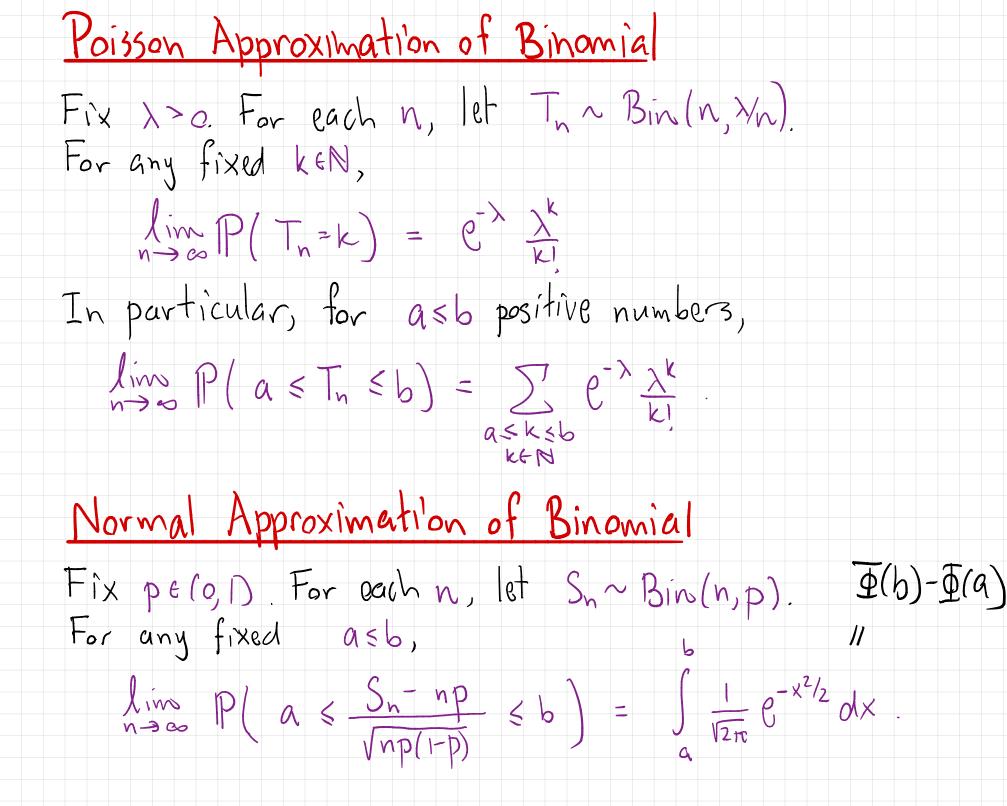# MATH 180A: INTRO TO PROBABILITY
## (FOR DATA SCIENCE)

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

Today:   § 4.1 - 4.3

Next:   § 4.3 - 4.4

Midterm grades released; regrade requests Wed 10/30 8am-11pm

Lab 3 grades released; regrade requests Thu 10/31 8am-11pm

HW4   due   Friday   by 11:59 pm

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

# Poisson Approximation of Binomial

Fix $\lambda > 0$. For each $n$, let $T_n \sim \text{Bin}(n, \lambda/n)$.
For any fixed $k \in \mathbb{N}$,

$$\lim_{n \to \infty} \mathbb{P}(T_n = k) = e^{-\lambda} \frac{\lambda^k}{k!}$$

In particular, for $a \leq b$ positive numbers,

$$\lim_{n \to \infty} \mathbb{P}(a \leq T_n \leq b) = \sum_{\substack{a \leq k \leq b \\ k \in \mathbb{N}}} e^{-\lambda} \frac{\lambda^k}{k!}.$$

# Normal Approximation of Binomial

Fix $p \in (0,1)$. For each $n$, let $S_n \sim \text{Bin}(n,p)$.

$$\Phi(b) - \Phi(a)$$
$$\parallel$$

For any fixed $a \leq b$,

$$\lim_{n \to \infty} \mathbb{P}\left(a \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq b\right) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \, dx.$$

# Question:

A fair 20-sided die is tossed 400 times. We want to calculate the probability that a 13 came up at least 25 times. We should use:
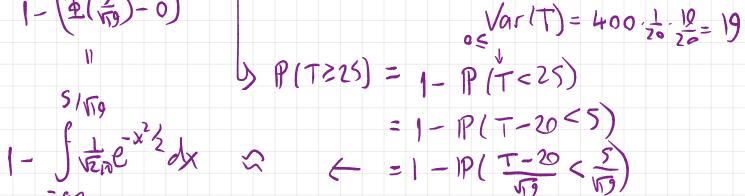
$T \sim Bin(400, \frac{1}{20})$
$Bin(n, \lambda/n)$

(a) Poisson Approximation.
(b) Normal Approximation.
(c) Either.
(d) Neither.

---

$T \sim Bin(400, \frac{1}{20})$

$P(T \geq 25) = 1 - P(T \leq 24)$

$= 1 - \sum_{k=0}^{24} \binom{400}{k} \left(\frac{1}{20}\right)^k \left(\frac{19}{20}\right)^{400-k}$

$= \boxed{15.10\%}$

12.57%

$1 - \left(\Phi\left(\frac{5}{\sqrt{19}}\right) - 0\right)$

$\parallel$

$5/\sqrt{19}$

$1 - \int_{-\infty}^{5/\sqrt{19}} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \quad \lessgtr$

---

## Poisson: $n = 400, \quad \lambda = 20$

$P(T \geq 25) = 1 - P(T < 25)$

$= 1 - \sum_{k=0}^{24} P(T = k)$

$\approx 1 - \sum_{k=0}^{24} e^{-20} \frac{(20)^k}{k!}$

$\doteq \boxed{15.68\%}$

### Normal Approx

$E(T) = 400 \cdot \frac{1}{20} = 20$

$Var(T) = 400 \cdot \frac{1}{20} \cdot \frac{19}{20} = 19$

$0 \leq$

$P(T \geq 25) = 1 - P(T < 25)$

$= 1 - P(T - 20 < 5)$

$= 1 - P\left(\frac{T-20}{\sqrt{19}} < \frac{5}{\sqrt{19}}\right)$

$T \sim Bin(400, \frac{1}{20})$
$Bin(n, \lambda/n)$

Eg. A fair die is rolled 720 times. What is the probability that exactly 113 sixes come up?

$$S = \# \text{ sixes} \sim \text{Bin}(720, \tfrac{1}{6}). \quad P(S=113) = \binom{720}{113}\left(\tfrac{1}{6}\right)^{113}\left(\tfrac{5}{6}\right)^{607} \approx 3.184\%$$

Normal Approximation:

$$\mathbb{E}(S) = 720 \cdot \tfrac{1}{6} = 120 \qquad \text{Var}(S) = 720 \cdot \tfrac{1}{6} \cdot \tfrac{5}{6} = 100$$

$$P(S=113) = P(113 \le S \le 113) = P\left(-0.7 \le \frac{S-120}{\sqrt{100}} \le -0.7\right) \stackrel{\sim}{=} \Phi(-0.7) - \Phi(-0.7)$$
$$= 0.$$

$$= P(112.5 \le S \le 113.5) = P\left(0.75 \le \frac{S-120}{10} \le -0.65\right)$$

$$= \Phi(-0.65) - \Phi(-0.75)$$

$$\doteq 3.122\%$$

"Continuity correction"

$$P(k_1 \le S \le k_2) = P\left(k_1 - \tfrac{1}{2} \le S \le k_2 + \tfrac{1}{2}\right) \approx \Phi\left(\frac{k_2 + \tfrac{1}{2} - np}{\sqrt{np(1-p)}}\right)$$
$$- \Phi\left(\left(k_1 - \tfrac{1}{2} - np\right)/\sqrt{np(1-p)}\right)$$

What is $\mathbb{E}$?

Precise definition: $\mathbb{E}(X) = \sum_{k} k \cdot \mathbb{P}(X = k)$

Intuition: Sample $X$ independently many times, compute the average.

Binomial $S_n \sim Bin(n,p)$ $\qquad \mathbb{E}(S_n) = np$

$\underbrace{X_1 + \cdots + X_n}, \quad X_j \sim Ber(p) \quad \mathbb{E}(X_j) = p. \qquad \mathbb{E}\left(\frac{S_n}{n}\right) = \frac{np}{n} = \boxed{p}$

## Theorem (Law of Large Numbers for Bernoulli Trials)

Let $X_1, X_2, \ldots, X_n$ be independent Bernoulli trials with success probability $p$. Then " $\underbrace{\frac{X_1 + \cdots + X_n}{n}}_{S_n/n} \longrightarrow p$ as $n \to \infty$ "

Precisely: For any $\varepsilon > 0$,

$$\lim_{n \to \infty} \mathbb{P}\left(\left|\frac{S_n}{n} - p\right| < \varepsilon\right) = 1.$$

If $S_n \sim \text{Bin}(n, p)$ (p fixed), and $\varepsilon > 0$,

$$\lim_{n \to \infty} \mathbb{P}\left(\left|\frac{S_n}{n} - p\right| > \varepsilon\right) = 0.$$

**Proof:**

$$\frac{S_n}{n} - p = \frac{S_n - np}{n} \frac{\sqrt{np(1-p)}}{\sqrt{np(1-p)}} = \frac{S_n - np}{\sqrt{np(1-p)}} \cdot \frac{\sqrt{n(p)(1-p)}}{n}$$

$$\downarrow \qquad\qquad \underbrace{\phantom{xx}}$$
$$\mathcal{N}(0,1) \qquad \frac{\sqrt{p(1-p)}}{\sqrt{n}}$$

$$\mathbb{P}\left(\left|\frac{S_n}{n} - p\right| < \varepsilon\right) = \mathbb{P}\left(\underbrace{\left|\frac{S_n - np}{\sqrt{np(1-p)}}\right| < \frac{\varepsilon}{\sqrt{p(1-p)}} \cdot \sqrt{n}}_{\text{large \#.}}\right)$$

Pick your favorite large $R$. Find some $n_0$ s.t.

$$\mathcal{N}(0,1) \qquad\qquad \frac{\varepsilon}{\sqrt{p(1-p)}}\sqrt{n_0} > R.$$
$$\downarrow$$

$$\approx \mathbb{P}(|X| < R) \to 1 \quad \text{as } R \to \infty. \qquad\qquad /\!/\!/$$

## Example

Flip a fair coin $n$ times. How does

$$\lim_{n \to \infty} \mathbb{P}\left(\frac{\# Heads}{n} \geq 50.01\%\right) = 0.$$

behave as $n \to \infty$?

Suppose after 10,000 flips, there are 5,001 Heads. Should we doubt that the coin is really fair?

What if, after 1,000,000 flips, there are 500,100 Heads. Now how confident should we be that the coin is really fair?

$$\mathbb{P}\left(\frac{S_n}{n} \geq \frac{1}{2} + \varepsilon\right) = \mathbb{P}\left(\frac{S_n - \frac{1}{2}n}{\sqrt{n} \cdot \frac{1}{2} \cdot \frac{1}{2}} \geq 2\varepsilon\sqrt{n}\right) \approx \mathbb{P}\left(X \geq 2\varepsilon\sqrt{n}\right)$$

$\uparrow$

$\varepsilon = 0.01$

$\overset{?}{\underset{}{N(0,1)}}$