

MATH 180A: INTRO TO PROBABILITY (FOR DATA SCIENCE)

www.math.ucsd.edu/~tkemp/180A

Today: § 4.4 - 4.5

Next: § 4.6

Lab 4 due Wednesday (Nov 6) by 11:59pm

HW 5 due Friday (Nov 8) by 11:59pm

Poisson vs. Normal Approximation - Quantitative

4.4

Theorem. Let $S_n \sim \text{Bin}(n, p)$
 $X \sim \text{Poisson}(np)$
 $Y \sim \mathcal{N}(0, 1)$

For any subset $A \subseteq \mathbb{N}$,

$$|\mathbb{P}(S_n \in A) - \mathbb{P}(X \in A)| \leq np^2$$

if $p = \frac{\lambda}{n^{0.51}}$ ($\lambda > 0$)
 $np^2 = n \left(\frac{\lambda}{n^{0.51}}\right)^2 = \frac{\lambda^2}{n^{1.02}}$
 $\rightarrow 0$
as $n \rightarrow \infty$

OTOH, for any $x \in \mathbb{R}$,

Berry-Essen Thm.

$$\left| \underbrace{\mathbb{P}\left(\frac{S_n - np}{\sqrt{np(1-p)}} \leq x\right)}_{\text{CDF of } \frac{S_n - np}{\sqrt{np(1-p)}}} - \underbrace{\mathbb{P}(Y \leq x)}_{\Phi(x)} \right| \leq \frac{3}{\sqrt{np(1-p)}} \leftarrow \text{optimal}$$

3 is not optimal

Upshot: if np^2 is small, use Poisson Approximation.

if $np(1-p)$ is quite large, use Normal Approximation.

Beyond independent trials:

- * The normal approximation breaks down quickly if the trials are dependent.
- * The Poisson approximation holds up well under "weak dependence"

Example. A factory experiences 3 accidents per month, on average.
What is the probability there will be 3 accidents this month?

$X =$ # accidents in a given month.

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

well modeled
by a Poisson.

$$3 = \mathbb{E}(X) = \lambda$$

$$P(X=3) = e^{-3} \frac{3^3}{3!} = 22.4\%$$

$$\frac{3^3}{3!} = \frac{3^2 \cdot \cancel{3}}{\cancel{3} \cdot 2 \cdot 1} = \frac{3^2}{2!}$$

$$P(X=2) = e^{-3} \frac{3^2}{2!} = 22.4\%$$

Wait Times

4.5

Question: You're tossing a fair die until you get a 6. It's been 12 tosses already, but no 6's so far. The time you have to wait until the first 6 from now is

- (a) Less than
- (b) The same as
- (c) Greater than

the time you would wait from the start.

"Gambler's Fallacy"

Time of first success $T \sim \text{Geom}(p)$.

$$P(T > t) = \sum_{k=t+1}^{\infty} (1-p)^{k-1} p = p(1-p)^t \sum_{l=0}^{\infty} (1-p)^l = p(1-p)^t \cdot \frac{1}{1-(1-p)}$$

$F_T(t) = P(T \leq t) = 1 - (1-p)^t$

Given that we've waited longer than t , what is the probability that we'll have to wait more than s more?

$$P(T > t+s \mid T > t)$$

$$\begin{aligned} &= \frac{P(T > t+s \ \& \ T > t)}{P(T > t)} \leftarrow = \frac{P(T > t+s)}{P(T > t)} = \frac{(1-p)^{t+s}}{(1-p)^t} \\ &= \frac{\cancel{(1-p)^t} (1-p)^s}{\cancel{(1-p)^t}} \\ &= (1-p)^s \\ &= P(T > s) \end{aligned}$$

"memoryless"

Continuous Wait Times

In the real world, most wait times are continuous random variables.

Eg. After a lull, the arrival time of the first customer at a post office.

Eg. The time until a radioactive particle decays.

These wait times are continuous, but have the same defining **memoryless** property:

$$G(t) = 1 - F_T(t) \quad \rightarrow \quad P(T > t+s) = P(T > s) P(T > t)$$

Set $P(T > t) = G(t)$. Thus $G(t+s) = G(t)G(s)$ (*)

Theorem: If $G: \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a ^{cont.} differentiable function satisfying (*), then for some $a \in \mathbb{R}$,

$$G(t) = e^{at}$$

$$\begin{aligned} & e^{a(t+s)} \\ &= e^{at} e^{as} \quad \checkmark \end{aligned}$$

Proof: $\frac{\partial}{\partial s} G(t+s) = \frac{\partial}{\partial s} [G(t)G(s)]$

$$\frac{\partial}{\partial s} G(t+s) \cdot G'(t+s) = G(t)G'(s)$$

Let $s \downarrow 0$.

$$G'(t) = G'(0)G(t)$$

$$G'(t) = aG(t)$$

$\therefore C^2 = C$ $C = 0$ or $\boxed{1}$

$$C e^{a(t+s)} = G(t+s) = G(t)G(s) \Rightarrow G(t) = C e^{at}$$

Now, if $G(t) = P(T > t) = 1 - P(T \leq t)$, we get

$$= C e^{at} - C e^{as} = C^2 e^{a(t+s)}$$

$$F_T(t) = P(T \leq t) = 1 - G(t) = 1 - e^{at} \rightarrow 1 \text{ as } t \rightarrow \infty$$

$\therefore a < 0$ Let $\lambda = -a$

Definition: The **exponential distribution** with parameter $\lambda > 0$

is given by CDF

$$F(t) = \begin{cases} 1 - e^{-\lambda t}, & t \geq 0 \\ 0, & t < 0 \end{cases}$$

$$f(t) = \begin{cases} \lambda e^{-\lambda t} & t \geq 0 \\ 0 & t < 0 \end{cases}$$

Mean and Variance of $T \sim \text{Exp}(\lambda)$

$$\mathbb{E}(T) = \int_0^{\infty} t f_T(t) = \int_0^{\infty} t \cdot \lambda e^{-\lambda t} dt = \text{CALCULUS} = \frac{1}{\lambda}.$$

$$\text{Var}(T) = \frac{1}{\lambda^2} \quad (\text{similar calculation})$$

E.g. The average phone call is 5 minutes in length. What is the probability your next phone call will be longer than 3 minutes?

$$T \sim \text{Exp}(\lambda), \quad \frac{1}{\lambda} = \mathbb{E}(T) = 5$$

$$\therefore P(T > 3) = \int_3^{\infty} \frac{1}{5} e^{-\frac{1}{5}t} dt = \left(-e^{-\frac{t}{5}} \right) \Big|_{t=3}^{t=\infty} = e^{-\frac{3}{5}} = 54.9\%$$

$$P(T > x) = e^{-x/5} = e^{-\lambda x}$$

E.g. On a forest road, cars come by Turtle Rock on average every 30 minutes. Tianyi the Turtle needs 10 minutes to cross the road. What is the probability she can cross safely?

$T =$ arrival time of next car, $T \sim \text{Exp}(\frac{1}{30})$

$$P(T > 10) = e^{-\frac{1}{30} \cdot 10} = e^{-\frac{1}{3}} \approx 71.7\%$$

Just before she starts to cross, Li-Tien the lemur tells her he's been hanging around for over 20 minutes and no cars have come by. Does this change Tianyi's mind about how safe it is to cross?

$$P(T > 20 + 10 \mid T > 20) = P(T > 10) = 71.7\%$$