

Section 7: Information Theory, Bootstrapping, and P-Values

Warmup: populations vs. samples

What is the difference between the population variance, σ^2 , and sample variance, S^2 ? What is the difference between sample variance, S^2 , and variance of the sample mean, $\text{Var}(\bar{X})$?

1 Song of the Quarter

This quarter in CS109 there were 167 songs that were voted on. For each song, we have a list of votes where each vote is an integer in the set $\{1, 2, 3, 4, 5\}$. We assume all votes for a song are IID samples from the “true” distribution of CS109 opinion on the song.

For each song i we have m_i votes stored in a list $\text{votes}[i] = [x_1, x_2, \dots, x_{m_i}]$. We have already calculated:

$$\begin{aligned}\mu_i &= \frac{1}{m_i} \sum_{j=1}^{m_i} x_j && \text{using } \text{np.mean}(\text{votes}[i]) \\ \text{var}_i &= \frac{1}{m_i} \sum_{j=1}^{m_i} (x_j - \mu_i)^2 && \text{using } \text{np.var}(\text{votes}[i]) \\ \text{svar}_i &= \frac{1}{m_i - 1} \sum_{j=1}^{m_i} (x_j - \mu_i)^2 && \text{using } \text{np.var}(\text{votes}[i], \text{ddof}=1)\end{aligned}$$

(a) Song 1 has $m_1 = 45$ votes. We have calculated:

$$\mu_1 = 3.82 \quad \text{var}_1 = 1.4 \quad \text{svar}_1 = 1.5$$

Estimate the probability that the true average rating for song 1 is less than 3.

(b) Song 1 has $m_1 = 45$ votes. Song 2 has $m_2 = 36$ votes. We have calculated:

$$\text{Song 1: } \mu_1 = 3.82 \quad \text{var}_1 = 1.4 \quad \text{svar}_1 = 1.5$$

$$\text{Song 2: } \mu_2 = 3.79 \quad \text{var}_2 = 1.7 \quad \text{svar}_2 = 1.8$$

What is the probability that the true average of Song 1 is greater than the true average for Song 2?

2 Entropy & Name2Age

Choosing a Diagnostic Test with Information Theory

A doctor is deciding which diagnostic test to administer to a patient. There are nine mutually exclusive possibilities: diseases A, B, C, D, E, F, G, H , or having *no disease* (labeled None).

You are given:

- **prior:** A prior distribution $P(x)$ over all diseases $x \in A, B, C, D, E, F, G, H, \text{None}$ stored in a dictionary called `prior`.
- A function `prob_pos_given_disease(test, x)` that returns the probability that a given test yields a positive result if the patient truly has disease x . (For the “None” case, this value represents the false-positive rate.)
- **tests:** A list of available tests, stored as `tests`.

When a test is run, it will return either a + or – result. However, the probability that a test is positive can vary depending on which disease the patient actually has. For example, a test designed to detect disease A may also occasionally return a positive result if the patient has disease B (a cross-reactivity), even though the true disease is not A . So while we only run one test at a time and get one result, that result provides evidence that updates our belief about *all* diseases.

The goal is to determine which test to run by choosing the one that is expected to reduce our uncertainty about the patients condition the most. To do this, write code to compute the expected uncertainty for each test.

3 Variance of Hemoglobin Levels

A medical researcher treats patients with dangerously low hemoglobin levels. She has formulated two slightly different drugs and is now testing them on patients. First, she administered drug A to one group of 50 patients and drug B to a separate group of 50 patients. Then, she measured all the patients' hemoglobin levels post-treatment. For simplicity, assume that all variation in the patient outcomes is due to their different reactions to treatment.

The researcher notes that the sample mean is similar between the two groups: both have mean hemoglobin levels around 10g/dL. However, drug B's group has a **sample variance** that is 3 (g/dL)² **greater** than drug A's group. The researcher thinks that patients respond to drugs A and B differently. Specifically, she wants to make the scientific claim that drug A's patients will end up with a significantly different spread of hemoglobin levels compared to drug B's.

You are skeptical. It is possible that the two drugs have practically identical effects and that the observed different in variance was a result of chance and a small sample size, i.e. the **null hypothesis**. Calculate the probability of the null hypothesis using bootstrapping. Here is the data. Each number is the level of an independently sampled patient:

Hemoglobin Levels of Drug A's Group ($S^2 = 6.0$): 13, 12, 7, 16, 9, 11, 7, 10, 9, 8, 9, 7, 16, 7, 9, 8, 13, 10, 11, 9, 13, 13, 10, 10, 9, 7, 7, 6, 7, 8, 12, 13, 9, 6, 9, 11, 10, 8, 12, 10, 9, 10, 8, 14, 13, 13, 10, 11, 12, 9

Hemoglobin Levels of Drug B's Group ($S^2 = 9.1$): 8, 8, 16, 16, 9, 13, 14, 13, 10, 12, 10, 6, 14, 8, 13, 14, 7, 13, 7, 8, 4, 11, 7, 12, 8, 9, 12, 8, 11, 10, 12, 6, 10, 15, 11, 12, 3, 8, 11, 10, 10, 8, 12, 8, 11, 6, 7, 10, 8, 5

How would this calculation be different if you were interested in looking at the statistical significance of the difference in sample mean? Or the 95th percentile?