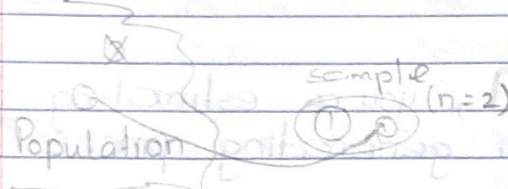


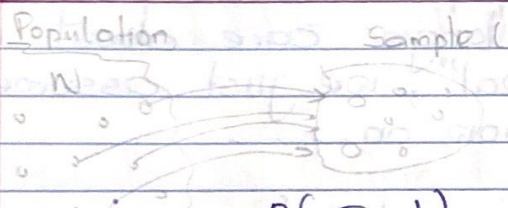
$$X_1 \sim \text{Bern}(\theta) = \text{Bern}\left(\frac{X}{N}\right)$$

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



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

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



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

Hypergeometric distribution

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

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

Let  $X, N \rightarrow \infty$  but  $\theta = \frac{X}{N}$  (make the ratio constant)  $\rightarrow$  Simplifying Assumption

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

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



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  $\theta$ , you're estimating  $\theta$  in a "process", i.e. a "data generating process" (DGP), i.i.d.  $\text{Bern}(\theta)$ .

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

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

$$\hat{\theta} = \bar{x} = \frac{1}{n} (x_1 + \dots + x_n). \quad x_1, \dots, x_n \text{ are random realizations from } X_1, \dots, X_n \text{ i.i.d. } \text{Bern}(\theta).$$

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

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

$$\hat{\theta} \text{ is a realization from the rv } \hat{\theta}_n := \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n \epsilon_i$$

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 because it tells us a lot about our estimates.

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

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

↑  
overall  $x_1, \dots, x_n$

in our iid Bern( $\theta$ ) setting.  
 $= \theta \Rightarrow \hat{\theta}_n$  is unbiased"

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

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

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

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

There are many loss functions. e.g

$$l(\hat{\theta}, \theta) := |\hat{\theta} - \theta| \quad \text{absolute error loss (L}_1 \text{ loss)}$$

$$* \quad l(\hat{\theta}, \theta) := |\hat{\theta} - \theta|^2 \quad \text{squared error loss (L}_2 \text{ loss)}$$

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

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



How far away on average are we?

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

Risk of an estimator 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)"

If the estimator is unbiased, does its MSE simplify?  $\text{MSE} = \text{Variance}$

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

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

For a biased estimator (ie the general case),

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

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

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

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

$$= \text{Var}[\hat{\theta}] + \text{Bias}[\hat{\theta}]^2 \quad \text{Bias-variance decomposition of MSE}$$

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