

Technical Manual

# Pragmastat: Pragmatic Statistical Toolkit

29.01.2026

Andrey Akinshin · [andrey.akinshin@gmail.com](mailto:andrey.akinshin@gmail.com)

Version 5.0.0 · DOI: [10.5281/zenodo.17236778](https://doi.org/10.5281/zenodo.17236778)

This manual presents a toolkit of statistical procedures that provide reliable results across diverse real-world distributions, with ready-to-use implementations and detailed explanations. The toolkit consists of renamed, recombined, and refined versions of existing methods.

<b>Documentation</b>	pragmastat-v5.0.0.pdf web-v5.0.0.zip
<b>Implementations</b>	py-v5.0.0.zip ts-v5.0.0.zip r-v5.0.0.zip cs-v5.0.0.zip kt-v5.0.0.zip rs-v5.0.0.zip go-v5.0.0.zip
<b>Reference data</b>	tests-v5.0.0.zip sim-v5.0.0.zip
<b>Source code</b>	pragmastat-5.0.0.zip

Andrey Akinshin  
andrey.akinshin@gmail.com  
DOI: 10.5281/zenodo.17236778

**Copyright © 2025–2026 Andrey Akinshin**

This manual is licensed under the **Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License** (CC BY-NC-SA 4.0). You are free to share and adapt this material for non-commercial purposes, provided you give appropriate credit, indicate if changes were made, and distribute your contributions under the same license.

The accompanying source code and software implementations are licensed under the **MIT License**. You are free to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the software, subject to the conditions stated in the license. For complete license terms, see the LICENSE file in the source repository.

While the information in this manual is believed to be accurate at the date of publication, the author makes no warranty, express or implied, with respect to the material contained herein. The author shall not be liable for any errors, omissions, or damages arising from the use of this information.

Source code and implementations are available at [github.com/AndreyAkinshin/pragmastat](https://github.com/AndreyAkinshin/pragmastat).

Typeset with Typst. Text refined with LLM assistance.

## Contents

<b>1. Toolkit</b>	<b>4</b>
1.1. Rng .....	4
1.2. Sample .....	5
1.3. Shuffle .....	6
1.4. Median .....	7
1.5. Center .....	8
1.6. Spread .....	9
1.7. RelSpread .....	10
1.8. Shift .....	11
1.9. Ratio .....	12
1.10. AvgSpread .....	13
1.11. Disparity .....	14
1.12. PairwiseMargin .....	15
1.13. ShiftBounds .....	16
<b>2. Distributions</b>	<b>17</b>
2.1. Additive ('Normal') .....	17
2.2. Multiplic ('LogNormal') .....	18
2.3. Exponential .....	19
2.4. Power ('Pareto') .....	19
2.5. Uniform .....	20
<b>3. Algorithms</b>	<b>22</b>
3.1. Pseudorandom Number Generation .....	22
3.2. Shuffle Algorithm .....	23
3.3. Selection Sampling .....	23
3.4. Distribution Sampling .....	23
3.5. Fast Center .....	24
3.6. Fast Spread .....	30
3.7. Fast Shift .....	36
3.8. Fast PairwiseMargin .....	40
<b>4. Studies</b>	<b>49</b>
4.1. Summary Estimator Properties .....	49
4.2. Breakdown .....	50
4.3. Drift .....	50
4.4. Invariance .....	57
4.5. From Statistical Efficiency to Drift .....	58
4.6. From Confidence Level to Misrate .....	60
4.7. From Mann-Whitney U-test to Pairwise Margin .....	61
4.8. Additive ('Normal') Distribution .....	64

<b>5. Reference Implementations</b>	<b>69</b>
5.1. Python .....	69
5.2. TypeScript .....	71
5.3. R .....	72
5.4. C# .....	74
5.5. Kotlin .....	76
5.6. Rust .....	78
5.7. Go .....	80
<b>6. Reference Tests</b>	<b>83</b>
6.1. Center Tests .....	83
6.2. Spread Tests .....	84
6.3. RelSpread Tests .....	86
6.4. Shift Tests .....	87
6.5. Ratio Tests .....	89
6.6. AvgSpread Tests .....	91
6.7. Disparity Tests .....	92
6.8. PairwiseMargin Tests .....	93
6.9. ShiftBounds Tests .....	95
6.10. Test Framework .....	98
<b>7. Methodology</b>	<b>100</b>
7.1. Pragmatic Philosophy .....	100
7.2. Procedure-First Empiricism .....	101
7.3. Epistemic Humility .....	102
7.4. The Pairwise Principle .....	103
7.5. Median as Universal Aggregator .....	104
7.6. Convergence Conventions .....	105
7.7. Structural Unity .....	106
7.8. Generative Naming .....	107
7.9. The Inversion Principle .....	108
7.10. Multi-Audience Design .....	109
7.11. Reference Tests as Specification .....	110
7.12. Cross-Language Determinism .....	111
7.13. Summary Principles .....	112
<b>Bibliography</b>	<b>113</b>

## 1. Toolkit

This chapter provides formal definitions and properties of each toolkit function.

### 1.1. Rng

$$\text{Rng}(s)$$

Deterministic pseudorandom number generator from seed  $s$ .

**Algorithm** — xoshiro256++ with SplitMix64 seeding (see Pseudorandom Number Generation)

**Seed types** — integer seed or string seed (hashed via FNV-1a)

**Determinism** — identical sequences across all supported languages

**Period** —  $2^{256} - 1$

#### Notation

$X, Y$       random variables (generators of real measurements)

#### Properties

**Reproducibility**       $\text{Rng}(s)$  produces identical sequence for same  $s$

**Independence**      different seeds produce uncorrelated sequences

#### Example

$\text{Rng}(1729)$  — numeric seed

$\text{Rng}(\text{"experiment-1"})$  — string seed for named experiments

Use `Rng` whenever you need random numbers that can be reproduced later. The same seed will produce exactly the same sequence of random values, and this works identically across Python, TypeScript, R, C#, Kotlin, Rust, and Go. You can pass a simple integer like 42, or a descriptive string like "experiment-1" to make your code self-documenting. Each time you draw from the generator, its internal state advances, so if you need independent random streams, create separate generators with different seeds.

## 1.2. Sample

$r.\text{Sample}(x, k)$

Select  $k$  elements from sample  $x$  without replacement using generator  $r$ .

**Algorithm** — selection sampling (Fan, Muller, Rezucha 1962), see Pseudorandom Number Generation

**Complexity** —  $O(n)$  time, single pass

**Output** — preserves original order of selected elements

**Domain** —  $1 \leq k \leq n$

### Notation

$x = (x_1, \dots, x_n), y = (y_1, \dots, y_m)$       samples ( $n, m \geq 1$ )

$x_i, y_j$       individual measurements

### Properties

**Simple random sample**      each  $k$ -subset has equal probability

**Order preservation**      selected elements appear in order of first occurrence

**Determinism**      same generator state produces same selection

### Example

`Sample([1, 2, 3, 4, 5], 3, Rng(1729))` — select 3 elements

`Sample(x, n, r) = x` — selecting all elements returns original order

Use `Sample` when you need to pick a random subset of your data without replacement. Common uses include bootstrap resampling, creating cross-validation splits, or reducing a large dataset to a manageable size. Every possible subset of size  $k$  has equal probability of being selected, and the selected elements keep their original order. To make your subsampling reproducible, combine it with a seeded generator: `Sample(data, 100, Rng("training-set"))` will always select the same 100 elements.

## 1.3. Shuffle

`r.Shuffle(x)`

Uniformly random permutation of sample `x` using generator `r`.

**Algorithm** — Fisher-Yates (Knuth shuffle), see Pseudorandom Number Generation

**Complexity** —  $O(n)$  time,  $O(1)$  additional space

**Output** — new array (does not modify input)

### Properties

**Uniformity**      each of  $n!$  permutations has equal probability

**Determinism**      same generator state produces same permutation

### Example

`Shuffle([1, 2, 3, 4, 5], Rng(1729))` — shuffled copy

`Shuffle(x, r)` preserves multiset (same elements, different order)

Use `Shuffle` when you need a random reordering of your data. This is essential for permutation tests and useful whenever you want to eliminate any bias that might come from the original ordering. Every possible arrangement has exactly equal probability, which is required for valid statistical inference. The function returns a new shuffled array and leaves your original data unchanged. For reproducible results, pass a seeded generator: `Shuffle(data, Rng(1729))` will always produce the same permutation.

## 1.4. Median

$$\text{Median}(\mathbf{x}) = \begin{cases} x_{((n+1)/2)} & \text{if } n \text{ is odd} \\ \frac{x_{(n/2)} + x_{(n/2+1)}}{2} & \text{if } n \text{ is even} \end{cases}$$

The value splitting a sorted sample into two equal parts.

**Also known as** — 50th percentile, second quartile (Q2)

**Asymptotic** — value where  $P(X \leq \text{Median}) = 0.5$

**Complexity** —  $O(n)$  with selection,  $O(n \log n)$  with sorting

**Domain** — any real numbers

**Unit** — same as measurements

### Notation

$x_{(1)}, \dots, x_{(n)}$       order statistics (sorted sample)

### Properties

**Shift equivariance**       $\text{Median}(\mathbf{x} + k) = \text{Median}(\mathbf{x}) + k$

**Scale equivariance**       $\text{Median}(k \cdot \mathbf{x}) = k \cdot \text{Median}(\mathbf{x})$

### Example

$\text{Median}([1, 2, 3, 4, 5]) = 3$

$\text{Median}([1, 2, 3, 4]) = 2.5$

Use Median when you need maximum protection against outliers and corrupted data. It achieves a 50% breakdown point, meaning that up to half of your data can be arbitrarily bad before the estimate becomes meaningless. However, this extreme robustness comes at a cost: the median is less precise than Center when your data is clean. For most practical applications, Center offers a better tradeoff (29% breakdown with 95% efficiency). Reserve Median for situations where you suspect contamination levels above 29% or need the strongest possible robustness guarantee.



## 1.5. Center

$$\text{Center}(\mathbf{x}) = \text{Median}_{1 \leq i \leq j \leq n} \frac{x_i + x_j}{2}$$

Robust measure of location (central tendency).

**Also known as** — Hodges-Lehmann estimator, pseudomedian

**Asymptotic** — median of the average of two random measurements from  $X$

**Complexity** —  $O(n^2 \log n)$  naive,  $O(n \log n)$  fast (see Fast Center)

**Domain** — any real numbers

**Unit** — same as measurements

### Properties

**Shift equivariance**       $\text{Center}(\mathbf{x} + k) = \text{Center}(\mathbf{x}) + k$

**Scale equivariance**       $\text{Center}(k \cdot \mathbf{x}) = k \cdot \text{Center}(\mathbf{x})$

### Example

$\text{Center}([0, 2, 4, 6, 8]) = 4$

$\text{Center}(\mathbf{x} + 10) = 14$      $\text{Center}(3\mathbf{x}) = 12$

### References

Hodges & Lehmann (1963)

Sen (1963)

Use Center as your default choice when you need a single number to represent “where the data is.” It works like the familiar mean but does not break when your data contains a few bad measurements or outliers. You can have up to 29% of your data corrupted before Center becomes unreliable. At the same time, when your data is clean, Center is nearly as precise as the mean (95% efficiency), so you pay almost no cost for the added protection. If you are unsure whether to use mean, median, or something else, start with Center.

## 1.6. Spread

$$\text{Spread}(\mathbf{x}) = \text{Median}_{1 \leq i < j \leq n} |x_i - x_j|$$

Robust measure of dispersion (variability, scatter).

**Also known as** — Shamos scale estimator

**Asymptotic** — median of the absolute difference between two random measurements from  $X$

**Complexity** —  $O(n^2 \log n)$  naive,  $O(n \log n)$  fast (see Fast Spread)

**Domain** — any real numbers

**Unit** — same as measurements

**Corner case** —  $\text{Spread}(\mathbf{x}) = 0$  when  $n = 1$

### Properties

**Shift invariance**       $\text{Spread}(\mathbf{x} + k) = \text{Spread}(\mathbf{x})$

**Scale equivariance**       $\text{Spread}(k \cdot \mathbf{x}) = |k| \cdot \text{Spread}(\mathbf{x})$

**Non-negativity**       $\text{Spread}(\mathbf{x}) \geq 0$

### Example

$\text{Spread}([0, 2, 4, 6, 8]) = 4$

$\text{Spread}(\mathbf{x} + 10) = 4$      $\text{Spread}(2\mathbf{x}) = 8$

### References

Shamos (1976)

Use Spread when you want to know how much your measurements vary from each other. It serves the same purpose as standard deviation but does not explode when you have outliers or heavy-tailed data. The result comes in the same units as your measurements, so if Spread is 5 milliseconds, that tells you how much your values typically differ. Like Center, it tolerates up to 29% corrupted data. When comparing variability across datasets, Spread gives you a reliable answer even when standard deviation would be misleading or infinite.

## 1.7. RelSpread

$$\text{RelSpread}(\mathbf{x}) = \frac{\text{Spread}(\mathbf{x})}{|\text{Center}(\mathbf{x})|}$$

Relative dispersion normalized by location.

**Also known as** — robust coefficient of variation

**Domain** —  $\text{Center}(\mathbf{x}) \neq 0$

**Unit** — dimensionless

### Properties

**Scale invariance**       $\text{RelSpread}(k \cdot \mathbf{x}) = \text{RelSpread}(\mathbf{x})$

**Non-negativity**       $\text{RelSpread}(\mathbf{x}) \geq 0$

### Example

$\text{RelSpread}([0, 2, 4, 6, 8]) = 1$

$\text{RelSpread}(5x) = 1$

Use RelSpread when you want to compare how “noisy” different datasets are, even if they have completely different scales or units. A dataset centered around 100 with spread of 10 has the same relative variability as one centered around 1000 with spread of 100. Both show 10% relative variation, and RelSpread captures exactly this. This makes it useful for comparing measurement quality across different experiments, instruments, or physical quantities where absolute numbers are not directly comparable.

## 1.8. Shift

$$\text{Shift}(\mathbf{x}, \mathbf{y}) = \text{Median}_{1 \leq i \leq n, 1 \leq j \leq m} (x_i - y_j)$$

Robust measure of location difference between two samples.

**Also known as** — Hodges-Lehmann estimator for two samples

**Asymptotic** — median of the difference between random measurements from  $X$  and  $Y$

**Complexity** —  $O(mn \log(mn))$  naive,  $O((m+n) \log L)$  fast (see Fast Shift)

**Domain** — any real numbers

**Unit** — same as measurements

### Properties

**Self-difference**       $\text{Shift}(\mathbf{x}, \mathbf{x}) = 0$

**Shift equivariance**       $\text{Shift}(\mathbf{x} + k_x, \mathbf{y} + k_y) = \text{Shift}(\mathbf{x}, \mathbf{y}) + k_x - k_y$

**Scale equivariance**       $\text{Shift}(k \cdot \mathbf{x}, k \cdot \mathbf{y}) = k \cdot \text{Shift}(\mathbf{x}, \mathbf{y})$

**Antisymmetry**       $\text{Shift}(\mathbf{x}, \mathbf{y}) = -\text{Shift}(\mathbf{y}, \mathbf{x})$

### Example

$\text{Shift}([0, 2, 4, 6, 8], [10, 12, 14, 16, 18]) = -10$

$\text{Shift}(\mathbf{y}, \mathbf{x}) = -\text{Shift}(\mathbf{x}, \mathbf{y})$

### References

Hodges & Lehmann (1963)

Sidak et al. (1999)

Use Shift when you have two groups and want to know how much one differs from the other. If you are comparing response times between version A and version B, Shift tells you by how many milliseconds A is faster or slower than B. A negative result means the first group tends to be lower; positive means it tends to be higher. Unlike comparing means, Shift handles outliers gracefully and works well with skewed data. The result comes in the same units as your measurements, making it easy to interpret.

## 1.9. Ratio

$$\text{Ratio}(\mathbf{x}, \mathbf{y}) = \text{Median}_{1 \leq i \leq n, 1 \leq j \leq m} \left( \frac{x_i}{y_j} \right)$$

Robust measure of scale ratio between two samples.

**Asymptotic** — median of the ratio of random measurements from  $X$  and  $Y$

**Domain** —  $x_i, y_j > 0$  or  $x_i, y_j < 0$  (exclude  $|y_j| \approx 0$ )

**Unit** — dimensionless

**Caveat** —  $\text{Ratio}(\mathbf{x}, \mathbf{y}) \neq \frac{1}{\text{Ratio}(\mathbf{y}, \mathbf{x})}$  in general

### Properties

**Self-ratio**       $\text{Ratio}(\mathbf{x}, \mathbf{x}) = 1$

**Scale equivariance**       $\text{Ratio}(k_x \cdot \mathbf{x}, k_y \cdot \mathbf{y}) = \left( \frac{k_x}{k_y} \right) \cdot \text{Ratio}(\mathbf{x}, \mathbf{y})$

### Example

$$\text{Ratio}([1, 2, 4, 8, 16], [2, 4, 8, 16, 32]) = 0.5$$

$$\text{Ratio}(\mathbf{x}, \mathbf{x}) = 1 \quad \text{Ratio}(2\mathbf{x}, 5\mathbf{y}) = 0.4 \cdot \text{Ratio}(\mathbf{x}, \mathbf{y})$$

Use Ratio when you care about multiplicative relationships rather than additive differences. If one system is “twice as fast” or prices are “30% lower,” you are thinking in ratios. A result of 0.5 means the first group is typically half the size of the second; 2.0 means twice as large. This estimator is appropriate for quantities like prices, response times, and concentrations where relative comparisons make more sense than absolute ones. All values must have the same sign, and you should avoid values close to zero in the denominator sample.

## 1.10. AvgSpread

$$\text{AvgSpread}(\mathbf{x}, \mathbf{y}) = \frac{n \cdot \text{Spread}(\mathbf{x}) + m \cdot \text{Spread}(\mathbf{y})}{n + m}$$

Weighted average of dispersions (pooled scale).

**Also known as** — robust pooled standard deviation

**Domain** — any real numbers

**Unit** — same as measurements

**Caveat** —  $\text{AvgSpread}(\mathbf{x}, \mathbf{y}) \neq \text{Spread}(\mathbf{x} \cup \mathbf{y})$  (pooled scale, not concatenated spread)

### Properties

**Self-average**       $\text{AvgSpread}(\mathbf{x}, \mathbf{x}) = \text{Spread}(\mathbf{x})$

**Symmetry**       $\text{AvgSpread}(\mathbf{x}, \mathbf{y}) = \text{AvgSpread}(\mathbf{y}, \mathbf{x})$

**Scale equivariance**       $\text{AvgSpread}(k \cdot \mathbf{x}, k \cdot \mathbf{y}) = |k| \cdot \text{AvgSpread}(\mathbf{x}, \mathbf{y})$

**Mixed scaling**       $\text{AvgSpread}(k_1 \cdot \mathbf{x}, k_2 \cdot \mathbf{x}) = \frac{|k_1| + |k_2|}{2} \cdot \text{Spread}(\mathbf{x})$

### Example

$\text{AvgSpread}(\mathbf{x}, \mathbf{y}) = 5$  where  $\text{Spread}(\mathbf{x}) = 6$ ,  $\text{Spread}(\mathbf{y}) = 4$ ,  $n = m$

$\text{AvgSpread}(\mathbf{x}, \mathbf{y}) = \text{AvgSpread}(\mathbf{y}, \mathbf{x})$

Use AvgSpread when you need a single number representing the typical variability across two groups. It combines the spread of both samples, giving more weight to larger samples since they provide more reliable estimates. This pooled spread serves as a common reference scale, which is essential when you want to express a difference in relative terms. Disparity uses AvgSpread internally to normalize the shift into a scale-free effect size.

## 1.11. Disparity

$$\text{Disparity}(\mathbf{x}, \mathbf{y}) = \frac{\text{Shift}(\mathbf{x}, \mathbf{y})}{\text{AvgSpread}(\mathbf{x}, \mathbf{y})}$$

Robust effect size (shift normalized by pooled dispersion).

**Also known as** — robust Cohen’s d (Cohen (1988); estimates differ due to robust construction)

**Domain** —  $\text{AvgSpread}(\mathbf{x}, \mathbf{y}) > 0$

**Unit** — spread units

### Properties

**Location invariance**       $\text{Disparity}(\mathbf{x} + k, \mathbf{y} + k) = \text{Disparity}(\mathbf{x}, \mathbf{y})$

**Scale invariance**       $\text{Disparity}(k \cdot \mathbf{x}, k \cdot \mathbf{y}) = \text{sign}(k) \cdot \text{Disparity}(\mathbf{x}, \mathbf{y})$

**Antisymmetry**       $\text{Disparity}(\mathbf{x}, \mathbf{y}) = -\text{Disparity}(\mathbf{y}, \mathbf{x})$

### Example

$\text{Disparity}(\mathbf{x}, \mathbf{y}) = 0.4$  where  $\text{Shift} = 2, \text{AvgSpread} = 5$

$\text{Disparity}(\mathbf{x} + c, \mathbf{y} + c) = \text{Disparity}(\mathbf{x}, \mathbf{y})$      $\text{Disparity}(k\mathbf{x}, k\mathbf{y}) = \text{Disparity}(\mathbf{x}, \mathbf{y})$

Use Disparity when you want to express a difference between groups in a way that does not depend on the original measurement units. A disparity of 0.5 means the groups differ by half a spread unit; 1.0 means one full spread unit. Because it is dimensionless, you can compare effect sizes across different studies, metrics, or measurement scales. What counts as a “large” or “small” disparity depends entirely on your domain and what matters practically in your application. Do not rely on universal thresholds; interpret the number in context.

## 1.12. PairwiseMargin

`PairwiseMargin( $n, m, \text{misrate}$ )`

Exclusion count for dominance-based bounds.

**Purpose** — determines extreme pairwise differences to exclude when constructing bounds

**Based on** — distribution of  $\text{Dominance}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}(x_i > y_j)$  under random sampling

**Returns** — total margin split evenly between lower and upper tails

**Used by** — `ShiftBounds` to select appropriate order statistics

**Complexity** — exact for small samples, approximated for large (see `Fast PairwiseMargin`)

**Domain** —  $n, m \geq 1, \text{misrate} \in (0, 1)$

**Unit** — count

### Properties

**Symmetry**       $\text{PairwiseMargin}(n, m, \text{misrate}) = \text{PairwiseMargin}(m, n, \text{misrate})$

**Bounds**         $0 \leq \text{PairwiseMargin}(n, m, \text{misrate}) \leq nm$

### Example

```
PairwiseMargin(30, 30, 1e-6) = 276
```

```
PairwiseMargin(30, 30, 1e-4) = 390
```

```
PairwiseMargin(30, 30, 1e-3) = 464
```

This is a supporting function that `ShiftBounds` uses internally, so most users do not need to call it directly. It calculates how many extreme pairwise differences should be excluded when constructing bounds, based on sample sizes and the desired error rate. When you request a lower misrate (higher confidence), the margin becomes smaller, which produces wider bounds. The function automatically chooses between exact computation for small samples and a fast approximation for large samples.



### 1.13. ShiftBounds

$$\text{ShiftBounds}(\mathbf{x}, \mathbf{y}, \text{misrate}) = [z_{(k_{\text{left}})}, z_{(k_{\text{right}})}]$$

where  $\mathbf{z} = \{x_i - y_j\}$  (sorted),  $k_{\text{left}} = \left\lfloor \frac{\text{PairwiseMargin}}{2} \right\rfloor + 1$ ,  $k_{\text{right}} = nm - \left\lfloor \frac{\text{PairwiseMargin}}{2} \right\rfloor$

Robust bounds on  $\text{Shift}(\mathbf{x}, \mathbf{y})$  with specified coverage.

**Also known as** — distribution-free confidence interval for Hodges-Lehmann

**Interpretation** — misrate is probability that true shift falls outside bounds

**Domain** — any real numbers

**Unit** — same as measurements

#### Properties

**Shift invariance**       $\text{ShiftBounds}(\mathbf{x} + k, \mathbf{y} + k, \text{misrate}) = \text{ShiftBounds}(\mathbf{x}, \mathbf{y}, \text{misrate})$

**Scale equivariance**       $\text{ShiftBounds}(k \cdot \mathbf{x}, k \cdot \mathbf{y}, \text{misrate}) = k \cdot \text{ShiftBounds}(\mathbf{x}, \mathbf{y}, \text{misrate})$

#### Example

`ShiftBounds([1..30], [21..50], 1e-4) = [-30, -10]` where `Shift = -20`

Bounds fail to cover true shift with probability  $\approx$  misrate

Use `ShiftBounds` when you want to know not just the estimated shift but also how uncertain that estimate is. The function returns an interval of plausible shift values given your data. Set `misrate` to control how often the bounds might fail to contain the true shift: use  $10^{-6}$  for critical decisions where errors are costly, or  $10^{-3}$  for everyday analysis. These bounds require no assumptions about your data distribution, so they remain valid for any continuous measurements. If the bounds exclude zero, that suggests a reliable difference between the two groups.

## 2. Distributions

Distributions are parametrized random generators with well-defined statistical properties. Each distribution describes a family of random variables characterized by specific parameters.

### Notation

$X \sim \text{Additive}(0, 1)$  —  $X$  is distributed as standard normal

$\text{Estimator}(\mathbf{x})$  — estimate computed from sample

$\text{Estimator}[X]$  — true value (asymptotic limit)

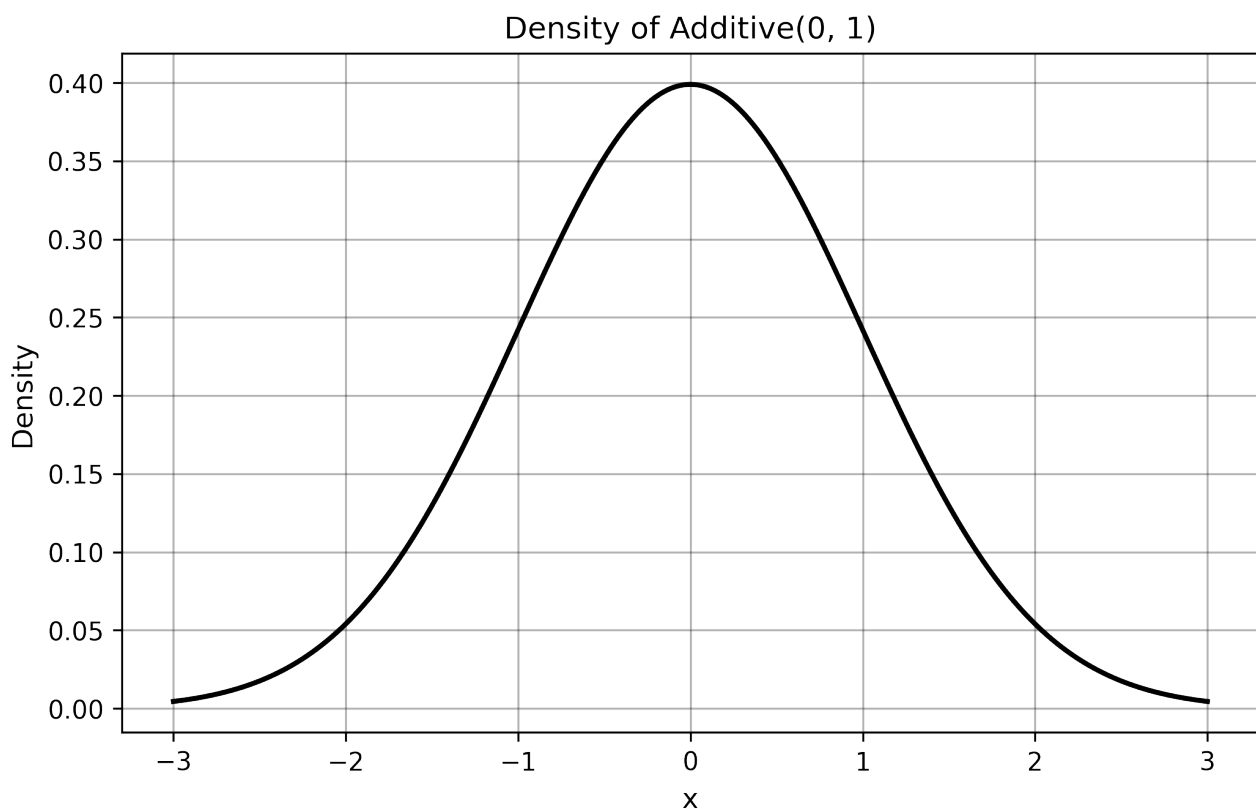
$n \rightarrow \infty$  — asymptotic case (large sample approximation)

### 2.1. Additive ('Normal')

Additive(mean, stdDev)

mean: location parameter (center of the distribution), consistent with Center

stdDev: scale parameter (standard deviation), can be rescaled to Spread



**Formation:** the sum of many variables  $X_1 + X_2 + \dots + X_n$  under mild CLT (Central Limit Theorem) conditions (e.g., Lindeberg-Feller).

**Origin:** historically called 'Normal' or 'Gaussian' distribution after Carl Friedrich Gauss and others.

**Rename Motivation:** renamed to Additive to reflect its formation mechanism through addition.

**Properties:** symmetric, bell-shaped, characterized by central limit theorem convergence.

**Applications:** measurement errors, heights and weights in populations, test scores, temperature variations.

**Characteristics:** symmetric around the mean, light tails, finite variance.

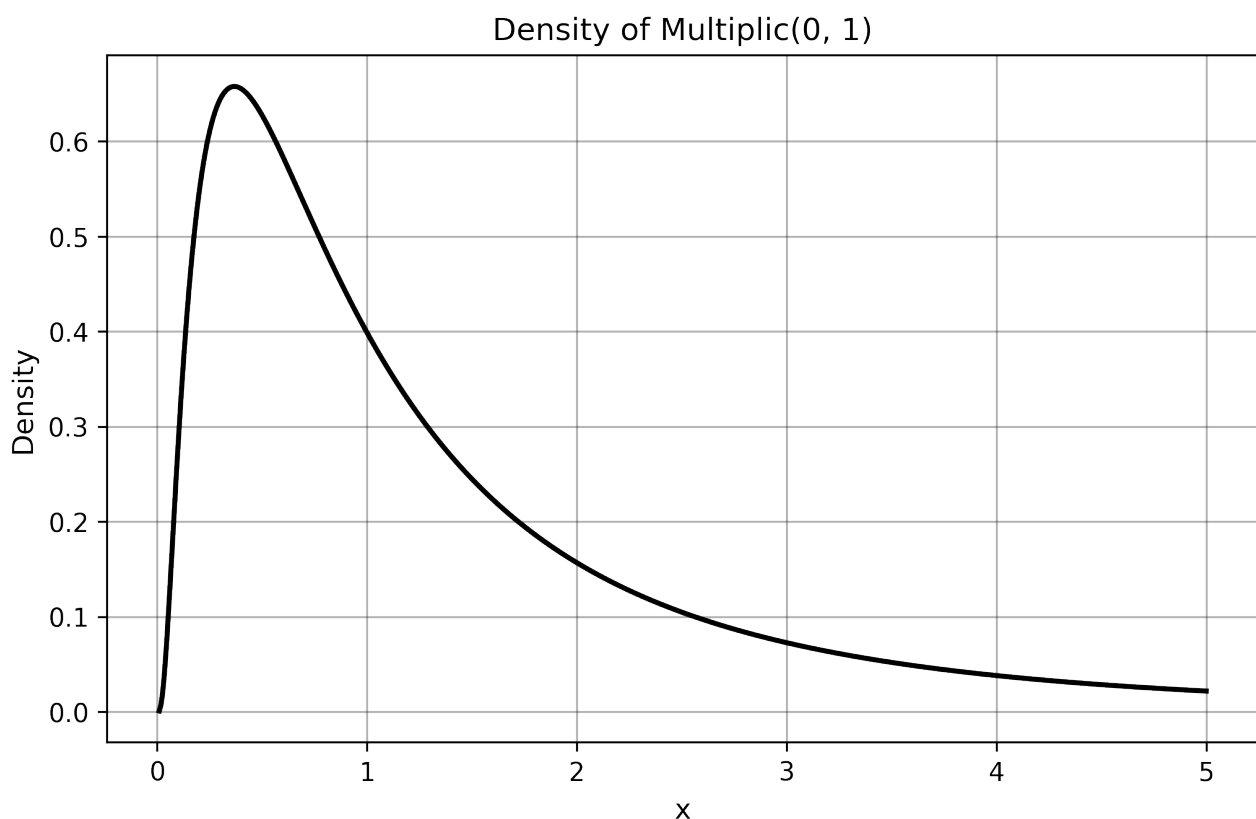
**Caution:** no perfectly additive distributions exist in real data; all real-world measurements contain deviations. Traditional estimators like Mean and StdDev lack robustness to outliers; use them only when strong evidence supports approximate additivity with no extreme measurements.

## 2.2. Multiplic ('LogNormal')

Multiplic(logMean, logStdDev)

logMean: mean of log values (location parameter;  $e^{\text{logMean}}$  equals the geometric mean)

logStdDev: standard deviation of log values (scale parameter; controls multiplicative spread)



**Formation:** the product of many positive variables  $X_1 \cdot X_2 \cdot \dots \cdot X_n$  with mild conditions (e.g., finite variance of  $\log X$ ).

**Origin:** historically called 'Log-Normal' or 'Galton' distribution after Francis Galton.

**Rename Motivation:** renamed to Multiplic to reflect its formation mechanism through multiplication.

**Properties:** logarithm of a Multiplic ('LogNormal') variable follows an Additive ('Normal') distribution.

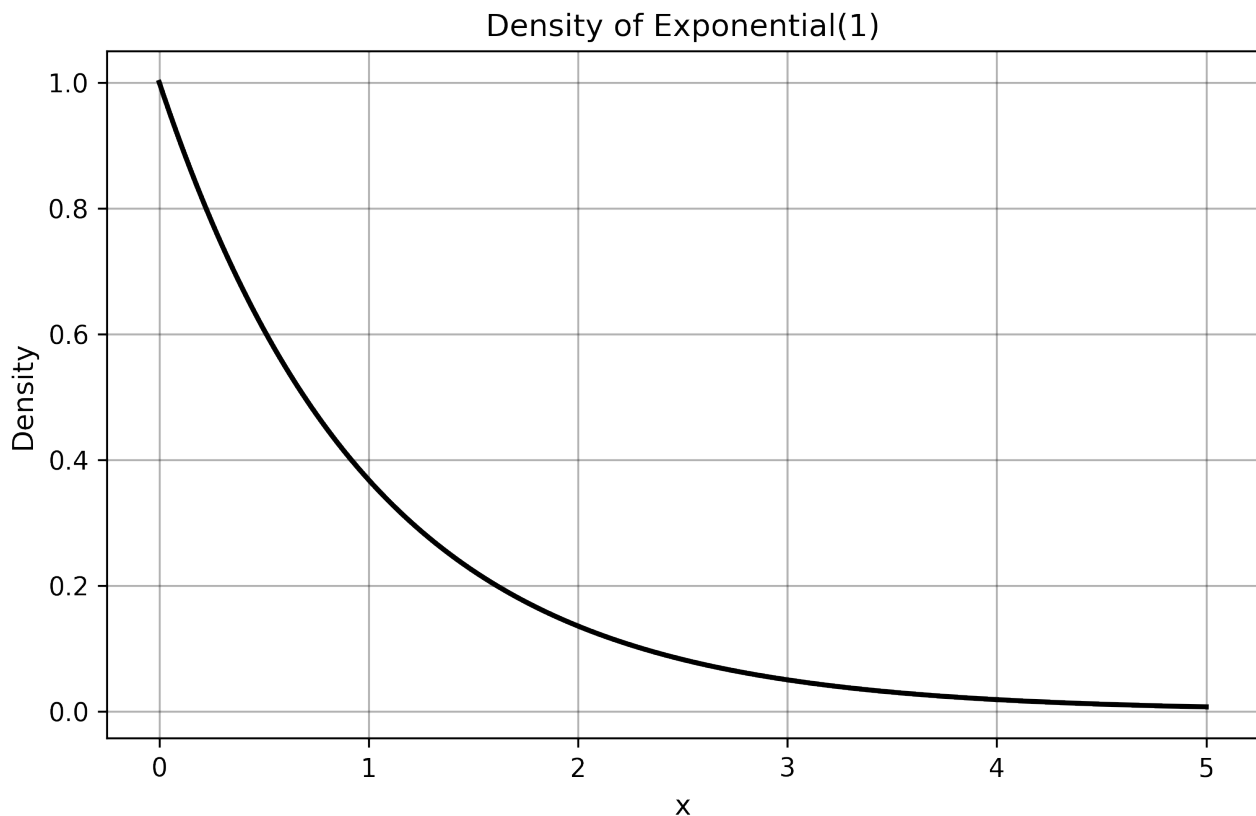
**Applications:** stock prices, file sizes, reaction times, income distributions, biological growth rates.

**Caution:** no perfectly multiplicative distributions exist in real data; all real-world measurements contain deviations. Traditional estimators struggle with the inherent skewness and heavy right tail.

## 2.3. Exponential

Exp(rate)

rate: rate parameter ( $\lambda > 0$ , controls decay speed; mean =  $1/\text{rate}$ )



**Formation:** the waiting time between events in a Poisson process.

**Origin:** naturally arises from memoryless processes where the probability of an event occurring is constant over time.

**Properties:** memoryless (past events do not affect future probabilities).

**Applications:** time between failures, waiting times in queues, radioactive decay, customer service times.

**Characteristics:** always positive, right-skewed with a light (exponential) tail.

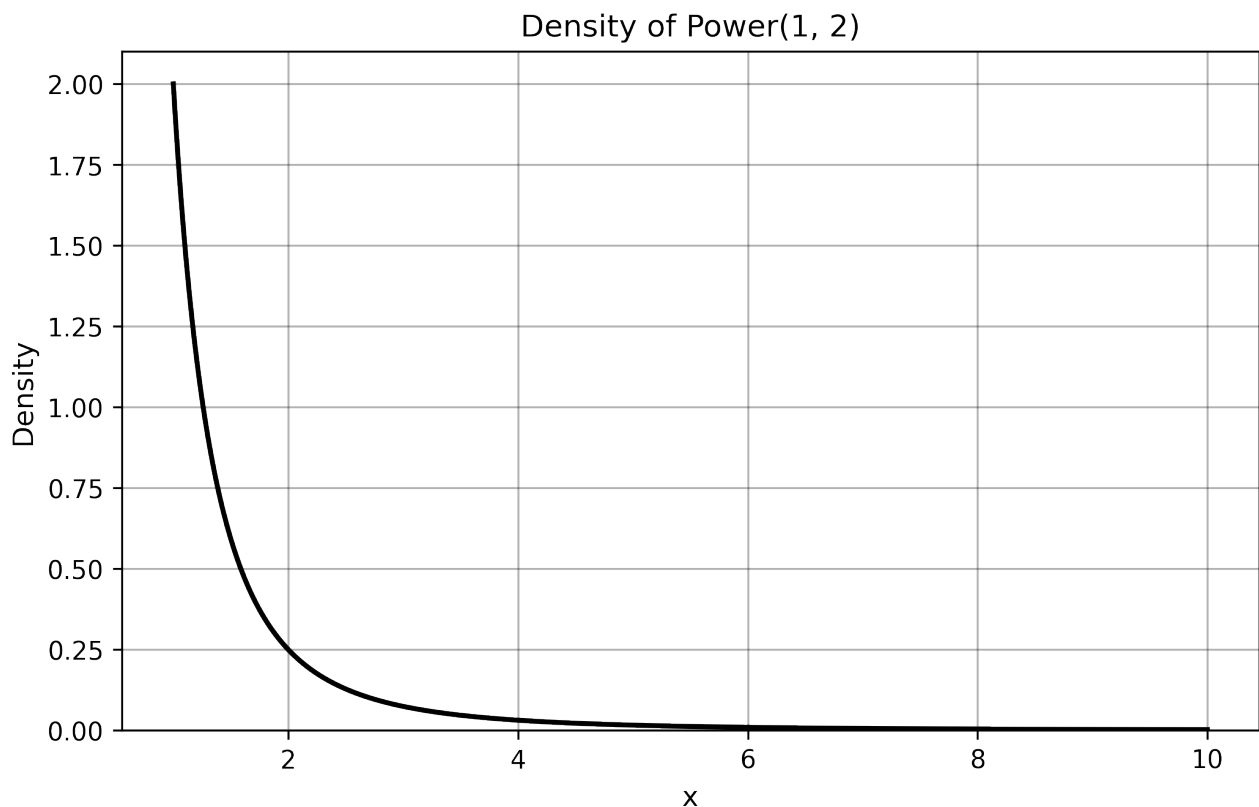
**Caution:** extreme skewness makes traditional location estimators like Mean unreliable; robust estimators provide more stable results.

## 2.4. Power ('Pareto')

Power(min, shape)

min: minimum value (lower bound,  $\text{min} > 0$ )

shape: shape parameter ( $\alpha > 0$ , controls tail heaviness; smaller values = heavier tails)



**Formation:** follows a power-law relationship where large values are rare but possible.

**Origin:** historically called 'Pareto' distribution after Vilfredo Pareto's work on wealth distribution.

**Rename Motivation:** renamed to Power to reflect its connection with power-law.

**Properties:** exhibits scale invariance and extremely heavy tails.

**Applications:** wealth distribution, city population sizes, word frequencies, earthquake magnitudes, website traffic.

**Characteristics:** infinite variance for many parameter values; extreme outliers are common.

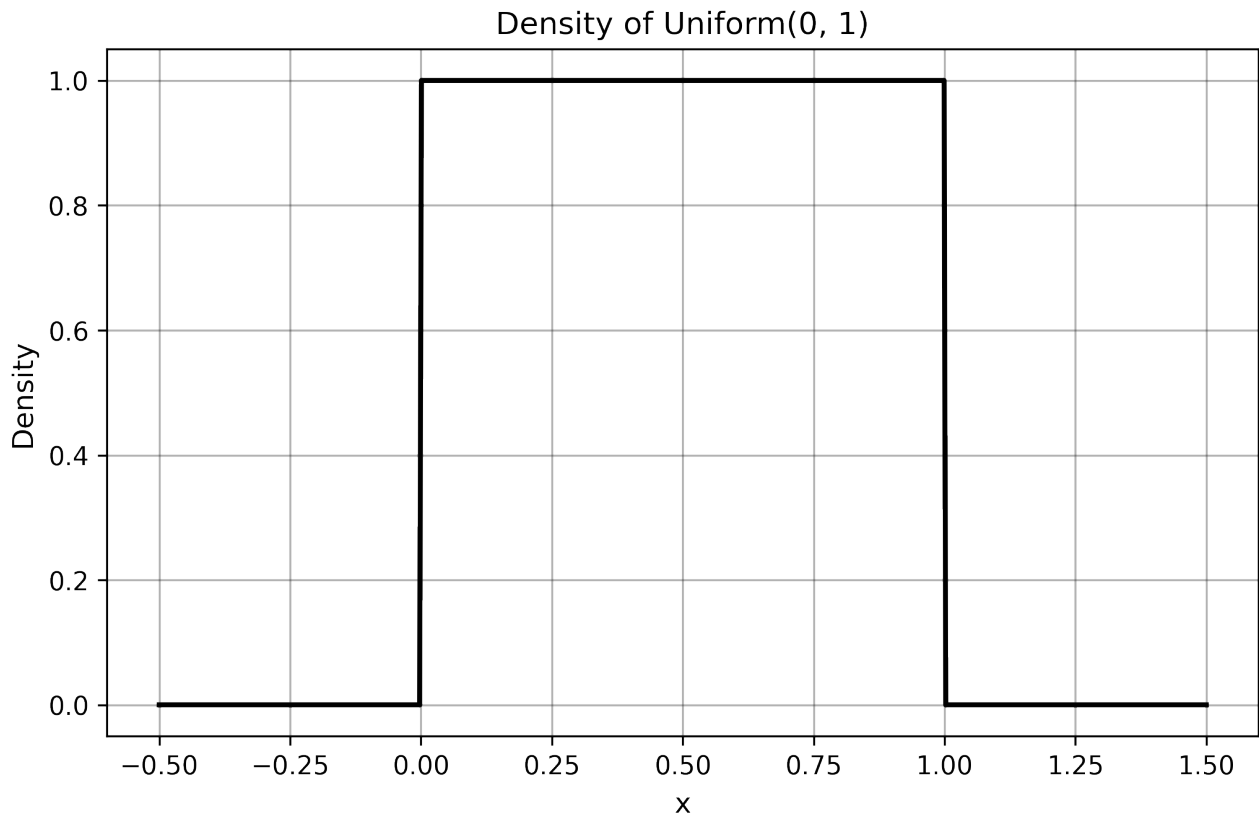
**Caution:** traditional variance-based estimators completely fail; robust estimators are essential for reliable analysis.

## 2.5. Uniform

Uniform(min, max)

min: lower bound of the support interval

max: upper bound of the support interval ( $\max > \min$ )



**Formation:** all values within a bounded interval have equal probability.

**Origin:** represents complete uncertainty within known bounds.

**Properties:** rectangular probability density, finite support with hard boundaries.

**Applications:** random number generation, round-off errors, arrival times within known intervals.

**Characteristics:** symmetric, bounded, no tail behavior.

**Note:** traditional estimators work reasonably well due to symmetry and bounded nature.

## 3. Algorithms

This chapter describes the core algorithms that power the toolkit.

### 3.1. Pseudorandom Number Generation

The toolkit provides deterministic randomization utilities that produce identical sequences across all supported programming languages. This cross-language reproducibility enables deterministic experiments and cross-validation of statistical analyses.

The core random number generator uses the xoshiro256++ algorithm (see Blackman & Vigna (2021)), a member of the xoshiro/xoroshiro family developed by David Blackman and Sebastiano Vigna. This algorithm was selected for several reasons:

**Quality:** passes all tests in the BigCrush test suite from TestU01

**Speed:** extremely fast due to simple bitwise operations (shifts, rotations, XORs)

**Period:** period of  $2^{256} - 1$ , sufficient for parallel simulations

**Adoption:** used by .NET 6+, Julia, and Rust's rand crate

The algorithm maintains a 256-bit state  $(s_0, s_1, s_2, s_3)$  and produces 64-bit outputs. Each step updates the state through a combination of XOR, shift, and rotate operations.

#### Seed Initialization

Converting a single seed value into the full 256-bit state requires a seeding algorithm. The toolkit uses SplitMix64 (see Steele et al. (2014)) for this purpose:

$$\begin{aligned} x &\leftarrow x + 0x9e3779b97f4a7c15 \\ z &\leftarrow (x \oplus (x \gg 30)) \times 0xbf58476d1ce4e5b9 \\ z &\leftarrow (z \oplus (z \gg 27)) \times 0x94d049bb133111eb \\ \text{output} &\leftarrow z \oplus (z \gg 31) \end{aligned}$$

Four consecutive outputs from SplitMix64 initialize the xoshiro256++ state  $(s_0, s_1, s_2, s_3)$ . This approach provides high-quality initial states from simple integer seeds.

#### String Seeds

For named experiments (e.g., `Rng(experiment-1)`), string seeds are converted to integers using FNV-1a hash (see Fowler et al. (1991)):

$$\begin{aligned} \text{hash} &\leftarrow 0xcbf29ce484222325 \quad (\text{offset basis}) \\ \text{for each byte } b : \quad \text{hash} &\leftarrow (\text{hash} \oplus b) \times 0x00000100000001b3 \quad (\text{FNV prime}) \end{aligned}$$

This enables meaningful experiment identifiers while maintaining determinism.

#### Uniform Floating-Point Generation

To generate uniform values in  $[0, 1)$ , the upper 53 bits of a 64-bit output are used:

$$\text{uniform}() = (\text{next}() \gg 11) \times 2^{-53}$$

The 53-bit mantissa of IEEE 754 double precision ensures all representable values in  $[0, 1)$  are reachable.

### 3.2. Shuffle Algorithm

The Shuffle function uses the Fisher-Yates algorithm (see Fisher & Yates (1938), Knuth (1997)), also known as the Knuth shuffle:

```
for i from n-1 down to 1:
    j = uniform_int(0, i+1)
    swap(array[i], array[j])
```

This produces a uniformly random permutation in  $O(n)$  time with  $O(1)$  additional space. The algorithm is unbiased: each of the  $n!$  permutations has equal probability.

### 3.3. Selection Sampling

The Sample function uses selection sampling (see Fan et al. (1962)) to select  $k$  elements from  $n$  without replacement:

```
seen = 0, selected = 0
for each element x at position i:
    if uniform() < (k - selected) / (n - seen):
        output x
        selected += 1
    seen += 1
```

This algorithm preserves the original order of elements (order of first appearance) and requires only a single pass through the data. Each element is selected independently with the correct marginal probability, producing a simple random sample.

### 3.4. Distribution Sampling

#### Box-Muller Transform

The Additive (normal) distribution uses the Box-Muller transform (see Box & Muller (1958)) to convert uniform random values to normally distributed values. Given two independent uniform values  $U_1, U_2 \in [0, 1)$ :

$$Z_0 = \sqrt{-2 \ln(U_1)} \cos(2\pi U_2)$$

$$Z_1 = \sqrt{-2 \ln(U_1)} \sin(2\pi U_2)$$

Both  $Z_0$  and  $Z_1$  are independent standard normal values. The implementation uses only  $Z_0$  to maintain cross-language determinism.

#### Numerical Constants



Distribution sampling requires handling edge cases where floating-point operations would be undefined. Two constants are used across all language implementations:

**Machine Epsilon** ( $\varepsilon_{\text{mach}}$ ): The smallest  $\varepsilon$  such that  $1 + \varepsilon \neq 1$  in float64 arithmetic.

$$\varepsilon_{\text{mach}} = 2^{-52} \approx 2.22 \times 10^{-16}$$

Used when  $U = 1$  to avoid  $\ln(0)$  or division by zero in inverse transform sampling.

**Smallest Positive Subnormal** ( $\varepsilon_{\text{sub}}$ ): The smallest positive value representable in IEEE 754 binary64.

$$\varepsilon_{\text{sub}} = 2^{-1074} \approx 4.94 \times 10^{-324}$$

Used when  $U = 0$  to avoid  $\ln(0)$  in the Box-Muller transform.

All language implementations use the same literal values for these constants (not language-specific builtins like `Number.EPSILON` or `f64::EPSILON`) to ensure bit-identical outputs across languages.

### Inverse Transform Sampling

Other distributions use inverse transform sampling. Given a uniform value  $U \in [0, 1)$  and the inverse CDF  $F^{-1}$ :

**Exponential:**  $X = -\ln(1 - U)/\lambda$

**Pareto (Power):**  $X = x_{\min}/(1 - U)^{1/\alpha}$

**Uniform:**  $X = a + U(b - a)$

The log-normal (Multiplic) distribution samples from Additive and applies the exponential transform:  $X = \exp(\text{Additive}(\mu, \sigma))$ .

## 3.5. Fast Center

The Center estimator computes the median of all pairwise averages from a sample. Given a dataset  $x = (x_1, x_2, \dots, x_n)$ , this estimator is defined as:

$$\text{Center}(\mathbf{x}) = \text{median}_{1 \leq i \leq j \leq n} \frac{x_i + x_j}{2}$$

A direct implementation would generate all  $\frac{n(n+1)}{2}$  pairwise averages and sort them. With  $n = 10000$ , this approach creates approximately 50 million values, requiring quadratic memory and  $O(n^2 \log n)$  time.

The breakthrough came in 1984 when John Monahan developed an algorithm that reduces expected complexity to  $O(n \log n)$  while using only linear memory (see Monahan (1984)). The algorithm exploits the inherent structure in pairwise sums rather than computing them explicitly. After sorting the values  $x_1 \leq x_2 \leq \dots \leq x_n$ , the algorithm considers the implicit upper triangular matrix  $M$  where  $M_{i,j} = x_i + x_j$  for  $i \leq j$ . This matrix has a crucial property: each row and column is sorted in non-decreasing order, enabling efficient median selection without storing the matrix.

Rather than sorting all pairwise sums, the algorithm uses a selection approach similar to quickselect. It maintains search bounds for each matrix row and iteratively narrows the search space. For each row  $i$ , the algorithm tracks active column indices from  $i + 1$  to  $n$ , defining which pairwise sums remain candidates for the median. It selects a candidate sum as a pivot using randomized selection from active matrix elements, then counts how many pairwise sums fall below the pivot. Because both rows and columns are sorted, this counting takes only  $O(n)$  time using a two-pointer sweep from the matrix's upper-right corner.

The median corresponds to rank  $k = \lfloor \frac{N+1}{2} \rfloor$  where  $N = \frac{n(n+1)}{2}$ . If fewer than  $k$  sums lie below the pivot, the median must be larger; if more than  $k$  sums lie below the pivot, the median must be smaller. Based on this comparison, the algorithm eliminates portions of each row that cannot contain the median, shrinking the active search space while preserving the true median.

Real data often contain repeated values, which can cause the selection process to stall. When the algorithm detects no progress between iterations, it switches to a midrange strategy: find the smallest and largest pairwise sums still in the search space, then use their average as the next pivot. If the minimum equals the maximum, all remaining candidates are identical and the algorithm terminates. This tie-breaking mechanism ensures reliable convergence with discrete or duplicated data.

The algorithm achieves  $O(n \log n)$  time complexity through linear partitioning (each pivot evaluation requires only  $O(n)$  operations) and logarithmic iterations (randomized pivot selection leads to expected  $O(\log n)$  iterations, similar to quickselect). The algorithm maintains only row bounds and counters, using  $O(n)$  additional space. This matches the complexity of sorting a single array while avoiding the quadratic memory and time explosion of computing all pairwise combinations.

```
namespace Pragmastat.Algorithms;
```

```
internal static class FastCenter
```

```
{
```

```
    /// <summary>
```

```
    /// ACM Algorithm 616: fast computation of the Hodges-Lehmann location estimator
```

```
    /// </summary>
```

```
    /// <remarks>
```

```
    /// Computes the median of all pairwise averages (xi + xj)/2 efficiently.
```

```
    /// See: John F Monahan, "Algorithm 616: fast computation of the Hodges-Lehmann location estimator"
```

```
    /// (1984) DOI: 10.1145/1271.319414
```

```
    /// </remarks>
```

```
    /// <param name="values">A sorted sample of values</param>
```

```
    /// <param name="random">Random number generator</param>
```

```
    /// <param name="isSorted">If values are sorted</param>
```

```
    /// <returns>Exact center value (Hodges-Lehmann estimator)</returns>
```

```
    public static double Estimate(IReadOnlyList<double> values, Random? random = null, bool isSorted = false)
```

```
    {
```

```
        int n = values.Count;
```

```
        if (n == 1) return values[0];
```

```
        if (n == 2) return (values[0] + values[1]) / 2;
```

```
        random ??= new Random();
```

```
        if (!isSorted)
```

```

    values = values.OrderBy(x => x).ToList();

    // Calculate target median rank(s) among all pairwise sums
    long totalPairs = (long)n * (n + 1) / 2;
    long medianRankLow = (totalPairs + 1) / 2; // For odd totalPairs, this is the median
    long medianRankHigh =
        (totalPairs + 2) / 2; // For even totalPairs, average of ranks medianRankLow and
medianRankHigh

    // Initialize search bounds for each row in the implicit matrix
    long[] leftBounds = new long[n];
    long[] rightBounds = new long[n];
    long[] partitionCounts = new long[n];

    for (int i = 0; i < n; i++)
    {
        leftBounds[i] = i + 1; // Row i can pair with columns [i+1..n] (1-based indexing)
        rightBounds[i] = n; // Initially, all columns are available
    }

    // Start with a good pivot: sum of middle elements (handles both odd and even n)
    double pivot = values[(n - 1) / 2] + values[n / 2];
    long activeSetSize = totalPairs;
    long previousCount = 0;

    while (true)
    {
        // === PARTITION STEP ===
        // Count pairwise sums less than current pivot
        long countBelowPivot = 0;
        long currentColumn = n;

        for (int row = 1; row <= n; row++)
        {
            partitionCounts[row - 1] = 0;

            // Move left from current column until we find sums < pivot
            // This exploits the sorted nature of the matrix
            while (currentColumn >= row && values[row - 1] + values[(int)currentColumn - 1] >=
pivot)
                currentColumn--;

            // Count elements in this row that are < pivot
            if (currentColumn >= row)
            {
                long elementsBelow = currentColumn - row + 1;
                partitionCounts[row - 1] = elementsBelow;
                countBelowPivot += elementsBelow;
            }
        }

        // === CONVERGENCE CHECK ===
        // If no progress, we have ties - break them using midrange strategy

```

```

    if (countBelowPivot == previousCount)
    {
        double minActiveSum = double.MaxValue;
        double maxActiveSum = double.MinValue;

        // Find the range of sums still in the active search space
        for (int i = 0; i < n; i++)
        {
            if (leftBounds[i] > rightBounds[i]) continue; // Skip empty rows

            double rowValue = values[i];
            double smallestInRow = values[(int)leftBounds[i] - 1] + rowValue;
            double largestInRow = values[(int)rightBounds[i] - 1] + rowValue;

            minActiveSum = Min(minActiveSum, smallestInRow);
            maxActiveSum = Max(maxActiveSum, largestInRow);
        }

        pivot = (minActiveSum + maxActiveSum) / 2;
        if (pivot <= minActiveSum || pivot > maxActiveSum) pivot = maxActiveSum;

        // If all remaining values are identical, we're done
        if (minActiveSum == maxActiveSum || activeSetSize <= 2)
            return pivot / 2;

        continue;
    }

    // === TARGET CHECK ===
    // Check if we've found the median rank(s)
    bool atTargetRank = countBelowPivot == medianRankLow || countBelowPivot ==
medianRankHigh - 1;
    if (atTargetRank)
    {
        // Find the boundary values: largest < pivot and smallest >= pivot
        double largestBelowPivot = double.MinValue;
        double smallestAtOrAbovePivot = double.MaxValue;

        for (int i = 1; i <= n; i++)
        {
            long countInRow = partitionCounts[i - 1];
            double rowValue = values[i - 1];
            long totalInRow = n - i + 1;

            // Find largest sum in this row that's < pivot
            if (countInRow > 0)
            {
                long lastBelowIndex = i + countInRow - 1;
                double lastBelowValue = rowValue + values[(int)lastBelowIndex - 1];
                largestBelowPivot = Max(largestBelowPivot, lastBelowValue);
            }

            // Find smallest sum in this row that's >= pivot

```

```

        if (countInRow < totalInRow)
        {
            long firstAtOrAboveIndex = i + countInRow;
            double firstAtOrAboveValue = rowValue + values[(int)firstAtOrAboveIndex - 1];
            smallestAtOrAbovePivot = Min(smallestAtOrAbovePivot, firstAtOrAboveValue);
        }
    }

    // Calculate final result based on whether we have odd or even number of pairs
    if (medianRankLow < medianRankHigh)
    {
        // Even total: average the two middle values
        return (smallestAtOrAbovePivot + largestBelowPivot) / 4;
    }
    else
    {
        // Odd total: return the single middle value
        bool needLargest = countBelowPivot == medianRankLow;
        return (needLargest ? largestBelowPivot : smallestAtOrAbovePivot) / 2;
    }
}

// === UPDATE BOUNDS ===
// Narrow the search space based on partition result
if (countBelowPivot < medianRankLow)
{
    // Too few values below pivot - eliminate smaller values, search higher
    for (int i = 0; i < n; i++)
        leftBounds[i] = i + partitionCounts[i] + 1;
}
else
{
    // Too many values below pivot - eliminate larger values, search lower
    for (int i = 0; i < n; i++)
        rightBounds[i] = i + partitionCounts[i];
}

// === PREPARE NEXT ITERATION ===
previousCount = countBelowPivot;

// Recalculate how many elements remain in the active search space
activeSetSize = 0;
for (int i = 0; i < n; i++)
{
    long rowSize = rightBounds[i] - leftBounds[i] + 1;
    activeSetSize += Max(0, rowSize);
}

// Choose next pivot based on remaining active set size
if (activeSetSize > 2)
{
    // Use randomized row median strategy for efficiency
    // Handle large activeSetSize by using double precision for random selection

```

```

    double randomFraction = random.NextDouble();
    long targetIndex = (long)(randomFraction * activeSetSize);
    int selectedRow = 0;

    // Find which row contains the target index
    long cumulativeSize = 0;
    for (int i = 0; i < n; i++)
    {
        long rowSize = Max(0, rightBounds[i] - leftBounds[i] + 1);
        if (targetIndex < cumulativeSize + rowSize)
        {
            selectedRow = i;
            break;
        }

        cumulativeSize += rowSize;
    }

    // Use median element of the selected row as pivot
    long medianColumnInRow = (leftBounds[selectedRow] + rightBounds[selectedRow]) / 2;
    pivot = values[selectedRow] + values[(int)medianColumnInRow - 1];
}
else
{
    // Few elements remain - use midrange strategy
    double minRemainingSum = double.MaxValue;
    double maxRemainingSum = double.MinValue;

    for (int i = 0; i < n; i++)
    {
        if (leftBounds[i] > rightBounds[i]) continue; // Skip empty rows

        double rowValue = values[i];
        double minInRow = values[(int)leftBounds[i] - 1] + rowValue;
        double maxInRow = values[(int)rightBounds[i] - 1] + rowValue;

        minRemainingSum = Min(minRemainingSum, minInRow);
        maxRemainingSum = Max(maxRemainingSum, maxInRow);
    }

    pivot = (minRemainingSum + maxRemainingSum) / 2;
    if (pivot <= minRemainingSum || pivot > maxRemainingSum)
        pivot = maxRemainingSum;

    if (minRemainingSum == maxRemainingSum)
        return pivot / 2;
}
}
}
}

```

### 3.6. Fast Spread

The Spread estimator computes the median of all pairwise absolute differences. Given a sample  $x = (x_1, x_2, \dots, x_n)$ , this estimator is defined as:

$$\text{Spread}(\mathbf{x}) = \text{median}_{1 \leq i < j \leq n} |x_i - x_j|$$

Like Center, computing Spread naively requires generating all  $\frac{n(n-1)}{2}$  pairwise differences, sorting them, and finding the median — a quadratic approach that becomes computationally prohibitive for large datasets.

The same structural principles that accelerate Center computation also apply to pairwise differences, yielding an exact  $O(n \log n)$  algorithm. After sorting the input to obtain  $y_1 \leq y_2 \leq \dots \leq y_n$ , all pairwise absolute differences  $|x_i - x_j|$  with  $i < j$  become positive differences  $y_j - y_i$ . This allows considering the implicit upper triangular matrix  $D$  where  $D_{i,j} = y_j - y_i$  for  $i < j$ . This matrix has a crucial structural property: for a fixed row  $i$ , differences increase monotonically, while for a fixed column  $j$ , differences decrease as  $i$  increases. This sorted structure enables linear-time counting of elements below any threshold.

The algorithm applies Monahan’s selection strategy (Monahan (1984)), adapted for differences rather than sums. For each row  $i$ , it tracks active column indices representing differences still under consideration, initially spanning columns  $i + 1$  through  $n$ . It chooses candidate differences from the active set using weighted random row selection, maintaining expected logarithmic convergence while avoiding expensive pivot computations. For any pivot value  $p$ , the number of differences falling below  $p$  is counted using a single sweep; the monotonic structure ensures this counting requires only  $O(n)$  operations. While counting, the largest difference below  $p$  and the smallest difference at or above  $p$  are maintained — these boundary values become the exact answer when the target rank is reached.

The algorithm naturally handles both odd and even cases. For an odd number of differences, it returns the single middle element when the count exactly hits the median rank. For an even number of differences, it returns the average of the two middle elements; boundary tracking during counting provides both values simultaneously. Unlike approximation methods, this algorithm returns the precise median of all pairwise differences; randomness affects only performance, not correctness.

The algorithm includes the same stall-handling mechanisms as the center algorithm. It tracks whether the count below the pivot changes between iterations; when progress stalls due to tied values, it computes the range of remaining active differences and pivots to their midrange. This midrange strategy ensures convergence even with highly discrete data or datasets with many identical values.

Several optimizations make the implementation practical for production use. A global column pointer that never moves backward during counting exploits the matrix structure to avoid redundant comparisons. Exact boundary values are captured during each counting pass, eliminating the need for additional searches when the target rank is reached. Using only  $O(n)$  additional space for row bounds and counters, the algorithm achieves  $O(n \log n)$  time complexity with minimal memory overhead, making robust scale estimation practical for large datasets.

```
namespace Pragmestat.Algorithms;
```

```

internal static class FastSpread
{
    /// <summary>
    /// Shamos "Spread". Expected  $O(n \log n)$  time,  $O(n)$  extra space. Exact.
    /// </summary>
    public static double Estimate(IReadOnlyList<double> values, Random? random = null, bool
isSorted = false)
    {
        int n = values.Count;
        if (n <= 1) return 0;
        if (n == 2) return Abs(values[1] - values[0]);
        random ??= new Random();

        // Prepare a sorted working copy.
        double[] a = isSorted ? CopySorted(values) : EnsureSorted(values);

        // Total number of pairwise differences with  $i < j$ 
        long N = (long)n * (n - 1) / 2;
        long kLow = (N + 1) / 2; // 1-based rank of lower middle
        long kHigh = (N + 2) / 2; // 1-based rank of upper middle

        // Per-row active bounds over columns j (0-based indices).
        // Row i allows j in [i+1, n-1] initially.
        int[] L = new int[n];
        int[] R = new int[n];
        long[] rowCounts = new long[n]; // # of elements in row i that are < pivot (for
current partition)

        for (int i = 0; i < n; i++)
        {
            L[i] = Min(i + 1, n); // n means empty
            R[i] = n - 1; // inclusive
            if (L[i] > R[i])
            {
                L[i] = 1;
                R[i] = 0;
            } // mark empty
        }

        // A reasonable initial pivot: a central gap
        double pivot = a[n / 2] - a[(n - 1) / 2];

        long prevCountBelow = -1;

        while (true)
        {
            // === PARTITION: count how many differences are < pivot; also track boundary
neighbors ===
            long countBelow = 0;
            double largestBelow = double.NegativeInfinity; // max difference < pivot
            double smallestAtOrAbove = double.PositiveInfinity; // min difference >= pivot

            int j = 1; // global two-pointer (non-decreasing across rows)

```



```

for (int i = 0; i < n - 1; i++)
{
    if (j < i + 1) j = i + 1;
    while (j < n && a[j] - a[i] < pivot) j++;

    long cntRow = j - (i + 1);
    if (cntRow < 0) cntRow = 0;
    rowCounts[i] = cntRow;
    countBelow += cntRow;

    // boundary elements for this row
    if (cntRow > 0)
    {
        // last < pivot in this row is (j-1)
        double candBelow = a[j - 1] - a[i];
        if (candBelow > largestBelow) largestBelow = candBelow;
    }

    if (j < n)
    {
        double candAtOrAbove = a[j] - a[i];
        if (candAtOrAbove < smallestAtOrAbove) smallestAtOrAbove = candAtOrAbove;
    }
}

// === TARGET CHECK ===
// If we've split exactly at the middle, we can return using the boundaries we just
found.
bool atTarget =
    (countBelow == kLow) || // lower middle is the largest < pivot
    (countBelow == (kHigh - 1)); // upper middle is the smallest >= pivot

if (atTarget)
{
    if (kLow < kHigh)
    {
        // Even N: average the two central order stats.
        return 0.5 * (largestBelow + smallestAtOrAbove);
    }
    else
    {
        // Odd N: pick the single middle depending on which side we hit.
        bool needLargest = (countBelow == kLow);
        return needLargest ? largestBelow : smallestAtOrAbove;
    }
}

// === STALL HANDLING (ties / no progress) ===
if (countBelow == prevCountBelow)
{
    // Compute min/max remaining difference in the ACTIVE set and pivot to their
midrange.
    double minActive = double.PositiveInfinity;

```

```

double maxActive = double.NegativeInfinity;
long active = 0;

for (int i = 0; i < n - 1; i++)
{
    int Li = L[i], Ri = R[i];
    if (Li > Ri) continue;

    double rowMin = a[Li] - a[i];
    double rowMax = a[Ri] - a[i];
    if (rowMin < minActive) minActive = rowMin;
    if (rowMax > maxActive) maxActive = rowMax;
    active += (Ri - Li + 1);
}

if (active <= 0)
{
    // No active candidates left: the only consistent answer is the boundary implied
    by counts.
    // Fall back to neighbors from this partition.
    if (kLow < kHigh) return 0.5 * (largestBelow + smallestAtOrAbove);
    return (countBelow >= kLow) ? largestBelow : smallestAtOrAbove;
}

if (maxActive <= minActive) return minActive; // all remaining equal

double mid = 0.5 * (minActive + maxActive);
pivot = (mid > minActive && mid <= maxActive) ? mid : maxActive;
prevCountBelow = countBelow;
continue;
}

// === SHRINK ACTIVE WINDOW ===
// --- SHRINK ACTIVE WINDOW (fixed) ---
if (countBelow < kLow)
{
    // Need larger differences: discard all strictly below pivot.
    for (int i = 0; i < n - 1; i++)
    {
        // First j with a[j] - a[i] >= pivot is j = i + 1 + cntRow (may be n => empty
row)

        int newL = i + 1 + (int)rowCounts[i];
        if (newL > L[i]) L[i] = newL; // do NOT clamp; allow L[i] == n to mean empty
        if (L[i] > R[i])
        {
            L[i] = 1;
            R[i] = 0;
        } // mark empty
    }
}
else
{
    // Too many below: keep only those strictly below pivot.

```

```

    for (int i = 0; i < n - 1; i++)
    {
        // Last j with a[j] - a[i] < pivot is j = i + cntRow (not cntRow-1!)
        int newR = i + (int)rowCounts[i];
        if (newR < R[i]) R[i] = newR; // shrink downward to the true last-below
        if (R[i] < i + 1)
        {
            L[i] = 1;
            R[i] = 0;
        } // empty row if none remain
    }
}

prevCountBelow = countBelow;

// === CHOOSE NEXT PIVOT FROM ACTIVE SET (weighted random row, then row median) ===
long activeSize = 0;
for (int i = 0; i < n - 1; i++)
{
    if (L[i] <= R[i]) activeSize += (R[i] - L[i] + 1);
}

if (activeSize <= 2)
{
    // Few candidates left: return midrange of remaining exactly.
    double minRem = double.PositiveInfinity, maxRem = double.NegativeInfinity;
    for (int i = 0; i < n - 1; i++)
    {
        if (L[i] > R[i]) continue;
        double lo = a[L[i]] - a[i];
        double hi = a[R[i]] - a[i];
        if (lo < minRem) minRem = lo;
        if (hi > maxRem) maxRem = hi;
    }

    if (activeSize <= 0) // safety net; fall back to boundary from last partition
    {
        if (kLow < kHigh) return 0.5 * (largestBelow + smallestAtOrAbove);
        return (countBelow >= kLow) ? largestBelow : smallestAtOrAbove;
    }

    if (kLow < kHigh) return 0.5 * (minRem + maxRem);
    return (Abs((kLow - 1) - countBelow) <= Abs(countBelow - kLow)) ? minRem : maxRem;
}
else
{
    long t = NextIndex(random, activeSize); // 0..activeSize-1
    long acc = 0;
    int row = 0;
    for (; row < n - 1; row++)
    {
        if (L[row] > R[row]) continue;
        long size = R[row] - L[row] + 1;

```

```

        if (t < acc + size) break;
        acc += size;
    }

    // Median column of the selected row
    int col = (L[row] + R[row]) >> 1;
    pivot = a[col] - a[row];
    }
}
}
// --- Helpers ---

private static double[] CopySorted(IReadOnlyList<double> values)
{
    var a = new double[values.Count];
    for (int i = 0; i < a.Length; i++)
    {
        double v = values[i];
        if (double.IsNaN(v)) throw new ArgumentException("NaN not allowed.",
nameof(values));
        a[i] = v;
    }

    Array.Sort(a);
    return a;
}

private static double[] EnsureSorted(IReadOnlyList<double> values)
{
    // Trust caller; still copy to array for fast indexed access.
    var a = new double[values.Count];
    for (int i = 0; i < a.Length; i++)
    {
        double v = values[i];
        if (double.IsNaN(v)) throw new ArgumentException("NaN not allowed.",
nameof(values));
        a[i] = v;
    }

    return a;
}

private static long NextIndex(Random rng, long limitExclusive)
{
    // Uniform 0..limitExclusive-1 even for large ranges.
    // Use rejection sampling for correctness.
    ulong uLimit = (ulong)limitExclusive;
    if (uLimit <= int.MaxValue)
    {
        return rng.Next((int)uLimit);
    }

    while (true)

```

```

{
    ulong u = ((ulong)(uint)rng.Next() << 32) | (uint)rng.Next();
    ulong r = u % uLimit;
    if (u - r <= ulong.MaxValue - (ulong.MaxValue % uLimit)) return (long)r;
}
}
}

```

### 3.7. Fast Shift

The Shift estimator measures the median of all pairwise differences between elements of two samples. Given samples  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  and  $\mathbf{y} = (y_1, y_2, \dots, y_m)$ , this estimator is defined as:

$$\text{Shift}(\mathbf{x}, \mathbf{y}) = \text{median}_{1 \leq i \leq n, 1 \leq j \leq m} (x_i - y_j)$$

This definition represents a special case of a more general problem: computing arbitrary quantiles of all pairwise differences. For samples of size  $n$  and  $m$ , the total number of pairwise differences is  $n \times m$ . A naive approach would materialize all differences, sort them, and extract the desired quantile. With  $n = m = 10000$ , this approach creates 100 million values, requiring quadratic memory and  $O(nm \log(nm))$  time.

The presented algorithm avoids materializing pairwise differences by exploiting the sorted structure of the samples. After sorting both samples to obtain  $x_1 \leq x_2 \leq \dots \leq x_n$  and  $y_1 \leq y_2 \leq \dots \leq y_m$ , the key insight is that it's possible to count how many pairwise differences fall below any threshold without computing them explicitly. This counting operation enables a binary search over the continuous space of possible difference values, iteratively narrowing the search range until it converges to the exact quantile.

The algorithm operates through a value-space search rather than index-space selection. It maintains a search interval  $[\text{searchMin}, \text{searchMax}]$ , initialized to the range of all possible differences:  $[x_1 - y_m, x_n - y_1]$ . At each iteration, it selects a candidate value within this interval and counts how many pairwise differences are less than or equal to this threshold. For the median (quantile  $p = 0.5$ ), if fewer than half the differences lie below the threshold, the median must be larger; if more than half lie below, the median must be smaller. Based on this comparison, the search space is reduced by eliminating portions that cannot contain the target quantile.

The counting operation achieves linear complexity via a two-pointer sweep. For a given threshold  $t$ , the number of pairs  $(i, j)$  satisfying  $x_i - y_j \leq t$  is counted. This is equivalent to counting pairs where  $y_j \geq x_i - t$ . For each row  $i$  in the implicit matrix of differences, a column pointer advances through the sorted  $y$  array while  $x_i - y_j > t$ , stopping at the first position where  $x_i - y_j \leq t$ . All subsequent positions in that row satisfy the condition, contributing  $(m - j)$  pairs to the count for row  $i$ . Because both samples are sorted, the column pointer advances monotonically and never backtracks across rows, making each counting pass  $O(n + m)$  regardless of the total number of differences.

During each counting pass, the algorithm tracks boundary values: the largest difference at or below the threshold and the smallest difference above it. When the count exactly matches the target rank (or the two middle ranks for even-length samples), these boundary values provide the exact answer without additional searches. For Type-7 quantile computation (Hyndman & Fan (1996)), which inter-

polates between order statistics, the algorithm collects the necessary boundary values in a single pass and performs linear interpolation:  $(1 - w) \cdot \text{lower} + w \cdot \text{upper}$ .

Real datasets often contain discrete or repeated values that can cause search stagnation. The algorithm detects when the search interval stops shrinking between iterations, indicating that multiple pairwise differences share the same value. When the closest difference below the threshold equals the closest above, all remaining candidates are identical and the algorithm terminates immediately. Otherwise, it uses the boundary values from the counting pass to snap the search interval to actual difference values, ensuring reliable convergence even with highly discrete data.

The binary search employs numerically stable midpoint calculations and terminates when the search interval collapses to a single value or when boundary tracking confirms convergence. Iteration limits are included as a safety mechanism, though convergence typically occurs much earlier due to the exponential narrowing of the search space.

The algorithm naturally generalizes to multiple quantiles by computing each one independently. For  $k$  quantiles with samples of size  $n$  and  $m$ , the total complexity is  $O(k(n + m) \log L)$ , where  $L$  represents the convergence precision. This is dramatically more efficient than the naive  $O(nm \log(nm))$  approach, especially for large  $n$  and  $m$  with small  $k$ . The algorithm requires only  $O(1)$  additional space beyond the input arrays, making it practical for large-scale statistical analysis where memory constraints prohibit materializing quadratic data structures.

```
namespace Pragmatstat.Algorithms;
```

```
using System;
```

```
using System.Collections.Generic;
```

```
using System.Linq;
```

```
public static class FastShift
```

```
{
```

```
    /// <summary>
```

```
    /// Computes quantiles of all pairwise differences { x_i - y_j }.
```

```
    /// Time:  $O((m + n) \cdot \log(\text{precision}))$  per quantile. Space:  $O(1)$ .
```

```
    /// </summary>
```

```
    /// <param name="p">Probabilities in [0, 1].</param>
```

```
    /// <param name="assumeSorted">If false, collections will be sorted.</param>
```

```
    public static double[] Estimate(IReadOnlyList<double> x, IReadOnlyList<double> y,
    double[] p, bool assumeSorted = false)
```

```
    {
```

```
        if (x == null || y == null || p == null)
```

```
            throw new ArgumentNullException();
```

```
        if (x.Count == 0 || y.Count == 0)
```

```
            throw new ArgumentException("x and y must be non-empty.");
```

```
        foreach (double pk in p)
```

```
            if (double.IsNaN(pk) || pk < 0.0 || pk > 1.0)
```

```
                throw new ArgumentOutOfRangeException(nameof(p), "Probabilities must be within [0, 1].");
```

```
        double[] xs, ys;
```

```

    if (assumeSorted)
    {
        xs = x as double[] ?? x.ToArray();
        ys = y as double[] ?? y.ToArray();
    }
    else
    {
        xs = x.OrderBy(v => v).ToArray();
        ys = y.OrderBy(v => v).ToArray();
    }

    int m = xs.Length;
    int n = ys.Length;
    long total = (long)m * n;

    // Type-7 quantile:  $h = 1 + (n-1)*p$ , then interpolate between floor(h) and ceil(h)
    var requiredRanks = new SortedSet<long>();
    var interpolationParams = new (long lowerRank, long upperRank, double weight)
[p.Length];

    for (int i = 0; i < p.Length; i++)
    {
        double h = 1.0 + (total - 1) * p[i];
        long lowerRank = (long)Math.Floor(h);
        long upperRank = (long)Math.Ceiling(h);
        double weight = h - lowerRank;
        if (lowerRank < 1) lowerRank = 1;
        if (upperRank > total) upperRank = total;
        interpolationParams[i] = (lowerRank, upperRank, weight);
        requiredRanks.Add(lowerRank);
        requiredRanks.Add(upperRank);
    }

    var rankValues = new Dictionary<long, double>();
    foreach (long rank in requiredRanks)
        rankValues[rank] = SelectKthPairwiseDiff(xs, ys, rank);

    var result = new double[p.Length];
    for (int i = 0; i < p.Length; i++)
    {
        var (lowerRank, upperRank, weight) = interpolationParams[i];
        double lower = rankValues[lowerRank];
        double upper = rankValues[upperRank];
        result[i] = weight == 0.0 ? lower : (1.0 - weight) * lower + weight * upper;
    }

    return result;
}

// Binary search in [min_diff, max_diff] that snaps to actual discrete values.
// Avoids materializing all m*n differences.
private static double SelectKthPairwiseDiff(double[] x, double[] y, long k)
{

```

```

    int m = x.Length;
    int n = y.Length;
    long total = (long)m * n;

    if (k < 1 || k > total)
        throw new ArgumentOutOfRangeException(nameof(k));

    double searchMin = x[0] - y[n - 1];
    double searchMax = x[m - 1] - y[0];

    if (double.IsNaN(searchMin) || double.IsNaN(searchMax))
        throw new InvalidOperationException("NaN in input values.");

    const int maxIterations = 128; // Sufficient for double precision convergence
    double prevMin = double.NegativeInfinity;
    double prevMax = double.PositiveInfinity;

    for (int iter = 0; iter < maxIterations && searchMin != searchMax; iter++)
    {
        double mid = Midpoint(searchMin, searchMax);
        CountAndNeighbors(x, y, mid, out long countLessOrEqual, out double closestBelow, out
double closestAbove);

        if (closestBelow == closestAbove)
            return closestBelow;

        // No progress means we're stuck between two discrete values
        if (searchMin == prevMin && searchMax == prevMax)
            return countLessOrEqual >= k ? closestBelow : closestAbove;

        prevMin = searchMin;
        prevMax = searchMax;

        if (countLessOrEqual >= k)
            searchMax = closestBelow;
        else
            searchMin = closestAbove;
    }

    if (searchMin != searchMax)
        throw new InvalidOperationException("Convergence failure (pathological input).");

    return searchMin;
}

// Two-pointer algorithm: counts pairs where x[i] - y[j] <= threshold, and tracks
// the closest actual differences on either side of threshold.
private static void CountAndNeighbors(
    double[] x, double[] y, double threshold,
    out long countLessOrEqual, out double closestBelow, out double closestAbove)
{
    int m = x.Length, n = y.Length;
    long count = 0;

```



```

double maxBelow = double.NegativeInfinity;
double minAbove = double.PositiveInfinity;

int j = 0;
for (int i = 0; i < m; i++)
{
    while (j < n && x[i] - y[j] > threshold)
        j++;

    count += (n - j);

    if (j < n)
    {
        double diff = x[i] - y[j];
        if (diff > maxBelow) maxBelow = diff;
    }

    if (j > 0)
    {
        double diff = x[i] - y[j - 1];
        if (diff < minAbove) minAbove = diff;
    }
}

// Fallback to actual min/max if no boundaries found (shouldn't happen in normal
operation)
if (double.IsNegativeInfinity(maxBelow))
    maxBelow = x[0] - y[n - 1];
if (double.IsPositiveInfinity(minAbove))
    minAbove = x[m - 1] - y[0];

countLessOrEqual = count;
closestBelow = maxBelow;
closestAbove = minAbove;
}

private static double Midpoint(double a, double b) => a + (b - a) * 0.5;
}

```

### 3.8. Fast PairwiseMargin

The PairwiseMargin function determines how many extreme pairwise differences to exclude when constructing bounds around  $\text{Shift}(\mathbf{x}, \mathbf{y})$ . Given samples  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_m)$ , the ShiftBounds estimator computes all  $nm$  pairwise differences  $z_{ij} = x_i - y_j$  and sorts them. The bounds select specific order statistics from this sorted sequence:  $[z_{(k_{\text{left}})}, z_{(k_{\text{right}})}]$ . The challenge lies in determining which order statistics produce bounds that contain the true shift  $\text{Shift}[X, Y]$  with probability  $1 - \text{misrate}$ .

Random sampling creates natural variation in pairwise differences. Even when populations have identical distributions, sampling variation produces both positive and negative differences. The margin function quantifies this sampling variability: it specifies how many extreme pairwise differences could

occur by chance with probability misrate. For symmetric bounds, this margin splits evenly between the tails, giving  $k_{\text{left}} = \left\lfloor \frac{\text{PairwiseMargin}(n, m, \text{misrate})}{2} \right\rfloor + 1$  and  $k_{\text{right}} = nm - \left\lfloor \frac{\text{PairwiseMargin}(n, m, \text{misrate})}{2} \right\rfloor$ .

Computing the margin requires understanding the distribution of pairwise comparisons. Each pairwise difference corresponds to a comparison:  $x_i - y_j > 0$  exactly when  $x_i > y_j$ . This connection motivates the dominance function:

$$\text{Dominance}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}(x_i > y_j)$$

The dominance function counts how many pairwise comparisons favor  $\mathbf{x}$  over  $\mathbf{y}$ . Both Shift and Dominance operate on the same collection of  $nm$  pairwise differences. The Shift estimator examines difference values, returning the median as a location estimate. The Dominance function examines difference signs, counting how many comparisons produce positive differences. While Shift provides the estimate itself, Dominance determines which order statistics form reliable bounds around that estimate.

When populations have equivalent distributions, Dominance concentrates near  $nm/2$  by symmetry. The distribution of Dominance across all possible sample orderings determines reliable bounds. If Dominance deviates from  $nm/2$  by at least  $k/2$  with probability misrate, then the interval excluding the  $k$  most extreme pairwise differences contains zero with probability  $1 - \text{misrate}$ . Translation invariance extends this relationship to arbitrary shifts: the margin computed from the comparison distribution applies regardless of the true shift value.

Two computational approaches provide the distribution of Dominance: exact calculation for small samples and approximation for large samples.

### Exact method

Small sample sizes allow exact computation without approximation. The exact approach exploits a fundamental symmetry: under equivalent populations, all  $\binom{n+m}{n}$  orderings of the combined measurements occur with equal probability. This symmetry enables direct calculation of how many orderings produce each comparison count.

Direct computation faces a combinatorial challenge. Enumerating all orderings to count comparison outcomes requires substantial resources. For samples beyond a few dozen measurements, naive implementation becomes impractical.

Löffler's recurrence relation (Löffler (1982)) resolves this through algebraic structure. The recurrence exploits cycle properties in the comparison distribution, reducing memory requirements while maintaining exact calculation. The algorithm builds cumulative probabilities sequentially until reaching the threshold corresponding to the desired error rate. This approach extends practical exact computation to combined sample sizes of several hundred.

Define  $p_{n,m}(c)$  as the number of orderings producing exactly  $c$  comparisons favoring  $\mathbf{x}$ . The probability mass function becomes:

$$\Pr(\text{Dominance} = c) = p_{n,m} \frac{c}{\binom{n+m}{n}}$$

A direct recurrence follows from considering the largest measurement. The rightmost element comes from either  $x$  (contributing  $m$  comparisons) or  $y$  (contributing zero):

$$p_{n,m}(c) = p_{n-1,m}(c-m) + p_{n,m-1}(c)$$

with base cases  $p_{n,0}(0) = 1$  and  $p_{0,m}(0) = 1$ .

Direct implementation requires  $O(n \cdot m \cdot nm)$  time and  $O(nm)$  memory. An alternative recurrence (Löffler (1982)) exploits cycle structure:

$$p_{n,m}(c) = \frac{1}{c} \sum_{i=0}^{c-1} p_{n,m}(i) \cdot \sigma_{n,m}(c-i)$$

where  $\sigma_{n,m}(d)$  captures structural properties through divisors:

$$\sigma_{n,m}(d) = \sum_{k|d} \varepsilon_k \cdot k, \quad \varepsilon_k = \begin{cases} 1 & \text{if } 1 \leq k \leq n \\ -1 & \text{if } m+1 \leq k \leq m+n \\ 0 & \text{otherwise} \end{cases}$$

This reduces memory to  $O(nm)$  and enables efficient computation through  $c = nm$ .

The algorithm computes cumulative probabilities  $\Pr(\text{Dominance} \leq c)$  sequentially until the threshold  $\text{misrate}/2$  is exceeded. By symmetry, the lower and upper thresholds determine the total margin  $\text{PairwiseMargin} = 2c$ .

The sequential computation proceeds incrementally. Starting from  $u = 0$  with base probability  $p_{n,m}(0) = 1$ , the algorithm computes  $p_{n,m}(1)$ , then  $p_{n,m}(2)$ , and so on, accumulating the cumulative distribution function with each step. The loop terminates as soon as  $\Pr(\text{Dominance} \leq u)$  reaches  $\text{misrate}/2$ , returning the threshold value  $u$  without computing further probabilities.

This sequential approach performs particularly well for small misrates. For  $\text{misrate} = 10^{-6}$ , the threshold  $u$  typically remains small even with large sample sizes, requiring only a few iterations regardless of whether  $n$  and  $m$  equal 50 or 200. The algorithm computes only the extreme tail probabilities needed to reach the threshold, never touching the vast majority of probability mass concentrated near  $nm/2$ . This efficiency advantage grows as misrates decrease: stricter bounds require fewer computed values, making exact calculation particularly attractive for high-confidence applications.

### Approximate method

Large samples make exact computation impractical. The dominance count  $\text{Dominance}$  concentrates near  $nm/2$  with variance  $nm(n+m+1)/12$ . A basic Additive (Normal) approximation suffices asymptotically:

$$\text{Dominance} \approx \text{Additive}\left(nm/2, \sqrt{nm(n+m+1)/12}\right)$$

This approximation underestimates tail probabilities for moderate sample sizes. The Additive (Normal) approximation provides a baseline but fails to capture the true distribution shape in the tails, producing mis-calibrated probabilities that become problematic for small error rates.

The Edgeworth expansion refines this approximation through moment-based corrections (Fix & Hodges (1955)). The expansion starts with the Additive (Normal) cumulative distribution as a baseline, then adds correction terms that account for the distribution's asymmetry (skewness) and tail weight (kurtosis). These corrections use Hermite polynomials to adjust the baseline curve where the Additive (Normal) approximation deviates most from the true distribution. The first few correction terms typically achieve the practical balance between accuracy and computational cost, substantially improving tail probability estimates compared to the basic approximation.

The standardized comparison count:

$$z = \frac{c - nm/2}{\sqrt{nm(n+m+1)/12}}$$

produces the approximated cumulative distribution:

$$\Pr(\text{Dominance} \leq c) \approx \Phi(z) + e_3\varphi^{(3)}(z) + e_5\varphi^{(5)}(z) + e_7\varphi^{(7)}(z)$$

where  $\Phi$  denotes the standard Additive (Normal) CDF.

The correction coefficients depend on standardized moments:

$$e_3 = \frac{1}{24} \left( \frac{\mu_4}{\mu_2^2} - 3 \right), \quad e_5 = \frac{1}{720} \left( \frac{\mu_6}{\mu_2^3} - 15 \frac{\mu_4}{\mu_2^2} + 30 \right), \quad e_7 = \frac{35}{40320} \left( \frac{\mu_4}{\mu_2^2} - 3 \right)^2$$

The moments  $\mu_2, \mu_4, \mu_6$  are computed from sample sizes:

$$\begin{aligned} \mu_2 &= \frac{nm(n+m+1)}{12} \\ \mu_4 &= \frac{nm(n+m+1)}{240} (5nm(n+m) - 2(n^2 + m^2) + 3nm - 2(n+m)) \\ \mu_6 &= \frac{nm(n+m+1)}{4032} (35n^2m^2(n^2 + m^2) + 70n^3m^3 - 42nm(n^3 + m^3) \\ &\quad - 14n^2m^2(n+m) + 16(n^4 + m^4) - 52nm(n^2 + m^2) \\ &\quad - 43n^2m^2 + 32(n^3 + m^3) + 14nm(n+m) \\ &\quad + 8(n^2 + m^2) + 16nm - 8(n+m)) \end{aligned}$$

The correction terms use Hermite polynomials:

$$\varphi^{(k)}(z) = -\varphi(z)H_{k(z)}$$

$$H_3(z) = z^3 - 3z, \quad H_5(z) = z^5 - 10z^3 + 15z, \quad H_7(z) = z^7 - 21z^5 + 105z^3 - 105z$$

Binary search locates the threshold value efficiently. The algorithm maintains a search interval  $[a, b]$  initialized to  $[0, nm]$ . Each iteration computes the midpoint  $c = (a + b)/2$  and evaluates the Edgeworth CDF at  $c$ . If  $\Pr(\text{Dominance} \leq c) < \text{misrate}/2$ , the threshold lies above  $c$  and the search continues in  $[c, b]$ . If  $\Pr(\text{Dominance} \leq c) \geq \text{misrate}/2$ , the threshold lies below  $c$  and the search continues in  $[a, c]$ . The loop terminates when  $a$  and  $b$  become adjacent, requiring  $O(\log(nm))$  CDF evaluations.

This binary search exhibits uniform performance across misrate values. Whether computing bounds for  $\text{misrate} = 10^{-6}$  or  $\text{misrate} = 0.05$ , the algorithm performs the same number of iterations determined solely by the sample sizes. Each CDF evaluation costs constant time regardless of the threshold location, making the approximate method particularly efficient for large samples where exact computation becomes impractical. The logarithmic scaling ensures that doubling the sample size adds only one additional iteration, enabling practical computation for samples in the thousands or tens of thousands.

The toolkit selects between exact and approximate computation based on combined sample size: exact method for  $n + m \leq 400$ , approximate method for  $n + m > 400$ . The exact method guarantees correctness but scales as  $O(nm)$  memory and  $O((nm)^2)$  time. For  $n = m = 200$ , this requires 40,000 memory locations. The approximate method achieves 1% accuracy with  $O(\log(nm))$  constant-time evaluations. For  $n = m = 10000$ , the approximate method completes in milliseconds versus minutes for exact computation.

Both methods handle discrete data. Repeated measurements produce tied pairwise differences, creating plateaus in the sorted sequence. The exact method counts orderings without assuming continuity. The approximate method's moment-based corrections capture the actual distribution shape regardless of discreteness.

### Minimum reasonable misrate

The misrate parameter controls how many extreme pairwise differences the bounds exclude. Lower misrate produces narrower bounds with higher confidence but requires excluding fewer extreme values. However, sample size limits how small misrate can meaningfully become.

Consider the most extreme configuration: all measurements from  $x$  exceed all measurements from  $y$ , giving  $x_1, \dots, x_n > y_1, \dots, y_m$ . Under equivalent populations, this arrangement occurs purely by chance. The probability equals the chance of having all  $n$  elements from  $x$  occupy the top  $n$  positions among  $n + m$  total measurements:

$$\text{misrate}_{\min} = \frac{1}{\binom{n+m}{n}} = \frac{n! \cdot m!}{(n+m)!}$$

This represents the minimum probability of the most extreme ordering under random sampling. Setting  $\text{misrate} < \text{misrate}_{\min}$  makes bounds construction problematic. The exact distribution of Dominance cannot support misrates smaller than the probability of its most extreme realization. Attempting to construct bounds with  $\text{misrate} < \text{misrate}_{\min}$  forces the algorithm to exclude zero pairwise differences from the tails, making  $\text{PairwiseMargin} = 0$ . The resulting bounds span all  $nm$  pairwise differences, returning  $[z_{(1)}, z_{(nm)}]$  regardless of the desired confidence level. These bounds convey no useful information beyond the range of observed pairwise differences.

For small samples,  $\text{misrate}_{\min}$  can exceed commonly used values. With  $n = m = 6$ , the minimum misrate equals  $1/\binom{12}{6} \approx 0.00108$ , making the typical choice of  $\text{misrate} = 10^{-3}$  impossible. With  $n = m = 4$ , the minimum becomes  $1/\binom{8}{4} \approx 0.0143$ , exceeding even  $\text{misrate} = 0.01$ .

The table below shows  $\text{misrate}_{\min}$  for small sample sizes:

	1	2	3	4	5	6	7	8	9	10
1	0.500000	0.333333	0.250000	0.200000	0.166667	0.142857	0.125000	0.111111	0.100000	0.090909
2	0.333333	0.166667	0.100000	0.066667	0.047619	0.035714	0.027778	0.022222	0.018182	0.015152
3	0.250000	0.100000	0.050000	0.028571	0.017857	0.011905	0.008333	0.006061	0.004545	0.003497
4	0.200000	0.066667	0.028571	0.014286	0.007937	0.004762	0.003030	0.002020	0.001399	0.000999
5	0.166667	0.047619	0.017857	0.007937	0.003968	0.002165	0.001263	0.000777	0.000500	0.000333
6	0.142857	0.035714	0.011905	0.004762	0.002165	0.001082	0.000583	0.000333	0.000200	0.000125
7	0.125000	0.027778	0.008333	0.003030	0.001263	0.000583	0.000291	0.000155	0.000087	0.000051
8	0.111111	0.022222	0.006061	0.002020	0.000777	0.000333	0.000155	0.000078	0.000041	0.000023
9	0.100000	0.018182	0.004545	0.001399	0.000500	0.000200	0.000087	0.000041	0.000021	0.000011
10	0.090909	0.015152	0.003497	0.000999	0.000333	0.000125	0.000051	0.000023	0.000011	0.000005

For meaningful bounds construction, choose  $\text{misrate} > \text{misrate}_{\min}$ . This ensures the margin function excludes at least some extreme pairwise differences, producing bounds narrower than the full range. When working with small samples, verify that the desired misrate exceeds  $\text{misrate}_{\min}$  for the given sample sizes. With moderate sample sizes ( $n, m \geq 15$ ),  $\text{misrate}_{\min}$  drops below  $10^{-8}$ , making standard choices like  $\text{misrate} = 10^{-6}$  feasible.

```
using System;
using JetBrains.Annotations;
using Pragmatat.Internal;

namespace Pragmatat.Functions;

/// <summary>
/// PairwiseMargin function
/// </summary>
/// <param name="threshold">The maximum value for n+m, after which implementation switches
/// from exact to approx</param>
public class PairwiseMargin(int threshold = PairwiseMargin.MaxExactSize)
{
    public static readonly PairwiseMargin Instance = new();

    private const int MaxExactSize = 400;

    [PublicAPI]
    public int Calc(int n, int m, double misrate)
    {
        Assertion.MoreThan(nameof(n), n, 0);
        Assertion.MoreThan(nameof(m), m, 0);
        Assertion.InRangeInclusive(nameof(misrate), misrate, 0, 1);

        return n + m <= threshold
    }
}
```

```

        ? CalcExact(n, m, misrate)
        : CalcApprox(n, m, misrate);
    }

    [PublicAPI]
    public int CalcExact(int n, int m, double misrate)
    {
        int raw = CalcExactRaw(n, m, misrate / 2);
        return checked(raw * 2);
    }

    [PublicAPI]
    public int CalcApprox(int n, int m, double misrate)
    {
        long raw = CalcApproxRaw(n, m, misrate / 2);
        long margin = raw * 2;
        if (margin > int.MaxValue)
            throw new OverflowException($"Pairwise margin exceeds supported range for n={n},
m={m}");
        return (int)margin;
    }

    // Inversed implementation of Andreas Löffler's (1982) "Über eine Partition der nat.
    // Zahlen und ihre Anwendung beim U-Test"
    private static int CalcExactRaw(int n, int m, double p)
    {
        double total = n + m < BinomialCoefficientFunction.MaxAcceptableN
            ? BinomialCoefficientFunction.BinomialCoefficient(n + m, m)
            : BinomialCoefficientFunction.BinomialCoefficient(n + m, m * 1.0);

        var pmf = new List<double> { 1 }; // pmf[0] = 1
        var sigma = new List<double> { 0 }; // sigma[0] is unused

        int u = 0;
        double cdf = 1.0 / total;

        if (cdf >= p)
            return 0;

        while (true)
        {
            u++;
            // Ensure sigma has entry for u
            if (sigma.Count <= u)
            {
                int value = 0;
                for (int d = 1; d <= n; d++)
                    if (u % d == 0 && u >= d)
                        value += d;
                for (int d = m + 1; d <= m + n; d++)
                    if (u % d == 0 && u >= d)
                        value -= d;
                sigma.Add(value);
            }
        }
    }

```

```

    }

    // Compute pmf[u] using Loeffler recurrence
    double sum = 0;
    for (int i = 0; i < u; i++)
        sum += pmf[i] * sigma[u - i];
    sum /= u;
    pmf.Add(sum);

    cdf += sum / total;
    if (cdf >= p)
        return u;
    if (sum == 0)
        break;
}

return pmf.Count - 1;
}

// Inverse Edgeworth Approximation
private static long CalcApproxRaw(int n, int m, double misrate)
{
    long a = 0;
    long b = (long)n * m;
    while (a < b - 1)
    {
        long c = (a + b) / 2;
        double p = EdgeworthCdf(n, m, c);
        if (p < misrate)
            a = c;
        else
            b = c;
    }

    return EdgeworthCdf(n, m, b) < misrate ? b : a;
}

private static double EdgeworthCdf(int n, int m, long u)
{
    double nm = (double)n * m;
    double mu = nm / 2.0;
    double su = Sqrt(nm * (n + m + 1) / 12.0);
    double z = (u - mu - 0.5) / su;
    double phi = Exp(-z.Sqr() / 2) / Sqrt(2 * PI);
    double Phi = AcmlAlgorithm209.Gauss(z);

    // Pre-compute powers of n and m for efficiency
    double n2 = n * n;
    double n3 = n2 * n;
    double n4 = n2 * n2;
    double m2 = m * m;
    double m3 = m2 * m;
    double m4 = m2 * m2;

```



```

double mu2 = n * m * (n + m + 1) / 12.0;
double mu4 =
    n * m * (n + m + 1) *
    (0
    + 5 * m * n * (m + n)
    - 2 * (m2 + n2)
    + 3 * m * n
    - 2 * (n + m)
    ) / 240.0;

double mu6 =
    n * m * (n + m + 1) *
    (0
    + 35 * m2 * n2 * (m2 + n2)
    + 70 * m3 * n3
    - 42 * m * n * (m3 + n3)
    - 14 * m2 * n2 * (n + m)
    + 16 * (n4 + m4)
    - 52 * n * m * (n2 + m2)
    - 43 * n2 * m2
    + 32 * (m3 + n3)
    + 14 * m * n * (n + m)
    + 8 * (n2 + m2)
    + 16 * n * m
    - 8 * (n + m)
    ) / 4032.0;

// Pre-compute powers of mu2 and related terms
double mu2_2 = mu2 * mu2;
double mu2_3 = mu2_2 * mu2;
double mu4_mu2_2 = mu4 / mu2_2;

// Factorial constants: 4! = 24, 6! = 720, 8! = 40320
double e3 = (mu4_mu2_2 - 3) / 24.0;
double e5 = (mu6 / mu2_3 - 15 * mu4_mu2_2 + 30) / 720.0;
double e7 = 35 * (mu4_mu2_2 - 3) * (mu4_mu2_2 - 3) / 40320.0;

// Pre-compute powers of z for Hermite polynomials
double z2 = z * z;
double z3 = z2 * z;
double z5 = z3 * z2;
double z7 = z5 * z2;

double f3 = -phi * (z3 - 3 * z);
double f5 = -phi * (z5 - 10 * z3 + 15 * z);
double f7 = -phi * (z7 - 21 * z5 + 105 * z3 - 105 * z);

double edgeworth = Phi + e3 * f3 + e5 * f5 + e7 * f7;
return Min(Max(edgeworth, 0), 1);
}
}

```

## 4. Studies

This section analyzes the estimators' properties using mathematical proofs. Most proofs are adapted from various textbooks and research papers, but only essential references are provided.

Unlike the main part of the manual, the studies require knowledge of classic statistical methods. Well-known facts and commonly accepted notation are used without special introduction. The studies provide detailed analyses of estimator properties for practitioners interested in rigorous proofs and numerical simulation results.

### 4.1. Summary Estimator Properties

This section compares the toolkit's robust estimators against traditional statistical methods to demonstrate their advantages across diverse conditions. While traditional estimators often work well under ideal conditions, the toolkit's estimators maintain reliable performance across diverse real-world scenarios.

Average Estimators:

**Mean** (arithmetic average):

$$\text{Mean}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n x_i$$

**Median:**

$$\text{Median}(\mathbf{x}) = \begin{cases} x_{((n+1)/2)} & \text{if } n \text{ is odd} \\ \frac{x_{(n/2)} + x_{(n/2+1)}}{2} & \text{if } n \text{ is even} \end{cases}$$

**Center** (Hodges-Lehmann estimator):

$$\text{Center}(\mathbf{x}) = \text{Median}_{1 \leq i \leq j \leq n} \left( \frac{x_i + x_j}{2} \right)$$

Dispersion Estimators:

**Standard Deviation:**

$$\text{StdDev}(\mathbf{x}) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \text{Mean}(\mathbf{x}))^2}$$

**Median Absolute Deviation** (around the median):

$$\text{MAD}(\mathbf{x}) = \text{Median}(|x_i - \text{Median}(\mathbf{x})|)$$

**Spread** (Shamos scale estimator):

$$\text{Spread}(\mathbf{x}) = \text{Median}_{1 \leq i < j \leq n} |x_i - x_j|$$

## 4.2. Breakdown

Heavy-tailed distributions naturally produce extreme outliers that completely distort traditional estimators. A single extreme measurement from the Power distribution can make the sample mean arbitrarily large. Real-world data can also contain corrupted measurements from instrument failures, recording errors, or transmission problems. Both natural extremes and data corruption create the same challenge: extracting reliable information when some measurements are too influential.

The breakdown point (Huber (2009)) is the fraction of a sample that can be replaced by arbitrarily large values without making an estimator arbitrarily large. The theoretical maximum is 50%; no estimator can guarantee reliable results when more than half the measurements are extreme or corrupted. In such cases, summary estimators are not applicable, and a more sophisticated approach is needed.

A 50% breakdown point is rarely needed in practice, as more conservative values also cover practical needs. Additionally, a high breakdown point corresponds to low precision (information is lost by neglecting part of the data). The optimal practical breakdown point should be between 0% (no robustness) and 50% (low precision).

The Center and Spread estimators achieve 29% breakdown points, providing substantial protection against realistic contamination levels while maintaining good precision. Below is a comparison with traditional estimators.

**Asymptotic breakdown points** for average estimators:

Mean	Median	Center
0%	50%	29%

**Asymptotic breakdown points** for dispersion estimators:

StdDev	MAD	Spread
0%	50%	29%

## 4.3. Drift

Drift measures estimator precision by quantifying how much estimates scatter across repeated samples. It is based on the Spread of estimates and therefore has a breakdown point of approximately 29%.

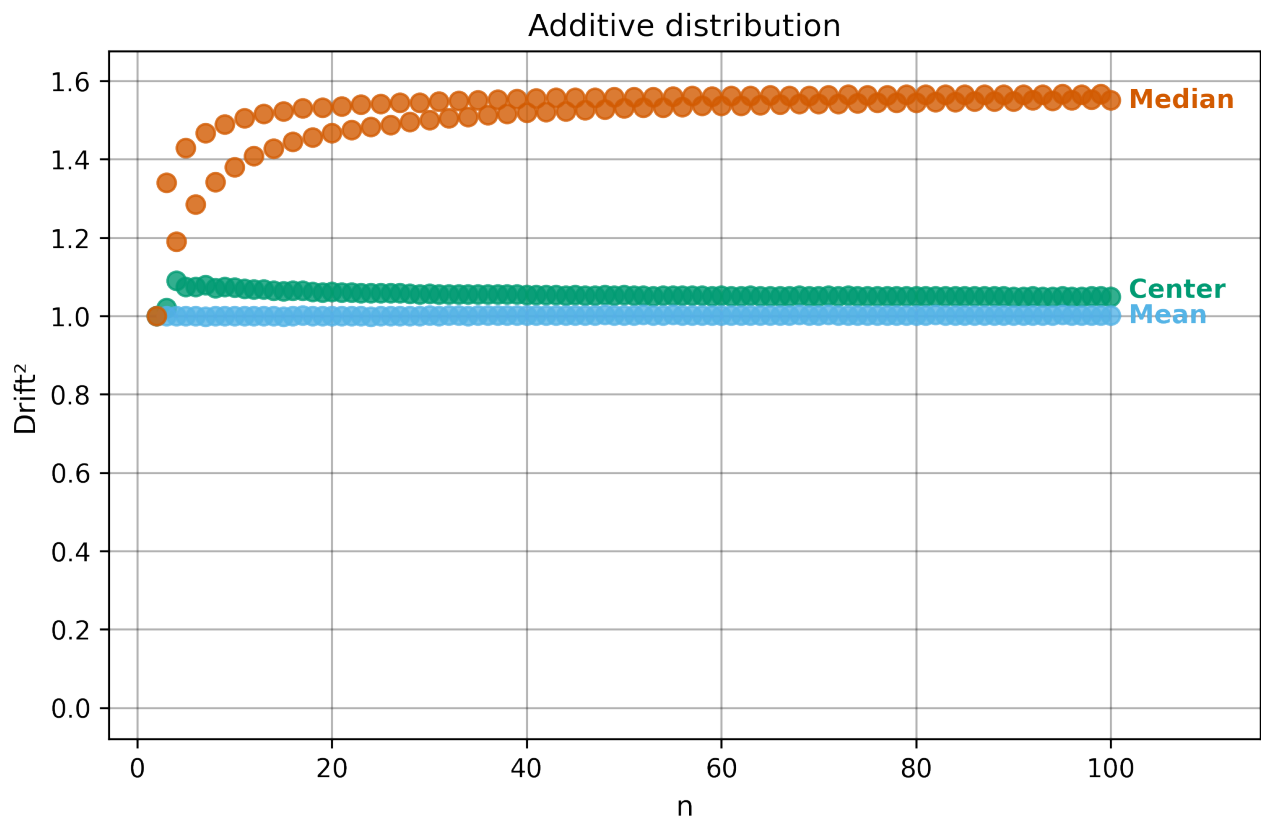
Drift is useful for comparing the precision of several estimators. To simplify the comparison, one of the estimators can be chosen as a baseline. A table of squared drift values, normalized by the baseline, shows the required sample size adjustment factor for switching from the baseline to another estimator. For example, if Center is the baseline and the rescaled drift square of Median is 1.5, this means that Median requires 1.5 times more data than Center to achieve the same precision. See the “From Statistical Efficiency to Drift” section for details.

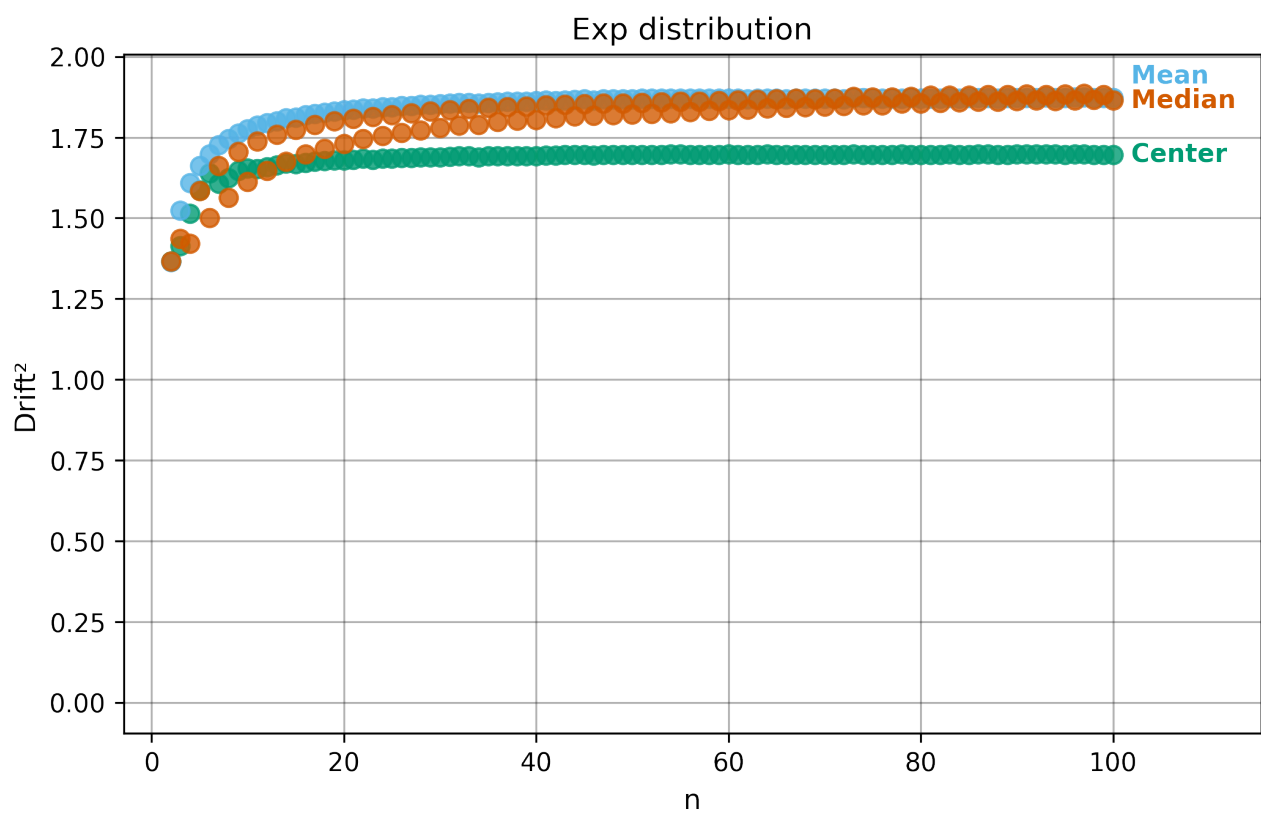
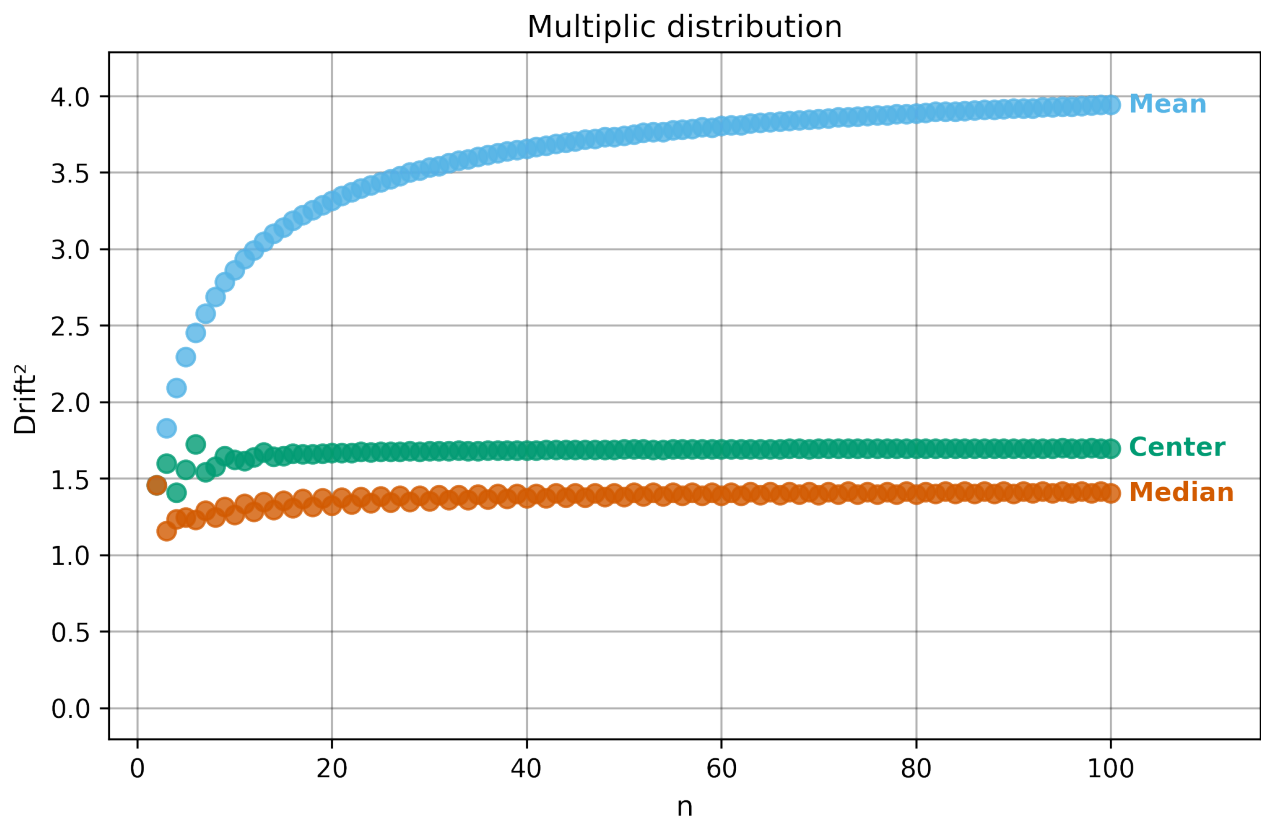
**Squared Asymptotic Drift of Average Estimators** (values are approximated):

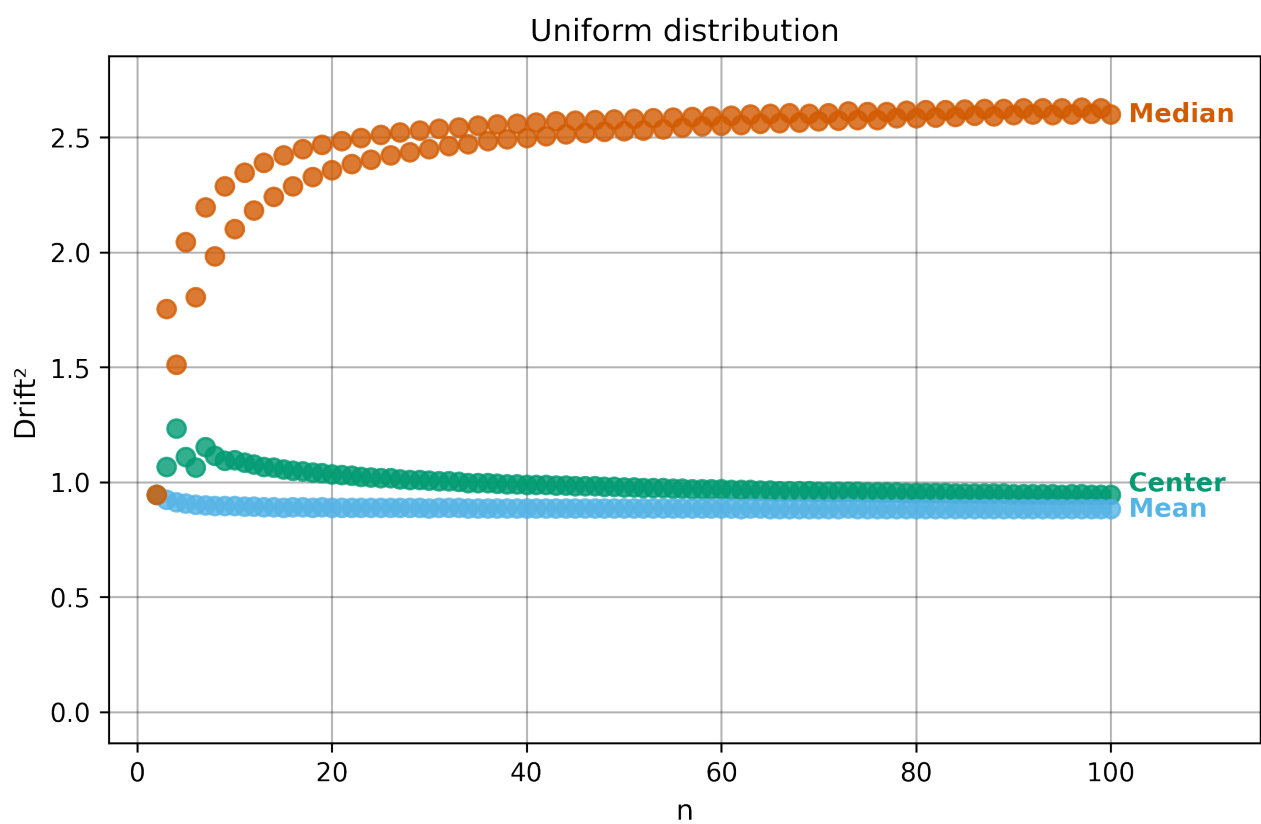
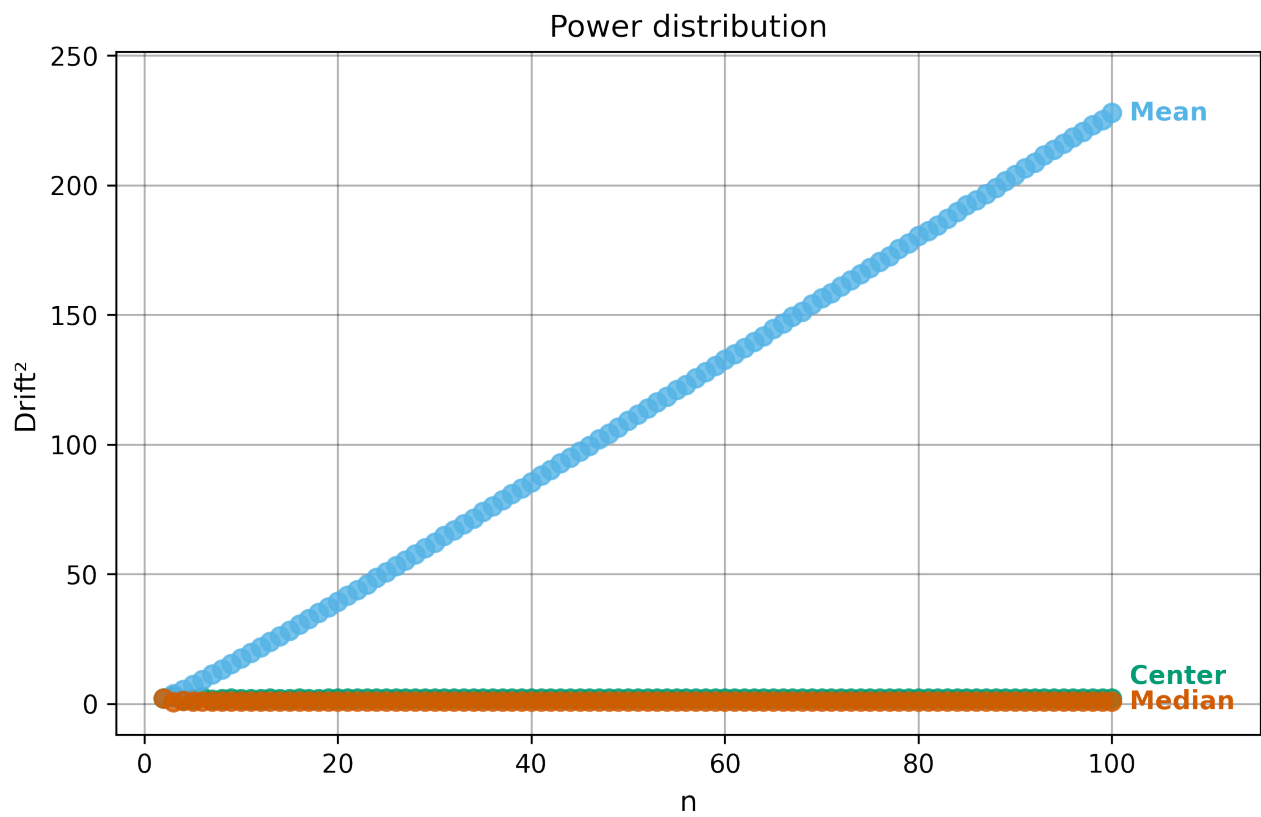
	Mean	Median	Center
<u>Additive</u>	1.0	1.571	1.047
<u>Multiplic</u>	3.95	1.40	1.7
<u>Exp</u>	1.88	1.88	1.69
<u>Power</u>	$\infty$	0.9	2.1
<u>Uniform</u>	0.88	2.60	0.94

Rescaled to Center (sample size adjustment factors):

	Mean	Median	Center
<u>Additive</u>	0.96	1.50	1.0
<u>Multiplic</u>	2.32	0.82	1.0
<u>Exp</u>	1.11	1.11	1.0
<u>Power</u>	$\infty$	0.43	1.0
<u>Uniform</u>	0.936	2.77	1.0





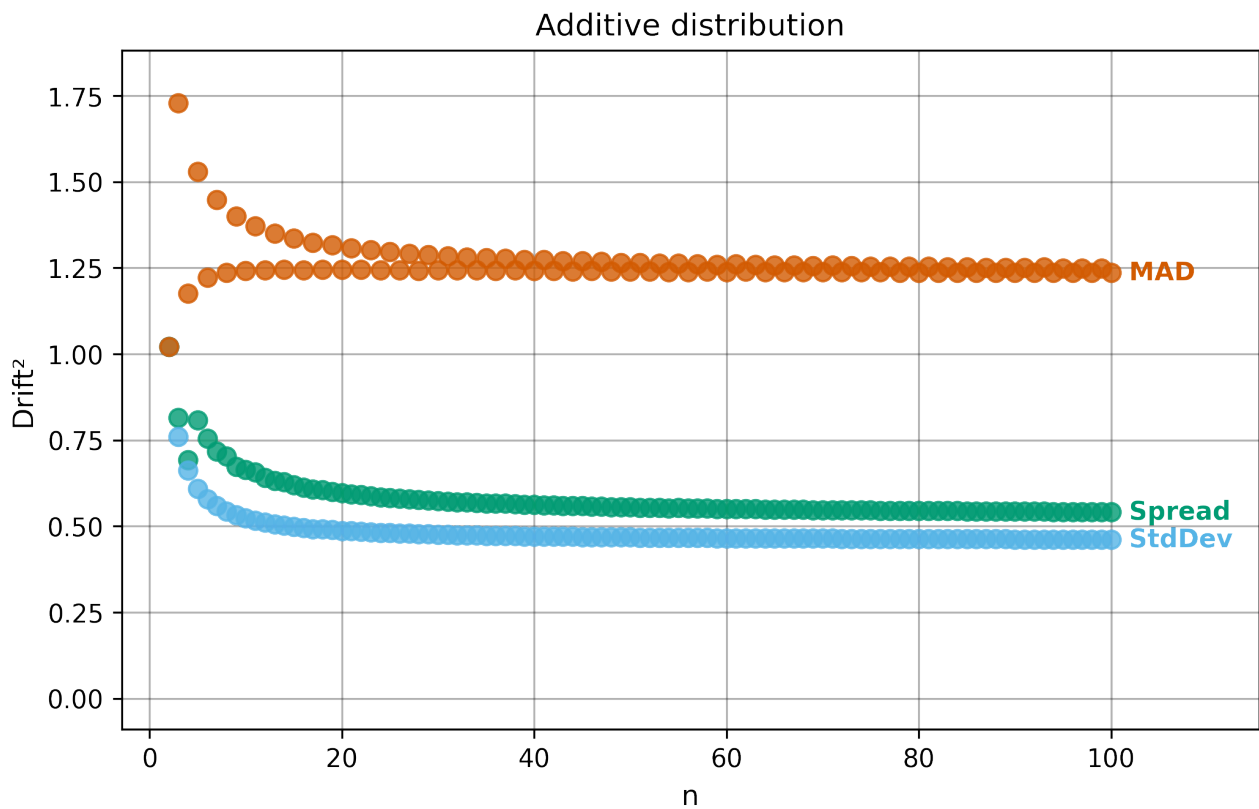


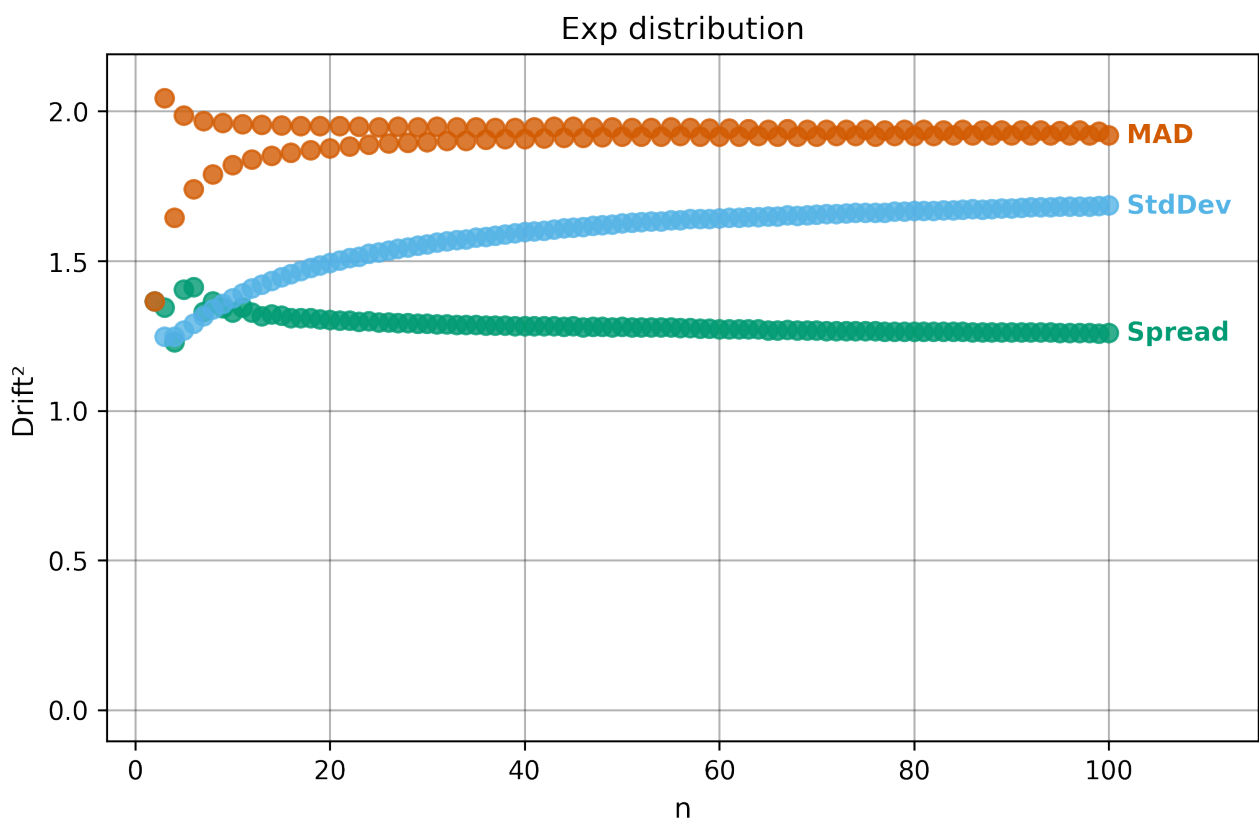
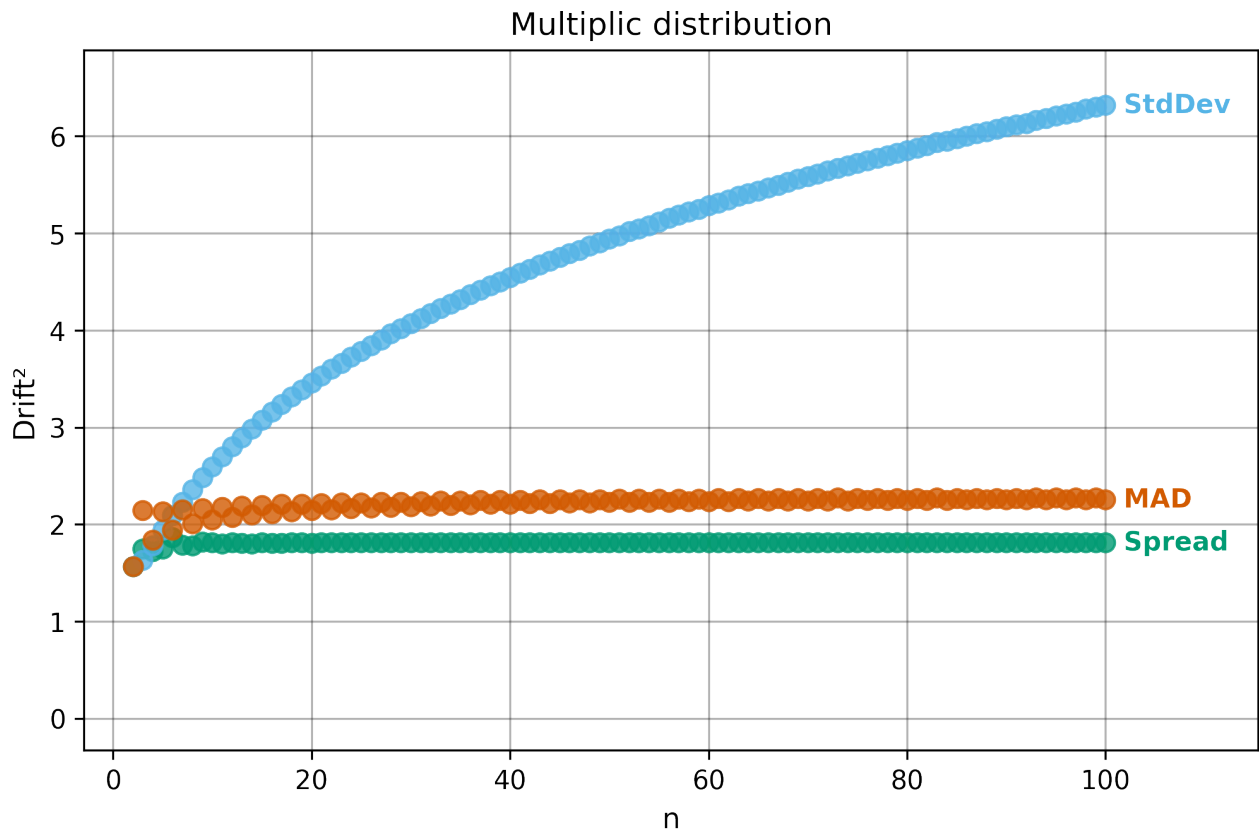
**Squared Asymptotic Drift of Dispersion Estimators** (values are approximated):

	StdDev	MAD	Spread
<u>Additive</u>	0.45	1.22	0.52
<u>Multiplic</u>	$\infty$	2.26	1.81
<u>Exp</u>	1.69	1.92	1.26
<u>Power</u>	$\infty$	3.5	4.4
<u>Uniform</u>	0.18	0.90	0.43

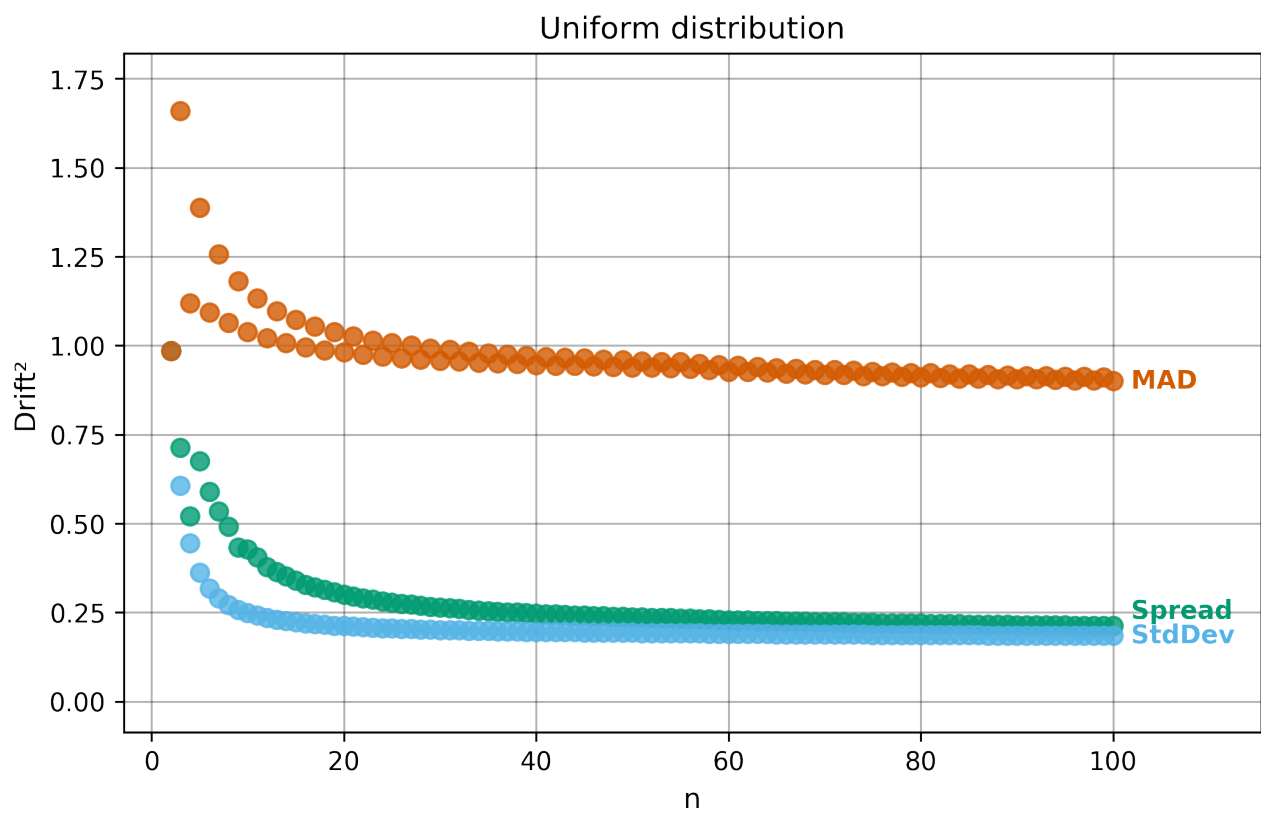
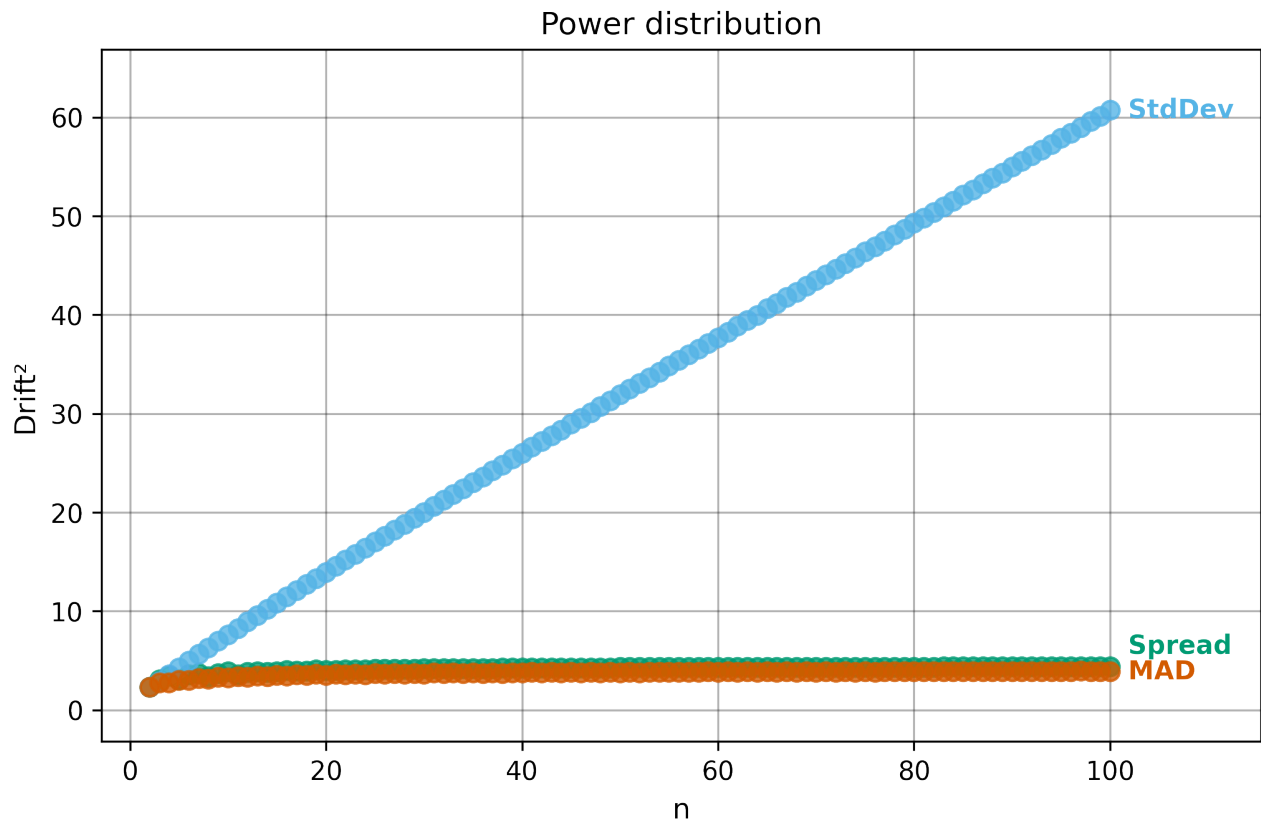
Rescaled to Spread (sample size adjustment factors):

	StdDev	MAD	Spread
<u>Additive</u>	0.87	2.35	1.0
<u>Multiplic</u>	$\infty$	1.25	1.0
<u>Exp</u>	1.34	1.52	1.0
<u>Power</u>	$\infty$	0.80	1.0
<u>Uniform</u>	0.42	2.09	1.0









## 4.4. Invariance

Invariance properties determine how estimators respond to data transformations. These properties are crucial for analysis design and interpretation:

**Location-invariant** estimators are invariant to additive shifts:  $T(\mathbf{x} + k) = T(\mathbf{x})$

**Scale-invariant** estimators are invariant to positive rescaling:  $T(k \cdot \mathbf{x}) = T(\mathbf{x})$  for  $k > 0$

**Equivariant** estimators change predictably with transformations, maintaining relative relationships

Choosing estimators with appropriate invariance properties ensures that results remain meaningful across different measurement scales, units, and data transformations. For example, when comparing datasets collected with different instruments or protocols, location-invariant estimators eliminate the need for data centering, while scale-invariant estimators eliminate the need for normalization.

**Location-invariance:** An estimator  $T$  is location-invariant if adding a constant to the measurements leaves the result unchanged:

$$T(\mathbf{x} + k) = T(\mathbf{x})$$

$$T(\mathbf{x} + k, \mathbf{y} + k) = T(\mathbf{x}, \mathbf{y})$$

**Location-equivariance:** An estimator  $T$  is location-equivariant if it shifts with the data:

$$T(\mathbf{x} + k) = T(\mathbf{x}) + k$$

$$T(\mathbf{x} + k_1, \mathbf{y} + k_2) = T(\mathbf{x}, \mathbf{y}) + f(k_1, k_2)$$

**Scale-invariance:** An estimator  $T$  is scale-invariant if multiplying by a positive constant leaves the result unchanged:

$$T(k \cdot \mathbf{x}) = T(\mathbf{x}) \quad \text{for } k > 0$$

$$T(k \cdot \mathbf{x}, k \cdot \mathbf{y}) = T(\mathbf{x}, \mathbf{y}) \quad \text{for } k > 0$$

**Scale-equivariance:** An estimator  $T$  is scale-equivariant if it scales proportionally with the data:

$$T(k \cdot \mathbf{x}) = k \cdot T(\mathbf{x}) \text{ or } |k| \cdot T(\mathbf{x}) \quad \text{for } k \neq 0$$

$$T(k \cdot \mathbf{x}, k \cdot \mathbf{y}) = k \cdot T(\mathbf{x}, \mathbf{y}) \text{ or } |k| \cdot T(\mathbf{x}, \mathbf{y}) \quad \text{for } k \neq 0$$

	Location	Scale
Center	Equivariant	Equivariant
Spread	Invariant	Equivariant
RelSpread	–	Invariant
Shift	Invariant	Equivariant
Ratio	–	Invariant
AvgSpread	Invariant	Equivariant
Disparity	Invariant	Invariant

## 4.5. From Statistical Efficiency to Drift

Statistical efficiency measures estimator precision (Serfling (2009)). When multiple estimators target the same quantity, efficiency determines which provides more reliable results.

Efficiency measures how tightly estimates cluster around the true value across repeated samples. For an estimator  $T$  applied to samples from distribution  $X$ , absolute efficiency is defined relative to the optimal estimator  $T^*$ :

$$\text{Efficiency}(T, X) = \frac{\text{Var}[T^*(X_1, \dots, X_n)]}{\text{Var}[T(X_1, \dots, X_n)]}$$

Relative efficiency compares two estimators by taking the ratio of their variances:

$$\text{RelativeEfficiency}(T_1, T_2, X) = \frac{\text{Var}[T_2(X_1, \dots, X_n)]}{\text{Var}[T_1(X_1, \dots, X_n)]}$$

Under Additive (Normal) distributions, this approach works well. The sample mean achieves optimal efficiency, while the median operates at roughly 64% efficiency.

However, this variance-based definition creates four critical limitations:

Absolute efficiency requires knowing the optimal estimator, which is difficult to determine. For many distributions, deriving the minimum-variance unbiased estimator requires complex mathematical analysis. Without this reference point, absolute efficiency cannot be computed.

Relative efficiency only compares estimator pairs, preventing systematic evaluation. This limits understanding of how multiple estimators perform relative to each other. Practitioners cannot rank estimators comprehensively or evaluate individual performance in isolation.

The approach depends on variance calculations that break down when variance becomes infinite or when distributions have heavy tails. Many real-world distributions, such as those with power-law tails, exhibit infinite variance. When the variance is undefined, efficiency comparisons become impossible.

Variance is not robust to outliers, which can corrupt efficiency calculations. A single extreme observation can greatly inflate variance estimates. This sensitivity can make efficient estimators look inefficient and vice versa.

The Drift concept provides a robust alternative. Drift measures estimator precision using Spread instead of variance, providing reliable comparisons across a wide range of distributions.

For an average estimator  $T$ , random variable  $X$ , and sample size  $n$ :

$$\text{AvgDrift}(T, X, n) = \frac{\sqrt{n} \cdot \text{Spread}[T(X_1, \dots, X_n)]}{\text{Spread}[X]}$$

This formula measures estimator variability compared to data variability.  $\text{Spread}[T(X_1, \dots, X_n)]$  captures the median absolute difference between estimates across repeated samples. Multiplying by  $\sqrt{n}$  removes sample size dependency, making drift values comparable across different sample sizes. Dividing by  $\text{Spread}[X]$  creates a scale-free measure that provides consistent drift values across different distribution parameters and measurement units.

Dispersion estimators use a parallel formulation:

$$\text{DispDrift}(T, X, n) = \sqrt{n} \cdot \text{RelSpread}[T(X_1, \dots, X_n)]$$

Here  $\text{RelSpread}$  normalizes by the estimator's typical value for fair comparison.

Drift offers four key advantages:

For estimators with  $\sqrt{n}$  convergence rates, drift remains finite and comparable across distributions; for heavier tails, drift may diverge, flagging estimator instability.

It provides absolute precision measures rather than only pairwise comparisons.

The robust  $\text{Spread}$  foundation resists outlier distortion that corrupts variance-based calculations.

The  $\sqrt{n}$  normalization makes drift values comparable across different sample sizes, enabling direct comparison of estimator performance regardless of sample size.

Under Additive (Normal) conditions, drift matches traditional efficiency. The sample mean achieves drift near 1.0; the median achieves drift around 1.25. This consistency validates drift as a proper generalization of efficiency that extends to realistic data conditions where traditional efficiency fails.

When switching from one estimator to another while maintaining the same precision, the required sample size adjustment follows:

$$n_{\text{new}} = n_{\text{original}} \cdot \text{Drift}^2 \frac{T_2, X}{\text{Drift}^2(T_1, X)}$$

This applies when estimator  $T_1$  has lower drift than  $T_2$ .

The ratio of squared drifts determines the data requirement change. If  $T_2$  has drift 1.5 times higher than  $T_1$ , then  $T_2$  requires  $(1.5)^2 = 2.25$  times more data to match  $T_1$ 's precision. Conversely, switching to a more precise estimator allows smaller sample sizes.

For asymptotic analysis,  $\text{Drift}(T, X)$  denotes the limiting value as  $n \rightarrow \infty$ . With a baseline estimator, rescaled drift values enable direct comparisons:

$$\text{Drift}_{\text{baseline}(T, X)} = \frac{\text{Drift}(T, X)}{\text{Drift}(T_{\text{baseline}}, X)}$$

The standard drift definition assumes  $\sqrt{n}$  convergence rates typical under Additive (Normal) conditions. For broader applicability, drift generalizes to:

$$\text{AvgDrift}(T, X, n) = \frac{n^{\text{instability}} \cdot \text{Spread}[T(X_1, \dots, X_n)]}{\text{Spread}[X]}$$

$$\text{DispDrift}(T, X, n) = n^{\text{instability}} \cdot \text{RelSpread}[T(X_1, \dots, X_n)]$$

The instability parameter adapts to estimator convergence rates. The toolkit uses  $\text{instability} = 1/2$  throughout because this choice provides natural intuition and mental representation for the Additive (Normal) distribution. Rather than introduce additional complexity through variable instability parameters, the fixed  $\sqrt{n}$  scaling offers practical convenience while maintaining theoretical rigor for the distribution classes most common in applications.

## 4.6. From Confidence Level to Misrate

Traditional statistics expresses uncertainty through confidence levels: “95% confidence interval”, “99% confidence”, “99.9% confidence”. This convention emerged from early statistical practice when tables printed confidence intervals for common levels like 90%, 95%, and 99%.

The confidence level approach creates practical problems:

**Cognitive difficulty with high confidence.** Distinguishing between 99.999% and 99.9999% confidence requires mental effort. The difference matters — one represents a 1-in-100,000 error rate, the other 1-in-1,000,000 — but the representation obscures this distinction.

**Asymmetric scale.** The confidence level scale compresses near 100%, where most practical values cluster. Moving from 90% to 95% represents a 2× change in error rate, while moving from 99% to 99.9% represents a 10× change, despite similar visual spacing.

**Indirect interpretation.** Practitioners care about error rates, not success rates. “What’s the chance I’m wrong?” matters more than “What’s the chance I’m right?” Confidence level forces mental subtraction to answer the natural question.

**Unclear defaults.** Traditional practice offers no clear default confidence level. Different fields use different conventions (95%, 99%, 99.9%), creating inconsistency and requiring arbitrary choices.

The misrate parameter provides a more natural representation. Misrate expresses the probability that computed bounds fail to contain the true value:

$$\text{misrate} = 1 - \text{confidence level}$$

This simple inversion provides several advantages:

**Direct interpretation.**  $\text{misrate} = 0.01$  means “1% chance of error” or “wrong 1 time in 100”.  $\text{misrate} = 10^{-6}$  means “wrong 1 time in a million”. No mental arithmetic required.

**Linear scale for practical values.**  $\text{misrate} = 0.1$  (10%),  $\text{misrate} = 0.01$  (1%),  $\text{misrate} = 0.001$  (0.1%) form a natural sequence. Scientific notation handles extreme values cleanly:  $10^{-3}$ ,  $10^{-6}$ ,  $10^{-9}$ .

**Clear comparisons.**  $10^{-5}$  versus  $10^{-6}$  immediately shows a 10× difference in error tolerance. 99.999% versus 99.9999% confidence obscures this same relationship.

**Pragmatic default.** The toolkit recommends  $\text{misrate} = 10^{-6}$  (one-in-a-million error rate) as a reasonable default for most applications. This represents extremely high confidence (99.9999%) while remaining computationally practical and conceptually clear.

The terminology shift from “confidence level” to “misrate” parallels other clarifying renames in this toolkit. Just as Additive better describes the distribution’s formation than ‘Normal’, and Center better describes the estimator’s purpose than ‘Hodges-Lehmann’, misrate better describes the quantity practitioners actually reason about: the probability of error.

Traditional confidence intervals become “bounds” in this framework, eliminating statistical jargon in favor of descriptive terminology. `ShiftBounds(x, y, misrate)` clearly indicates: it provides bounds on the shift, with a specified error rate. No background in classical statistics required to understand the concept.

## 4.7. From Mann-Whitney U-test to Pairwise Margin

The Mann-Whitney  $U$  test (also known as the Wilcoxon rank-sum test) ranks among the most widely used non-parametric statistical tests, testing whether two independent samples come from the same distribution. Under Additive (Normal) conditions, it achieves nearly the same precision as the Student’s  $t$ -test, while maintaining reliability under diverse distributional conditions where the  $t$ -test fails.

The test operates by comparing all pairs of measurements between the two samples. Given samples  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_m)$ , the Mann-Whitney  $U$  statistic counts how many pairs satisfy  $x_i > y_j$ :

$$U = \sum_{i=1}^n \sum_{j=1}^m \mathbb{1}(x_i > y_j)$$

If the samples come from the same distribution,  $U$  should be near  $nm/2$  (roughly half the pairs favor  $\mathbf{x}$ , half favor  $\mathbf{y}$ ). Large deviations from  $nm/2$  suggest the distributions differ.

The test answers: “Could this  $U$  value arise by chance if the samples were truly equivalent?” The  $p$ -value quantifies this probability. If  $p < 0.05$ , traditional practice declares the difference “statistically significant”.

This approach creates several problems for practitioners:

**Binary thinking.** The test produces a yes/no answer: reject or fail to reject the null hypothesis. Practitioners typically want to know the magnitude of difference, not just whether one exists.

**Arbitrary thresholds.** The 0.05 threshold has no universal justification, yet it dominates practice and creates a false dichotomy between  $p = 0.049$  and  $p = 0.051$ .

**Hypothesis-centric framework.** The test assumes a null hypothesis of “no difference” and evaluates evidence against it. Real questions rarely concern exact equality; practitioners want to know “how different?” rather than “different or not?”

**Inverted logic.** The natural question is “what shifts are consistent with my data?” The test answers “is this specific shift (zero) consistent with my data?”

The toolkit inverts this framework. Instead of testing whether a hypothesized shift is plausible, we compute which shifts are plausible given the data. This inversion transforms hypothesis testing into bounds estimation.

The mathematical foundation remains the same. The distribution of pairwise comparisons under random sampling determines which order statistics of pairwise differences form reliable bounds. The Mann-Whitney  $U$  statistic measures pairwise comparisons ( $x_i > y_j$ ). The Shift estimator uses pairwise differences ( $x_i - y_j$ ). These quantities are mathematically related: a pairwise difference  $x_i - y_j$  is positive exactly when  $x_i > y_j$ . The toolkit renames this comparison count from  $U$  to  $\text{Dominance}(\mathbf{x}, \mathbf{y})$ , clarifying its purpose: measuring how often one sample dominates the other in pairwise comparisons.

The distribution of  $\text{Dominance}$  determines which order statistics form reliable bounds. Define the margin function:

$$\text{PairwiseMargin}(n, m, \text{misrate}) = \text{number of pairwise differences to exclude from bounds}$$

This function computes how many extreme pairwise differences could occur by chance with probability  $\text{misrate}$ , based on the distribution of pairwise comparisons.

The  $\text{PairwiseMargin}$  function requires knowing the distribution of pairwise comparisons under sampling. Two computational approaches exist:

**Exact computation** (Löffler’s algorithm, 1982). Uses a recurrence relation to compute the exact distribution of pairwise comparisons for small samples. Practical for combined sample sizes up to several hundred.

**Approximation** (Edgeworth expansion, 1955). Refines the normal approximation with correction terms based on higher moments of the distribution. Provides accurate results for large samples where exact computation becomes impractical.

The toolkit automatically selects the appropriate method based on sample sizes, ensuring both accuracy and computational efficiency.

This approach naturally complements Center and Spread:

$\text{Center}(\mathbf{x})$  uses the median of pairwise averages  $(x_i + x_j)/2$

$\text{Spread}(\mathbf{x})$  uses the median of pairwise differences  $|x_i - x_j|$

$\text{Shift}(\mathbf{x}, \mathbf{y})$  uses the median of pairwise differences  $x_i - y_j$

$\text{ShiftBounds}(\mathbf{x}, \mathbf{y}, \text{misrate})$  uses order statistics of the same pairwise differences

All procedures build on pairwise operations. This structural consistency reflects the mathematical unity underlying robust statistics: pairwise operations provide natural robustness while maintaining computational feasibility and statistical efficiency.

The inversion from hypothesis testing to bounds estimation represents a philosophical shift in statistical practice. Traditional methods ask “should I believe this specific hypothesis?” Pragmatic methods ask “what should I believe, given this data?” Bounds provide actionable answers: they tell practitioners which values are plausible, enabling informed decisions without arbitrary significance thresholds.

Traditional Mann-Whitney implementations apply tie correction when samples contain repeated values. This correction modifies variance calculations to account for tied observations, changing  $p$ -values and confidence intervals in ways that depend on measurement precision. The toolkit deliberately omits tie correction. Continuous distributions produce theoretically distinct values; observed ties result from finite measurement precision and digital representation. When measurements appear identical, this reflects rounding of underlying continuous variation, not true equality in the measured quantity. Treating ties as artifacts of discretization rather than distributional features simplifies computation while maintaining accuracy. The exact and approximate methods compute comparison distributions without requiring adjustments for tied values, eliminating a source of complexity and potential inconsistency in statistical practice.

### Historical Development

The mathematical foundations emerged through decades of refinement. Mann and Whitney (1947) established the distribution of pairwise comparisons under random sampling, creating the theoretical basis for comparing samples through rank-based methods. Their work demonstrated that comparison counts follow predictable patterns regardless of the underlying population distributions.

The original computational approaches suffered from severe limitations. Mann and Whitney proposed a slow exact method requiring exponential resources and a normal approximation that proved grossly inaccurate for practical use. The approximation works reasonably in distribution centers but fails catastrophically in the tails where practitioners most need accuracy. For moderate sample sizes, approximation errors can exceed factors of  $10^{11}$ .

Fix and Hodges (1955) addressed the approximation problem through higher-order corrections. Their expansion adds terms based on the distribution's actual moments rather than assuming perfect normality. This refinement reduces tail probability errors from orders of magnitude to roughly 1%, making approximation practical for large samples where exact computation becomes infeasible.

Löffler (1982) solved the exact computation problem through algorithmic innovation. The naive recurrence requires quadratic memory— infeasible for samples beyond a few dozen measurements. Löffler discovered a reformulation that reduces memory to linear scale, making exact computation practical for combined sample sizes up to several hundred.

Despite these advances, most statistical software continues using the 1947 approximation. The computational literature contains the solutions, but software implementations lag decades behind theoretical developments. This toolkit implements both the exact method for small samples and the refined approximation for large samples, automatically selecting the appropriate approach based on sample sizes.

The shift from hypothesis testing to bounds estimation requires no new mathematics. The same comparison distributions that enable hypothesis tests also determine which order statistics form reliable bounds. Traditional applications ask “is zero plausible?” and answer yes or no. This toolkit asks “which values are plausible?” and answers with an interval. The perspective inverts while the mathematical foundation remains identical.



## 4.8. Additive ('Normal') Distribution

The Additive ('Normal') distribution has two parameters: the mean and the standard deviation, written as Additive(mean, stdDev).

### Asymptotic Spread Value

Consider two independent draws  $X$  and  $Y$  from the Additive(mean, stdDev) distribution. The goal is to find the median of their absolute difference  $|X - Y|$ . Define the difference  $D = X - Y$ . By linearity of expectation,  $\mathbb{E}[D] = 0$ . By independence,  $\text{Var}[D] = 2 \cdot \text{stdDev}^2$ . Thus  $D$  has distribution Additive(0,  $\sqrt{2} \cdot \text{stdDev}$ ), and the problem reduces to finding the median of  $|D|$ . The location parameter mean disappears, as expected, because absolute differences are invariant to shifts.

Let  $\tau = \sqrt{2} \cdot \text{stdDev}$ , so that  $D \sim \text{Additive}(0, \tau)$ . The random variable  $|D|$  then follows the Half-Additive ('Folded Normal') distribution with scale  $\tau$ . Its cumulative distribution function for  $z \geq 0$  becomes

$$F_{|D|}(z) = \Pr(|D| \leq z) = 2\Phi(z/\tau) - 1$$

where  $\Phi$  denotes the standard Additive ('Normal') CDF.

The median  $m$  is the point at which this cdf equals 1/2. Setting  $F_{|D|}(m) = 1/2$  gives

$$2\Phi(m/\tau) - 1 = 1/2 \Rightarrow \Phi(m/\tau) = 3/4$$

Applying the inverse cdf yields  $m/\tau = z_{0.75}$ . Substituting back  $\tau = \sqrt{2} \cdot \text{stdDev}$  produces

$$\text{Median}(|X - Y|) = \sqrt{2} \cdot z_{0.75} \cdot \text{stdDev}$$

Define  $z_{0.75} := \Phi^{-1}(0.75) \approx 0.6744897502$ . Numerically, the median absolute difference is approximately  $\sqrt{2} \cdot z_{0.75} \cdot \text{stdDev} \approx 0.9538725524 \cdot \text{stdDev}$ . This expression depends only on the scale parameter stdDev, not on the mean, reflecting the translation invariance of the problem.

### Lemma: Average Estimator Drift Formula

For average estimators  $T_n$  with asymptotic standard deviation  $a \cdot \text{stdDev} / \sqrt{n}$  around the mean  $\mu$ , define  $\text{RelSpread}[T_n] := \text{Spread}[T_n] / \text{Spread}[X]$ . In the Additive ('Normal') case,  $\text{Spread}[X] = \sqrt{2} \cdot z_{0.75} \cdot \text{stdDev}$ .

For any average estimator  $T_n$  with asymptotic standard deviation  $a \cdot \text{stdDev} / \sqrt{n}$  around the mean  $\mu$ , the drift calculation follows:

The spread of two independent estimates:  $\text{Spread}[T_n] = \sqrt{2} \cdot z_{0.75} \cdot a \cdot \text{stdDev} / \sqrt{n}$

The relative spread:  $\text{RelSpread}[T_n] = a / \sqrt{n}$

The asymptotic drift:  $\text{Drift}(T, X) = a$

### Asymptotic Mean Drift

For the sample mean  $\text{Mean}(\mathbf{x}) = 1/n \sum_{i=1}^n x_i$  applied to samples from Additive(mean, stdDev), the sampling distribution of Mean is also additive with mean mean and standard deviation  $\text{stdDev}/\sqrt{n}$ .

Using the lemma with  $a = 1$  (since the standard deviation is  $\text{stdDev}/\sqrt{n}$ ):

$$\text{Drift}(\text{Mean}, X) = 1$$

Mean achieves unit drift under the Additive ('Normal') distribution, serving as the natural baseline for comparison. Mean is the optimal estimator under the Additive ('Normal') distribution: no other estimator achieves lower Drift.

### Asymptotic Median Drift

For the sample median  $\text{Median}(\mathbf{x})$  applied to samples from Additive(mean, stdDev), the asymptotic sampling distribution of Median is approximately Additive ('Normal') with mean mean and standard deviation  $\sqrt{\pi/2} \cdot \text{stdDev}/\sqrt{n}$ .

This result follows from the asymptotic theory of order statistics. For the median of a sample from a continuous distribution with density  $f$  and cumulative distribution  $F$ , the asymptotic variance is  $1/\left(4n[f(F^{-1}(0.5))]^2\right)$ . For the Additive ('Normal') distribution with standard deviation  $\text{stdDev}$ , the density at the median (which equals the mean) is  $1/(\text{stdDev} \sqrt{2\pi})$ . Thus the asymptotic variance becomes  $\pi \cdot \text{stdDev}^2/(2n)$ .

Using the lemma with  $a = \sqrt{\pi/2}$ :

$$\text{Drift}(\text{Median}, X) = \sqrt{\pi/2}$$

Numerically,  $\sqrt{\pi/2} \approx 1.2533$ , so the median has approximately 25% higher drift than the mean under the Additive ('Normal') distribution.

### Asymptotic Center Drift

For the sample center  $\text{Center}(\mathbf{x}) = \text{Median}_{1 \leq i \leq j \leq n} (x_i + x_j)/2$  applied to samples from Additive(mean, stdDev), its asymptotic sampling distribution must be determined.

The center estimator computes all pairwise averages (including  $i = j$ ) and takes their median. For the Additive ('Normal') distribution, asymptotic theory shows that the center estimator is asymptotically Additive ('Normal') with mean mean.

The exact asymptotic variance of the center estimator for the Additive ('Normal') distribution is:

$$\text{Var}[\text{Center}(X_{1:n})] = (\pi \cdot \text{stdDev}^2)/(3n)$$

This gives an asymptotic standard deviation of:

$$\text{StdDev}[\text{Center}(X_{1:n})] = \sqrt{\pi/3} \cdot \text{stdDev}/\sqrt{n}$$

Using the lemma with  $a = \sqrt{\pi/3}$ :

$$\text{Drift}(\text{Center}, X) = \sqrt{\pi/3}$$

Numerically,  $\sqrt{\pi/3} \approx 1.0233$ , so the center estimator achieves a drift very close to 1 under the Additive ('Normal') distribution, performing nearly as well as the mean while offering greater robustness to outliers.

### Lemma: Dispersion Estimator Drift Formula

For dispersion estimators  $T_n$  with asymptotic center  $b \cdot \text{stdDev}$  and standard deviation  $a \cdot \text{stdDev} / \sqrt{n}$ , define  $\text{RelSpread}[T_n] := \text{Spread}[T_n] / (b \cdot \text{stdDev})$ .

For any dispersion estimator  $T_n$  with asymptotic distribution  $T_n \sim \text{approx Additive}(b \cdot \text{stdDev}, (a \cdot \text{stdDev})^2/n)$ , the drift calculation follows:

The spread of two independent estimates:  $\text{Spread}[T_n] = \sqrt{2} \cdot z_{0.75} \cdot a \cdot \text{stdDev} / \sqrt{n}$

The relative spread:  $\text{RelSpread}[T_n] = \sqrt{2} \cdot z_{0.75} \cdot a / (b\sqrt{n})$

The asymptotic drift:  $\text{Drift}(T, X) = \sqrt{2} \cdot z_{0.75} \cdot a/b$

Note: The  $\sqrt{2}$  factor comes from the standard deviation of the difference  $D = T_1 - T_2$  of two independent estimates, and the  $z_{0.75}$  factor converts this standard deviation to the median absolute difference.

### Asymptotic StdDev Drift

For the sample standard deviation  $\text{StdDev}(\mathbf{x}) = \sqrt{1/(n-1) \sum_{i=1}^n (x_i - \text{Mean}(\mathbf{x}))^2}$  applied to samples from Additive(mean, stdDev), the sampling distribution of StdDev is approximately Additive ('Normal') for large  $n$  with mean stdDev and standard deviation  $\text{stdDev}/\sqrt{2n}$ .

Applying the lemma with  $a = 1/\sqrt{2}$  and  $b = 1$ :

$$\text{Spread}[\text{StdDev}(X_{1:n})] = \sqrt{2} \cdot z_{0.75} \cdot 1/\sqrt{2} \cdot \text{stdDev}/\sqrt{n} = z_{0.75} \cdot \text{stdDev}/\sqrt{n}$$

For the dispersion drift, we use the relative spread formula:

$$\text{RelSpread}[\text{StdDev}(X_{1:n})] = \text{Spread}[\text{StdDev}(X_{1:n})] / \text{Center}[\text{StdDev}(X_{1:n})]$$

Since  $\text{Center}[\text{StdDev}(X_{1:n})] \approx \text{stdDev}$  asymptotically:

$$\text{RelSpread}[\text{StdDev}(X_{1:n})] = (z_{0.75} \cdot \text{stdDev}/\sqrt{n}) / \text{stdDev} = z_{0.75}/\sqrt{n}$$

Therefore:

$$\text{Drift}(\text{StdDev}, X) = \lim_{n \rightarrow \infty} \sqrt{n} \cdot \text{RelSpread}[\text{StdDev}(X_{1:n})] = z_{0.75}$$

Numerically,  $z_{0.75} \approx 0.67449$ .

### Asymptotic MAD Drift

For the median absolute deviation  $\text{MAD}(\mathbf{x}) = \text{Median}(|x_i - \text{Median}(\mathbf{x})|)$  applied to samples from Additive(mean, stdDev), the asymptotic distribution is approximately Additive ('Normal').

For the Additive ('Normal') distribution, the population MAD equals  $z_{0.75} \cdot \text{stdDev}$ . The asymptotic standard deviation of the sample MAD is:

$$\text{StdDev}[\text{MAD}(X_{1:n})] = c_{\text{mad}} \cdot \text{stdDev} / \sqrt{n}$$

where  $c_{\text{mad}} \approx 0.78$ .

Applying the lemma with  $a = c_{\text{mad}}$  and  $b = z_{0.75}$ :

$$\text{Spread}[\text{MAD}(X_{1:n})] = \sqrt{2} \cdot z_{0.75} \cdot c_{\text{mad}} \cdot \text{stdDev} / \sqrt{n}$$

Since  $\text{Center}[\text{MAD}(X_{1:n})] \approx z_{0.75} \cdot \text{stdDev}$  asymptotically:

$$\text{RelSpread}[\text{MAD}(X_{1:n})] = (\sqrt{2} \cdot z_{0.75} \cdot c_{\text{mad}} \cdot \text{stdDev} / \sqrt{n}) / (z_{0.75} \cdot \text{stdDev}) = (\sqrt{2} \cdot c_{\text{mad}}) / \sqrt{n}$$

Therefore:

$$\text{Drift}(\text{MAD}, X) = \lim_{n \rightarrow \infty} \sqrt{n} \cdot \text{RelSpread}[\text{MAD}(X_{1:n})] = \sqrt{2} \cdot c_{\text{mad}}$$

Numerically,  $\sqrt{2} \cdot c_{\text{mad}} \approx \sqrt{2} \cdot 0.78 \approx 1.10$ .

### Asymptotic Spread Drift

For the sample spread  $\text{Spread}(\mathbf{x}) = \text{Median}_{1 \leq i < j \leq n} |x_i - x_j|$  applied to samples from Additive(mean, stdDev), the asymptotic distribution is approximately Additive ('Normal').

The spread estimator computes all pairwise absolute differences and takes their median. For the Additive ('Normal') distribution, the population spread equals  $\sqrt{2} \cdot z_{0.75} \cdot \text{stdDev}$  as derived in the Asymptotic Spread Value section.

The asymptotic standard deviation of the sample spread for the Additive ('Normal') distribution is:

$$\text{StdDev}[\text{Spread}(X_{1:n})] = c_{\text{spr}} \cdot \text{stdDev} / \sqrt{n}$$

where  $c_{\text{spr}} \approx 0.72$ .

Applying the lemma with  $a = c_{\text{spr}}$  and  $b = \sqrt{2} \cdot z_{0.75}$ :

$$\text{Spread}[\text{Spread}(X_{1:n})] = \sqrt{2} \cdot z_{0.75} \cdot c_{\text{spr}} \cdot \text{stdDev} / \sqrt{n}$$

Since  $\text{Center}[\text{Spread}(X_{1:n})] \approx \sqrt{2} \cdot z_{0.75} \cdot \text{stdDev}$  asymptotically:

$$\text{RelSpread}[\text{Spread}(X_{1:n})] = (\sqrt{2} \cdot z_{0.75} \cdot c_{\text{spr}} \cdot \text{stdDev} / \sqrt{n}) / (\sqrt{2} \cdot z_{0.75} \cdot \text{stdDev}) = c_{\text{spr}} / \sqrt{n}$$

Therefore:

$$\text{Drift}(\text{Spread}, X) = \lim_{n \rightarrow \infty} \sqrt{n} \cdot \text{RelSpread}[\text{Spread}(X_{1:n})] = c_{\text{spr}}$$

Numerically,  $c_{\text{spr}} \approx 0.72$ .

## Summary

### Summary for average estimators:

Estimator	$\text{Drift}(E, X)$	$\text{Drift}^2(E, X)$	$1 / \text{Drift}^2(E, X)$
Mean	1	1	1
Median	$\approx 1.253$	$\pi/2 \approx 1.571$	$2/\pi \approx 0.637$
Center	$\approx 1.023$	$\pi/3 \approx 1.047$	$3/\pi \approx 0.955$

The squared drift values indicate the sample size adjustment needed when switching estimators. For instance, switching from Mean to Median while maintaining the same precision requires increasing the sample size by a factor of  $\pi/2 \approx 1.571$  (about 57% more observations). Similarly, switching from Mean to Center requires only about 5% more observations.

The inverse squared drift (rightmost column) equals the classical statistical efficiency relative to the Mean. The Mean achieves optimal performance (unit efficiency) for the Additive ('Normal') distribution, as expected from classical theory. The Center maintains 95.5% efficiency while offering greater robustness to outliers, making it an attractive alternative when some contamination is possible. The Median, while most robust, operates at only 63.7% efficiency under purely Additive ('Normal') conditions.

### Summary for dispersion estimators:

For the Additive ('Normal') distribution, the asymptotic drift values reveal the relative precision of different dispersion estimators:

Estimator	$\text{Drift}(E, X)$	$\text{Drift}^2(E, X)$	$1 / \text{Drift}^2(E, X)$
StdDev	$\approx 0.67$	$\approx 0.45$	$\approx 2.22$
MAD	$\approx 1.10$	$\approx 1.22$	$\approx 0.82$
Spread	$\approx 0.72$	$\approx 0.52$	$\approx 1.92$

The squared drift values indicate the sample size adjustment needed when switching estimators. For instance, switching from StdDev to MAD while maintaining the same precision requires increasing the sample size by a factor of  $1.22/0.45 \approx 2.71$  (more than doubling the observations). Similarly, switching from StdDev to Spread requires a factor of  $0.52/0.45 \approx 1.16$ .

The StdDev achieves optimal performance for the Additive ('Normal') distribution. The MAD requires about 2.7 times more data to match the StdDev precision while offering greater robustness to outliers. The Spread requires about 1.16 times more data to match the StdDev precision under purely Additive ('Normal') conditions while maintaining robustness.

## 5. Reference Implementations

### 5.1. Python

Install from PyPI:

```
pip install pragmastat==5.0.0
```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/py>

Pragmastat on PyPI: <https://pypi.org/project/pragmastat/>

Demo:

```
from pragmastat import (
    Rng,
    median,
    center,
    spread,
    rel_spread,
    shift,
    ratio,
    avg_spread,
    disparity,
    pairwise_margin,
    shift_bounds,
)
from pragmastat.distributions import Additive, Exp, Multiplic, Power, Uniform

def main():
    # --- Randomization ---

    rng = Rng(1729)
    print(rng.uniform()) # 0.3943034703296536
    print(rng.uniform()) # 0.5730893757071377

    rng = Rng("experiment-1")
    print(rng.uniform()) # 0.9535207726895857

    rng = Rng(1729)
    print(rng.sample([0, 1, 2, 3, 4, 5, 6, 7, 8, 9], 3)) # [6, 8, 9]

    rng = Rng(1729)
    print(rng.shuffle([1, 2, 3, 4, 5])) # [4, 2, 3, 5, 1]

    # --- Distribution Sampling ---

    rng = Rng(1729)
    dist = Uniform(0, 10)
    print(dist.sample(rng)) # 3.9430347032965365

    rng = Rng(1729)
```

```
dist = Additive(0, 1)
print(dist.sample(rng)) # -1.222932972163442

rng = Rng(1729)
dist = Exp(1)
print(dist.sample(rng)) # 0.5013761944646019

rng = Rng(1729)
dist = Power(1, 2)
print(dist.sample(rng)) # 1.284909255071668

rng = Rng(1729)
dist = Multiplic(0, 1)
print(dist.sample(rng)) # 0.2943655336550937

# --- Single-Sample Statistics ---

x = [0, 2, 4, 6, 8]

print(median(x)) # 4
print(center(x)) # 4
print(spread(x)) # 4
print(spread([v + 10 for v in x])) # 4
print(spread([v * 2 for v in x])) # 8
print(rel_spread(x)) # 1

# --- Two-Sample Comparison ---

x = [0, 3, 6, 9, 12]
y = [0, 2, 4, 6, 8]

print(shift(x, y)) # 2
print(shift(y, x)) # -2
print(avg_spread(x, y)) # 5
print(disparity(x, y)) # 0.4
print(disparity(y, x)) # -0.4

x = [1, 2, 4, 8, 16]
y = [2, 4, 8, 16, 32]
print(ratio(x, y)) # 0.5

# --- Confidence Bounds ---

x = list(range(1, 31))
y = list(range(21, 51))

print(pairwise_margin(30, 30, 1e-4)) # 390
print(shift(x, y)) # -20
bounds = shift_bounds(x, y, 1e-4) # [-30, -10]
print(f"Bounds(lower={bounds.lower}, upper={bounds.upper})")
```

```
if __name__ == "__main__":  
    main()
```

## 5.2. TypeScript

Install from npm:

```
npm i pragmastat@5.0.0
```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/ts>

Pragmastat on npm: <https://www.npmjs.com/package/pragmastat>

Demo:

```
import {  
    median, center, spread, relSpread, shift, ratio, avgSpread, disparity, shiftBounds,  
    pairwiseMargin,  
    Rng, Uniform, Additive, Exp, Power, Multiplic  
} from '../src';  
  
function main() {  
    // --- Randomization ---  
  
    let rng = new Rng(1729);  
    console.log(rng.uniform()); // 0.3943034703296536  
    console.log(rng.uniform()); // 0.5730893757071377  
  
    rng = new Rng("experiment-1");  
    console.log(rng.uniform()); // 0.9535207726895857  
  
    rng = new Rng(1729);  
    console.log(rng.sample([0, 1, 2, 3, 4, 5, 6, 7, 8, 9], 3)); // [6, 8, 9]  
  
    rng = new Rng(1729);  
    console.log(rng.shuffle([1, 2, 3, 4, 5])); // [4, 2, 3, 5, 1]  
  
    // --- Distribution Sampling ---  
  
    rng = new Rng(1729);  
    let dist = new Uniform(0, 10);  
    console.log(dist.sample(rng)); // 3.9430347032965365  
  
    rng = new Rng(1729);  
    let addDist = new Additive(0, 1);  
    console.log(addDist.sample(rng)); // -1.222932972163442  
  
    rng = new Rng(1729);  
    let expDist = new Exp(1);  
    console.log(expDist.sample(rng)); // 0.5013761944646019  
  
    rng = new Rng(1729);  
    let powDist = new Power(1, 2);
```



```

    console.log(powDist.sample(rng)); // 1.284909255071668

    rng = new Rng(1729);
    let mulDist = new Multiplic(0, 1);
    console.log(mulDist.sample(rng)); // 0.2943655336550937

    // --- Single-Sample Statistics ---

    let x = [0, 2, 4, 6, 8];

    console.log(median(x)); // 4
    console.log(center(x)); // 4
    console.log(spread(x)); // 4
    console.log(spread(x.map(v => v + 10))); // 4
    console.log(spread(x.map(v => v * 2))); // 8
    console.log(relSpread(x)); // 1

    // --- Two-Sample Comparison ---

    x = [0, 3, 6, 9, 12];
    let y = [0, 2, 4, 6, 8];

    console.log(shift(x, y)); // 2
    console.log(shift(y, x)); // -2
    console.log(avgSpread(x, y)); // 5
    console.log(disparity(x, y)); // 0.4
    console.log(disparity(y, x)); // -0.4

    x = [1, 2, 4, 8, 16];
    y = [2, 4, 8, 16, 32];
    console.log(ratio(x, y)); // 0.5

    // --- Confidence Bounds ---

    x = Array.from({ length: 30 }, (_, i) => i + 1);
    y = Array.from({ length: 30 }, (_, i) => i + 21);

    console.log(pairwiseMargin(30, 30, 1e-4)); // 390
    console.log(shift(x, y)); // -20
    console.log(shiftBounds(x, y, 1e-4)); // [-30, -10]
  }

  main();

```

### 5.3. R

Install from GitHub:

```

install.packages("remotes") # If 'remotes' is not installed
remotes::install_github("AndreyAkinshin/pragmastat",
                        subdir = "r/pragmastat", ref = "v5.0.0")
library(pragmastat)

```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/r>

Demo:

```
library(pragmastat)

# --- Randomization ---

r <- rng(1729)
print(r$uniform()) # 0.3943034703296536
print(r$uniform()) # 0.5730893757071377

r <- rng("experiment-1")
print(r$uniform()) # 0.9535207726895857

r <- rng(1729)
print(r$sample(0:9, 3)) # [6, 8, 9]

r <- rng(1729)
print(r$shuffle(c(1, 2, 3, 4, 5))) # [4, 2, 3, 5, 1]

# --- Distribution Sampling ---

r <- rng(1729)
dist <- dist_uniform(0, 10)
print(dist$sample(r)) # 3.9430347032965365

r <- rng(1729)
dist <- dist_additive(0, 1)
print(dist$sample(r)) # -1.222932972163442

r <- rng(1729)
dist <- dist_exp(1)
print(dist$sample(r)) # 0.5013761944646019

r <- rng(1729)
dist <- dist_power(1, 2)
print(dist$sample(r)) # 1.284909255071668

r <- rng(1729)
dist <- dist_multiplic(0, 1)
print(dist$sample(r)) # 0.2943655336550937

# --- Single-Sample Statistics ---

x <- c(0, 2, 4, 6, 8)

print(median(x)) # 4
print(center(x)) # 4
print(spread(x)) # 4
print(spread(x + 10)) # 4
print(spread(x * 2)) # 8
print(rel_spread(x)) # 1
```

```
# --- Two-Sample Comparison ---

x <- c(0, 3, 6, 9, 12)
y <- c(0, 2, 4, 6, 8)

print(shift(x, y)) # 2
print(shift(y, x)) # -2
print(avg_spread(x, y)) # 5
print(disparity(x, y)) # 0.4
print(disparity(y, x)) # -0.4

x <- c(1, 2, 4, 8, 16)
y <- c(2, 4, 8, 16, 32)
print(ratio(x, y)) # 0.5

# --- Confidence Bounds ---

x <- 1:30
y <- 21:50

print(pairwise_margin(30, 30, 1e-4)) # 390
print(shift(x, y)) # -20
bounds <- shift_bounds(x, y, 1e-4) # [-30, -10]
print(paste("[", bounds$lower, ", ", bounds$upper, "]", sep=""))
```

## 5.4. C#

Install from NuGet via .NET CLI:

```
dotnet add package Pragmastat --version 5.0.0
```

Install from NuGet via Package Manager Console:

```
NuGet\Install-Package Pragmastat -Version 5.0.0
```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/cs>

Pragmastat on NuGet: <https://www.nuget.org/packages/Pragmastat/>

Demo:

```
using static System.Console;
using Pragmastat.Distributions;
using Pragmastat.Functions;
using Pragmastat.Randomization;

namespace Pragmastat.Demo;

class Program
{
    static void Main()
    {
```

```
// --- Randomization ---

var rng = new Rng(1729);
WriteLine(rng.Uniform()); // 0.3943034703296536
WriteLine(rng.Uniform()); // 0.5730893757071377

rng = new Rng("experiment-1");
WriteLine(rng.Uniform()); // 0.9535207726895857

rng = new Rng(1729);
WriteLine(string.Join(", ", rng.Sample([0.0, 1, 2, 3, 4, 5, 6, 7, 8, 9], 3))); // 6,
8, 9

rng = new Rng(1729);
WriteLine(string.Join(", ", rng.Shuffle([1.0, 2, 3, 4, 5]))); // 4, 2, 3, 5, 1

// --- Distribution Sampling ---

rng = new Rng(1729);
IDistribution dist = new Uniform(0, 10);
WriteLine(dist.Sample(rng)); // 3.9430347032965365

rng = new Rng(1729);
dist = new Additive(0, 1);
WriteLine(dist.Sample(rng)); // -1.222932972163442

rng = new Rng(1729);
dist = new Exp(1);
WriteLine(dist.Sample(rng)); // 0.5013761944646019

rng = new Rng(1729);
dist = new Power(1, 2);
WriteLine(dist.Sample(rng)); // 1.284909255071668

rng = new Rng(1729);
dist = new Multiplic(0, 1);
WriteLine(dist.Sample(rng)); // 0.2943655336550937

// --- Single-Sample Statistics ---

var x = new Sample(0, 2, 4, 6, 8);

WriteLine(Toolkit.Median(x)); // 4
WriteLine(x.Center()); // 4
WriteLine(x.Spread()); // 4
WriteLine((x + 10).Spread()); // 4
WriteLine((x * 2).Spread()); // 8
WriteLine(x.RelSpread()); // 1

// --- Two-Sample Comparison ---

x = new Sample(0, 3, 6, 9, 12);
var y = new Sample(0, 2, 4, 6, 8);
```

```

WriteLine(Toolkit.Shift(x, y)); // 2
WriteLine(Toolkit.Shift(y, x)); // -2
WriteLine(Toolkit.AvgSpread(x, y)); // 5
WriteLine(Toolkit.Disparity(x, y)); // 0.4
WriteLine(Toolkit.Disparity(y, x)); // -0.4

x = new Sample(1, 2, 4, 8, 16);
y = new Sample(2, 4, 8, 16, 32);
WriteLine(Toolkit.Ratio(x, y)); // 0.5

// --- Confidence Bounds ---

x = new Sample(
    1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15,
    16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30);
y = new Sample(
    21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35,
    36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50);

WriteLine(PairwiseMargin.Instance.Calc(30, 30, 1e-4)); // 390
WriteLine(Toolkit.Shift(x, y)); // -20
WriteLine(Toolkit.ShiftBounds(x, y, 1e-4)); // [-30, -10]
}
}

```

## 5.5. Kotlin

Install from Maven Central Repository via Apache Maven:

```

<dependency>
  <groupId>dev.pragmastat</groupId>
  <artifactId>pragmastat</artifactId>
  <version>5.0.0</version>
</dependency>

```

Install from Maven Central Repository via Gradle:

```
implementation 'dev.pragmastat:pragmastat:5.0.0'
```

Install from Maven Central Repository via Gradle (Kotlin):

```
implementation("dev.pragmastat:pragmastat:5.0.0")
```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/kt>

Pragmastat on Maven Central Repository: <https://central.sonatype.com/artifact/dev.pragmastat/pragmastat/overview>

Demo:

```
package dev.pragmastat.demo
```

```
import dev.pragmastat.*
```

```
import dev.pragmastat.distributions.*

fun main() {
    // --- Randomization ---

    var rng = Rng(1729)
    println(rng.uniform()) // 0.3943034703296536
    println(rng.uniform()) // 0.5730893757071377

    rng = Rng("experiment-1")
    println(rng.uniform()) // 0.9535207726895857

    rng = Rng(1729)
    println(rng.sample(listOf(0.0, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0), 3)) //
[6, 8, 9]

    rng = Rng(1729)
    println(rng.shuffle(listOf(1.0, 2.0, 3.0, 4.0, 5.0))) // [4, 2, 3, 5, 1]

    // --- Distribution Sampling ---

    rng = Rng(1729)
    var dist: Distribution = Uniform(0.0, 10.0)
    println(dist.sample(rng)) // 3.9430347032965365

    rng = Rng(1729)
    dist = Additive(0.0, 1.0)
    println(dist.sample(rng)) // -1.222932972163442

    rng = Rng(1729)
    dist = Exp(1.0)
    println(dist.sample(rng)) // 0.5013761944646019

    rng = Rng(1729)
    dist = Power(1.0, 2.0)
    println(dist.sample(rng)) // 1.284909255071668

    rng = Rng(1729)
    dist = Multiplic(0.0, 1.0)
    println(dist.sample(rng)) // 0.2943655336550937

    // --- Single-Sample Statistics ---

    var x = listOf(0.0, 2.0, 4.0, 6.0, 8.0)

    println(median(x)) // 4
    println(center(x)) // 4
    println(spread(x)) // 4
    println(spread(x.map { it + 10 }))) // 4
    println(spread(x.map { it * 2 }))) // 8
    println(relSpread(x)) // 1

    // --- Two-Sample Comparison ---
```

```

x = listOf(0.0, 3.0, 6.0, 9.0, 12.0)
var y = listOf(0.0, 2.0, 4.0, 6.0, 8.0)

println(shift(x, y)) // 2
println(shift(y, x)) // -2
println(avgSpread(x, y)) // 5
println(disparity(x, y)) // 0.4
println(disparity(y, x)) // -0.4

x = listOf(1.0, 2.0, 4.0, 8.0, 16.0)
y = listOf(2.0, 4.0, 8.0, 16.0, 32.0)
println(ratio(x, y)) // 0.5

// --- Confidence Bounds ---

x = (1..30).map { it.toDouble() }
y = (21..50).map { it.toDouble() }

println(pairwiseMargin(30, 30, 1e-4)) // 390
println(shift(x, y)) // -20
println(shiftBounds(x, y, 1e-4)) // [-30, -10]
}

```

## 5.6. Rust

Install from crates.io via cargo:

```
cargo add pragmastat@5.0.0
```

Install from crates.io via Cargo.toml:

```

[dependencies]
pragmastat = "5.0.0"

```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/rs>

Pragmastat on crates.io: <https://crates.io/crates/pragmastat>

Demo:

```

use pragmastat::distributions::{Additive, Distribution, Exp, Multiplic, Power, Uniform};
use pragmastat::*;

fn print(result: Result<f64, &str>) {
    println!("{}", result.unwrap());
}

fn add(x: &[f64], val: f64) -> Vec<f64> {
    x.iter().map(|v| v + val).collect()
}

fn multiply(x: &[f64], val: f64) -> Vec<f64> {

```

```

    x.iter().map(|v| v * val).collect()
}

fn main() {
    // --- Randomization ---

    let mut rng = Rng::from_seed(1729);
    println!("{}", rng.uniform()); // 0.3943034703296536
    println!("{}", rng.uniform()); // 0.5730893757071377

    let mut rng = Rng::from_string("experiment-1");
    println!("{}", rng.uniform()); // 0.9535207726895857

    let mut rng = Rng::from_seed(1729);
    println!(
        "{:?}",
        rng.sample(&[0.0, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0, 9.0], 3)
    ); // [6, 8, 9]

    let mut rng = Rng::from_seed(1729);
    println!("{:?}", rng.shuffle(&[1.0, 2.0, 3.0, 4.0, 5.0])); // [4, 2, 3, 5, 1]

    // --- Distribution Sampling ---

    let mut rng = Rng::from_seed(1729);
    let dist = Uniform::new(0.0, 10.0);
    println!("{}", dist.sample(&mut rng)); // 3.9430347032965365

    let mut rng = Rng::from_seed(1729);
    let dist = Additive::new(0.0, 1.0);
    println!("{}", dist.sample(&mut rng)); // -1.222932972163442

    let mut rng = Rng::from_seed(1729);
    let dist = Exp::new(1.0);
    println!("{}", dist.sample(&mut rng)); // 0.5013761944646019

    let mut rng = Rng::from_seed(1729);
    let dist = Power::new(1.0, 2.0);
    println!("{}", dist.sample(&mut rng)); // 1.284909255071668

    let mut rng = Rng::from_seed(1729);
    let dist = Multiplic::new(0.0, 1.0);
    println!("{}", dist.sample(&mut rng)); // 0.2943655336550937

    // --- Single-Sample Statistics ---

    let x = vec![0.0, 2.0, 4.0, 6.0, 8.0];

    print(median(&x)); // 4
    print(center(&x)); // 4
    print(spread(&x)); // 4
    print(spread(&add(&x, 10.0))); // 4
    print(spread(&multiply(&x, 2.0))); // 8

```



```

print(rel_spread(&x)); // 1

// --- Two-Sample Comparison ---

let x = vec![0.0, 3.0, 6.0, 9.0, 12.0];
let y = vec![0.0, 2.0, 4.0, 6.0, 8.0];

print(shift(&x, &y)); // 2
print(shift(&y, &x)); // -2
print(avg_spread(&x, &y)); // 5
print(disparity(&x, &y)); // 0.4
print(disparity(&y, &x)); // -0.4

let x = vec![1.0, 2.0, 4.0, 8.0, 16.0];
let y = vec![2.0, 4.0, 8.0, 16.0, 32.0];
print(ratio(&x, &y)); // 0.5

// --- Confidence Bounds ---

let x: Vec<f64> = (1..=30).map(|i| i as f64).collect();
let y: Vec<f64> = (21..=50).map(|i| i as f64).collect();

println!("{}", pairwise_margin(30, 30, 1e-4).unwrap()); // 390
print(shift(&x, &y)); // -20

let bounds = shift_bounds(&x, &y, 1e-4).unwrap(); // [-30, -10]
println!("{lower: {}, upper: {}}", bounds.lower, bounds.upper);
}

```

## 5.7. Go

Install from GitHub:

```
go get github.com/AndreyAkinshin/pragmastat/go/v4@v5.0.0
```

Source code: <https://github.com/AndreyAkinshin/pragmastat/tree/v5.0.0/go>

Demo:

```

package main

import (
    "fmt"
    "log"

    pragmastat "github.com/AndreyAkinshin/pragmastat/go/v4"
)

func must[T any](val T, err error) T {
    if err != nil {
        log.Fatal(err)
    }
    return val
}

```

```
}

func print(val float64, err error) {
    fmt.Println(must(val, err))
}

func add(x []float64, val float64) []float64 {
    result := make([]float64, len(x))
    for i, v := range x {
        result[i] = v + val
    }
    return result
}

func multiply(x []float64, val float64) []float64 {
    result := make([]float64, len(x))
    for i, v := range x {
        result[i] = v * val
    }
    return result
}

func main() {
    // --- Randomization ---

    rng := prng.NewRngFromSeed(1729)
    fmt.Println(rng.Uniform()) // 0.3943034703296536
    fmt.Println(rng.Uniform()) // 0.5730893757071377

    rng = prng.NewRngFromString("experiment-1")
    fmt.Println(rng.Uniform()) // 0.9535207726895857

    rng = prng.NewRngFromSeed(1729)
    fmt.Println(prng.Sample(rng, []float64{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}, 3)) // [6 8
9]

    rng = prng.NewRngFromSeed(1729)
    fmt.Println(prng.Shuffle(rng, []float64{1, 2, 3, 4, 5})) // [4 2 3 5 1]

    // --- Distribution Sampling ---

    rng = prng.NewRngFromSeed(1729)
    dist := prng.NewUniform(0, 10)
    fmt.Println(dist.Sample(rng)) // 3.9430347032965365

    rng = prng.NewRngFromSeed(1729)
    addDist := prng.NewAdditive(0, 1)
    fmt.Println(addDist.Sample(rng)) // -1.222932972163442

    rng = prng.NewRngFromSeed(1729)
    expDist := prng.NewExp(1)
    fmt.Println(expDist.Sample(rng)) // 0.5013761944646019
```

```

rng = pragmastat.NewRngFromSeed(1729)
powDist := pragmastat.NewPower(1, 2)
fmt.Println(powDist.Sample(rng)) // 1.284909255071668

rng = pragmastat.NewRngFromSeed(1729)
mulDist := pragmastat.NewMultiplic(0, 1)
fmt.Println(mulDist.Sample(rng)) // 0.2943655336550937

// --- Single-Sample Statistics ---

x := []float64{0, 2, 4, 6, 8}

print(pragmastat.Median(x))           // 4
print(pragmastat.Center(x))           // 4
print(pragmastat.Spread(x))           // 4
print(pragmastat.Spread(add(x, 10)))  // 4
print(pragmastat.Spread(multiply(x, 2))) // 8
print(pragmastat.RelSpread(x))        // 1

// --- Two-Sample Comparison ---

x = []float64{0, 3, 6, 9, 12}
y := []float64{0, 2, 4, 6, 8}

print(pragmastat.Shift(x, y))         // 2
print(pragmastat.Shift(y, x))         // -2
print(pragmastat.AvgSpread(x, y))     // 5
print(pragmastat.Disparity(x, y))     // 0.4
print(pragmastat.Disparity(y, x))     // -0.4

x = []float64{1, 2, 4, 8, 16}
y = []float64{2, 4, 8, 16, 32}
print(pragmastat.Ratio(x, y)) // 0.5

// --- Confidence Bounds ---

x = []float64{
    1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15,
    16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30}
y = []float64{
    21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35,
    36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50}

margin, _ := pragmastat.PairwiseMargin(30, 30, 1e-4)
fmt.Println(margin) // 390
print(pragmastat.Shift(x, y)) // -20
fmt.Println(must(pragmastat.ShiftBounds(x, y, 1e-4))) // [-30, -10]
}

```

## 6. Reference Tests

### 6.1. Center Tests

$$\text{Center}(\mathbf{x}) = \text{median}_{1 \leq i \leq j \leq n} \frac{x_i + x_j}{2}$$

The Center test suite contains 38 correctness test cases stored in the repository (24 original + 14 unsorted), plus 1 performance test that should be implemented manually (see Test Framework).

**Demo examples** ( $n = 5$ ) — from manual introduction, validating properties:

demo-1:  $\mathbf{x} = (0, 2, 4, 6, 8)$ , expected output: 4 (base case)

demo-2:  $\mathbf{x} = (10, 12, 14, 16, 18)$  (= demo-1 + 10), expected output: 14 (location equivariance)

demo-3:  $\mathbf{x} = (0, 6, 12, 18, 24)$  ( $= 3 \times \text{demo-1}$ ), expected output: 12 (scale equivariance)

**Natural sequences** ( $n = 1, 2, 3, 4$ ) — canonical happy path examples:

natural-1:  $\mathbf{x} = (1)$ , expected output: 1

natural-2:  $\mathbf{x} = (1, 2)$ , expected output: 1.5

natural-3:  $\mathbf{x} = (1, 2, 3)$ , expected output: 2

natural-4:  $\mathbf{x} = (1, 2, 3, 4)$ , expected output: 2.5 (smallest even size with rich structure)

**Negative values** ( $n = 3$ ) — sign handling validation:

negative-3:  $\mathbf{x} = (-3, -2, -1)$ , expected output:  $-2$

**Zero values** ( $n = 1, 2$ ) — edge case testing with zeros:

zeros-1:  $\mathbf{x} = (0)$ , expected output: 0

zeros-2:  $\mathbf{x} = (0, 0)$ , expected output: 0

**Additive distribution** ( $n = 5, 10, 30$ ) — fuzzy testing with Additive(10, 1):

additive-5, additive-10, additive-30: random samples generated with seed 0

**Uniform distribution** ( $n = 5, 100$ ) — fuzzy testing with Uniform(0, 1):

uniform-5, uniform-100: random samples generated with seed 1

The random samples validate that Center performs correctly on realistic distributions at various sample sizes. The progression from small ( $n = 5$ ) to large ( $n = 100$ ) samples helps identify issues that only manifest at specific scales.

**Algorithm stress tests** — edge cases for fast algorithm implementation:

duplicates-5:  $\mathbf{x} = (3, 3, 3, 3, 3)$  (all identical, stress stall handling)

duplicates-10:  $\mathbf{x} = (1, 1, 1, 2, 2, 2, 3, 3, 3, 3)$  (many duplicates, stress tie-breaking)

parity-odd-7:  $\mathbf{x} = (1, 2, 3, 4, 5, 6, 7)$  (odd sample size for odd total pairs)

parity-even-6:  $\mathbf{x} = (1, 2, 3, 4, 5, 6)$  (even sample size for even total pairs)

parity-odd-49: 49-element sequence  $(1, 2, \dots, 49)$  (large odd, 1225 pairs)

parity-even-50: 50-element sequence (1, 2, ..., 50) (large even, 1275 pairs)

**Extreme values** — numerical stability and range tests:

extreme-large-5:  $x = (10^8, 2 \cdot 10^8, 3 \cdot 10^8, 4 \cdot 10^8, 5 \cdot 10^8)$  (very large values)

extreme-small-5:  $x = (10^{-8}, 2 \cdot 10^{-8}, 3 \cdot 10^{-8}, 4 \cdot 10^{-8}, 5 \cdot 10^{-8})$  (very small positive values)

extreme-wide-5:  $x = (0.001, 1, 100, 1000, 1000000)$  (wide range, tests precision)

**Unsorted tests** — verify sorting correctness (14 tests):

unsorted-reverse- $\{n\}$  for  $n \in \{2, 3, 4, 5, 7\}$ : reverse sorted natural sequences (5 tests)

unsorted-shuffle-3:  $x = (2, 1, 3)$  (middle element first)

unsorted-shuffle-4:  $x = (3, 1, 4, 2)$  (interleaved)

unsorted-shuffle-5:  $x = (5, 2, 4, 1, 3)$  (complex shuffle)

unsorted-last-first-5:  $x = (5, 1, 2, 3, 4)$  (last moved to first)

unsorted-first-last-5:  $x = (2, 3, 4, 5, 1)$  (first moved to last)

unsorted-duplicates-mixed-5:  $x = (3, 3, 3, 3, 3)$  (all identical, any order)

unsorted-duplicates-unsorted-10:  $x = (3, 1, 2, 3, 1, 3, 2, 1, 3, 2)$  (duplicates mixed)

unsorted-extreme-large-unsorted-5:  $x = (5 \cdot 10^8, 10^8, 4 \cdot 10^8, 2 \cdot 10^8, 3 \cdot 10^8)$  (large values unsorted)

unsorted-parity-odd-reverse-7:  $x = (7, 6, 5, 4, 3, 2, 1)$  (odd size reverse)

These tests ensure implementations correctly sort input data before computing pairwise averages. The variety of shuffle patterns (reverse, rotation, interleaving, single element displacement) catches common sorting bugs.

**Performance test** — validates the fast  $O(n \log n)$  algorithm:

**Input:**  $x = (1, 2, 3, \dots, 100000)$

**Expected output:** 50000.5

**Time constraint:** Must complete in under 5 seconds

**Purpose:** Ensures that the implementation uses the efficient algorithm rather than materializing all  $\binom{n+1}{2} \approx 5$  billion pairwise averages

This test case is not stored in the repository because it generates a large JSON file (approximately 1.5 MB). Each language implementation should manually implement this test with the hardcoded expected result.

## 6.2. Spread Tests

$$\text{Spread}(x) = \text{median}_{1 \leq i < j \leq n} |x_i - x_j|$$

The Spread test suite contains 38 correctness test cases stored in the repository (24 original + 14 unsorted), plus 1 performance test that should be implemented manually (see Test Framework).

**Demo examples** ( $n = 5$ ) — from manual introduction, validating properties:

demo-1:  $x = (0, 2, 4, 6, 8)$ , expected output: 4 (base case)

demo-2:  $x = (10, 12, 14, 16, 18)$  (= demo-1 + 10), expected output: 4 (location invariance)

demo-3:  $x = (0, 4, 8, 12, 16)$  ( $= 2 \times \text{demo-1}$ ), expected output: 8 (scale equivariance)

**Natural sequences** ( $n = 1, 2, 3, 4$ ):

natural-1:  $x = (1)$ , expected output: 0 (single element has zero dispersion)

natural-2:  $x = (1, 2)$ , expected output: 1

natural-3:  $x = (1, 2, 3)$ , expected output: 1

natural-4:  $x = (1, 2, 3, 4)$ , expected output: 1.5 (smallest even size with rich structure)

**Negative values** ( $n = 3$ ) — sign handling validation:

negative-3:  $x = (-3, -2, -1)$ , expected output: 1

**Zero values** ( $n = 1, 2$ ):

zeros-1:  $x = (0)$ , expected output: 0

zeros-2:  $x = (0, 0)$ , expected output: 0

**Additive distribution** ( $n = 5, 10, 30$ ) — Additive(10, 1):

additive-5, additive-10, additive-30: random samples generated with seed 0

**Uniform distribution** ( $n = 5, 100$ ) — Uniform(0, 1):

uniform-5, uniform-100: random samples generated with seed 1

The natural sequence cases validate the basic pairwise difference calculation. The zero cases confirm that constant samples correctly produce zero spread.

**Algorithm stress tests** — edge cases for fast algorithm implementation:

duplicates-5:  $x = (3, 3, 3, 3, 3)$  (all identical, expected output: 0)

duplicates-10:  $x = (1, 1, 1, 2, 2, 2, 3, 3, 3, 3)$  (many duplicates, stress tie-breaking)

parity-odd-7:  $x = (1, 2, 3, 4, 5, 6, 7)$  (odd sample size, 21 differences)

parity-even-6:  $x = (1, 2, 3, 4, 5, 6)$  (even sample size, 15 differences)

parity-odd-49: 49-element sequence  $(1, 2, \dots, 49)$  (large odd, 1176 differences)

parity-even-50: 50-element sequence  $(1, 2, \dots, 50)$  (large even, 1225 differences)

**Extreme values** — numerical stability and range tests:

extreme-large-5:  $x = (10^8, 2 \cdot 10^8, 3 \cdot 10^8, 4 \cdot 10^8, 5 \cdot 10^8)$  (very large values)

extreme-small-5:  $x = (10^{-8}, 2 \cdot 10^{-8}, 3 \cdot 10^{-8}, 4 \cdot 10^{-8}, 5 \cdot 10^{-8})$  (very small positive values)

extreme-wide-5:  $x = (0.001, 1, 100, 1000, 1000000)$  (wide range, tests precision)

**Unsorted tests** — verify sorting correctness (14 tests):

unsorted-reverse- $\{n\}$  for  $n \in \{2, 3, 4, 5, 7\}$ : reverse sorted natural sequences (5 tests)

unsorted-shuffle-3:  $x = (3, 1, 2)$  (rotated)

unsorted-shuffle-4:  $x = (4, 2, 1, 3)$  (mixed order)

unsorted-shuffle-5:  $x = (5, 1, 3, 2, 4)$  (partial shuffle)

unsorted-last-first-5:  $x = (5, 1, 2, 3, 4)$  (last moved to first)

unsorted-first-last-5:  $x = (2, 3, 4, 5, 1)$  (first moved to last)  
 unsorted-duplicates-mixed-5:  $x = (3, 3, 3, 3, 3)$  (all identical)  
 unsorted-duplicates-unsorted-10:  $x = (2, 3, 1, 3, 2, 1, 2, 3, 1, 3)$  (duplicates mixed)  
 unsorted-extreme-wide-unsorted-5:  $x = (1000, 0.001, 1000000, 100, 1)$  (wide range unsorted)  
 unsorted-negative-unsorted-5:  $x = (-1, -5, -2, -4, -3)$  (negative unsorted)

These tests verify that implementations correctly sort input before computing pairwise differences. Since Spread uses absolute differences, order-dependent bugs would manifest differently than in Center.

**Performance test** — validates the fast  $O(n \log n)$  algorithm:

**Input:**  $x = (1, 2, 3, \dots, 100000)$

**Expected output:** 29290

**Time constraint:** Must complete in under 5 seconds

**Purpose:** Ensures that the implementation uses the efficient algorithm rather than materializing all  $\binom{n}{2} \approx 5$  billion pairwise differences

This test case is not stored in the repository because it generates a large JSON file (approximately 1.5 MB). Each language implementation should manually implement this test with the hardcoded expected result.

### 6.3. RelSpread Tests

$$\text{RelSpread}(x) = \frac{\text{Spread}(x)}{|\text{Center}(x)|}$$

The RelSpread test suite contains 25 test cases (15 original + 10 unsorted) focusing on relative dispersion.

**Demo examples** ( $n = 5$ ) — from manual introduction, validating properties:

demo-1:  $x = (0, 2, 4, 6, 8)$ , expected output: 1 (base case)

demo-2:  $x = (0, 10, 20, 30, 40)$  ( $= 5 \times \text{demo-1}$ ), expected output: 1 (scale invariance)

**Natural sequences** ( $n = 1, 2, 3, 4$ ):

natural-1:  $x = (1)$ , expected output: 0

natural-2:  $x = (1, 2)$ , expected output:  $\approx 0.667$

natural-3:  $x = (1, 2, 3)$ , expected output: 0.5

natural-4:  $x = (1, 2, 3, 4)$ , expected output: 0.6 (validates composite with even size)

**Negative values** ( $n = 3$ ) — validates absolute value in denominator:

negative-3:  $x = (-3, -2, -1)$ , expected output: 0.5

**Uniform distribution** ( $n = 5, 10, 20, 30, 100$ ) — Uniform(0, 1):

uniform-5, uniform-10, uniform-20, uniform-30, uniform-100: random samples generated with seed 0

The uniform distribution tests span multiple sample sizes to verify that RelSpread correctly normalizes dispersion. The absence of zero-value tests reflects the domain constraint requiring  $\text{Center}(\mathbf{x}) \neq 0$ .

**Composite estimator stress tests** — edge cases specific to division operation:

composite-small-center:  $\mathbf{x} = (0.001, 0.002, 0.003, 0.004, 0.005)$  (small center, tests division stability)  
 composite-large-spread:  $\mathbf{x} = (1, 100, 200, 300, 1000)$  (large spread relative to center)  
 composite-extreme-ratio:  $\mathbf{x} = (1, 1.0001, 1.0002, 1.0003, 1.0004)$  (tiny spread, tests precision)

**Unsorted tests** — verify sorting for composite estimator (10 tests):

unsorted-reverse- $\{n\}$  for  $n \in \{3, 4, 5\}$ : reverse sorted natural sequences (3 tests)  
 unsorted-shuffle-4:  $\mathbf{x} = (4, 1, 3, 2)$  (mixed order)  
 unsorted-shuffle-5:  $\mathbf{x} = (5, 3, 1, 4, 2)$  (complex shuffle)  
 unsorted-negative-unsorted-3:  $\mathbf{x} = (-1, -3, -2)$  (negative unsorted)  
 unsorted-demo-unsorted-5:  $\mathbf{x} = (8, 0, 4, 2, 6)$  (demo case unsorted)  
 unsorted-composite-small-unsorted:  $\mathbf{x} = (0.005, 0.001, 0.003, 0.002, 0.004)$  (small center unsorted)  
 unsorted-composite-large-unsorted:  $\mathbf{x} = (1000, 1, 300, 100, 200)$  (large spread unsorted)  
 unsorted-extreme-ratio-unsorted-4:  $\mathbf{x} = (1.0003, 1, 1.0002, 1.0001)$  (extreme ratio unsorted)

Since RelSpread combines both Center and Spread, these tests verify that sorting works correctly for composite estimators.

## 6.4. Shift Tests

$$\text{Shift}(\mathbf{x}, \mathbf{y}) = \text{median}_{1 \leq i \leq n, 1 \leq j \leq m} (x_i - y_j)$$

The Shift test suite contains 60 correctness test cases stored in the repository (42 original + 18 unsorted), plus 1 performance test that should be implemented manually (see Test Framework).

**Demo examples** ( $n = m = 5$ ) — from manual introduction, validating properties:

demo-1:  $\mathbf{x} = (0, 2, 4, 6, 8)$ ,  $\mathbf{y} = (10, 12, 14, 16, 18)$ , expected output:  $-10$  (base case)  
 demo-2:  $\mathbf{x} = (0, 2, 4, 6, 8)$ ,  $\mathbf{y} = (0, 2, 4, 6, 8)$ , expected output:  $0$  (identity property)  
 demo-3:  $\mathbf{x} = (7, 9, 11, 13, 15)$ ,  $\mathbf{y} = (13, 15, 17, 19, 21)$  ( $= \text{demo-1} + [7, 3]$ ), expected output:  $-6$  (location equivariance)  
 demo-4:  $\mathbf{x} = (0, 4, 8, 12, 16)$ ,  $\mathbf{y} = (20, 24, 28, 32, 36)$  ( $= 2 \times \text{demo-1}$ ), expected output:  $-20$  (scale equivariance)  
 demo-5:  $\mathbf{x} = (10, 12, 14, 16, 18)$ ,  $\mathbf{y} = (0, 2, 4, 6, 8)$  ( $= \text{reversed demo-1}$ ), expected output:  $10$  (anti-symmetry)

**Natural sequences** ( $[n, m] \in \{1, 2, 3\} \times \{1, 2, 3\}$ ) — 9 combinations:

natural-1-1:  $\mathbf{x} = (1)$ ,  $\mathbf{y} = (1)$ , expected output:  $0$   
 natural-1-2:  $\mathbf{x} = (1)$ ,  $\mathbf{y} = (1, 2)$ , expected output:  $-0.5$   
 natural-1-3:  $\mathbf{x} = (1)$ ,  $\mathbf{y} = (1, 2, 3)$ , expected output:  $-1$



natural-2-1:  $x = (1, 2), y = (1)$ , expected output: 0.5  
 natural-2-2:  $x = (1, 2), y = (1, 2)$ , expected output: 0  
 natural-2-3:  $x = (1, 2), y = (1, 2, 3)$ , expected output: -0.5  
 natural-3-1:  $x = (1, 2, 3), y = (1)$ , expected output: 1  
 natural-3-2:  $x = (1, 2, 3), y = (1, 2)$ , expected output: 0.5  
 natural-3-3:  $x = (1, 2, 3), y = (1, 2, 3)$ , expected output: 0

**Negative values** ( $[n, m] = [2, 2]$ ) — sign handling validation:

negative-2-2:  $x = (-2, -1), y = (-2, -1)$ , expected output: 0

**Mixed-sign values** ( $[n, m] = [2, 2]$ ) — validates anti-symmetry across zero:

mixed-2-2:  $x = (-1, 1), y = (-1, 1)$ , expected output: 0

**Zero values** ( $[n, m] \in \{1, 2\} \times \{1, 2\}$ ) — 4 combinations:

zeros-1-1, zeros-1-2, zeros-2-1, zeros-2-2: all produce output 0

**Additive distribution** ( $[n, m] \in \{5, 10, 30\} \times \{5, 10, 30\}$ ) — 9 combinations with Additive(10, 1):

additive-5-5, additive-5-10, additive-5-30  
 additive-10-5, additive-10-10, additive-10-30  
 additive-30-5, additive-30-10, additive-30-30  
 Random generation: x uses seed 0, y uses seed 1

**Uniform distribution** ( $[n, m] \in \{5, 100\} \times \{5, 100\}$ ) — 4 combinations with Uniform(0, 1):

uniform-5-5, uniform-5-100, uniform-100-5, uniform-100-100  
 Random generation: x uses seed 2, y uses seed 3

The natural sequences validate anti-symmetry ( $\text{Shift}(x, y) = -\text{Shift}(y, x)$ ) and the identity property ( $\text{Shift}(x, x) = 0$ ). The asymmetric size combinations test the two-sample algorithm with unbalanced inputs.

**Algorithm stress tests** — edge cases for fast binary search algorithm:

duplicates-5-5:  $x = (3, 3, 3, 3, 3), y = (3, 3, 3, 3, 3)$  (all identical, expected output: 0)  
 duplicates-10-10:  $x = (1, 1, 2, 2, 3, 3, 4, 4, 5, 5), y = (1, 1, 2, 2, 3, 3, 4, 4, 5, 5)$  (many duplicates)  
 parity-odd-7-7:  $x = (1, 2, 3, 4, 5, 6, 7), y = (1, 2, 3, 4, 5, 6, 7)$  (odd sizes, 49 differences, expected output: 0)  
 parity-even-6-6:  $x = (1, 2, 3, 4, 5, 6), y = (1, 2, 3, 4, 5, 6)$  (even sizes, 36 differences, expected output: 0)  
 parity-asymmetric-7-6:  $x = (1, 2, 3, 4, 5, 6, 7), y = (1, 2, 3, 4, 5, 6)$  (mixed parity, 42 differences)  
 parity-large-49-50:  $x = (1, 2, \dots, 49), y = (1, 2, \dots, 50)$  (large asymmetric, 2450 differences)

**Extreme asymmetry** — tests with very unbalanced sample sizes:

asymmetry-1-100:  $x = (50), y = (1, 2, \dots, 100)$  (single vs many, 100 differences)  
 asymmetry-2-50:  $x = (10, 20), y = (1, 2, \dots, 50)$  (tiny vs medium, 100 differences)

asymmetry-constant-varied:  $x = (5, 5, 5, 5, 5)$ ,  $y = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$  (constant vs varied)

**Unsorted tests** — verify independent sorting of each sample (18 tests):

unsorted-x-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4), (5, 5)\}$ : X unsorted (reversed), Y sorted (3 tests)  
 unsorted-y-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4), (5, 5)\}$ : X sorted, Y unsorted (reversed) (3 tests)  
 unsorted-both-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4), (5, 5)\}$ : both unsorted (reversed) (3 tests)  
 unsorted-reverse-3-3:  $x = (3, 2, 1)$ ,  $y = (3, 2, 1)$  (both reversed)  
 unsorted-x-shuffle-3-3:  $x = (2, 1, 3)$ ,  $y = (1, 2, 3)$  (X shuffled, Y sorted)  
 unsorted-y-shuffle-3-3:  $x = (1, 2, 3)$ ,  $y = (3, 1, 2)$  (X sorted, Y shuffled)  
 unsorted-both-shuffle-4-4:  $x = (3, 1, 4, 2)$ ,  $y = (4, 2, 1, 3)$  (both shuffled)  
 unsorted-duplicates-mixed-5-5:  $x = (3, 3, 3, 3, 3)$ ,  $y = (3, 3, 3, 3, 3)$  (all identical)  
 unsorted-x-unsorted-duplicates:  $x = (2, 1, 3, 2, 1)$ ,  $y = (1, 1, 2, 2, 3)$  (X has unsorted duplicates)  
 unsorted-y-unsorted-duplicates:  $x = (1, 1, 2, 2, 3)$ ,  $y = (3, 2, 1, 3, 2)$  (Y has unsorted duplicates)  
 unsorted-asymmetric-unsorted-2-5:  $x = (2, 1)$ ,  $y = (5, 2, 4, 1, 3)$  (asymmetric sizes, both unsorted)  
 unsorted-negative-unsorted-3-3:  $x = (-1, -3, -2)$ ,  $y = (-2, -3, -1)$  (negative unsorted)

These tests are critical for two-sample estimators because they verify that  $x$  and  $y$  are sorted **independently**. The variety includes cases where only one sample is unsorted, ensuring implementations don't incorrectly assume pre-sorted input or sort samples together.

**Performance test** — validates the fast  $O((m + n) \log L)$  binary search algorithm:

**Input:**  $x = (1, 2, 3, \dots, 100000)$ ,  $y = (1, 2, 3, \dots, 100000)$

**Expected output:** 0

**Time constraint:** Must complete in under 5 seconds

**Purpose:** Ensures that the implementation uses the efficient algorithm rather than materializing all  $mn = 10$  billion pairwise differences

This test case is not stored in the repository because it generates a large JSON file (approximately 1.5 MB). Each language implementation should manually implement this test with the hardcoded expected result.

## 6.5. Ratio Tests

$$\text{Ratio}(x, y) = \text{median}_{1 \leq i \leq n, 1 \leq j \leq m} \left( \frac{x_i}{y_j} \right)$$

The Ratio test suite contains 37 test cases (25 original + 12 unsorted), excluding zero values due to division constraints.

**Demo examples** ( $n = m = 5$ ) — from manual introduction, validating properties:

demo-1:  $x = (1, 2, 4, 8, 16)$ ,  $y = (2, 4, 8, 16, 32)$ , expected output: 0.5 (base case)

demo-2:  $x = (1, 2, 4, 8, 16)$ ,  $y = (1, 2, 4, 8, 16)$ , expected output: 1 (identity property)

demo-3:  $x = (2, 4, 8, 16, 32)$ ,  $y = (10, 20, 40, 80, 160) (= [2 \times \text{demo-1.x}, 5 \times \text{demo-1.y}])$ , expected output: 0.2 (scale property)

**Natural sequences** ( $[n, m] \in \{1, 2, 3\} \times \{1, 2, 3\}$ ) — 9 combinations:

natural-1-1:  $x = (1)$ ,  $y = (1)$ , expected output: 1

natural-1-2:  $x = (1)$ ,  $y = (1, 2)$ , expected output:  $\approx 0.667$

natural-1-3:  $x = (1)$ ,  $y = (1, 2, 3)$ , expected output: 0.5

natural-2-1:  $x = (1, 2)$ ,  $y = (1)$ , expected output: 1.5

natural-2-2:  $x = (1, 2)$ ,  $y = (1, 2)$ , expected output: 1

natural-2-3:  $x = (1, 2)$ ,  $y = (1, 2, 3)$ , expected output:  $\approx 0.833$

natural-3-1:  $x = (1, 2, 3)$ ,  $y = (1)$ , expected output: 2

natural-3-2:  $x = (1, 2, 3)$ ,  $y = (1, 2)$ , expected output: 1.5

natural-3-3:  $x = (1, 2, 3)$ ,  $y = (1, 2, 3)$ , expected output: 1

**Additive distribution** ( $[n, m] \in \{5, 10, 30\} \times \{5, 10, 30\}$ ) — 9 combinations with Additive(10, 1):

additive-5-5, additive-5-10, additive-5-30

additive-10-5, additive-10-10, additive-10-30

additive-30-5, additive-30-10, additive-30-30

Random generation: x uses seed 0, y uses seed 1

**Uniform distribution** ( $[n, m] \in \{5, 100\} \times \{5, 100\}$ ) — 4 combinations with Uniform(0, 1):

uniform-5-5, uniform-5-100, uniform-100-5, uniform-100-100

Random generation: x uses seed 2, y uses seed 3

The natural sequences verify the identity property ( $\text{Ratio}(x, x) = 1$ ) and validate ratio calculations with simple integer inputs. Note that implementations should handle the practical constraint of avoiding division by values near zero.

**Unsorted tests** — verify independent sorting for ratio calculation (12 tests):

unsorted-x-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : X unsorted (reversed), Y sorted (2 tests)

unsorted-y-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : X sorted, Y unsorted (reversed) (2 tests)

unsorted-both-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : both unsorted (reversed) (2 tests)

unsorted-demo-unsorted-x:  $x = (16, 1, 8, 2, 4)$ ,  $y = (2, 4, 8, 16, 32)$  (demo-1 with X unsorted)

unsorted-demo-unsorted-y:  $x = (1, 2, 4, 8, 16)$ ,  $y = (32, 2, 16, 4, 8)$  (demo-1 with Y unsorted)

unsorted-demo-both-unsorted:  $x = (8, 1, 16, 4, 2)$ ,  $y = (16, 32, 2, 8, 4)$  (demo-1 both unsorted)

unsorted-identity-unsorted:  $x = (4, 1, 8, 2, 16)$ ,  $y = (16, 1, 8, 4, 2)$  (identity property, both unsorted)

unsorted-asymmetric-unsorted-2-3:  $x = (2, 1)$ ,  $y = (3, 1, 2)$  (asymmetric, both unsorted)

unsorted-power-unsorted-5:  $x = (16, 2, 8, 1, 4)$ ,  $y = (32, 4, 16, 2, 8)$  (powers of 2 unsorted)

## 6.6. AvgSpread Tests

$$\text{AvgSpread}(\mathbf{x}, \mathbf{y}) = \frac{n \cdot \text{Spread}(\mathbf{x}) + m \cdot \text{Spread}(\mathbf{y})}{n + m}$$

The AvgSpread test suite contains 49 test cases (35 original + 14 unsorted). Since AvgSpread computes  $\text{Spread}(\mathbf{x})$  and  $\text{Spread}(\mathbf{y})$  independently, unsorted tests are critical to verify that both samples are sorted independently before computing their spreads.

**Demo examples** ( $n = m = 5$ ) — from manual introduction, validating properties:

demo-1:  $\mathbf{x} = (0, 3, 6, 9, 12)$ ,  $\mathbf{y} = (0, 2, 4, 6, 8)$ , expected output: 5 (base case:  $(5 \cdot 6 + 5 \cdot 4)/10$ )

demo-2:  $\mathbf{x} = (0, 3, 6, 9, 12)$ ,  $\mathbf{y} = (0, 3, 6, 9, 12)$ , expected output: 6 (identity case)

demo-3:  $\mathbf{x} = (0, 6, 12, 18, 24)$ ,  $\mathbf{y} = (0, 9, 18, 27, 36)$  ( $= [2 \times \text{demo-1.x}, 3 \times \text{demo-1.y}]$ ), expected output: 15 (scale equivariance)

demo-4:  $\mathbf{x} = (0, 2, 4, 6, 8)$ ,  $\mathbf{y} = (0, 3, 6, 9, 12)$  ( $=$  reversed demo-1), expected output: 5 (symmetry)

demo-5:  $\mathbf{x} = (0, 6, 12, 18, 24)$ ,  $\mathbf{y} = (0, 4, 8, 12, 16)$  ( $= 2 \times \text{demo-1}$ ), expected output: 10 (uniform scaling)

**Natural sequences** ( $[n, m] \in \{1, 2, 3\} \times \{1, 2, 3\}$ ) — 9 combinations:

All combinations from single-element to three-element samples, validating the weighted average calculation

**Negative values** ( $[n, m] = [2, 2]$ ) — validates spread calculation with negative values:

negative-2-2:  $\mathbf{x} = (-2, -1)$ ,  $\mathbf{y} = (-2, -1)$ , expected output: 1

**Zero values** ( $[n, m] \in \{1, 2\} \times \{1, 2\}$ ) — 4 combinations:

All produce output 0 since Spread of constant samples is zero

**Additive distribution** ( $[n, m] \in \{5, 10, 30\} \times \{5, 10, 30\}$ ) — 9 combinations with Additive(10, 1):

Tests pooled dispersion across different sample size combinations

Random generation:  $\mathbf{x}$  uses seed 0,  $\mathbf{y}$  uses seed 1

**Uniform distribution** ( $[n, m] \in \{5, 100\} \times \{5, 100\}$ ) — 4 combinations with Uniform(0, 1):

Validates correct weighting when sample sizes differ substantially

Random generation:  $\mathbf{x}$  uses seed 2,  $\mathbf{y}$  uses seed 3

The asymmetric size combinations are particularly important for AvgSpread because the estimator must correctly weight each sample's contribution by its size.

**Composite estimator stress tests** — edge cases for weighted averaging:

composite-asymmetric-weights:  $\mathbf{x} = (1, 2)$ ,  $\mathbf{y} = (3, 4, 5, 6, 7, 8, 9, 10)$  (2 vs 8, tests weighting formula)

composite-zero-spread-one:  $\mathbf{x} = (5, 5, 5)$ ,  $\mathbf{y} = (1, 2, 3, 4, 5)$  (one zero spread, tests edge case)

composite-extreme-sizes:  $\mathbf{x} = (10)$ ,  $\mathbf{y} = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$  (1 vs 10, extreme weighting)

**Unsorted tests** — critical for verifying independent sorting (14 tests):

unsorted-x-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : X unsorted (reversed), Y sorted (2 tests)  
 unsorted-y-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : X sorted, Y unsorted (reversed) (2 tests)  
 unsorted-both-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : both unsorted (reversed) (2 tests)  
 unsorted-demo-unsorted-x:  $x = (12, 0, 6, 3, 9)$ ,  $y = (0, 2, 4, 6, 8)$  (demo-1 with X unsorted)  
 unsorted-demo-unsorted-y:  $x = (0, 3, 6, 9, 12)$ ,  $y = (8, 0, 4, 2, 6)$  (demo-1 with Y unsorted)  
 unsorted-demo-both-unsorted:  $x = (9, 0, 12, 3, 6)$ ,  $y = (6, 0, 8, 2, 4)$  (demo-1 both unsorted)  
 unsorted-identity-unsorted:  $x = (6, 0, 12, 3, 9)$ ,  $y = (9, 0, 12, 6, 3)$  (demo-2 unsorted)  
 unsorted-negative-unsorted:  $x = (-1, -2)$ ,  $y = (-1, -2)$  (negative unsorted)  
 unsorted-zero-unsorted-2-2:  $x = (0, 0)$ ,  $y = (0, 0)$  (zeros, any order)  
 unsorted-asymmetric-weights-unsorted:  $x = (2, 1)$ ,  $y = (8, 3, 6, 4, 10, 5, 9, 7)$  (asymmetric unsorted)  
 unsorted-zero-spread-x-unsorted:  $x = (5, 5, 5)$ ,  $y = (5, 1, 4, 2, 3)$  (zero spread X, Y unsorted)

These tests verify that implementations compute  $\text{Spread}(x)$  and  $\text{Spread}(y)$  with properly sorted samples.

## 6.7. Disparity Tests

$$\text{Disparity}(x, y) = \frac{\text{Shift}(x, y)}{\text{AvgSpread}(x, y)}$$

The Disparity test suite contains 28 test cases (16 original + 12 unsorted). Since Disparity combines Shift and AvgSpread, unsorted tests verify both components handle sorting correctly.

**Demo examples** ( $n = m = 5$ ) — from manual introduction, validating properties:

demo-1:  $x = (0, 3, 6, 9, 12)$ ,  $y = (0, 2, 4, 6, 8)$ , expected output: 0.4 (base case:  $2/5$ )  
 demo-2:  $x = (5, 8, 11, 14, 17)$ ,  $y = (5, 7, 9, 11, 13)$  (= demo-1 + 5), expected output: 0.4 (location invariance)  
 demo-3:  $x = (0, 6, 12, 18, 24)$ ,  $y = (0, 4, 8, 12, 16)$  ( $= 2 \times \text{demo-1}$ ), expected output: 0.4 (scale invariance)  
 demo-4:  $x = (0, 2, 4, 6, 8)$ ,  $y = (0, 3, 6, 9, 12)$  (= reversed demo-1), expected output:  $-0.4$  (anti-symmetry)

**Natural sequences** ( $[n, m] \in \{2, 3\} \times \{2, 3\}$ ) — 4 combinations:

natural-2-2, natural-2-3, natural-3-2, natural-3-3

Minimum size  $n, m \geq 2$  required for meaningful dispersion calculations

**Negative values** ( $[n, m] = [2, 2]$ ) — end-to-end validation with negative values:

negative-2-2:  $x = (-2, -1)$ ,  $y = (-2, -1)$ , expected output: 0

**Uniform distribution** ( $[n, m] \in \{5, 100\} \times \{5, 100\}$ ) — 4 combinations with Uniform(0, 1):

uniform-5-5, uniform-5-100, uniform-100-5, uniform-100-100

Random generation: x uses seed 0, y uses seed 1

The smaller test set for Disparity reflects implementation confidence. Since Disparity combines Shift and AvgSpread, correct implementation of those components ensures Disparity correctness. The test cases validate the division operation and confirm scale-free properties.

**Composite estimator stress tests** — edge cases for effect size calculation:

composite-small-avgspread:  $x = (10.001, 10.002, 10.003)$ ,  $y = (10.004, 10.005, 10.006)$  (tiny spread, large shift)  
 composite-large-avgspread:  $x = (1, 100, 200)$ ,  $y = (50, 150, 250)$  (large spread, small shift)  
 composite-extreme-disparity:  $x = (1, 1.001)$ ,  $y = (100, 100.001)$  (extreme ratio, tests precision)

**Unsorted tests** — verify both Shift and AvgSpread handle sorting (12 tests):

unsorted-x-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : X unsorted (reversed), Y sorted (2 tests)  
 unsorted-y-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : X sorted, Y unsorted (reversed) (2 tests)  
 unsorted-both-natural- $\{n\}$ - $\{m\}$  for  $(n, m) \in \{(3, 3), (4, 4)\}$ : both unsorted (reversed) (2 tests)  
 unsorted-demo-unsorted-x:  $x = (12, 0, 6, 3, 9)$ ,  $y = (0, 2, 4, 6, 8)$  (demo-1 with X unsorted)  
 unsorted-demo-unsorted-y:  $x = (0, 3, 6, 9, 12)$ ,  $y = (8, 0, 4, 2, 6)$  (demo-1 with Y unsorted)  
 unsorted-demo-both-unsorted:  $x = (9, 0, 12, 3, 6)$ ,  $y = (6, 0, 8, 2, 4)$  (demo-1 both unsorted)  
 unsorted-location-invariance-unsorted:  $x = (17, 5, 11, 8, 14)$ ,  $y = (13, 5, 9, 7, 11)$  (demo-2 unsorted)  
 unsorted-scale-invariance-unsorted:  $x = (24, 0, 12, 6, 18)$ ,  $y = (16, 0, 8, 4, 12)$  (demo-3 unsorted)  
 unsorted-anti-symmetry-unsorted:  $x = (8, 0, 4, 2, 6)$ ,  $y = (12, 0, 6, 3, 9)$  (demo-4 reversed and unsorted)

As a composite estimator, Disparity tests both the numerator (Shift) and denominator (AvgSpread). Unsorted variants verify end-to-end correctness including invariance properties.

## 6.8. PairwiseMargin Tests

$\text{PairwiseMargin}(n, m, \text{misrate})$

The PairwiseMargin test suite contains 346 correctness test cases (4 demo + 32 natural + 10 edge + 300 comprehensive grid).

**Demo examples** ( $n = m = 30$ ) — from manual introduction:

demo-1:  $n = 30, m = 30, \text{misrate} = 10^{-6}$ , expected output: 276  
 demo-2:  $n = 30, m = 30, \text{misrate} = 10^{-5}$ , expected output: 328  
 demo-3:  $n = 30, m = 30, \text{misrate} = 10^{-4}$ , expected output: 390  
 demo-4:  $n = 30, m = 30, \text{misrate} = 10^{-3}$ , expected output: 464

These demo cases match the reference values used throughout the manual to illustrate ShiftBounds construction.

**Natural sequences** ( $[n, m] \in \{1, 2, 3, 4\} \times \{1, 2, 3, 4\} \times 2$  misrates) — 32 tests:

Misrate values:  $\text{misrate} \in \{10^{-1}, 10^{-2}\}$   
 Test naming: natural- $\{n\}$ - $\{m\}$ -mr $\{k\}$  where  $k$  is the negative log10 of misrate

**Examples:**

natural-1-1-mr1:  $n = 1, m = 1$ , misrate = 0.1, expected output: 0  
 natural-2-2-mr1:  $n = 2, m = 2$ , misrate = 0.1, expected output: 0  
 natural-3-3-mr2:  $n = 3, m = 3$ , misrate = 0.01, expected output: 0  
 natural-4-4-mr1:  $n = 4, m = 4$ , misrate = 0.1, expected output: 4

The natural sequences provide canonical examples with small, easily verified parameter values.

**Edge cases** — boundary condition validation:

boundary-min:  $n = 1, m = 1$ , misrate = 0.5 (minimum samples, expected output: 0)  
 boundary-zero-margin-small:  $n = 2, m = 2$ , misrate =  $10^{-6}$  (misrate too strict, expected output: 0)  
 boundary-loose:  $n = 5, m = 5$ , misrate = 0.9 (very permissive misrate)  
 symmetry-2-5:  $n = 2, m = 5$ , misrate = 0.1 (tests symmetry property)  
 symmetry-5-2:  $n = 5, m = 2$ , misrate = 0.1 (symmetric counterpart, same output as above)  
 symmetry-3-7:  $n = 3, m = 7$ , misrate = 0.05 (asymmetric sizes)  
 symmetry-7-3:  $n = 7, m = 3$ , misrate = 0.05 (symmetric counterpart)  
 asymmetry-extreme-1-100:  $n = 1, m = 100$ , misrate = 0.1 (extreme size difference)  
 asymmetry-extreme-100-1:  $n = 100, m = 1$ , misrate = 0.1 (reversed extreme)  
 asymmetry-extreme-2-50:  $n = 2, m = 50$ , misrate = 0.05 (highly unbalanced)

These edge cases validate correct handling of boundary conditions, the symmetry property  $\text{PairwiseMargin}(n, m, \text{misrate}) = \text{PairwiseMargin}(m, n, \text{misrate})$ , and extreme asymmetry in sample sizes.

**Comprehensive grid** — systematic coverage for thorough validation:

Small sample combinations ( $[n, m] \in \{1, 2, 3, 4, 5\} \times \{1, 2, 3, 4, 5\} \times 6$  misrates) — 150 tests:

Misrate values:  $\text{misrate} \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$   
 Test naming:  $n\{n\}_m\{m\}_{mr}\{k\}$  where  $k$  is the negative log<sub>10</sub> of misrate  
 Examples:

n1\_m1\_mr1:  $n = 1, m = 1$ , misrate = 0.1, expected output: 0  
 n5\_m5\_mr1:  $n = 5, m = 5$ , misrate = 0.1, expected output: 10  
 n5\_m5\_mr3:  $n = 5, m = 5$ , misrate = 0.001, expected output: 0

Large sample combinations ( $[n, m] \in \{10, 20, 30, 50, 100\} \times \{10, 20, 30, 50, 100\} \times 6$  misrates) — 150 tests:

Misrate values: same as small samples  
 Test naming:  $n\{n\}_m\{m\}_r\{k\}$  where  $k$  is the negative log<sub>10</sub> of misrate  
 Examples:

n10\_m10\_r1:  $n = 10, m = 10$ , misrate = 0.1, expected output: 56  
 n10\_m10\_r6:  $n = 10, m = 10$ , misrate =  $10^{-6}$ , expected output: 0  
 n50\_m50\_r3:  $n = 50, m = 50$ , misrate = 0.001, expected output: 1556  
 n100\_m100\_r6:  $n = 100, m = 100$ , misrate =  $10^{-6}$ , expected output: 6060

The comprehensive grid validates both symmetric ( $n = m$ ) and asymmetric sample size combinations across six orders of magnitude in misrate, ensuring robust coverage of the parameter space.

## 6.9. ShiftBounds Tests

$$\text{ShiftBounds}(\mathbf{x}, \mathbf{y}, \text{misrate}) = [z_{(k_{\text{left}})}, z_{(k_{\text{right}})}]$$

where

$$\mathbf{z} = \{x_i - y_j\}_{1 \leq i \leq n, 1 \leq j \leq m} \quad (\text{sorted})$$

$$k_{\text{left}} = \left\lfloor \frac{\text{PairwiseMargin}(n, m, \text{misrate})}{2} \right\rfloor + 1$$

$$k_{\text{right}} = nm - \left\lfloor \frac{\text{PairwiseMargin}(n, m, \text{misrate})}{2} \right\rfloor$$

The ShiftBounds test suite contains 61 correctness test cases (3 demo + 9 natural + 6 property + 10 edge + 9 additive + 4 uniform + 5 misrate + 15 unsorted). Since ShiftBounds returns bounds rather than a point estimate, tests validate that the bounds contain Shift( $\mathbf{x}, \mathbf{y}$ ) and satisfy equivariance properties. Each test case output is a JSON object with lower and upper fields representing the interval bounds.

**Demo examples** ( $n = m = 5$ ) — from manual introduction, validating basic bounds:

demo-1:  $\mathbf{x} = (1, 2, 3, 4, 5), \mathbf{y} = (3, 4, 5, 6, 7), \text{misrate} = 0.05$ , expected output:  $[-4, 0]$

demo-2:  $\mathbf{x} = (1, 2, 3, 4, 5), \mathbf{y} = (3, 4, 5, 6, 7), \text{misrate} = 0.01$ , expected output:  $[-5, 1]$

demo-3:  $\mathbf{x} = (3, 4, 5, 6, 7), \mathbf{y} = (3, 4, 5, 6, 7), \text{misrate} = 0.05$ , expected output: bounds containing 0 (identity case)

These cases illustrate how tighter misrates produce wider bounds and validate the identity property where identical samples yield bounds containing zero.

**Natural sequences** ( $[n, m] \in \{1, 2, 3\} \times \{1, 2, 3\}, \text{misrate} = 10^{-2}$ ) — 9 combinations:

natural-1-1:  $\mathbf{x} = (1), \mathbf{y} = (1)$ , expected bounds containing 0

natural-1-2:  $\mathbf{x} = (1), \mathbf{y} = (1, 2)$ , expected bounds containing  $-0.5$

natural-1-3:  $\mathbf{x} = (1), \mathbf{y} = (1, 2, 3)$ , expected bounds containing  $-1$

natural-2-1:  $\mathbf{x} = (1, 2), \mathbf{y} = (1)$ , expected bounds containing  $0.5$

natural-2-2:  $\mathbf{x} = (1, 2), \mathbf{y} = (1, 2)$ , expected bounds containing 0

natural-2-3:  $\mathbf{x} = (1, 2), \mathbf{y} = (1, 2, 3)$ , expected bounds containing  $-0.5$

natural-3-1:  $\mathbf{x} = (1, 2, 3), \mathbf{y} = (1)$ , expected bounds containing 1

natural-3-2:  $\mathbf{x} = (1, 2, 3), \mathbf{y} = (1, 2)$ , expected bounds containing  $0.5$

natural-3-3:  $\mathbf{x} = (1, 2, 3), \mathbf{y} = (1, 2, 3)$ , expected bounds containing 0

These canonical cases validate that bounds properly contain the corresponding Shift values and handle small sample sizes correctly.

**Property validation** ( $n = m = 5, \text{misrate} = 10^{-3}$ ) — 6 tests:



property-identity:  $x = (0, 2, 4, 6, 8)$ ,  $y = (0, 2, 4, 6, 8)$ , bounds must contain 0  
 property-location-shift:  $x = (7, 9, 11, 13, 15)$ ,  $y = (13, 15, 17, 19, 21)$  (= demo-1 + [7, 3])  
 Must produce same bounds as base case (location invariance)  
 property-scale-2x:  $x = (2, 4, 6, 8, 10)$ ,  $y = (6, 8, 10, 12, 14)$  (= 2 × demo-1)  
 Bounds must be 2× the base case bounds (scale equivariance)  
 property-antisymmetry:  $x = (3, 4, 5, 6, 7)$ ,  $y = (1, 2, 3, 4, 5)$  (= reversed demo-1)  
 Bounds must be negated: if original is  $[a, b]$ , this yields  $[-b, -a]$   
 property-negative:  $x = (-5, -4, -3, -2, -1)$ ,  $y = (-7, -6, -5, -4, -3)$   
 Validates sign handling with all negative values  
 property-mixed-signs:  $x = (-2, -1, 0, 1, 2)$ ,  $y = (-1, 0, 1, 2, 3)$   
 Validates bounds crossing zero with mixed-sign samples

### Edge cases — boundary conditions and extreme scenarios (10 tests):

edge-min-samples:  $x = (1)$ ,  $y = (2)$ , misrate =  $10^{-2}$  (minimum samples, single difference)  
 edge-permissive-misrate:  $x = (1, 2, 3, 4, 5)$ ,  $y = (3, 4, 5, 6, 7)$ , misrate = 0.5 (very wide bounds)  
 edge-strict-misrate:  $x = (1, 2, 3, 4, 5)$ ,  $y = (3, 4, 5, 6, 7)$ , misrate =  $10^{-6}$  (very narrow bounds)  
 edge-zero-shift:  $x = (5, 5, 5)$ ,  $y = (5, 5, 5)$ , misrate =  $10^{-3}$  (all identical, bounds around 0)  
 edge-asymmetric-1-100:  $x = (50)$ ,  $y = (1, 2, \dots, 100)$ , misrate =  $10^{-2}$  (extreme size difference)  
 edge-asymmetric-2-50:  $x = (25, 26)$ ,  $y = (1, 2, \dots, 50)$ , misrate =  $10^{-3}$  (highly unbalanced)  
 edge-duplicates:  $x = (3, 3, 3, 3, 3)$ ,  $y = (5, 5, 5, 5, 5)$ , misrate =  $10^{-2}$  (all duplicates, bounds around -2)  
 edge-wide-range:  $x = (0.001, 1, 100, 1000, 10000)$ ,  $y = (0.1, 10, 1000, 100000)$ , misrate =  $10^{-3}$  (extreme value range)  
 edge-tiny-values:  $x = (10^{-8}, 2 \cdot 10^{-8}, 3 \cdot 10^{-8})$ ,  $y = (2 \cdot 10^{-8}, 3 \cdot 10^{-8}, 4 \cdot 10^{-8})$ , misrate =  $10^{-2}$  (numerical precision)  
 edge-large-values:  $x = (10^8, 2 \cdot 10^8, 3 \cdot 10^8)$ ,  $y = (2 \cdot 10^8, 3 \cdot 10^8, 4 \cdot 10^8)$ , misrate =  $10^{-2}$  (large magnitude)

These edge cases stress-test boundary conditions, numerical stability, and the margin calculation with extreme parameters.

**Additive distribution** ( $[n, m] \in \{5, 10, 30\} \times \{5, 10, 30\}$ , misrate =  $10^{-3}$ ) — 9 combinations with Additive(10, 1):

additive-5-5, additive-5-10, additive-5-30  
 additive-10-5, additive-10-10, additive-10-30  
 additive-30-5, additive-30-10, additive-30-30  
 Random generation:  $x$  uses seed 0,  $y$  uses seed 1

These fuzzy tests validate that bounds properly encompass the shift estimate for realistic normally-distributed data at various sample sizes.

**Uniform distribution** ( $[n, m] \in \{5, 100\} \times \{5, 100\}$ , misrate =  $10^{-4}$ ) — 4 combinations with Uniform(0, 1):

uniform-5-5, uniform-5-100, uniform-100-5, uniform-100-100

Random generation: x uses seed 2, y uses seed 3

The asymmetric size combinations are particularly important for testing margin calculation with unbalanced samples.

**Misrate variation** ( $x = (0, 2, 4, 6, 8)$ ,  $y = (10, 12, 14, 16, 18)$ ) — 5 tests with varying misrates:

```
misrate-1e-2: misrate =  $10^{-2}$ 
misrate-1e-3: misrate =  $10^{-3}$ 
misrate-1e-4: misrate =  $10^{-4}$ 
misrate-1e-5: misrate =  $10^{-5}$ 
misrate-1e-6: misrate =  $10^{-6}$ 
```

These tests use identical samples with varying misrates to validate the monotonicity property: smaller misrates (higher confidence) produce wider bounds. The sequence demonstrates how bound width increases as misrate decreases, helping implementations verify correct margin calculation.

**Unsorted tests** — verify independent sorting of x and y (15 tests):

```
unsorted-x-natural-3-3: x = (3, 2, 1), y = (1, 2, 3), misrate =  $10^{-2}$  (X reversed, Y sorted)
unsorted-y-natural-3-3: x = (1, 2, 3), y = (3, 2, 1), misrate =  $10^{-2}$  (X sorted, Y reversed)
unsorted-both-natural-3-3: x = (3, 2, 1), y = (3, 2, 1), misrate =  $10^{-2}$  (both reversed)
unsorted-x-shuffle-4-4: x = (3, 1, 4, 2), y = (1, 2, 3, 4), misrate =  $10^{-3}$  (X shuffled)
unsorted-y-shuffle-4-4: x = (1, 2, 3, 4), y = (4, 2, 1, 3), misrate =  $10^{-3}$  (Y shuffled)
unsorted-both-shuffle-4-4: x = (3, 1, 4, 2), y = (2, 4, 1, 3), misrate =  $10^{-3}$  (both shuffled)
unsorted-demo-unsorted-x: x = (5, 1, 4, 2, 3), y = (3, 4, 5, 6, 7), misrate = 0.05 (demo-I X unsorted)
unsorted-demo-unsorted-y: x = (1, 2, 3, 4, 5), y = (7, 3, 6, 4, 5), misrate = 0.05 (demo-I Y unsorted)
unsorted-demo-both-unsorted: x = (4, 1, 5, 2, 3), y = (6, 3, 7, 4, 5), misrate = 0.05 (demo-I both unsorted)
unsorted-identity-unsorted: x = (4, 1, 5, 2, 3), y = (5, 1, 4, 3, 2), misrate =  $10^{-2}$  (identity property, both unsorted)
unsorted-negative-unsorted: x = (-1, -3, -2), y = (-2, -3, -1), misrate =  $10^{-2}$  (negative values unsorted)
unsorted-asymmetric-2-5: x = (2, 1), y = (5, 2, 4, 1, 3), misrate =  $10^{-3}$  (asymmetric sizes, both unsorted)
unsorted-duplicates: x = (3, 3, 3, 3, 3), y = (5, 5, 5, 5, 5), misrate =  $10^{-2}$  (all duplicates, any order)
unsorted-mixed-duplicates-x: x = (2, 1, 3, 2, 1), y = (1, 1, 2, 2, 3), misrate =  $10^{-3}$  (X has unsorted duplicates)
unsorted-mixed-duplicates-y: x = (1, 1, 2, 2, 3), y = (3, 2, 1, 3, 2), misrate =  $10^{-3}$  (Y has unsorted duplicates)
```

These unsorted tests are critical because ShiftBounds computes bounds from pairwise differences, requiring both samples to be sorted independently. The variety ensures implementations don't incorrectly assume pre-sorted input or sort samples together. Each test must produce identical output to its sorted counterpart, validating that the implementation correctly handles the sorting step.

**No performance test** — ShiftBounds uses the FastShift algorithm internally, which is already validated by the Shift performance test. Since bounds computation involves only two quantile calculations from the pairwise differences (at positions determined by PairwiseMargin), the performance characteristics are equivalent to computing two Shift estimates, which completes efficiently for large samples.

## 6.10. Test Framework

The reference test framework consists of three components:

**Test generation** — The C# implementation defines test inputs programmatically using builder patterns. For deterministic cases, inputs are explicitly specified. For random cases, the framework uses controlled seeds with `System.Random` to ensure reproducibility across all platforms.

The random generation mechanism works as follows:

Each test suite builder maintains a seed counter initialized to zero.

For one-sample estimators, each distribution type receives the next available seed. The same random generator produces all samples for all sizes within that distribution.

For two-sample estimators, each pair of distributions receives two consecutive seeds: one for the x sample generator and one for the y sample generator.

The seed counter increments with each random generator creation, ensuring deterministic test data generation.

For Additive distributions, random values are generated using the Box-Müller transform, which converts pairs of uniform random values into normally distributed values. The transform applies the formula:

$$X = \mu + \sigma \sqrt{-2 \ln(U_1)} \sin(2\pi U_2)$$

where  $U_1, U_2$  are uniform random values from Uniform(0, 1),  $\mu$  is the mean, and  $\sigma$  is the standard deviation.

For Uniform distributions, random values are generated directly using the quantile function:

$$X = \min + U \cdot (\max - \min)$$

where  $U$  is a uniform random value from Uniform(0, 1).

The framework executes the reference implementation on all generated inputs and serializes input-output pairs to JSON format.

**Test validation** — Each language implementation loads the JSON test cases and executes them against its local estimator implementation. Assertions verify that outputs match expected values within a given numerical tolerance (typically  $10^{-10}$  for relative error).

**Test data format** — Each test case is a JSON file containing input and output fields. For one-sample estimators, the input contains array x and optional parameters. For two-sample estimators, input contains arrays x and y. Output is a single numeric value.

**Performance testing** — The toolkit provides  $O(n \log n)$  fast algorithms for Center, Spread, and Shift estimators, dramatically more efficient than naive implementations that materialize all pairwise combinations. Performance tests use sample size  $n = 100,000$  (for one-sample) or  $n = m = 100,000$  (for two-sample). This specific size creates a clear performance distinction: fast implementations ( $O(n \log n)$  or  $O((m + n) \log L)$ ) complete in under 5 seconds on modern hardware across all supported languages, while naive implementations ( $O(n^2 \log n)$  or  $O(mn \log(mn))$ ) would be prohibitively slow (taking hours or failing due to memory exhaustion). With  $n = 100,000$ , naive approaches would need to materialize approximately 5 billion pairwise values for Center/Spread or 10 billion for Shift, whereas fast algorithms require only  $O(n)$  additional memory. Performance tests serve dual purposes: correctness validation at scale and performance regression detection, ensuring implementations use the efficient algorithms and remain practical for real-world datasets with hundreds of thousands of observations. Performance test specifications are provided in the respective estimator sections above.

This framework ensures that all seven language implementations maintain strict numerical agreement across the full test suite.

## 7. Methodology

This chapter examines the methodological principles that guide the pragmatist's design and application.

### 7.1. Pragmatic Philosophy

The toolkit's foundations rest on pragmatist epistemology: truth is determined by practical consequences, not abstract correspondence with reality.

**Truth is what works** — An estimator is “correct” if it produces useful results across realistic conditions

**Meaning from consequences** — The value of a statistical method lies in what it enables, not its theoretical elegance

**Theory serves practice** — Mathematical analysis provides insight, but empirical validation determines adoption

**Utility as criterion** — When methods conflict, prefer the one that solves more real problems

This stance inverts the traditional relationship between theory and practice. Rather than deriving methods from first principles and hoping they apply, we evaluate methods by their performance and seek theoretical understanding afterward.

## 7.2. Procedure-First Empiricism

Traditional statistical practice follows an assumptions-first methodology:

1. Assume a data-generating model (e.g., “observations are normally distributed”)
2. Derive the optimal procedure under those assumptions
3. Apply the procedure to data, hoping assumptions approximately hold

This toolkit inverts the process:

1. Select procedures based on desired properties (robustness, equivariance, interpretability)
2. Empirically measure performance across a wide range of conditions
3. Use theory to explain and predict observed behavior

Monte Carlo simulation serves as the primary instrument of knowledge. Rather than deriving asymptotic formulas for estimator variance, we measure actual variance across thousands of simulated samples. Drift tables in this manual are empirically measured, not analytically derived.

This approach has practical advantages: simulations can explore conditions that resist closed-form analysis, and empirical results are self-validating—they show what actually happens, not what theory predicts should happen.

### 7.3. Epistemic Humility

No perfectly Gaussian, log-normal, or Pareto distributions exist in real data. Every distribution we name is a useful fiction — a model we employ because it approximates reality well enough for our purposes, while knowing it cannot be exactly correct.

**Models are approximations** — They capture essential structure while ignoring irrelevant details

**Approximations fail at boundaries** — Edge cases, extreme values, and distribution tails often violate assumptions

**Graceful degradation** — Methods should produce sensible (if less precise) results when assumptions weaken

The toolkit embodies this humility by choosing estimators that remain interpretable and bounded even when distributional assumptions break down. A robust estimator may sacrifice some efficiency under ideal conditions in exchange for reliable behavior when conditions degrade.

## 7.4. The Pairwise Principle

A structural insight unifies all primary robust estimators in this toolkit: they are medians of pairwise operations.

Estimator	Pairwise Operation	Result
Center	$(x_i + x_j)/2$	Median of pairwise averages
Spread	$ x_i - x_j $	Median of pairwise differences
Shift	$x_i - y_j$	Median of cross-sample differences
Dominance	$1(x_i > y_j)$	Proportion of pairwise comparisons

This pairwise structure provides three benefits:

**Natural robustness** — Comparing measurements to each other, not to external references, limits outlier influence

**Self-calibration** — The sample serves as its own reference distribution, requiring no external assumptions

**Algebraic closure** — Pairwise operations preserve symmetry and equivariance properties

The pairwise principle also enables efficient computation. Matrices of pairwise operations have structural properties (sorted rows and columns) that fast algorithms exploit to achieve  $O(n \log n)$  complexity.



## 7.5. Median as Universal Aggregator

The median is the final step in each pairwise estimator. Why median specifically?

The median achieves the maximum possible breakdown point (50%) among all translation-equivariant location estimators. Up to half the data can be arbitrarily corrupted before the median becomes unbounded.

However, Center and Spread achieve only 29% breakdown — not 50%. This is deliberate: a tradeoff between robustness and precision.

Breakdown	Robustness	Precision	Estimators
0%	None	Optimal under assumptions	Mean, StdDev
29%	Substantial	Near-optimal	Center, Spread
50%	Maximum	Reduced	Median, MAD

The 29% breakdown point survives approximately one corrupted measurement in four while maintaining roughly 95% asymptotic efficiency under ideal Gaussian conditions. This represents the practical optimum: enough robustness for realistic contamination levels, enough efficiency to compete with traditional methods when data is clean.

## 7.6. Convergence Conventions

Drift normalizes estimator variability by  $\sqrt{n}$ , making precision comparable across sample sizes:

$$\text{Drift} = \text{Spread}(\text{estimates}) \times \sqrt{n}$$

This normalization embeds a deliberate assumption: most useful estimators converge at the  $\sqrt{n}$  rate. The Central Limit Theorem guarantees this rate for means under mild conditions, and median-based estimators inherit similar convergence behavior.

**Common case default** —  $\sqrt{n}$  convergence covers the vast majority of practical estimators

**Intuitive interpretation** — Drift represents “effective standard deviation at  $n = 1$ ”

**Mental calculation** — Expected precision at any  $n$  is simply  $\text{Drift} / \sqrt{n}$

For estimators with non-standard convergence (e.g., extreme value statistics), drift generalizes to  $n^{\text{instability}}$  where instability differs from 0.5. But the toolkit deliberately uses  $\sqrt{n}$  throughout because it matches the common case and provides intuitive interpretation without complicating the universal mechanism.

This is pragmatic universalism: adopt the common case as default, acknowledge exceptions exist, and handle them explicitly rather than burdening the common case with unnecessary generality.

## 7.7. Structural Unity

All robust estimators in this toolkit share a common mathematical structure:

$$\text{Estimator} = \text{Median}(\text{Pairwise Operations})$$

This structural unity is not merely aesthetic — it enables unified algorithmic optimization.

**Sorted structure** — Matrices of pairwise operations have sorted rows and columns

**Monahan's algorithm** — Exploits sorted structure for  $O(n \log n)$  Center/Spread

**Fast shift** — Exploits cross-sample matrix structure for efficient two-sample comparison

Because all estimators share the same “median of pairwise” form, insights that accelerate one can often be adapted to accelerate others. A single theoretical framework covers all primary estimators.

## 7.8. Generative Naming

Names in this toolkit encode operational knowledge rather than historical provenance.

Traditional	Pragmastat	What's Encoded
Gaussian / Normal	<u>Additive</u>	Formation: sum of independent factors (CLT)
Log-normal / Galton	<u>Multiplic</u>	Formation: product of independent factors
Pareto	<u>Power</u>	Behavior: power-law relationship
Hodges-Lehmann	Center	Function: measures central tendency
Shamos	Spread	Function: measures variability

When you read “Additive”, your mind activates a generative model: this distribution arises when many independent factors add together. When you read “Gaussian”, you must recall an association with Carl Friedrich Gauss, then remember what properties that name implies.

Generative names create immediate intuition about when a model applies. Additive distributions arise from additive processes. Multiplic distributions arise from multiplicative processes. The name itself encodes the formation mechanism.

## 7.9. The Inversion Principle

Traditional statistical outputs often require mental transformation before use. This toolkit inverts such framings to present information in directly actionable form, following principles of user-centered design (Norman (2013)).

Traditional	Pragmastat	Reason for Inversion
Confidence level (95%)	misrate (0.05)	Direct error interpretation
Confidence interval	Bounds	Plain language, no jargon
Hypothesis test (p-value)	Bounds estimation	“What’s plausible?” not “Is zero plausible?”
Efficiency (variance ratio)	Drift (spread-based)	Works with heavy tails

Consider the confidence level vs. misrate inversion. A “95% confidence interval” requires understanding: “If I repeated this procedure infinitely, 95% of intervals would contain the true value.” A “5% misrate” states directly: “This procedure errs about 5% of the time.”

The shift from confidence intervals to bounds, and from hypothesis testing to interval estimation, moves from frequentist theology toward decision-relevant inference. The practitioner asks “What values are plausible for this parameter?” rather than “Can I reject the hypothesis that this parameter equals zero?”

## 7.10. Multi-Audience Design

This manual serves readers with diverse backgrounds and conflicting preferences:

Audience	Priorities	Challenges
Experienced academics	Rigor, derivation, formalism, citations	May find practical focus too shallow
Professional developers	Examples, APIs, searchability, minimalism	May find theory intimidating
Students and beginners	Clarity, intuition, progressive disclosure	Need both theory and practice
Large language models	Structure, consistency, unambiguous definitions	Need form-independent content

These audiences have conflicting needs. Academics want complete derivations; developers want quick answers. Beginners need gentle introductions; experts need dense references. LLMs need predictable structure; humans appreciate variety.

The manual targets a “neutral zone” where all audiences find acceptable content:

**Signature first** — Mathematical definition immediately visible

**Example second** — Concrete computation before abstract explanation

**Detail optional** — Properties, corner cases, and theory follow for those who need them

**Every sentence earns its place** — No filler prose, no redundant explanation

### Structural Principles

**Concrete over abstract** — Numbers and examples before symbols and theory

**Precision without verbosity** — Mathematical rigor in minimal words

**Consistent layout** — Same structure across all toolkit items enables scanning

**Self-contained sections** — Each section readable independently

### LLM-Friendliness

The manual’s structure also serves machine readers:

**Predictable patterns** — Consistent section ordering aids extraction

**Explicit definitions** — No implicit knowledge assumed

**Tabular data** — Structured information in tables, not prose

**Short paragraphs** — Content chunks cleanly for context windows

This multi-audience optimization forces elimination of audience-specific conventions, revealing form-independent essential content that serves everyone adequately rather than serving one group perfectly and others poorly.

### 7.11. Reference Tests as Specification

The toolkit maintains seven implementations across different programming languages: Python, TypeScript, R, C#, Kotlin, Rust, and Go. Each implementation must produce identical numerical results for all estimators.

This cross-language consistency is achieved through executable specifications:

Manual (definitions) ↔ C# (reference) → JSON (tests) → All languages (validation)

The specification IS the test suite. Reference tests serve three critical purposes:

**Cross-language validation** — All implementations pass identical test cases

**Regression prevention** — Changes validated against known outputs

**Implementation guidance** — Concrete examples for porting to new languages

#### Test Design Principles

**Minimal sufficiency** — Smallest test set providing high confidence in correctness

**Comprehensive coverage** — Both typical cases and edge cases that expose errors

**Deterministic reproducibility** — Fixed seeds for all random tests

#### Test Categories

**Canonical cases** — Deterministic inputs like natural number sequences where outputs are easily verified

**Edge cases** — Boundary conditions: single element, zeros, minimum viable sample sizes

**Fuzzy tests** — Controlled random exploration beyond hand-crafted examples

The C# implementation serves as the reference generator. All test cases are defined programmatically, executed to produce expected outputs, and serialized to JSON. Other implementations load these JSON files and verify their outputs match within numerical tolerance.

## 7.12. Cross-Language Determinism

Reproducibility requires determinism at every layer. When a simulation in Python produces a result, the same simulation in Rust, Go, or any other supported language must produce the identical result.

**Portable RNG** — `Rng(experiment-1)` produces identical sequences in all languages

**Specified algorithms** — xoshiro256++ for generation, SplitMix64 for seeding, FNV-1a for string hashing

**No implementation-dependent behavior** — Floating-point operations follow IEEE 754

### Unified API

Beyond numerical determinism, the toolkit maintains a consistent API across all implementations. Function names, parameter orders, and return types follow the same conventions in every language.

**Same vocabulary** — Center, Spread, Shift mean the same thing everywhere

**Same signatures** — `Center(x)` in Python, `Center(x)` in Rust, `Center(x)` in Go

**Same behavior** — Edge cases, error conditions, and defaults are identical

This unified API enables frictionless language switching. A practitioner prototyping in Python can port to Rust for production without learning new abstractions or revalidating statistical assumptions. The mental model transfers directly; only syntax changes.

### Benefits of Unification

**Debugging across languages** — A failing test in TypeScript can be debugged in C#

**Verified ports** — New implementations can be validated against existing ones

**Reproducible research** — Results can be reproduced in any supported language

**Team flexibility** — Different team members can use preferred languages on the same analysis

**Migration paths** — Move from prototype to production without statistical revalidation



### 7.13. Summary Principles

The methodology of this toolkit can be distilled into twelve guiding principles:

1. **Name things by what they do, not who discovered them** — Generative names encode operational knowledge
2. **All models are wrong; design for graceful degradation** — Robust methods fail gently
3. **Evaluate empirically, organize theoretically** — Simulation before derivation
4. **Self-reference provides robustness** — Pairwise operations compare data to itself
5. **29% breakdown is the practical optimum** — Balance robustness and precision
6. **Invert framings that require mental transformation** — Present directly actionable information
7. **Default to the common case** — Use  $\sqrt{n}$  convergence; handle exceptions explicitly
8. **Multi-audience optimization reveals essential content** — Serve everyone adequately, not one group perfectly
9. **Executable specifications are reliable specifications** — Tests define correctness
10. **Reproducibility requires portable determinism** — Same seeds, same results, any language
11. **Structural unity enables unified optimization** — “Median of pairwise” admits fast algorithms
12. **Utility is the ultimate criterion** — Methods that solve real problems are correct methods

## Bibliography

- Blackman, D., & Vigna, S. (2021). Scrambled Linear Pseudorandom Number Generators. *ACM Transactions on Mathematical Software*, 47(4), 1–32. <https://dl.acm.org/doi/10.1145/3460772>
- Box, G. E. P., & Muller, M. E. (1958). A Note on the Generation of Random Normal Deviates. *The Annals of Mathematical Statistics*, 29(2), 610–611. <https://projecteuclid.org/euclid.aoms/1177706645>
- Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences* (2nd ed.). Lawrence Erlbaum Associates.
- Fan, C. T., Muller, M. E., & Rezucha, I. (1962). Development of Sampling Plans by Using Sequential (Item by Item) Selection Techniques and Digital Computers. *Journal of the American Statistical Association*, 57(298), 387–402. <https://www.tandfonline.com/doi/abs/10.1080/01621459.1962.10480667>
- Fisher, R. A., & Yates, F. (1938). *Statistical Tables for Biological, Agricultural and Medical Research*. Oliver and Boyd.
- Fix, E., & Hodges, J. L. (1955). Significance probabilities of the Wilcoxon test. *The Annals of Mathematical Statistics*, 26(2), 301–312. <https://projecteuclid.org/euclid.aoms/1177728547>
- Fowler, G., Noll, L. C., & Vo, K.-P. (1991, ). *FNV Hash*. <http://www.isthe.com/chongo/tech/comp/fnv/>
- Hodges, J. L., & Lehmann, E. L. (1963). Estimates of Location Based on Rank Tests. *The Annals of Mathematical Statistics*, 34(2), 598–611. <http://projecteuclid.org/euclid.aoms/1177704172>
- Huber, P. J. (2009). Robust statistics. In *International encyclopedia of statistical science* (pp. 1248–1251). Springer.
- Hyndman, R. J., & Fan, Y. (1996). Sample Quantiles in Statistical Packages. *The American Statistician*, 50(4), 361. <https://www.jstor.org/stable/2684934>
- Knuth, D. E. (1997). *The Art of Computer Programming, Volume 2: Seminumerical Algorithms* (3rd ed.). Addison-Wesley.
- Löffler, A. (1982). *On the calculation of the exact null-distribution of the Wilcoxon-Mann-Whitney U-statistic*.
- Monahan, J. F. (1984). Algorithm 616: fast computation of the Hodges-Lehmann location estimator. *ACM Transactions on Mathematical Software*, 10(3), 265–270. <https://dl.acm.org/doi/10.1145/1271.319414>
- Norman, D. A. (2013). *The Design of Everyday Things: Revised and Expanded Edition*. Basic Books.
- Sen, P. K. (1963). On the Estimation of Relative Potency in Dilution (–Direct) Assays by Distribution-Free Methods. *Biometrics*, 19(4), 532. <https://www.jstor.org/stable/2527532>
- Serfling, R. J. (2009). *Approximation theorems of mathematical statistics*. John Wiley & Sons.
- Shamos, M. I. (1976). *Geometry and Statistics: Problems at the Interface*.
- Sidak, Z., Sen, P. K., & Hajek, J. (1999). *Theory of rank tests*. Elsevier.

---

Steele, G. L., Lea, D., & Flood, C. H. (2014). *Fast Splittable Pseudorandom Number Generators* (pp. 453–472). ACM. <https://dl.acm.org/doi/10.1145/2660193.2660195>