

## Quantiles in Machine Learning

**Quantiles** offers valuable insights into data distribution and helping in various aspects of analysis. This article describes quantiles, looks at **how to calculate** them, and talks about **how important they are for machine learning applications**. We also discuss the problems with quantiles and how box plots may be used to represent them. For anybody dealing with data in the field of machine learning, having a firm understanding of quantiles is crucial.

### *What are Quantiles?*

Quantiles divide the dataset into equal parts based on rank or percentile. They represent the values at certain points in a dataset sorted in increasing order. General quantiles include the [median](#) (50th percentile), quartiles (25th, 50th, and 75th percentiles), and percentiles (values ranging from 0 to 100).

In machine learning and data science, quantiles play an important role in understanding the data, detecting outliers and evaluating model performance.

### **Types of Quantiles**

- **Quartiles:** [Quartiles](#) divide a dataset into four equal parts, representing the 25th, 50th (median), and 75th percentiles.
- **Quintiles:** Quintiles divide a dataset into five equal parts, each representing 20% of the data.
- **Deciles:** Deciles divide a dataset into ten equal parts, with each decile representing 10% of the data.

- **Percentiles:** Percentiles divide a dataset into 100 equal parts, with each percentile representing 1% of the data.

### *Steps to Calculate Quantiles*

The steps for calculating quantiles involve:

1. **Sorting the Data:** Arrange the dataset in increasing order.
2. **Determine the Position:** Calculate the position of the desired quantile based on the given formula: “**Position**=(quantile×(n+1))/100”, where n is the total number of observations.
3. **Interpolation (if needed):** Interpolate between two adjacent values to find the quantile if the position is not an integer.

### **Example with Mathematical Imputation:**

Let's consider a dataset: [5, 10, 15, 20, 25, 30, 35, 40, 45, 50].

1. **Median (Q2):** There are 10 observations, so the median position is  $(2 \times (10+1))/2 = 5.5$ . Since, 5.5 is not an integer, we interpolate between the 5th and 6th observations: Median =  $(25+30)/2 = 27.5$ .
2. **First Quartile (Q1):**  $(25 \times (10+1))/4 = 13.75$ . Interpolating between the 13th and 14th observations: Q1 =  $(15+20)/2 = 17.5$ .
3. **Third Quartile (Q3):**  $(75 \times (10+1))/4 = 41.25$ . Interpolating between the 41st and 42nd observations: Q3 =  $(40+45)/2 = 42.5$ .

### *Uses of Quantiles in Machine Learning*

Quantiles play a crucial role in various aspects of machine learning and [data analysis](#). Here are some key uses:

1. **Descriptive Statistics:** Quantiles help summarize the distribution of a dataset, providing insights into its spread and central tendency.
2. **Outlier Detection:** Observations that fall far from certain quantiles may be considered outliers, aiding in anomaly detection.

3. **Probability Distributions:** Quantiles are used to describe the distribution of random variables, facilitating the analysis of probability distributions in machine learning models.
4. **Comparative Analysis:** By comparing quantiles across different datasets, analysts can make informed decisions about the relative standing and characteristics of the datasets.
5. **Risk Assessment:** In finance and other fields, quantiles are used to assess the risk of investments by determining the potential for loss or gain based on the distribution of data.

Understanding these uses is essential for effectively utilizing quantiles in machine learning and data analysis tasks.

### *Challenges and Limitations of Quantiles*

1. **Influence of Outliers:** Quantiles can be sensitive to outliers, especially when calculating quartiles. Outliers can significantly affect the position of quantiles, potentially leading to a misrepresentation of the data's central tendency and spread.
2. **Skewed Distributions:** Quantiles may not fully capture the characteristics of skewed distributions. For highly skewed datasets, the quantiles may not provide a complete picture of the data distribution, especially in the tails.
3. **Variability in Calculations:** Different methods and software packages may use different algorithms for calculating quantiles, leading to variability in results. This can be a challenge when comparing quantiles across different datasets or when using quantiles for decision-making.