

# Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#231 Statistical Inference

été 2020

# Average

Consider observations  $\{x_1, \dots, x_n\}$  with true mean  $\mu$ , and standard deviation  $\sigma$ .

$$\mathbb{E}(\bar{X}_n) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \stackrel{*}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} n\mu = \mu$$

\* since the expected value is linear

$\bar{x}$  is an unbiased estimator of the (true) mean  $\mu$

$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \stackrel{*}{=} \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}$$

\* because variables are independent, and variance is a quadratic function

The variance of  $\bar{x}$  is  $\frac{\sigma^2}{n}$ , and  $\text{Var}(\bar{x}) \rightarrow 0$  as  $n \rightarrow \infty$

# Average

Problem:  $\sigma$  is unknown...

$$s_n^2 = \frac{1}{n-1} \sum_{i=1}^n [X_i - \bar{X}_n]^2.$$

$$\mathbb{E}(S_n^2) = \mathbb{E} \left( \frac{1}{n-1} \sum_{i=1}^n [X_i - \bar{X}_n]^2 \right) \stackrel{*}{=} \mathbb{E} \left( \frac{1}{n-1} \left[ \sum_{i=1}^n X_i^2 - n\bar{X}_n^2 \right] \right)$$

\* from the same property as before

$$\mathbb{E}(S_n^2) = \frac{1}{n-1} [n\mathbb{E}(X_i^2) - n\mathbb{E}(\bar{X}^2)] \stackrel{*}{=} \frac{1}{n-1} \left[ n(\sigma^2 + \mu^2) - n \left( \frac{\sigma^2}{n} + \mu^2 \right) \right]$$

\* since  $\text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$

An unbiased estimator of the variance of  $\bar{x}$  is  $\frac{s^2}{n}$

# Average

Finally, from the central limit theorem (#431),

$$\bar{X}_n \approx \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

or

$$\bar{X}_n \approx \mathcal{N}\left(\mu, \frac{s^2}{n}\right)$$

so, a 95% confidence interval for  $\mu$  is

$$\left[ \bar{x} \pm 2 \frac{s}{\sqrt{n}} \right]$$

# Average

In a given city, a random sample of  $n = 400$  people is taken. The average years of schooling of this sample is  $\bar{x} = 11.6$  years, with a standard deviation of  $\hat{s} = 4.1$ . Find a 95% confidence interval for the average educational level of people in this city.

Let  $\mu$  be the true mean educational level,  $\sigma$  the standard deviation.

From the central limit theorem (#431),  $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

95% confidence interval for  $\mu$  is  $\left[\bar{x} \pm 2 \frac{\sigma}{\sqrt{n}}\right]$

Standard deviation of  $\bar{X}$  can be estimated by  $\hat{s}/\sqrt{n} \simeq 0.2$ ,

Therefore, 95% confidence interval is  $[11.6 \pm 0.4]$

# Proportion

In a given city, a random sample of  $n = 400$  people is taken. The number of people who intend to vote for a candidate is 212. Find a 95% confidence interval for the probability to vote for that candidate, in this city.

Let  $p$  be the true probability, and  $\hat{p}$  denote the (empirical) frequency, i.e.  $\hat{p} = \frac{212}{400}$

From the central limit theorem ([#431](#)),  $\hat{p} \sim \mathcal{N}\left(p, \frac{p(1-p)}{n}\right)$

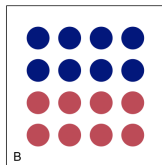
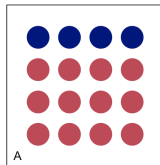
95% confidence interval for  $p$  is  $\left[\hat{p} \pm 2 \frac{p(1-p)}{\sqrt{n}}\right]$

Standard deviation of  $\bar{X}$  can be estimated by  $\hat{s}/\sqrt{n} \simeq 0.2$ ,  
Therefore, 95% confidence interval is  $[11.6 \pm 0.4]$

# Bayesian Approach

Consider two bags, A (4 blue balls, 12 red) and B (8 blue balls, 8 red).

We pick randomly a bag, and randomly a ball. The ball is red. What probability that bag A was selected ?



The prior estimate is  $\mathbb{P}(A) = \mathbb{P}(B)$

Given that the ball is red, the posterior estimate is

$$\mathbb{P}(A|\bullet) = \frac{\mathbb{P}(A)\mathbb{P}(\bullet|A)}{\mathbb{P}(\bullet)} = \frac{\mathbb{P}(A)\mathbb{P}(\bullet|A)}{\mathbb{P}(\bullet|A)\mathbb{P}(A) + \mathbb{P}(\bullet|B)\mathbb{P}(B)}$$

$$\text{so, numerically, } \mathbb{P}(A|\bullet) = \frac{\frac{1}{2} \cdot \frac{12}{16}}{\frac{12}{16} \cdot \frac{1}{2} + \frac{8}{16} \cdot \frac{1}{2}} = \frac{3}{5} = 60\%$$

$$\text{et } \mathbb{P}(B|\bullet) = 40\%.$$

# Bayesian Approach

Application to statistical inference:

Assume that  $x$  has a density  $f(x|\theta)$ .

Probability to observe sample  $\mathbf{x} = \{x_1, \dots, x_n\}$  (independent) is

$$f(\mathbf{x}|\theta) = f(x_1|\theta) \cdot f(x_2|\theta) \cdots f(x_n|\theta) \quad (\text{called likelihood})$$

Assume a prior distribution for  $\theta$ , density  $\pi$ .

Its posterior distribution is

$$\pi(\theta|\mathbf{x}) = \frac{\pi(\theta) \cdot f(\mathbf{x}|\theta)}{f(\mathbf{x})} = \frac{\pi(\theta) \cdot f(\mathbf{x}|\theta)}{\int f(\mathbf{x}|\tau)\pi(\tau)d\tau}$$

which is the general form of

$$\mathbb{P}(A|\bullet) = \frac{\mathbb{P}(A)\mathbb{P}(\bullet|A)}{\mathbb{P}(\bullet)} = \frac{\mathbb{P}(A)\mathbb{P}(\bullet|A)}{\mathbb{P}(\bullet|A)\mathbb{P}(A) + \mathbb{P}(\bullet|B)\mathbb{P}(B)}$$



# Proportion

Over five years, no student got caught cheating in a course.  
Estimate the yearly probability to have a student cheating

$$\hat{p} = \frac{0}{5} = 0.00\%$$

Give a 95% confidence interval

Let  $X_i = \mathbf{1}(\text{cheated on year } i)$ ,

$$\mathbb{P}(X_i = x) = \begin{cases} \theta & \text{if } x = 1 \text{ cheated on year } i \\ 1 - \theta & \text{if } x = 0 \text{ did not cheat on year } i \end{cases}$$

then, the probability to have observed  $(0, 1, 0)$  over 3 years

$$\mathbb{P}(X_1 = 0, X_2 = 1, X_3 = 0) = \underbrace{\mathbb{P}(X_1 = 0)}_{=(1-\theta)} \cdot \underbrace{\mathbb{P}(X_2 = 1)}_{=\theta} \cdot \underbrace{\mathbb{P}(X_3 = 0)}_{=(1-\theta)}$$

# Proportion

Consider sample  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$ .

$$\begin{cases} p(x_i|\theta) = \theta^{x_i}[1 - \theta]^{1-x_i} \\ p(\mathbf{x}|\theta) = \theta^{\mathbf{x}^T \mathbf{1}}[1 - \theta]^{n - \mathbf{x}^T \mathbf{1}} \end{cases}$$

Assume a prior uniform distribution for  $\theta$ ,  $\pi(\theta) = 1$  on  $[0, 1]$ , so that the posterior distribution is

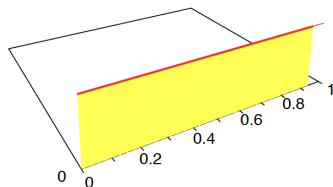
$$\pi(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)\pi(\theta)}{p(\mathbf{x})}$$

i.e.

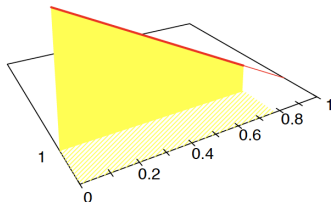
$$\pi(\theta|\mathbf{x}) \propto \theta^{\mathbf{x}^T \mathbf{1}}[1 - \theta]^{n - \mathbf{x}^T \mathbf{1}}$$

(called a Beta distribution)

$$\pi(\theta) = 1 \text{ on } [0, 1]$$



$$\pi(\theta|x_1 = 0) = 2(1 - \theta)$$

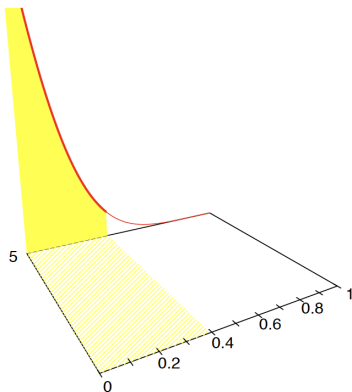
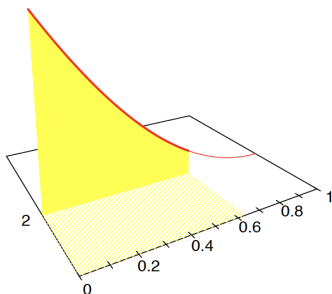


# Proportion

$$\pi(\theta|\mathbf{x}) = \frac{\theta^{\alpha+\mathbf{x}^\top \mathbf{1}} [1-\theta]^{\beta+n-\mathbf{x}^\top \mathbf{1}}}{B(\alpha+\mathbf{x}^\top \mathbf{1}, \beta+n-\mathbf{x}^\top \mathbf{1})}$$

$$\pi(\theta|(0,0,0)) = 4(1-\theta)^3$$

$$\pi(\theta|(0,0,0,0,0)) = 6(1-\theta)^5$$



# Proportion

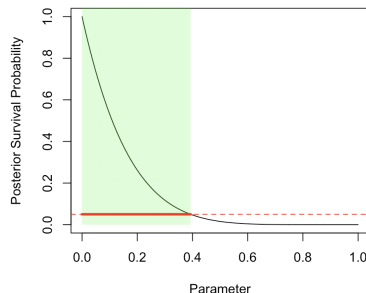
we want  $\mathbb{P}(\theta > u | \mathbf{x}) = 5\%$

or  $\mathbb{P}(\theta \in [0, u] | \mathbf{x}) = 95\%$

i.e.  $\int_0^u 6(1 - \theta)^5 d\theta = 95\%$

(i.e. quantile of Beta distribution)

Numerically, we have



```
1 > import scipy.integrate as integrate
2 > integrate.quad(lambda x: 6*(1-x)**5, 0, .393)
3 (0.9499813292139772, 1.0546911446573683e-14)
```

i.e.  $\mathbb{P}(\theta \in [0, 0.393] | \mathbf{x}) \sim 95\%$

# Estimating a proportion?

## RANDOMIZED RESPONSE: A SURVEY TECHNIQUE FOR ELIMINATING EVASIVE ANSWER BIAS

STANLEY L. WARNER  
*Claremont Graduate School*

For various reasons individuals in a sample survey may prefer not to confide to the interviewer the correct answers to certain questions. In such cases the individuals may elect not to reply at all or to reply with incorrect answers. The resulting evasive answer bias is ordinarily difficult to assess. In this paper it is argued that such bias is potentially removable through allowing the interviewee to maintain privacy through the device of randomizing his response. A randomized response method for estimating a population proportion is presented as an example. Unbiased maximum likelihood estimates are obtained and their mean square errors are compared with the mean square errors of conventional estimates under various assumptions about the underlying population.

## Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias

## Estimating a proportion?

Let  $\theta$  denote the probability to have ever cheated on your husband/wife

Consider the following strategy : flip a coin, do not show me

- ▶ head (prob 50%) : answer *have you ever seen Titanic ?*
- ▶ tail (prob 50%) : answer *have you ever seen cheated on your wife ?*

$$\mathbb{P}[\text{answer yes}] = \underbrace{\mathbb{P}[\text{answer yes}|\text{head}]}_{=85\%} \cdot \underbrace{\mathbb{P}(\text{head})}_{=50\%} + \underbrace{\mathbb{P}[\text{answer yes}|\text{tail}]}_{=\theta} \cdot \underbrace{\mathbb{P}(\text{tail})}_{=50\%}$$

(see <http://fivethirtyeight.com> for the Titanic probability)

# More Formally... Accuracy & Precision - Mean & Variance

Consider data  $\{x_1, \dots, x_n\}$  with (unknown) distribution  $F_\theta$

Let  $\hat{\theta}$  denote an estimator of (unknown)  $\theta$

$\hat{\theta}$  is **unbiased** if  $\mathbb{E}(\hat{\theta}) = \theta$

Its **precision** is given by  $\text{Var}(\hat{\theta})$

