

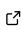


Tidy-TS: A Type-Safe Framework for Statistical Data Analysis in the TypeScript Ecosystem

John Thomas Menchaca MD ^{1,2}

¹ University of Utah, Department of Biomedical Informatics ² University of Utah, Department of Internal Medicine

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Tidy-TS is a type-safe library for statistical computing, data transformation, and visualization in the TypeScript ecosystem. It introduces a functional grammar for manipulating tabular data using arrays of objects, which are the idiomatic format in JavaScript runtimes. These are organized and operated on as DataFrames, following conventions established in statistical computing languages like R and Python. Inspired by the tidyverse philosophy ([Wickham et al., 2019](#)), it brings static typing, schema validation, and compile time guarantees to data workflows in TypeScript.

The library supports pipelines that load data from files, APIs, or databases, apply transformations that preserve type information, run statistical analyses, and produce interactive visualizations using Vega-Lite ([Satyanarayan et al., 2016, 2017](#)). These workflows remain entirely within TypeScript and avoid the need to switch languages for analysis or presentation.

Tidy-TS is designed for use in browsers, servers, and notebooks. It targets teams working fully in TypeScript who face fragmented toolchains and inconsistent data handling. Features such as asynchronous row-wise operations, concurrency controls, columnar storage, and WebAssembly acceleration support analytics workflows on millions of records in both research and production environments. The addition of a type system extends this paradigm to include substantive compile-time guarantees, which existing literature suggests can prevent 15–38% of production bugs ([Bunge, 2019](#); [Gao et al., 2017](#); [Khan et al., 2021](#)).

Statement of need

JavaScript and TypeScript are widely used for building interactive data applications, but they lack a type-aware and expressive toolkit for statistical computing. Existing libraries like Arquero and Danfo.js offer partial solutions for DataFrame-style operations, but they do not provide a comprehensive type system for a full data analysis pipeline.

Without a consistent type model, developers are often left writing brittle glue code across loosely coupled tools. Many teams using TypeScript for data ingestion and application logic still rely on Python or R for analysis and modeling. This workflow involves exporting and reimporting data with assumptions about column names, types, and missing values. Subtle errors such as incorrect joins, inconsistent handling of null and undefined values, or mismatched column references often go undetected until late stages. These issues create avoidable risk in data pipelines. Tidy-TS is designed to minimize this class of error.

This library implements and builds upon a tidyverse-style grammar of data manipulation using TypeScript's type system. At its core, Tidy-TS provides a strongly typed DataFrame abstraction that models tabular data as arrays of objects with explicit column types. Each operation in a pipeline that transforms a DataFrame also updates the underlying inferred structure for

the compiler, allowing it to detect potential errors in downstream transformations, tests, and visualizations. Operations such as mutate, select, drop, join, pivot, and summarize all preserve and refine type information as data evolve through the workflow. Descriptive statistics follow conventions from R but enforce type correctness at compile time rather than through runtime checks. For example, functions that summarize over potentially null values require the developer to explicitly handle missing data before returning numeric output. This type-safe approach is illustrated in Example 1, which demonstrates how a DataFrame-based analysis carries type information through transformations, grouping, and aggregation.

Example 1: Type-Safe DataFrame Transformations

```
import { createDataFrame, stats as s } from "@tidy-ts/dataframe";

const sales = createDataFrame([
  { region: "North", product: "Widget", quantity: 10, price: 100 },
  { region: "North", product: "Gizmo", quantity: 20, price: 50 },
  { region: "South", product: "Widget", quantity: 20, price: 100 },
  { region: "East", product: "Widget", quantity: 8, price: 100 },
]);

const analysis = sales
  .mutate({
    revenue: (row) => row.quantity * row.price,
    moreQuantityThanAvg: (row, _i, df) => row.quantity > s.mean(df.quantity)
  })
  .groupBy("region")
  .summarize({
    total_revenue: (group) => s.sum(group.revenue),
    avg_quantity: (group) => s.mean(group.quantity)
  })
  .arrange("total_revenue", "desc");
```

Tidy-TS also includes support for statistical hypothesis testing (Example 2). These functions are validated against R using randomized test suites to ensure parity in results. The validation framework compares test statistics and p-values across implementations, requiring differences to be within 1e-4 tolerance. Output from statistical tests is returned in typed objects that contain test names, effect sizes, p-values, confidence intervals, and other relevant statistics in a structured format. Providing these features natively is a priority for this work, particularly as research applications increasingly converge data collection and data analysis through technologies like TypeScript and interactive applications. Example 2 demonstrates how Tidy-TS exposes statistical tests directly and via an API that guides users to the correct test based on their intention - invoking the tidyverse philosophy of human-centered design. Both methods return the same structured, typed results.

Example 2: Statistical Hypothesis Testing with Type Safety

```
import { stats as s } from "@tidy-ts/dataframe";

// Create test "height" data for 6 individuals
const heights = [170, 165, 180, 175, 172, 168];

// Direct test API
// Access specific statistical tests
const directTest = s.test.t.oneSample({
  data: heights,
  mu: 170,
  alternative: "two-sided",
```

```

    alpha: 0.05
  });

  // Compare API
  // Intent-driven hypothesis testing api
  const compareAPI = s.compare.oneGroup.centralTendency.toValue({
    data: heights,
    comparator: "not equal to" // or "less than" | "greater than"
    hypothesizedValue: 170,
    parametric: "parametric", // or "nonparametric" | "auto"
    alpha: 0.05
  });

  // Both return the same typed result:
  // {
  //   test_name: "One-sample t-test",
  //   p_value: 0.47...,
  //   effect_size: { value: 0.31..., name: "Cohen's D" },
  //   test_statistic: { value: 0.76..., name: "T-Statistic" },
  //   confidence_interval: { lower: 166.08..., upper: 177.24...,
  //     confidence_level: 0.95 },
  //   degrees_of_freedom: 5,
  //   alpha: 0.05,
  //   alternative: "Two-Sided"
  // }

```

62 Many modern analytics workflows also rely on remote API calls or external services for validation
 63 and analysis. Tidy-TS supports asynchronous transformations, permitting row-wise async
 64 operations within mutate, filter, and summarize operations. It also includes both dataframe-
 65 level and operation-level controls for concurrency and retries. This eases the task of building
 66 analytics pipelines dependent on external services or subject to API rate limits. Applications
 67 that call remote artificial intelligence (AI) models or services are a key use case.

68 Modern workflows also ingest data from various sources. Tidy-TS provides tools to import
 69 CSV, Parquet, Arrow, and JSON in a type-safe manner with the help of Zod schema. Likewise,
 70 databases queries made with raw SQL can be made type-safe with schema validation. For
 71 data queried via popular type-safe Object-Relational Mappers, Tidy-TS can create dataframes
 72 using their provided types directly.

73 Tidy-TS ultimately builds on lessons from tools like pandas, dplyr, and Polars and adapts them
 74 for modern Typescript development needs. Tidy-TS is tested on each commit with a suite of
 75 over 900 tests and across Node, Bun, Deno, and browser targets.

76 Research applications

77 Tidy-TS supports many common data analysis workflows in research environments. At the
 78 University of Utah's Department of Biomedical Informatics, the library is used for healthcare
 79 data analysis where the added type safety aids in integrating multiple data sources. These
 80 projects include both ad hoc quality improvement analysis and real-time evaluation of clinical
 81 data streams intended for integration into the electronic health record.

82 The ability to perform asynchronous transformations can also help researchers more easily
 83 incorporate external data and tooling. One can fetch data from a repository, process and clean
 84 it, invoke async AI tools, perform statistical testing, and visualize the results (see Figure 1) all
 85 within a single TypeScript workflow without requiring a switch to Python or R for analysis.

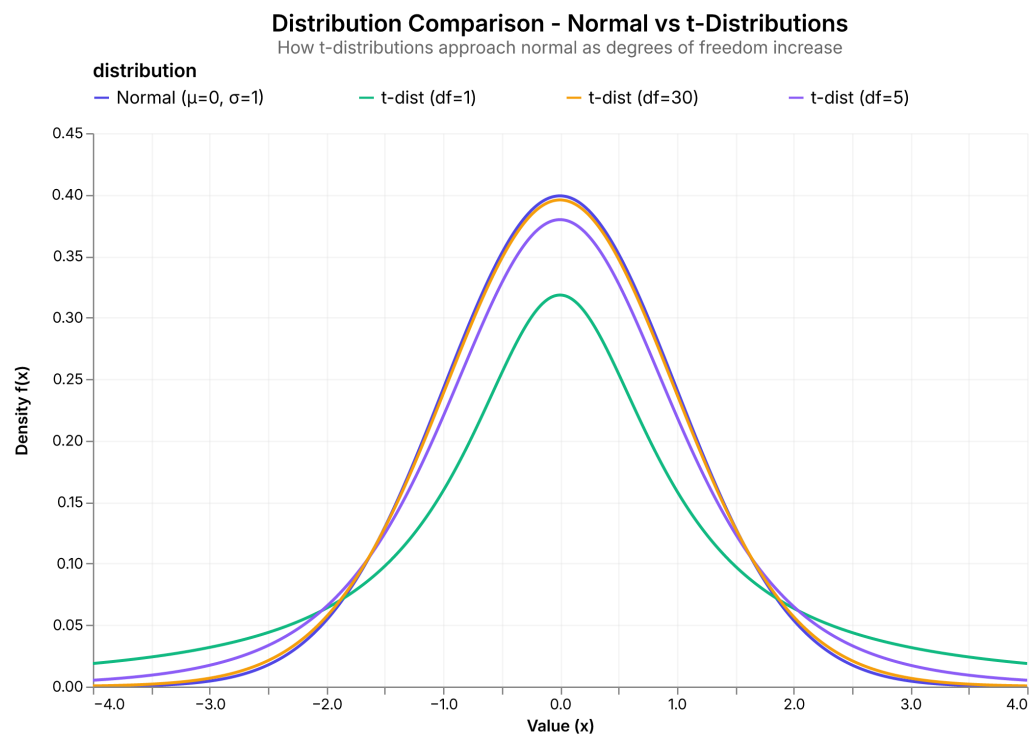


Figure 1: Publication-quality statistical visualization generated with Tidy-TS, showing normal and t-distributions with varying degrees of freedom.

86 Tidy-TS runs in multiple environments including Node.js, Deno, Bun, and modern browsers.
87 It can also be used in Jupyter notebooks through the Deno kernel, enabling interactive data
88 analyses and exploration with minimal configuration in a common workflow in data science
89 and research.

90 Acknowledgements

91 No external funding was received in support of this work. The author is grateful to Dr. Kensaku
92 Kawamoto for his mentorship and support, without which this work would not be possible.

93 Tidy-TS is also only possible due to numerous open source projects. Beyond the pioneering
94 work by the tidyverse teams, the author would like to specifically acknowledge and give thanks
95 to 1) the [University of Washington Interactive Data Lab](#) for their transformative work with
96 Vega-Lite, Vega, and the [arquero javascript](#) library, 2) the [hyparam team](#) for their [parquet](#)
97 processing libraries and 3) the [cross-org team](#) for their cross-runtime testing tooling.

98 Bunge, B. (2019). *Adopting TypeScript at scale*. JSConf Hawaii 2019. <https://www.youtube.com/watch?v=P-J9Eg7hJwE&feature=youtu.be&t=702>

100 Gao, Z., Bird, C., & Barr, E. T. (2017). To type or not to type: Quantifying detectable bugs
101 in JavaScript. *Proceedings of the 39th International Conference on Software Engineering*,
102 758–769.

103 Khan, S., Uddin, G., & Niazi, M. (2021). A systematic literature review on the empirical
104 evidence for the effectiveness of type systems. *IEEE Transactions on Software Engineering*,
105 47(12), 2670–2691. <https://doi.org/10.1109/TSE.2019.2961751>

106 Satyanarayan, A., Moritz, D., Wongsuphasawat, K., & Heer, J. (2017). Vega-lite: A grammar
107 of interactive graphics. *IEEE Transactions on Visualization and Computer Graphics*, 23(1),

- 108 341–350. <https://doi.org/10.1109/TVCG.2016.2599030>
- 109 Satyanarayan, A., Russell, R., Hoffswell, J., & Heer, J. (2016). Reactive vega: A streaming
110 dataflow architecture for declarative interactive visualization. *IEEE Transactions on Visual-*
111 *ization and Computer Graphics*, 22(1), 659–668. [https://doi.org/10.1109/TVCG.2015.](https://doi.org/10.1109/TVCG.2015.2467091)
112 [2467091](https://doi.org/10.1109/TVCG.2015.2467091)
- 113 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Golemund,
114 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Lin Pedersen, T., Miller, E., Milton
115 Bache, S., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., ... Yutani, H.
116 (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686.
117 <https://doi.org/10.21105/joss.01686>

DRAFT