

[Jan 26]

Remark Use Poisson distr. for rare successes, e.g.

typos on a page of a book

software failures on a given day, ...

Ex (in Ex. p. 34): of 40 students, ~~two~~ 2 are sick on ave. $X \doteq \# \text{ sick}$.

$$\Rightarrow X \approx \text{Poisson}(2) \Rightarrow P\{X=2\} \approx e^{-2} \cdot \frac{2^2}{2!} = 0.0902 \quad (\text{compare to } 0.0901).$$

~~Ex~~ • $P\{\text{no sick student}\} = P\{X=0\} = e^{-2} = 0.14$.

• $P\{\text{at least two sick}\} = 1 - P\{X=0\} - P\{X=1\} = 1 - e^{-2} - 2e^{-2} = 0.56$.

Prop (Sum of Bernoulli = Binomial) $X \sim \text{Binom}(n, p)$, $Y \sim \text{Binom}(m, p)$ indep \Rightarrow $X+Y \sim \text{Binom}(n+m, p)$

Prop (Sum of Poisson = Poisson)

(Let $X \sim \text{Pois}(\lambda)$, $Y \sim \text{Pois}(\mu)$ be independent r.v. Then

~~QED~~ $X+Y \sim \text{Pois}(\lambda+\mu)$

~~for Poisson~~ Poisson. (if $X \sim \text{Binom}(n, p)$

$$* P\{X+Y=n\} = \sum_{k=0}^n P\{X+Y=n, X=k\} \\ = \sum_{k=0}^n P\{X=k, Y=n-k\} = \sum_{k=0}^n P\{X=k\} \cdot P\{Y=n-k\} \quad (\text{indep.})$$

$$= \sum_{k=0}^n e^{-\lambda} \frac{\lambda^k}{k!} \cdot e^{-\mu} \frac{\mu^{n-k}}{(n-k)!}$$

$$= \frac{e^{-(\lambda+\mu)}}{n!} \sum_{k=0}^n \binom{n}{k} \lambda^k \mu^{n-k}$$

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

$$= \frac{e^{-(\lambda+\mu)}}{n!} (\lambda+\mu)^n \quad (\text{Binomial Thm}).$$

2.8. Geometric Distribution.

"Time until first success"

- Consider ∞ sequence of indep. trials, each resulting in a success with prob. p .
Let $X = \# \text{ trials} \rightarrow$ that is the first success.

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

- PMF: $P\{X=k\} = \underbrace{(1-p)^{k-1}}_{\substack{\text{first } k-1 \text{ trials} \\ \text{are failures}}} \cdot \underbrace{p}_{\substack{k-\text{th is} \\ \text{a success}}}, \quad k=1, 2, 3, \dots$

Def $X \sim \text{Geom}(p)$, $p \in (0, 1)$ is a r.v. that takes values $1, 2, \dots$, and
 $P\{X=k\} = (1-p)^{k-1} p, \quad k=1, 2, \dots$

Prop $E[X] = \frac{1}{p}$; $\text{Var}(X) = \frac{1-p}{p^2}$

$$\begin{aligned} E[X] &= \sum_{k=1}^{\infty} k(1-p)^{k-1} p = P \cdot \sum_{k=1}^{\infty} k q^{k-1} \quad (\text{where } q = 1-p) \\ &\quad \leftarrow \sum_{k=1}^{\infty} q^{k-1} + \sum_{k=2}^{\infty} q^{k-1} + \sum_{k=3}^{\infty} q^{k-1} + \dots \\ &= (1+q+q^2+\dots) \sum_{k=1}^{\infty} q^{k-1} = \left(\sum_{k=1}^{\infty} q^{k-1} \right)^2 = \frac{1}{(1-q)^2} = \frac{1}{p^2}. \end{aligned}$$

$E[X] = \sum_{k=1}^{\infty} P\{X \geq k\}$ (by KW Exercise 237)

~~$\Rightarrow P\{X \geq k\} = P\{X \geq k\}$~~ Now, $X \geq k \Leftrightarrow$ the first $k-1$ trials were failures
 $\Rightarrow P\{X \geq k\} = (1-p)^{k-1}$.

$\Rightarrow E[X] = \sum_{k=1}^{\infty} (1-p)^{k-1} = \frac{1}{p}$ (geometric series).

Ex Research teams A and B are competing for the discovery of a new particle. Each day, team A may discover the particle with prob. P , and team B with prob. q (independently). What is the prob. that team A discovers the particle before team B?

$$\left[\begin{array}{l} X = \text{day team A discovers particle}, \quad X \sim \text{Geom}(P) \\ Y = " \text{---" B } \quad \text{---} . \quad Y \sim \text{Geom}(q) \end{array} \right] \text{indep.}$$

$$\begin{aligned} P\{X < Y\} &= \sum_{k=1}^{\infty} P\{X < Y, X = k\} \\ &= \sum_{k=1}^{\infty} P\{Y > k, X = k\} \\ &= \sum_{k=1}^{\infty} P\{Y > k\} \cdot P\{X = k\} \quad (\text{indep.}) \end{aligned}$$

$$= \sum_{k=1}^{\infty} \underbrace{(1-q)^k}_{\substack{\text{failures on} \\ \text{first } k \text{ days}}} \cdot \underbrace{(1-p)^{k-1} p}_{\substack{}}.$$

$$= \boxed{\frac{p(1-q)}{1 - (1-p)(1-q)}} \quad \left(\approx \frac{p}{p+q} \text{ if } p, q \text{ small} \right)$$

Ex (Coupon Collecting Problem I) (Euler)

There are n coupons different types of coupons.

Each time one obtains a coupon, it is equally likely to be of ^{each} type.

Compute Expected of number of the different coupons among N collected.

Applications: Clinical trials - collecting info on side effects of drug.

$$Y = N - X, \text{ # of uncollected coupons.}$$

$$E(Y) = ?$$

$$Y = Y_1 + Y_2 + \dots + Y_n, \text{ where } Y_i = \begin{cases} 1, & \text{coupon of } i^{\text{th}} \text{ type is not collected} \\ 0, & \text{---} \end{cases}$$

~~$$Y \sim \text{Bernoulli}(p)$$~~
$$E(Y) = \sum_{i=1}^n E(Y_i) = E(Y_i) = 1 \cdot p + 0 \cdot (1-p) = p, \text{ where } p =$$

$$\boxed{p = P\{\text{coupon of } i^{\text{th}} \text{ type is not collected}\} = \left(1 - \frac{1}{n}\right)^N}$$

$$\Rightarrow E(Y) = \underbrace{\sum_{i=1}^n p}_{\substack{\text{collected at one trial.} \\ \text{all } Y_i \text{ are independent}}} = n \left(1 - \frac{1}{n}\right)^N.$$

$$E(X) = \boxed{n - n \left(1 - \frac{1}{n}\right)^N}$$

$$\text{E.g. if } N=n \Rightarrow E(X) = \underbrace{(1 - \frac{1}{e})}_\infty n = 0.63n.$$

Asympt. Analysis: $n \rightarrow \infty, N = tn$

$$E(Y) \approx n e^{-N/n} = \cancel{n} \boxed{ne^{-t}}$$

Thus: 63% of collection
| after collecting n coupons.

$$E(Y) \leq 1 \text{ for } t \approx \log n \Rightarrow \text{for } N \approx n \log n.$$

Should Expect a complete collection in time $N \approx n \log n$.

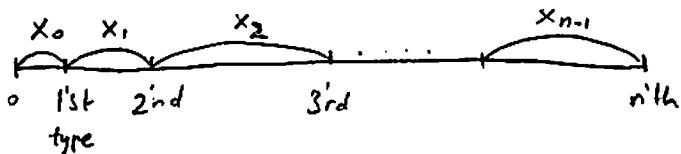
Let's verify this:

Ex (Coupon Collector's Problem II). (Ex. 2i) $\Leftarrow X$

Lec 33 03/30

What is the expected number of coupons one needs to collect before obtaining a complete set of all n types of coupons?

- $X = X_0 + X_1 + \dots + X_{n-1}$, where X_i = number of additional coupons (after i types have been collected) in order to obtain a new type



- $E[X] = \sum_{i=0}^{n-1} E[X_i]$.

- $X_i \sim ?$

When i types of coupons have already been collected, we are waiting for a new type. The coupons that we get now are of a new type with prob.

$$P_i = \frac{n-i}{n} \quad \begin{matrix} \leftarrow \# \text{ new types} \\ \rightarrow \# \text{ all types} \end{matrix}$$

Hence

$$X_i \sim \text{Geom}(P_i) \Rightarrow E[X_i] = \frac{1}{P_i}$$

- $\Rightarrow E[X] = \sum_{i=0}^{n-1} \frac{1}{P_i} = \frac{n}{n} + \frac{n}{n-1} + \dots + \frac{n}{1} = n \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \right)$

Asymptotic analysis $\xrightarrow{\text{as } n \rightarrow \infty}$ $1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \approx \ln n$ (harmonic series).

$\Rightarrow E[X] \approx n \ln n$ (log-oversampling)

More precisely, harmonic number $H_n = \sum_{k=1}^n \frac{1}{k} = \ln n + \gamma + \frac{1}{2} n^{-1} + \tilde{o}(n^{-2})$

0.58 (Euler-Mascheroni constant)

$$\Rightarrow E[X] \approx n \ln n + \gamma n + \frac{1}{2} + \tilde{o}(1)$$

Example: $n=50$, $E[X] \approx 225$

Remark. (Erdős-Renyi '61). $P\{X < n \ln n + \epsilon n\} \rightarrow \exp(-e^{-t})$, $n \rightarrow \infty$ ($\forall t > 0$)