

Victoria Lombardi

①

Math 369

8/31/20

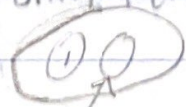
Lecture 2

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

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

sample(n=2)

pop



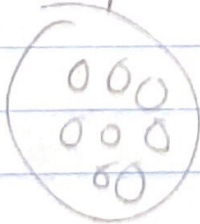
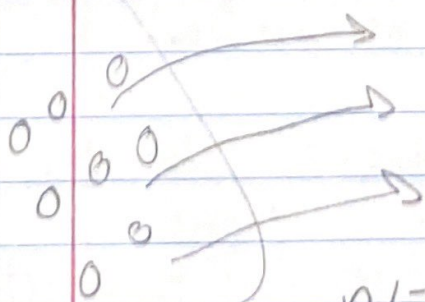
$$P(X_2 = 1 | X_1 = 1)$$

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

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

sample(n)

$$T_n = X_1 + \dots + X_n \rightarrow \text{Hyper}(n, X, N)$$



Hypergeometric distribution

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

(2)

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

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

Simplifying assumption

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

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

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", i.e. 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.

(3)

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

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

$\hat{\theta}$
 \hat{p}

e.g. $\vec{x} [1 0 0 1 0] \Rightarrow \hat{\theta} = 0.4$

but e.g. $\vec{x} [1 1 1 0 1] \Rightarrow \hat{\theta} = 0.8 \Rightarrow \hat{\theta} \text{ random}$

$\hat{\theta}$ is a realization from the random variable
 $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$ called a "statistical estimator"
or just "estimator". The statistic (estimate, statistical 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 because it tells us a lot about our estimates

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

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

over all x_1, \dots, x_n

$= \frac{1}{n} E[X_i] = \theta \Rightarrow \hat{\theta} \text{ is unbiased}$

\leftarrow in our iid Bern(θ) setting

expectation
↓

(4)

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

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

How far is $\hat{\theta}$ from θ ?

We define a distance function AKA "loss function", ("error function"), parameter space

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

$l=0$ only if $\hat{\theta} = \theta$

There are many loss function e.g.

① $l(\hat{\theta}, \theta) := |\hat{\theta} - \theta|$ absolute error loss (L_1 loss)

default

*

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

③ $l(\hat{\theta}, \theta) := |\hat{\theta} - \theta|^p, p > 0$ L_p loss

④ $l(\hat{\theta}, \theta) := \int \ln\left(\frac{f(x; \theta)}{f(x; \hat{\theta})}\right) f(x; \theta) d\tilde{x}$ *

Kullback-Leibler (KL) loss for continuous rv's

How far away on average are we?

ave.

x_1, \dots, x_n

$$\rightarrow E[l(\theta, \hat{\theta})]$$

(5)

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

Risk of an estimator

If we use squared error loss,

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

"mean squared error" (MSE)

If the estimator is unbiased, does its MSE simplify?

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

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

MSE = Variance

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

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

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

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

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

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

$$= \text{Var}[\hat{\theta}] + \text{Bias}[\hat{\theta}]^2$$

Bias-Variance
decomposition of MSE

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