

Math 3501 - Probability and Statistics I

2.4 - The binomial distribution

Binomial distribution

Let X equal the number of observed successes in n Bernoulli trials.

The pmf of X is

$$P(X=x) \leftarrow \text{probability of obtaining } x \text{ success in } n \text{ Bernoulli trial}$$
$$\rightarrow f(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, 2, \dots, n$$

total # of sequences of 1s and 0s of length n with x 1s and $n-x$ 0s

These probabilities are called binomial probabilities, and the random variable X is said to have a binomial distribution.

Notation and Terminology:

- A binomial distribution is denoted by $b(n, p)$ and we sometimes write

n = # of repetition of Bernoulli trial

p = prob. of success in each Bernoulli trials

x : # of success observed

to denote that the distribution of X is $b(n, p)$

- The constants n and p are called the parameters of the binomial distribution

$$X \sim b(n, p) \leftarrow \text{NOTATION}$$

repetition *probability of success*

Summary

A binomial experiment satisfies the following properties:

1. A Bernoulli (success-failure) experiment is performed n times, where n is a (non-random) constant.
2. The Bernoulli trials are independent.
3. The probability of success on each trial is a constant p ; the probability of failure is $q = 1 - p$.
4. The random variable X equals the number of successes in the n trials.

$X \sim bi(n, p)$ means $X = \# \text{successes in } n \text{ independent Bernoulli trials, each with prob. of success } p$

Remark

Let $X \sim b(n, p)$. Then

$X = \# \text{ successes in}$
 $n \text{ Bernoulli trials}$

$\# \text{ failures in } n \text{ Bernoulli trials}$

$$Y = n - X \sim b(n, 1 - p) .$$

$\left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Convenient for certain evaluations}$

Interpretation: If $X \sim b(n, p)$ represents the number of observed successes in n independent Bernoulli trials, then $Y = n - X$ is the corresponding number of observed failures.

Example

Suppose that the probability of germination of a beet seed is 0.8 and the germination of a seed is called a success.

Suppose that the germination of one seed is independent of the germination of another seed.

If we plant 10 seeds, find the probability that:

- a) • exactly 4 germinate;
- b) • at most 8 germinate;
- c) • at least 3 germinate.

Define the r.v. $X = \# \text{ seeds (out of 10) that germinate}$

We know that $X \sim b(10, 0.8)$

$$\left\{ \begin{array}{l} m=10 \text{ (# of repetitions)} \\ p=0.8 \text{ (prob of success in each repetition)} \end{array} \right.$$

a) $P(X=4) = \binom{10}{4} (0.8)^4 \cdot (0.2)^6 \approx \dots \text{ (calculator)}$

b) $P(X \leq 8) = ?$

We can find this number in multiple ways:

$$P(X \leq 8) = \sum_{x=0}^8 \binom{10}{x} \cdot (0.8)^x \cdot (0.2)^{10-x}$$

$X \leq 8$ means $X = 0, 1, 2, 3, \dots, 8$ (at most 8 successes)

is equivalent to

$\underbrace{Y = 10 - X}_{\text{number of failures}} \geq 2$, i.e., we observe at least two failures

\downarrow prob. of failure

But we know that $Y = 10 - X \sim \text{bi}(10, 0.2)$

\uparrow read from table

$$P(X \leq 8) = P(Y \geq 2) = 1 - \overbrace{P(Y \leq 1)}^{m=10, p=0.2, x=1} = 1 - 0.3758 = \dots$$

$$e) P(X > 3) = P(Y \leq 7) = \sum_{y=0}^7 \binom{10}{y} (0.2)^y (0.8)^{10-y} = \dots$$

$X \sim b(10, 0.2)$ ← there is no table

$Y = 10 - X \sim b(10, 0.2)$ ← we can use table

failures

$$X = 3, 4, 5, \dots, 10 \Leftrightarrow Y = 0, 1, 2, \dots, 7$$

$$X > 3 \Leftrightarrow -X \leq -3 \Leftrightarrow 10 - X \leq 10 - 3 \Leftrightarrow Y \leq 7$$

table

$$\begin{aligned} n &= 10 \\ x &= 7 \\ p &= 0.2 \end{aligned}$$

0.9999

Moment generating function of binomial distribution

Let $X \sim b(n, p)$. The mgf of X is given by

$$M(t) = E(e^{tX}) = [(1-p) + pe^t]^n, \quad -\infty < t < \infty.$$

Proof :

$$\begin{aligned}
 M(t) &= E[e^{tX}] = \sum_{x=0}^n e^{tx} \cdot f(x) = \sum_{x=0}^n \cancel{x^{tx}} \cdot \binom{n}{x} \cdot \cancel{p^x} \cdot (1-p)^{n-x} = \\
 &\quad \underset{\substack{x \text{ takes values on } \{0, 1, 2, \dots, n\}}}{=} \\
 &= \sum_{x=0}^n \binom{n}{x} (pe^t)^x \cdot (1-p)^{n-x} \quad \begin{aligned} a &= pe^t \\ b &= 1-p \end{aligned} \\
 &= (pe^t + 1-p)^n = (1-p + pe^t)^n
 \end{aligned}$$

$$\text{Recall : } (a+b)^n = \sum_{x=0}^n \binom{n}{x} \cdot a^x \cdot b^{n-x}$$

converges for all $t \in \mathbb{R}$
(finite # of terms)

we can use this to find
 $E[X] = M'(0)$
 $E[X^2] = M''(0)$

Binomial distribution mean and variance

Differentiating

we obtain *chain rule*

$$M(t) = [(1-p) + pe^t]^n \quad \leftarrow \text{from previous slide}$$

$$M'(t) = n \underbrace{[(1-p) + pe^t]^{n-1}}_{\text{product rule}} \underbrace{(pe^t)}_{\text{chain rule}}$$

and

$$M''(t) = n(n-1) [(1-p) + pe^t]^{n-2} (pe^t)^2 + n [(1-p) + pe^t]^{n-1} (pe^t) \quad \leftarrow$$

$$M''(0) = n(n-1) \underbrace{[1-p+p]}_{=1}^{n-2} \cdot p^2 + n \underbrace{[1-p+p]}_{=1}^{n-1} \cdot p = n(n-1)p^2 + np$$

Evaluating at $t = 0$, we find that

$$\mu = E(X) = M'(0) = np \quad \leftarrow \quad \left\{ \begin{array}{l} M'(0) = n \underbrace{[1-p+p]}_{=1}^{n-1} \cdot p \end{array} \right.$$

and

$$\sigma^2 = E(X^2) - [E(X)]^2 = M''(0) - [M'(0)]^2 \quad \leftarrow E[X^2] = M''(0)$$

$$= n(n-1)p^2 + np - (np)^2 = np(1-p) \quad \leftarrow \underbrace{E[X^2]}_{= n^2 p^2}, \underbrace{(E[X])^2}_{= n^2 p^2}$$

$$\quad \quad \quad (np(1-p)) = np - np^2 + np = np(1-p)$$

Conclusion: If $X \sim b(n, p)$:

- $\mu = E[X] = np$
- $\sigma^2 = \text{Var}(X) = npq$

Intuition: when p is the probability of success on each trial, the expected number of successes in n trials is np .

Special case: when $n = 1$, X has a Bernoulli distribution and we have

- $M(t) = (1 - p) + pe^t$
- $\mu = E[X] = p$
- $\sigma^2 = \text{Var}(X) = pq$

at $n=1$

→ Bernoulli distribution is the
same as $b(1, p)$

↑
 $n=1$

Example

Suppose that observation over a long period of time has disclosed that, on the average, 1 out of 10 items produced by a certain process is defective.

Select five items independently from the production line and test them. Let X denote the number of defective items among the $n = 5$ items.

Determine the mean and variance of X .

5 repetitions of independent Bernoulli trials.
we're counting the # of defective item.

We're told 1 out of 10 items, on average, $\underbrace{\quad}_{m}$ are defective

We're told that probability of having 1 defective item is $P = 0.1$

$$X = \# \text{ defective items (out of 5 inspected)} \quad X \sim b(5, 0.1)$$

$\uparrow \quad \uparrow$
 $m \quad p$

$$\text{mean } \mu = E[X] = m \cdot p = 5 \cdot (0.1) = 0.5$$

$$\text{variance } \sigma^2 = \text{Var}(X) = m \cdot p \cdot q = 5 \cdot (0.1) \cdot (0.9) = 0.45$$

Math 3501 - Probability and Statistics I

2.5 - The hypergeometric distribution

Random experiment

total # balls

red balls
blue balls

Consider a collection of $N = N_1 + N_2$ similar objects split into two dichotomous classes (e.g. red balls vs blue balls):

- $N_1 > 0$ belong to one class (red balls)
- $N_2 > 0$ belong to the other class (blue balls)

A collection of n objects, where $1 \leq n \leq \overbrace{N_1 + N_2}^N$, is selected from these N objects at random and without replacement.

Define the random variable:

X = number of objects selected that belong to the first class.

p.m.f of r.v. X

Question: Determine $P(X = x)$, where the nonnegative integer x satisfies

$$\underline{x \leq n}, \quad \underline{x \leq N_1}, \quad \text{and} \quad \underline{n - x \leq N_2},$$

objects drawn from 2nd class

"analogue" to number of successes in Bernoulli trials

that is, the probability that exactly x of these n objects belong to the first class and $n - x$ belong to the second.

Strategy: We can select:

- x objects from the first class in $\binom{N_1}{x}$ ways;
- $n - x$ objects from the second class in $\binom{N_2}{n-x}$ ways;
- By the multiplication principle, the product

$$\binom{N_1}{x} \binom{N_2}{n-x}$$

equals the number of ways the joint selection operation can be performed;

- there are $\binom{N}{n}$ ways of selecting n objects from $N = N_1 + N_2$ objects.

Conclusion:

$$P(X=x) = \frac{\binom{N_1}{x} \binom{N_2}{n-x}}{\binom{N}{n}},$$

where the space of X is the set S of nonnegative integers x such that

$$x \leq n, \quad x \leq N_1, \quad \text{and} \quad n - x \leq N_2.$$

Hypergeometric distribution

We say that a random variable X has a hypergeometric distribution with parameters N_1 , N_2 , and n , if its pmf is of the form

$$f(x) = \frac{\binom{N_1}{x} \binom{N_2}{n-x}}{\binom{N}{n}},$$

*# of objects
of type 1 # of objects
of type 2 sample size*

$N = N_1 + N_2$

where x is an integer such that

$$\rightarrow \max\{n - N_2, 0\} \leq x \leq \min\{n, N_1\}.$$

Notation: $X \sim HG(N_1, N_2, n)$.

We do not need to specify N

since we know that $N = N_1 + N_2$

Example

A small pond has 50 fish, ten of which have been tagged. If a fisherman's catch consists of seven fish selected at random and without replacement, determine the probability that exactly two tagged fish are caught.

$\frac{50}{N}$ fish
→ 10 are tagged
→ 40 are not tagged

Two types of objects:
 $N_1 = 10$ tagged fish
 $N_2 = 40$ non tagged fish } 50 fish in total

Sample size $n = 7$

Define the r.v.: $X = \# \text{ of tagged fish caught (out of 7)}$

The pmf of X is $f(x) = P(X=x) = \frac{\binom{10}{x} \cdot \binom{40}{7-x}}{\binom{50}{7}}$,
 $0 \leq x \leq 7$

$$f(2) = P(X=2) = \frac{\binom{10}{2} \cdot \binom{40}{5}}{\binom{50}{7}}$$

Example

A lot consisting of 100 fuses is inspected by the following procedure:

Five fuses are chosen at random and tested; if all five blow at the correct amperage, the lot is accepted.

Suppose that the lot contains 20 defective fuses. Determine the probability of accepting the lot.

Lot has 100 fuses
↑
 $N=100$

20 defective $N_1 = 20$
30 are not defective $N_2 = 80$

Pick a sample of $n = 5$ fuses (sample size is $n = 5$)

Define the r.v. $X = \#$ of defective fuses out of the 5 selected (without replacement)

Lot is accepted if $X = 0$. We want to find $P(X=0)$ knowing that $X \sim HG(20, 80, 5)$

$$\text{and so } P(X=0) = \frac{\binom{20}{0} \cdot \binom{80}{5}}{\binom{100}{5}}$$

\uparrow
 \uparrow
 \uparrow
 n_1
 n_2
 n

Hypergeometric distribution mean and variance

Let $X \sim HG(N_1, N_2, n)$. Then:

$$\bullet \mu = E[X] = n \cdot \frac{N_1}{N} \quad \leftarrow \text{identical to the mean of } b(m, p) \text{ with } p \text{ replaced by } \frac{N_1}{N}$$

$$\bullet \sigma^2 = \text{Var}(X) = n \cdot \frac{N_1}{N} \cdot \frac{N_2}{N} \cdot \frac{N-n}{N-1}$$

\uparrow \uparrow \uparrow \uparrow

of objects selected like p like q

Recall variance of $b(m, p)$ is
 $\sigma^2 = m \cdot p \cdot q$

last factor adjusts the variance to sampling without replacement

$$\text{with } \frac{N-m}{N-1} \leq 1$$

\Rightarrow "no replacement" reduces variance

Examples (previous example continued)

$N = 100$ firms

$N_1 = 20$ defective

$N_2 = 80$ non defective

$n = 5$ objects selected without replacement



$X \sim HG(20, 80, 5)$

of defective items found
in a sample of ≈ 5
obtained without replacement

Mean of X is

$$\mu = E[X] = m \cdot \frac{N_1}{N} = 5 \cdot \frac{20}{100} = 1$$

Variance of X is

$$\sigma^2 = \text{Var}(X) = m \cdot \frac{N_1}{N} \cdot \frac{N_2}{N} \cdot \frac{N-m}{N-1} = 5 \cdot (0.2) (0.8) \frac{\frac{95}{99}}{\frac{95}{99}} \approx \dots$$