

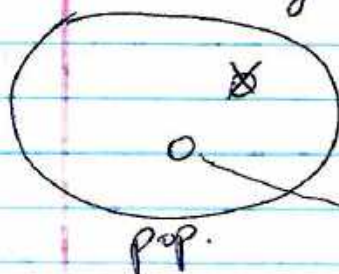
MA369

08/31/2020.

Lecture 02

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

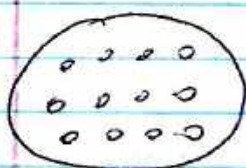
Let's draw a second sample from the population assuming $X_1 = 1$.

Sample ($n=2$)

Not Independent!

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

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



$$T_n = X_1 + \dots + X_n$$

$$\sim \text{Hyper}(n, x, N)$$

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

Hypergeometric
Distribution.

Dealing with the hypergeometric is complicated (but doable).

What can we assume to make this go away?

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

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

→ $\theta = \frac{x}{N}$ Simplifying assumption → $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p)$

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 $X_1 = 1$, $X_2 = 1$, ..., $X_{100} = 1$

What population are you sampling from? What is N ?

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

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,

$$\hat{\theta} = \frac{1}{n} (X_1 + \dots + X_n)$$

$\hat{\theta}$
 \bar{x} \hat{p}

X_1, \dots, X_n are random realizations from

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bern}(\theta)$

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

but e.g. $\vec{x} = [11101] \Rightarrow \hat{\theta} = 0.8$

$\Rightarrow \hat{\theta}$ random.

$\hat{\theta}$ is a realization from the r.v. $\hat{\theta} := \frac{1}{n} \sum_{i=1}^n X_i$ = $w(X_1, \dots, X_n)$
 called a "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 are very
 important b/c ~~they~~ it tells us a lot about our
 estimates.

One property is the estimator's expectation, the mean.
It would be nice if $E[\hat{\theta}] = \theta$. over all sample of
 size n .

$$\underset{\substack{\uparrow \\ \text{over all}}} E[\hat{\theta}] = E\left[\frac{1}{n}(X_1 + \dots + X_n)\right]$$

$$X_1, \dots, X_n = \frac{1}{n} \sum E[X_i] = \frac{1}{n} \cdot n E[X_i] = \theta.$$

in θ iid Bern(θ) setting

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

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

If $\text{Bias}[\hat{\theta}] \neq 0 \Rightarrow \hat{\theta}$ is biased.

How far is $\hat{\theta}$ from θ ?

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

$$\ell(\hat{\theta}, \theta). \quad \ell: \mathcal{H} \times \mathcal{H} \rightarrow [0, \infty)$$

$\ell = 0$ only if $\hat{\theta} = \theta$.

There are many loss functions e.g.

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

Default *

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

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

$$\ell(\hat{\theta}, \theta) := \int \ln\left(\frac{f(x; \theta)}{f(x; \hat{\theta})}\right) f(x; \theta) d\mathbb{P}$$

$\mathbb{P} \in \mathcal{P}$.

Kullback-Leibler (KL) loss
for continuous RV's.

How far away on average are we?

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

↑
over X_1, \dots, X_n .

Risk of an estimator.

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

~~Under squared error loss~~

If the estimator is unbiased, does its MSE simplify?

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

if θ is unbiased, $E[\hat{\theta}] = \theta$
 $= \text{Var}[\hat{\theta}]$.

For a biased estimator (i.e. the general case)

$$\begin{aligned} \text{MSE}[\hat{\theta}] &= E[(\hat{\theta} - \theta)^2] = E[\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2] \\ &= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \end{aligned}$$

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

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

Bias - variance decomposition of MSE.

$$\text{SE}[\hat{\theta}] := \sqrt{\text{Var}[\hat{\theta}]} \quad \text{"standard error of the estimator"}$$

Note

r.v. $X \sim P(X)$.

$$\mu = E[X] = \sum_{X \in \text{Supp}[X]} X P(X).$$

$$\sigma^2 = \text{Var}[X] = E[(X - \mu)^2].$$

the square of the distance.

that how far an X away from μ .

If $\hat{\theta} = \bar{X}$,

$$\theta = E[\theta] = \sum_{\hat{\theta} \in \text{Supp}[\hat{\theta}]} \hat{\theta} P(\hat{\theta}).$$

$\hat{\theta}$ is unbiased.

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

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

$$\ell(\hat{\theta}, \theta)$$

$$\hat{\theta} = \bar{X} = \frac{X_1 + \dots + X_n}{n}$$

$\hat{\theta} = \bar{x}$ realization from \bar{X}