# Overview

## Dataset statistics

| | |
|---|---|
| **Number of variables** | 12 |
| **Number of observations** | 11127 |
| **Missing cells** | 0 |
| **Missing cells (%)** | 0.0% |
| **Duplicate rows** | 0 |
| **Duplicate rows (%)** | 0.0% |
| **Total size in memory** | 1.0 MiB |
| **Average record size in memory** | 96.0 B |

## Variable types

| | |
|---|---|
| **Numeric** | 6 |
| **Text** | 5 |
| **Categorical** | 1 |

## Alerts

| | |
|---|---|
| `ratings_count` is highly overall correlated with `text_reviews_count` | High correlation |
| `text_reviews_count` is highly overall correlated with `ratings_count` | High correlation |
| `language_code` is highly imbalanced (76.6%) | Imbalance |
| `isbn13` is highly skewed ($\gamma_1 = -21.07028799$) | Skewed |
| `bookID` has unique values | Unique |
| `isbn` has unique values | Unique |
| `text_reviews_count` has 625 (5.6%) zeros | Zeros |

## Reproduction

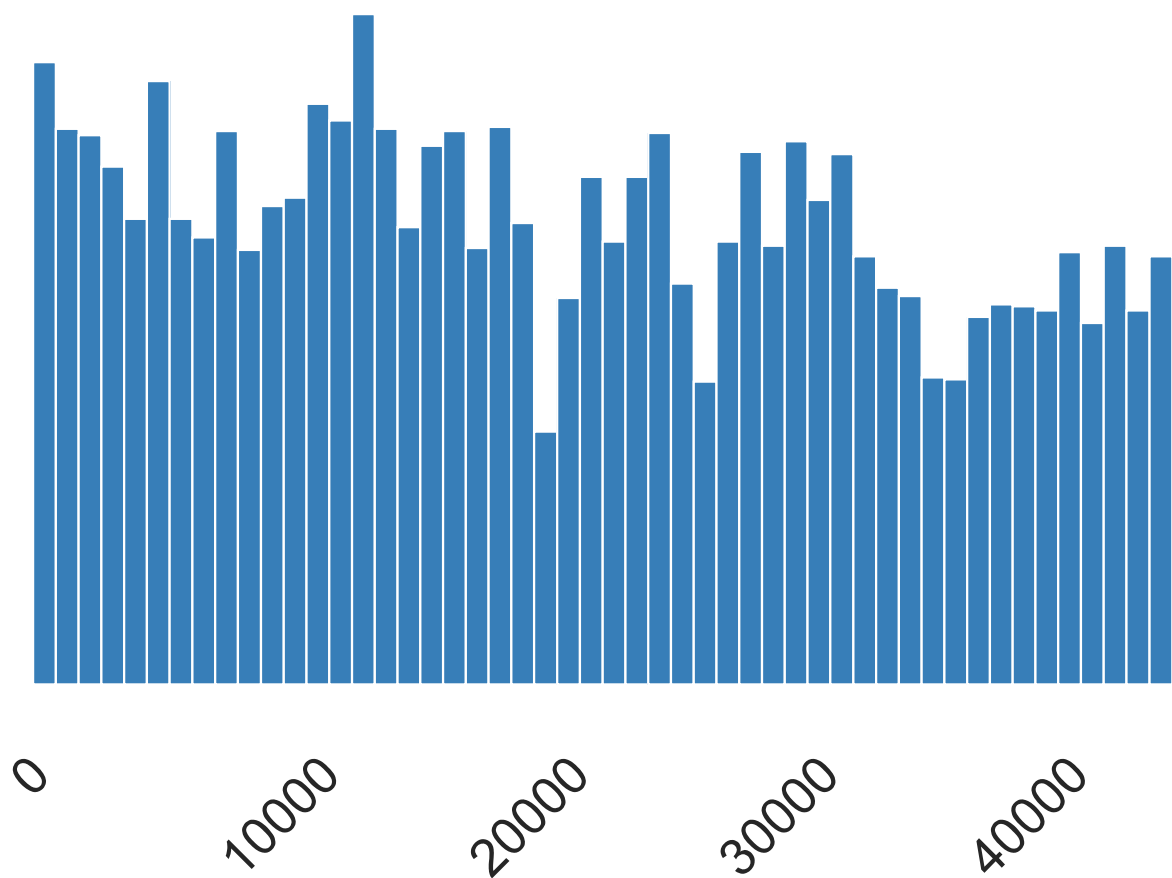| | |
|---|---|
| **Analysis started** | 2024-03-03 14:58:51.657057 |
| **Analysis finished** | 2024-03-03 14:58:58.137990 |
| **Duration** | 6.48 seconds |
| **Software version** | ydata-profiling vv4.6.5 (https://github.com/ydataai/ydata-profiling) |
| **Download configuration** | config.json (data:text/plain;charset=utf-8,%7B%22title%22%3A%20%22Profiling%20Report%22%2C%20%22dataset%22%3A%20%7B%22description%22%3A%20%22%22%2C%20%22creato |

# Variables

Select Columns ⌄

### bookID
Real number (ℝ)

UNIQUE

| | |
|---|---|
| **Distinct** | 11127 |
| **Distinct (%)** | 100.0% |
| **Missing** | 0 |
| **Missing (%)** | 0.0% |
| **Infinite** | 0 |
| **Infinite (%)** | 0.0% |
| **Mean** | 21310.939 |

| | |
|---|---|
| **Minimum** | 1 |
| **Maximum** | 45641 |
| **Zeros** | 0 |
| **Zeros (%)** | 0.0% |
| **Negative** | 0 |
| **Negative (%)** | 0.0% |
| **Memory size** | 87.1 KiB |

## Quantile statistics

| | |
|---|---|
| **Minimum** | 1 |
| **5-th percentile** | 1800.3 |
| **Q1** | 10287 |
| **median** | 20287 |
| **Q3** | 32104.5 |
| **95-th percentile** | 43066.5 |
| **Maximum** | 45641 |
| **Range** | 45640 |
| **Interquartile range (IQR)** | 21817.5 |

## Descriptive statistics

| | |
|---|---|
| **Standard deviation** | 13093.358 |
| **Coefficient of variation (CV)** | 0.61439611 |
| **Kurtosis** | -1.1463568 |
| **Mean** | 21310.939 |
| **Median Absolute Deviation (MAD)** | 10879 |
| **Skewness** | 0.14405166 |
| **Sum** | $2.3712682 \times 10^8$ |
| **Variance** | $1.7143602 \times 10^8$ |

**Monotonicity**                                                                                 Not monotonic

**Histogram with fixed size bins** (bins=50)

| Value | Count | Frequency (%) |
|---|---|---|
| 34889 | 1 | < 0.1% |
| 28532 | 1 | < 0.1% |
| 28510 | 1 | < 0.1% |
| 28511 | 1 | < 0.1% |
| 28514 | 1 | < 0.1% |
| 28522 | 1 | < 0.1% |
| 28524 | 1 | < 0.1% |
| 28529 | 1 | < 0.1% |
| 28530 | 1 | < 0.1% |
| 28531 | 1 | < 0.1% |
| Other values (11117) | 11117 | 99.9% |

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| 1 | 1 | < 0.1% |
| 2 | 1 | < 0.1% |
| 4 | 1 | < 0.1% |
| 5 | 1 | < 0.1% |
| 8 | 1 | < 0.1% |
| 9 | 1 | < 0.1% |
| 10 | 1 | < 0.1% |
| 12 | 1 | < 0.1% |
| 13 | 1 | < 0.1% |
| 14 | 1 | < 0.1% |

| Value | Count | Frequency (%) |
|---|---|---|
| 45641 | 1 | < 0.1% |
| 45639 | 1 | < 0.1% |
| 45634 | 1 | < 0.1% |
| 45633 | 1 | < 0.1% |
| 45631 | 1 | < 0.1% |
| 45630 | 1 | < 0.1% |
| 45626 | 1 | < 0.1% |
| 45625 | 1 | < 0.1% |
| 45623 | 1 | < 0.1% |
| 45617 | 1 | < 0.1% |

## title
Text

| Distinct | 10352 |
|---|---|
| Distinct (%) | 93.0% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Memory size | 87.1 KiB |



## Length

| Max length | 254 |
|---|---|
| Median length | 141 |
| Mean length | 35.749348 |
| Min length | 2 |

## Characters and Unicode

| Total characters | 397783 | |
|---|---|---|
| Distinct characters | 296 | |
| Distinct categories | 17 (https://en.wikipedia.org/wiki/Unicode_character_property#General_Category) | ? |
| Distinct scripts | 8 (https://en.wikipedia.org/wiki/Script_(Unicode)#List_of_scripts_in_Unicode) | ? |
| Distinct blocks | 9 (https://en.wikipedia.org/wiki/Unicode_block) | ? |

The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.

## Unique

| Unique | 9865 | ? |
|---|---|---|
| Unique (%) | 88.7% | |

## Sample

| 1st row | Brown's Star Atlas: Showing All The Bright Stars With Full Instructions How To Find And Use Them For Navigational Purposes And Department Of Trade Examinations. |
|---|---|
| 2nd row | The Tolkien Fan's Medieval Reader |
| 3rd row | Streetcar Suburbs: The Process of Growth in Boston 1870-1900 |
| 4th row | Patriots (The Coming Collapse) |
| 5th row | Harry Potter and the Half-Blood Prince (Harry Potter #6) |

| Value | Count | Frequency (%) |
|---|---|---|
| the | 6692 | 10.1% |
| of | 3336 | 5.0% |
| and | 1653 | 2.5% |
| a | 1335 | 2.0% |
| 1 | 796 | 1.2% |
| in | 778 | 1.2% |
| to | 698 | 1.1% |
| | 588 | 0.9% |
| 2 | 519 | 0.8% |
| 3 | 399 | 0.6% |
| Other values (12076) | 49535 | 74.7% |

## Most occurring characters

| Value | Count | Frequency (%) |
|---|---|---|
|  | 58892 | 14.8% |
| e | 36631 | 9.2% |
| o | 23592 | 5.9% |
| a | 22308 | 5.6% |
| i | 20615 | 5.2% |
| r | 20216 | 5.1% |
| n | 20032 | 5.0% |
| t | 19155 | 4.8% |
| s | 16700 | 4.2% |
| h | 13698 | 3.4% |
| Other values (286) | 145944 | 36.7% |

## Most occurring categories

| Value | Count | Frequency (%) |
|---|---|---|
| Lowercase Letter | 265634 | 66.8% |
| Space Separator | 58892 | 14.8% |
| Uppercase Letter | 52413 | 13.2% |
| Other Punctuation | 8579 | 2.2% |
| Decimal Number | 5487 | 1.4% |
| Close Punctuation | 2765 | 0.7% |
| Open Punctuation | 2764 | 0.7% |
| Dash Punctuation | 808 | 0.2% |
| Other Letter | 373 | 0.1% |
| Math Symbol | 27 | < 0.1% |
| Other values (7) | 41 | < 0.1% |

## Most frequent character per category

*Other Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| の | 16 | 4.3% |
| 夜 | 13 | 3.5% |
| 犬 | 13 | 3.5% |
| 師 | 13 | 3.5% |
| 術 | 13 | 3.5% |
| 金 | 13 | 3.5% |
| 鋼 | 13 | 3.5% |
| 叉 | 13 | 3.5% |
| 碁 | 11 | 2.9% |
| ヒ | 11 | 2.9% |
| Other values (144) | 244 | 65.4% |

*Lowercase Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| e | 36631 | 13.8% |
| o | 23592 | 8.9% |
| a | 22308 | 8.4% |
| i | 20615 | 7.8% |
| r | 20216 | 7.6% |
| n | 20032 | 7.5% |
| t | 19155 | 7.2% |
| s | 16700 | 6.3% |

| Value | Count | Frequency (%) |
|---|---|---|
| h | 13698 | 5.2% |
| I | 12851 | 4.8% |
| Other values (48) | 59836 | 22.5% |

*Uppercase Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| T | 7023 | 13.4% |
| S | 4510 | 8.6% |
| A | 3662 | 7.0% |
| C | 3562 | 6.8% |
| M | 3165 | 6.0% |
| B | 2677 | 5.1% |
| W | 2612 | 5.0% |
| P | 2599 | 5.0% |
| L | 2507 | 4.8% |
| D | 2490 | 4.8% |
| Other values (23) | 17606 | 33.6% |

*Other Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
| : | 3025 | 35.3% |
| # | 2431 | 28.3% |
| ' | 1397 | 16.3% |
| . | 709 | 8.3% |
| / | 414 | 4.8% |
| & | 258 | 3.0% |
| ! | 135 | 1.6% |
| ? | 70 | 0.8% |
| ; | 59 | 0.7% |
| " | 50 | 0.6% |
| Other values (8) | 31 | 0.4% |

*Decimal Number*

| Value | Count | Frequency (%) |
|---|---|---|
| 1 | 1724 | 31.4% |
| 2 | 855 | 15.6% |
| 3 | 645 | 11.8% |
| 4 | 415 | 7.6% |
| 0 | 396 | 7.2% |
| 9 | 371 | 6.8% |

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| 5 | 347 | 6.3% |
| 6 | 282 | 5.1% |
| 8 | 233 | 4.2% |
| 7 | 219 | 4.0% |

*Dash Punctuation*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| - | 773 | 95.7% |
| – | 18 | 2.2% |
| — | 15 | 1.9% |
| ― | 2 | 0.2% |

*Nonspacing Mark*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| ´ | 3 | 50.0% |
| ¨ | 2 | 33.3% |
| ˇ | 1 | 16.7% |

*Close Punctuation*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| ) | 2756 | 99.7% |
| ] | 9 | 0.3% |

*Open Punctuation*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| ( | 2755 | 99.7% |
| [ | 9 | 0.3% |

*Math Symbol*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| + | 18 | 66.7% |
| = | 9 | 33.3% |

*Final Punctuation*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| ' | 12 | 92.3% |
| " | 1 | 7.7% |

*Other Number*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| ½ | 11 | 91.7% |

| Value | Count | Frequency (%) |
|---|---|---|
| ² | 1 | 8.3% |

*Initial Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
| ' | 1 | 50.0% |
| " | 1 | 50.0% |

*Space Separator*

| Value | Count | Frequency (%) |
|---|---|---|
| | 58892 | 100.0% |

*Modifier Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| ― | 3 | 100.0% |

*Currency Symbol*

| Value | Count | Frequency (%) |
|---|---|---|
| $ | 3 | 100.0% |

*Connector Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
| _ | 2 | 100.0% |

## Most occurring scripts

| Value | Count | Frequency (%) |
|---|---|---|
| Latin | 318031 | 80.0% |
| Common | 79357 | 19.9% |
| Han | 232 | 0.1% |
| Katakana | 72 | < 0.1% |
| Hiragana | 61 | < 0.1% |
| Cyrillic | 16 | < 0.1% |
| Arabic | 8 | < 0.1% |
| Inherited | 6 | < 0.1% |

## Most frequent character per script

*Han*

| Value | Count | Frequency (%) |
|---|---|---|
| 夜 | 13 | 5.6% |
| 犬 | 13 | 5.6% |
| 師 | 13 | 5.6% |
| 術 | 13 | 5.6% |
| 金 | 13 | 5.6% |
| 鋼 | 13 | 5.6% |
| 叉 | 13 | 5.6% |
| 碁 | 11 | 4.7% |
| 之 | 8 | 3.4% |
| 鍊 | 8 | 3.4% |
| Other values (90) | 114 | 49.1% |

*Latin*

| Value | Count | Frequency (%) |
|---|---|---|
| e | 36631 | 11.5% |
| o | 23592 | 7.4% |
| a | 22308 | 7.0% |
| i | 20615 | 6.5% |
| r | 20216 | 6.4% |
| n | 20032 | 6.3% |
| t | 19155 | 6.0% |
| s | 16700 | 5.3% |
| h | 13698 | 4.3% |
| l | 12851 | 4.0% |
| Other values (73) | 112233 | 35.3% |

*Common*

| Value | Count | Frequency (%) |
|---|---|---|
|  | 58892 | 74.2% |
| : | 3025 | 3.8% |
| ) | 2756 | 3.5% |
| ( | 2755 | 3.5% |
| # | 2431 | 3.1% |
| 1 | 1724 | 2.2% |
| ' | 1397 | 1.8% |
| 2 | 855 | 1.1% |
| - | 773 | 1.0% |
| . | 709 | 0.9% |
| Other values (38) | 4040 | 5.1% |

*Hiragana*

| Value | Count | Frequency (%) |
|---|---|---|
| の | 16 | 26.2% |
| ん | 5 | 8.2% |
| か | 3 | 4.9% |
| る | 3 | 4.9% |
| て | 3 | 4.9% |
| き | 3 | 4.9% |
| た | 3 | 4.9% |
| ら | 3 | 4.9% |
| な | 2 | 3.3% |
| ぜ | 2 | 3.3% |
| Other values (14) | 18 | 29.5% |

*Katakana*

| Value | Count | Frequency (%) |
|---|---|---|
| ヒ | 11 | 15.3% |
| ル | 11 | 15.3% |
| カ | 11 | 15.3% |
| ツ | 5 | 6.9% |
| サ | 5 | 6.9% |
| バ | 5 | 6.9% |
| リ | 2 | 2.8% |
| ト | 2 | 2.8% |
| ス | 2 | 2.8% |
| ャ | 2 | 2.8% |

| Value | Count | Frequency (%) |
|---|---|---|
| Other values (14) | 16 | 22.2% |

*Cyrillic*

| Value | Count | Frequency (%) |
|---|---|---|
| а | 4 | 25.0% |
| р | 3 | 18.8% |
| м | 2 | 12.5% |
| т | 2 | 12.5% |
| и | 2 | 12.5% |
| с | 1 | 6.2% |
| е | 1 | 6.2% |
| г | 1 | 6.2% |

*Arabic*

| Value | Count | Frequency (%) |
|---|---|---|
| م | 2 | 25.0% |
| ل | 2 | 25.0% |
| ح | 1 | 12.5% |
| ا | 1 | 12.5% |
| ن | 1 | 12.5% |
| د | 1 | 12.5% |

*Inherited*

| Value | Count | Frequency (%) |
|---|---|---|
| ́ | 3 | 50.0% |
| ̈ | 2 | 33.3% |
| ̌ | 1 | 16.7% |

## Most occurring blocks

| Value | Count | Frequency (%) |
|---|---|---|
| ASCII | 397016 | 99.8% |
| None | 319 | 0.1% |
| CJK | 232 | 0.1% |
| Katakana | 75 | < 0.1% |
| Hiragana | 61 | < 0.1% |
| Punctuation | 50 | < 0.1% |
| Cyrillic | 16 | < 0.1% |
| Arabic | 8 | < 0.1% |
| Diacriticals | 6 | < 0.1% |

## Most frequent character per block

*ASCII*

| Value | Count | Frequency (%) |
|---|---|---|
|  | 58892 | 14.8% |
| e | 36631 | 9.2% |
| o | 23592 | 5.9% |
| a | 22308 | 5.6% |
| i | 20615 | 5.2% |
| r | 20216 | 5.1% |
| n | 20032 | 5.0% |
| t | 19155 | 4.8% |
| s | 16700 | 4.2% |
| h | 13698 | 3.5% |
| Other values (75) | 145177 | 36.6% |

*None*

| Value | Count | Frequency (%) |
|---|---|---|
| é | 62 | 19.4% |
| á | 35 | 11.0% |
| ó | 31 | 9.7% |
| í | 28 | 8.8% |
| ä | 19 | 6.0% |
| ü | 17 | 5.3% |
| ñ | 14 | 4.4% |
| ½ | 11 | 3.4% |
| ` | 11 | 3.4% |
| è | 11 | 3.4% |

| Value | Count | Frequency (%) |
|---|---|---|
| Other values (28) | 80 | 25.1% |

*Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
| 一 | 18 | 36.0% |
| — | 15 | 30.0% |
| ' | 12 | 24.0% |
| — | 2 | 4.0% |
| ' | 1 | 2.0% |
| " | 1 | 2.0% |
| " | 1 | 2.0% |

*Hiragana*

| Value | Count | Frequency (%) |
|---|---|---|
| の | 16 | 26.2% |
| ん | 5 | 8.2% |
| か | 3 | 4.9% |
| る | 3 | 4.9% |
| て | 3 | 4.9% |
| き | 3 | 4.9% |
| た | 3 | 4.9% |
| ら | 3 | 4.9% |
| な | 2 | 3.3% |
| ぜ | 2 | 3.3% |
| Other values (14) | 18 | 29.5% |

*CJK*

| Value | Count | Frequency (%) |
|---|---|---|
| 夜 | 13 | 5.6% |
| 犬 | 13 | 5.6% |
| 師 | 13 | 5.6% |
| 術 | 13 | 5.6% |
| 金 | 13 | 5.6% |
| 鋼 | 13 | 5.6% |
| 又 | 13 | 5.6% |
| 碁 | 11 | 4.7% |
| 之 | 8 | 3.4% |
| 錬 | 8 | 3.4% |
| Other values (90) | 114 | 49.1% |

*Katakana*

| Value | Count | Frequency (%) |
|---|---|---|
| ヒ | 11 | 14.7% |
| ル | 11 | 14.7% |
| カ | 11 | 14.7% |
| ツ | 5 | 6.7% |
| サ | 5 | 6.7% |
| バ | 5 | 6.7% |
| ー | 3 | 4.0% |
| リ | 2 | 2.7% |
| ト | 2 | 2.7% |
| ス | 2 | 2.7% |
| Other values (15) | 18 | 24.0% |

*Cyrillic*

| Value | Count | Frequency (%) |
|---|---|---|
| а | 4 | 25.0% |
| р | 3 | 18.8% |
| м | 2 | 12.5% |
| т | 2 | 12.5% |
| и | 2 | 12.5% |
| с | 1 | 6.2% |
| е | 1 | 6.2% |
| г | 1 | 6.2% |

*Diacriticals*

| Value | Count | Frequency (%) |
|---|---|---|
| ´ | 3 | 50.0% |
| ¨ | 2 | 33.3% |
| ˇ | 1 | 16.7% |

*Arabic*

| Value | Count | Frequency (%) |
|---|---|---|
| م | 2 | 25.0% |
| ل | 2 | 25.0% |
| ح | 1 | 12.5% |
| ا | 1 | 12.5% |
| ن | 1 | 12.5% |
| د | 1 | 12.5% |

## authors
Text

| | |
|---|---|
| **Distinct** | 6643 |
| **Distinct (%)** | 59.7% |
| **Missing** | 0 |
| **Missing (%)** | 0.0% |
| **Memory size** | 87.1 KiB |



## Length

| | |
|---|---|
| **Max length** | 750 |
| **Median length** | 372 |
| **Mean length** | 24.724005 |
| **Min length** | 3 |

## Characters and Unicode

| | | |
|---|---|---|
| **Total characters** | 275104 | |
| **Distinct characters** | 267 | |
| **Distinct categories** | 11 (https://en.wikipedia.org/wiki/Unicode_character_property#General_Category) | ? |
| **Distinct scripts** | 8 (https://en.wikipedia.org/wiki/Script_(Unicode)#List_of_scripts_in_Unicode) | ? |
| **Distinct blocks** | 8 (https://en.wikipedia.org/wiki/Unicode_block) | ? |

The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.
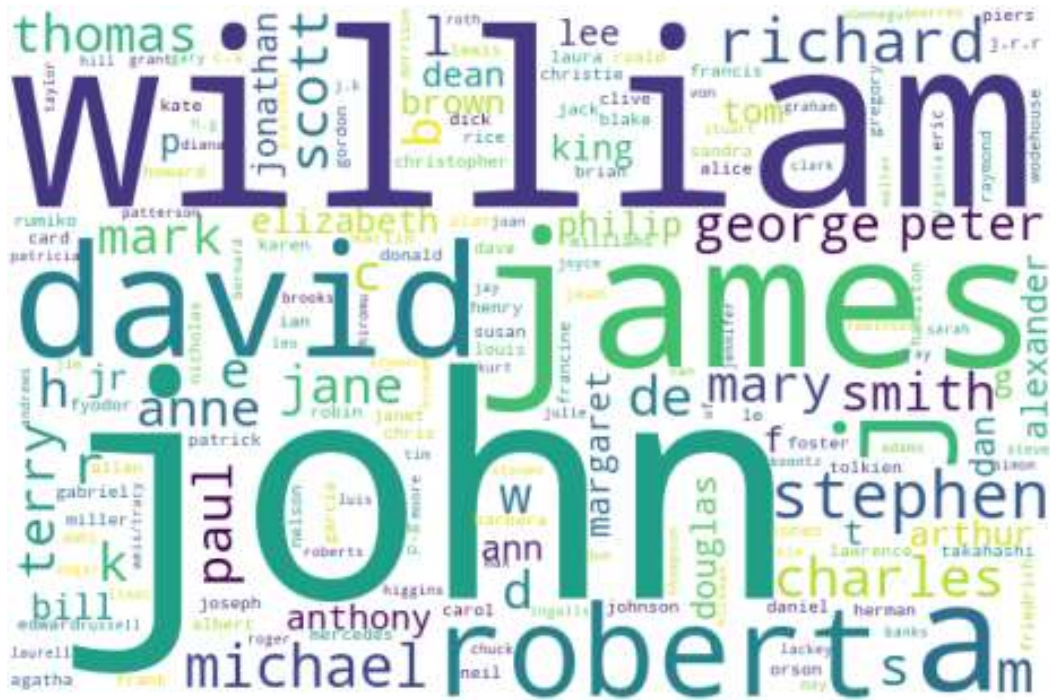
## Unique

| Unique | 5282 | ? |
|---|---|---|
| Unique (%) | 47.5% | |

## Sample

| | |
|---|---|
| **1st row** | Brown Son & Ferguson |
| **2nd row** | David E. Smith (Turgon of TheOneRing.net one of the founding members of this Tolkien website)/Verlyn Flieger/Turgon (=David E. Smith) |
| **3rd row** | Sam Bass Warner Jr./Sam B. Warner |
| **4th row** | James Wesley Rawles |
| **5th row** | J.K. Rowling/Mary GrandPré |

| Value | Count | Frequency (%) |
|---|---|---|
| john | 279 | 0.8% |
| william | 262 | 0.8% |
| james | 228 | 0.7% |
| david | 203 | 0.6% |
| a | 191 | 0.6% |
| robert | 185 | 0.5% |
| j | 181 | 0.5% |
| stephen | 176 | 0.5% |
| richard | 157 | 0.5% |
| m | 155 | 0.5% |
| Other values (12644) | 31795 | 94.0% |

## Most occurring characters

| Value | Count | Frequency (%) |
|---|---|---|
| e | 23758 | 8.6% |
|  | 23425 | 8.5% |
| a | 22547 | 8.2% |
| r | 18119 | 6.6% |
| n | 17426 | 6.3% |
| i | 15677 | 5.7% |
| o | 14415 | 5.2% |
| l | 13273 | 4.8% |
| s | 9822 | 3.6% |
| t | 9527 | 3.5% |
| Other values (257) | 107115 | 38.9% |

## Most occurring categories

| Value | Count | Frequency (%) |
|---|---|---|
| Lowercase Letter | 195795 | 71.2% |
| Uppercase Letter | 43440 | 15.8% |
| Space Separator | 23425 | 8.5% |
| Other Punctuation | 11849 | 4.3% |
| Other Letter | 378 | 0.1% |
| Dash Punctuation | 200 | 0.1% |
| Close Punctuation | 5 | < 0.1% |
| Open Punctuation | 5 | < 0.1% |
| Decimal Number | 4 | < 0.1% |
| Format | 2 | < 0.1% |

## Most frequent character per category

*Other Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| ا | 25 | 6.6% |
| ر | 22 | 5.8% |
| ن | 19 | 5.0% |
| ل | 18 | 4.8% |
| ب | 15 | 4.0% |
| م | 13 | 3.4% |
| ج | 12 | 3.2% |
| 川 | 9 | 2.4% |
| ي | 9 | 2.4% |
| 方 | 8 | 2.1% |
| Other values (99) | 228 | 60.3% |

*Lowercase Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| e | 23758 | 12.1% |
| a | 22547 | 11.5% |
| r | 18119 | 9.3% |
| n | 17426 | 8.9% |
| i | 15677 | 8.0% |
| o | 14415 | 7.4% |
| l | 13273 | 6.8% |
| s | 9822 | 5.0% |
| t | 9527 | 4.9% |

| Value | Count | Frequency (%) |
|---|---|---|
| h | 7628 | 3.9% |
| Other values (86) | 43603 | 22.3% |

*Uppercase Letter*

| Value | Count | Frequency (%) |
|---|---|---|
| M | 3731 | 8.6% |
| S | 3497 | 8.1% |
| J | 3286 | 7.6% |
| C | 3074 | 7.1% |
| R | 2655 | 6.1% |
| A | 2611 | 6.0% |
| B | 2520 | 5.8% |
| D | 2490 | 5.7% |
| H | 2232 | 5.1% |
| P | 2200 | 5.1% |
| Other values (37) | 15144 | 34.9% |

*Other Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
| / | 8117 | 68.5% |
| . | 3577 | 30.2% |
| ' | 148 | 1.2% |
| ! | 4 | < 0.1% |
| " | 2 | < 0.1% |
| & | 1 | < 0.1% |

*Decimal Number*

| Value | Count | Frequency (%) |
|---|---|---|
| 9 | 2 | 50.0% |
| 1 | 1 | 25.0% |
| 2 | 1 | 25.0% |

*Space Separator*

| Value | Count | Frequency (%) |
|---|---|---|
| | 23425 | 100.0% |

*Dash Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
| - | 200 | 100.0% |

*Close Punctuation*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| )     | 5     | 100.0%        |

*Open Punctuation*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| (     | 5     | 100.0%        |

*Format*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
|       | 2     | 100.0%        |

*Math Symbol*

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| =     | 1     | 100.0%        |

## Most occurring scripts

| Value | Count | Frequency (%) |
|---|---|---|
| Latin | 239152 | 86.9% |
| Common | 35489 | 12.9% |
| Arabic | 187 | 0.1% |
| Han | 173 | 0.1% |
| Greek | 53 | < 0.1% |
| Cyrillic | 30 | < 0.1% |
| Hiragana | 18 | < 0.1% |
| Inherited | 2 | < 0.1% |

## Most frequent character per script

*Latin*

| Value | Count | Frequency (%) |
|---|---|---|
| e | 23758 | 9.9% |
| a | 22547 | 9.4% |
| r | 18119 | 7.6% |
| n | 17426 | 7.3% |
| i | 15677 | 6.6% |
| o | 14415 | 6.0% |
| l | 13273 | 5.6% |
| s | 9822 | 4.1% |
| t | 9527 | 4.0% |
| h | 7628 | 3.2% |
| Other values (88) | 86960 | 36.4% |

*Han*

| Value | Count | Frequency (%) |
|---|---|---|
| 川 | 9 | 5.2% |
| 方 | 8 | 4.6% |
| 荒 | 8 | 4.6% |
| 二 | 8 | 4.6% |
| 郁 | 8 | 4.6% |
| 仁 | 8 | 4.6% |
| 弘 | 8 | 4.6% |
| 伊 | 8 | 4.6% |
| 藤 | 8 | 4.6% |
| 潤 | 8 | 4.6% |
| Other values (56) | 92 | 53.2% |

*Arabic*

| Value | Count | Frequency (%) |
|---|---|---|
| ا | 25 | 13.4% |
| ر | 22 | 11.8% |
| ن | 19 | 10.2% |
| ل | 18 | 9.6% |
| ب | 15 | 8.0% |
| م | 13 | 7.0% |
| ج | 12 | 6.4% |
| ي | 9 | 4.8% |
| ى | 7 | 3.7% |
| خ | 6 | 3.2% |
| Other values (19) | 41 | 21.9% |

*Greek*

| Value | Count | Frequency (%) |
|---|---|---|
| ο | 7 | 13.2% |
| α | 5 | 9.4% |
| υ | 4 | 7.5% |
| ς | 4 | 7.5% |
| λ | 4 | 7.5% |
| ί | 3 | 5.7% |
| ι | 3 | 5.7% |
| κ | 2 | 3.8% |
| τ | 2 | 3.8% |
| π | 2 | 3.8% |
| Other values (17) | 17 | 32.1% |

*Cyrillic*

| Value | Count | Frequency (%) |
|---|---|---|
| а | 5 | 16.7% |
| л | 4 | 13.3% |
| и | 3 | 10.0% |
| н | 2 | 6.7% |
| о | 2 | 6.7% |
| в | 2 | 6.7% |
| А | 1 | 3.3% |
| В | 1 | 3.3% |
| Ь | 1 | 3.3% |
| е | 1 | 3.3% |

| Value | Count | Frequency (%) |
|---|---|---|
| Other values (8) | 8 | 26.7% |

*Common*

| Value | Count | Frequency (%) |
|---|---|---|
| | 23425 | 66.0% |
| / | 8117 | 22.9% |
| . | 3577 | 10.1% |
| - | 200 | 0.6% |
| ' | 148 | 0.4% |
| ) | 5 | < 0.1% |
| ( | 5 | < 0.1% |
| ! | 4 | < 0.1% |
| 9 | 2 | < 0.1% |
| " | 2 | < 0.1% |
| Other values (4) | 4 | < 0.1% |

*Hiragana*

| Value | Count | Frequency (%) |
|---|---|---|
| き | 2 | 11.1% |
| た | 2 | 11.1% |
| し | 2 | 11.1% |
| か | 2 | 11.1% |
| ぐ | 1 | 5.6% |
| つ | 1 | 5.6% |
| み | 1 | 5.6% |
| ず | 1 | 5.6% |
| あ | 1 | 5.6% |
| ゆ | 1 | 5.6% |
| Other values (4) | 4 | 22.2% |

*Inherited*

| Value | Count | Frequency (%) |
|---|---|---|
| | 2 | 100.0% |

## Most occurring blocks

| Value | Count | Frequency (%) |
| --- | --- | --- |
| ASCII | 274080 | 99.6% |
| None | 613 | 0.2% |
| Arabic | 187 | 0.1% |
| CJK | 173 | 0.1% |
| Cyrillic | 30 | < 0.1% |
| Hiragana | 18 | < 0.1% |
| Punctuation | 2 | < 0.1% |
| Latin Ext Additional | 1 | < 0.1% |

## Most frequent character per block

*ASCII*

| Value | Count | Frequency (%) |
| --- | --- | --- |
| e | 23758 | 8.7% |
|  | 23425 | 8.5% |
| a | 22547 | 8.2% |
| r | 18119 | 6.6% |
| n | 17426 | 6.4% |
| i | 15677 | 5.7% |
| o | 14415 | 5.3% |
| l | 13273 | 4.8% |
| s | 9822 | 3.6% |
| t | 9527 | 3.5% |
| Other values (56) | 106091 | 38.7% |

*None*

| Value | Count | Frequency (%) |
| --- | --- | --- |
| é | 115 | 18.8% |
| í | 85 | 13.9% |
| á | 57 | 9.3% |
| ō | 45 | 7.3% |
| ó | 34 | 5.5% |
| ë | 20 | 3.3% |
| è | 20 | 3.3% |
| ü | 19 | 3.1% |
| ł | 17 | 2.8% |
| ï | 15 | 2.4% |
| Other values (62) | 186 | 30.3% |

*Arabic*

| Value | Count | Frequency (%) |
|---|---|---|
| ا | 25 | 13.4% |
| ر | 22 | 11.8% |
| ن | 19 | 10.2% |
| ل | 18 | 9.6% |
| ب | 15 | 8.0% |
| م | 13 | 7.0% |
| ح | 12 | 6.4% |
| ي | 9 | 4.8% |
| ى | 7 | 3.7% |
| خ | 6 | 3.2% |
| Other values (19) | 41 | 21.9% |

*CJK*

| Value | Count | Frequency (%) |
|---|---|---|
| 川 | 9 | 5.2% |
| 方 | 8 | 4.6% |
| 荒 | 8 | 4.6% |
| 二 | 8 | 4.6% |
| 郁 | 8 | 4.6% |
| 仁 | 8 | 4.6% |
| 弘 | 8 | 4.6% |
| 伊 | 8 | 4.6% |
| 藤 | 8 | 4.6% |
| 潤 | 8 | 4.6% |
| Other values (56) | 92 | 53.2% |

*Cyrillic*

| Value | Count | Frequency (%) |
|---|---|---|
| а | 5 | 16.7% |
| л | 4 | 13.3% |
| и | 3 | 10.0% |
| н | 2 | 6.7% |
| о | 2 | 6.7% |
| в | 2 | 6.7% |
| А | 1 | 3.3% |
| В | 1 | 3.3% |
| ь | 1 | 3.3% |
| е | 1 | 3.3% |

| Value | Count | Frequency (%) |
|---|---|---|
| Other values (8) | 8 | 26.7% |

*Hiragana*

| Value | Count | Frequency (%) |
|---|---|---|
| き | 2 | 11.1% |
| た | 2 | 11.1% |
| し | 2 | 11.1% |
| か | 2 | 11.1% |
| ぐ | 1 | 5.6% |
| つ | 1 | 5.6% |
| み | 1 | 5.6% |
| ず | 1 | 5.6% |
| あ | 1 | 5.6% |
| ゆ | 1 | 5.6% |
| Other values (4) | 4 | 22.2% |

*Punctuation*

| Value | Count | Frequency (%) |
|---|---|---|
|  | 2 | 100.0% |

*Latin Ext Additional*

| Value | Count | Frequency (%) |
|---|---|---|
| ệ | 1 | 100.0% |

## average_rating
Real number ($\mathbb{R}$)

| | |
|---|---|
| **Distinct** | 209 |
| **Distinct (%)** | 1.9% |
| **Missing** | 0 |
| **Missing (%)** | 0.0% |
| **Infinite** | 0 |
| **Infinite (%)** | 0.0% |
| **Mean** | 3.9336308 |
| | |
| **Minimum** | 0 |
| **Maximum** | 5 |
| **Zeros** | 26 |
| **Zeros (%)** | 0.2% |
| **Negative** | 0 |
| **Negative (%)** | 0.0% |
| **Memory size** | 87.1 KiB |



Quantile statistics

| Minimum | 0 |
|---|---|
| 5-th percentile | 3.44 |
| Q1 | 3.77 |
| median | 3.96 |
| Q3 | 4.135 |
| 95-th percentile | 4.38 |
| Maximum | 5 |
| Range | 5 |
| Interquartile range (IQR) | 0.365 |

## Descriptive statistics

| Standard deviation | 0.35244503 |
|---|---|
| Coefficient of variation (CV) | 0.089597893 |
| Kurtosis | 36.721777 |
| Mean | 3.9336308 |
| Median Absolute Deviation (MAD) | 0.18 |
| Skewness | -3.6383114 |
| Sum | 43769.51 |
| Variance | 0.1242175 |
| Monotonicity | Not monotonic |

**Histogram with fixed size bins** (bins=50)

| Value | Count | Frequency (%) |
|---|---|---|
| 4 | 219 | 2.0% |
| 3.96 | 195 | 1.8% |
| 4.02 | 178 | 1.6% |
| 3.94 | 176 | 1.6% |
| 4.07 | 172 | 1.5% |
| 3.92 | 168 | 1.5% |
| 3.93 | 168 | 1.5% |
| 4.05 | 168 | 1.5% |
| 3.83 | 166 | 1.5% |
| 3.89 | 166 | 1.5% |
| Other values (199) | 9351 | 84.0% |

| Value | Count | Frequency (%) |
|---|---|---|
| 0 | 26 | 0.2% |
| 1 | 2 | < 0.1% |
| 1.67 | 1 | < 0.1% |
| 2 | 6 | 0.1% |
| 2.33 | 1 | < 0.1% |
| 2.4 | 1 | < 0.1% |
| 2.5 | 1 | < 0.1% |
| 2.55 | 1 | < 0.1% |
| 2.61 | 1 | < 0.1% |
| 2.62 | 3 | < 0.1% |

| Value | Count | Frequency (%) |
|-------|-------|---------------|
| 5 | 22 | 0.2% |
| 4.91 | 1 | < 0.1% |
| 4.88 | 1 | < 0.1% |
| 4.86 | 1 | < 0.1% |
| 4.83 | 1 | < 0.1% |
| 4.82 | 1 | < 0.1% |
| 4.8 | 1 | < 0.1% |
| 4.78 | 2 | < 0.1% |
| 4.76 | 1 | < 0.1% |
| 4.75 | 2 | < 0.1% |