

# **Statistical Methods for High Dimensional Biology**

## **STAT/BIOF/GSAT 540**

Lecture 3 – Review of probability and  
statistical inference

Sara Mostafavi

January 11 2016

**\*\*Lectures prepared by Dr. Jenny Bryan\*\***

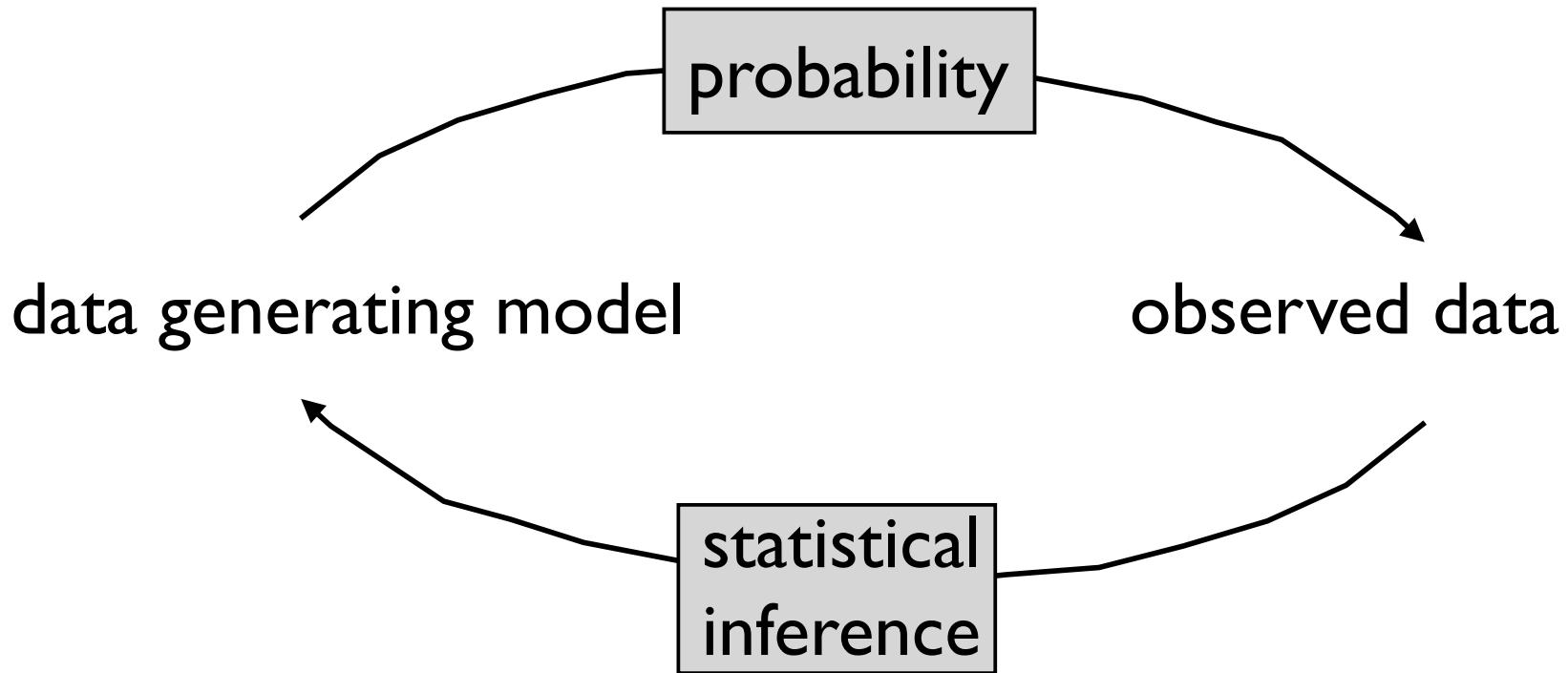
# So far we have reviewed:

- estimate/use data **generating model** to understand/describe an observed *sample*
- rv's and their distributions
- Importance of variance in hypothesis testing
- Intro hypothesis testing

# Today

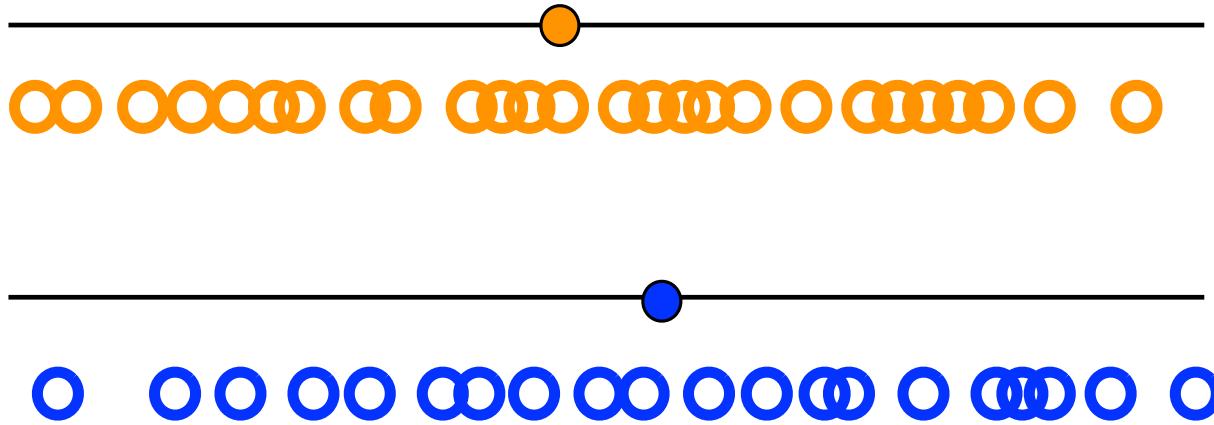
- IID
- Parameters of a distribution
- Law of large numbers
- Central Limit Theorem
- Hypothesis testing and parameter estimation
  - Method of maximum likelihood
- Types of errors in hypothesis testing

Going from data to model (vs model to data) requires lots of assumptions and simplifications.



Adapted from Figure 1 of "All of Statistics".

## Why we care about IID observations?



Regard the data as iid observations of random variables that have certain (unknown) distributions.

$$Y_1, \dots, Y_i, \dots, Y_{n_y} \sim \text{iid } F$$

$$Z_1, \dots, Z_i, \dots, Z_{n_z} \sim \text{iid } G$$

What do we mean by iid?

**iid**

independent  
identically  
distributed

$$Y_1, \dots, Y_i, \dots, Y_{n_y} \sim \text{iid } F$$

$$Z_1, \dots, Z_i, \dots, Z_{n_z} \sim \text{iid } G$$

But let's cut to the chase: independence of events or rvs makes it much easier to write down the probability of joint events or the joint distribution. It allows you to write these as a *simple product*.

Toss a fair coin 10 times.

A = at least one head

$T_j$  = toss  $j$  yields tails,  $j \in \{1, 2, \dots, 10\}$

What's the probability of A if you toss a fair coin 10 times?

Toss a fair coin 10 times.

A = at least one head

$T_j$  = toss  $j$  yields tails,  $j \in \{1, 2, \dots, 10\}$

$$\begin{aligned} P(A) &= 1 - P(\text{not } A) \\ &= 1 - P(\text{all tosses yield tails}) \\ &= 1 - P(T_1 T_2 \dots T_{10}) \\ &= 1 - P(T_1) P(T_2) \dots P(T_{10}) \quad * \\ &= 1 - 0.5^{10} \approx 0.999 \end{aligned}$$

\*Independence of the events  $T_j$  is critical to making this such a simple calculation!

**iid**

independent  
identically  
distributed

$$Y_1, \dots Y_i, \dots, Y_{n_y} \sim \text{iid } F$$

$$Z_1, \dots Z_i, \dots, Z_{n_z} \sim \text{iid } G$$

Independence of events or rvs makes it much easier to write down the probability of joint events or the joint distribution. Be aware of assumptions.!

$$P(\text{all tosses yield tails})$$

$$= P(T_1 T_2 \cdots T_{10})$$

$$= P(T_1) P(T_2) \cdots P(T_{10})$$

$$= \prod_{j=1}^{10} P(T_j)$$

**events**  $\longrightarrow T_j$  : toss  $j$  is a head

**rv**  $\longrightarrow X_j$  : number of heads in toss  $j$

**iid**  $X \sim \text{Bernoulli}(0.5)$

$$P(X = 1) = 0.5$$

$$P(X = 0) = 1 - 0.5$$

## Increasing abstraction .....

Coin comes up heads with probability  $p$ .  parameter

Toss it 10 times.

$A$  = at least one head

$T_j$  = toss  $j$  yields tails,  $j \in \{1, 2, \dots, 10\}$

$$P(T_j) = 1 - p$$

$$P(A) = 1 - P(\text{not } A)$$

$$= 1 - P(T_1 T_2 \cdots T_{10})$$

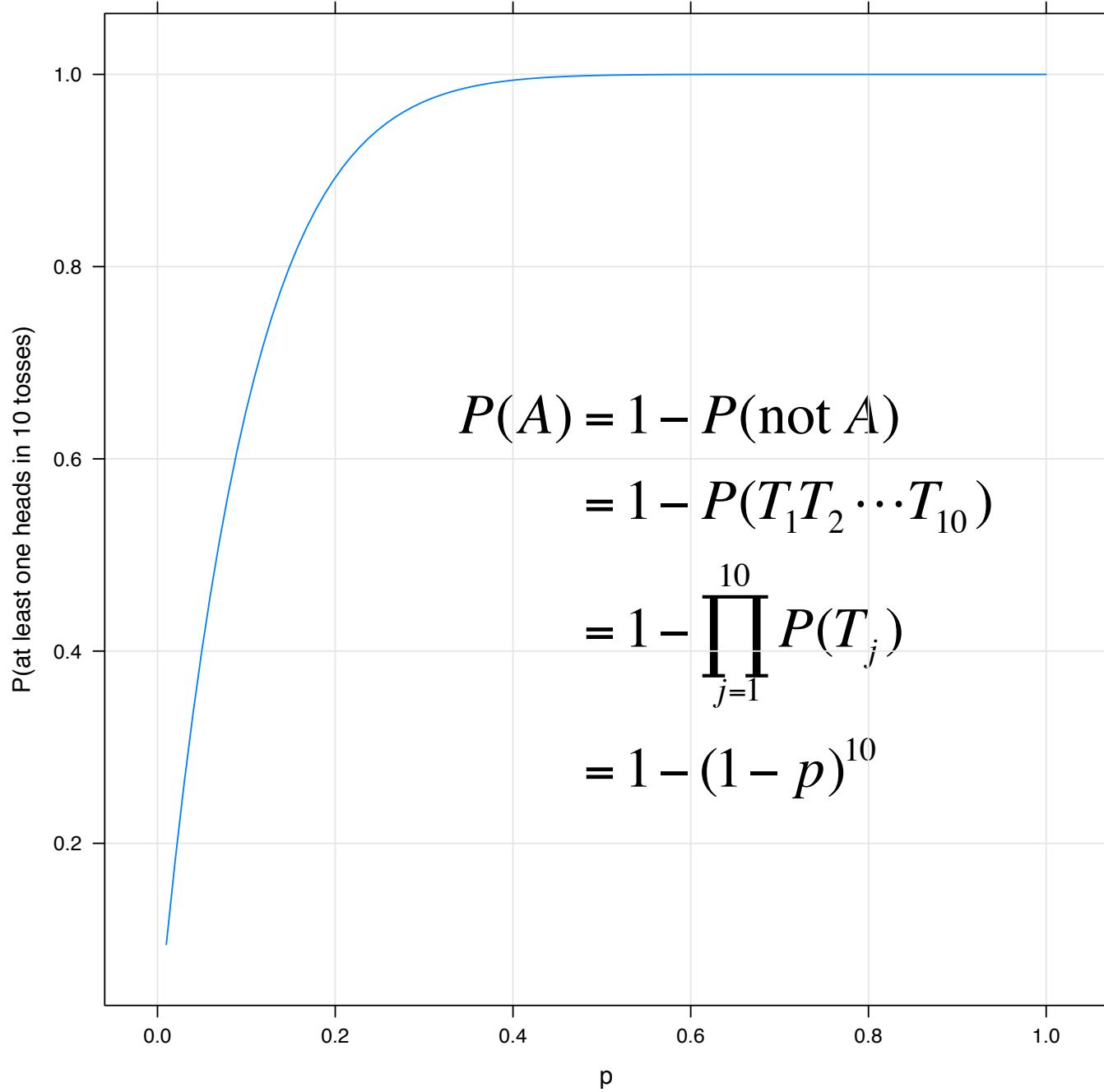
$$= 1 - \prod_{j=1}^{10} P(T_j)$$

$$= 1 - (1 - p)^{10}$$

$X$  : number of heads in 10 tosses

$$X \sim Bin(10, p)$$

$$P(X = 10) = (1 - p)^{10}$$



Increasing abstraction and sneaky foreshadowing of the incredible multiple testing problems faced in genomics.....

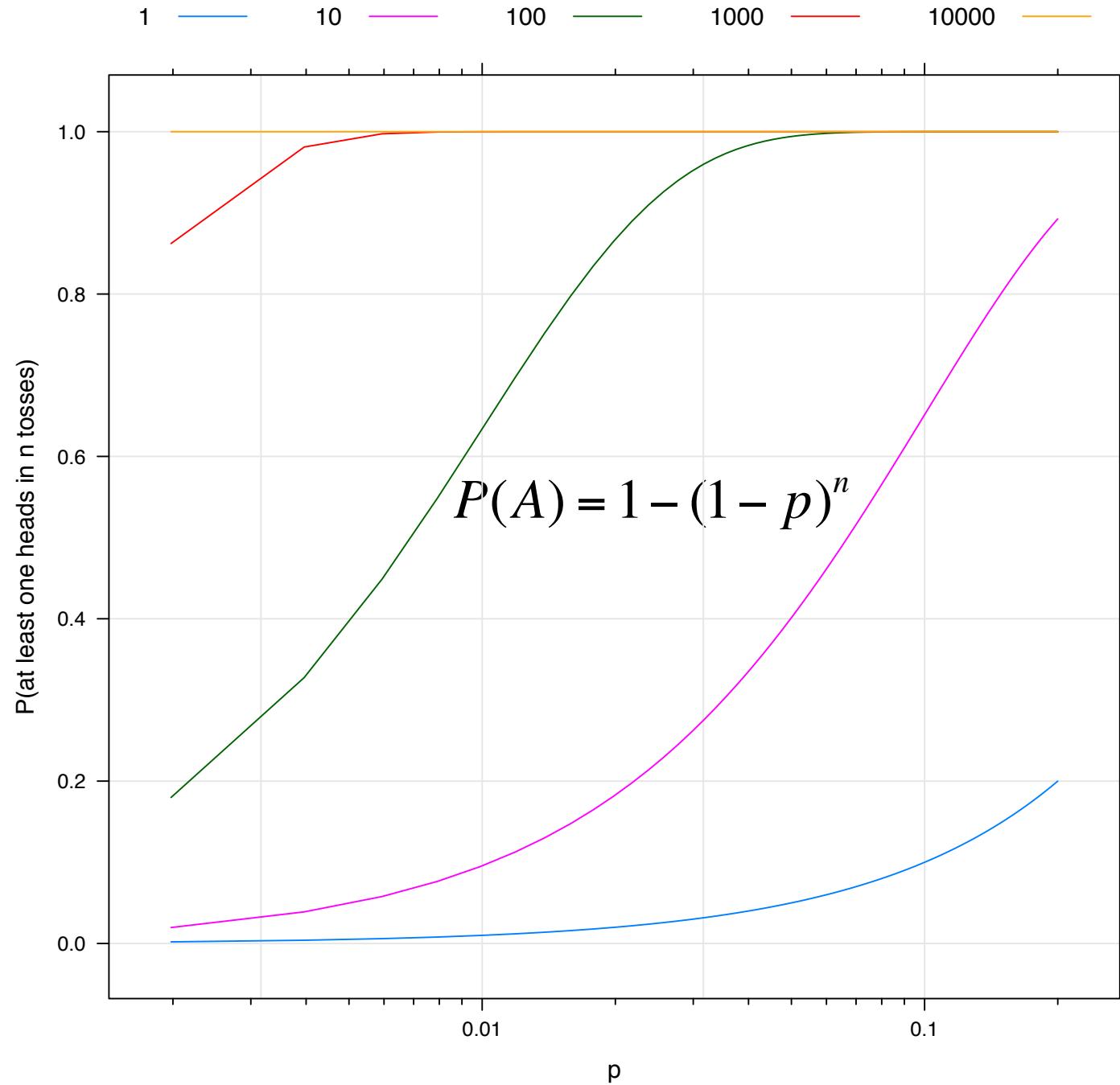
Coin comes up heads with probability  $p$ .  
Toss it  $n$  times.

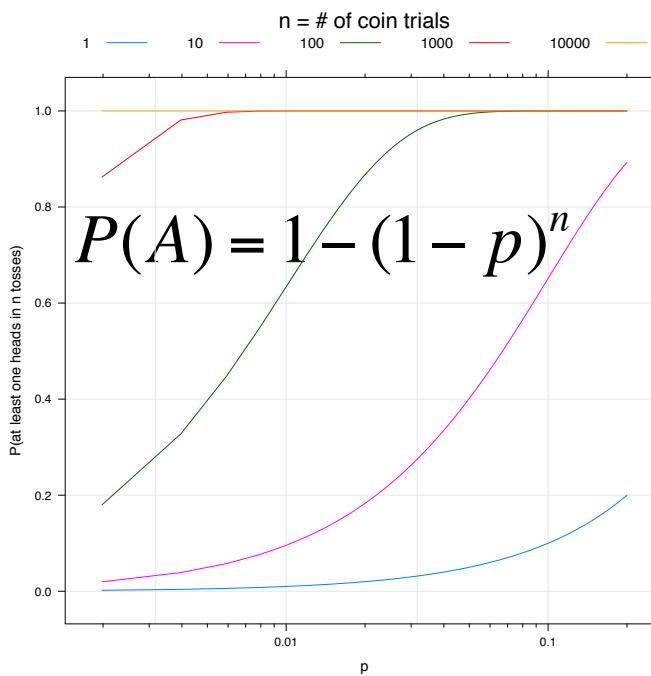
$A$  = at least one head

$$P(A) = \text{<same stuff as before, really>}$$

$$= 1 - (1 - p)^n$$

$n = \#$  of coin trials





In a genomics experiment...

What if “head” = false positive = false “significant” gene

Doing lots of tests today? Then I guarantee you’ll get a false positive. In fact, you’ll get *LOTS*.

This is the multiple testing problem and it is almost crippling in genomics. More on that later.

# Random variables can be characterized by a distribution

Following previous example...

$X$  : number of heads in  $n$  tosses

Variable

$$X \sim Bin(n, p)$$
$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

parameter

probability distribution

$$F_X(a) = P(X \leq a) = \sum_{x \leq a} p_X(x) \quad (\text{for a discrete } X)$$

$$\sum_{x=0}^a \binom{n}{x} p^x (1 - p)^{n-x}$$

cumulative distribution

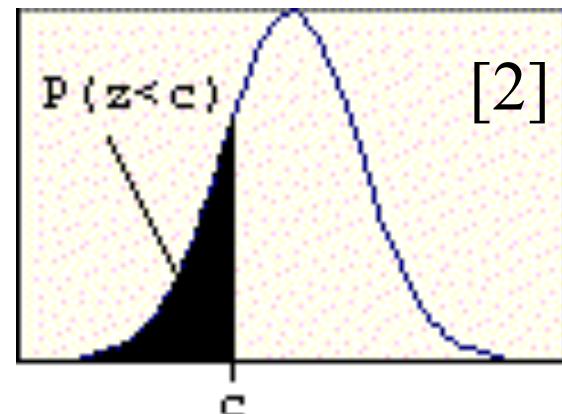
# how to get a probability from a density

$$[1] P(a < Y < b) = \int_a^b f_Y(y) dy$$

$$[2] P(Y \leq a) = \int_{-\infty}^a f_Y(y) dy$$

$$[3] P(Y \geq a) = \int_a^{\infty} f_Y(y) dy$$

$$[4] P(|Y| \geq a) = \int_{-\infty}^{-a} f_Y(y) dy + \int_a^{\infty} f_Y(y) dy$$



“cumulative distribution function”

## “cumulative distribution function (CDF)”

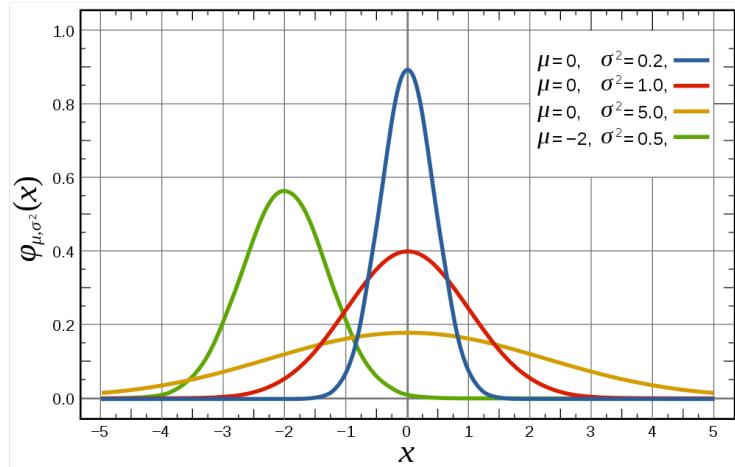
$$F_Y(a) = P(Y \leq a) = \int_{-\infty}^a f_Y(y) dy \text{ (for a continuous Y)}$$

$$F_Y(a) = P(Y \leq a) = \sum_{y_i \leq a} p_Y(y_i) \text{ (for a discrete Y)}$$

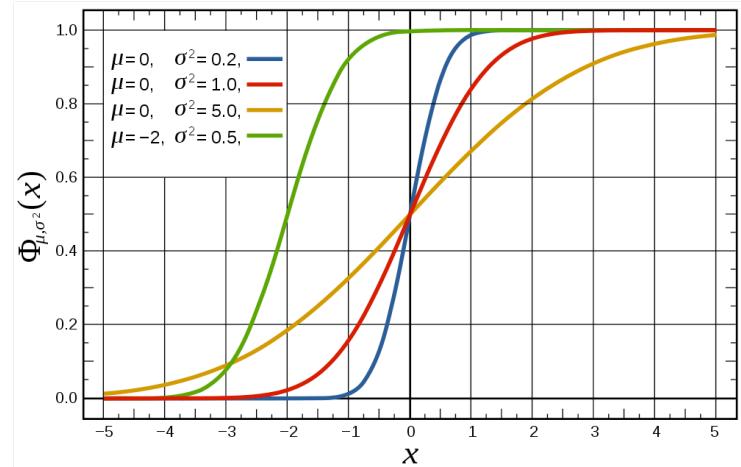
yes, we really do distinguish the density function  
(continuous rv) from the CDF with the  
deceptively subtle lowercase “*f*” vs. uppercase “*F*”

# density or prob. mass function

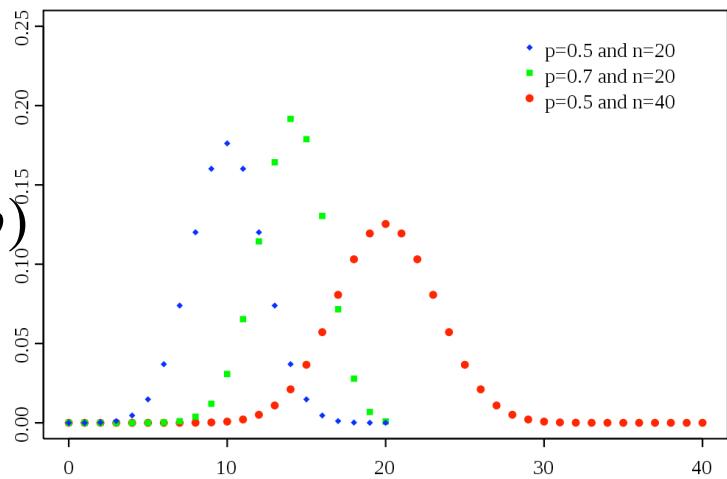
$N(\mu, \sigma^2)$



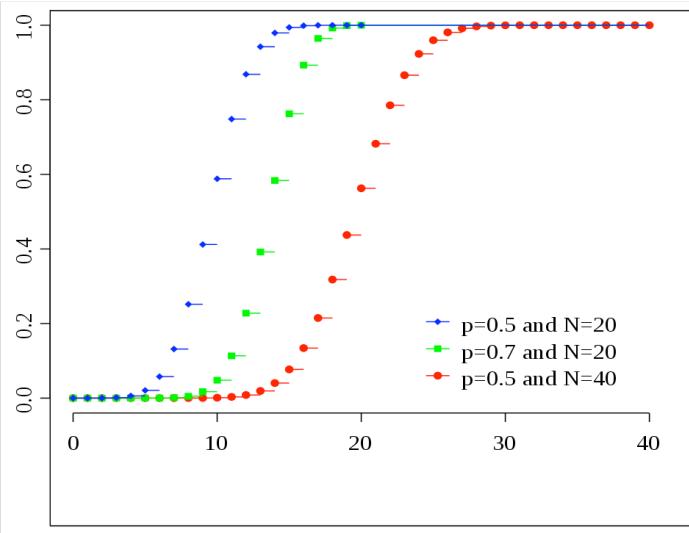
CDF



$Binom(n, p)$



$\Phi_{\mu, \sigma^2}$



# **sources of images on previous page**

[http://en.wikipedia.org/wiki/File:Normal\\_Distribution\\_PDF.svg](http://en.wikipedia.org/wiki/File:Normal_Distribution_PDF.svg)

[http://en.wikipedia.org/wiki/File:Normal\\_Distribution\\_CDF.svg](http://en.wikipedia.org/wiki/File:Normal_Distribution_CDF.svg)

[http://en.wikipedia.org/wiki/File:Binomial\\_distribution\\_pmf.svg](http://en.wikipedia.org/wiki/File:Binomial_distribution_pmf.svg)

[http://en.wikipedia.org/wiki/File:Binomial\\_distribution\\_cdf.svg](http://en.wikipedia.org/wiki/File:Binomial_distribution_cdf.svg)

# Why is it important to learn about probability distributions?

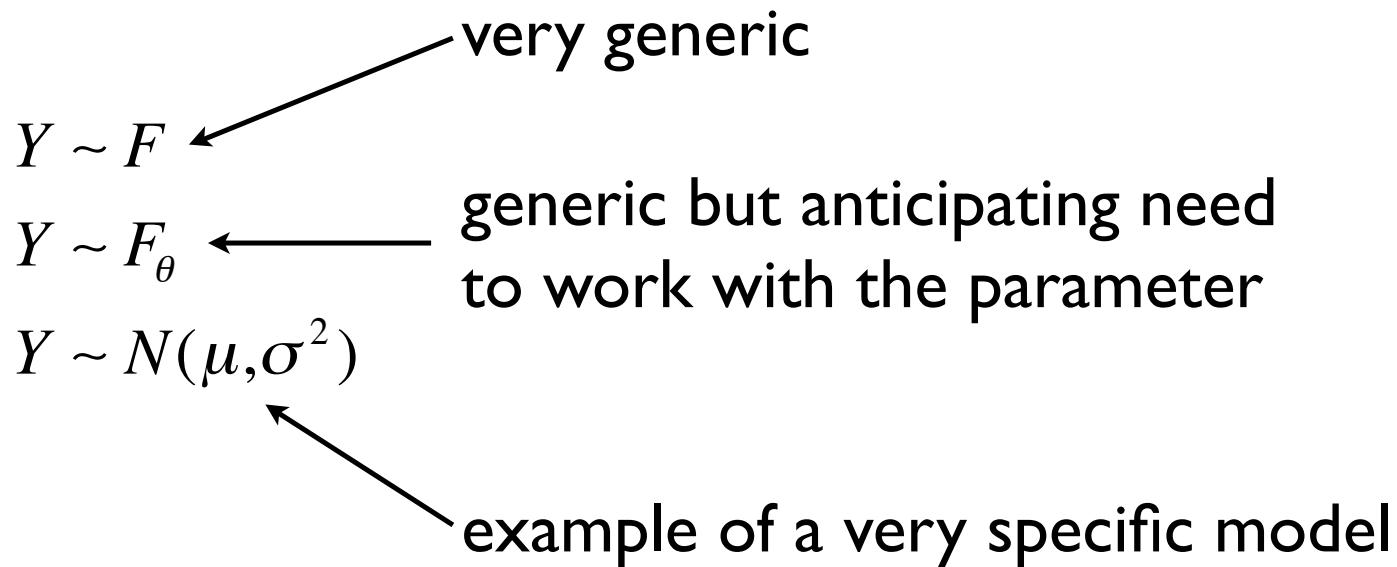
Given pmf/pdf of an r.v.  $X$ , we can:

- Compute the probability of various events, mean/variance of  $X$ , without having to perform experiments!
- We can simulate real systems and get the data.

we're starting to leave basic probability and  
transition into statistical inference .....

# First some vocabulary ...

## statistical model



a statistician doesn't mean much when they say  
“model” ... nothing terribly specific or mechanistic ...  
just specifying a probability distribution and, optionally,  
more details about the parameter(s)

## statistical model

the parameter space is the set of all possible values for the parameter

to say a model is “parametric” means the parameter space is a nice friendly Euclidean space

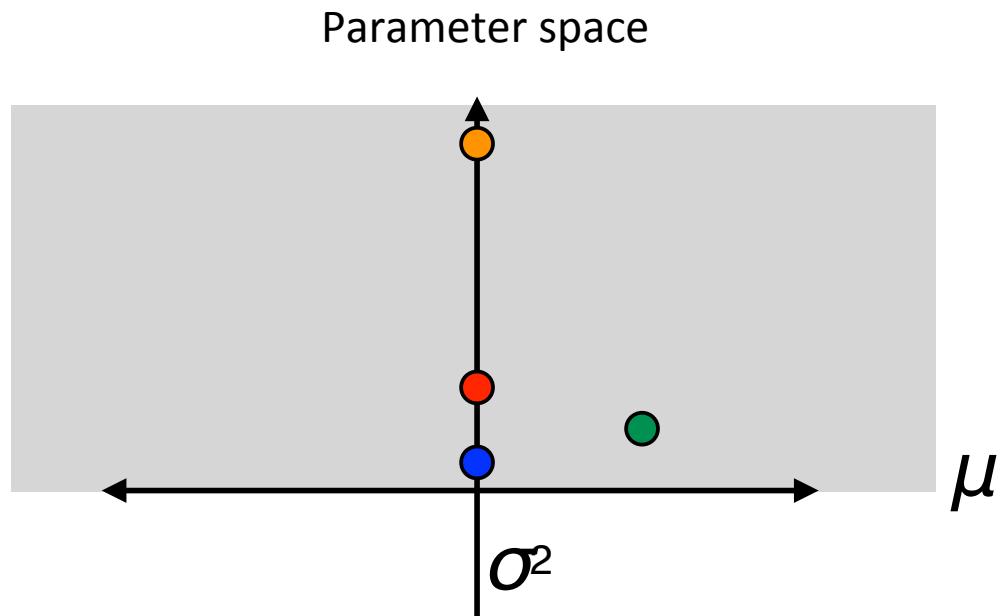
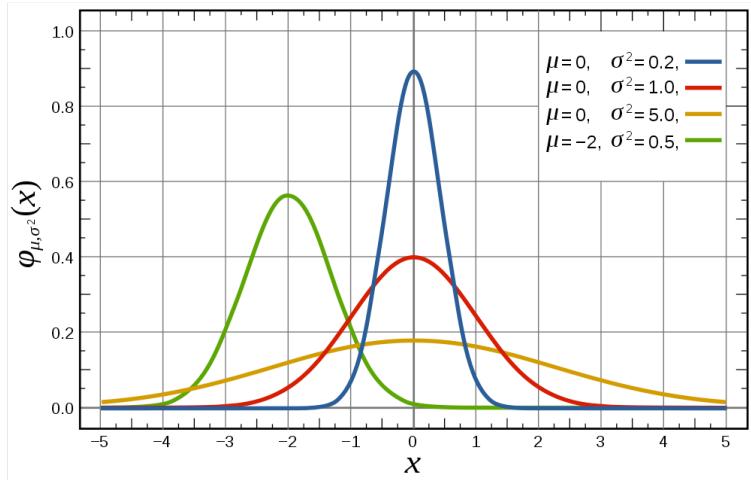
when we assume data is normally distributed about its mean ... we’re doing parametric inference; the parameter space is a nice friendly half-plane in  $\mathbb{R}^2$

# world's favorite parametric model

$$Y_1, \dots, Y_i, \dots, Y_n \sim F_{\theta} = N(\mu, \sigma^2)$$

$$\theta = (\mu, \sigma^2)$$

the parameter space, i.e. all possible values of  $\theta = (\mu, \sigma^2)$



parameter space = set of all possible values for the parameter

“model is parametric”  $\Leftrightarrow$  parameter space is a nice friendly Euclidean space

family	typical notation	parameter $\theta$
<generic>	$Y \sim F_\theta$	$\theta$
Bernoulli	$Y \sim \text{Bern}(p)$	$\theta = p$
binomial	$Y \sim \text{Bin}(n, p)$	$\theta = (n, p)$
uniform	$Y \sim \text{Unif}[a, b]$	$\theta = (a, b)$
Normal	$Y \sim N(\mu, \sigma^2)$	$\theta = (\mu, \sigma^2)$
Student's t	$Y \sim t_{df}$	$\theta = df$

*Parametric  
models we've  
reviewed ....*

“semi-parametric” and “nonparametric” imply the parameter space isn’t a simple Euclidean space

means the parameter space is more exotic, e.g. at least partially a function space, an infinite dimensional space

BUT one does not have to feel comfortable with, say, function spaces, to *apply* nonparametric statistical methods (e.g. rank based procedures like the Wilcoxon test) responsibly

"Let  $(Y_1, Y_2, \dots, Y_n)$  be independent, identically distributed random variables."

or

" $Y_i \sim F$ "

the two parameters of any distribution  $F$  you're mostly like to care about

- #1: it's expected value (or expectation or mean)
- #2: it's variance

expectation, expected value, the mean

it is a parameter

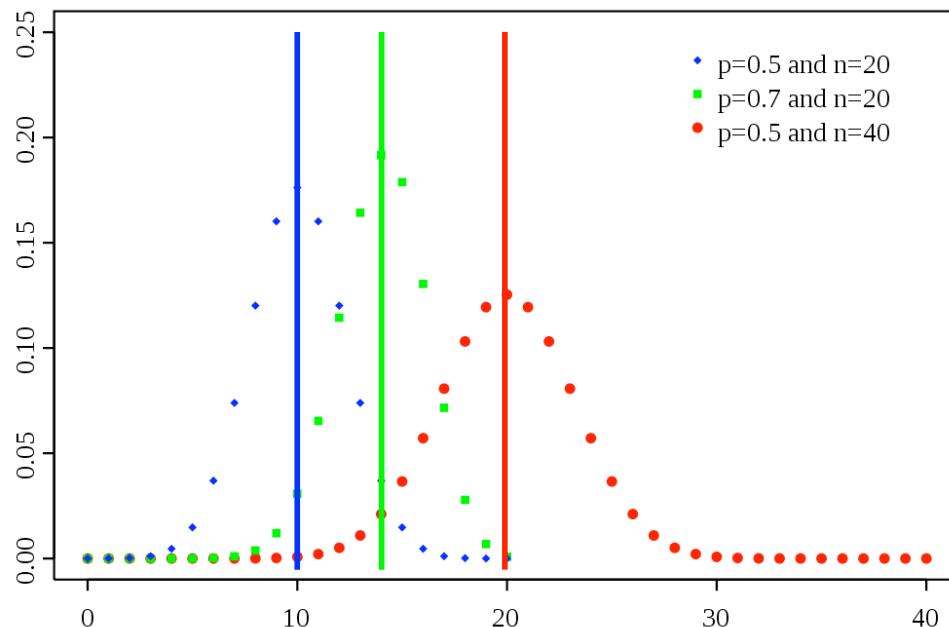
often denoted  $E(Y)$  or  $\mu$  or  $\mu_Y$

common sense “definition”: a long-run average  
 $E(Y)$  approx equal to (sum of  $Y_i$ ’s)/ $n$   
the bigger  $n$  is, the better the “approximation”

# expectation, expected value, the mean

$$E(Y) = \sum_y y p_Y(y) \text{ for discrete rv } Y$$

$$E(Y) = \int y f_Y(y) dy \text{ for continuous rv } Y$$



binomial example:

$$Y \sim Binom(n, p)$$

$$E(Y) = np$$

the mean is a measure of “location”

often is one of the “obvious” parameters (e.g. normal)  
or is easily computed from them (e.g. binomial)

and now something related *but different*

the average, the sample mean

it is a random variable!!!! not a parameter!!!

often denoted  $\bar{Y}$  or  $\bar{Y}_n$  or  $\hat{\mu}$  or  $\hat{\mu}_Y$

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$$

main usage:

as a point estimator of the true mean or as a test statistic -- or part of a test statistic -- for hypothesis tests re: the mean

$$\bar{Y} \text{ or } \bar{Y}_n \text{ or } \hat{\mu} \text{ or } \hat{\mu}_Y$$

notational sidebar:

statisticians LOVE to put hats on Greek letters

a constant visual reminder of what's random (the thing with the hat) and which parameter it is an estimator for (the Greek letter without the hat)

sometimes we put the sample size  $n$  in the subscript to reinforce that something is random and that its distribution depends on the sample size

the expected value of the sample mean is the true mean:

$$E(\bar{Y}_n) = \mu$$

“the sample mean is unbiased”

the variance of the sample mean is:

$$V(\bar{Y}_n) = \frac{\sigma^2}{n}$$

i.e., the variance of the sample mean is fundamentally determined by the underlying variance of the data --it is also affected by the sample size

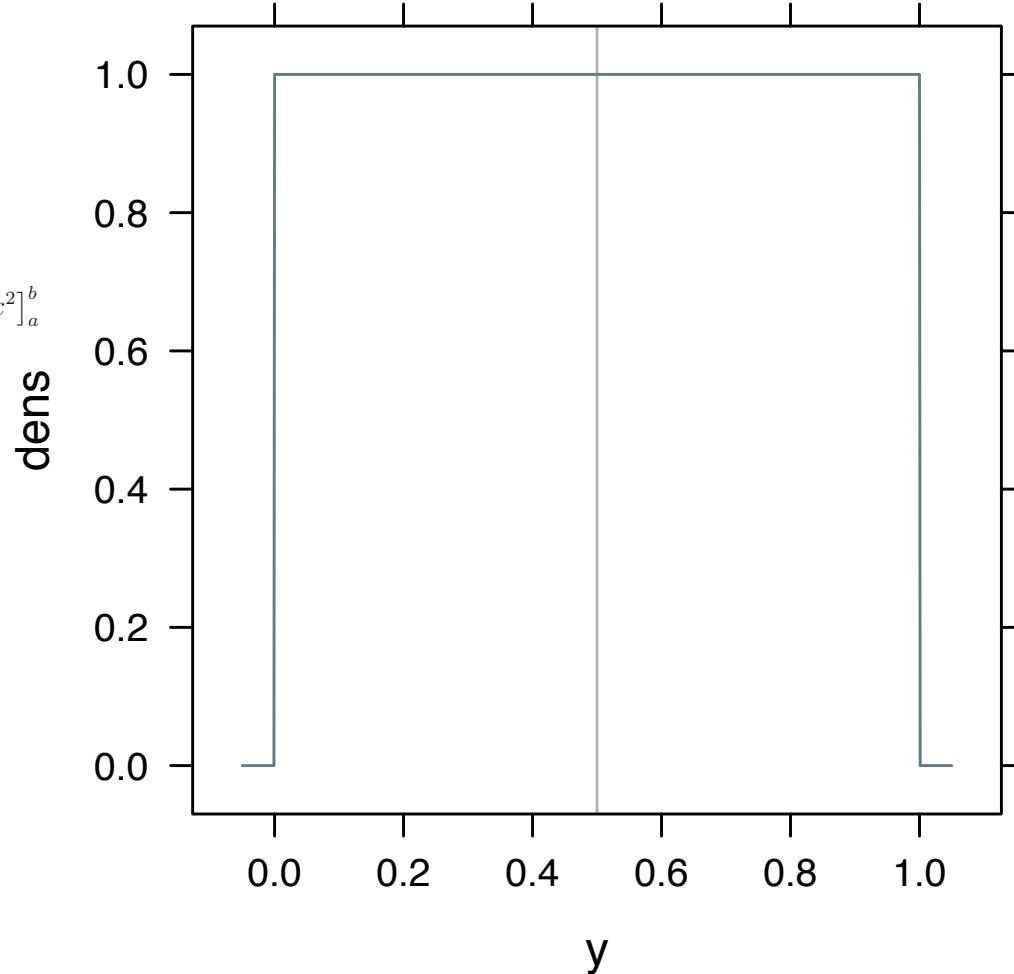
# the average, the sample mean: an empirical case

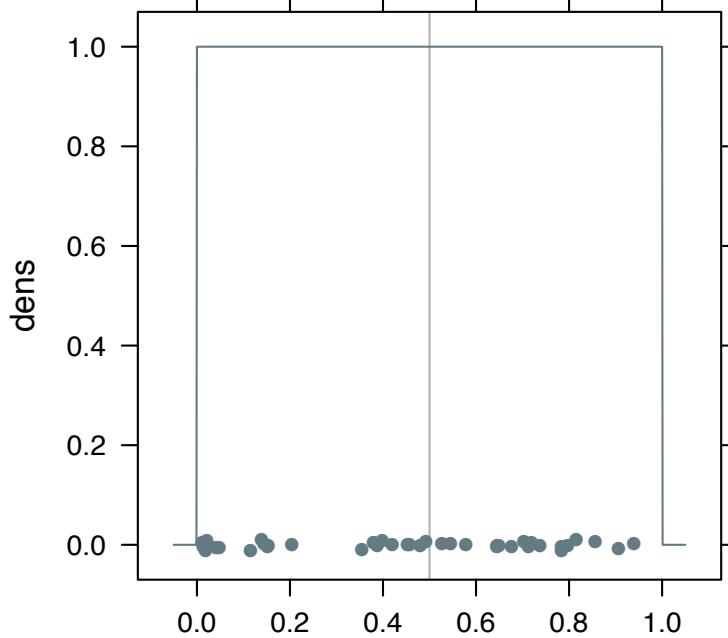
consider  $Y \sim \text{Unif}(0, 1)$

$$E(Y) = 0.5$$

$$E(X) = \int xf_x(x)dx$$

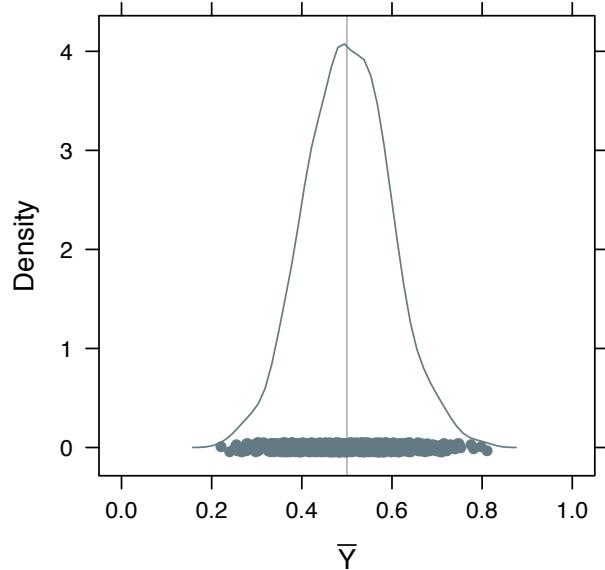
$$\begin{aligned} E(X) &= \int_{-\infty}^{\infty} xf(x)dx = \int_a^b x \frac{1}{b-a} dx = \frac{1}{2(b-a)} [x^2]_a^b \\ &= \frac{b^2 - a^2}{2(b-a)} \\ &= \frac{b+a}{2} \end{aligned}$$





take a sample of size n

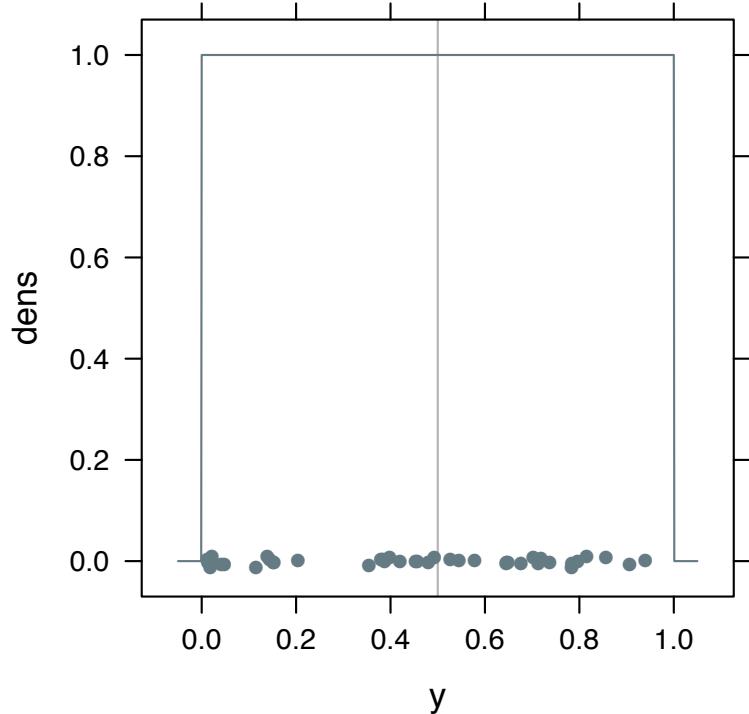
e.g. (0.3365 , 0.1733 , 0.0861, 0.3933 , 0.8044 , 0.0111, 0.2331, 0.9339, 0.2268, 0.7859)  $\bar{Y} = 0.3984$



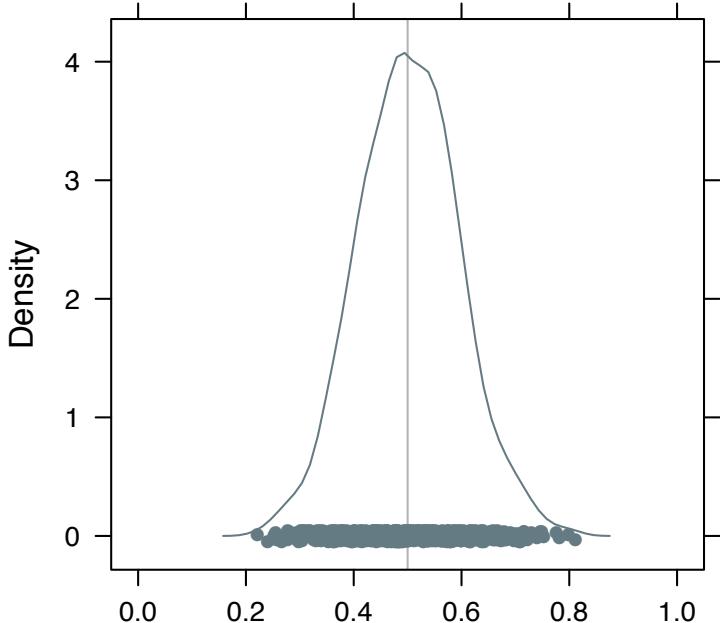
take the average

... now do that lots of times ...

what does that distribution look like?



visual confirmation of  
 $E(\bar{X}_n) = \mu$   
= 0.5 in this case

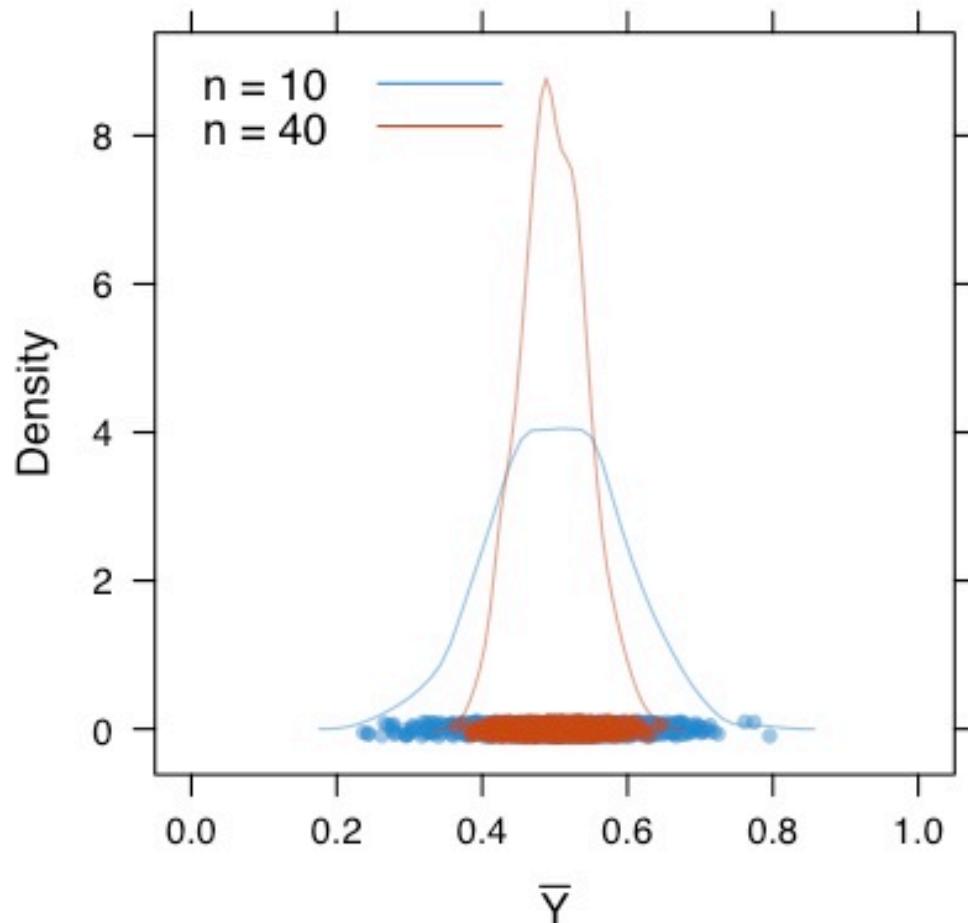


notice that the distribution of  
sample means doesn't look  
uniform .... Central Limit  
Theorem coming soon!

```
> n <- 10
> numSamp <- 1000
> xBar <- rowMeans(matrix(runif(n * numSamp), nrow = numSamp))
> n2 <- 40
> xBar2 <- rowMeans(matrix(runif(n2 * numSamp), nrow = numSamp))
> densityplot(~ xBar + xBar2, ...)
```

visual confirmation of

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$



the variance of the sample mean is

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$

i.e., the variance of the sample mean is fundamentally determined by the underlying variance of the data --it is also affected by the sample size

this is why it is nonsensical to ask if a sample size of  $n = 3$  (or 20 or whatever) is “enough” to perform statistical inference, in the absence of some info on  $\sigma^2$  (and specific discovery goals)

discouraging that the variance of the sample mean involves  $\sigma^2$  which we generally don’t know ....

what if you want to know more about  $Y$ ?

the mean gives a sense of what's "typical" or where the center of observed values of  $Y$  will lie, but what sort of spread will those observed values have?

the sample mean is obviously a good guess at  $E(Y) = \mu$ , but how good is it?

variance

standard deviation =  $\sqrt{\text{variance}}$

it is a parameter

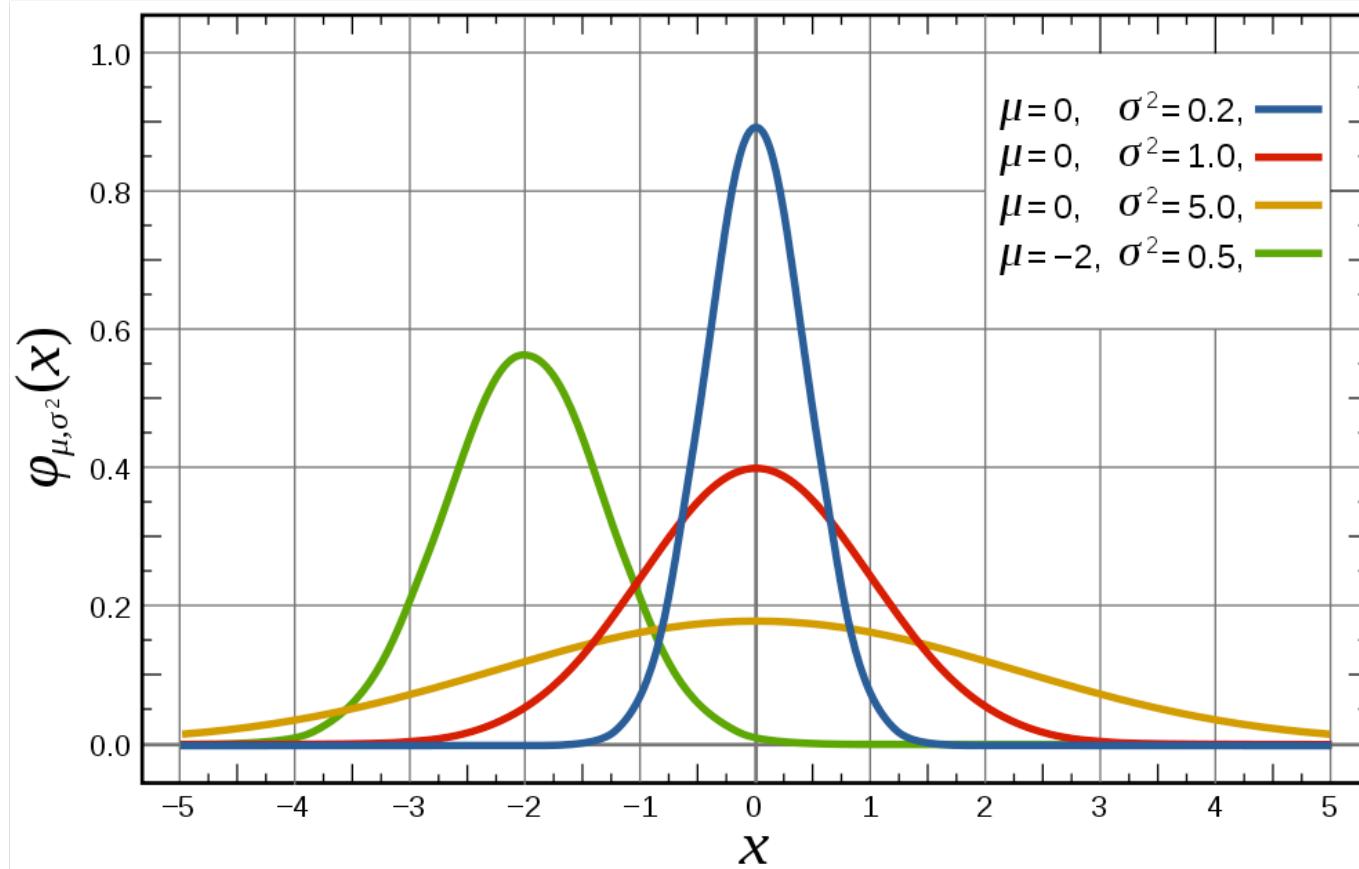
usually denoted  $V(X)$  or  $\sigma^2$  (variance) and  $\sigma$  (sd)

$$V(Y) = E(Y - \mu)^2$$

common sense “definition” of variance:  
a long-run average of the squared differences  
between obs vals  $Y = y$  and the true mean  $\mu$

variance

standard deviation =  $\sqrt{\text{variance}}$



normal as example; bigger  $\sigma^2 \leftrightarrow$  bigger “spread”

and now something related *but different*

## the sample variance

it is a random variable!!!! not a parameter!!!

often denoted  $s^2$  or  $\hat{\sigma}^2$  \*

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 \quad E(s^2) = \sigma^2$$

main usage:

I'm using the sample mean to infer something about the true mean and, to my horror, the quality of that guess depends on the variance. So I'm forced to worry about the variance. ("nuisance parameter")

\*  $n$  vs.  $n - 1$  sidebar

a “statistic” is a rv that’s a function of the data

classic examples:

- the sample mean
- the sample variance

two main reasons we love them:

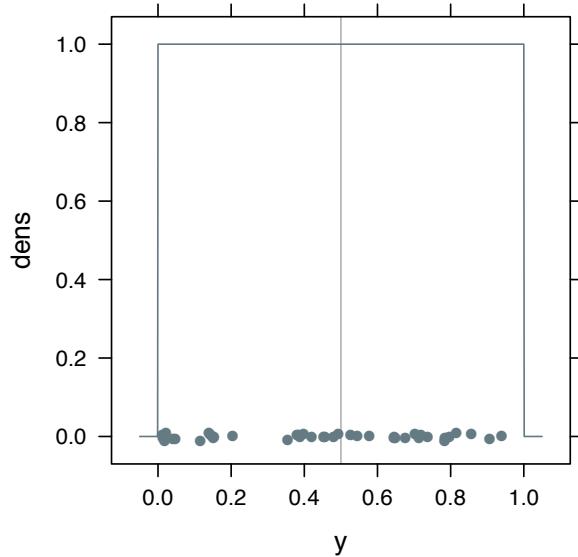
[1] sometimes they are estimators for parameters we care about

[2] sometimes they are the basis for a hypothesis test

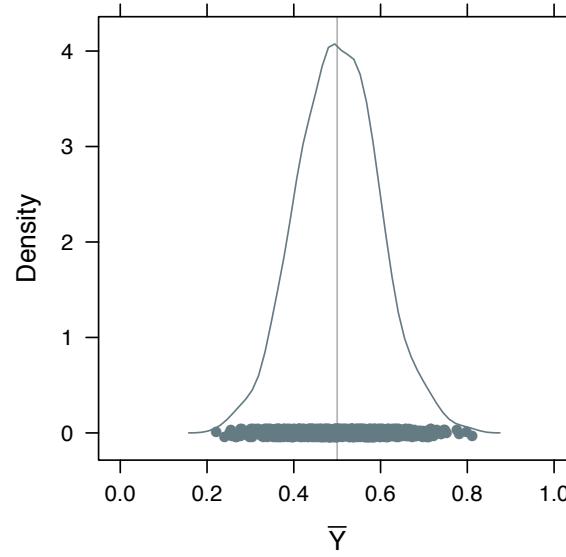
the distribution of a statistic is called its sampling distribution. it’s related to the distribution of the data but it is not the same

visual confirmation that distribution of a statistic -- such as sample mean here -- is NOT the same as the distribution of the underlying data

dist'n of data  
 $\text{Unif}(0, 1)$



dist'n of the sample average  
... NOT  $\text{Unif}(0, 1)$



... looks kind of like the normal dist'n, no?

we generally know more about a statistic's sampling distribution as  $n$  gets large

“large sample theory”, “limit theory”, “asymptotic theory”

# the law of large numbers

common sense “statement”:  
the average of a large, iid sample will be close to the true mean

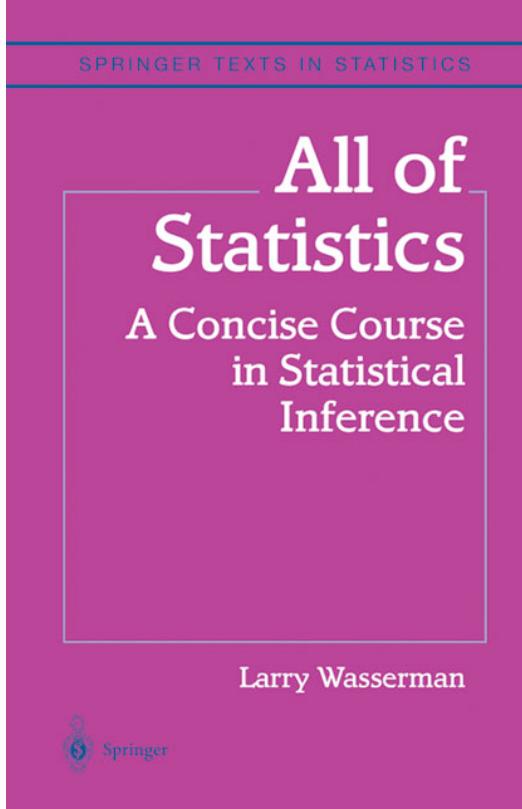
## the law of large numbers (formally)

Let  $X_1, X_2, \dots$  be an IID sample, let  $\mu = \mathbb{E}(X_1)$  and  $\sigma^2 = \mathbb{V}(X_1)$ . Recall that the sample mean is defined as  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$  and that  $\mathbb{E}(\bar{X}_n) = \mu$  and  $\mathbb{V}(\bar{X}_n) = \sigma^2/n$ .

**5.6 Theorem** (The Weak Law of Large Numbers (WLLN)).<sup>3</sup>

If  $X_1, \dots, X_n$  are IID, then  $\bar{X}_n \xrightarrow{\text{P}} \mu$ .

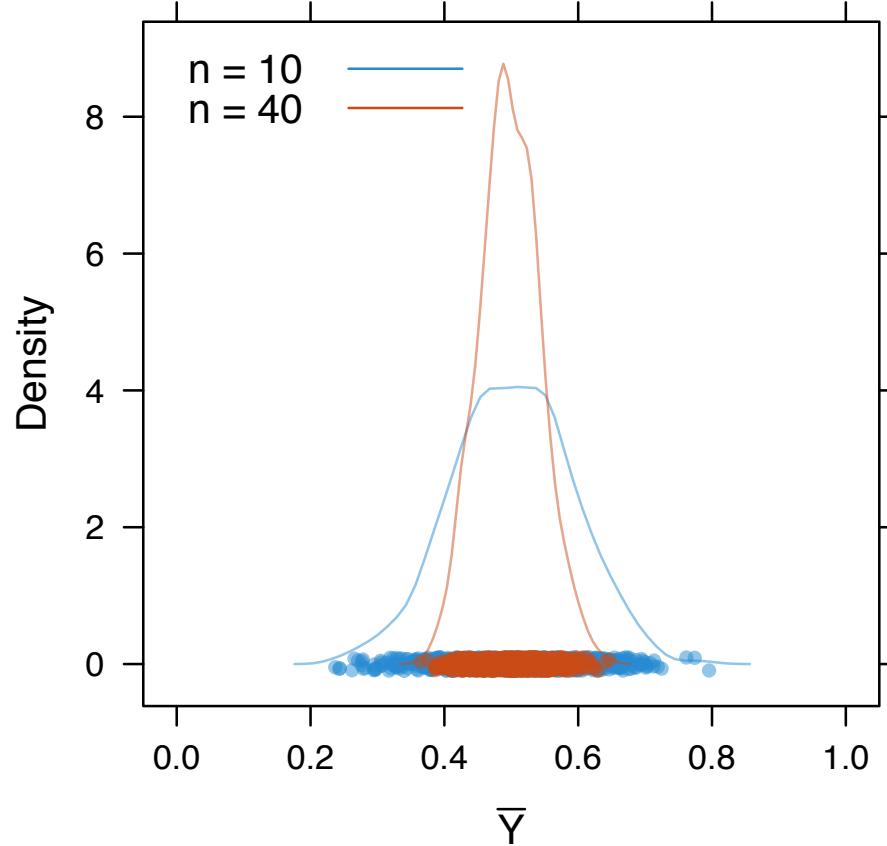
Interpretation of the WLLN: The distribution of  $\bar{X}_n$  becomes more concentrated around  $\mu$  as  $n$  gets large.



Available via SpringerLink!

Find bibliographic info and links  
on my resources page.

Imagine this trend continuing as  $n$  gets bigger and bigger ..... the sample mean sampling dist'n gets more and more concentrated around  $\mu = 0.5$



# the central limit theorem

common sense “statement”:

the sampling distribution for the average of a large, iid sample will be approximately a normal distribution

**5.8 Theorem** (The Central Limit Theorem (CLT)). *Let  $X_1, \dots, X_n$  be IID with mean  $\mu$  and variance  $\sigma^2$ . Let  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ . Then*

$$Z_n \equiv \frac{\bar{X}_n - \mu}{\sqrt{\mathbb{V}(\bar{X}_n)}} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow Z$$

where  $Z \sim N(0, 1)$ . In other words,

$$\lim_{n \rightarrow \infty} \mathbb{P}(Z_n \leq z) = \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

rookie misconception re: law of large numbers:

“If I can just make my sample big enough, I won’t have to worry about error.”

there is no sample that is “big enough” in an unqualified sense

in stats, there are precious few fundamental constants, like there are in math (think:  $\pi$  and e) or physics (think: speed of light)

context and goals always matter

rookie misconception re: central limit theorem:

“I can average any large-ish bunch of numbers and divide by the sd and call it a z-score. Then I can compare it to a  $N(0,1)$  to determine statistical significance. I’ve got a hit if the number’s greater than 1.96!”

the CLT assumes you’re averaging observations that are iid!

averaging gene expression for 1 gene across exchangeable subjects ... yeah, CLT applies

averaging gene expression for 1 subject across genes ... no, CLT does not apply (or, at least, you’ll have to convince me)

# Exercise: solve for expected value and variance of uniform distribution $\text{Unif}(a,b)$

First: recall that

$$E(X) = \int xf_x(x)dx$$

$$\begin{aligned}\text{Var}(X) &= E[(X - E[X])^2] \\ &= E[X^2 - 2X E[X] + (E[X])^2] \\ &= E[X^2] - 2E[X]E[X] + (E[X])^2 \\ &= E[X^2] - (E[X])^2\end{aligned}$$

Answers:

$$\begin{aligned}E(X) &= \int_{-\infty}^{\infty} xf(x)dx = \int_a^b x \frac{1}{b-a} dx = \frac{1}{2(b-a)} [x^2]_a^b \\ &= \frac{b^2 - a^2}{2(b-a)} \\ &= \frac{b+a}{2}\end{aligned}$$

$$\begin{aligned}V(X) &= E(X^2) - [E(X)]^2 \\ &= \int_a^b x^2 \cdot \frac{1}{b-a} dx - \left(\frac{b+a}{2}\right)^2 = \frac{1}{3(b-a)} [x^3]_a^b - \left(\frac{b+a}{2}\right)^2 \\ &= \frac{b^3 - a^3}{3(b-a)} - \left(\frac{b+a}{2}\right)^2 \\ &= \frac{b^2 + ab + a^2}{3} - \frac{b^2 + 2ab + a^2}{4} \\ &= \frac{(b-a)^2}{12}\end{aligned}$$

we have completely arrived at statistical inference  
now (vs. building our probability foundation)

canonical breakdown of typical statistical inference activities:

## hypothesis testing vs. estimation

in either case, you are trying to say something intelligent about a parameter

hyp testing: does the true value of the parameter lie in an exciting or boring part of the parameter space?

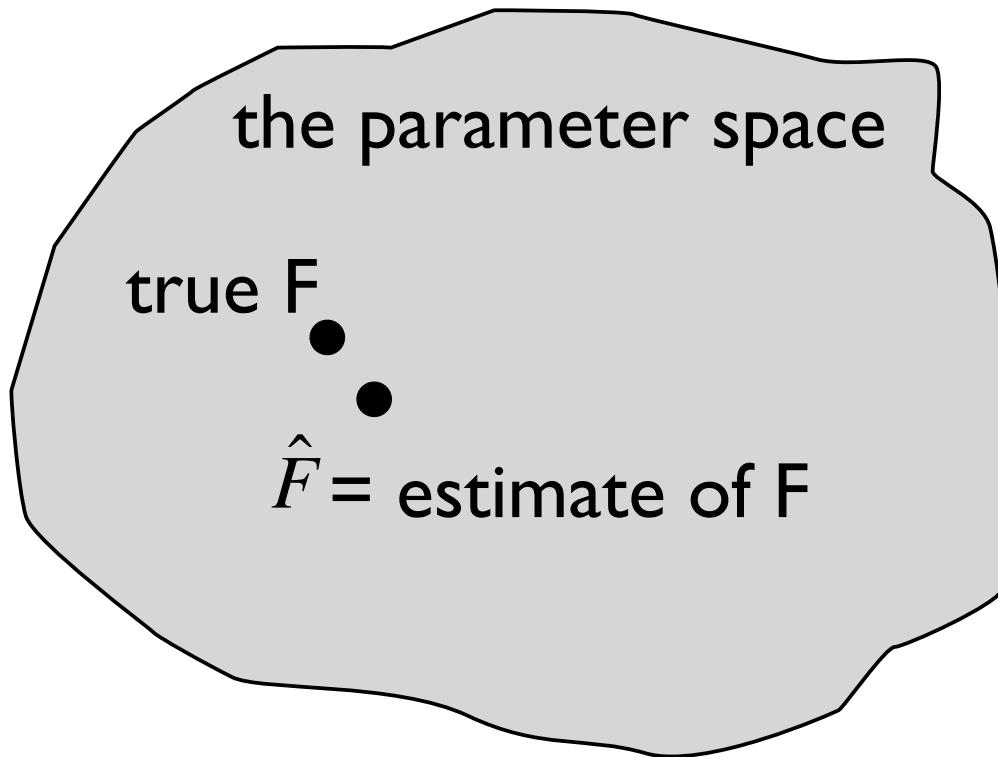
estimation: what's your best guess at the true value of the parameter?

# estimation in generic statistical model

$$Y_1, \dots Y_i, \dots, Y_n \sim F$$

Observe data ( $Y_1 = y_1, \dots Y_i = y_i, \dots Y_n = y_n$ ).

Estimate  $F$  with  $\hat{F}$ .



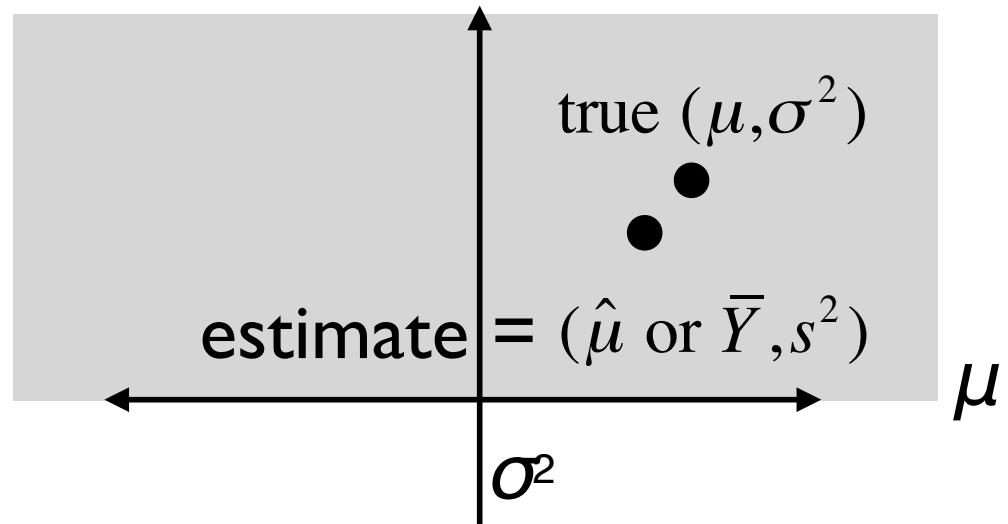
# estimation in very specific statistical model

$$Y_1, \dots, Y_i, \dots, Y_n \sim F = N(\mu, \sigma^2)$$

Observe data ( $Y_1 = y_1, \dots, Y_i = y_i, \dots, Y_n = y_n$ ).

Estimate  $F$ , i.e. estimate the mean  $\mu$  and the variance  $\sigma^2$ .

the parameter space



# hypothesis testing in high-throughput experiments

~thousands of individual “cases” being studied in a massively parallel fashion

e.g., expression level of each individual gene in a genome under two different conditions, A and B

some genes -- presumably a small minority -- are truly “interesting” (Efron) or “alternative”, i.e. expression levels are different in condition A vs. condition B

the rest -- presumably most genes -- are truly boring (?) or “null”

# hypothesis testing in high-throughput experiments

typical analytical goal:

based on observed, messy data, guess which genes are interesting and which are not and characterize the quality of your guessing

there's no magic from the “high-throughput” nature of this data (hurts more than helps, actually)

must begin with a clear understanding of how to do this for one gene and two conditions

then ... extend to more genes, more conditions

$$Y_1, \dots, Y_i, \dots, Y_{n_y} \sim \text{iid } F$$
$$Z_1, \dots, Z_i, \dots, Z_{n_z} \sim \text{iid } G$$
**testing**

Observe data  $(Y_1 = y_1, \dots, Y_i = y_i, \dots, Y_n = y_{n_y})$  and

$(Z_1 = z_1, \dots, Z_i = z_i, \dots, Z_n = z_{n_y})$ .

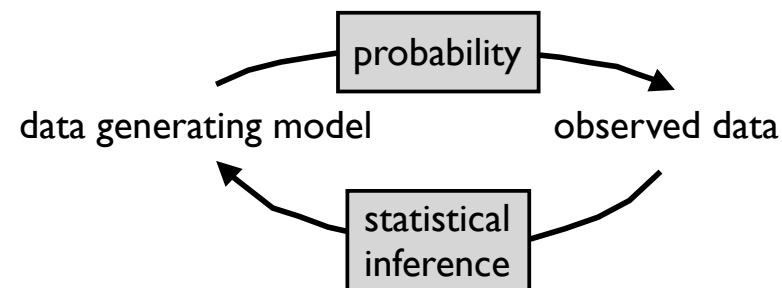
Does  $F = G$ ? OK, I'll settle for ...

does  $E_F(Y) = \mu_Y = \mu_Z = E_G(Z)$ ?

Call this statement the null hypothesis  $H_0$ :

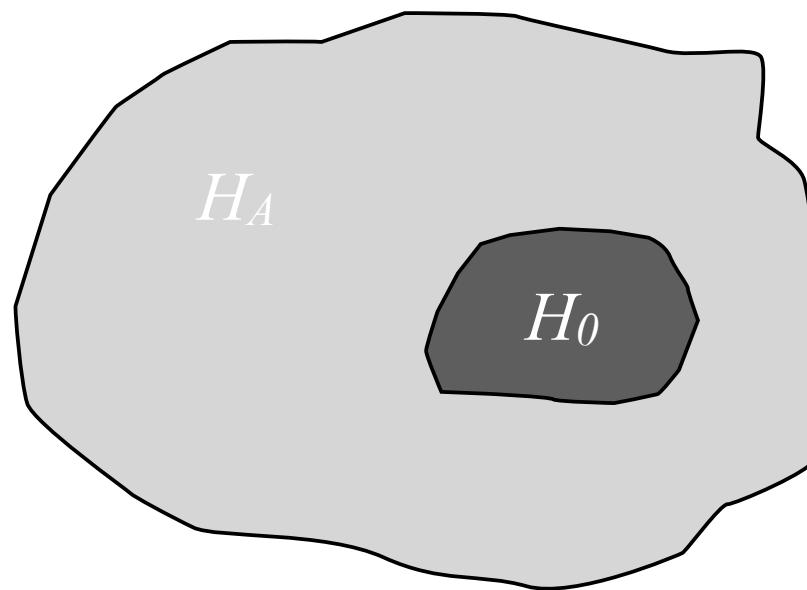
$$H_0 : \mu_Y = \mu_Z$$

Or, equivalently:

$$H_0 : \mu_Z - \mu_Y = 0$$


## statistical model

the parameter space



In formal hypothesis testing:

Define a “null (boring) region” for the parameter --  
the dark gray area.

Ask whether the true value lies in that region or  
outside, in the “alternative (interesting) region” --  
the light gray area.

# testing in world's favorite statistical model

$$Y_1, \dots, Y_i, \dots, Y_{n_y} \sim F = N(\mu_Y, \sigma^2)$$

$$Z_1, \dots, Z_i, \dots, Z_{n_z} \sim G = N(\mu_Z, \sigma^2)$$

Observe data  $(Y_1 = y_1, \dots, Y_i = y_i, \dots, Y_n = y_{n_y})$  and  
 $(Z_1 = z_1, \dots, Z_i = z_i, \dots, Z_n = z_{n_z})$ .

Does  $F = G$ ? OK, I'll settle for ...

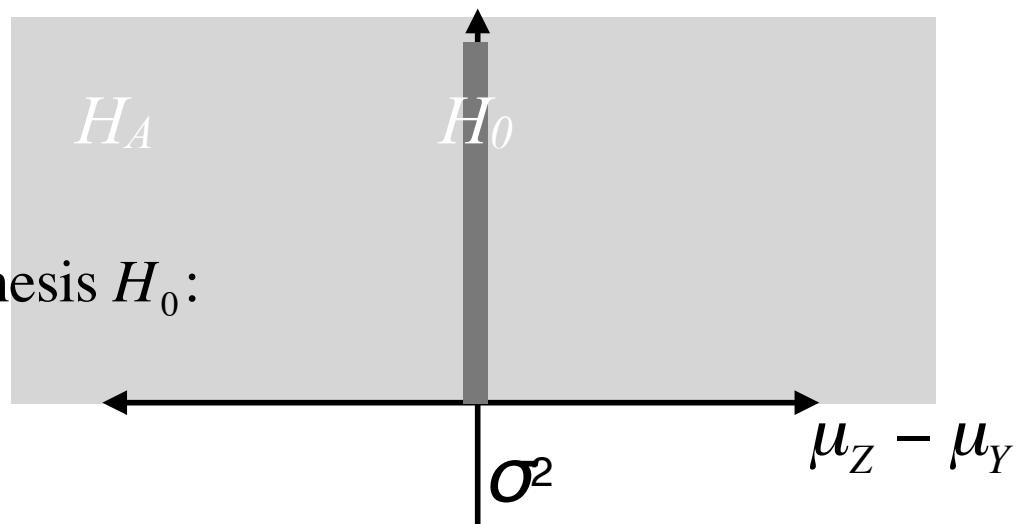
does  $E_F(Y) = \mu_Y = \mu_Z = E_G(Z)$ ?

Call this statement the null hypothesis  $H_0$ :

$$H_0 : \mu_Y = \mu_Z$$

Or, equivalently:

$$H_0 : \mu_Z - \mu_Y = 0$$

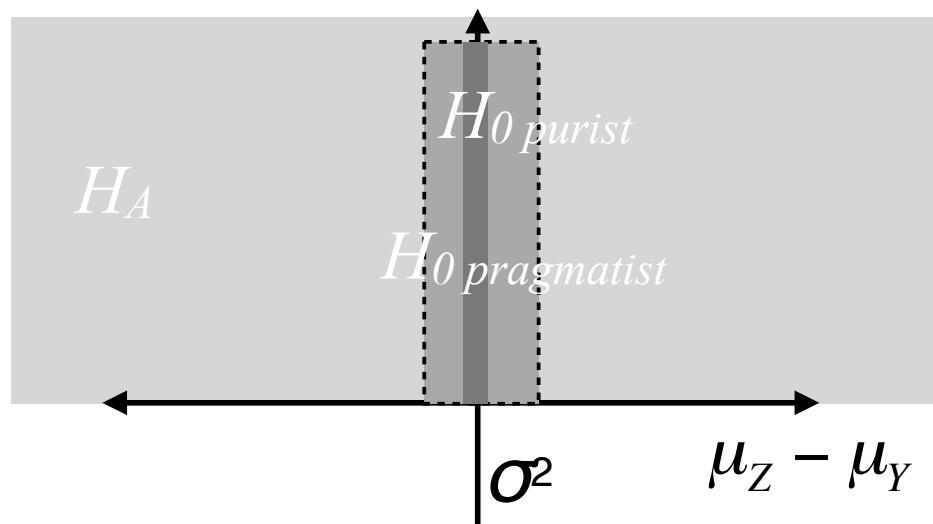


# reality check re: null and alternative regions/hypotheses

“purist” defines null region as half-line where  $\mu_Z - \mu_Y$  equals exactly zero

“realist” knows that the null region is a *neighborhood* around zero -- there are some differences too small to care about

“pragmatist” usually defines the null region like the “purist”, because the math is so much more tractable and then accounts for concerns of “realist” when interpreting results (or, e.g., does a post hoc filter on observed difference in sample means)



# Parameter estimation

# Parameters determine distributions

- When sampling from a population described by a pmf/pdf  $f(X|\theta)$ , then knowledge of  $\theta$  yields knowledge of the entire population.
- This is why parameter estimation is useful:
  - If we are tossing a coin, we would like to estimate the parameter  $p$

# Parameter estimation

- **Estimator:** rule/function whose calculated value is used to estimate the parameter
- **Estimate:** A particular realization of the estimator
- **Types of estimators:**
  - Point estimate: single number that can be regarded as the most plausible value of the parameter
  - Interval estimate: range of numbers, likely contain the true value of the parameter

# Methods of point estimation

- (Methods of moments)
- Maximum likelihood estimation (MLE)
- Bayesian Inference

# What are the properties of a good estimator?

- How well does the resulting estimate *explain* the “real world”?
- Proposed by geneticists/statisticians: Sir Ronald A Fisher in 1922
- Idea: we attempt to find the values of the parameters which would most likely produced the data that we in fact observed.

# What is *Likelihood*?

- **Before** we perform an experiment, the outcome is unknown. Probability density function allows us to predict the probability of any outcome based on known parameters:
  - $P(\text{Data} \mid \theta)$
- For example, say we know the probability of getting a head in a coin toss is  $p=0.6$ 
  - Then we can calculate the probability of any outcome:

$$D_1 = \{\text{HTHHHTHHHT}\} \quad P(D \mid \theta) = p^7(1-p)^3$$

$$D_2 = \{\text{HTH}\} \quad P(D \mid \theta) = p^2(1-p)$$

$$D_3 = \{\text{TTTH}\} \quad P(D \mid \theta) = p^3(1-p)$$

.

# What is *Likelihood*?

- **After** the experiment is done, we know the outcome. Now we want to know the *likelihood* that a given parameter value would generate the outcome:  
 $L(\text{Data} \mid \theta) : p(\text{Data} \mid \theta)$
- **Estimate**  $\theta$  by finding the value of  $\theta$  that makes the data most *likely* (our estimate:  $\hat{\theta}$ )

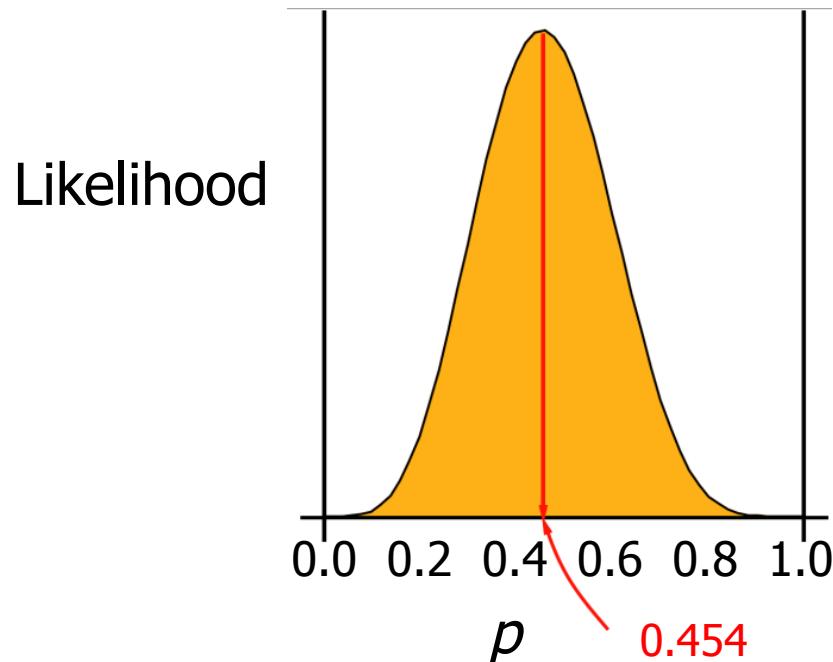
# The coin example

- We have data from 11 tosses of a coin (we don't know  $p$  the probability of head)
  - RV = outcome is head
- Outcome of the experiment: {HHTHTTTHTTH}
- Probability of the outcome of the experiment:  
$$pp(1-p)p(1-p)(1-p)(1-p)p(1-p)(1-p)p$$
- The likelihood is  $L(\text{Data} \mid p) = p^5(1-p)^6$

- We can plot the data against it's likelihood to figure out when we reach the maximum of the likelihood function.

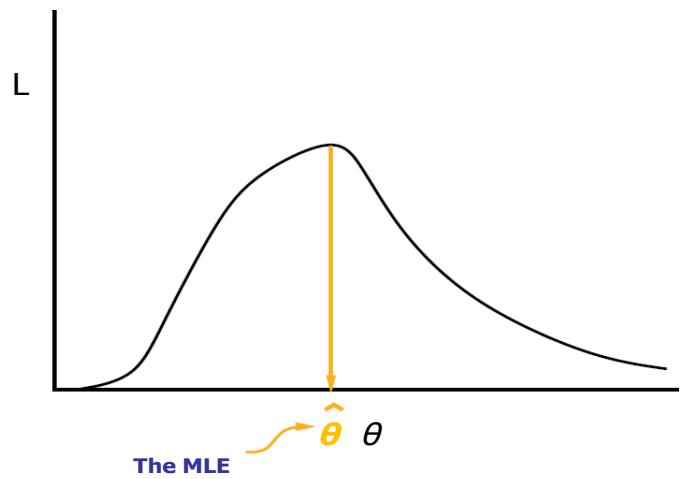
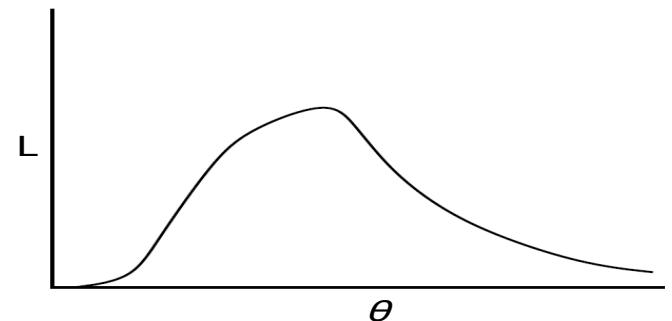
### Likelihood function

The likelihood is  $L(\text{Data} \mid p) = p^5(1-p)^6$



# The likelihood function

- A function of the parameter(s) of our model for the observed data.
- We want to find parameters that result in the maximum of the likelihood function.



- Often the math is easier to deal with if we take the log of the likelihood function
  - $\text{Log}(L)$  achieves its maximum at the same parameter values as  $L$
- Note that “simple” (i.e., convex) likelihood function achieve their maximum at one parameter setting; non-convex likelihood functions have multiple local maxima

# Solving for the solution of the maximum likelihood problem:

- General problem: we want to find the parameter settings that maximize some function given our data.
  - $\text{Log } L = \text{Log} ( p^5 \cdot (1-p)^6 ) = 5x \log(p) + 6 \log(1-p)$
- Differentiate the  $\text{log } L$  function and set derivative to zero.
- We will arrive at  $p = 5/11$

# World view according to Bayesians

- Classic philosophy (frequentist) assumes that parameters are *fixed* quantities that we want to estimate as precisely as possible.
- Bayesian perspective is different: parameters are random variables with probability assigned to particular values of parameters to reflect the degree of evidence for that value.

# Bayesian estimation

- In order to make probability statements about  $\theta$  we make use of Baye's rule:

$$P(\theta | D) = \frac{P(\theta)P(D | \theta)}{P(D)}$$

$$P(\theta | D) \propto P(\theta)P(D | \theta)$$

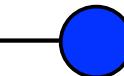
**Posterior**  $\propto$  **Prior**  $\times$  **Likelihood**

- Find  $\theta$ , such as posterior is maximized

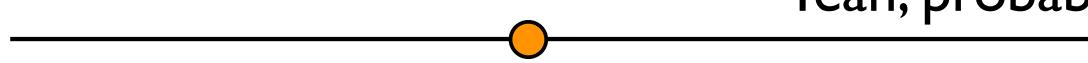
# Back to hypothesis testing ...

$$H_0 : \mu_Y - \mu_Z = 0 ?$$

I seriously doubt it.



Yeah, probably.



# Properties of a test statistic

- When observed value (based on our sample) is “big” or “extreme”, suggests that observed data is very unexpected under the null hypothesis  $H_0$
- We know the distribution of the test statistic under the null model: so we can compute a pvalue quantifying the incompatibility between observed value of test statistic and  $H_0$

- Point estimate: single best guess of the parameter
- Interval estimate (e.g., confidence interval) provides a range of possible values for the parameters.
- Constructing the interval estimator requires knowledge of the estimator's distribution

- Therefore ...
- To complete a hypothesis test, we need a statistic's sampling distribution
- “sampling” -- “hypothetical long repeats of the experiment”
- Standard error: standard deviation of the sampling distribution of an estimator.
- E.g., The standard error of the mean (SEM) (i.e., of using the sample mean as a method of estimating the population mean) is the standard deviation of those sample means over all possible samples (of a given size) drawn from the population

- p-value ...
- The probability under the null  $H_0$  of observing a test statistic *value* as or more extreme than the one computed from the data.
- Two-sided test: both very small and very large values are considered extreme

$$\text{p-value(obs. test stat.)} = P(|\text{test statistic rv}| \geq |\text{obs. test stat}|)$$

- Musing on p-values
- In some sense, it's laziness to work this way: easy because we only need to characterize the distribution of the test statistic under the null
- downside: an indirect measure of how “interesting” the data is
- Just saying something is not “null” or not “boring” is not exactly equivalent to saying what’s truly “exciting” about it.

# Errors in hypothesis testing

- p-values will eventually be thresholded to make decisions

p-value exceeds threshold	... does not
hit	not hit
statistically significant	not statistically significant
discovery!	?
reject $H_0$	accept $H_0$ (wince) fail to reject $H_0$ (roll eyes)

# confusion matrix

“call” based on obs. data true state of nature	“not hit”	reject $H_0$ “hit”	
$H_0$ holds	true negatives	false positives	# nulls
$H_A$ holds “interesting”	false negatives	true positives	# alts
		discoveries	# genes

“call” based on obs. data true state of nature	“not hit”	reject $H_0$ “hit”	
$H_0$ holds	true negatives	false positives Type I errors	# nulls
$H_A$ holds “interesting”	false negatives Type II errors	true positives	# alts
		discoveries	# genes

# Should you care about false positive rate or false negative rate?

- setting of alpha allows us to trade-off between FN rate and FP rate.
- False negative is preferred over false positive:
  - e.g., legal proceeding
- False positive is preferred over false negative:
  - E.g., quarantining people that are suspected to have acquired an infectious disease.