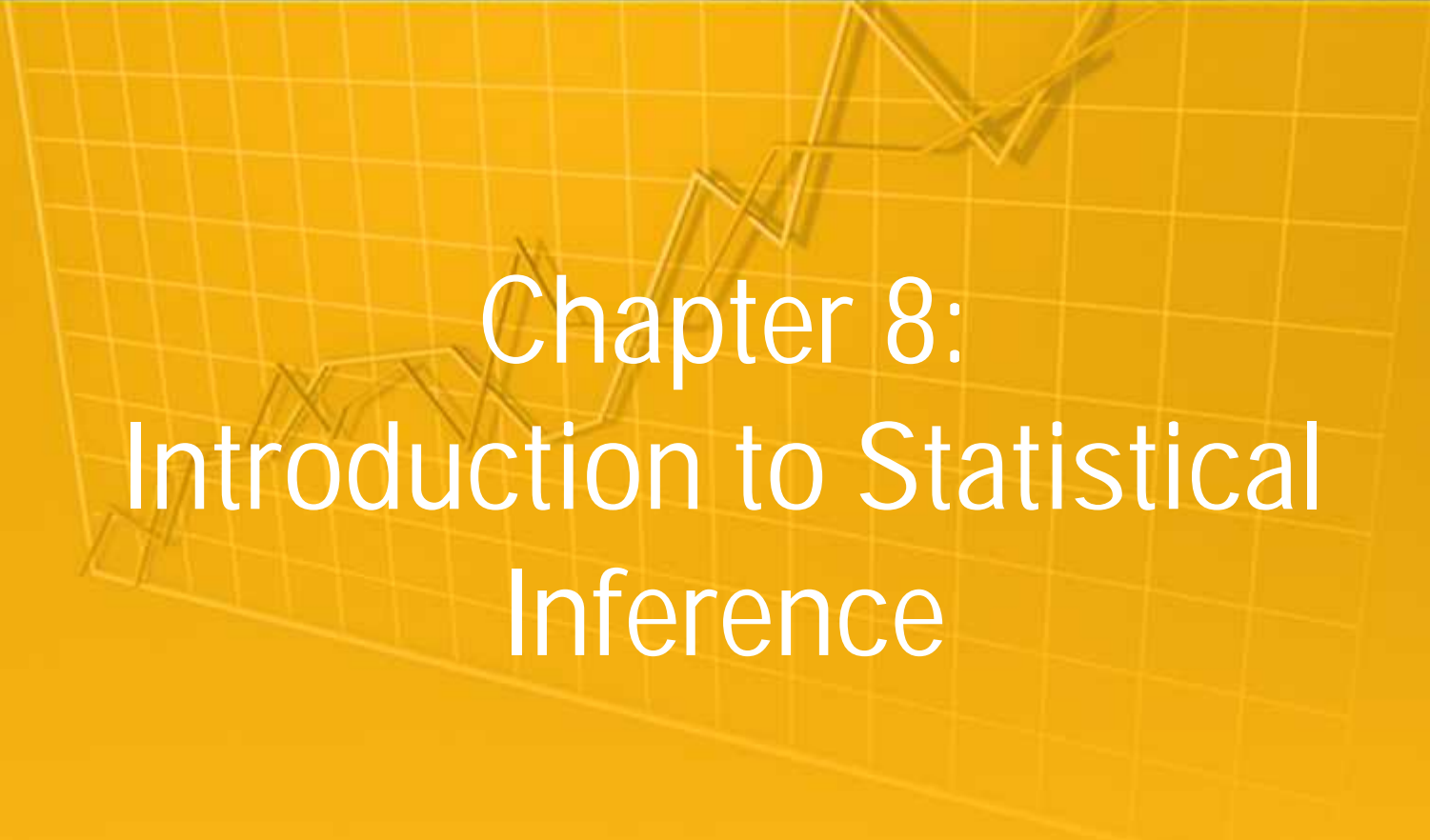


Basic Biostatistics

Statistics for Public Health Practice

B. Burt Gerstman



Chapter 8: Introduction to Statistical Inference

In Chapter 8:

8.1 Concepts

8.2 Sampling Behavior of a Mean

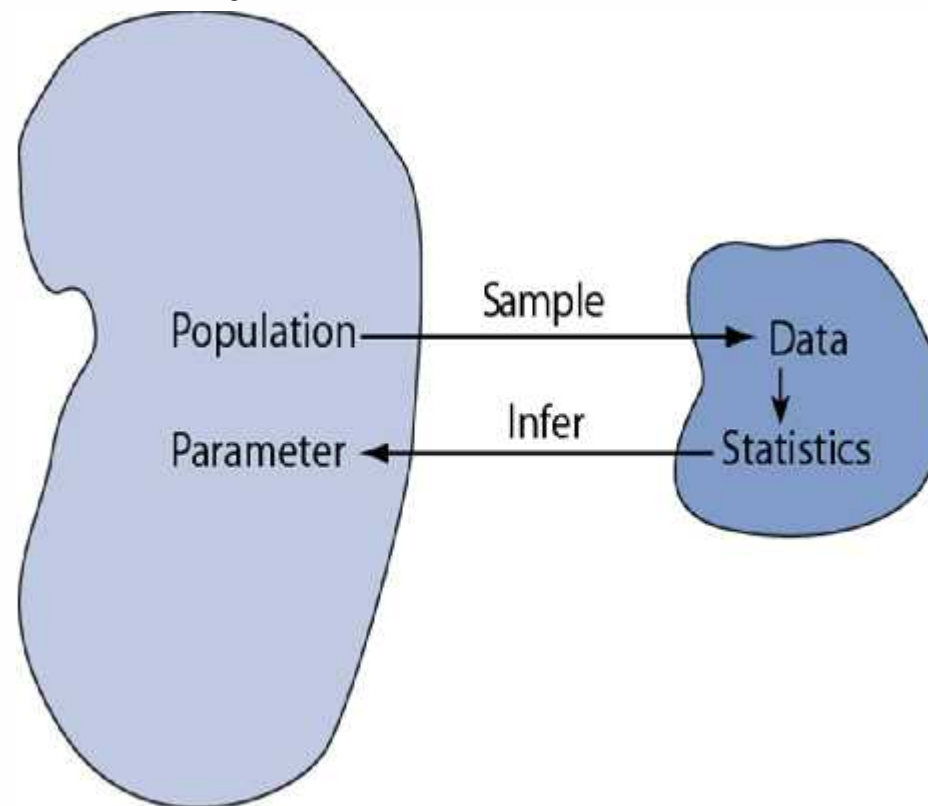
8.3 Sampling Behavior of a Count and Proportion

§8.1: Concepts

Statistical inference is the act of generalizing from a **sample** to a **population** with calculated degree of certainty.

We want to learn about population *parameters*

...



...but we can only calculate *sample statistics*

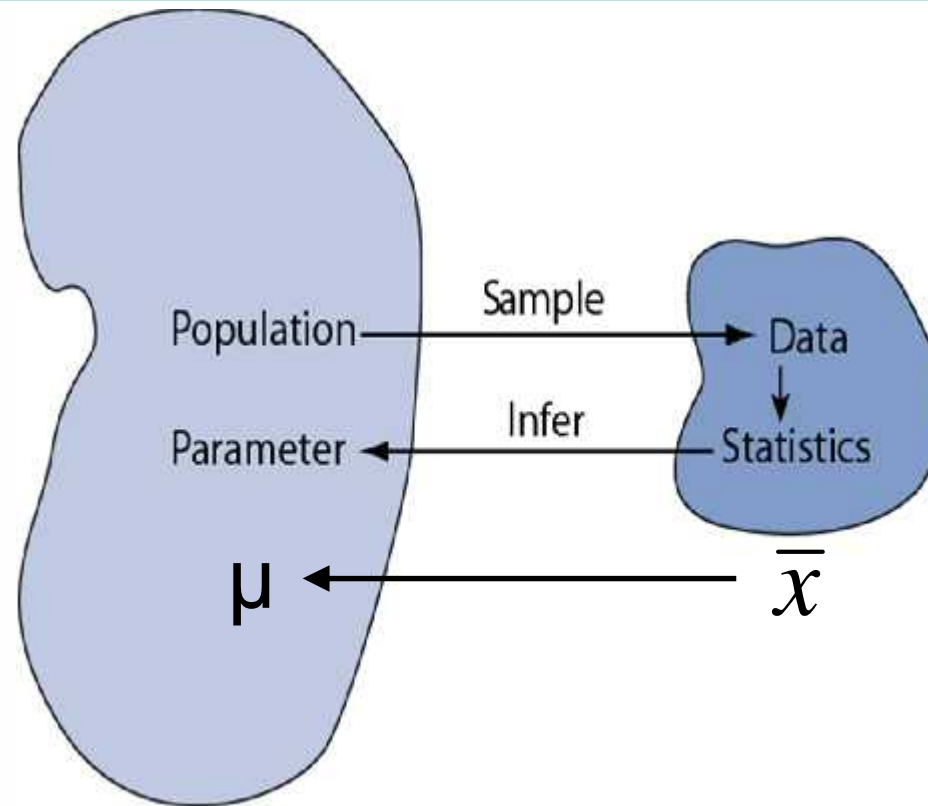
Parameters and Statistics

It is essential that we draw distinctions between parameters and statistics

	Parameters	Statistics
Source	Population	Sample
Calculated?	No	Yes
Constants?	Yes	No
Examples	μ, σ, ρ	\bar{x}, s, \hat{p}

Parameters and Statistics

We are going to illustrate inferential concept by considering how well a given sample mean “x-bar” reflects an underlying population mean μ



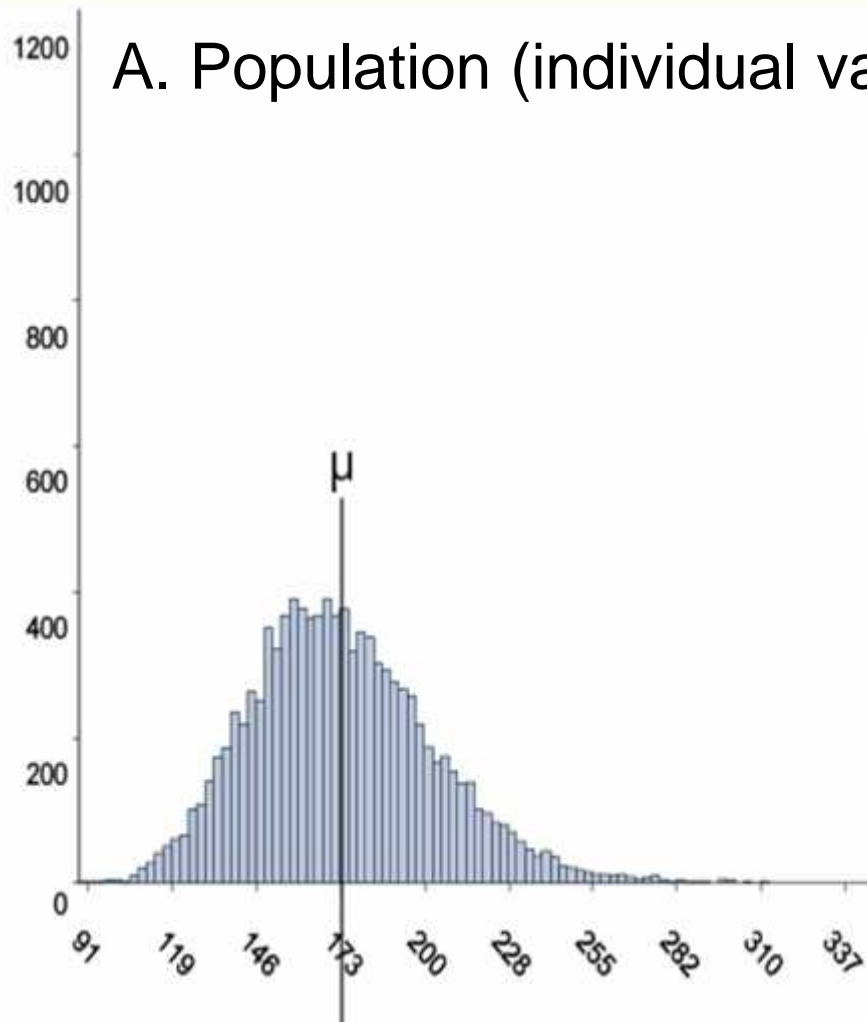
Precision and reliability

- How **precisely** does a given sample mean (\bar{x}) reflect underlying population mean (μ)? How reliable are our inferences?
- To answer these questions, we consider a **simulation experiment** in which we take all possible samples of size n taken from the population

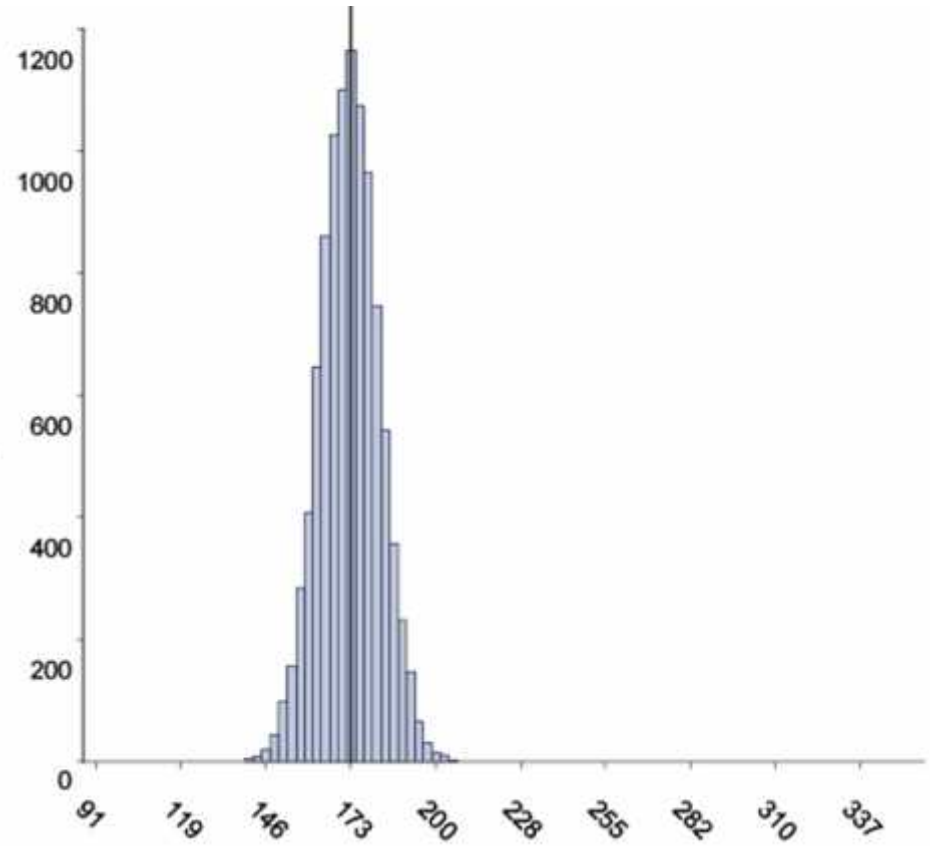
Simulation Experiment

- Population (Figure A, next slide)
 $N = 10,000$
Lognormal shape (positive skew)
 $\mu = 173$
 $\sigma = 30$
- Take repeated SRSs, each of $n = 10$
- Calculate \bar{x} in each sample
- Plot \bar{x} -bars (Figure B , next slide)

A. Population (individual values)

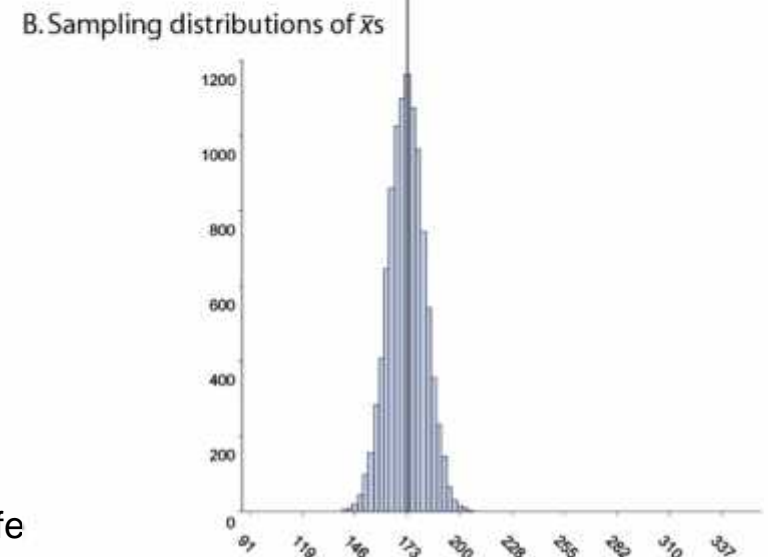
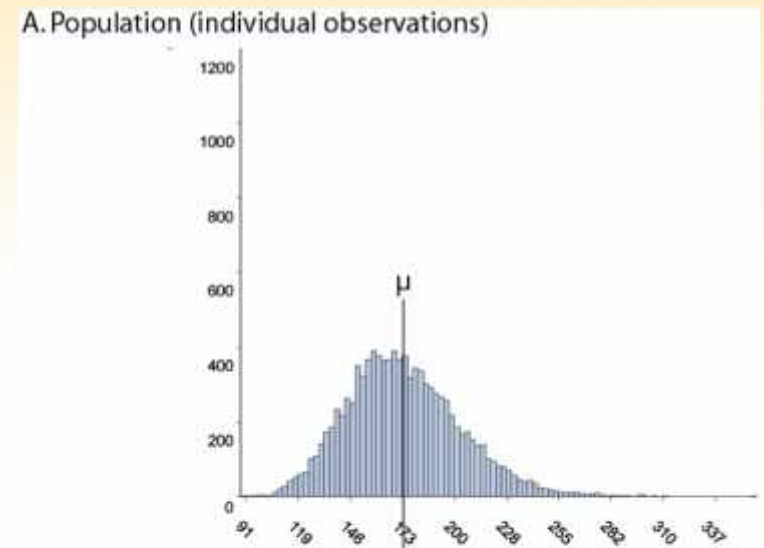


B. Sampling distribution of x-bars



Simulation Experiment Results

1. Distribution B is more Normal than distribution A $\tilde{\sigma}$ **Central Limit Theorem**
2. Both distributions centered on $\mu \tilde{\sigma}$ **\bar{x} is unbiased estimator of μ**
3. Distribution B is skinnier than distribution A $\tilde{\sigma}$ related to “**square root law**”



Reiteration of Key Findings

- **Finding 1 (central limit theorem):** the sampling distribution of \bar{x} tends toward Normality even when the population is not Normal (esp. strong in large samples).
- **Finding 2 (unbiasedness):** the expected value of \bar{x} is μ
- **Finding 3** is related to the **square root law**, which says:

$$\dagger_{\bar{x}} = \frac{\dagger}{\sqrt{n}}$$

Standard Deviation of the Mean

- The standard deviation of the sampling distribution of the mean has a *special name*: **standard error of the mean** (denoted $SE_{\bar{x}}$ or $SE_{\bar{x}}$)
- The square root law says:

$$SE_{\bar{x}} \equiv SE_{\bar{x}} = \frac{s}{\sqrt{n}}$$

Square Root Law

Example: $\sigma = 15$

$$\text{For } n = 1 \Rightarrow SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{1}} = 15$$

$$\text{For } n = 4 \Rightarrow SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{4}} = 7.5$$

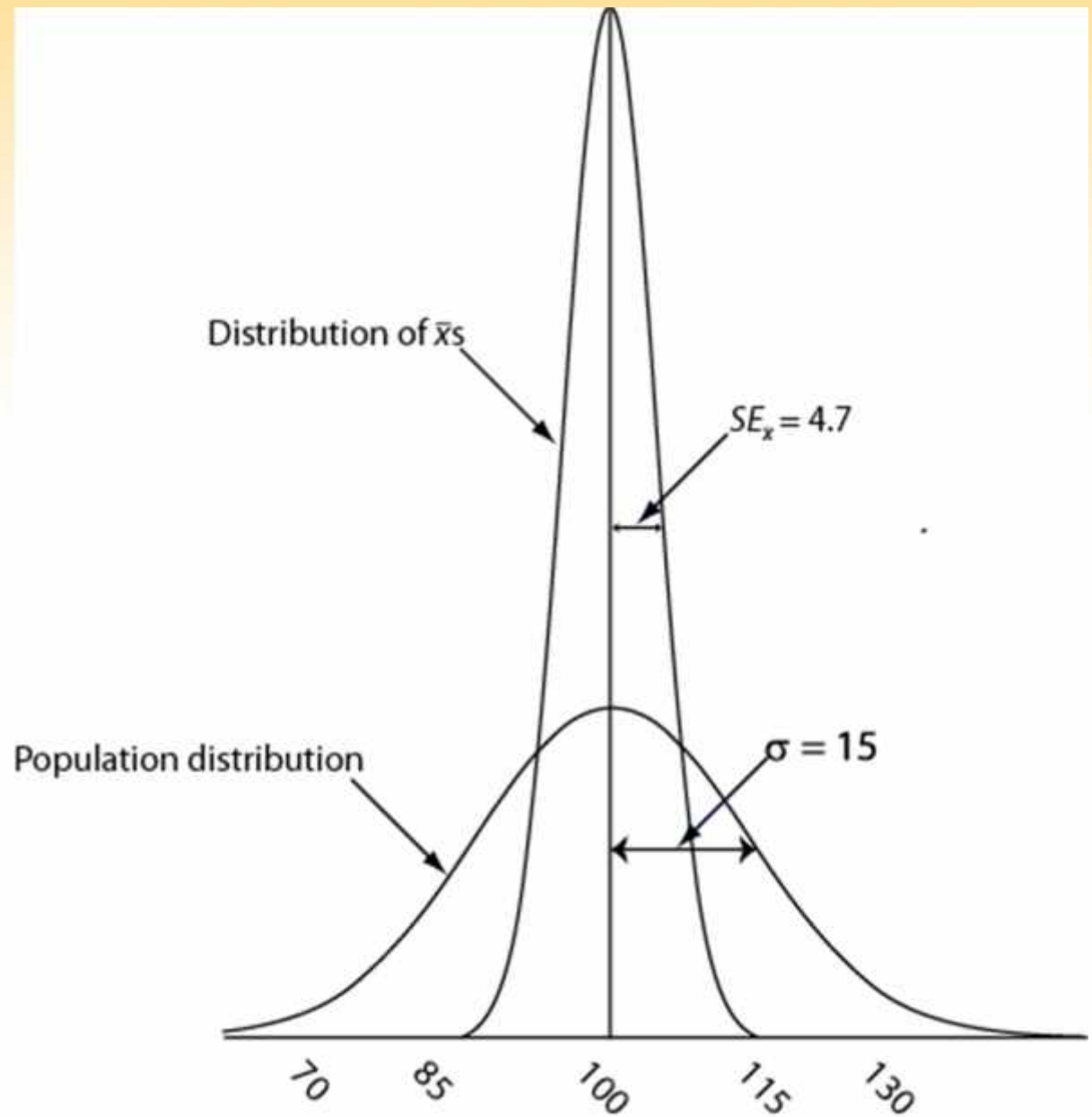
$$\text{For } n = 16 \Rightarrow SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{16}} = 3.75$$

Quadrupling the sample size cuts *the standard error of the mean* in half

Putting it together: $\bar{x} \sim N(\mu, SE)$

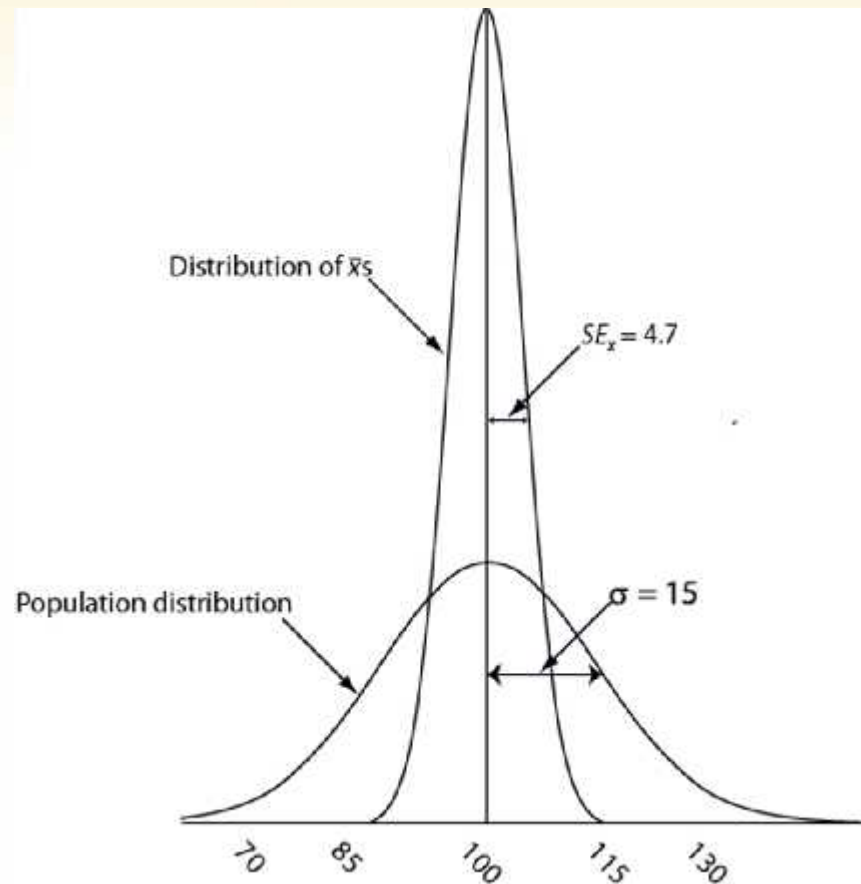
- The sampling distribution of \bar{x} tends to be Normal with mean μ and $SE_{\bar{x}} = \sigma / \sqrt{n}$
- Example: Let X represent Weschler Adult Intelligence Scores; $X \sim N(100, 15)$.
 - Take an SRS of $n = 10$
 - $SE_{\bar{x}} = \sigma / \sqrt{n} = 15 / \sqrt{10} = 4.7$
 - Thus, $\bar{x} \sim N(100, 4.7)$

Individual
WAIS
(population)
and mean
WAIS when
 $n = 10$



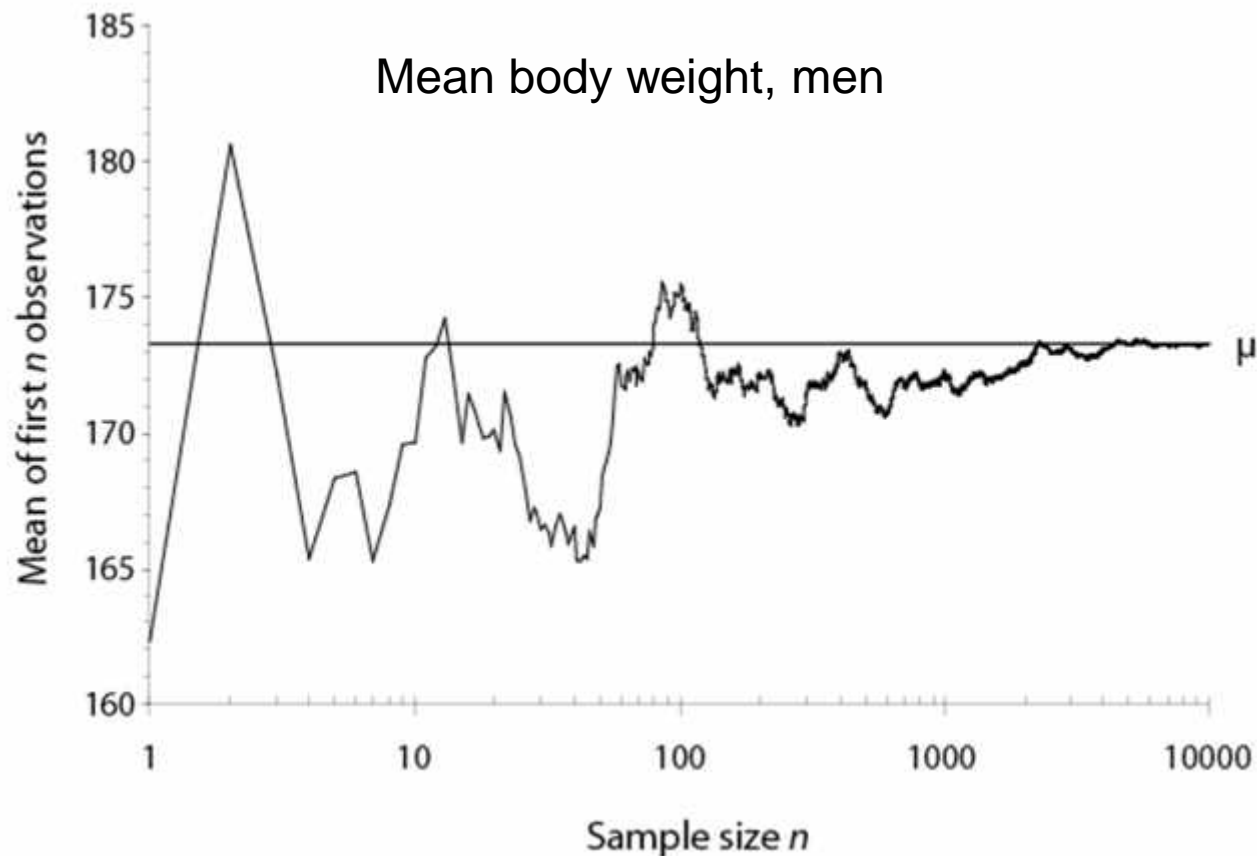
68-95-99.7 rule applied to the SDM

- We've established $\bar{x} \sim N(100, 4.7)$.
Therefore,
- 68% of x-bars within
 $\mu \pm \bar{x}$
 $= 100 \pm 4.7$
 $= 95.3 \text{ to } 104.7$
- 95% of x-bars within
 $\mu \pm 2 \cdot \bar{x}$
 $= 100 \pm (2 \cdot 4.7)$
 $= 90.6 \text{ to } 109.4$



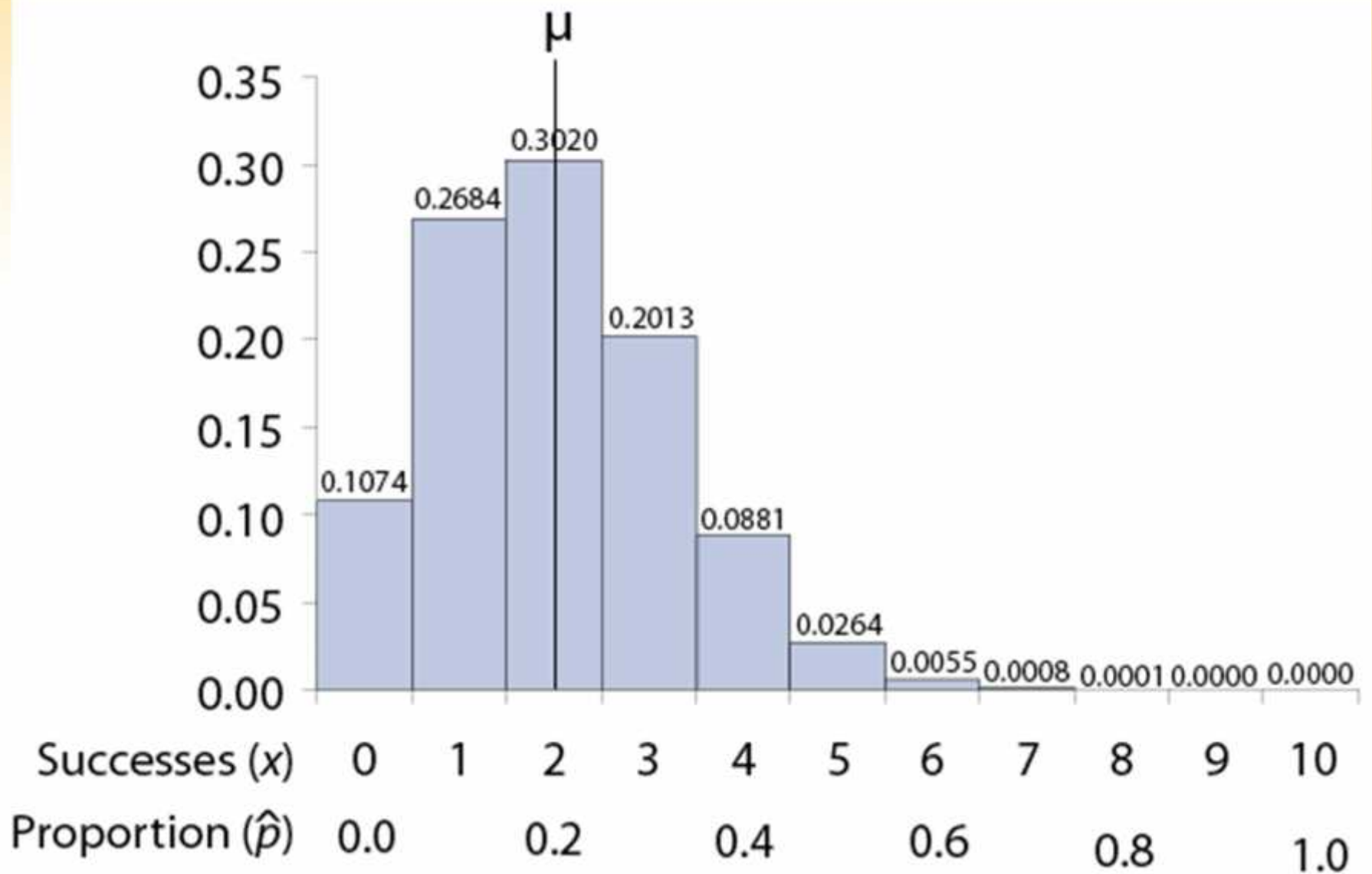
Law of Large Numbers

As a sample gets larger and larger, the \bar{x} approaches μ . Figure demonstrates results from an experiment done in a population with $\mu = 173.3$



8.3 Sampling Behavior of Counts and Proportions

- Recall Chapter: binomial random variable represents the random number of successes in n independent Bernoulli trials each with probability of success p ; notation $X \sim b(n, p)$
- $X \sim b(10, 0.2)$ is shown on the next slide. Note that $\mu = 2$
- Reexpress the counts of success as proportion $\hat{p} = x / n$. For this re-expression, $\mu = 0.2$



Normal Approximation to the Binomial (“ npq rule”)

- When n is large, the binomial distribution approximates a Normal distribution (“the Normal Approximation”)
- How large does the sample have to be to apply the Normal approximation? \Rightarrow One rule says that the Normal approximation applies when $npq \geq 5$

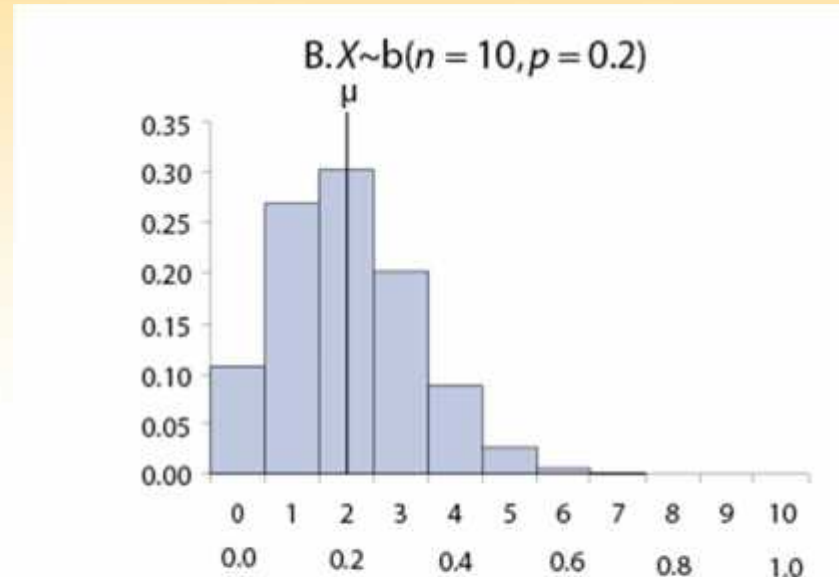
Top figure:

$$X \sim b(10, 0.2)$$

$$npq = 10 \cdot 0.2 \cdot (1 - 0.2)$$

$$= 1.6 \Rightarrow \text{Normal}$$

approximation does **not**
apply



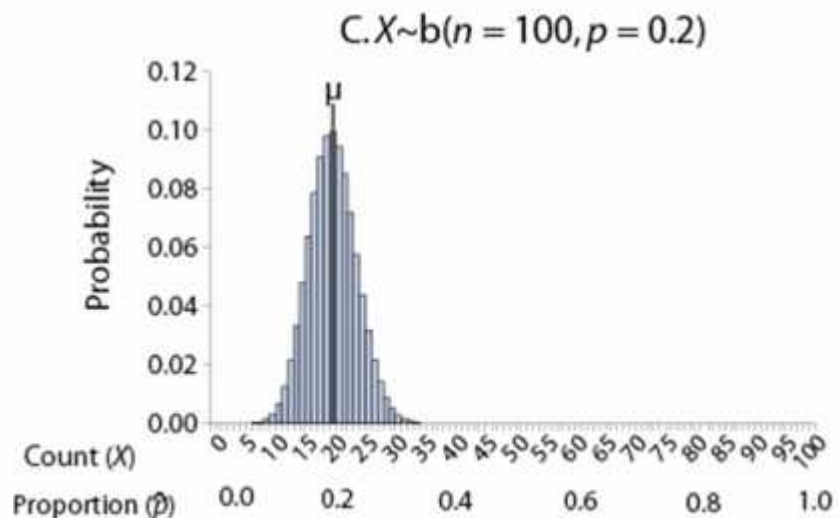
Bottom figure:

$$X \sim b(100, 0.2)$$

$$npq = 100 \cdot 0.2 \cdot (1 - 0.2)$$

$$= 16 \Rightarrow \text{Normal}$$

approximation applies



Normal Approximation for a Binomial Count

$$\sim = np \text{ and } \dagger = \sqrt{npq}$$

When Normal approximation applies:

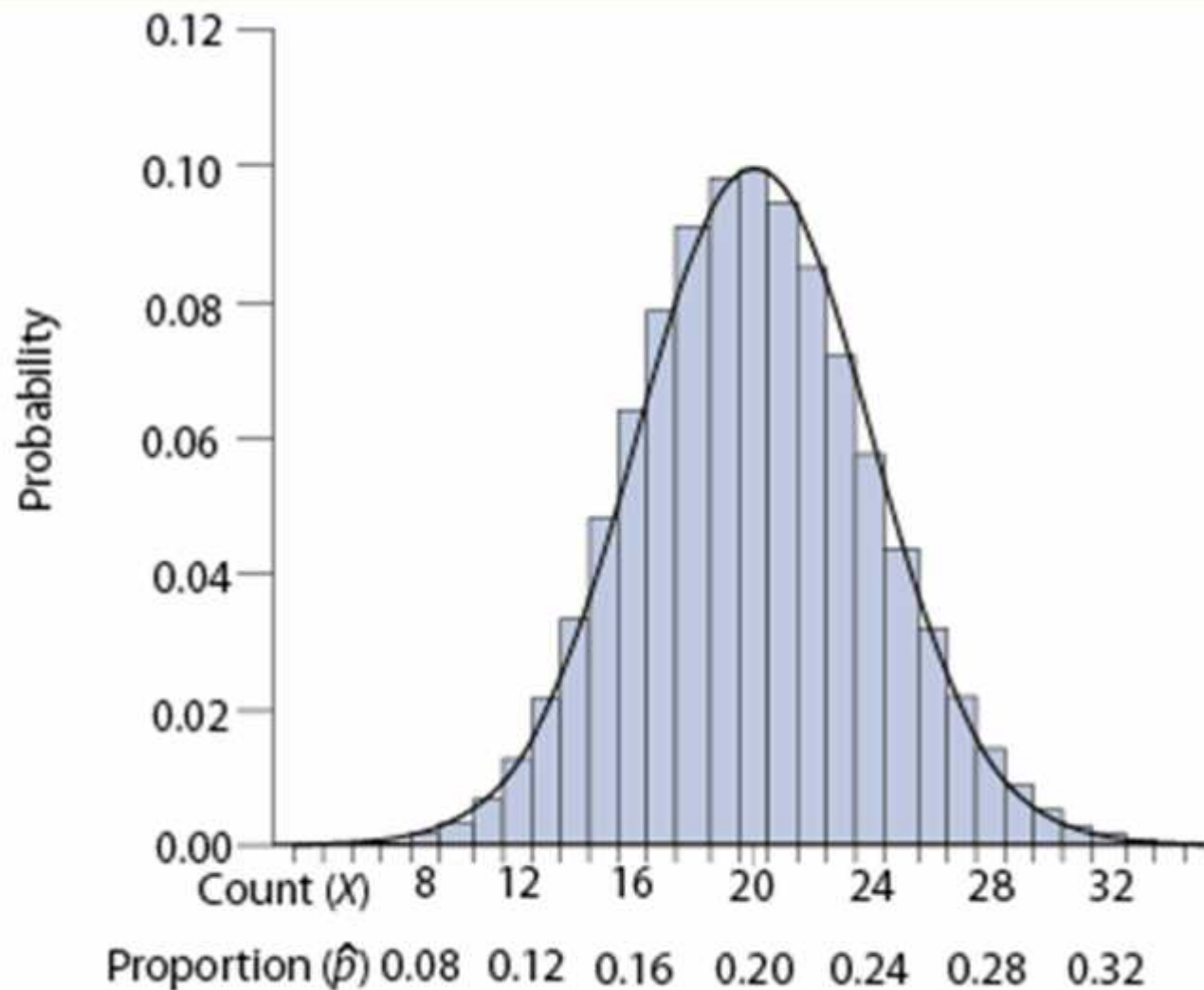
$$X \sim N\left(np, \sqrt{npq}\right)$$

Normal Approximation for a Binomial Proportion

$$\hat{p} \sim p \text{ and } \dagger = \sqrt{\frac{pq}{n}}$$

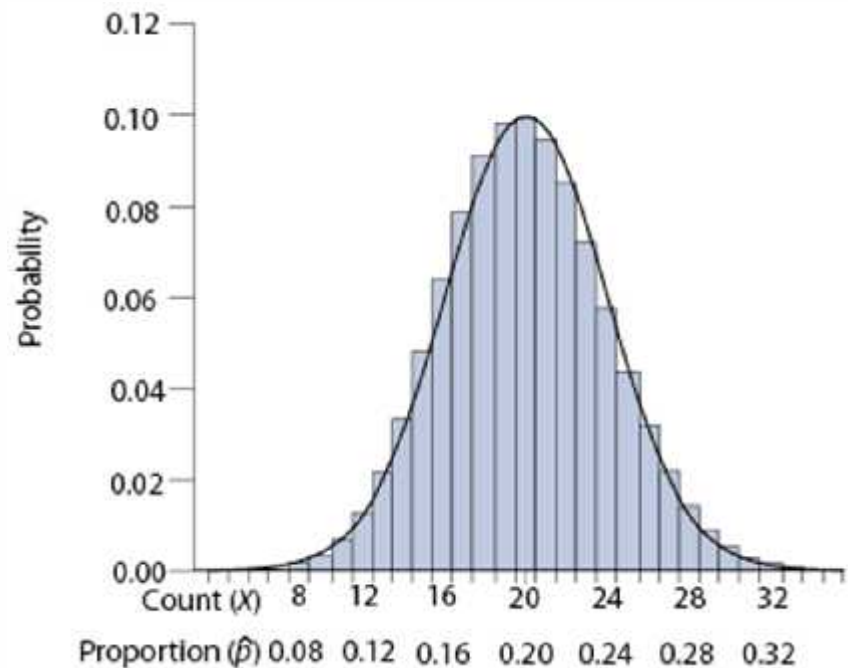
$$\hat{p} \sim N\left(p, \sqrt{\frac{pq}{n}}\right)$$

“p-hat” represents the sample proportion



Illustrative Example: Normal Approximation to the Binomial

- Suppose the prevalence of a risk factor in a population is 20%
- Take an SRS of $n = 100$ from population
- A variable number of cases in a sample will follow a binomial distribution with $n = 20$ and $p = .2$



Illustrative Example, cont.

The Normal approximation for the **count** is:

$$\sim = np = 100 \cdot .2 = 20$$

$$\text{and } \dagger = \sqrt{npq} = \sqrt{100 \cdot .2 \cdot .8} = 4$$

$$X \sim N(20, 4)$$

The Normal approximation for the **proportion** is:

$$\sim = p = .2$$

$$\dagger = \sqrt{\frac{pq}{n}} = \sqrt{\frac{.2 \cdot .8}{100}} = 0.04$$

$$\hat{p} \sim N(0.2, 0.04)$$

Illustrative Example, cont.

1. Statement of problem: Recall $X \sim N(20, 4)$
Suppose we observe 30 cases in a sample. What is the probability of observing at least 30 cases under these circumstance, i.e., $\Pr(X \geq 30) = ?$

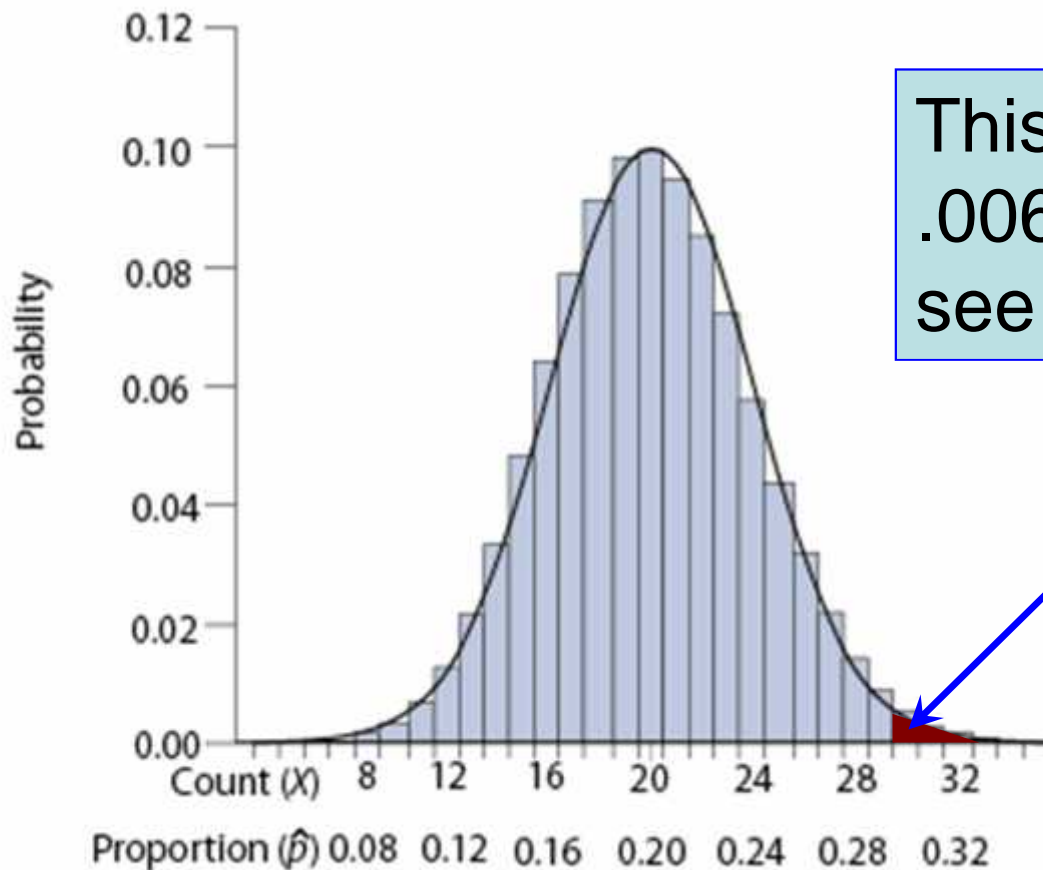
2. Standardize: $z = (30 - 20) / 4 = 2.5$

3. Sketch: next slide

4. Table B: $\Pr(Z \geq 2.5) = 0.0062$

Illustrative Example, cont.

Binomial and superimposed Normal distributions



This model suggests .0062 of samples will see 30 or more cases.