem normalize the result qualities they obtain on benchmark instances with the *Best Known Solutions* (BKS).

Normalized Data

- Often, our data is somehow **normalized**.
- The geometric mean is the only meaningful average if we have **normalized** data!²²
- And we very often have normalized data.
- For example, at least half of the papers on the Job Shop Scheduling Problem normalize the result qualities they obtain on benchmark instances with the *Best Known Solutions* (BKS) and then compute the arithmetic mean.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median:
 - If the arithmetic mean is much worse than the median, then

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median:
 - If the arithmetic mean is much worse than the median, then
 - maybe we have a bug in our code that only sometimes has an impact.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median:
 - If the arithmetic mean is much worse than the median, then
 - maybe we have a bug in our code that only sometimes has an impact or
 - our algorithm has a bad worst-case behavior (which is also good to know).

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median:
 - If the arithmetic mean is much worse than the median, then
 - maybe we have a bug in our code that only sometimes has an impact or
 - our algorithm has a bad worst-case behavior (which is also good to know).
 - If the median is much worse than the mean, then the mean is too optimistic, i.e., most of the time we should expect worse results.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median:
 - If the arithmetic mean is much worse than the median, then
 - maybe we have a bug in our code that only sometimes has an impact or
 - our algorithm has a bad worst-case behavior (which is also good to know).
 - If the median is much worse than the mean, then the mean is too optimistic, i.e., most of the time we should expect worse results.
- If there are outliers, the value of the arithmetic mean itself may be very different from any actually observed value, while the median is (almost always) similar to some actual measurements.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized, e.g.,
 - if we divide result qualities by results of well-known heuristics or “Best-Known Solutions” or

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized, e.g.,
 - if we divide result qualities by results of well-known heuristics or “Best-Known Solutions” or
 - if we normalize the runtime using another algorithm as standard.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized, e.g.,
 - if we divide result qualities by results of well-known heuristics or “Best-Known Solutions” or
 - if we normalize the runtime using another algorithm as standard.
- Then, the arithmetic mean and median can be very misleading and the geometric mean must be computed.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized, e.g.,
 - if we divide result qualities by results of well-known heuristics or “Best-Known Solutions” or
 - if we normalize the runtime using another algorithm as standard.
- Then, the arithmetic mean and median can be very misleading and the geometric mean must be computed.
- I think: On raw data, compute all three measures of average, and pay special attention to the one looking the worst.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized, e.g.,
 - if we divide result qualities by results of well-known heuristics or “Best-Known Solutions” or
 - if we normalize the runtime using another algorithm as standard.
- Then, the arithmetic mean and median can be very misleading and the geometric mean must be computed.
- I think: On raw data, compute all three measures of average, and pay special attention to the one looking the worst. On normalized data, compute the geometric mean.

Arithmetic Mean vs. Median vs. Geometric Mean

- Most publications report arithmetic mean results, many report median results, almost none report geometric means.
- The median is more robust against outliers compared to the arithmetic mean, however, in normal application scenarios, there are very few acceptable reasons for outliers.
- We therefore want to know both the arithmetic mean and the median.
- Often, our data is implicitly or explicitly normalized, e.g.,
 - if we divide result qualities by results of well-known heuristics or “Best-Known Solutions” or
 - if we normalize the runtime using another algorithm as standard.
- Then, the arithmetic mean and median can be very misleading and the geometric mean must be computed.
- I think: On raw data, compute all three measures of average, and pay special attention to the one looking the worst. On normalized data, compute the geometric mean, but also consider the arithmetic mean and median *if and only if they make **your** algorithm look worse.*

Measures of the Spread

- The average gives us a good impression about the central value or location of a distribution.

Measures of the Spread

- The average gives us a good impression about the central value or location of a distribution.
- It does not tell us much about the range of the data.

Measures of the Spread

- The average gives us a good impression about the central value or location of a distribution.
- It does not tell us much about the range of the data.
- We do not know whether the data we have measured is very similar to the median or whether it may differ very much from the mean.

Measures of the Spread

- The average gives us a good impression about the central value or location of a distribution.
- It does not tell us much about the range of the data.
- We do not know whether the data we have measured is very similar to the median or whether it may differ very much from the mean.
- An average alone is not very meaningful – if we known nothing about the range of the data.

Measures of the Spread

- The average gives us a good impression about the central value or location of a distribution.
- It does not tell us much about the range of the data.
- We do not know whether the data we have measured is very similar to the median or whether it may differ very much from the mean.
- An average alone is not very meaningful – if we know nothing about the range of the data.
- We can therefore compute a measure of dispersion, i.e., a value that tells us whether the observations are stretched and spread far or squeezed tight around the center.

Variance

Definition (Variance)

The variance is the expectation of the squared deviation of a random variable from its mean.

Variance

Definition (Variance)

The variance is the expectation of the squared deviation of a random variable from its mean. The variance $\text{var}(A)$ of a data sample $A = (a_0, a_1, \dots, a_{n-1})$ with n observations can be estimated as:

$$\text{var}(A) = \frac{1}{n-1} \sum_{i=0}^{n-1} (a_i - \text{mean}(A))^2$$

Standard Deviation

Definition (Standard Deviation)

The statistical estimate $\text{sd}(A)$ of the standard deviation of a data sample $A = (a_0, a_1, \dots, a_{n-1})$ with n observations is the square root of the estimated variance $\text{var}(A)$.

$$\text{sd}(A) = \sqrt{\text{var}(A)}$$

Standard Deviation

- Small standard deviations indicate that the observations tend to be similar to the mean.

Standard Deviation

- Small standard deviations indicate that the observations tend to be similar to the mean.
- Large standard deviations indicate that they tend to be far from the mean.

Standard Deviation

- Small standard deviations indicate that the observations tend to be similar to the mean.
- Large standard deviations indicate that they tend to be far from the mean.
- Small standard deviations in optimization results and runtime indicate that the algorithm is reliable.

Standard Deviation

- Small standard deviations indicate that the observations tend to be similar to the mean.
- Large standard deviations indicate that they tend to be far from the mean.
- Small standard deviations in optimization results and runtime indicate that the algorithm is reliable.
- Large standard deviations indicate unreliable algorithms.

Standard Deviation

- Small standard deviations indicate that the observations tend to be similar to the mean.
- Large standard deviations indicate that they tend to be far from the mean.
- Small standard deviations in optimization results and runtime indicate that the algorithm is reliable.
- Large standard deviations indicate unreliable algorithms, but may also offer a potential that could be exploited.

Standard Deviation

- Small standard deviations indicate that the observations tend to be similar to the mean.
- Large standard deviations indicate that they tend to be far from the mean.
- Small standard deviations in optimization results and runtime indicate that the algorithm is reliable.
- Large standard deviations indicate unreliable algorithms, but may also offer a potential that could be exploited: Given enough time, we can restart algorithms several times and expect to get different (and thus sometimes better) solutions.

Quantiles

Definition (Quantile)

The q -quantiles are the cut points that divide a sorted data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \forall i \in 1 \dots (n-1)$ into q -equally sized parts.

Quantiles

Definition (Quantile)

The q -quantiles are the cut points that divide a sorted data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \forall i \in 1 \dots (n-1)$ into q -equally sized parts. quantile_q^k be the k^{th} q -quantile, with $k \in 1 \dots (q-1)$, i.e., there are $q-1$ of the q -quantiles.

$$\begin{aligned} h &= (n-1) \frac{k}{q} \\ \text{quantile}_q^k(A) &= \begin{cases} a_h & \text{if } h \text{ is integer} \\ a_{\lfloor h \rfloor} + (h - \lfloor h \rfloor) * (a_{\lfloor h \rfloor + 1} - a_{\lfloor h \rfloor}) & \text{otherwise} \end{cases} \end{aligned}$$

Quantiles

Definition (Quantile)

The q -quantiles are the cut points that divide a sorted data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \forall i \in 1 \dots (n-1)$ into q -equally sized parts. quantile_q^k be the k^{th} q -quantile, with $k \in 1 \dots (q-1)$, i.e., there are $q-1$ of the q -quantiles.

$$\begin{aligned} h &= (n-1) \frac{k}{q} \\ \text{quantile}_q^k(A) &= \begin{cases} a_h & \text{if } h \text{ is integer} \\ a_{\lfloor h \rfloor} + (h - \lfloor h \rfloor) * (a_{\lfloor h \rfloor + 1} - a_{\lfloor h \rfloor}) & \text{otherwise} \end{cases} \end{aligned}$$

- Quantiles are a generalized form of the median.

Quantiles

Definition (Quantile)

The q -quantiles are the cut points that divide a sorted data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \forall i \in 1 \dots (n-1)$ into q -equally sized parts. quantile_q^k be the k^{th} q -quantile, with $k \in 1 \dots (q-1)$, i.e., there are $q-1$ of the q -quantiles.

$$\begin{aligned} h &= (n-1) \frac{k}{q} \\ \text{quantile}_q^k(A) &= \begin{cases} a_h & \text{if } h \text{ is integer} \\ a_{\lfloor h \rfloor} + (h - \lfloor h \rfloor) * (a_{\lfloor h \rfloor + 1} - a_{\lfloor h \rfloor}) & \text{otherwise} \end{cases} \end{aligned}$$

- Quantiles are a generalized form of the median.
- The $\text{quantile}_1^2(A)$ is the median of A

Quantiles

Definition (Quantile)

The q -quantiles are the cut points that divide a sorted data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \forall i \in 1 \dots (n-1)$ into q -equally sized parts. quantile_q^k be the k^{th} q -quantile, with $k \in 1 \dots (q-1)$, i.e., there are $q-1$ of the q -quantiles.

$$\begin{aligned} h &= (n-1) \frac{k}{q} \\ \text{quantile}_q^k(A) &= \begin{cases} a_h & \text{if } h \text{ is integer} \\ a_{\lfloor h \rfloor} + (h - \lfloor h \rfloor) * (a_{\lfloor h \rfloor + 1} - a_{\lfloor h \rfloor}) & \text{otherwise} \end{cases} \end{aligned}$$

- Quantiles are a generalized form of the median.
- The $\text{quantile}_1^2(A)$ is the median of A
- 4-quantiles are called *quartiles*.

Quantiles

Definition (Quantile)

The q -quantiles are the cut points that divide a sorted data sample $A = (a_0, a_1, \dots, a_{n-1})$ where $a_{i-1} \leq a_i \forall i \in 1 \dots (n-1)$ into q -equally sized parts. quantile_q^k be the k^{th} q -quantile, with $k \in 1 \dots (q-1)$, i.e., there are $q-1$ of the q -quantiles.

$$\begin{aligned} h &= (n-1) \frac{k}{q} \\ \text{quantile}_q^k(A) &= \begin{cases} a_h & \text{if } h \text{ is integer} \\ a_{\lfloor h \rfloor} + (h - \lfloor h \rfloor) * (a_{\lfloor h \rfloor + 1} - a_{\lfloor h \rfloor}) & \text{otherwise} \end{cases} \end{aligned}$$

- Quantiles are a generalized form of the median.
- The $\text{quantile}_1^2(A)$ is the median of A
- 4-quantiles are called *quartiles*.
- We often consider *percentiles* or write things like “98% quantile” or “0.98 percentile” or “98% percentile” meaning $\text{quantile}_{100}^{98}$.

Standard Deviation: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

Standard Deviation: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$\text{mean}(A) = 7$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

$$\text{mean}(B) = 533$$

Standard Deviation: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$\text{mean}(A) = 7$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

$$\text{mean}(B) = 533$$

$$\text{var}(A) = \frac{1}{19-1} \sum_{i=1}^{19} (a_i - \text{mean}(A))^2 = \frac{198}{18} = 11$$

$$\text{var}(B) = \frac{1}{19-1} \sum_{i=1}^{19} (b_i - \text{mean}(B))^2 = \frac{94'763'306}{18} \approx 5'264'628$$

Standard Deviation: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$\text{mean}(A) = 7$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

$$\text{mean}(B) = 533$$

$$\text{var}(A) = \frac{1}{19-1} \sum_{i=1}^{19} (a_i - \text{mean}(A))^2 = \frac{198}{18} = 11$$

$$\text{var}(B) = \frac{1}{19-1} \sum_{i=1}^{19} (b_i - \text{mean}(B))^2 = \frac{94'763'306}{18} \approx 5'264'628$$

$$\text{sd}(A) = \sqrt{\text{var}A} = \sqrt{11} \approx 3.3$$

$$\text{sd}(B) = \sqrt{\text{var}B} = \sqrt{\frac{94'763'306}{18}} \approx 2294$$

Standard Deviation: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

$$\text{var}(A) = \frac{1}{19-1} \sum_{i=1}^{19} (a_i - \text{mean}(A))^2 = \frac{198}{18} = 11$$

$$\text{var}(B) = \frac{1}{19-1} \sum_{i=1}^{19} (b_i - \text{mean}(B))^2 = \frac{94'763'306}{18} \approx 5'264'628$$

$$\text{sd}(A) = \sqrt{\text{var}A} = \sqrt{11} \approx 3.3$$

$$\text{sd}(B) = \sqrt{\text{var}B} = \sqrt{\frac{94'763'306}{18}} \approx 2294$$

- Being based on the arithmetic mean, the variance and standard deviation are heavily influenced by outliers – with all pros and cons coming with that. . .

Quantiles: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

Quantiles: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$

$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$

Quantiles: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

$$\text{quantile}_4^1(A) = \text{quantile}_4^1(B) = 4.5$$

$$\text{quantile}_4^3(A) = \text{quantile}_4^3(B) = 9$$

Quantiles: Example

- Two data samples A and B with $n_a = n_b = 19$ values.

$$A = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 14)$$

$$B = (1, 3, 4, 4, 4, 5, 6, 6, 6, 6, 7, 7, 9, 9, 9, 10, 11, 12, 10'008)$$

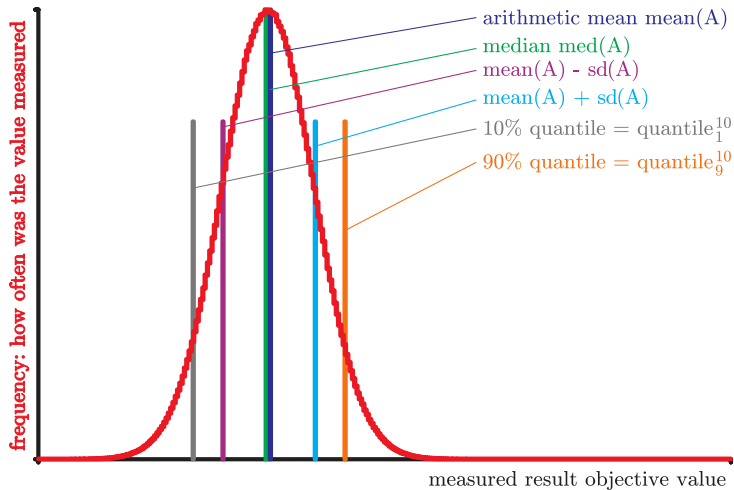
$$\text{quantile}_4^1(A) = \text{quantile}_4^1(B) = 4.5$$

$$\text{quantile}_4^3(A) = \text{quantile}_4^3(B) = 9$$

- Being generalizations of the median, the quantiles are little influenced by outliers – with all pros and cons coming with that. . .

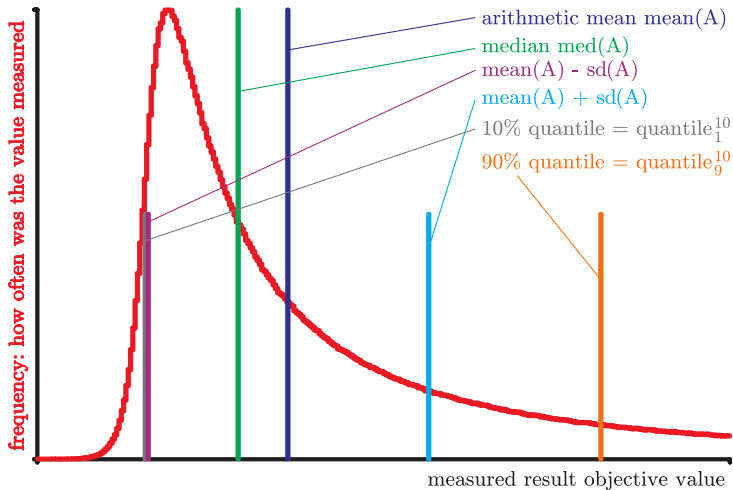
Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!



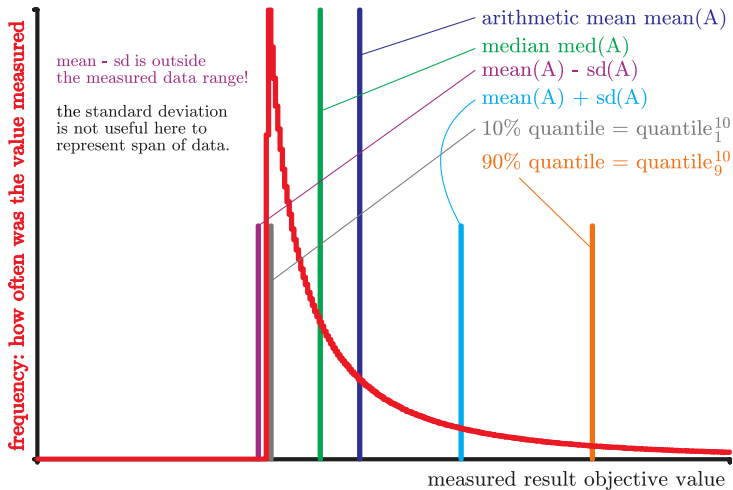
Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!



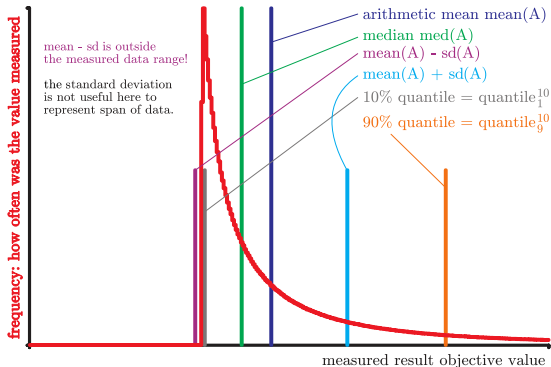
Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!



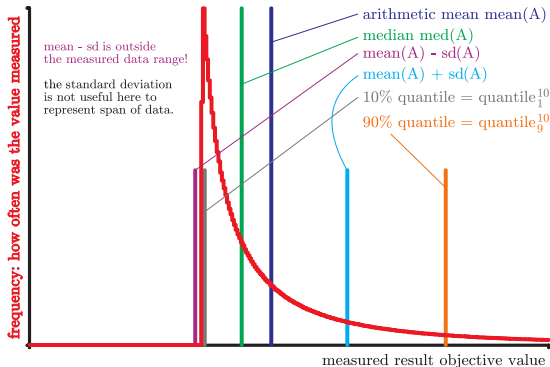
Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!
- Such a shape is possible in optimization!



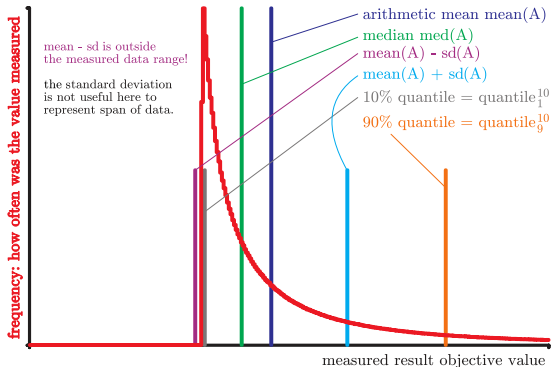
Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!
- Such a shape is possible in optimization:
 - The global optimum marks a lower bound for the possible objective values.



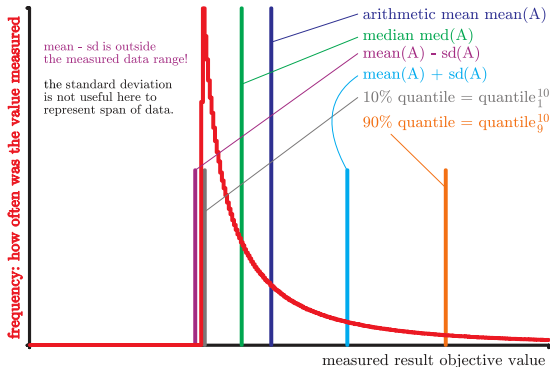
Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!
- Such a shape is possible in optimization:
 - The global optimum marks a lower bound for the possible objective values.
 - A good algorithm often returns results which are close-to-optimal.



Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!
- Such a shape is possible in optimization:
 - The global optimum marks a lower bound for the possible objective values.
 - A good algorithm often returns results which are close-to-optimal.
 - There may be a long tail of few but significantly worse runs.



Further Example

- The implicit assumption that $\text{mean} \pm \text{sd}$ is a meaningful range is not always true!
- Such a shape is possible in optimization:
 - The global optimum marks a lower bound for the possible objective values.
 - A good algorithm often returns results which are close-to-optimal.
 - There may be a long tail of few but significantly worse runs.
 - A statement such as “For this TSP instance, our algorithm can find tours for with a length of 100 ± 120 km.” makes little sense. . .

Statistical Comparisons



Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.
- For one of the two algorithms, the results will be better.

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.
- For one of the two algorithms, the results will be better.
- What does this mean?

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.
- For one of the two algorithms, the results will be better.
- What does this mean?
- It means that one of the two algorithms is better

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.
- For one of the two algorithms, the results will be better.
- What does this mean?
- It means that one of the two algorithms is better with a certain probability

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.
- For one of the two algorithms, the results will be better.
- What does this mean?
- It means that one of the two algorithms is better with a certain probability
- If we say "*A is better than B,*" we have a certain probability p to be wrong.

Introduction

- We can now, e.g., perform 20 runs each with two different optimization algorithms on one problem instance and compute the medians of a performance indicator.
- Likely, they will be different.
- For one of the two algorithms, the results will be better.
- What does this mean?
- It means that one of the two algorithms is better with a certain probability
- If we say “*A is better than B*,” we have a certain probability p to be wrong.
- The statement “*A is better than B*” makes only sense if we can give an upper bound α for the error probability p !

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and
- get a result (e.g., *“The median of A is bigger than the median of B ”*) together with an error probability p that the conclusion is wrong.

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and
- get a result (e.g., “*The median of A is bigger than the median of B*”) together with an error probability p that the conclusion is wrong.
- If p is less than a significance level (upper bound) α , we can accept the conclusion.

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and
- get a result (e.g., “*The median of A is bigger than the median of B*”) together with an error probability p that the conclusion is wrong.
- If p is less than a significance level (upper bound) α , we can accept the conclusion.
- Otherwise, the observation is not significant.

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and
- get a result (e.g., “*The median of A is bigger than the median of B*”) together with an error probability p that the conclusion is wrong.
- If p is less than a significance level (upper bound) α , we can accept the conclusion.
- Otherwise, the observation is not significant and must be ignored.

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and
- get a result (e.g., “*The median of A is bigger than the median of B*”) together with an error probability p that the conclusion is wrong.
- If p is less than a significance level (upper bound) α , we can accept the conclusion.
- Otherwise, the observation is not significant and must be ignored.
- But how can we arrive at such statements? How can we even estimate a probability to be wrong?

Statistical Tests

- Compare two data samples $A = (a_1, a_2, \dots)$ and $B = (b_1, b_2, \dots)$ and
- get a result (e.g., "*The median of A is bigger than the median of B*") together with an error probability p that the conclusion is wrong.
- If p is less than a significance level (upper bound) α , we can accept the conclusion.
- Otherwise, the observation is not significant and must be ignored.
- But how can we arrive at such statements? How can we even estimate a probability to be wrong?
- Disclaimer: I am not a mathematician. What follows are simplified explanations of concepts.

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.



heads



tails

Example for Underlying Idea

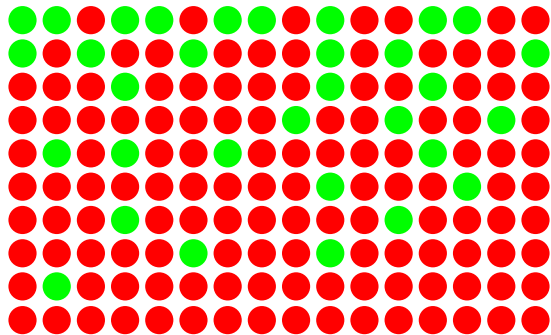
- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.



Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.
- Did I cheat? Is my coin “fixed?” (i.e., is your chance to win $\neq 0.5$)

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.
- Did I cheat? Is my coin “fixed?” (i.e., is your chance to win $\neq 0.5$)
- Assumption: I cheat. (alternative hypothesis H_1)

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.
- Did I cheat? Is my coin “fixed?” (i.e., is your chance to win $\neq 0.5$)
- Assumption: I cheat. (alternative hypothesis H_1)
- It is impossible to compute my winning probability if I cheated. . .

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.
- Did I cheat? Is my coin “fixed?” (i.e., is your chance to win $\neq 0.5$)
- Assumption: I cheat. (alternative hypothesis H_1)
- It is impossible to compute my winning probability if I cheated. . .
- Counter-Assumption: I did not cheat. (null hypothesis H_0)

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.
- Did I cheat? Is my coin “fixed?” (i.e., is your chance to win $\neq 0.5$)
- Assumption: I cheat. (alternative hypothesis H_1)
- It is impossible to compute my winning probability if I cheated. . .
- Counter-Assumption: I did not cheat. (null hypothesis H_0)
- How likely is it that I win **at least** 128 times if I did **not** cheat?

Example for Underlying Idea

- Coin flip game: We flip a coin. If it is heads, I give you 1 RMB, if it is tails, you give me 1 RMB.
- We play 160 times.
- I win 128 times. You win 32 times.
- Did I cheat? Is my coin “fixed?” (i.e., is your chance to win $\neq 0.5$)
- Assumption: I cheat. (alternative hypothesis H_1)
- It is impossible to compute my winning probability if I cheated. . .
- Counter-Assumption: I did not cheat. (null hypothesis H_0)
- How likely is it that I win **at least** 128 times if I did **not** cheat?
- (What we will do right now is called *binomial test*.)

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$P(k \geq z|n) = \sum_{i=z}^n P(i|n)$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$P(k \geq z|n) = \sum_{i=z}^n P(i|n) = \sum_{i=128}^{160} P(i|160) = \sum_{i=128}^{160} \left[\binom{160}{i} \frac{1}{2^{160}} \right]$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$P(k \geq z|n) = \sum_{i=z}^n P(i|n) = \frac{1}{2^{160}} \sum_{i=128}^{160} \binom{160}{i}$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$\begin{aligned} P(k \geq z|n) &= \sum_{i=z}^n P(i|n) = \frac{1}{2^{160}} \sum_{i=128}^{160} \binom{160}{i} \\ &= \frac{1'538'590'628'148'134'280'316'221'828'039'113}{365'375'409'332'725'729'550'921'208'179'070'754'913'983'135'744} \end{aligned}$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$P(k \geq z|n) = \sum_{i=z}^n P(i|n) = \frac{1}{2^{160}} \sum_{i=128}^{160} \binom{160}{i} \approx \frac{1.539 * 10^{33}}{3.654 * 10^{47}}$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$\begin{aligned} P(k \geq z|n) &= \sum_{i=z}^n P(i|n) = \frac{1}{2^{160}} \sum_{i=128}^{160} \binom{160}{i} \approx \frac{1.539 * 10^{33}}{3.654 * 10^{47}} \\ &\approx 0.0000000000000000421098571 \end{aligned}$$

Example for Underlying Idea

- How likely is it that I win **at least** 128 times if I did not cheat?
- Then, the probabilities for heads and tails are $q = P(\text{head}) = P(\text{tail}) = 0.5$.
- Flipping a coin n times is a Bernoulli Process
- The probability $P(k|n)$ to flip $k \in 0..n$ times heads (or tails) is thus:

$$P(k|n) = \binom{n}{k} 0.5^k * (1 - 0.5)^{n-k} = \binom{n}{k} 0.5^k * 0.5^{n-k} = \binom{n}{k} \frac{1}{2^n}$$

- For winning **at least** $z = 128$ times, we need to compute:

$$\begin{aligned} P(k \geq z|n) &= \sum_{i=z}^n P(i|n) = \frac{1}{2^{160}} \sum_{i=128}^{160} \binom{160}{i} \approx \frac{1.539 * 10^{33}}{3.654 * 10^{47}} \\ &\approx 4.211 \cdot 10^{-15} \end{aligned}$$

Example for Underlying Idea

- Question: How likely is it that I win at least 128 times if I did not cheat?

Example for Underlying Idea

- Question: How likely is it that I win at least 128 times if I did not cheat?
- If the coin was an ideal coin, the chance that I win at least 128 out of 160 times is about $4 \cdot 10^{-15}$.

Example for Underlying Idea

- Question: How likely is it that I win at least 128 times if I did not cheat?
- If the coin was an ideal coin, the chance that I win at least 128 out of 160 times is about $4 \cdot 10^{-15}$.
- If you claim that I cheat, your chance to be wrong is about $4 \cdot 10^{-15}$.

Example for Underlying Idea

- Question: How likely is it that I win at least 128 times if I did not cheat?
- If the coin was an ideal coin, the chance that I win at least 128 out of 160 times is about $4 \cdot 10^{-15}$.
- If you claim that I cheat, your chance to be wrong is about $4 \cdot 10^{-15}$.
- Thus, if we cannot accept a chance p to be wrong higher than a significance level $\alpha = 1\%$, we can still say:

The observation is significant, I did likely cheat.

A More Specific Example for Tests

- We want to compare two algorithms \mathcal{A} and \mathcal{B} on a given problem instance.

A More Specific Example for Tests

- We want to compare two algorithms \mathcal{A} and \mathcal{B} on a given problem instance.
- We have conducted a small experiment and measured objective values of their final results in a few runs in form of the two data sets A and B , respectively:

$$A = (2, 5, 6, 7, 9, 10)$$

$$B = (1, 3, 4, 8)$$

A More Specific Example for Tests

- We want to compare two algorithms \mathcal{A} and \mathcal{B} on a given problem instance.
- We have conducted a small experiment and measured objective values of their final results in a few runs in form of the two data sets A and B , respectively:

$$A = (2, 5, 6, 7, 9, 10)$$

$$B = (1, 3, 4, 8)$$

- From this, we can estimate the arithmetic means:

A More Specific Example for Tests

- We want to compare two algorithms \mathcal{A} and \mathcal{B} on a given problem instance.
- We have conducted a small experiment and measured objective values of their final results in a few runs in form of the two data sets A and B , respectively:

$$A = (2, 5, 6, 7, 9, 10)$$

$$B = (1, 3, 4, 8)$$

- From this, we can estimate the arithmetic means:

$$\text{mean}(A) = \frac{39}{6} = 6.5$$

$$\text{mean}(B) = \frac{16}{4} = 4$$

A More Specific Example

$$\begin{aligned}\text{mean}(A) &= \frac{39}{6} = 6.5 \\ \text{mean}(B) &= \frac{16}{4} = 4\end{aligned}$$

- It looks like algorithm \mathcal{B} may produce the smaller objective values.

A More Specific Example

$$\begin{aligned}\text{mean}(A) &= \frac{39}{6} = 6.5 \\ \text{mean}(B) &= \frac{16}{4} = 4\end{aligned}$$

- It looks like algorithm \mathcal{B} may produce the smaller objective values.
- But is this assumption justified based on the data we have?

A More Specific Example

$$\begin{aligned}\text{mean}(A) &= \frac{39}{6} = 6.5 \\ \text{mean}(B) &= \frac{16}{4} = 4\end{aligned}$$

- It looks like algorithm \mathcal{B} may produce the smaller objective values.
- But is this assumption justified based on the data we have?
- Is the difference between $\text{mean}(A)$ and $\text{mean}(B)$ significant at a threshold of $\alpha = 2\%$?

A More Specific Example

- If \mathcal{B} is truly better than \mathcal{A} , which is our hypothesis H_1 , then we cannot calculate anything.

A More Specific Example

- If \mathcal{B} is truly better than \mathcal{A} , which is our hypothesis H_1 , then we cannot calculate anything.
- Let us therefore assume as null hypothesis H_0 the observed difference did just happen by chance and, well, $\mathcal{A} \equiv \mathcal{B}$.

A More Specific Example

- If \mathcal{B} is truly better than \mathcal{A} , which is our hypothesis H_1 , then we cannot calculate anything.
- Let us therefore assume as null hypothesis H_0 the observed difference did just happen by chance and, well, $\mathcal{A} \equiv \mathcal{B}$.
- Then, this would mean that the data samples A and B stem from the **same** algorithm (as $\mathcal{A} \equiv \mathcal{B}$).

A More Specific Example

- If \mathcal{B} is truly better than \mathcal{A} , which is our hypothesis H_1 , then we cannot calculate anything.
- Let us therefore assume as null hypothesis H_0 the observed difference did just happen by chance and, well, $\mathcal{A} \equiv \mathcal{B}$.
- Then, this would mean that the data samples A and B stem from the **same** algorithm (as $\mathcal{A} \equiv \mathcal{B}$).
- The division into the two sets would only be artificial, an artifact of our experimental design.

A More Specific Example

- If \mathcal{B} is truly better than \mathcal{A} , which is our hypothesis H_1 , then we cannot calculate anything.
- Let us therefore assume as null hypothesis H_0 the observed difference did just happen by chance and, well, $\mathcal{A} \equiv \mathcal{B}$.
- Then, this would mean that the data samples A and B stem from the **same** algorithm (as $\mathcal{A} \equiv \mathcal{B}$).
- The division into the two sets would only be artificial, an artifact of our experimental design.
- Instead of having two data samples, we only have one, namely the union set O with 10 elements:

A More Specific Example

- If \mathcal{B} is truly better than \mathcal{A} , which is our hypothesis H_1 , then we cannot calculate anything.
- Let us therefore assume as null hypothesis H_0 the observed difference did just happen by chance and, well, $\mathcal{A} \equiv \mathcal{B}$.
- Then, this would mean that the data samples A and B stem from the **same** algorithm (as $\mathcal{A} \equiv \mathcal{B}$).
- The division into the two sets would only be artificial, an artifact of our experimental design.
- Instead of having two data samples, we only have one, namely the union set O with 10 elements:

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$
- There are $\binom{10}{4} = 210$ different ways to draw 4 (or 6) elements from O

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$
- There are $\binom{10}{4} = 210$ different ways to draw 4 (or 6) elements from O
- If H_0 holds, all have the same probability

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$
- There are $\binom{10}{4} = 210$ different ways to draw 4 (or 6) elements from O
- If H_0 holds, all have the same probability
- Use a program to test the combinations

A More Specific Example

```
/** an example class enumerating all combinations */
public class EnumerateAtLeastAsExtremeScenarios {
    public static void main(String[] args) {
        int meanLowerOrEqualTo4 = 0; // how often did we find a mean <= 4
        int totalCombinations = 0; // total number of tested combinations

        for (int i = 10; i > 0; i--) { // as 0 = numbers from 1 to 10
            for (int j = (i - 1); j > 0; j--) { // we can conveniently iterate
                for (int k = (j - 1); k > 0; k--) { // over all 4-element combos
                    for (int l = (k - 1); l > 0; l--) { // with 4 such nested loops
                        if (((i + j + k + l) / 4.0) <= 4) { // check for the extreme cases
                            meanLowerOrEqualTo4++; // count the extreme case
                            totalCombinations++; // add up combos, to verify
                        }
                    }
                }
            }
        }

        System.out.println(meanLowerOrEqualTo4 + "␣" + totalCombinations);
    }
}
```

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$
- There are $\binom{10}{4} = 210$ different ways to draw 4 (or 6) elements from O
- If H_0 holds, all have the same probability
- There are 27 such combinations with a mean of **less or equal 4**.

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$
- There are $\binom{10}{4} = 210$ different ways to draw 4 (or 6) elements from O
- If H_0 holds, all have the same probability
- There are 27 such combinations with a mean of less or equal 4.
- The probability p to observe a situation at least as extreme as A and B under H_0 is thus:

A More Specific Example

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

- Any division C into two sets with 4 and 6 elements has the same probability
- $|O| = 10$
- There are $\binom{10}{4} = 210$ different ways to draw 4 (or 6) elements from O
- If H_0 holds, all have the same probability
- There are 27 such combinations with a mean of less or equal 4.
- The probability p to observe a situation at least as extreme as A and B under H_0 is thus:

$$p = \frac{\# \text{cases } C : \text{mean}(c) \leq \text{mean}(b)}{\# \text{all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa.

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$
$$\sum_{\forall o \in O} o = \sum_{o=1}^{10} o = \frac{10(10+1)}{2} = 55$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

$$\sum_{\forall o \in O} o = \sum_{o=1}^{10} o = \frac{10(10+1)}{2} = 55$$

$$\text{mean}(b) = \left(\frac{1}{4} \sum_{\forall b \in B} b \right) \leq 4 \implies \left(\sum_{\forall b \in B} b \right) \leq 4 * 4 \leq 16$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

$$\sum_{\forall o \in O} o = \sum_{o=1}^{10} o = \frac{10(10+1)}{2} = 55$$

$$\text{mean}(b) = \left(\frac{1}{4} \sum_{\forall b \in B} b \right) \leq 4 \implies \left(\sum_{\forall b \in B} b \right) \leq 4 * 4 \leq 16$$

$$O = A \cup B \implies \sum_{\forall a \in A} a = \left(\sum_{\forall o \in O} o \right) - \left(\sum_{\forall b \in B} b \right)$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$O = A \cup B = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$$

$$\sum_{\forall o \in O} o = \sum_{o=1}^{10} o = \frac{10(10+1)}{2} = 55$$

$$\text{mean}(b) = \left(\frac{1}{4} \sum_{\forall b \in B} b \right) \leq 4 \implies \left(\sum_{\forall b \in B} b \right) \leq 4 * 4 \leq 16$$

$$O = A \cup B \implies \sum_{\forall a \in A} a = \left(\sum_{\forall o \in O} o \right) - \left(\sum_{\forall b \in B} b \right)$$

$$\sum_{\forall b \in B} b \leq 16 \implies \left(\sum_{\forall a \in A} a \right) \geq 55 - 16 \geq 39$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$\sum_{\forall o \in O} o = \sum_{o=1}^{10} o = \frac{10(10+1)}{2} = 55$$

$$\text{mean}(b) = \left(\frac{1}{4} \sum_{\forall b \in B} b \right) \leq 4 \implies \left(\sum_{\forall b \in B} b \right) \leq 4 * 4 \leq 16$$

$$O = A \cup B \implies \sum_{\forall a \in A} a = \left(\sum_{\forall o \in O} o \right) - \left(\sum_{\forall b \in B} b \right)$$

$$\sum_{\forall b \in B} b \leq 16 \implies \left(\sum_{\forall a \in A} a \right) \geq 55 - 16 \geq 39$$

$$\text{mean}(A) = \frac{1}{6} \left(\sum_{\forall a \in A} a \right)$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$\text{mean}(b) = \left(\frac{1}{4} \sum_{\forall b \in B} b \right) \leq 4 \implies \left(\sum_{\forall b \in B} b \right) \leq 4 * 4 \leq 16$$

$$O = A \cup B \implies \sum_{\forall a \in A} a = \left(\sum_{\forall o \in O} o \right) - \left(\sum_{\forall b \in B} b \right)$$

$$\sum_{\forall b \in B} b \leq 16 \implies \left(\sum_{\forall a \in A} a \right) \geq 55 - 16 \geq 39$$

$$\text{mean}(A) = \frac{1}{6} \left(\sum_{\forall a \in A} a \right)$$

$$\text{mean}(B) \leq 4 \implies \text{mean}(A) \geq \frac{39}{6} \geq 6.5$$

A More Specific Example

- Extreme cases into the other direction are the same, because if $\text{mean}(B) \leq 4$ then $\text{mean}(A) \geq 6.5$ for any division $A \cup B = O$ and vice versa:

$$\text{mean}(b) = \left(\frac{1}{4} \sum_{\forall b \in B} b \right) \leq 4 \implies \left(\sum_{\forall b \in B} b \right) \leq 4 * 4 \leq 16$$

$$O = A \cup B \implies \sum_{\forall a \in A} a = \left(\sum_{\forall o \in O} o \right) - \left(\sum_{\forall b \in B} b \right)$$

$$\sum_{\forall b \in B} b \leq 16 \implies \left(\sum_{\forall a \in A} a \right) \geq 55 - 16 \geq 39$$

$$\text{mean}(A) = \frac{1}{6} \left(\sum_{\forall a \in A} a \right)$$

$$\text{mean}(B) \leq 4 \implies \text{mean}(A) \geq \frac{39}{6} \geq 6.5$$

- So – of course – we could have also done the test the other way around with the same result!

A More Specific Example

- The probability p to observe a constellation at least as extreme as A or B under H_0 is thus:

$$p = \frac{\#\text{cases } C : \text{mean}(c) \leq \text{mean}(b)}{\#\text{all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

A More Specific Example

- The probability p to observe a constellation at least as extreme as A or B under H_0 is thus:

$$p = \frac{\#\text{cases } C : \text{mean}(c) \leq \text{mean}(b)}{\#\text{all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

- If we claim that A and B are from distributions with different means...

A More Specific Example

- The probability p to observe a constellation at least as extreme as A or B under H_0 is thus:

$$p = \frac{\#\text{cases } C : \text{mean}(c) \leq \text{mean}(b)}{\#\text{all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

- If we claim that A and B are from distributions with different means. . .
- . . . we are wrong with probability $p \approx 0.13$

A More Specific Example

- The probability p to observe a constellation at least as extreme as A or B under H_0 is thus:

$$p = \frac{\#\text{cases } C : \text{mean}(c) \leq \text{mean}(b)}{\#\text{all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

- If we claim that A and B are from distributions with different means. . .
- . . . we are wrong with probability $p \approx 0.13$
- At a significance level of $\alpha = 2\%$, the **means** of A and B are not significantly different! ($2\% < 0.13$)

A More Specific Example

- The probability p to observe a constellation at least as extreme as A or B under H_0 is thus:

$$p = \frac{\#\text{cases } C : \text{mean}(c) \leq \text{mean}(b)}{\#\text{all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

- If we claim that A and B are from distributions with different means. . .
- . . . we are wrong with probability $p \approx 0.13$
- At a significance level of $\alpha = 2\%$, the **means** of A and B are not significantly different! ($2\% < 0.13$)
- Actually: This here is an example for an *Randomization Test*^{23 24}.

A More Specific Example

- The probability p to observe a constellation at least as extreme as A or B under H_0 is thus:

$$p = \frac{\text{\#cases } C : \text{mean}(c) \leq \text{mean}(b)}{\text{\#all cases}} = \frac{27}{210} = \frac{9}{70} \approx 0.1286$$

- If we claim that A and B are from distributions with different means. . .
- . . . we are wrong with probability $p \approx 0.13$
- At a significance level of $\alpha = 2\%$, the **means** of A and B are not significantly different! ($2\% < 0.13$)
- Actually: This here is an example for an *Randomization Test*^{23 24}.
- The method here is only feasible for small sample sets, real tests are more sophisticated

A fair warning

- There are many algorithms and even more configuration parameters.

A fair warning

- There are many algorithms and even more configuration parameters.
- All kinds of algorithm modules and parameters have some kind of impact on the performance.

A fair warning

- There are many algorithms and even more configuration parameters.
- All kinds of algorithm modules and parameters have some kind of impact on the performance.
- If I have two different algorithms \mathcal{A} and \mathcal{B} , logic dictates that their performance is also different.

A fair warning

- There are many algorithms and even more configuration parameters.
- All kinds of algorithm modules and parameters have some kind of impact on the performance.
- If I have two different algorithms \mathcal{A} and \mathcal{B} , logic dictates that their performance is also different.
- But is this difference usually significant?

A fair warning

- There are many algorithms and even more configuration parameters.
- All kinds of algorithm modules and parameters have some kind of impact on the performance.
- If I have two different algorithms \mathcal{A} and \mathcal{B} , logic dictates that their performance is also different.
- But is this difference usually significant?
- From the viewpoint of statistics: Probably **yes**.

A fair warning

- There are many algorithms and even more configuration parameters.
- All kinds of algorithm modules and parameters have some kind of impact on the performance.
- If I have two different algorithms \mathcal{A} and \mathcal{B} , logic dictates that their performance is also different.
- But is this difference usually significant?
- From the viewpoint of statistics: Probably yes.
- If I just conduct enough runs, maybe thousands, or millions, than even a difference of 0.001% in performance will pass a test as **significant**.

A fair warning

- There are many algorithms and even more configuration parameters.
- All kinds of algorithm modules and parameters have some kind of impact on the performance.
- If I have two different algorithms \mathcal{A} and \mathcal{B} , logic dictates that their performance is also different.
- But is this difference usually significant?
- From the viewpoint of statistics: Probably yes.
- If I just conduct enough runs, maybe thousands, or millions, than even a difference of 0.001% in performance will pass a test as significant.
- To be **practically significant**, the measured difference of results should be statistically significant already with few runs, say, 11 or 21, not just with ≥ 100 runs.

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before.

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 - Assume that the data samples follow a certain distribution

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 - Assume that the data samples follow a certain distribution
 - Examples²⁵: t -test (assumes normal distribution)

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 - Assume that the data samples follow a certain distribution
 - Examples²⁵: t -test (assumes normal distribution)
 - The distribution of the data we measure is unknown. . .

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 - Assume that the data samples follow a certain distribution
 - Examples²⁵: t -test (assumes normal distribution)
 - The distribution of the data we measure is unknown. . .
 - . . . and usually not normal nor symmetric (see the further quantiles/stddev plot example).

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 - Assume that the data samples follow a certain distribution
 - Examples²⁵: t -test (assumes normal distribution)
 - The distribution of the data we measure is unknown. . .
 - . . . and usually not normal nor symmetric (see the further quantiles/stddev plot example).
 - The condition for using such tests often cannot be met (known distribution)

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 - Assume that the data samples follow a certain distribution
 - Examples²⁵: t -test (assumes normal distribution)
 - The distribution of the data we measure is unknown. . .
 - . . . and usually not normal nor symmetric (see the further quantiles/stddev plot example).
 - The condition for using such tests often cannot be met (known distribution)
 - **Parametric tests should usually not be used here!**

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests
 - Make few assumption about the distribution from which the data was sampled.

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests
 - Make few assumption about the distribution from which the data was sampled.
 - Examples²: the Wilcoxon rank sum test with continuity correction (also called Mann-Whitney U test²⁶⁻²⁹), Fisher's Exact Test³⁰, the Sign Test^{27 31}, the Randomization Test^{23 24}, and Wilcoxon's Signed Rank Test³².

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests
 - Make few assumption about the distribution from which the data was sampled.
 - Examples²: the Wilcoxon rank sum test with continuity correction (also called Mann-Whitney U test²⁶⁻²⁹), Fisher's Exact Test³⁰, the Sign Test^{27 31}, the Randomization Test^{23 24}, and Wilcoxon's Signed Rank Test³².
 - These tests are more **robust** (less assumptions)

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests
 - Make few assumption about the distribution from which the data was sampled.
 - Examples²: the Wilcoxon rank sum test with continuity correction (also called Mann-Whitney U test²⁶⁻²⁹), Fisher's Exact Test³⁰, the Sign Test^{27 31}, the Randomization Test^{23 24}, and Wilcoxon's Signed Rank Test³².
 - These tests are more **robust** (less assumptions)
 - **This usually is the kind of test we want to use.**

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests
 - Make few assumption about the distribution from which the data was sampled.
 - Examples²: the Wilcoxon rank sum test with continuity correction (also called Mann-Whitney U test²⁶⁻²⁹), Fisher's Exact Test³⁰, the Sign Test^{27 31}, the Randomization Test^{23 24}, and Wilcoxon's Signed Rank Test³².
 - These tests are more **robust** (less assumptions)
 - This usually is the kind of test we want to use.
 - They work similar to the previous test example, but with larger sample sizes

Statistical Tests: Types

- Statistical tests are more elegant mathematical approaches than the example shown before. In order to work, they have preconditions, they make certain assumptions.
- There are two types of tests:
 1. Parametric Tests
 2. Non-Parametric Tests
 - Make few assumption about the distribution from which the data was sampled.
 - Examples²: the Wilcoxon rank sum test with continuity correction (also called Mann-Whitney U test²⁶⁻²⁹), Fisher's Exact Test³⁰, the Sign Test^{27 31}, the Randomization Test^{23 24}, and Wilcoxon's Signed Rank Test³².
 - These tests are more **robust** (less assumptions)
 - This usually is the kind of test we want to use.
 - They work similar to the previous test example, but with larger sample sizes
 - **Often, the most suitable test is the Mann-Whitney U test.**

Compare $N \geq 2$ Algorithms

- For comparing $N \geq 2$ algorithms, we can compare any two algorithms with each other

Compare $N \geq 2$ Algorithms

- For comparing $N \geq 2$ algorithms, we can compare any two algorithms with each other
- N Algorithms $\Rightarrow k = N(N - 1)/2$ statistical tests (e.g., Mann-Whitney U)

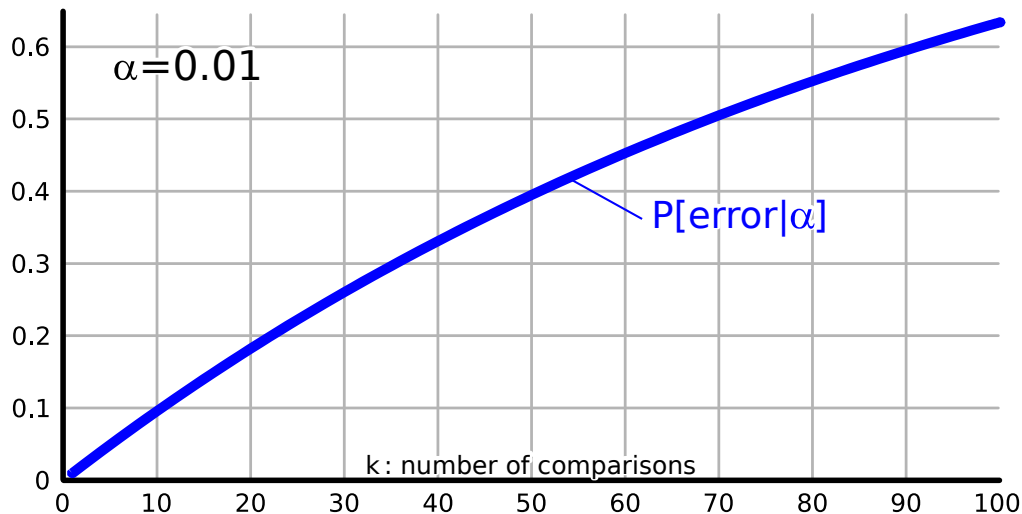
Compare $N \geq 2$ Algorithms

- For comparing $N \geq 2$ algorithms, we can compare any two algorithms with each other
- N Algorithms $\Rightarrow k = N(N - 1)/2$ statistical tests (e.g., Mann-Whitney U)
- k tests and each with error probability α

Compare $N \geq 2$ Algorithms

- For comparing $N \geq 2$ algorithms, we can compare any two algorithms with each other
- N Algorithms $\Rightarrow k = N(N - 1)/2$ statistical tests (e.g., Mann-Whitney U)
- k tests and each with error probability $\alpha \Rightarrow$ total probability E to make error
 $E = 1 - ((1 - \alpha)^k)$

Compare $N \geq 2$ Algorithms



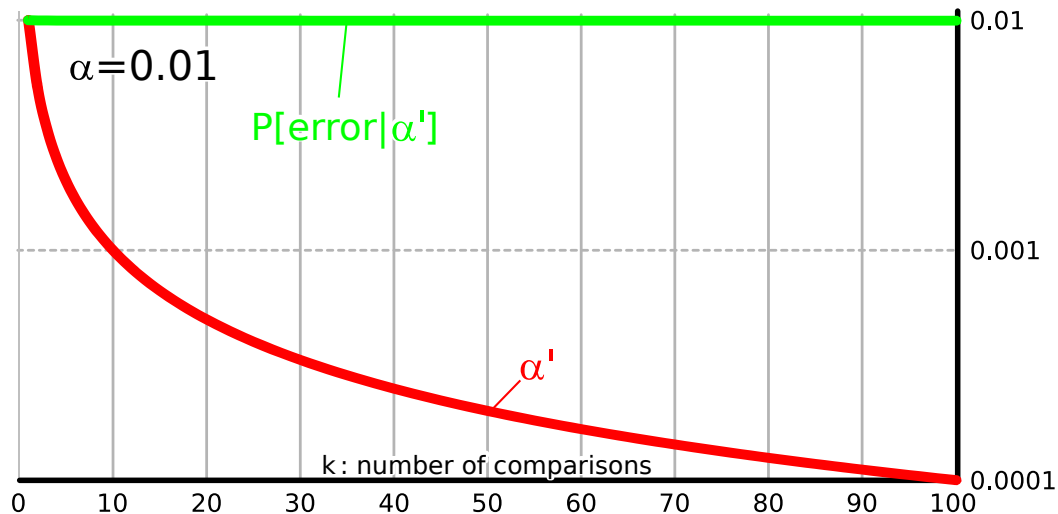
Compare $N \geq 2$ Algorithms

- For comparing $N \geq 2$ algorithms, we can compare any two algorithms with each other
- N Algorithms $\Rightarrow k = N(N - 1)/2$ statistical tests
- k tests and each with error probability $\alpha \Rightarrow$ total probability E to make error
 $E = 1 - ((1 - \alpha)^k)$
- Correction needed

Compare $N \geq 2$ Algorithms

- For comparing $N \geq 2$ algorithms, we can compare any two algorithms with each other
- N Algorithms $\Rightarrow k = N(N - 1)/2$ statistical tests
- k tests and each with error probability $\alpha \Rightarrow$ total probability E to make error
 $E = 1 - ((1 - \alpha)^k)$
- Correction needed: Bonferroni correction³³: Use $\alpha' = \alpha/k$ as significance level instead of α , then the overall probability E to make an error will remain $E \leq \alpha$.

Compare $N \geq 2$ Algorithms



Testing is Not Enough



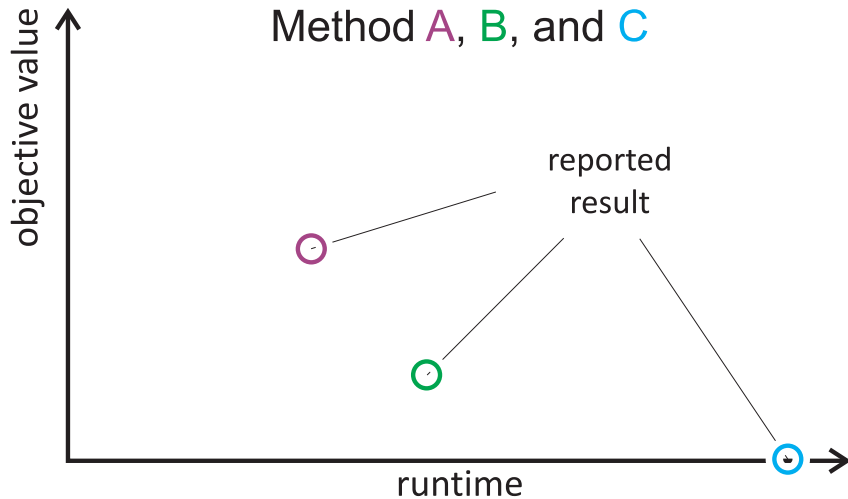
The question of termination

- Literature usually reports tuples “(instance, result, runtime)”

The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Papers often use different termination criteria

The question of termination



The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- **Problem**: Papers often use different termination criteria

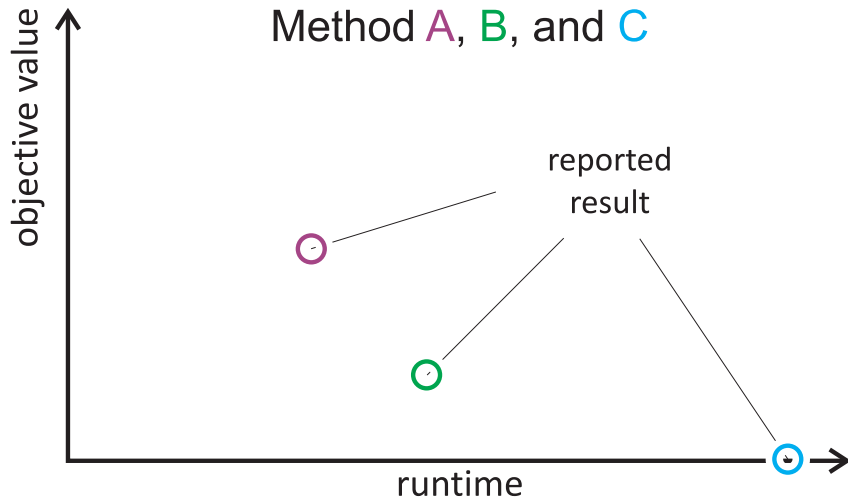
The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Problem: Papers often use different termination criteria
- Anytime Algorithms³⁴

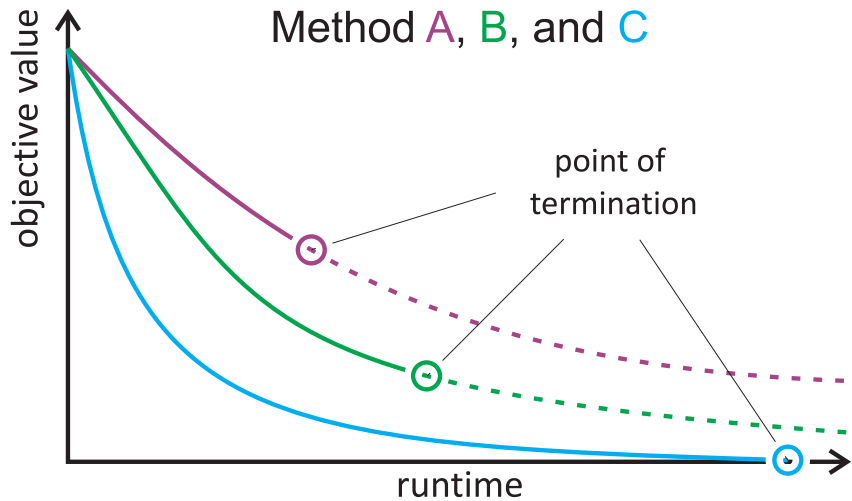
The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Problem: Papers often use different termination criteria
- Anytime Algorithms³⁴: Always have approximate solution, refine it iteratively

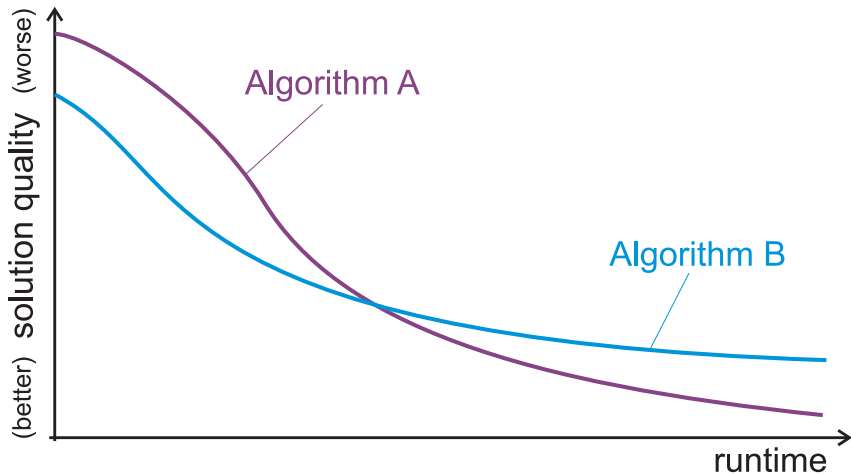
The question of termination



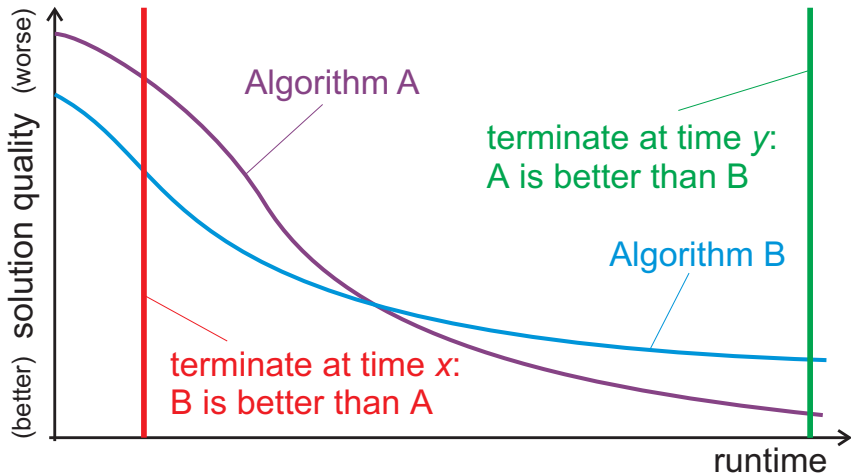
The question of termination



The question of termination



The question of termination



The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Problem: Papers often use different termination criteria
- Anytime Algorithms³⁴: Always have approximate solution, refine it iteratively
- One measure point per run or instance does not tell the whole story!

The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Problem: Papers often use different termination criteria
- Anytime Algorithms³⁴: Always have approximate solution, refine it iteratively
- One measure point per run or instance does not tell the whole story!
- Using statistical tests cannot solve this issue (still: at one point in time).

The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Problem: Papers often use different termination criteria
- Anytime Algorithms³⁴: Always have approximate solution, refine it iteratively
- One measure point per run or instance does not tell the whole story!
- Using statistical tests cannot solve this issue (still: at one point in time).
- We should have the “whole performance curves!”

The question of termination

- Literature usually reports tuples “(instance, result, runtime)”
- Problem: Papers often use different termination criteria
- Anytime Algorithms³⁴: Always have approximate solution, refine it iteratively
- One measure point per run or instance does not tell the whole story!
- Using statistical tests cannot solve this issue (still: at one point in time).
- We should have the “whole performance curves!” ... ideally mean or median curves over several runs!

Other Stuff



New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. I think this is not a good idea.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. I think this is not a good idea.
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.
- If we introduce a new optimization problem, we should apply well-known and well-established algorithms to it.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.
- If we introduce a new optimization problem, we should apply well-known and well-established algorithms to it. This way we get proper baseline results and can understand whether the problem is hard for the state-of-the-art and/or how far this state-of-the-art allows us to go.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.
- If we introduce a new optimization problem, we should apply well-known and well-established algorithms to it. This way we get proper baseline results and can understand whether the problem is hard for the state-of-the-art and/or how far this state-of-the-art allows us to go.
- If we have two “moving parts” at the same time, it is hard to understand whether an algorithm is good and whether a problem is hard.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.
- If we introduce a new optimization problem, we should apply well-known and well-established algorithms to it. This way we get proper baseline results and can understand whether the problem is hard for the state-of-the-art and/or how far this state-of-the-art allows us to go.
- If we have two “moving parts” at the same time, it is hard to understand whether an algorithm is good and whether a problem is hard.
- If you have an own new algorithm on a new problem and use other algorithms for comparision, you might be tempted to just use the most basic configurations of these algorithms.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.
- If we introduce a new optimization problem, we should apply well-known and well-established algorithms to it. This way we get proper baseline results and can understand whether the problem is hard for the state-of-the-art and/or how far this state-of-the-art allows us to go.
- If we have two “moving parts” at the same time, it is hard to understand whether an algorithm is good and whether a problem is hard.
- If you have an own new algorithm on a new problem and use other algorithms for comparision, you might be tempted to just use the most basic configurations of these algorithms. Then your algorithm might look good, while it actually is not.

New Algorithms and Problems

- There are many papers that introduce two things at the same time: a new optimization problem and a new algorithm. **I think this is not a good idea.**
- If we introduce a new optimization algorithm, we should test it on well-known, well-established benchmark problems. For such problems, results from other well-known and well-established algorithms exist – so we can compare our algorithm to them and investigate its performance objectively.
- If we introduce a new optimization problem, we should apply well-known and well-established algorithms to it. This way we get proper baseline results and can understand whether the problem is hard for the state-of-the-art and/or how far this state-of-the-art allows us to go.
- If we have two “moving parts” at the same time, it is hard to understand whether an algorithm is good and whether a problem is hard.
- If you have an own new algorithm on a new problem and use other algorithms for comparision, you might be tempted to just use the most basic configurations of these algorithms. Then your algorithm might look good, while it actually is not.
- **Know the standard benchmark instances for your field!**

Reproducibility

- Your experiments should be well-documented and reproducible.

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.
- For this purpose, you should make your code available, e.g., put it on GitHub or zenodo.org.

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.
- For this purpose, you should make your code available, e.g., put it on GitHub or zenodo.org.
- If your experiments are time-consuming, also make sure to properly store all your results in human- and machine-readable form (ideally in a CSV format).

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.
- For this purpose, you should make your code available, e.g., put it on GitHub or zenodo.org.
- If your experiments are time-consuming, also make sure to properly store all your results in human- and machine-readable form (ideally in a CSV format).
- You should make an archive such that a) I can directly run the same experiments that you did and b) also have all the data and tools to create the same statistics and figures.

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.
- For this purpose, you should make your code available, e.g., put it on GitHub or zenodo.org.
- If your experiments are time-consuming, also make sure to properly store all your results in human- and machine-readable form (ideally in a CSV format).
- You should make an archive such that a) I can directly run the same experiments that you did and b) also have all the data and tools to create the same statistics and figures.
- But what if someone finds an error in work?

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.
- For this purpose, you should make your code available, e.g., put it on GitHub or zenodo.org.
- If your experiments are time-consuming, also make sure to properly store all your results in human- and machine-readable form (ideally in a CSV format).
- You should make an archive such that a) I can directly run the same experiments that you did and b) also have all the data and tools to create the same statistics and figures.
- But what if someone finds an error in work?
- That is OK.

Reproducibility

- Your experiments should be well-documented and reproducible.
- In the ideal case, someone else can run your code and get the same results.
- For this purpose, you should make your code available, e.g., put it on GitHub or zenodo.org.
- If your experiments are time-consuming, also make sure to properly store all your results in human- and machine-readable form (ideally in a CSV format).
- You should make an archive such that a) I can directly run the same experiments that you did and b) also have all the data and tools to create the same statistics and figures.
- But what if someone finds an error in work?
- That is OK.
- Better they find it in your code that you voluntarily provided than after going through significant re-implementation effort. . .

Cheating

What are typical bad / cheating behavior in research on optimization?

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported. This can be prevented by generating the sequence of random seeds with a deterministic algorithm and reporting both.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported. This can be prevented by generating the sequence of random seeds with a deterministic algorithm and reporting both.
 - Only the benchmark instances where the algorithm performs well are chosen.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported. This can be prevented by generating the sequence of random seeds with a deterministic algorithm and reporting both.
 - Only the benchmark instances where the algorithm performs well are chosen. Be wary of statements such as “We now present the results of our algorithm on 10 of the TSPLib instances.” (TSPLib has more than 100...)

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported. This can be prevented by generating the sequence of random seeds with a deterministic algorithm and reporting both.
 - Only the benchmark instances where the algorithm performs well are chosen. Be wary of statements such as “We now present the results of our algorithm on 10 of the TSPLib instances.” (TSPLib has more than 100. . .)
 - Weak algorithms are chosen for comparison.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported. This can be prevented by generating the sequence of random seeds with a deterministic algorithm and reporting both.
 - Only the benchmark instances where the algorithm performs well are chosen. Be wary of statements such as “We now present the results of our algorithm on 10 of the TSPLib instances.” (TSPLib has more than 100. . .)
 - Weak algorithms are chosen for comparison. Comparison must always be done with the state-of-the-art on the specific problem at hand.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking:
 - On a benchmark instance, many runs are conducted with different random seeds. But only the 10 with the best results are reported. This can be prevented by generating the sequence of random seeds with a deterministic algorithm and reporting both.
 - Only the benchmark instances where the algorithm performs well are chosen. Be wary of statements such as “We now present the results of our algorithm on 10 of the TSPLib instances.” (TSPLib has more than 100. . .)
 - Weak algorithms are chosen for comparison. Comparison must always be done with the state-of-the-art on the specific problem at hand. Be wary of statements such as “We compare our algorithm with the standard Genetic Algorithm.” (because the SGA is usually not the state-of-the-art)

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated. Algorithm must be clearly specified and ideally the source code is available to prevent this.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated.
- Misleading statistics are reported (see our discussion on normalization).

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated.
- Misleading statistics are reported
- Uneven configuration effort: Much effort is spent on configuring the own algorithm, the algorithms used for comparison are used with bad settings.

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated.
- Misleading statistics are reported
- Uneven configuration effort.
- Incomparable results are reported (see our discussion on why testing is not enough).

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated.
- Misleading statistics are reported
- Uneven configuration effort.
- Incomparable results are reported.
- Misleading significance in test results (high α , many runs, no corrections).

Cheating

What are typical bad / cheating behavior in research on optimization?

- Cherry-Picking
- Sometimes, results may be straight up fabricated.
- Misleading statistics are reported
- Uneven configuration effort.
- Incomparable results are reported.
- Misleading significance in test results (high α , many runs, no corrections).

Reproducibility prevents cheating and misunderstandings!

Summary



Summary

- The optimization algorithms we consider in this lecture are **randomized**.

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a **statistical way** using data from **multiple runs**

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators
 2. quartiles or top/bottom 1% quantile to get a feeling for the usual range of values

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators
 2. quartiles or top/bottom 1% quantile to get a feeling for the usual range of values
 3. don't trust just arithmetic mean or standard deviation alone

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators
 2. quartiles or top/bottom 1% quantile to get a feeling for the usual range of values
 3. don't trust just arithmetic mean or standard deviation alone
 4. geometric mean if the data is normalized

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators
 2. quartiles or top/bottom 1% quantile to get a feeling for the usual range of values
 3. don't trust just arithmetic mean or standard deviation alone
 4. geometric mean if the data is normalized
- Use non-parametric statistical tests with corrections for multiple comparisons.

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators
 2. quartiles or top/bottom 1% quantile to get a feeling for the usual range of values
 3. don't trust just arithmetic mean or standard deviation alone
 4. geometric mean if the data is normalized
- Use non-parametric statistical tests with corrections for multiple comparisons.
- Do not only collect one data sample per run, try to plot progress curves.

Summary

- The optimization algorithms we consider in this lecture are randomized.
- Comparing them must be done in a statistical way using data from multiple runs
- Two views on performance:
 1. best result after fixed number of FEs/runtime
 2. number of FEs/runtime needed to get certain result
- For every single algorithm/configuration, compute:
 1. arithmetic and geometric mean and median of key performance indicators
 2. quartiles or top/bottom 1% quantile to get a feeling for the usual range of values
 3. don't trust just arithmetic mean or standard deviation alone
 4. geometric mean if the data is normalized
- Use non-parametric statistical tests with corrections for multiple comparisons.
- Do not only collect one data sample per run, try to plot progress curves.
- Use well-known benchmarks, provide your source code!

谢谢

Thank you



References I

1. Thomas Weise. *An Introduction to Optimization Algorithms*. Institute of Applied Optimization (IAO) [应用优化研究所] of the School of Artificial Intelligence and Big Data [人工智能与大数据学院] of Hefei University [合肥学院], Hefei [合肥市], Anhui [安徽省], China [中国], 2018–2020. URL <http://thomasweise.github.io/aitoa/>.
2. Thomas Weise. *Global Optimization Algorithms – Theory and Application*. it-weise.de (self-published), Germany, 2009. URL <http://www.it-weise.de/projects/book.pdf>.
3. Thomas Bartz-Beielstein, Carola Doerr, Daan van den Berg, Jakob Bossek, Sowmya Chandrasekaran, Tome Eftimov, Andreas Fischbach, Pascal Kerschke, William La Cava, Manuel López-Ibáñez, Katherine M. Malan, Jason H. Moore, Boris Naujoks, Patryk Orzechowski, Vanessa Volz, Markus Wagner, and Thomas Weise. Benchmarking in optimization: Best practice and open issues, July 2000. URL <https://arxiv.org/pdf/2007.03488.pdf>. arXiv:2007.03488v2 [cs.NE] 16 Dec 2020.
4. Eugene Leighton Lawler, Jan Karel Lenstra, Alexander Hendrik George Rinnooy Kan, and David B. Shmoys. Sequencing and scheduling: Algorithms and complexity. In Stephen C. Graves, Alexander Hendrik George Rinnooy Kan, and Paul H. Zipkin, editors, *Handbook of Operations Research and Management Science*, volume IV: Production Planning and Inventory, chapter 9, pages 445–522. North-Holland Scientific Publishers Ltd., Amsterdam, The Netherlands, 1993. doi:10.1016/S0927-0507(05)80189-6.
5. Bo Chen, Chris N. Potts, and Gerhard J. Woeginger. A review of machine scheduling: Complexity, algorithms and approximability. In Ding-Zhu Du and Panos M. Pardalos, editors, *Handbook of Combinatorial Optimization*, pages 1493–1641. Springer-Verlag US, Boston, MA, USA, 1998. ISBN 978-1-4613-7987-4. doi:10.1007/978-1-4613-0303-9_25. also pages 21–169 in volume 3/3 by Kluwer Academic Publishers.
6. Stephen Arthur Cook. The complexity of theorem-proving procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing (STOC'71)*, May 3–5, 1971, Shaker Heights, OH, USA, pages 151–158, New York, NY, USA, 1971. ACM. doi:10.1145/800157.805047.
7. Richard M. Karp. Reducibility among combinatorial problems. In Raymond E. Miller and James W. Thatcher, editors, *Complexity of Computer Computations. The IBM Research Symposia Series.*, pages 85–103. Springer, Boston, MA, USA, 1972. ISBN 978-1-4684-2003-6. doi:10.1007/978-1-4684-2001-2_9.
8. Scott Aaronson. The limits of quantum computers. *Scientific American*, 298(3):62–69, March 2008. doi:10.1038/scientificamerican0308-62. URL http://www.cs.virginia.edu/~robins/The_Limits_of_Quantum_Computers.pdf.
9. Nikolaus Hansen, Anne Auger, Steffen Finck, and Raymond Ros. Real-parameter black-box optimization benchmarking 2010: Experimental setup. Rapports de Recherche RR-7215, Institut National de Recherche en Informatique et en Automatique (INRIA), March 9 2010. URL <http://hal.inria.fr/inria-00462481>. inria-00462481.
10. Steffen Finck, Nikolaus Hansen, Raymond Ros, and Anne Auger. Coco documentation, release 15.03, November 17 2015. URL <http://coco.lri.fr/COC0doc/COC0.pdf>.

References II

11. Thomas Weise, Li Niu, and Ke Tang. AOAB – automated optimization algorithm benchmarking. In *Proceedings of the 12th Annual Conference Companion on Genetic and Evolutionary Computation (GECCO'10), July 7–11, 2010, Portland, OR, USA*, pages 1479–1486, New York, NY, USA, 2010. ACM Press. doi:10.1145/1830761.1830763.
12. Thomas Weise, Xiaofeng Wang, Qi Qi, Bin Li, and Ke Tang. Automatically discovering clusters of algorithm and problem instance behaviors as well as their causes from experimental data, algorithm setups, and instance features. *Applied Soft Computing Journal (ASOC)*, 73:366–382, December 2018. doi:10.1016/j.asoc.2018.08.030.
13. Thomas Weise, Raymond Chiong, Ke Tang, Jörg Lässig, Shigeyoshi Tsutsui, Wenxiang Chen, Zbigniew Michalewicz, and Xin Yao. Benchmarking optimization algorithms: An open source framework for the traveling salesman problem. *IEEE Computational Intelligence Magazine (CIM)*, 9:40–52, August 2014. doi:10.1109/MCI.2014.2326101.
14. Kenneth V. Price. Differential evolution vs. the functions of the 2nd ICEO. In Russ Eberhart, Peter Angeline, Thomas Bäck, Zbigniew Michalewicz, and Xin Yao, editors, *IEEE International Conference on Evolutionary Computation, April 13–16, 1997, Indianapolis, IN, USA*, pages 153–157. IEEE Computational Intelligence Society, 1997.
15. Anne Auger and Nikolaus Hansen. Performance evaluation of an advanced local search evolutionary algorithm. In *Proceedings of the IEEE Congress on Evolutionary Computation (CEC'05), September 2–4, 2005, Edinburgh, UK*, pages 1777–1784. IEEE, 2005. ISBN 0-7803-9363-5. doi:10.1109/CEC.2005.1554903. URL <http://www.cmap.polytechnique.fr/~nikolaus.hansen/cec2005localcmaes.pdf>.
16. Bernd Bischl, Pascal Kerschke, Lars Kotthoff, Thomas Marius Lindauer, Yuri Malitsky, Alexandre Fréchette, Holger H. Hoos, Frank Hutter, Kevin Leyton-Brown, Kevin Tierney, and Joaquin Vanschoren. ASlib: A benchmark library for algorithm selection. *Applied Intelligence – The International Journal of Research on Intelligent Systems for Real Life Complex Problems*, 237:41 – 58, 2016. doi:10.1016/j.artint.2016.04.003.
17. Pascal Kerschke, Jakob Bossek, and Heike Trautmann. Parameterization of state-of-the-art performance indicators: A robustness study based on inexact TSP solvers. In Hernán E. Aguirre and Keiki Takadama, editors, *Proceedings of the Genetic and Evolutionary Computation Conference Companion (GECCO'18), July 15–19, 2018, Kyoto, Japan*, pages 1737–1744. ACM, 2018. doi:10.1145/3205651.3208233.
18. Ke Tang, Xiaodong Li, Ponnuthurai Nagarathnam Suganthan, Zhenyu Yang, and Thomas Weise. Benchmark functions for the cec'2010 special session and competition on large-scale global optimization. Technical report, University of Science and Technology of China (USTC), School of Computer Science and Technology, Nature Inspired Computation and Applications Laboratory (NICAL), Hefei, Anhui, China, January 8 2010.
19. Frank E. Grubbs. Procedures for detecting outlying observations in samples. *Technometrics*, 11:1–21, 1969. doi:10.1080/00401706.1969.10490657.
20. Gangadharrao Soundalayarao Maddala. *Introduction to Econometrics*. MacMillan, New York, NY, USA, second edition, 1992. ISBN 978-0-02-374545-4.
21. Holger H. Hoos and Thomas Stützle. Local search algorithms for SAT: An empirical evaluation. *Journal of Automated Reasoning*, 24(4):421–481, May 2000. doi:10.1023/A:1006350622830. URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.29.978&type=pdf>.

References III

22. Philip J. Fleming and John J. Wallace. How not to lie with statistics: The correct way to summarize benchmark results. *Communications of the ACM*, 29(3): 218–221, 1986. doi:[10.1145/5666.5673](https://doi.org/10.1145/5666.5673).
23. Jürgen Bortz, Gustav Adolf Lienert, and Klaus Boehnke. *Verteilungsfreie Methoden in der Biostatistik*. Springer-Lehrbuch. Springer Medizin Verlag, Heidelberg, Germany, 3 edition, 2008. ISBN 3445110344. doi:[10.1007/978-3-540-74707-9](https://doi.org/10.1007/978-3-540-74707-9).
24. Eugene S. Edgington. *Randomization Tests*. CRC Press, Inc., Boca Raton, FL, USA, 3 edition, 1995. ISBN 0824796691 and 9780824796693.
25. Shaun Burke. Missing values, outliers, robust statistics & non-parametric methods. *LC.GC Europe Online Supplement*, 1(2):19–24, January 2001.
26. Daniel F. Bauer. Constructing confidence sets using rank statistics. *Journal of the American Statistical Association (J AM STAT ASSOC)*, 67:687–690, September 1972. doi:[10.1080/01621459.1972.10481279](https://doi.org/10.1080/01621459.1972.10481279).
27. Sidney Siegel and N. John Castellan Jr. *Nonparametric Statistics for The Behavioral Sciences*. Humanities/Social Sciences/Languages. McGraw-Hill, New York, NY, USA, 1988. ISBN 0-07-057357-3.
28. Myles Hollander and Douglas Alan Wolfe. *Nonparametric Statistical Methods*. John Wiley and Sons Ltd., New York, USA, 1973. ISBN 047140635X.
29. Henry B. Mann and Donald R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *The Annals of Mathematical Statistics (AOMS)*, 18(1):50–60, March 1947. doi:[10.1214/aoms/1177730491](https://doi.org/10.1214/aoms/1177730491). URL <http://projecteuclid.org/euclid.aoms/1177730491>.
30. Sir Ronald Aylmer Fisher. On the interpretation of χ^2 from contingency tables, and the calculation of p. *Journal of the Royal Statistical Society*, 85:87–94, 1922. URL <http://hdl.handle.net/2440/15173>.
31. Lorenz Gygax. Statistik für Nutztierethologen – Einführung in die statistische Denkweise: Was ist, was macht ein statistischer Test?, June 2003. URL <http://www.proximate-biology.ch/documents/introEtho.pdf>.
32. Frank Wilcoxon. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83, December 1945. URL <http://sci2s.ugr.es/keel/pdf/algorithm/articulo/wilcoxon1945.pdf>.
33. Olive Jean Dunn. Multiple comparisons among means. *Journal of the American Statistical Association (J AM STAT ASSOC)*, 56(293):52–64, March 1961. doi:[10.1080/01621459.1961.10482090](https://doi.org/10.1080/01621459.1961.10482090).
34. Mark S. Boddy and Thomas L. Dean. Solving time-dependent planning problems. Technical Report CS-89-03, Brown University, Department of Computer Science, Providence, RI, USA, February 1989. URL <ftp://ftp.cs.brown.edu/pub/techreports/89/cs89-03.pdf>.