# CSCI 3022
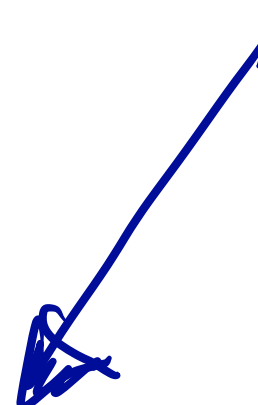
# intro to data science
# with probability & statistics

Lecture 13
February 28, 2018

1. The Central Limit Theorem

Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**

Dan Larremore

# Stuff & Things

- **Midterm** tonight.
  - 6:30 PM in HUMN 1B50.
    - Mix of multiple choice (13) and free-response (3)
    - You + one 8.5 x 11 inch sheet of handwritten notes; no magnif. glass.
    - You may bring a simple calculator that can't connect to the internet or store large data.
- **Office Hrs** as usual today: 11 to 1, Fleming 417.
- **HW3** due Friday - Answer only 4 out of the 5.

# Last time on CSCI 3022

mean     variance

- **Def**: A continuous random variable has a normal (or Gaussian) distribution with parameters $\mu$ and $\sigma^2$ if its probability density function is given by the following. We say $X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \longleftarrow \text{PDF}$$

But Muller

- **Proposition**: If X is a normally distributed random variable with mean $\mu$ and standard deviation $\sigma$, then Z is a standard normal distribution if

$$Z = \frac{X - \mu}{\sigma} \quad \text{or} \quad X = \sigma Z + \mu$$

- **Fact**: If Z is a standard normal random variable, then we can compute probabilities using the standard normal CDF

$$P(Z \leq z) = \int_{-\infty}^{z} f(x)dx = \Phi(z)$$

# Motivating example

- Soon, we'll be talking about *statistical inference* where we'll try to infer (learn) things about the true mean of a population using sample datasets

- **Examples**:
  - CU AERO student mean GPA // sample 30 students
  - Do all zebras have stripes? // sample 20 zebras in field
  - ...

# Random samples

- The random variables $X_1$, $X_2$, …, $X_n$ are said to form a (sample) random sample of size *n* if:

  - all $X_k$ $k = 1, 2, … n$ are independent
  - all $X_k$ $k = 1, 2, … n$ have the same distribution

- We say that these $X_k$'s are $i.i.d$

  stands for independent and identically distributed.

# Estimators and their distributions

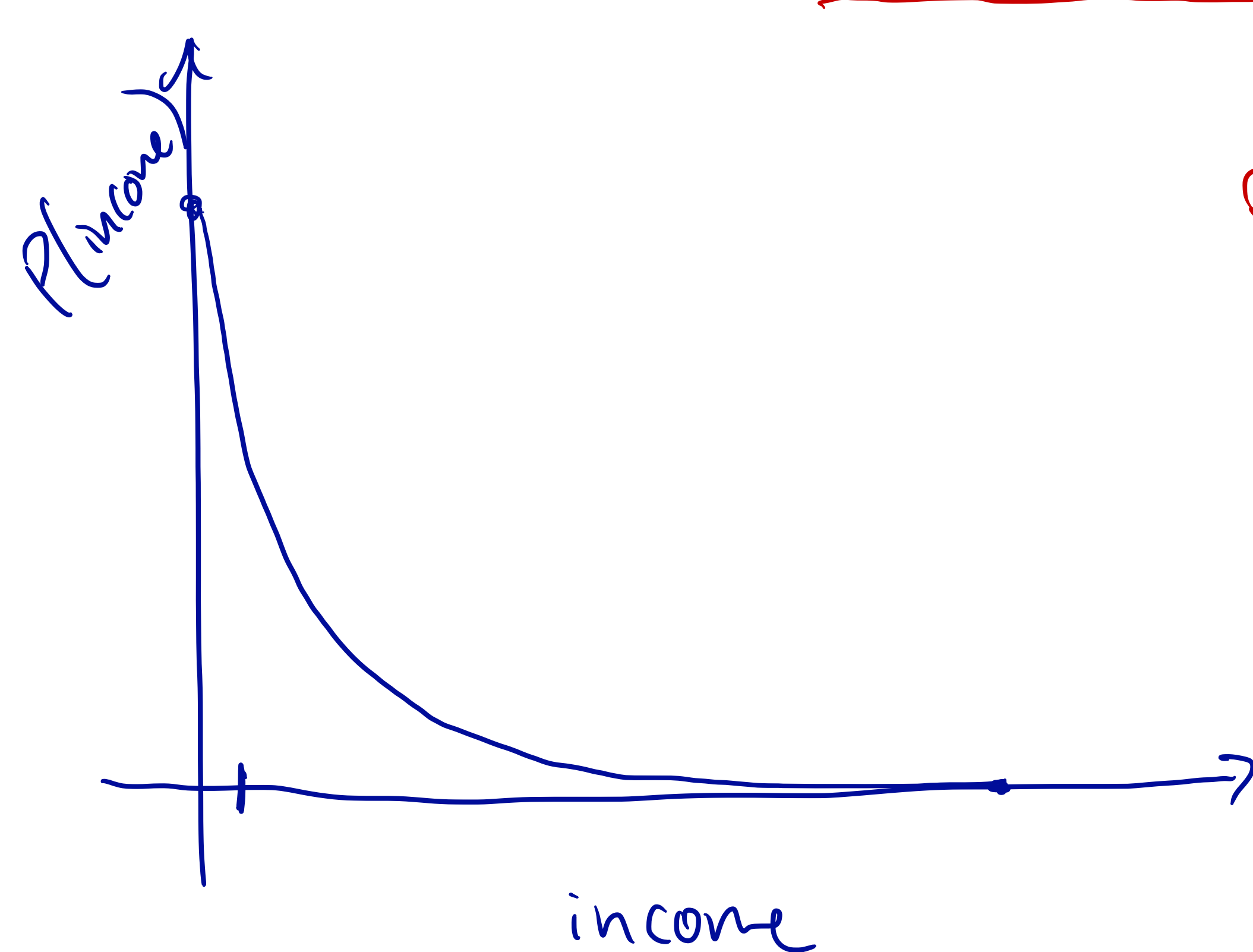- We use **estimators** to summarize our i.i.d. sample

- **Examples**:

  ① $\bar{X}$  Sample mean $\longrightarrow$ use to estimate true mean $\mu$

  ② $\hat{p}$  sample proportion $\longrightarrow$ use to estimate true proportion $P$

  ③ $s^2$  sample variance $\longrightarrow$ use to estimate true variance $\sigma^2$

# Estimators and their distributions

- We use **estimators** to summarize our i.i.d. sample *R.V.*

- Any estimator, including the **sample mean**, $\bar{X}$, is a random variable. Why? Because it's based on a random sample.

- This means that $\bar{X}$ has a distribution of its own, which is referred to as the **sampling distribution of the sample mean**.

- The sampling distribution depends on:

  - What is the underlying (true) distribution
  - number of samples, $n$
  - method of sampling.

# Distribution of the Sample Mean

- What does the distribution of the sample mean actually look like?

- For example, does it look like the distribution that it's sampling? *Not really...*

P(income)

income

P(mean income)

mean income

# Distribution of the Sample Mean

- What does the distribution of the sample mean actually look like?

- **Proposition**: Let $X_1, X_2, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$ . Then for any $n$,

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

*not just $\sigma^2$ variance, it's also divided by $n$, the # of samples*

*estimates of mean from samples*

*same mean as the iid variables*

- We know everything there is to know about the distribution of the sample mean when the population distribution is normal!

# Distribution of the Sample Mean

- If the population is normally distributed, then:

$\bar{X}$ distribution when $n = 10$

var $\dfrac{\sigma^2}{10}$

Coincidence

$\bar{X}$ distribution when $n = 4$

var $\sigma^2/4$

Population distribution
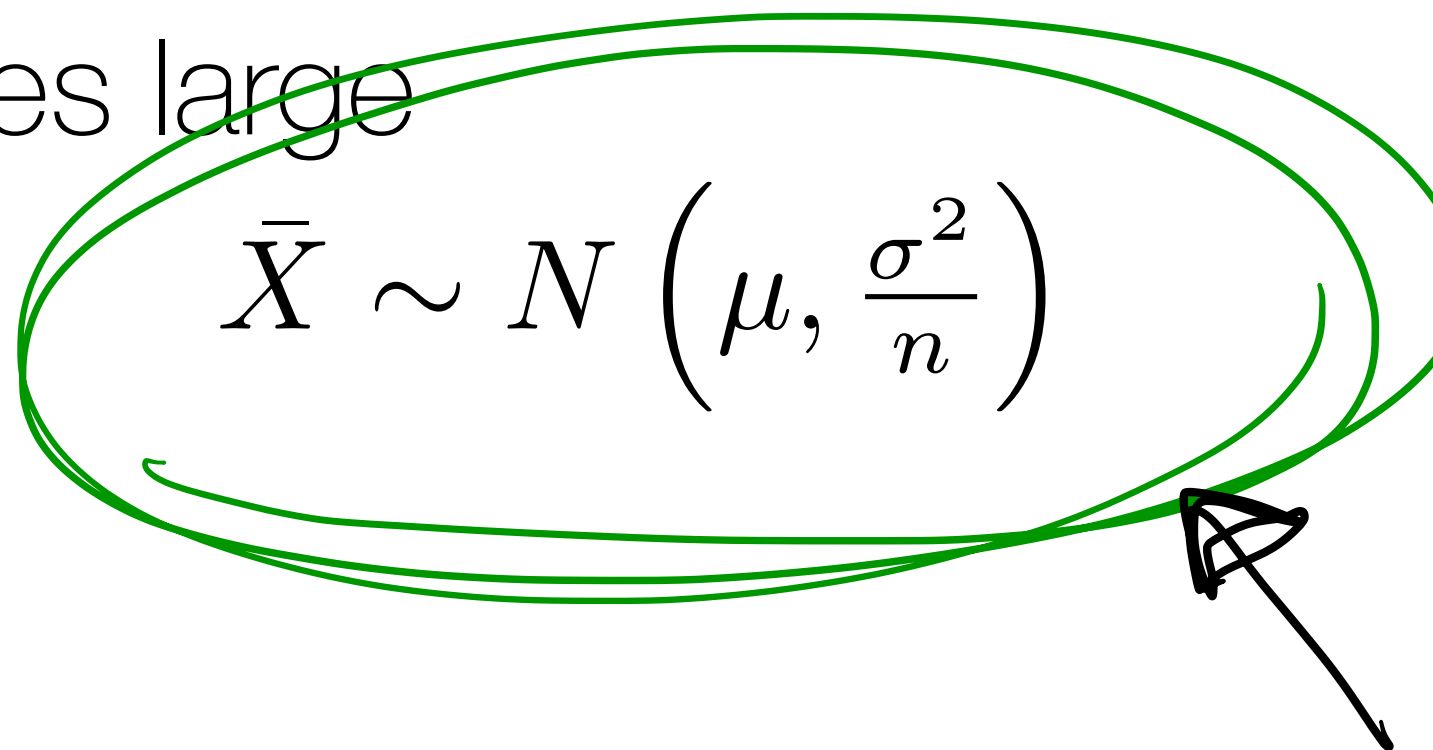
var $\sigma^2$

$\bar{X}$

# Distribution of the Sample Mean

- What if the population is *not* normally distributed?
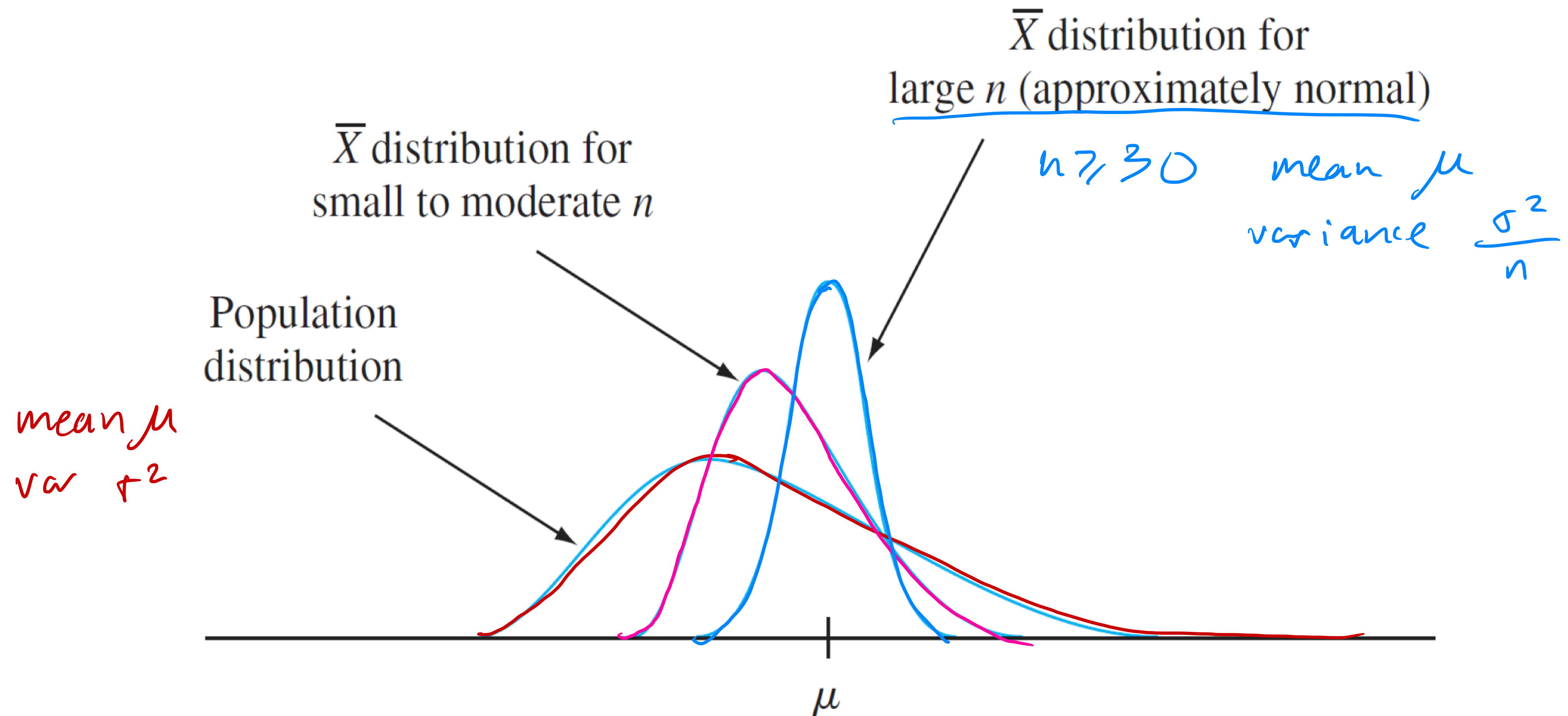
# The Central Limit Theorem

- What if the population is *not* normally distributed?

- **Important**: When the population distribution is non-normal, averaging produces a distribution more bell-shaped than the one being sampled.

- A reasonable assumption is that *if n is large*, a suitable normal curve will well-approximate the actual distribution of the sample mean.

- **The Central Limit Theorem**: Let $X_1, X_2, \ldots, X_n$ be i.i.d. draws from some distribution. Then as n becomes large

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- Rule of Thumb: $n \geq 30$

# Distribution of the sample mean:

- If the population is *not* normally distributed



$\overline{X}$ distribution for
large $n$ (approximately normal)

$n \geqslant 30$    mean $\mu$

variance $\dfrac{\sigma^2}{n}$

$\overline{X}$ distribution for
small to moderate $n$

Population
distribution

mean $\mu$
var $\tau^2$

$\mu$

# Examples:

- **Example 1**: A hardware store receives a shipment of bolts that are supposed to be 12cm long. The mean is indeed 12cm, and the standard deviation is 0.2cm. For quality control, the hardware store chooses 100 bolts at random to measure. They will call the shipment defective and return it to the manufacturer if the average length of the 100 bolts is less than 11.97cm or greater than 12.04cm. Find the probability that the shipment is found satisfactory.

① What's info about pop: $\mu = 12$, $\sigma = 0.2$

What's info about sample: $n = 100$    $P(11.97 \leq \bar{X} \leq 12.04)$ ?

② How is $\bar{X}$ distributed?

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

plug in!    $\bar{X} \sim N\left(12, \frac{0.2^2}{100}\right)$

# Examples:

"side" calc: if $\sigma^2 = \dfrac{0.2^2}{100}$ $\sigma = \dfrac{\sqrt{0.2^2}}{\sqrt{100}}$ $\sigma = \dfrac{0.2}{10}$ $\sigma = 0.02$

$\displaystyle\int_{-1.5}^{2} f(z)\,dz = F(2) - F(-1.5)$

$\Phi(2) - \Phi(-1.5)$

- **Example 1**: A hardware store receives a shipment of bolts that are supposed to be 12cm long.  The mean is indeed 12cm, and the standard deviation is 0.2cm.  For quality control, the hardware store chooses 100 bolts at random to measure.  They will call the shipment defective and return it to the manufacturer if the average length of the 100 bolts is less than 11.97cm or greater than 12.04cm.  Find the probability that the shipment is found satisfactory.

Goal: $P\left(11.97 \le \bar{X} \le 12.04\right)$

Know, by CLT: $\bar{X} \sim N\left(12, \dfrac{0.2^2}{100}\right)$

Box Muller $\quad Z = \dfrac{\bar{X} - 12}{0.02}$

New Goal: $P\left(\dfrac{11.97 - 12}{0.02} \le Z \le \dfrac{12.04 - 12}{0.02}\right)$

$= P\left(-1.5 \le Z \le 2\right)$

$= \Phi(2) - \Phi(-1.5)$

$\boxed{= 0.91}$

# Examples:

- **Example 2**: Suppose you have a jar of lemon and banana jelly beans where it is known that the <u>true proportion of lemon jelly beans is 0.5</u>. You try to estimate the proportion of lemon beans by reaching in and <u>drawing 50 jelly beans</u> and testing them (by eating them). What is the <u>probability that your sample is 75% or more lemon jelly beans?</u>

- Note: this is a little different because we're estimating a *proportion*. What changes?

Population: $p = 0.5$

Sample: $n = 50$    $P(X \geq 0.75) = 1 - P(X \leq 0.75)$

Proportions are different. $\bar{X} = \hat{p} = \dfrac{Bin(n,p)}{n}$

What is variance of $\hat{p}$? $Var(\hat{p}) = Var\left(\dfrac{Bin(n,p)}{n}\right) = \dfrac{1}{n^2} Var(Bin(n,p))$

$$= \dfrac{p(1-p)}{n}$$

$$= \dfrac{1}{n^2} n p(1-p)$$

# Examples:

- **Example 2**: Suppose you have a jar of lemon and banana jelly beans where it is known that the true proportion of lemon jelly beans is 0.5. You try to estimate the proportion of lemon beans by reaching in and drawing 50 jelly beans and testing them (by eating them). What is the probability that your sample is 75% or more lemon jelly beans?

# Problem-solving hints:

- **First**, identify the *population* and identify the *sample*.

- **Second**, is the problem about *means* or *proportions*?

- **Then**, we're off to the races using the CLT, the Box-Muller transform, and our standard normal distribution!