Benjamin Nguyen
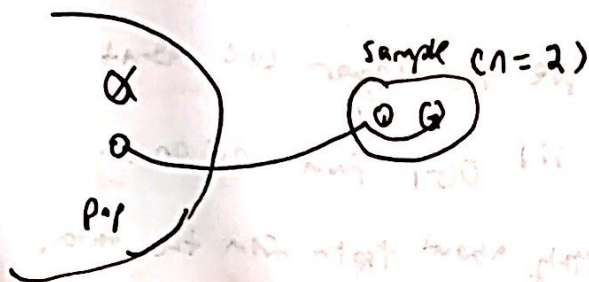MHTA 369
8/31/20
Lecture #2
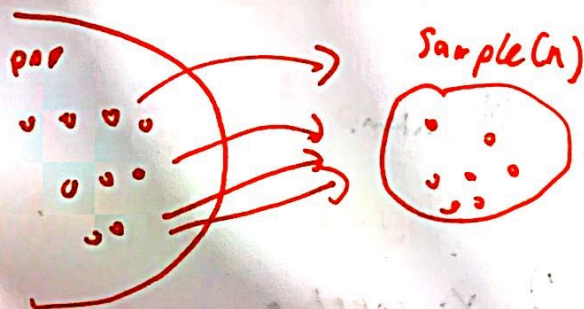
$X_1 \sim Bern(\theta) = Bern\left(\frac{K}{N}\right)$

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



sample $(n=2)$

$P(X_2 = 1 \mid X_1 = 1)$

$= \frac{K-1}{N-1} < \frac{K}{N} = \theta$

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



pop          Sample $(n)$

$T_1 = X_1 + \dots + X_n \sim Hyperge\!\cdot\!(n, K, N)$

$P(T_n = t) = \dfrac{\binom{K}{t}\binom{N-K}{n-t}}{\binom{N}{n}}$

Dealing with the hypergeo is complicated. What can we assume to get iid or simplifying assumption?

this problem?

Let $K, N \to \infty$ but $\theta = \frac{K}{N}$  $\Rightarrow$  $X_1 \dots X_n \overset{iid}{\sim} Bern\,\theta$

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

Pretend you work at the Iphone factory, they sample new Iphones to ensure they win to ensure the manufacturing is working properly. You check the first one $x_1 = 1$, $x_2 = 2$.

What population are you sampling from? What is N?

When you estimate theta, you're estimating theta in a "process" i.e a "data generating process" (DGP), iid Bern (theta)...

DGPS are infinite population Sampling it 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 theta from the data.

First subgoal: point estimator. Recall,

$$\hat{\theta} = \frac{1}{n}(X_1 + \dots + X_n).$$
$x_1, \dots x_n$ are random realizations from $X_1, \dots, X_n \overset{iid}{\sim}$ Bern(θ)

$\bar{x}$, $\hat{p}$    e.g $\vec{x} = [10010] \Rightarrow \hat{\theta} = 0.4$    $\Rightarrow \hat{\theta}$ random.

$$\vec{x} = [11101] \Rightarrow \hat{\theta} = 0.8$$

$\hat{\theta}$ is a realization from the r.v $\hat{\theta} := \frac{1}{n} \sum_{i=1}^{n} X_i$ called $\hat{q}$

"Statistical estimator" or just "estimator". The statistic (statistical estimate, estimate) is a realization from the estimator. The distribution of the estimator, $\hat{\theta}$ is called the "sampling distribution". This sampling distribution and its properties one very important because it tells us a lot about our estimates.

One property is the estimator's expectation the mean often over all samples of size n.

$$E[\hat{\theta}] = \theta$$
why would it be nice if
↗
overall
$x_1, \dots, x_n$

$$E[\hat{\theta}] = E\left[\frac{1}{n}(x_1 + \dots + \frac{1}{n})\right] = \frac{1}{n}\sum E[x_i] = \frac{1}{n}\cdot n \cdot E[x_i] = \theta$$

In our iid $Bern(\theta)$ sett.

$\Rightarrow \hat{\theta}$ is unbiased.

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

---

How far is $\hat{\theta}$ from $\theta$?

We define a distance function AKA "loss function". ("error function")

$$\ell(\hat{\theta}, \theta), \quad \ell: \textcircled{H} \times \textcircled{H} \to [0, \infty). \quad \ell = 0 \text{ only if } \hat{\theta} = \theta$$

There are many ways to define a loss function e.g.

$$\ell(\hat{\theta}, \theta) := |\hat{\theta} - \theta| \quad \text{absolute error loss } (L_1, Loss)$$

$$\star \quad \ell(\hat{\theta}, \theta) := |\hat{\theta} - \theta|^2 \quad \text{squared error loss } (L_2, loss)$$

$$\ell(\hat{\theta}, \theta) := |\hat{\theta} - \theta|^p, \quad p > 0 \quad L_p \text{ loss}$$

$$\ell(\hat{\theta}, \theta) := \int_{x \in \mathcal{X}} \ln\left(\frac{f(x;\theta)}{f(x;\hat{\theta})}\right) f(x;\theta) \, d\lambda \quad \text{Kullback-Leibler (KL) loss for continuous r.v.'s}$$

---

How far away on average are we?

Risk: $R(\theta, \hat{\theta}) = E[\ell(\theta, \hat{\theta})]$

Risk of an estimator.

If we use squared error loss,

$$R(\hat{\theta}, \theta) = MSE[\hat{\theta}] = E[(\hat{\theta} - \theta)^2]$$

"mean squared error"

Scanned with CamScanner

If the estimator is unbiased, does its MSE simplify?

$$\text{MSE}[\hat{\theta}] = E\left[(\hat{\theta} - \theta)^2\right] = E\left[(\hat{\theta} - E[\hat{\theta}])^2\right] = \text{variance or } \text{Var}[\hat{\theta}].$$

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

---

For a biased estimator (ie the general case),

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

$$= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \qquad \text{Recall}$$
$$\text{var}[\hat{\theta}] = E[\hat{\theta}^2] - E[\hat{\theta}]^2.$$

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

$$= \text{var}[\hat{\theta}] + (E[\hat{\theta}] - \theta)^2$$

$$= \text{var}[\hat{\theta}] + \text{Bias}[\hat{\theta}]^2 \qquad \text{Bias-variance decomp of MSE.}$$

---

$$\text{SE}[\hat{\theta}] := \sqrt{\text{var}[\hat{\theta}]} \qquad \text{"Standard error of the estimator"}$$