

## LECTURE 20

- **Readings:** Section 5.4

## Lecture outline

- Review
- Central limit theorem
- Normal approximations
- De Moivre–Laplace binomial approximation

## Review: Convergence in probability

- Sequence of random variables  $Y_1, Y_2, \dots$  **converges in probability to a number**  $a$  when

$$\text{for any } \epsilon > 0, \quad \lim_{n \rightarrow \infty} \mathbf{P}(|Y_n - a| \geq \epsilon) = 0$$

- Examples of  $Y_1, Y_2, \dots$  converging in probability to 0:

–  $Y_n$  continuous uniform over  $[0, 1/n]$

–  $Y_n$  continuous uniform over  $[-1/n, 1/n]$

$$- p_{Y_n}(y) = \begin{cases} 1 - \frac{1}{n}, & y = 0 \\ \frac{1}{n}, & y = 1 \end{cases}$$

$$- p_{Y_n}(y) = \begin{cases} 1 - \frac{1}{n}, & y = 0 \\ \frac{1}{n}, & y = n^2 \end{cases} \quad - p_{Y_n}(y) = \begin{cases} 1 - \frac{1}{n}, & y = e^{-n} \\ \frac{1}{n}, & y = e^n \end{cases}$$

## Central limit theorem

- Let  $Z$  be a standard normal r.v. (zero mean, unit variance),  $X_1, X_2, \dots$  i.i.d. with finite mean  $\mu$  and variance  $\sigma^2$ , and

$$Z_n = \frac{(X_1 + X_2 + \dots + X_n) - n\mu}{\sqrt{n}}.$$

$$\mathbf{E}[Z_n] = 0$$

$$\text{var}(Z_n) = 1$$

For every  $z$ ,

$$\lim_{n \rightarrow \infty} \mathbf{P}(Z_n \leq z) = \mathbf{P}(Z \leq z) = \Phi(z).$$

- Asymptotic equality of CDFs
- Not a statement about PDFs or CDFs
- Standard normal CDF  $\Phi(z)$  is built in to many packages and available in tables
- Proof is not too difficult with transforms

## Normal approximation

- Treat  $Z_n$  as if standard normal
  - Makes  $S_n$  have the  $\mathcal{N}(n\mu, n\sigma^2)$  distribution
- How good is it when  $n$  is “moderate”?
  - Never exactly correct (unless  $X_i$ s are normal)
  - Speed of convergence depends on  $X_i$  distribution
  - Convergence is faster with symmetric distributions
  - Convergence is faster with unimodal distributions

## The pollster’s problem revisited

- $f$ : fraction of population that “...”
- $i$ th (randomly selected) person polled:  $X_i = \begin{cases} 1, & \text{if yes;} \\ 0, & \text{if no.} \end{cases}$
- $M_n = (X_1 + \dots + X_n)/n$  is fraction of “yes” in our sample
- Goal: “95% confidence in being within 1% error”

$$\mathbf{P}(|M_n - f| \geq 0.01) \leq 0.05$$

- Use normal approximation: treat  $M_n$  as normal

## The pollster’s problem revisited (2)

- Event of interest (to be made unlikely):  $|M_n - f| \geq 0.01$

$$\left| \frac{X_1 + \dots + X_n - nf}{n} \right| \geq 0.01$$

$$\left| \frac{X_1 + \dots + X_n - nf}{\sqrt{n}\sigma} \right| \geq \frac{0.01\sqrt{n}}{\sigma}$$

$$\mathbf{P}(|M_n - f| \geq 0.01) \approx \mathbf{P}(|Z| \geq 0.01\sqrt{n}/\sigma)$$

$$\leq \mathbf{P}(|Z| \geq 0.02\sqrt{n})$$

$$= 2\mathbf{P}(Z < -0.02\sqrt{n}) = 2(1 - \Phi(0.02\sqrt{n}))$$

- Pick  $n$  large enough that  $\Phi(0.02\sqrt{n}) \geq 1 - \frac{0.05}{2} = 0.975$ 
  - $\Phi(1.96) \approx 0.975$ , so  $n = 9604$  suffices

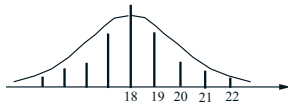
### The pollster's problem revisited (3)

- Why is Chebyshev-inspired pollster so conservative?
- Fix  $p \in (0, 1)$  and let  $X_1, X_2, \dots$  be independent Bernoulli( $p$ )
- $S_n = X_1 + \dots + X_n$  is binomial( $n, p$ )
  - mean  $np$ , variance  $np(1-p)$
- CDF of  $\frac{S_n - np}{\sqrt{np(1-p)}}$  approaches standard normal CDF

### Approximating a binomial: Example

- $n = 36, p = 0.5$ ; find  $P(S_n \leq 21)$
- Exact answer:
 
$$P(S_n \leq 21) = \sum_{k=0}^{21} \binom{36}{k} \left(\frac{1}{2}\right)^{36} \approx 0.8785$$
- Normal approximation:
 
$$P(S_n \leq 21) = P\left(\frac{S_n - 36 \cdot 0.5}{\sqrt{36 \cdot 0.5 \cdot 0.5}} \leq \frac{21 - 36 \cdot 0.5}{\sqrt{36 \cdot 0.5 \cdot 0.5}}\right) \approx \Phi(1) \approx 0.8413$$
- But also:
 
$$\begin{aligned} P(S_n \leq 21) &= P(S_n < 22) \\ &= P\left(\frac{S_n - 36 \cdot 0.5}{\sqrt{36 \cdot 0.5 \cdot 0.5}} < \frac{22 - 36 \cdot 0.5}{\sqrt{36 \cdot 0.5 \cdot 0.5}}\right) \approx \Phi\left(\frac{4}{3}\right) \approx 0.9088 \end{aligned}$$

### De Moivre–Laplace approximation



- Using the midpoint of discrete values is called the De Moivre–Laplace approximation:

$$\begin{aligned} P(S_n \leq 21) &= P(S_n \leq 21.5) \\ &= P\left(\frac{S_n - 36 \cdot 0.5}{\sqrt{36 \cdot 0.5 \cdot 0.5}} \leq \frac{21.5 - 36 \cdot 0.5}{\sqrt{36 \cdot 0.5 \cdot 0.5}}\right) \approx \Phi\left(\frac{7}{6}\right) \approx 0.8783 \end{aligned}$$

### De Moivre–Laplace approximation for binomial PMF

- “1/2 correction” enables a reasonable PMF approximation:
 
$$P(S_n = 19) = P(18.5 \leq S_n \leq 19.5)$$

$$18.5 \leq S_n \leq 19.5 \iff \frac{18.5 - 18}{3} \leq \frac{S_n - 18}{3} \leq \frac{19.5 - 18}{3}$$

$$\iff \frac{1}{6} \leq Z_n \leq \frac{1}{2}$$

$$\begin{aligned} P(S_n = 19) &\approx P\left(\frac{1}{6} \leq Z \leq \frac{1}{2}\right) \\ &= \Phi\left(\frac{1}{2}\right) - \Phi\left(\frac{1}{6}\right) \approx 0.6915 - 0.5675 = 0.124 \end{aligned}$$
- Exact answer:  $\binom{36}{19} \left(\frac{1}{2}\right)^{36} \approx 0.1251$

### Poisson vs. normal approximations of the binomial

- Poisson arrivals during unit interval equals: sum of  $n$  (independent) Poisson arrivals during  $n$  intervals of length  $1/n$ 
  - Let  $n \rightarrow \infty$ , apply CLT (?)
  - Poisson  $\stackrel{?}{=}$  normal
- Binomial( $n, p$ )
  - $p$  fixed,  $n \rightarrow \infty$ : normal
  - $np$  fixed,  $n \rightarrow \infty, p \rightarrow 0$ : Poisson
- Single instance: use Poisson when  $p$  is very small
  - $p = 1/100, n = 100$ : use Poisson approximation
  - $p = 1/10, n = 500$ : use normal approximation

### Application of the CLT

- Easy to apply: Uses only means and variances
  - Useful computational shortcut, even if the distribution of  $S_n$  is known
  - Justifies noise models involving normal random variables
    - Accumulation of many small, approx.-independent effects
  - Normal distribution has many extremal properties
    - maximizing entropy, minimizing Fisher information, Heisenberg uncertainty principle, ...
- “Physicists believe that the Gaussian law has been proved in mathematics while mathematicians think that it was experimentally established in physics.”  
— Henri Poincaré [*Science et Hypothèse*, 1904]