

Learning to Calculate Grouped Quantiles with Pandas

Authored by
Mohammed loot

November 3, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning to Calculate Grouped Quantiles with Pandas*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=8934>

Introduction to Grouped Quantile Analysis

In the vast landscape of data analysis, deriving meaningful insights often requires looking beyond simple averages. While aggregate statistics provide a broad overview, true understanding of data distribution necessitates the calculation of metrics within specific subgroups. This process, known as grouped quantile calculation, is a fundamental technique in modern data science, offering precise views into how data is distributed across distinct categories. This tutorial focuses on demonstrating how to execute this operation efficiently using the highly regarded [Pandas](#) library in [Python](#).

[Quantiles](#)--which include essential metrics like the median (the 50th percentile) and [quartiles](#)--serve to divide a variable's probability distribution into continuous intervals that hold equal probabilities. By combining these statistical measures with grouping capabilities, data analysts can effectively compare key distribution points. For instance, this allows for robust comparison of test scores between different educational institutions or the nuanced evaluation of sales performance across various geographical regions, moving beyond mere means or sums.

Mastering the application of [Pandas](#) for this task is crucial for anyone working with structured data. When dealing with large datasets, the ability to rapidly calculate multiple [percentiles](#) across different dimensions of a [DataFrame](#) is invaluable for exploratory data analysis (EDA) and subsequent model building.

The Core Pandas Syntax for Grouped Metrics

The most common and straightforward method for performing group-wise quantile calculation involves chaining two primary methods: **groupby()** followed by **quantile()**. This chain elegantly handles the segmentation of the data before applying the desired statistical function to each segment independently. The syntax below illustrates the general structure necessary to calculate a specified [quantile](#) (here represented by 0.5, which corresponds to the median) based on a designated grouping column within a [DataFrame](#).

```
df.groupby('grouping_variable').quantile(.5)
```

Understanding the logic behind this foundational syntax is key to efficient data manipulation. We first must identify the column that will be used to segment the data (the `grouping_variable`), which effectively creates temporary subgroups. Next, we apply the desired statistical function, in this case, **quantile()**, while specifying the desired probability level (e.g., 0.5 for the 50th [percentile](#) or 0.90 for the 90th [percentile](#)). The subsequent practical examples will demonstrate how to implement this powerful syntax effectively in realistic data scenarios.

The result of this operation is a new [DataFrame](#) where the grouping column becomes the index, and the remaining columns show the calculated quantile values for the numerical variables. This provides an immediate, aggregated view of the distribution statistics across all defined groups.

Example 1: Calculating a Single Percentile by Group

For our first practical application, we will define a sample **Pandas DataFrame** that simulates scores achieved by two distinct teams. This simple dataset will allow us to calculate specific high-level distribution points, such as the 90th [percentile](#), for each team independently, providing a clear benchmark of upper-end performance.

```
import pandas as pd
```

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

```
#view first five rows
df.head()
```

```
team score
0 1 3
1 1 4
2 1 4
3 1 5
4 1 5
```

Having initialized our data, we now apply the combination of the [groupby\(\)](#) and [quantile\(\)](#) functions. Our objective is to determine the 90th [percentile](#) (represented by the probability 0.90) of the `score` column, segmented by the `team` column. This computation provides a critical high-level performance benchmark for each group, enabling analysts to quickly understand the upper limits of the scores achieved by the respective teams.

```
df.groupby('team').quantile(.90)
```

```
score
team
1 6.5
2 4.0
```

This concise output demonstrates the immediate power of Pandas for segmented analysis. The

index now represents the unique teams, and the `score` column contains the calculated 90th [percentile](#) value for that group.

Interpreting High-End Performance Distribution Points

The resulting output is a simplified [DataFrame](#), indexed by the grouping variable (`team`), containing the calculated quantile value for the specified numerical column (`score`). Interpreting these results correctly is essential for translating raw numbers into actionable business or statistical insights. The 90th [percentile](#) is defined as the score value below which 90% of the observations fall for that specific group.

Our results clearly show a substantial difference in high-end performance between the two groups. Team 1 exhibits superior top-tier scores compared to Team 2, which has a much tighter distribution of high scores. This type of analysis is far more informative than simply comparing the average score, which could be skewed by outliers.

We can summarize the performance difference based on the 90th [percentile](#) as follows:

The 90th [percentile](#) of scores for **team 1** is **6.5**. This signifies that 90% of Team 1's recorded scores were 6.5 or below.

The 90th [percentile](#) of scores for **team 2** is **4.0**. This indicates that 90% of Team 2's scores were 4.0 or below, highlighting a lower ceiling on their performance within the top decile.

Example 2: Calculating Multiple Quantiles Simultaneously with Aggregation

In many sophisticated analyses, reliance on a single quantile is insufficient; often, a complete picture requires multiple distribution points, such as the first, second (median), and third [quantiles](#), typically used to construct a comprehensive five-number summary or to calculate the interquartile range (IQR). [Pandas](#) provides an efficient pathway to calculate multiple metrics in a single operation using the [agg\(\)](#) (aggregation) method in conjunction with the [groupby\(\)](#) operation.

When leveraging the powerful [agg\(\)](#) function, it is necessary to first define custom functions for the specific [quantiles](#) we intend to calculate, especially if we wish to rename the resulting columns clearly. In this scenario, we create simple functions named `q1` (for the 25th [quantile](#)) and `q3` (for the 75th [quantile](#)). We then pass these functions to the aggregation dictionary.

```
import pandas as pd
```

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

```
#create functions to calculate 1st and 3rd quartiles
def q1(x):
    return x.quantile(0.25)

def q3(x):
    return x.quantile(0.75)

#calculate 1st and 3rd quartiles by group
vals = {'score': }

df.groupby('team').agg(vals)

score
q1 q3
team
1 4.0 5.0
2 2.0 3.0
```

Advanced Aggregation and Interpreting Quartile Spreads

The mechanism employed in the previous example involves passing a dictionary (`vals`) to the `agg()` method. The structure of this dictionary is defined such that the keys represent the column names we wish to operate on (in this case, `score`), and the values are lists of functions (`q1`, `q3`) that Pandas is instructed to apply to the specified column after the data has been grouped. This approach is exceptionally flexible for generating complex statistical summaries involving multiple custom metrics simultaneously within a single pass.

The resulting output is a hierarchical [DataFrame](#) structure, where the original `score` column now contains two sub-columns: `q1` and `q3`. These values represent the lower and upper bounds of the central 50% of the data, respectively, providing a clear and concise summary of the data spread (variability) for each team. The difference between Q3 and Q1 is the Interquartile Range (IQR), a measure of statistical dispersion.

Analyzing these quartile ranges allows us to draw conclusions about the consistency of performance:

The **first quartile (Q1)** and **third quartile (Q3)** of scores for **team 1** are **4.0** and **5.0**, respectively. This results in an IQR of 1.0, indicating a relatively tight clustering of scores around the median.

The **first quartile (Q1)** and **third quartile (Q3)** of scores for **team 2** are **2.0** and **3.0**, respectively. This also yields an IQR of 1.0, but the entire distribution is shifted lower compared to Team 1.

Conclusion and Further Resources

Calculating [quantiles](#) by group is a straightforward yet immensely powerful technique available in **Pandas**. It is essential for accurately understanding the spread, skewness, and distribution of data within specific subsets, providing insights that simple mean calculations often mask. Whether the task involves calculating a single [percentile](#) using the concise **groupby().quantile()** chain or generating multiple [quartiles](#) via the flexible **groupby().agg()** method, these operations form the bedrock for robust statistical analysis and effective exploratory data visualization.

To further solidify your expertise in data manipulation and advanced aggregation techniques, it is highly recommended to explore related functionalities within the expansive **Pandas** ecosystem. Understanding how to apply custom functions using [agg\(\)](#) opens up possibilities for complex, specialized metric calculation beyond standard statistical functions.

The official [Pandas](#) documentation provides comprehensive guidance on performing other common statistical and grouping operations necessary for professional data processing tasks.