Let's do a "survey". who has an iphone? Let's begin with me.

I do not          $X_1 = 0$ ← code for "No". Yes will be 1.

"raise your    standard        I'm the first survey element.
hand"          symbol          I'm numero uno.
               for datum

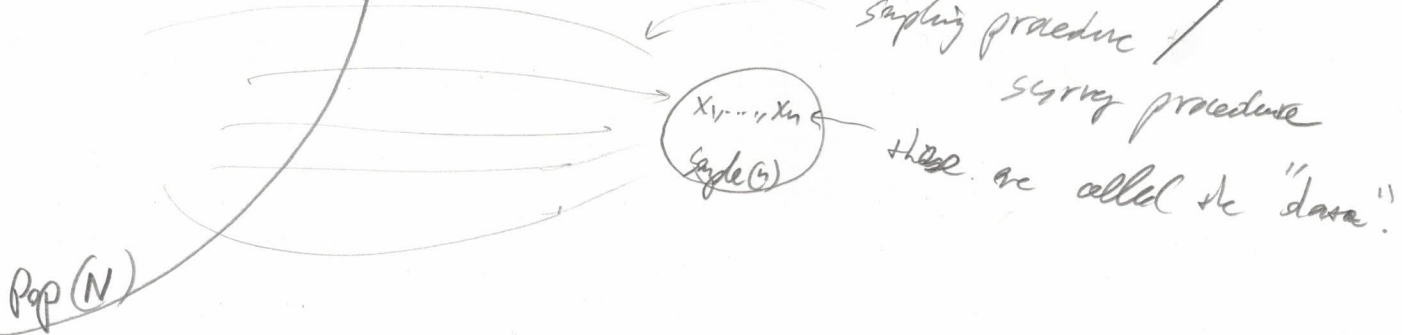$X_1 = 0, X_1 = 1, x_2 = 1, \ldots \ldots \ldots, X_{20} = 0$       (ie. a subset)

Do we believe this survey has a "sample" of $n = 20$ elements from a supposed allof the population? This is the "population model sampling assumption". Let's assume it. Where is the population?

- All people on Earth?
- All people in America?
- All college students?
- _____ in NYC?
- All public college students in NYC?

This is a typical situation. Given a sample, assume population model, then identify the population. This happens in data science all the time.

The more typical situation in classical statistics, is you start with a conception of a population e.g. Pop = All Americans.
Then you take a sample of the population elements, size n, and survey those.

The pop. has $N$ total elements. You should have some idea as to what $N$ is. You defined the population!

$N \approx 333$ million in this case.

$X_1, \ldots, X_n$
Sample (n)

Pop (N)

sampling procedure / survey procedure
these are called the "data"

We see the known in the sample, but not the population. Can we use the sample to tell us something about the population?

Yes, the sample data is used to "infer" properties about the population. Numeric properties are called "population parameters."

"Infer" means to make an educated guess from the particular → the universal. A synonym is "induction."

The opposite is "deduction" which goes from universal → particular.

The process of "inference" is difficult and because is a guess, you can <u>never</u> be sure your inference is correct.

Assume: all swans white $\xrightarrow{deduce}$ thesis: 5 swans are white. | Observe: white 5 swans $\xrightarrow{infer}$ all swans are white. this inference may be wrong.

How is inference done? We generate "Statistics" by running functions on the data.

double-line them

is my → $\hat{\theta} = W(X_1, \ldots, X_n)$ where $\hat{\theta}$ is usually a scalar

own
notation

e.g.

in our iphone survey

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} X_i = .527$$

sample    $\hat{p}$     $\bar{X}$   sample
proportion                average

What can you infer from this scenario? The true, population proportion parameter $\theta$. "Statistical Inference": using scenarios to make inferences.

What is $\theta$?

$$\theta_1 = \frac{x}{N}$$

$x \leftarrow$ # of people with iphones in pop (unknown)

$N \leftarrow$ pop size (somewhat known)
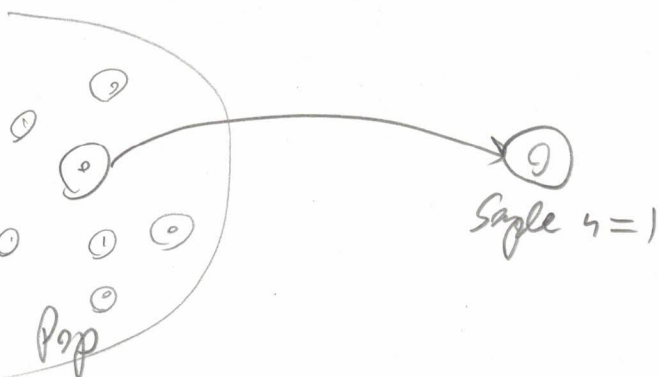
$\theta \in \Theta$ All possible values

$$\Theta = \{0, \tfrac{1}{N}, \tfrac{2}{N}, \ldots \tfrac{N-1}{N}, 1\}$$

param space

Convention: Greek letters are used for unknowable quantities and Roman letters are used for knowable quantities.

$\hat{\theta}$ is a "point estimate" for $\theta$. "Point" meaning one single value that you believe is a "good" guess for the value of $\theta$.

"Point estimation" is one type of statistical inference. The two common other goals are "confidence set creation" (giving an interval of possible values of $\theta$ at a "certainty level" $1-\alpha$) and "theory testing" (testing a theory about the true value $\theta$ at a "certainty level" $\alpha$).

---

Let's discuss sampling and surveying more. Let's take one sample + survey.



Pop    Sample $n=1$

How is this element chosen? Randomly. Technically uniformly sampled i.e. every element of the pop. has prob. $\frac{1}{N}$ of being chosen.

Representative sample: a sample that faithfully reflects the population. Untruly random samples are representative.

What is the probability that $X_1 = 1$?

uses the naive definition of probability

$$P(X_1 = x_1 = 1) = \frac{x}{N} = \theta$$

$$\frac{\text{\# elements satisfying even}}{\text{\# of total elements}}$$

the r.v. modeling the first survey datum (capital letter)

its realization (lowercase letter)

a possible value of its realization

data: realizations of r.v.'s. Surveying: forcing a r.v. to realize.

Bernoulli r.v. with parameter $\theta$

$$\Rightarrow X_1 \sim \text{Bern}\left(\frac{x}{N}\right) = \text{Bern}(\theta) = \theta^{x_1}(1-\theta)^{1-x_1} \ldots$$

prob. mass function (PMF).

$$\text{Supp}[X_1] = \{0, 1\}$$ the support of the r.v. is the set of all possible realization values.

Note: the parameter of the survey r.v. is the same as the population parameter we would like to draw inference about.

---

Let's draw a second sample assuming $X_1 = 1$. And ask some questions
How? Each remaining element has $\frac{1}{N-1}$ prob. of being drawn.

sample $n=2$



Pop

$$P(X_2 = 1 \mid X_1 = 1) = \frac{x-1}{N-1} < \theta = P(X_1 = 1)$$

$$X_2 \mid X_1 = 1 \sim \text{Bern}\left(\frac{x-1}{N-1}\right)$$

Conditional prob.

cond. r.v. model

$$\Rightarrow X_1, X_2 \text{ are dependent r.v.'s}$$

If $X_1 = 0 \Rightarrow P(X_2=1 \mid X_1=0) = \frac{x}{N-1} > \theta$

either way dependent

what is $P(X_2=1) = \frac{x}{N} = \theta \Rightarrow X_2 \sim bern(\theta)$

↑ conditional prob. $X_1$ was realized but... you don't know who it is thus you pretend it doesn't exist

$\Rightarrow X_1 \overset{d}{=} X_2$ Hence "identically distributed" since they have the same PMF.

Let's sample all $n$. Let $T_n = X_1 + \dots + X_n$ i.e. the r.v. that tallies the total # of 1's.

$$P(T_n=t) = \frac{\binom{x}{t}\binom{N-x}{n-t}}{\binom{N}{n}} = Hyper(n, x, N)$$

sample size, total # of 1's in pop, total pop.

# of unique samples!

Hypergeometric r.v. model

How did it get this complicated? Because $\frac{x}{N} \neq \frac{x-1}{N-1}$ !

Let's make a simplifying assumption. Let $x, N \to \infty$ with $\frac{x}{N} = \theta$.

$$\Rightarrow \frac{x}{N} \approx \frac{x-1}{N-1} = \frac{x-k_1}{N-k_2} = \theta \text{ for all } k_1, k_2.$$

with this limit, we assume $\Theta = [0,1]$ all possible real #'s

Now $P(X_2=1 \mid X_1=1) = \frac{x-1}{N-1} \overset{assume}{=} \theta = P(X_2=1)$

$\Rightarrow X_1, X_2$ are independent $\Rightarrow X_1, X_2 \overset{iid}{\sim} bern(\theta)$

For all $n$ samples, $X_1, X_2, \dots, X_n \overset{iid}{\sim} bern(\theta)$

$$\Rightarrow T_n = X_1 + \dots + X_n \sim Binomial(n, \theta) = \binom{n}{t}\theta^t(1-\theta)^{n-t}$$

Markdown:

Let's consider a new sampling problem. At the iphone factory, they check every new iphone to make sure it works. Let's say they check the first one, $X_1 = 1$ ← it works, the second $X_2 = 0$ ← it doesn't work, ..., $X_{100} = 1$. What population is this sample from? All iphones? $N = ?$

Are you drawing one sample of $n$ from $\binom{N}{n}$? Not really. What is $\Theta$? Is it a "population" parameter?

Would you agree $X_1, ..., X_n \overset{iid}{\sim} Bern(\Theta)$?

Is it a "process" parameter? Process or infinite population...

we still have a r.v. model that describes the survey. At this point we no longer care whether the pop. is real or if its a process, we just need an iid r.v. model assumption called the "data generating process."

Let's return to our main goal: inference. Specifically: parameter estimation of a parameter $\theta$.

$$\hat{\theta} = \frac{1}{n}(x_1 + \dots + x_n) \approx \theta. \quad \text{How approximate is it?}$$

Since $x_1, \dots, x_n$ here random realization of $X_1, \dots, X_n$, $\hat{\theta}$ could have been different. e.g if $\vec{x} = [1 \ 0 \ 0 \ 1 \ 0] \Rightarrow \hat{\theta} = 0.4$

but if $\vec{x} = [1 1 1 0 1] \Rightarrow \hat{\theta} = 0.8$. Thus $\hat{\theta}$ is a realization itself from a r.v. $\hat{\theta}_n = \frac{1}{n}(X_1 + \dots + X_n)$ called a "statistical estimator" or just "estimator". So... $\hat{\theta}$ is a realization from $\hat{\theta}_n$.

(the distr. of $\hat{\theta}_n$ is called the "sampling distr.")

The properties of the estimator are very important because they tell us a lot about our estimate.

One property is the expectation,

$$E[\hat{\theta}_n] = E\left[\frac{1}{n}T_n\right] = \frac{1}{n}E[T_n] = \frac{1}{n}n\theta = \theta \quad \forall n$$

$\uparrow$ $X_1, \dots, X_n$   $\uparrow$ over all $t$

$\Rightarrow E[\hat{\theta}_n] = \theta$ this is special.

It means $\hat{\theta}$ is "unbiased"

this expectation is taken over the $Supp[X_1], \dots, Supp[X_n] := \mathcal{X}$ weighted by the joint mass function. Now again our distr. for $\theta$: $\theta$ is one #!

$p(x_1, \dots, x_n)$. (HW)

In general, for any estimator and any pop. param.

$$Bias[\hat{\theta}_n] := E[\hat{\theta}_n] - \theta$$

If $Bias[\hat{\theta}_n] = 0$, $\hat{\theta}$ is "unbiased".
If $Bias(\hat{\theta}_n) \neq 0$, $\hat{\theta}$ is "biased".

Across every possible sample of any size $n$, the average estimate will be the pop. parameter $\theta$.

This is certainly reasonable. How "far" is $\hat{\theta}$ away from $\theta$?

Let's define "far" by a loss function $l(\hat{\theta}, \theta)$ where

$$l: \Theta \times \Theta \longrightarrow [0, \infty) \text{ and } l(\hat{\theta}, \theta) = 0 \text{ only when } \hat{\theta} = \theta.$$

Some examples...

**default**

$$l(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2 \quad \text{"squared error loss" or "} L_2 \text{ loss"}$$

$$l(\hat{\theta}, \theta) = |\hat{\theta} - \theta| \quad \text{"absolute error loss" or "} L_1 \text{ loss"}$$

$$l(\hat{\theta}, \theta) = |\hat{\theta} - \theta|^p \quad \text{"} L_p \text{ loss" for } p > 0.$$

$$l(\hat{\theta}, \theta) = \int_{\mathcal{X}} \ln\left(\frac{f(x;\theta)}{f(x;\hat{\theta})}\right) f(x;\theta) \, d\vec{x} \quad \text{Kullback-Leibler loss}$$

for cont. r.v.'s $X_1, \ldots, X_n$

estimator $\Rightarrow$ loss has a distribution

$$R(\hat{\theta}, \theta) := \underset{\mathcal{X}}{E}\left[l(\hat{\theta}, \theta)\right]$$

Risk of an estimator; what is the average loss accd. to our loss function. If it's sqd error loss...

mean squared error

$$R(\hat{\theta}, \theta) = E_{\mathcal{X}}\left[(\hat{\theta} - \theta)^2\right] = MSE(\hat{\theta}) = E_{\mathcal{X}}\left[(\hat{\theta} - \underset{\mathcal{X}}{E}(\hat{\theta}))^2\right] = \underset{\mathcal{X}}{Var}(\hat{\theta}).$$

If $\hat{\theta}$ is unbiased $\Rightarrow E(\hat{\theta}) = \theta$

$$Risk = MSE = Variance \text{ for an unbiased estimator under } L_2 \text{ loss.}$$

under $L_2$ loss

For a biased estimator,

$$MSE = E_X\left[(\hat{\theta}-\theta)^2\right] = E_X\left[\hat{\theta}^2 - 2\theta\hat{\theta} + \theta^2\right]$$

$$= E_X[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2$$
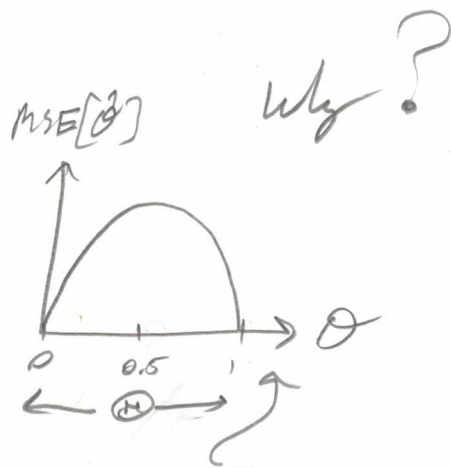
Recall: $Var(\hat{\theta}) = E[\hat{\theta}^2] - E(\hat{\theta})^2$

$$= Var[\hat{\theta}] + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2$$

$$= Var[\hat{\theta}] + \left(E[\hat{\theta}] - \theta\right)^2$$

$$= Var[\hat{\theta}] + bias[\hat{\theta}]^2 \qquad \text{"Bias-Variance decomposition of MSE."}$$


$MSE[\hat{\theta}]$    Why?

$$\boxed{SE(\hat{\theta}) := \sqrt{Var(\hat{\theta})}}$$

estimator
Standard error is the standard deviation of an estimator.

$E[\hat{\theta}_n] = \theta$ for all iid DGP's $X_1,...X_n$ with mean $\theta$, variance $\sigma^2$

Back to our example.... $\hat{\theta} = \bar{X}$

variance law for iid r.v.'s        for all iid DGP's

$$E(\hat{\theta}) = \sqrt{Var[\hat{\theta}_n]} = \sqrt{Var\left[\tfrac{1}{n}T_n\right]} = \sqrt{\tfrac{1}{n^2}Var[T_n]} = \sqrt{\tfrac{1}{n^2}\,n\,Var(X)} = \frac{SD(X)}{\sqrt{n}} = \sqrt{\frac{\theta(1-\theta)}{n}}$$

$\sigma$
Bernoulli

Law of Large #'s  for $\hat{\theta}_i = \bar{X}$,

$\Rightarrow n \to \infty \Rightarrow SE(\hat{\theta}_n) \longrightarrow 0 \Rightarrow \hat{\theta}_n \to \theta$ wp1

this limit isn't precisely defined.

r.v.   #

But ignore this detail....

For all estimators which are sample averages of iid r.v.'s, we're kind of done with Goal #1: Pt. Est.

Now that we have an idea about how variable the

standard inference

estimator is we can move to $\land$ goal #3: testing. Turns out, which $\land$ is

goal #2 is harder!

confidence sets