# Exploratory Data Analysis (EDA)

Statsomat.com

19 April 2021

# Basic Information

Automatic statistics for the file:

| File |
| --- |
| Baseball_small.csv |

Your selection for the encoding: Auto
Your selection for the decimal character: Auto
Observations (rows with at least one non-missing value): 18
Variables (columns with at least one non-missing value): 23
Variables considered continuous: 16

| Variables considered continuous |
| --- |
| assist86 |
| atbat |
| atbat86 |
| error86 |
| hits |
| hits86 |
| homer86 |
| homeruns |
| outs86 |
| rbi |
| rbi86 |
| runs |
| runs86 |
| V1 |
| walks |
| walks86 |

Variables considered categorical: 7

| Variables considered categorical |
| --- |
| div86 |
| league86 |
| name1 |
| name2 |
| posit86 |
| team86 |
| years |

# Results for Numerical Variables

## Descriptive Statistics

Variables are sorted alphabetically. Missings are omitted in stats. CV only for positive variables.

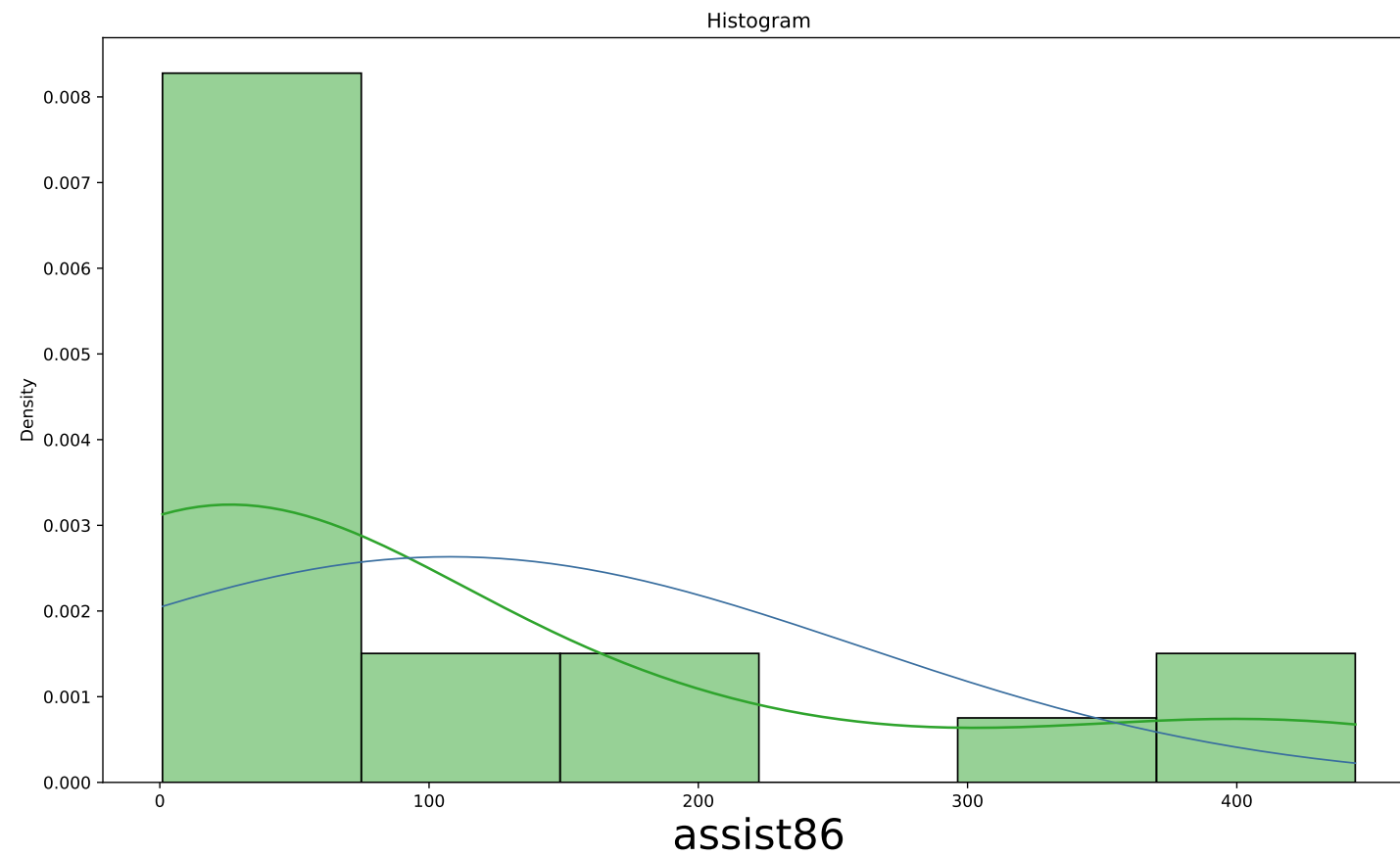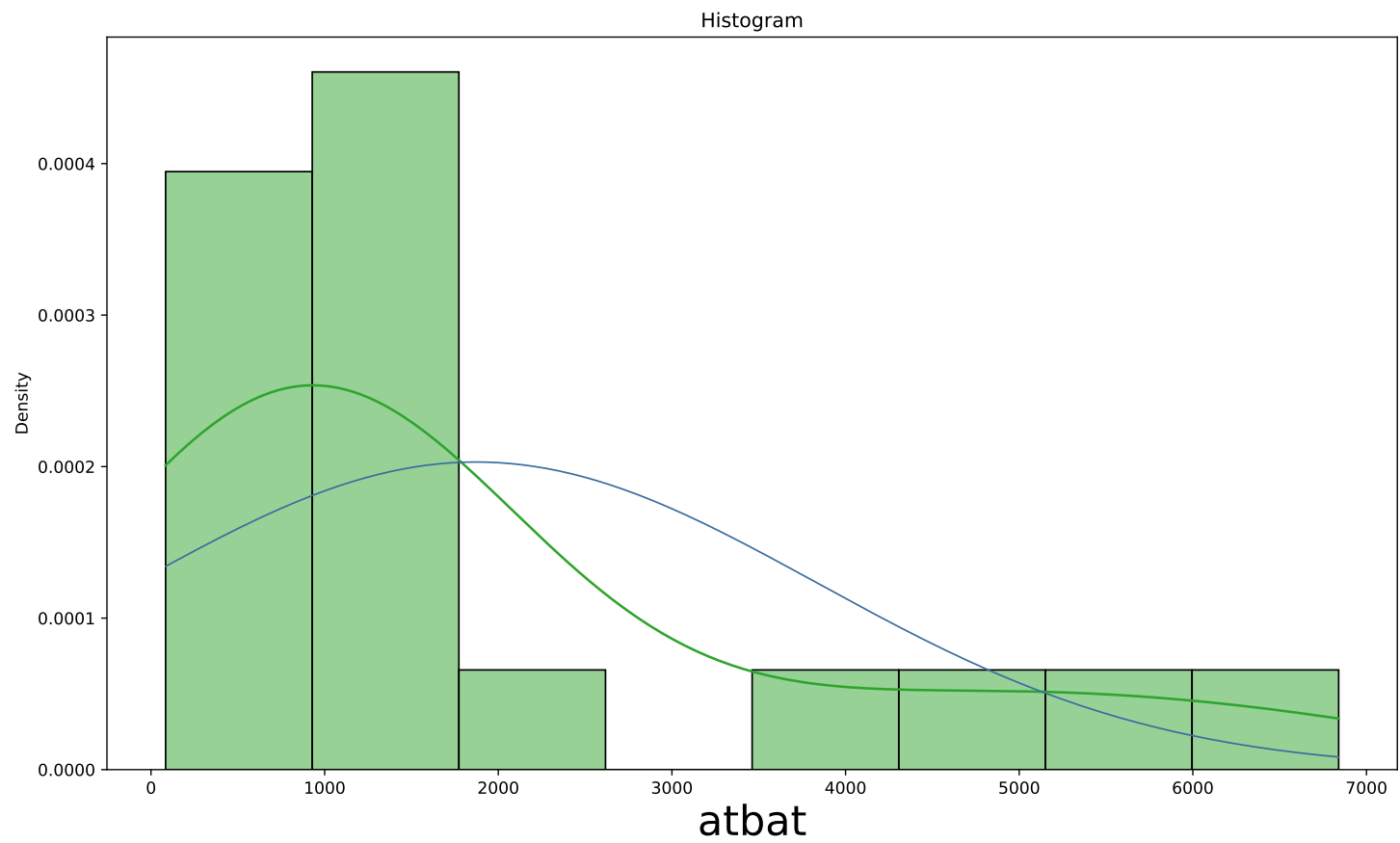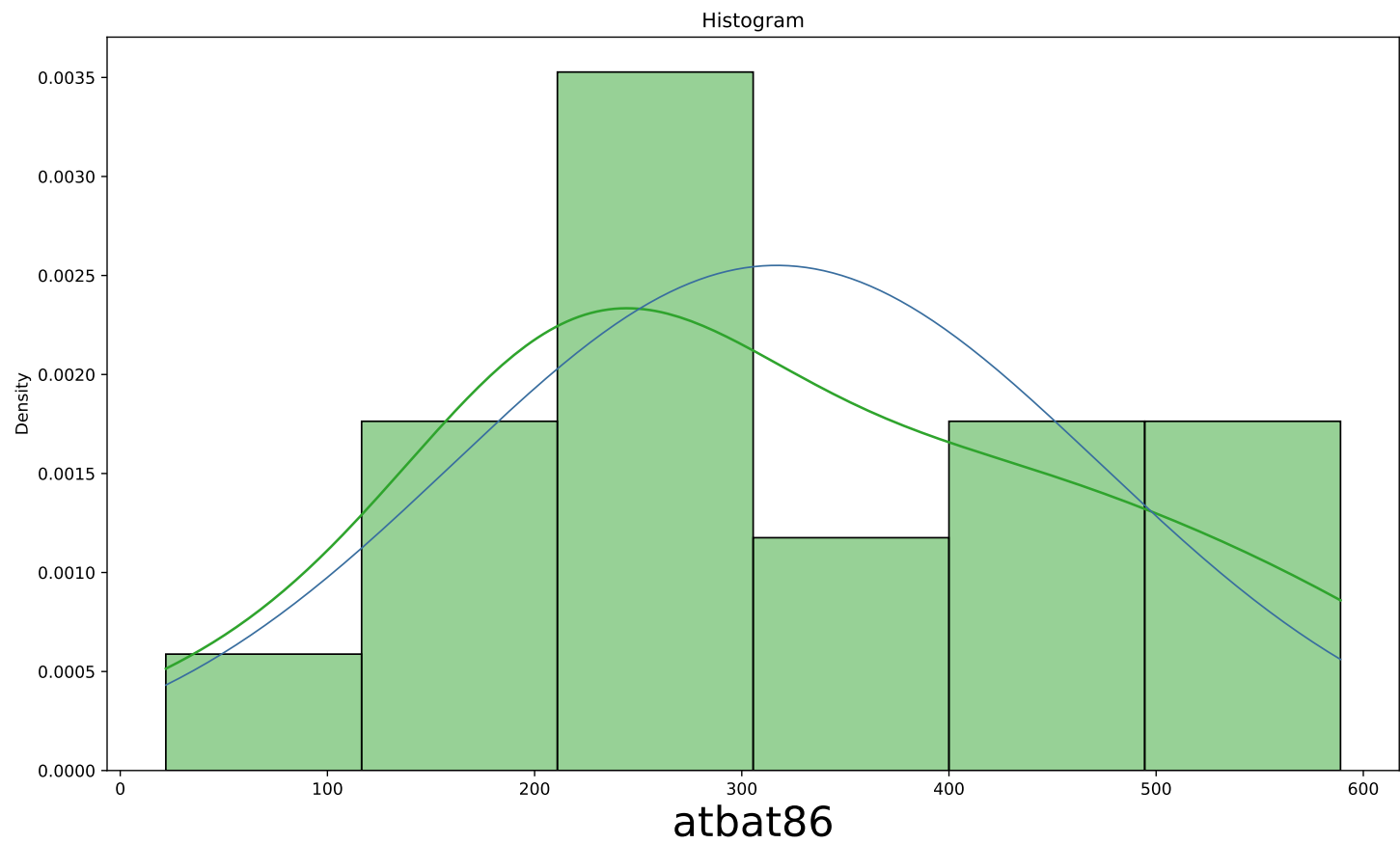|  | N Obs | N Missing | N Valid | % Complete | N Unique | Mean | SD | Median | MAD | Min | Max | Skewness | Kurtosis | CV |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| assist86 | 18 | 0 | 18 | 100 | 16 | 107.78 | 151.52 | 23.5 | 28.91 | 1 | 444 | 1.44 | 0.75 | 1.41 |
| atbat | 18 | 0 | 18 | 100 | 18 | 1872.39 | 1964.87 | 1160.5 | 825.07 | 84 | 6840 | 1.55 | 1.45 | 1.05 |
| atbat86 | 18 | 0 | 18 | 100 | 18 | 316.83 | 156.39 | 268.5 | 136.4 | 22 | 589 | 0.23 | -0.51 | 0.49 |
| error86 | 18 | 0 | 18 | 100 | 12 | 7.28 | 4.98 | 5.5 | 4.45 | 1 | 20 | 1.18 | 1.26 | 0.68 |
| hits | 18 | 0 | 18 | 100 | 18 | 504.56 | 556.62 | 307.5 | 248.34 | 26 | 1910 | 1.59 | 1.53 | 1.1 |
| hits86 | 18 | 0 | 18 | 100 | 18 | 82.33 | 41.53 | 72.5 | 37.07 | 10 | 157 | 0.28 | -0.77 | 0.5 |
| homer86 | 18 | 0 | 18 | 100 | 14 | 8.06 | 6.55 | 5.5 | 5.19 | 1 | 21 | 0.9 | -0.49 | 0.81 |
| homeruns | 18 | 0 | 18 | 100 | 17 | 42.17 | 61.25 | 24 | 25.95 | 2 | 259 | 2.93 | 9.74 | 1.45 |
| outs86 | 18 | 0 | 18 | 100 | 17 | 241.89 | 168 | 214.5 | 101.56 | 59 | 812 | 2.42 | 7.81 | 0.69 |
| rbi | 18 | 0 | 18 | 100 | 18 | 222.61 | 269.25 | 119.5 | 120.83 | 9 | 1067 | 2.12 | 4.94 | 1.21 |
| rbi86 | 18 | 0 | 18 | 100 | 16 | 37.83 | 23.48 | 28.5 | 21.5 | 2 | 86 | 0.71 | -0.33 | 0.62 |
| runs | 18 | 0 | 18 | 100 | 17 | 248.33 | 274.21 | 144.5 | 122.31 | 9 | 915 | 1.46 | 0.97 | 1.1 |
| runs86 | 18 | 0 | 18 | 100 | 16 | 42.56 | 23.5 | 41 | 19.27 | 4 | 95 | 0.84 | 0.74 | 0.55 |
| V1 | 18 | 0 | 18 | 100 | 18 | 132.94 | 76.2 | 114.5 | 97.11 | 22 | 273 | 0.19 | -1.16 | 0.57 |
| walks | 18 | 0 | 18 | 100 | 18 | 154.83 | 173.99 | 98 | 92.66 | 3 | 576 | 1.6 | 1.74 | 1.12 |
| walks86 | 18 | 0 | 18 | 100 | 16 | 27.22 | 18.58 | 21.5 | 13.34 | 1 | 64 | 1.04 | 0.26 | 0.68 |

# Graphics

## Histograms

Details: Density Histograms. One large figure per page for each variable, sorted alphabetically. The blue line represents the normal density approximation. The green line represents a special kernel density approximation.
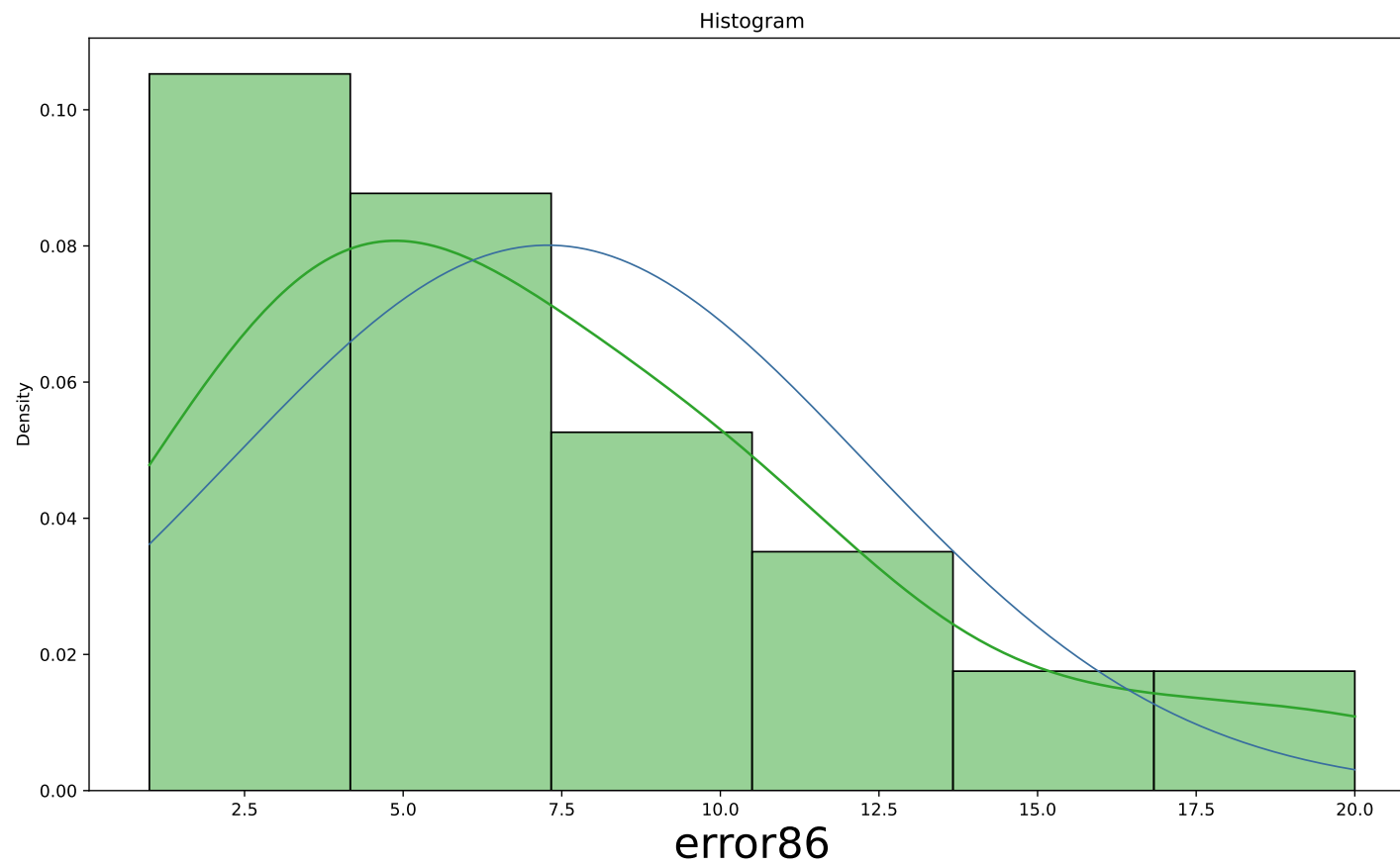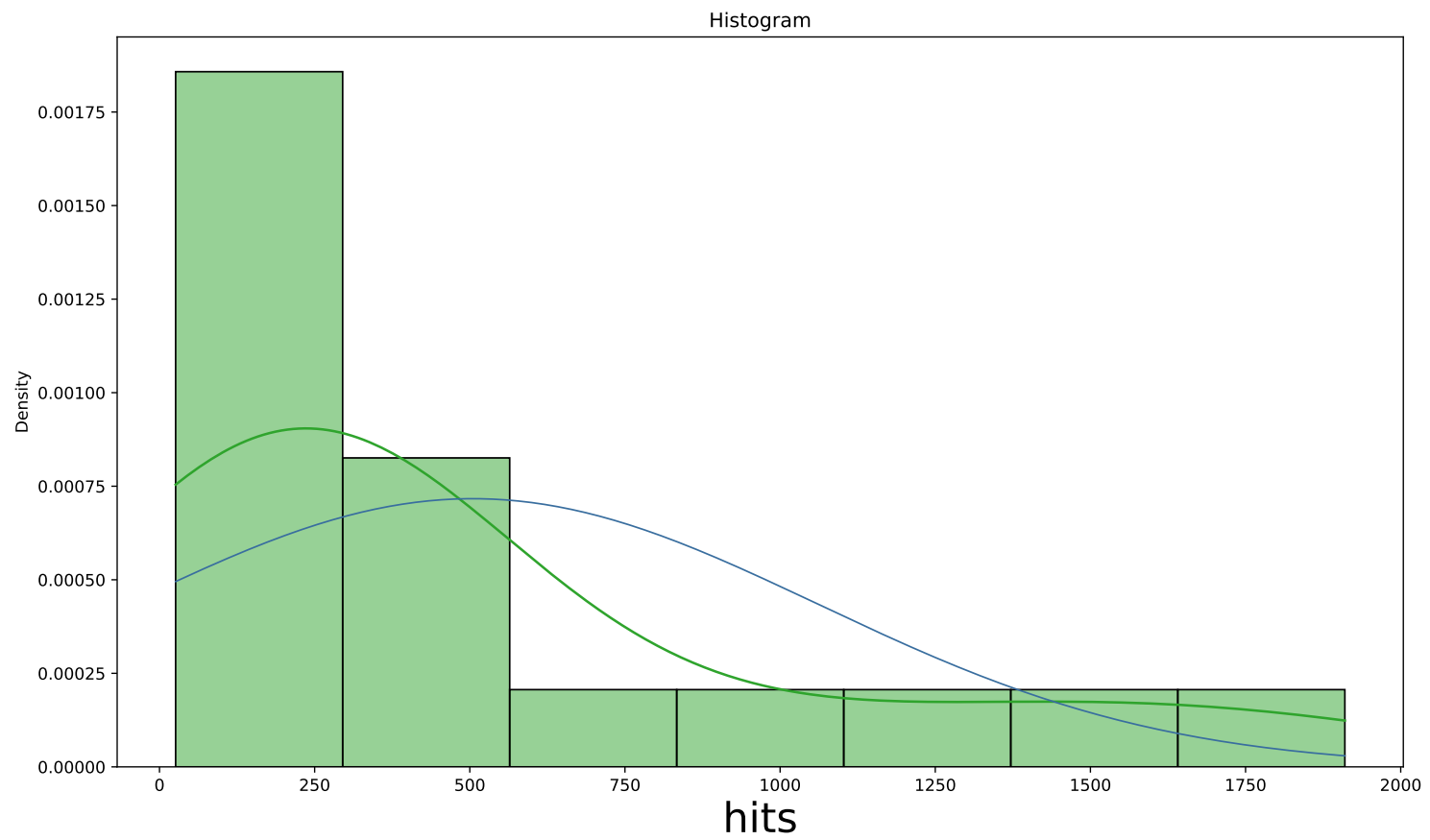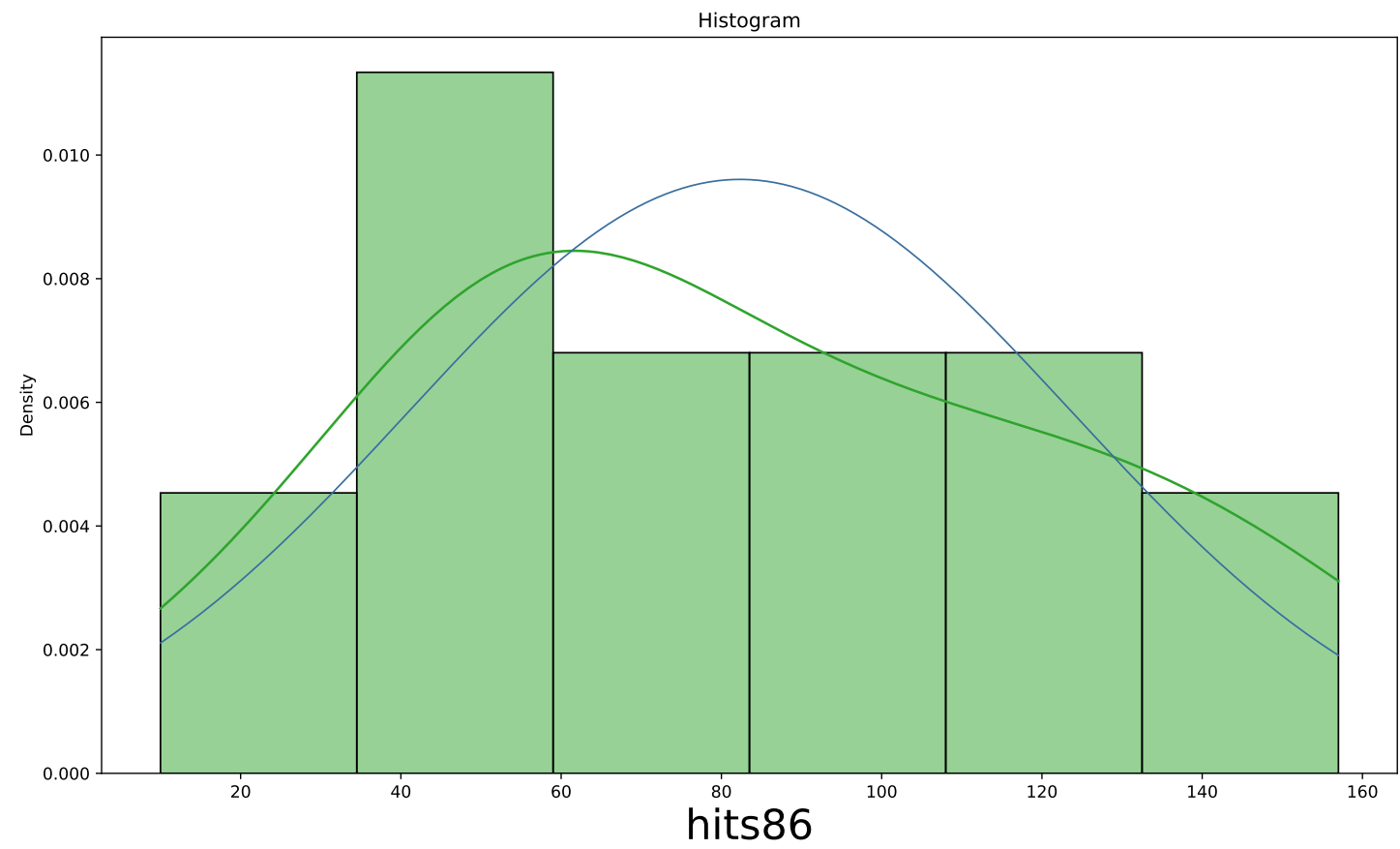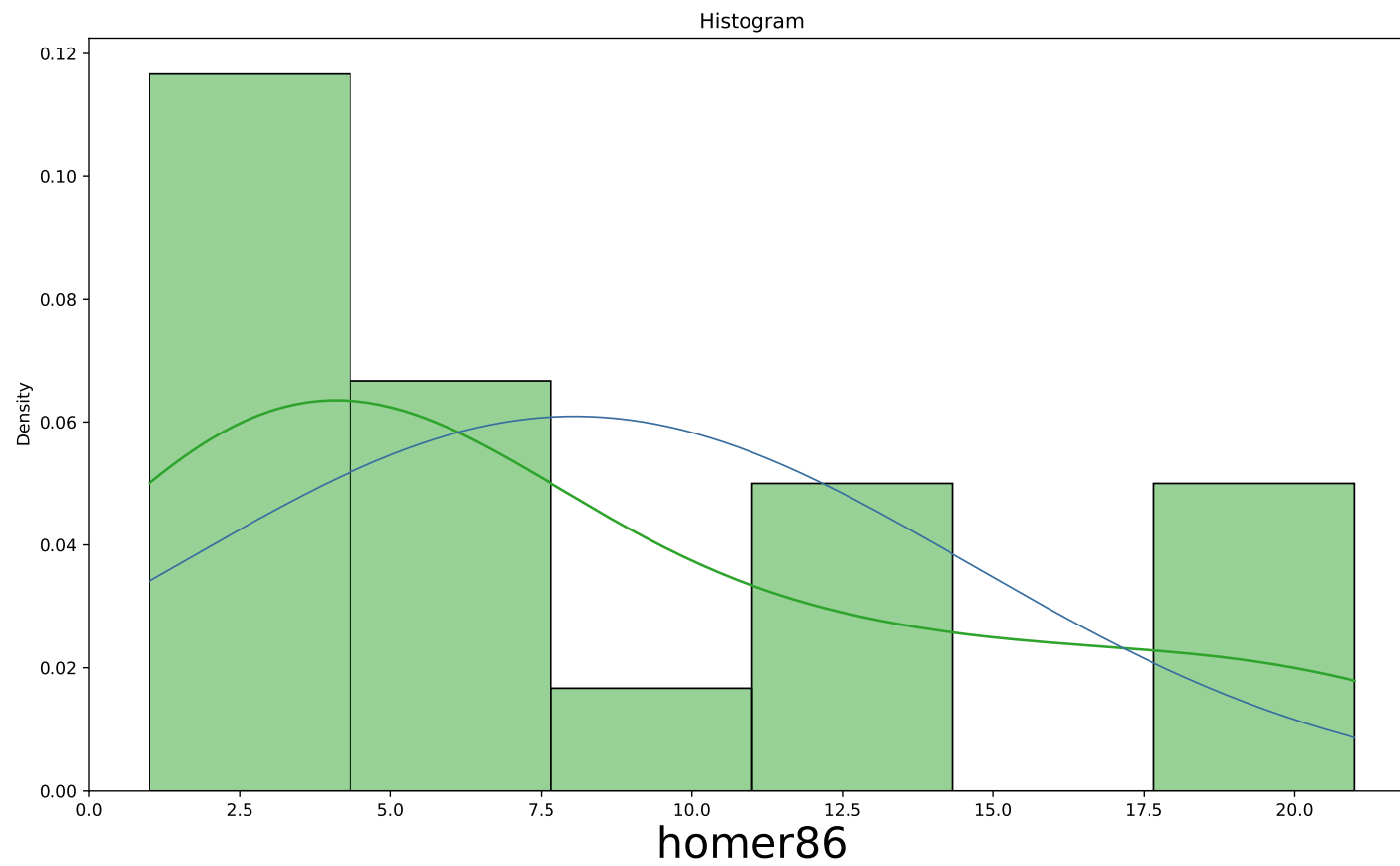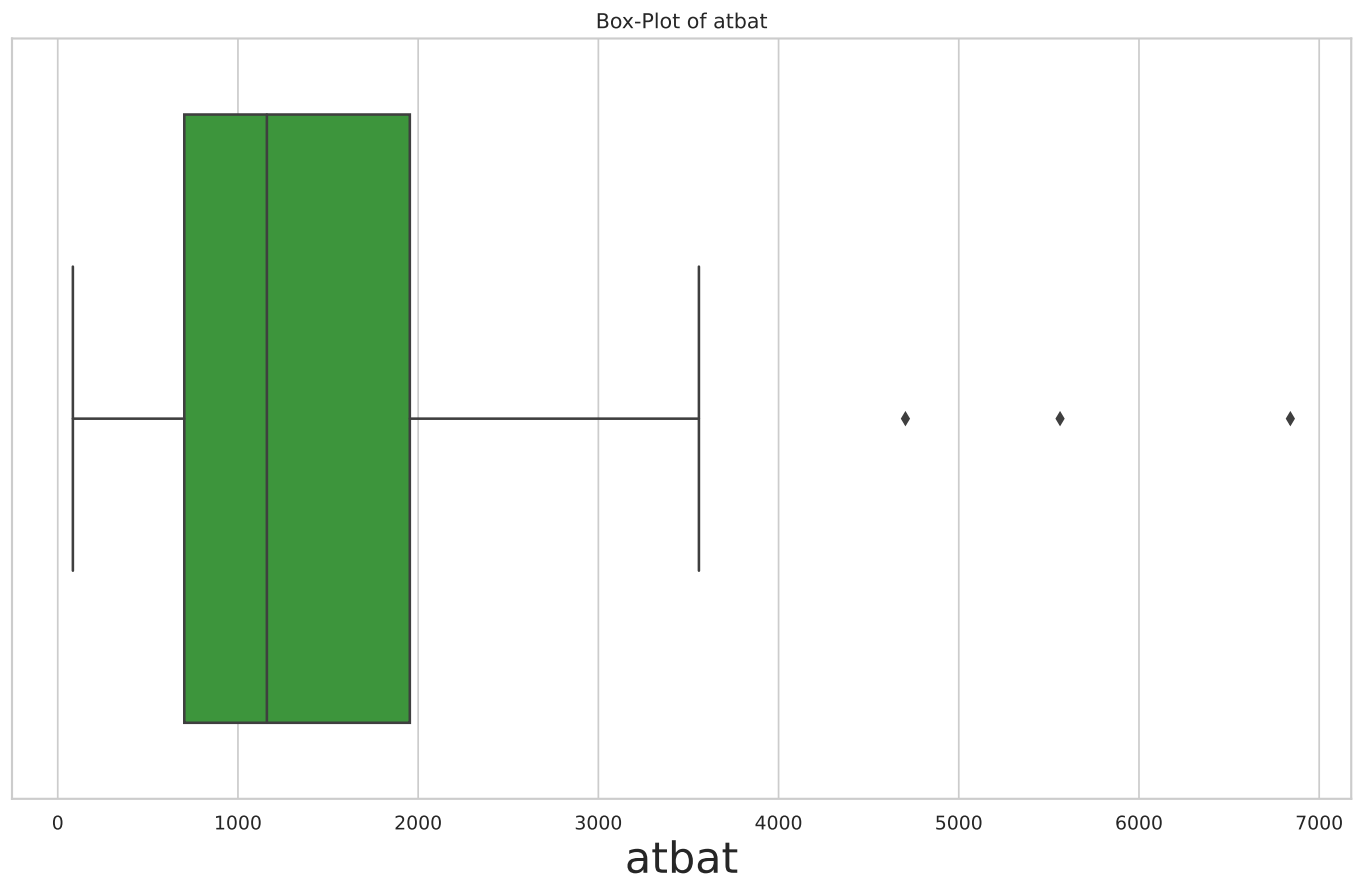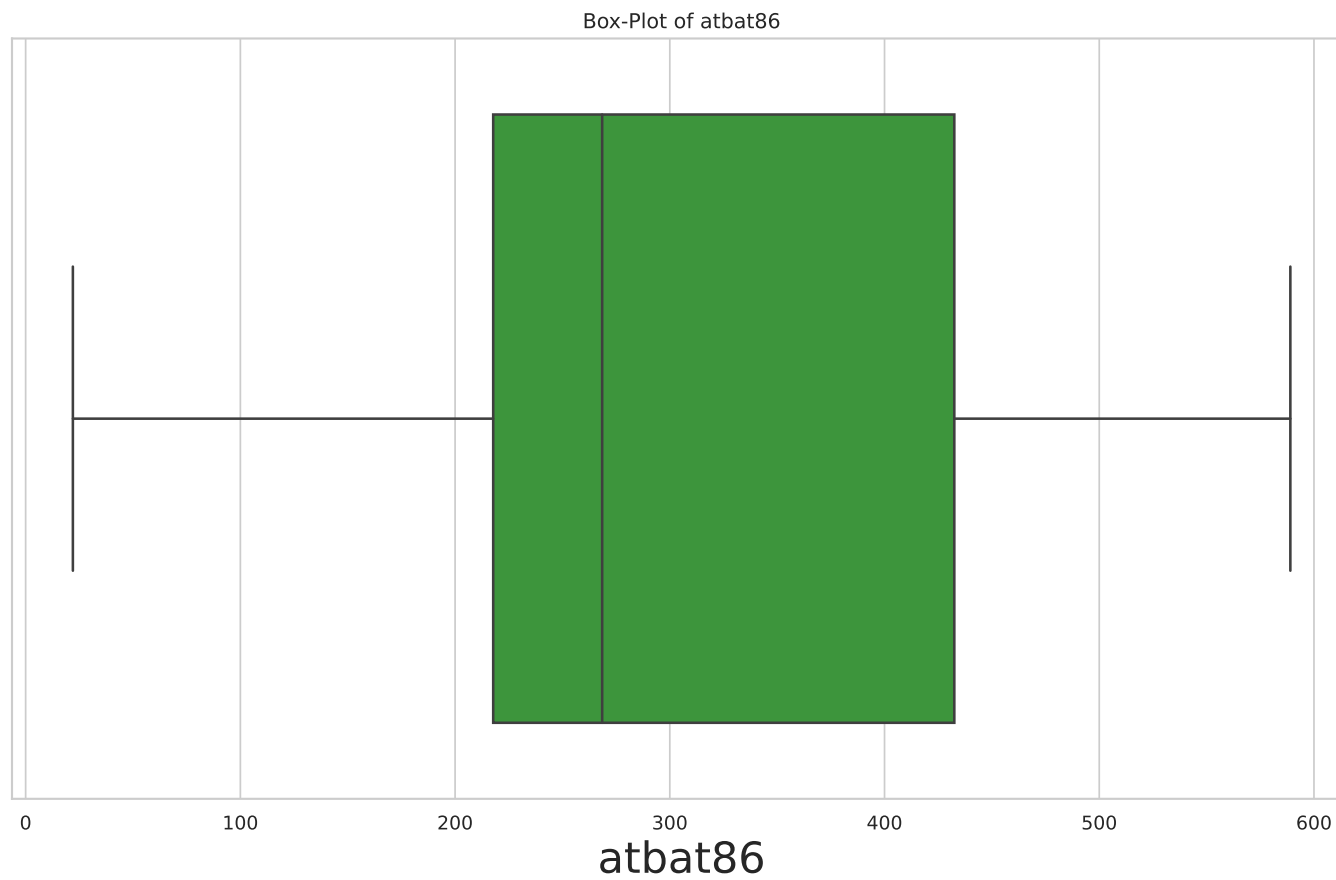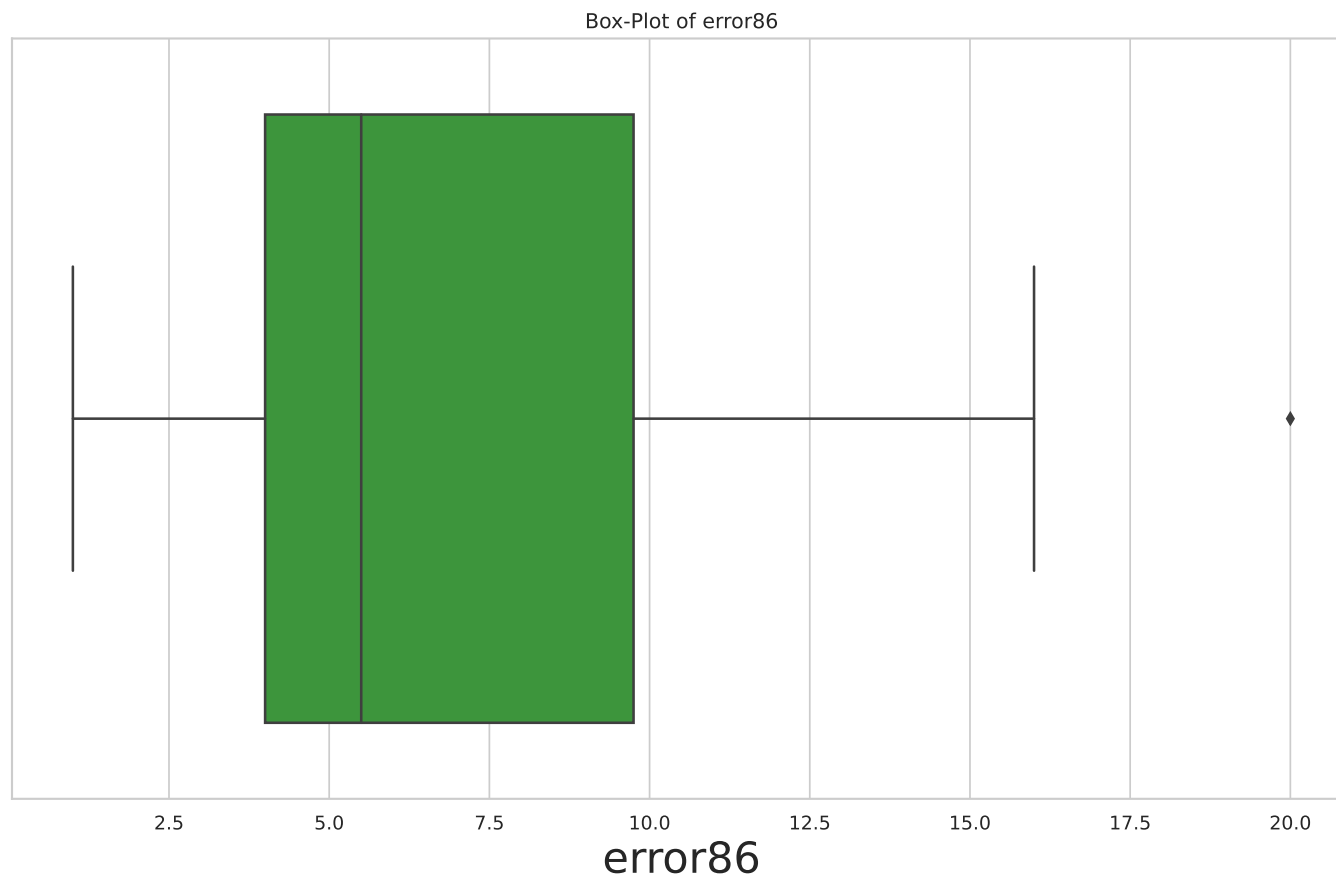
See figures on next page.

[]

Histogram

assist86

Histogram

Density

atbat

Histogram

Density

atbat86

Histogram

Density

error86

Histogram

Histogram

Density

hits86

Histogram

homer86

Histogram

Density

homeruns

Histogram

Density

outs86

Histogram

Histogram

Density

rbi86

Histogram

Histogram

runs86

Histogram

Histogram

walks

Histogram

walks86

20

## Histograms Summary

Multiple Relative Frequency Histogram in one figure. Variables are sorted alphabetically. The blue line represents the normal density approximation. The green line represents a special kernel density approximation.

**Box-Plots**

One Box-Plot per page for each variable. Variables are sorted alphabetically.

[]

Box-Plot of assist86



assist86

Box-Plot of atbat

atbat

Box-Plot of atbat86



atbat86

24

Box-Plot of error86

error86

Box-Plot of hits

hits

Box-Plot of hits86

hits86

Box-Plot of homer86

homer86

Box-Plot of homeruns

homeruns

Box-Plot of outs86



outs86

Box-Plot of rbi

rbi

Box-Plot of rbi86



rbi86

Box-Plot of runs

runs

Box-Plot of runs86



runs86

Box-Plot of V1

V1

Box-Plot of walks



walks

Box-Plot of walks86



walks86

## Box-Plots Summary

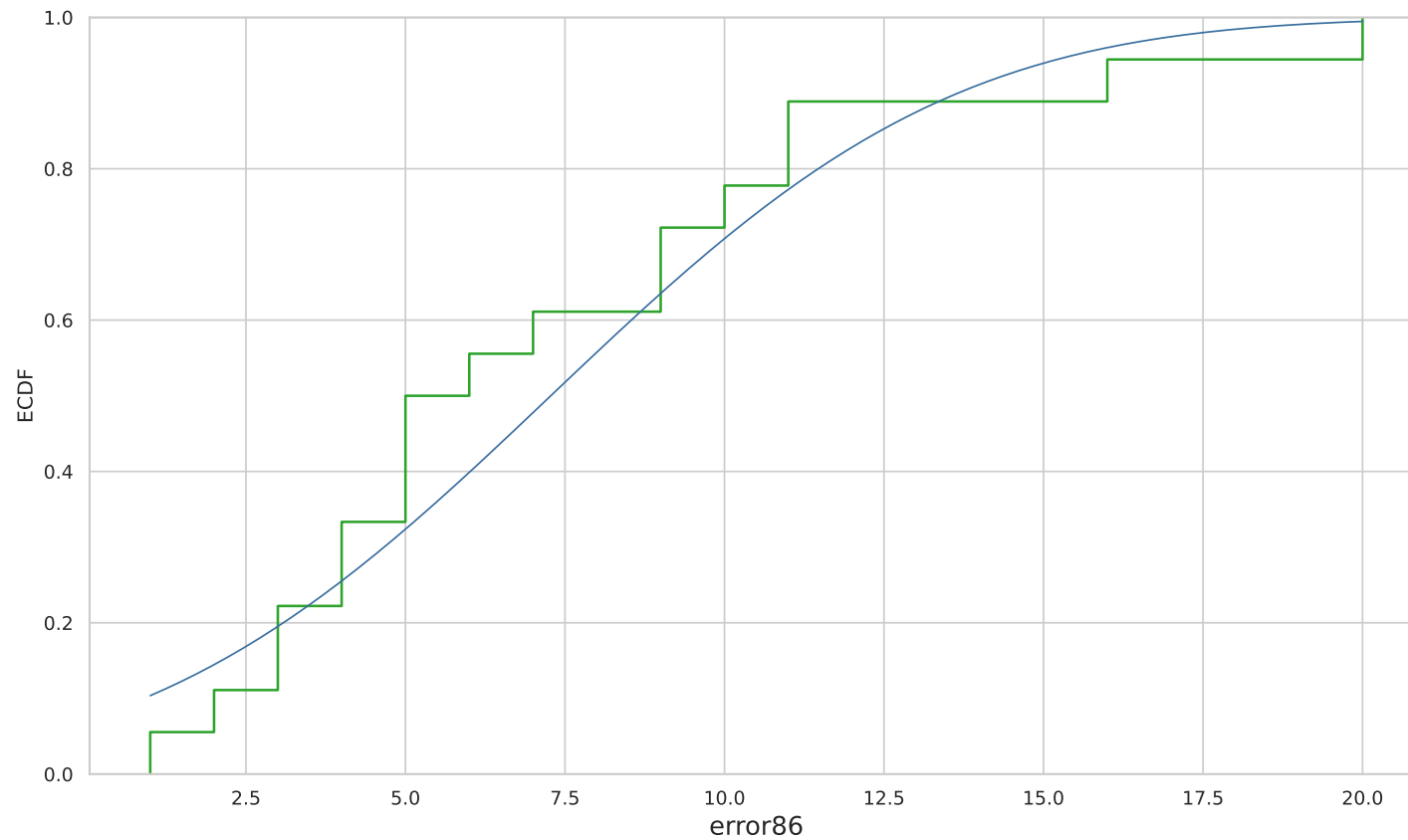Multiple Box-Plots of variables in one figure. Variables are sorted alphabetically.

## ECDF Plots

One ECDF (Empirical Cumulative Distribution Function) Plot per page for each variable. Variables are sorted alphabetically. The blue line represents the CDF of a normal distribution. If the variable is normally distributed, the blue line approximates well the ECDF.
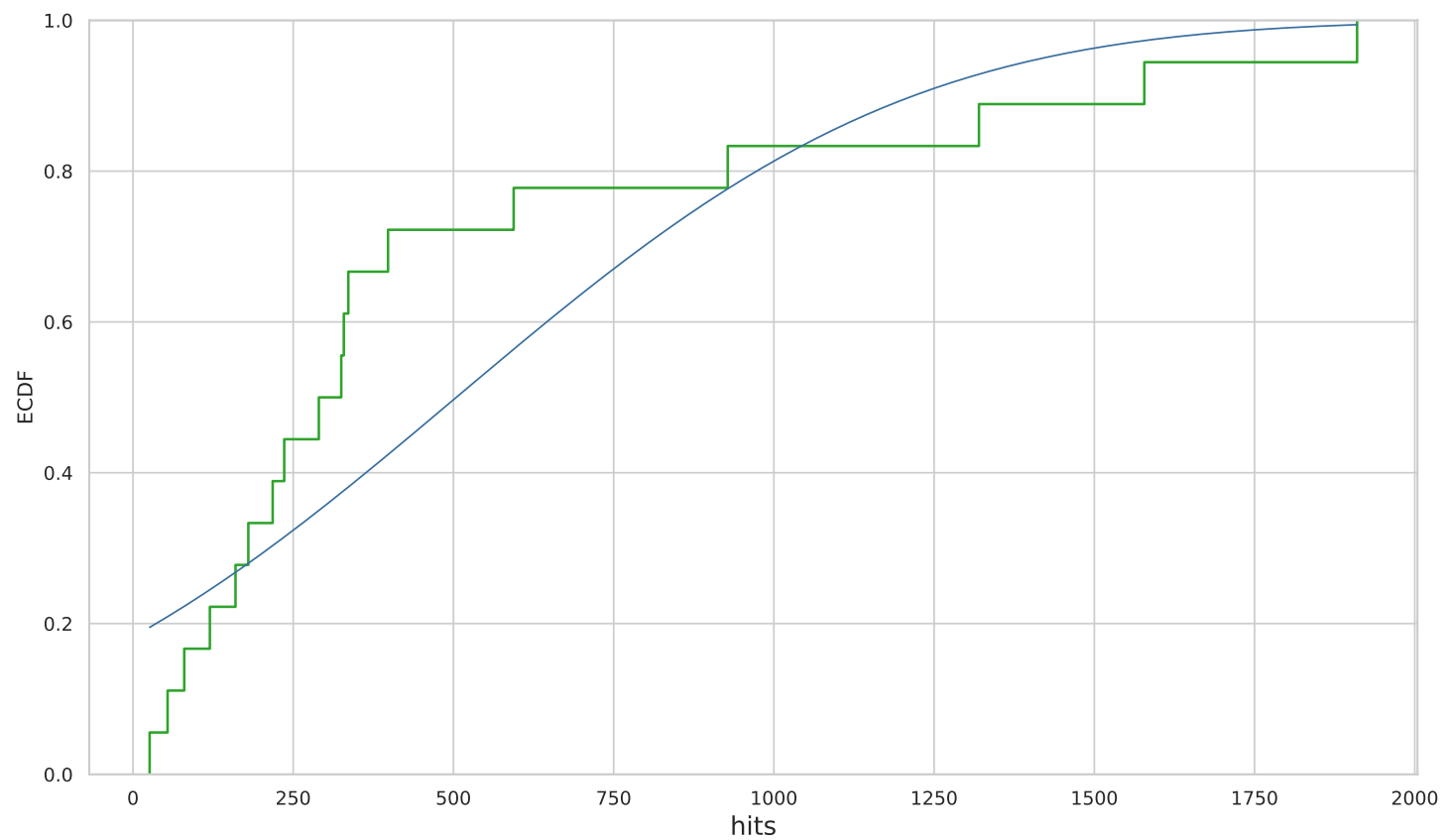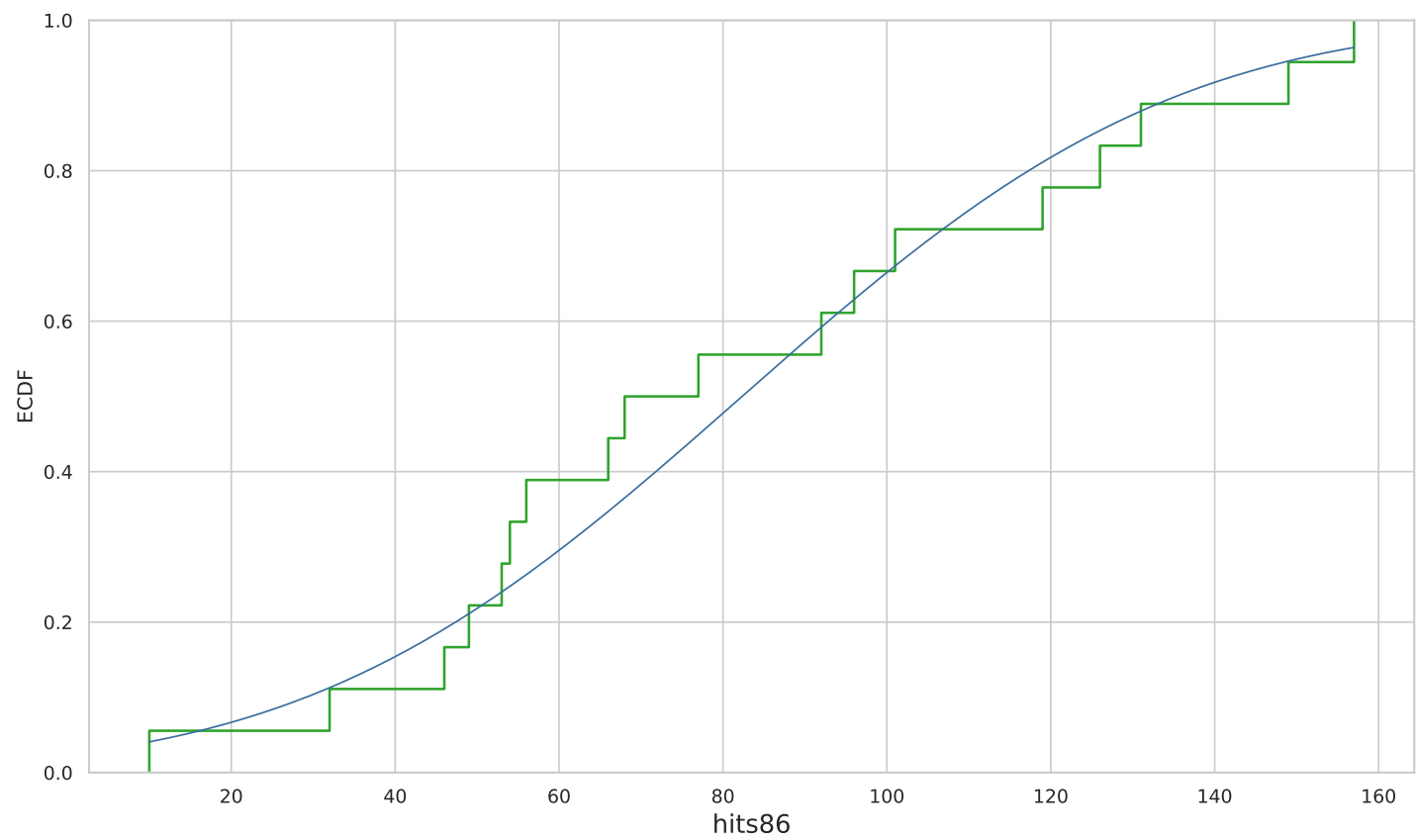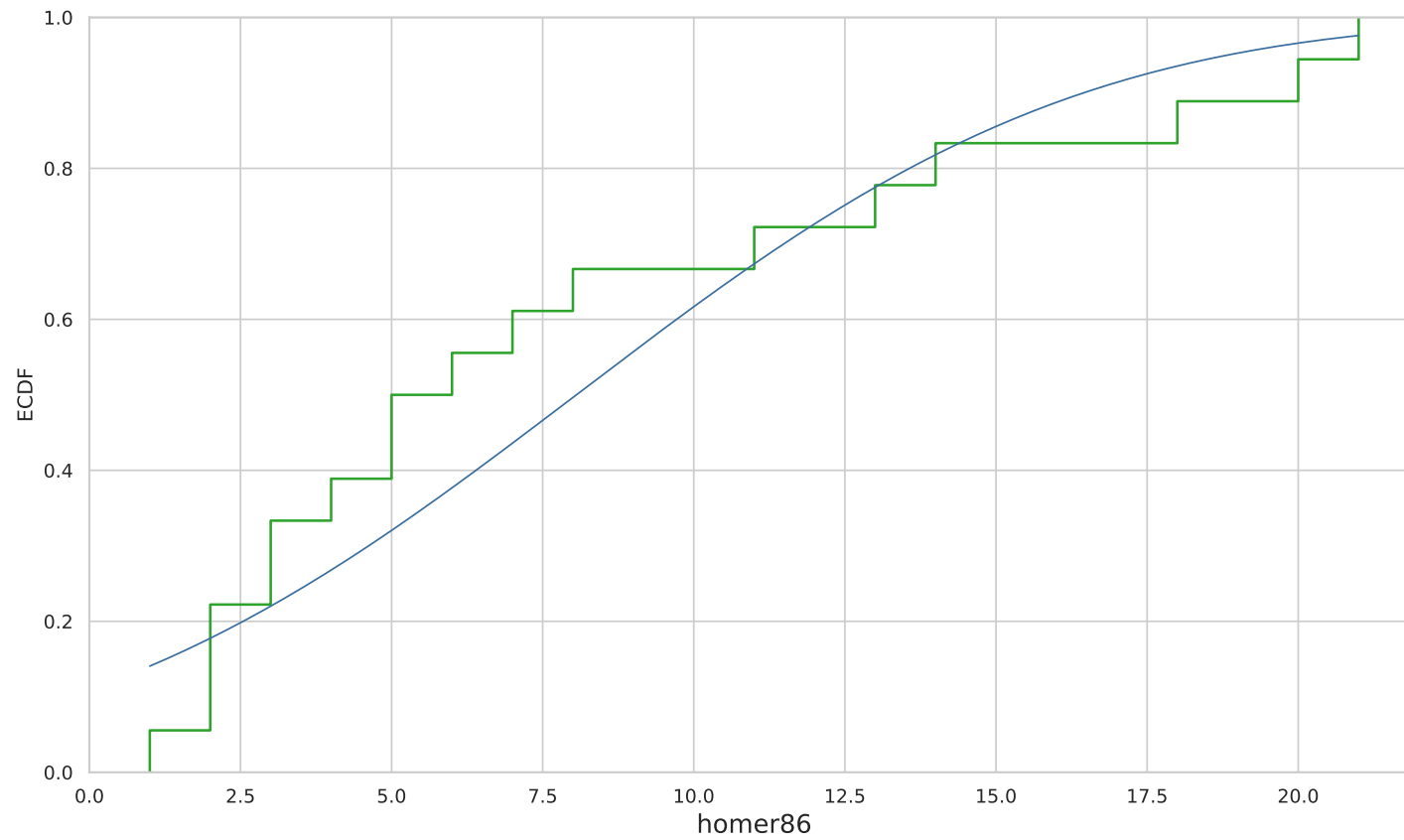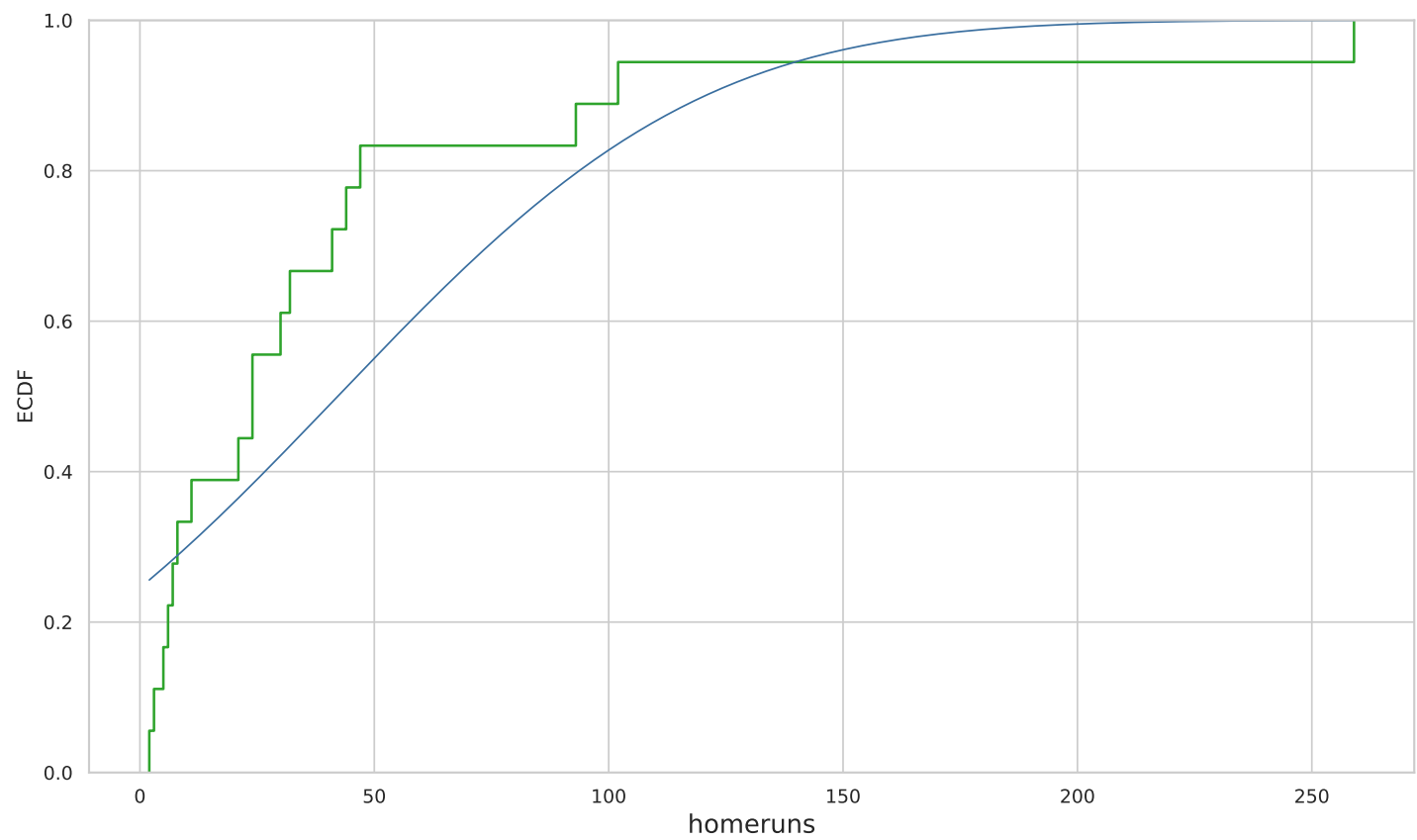
[]

## ECDF Plots Summary

Multiple ECDF Plots of variables in one figure. Variables are sorted alphabetically.

**QQ-Plots**

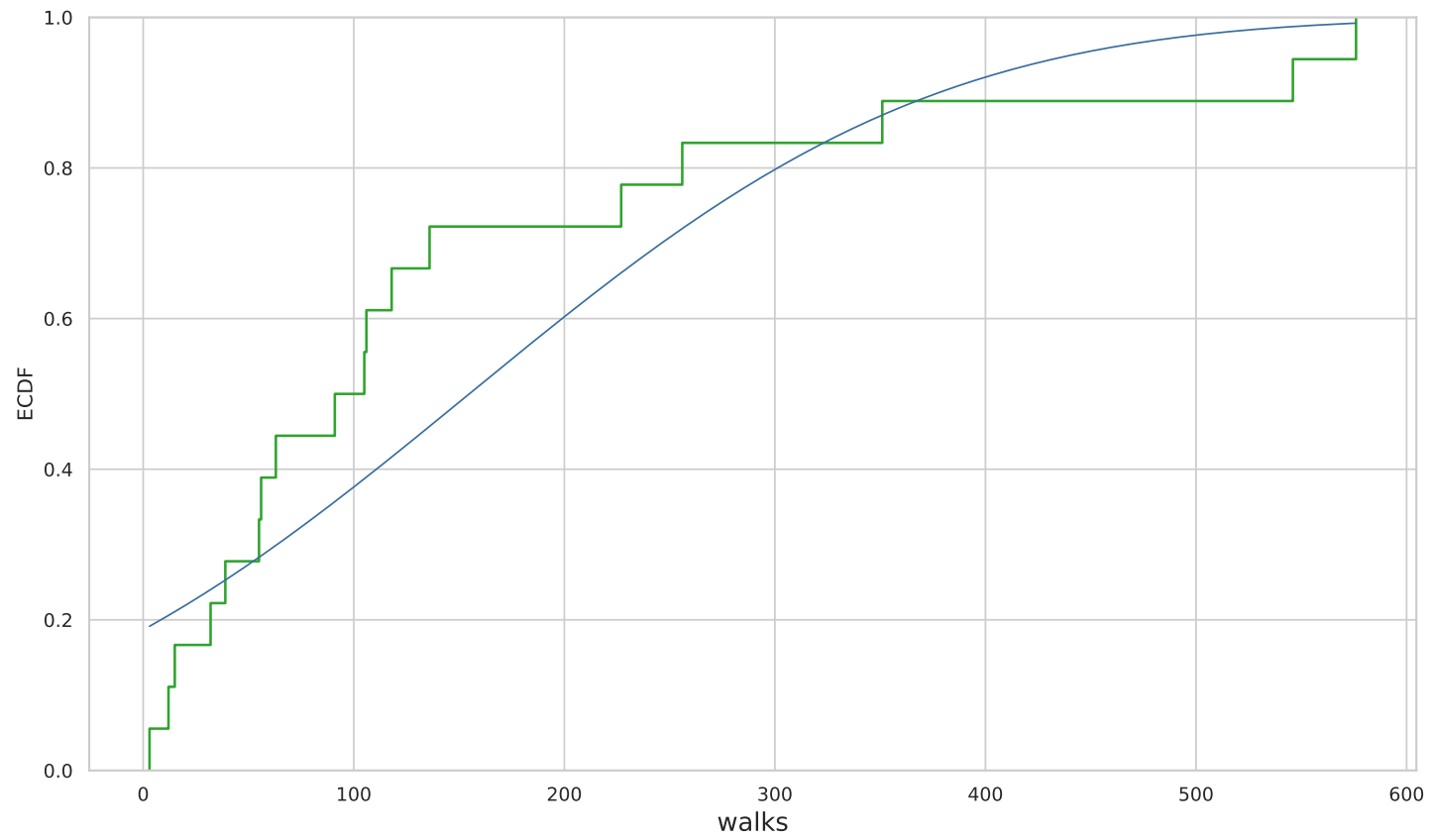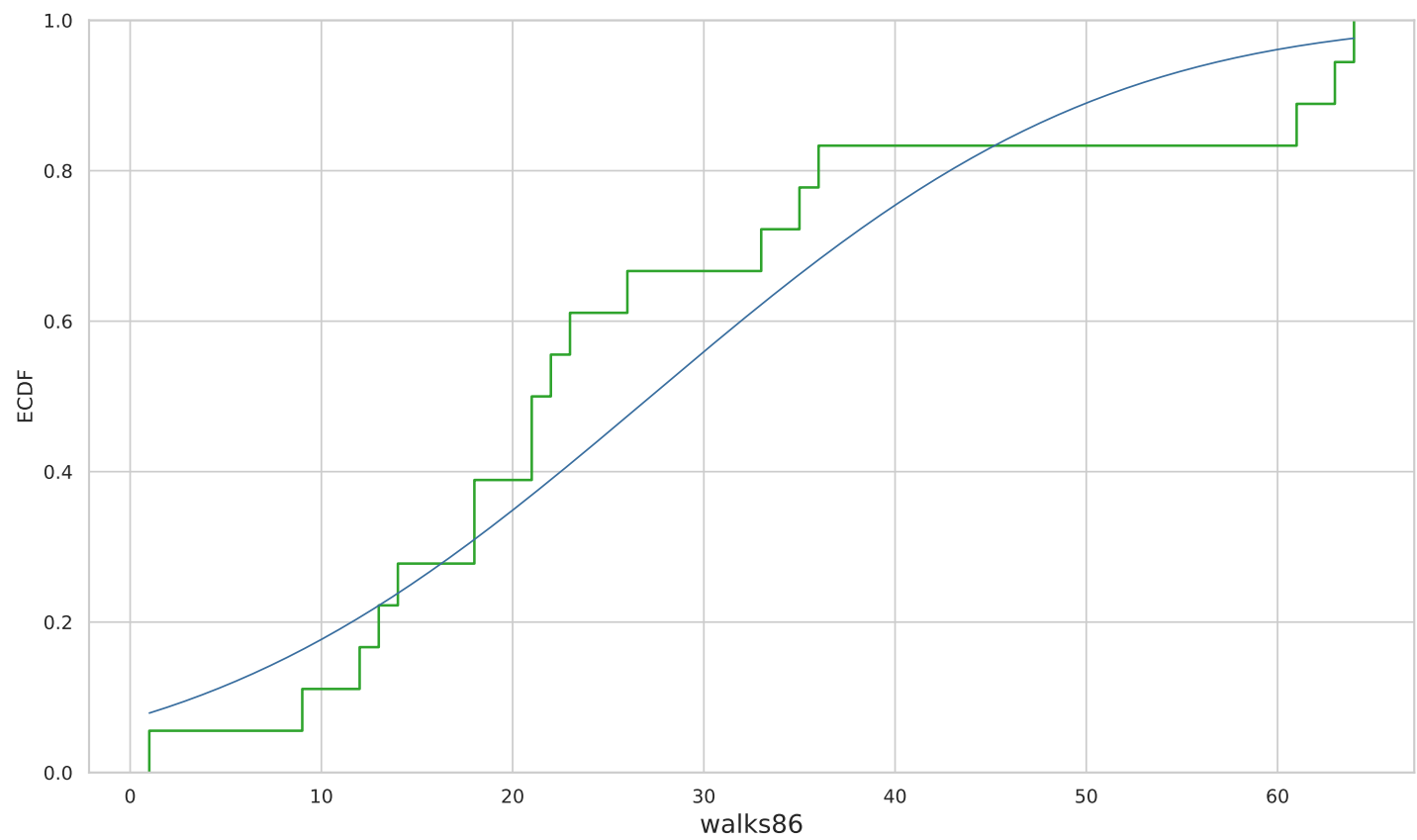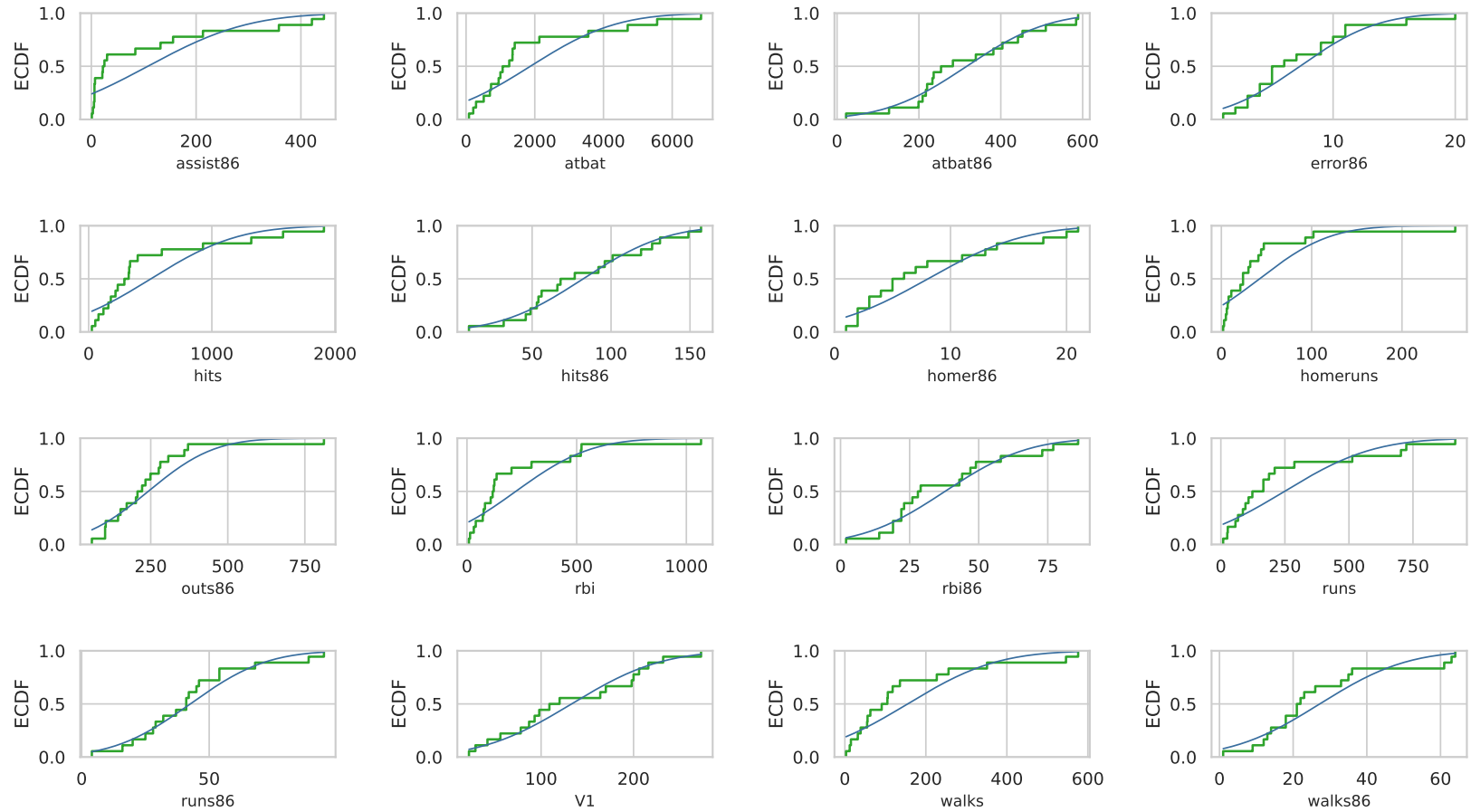One QQ-Plot per page for each variable. Variables are sorted alphabetically.

[]

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot
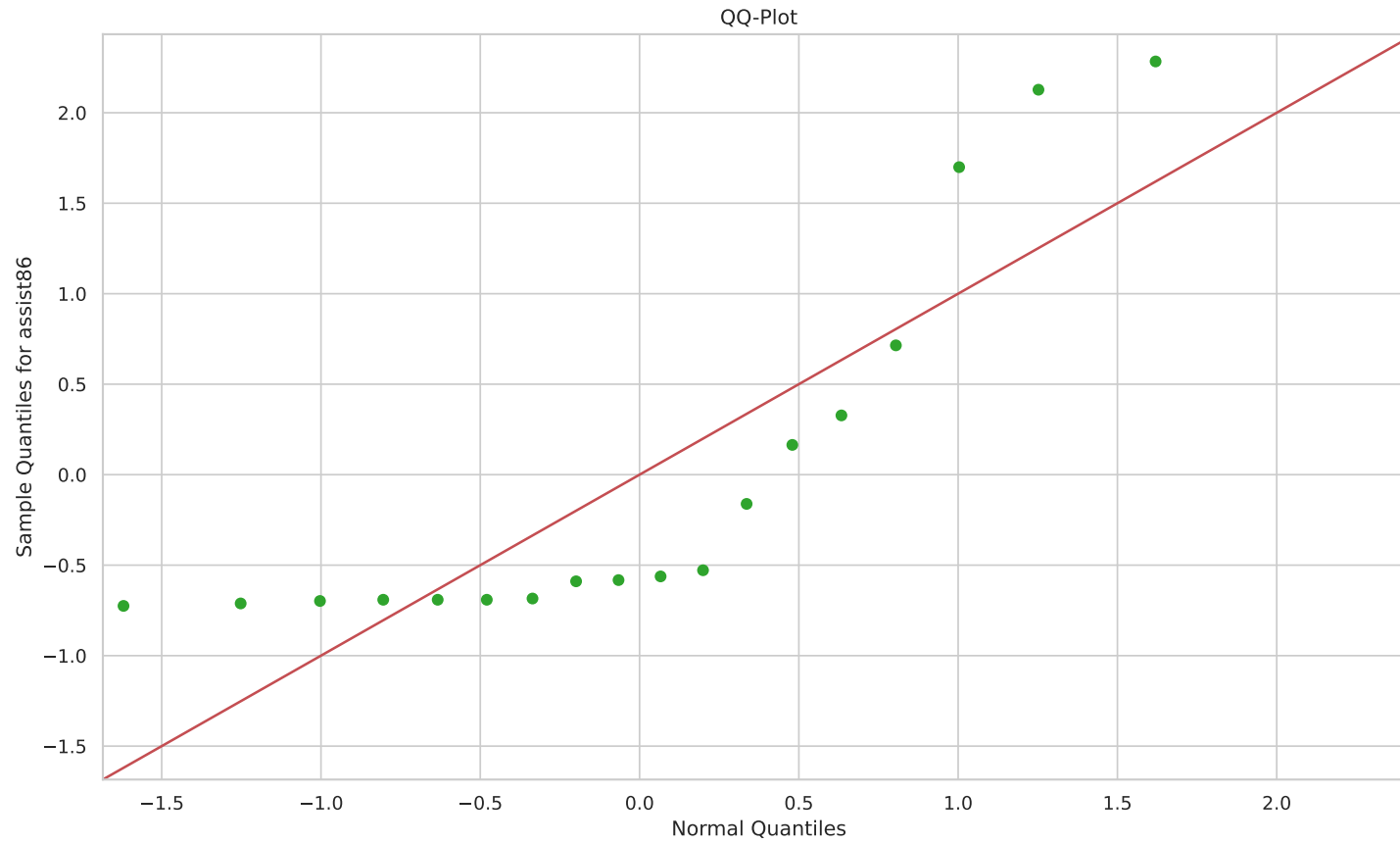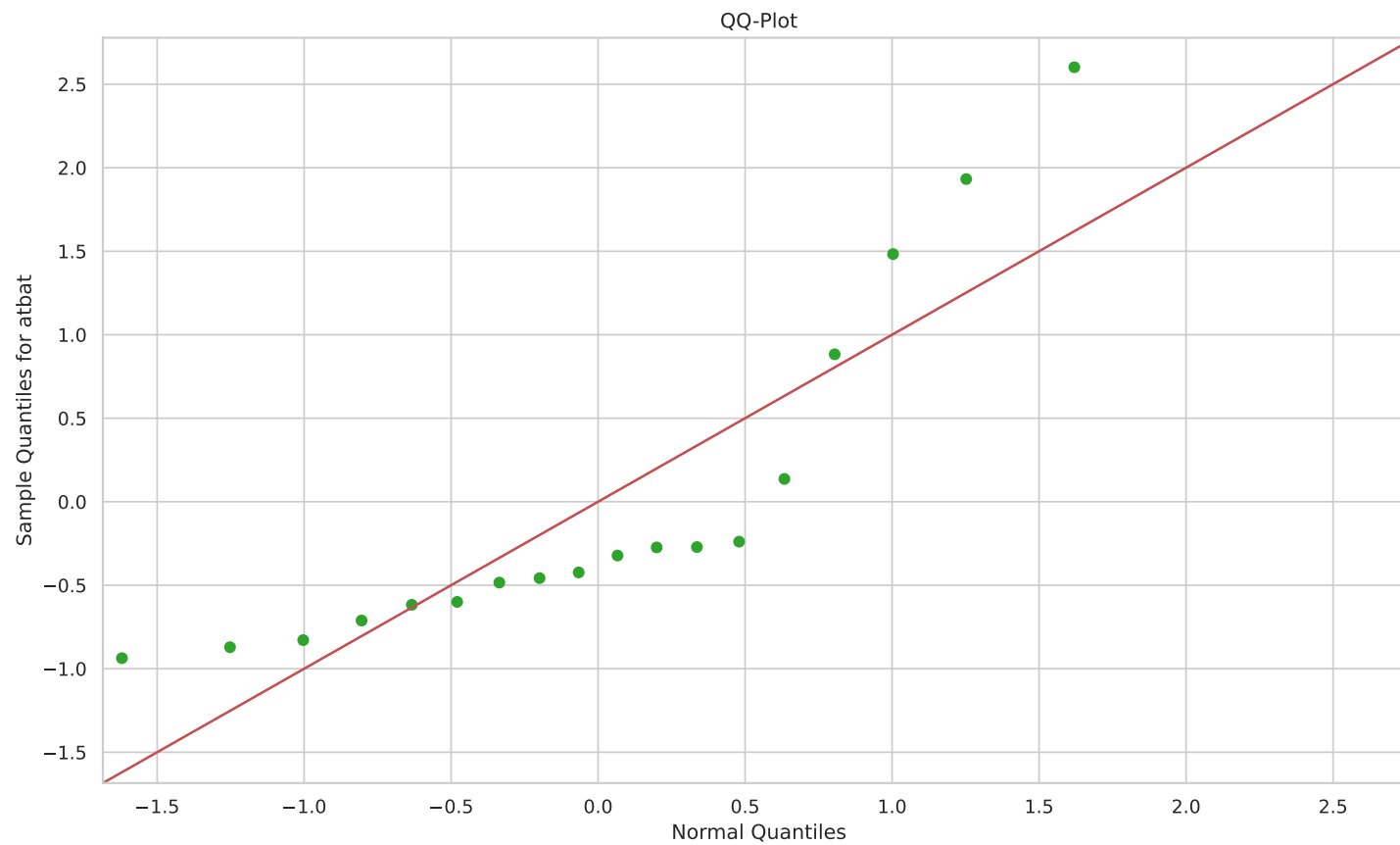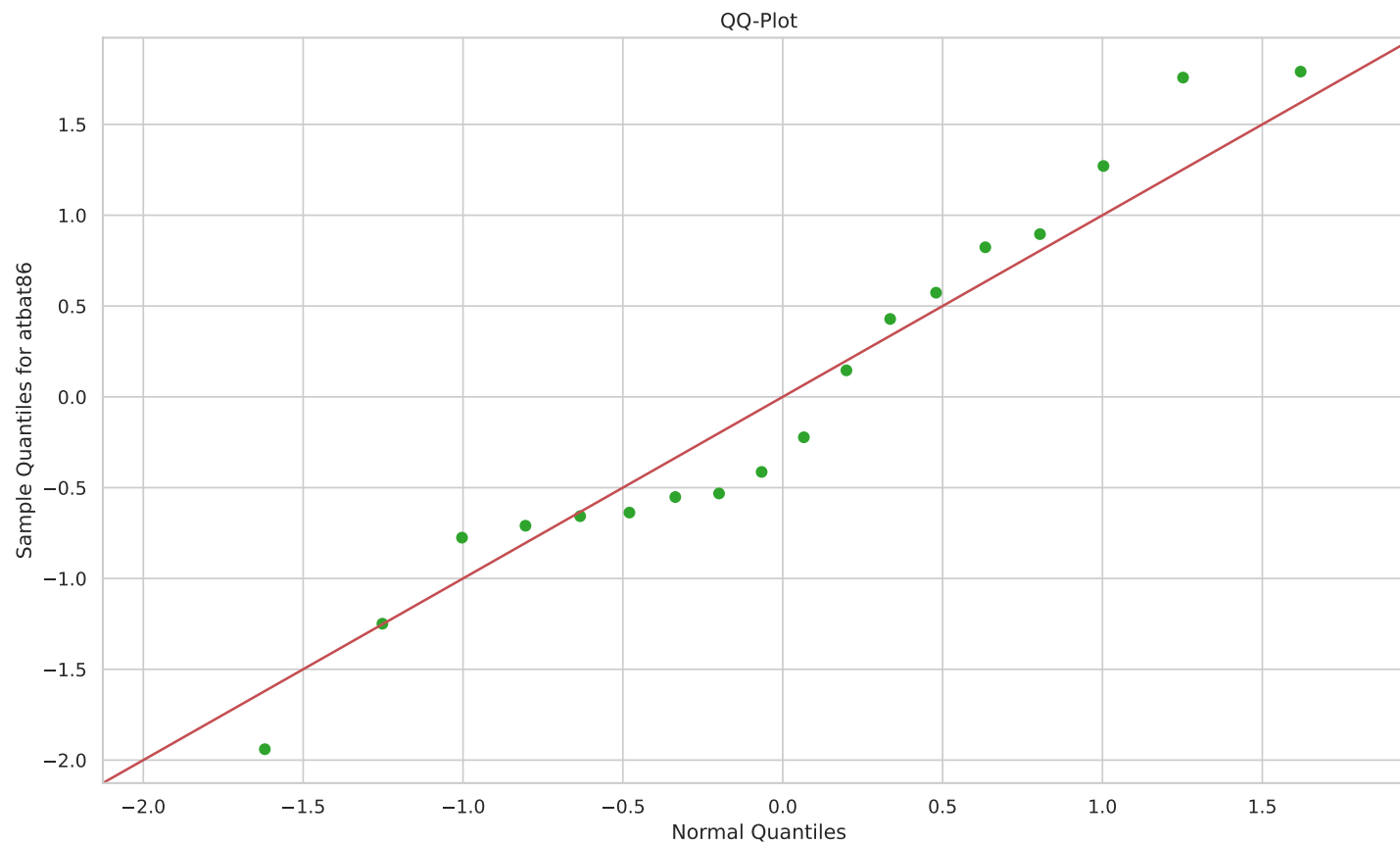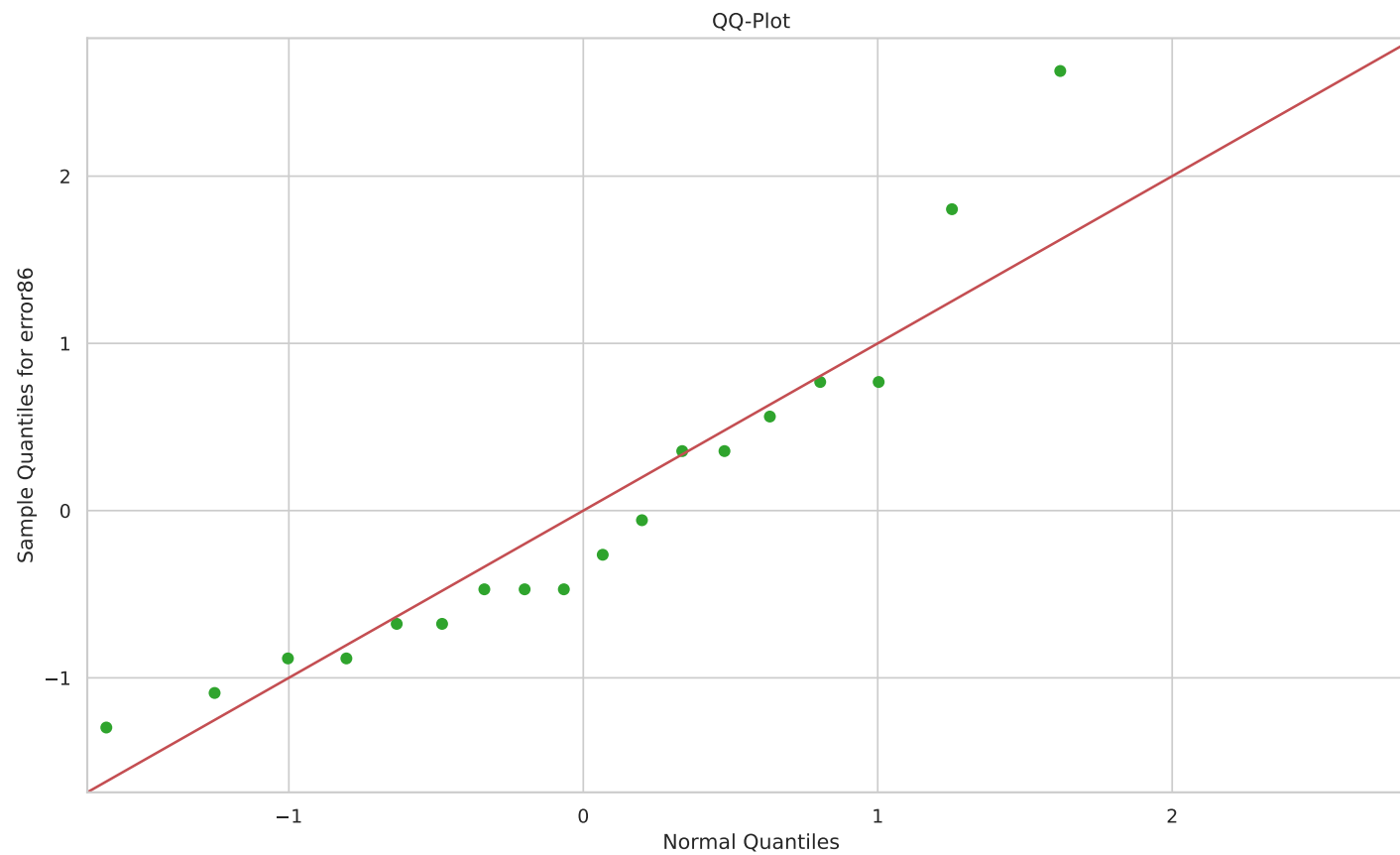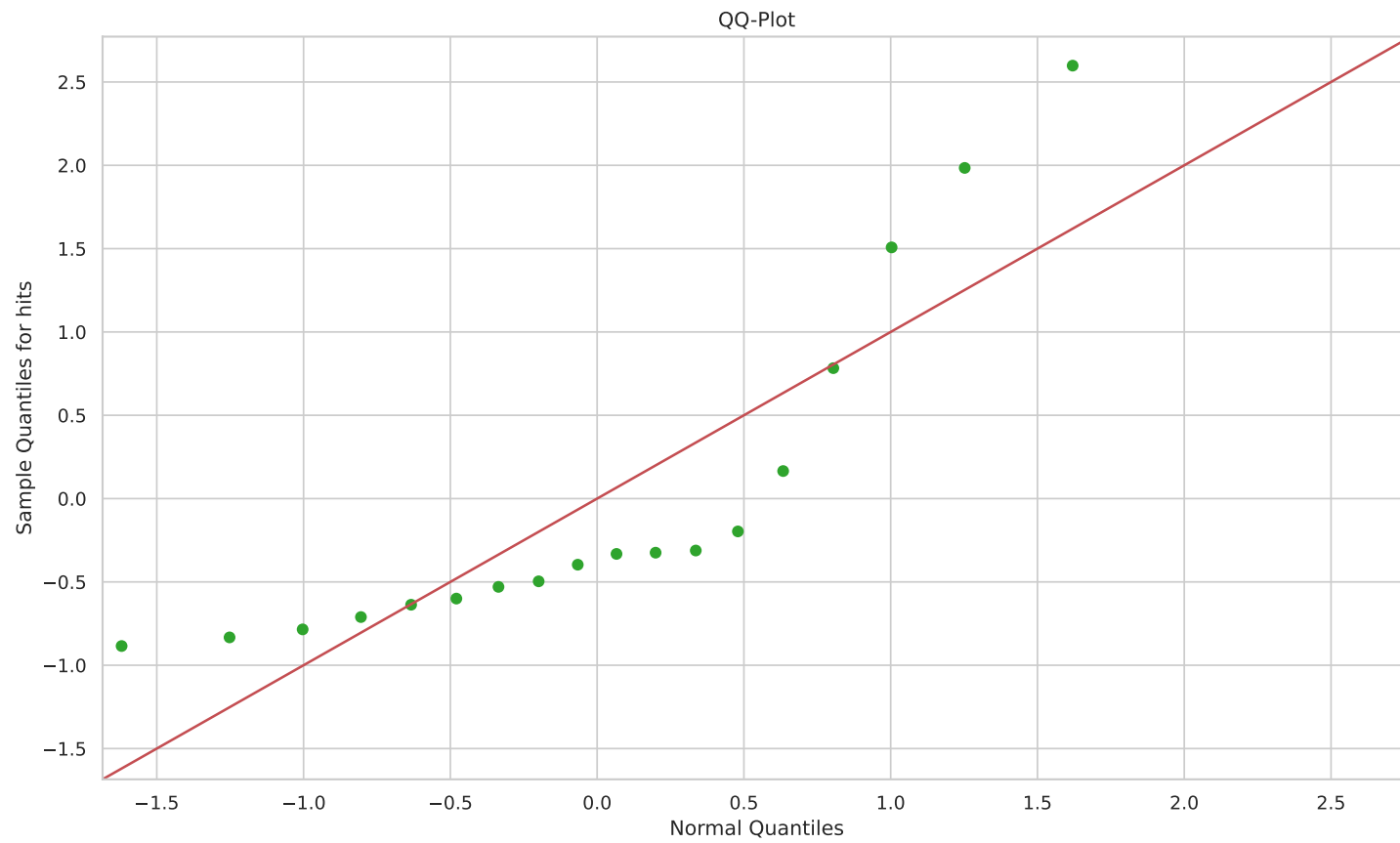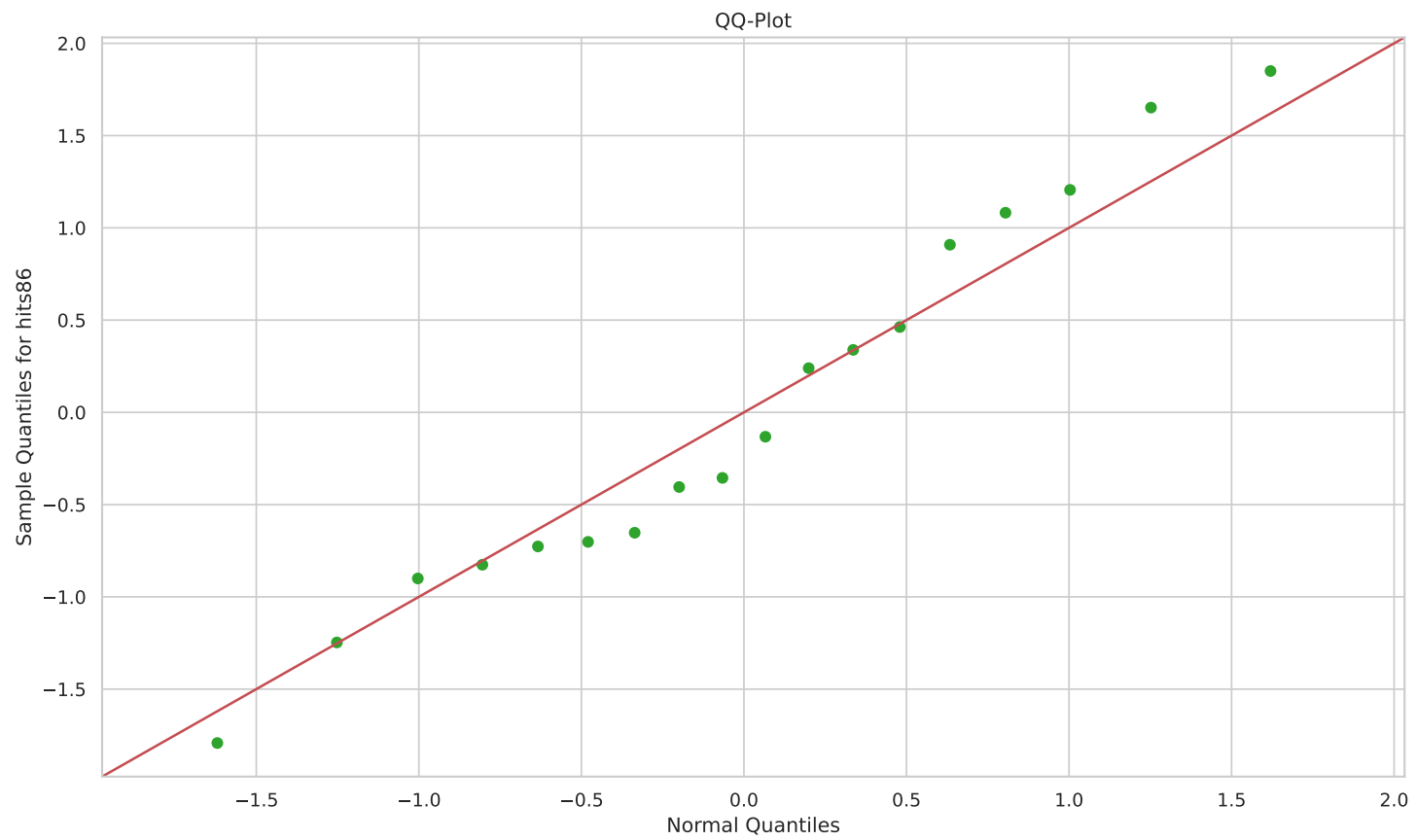
QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

QQ-Plot

## QQ-Plots Summary

Multiple QQ-Plots of variables in one figure. Variables are sorted alphabetically.

# Results for Discrete Variables

## Descriptive Statistics

### Totals

The table is sorted by the variable name. If any, N Unique contains the missing category.

|  | N Obs | N Missing | N Valid | % Complete | N Unique |
|---|---|---|---|---|---|
| div86 | 18 | 0 | 18 | 100 | 2 |
| league86 | 18 | 0 | 18 | 100 | 2 |
| name1 | 18 | 0 | 18 | 100 | 18 |
| name2 | 18 | 0 | 18 | 100 | 18 |
| posit86 | 18 | 0 | 18 | 100 | 9 |
| team86 | 18 | 0 | 18 | 100 | 14 |
| years | 18 | 0 | 18 | 100 | 11 |

## Frequencies

The table is sorted by the variable name. For each variable, a maximum of 20 unique values are considered, sorted in decreasing order of their frequency. If any, missings are counted as a category.

| Variable | Category | Frequency | Percent |
|---|---|---|---|
| div86 | W | 10 | 0.555556 |
| div86 | E | 8 | 0.444444 |
| league86 | A | 14 | 0.777778 |
| league86 | N | 4 | 0.222222 |
| name1 | Chris | 1 | 0.0555556 |
| name1 | Jose | 1 | 0.0555556 |
| name1 | Garry | 1 | 0.0555556 |
| name1 | Scott | 1 | 0.0555556 |
| name1 | John | 1 | 0.0555556 |
| name1 | Lloyd | 1 | 0.0555556 |
| name1 | Denny | 1 | 0.0555556 |
| name1 | Mel | 1 | 0.0555556 |
| name1 | Lou | 1 | 0.0555556 |
| name1 | Donnie | 1 | 0.0555556 |
| name1 | Mike | 1 | 0.0555556 |
| name1 | Bruce | 1 | 0.0555556 |
| name1 | Bob | 1 | 0.0555556 |
| name1 | George | 1 | 0.0555556 |
| name1 | Ed | 1 | 0.0555556 |
| name1 | Pat | 1 | 0.0555556 |
| name1 | Daryl | 1 | 0.0555556 |
| name1 | Bill | 1 | 0.0555556 |
| name2 | Moseby | 1 | 0.0555556 |
| name2 | Hill | 1 | 0.0555556 |
| name2 | Bochy | 1 | 0.0555556 |
| name2 | Shelby | 1 | 0.0555556 |
| name2 | Templeton | 1 | 0.0555556 |
| name2 | Walling | 1 | 0.0555556 |
| name2 | Boone | 1 | 0.0555556 |
| name2 | Kingery | 1 | 0.0555556 |
| name2 | Boston | 1 | 0.0555556 |
| name2 | Hall | 1 | 0.0555556 |
| name2 | Romero | 1 | 0.0555556 |
| name2 | Bando | 1 | 0.0555556 |
| name2 | Whitaker | 1 | 0.0555556 |
| name2 | Bradley | 1 | 0.0555556 |
| name2 | Uribe | 1 | 0.0555556 |

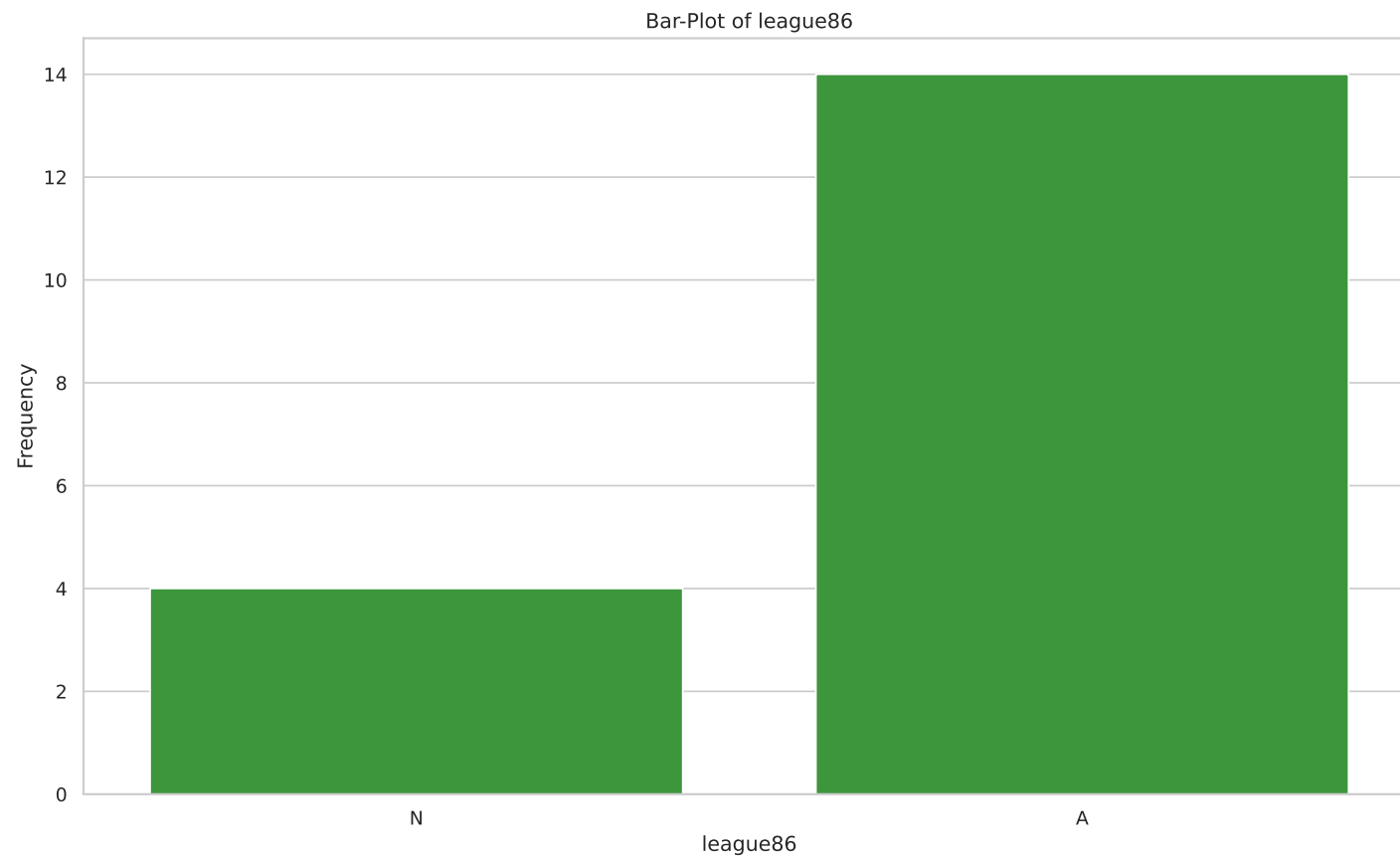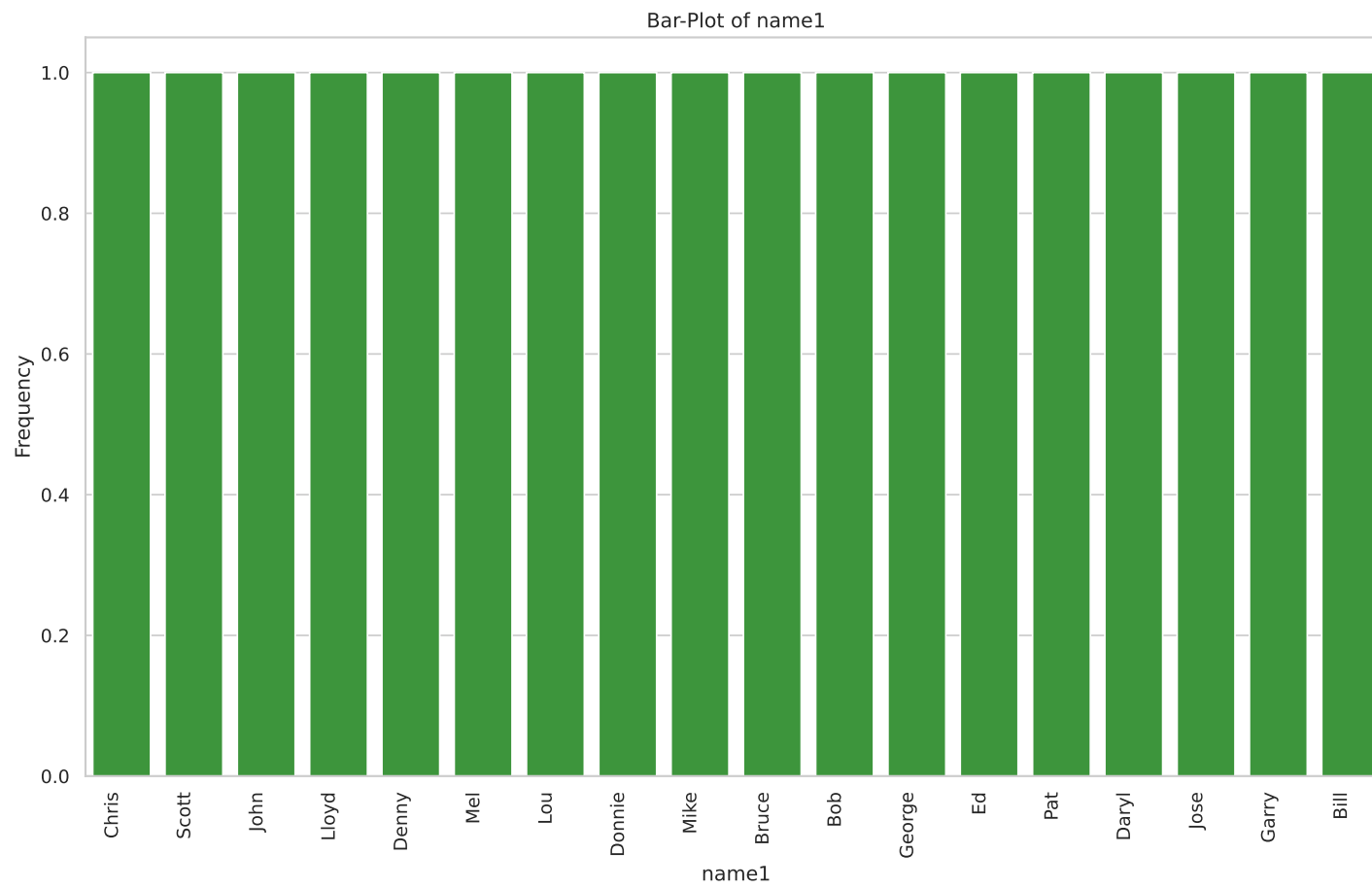| Variable | Category | Frequency | Percent |
|---|---|---|---|
| name2 | Schroeder | 1 | 0.0555556 |
| name2 | Hendrick | 1 | 0.0555556 |
| name2 | Sheridan | 1 | 0.0555556 |
| posit86 | C | 4 | 0.222222 |
| posit86 | OF | 4 | 0.222222 |
| posit86 | SS | 3 | 0.166667 |
| posit86 | CF | 2 | 0.111111 |
| posit86 | 2B | 1 | 0.0555556 |
| posit86 | 3B | 1 | 0.0555556 |
| posit86 | LF | 1 | 0.0555556 |
| posit86 | UT | 1 | 0.0555556 |
| posit86 | 23 | 1 | 0.0555556 |
| team86 | Cal | 2 | 0.111111 |
| team86 | SD | 2 | 0.111111 |
| team86 | Det | 2 | 0.111111 |
| team86 | Cle | 2 | 0.111111 |
| team86 | Chi | 1 | 0.0555556 |
| team86 | Bal | 1 | 0.0555556 |
| team86 | Bos | 1 | 0.0555556 |
| team86 | SF | 1 | 0.0555556 |
| team86 | Sea | 1 | 0.0555556 |
| team86 | Tor | 1 | 0.0555556 |
| team86 | Hou | 1 | 0.0555556 |
| team86 | Oak | 1 | 0.0555556 |
| team86 | Mil | 1 | 0.0555556 |
| team86 | KC | 1 | 0.0555556 |
| years | 6 | 4 | 0.222222 |
| years | 3 | 3 | 0.166667 |
| years | 8 | 2 | 0.111111 |
| years | 4 | 2 | 0.111111 |
| years | 10 | 1 | 0.0555556 |
| years | 7 | 1 | 0.0555556 |
| years | 1 | 1 | 0.0555556 |
| years | 16 | 1 | 0.0555556 |
| years | 11 | 1 | 0.0555556 |
| years | 5 | 1 | 0.0555556 |
| years | 12 | 1 | 0.0555556 |

# Graphics

## Bar-Plots

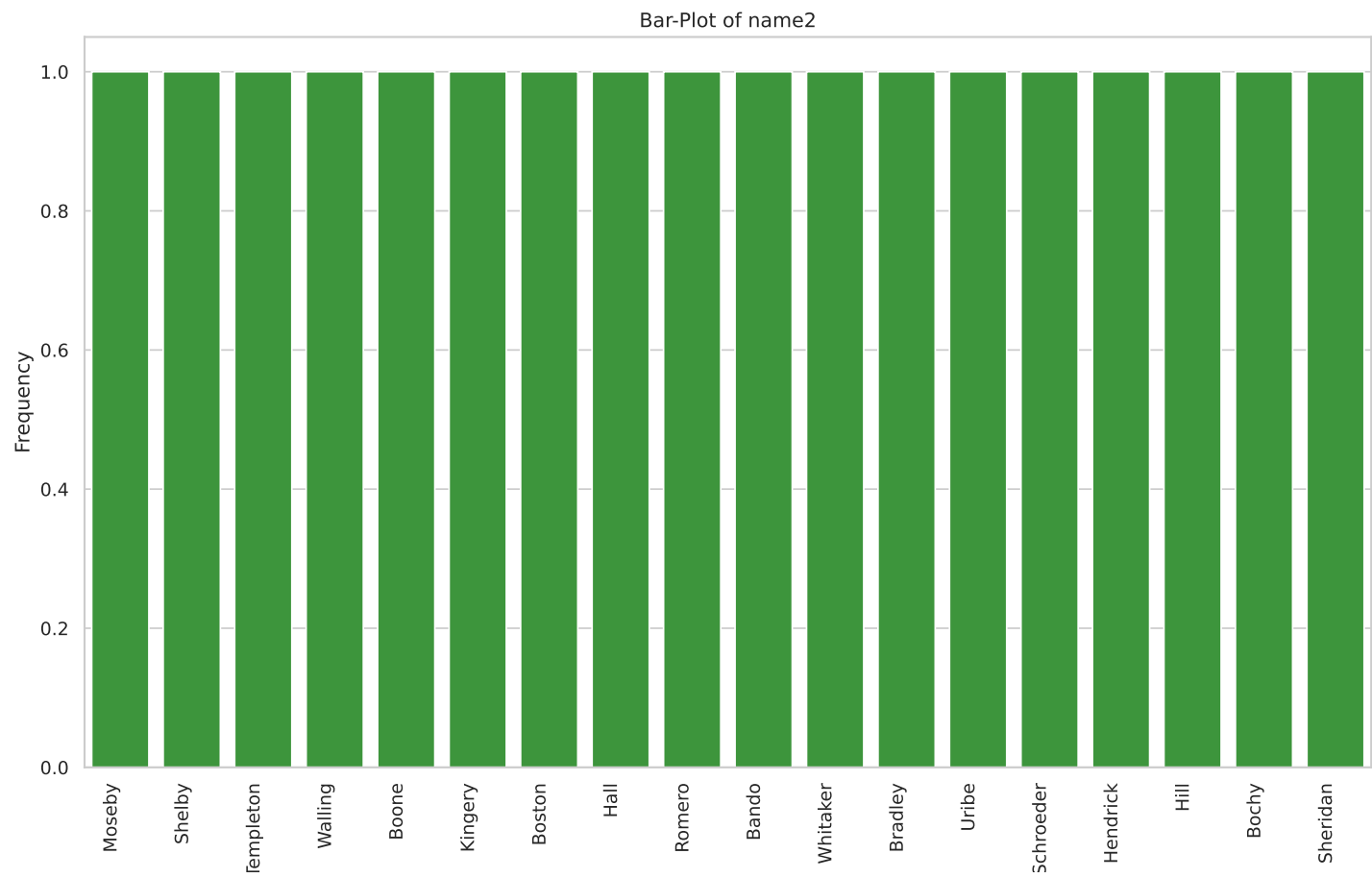One Bar-Plot per page for each variable. Variables are sorted alphabetically. No labels for variables with more than 40 categories.
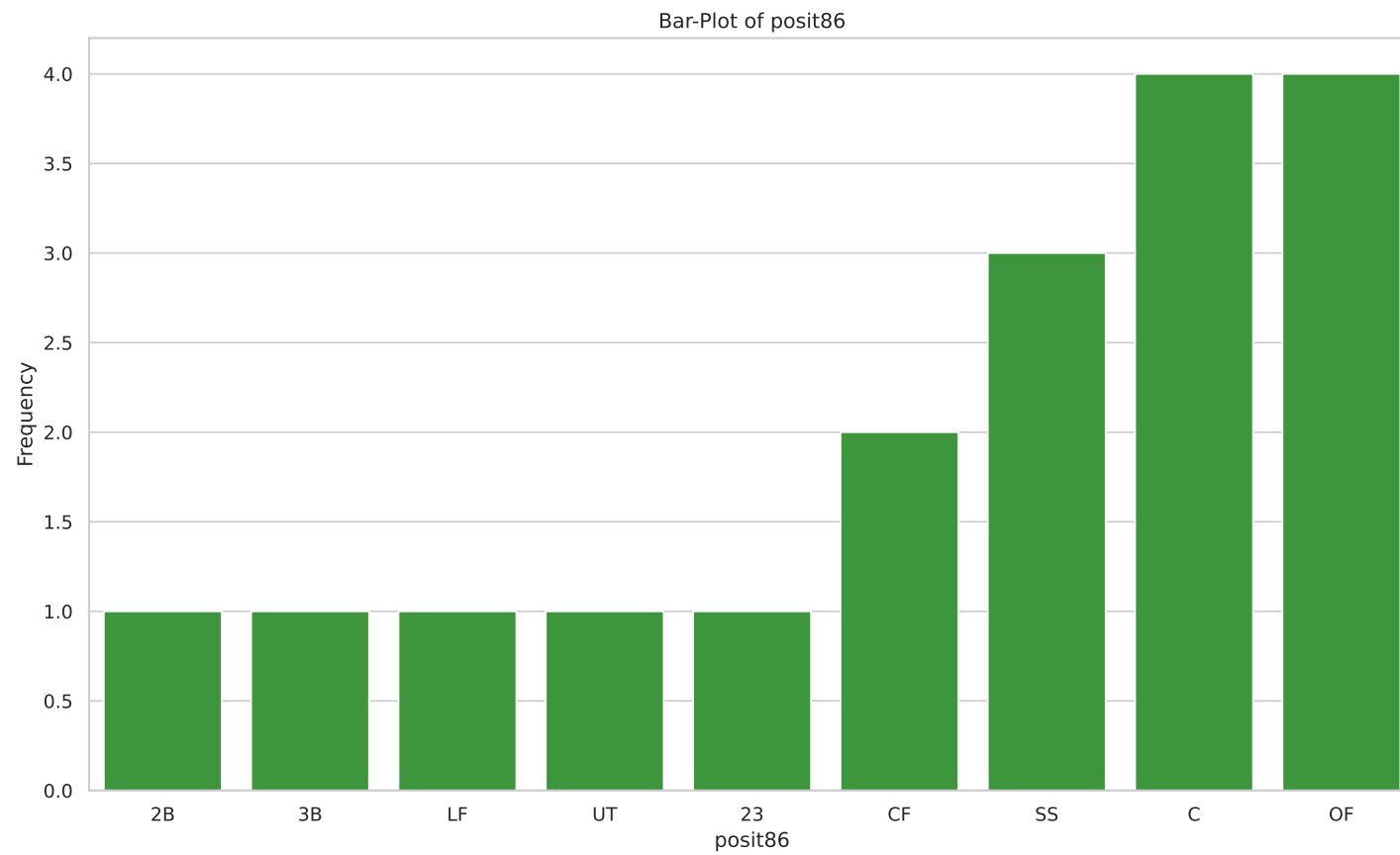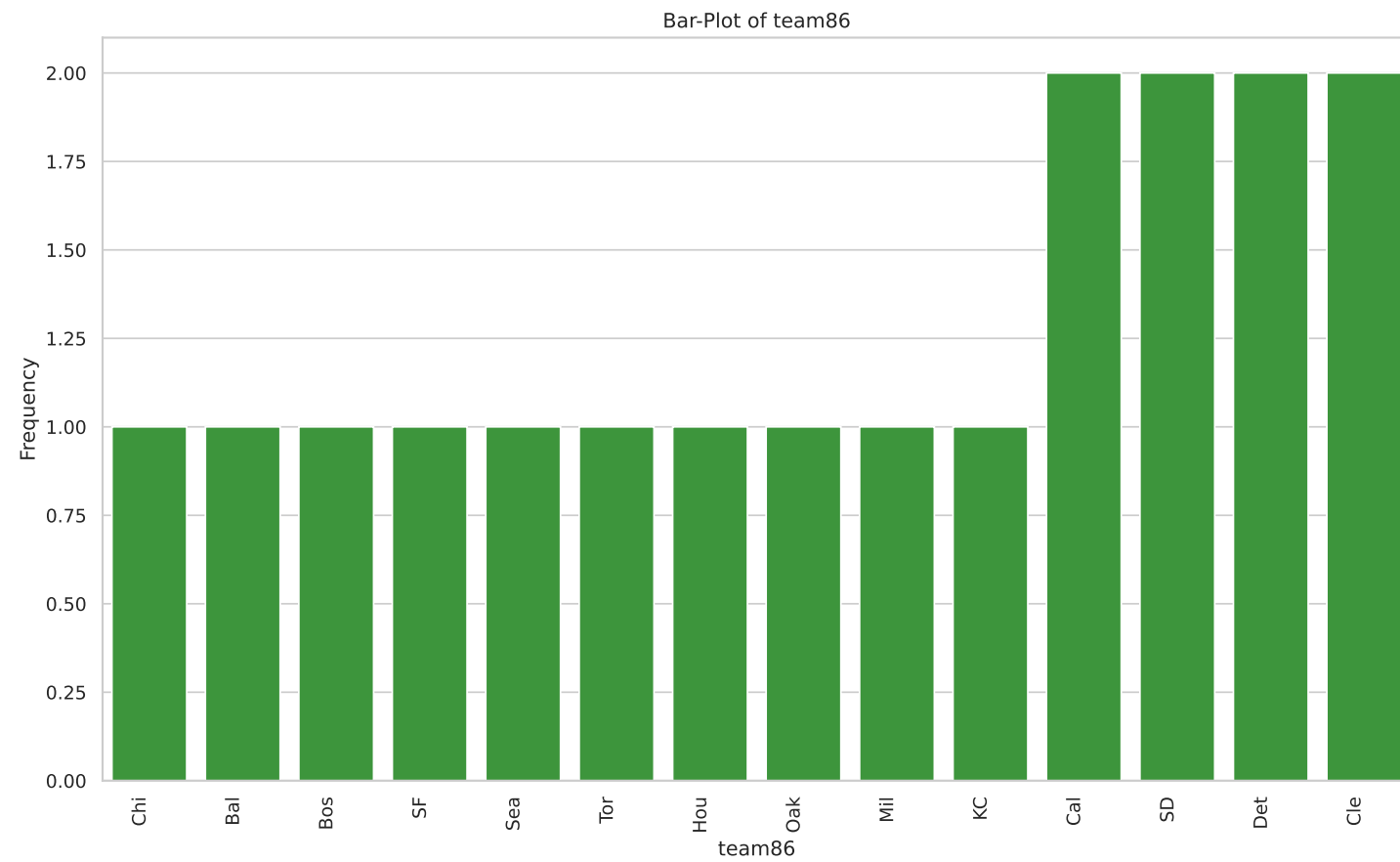
[]



Bar-Plot of div86
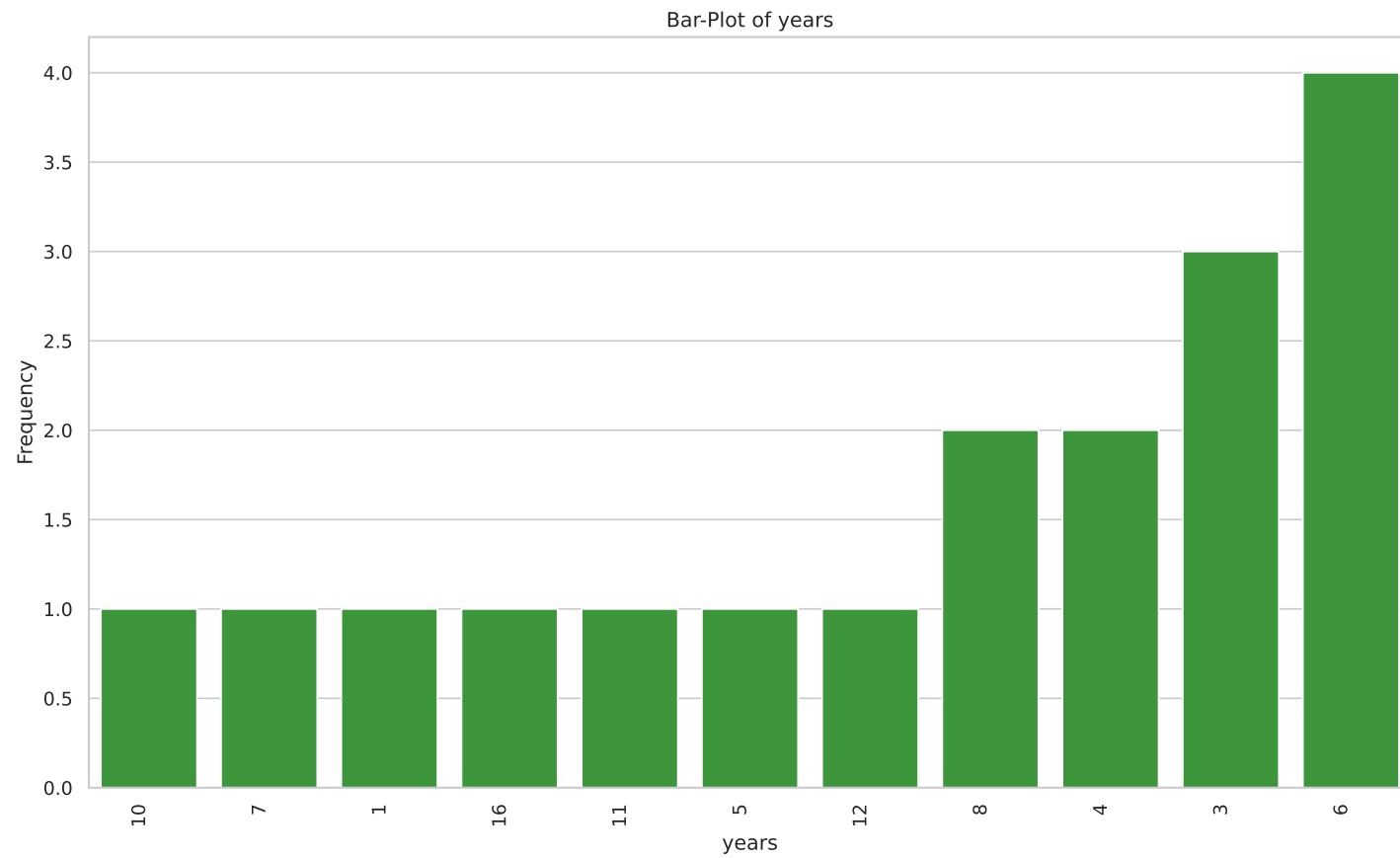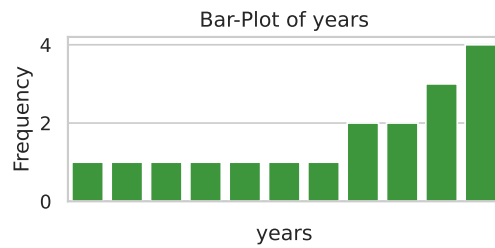
Bar-Plot of league86

Bar-Plot of name1

Bar-Plot of name2

Bar-Plot of posit86

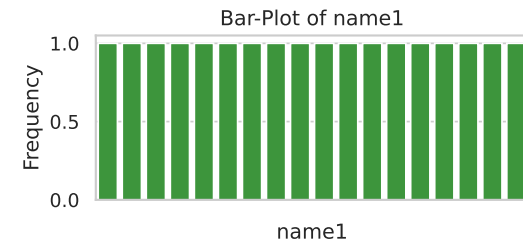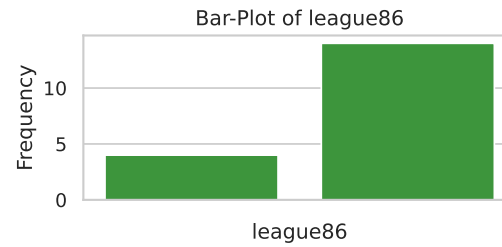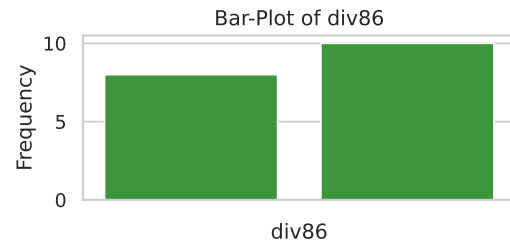Bar-Plot of team86

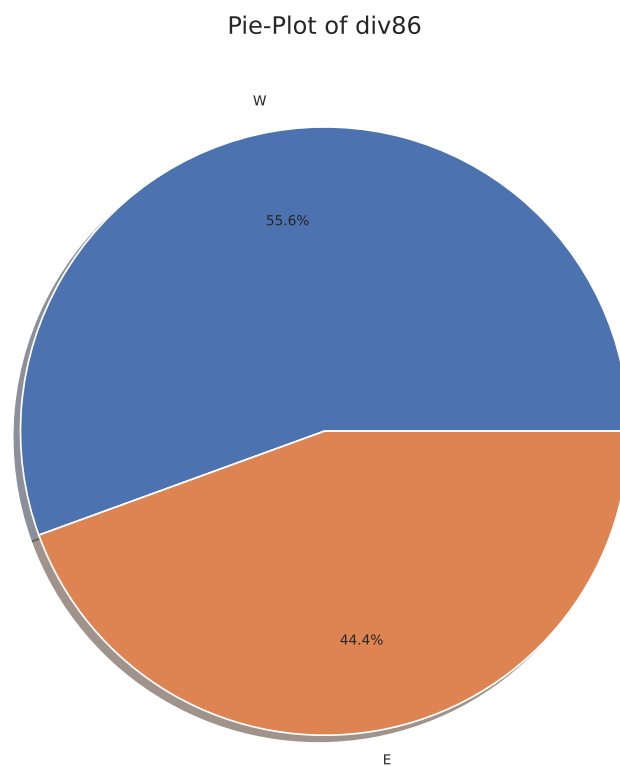Bar-Plot of years

## Bar-Plots Summary

Multiple Bar-Plots of variables in one figure. Variables are sorted alphabetically. No labels displayed.
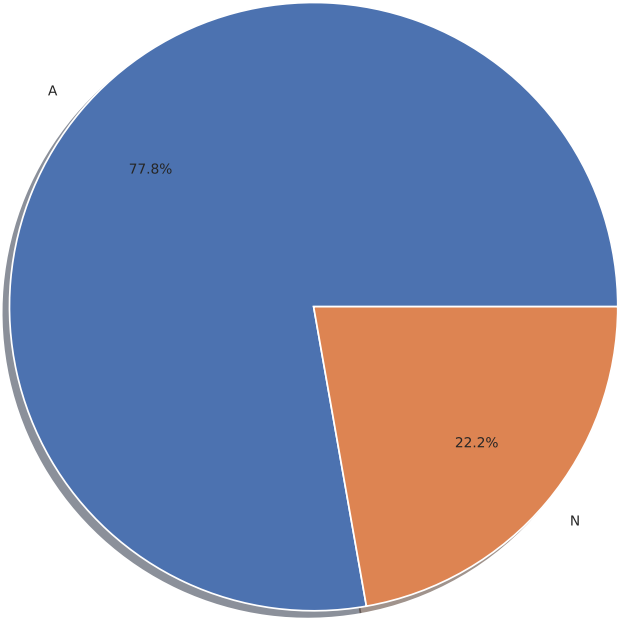
**Pie Plots**

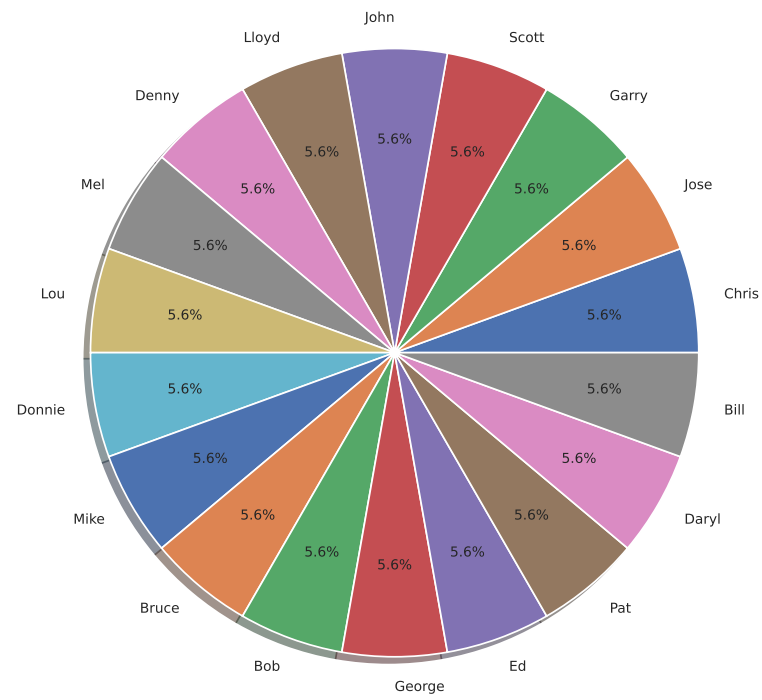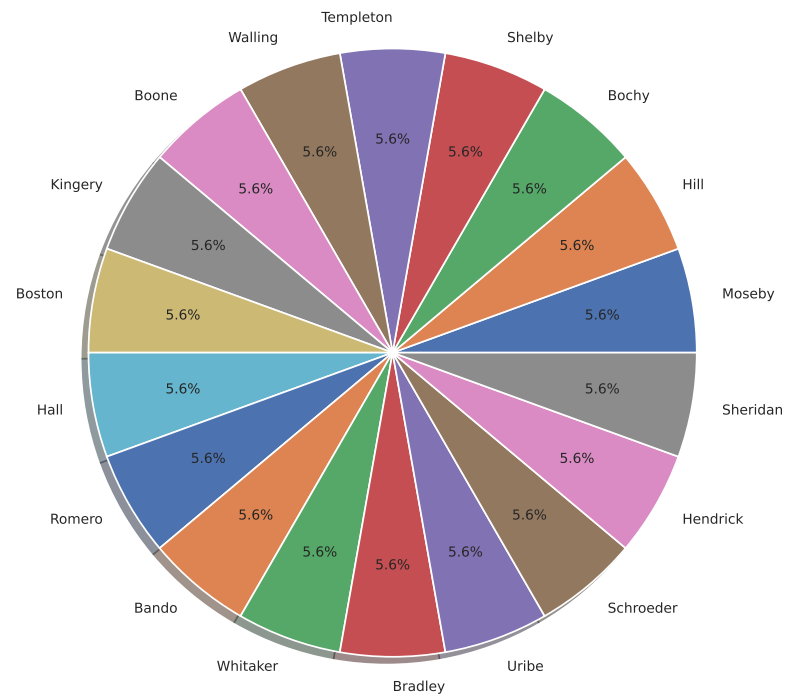One Pie Plot per page for each variable. Variables are sorted alphabetically.

[]

Pie-Plot of div86
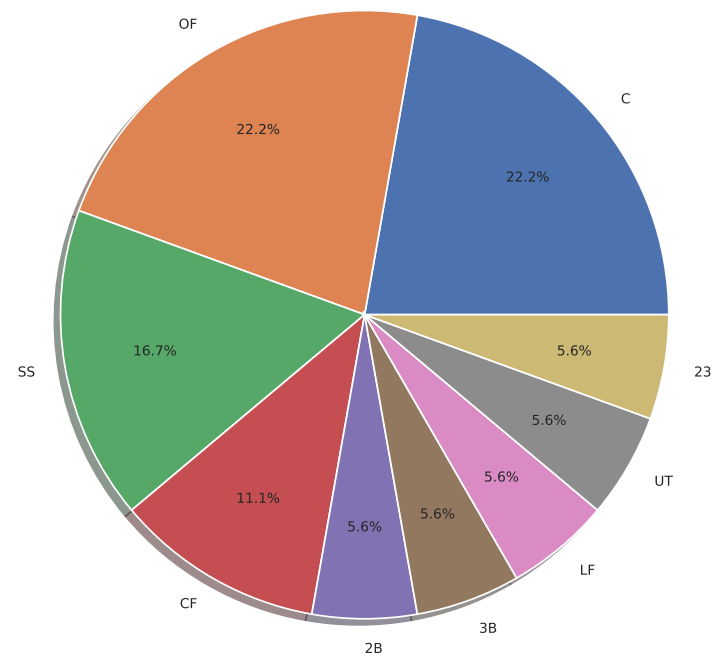
Pie-Plot of league86
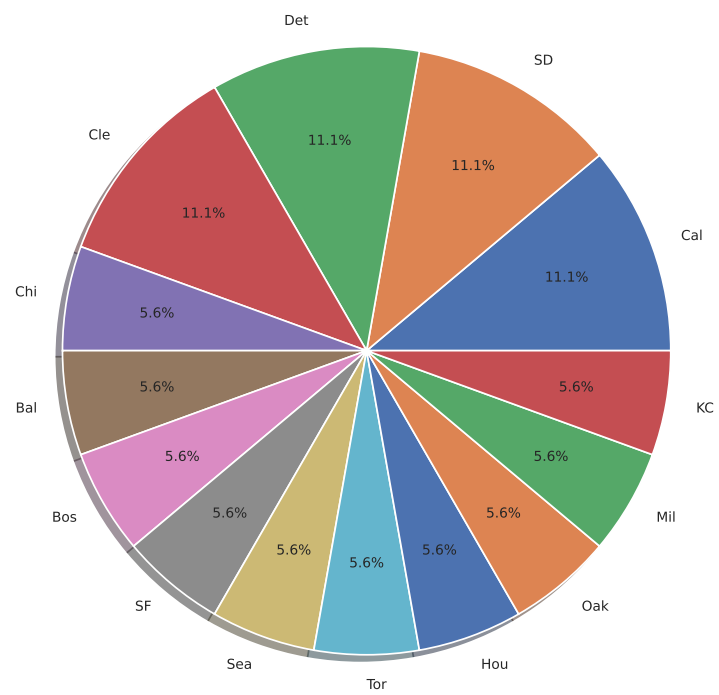
A

77.8%

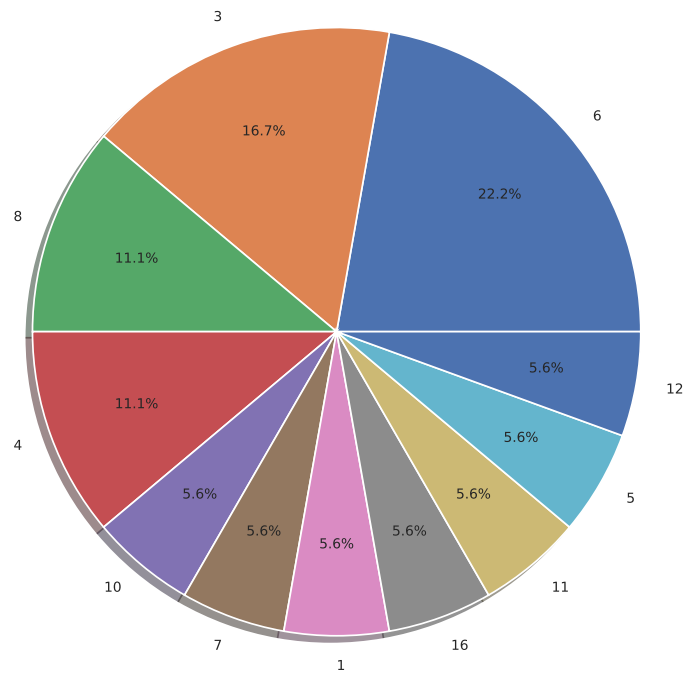22.2%

N

Pie-Plot of name1

Pie-Plot of name2

# Pie-Plot of posit86

Pie-Plot of team86

# Pie-Plot of years

**Pie Plots Summary**

Multiple Pie Plots of variables in one figure. Variables are sorted alphabetically.

See figures on next page.