

Learn About the Hypergeometric Distribution: Definition, Formula, and Examples

Authored by
Mohammed looti

November 8, 2025

RECOMMENDED CITATION

Mohammed looti (2025). *Learn About the Hypergeometric Distribution: Definition, Formula, and Examples*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=13106>

The [hypergeometric distribution](#) is a fundamental concept within [probability](#) theory and **statistics**, meticulously designed to model sampling processes derived from a **finite population**. It stands distinct from other common models, such as the Binomial distribution, because it applies exclusively to situations where sampling is conducted [without replacement](#). This critical distinction means that once an item is selected from the population, it is permanently removed, fundamentally changing the composition of the remaining pool and thereby altering the probability of subsequent selections.

The primary function of the hypergeometric model is to calculate the precise [probability](#) of observing exactly k successes--defined as objects possessing a specific feature--within a total sample of size n . This calculation is conditional on the total population size N containing exactly K such successful objects. This distribution proves indispensable across various fields, including **quality control testing**, where defective items are removed; **evolutionary biology**; and any form of statistical analysis where the population is limited, ensuring that the act of sampling tangibly impacts the remaining probabilities.

Introduction: Defining the Hypergeometric Model

In the realm of probability, defining the context of the experiment is paramount. The hypergeometric distribution establishes the required analytical framework precisely when the assumption of **independence**--a cornerstone of models like the Binomial--is invalidated. This violation occurs inherently through the process of sampling **without replacement**. For instance, consider a routine manufacturing inspection: if a quality control officer identifies a defective component and permanently removes it from the lot, the probability of finding a subsequent defective item is instantly modified, as the overall proportion of defective items remaining in the batch has decreased.

The application of this distribution relies upon the accurate identification of four interconnected parameters that define the sampling scenario: the size of the total [finite population](#) (**N**), the total count of items in that population designated as "successes" (**K**), the size of the sample being drawn (**n**), and the specific number of successes desired within that sample (**k**). The initial and most critical step in successful modeling is ensuring that these parameters are unambiguously and correctly assigned based on the problem statement.

The final calculation yields the likelihood of achieving exactly k successes from n attempts. This is achieved by carefully balancing the number of favorable outcomes (selecting k successes and $n-k$ failures) against the total number of possible outcomes (choosing n items from N). This complex but rigorous accounting process is rooted firmly in **combinatorial mathematics**, which ensures that every possibility is weighted accurately according to the constraints imposed by the limited, non-replenishing population.

The Core Formula and Combinatorial Notation

When a random variable X adheres to the rules of a **hypergeometric distribution**, the likelihood of selecting exactly k objects that possess the feature of interest is determined by a ratio derived from the concept of **combinations**. This formula elegantly balances the desired successful outcomes against the comprehensive total of possible outcomes:

$$P(X=k) = \frac{K C_k (N-K) C_{n-k}}{N C_n}$$

This mathematical structure is fundamentally rooted in the basic principles of counting, often referred to as combinatorics. Specifically, the denominator, $N C_n$, calculates the total number of unique ways one can select a sample of size n from the complete population N . The numerator quantifies the desired number of favorable outcomes, which is the product of two distinct **combinations**: first, the number of ways to choose k successful items from the total available K successes ($K C_k$); and second, the number of ways to select the remaining $n-k$ "failures" from the total pool of $N-K$ failures ($(N-K) C_{n-k}$).

Understanding the notation is crucial for accurate calculation. We define the parameters as follows:

N: Represents the total **population size** (the pool from which items are drawn).

K: Represents the number of "success objects" available in the population, meaning items possessing the **specific feature of interest**.

n: Represents the total **sample size**, corresponding to the number of trials or selections made.

k: Represents the desired **number of successes** to be observed within the drawn sample.

A C_b: Denotes the number of **combinations** of selecting b items from a group of A items.

When applying this powerful formula, strict attention must be paid to the physical constraints of the parameters. The value of k (sample successes) cannot exceed n (the sample size) or K (the total successes available in the population). Similarly, the number of failures selected in the sample ($n-k$) must not exceed the total number of failures available in the population ($N-K$). Should any of these logical constraints be violated, the resultant probability calculation will correctly and automatically yield a value of zero.

Distinguishing Hypergeometric from Binomial Sampling

A frequent conceptual hurdle for those studying **probability** theory is accurately differentiating between the **hypergeometric distribution** and the **binomial distribution**. The sole determinant separating these two models is the selection methodology: whether sampling is performed **with replacement** or **without replacement**. In the Binomial model, every trial is considered independent because the selected item is returned to the pool, guaranteeing that the underlying success probability remains perfectly constant across all draws.

In sharp contrast, the hypergeometric distribution is required when trials are inherently dependent. The removal of an item **without replacement** instantly alters the makeup of the remaining [finite population](#). For example, if the first selection is a success, the ratio of successes to failures remaining in the pool drops, leading to a decreased likelihood of success on the subsequent draw. It is worth noting, however, that if the population size (N) is vastly larger than the sample size (n), the impact of removing a single item becomes statistically insignificant, allowing the **hypergeometric distribution** to be accurately approximated by the [binomial distribution](#) for ease of calculation.

Ultimately, the decision to employ the correct distribution rests solely on whether the process ensures a constant probability of success from trial to trial. If items are replenished, maintaining independence, the Binomial model is appropriate. If the population is depleted, creating statistical dependence and requiring continuous probability adjustments, the **hypergeometric distribution** is the necessary tool for precise and accurate modeling of the outcome.

Practical Application: A Detailed Card Example

To truly appreciate the power and necessity of this distribution, let us examine a classic scenario involving a standard deck of 52 playing cards, which contains 4 Queens. Imagine randomly selecting one card, and then immediately selecting a second card **without replacement**. Our goal is to determine the exact [probability](#) that both cards drawn are Queens.

To solve this using the **hypergeometric distribution**, we must first establish the four critical parameters by mapping the scenario constraints to our variables:

N: Population size = 52 cards (total deck size)

K: Number of objects with the feature = 4 Queens (total successes available)

n: Sample size = 2 draws (total cards selected)

k: Number of successes in sample = 2 Queens (desired outcome)

By substituting these established values into the hypergeometric formula, we can calculate the exact probability. Conceptually, we are determining the ratio of the number of unique ways to select 2 Queens and 0 non-Queens, divided by the total number of unique ways possible to draw 2 cards from the 52-card deck.

$$P(X=2) = \frac{4C2 (52-4)C(2-2)}{52C2}$$

The calculation proceeds by evaluating the [combinations](#): First, $4C2$ (the ways to choose 2 Queens from the 4 available) equals 6. Second, $52-4)C(2-2)$ simplifies to $48C0$ (the ways to choose 0 non-Queens from the 48 available non-Queens), which equals 1. Finally, the denominator $52C2$ (the total ways to choose 2 cards from 52) equals 1326.

The resulting computation is therefore $(6 * 1) / 1326$, which calculates to **0.00452**. This extremely low probability provides intuitive confirmation that drawing two highly specific items sequentially from a large population, especially when sampling **without replacement**, is statistically a rare occurrence.

Essential Mathematical Properties

Beyond determining the probability of single events, the [hypergeometric distribution](#) is characterized by specific statistical properties relating to its **central tendency** and **dispersion**. These properties--namely the **mean** (expected value) and the **variance**--are fundamental tools for theoretical modeling and for forecasting the long-term results of iterative sampling processes conducted under the constraints of this model.

The **mean**, denoted as the expected value (E), quantifies the average number of successes that one would statistically anticipate observing if the sampling experiment were to be repeated indefinitely. This value is logically derived by multiplying the overall rate of success in the population (K/N) by the size of the sample (n). The resulting formula is highly concise and intuitive, directly reflecting the expected proportional yield of successes within the sample:

The mean of the distribution is $(nK) / N$

The **variance** (Var) serves to measure the distribution's spread or variability--that is, the extent to which the observed number of successes typically deviates from the calculated expected mean. Crucially, because the process involves sampling **without replacement** from a [finite population](#), the variance calculation must incorporate a specific multiplier known as the **Finite Population Correction Factor (FPCF)**. This factor, represented by the term $(N-n) / (N-1)$, mathematically compensates for the dependency introduced between sequential trials.

The variance of the distribution is $(nK)(N-K)(N-n) / (N^2(N-1))$

It should be noted that the denominator in the variance formula is occasionally misstated in less rigorous texts. For maximum accuracy and theoretical compliance, the standard formula explicitly includes the term $N-1$ in the denominator. This inclusion ensures that the correction factor properly addresses the dependency inherent in drawing from a **finite population**, which logically forces the variance to decrease as the sample size (n) grows closer to the population size (N).

Hypergeometric Distribution Practice Problems

To test and solidify your comprehension of the [hypergeometric distribution](#), the following practice problems illustrate its application across a range of classic statistical scenarios.

Problem 1: Card Draws (Extended Sample)

Question: Assume four cards are randomly selected from a standard 52-card deck **without replacement**. What is the corresponding [probability](#) of observing exactly two Queens among the four drawn cards?

This problem utilizes the same population parameters as the initial example but increases the sample size (n) from 2 to 4, while keeping the desired number of successes (k) at 2.

N: Population size = 52 cards

K: Number of objects with a certain feature = 4 Queens

n: Sample size = 4 draws

k: Number of objects in sample with a certain feature = 2 Queens

When these values are inserted into the hypergeometric formula, the resulting probability is calculated as approximately **0.025**. This probability is notably higher than the outcome of the two-card draw example, which logically follows, as increasing the sample size (n) provides more opportunities to achieve the target number of successes ($k=2$).

Problem 2: Urn Selection

Question: An urn contains a total of 3 red balls and 5 green balls. If you randomly select 4 balls **without replacement**, what is the [probability](#) of choosing exactly 2 red balls?

In this scenario, the total **population size** (N) is the sum of all balls ($3 + 5 = 8$). The success group (K) is defined by the number of red balls available in the urn.

N: Population size = 8 balls

K: Number of objects in population with a certain feature = 3 red balls

n: Sample size = 4 draws

k: Number of objects in sample with a certain feature = 2 red balls

Applying the formula requires calculating the number of ways to choose 2 red balls from the 3 available, multiplied by the ways to choose 2 green balls (failures) from the 5 available, all divided by the total number of ways to choose 4 balls from 8. The calculated probability for this outcome is high, equaling **0.42857**.

Problem 3: Marbles in a Basket

Question: A basket holds 7 purple marbles and 3 pink marbles. If you randomly select 6 marbles **without replacement**, what is the probability that your selection includes exactly 3 pink marbles?

This final scenario tests the limits of the [hypergeometric distribution](#) by analyzing a large sample size ($n=6$) relative to a very small [finite population](#) ($N=10$). The desired successful outcome (K) is defined by drawing pink marbles.

N: Population size = 10 marbles

K: Number of objects in population with a certain feature = 3 pink marbles

n: Sample size = 6 draws

k: Number of objects in sample with a certain feature = 3 pink marbles

Since only 3 pink marbles are available, achieving a sample of exactly 3 pink marbles requires that the remaining 3 marbles in the sample must be purple (failures). The calculation involves using [combinations](#) to determine the ratio of favorable outcomes, yielding a probability of **0.16667**.