

## 15. Commonly used Discrete Distribution

- Review r.v., p.m.f., p.d.f., c.d.f. and expected value
- Population variance and standard deviation
- Commonly use discrete distributions
  - Bernoulli
  - Binomial
  - Poisson

### 1. Review:

Random variable  $X$ : a function that numerically records the outcomes of a random experiment

p.m.f. (Discrete) and p.d.f. (continuous): functions

that define the values and probability distribution of the values of  $X$ .

- p.m.f. :  $f(x) = P(X = x), x \in S$

- p.d.f. :  $f(x) = F'(x), x \in S$

c.d.f. : another function that obtains the information of the distribution of  $X$

$$F(x) = P(X \leq x)$$

Expected value : the weighted average of  $X$  that shows

$\mu$

the "center" of the population

Discrete:  $E(X) = \sum_{x \in S} x \cdot f(x)$

continuous:  $E(X) = \int_{x \in S} x \cdot f(x) dx$

In real life, working on estimating the distribution of the whole population is rather difficult, so most statistical methods

focus on some key characteristics of the distribution: parameters

Two most commonly used parameters:

- mean (expected value) (center)

- variance (spread-out)

For ex: a data, by nature, is close to a Normal distribution behavior, then knowing the mean and variance can fully define this Normal distribution.

## 2. Population variance $\sigma^2$ or $\text{Var}(X)$

What is variance? : it measures the variability in the outcomes. The larger it is, the more spread-out the outcome values are.

How is it calculated? : standard deviation measures the "average" distance of the values from the mean; variance is the squared version of that.

Def. The variance of a r.v.  $X$  (or its distribution) is

given by  $\sigma_x^2 = \text{Var}(X) = E\left[\underbrace{(X - \mu_x)^2}_{\substack{\text{"average"} \\ \nearrow \\ \text{distance from the mean}}}\right]$

[Thm]

$$\sigma_x^2 = E[(X - \mu_x)^2] = E(X^2) - \mu_x^2$$

$$\text{p.f. } E[(X - \mu_x)^2]$$

$$= E[X^2 - 2\mu_x \cdot X + \mu_x^2]$$

$$= E(X^2) - 2\mu_x \cdot \underbrace{E(X)}_{= \mu_x} + \mu_x^2$$

↑  
remark: for any r.v.  $X$ ,

$\mu_x$  is a fixed constant

(no matter unknown or known)

$$= E(X^2) - 2\mu_x^2 + \mu_x^2$$

$$= E(X^2) - \mu_x^2$$

[Ex 1]

Bernoulli cp)

$$\text{Var}(X) = E[(X - \mu_x)^2]$$

$$= E[(X - p)^2]$$

$$= \sum (x-p)^2 \cdot f(x)$$

$$= (0-p)^2 \cdot (1-p) + (1-p)^2 \cdot p$$

$$= p^2(1-p) + (1-2p+p^2) \cdot p$$

$$= p^2 - p^3 + p - 2p^2 + p^3$$

$$= p - p^2$$

$$= p(1-p)$$

$$= E(X^2) - \mu^2$$

$$= \sum x^2 f(x) - p^2$$

$$= 0^2 \cdot (1-p) + 1^2 \cdot p - p^2$$

$$= p - p^2$$

$$= p(1-p)$$

$X \sim \text{Bernoulli}(p)$

$$\text{p.m.f. : } f(x) = \begin{cases} p, & x=1 \\ 1-p, & x=0 \end{cases}$$

$$\mu = p, \quad \sigma^2 = p(1-p)$$

Take away: when we have a binary data,  
 we can assume it's from a Bernoulli distribution.  
 if we can estimate  $p$ , we get all the facts  
 about the distribution.

Ex 2 Uniform distribution on  $[a, b]$

$$\sigma^2 = E[(X - \mu_X)^2]$$

$$= \int_a^b \left(X - \frac{a+b}{2}\right)^2 \cdot \frac{1}{b-a} dx$$

$$= \left. \frac{\left(X - \frac{a+b}{2}\right)^3}{3(b-a)} \right|_a^b$$

$$= \frac{\left(b - \frac{a+b}{2}\right)^3}{3(b-a)} - \frac{\left(a - \frac{a+b}{2}\right)^3}{3(b-a)}$$

$$= \frac{\left(\frac{b-a}{2}\right)^3}{3(b-a)} - \frac{\left(\frac{a-b}{2}\right)^3}{3(b-a)}$$

$$= \frac{(b-a)^2}{24} + \frac{(b-a)^2}{24} = \frac{(b-a)^2}{12}$$

Try the  
other way?

$$X \sim \text{Uniform}[a, b]$$

$$f(x) = \frac{1}{b-a}, \quad x \in [a, b]$$

$$\mu = \frac{a+b}{2}, \quad \sigma^2 = \frac{(b-a)^2}{12}$$

### Properties of Variance

$$\begin{aligned} (1) \quad \text{Var}(h(X)) &= E[(h(X))^2 - (E[h(X)])^2] \\ &= E[(h(X))^2] - [E(h(X))]^2 \end{aligned}$$

$$(2) \quad \text{Var}(c) = 0 \text{ for a constant } c$$

$$\begin{aligned} \text{proof: } \text{Var}(c) &= E[c^2] - [E(c)]^2 \\ &= c^2 - c^2 = 0 \end{aligned}$$

$$(3) \quad \text{Var}(aX+b) = a^2 \text{Var}(X)$$

proof?

$$\text{Do we have } \text{Var}(h_1(x) + h_2(x)) = \text{Var}(h_1(x)) + \text{Var}(h_2(x))?$$

[optional] The Moment-Generating Function (m.g.f.)

- A lot of times it's hard to calculate the mean and variance of a distribution, or the higher moments:

$$E(X^3), E(X^4), \text{ etc.}$$

- m.g.f. helps with the calculation

[Def] : the moment-generating function of  $X$  is defined

$$\text{as } M(t) = E(e^{tx})$$

$$[\text{Thm}] : M^{(r)}(0) = E(X^r)$$

$$\text{particularly, } \mu = M'(0)$$

$$\sigma^2 = M''(0) - [M'(0)]^2$$

$$[\text{Ex}] \quad X \sim \text{Bernoulli}(p)$$

$$M(t) = E(e^{tx})$$

$$= \sum_x e^{tx} \cdot f(x)$$



$$= e^{t \cdot 0} \cdot (1-p) + e^{t \cdot 1} \cdot p$$

$$= (1-p) + p \cdot e^t$$

$$M'(t) = p \cdot e^t$$

$$M''(t) = p \cdot e^t$$

$$M'(0) = p \cdot e^0 = p = \mu$$

$$M''(0) = p \cdot e^0 = p = E(X^2)$$

$$M''(0) - (M'(0))^2 = p - p^2 = p(1-p) = \sigma^2$$

### 3. Commonly used discrete distribution

When given a data, most likely we don't know the "exact" distribution. However, based on the behavior and nature of the data, we can make reasonable assumptions of the distribution the data

is drawn from. Based on the assumption, we can then choose appropriate analysis methods.

a. Bernoulli ( $p$ )  $\longleftrightarrow$  Binary outcomes  
 $\quad \quad \quad = \text{parameter } p$

A random experiment has two mutually exclusive outcomes.

Let r.v.  $X$  denote the outcome of the experiment,

then  $X \sim \text{Bernoulli}(p)$ .

$$f(x) = \begin{cases} p & , x=1 \\ 1-p & , x=0 \end{cases}$$

$$\mu = p, \sigma^2 = p(1-p)$$

Application: this is the mostly used distribution in all classification models.

b. Binomial  $(n, p)$  or  $b(\underline{n}, \underline{p})$  parameters  $(n, p)$

— Perform the Bernoulli( $p$ ) experiment  $n$  times.

— The trials are independent

— Let r.v.  $X$  equals the number of successes

$(1s)$  in the  $n$  trials

$X \sim b(n, p)$  s.t.

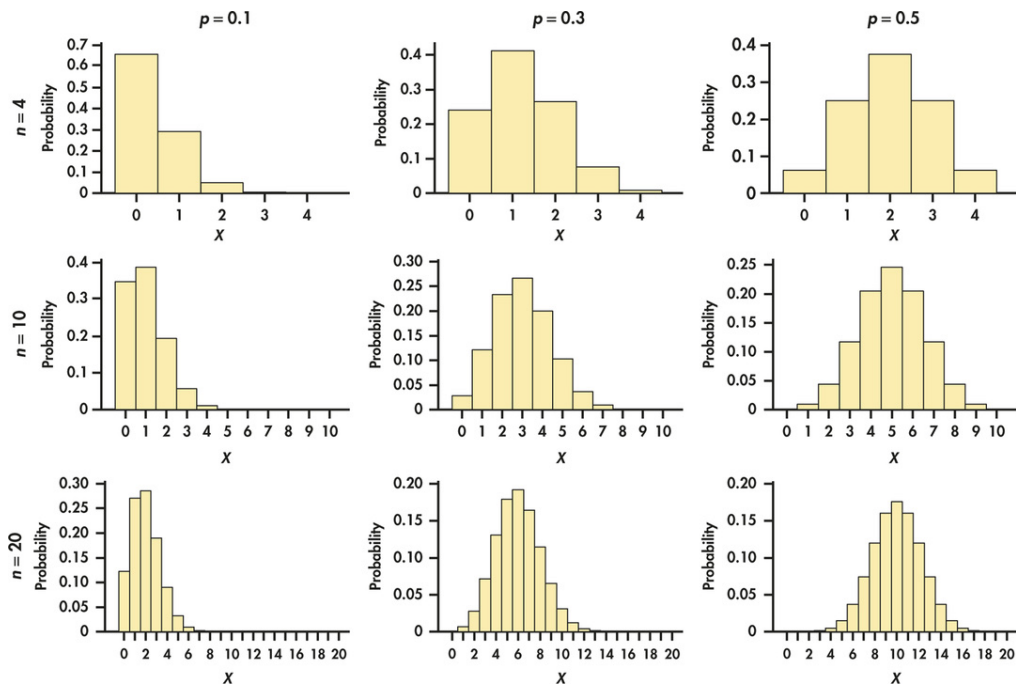
$$f(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x=0, 1, \dots, n$$

$$\mu = np$$

$$\sigma^2 = np(1-p)$$

How does binomial distribution look like?

when  $n$  is large,  $\overset{\text{Appr}}{\sim}$  Normal  $(np, np(1-p))$   
used in  $z$ -test for  $p$  later)



**Ex 3** Lottery! Match 3 Numbers (1-10)

I have 10 people bought tickets (one for each) this time, what's chance that nobody won?

For a single ticket:

$$P_{\text{win}} = \frac{1}{10^3} = 0.001$$

Let  $X$  be the number of winning tickets out of 10.

$$X \sim b(10, 0.001)$$

$$\begin{aligned} P(X=0) &= \binom{10}{0} 0.001^0 (1-0.001)^{10-0} \\ &= 0.999^{10} \\ &\approx 0.99 \end{aligned}$$

What is  $\mu = np = 10 \times 0.001 = 0.01$ ?

the expected number of people will win!

C. Poisson ( $\lambda$ ) parameter  $\lambda$

Let  $X$  be the number of occurrence of an event in a given continuous interval, knowing that the average number of occurrence during

the same interval is  $\lambda$ . Then

$$X \sim \text{Poisson}(\lambda)$$

Ex 4 Earthquake!

It's known that the average number of earthquakes in the area is 5/year.

Let  $X$  be the number of earthquakes

next year.

units match

$$X = 0, 1, 2, \dots, \infty$$

[optional] How to get p.m.f. of this process?



- split the interval into many small intervals, so small that at most one occurrence could happen
- count how many intervals that event happened. out of  $n$ .
- $X \overset{\text{APPR}}{\sim} b(n, p)$
- $p \approx \lambda \cdot \frac{1}{n}$  (this is the trick)

$$P(X = x) = \lim_{n \rightarrow \infty} \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$= \frac{\lambda^x e^{-\lambda}}{x!}$$

Def a r.v.  $X \sim \text{Poisson}(\lambda)$  when

$$f(x) = \frac{\lambda^x e^{-\lambda}}{\lambda!}, \quad x = 0, 1, 2, \dots$$

correction:  $f(x) = \frac{\lambda^x e^{-\lambda}}{x!}$

$$E(X) = \lambda$$

$$\text{Var}(X) = \lambda$$

} m.g.f.

Reist Ex 4

$$X \sim \text{Poisson}(5)$$



$$P(X=0) = \frac{5^0 e^{-5}}{0!} \approx 0.007$$

### Applications

- Poisson regression : predict the number of events occurring within an given interval of time.

- Bridge between discrete and continuous :

Let  $T$  be the time between two occurrences of the events.

$T \sim \text{Exponential}(\lambda)$

key distribution used in survival analysis.

