

How to Understand and Handle Truncated and Censored Data

Authored by
stats writer

December 6, 2025

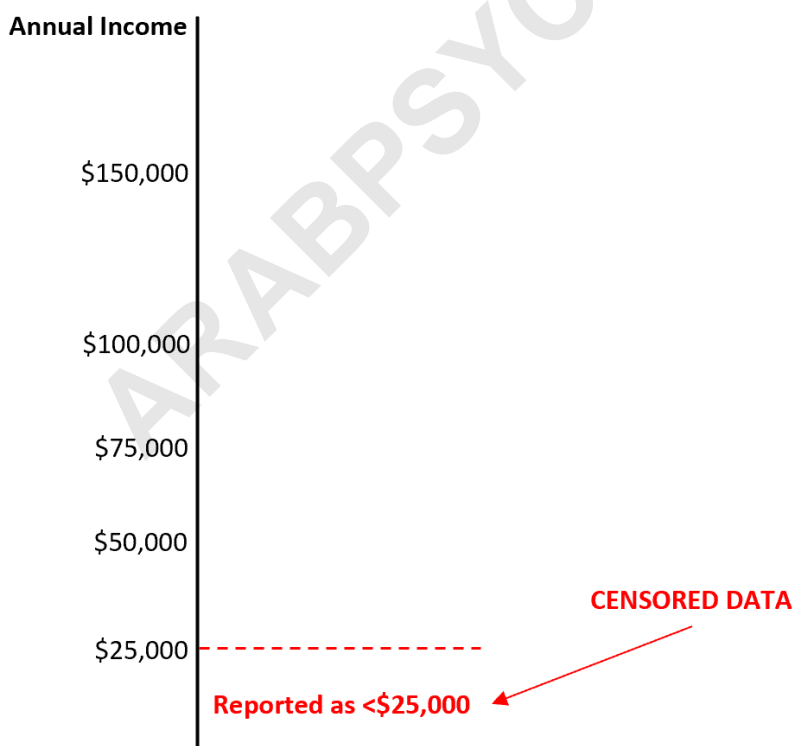
RECOMMENDED CITATION

stats writer (2025). *How to Understand and Handle Truncated and Censored Data*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=106474>

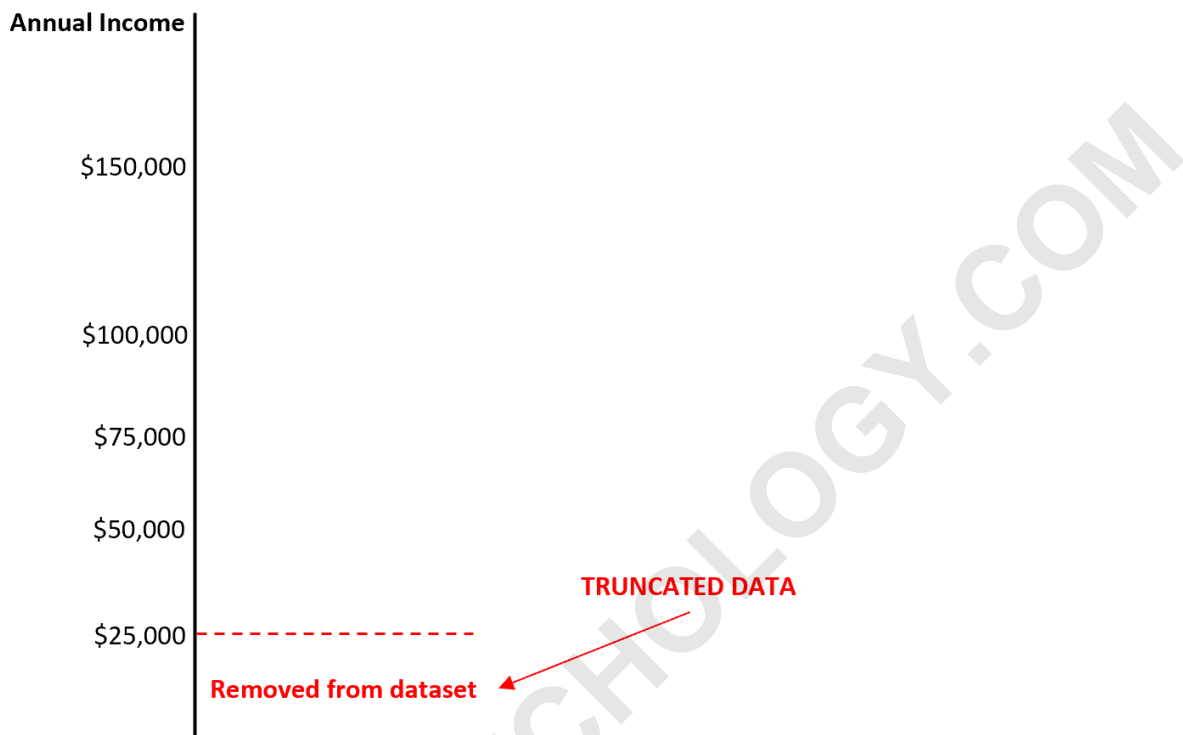
In the field of statistics and data collection, researchers frequently encounter scenarios where complete information about every observation is unavailable or intentionally excluded. This necessity often leads to the use of techniques known as data censoring and data truncation. Understanding these concepts is fundamental because they profoundly influence the interpretation and validity of statistical models, potentially introducing statistical bias if not properly accounted for. While both methods involve incomplete data, the mechanism by which information is lost differentiates them, leading to distinct analytical challenges.

To **censor** data means that we only collect partial information regarding data points that fall below or above a specific threshold value. In this scenario, we know the observation exists and we know its boundary (e.g., greater than X or less than Y), but we lack the precise, exact measurement. This typically occurs due to limits in measurement instruments, constraints in privacy reporting, or the design of the study protocol itself.

Conversely, to **truncate** data values involves the complete removal or exclusion of observations from the dataset that fall outside defined limits. When data is truncated, the researcher is not simply recording incomplete information; they are designing the study or filtering the results such that certain parts of the underlying population distribution are never observed or analyzed. This total exclusion results in a fundamental alteration of the sample distribution compared to the population distribution.



The illustration above visually represents the concept of censoring, where the values outside a measurable range are known only by their boundary condition, not their true magnitude. This distinction is critical for statistical modeling, as truncated data requires entirely different approaches than censored data for accurate inference.



The figure above, depicting truncation, shows how data points are entirely missing from the observed sample space, affecting the shape and properties of the resulting distribution. The remainder of this tutorial will provide detailed examples and discuss the implications of both data handling methods.

Understanding Data Censoring and Its Types

Data censoring occurs when the exact time, value, or magnitude of an event is not fully known, but we know it falls within a certain range or beyond a certain point. It is a highly common phenomenon in various fields, particularly in medical trials, reliability engineering, and economic studies, often being an inherent feature of the data collection process itself. Unlike missing data where the value is simply absent, censored data provides valuable boundary information that must be incorporated into the statistical estimation process. Failing to acknowledge and model censored observations correctly results in biased estimates of parameters such as means, variances, and regression coefficients.

Statisticians typically categorize censoring based on where the unknown value lies relative to the

observable range. The three main types are left-censoring, right-censoring, and interval censoring. In **left-censoring**, the true value is known to be less than or equal to the recorded value (e.g., an instrument can only measure above 0.05 ppm). In **right-censoring**, the true value is known to be greater than or equal to the recorded value, often occurring in longitudinal studies where the study ends before the event of interest happens for some subjects (e.g., a patient survives past the study endpoint). This specific type of censoring is particularly prevalent in Survival Analysis.

Interval censoring is the most general form, where the event of interest is known to have occurred between two specific time points, but the exact time is unknown. For instance, if a medical check-up happens annually, a disease might be detected at the second check-up, meaning its onset occurred sometime between the first and second visits. The presence of censored data necessitates specialized statistical methodologies, such as the use of likelihood functions tailored to incorporate the partial information, ensuring that these incomplete observations contribute meaningfully to the final model without distorting the underlying population characteristics.

Detailed Examples of Data Censoring in Practice

Example 1: Annual Income Reporting (Left-Censoring) - Suppose a researcher is conducting a nationwide survey requiring participants to report their yearly earnings. To protect the privacy of individuals earning very low incomes or to standardize reporting for administrative ease, the survey might specify that if an individual earns less than \$25,000 per year, they simply select the category "<\$25,000," rather than providing the precise dollar amount. This decision results in left-censored data. We know these individuals fall into the bottom segment of the income distribution, meaning their actual income might be \$10,000 or \$24,999, but this level of detail is unavailable. This partial information is still valuable, informing the shape of the lower tail of the income distribution, but requires specialized statistical techniques, like Tobit regression, to accurately estimate the average income across the entire population.

Example 2: Environmental Pollution Measurement (Limit of Detection) - Consider a biologist measuring pollution levels in various water sources using a highly sensitive instrument. Every measuring tool has a defined limit of detection (LOD); suppose this tool cannot reliably measure pollution below 0.002 parts per million (ppm). If a water sample has pollution levels below this threshold, the instrument simply registers the result as "non-detectable" or "<0.002 ppm." This is a classic instance of technical censoring inherent to the measurement process. We are certain that the true pollution level is non-negative and below the LOD, but the exact value remains unknown. If the researcher simply replaced these censored values with zero or the LOD, it would severely skew the estimate of mean pollution and potentially introduce statistical bias, especially if a large proportion of samples fall below the detection limit.

Example 3: Time-to-Event Studies (Right-Censoring) - Right-censoring is commonly encountered in clinical trials or reliability testing, often referred to as 'loss to follow-up.' Imagine a drug trial tracking how long it takes for a cancer patient to relapse. If the trial is scheduled to run for five years, any patient who has not relapsed by the five-year mark is right-censored. We know their time-to-relapse is greater than five years, but we do not know the exact time. The same applies in engineering when testing the lifespan of a component: if the component has not failed when the testing budget runs out, its failure time is censored. Specialized techniques like the Kaplan-Meier estimator are employed in survival analysis to properly incorporate these right-censored observations, recognizing that they contribute information about the longevity or durability of the subjects.

Defining Data Truncation and Sampling Restrictions

Data truncation represents a more severe form of information loss compared to censoring. Truncation occurs when entire subsets of the population are systematically excluded from the observed sample because they do not meet a specific eligibility criterion. If an observation falls outside the defined range, it is not merely recorded partially; it is completely absent from the dataset. This removal fundamentally alters the underlying sample distribution, meaning the observed sample is no longer representative of the original population from which the data was theoretically drawn.

Truncation is often a consequence of the study design or the data collection mechanism, where the ability to observe a value depends directly on the value itself. For example, if we study salaries by only surveying union members earning above \$50,000, we have truncated the data from below. The key analytical challenge with truncation is that the probability of observing a data point is conditional on it satisfying the truncation criteria. Therefore, standard statistical methods that assume a random sample from the entire population will yield highly inaccurate and statistically biased results.

Similar to censoring, truncation can be categorized as left-truncated, right-truncated, or doubly truncated. **Left-truncation** occurs when only observations above a certain threshold are included in the sample (e.g., studying only individuals over 6 feet tall). **Right-truncation** occurs when only observations below a threshold are included. If we only sample data points within an upper and lower bound, we have **doubly truncated** the sample. Specialized likelihood functions that account for the conditional nature of the observed data are required to estimate the parameters of the original, untruncated population distribution accurately.

Detailed Examples of Data Truncation Mechanisms

Example 1: Law Enforcement Crime Research (Left-Truncation) - Consider a criminologist

studying the characteristics and types of offenses committed by individuals in a specific metropolitan area. The dataset compiled by the police department might inherently only contain records for individuals who have committed at least one crime (i.e., Number of Crimes > 0). By definition, any individual who has committed zero crimes is excluded entirely from the dataset; their existence is simply not recorded within the context of this study's variable of interest. This situation represents left-truncation at zero. If the researcher attempts to calculate the average number of crimes committed in the entire population (including non-offenders) based solely on this truncated sample, the estimate will be severely inflated, yielding a highly misleading picture of community crime rates.

Example 2: Education Level and Program Eligibility (Left-Truncation) - Suppose a university professor is evaluating the efficacy of a new, rigorous academic program and decides to only monitor students who demonstrate a high likelihood of success upon entry. The selection criterion is set as a minimum cumulative Grade Point Average (GPA) of 3.5. Consequently, any student applying to the program who has a GPA less than 3.5 is immediately excluded and never enters the research sample. This administrative restriction results in a left-truncated dataset. The resulting analysis comparing the study program's effectiveness will only reflect outcomes for high-achieving students, meaning any inferences drawn about the program's general success or failure cannot be extrapolated back to the broader student body that includes those with lower GPAs.

Example 3: Astronomical Observations (Right-Truncation) - Truncation can also occur in less obvious fields, such as astronomy. When astronomers study stellar populations, they may use a telescope with a limited range, meaning they can only detect objects (stars, galaxies) that are within a certain distance or possess a minimum observable brightness. If researchers are studying the lifespan of stars and can only observe stars whose lifespan is below the current age of the universe (due to detection limits for extremely old, faint stars), this constitutes a form of truncation. Similarly, if studying the size of objects and the observation method only captures objects smaller than a certain measurable diameter, the data is effectively right-truncated, severely distorting the estimation of the true distribution of object sizes in the cosmos.

Key Differences and Statistical Implications

While both censoring and truncation introduce incomplete data, the crucial difference lies in the availability of the observation itself. With **censoring**, the observation is still present in the dataset, and we know that the true value lies within a certain range (e.g., $X_i > 50$). We have partial information about the event's magnitude or timing. With **truncation**, the observation is entirely missing; we never even observe the presence of the data point because it failed to meet the selection criterion necessary for inclusion in the sample. This distinction means that censoring deals with missing precision, while truncation deals with missing subjects entirely.

The severity of the impact on statistical analysis is also different. Truncation results in a greater loss of information because it removes entire segments of the distribution, leading to a sampling frame that is inherently biased toward the observed range. This conditional sampling process requires the use of truncated distribution models for accurate parameter estimation. Conversely, censored data retains the identity of the subject, allowing researchers to incorporate the known boundaries into the likelihood function. Although challenging, techniques for handling censored data (like those used in survival analysis) are often less prone to selection bias than those required for severely truncated data, assuming the mechanism of censoring is properly understood.

Failure to address either phenomenon appropriately invariably introduces statistical bias. If a researcher treats censored data as if the boundary value were the true value (e.g., treating "<\$25,000" as exactly \$25,000), they underestimate the true variance and distort the mean. If a researcher treats truncated data as a random sample from the original population, they fundamentally misunderstand the underlying probability distribution, leading to systematic errors in hypothesis testing and model predictions. Therefore, accurately identifying whether a dataset is censored or truncated is the prerequisite for selecting the correct statistical model.

Mitigating Bias in Incomplete Datasets

Mitigating the bias introduced by censored or truncated data requires employing advanced econometric or statistical techniques. For censored data, methods revolve around maximizing the likelihood function conditional on the available information. For instance, in right-censored time-to-event data, the likelihood contribution for a censored observation involves the probability of survival beyond the censoring time, rather than the probability density function used for uncensored observations. Similarly, the Tobit model is often used when dealing with censored variables, providing estimates of both the observed and the latent (unobserved) variables.

For truncated data, the correction must account for the fact that the sample is derived from a conditional distribution. This usually involves adjusting the probability density function (PDF) or probability mass function (PMF) of the distribution by dividing it by the probability that an observation falls within the observation window. For example, if data is left-truncated at value 'a', the density function used for inference must be $f(x) / P(X \geq a)$. This ensures that the analytical model accurately reflects the conditional nature of the observed data, thereby recovering the parameters of the original, untruncated population distribution and minimizing bias.

Ultimately, the best strategy is proactive study design. When planning data collection, researchers should strive to minimize the need for both censoring and truncation by using instruments with lower detection limits, extending follow-up periods in longitudinal studies, and ensuring sampling methods cover the full range of the population. When these limitations are unavoidable due to ethical, practical, or technical constraints, comprehensive documentation of the censoring or

truncation mechanism is essential for accurate subsequent statistical modeling.

Summary of Key Differences

To summarize the critical distinctions, **censoring** data means retaining the observation but knowing only its boundary, providing partial information about the data value (e.g., failure occurred after day 100). Conversely, to **truncate** data means removing the observation entirely from the dataset because it failed an inclusion criterion, resulting in a conditional and inherently modified sample distribution.

Both phenomena lead to a loss of complete information, making standard statistical inference techniques inadequate. However, truncation generally results in a greater distortion of the sample distribution relative to the target population, requiring more complex adjustments to avoid significant parameter bias in the final analysis.