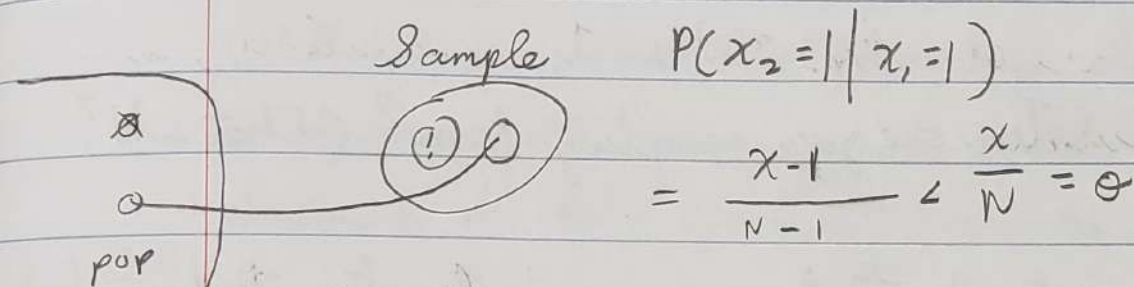


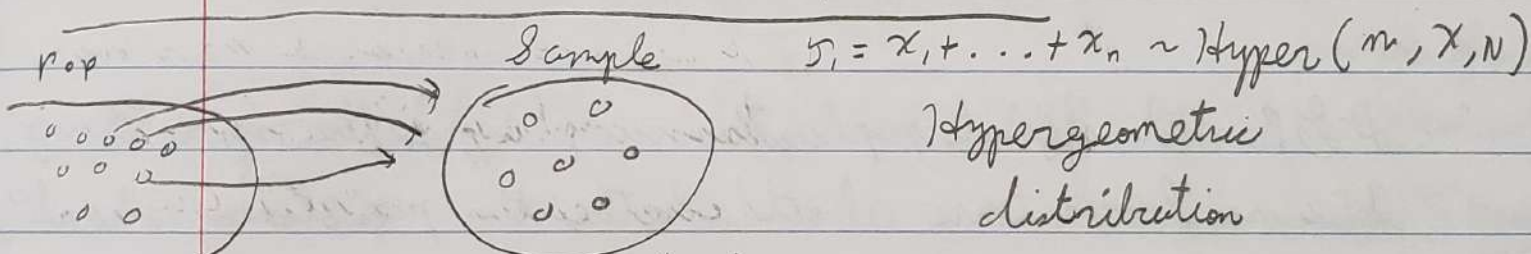
8/31/20

$$X_1 \sim \text{Bern}(\theta) = \text{Bern}\left(\frac{x}{N}\right)$$

Let's draw a second sample from the population assuming $X_1 = 1$



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



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

Dealing with the hypergeometric is complicated (but doable) What can we assume to make this go away?

Simplifying assumption

Let $x, N \rightarrow \infty$ but $\theta = \frac{x}{N}$

$$\lim P(X_2 = 1 | X_1 = 1) = \lim \frac{x-1}{N-1} = \theta \Rightarrow X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bern}(\theta)$$

Pretend you work at the iPhone factory, they sample new iPhones to ensure they work to ensure the manufacturing is working properly. You check the first one $x_1=1, x_2=1, \dots, x_{100}=1$.

What population are you sampling from? What is N ?

When you estimate θ , you're estimating θ in a "process", is a "data generating process" (DGP), iid $\text{Bern}(\theta)$.

DGPs and infinite population sampling is the same thing. We no longer care about whether the population is "real", we just assume an iid DGP from now on...

Returning to our main goal: inference i.e. knowing something about θ from the data. First subgoal: point estimation. Recall,

$$\hat{\theta} = \frac{1}{n} (x_1 + \dots + x_n) \quad x_1, \dots, x_n \text{ are random realizations from } x_1, \dots, x_n \stackrel{\text{iid}}{\sim} \text{Bern}(\theta).$$

$$\text{eg } \bar{x} = [10010] \Rightarrow \hat{\theta} = 0.4 \Rightarrow \hat{\theta} \text{ random}$$

$$\text{eg } \bar{x} = [11101] \Rightarrow \hat{\theta} = 0.8$$

$\hat{\theta}$ is a realization from the r.v. $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$ called "statistical estimator" or just "estimator"

The statistic (statistical estimate, estimate) is a realization from the estimator. The distribution of the estimator, that is called the "sampling distribution". The sampling distribution and its properties are very important because it tells us a lot about an estimator

think about
D, avg est.
or?

One property is the estimator's expectation, the mean over all

$$E[\hat{\theta}] = E\left[\frac{1}{n}(x_1 + \dots + x_n)\right] = \frac{1}{n} \sum E[x_i] = \frac{1}{n} n E[x_1] = \theta \Rightarrow \hat{\theta} \text{ is unbiased}$$

\uparrow over all x_1, \dots, x_n \uparrow in our iid Bern(θ) setting

Bias $[\hat{\theta}] := E[\hat{\theta}] - \theta$. If Bias $[\hat{\theta}] = 0 \Rightarrow \hat{\theta}$ is unbiased
Bias $[\hat{\theta}] \neq 0 \Rightarrow \hat{\theta}$ is biased

How far is that from theta? ("error function")

We define a distance function AKA "loss function"

$l(\hat{\theta}, \theta)$, $l: \mathbb{H} \times \mathbb{H} \rightarrow [0, \infty)$. $l = 0$ only if $\hat{\theta} = \theta$

There are many loss function e.g.

$l(\hat{\theta}, \theta) := |\hat{\theta} - \theta|$ absolute error loss (L_1 , loss) \rightarrow

default

* $l(\hat{\theta}, \theta) := |\hat{\theta} - \theta|^2$ square error loss (L_2 loss)

$l(\hat{\theta}, \theta) := |\hat{\theta} - \theta|^p, p \geq 0$ L_p loss

$l(\hat{\theta}, \theta) := \int \ln \left(\frac{f(x; \theta)}{f(x; \hat{\theta})} \right) f(x; \theta) dx$ Kullback-Leibler loss
for continuous r.v.s
 $\frac{1}{x} \in \mathcal{X}$

How far away on average are we?

Risk of an estimator
 $R(\hat{\theta}, \theta) := E[l(\theta, \hat{\theta})]$
 \uparrow
over x_1, \dots, x_n

If we use squared error loss,
 $R(\hat{\theta}, \theta) = \text{MSE}[\hat{\theta}] = E[(\hat{\theta} - \theta)^2]$
"mean squared error" (MSE)

under squared error loss and DGP: iid Bern(θ)

$R(\hat{\theta}, \theta) =$

If the estimator is unbiased, does its MSE simplify?

$$MSE[\hat{\theta}] = E[(\hat{\theta} - \theta)^2] \stackrel{\text{if } \hat{\theta} \text{ is unbiased, } E[\hat{\theta}] = \theta}{=} E[(\hat{\theta} - E[\hat{\theta}])^2] = \text{Var}[\hat{\theta}]$$

If $\hat{\theta}$ is unbiased, $E[\hat{\theta}] = \theta$

MSE = Variance

For a biased estimator (ie the general case),

$$\begin{aligned} MSE[\hat{\theta}] &= E[(\hat{\theta} - \theta)^2] = E[\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2] \\ &= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \quad \text{Recall } \text{var}[\hat{\theta}] = E[\hat{\theta}^2] - E[\hat{\theta}]^2 \\ &= \text{Var}[\hat{\theta}] + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2 \\ &= \text{Var}[\hat{\theta}] + (E[\hat{\theta}] - \theta)^2 \\ &= \text{Var}[\hat{\theta}] + \text{Bias}[\theta]^2 \quad \text{Bias-variance decomposition of MSE} \end{aligned}$$

$SE[\hat{\theta}] := \sqrt{\text{var}[\hat{\theta}]}$ "standard error of the estimation"