

L13.

We wish to model how many patients will be visiting the hospital from 10:am ~ 11:00 am during a day.

Suppose the average number of patients that would arrive is 2.

possible values: $0, 1, 2, \dots, \infty$

(Binomial(n, p), $1, 2, \dots, n$ (upper limit)).

average number is known and fixed.

(Bin(n, p), n known, p -pos).

of patients within a specified time interval 10:00 ~ 11:00.

(Bin(n, p), repeat n times).

Relations:

time interval = 1 hr.

cut the 1 hr interval into 6-sub intervals of 10 mins.

- each interval - the number of patients | 0. prob $\frac{2}{6}$.

- independent.

The total # of patients Binomial(6, $\frac{2}{6}$). $p = \frac{2}{6}$.

Cut into 60 - Sub intervals $\text{Bin}(60, \frac{2}{60})$

i

n -sub intervals. $\text{Bin}(n, \frac{2}{n})$

$n \rightarrow \infty \rightarrow$ poisson distribution

L4 10/24

More examples of Poisson Distribution.

- The number of accident claims per month by an auto insurance company.
- The number of flaws on a loom cable length internal --
- The # of typing errors on a page "area internal"
- The # of raindrops per square inch.
"area internal"

Def. Poisson Experiment.

Suppose we are given an interval (this could be time / length / area / volume) and we are interested in the number of "successes" in that interval. a certain event occurs.

Assume that the interval can be divided into very small subintervals.

- s.t. ① the probability of more than one success in any subinterval is 0
- ② the probability of success in a subinterval is constant for all subintervals and is proportional to the length of the interval.
- ③ the subintervals are independent

This is called a Poisson experiment.

Def. Poisson Distribution

We assume the following:

- The random variable X denotes the # of successes in the whole interval.
- λ is the mean number of successes in the whole interval.
- Then X has a Poisson distribution with parameter λ , denoted $X \sim \text{Poisson}(\lambda)$.

with pmf

$$P(X=k) = \frac{e^{-\lambda} \cdot \lambda^k}{k!} \quad k=0, 1, 2, \dots$$

E.g. 1. The number X of customers arriving at a service desk follows a Poisson Distribution with rate 2 per minute. Find the probability that in a given minute the # of customers will exceed 4.

- $X \sim \text{Poisson}(2)$.

$$\begin{aligned} P(X>4) &= P(X=5) + P(X=6) + \dots \\ &= 1 - P(X=0) - P(X=1) - P(X=2) - P(X=3) - P(X=4) \\ &= 1 - \frac{e^{-2} \cdot 2^0}{0!} - \frac{e^{-2} \cdot 2^1}{1!} - \frac{e^{-2} \cdot 2^2}{2 \times 1} - \frac{e^{-2} \cdot 2^3}{3 \times 2 \times 1} - \frac{e^{-2} \cdot 2^4}{4 \times 3 \times 2 \times 1} \\ &\approx 0.942 \quad = 9.42\% \end{aligned}$$

E.g. 2. The number of flaws in a fibre optic cable follows a Poisson distribution.

The average number of flaws in 50m of cables is 1.2.

(a) what is the probability of exactly 3 flaws in 150m of cable.

- X : # of flaws in 150m of cable.

$X \sim \text{Poisson}(3.6)$.

$$P(X=3) = \frac{e^{-3.6} \cdot 3.6^3}{3 \times 2 \times 1} \approx 0.212$$

(b) what is the probability of at least 2 flaws in 100m of cable.

- X : # of flaws in 100m of cable

$X \sim \text{Poisson}(2.4)$

$$P(X \geq 2) = 1 - P(X=0) - P(X=1)$$

$$= 1 - \frac{e^{-2.4} \cdot 2.4^0}{0!} - \frac{e^{-2.4} \cdot 2.4^1}{1!}$$

$$\approx 0.691$$

(c) what is the probability of exactly one flaw in the first 50m of cable
and exactly one flaw in the second 50m of cable.

- X : # of flaws in 50m of cable.

$X \sim \text{Poisson}(1.2)$.

$$P(0) = P(X=0) \cdot P(X=1)$$

$$= \left(\frac{e^{-1.2} \cdot 1.2^0}{1!} \right)^2$$

$$\approx 0.361$$

Remark (a) 150m

$$\text{Poisson } (3.6) \quad [\text{Poisson}(1.2)]$$
$$P(X=3) = \frac{e^{-3.6} \cdot 3.6^3}{3!} \neq 3 \times \frac{e^{-1.2} \cdot 1.2^3}{3!}$$

(c) 1st 50m and 2nd 50m
 \downarrow \downarrow \downarrow
 $\text{Poisson}(1.2)$ $\text{Poisson}(1.2)$ $\text{Poisson}(2.4)$

$$P(X_1=1, X_2=1) \neq P(Y=2).$$

Mean and Variance

If $X \sim \text{Poisson}(X)$, then

$$\textcircled{1} E(X) = \lambda$$

$$\textcircled{2} \text{Var}(X) = \lambda,$$

$$\textcircled{3} SD(X) = \sqrt{\lambda}.$$

$$\begin{aligned}
 \text{Proof. } \textcircled{1} \quad E(X) &\stackrel{\text{def}}{=} \sum_{k=0}^{\infty} k \cdot P(X=k) \\
 &= \sum_{k=0}^{\infty} k \frac{e^{-\lambda} \cdot \lambda^k}{k!} \\
 &= e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \cdot k \\
 &= e^{-\lambda} \cdot \lambda \underbrace{\sum_{k=0}^{\infty} \frac{\lambda^{k-1}}{(k-1)!}}_{= \frac{\lambda^0}{0!} + \frac{\lambda^1}{1!} + \frac{\lambda^2}{2!} \dots} \\
 &= e^{-\lambda} \cdot \lambda \cdot e^{\lambda} \\
 &= \lambda
 \end{aligned}$$

E.g. It's believed that the number of bookings taken per day at an online travel agency follows a Poisson Distribution. Past records indicate the daily mean number of bookings is 20 and the s.d is 6.

Comment on the suitability of the Poisson Distribution for this example.

- Poisson distribution?

$$\lambda = \text{mean} = 20 \neq \text{sd}^2 = 36$$

Sum of independent Poisson distribution r.v.s.

If X_1, X_2, \dots, X_j are independent Poisson r.v.s with parameters $\lambda_1, \lambda_2, \dots, \lambda_j$.

then $X_1 + X_2 + \dots + X_j$ also follows Poisson distribution with parameter

equals $\lambda_1 + \lambda_2 + \dots + \lambda_j$.

Proof: textbook §3.5 part.

E.g. Suppose X and Y are independent Poisson r.v.s each w/ mean 1.5.

Find (1) $P(X+Y=2)$

(2) $E[(X+Y)^2]$

(1) $X+Y \sim \text{Poisson}(3)$.

$$(i) P(X+Y=2) = \frac{\lambda^{-3} \cdot 3^2}{2!} \approx 0.2241.$$

$$(ii) P(X=0, Y=2) + P(X=1, Y=1) + P(X=2, Y=0).$$

$$(2)(i) E[(X+Y)^2] \stackrel{\text{def}}{=} \sum (x+y)^2 P(X=x, Y=y).$$

$$(ii) \text{Var}(X+Y) = E[(X+Y)^2] - [E(X+Y)]^2$$
$$= E[(X+Y)^2] - 3^2$$

$$E[(X+Y)^2] = 12$$

$$(iii) E[(X+Y)^2] = E(X^2 + Y^2 + 2XY).$$

$$= EX^2 + 2E(X)E(Y) + EY^2.$$

$$= V(X) + [E(X)]^2 + 2 \times 1.5^2 + \text{Var}(Y) + [E(Y)]^2$$

$$= 12$$

L15 10/29

Poisson Approximation (to Binomial dist (n, p)) §2.4.

Suppose

$$n \rightarrow \infty$$

$p \rightarrow 0$, with np staying constant

Thus, writing $\lambda = np$, it can be shown that the binomial probabilities can be approximated by poisson probabilities.

$$P(X=k) = \binom{n}{k} p^k q^{n-k} \approx \frac{e^{-\lambda} \cdot \lambda^k}{k!}, \text{ where } \lambda = np. \quad \forall k=1, 2, \dots, n.$$

Rule of Thumb

$n \geq 20$ and $p \leq .05 \rightarrow$ acceptable approximation.

$n \geq 100$, and $np \leq 10 \rightarrow$ excellent approximation.

E.g. Suppose we sample 100 items from a production line, on average, 2% defectives. Use Poisson approximation to estimate the probability of exactly 3 defectives.

X - # of defectives among 100 items.

$X \sim \text{Binomial}(100, 0.02) \quad \lambda = np = 100 \times 0.02 = 2$.

$$P(X=3) = \binom{100}{3} \times 0.02^3 \times 0.98^{97} \approx 0.1822759 \\ \approx \frac{e^{-2} \cdot 2^3}{3!} = 0.180477.$$

Geometric Distribution

Motivating example.

Suppose we toss a fair dice until the first six appear. Let x denote

the number of tosses. Find pmf of x .

possible value of $x: 1, 2, \dots$

$$P(x=k) = ?$$

$$P(x=1) = \frac{1}{6}$$

$$P(x=2) = \frac{5}{6} \times \frac{1}{6}$$

:

$$P(x=k) = \left(\frac{5}{6}\right)^{k-1} \cdot \frac{1}{6} \quad k=1, 2, \dots \Rightarrow \text{Geometric } \left(\frac{1}{6}\right).$$

Def: the experiment consists of a sequence of independent trials

Each trial results in either S or F.

$P(S)=p$, constant from $\stackrel{\downarrow}{\text{"event of interests occurs"}}$ trial to trial

the trials are performed until first success is observed.

Suppose X denote the total number of trials required, then X has the following pmf.

$$P(X=k) = q^{k-1} \cdot p, \quad k=1, 2, 3, \dots$$

and X is called a geometric distributed random variable

$X \sim \text{Geometric}(p)$.

$$\sum_{k=1}^{\infty} P(X=k) = 1 \quad (0 < p < 1).$$

$$\begin{aligned} \text{proof: } \sum_{k=1}^{\infty} q^{k-1} \cdot p &= p \sum_{k=1}^{\infty} q^{k-1} \\ &= p(q^0 + q^1 + \dots + q^{\infty}) \\ &= p \left(\frac{1}{1-q} \right) \\ &= 1 \quad \square \end{aligned}$$

Mean and Variance.

If $X \sim \text{Geometric}(p)$, then

- $E(X) = \frac{1}{p} \rightarrow$ on average, the # of trials required to achieve
- $\text{Var}(X) = \frac{1-p}{p^2}$ the first success.

proof. LHS = $E(X) = \sum_{k=1}^{\infty} k \cdot p(x=k) = \sum_{k=1}^{\infty} k \cdot q^{k-1} \cdot p$
 $= p \sum_{k=1}^{\infty} k q^{k-1}$
 $= p \cdot \Sigma.$

$$(1) \Sigma = 1 + 2q + 3q^2 + \dots$$

$$(2) q\Sigma = 0 + q + 2q^2 + \dots$$

$$(1-q)\Sigma = 1 + q + q^2 + \dots - q^\infty = \frac{1}{1-q}.$$

$$\Sigma = \frac{1}{(1-q)^2} = \frac{1}{p^2}.$$

$$E(X) = \frac{p}{p^2} = \frac{1}{p} = \text{RHS} \quad \square.$$

P213 Proof for Variance.

Now, suppose a fair dice is tossed until r sixes are observed.

Let X denote the # of tosses required until the r -th six.

Find the distribution of X .

→ say $T=2$.

Possible value of X $2, 3, \dots$

$$P(X=k) = \binom{k-1}{r-1} \left(\frac{1}{6}\right)^{r-1} \cdot \left(\frac{5}{6}\right)^{k-r} \cdot \frac{1}{6} = \binom{k-1}{r-1} \cdot \left(\frac{1}{6}\right)^r \cdot \left(\frac{5}{6}\right)^{k-r}.$$

↳ k tosses, r sixes.

a negative binomial dist
 $(r, \frac{1}{6})$

Def: the experiment consists of a sequence of independent trials

- each trial - S or F
- $P(S) = p$ - constant.
- trials are performed until r successes have been observed.

Let X be the total number of trials required. X has the following pmf.

$$P(X=k) = \binom{k-1}{r-1} p^r \cdot q^{k-r}, \quad k=r, r+1, r+2, \dots$$

* if $r=1$, the negative bin distribution is geometric distribution.

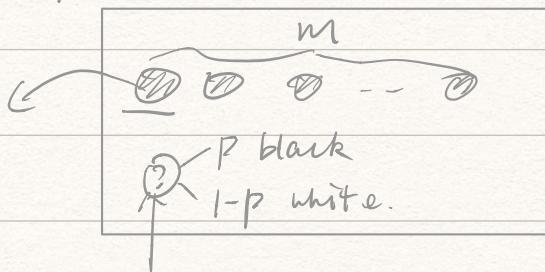
$$X \sim nb(r, p)$$

Mean and Variance.

$$\text{E}(X) = \frac{r}{p}$$

$$\text{Var}(X) = \frac{r(1-p)}{p^2}$$

E.g. A box contains m black balls. At each step, a black ball is removed and a new ball, with probability p to be black and $(1-p)$ to be white, is put into its place. Find the expected (number of steps) required until there are no black balls in the box



at least m steps: no black balls.

$X: m, m+1, \dots$

"Success": if a black ball is removed.

and a black ball is put in.

$$P(s) = 1-p.$$

n successes are required.

$$X \sim \text{nb}(n, 1-p)$$

$$E(X) = \frac{n}{1-p}.$$

e.g. $m=4$ black

$P=0.2$ black

0.8 white.

$$E(X) = \frac{4}{0.8} = 5.$$

E.g. the collector's problem.

Suppose each box of a particular brand of cereal contains one out of a set of n different coupons. Suppose that the coupon in each box is equally likely to be anyone of the set of n , independent of what coupon are in other boxes. Let X be the number of cereal boxes a collector must buy in order to obtain a complete set of n different coupons.

- Find $E(X)$ #of cereal boxes the collector must buy.
- Find the variance $\text{Var}(X)$.

X : possible value $n, n+1, \dots$

Say $n=2$.

1 box is needed to obtain the 1st coupon.

X_2 boxes required to obtain the "2nd coupon" after the 1st box.

$$P(\text{2nd}) = \frac{1}{2}, \quad X_2 \sim \text{Geometric}(\frac{1}{2}).$$

$$E(X) = 1 + \frac{1}{\frac{1}{2}} = 3.$$

Q6 10/31

In general:

X_1 : # of boxes to buy to obtain 1st coupon. $X_1 = 1$.

X_2 : # of boxes to buy to obtain the 2nd coupon after the 1st one is obtained. $P = \frac{n-1}{n}$ $X_2 \sim \text{Geometric}(\frac{n-1}{n})$

⋮

X_n : # of boxes to buy until the n^{th} coupon is observed after the first $(n-1)$ coupons are obtained.

$$P_n = \frac{1}{n}$$

$X_n \sim \text{Geometric}(\frac{1}{n})$.

$X_k \sim \text{Geometric}(\frac{n-(k-1)}{n})$, $k=1, 2, \dots, n$. independent.

$$X = X_1 + X_2 + \dots + X_n.$$

$$\begin{aligned} E(X) &= E(X_1) + \dots + E(X_n) \\ &= 1 + \frac{1}{\frac{n-1}{n}} + \frac{1}{\frac{n-2}{n}} + \dots + \frac{1}{\frac{1}{n}} \\ &= 1 + \frac{n}{n-1} + \frac{n}{n-2} + \dots + n. \\ &= n(\frac{1}{n} + \frac{1}{n-1} + \frac{1}{n-2} + \dots + \frac{1}{1} + 1). \end{aligned}$$

→ Harmonic Series H_n .

For $n=2$

$$E(X) = 2(\frac{1}{2} + 1) = 3.$$

$$\text{Var}(X) = \text{Var}(X_1) + \dots + \text{Var}(X_n) = \sum_{k=1}^n \frac{\frac{k-1}{n}}{(1 - \frac{k-1}{n})^2}$$

Remark $E(X) = nH_n$.

Euler's Constant

as $n \rightarrow \infty$, $H_n \approx \ln n + \gamma + \frac{1}{2n}$, where $\gamma = 0.57721$.

Thus, $E(X) = n \ln n = n\gamma + \frac{1}{2}$ as $n \rightarrow \infty$. by using Euler's approximation.

Hypergeometric Distribution

Def: - A finite population with N elements.

- Each element is a "success" or a "failure".

Total M successes.

- A sample of size n is drawn in such a way that each subset of size n is equally likely to be chosen. [w/o replacement]

Let $X =$ the number of successes in the sample.

X is a hypergeometric distributed r.v.

The probabilities $P(X=x) = h(x; n, m, N)$, will depend on n , M , N .

* Characteristics.

- 1) 2 possible outcomes
- 2) Probability of a success is not the same.
- 3) trials are dependent.
- 4) Population is finite ($N < \infty$).

Def: pmf of X is given by $P(X=x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}}$ for $x=0, 1, 2, \dots, n$.
also impose $x \leq M$.

$$n-x \leq N-M \rightarrow x \geq n-(N-M)$$

Thus, $\max(0, n-N+M) \leq x \leq \min(n, M)$.

E.g. Randomly draw 5 balls from a box containing 4 blue balls and 6 black balls. Find the pmf of x , where $x = \#$ of black balls in the sample.

$$N=10.$$

\textcircled{O}	\textcircled{O}	\textcircled{O}	\textcircled{D}
\textcircled{O}	\textcircled{O}	\textcircled{O}	
\textcircled{O}	\textcircled{O}		

\rightarrow select 5.
 $\frac{\text{''}}{n}.$

$$M=6$$

$$P(x=x) = \frac{\binom{6}{x} \binom{4}{5-x}}{\binom{10}{5}}, \quad x=1, 2, 3, 4, 5.$$

$$\max(0, 5-10+6) \leq x \leq \min(5, 6).$$

$$P(x=2) = \frac{\binom{6}{2} \binom{4}{3}}{\binom{10}{5}} = \frac{60}{252}.$$

Mean, Variance

If $X \sim h(n, M, N)$, then.

$$E(X) = n \cdot \frac{M}{N}.$$

$$\text{Var}(X) = n \cdot \frac{M}{N} \cdot \left(1 - \frac{M}{N}\right) \cdot \underbrace{\left(\frac{N-n}{N-1}\right)}_{\text{finite population correction factor.}} \nearrow \infty \rightarrow 1.$$

$$\text{v.s. bin}(n, p) \quad E(X) = np \quad \text{Var}(X) = npq.$$

* Sampling with replacement $X \sim \text{Bin}(n, p)$ $P = \frac{M}{N}$.
 w/o $X \sim \text{hyper}(n, M, N)$

- the expectation of binomial dist & hyper distribution the same.
- the variance of bin dist & hyperg dist is diff by a factor.
- if N is large, n is small. $\frac{n}{N}$ is very small. $P = \frac{m}{N}$ is among from 0 to 1. then, the 2 distributions are approximately equal.

E.g. A lake contains 500 fish, 60 of which have been tagged by scientists. A researcher randomly catches 10 fish from the lake, find a formula for the probability mass function of X , where X is the number of tagged fish in the researcher's sample.

$$P(X=x) = \frac{\binom{60}{x} \cdot \binom{500-60}{10-x}}{\binom{500}{10}}, \quad x=0, 1, 2, \dots, 10$$

Summary.

1) $E(X_1 + \dots + X_n) = E(X_1) + E(X_2) + \dots + E(X_n)$ without assumption.

$$E(X_1 + \dots + X_n) = E(X_1) E(X_2) \dots E(X_n) \quad \left. \right\} \text{ independent.}$$

$$\text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n)$$

$$\text{Var}(X_1 X_2 \dots X_n) \neq \text{Var}(X_1) \text{Var}(X_2) \dots \text{Var}(X_n).$$

2) If X_1, \dots, X_n are independent random variables with same distribution as r.v. X

- $E(X_i) = E(X), \quad i=1, \dots, n.$

- $\text{Var}(X_i) = \text{Var}(X), \quad i=1, \dots, n.$

Then $E(X_1 + X_2 + \dots + X_n) = nE(X)$

$$\text{Var}(X_1 + \dots + X_n) = n\text{Var}(X)$$

$$\text{SD}(X_1 + \dots + X_n) = \sqrt{n} \text{SD}(X)$$

$\underbrace{X_1, \dots, X_n}$ independent identically distributed r.v.s.
"iid" random variables.