

Chapter 4: Discrete Distribution

- A random variable assigns numbers to outcomes in the sample space.
 - Eg. # of children with retinitis pigmentosa
of individuals with leukemia
etc...
- Basically, an event that is a number
- Discrete random variable: Can count them (but no, be infinite)
 - Typically 0, 1, 2, 3, ...
- Continuous random variable: Cannot count them
 - Typically some interval $(-\infty, \infty)$, $[0, 1]$, etc..
- Discrete example: 0, $\frac{1}{2}$, 1, $\frac{3}{2}$, ...
 - Continuous example: $[0, 1] \cup (3, -]$
- Denote random variables with capital X, Y, Z , etc.
- A probability mass function (PMF) assigns a probability to a possible value r . Denote this probability by $\Pr(X = r)$.
 - ↑ function of r , not X . X is used to denote the random variable.

Hypergeometric

Ex.) Let $X = \#$ of patients in a trial of 4 wh. have improved blood pressure.

$$\Pr(X=0) = 0.008$$

$$\Pr(X=1) = 0.076$$

$$\Pr(X=2) = 0.265$$

$$\Pr(X=3) = 0.411$$

$$\Pr(X=4) = 0.240$$

- $0 \leq \Pr(X=r) \leq 1$ for all r

- $\sum_r \Pr(X=r) = 1$

↑
all possible r

- Expected Value: Measure of center of a PMF (aka mean)

$$E[X] = \sum_r r \Pr(X=r) = p$$

↑
all possible r 's

Hypothetical

- Ex.) $E[X] = 0 \cdot 0.008 + 1 \cdot 0.076 + 2 \cdot 0.265 + 3 \cdot 0.411 + 4 \cdot 0.240$
 $= 2.8$

↑
Average value of X across many trials

- Note: μ is a population parameter, not a statistic, which is a function of observed data.

E.g. Across many trials, we might see

x	freq
0	0
1	9
2	24
3	48
4	19

↑ theses seen

$$\bar{x} = 0 \cdot 0 + 1 \cdot \frac{9}{100} + 2 \cdot \frac{24}{100} + 3 \cdot \frac{48}{100} + 4 \cdot \frac{19}{100}$$

$$= 2.77$$

- Variance: Measure of spread

$$\text{Var}(x) = \sum_{\text{all } r} (r - \mu)^2 P_r(x=r) = \sigma^2$$

$$SD(x) = \sqrt{\sigma^2} = \sigma$$

- Note: $\text{Var}(x) = E(x^2) - E(x)^2$

$$= \sum_r r^2 P_r(x=r) - \sum_r r P_r(x=r)$$

↑ larger means more variable

- Cumulative Distribution Function (CDF)

$F(x) = P_r(X \leq x)$ = Probability X is less than or equal to x

- Ex.) Hypergeometric

$$F(0) = 0.008$$

$$F(1) = 0.008 + 0.076$$

$$F(2) = 0.008 + 0.076 + 0.265$$

$$F(3) = 0.008 + 0.076 + 0.265 + 0.441$$

$$F(4) = 0.008 + 0.076 + 0.265 + 0.441 + 0.240$$

- Useful for probability calculation

$$P_r(1 \leq X \leq 3) = P_r(X \leq 3) - P_r(X \leq 0)$$

- Exercise) $X = \# \text{ boys in a family of 4}$

r	$P_r(X=r)$
0	$\frac{1}{16}$
1	$\frac{1}{4}$
2	$\frac{3}{8}$
3	$\frac{1}{4}$
4	$\frac{1}{16}$

1) calculate $E(X)$, $SD(X)$, $F(x)$

- Some distributions are seen in real data over and over again.
 - ↳ Binomial = # out of n
 - ↳ Poisson = # during some time interval
- ↳ These have specific PDF's / CDF's,
- ↳ We need to know about permutations / combinations to understand them.
- Number of permutations of n -things taken k at a time

$$n P_k = n(n-1) \cdots (n-k+1)$$

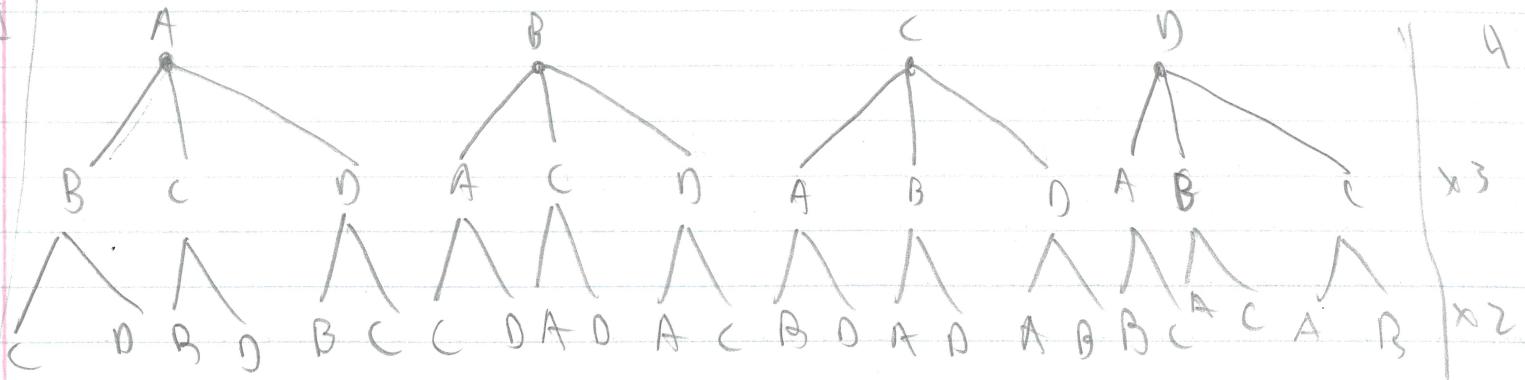
$$= \frac{n!}{(n-k)!}$$

Ex.) Individuals = A, B, C, D

$$4 P_2 = \frac{4!}{2!} = 4 \cdot 3 = 12$$

A, B	C, A
A, C	C, B
A, D	C, D
B, A	D, A
B, C	D, B
B, D	D, C

${}_4P_3$ Example



- What if order does not matter? E.g. $\{A, B\} = \{B, A\}$
- The number of combinations of n things taken k at a time

$${}^nC_k = \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

- There are $\frac{n!}{(n-k)!}$ permutations of size k

Each of those shares elements with $k! = {}^nC_n$

E.g. A, B, C, D $\quad {}^4C_3$

$$A, B, C = A, C, B = B, A, C = B, C, A = C, A, B = C, B, A = 3!$$

↑ Divide by $k!$ to get # combinations.

- Binomial Distribution

1.) n trials

2.) outcome of each trial is "success" or "failure"

3.) $\Pr(\text{"success"}) = p$ for each trial

4.) trials independent

$X = \# \text{ successes}$

then $X \sim \text{Bin}(n, p)$

↳ distributed as

- Ex) White blood count

$X = \# \text{ neutrophiles out of 100 white blood cells}$

$$\Pr(\text{neutrophile}) = 0.6$$

$$\Rightarrow X \sim \text{Bin}(100, 0.6)$$

- If $X \sim \text{Bin}(n, p)$ then

$$\Pr(X=r) = \binom{n}{r} p^r (1-p)^{n-r}$$

Ex) Suppose $X \sim \text{Bin}(3, 0.3)$

$$\Pr(X=2) = \Pr(\text{2 success and 1 failure})$$

$$\Pr(\text{SSF}) = \Pr(\text{SFS}) = \Pr(\text{FSS}) = p^2(1-p)$$

$$\text{So } \Pr(X=2) = 3 p^2(1-p) = \binom{3}{2} p^2(1-p)$$

• Claim: # of ways to order r successes and $n-r$ failures is $\binom{n}{r}$

Proof

Position 1, 2, ..., n

choose r of these to be S, rest are F //

Ex.) $X = \# \text{ boys out of 5 children, } p=0.51$

$$\Pr(X=2) = \binom{5}{2} 0.51^2 0.49^3 = 0.306$$

$$\binom{5}{2} = \frac{5 \cdot 4}{2 \cdot 1} = 10$$

- CDF:

$$\Pr(X \leq x) = \sum_{r=0}^x \Pr(X=r) = \sum_{r=0}^x \binom{n}{r} p^r (1-p)^{n-r}$$

↑ No simpler form

- Mean:

$$E(X) = \sum_{r=0}^n r \binom{n}{r} p^r (1-p)^{n-r} = np$$

↑ Expected # success = # trials • Pr(success)

- Variance

$$\text{Var}(X) = np(1-p)$$

↑ $n \uparrow \Rightarrow \text{Var} \uparrow$

↑ Variance highest @ $p=0.5$, smallest at $p=0$ or 1

- Binomial functions in R

$$\text{dbinom}() = \Pr(X=r)$$

$$\text{pbinom}() = \Pr(X \leq x)$$

$$\text{qbinom}() = \text{quantile}$$

$$\text{rbinom}() = \text{random generation.}$$

Poisson:

Counts of rare events over some period of time
or space

Eg.) # typhoid cases in a year

Eg.) # bacterial colonies on an agar plate

Assume

- 1.) For small time interval Δt , $\Pr(\text{"success"} \text{ in } \Delta t)$ is about $\lambda \Delta t$ (for some λ)
- 2.) $\Pr(\text{more than 2 "successes" in } \Delta t) \approx 0$
- 3.) Stationarity; $\Pr(\text{"success"})$ about the same for all time intervals
- 4.) Independence: One success has no bearing on any other success

↳ Violated, e.g., in epidemic

Then $X = \# \text{ "successes" in time } t$

$$X \sim \text{Pois}(\mu) \quad \text{s.t. } \mu = \lambda t$$

$$\Pr(X = k) = \frac{e^{-\mu} \mu^k}{k!}$$

Note

If $X \sim \text{Pois}(\mu)$ over the t

then $X \sim \text{Pois}(c\mu)$ over the ct

Ex) $X = \#$ typhoid deaths λ 1 year

$$X \sim \text{Pois}(4.6)$$

Let $Y = \#$ typhoid deaths in half a year

$$Y \sim \text{Pois}(4.6/2) = \text{Pois}(2.3)$$

Ex.) $X = \#$ bacteria colonies in 100 cm^2

$$X \sim \text{Pois}(2)$$

$Y = \#$ bacteria colonies in 1000 cm^2

$$Y \sim \text{Pois}(20)$$

• Mean: $X \sim \text{Pois}(\mu)$, $E(X) = \mu$

• Variance: $\text{Var}(X) = \mu$

• If $X \sim \text{Pois}(\lambda_1)$, $Y \sim \text{Pois}(\lambda_2)$, then $X+Y \sim \text{Pois}(\lambda_1+\lambda_2)$

Not generally true for other distributions
(E.g. Not for binomial).

Independent!

• Relation to binomial:

• If $X \sim \text{Bin}(n, p)$

n large (> 100)

p small (< 0.01)

np intermediate

then $X \approx \text{Pois}(np)$.

↑ approximate

• This is used to justify Poisson in cases where we know n is large, but we don't know it exactly.

Eg.) $X = \# \text{ RNA molecules at a gene observed}$
(on the order of 100)

- We don't know n , but know it is large

- we don't know p , but we know it is small
(because $X \sim \text{Bin}(\approx 100)$)

- use poisson to model X !

R functions

Exercises (4.24 - 4.29)

of episodes for 1 child to have OTIS media (ear disease) in 1 year is Pois(1.6)

4.24.) What is Prob of getting 3 or more episodes in the first 2 years of life?

$$\text{sol: } X = \# \text{ in 2 years} \sim \text{Pois}(3.2)$$

$$1 - \text{ppois}(q=2, \lambda=3.2) \\ = 0.6201$$

4.25.) What is Prob of not getting any in 1st year?

$$X = \# \text{ in first year} \sim \text{Pois}(1.6)$$

$$\text{dpois}(x=0, \lambda=1.6)$$

$$= 0.2019$$

4.26.) Prob 2 siblings will both have 3 or more episodes in first year of life?
Assume independence

$$= \Pr(\text{Sib1 has 2 or more}) \cdot \Pr(\text{Sib2 has 2 or more})$$

$$= 0.6201 \cdot 0.6201 = 0.3845$$

out of 2

4.27.) Prob exactly 1 sibling will have 3 or more,

$$Y = \# \text{ sibs}$$

$$\sim \text{Bin}(2, 0.6201)$$

$$\begin{aligned} p_i(Y=1) &= \text{dbinom}(1, \text{size}=2, \text{prob}=0.6201) \\ &= 2 \cdot 0.6201 \cdot (1-0.6201) \\ &\approx 0.4712 \end{aligned}$$

4.28.) Prob neither will have 3 or more episodes in first 2 years?

$$(1-0.6201)^2 \approx 0.1443$$

4.29.) Expected number of siblings in 2-sibling household who will have 3 or more episodes in first two years?

$$E[Y] = 2 \cdot 0.6201 = 1.24$$