

Central Limit Thm, Normal Approximations

Engineering Statistics
Section 5.4

Josh Engwer

TTU

23 March 2016

PART I:
CENTRAL LIMIT THEOREM

Expected Value & Variance of a Sum of iid rv's

Proposition

Let X_1, \dots, X_n be a random sample from a population and $c_1, \dots, c_n \neq 0$.
Then:

- $\mathbb{E}[c_1X_1 + \dots + c_nX_n] = c_1\mathbb{E}[X_1] + \dots + c_n\mathbb{E}[X_n]$
- $\mathbb{V}[c_1X_1 + \dots + c_nX_n] = c_1^2\mathbb{V}[X_1] + \dots + c_n^2\mathbb{V}[X_n]$

PROOF: CASE I ($n = 2$): $X_1, X_2 \stackrel{iid}{\sim}$ (discrete population)

$$\mathbb{E}[c_1X_1 + c_2X_2] := \sum_{(j,k) \in \text{Supp}(X_1, X_2)} (c_1j + c_2k) \cdot p_{(X_1, X_2)}(j, k)$$

$$\stackrel{iid}{=} \sum_{j \in \text{Supp}(X_1)} \sum_{k \in \text{Supp}(X_2)} (c_1j + c_2k) \cdot p_{X_1}(j) \cdot p_{X_2}(k)$$

$$= c_1 \sum_{j \in \text{Supp}(X_1)} j \cdot p_{X_1}(j) + c_2 \sum_{k \in \text{Supp}(X_2)} k \cdot p_{X_2}(k)$$

$$:= c_1\mathbb{E}[X_1] + c_2\mathbb{E}[X_2]$$

Expected Value & Variance of a Sum of iid rv's

Proposition

Let X_1, \dots, X_n be a random sample from a population and $c_1, \dots, c_n \neq 0$.
Then:

- $\mathbb{E}[c_1X_1 + \dots + c_nX_n] = c_1\mathbb{E}[X_1] + \dots + c_n\mathbb{E}[X_n]$
- $\mathbb{V}[c_1X_1 + \dots + c_nX_n] = c_1^2\mathbb{V}[X_1] + \dots + c_n^2\mathbb{V}[X_n]$ (requires **independence**)

PROOF: CASE II ($n = 2$): $X_1, X_2 \stackrel{iid}{\sim}$ (continuous population)

$$\begin{aligned}\mathbb{E}[c_1X_1 + c_2X_2] &:= \iint_{\text{Supp}(X_1, X_2)} (c_1x_1 + c_2x_2) \cdot f_{(X_1, X_2)}(x_1, x_2) dx_1 dx_2 \\ &\stackrel{iid}{=} \int_{\text{Supp}(X_1)} \int_{\text{Supp}(X_2)} (c_1x_1 + c_2x_2) \cdot f_{X_1}(x_1) \cdot f_{X_2}(x_2) dx_1 dx_2 \\ &= c_1 \int_{\text{Supp}(X_1)} x_1 \cdot f_{X_1}(x_1) dx_1 + c_2 \int_{\text{Supp}(X_2)} x_2 \cdot f_{X_2}(x_2) dx_2 \\ &:= c_1\mathbb{E}[X_1] + c_2\mathbb{E}[X_2]\end{aligned}$$

Properties of Sample Mean

Proposition

Let X_1, \dots, X_n be a random sample from a distribution with mean μ and variance σ^2 . Then:

- $\mu_{\bar{X}} = \mathbb{E}[\bar{X}] = \mu$
- $\sigma_{\bar{X}}^2 = \mathbb{V}[\bar{X}] = \sigma^2/n$
- $\sigma_{\bar{X}} = \sigma/\sqrt{n}$

PROOF:

$$\begin{aligned}\mathbb{E}[\bar{X}] &= \mathbb{E}\left[\frac{1}{n}(X_1 + \dots + X_n)\right] = \frac{1}{n} [\mathbb{E}[X_1] + \dots + \mathbb{E}[X_n]] \\ &= \frac{1}{n} [\mu + \dots + \mu] = \frac{1}{n} [n\mu] = \mu\end{aligned}$$

$$\begin{aligned}\mathbb{V}[\bar{X}] &= \mathbb{V}\left[\frac{1}{n}(X_1 + \dots + X_n)\right] = \frac{1}{n^2} [\mathbb{V}[X_1] + \dots + \mathbb{V}[X_n]] \\ &= \frac{1}{n^2} [\sigma^2 + \dots + \sigma^2] = \frac{1}{n^2} [n\sigma^2] = \sigma^2/n\end{aligned}$$

$$\sigma_{\bar{X}} = \sqrt{\mathbb{V}[\bar{X}]} = \sqrt{\sigma^2/n} = \sigma/\sqrt{n} \quad \square$$

Properties of Sample Total

Proposition

Let X_1, \dots, X_n be a random sample from a distribution with mean μ and variance σ^2 . Then:

- $\mu_{X_1 + \dots + X_n} = \mathbb{E}[X_1 + \dots + X_n] = n\mu$
- $\sigma_{X_1 + \dots + X_n}^2 = \mathbb{V}[X_1 + \dots + X_n] = n\sigma^2$
- $\sigma_{X_1 + \dots + X_n} = \sigma\sqrt{n}$

PROOF:

$$\mathbb{E}[X_1 + \dots + X_n] = \mathbb{E}[n\bar{X}] = n\mathbb{E}[\bar{X}] = n \cdot \mu = n\mu$$

$$\mathbb{V}[X_1 + \dots + X_n] = \mathbb{V}[n\bar{X}] = n^2\mathbb{V}[\bar{X}] = n^2[\sigma^2/n] = n\sigma^2$$

$$\sigma_{X_1 + \dots + X_n} = \sqrt{\mathbb{V}[X_1 + \dots + X_n]} = \sqrt{n\sigma^2} = \sigma\sqrt{n} \quad \square$$

Sample Mean & Sample Total of a Random Sample from a Normal Population

The proceeding properties of sample means & sample variances can be applied to particular population distributions:

Proposition

Let random sample $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2)$. Then:

- For any sample size $n > 1$, $\bar{X} \sim \text{Normal}(\mu, \sigma^2/n)$
- For any sample size $n > 1$, $X_1 + \dots + X_n \sim \text{Normal}(n\mu, n\sigma^2)$

PROOF: $X \sim \text{Normal}(\mu, \sigma^2) \implies \mathbb{E}[X] = \mu, \mathbb{V}[X] = \sigma^2$

$$\begin{aligned}\mathbb{E}[\bar{X}] &= \mathbb{E}\left[\frac{1}{n}(X_1 + \dots + X_n)\right] = \frac{1}{n}[\mathbb{E}[X_1] + \dots + \mathbb{E}[X_n]] \\ &= \frac{1}{n}[\mu + \dots + \mu] = \frac{1}{n}[n\mu] = \mu\end{aligned}$$

$$\begin{aligned}\mathbb{V}[\bar{X}] &= \mathbb{V}\left[\frac{1}{n}(X_1 + \dots + X_n)\right] = \frac{1}{n^2}[\mathbb{V}[X_1] + \dots + \mathbb{V}[X_n]] \\ &= \frac{1}{n^2}[\sigma^2 + \dots + \sigma^2] = \frac{1}{n^2}[n\sigma^2] = \sigma^2/n\end{aligned}$$

$\therefore \bar{X} \sim \text{Normal}(\mu, \sigma^2/n)$ \square

Sample Mean & Sample Total of a Random Sample from a Normal Population

The proceeding properties of sample means & sample variances can be applied to particular population distributions:

Proposition

Let random sample $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2)$. Then:

- For any sample size $n > 1$, $\bar{X} \sim \text{Normal}(\mu, \sigma^2/n)$
- For any sample size $n > 1$, $X_1 + \dots + X_n \sim \text{Normal}(n\mu, n\sigma^2)$

PROOF: $X \sim \text{Normal}(\mu, \sigma^2) \implies \mathbb{E}[X] = \mu, \mathbb{V}[X] = \sigma^2$

$$\mathbb{E}[X_1 + \dots + X_n] = \mathbb{E}[n\bar{X}] = n\mathbb{E}[\bar{X}] = n \cdot \mu = n\mu$$

$$\mathbb{V}[X_1 + \dots + X_n] = \mathbb{V}[n\bar{X}] = n^2\mathbb{V}[\bar{X}] = n^2[\sigma^2/n] = n\sigma^2$$

$$\therefore X_1 + \dots + X_n \sim \text{Normal}(n\mu, n\sigma^2) \quad \square$$

Central Limit Theorem (CLT)

The **Central Limit Thm** is considered a fundamental theorem of Statistics:

Theorem

(Central Limit Theorem)

Let X_1, \dots, X_n be a random sample from a non-normal distribution with mean μ and variance σ^2 . Then:

- The sample mean \bar{X} is approximately normal as follows:

$$\bar{X} \overset{\text{approx}}{\sim} \text{Normal}(\mu, \sigma^2/n)$$

- The sample total $X_1 + \dots + X_n$ is approximately normal as follows:

$$X_1 + \dots + X_n \overset{\text{approx}}{\sim} \text{Normal}(n\mu, n\sigma^2)$$

Requirement for this normal approximation to be valid: $n > 30$

The larger the sample size n , the better the approximation.

PROOF: Beyond the scope of this course.

PART II:

NORMAL APPROXIMATION TO THE BINOMIAL

The Sum of several iid Bernoulli(p) rv's

Proposition

Let random sample $X_1, \dots, X_m \stackrel{iid}{\sim} \text{Bernoulli}(p)$. Then:

- For any sample size $n > 1$, $X_1 + \dots + X_n \sim \text{Binomial}(n, p)$

PROOF: $X \sim \text{Bernoulli}(p) \implies \mathbb{E}[X] = p, \mathbb{V}[X] = pq \quad (q = 1 - p)$

$$\mathbb{E}[X_1 + \dots + X_n] = \mathbb{E}[X_1] + \dots + \mathbb{E}[X_n] \stackrel{iid}{=} n \cdot \mathbb{E}[X_1] = n \cdot p = np$$

$$\mathbb{V}[X_1 + \dots + X_n] \stackrel{iid}{=} \mathbb{V}[X_1] + \dots + \mathbb{V}[X_n] \stackrel{iid}{=} n \cdot \mathbb{V}[X_1] = n \cdot pq = npq$$

$\therefore X_1 + \dots + X_n \sim \text{Binomial}(n, p) \quad \square$

Normal Approximation to the Binomial

Corollary

Let $X \sim \text{Binomial}(n, p)$. Then:

$$X \stackrel{\text{approx}}{\sim} \text{Normal}(\mu = np, \sigma^2 = npq) \quad (\text{where } q = 1 - p)$$

Requirement for this normal approximation to be valid: $\min\{np, nq\} \geq 10$
i.e. It's required that both $np \geq 10$ and $nq \geq 10$.

NOTE: Remember that normal distributions are symmetric.
If $\min\{np, nq\} < 10$, then the binomial distribution is too skewed.

PROOF:

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p) \implies \mu_{X_1} = \mathbb{E}[X_1] = p$ and $\sigma_{X_1}^2 = \mathbb{V}[X_1] = pq$.

Then $X := X_1 + \dots + X_n \sim \text{Binomial}(n, p)$.

Moreover, the CLT asserts that $X_1 + \dots + X_n \stackrel{\text{approx}}{\sim} \text{Normal}(\mu = n\mu_{X_1}, \sigma^2 = n\sigma_{X_1}^2)$

$\therefore X \sim \text{Binomial}(n, p) \implies X \stackrel{\text{approx}}{\sim} \text{Normal}(\mu = np, \sigma^2 = npq) \quad \square$

Normal Approximation to Binomial Probability

Corollary

Let $X \sim \text{Binomial}(n, p)$. Then:

$$\mathbb{P}(X \leq x) = \text{Bi}(x; n, p) \approx \Phi\left(\frac{x + 0.5 - np}{\sqrt{npq}}\right) \quad (\text{where } q = 1 - p)$$

Requirement for this normal approximation to be valid: $\min\{np, nq\} \geq 10$

NOTE: The **continuity correction term** "+ 0.5" improves the approximation.

PROOF: Assume that the requirement $\min\{np, nq\} \geq 10$ is satisfied. Then:

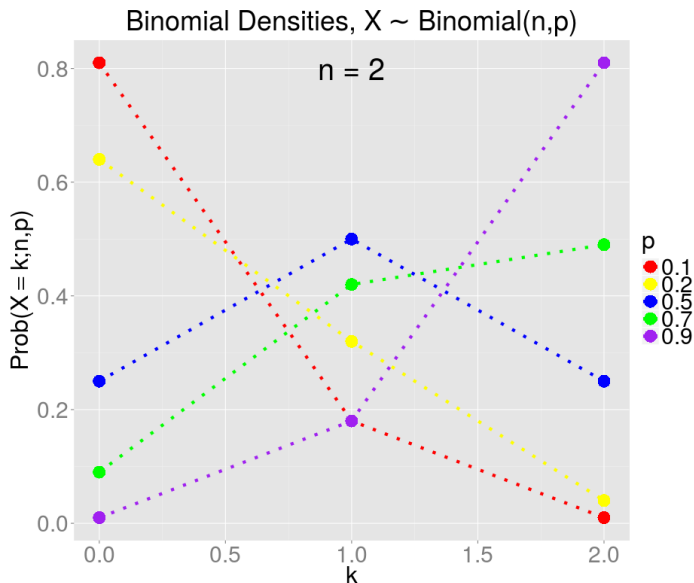
$$X \sim \text{Binomial}(n, p) \implies X \overset{\text{approx}}{\sim} \text{Normal}(\mu = np, \sigma^2 = npq)$$

$$\implies \text{Bi}(x; n, p) = \mathbb{P}(X \leq x) \approx \mathbb{P}\left(Z \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - np}{\sqrt{npq}}\right)$$

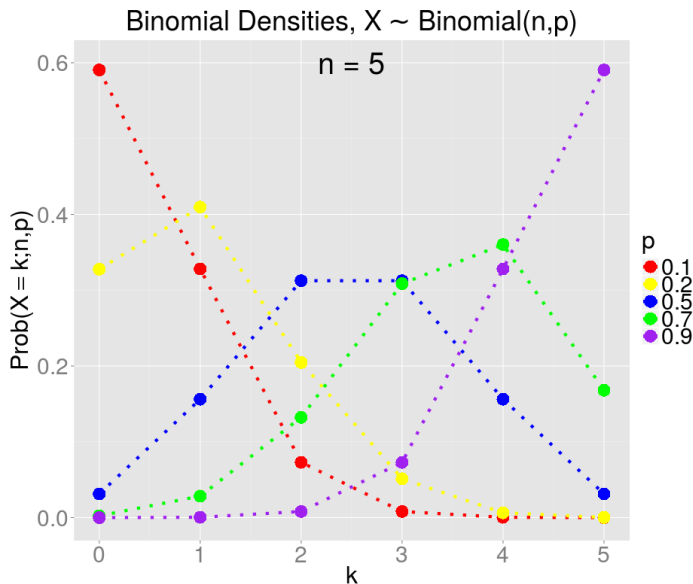
Add continuity correction "+ 0.5" in numerator of Φ to better the approximation.

□

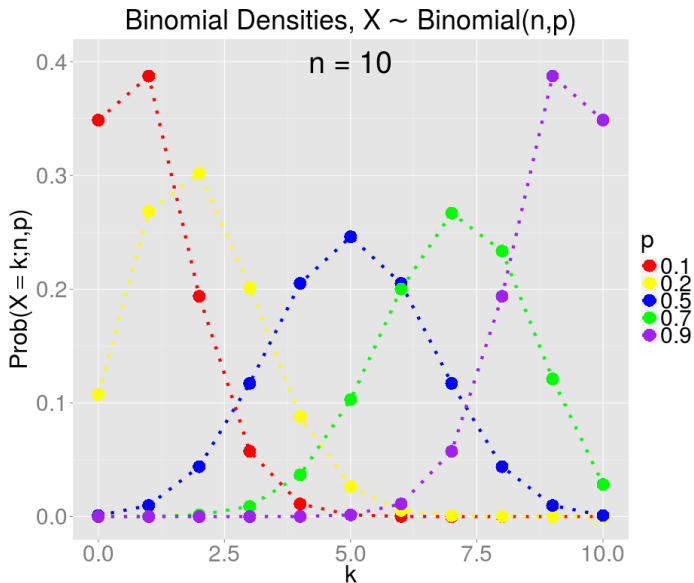
Binomial Density Plots (pmf's) for Sample Size $n = 2$



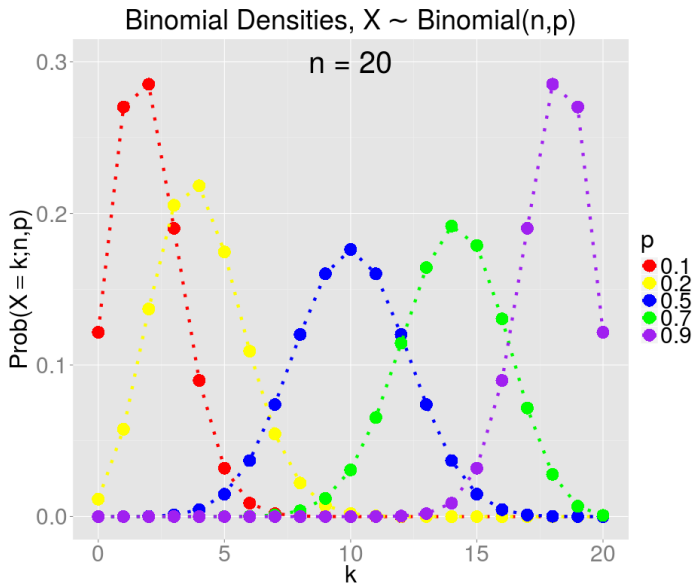
Binomial Density Plots (pmf's) for Sample Size $n = 5$



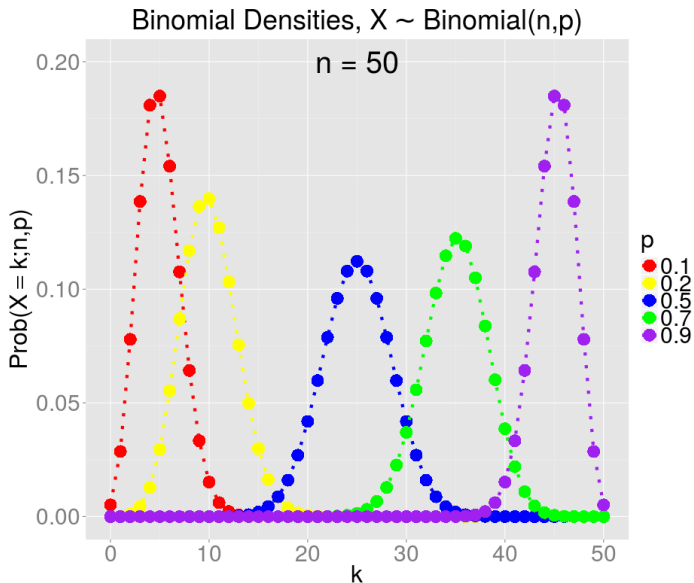
Binomial Density Plots (pmf's) for Sample Size $n = 10$



Binomial Density Plots (pmf's) for Sample Size $n = 20$



Binomial Density Plots (pmf's) for Sample Size $n = 50$



PART III:

NORMAL APPROXIMATION TO THE POISSON

The Sum of several iid $\text{Poisson}(\lambda)$ rv's

The proceeding properties of sample totals can be applied to particular population distributions:

Proposition

Let random sample $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$. Then:

- For any sample size $n > 1$, $X_1 + \dots + X_n \sim \text{Poisson}(n\lambda)$

PROOF: $X \sim \text{Poisson}(\lambda) \implies \mathbb{E}[X] = \lambda, \mathbb{V}[X] = \lambda$

$$\mathbb{E}[X_1 + \dots + X_n] = \mathbb{E}[X_1] + \dots + \mathbb{E}[X_n] \stackrel{iid}{=} n \cdot \mathbb{E}[X_1] = n \cdot \lambda = n\lambda$$

$$\mathbb{V}[X_1 + \dots + X_n] \stackrel{iid}{=} \mathbb{V}[X_1] + \dots + \mathbb{V}[X_n] \stackrel{iid}{=} n \cdot \mathbb{V}[X_1] = n \cdot \lambda = n\lambda$$

$\therefore X_1 + \dots + X_n \sim \text{Poisson}(n\lambda)$ \square

Normal Approximation to the Poisson

Corollary

Let $X \sim \text{Poisson}(\lambda)$. Then:

$$X \overset{\text{approx}}{\sim} \text{Normal}(\mu = \lambda, \sigma^2 = \lambda)$$

Requirement for this normal approximation to be valid: $\lambda > 20$

NOTE: Remember that normal distributions are symmetric.
If $\lambda \leq 20$, then the binomial distribution is too skewed.

PROOF:

Let $X_1, \dots, X_n \sim \text{Poisson}(\lambda/n) \implies \mu_{X_1} = \mathbb{E}[X_1] = \lambda/n$ and $\sigma_{X_1}^2 = \mathbb{V}[X_1] = \lambda/n$.

Then $X := X_1 + \dots + X_n \sim \text{Poisson}(\lambda)$.

Moreover, the CLT asserts that $X_1 + \dots + X_n \overset{\text{approx}}{\sim} \text{Normal}(\mu = n\mu_{X_1}, \sigma^2 = n\sigma_{X_1}^2)$

$\therefore X \sim \text{Poisson}(\lambda) \implies X \overset{\text{approx}}{\sim} \text{Normal}(\mu = \lambda, \sigma^2 = \lambda) \quad \square$

Normal Approximation to Poisson Probability

Corollary

Let $X \sim \text{Poisson}(\lambda)$. Then:

$$\mathbb{P}(X \leq x) = \text{Pois}(x; \lambda) \approx \Phi\left(\frac{x + 0.5 - \lambda}{\sqrt{\lambda}}\right)$$

Requirement for this normal approximation to be valid: $\lambda > 20$

NOTE: The **continuity correction term** "+ 0.5" improves the approximation.

PROOF: Assume that the requirement $\lambda > 20$ is satisfied. Then:

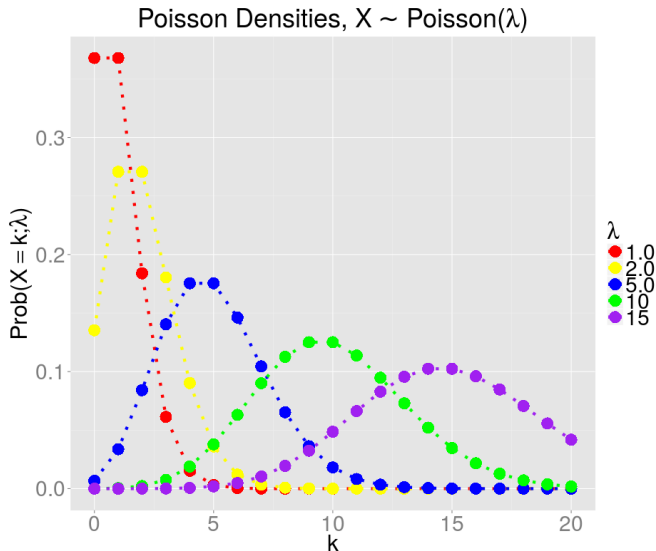
$$X \sim \text{Poisson}(\lambda) \implies X \overset{\text{approx}}{\sim} \text{Normal}(\mu = \lambda, \sigma^2 = \lambda)$$

$$\implies \text{Pois}(x; \lambda) = \mathbb{P}(X \leq x) \approx \mathbb{P}\left(Z \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \lambda}{\sqrt{\lambda}}\right)$$

Add continuity correction "+ 0.5" in numerator of Φ to better the approximation.

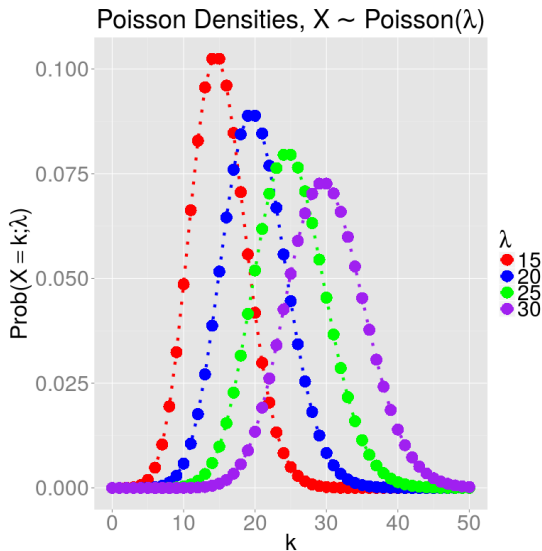


Poisson Density Plots (pmf's)



Notice the Poisson density curves for $\lambda = 1, 2, 5, 10, 15$ are skewed.

Poisson Density Plots (pmf's)



Notice the Poisson density curves for $\lambda \geq 20$ are nearly symmetric.

Textbook Logistics for Section 5.4

- Difference(s) in Notation:

CONCEPT	TEXTBOOK NOTATION	SLIDES/OUTLINE NOTATION
Probability of Event	$P(A)$	$\mathbb{P}(A)$
Support of a r.v.	"All possible values of X "	$\text{Supp}(X)$
pmf of a r.v.	$p_X(x)$	$p_X(k)$
Expected Value of r.v.	$E(X)$	$\mathbb{E}[X]$
Variance of r.v.	$V(X)$	$\mathbb{V}[X]$
Sample Total	T_o	$\sum X_k$
pmf of Sample Mean	$p_{\bar{X}}(\bar{x})$	$p_{\bar{X}}(k)$
pmf of Sample Variance	$p_{S^2}(s^2)$	$p_{S^2}(k)$

- Ignore Lognormal Approximation (bottom of pg 236)
 - The Lognormal distribution was part of Section 4.5 that was skipped

Fin.