

# Simulation wrap-up

---

# Last time

---

- Time-driven, event-driven
- “Simulation” from differential equations
- Cellular automata, microsimulation, agent-based simulation
  - see e.g.  
<http://www.microsimulation.org/IMA/What%20is%20microsimulation.htm>
- Example applications: SIR disease model, population genetics

# Simulation: Pros and Cons

---

- Pros:
  - Building model can be easy (easier) than other approaches
  - Outcomes can be easy to understand
  - Cheap, safe
  - Good for comparisons
- Cons:
  - Hard to debug
  - No guarantee of optimality
  - Hard to establish validity
  - Can't produce absolute numbers

# Simulation: Important Considerations

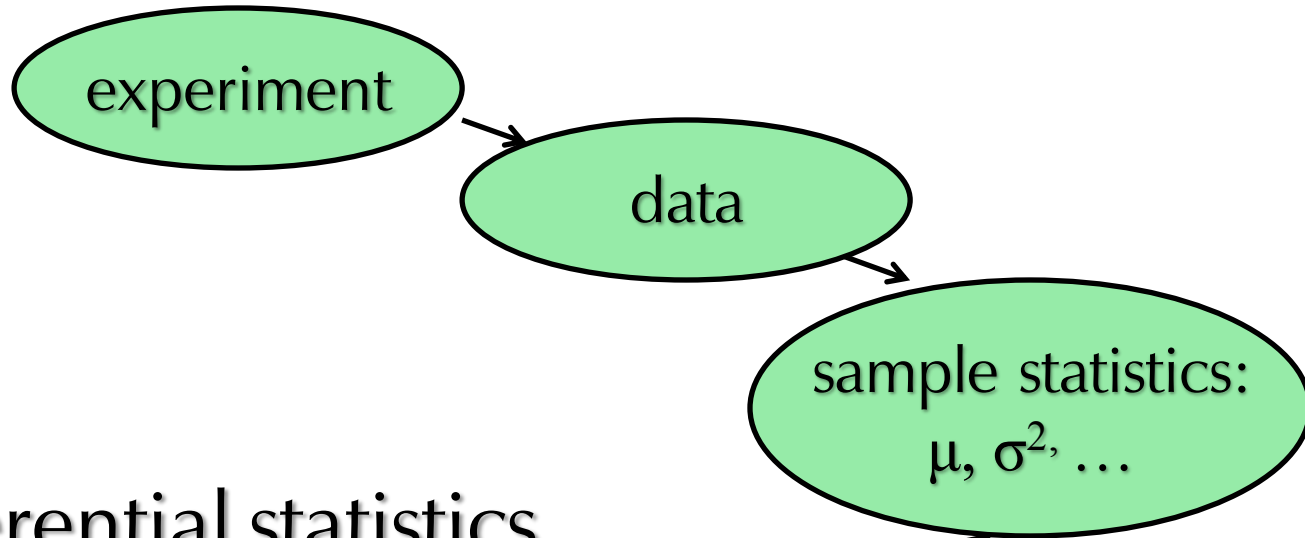
---

- Are outcomes statistically significant? (Need many simulation runs to assess this)
- What should initial state be?
- How long should the simulation run?
- Is the model realistic?
- How sensitive is the model to parameters, initial conditions?

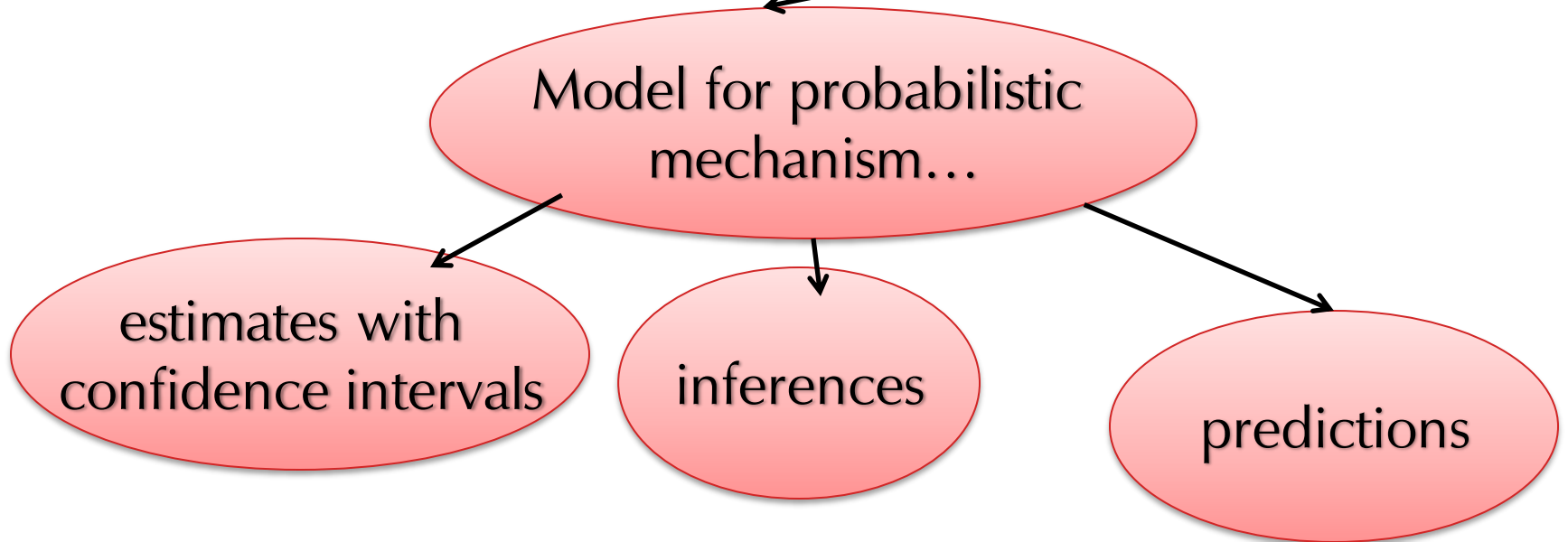
# Statistics Overview

---

# Descriptive statistics



# Inferential statistics



# Random Variables

---

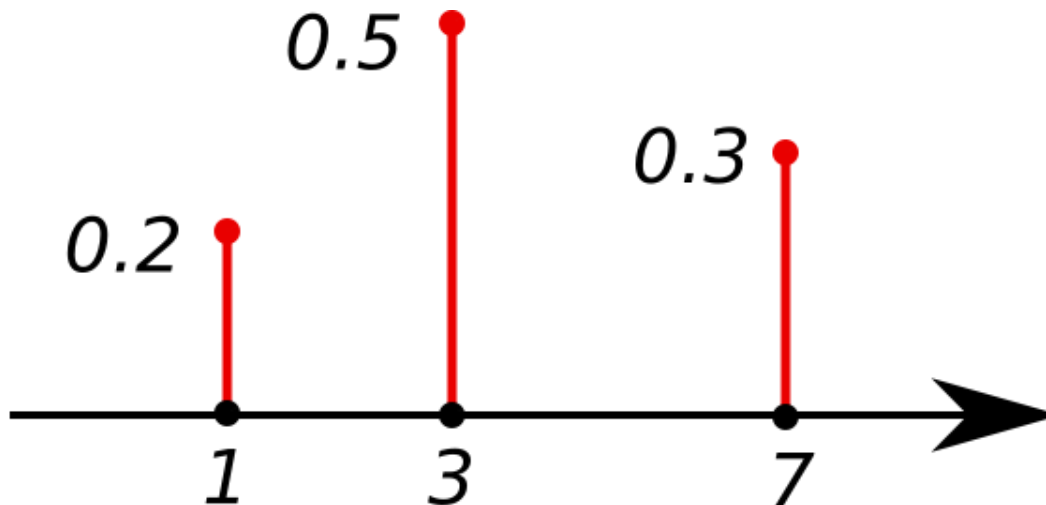
- A random variable is any “probabilistic outcome”
  - e.g., a coin flip, height of someone randomly chosen from a population
- A R.V. takes on a value in a *sample space*
  - space can be discrete, e.g., {H, T}
  - or continuous, e.g. height in (0, infinity)
- R.V. denoted with capital letter (X), a realization with lowercase letter (x)
  - e.g., X is a coin flip, x is the value (H or T) of that coin flip

# Probability Mass Function

---

- Describes probability for a discrete R.V.
- e.g.,

$$f_X(x) = \begin{cases} \frac{1}{2}, & x \in \{0, 1\}, \\ 0, & x \notin \{0, 1\}. \end{cases}$$

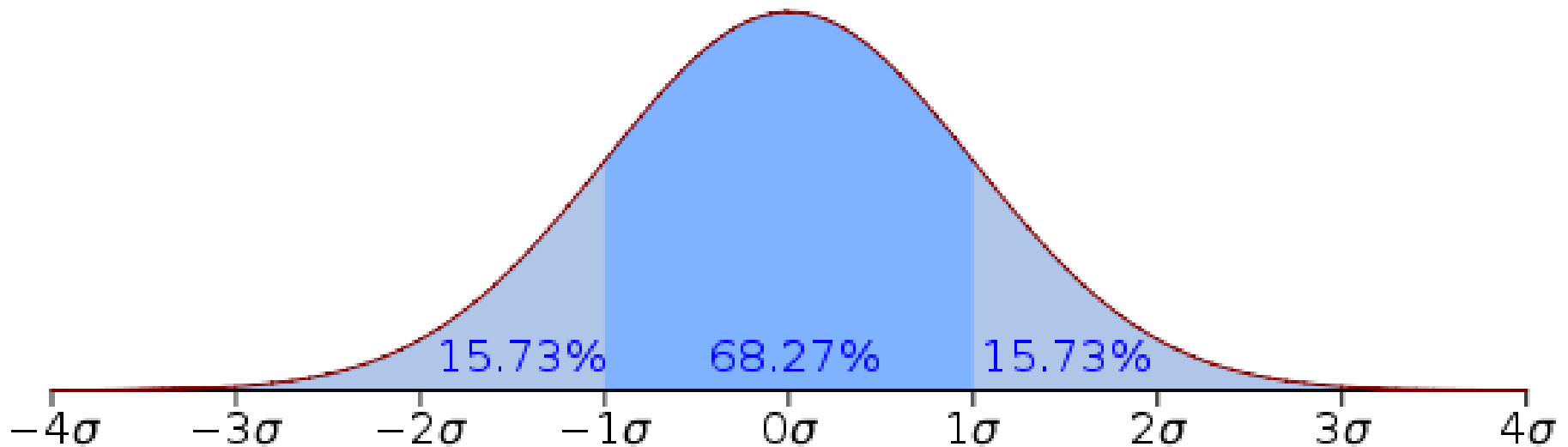




# Probability Density Function

---

- Describes probability for a continuous R.V.
- e.g.,



# [Population] Mean of a Random Variable

---

- aka expected value, first moment

- for discrete RV:  $E[X] = \mu = \sum_i x_i p_i$

(requires that  $\sum_i p_i = 1$ )

- for continuous RV:  $E[X] = \mu = \int_{-\infty}^{\infty} x p(x) dx$

(requires that  $\int_{-\infty}^{\infty} p(x) dx = 1$ )

## [Population] Variance

---

$$\begin{aligned}\sigma^2 &= \mathbf{E}[(X - \mu)^2] \\ &= \mathbf{E}[X^2 - 2X\mu + \mu^2] \\ &= \mathbf{E}[X^2] - \mu^2 \\ &= \mathbf{E}[X^2] - (\mathbf{E}[X])^2\end{aligned}$$

- for discrete RV: 
$$\sigma^2 = \sum_i p_i (x_i - \mu)^2$$

- for continuous RV: 
$$\sigma^2 = \int (x - \mu)^2 p(x) dx$$

# Sample mean and sample variance

---

- Suppose we have  $N$  independent observations of  $X$ :  $x_1, x_2, \dots, x_N$

- Sample mean:

$$\frac{1}{N} \sum_{i=1}^N x_i = \bar{x}$$

Unbiased:

$$E[\bar{x}] = \mu$$

- Sample variance:

$$\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 = s^2$$

$$E[s^2] = \sigma^2$$

# $1/(N-1)$ and the sample variance

---

- The  $N$  differences  $x_i - \bar{x}$  are not independent:

$$\sum (x_i - \bar{x}) = 0$$

- If you know  $N-1$  of these values, you can deduce the last one
  - i.e., only  $N-1$  degrees of freedom

- Could treat sample as population and compute population variance:

$$\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

- BUT this underestimates true population variance (especially bad if sample is small)

# Sample variance using $1/(N-1)$ is unbiased

---

$$\begin{aligned} E[s^2] &= E\left[\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2\right] \\ &= \frac{1}{N-1} E\left[\sum_{i=1}^N x_i^2 - N\bar{x}^2\right] \\ &= \frac{1}{N-1} \left[ N(\sigma^2 + \mu^2) - N\left(\frac{\sigma^2}{N} + \mu^2\right) \right] \\ &= \sigma^2 \end{aligned}$$

# Computing Sample Variance

---

- Can compute as

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

- Prefer:

$$s^2 = \frac{\left( \sum_{i=1}^N x_i^2 \right) - N(\bar{x})^2}{N-1}$$

(one pass, fewer operations, more accurate)

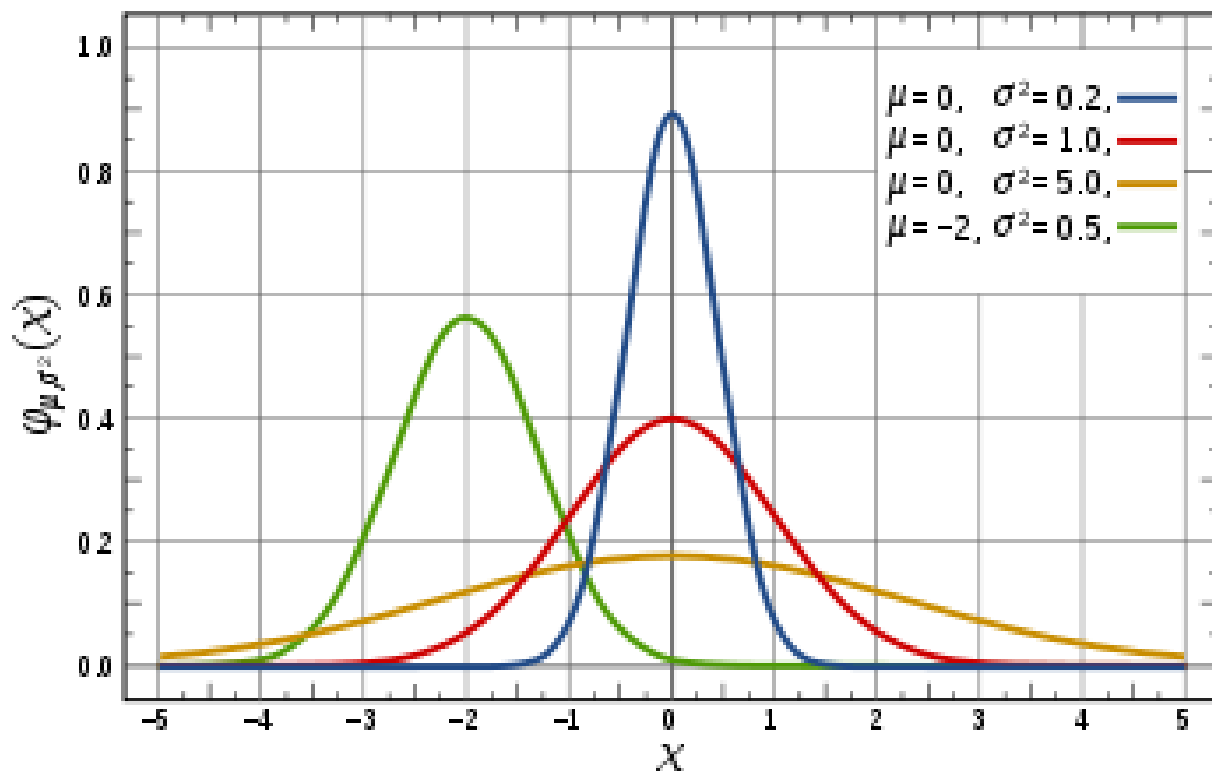
# The Gaussian Distribution

---

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$E[X] = \mu$$

$$\text{Var}[X] = \sigma^2$$





# Why so important?

---

- **sum** of independent observations of random variables converges to Gaussian \*(with some assumptions)
- **in nature**, events having variations resulting from many small, independent effects tend to have Gaussian distributions
  - demo: <http://www.mongrav.org/math/falling-balls-probability.htm>
  - e.g., measurement error
  - if effects are multiplicative, **logarithm** is often **normally** distributed

# Central Limit Theorem

---

- Suppose we sample  $x_1, x_2, \dots, x_N$  from a distribution with mean  $\mu$  and variance  $\sigma^2$

- Let 
$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

holds for \*(almost) any parent distribution!

- then

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{N}} \rightarrow N(0,1)$$

- i.e.,  $\bar{x}$  distributed normally with mean  $\mu$ , variance  $\sigma^2/N$

# Important Properties of Normal Distribution

---

1. Family of normal distributions closed under linear transformations:

if  $X \sim N(\mu, \sigma^2)$  then

$$(aX + b) \sim N(a\mu + b, a^2\sigma^2)$$

2. Linear combination of normals is also normal:

if  $X_1 \sim N(\mu_1, \sigma_1^2)$  and  $X_2 \sim N(\mu_2, \sigma_2^2)$  then

$$aX_1 + bX_2 \sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$$

# Important Properties of Normal Distribution

---

3. Of all distributions with mean and variance, normal has **maximum entropy**

Information theory: Entropy like “uninformativeness”

Principle of maximum entropy: choose to represent the world with as uninformative a distribution as possible, subject to “testable information”

**If** we know  $x$  is in  $[a, b]$ , then uniform distribution on  $[a, b]$  has least entropy

If we know distribution has mean  $\mu$ , variance  $\sigma^2$ , normal distribution  $N(\mu, \sigma^2)$  has least entropy

# Important Properties of Normal Distribution

---

4. If errors are normally distributed, a least-squares fit yields the maximum likelihood estimator

Finding least-squares  $x$  st  $Ax \approx b$  finds the value of  $x$  that maximizes the likelihood of data  $A$  under some model

$$P(\text{data}|\text{model}) \propto \prod_{i=1}^n \exp \left[ -\frac{1}{2} \left( \frac{y_i - y(x_i)}{\sigma_i} \right)^2 \right] \Delta y$$

$$P(\text{model}|\text{data}) \propto P(\text{data}|\text{model})P(\text{model})$$

# Important Properties of Normal Distribution

---

5. Many derived random variables have analytically-known densities  
e.g., sample mean, sample variance
6. Sample mean and variance of  $n$  identical independent **samples** are independent; sample mean is a normally-distributed random variable

$$\bar{X}_n \sim N(\mu, \sigma^2 / n)$$

# Distribution of Sample Variance

## (For Gaussian R.V. X)

---

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$\text{define } U = \frac{(n-1)s^2}{\sigma^2}$$

then  $U$  has a  $\chi^2$  distribution with  $(n-1)$  d.o.f.

$$p(x) = \left[ 2^{n/2} \Gamma\left(\frac{n}{2}\right) \right]^{-1} (x)^{\frac{n}{2}-1} e^{-x/2}, \quad x \geq 0$$

$$E[U] = n - 1, \quad \text{Var}[U] = 2(n - 1)$$

# The Gamma Function

---

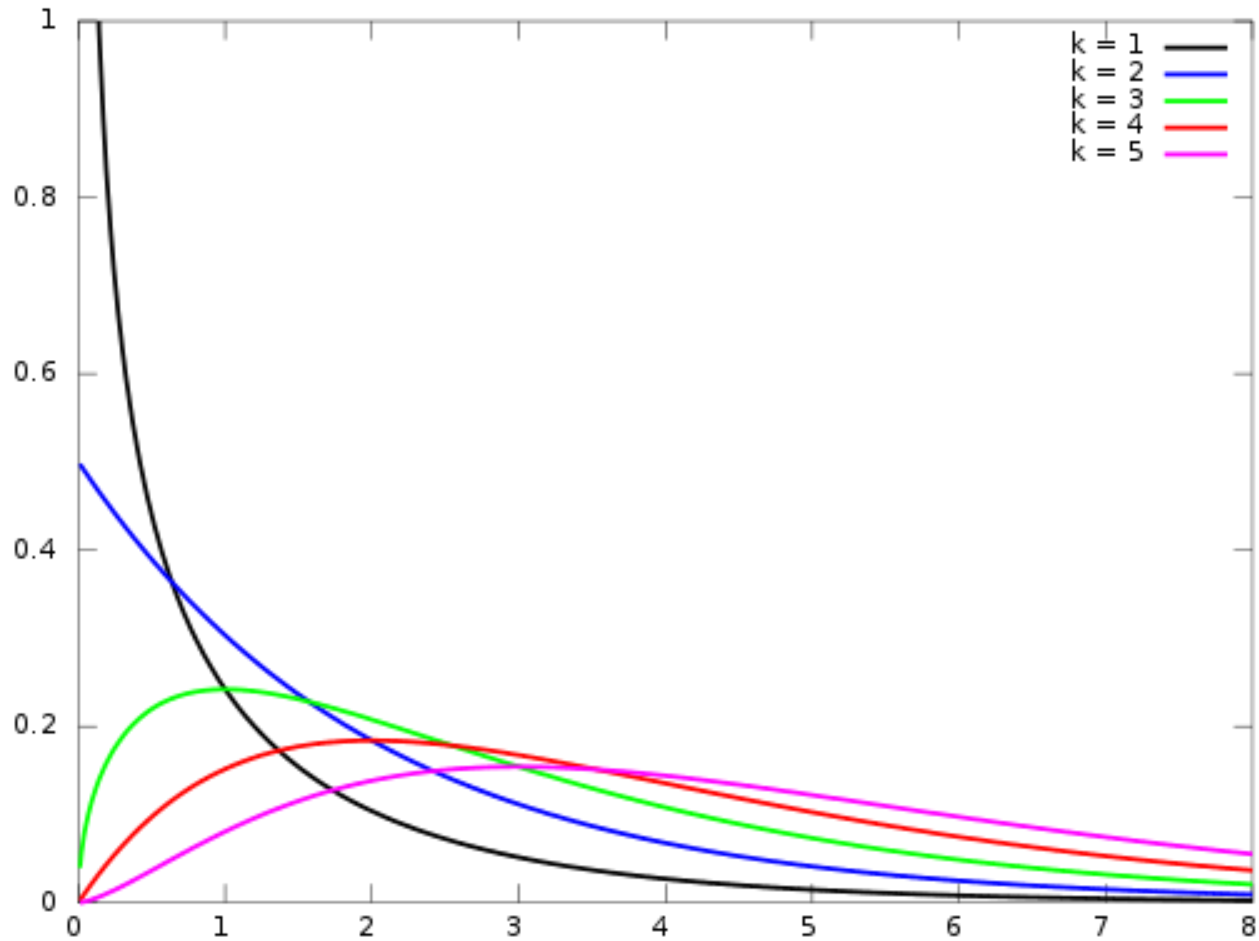
$$\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt.$$

$$\Gamma(n) = (n - 1)!$$



# The Chi-Squared Distribution

---



# What if we don't know true variance?

---

- Sample mean is normally distributed R.V.

$$\bar{X}_n \sim N(\mu, \sigma^2 / n)$$

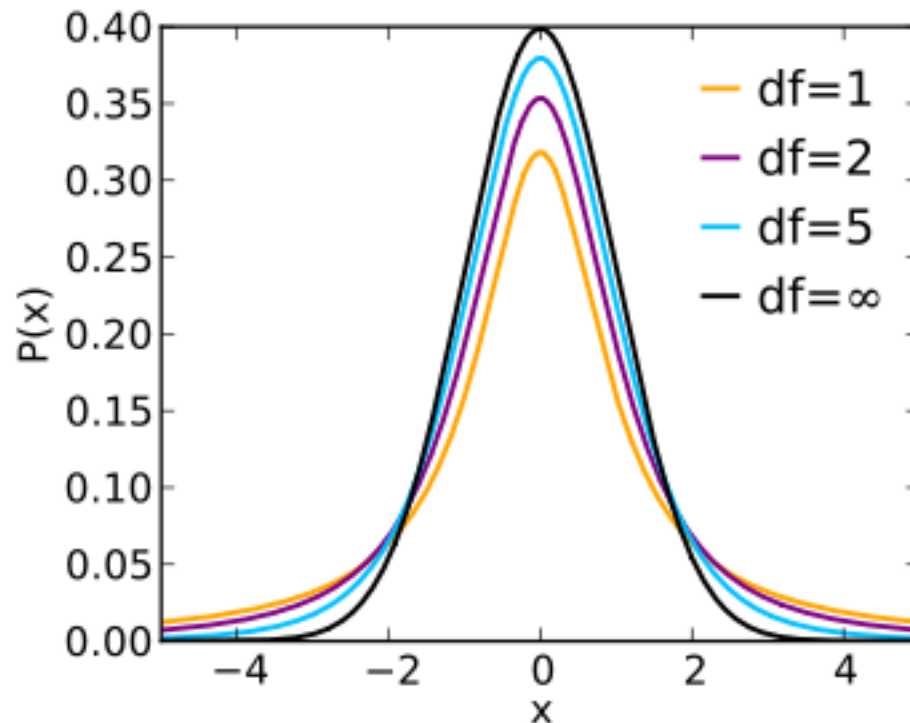
- Taking advantage of this presumes we know  $\sigma^2$

- $\frac{\bar{x} - \mu}{s_n / \sqrt{n}}$  has a  $t$  distribution with  $(n-1)$  d.o.f.

# [Student's] t-distribution

---

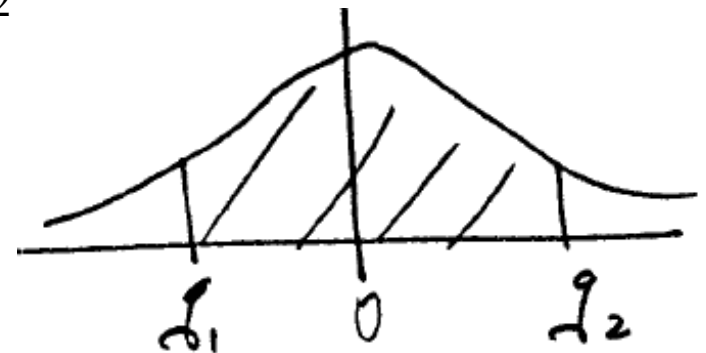
$$f(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi} \Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}},$$



# Forming a confidence interval

---

- e.g., given that I observed a sample mean of \_\_\_\_\_, I'm 99% confident that the true mean lies between \_\_\_\_\_ and \_\_\_\_\_.
- Know that  $\frac{\bar{x} - \mu}{s_n / \sqrt{n}}$  has  $t$  distribution
- Choose  $q_1, q_2$  such that student  $t$  with  $(n-1)$  dof has 99% probability of lying between  $q_1, q_2$



# Confidence interval for the mean

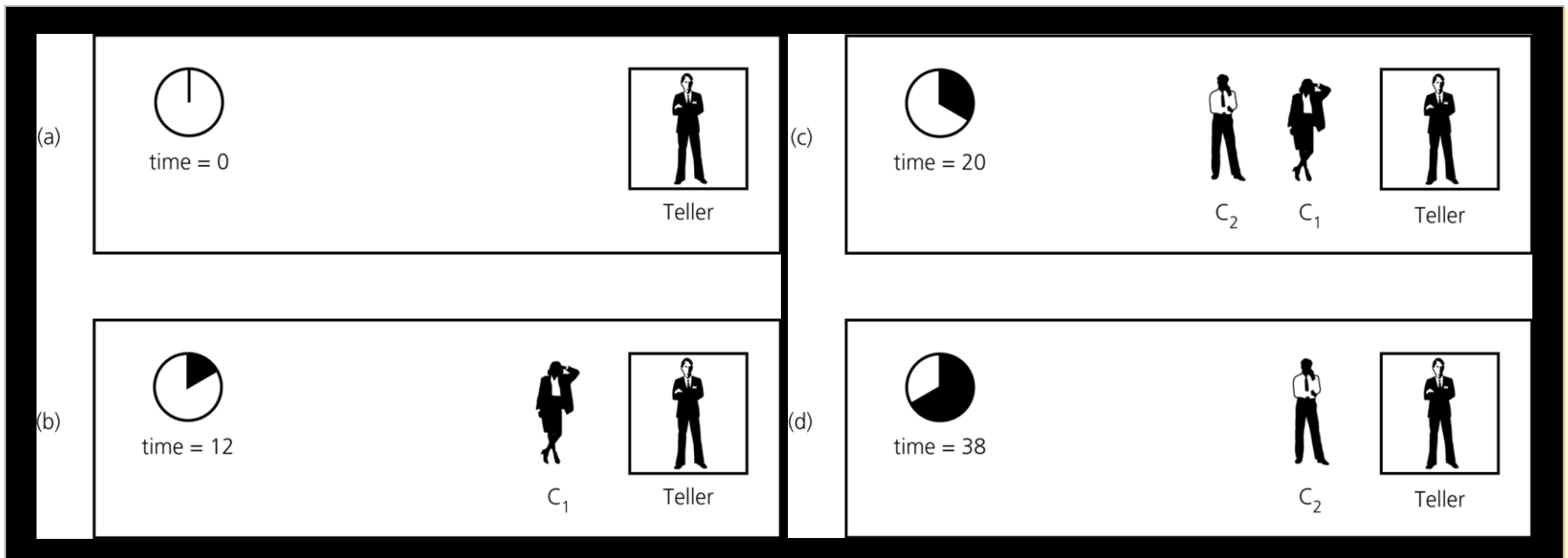
---

$$\text{if } P \left\{ q_1 < \frac{\bar{x}_n - \mu}{s_n / \sqrt{n}} < q_2 \right\} = 0.99$$

$$\text{then } P \left\{ \bar{x}_n - q_2 \frac{s_n}{\sqrt{n}} < \mu < \bar{x}_n - q_1 \frac{s_n}{\sqrt{n}} \right\} = 0.99$$

# Interpreting Simulation Outcomes

- How long will customers have to wait, on average?
  - e.g., for given # tellers, arrival rate, service time distribution, etc.



# Interpreting Simulation Outcomes

---

- Simulate bank for  $N$  customers
- Let  $x_i$  be the wait time of customer  $i$
- Is  $\text{mean}(x)$  a good estimate for  $\mu$ ?
- How to compute a 95% confidence interval for  $\mu$ ?
  - Problem:  $x_i$  are not independent!

# Replications

---

- Run simulation to get  $M$  observations
- Repeat simulation  $N$  times (different random numbers each time)
- Treat the sample mean of different runs as approximately uncorrelated

$$\bar{X}_i = \frac{1}{m} \sum_{j=1}^n x_{ij} \qquad \bar{\bar{X}} = \frac{1}{n} \sum_{i=1}^n \bar{X}_i$$

$$s^2 = \frac{1}{n-1} \sum_i (\bar{X}_i - \bar{\bar{X}})^2$$