

How to Easily Create Percentage-Based Crosstabs with Pandas

Authored by
stats writer

November 20, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Easily Create Percentage-Based Crosstabs with Pandas*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=98396>

Mastering Percentage Calculation in Pandas Crosstabs

Analyzing categorical data often requires transforming raw counts into proportional measures to better understand the distribution and relationship between variables. In the [Pandas](#) library, the powerful `pd.crosstab()` function is essential for constructing a **frequency table** or contingency table from two or more factors. While the default output provides raw counts, deriving percentages is crucial for comparative analysis, allowing data scientists to quickly assess relative contributions within the dataset.

The core mechanism for generating these percentages lies within the **normalize** parameter of the `crosstab()` function. This parameter governs how the normalization--the conversion from absolute counts to relative frequencies--is applied. By intelligently setting the value of `normalize`, we can calculate percentages based on the grand total, row totals, or column totals, providing three distinct perspectives on the underlying data relationships. This capability is fundamental for sophisticated data auditing and reporting when utilizing the [DataFrame](#) structure.

Furthermore, for enhanced clarity and completeness in reporting, Pandas offers the **margins** parameter. When set to `True`, `margins` calculates and includes aggregate sums for both rows and columns. When combined with normalization, the resulting marginal totals are also expressed as percentages, providing immediate confirmation of the calculation basis.

The Critical Role of the `normalize` Parameter

The calculation of percentages in a crosstab pivots entirely on the **normalize** argument. This argument dictates the denominator used in the percentage calculation. For instance, normalizing by the index (rows) means dividing each cell count by its corresponding row total, whereas normalizing by columns means dividing by the column total. Selecting the appropriate normalization method is vital, as it directly impacts the interpretation of the resulting statistical report.

You can use the **normalize** argument within the pandas `crosstab()` function to create a crosstab that displays percentage values instead of counts:

```
pd.crosstab(df.col1, df.col2, normalize='index')
```

The **normalize** argument accepts three distinct string values, each producing a unique type of percentage calculation relative to the overall table structure:

all: Display percentage relative to the **grand total** of all observations in the table.

index: Display percentage as a proportion of the corresponding **row total**. This is often used to show internal distribution within categories defined by the row variable.

columns: Display percentage as a proportion of the corresponding **column total**. This is ideal for comparing how row categories contribute to the totals of column categories.

Setting Up the Practical Data Example

To demonstrate the application and interpretation of the three normalization modes, we will utilize a sample `DataFrame` representing player data across various teams and positions. This dataset is simple yet effective for illustrating how different normalization techniques alter the relational view of the data. We first import the `Pandas` library and construct the required data structure.

import pandas as pd

```
#create DataFrame
df = pd.DataFrame({'team': ,
'position':,
'points': })
```

```
#view DataFrame
print(df)
```

```
team position points
0 A G 22
1 A G 25
2 A F 24
3 B G 39
4 B F 34
5 B F 20
6 B F 18
7 C G 17
8 C G 20
9 C F 19
10 C F 22
```

Before applying any normalization, it is instructive to examine the standard, count-based cross-tabulation. This default view provides the absolute count of players categorized by their `team` (index) and `position` (columns). This baseline table confirms that there are a total of 11 players across all teams and positions.

```
#create crosstab that displays count by team and position
pd.crosstab(df.team, df.position)
```

```
position F G
team
A 1 2
B 3 1
C 2 2
```

Example 1: Calculating Percentages Relative to the Grand Total (normalize='all')

The first and most straightforward method of normalization involves setting the `normalize` parameter to `'all'`. This calculates the percentage of each cell entry relative to the **grand total** number of observations in the dataset. Since our sample DataFrame contains 11 entries (players), every count is divided by 11 to obtain its proportional value. This approach is particularly useful for gauging the overall contribution of a specific subgroup to the entire population being studied.

When `normalize='all'` is used, the resulting values are represented as floating-point decimals between 0 and 1. The sum of all cells in this resulting frequency table will always equal 1.0.

#create crosstab that displays counts as percentage relative to total count
pd.crosstab(df.team, df.position, normalize='all')

```
position F G
team
A 0.090909 0.181818
B 0.272727 0.090909
C 0.181818 0.181818
```

Here is how to interpret the output:

Players on team A in position F (1 player) account for **9.09%** (1/11) of total players.

Players on team A in position G (2 players) account for **18.18%** (2/11) of total players.

The results clearly indicate that the largest single subgroup is 'Team B, Position F', which accounts for over 27% of all players examined. This 'all' normalization provides a bird's-eye view of demographic significance across the entire sample space.

Example 2: Analyzing Internal Distribution with Row Percentages (normalize='index')

When the goal is to understand the composition **within** each category defined by the row variable,

we use `normalize='index'`. In this mode, each cell count is divided by the sum of its respective row. If we consider the rows to represent the teams, this calculation shows what percentage of Team A's players are in Position F, and what percentage are in Position G. This is vital for internal comparisons and understanding the characteristics unique to each row group.

For Team A, which has 3 total players (1 F, 2 G), the calculation for Position F will be 1/3, and for Position G, it will be 2/3. Crucially, the percentages across each individual row will sum up to 1 (or 100%), allowing for straightforward interpretation of internal distribution.

#create crosstab that displays counts as percentage relative to row totals
pd.crosstab(df.team, df.position, normalize='index')

```
position F G
team
A 0.333333 0.666667
B 0.750000 0.250000
C 0.500000 0.500000
```

Players in position F account for **33.33%** of total players on team A, while position G accounts for 66.67%.

Players in position F account for **75%** of total players on team B, indicating a heavy specialization in this role within that team.

Players in position F account for **50%** of total players on team C, showing a perfectly balanced distribution between the two positions.

This normalization mode is highly effective for identifying internal biases or specialization among the different teams, as illustrated by Team B's clear preference for 'F' position players.

Example 3: Comparing Contributions Across Column Categories (`normalize='columns'`)

If the analytical focus shifts to understanding how different row categories contribute to a specific column total, we set `normalize='columns'`. In this scenario, each cell count is divided by the sum of its respective column. This normalization addresses the question: "Of all players in Position Y, what percentage come from Team X?" This is particularly useful for resource allocation or benchmarking where the column variable represents a constrained resource or role.

The column percentages are calculated independently for each column. For the 'F' position column (total count 6), every cell in that column will be divided by 6. For the 'G' position column (total count 5), those cells will be divided by 5. The sum of percentages down each individual column will be 1

(or 100%).

```
#create crosstab that displays counts as percentage relative to column totals
pd.crosstab(df.team, df.position, normalize='columns')
```

```
position F G
team
A 0.166667 0.4
B 0.500000 0.2
C 0.333333 0.4
```

Here is how to interpret the output:

Players on team A account for **16.67%** (1/6) of total players with a position of F.

Players on team B account for **50%** (3/6) of total players with a position of F.

Players on team C account for **33.33%** (2/6) of total players with a position of F.

This normalization method clearly highlights the dominance of Team B, which provides half of all the 'F' position players in the dataset.

Enhancing Readability with Marginal Totals

While normalization provides the proportional view of the data, including marginal totals is often necessary for comprehensive data presentation. The `pd.crosstab()` function supports the **margins** parameter, which, when set to `True`, adds a row and a column showing the subtotals. When used in conjunction with normalization, these marginal cells display the total percentages for the corresponding rows and columns.

If we were to apply `normalize='index'` along with `margins=True`, the new 'All' column (the row margin) would contain 1.0 (or 100%) for every row, confirming that the percentages were correctly calculated based on the row totals. The 'All' row (the column margin) would represent the total proportion of players in each position relative to the overall population, which is a calculation based on the frequency table's grand total.

Using `margins=True` is considered a best practice in statistical reporting, as it provides built-in validation for the normalization method applied. Furthermore, it allows the reader to quickly ascertain the overall distribution (like the percentage breakdown of positions across all teams) without running a separate calculation.

Summary of Normalization Techniques

Selecting the correct method for calculating percentages in a Pandas crosstab depends entirely on the question the analysis aims to answer. If the goal is to understand the significance of a data point relative to the entire population, `normalize='all'` is appropriate. If the objective is to assess the internal makeup of row categories, `normalize='index'` provides the necessary distribution. Finally, to evaluate the contribution of row categories toward the totals of column categories, `normalize='columns'` is the ideal choice.

By mastering the **normalize** and **margins** parameters, data analysts can leverage Pandas to generate highly detailed and statistically valid contingency tables that move beyond simple counts to provide profound relational insights.

Note: You can find the complete documentation for the pandas `crosstab()` function .