

Homework Stats 1

Week 1

6733172621 Patthadon Phengpinij

Collaborators. ChatGPT (for L^AT_EX styling and grammar checking)

1 Sampling

Sampling is a process that is very important for writing simulations. In this section, you will try to sample from some common distributions. You may implement them yourself or use the provided distribution from `scipy.stats`.

```
1 from scipy.stats import norm, bernoulli, binom, uniform, geom,
   expon
2
3 # Sample from Uniform(a, b)
4 def sample_uniform(sample_size, a, b):
5     # [YOUR CODE HERE]
6     dist = uniform(a, b)
7     return dist.rvs(sample_size)
8
9 def sample_normal(sample_size, mu, sigma):
10    # [YOUR CODE HERE]
11    dist = norm(mu, sigma)
12    return dist.rvs(sample_size)
13
14 def sample_bernoulli(sample_size, p):
15    # [YOUR CODE HERE]
16    dist = bernoulli(p)
17    return dist.rvs(sample_size)
18
19 def sample_binomial(sample_size, n, p):
20    # [YOUR CODE HERE]
21    dist = binom(n, p)
22    return dist.rvs(sample_size)
23
24 def sample_geometric(sample_size, p):
25    # [YOUR CODE HERE]
26    dist = geom(p)
27    return dist.rvs(sample_size)
28
29 def sample_exponential(sample_size, l):
30    # [YOUR CODE HERE]
31    dist = expon(l)
32    return dist.rvs(sample_size)
```

TO SUBMIT

Hamtaro and his friends are collecting sunflower seeds. The bigger the sunflower, the more seeds they can find! The probability of finding a sunflower of a certain height x (in cm, from 0 to 10) increases with its height, following the probability density function $f(x) = \frac{x}{50}$. Write a function `sample_increasing(sample_size)` to simulate the heights of the sunflowers the Ham-Hams find.

Solution. First, we need to find the cumulative distribution function (CDF) from the given probability density function (PDF) $f(x) = \frac{x}{50}$. To find the CDF, we integrate the PDF from 0 to x :

$$F(x) = \int_0^x f(t) dt = \int_0^x \frac{t}{50} dt = \left[\frac{t^2}{100} \right]_0^x = \frac{x^2}{100} \text{ for } 0 \leq x \leq 10.$$

Next, we need to find the inverse of the CDF, $F^{-1}(y)$, to use the inverse transform sampling method. Setting $y = F(x)$, we have:

$$y = \frac{x^2}{100} \implies x^2 = 100y \implies x = 10\sqrt{y}.$$

Thus, the inverse CDF is $F^{-1}(y) = 10\sqrt{y}$.

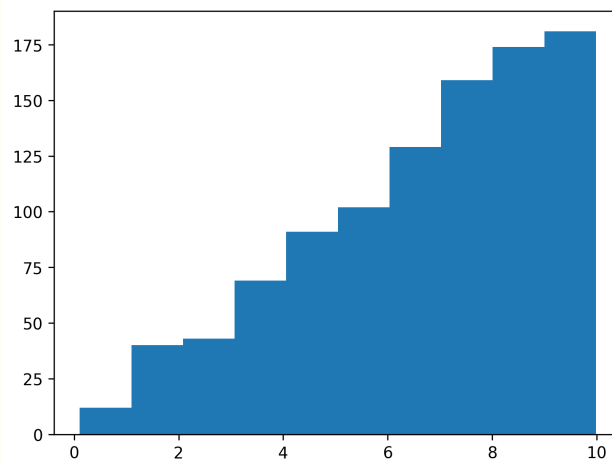
Now, we can implement the function `sample_increasing(sample_size)` to generate random samples from the distribution.

```

1 import numpy as np
2
3 np.random.seed(1) # For reproducibility
4
5 # sample from pdf f(x)=x/50, 0<=x<=10
6 def sample_increasing(sample_size):
7     # [YOUR CODE HERE]
8     u = np.random.rand(sample_size)
9     dist = 10 * np.sqrt(u)
10    return dist

```

We can plot the histogram of our samples. If the sample functions are implemented correctly, the histogram should look like our distribution.

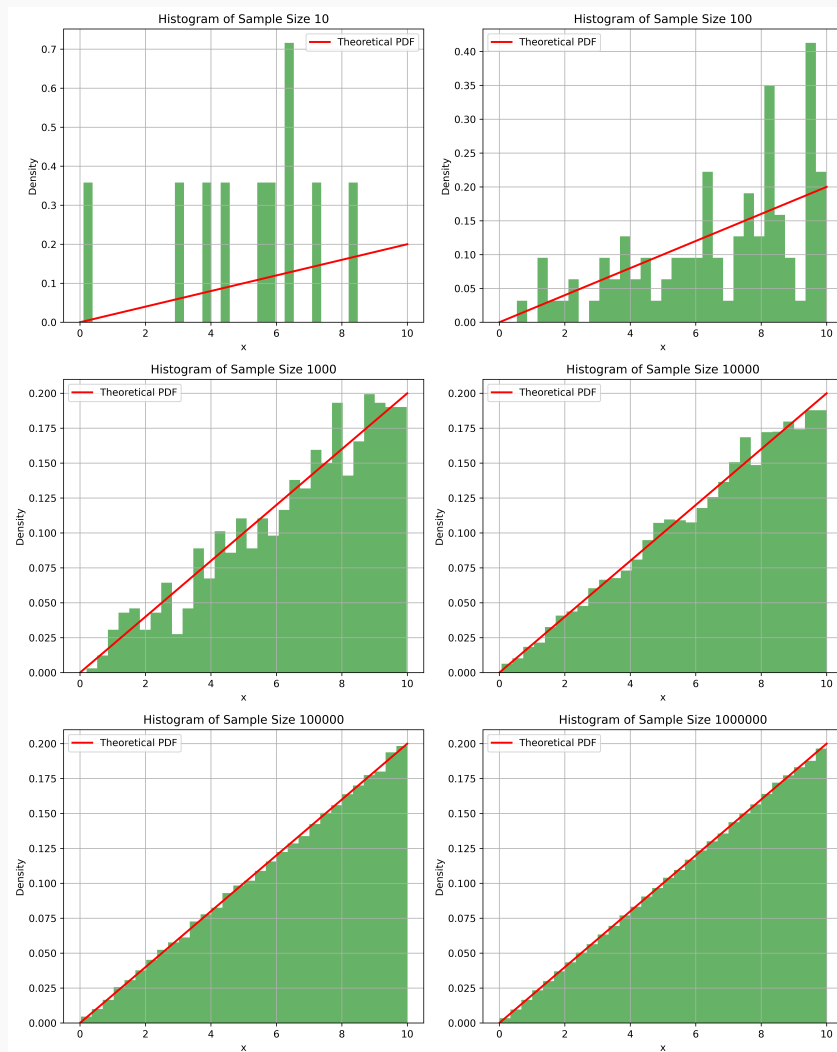


Problem 0. Try playing with the sample size and see how the histogram change with each run. Check if the result match what you think. Explain in detail.

Solution. When we increase the sample size, the histogram of the sampled data becomes smoother and more closely resembles the theoretical probability density function (PDF) $f(x) = \frac{x}{50}$. This is due to the Law of Large Numbers, which states that as the number of trials increases, the sample mean will converge to the expected value. Using Python code, we can visualize this effect by plotting histograms for different sample sizes.

```
1 import matplotlib.pyplot as plt
2
3 np.random.seed(1) # For reproducibility
4
5 Ns = [10, 100, 1000, 10000, 100000, 1000000]
6
7 # Create a 3x2 subplot grid
8 fig, axes = plt.subplots(3, 2, figsize=(12, 15))
9 axes = axes.flatten() # Flatten to make indexing easier
10
11 for i, N in enumerate(Ns):
12     data = sample_increasing(N)
13
14     # Plot histogram on the corresponding subplot
15     axes[i].hist(data, bins=30, density=True, alpha=0.6,
16                 color="g")
17
18     # Plot the theoretical PDF
19     x = np.linspace(0, 10, 100)
20     pdf = x / 50
21     axes[i].plot(x, pdf, "r-", lw=2, label="Theoretical PDF")
22
23     axes[i].set_title(f"Histogram of Sample Size {N}")
24     axes[i].set_xlabel("x")
25     axes[i].set_ylabel("Density")
26     axes[i].grid(True)
27     axes[i].legend()
28
29 plt.tight_layout()
30 plt.show()
```

The following figure shows the histograms for different sample sizes.



As we can see, with a small sample size (e.g., $N = 10$), the histogram is quite irregular and does not closely follow the theoretical PDF. However, as the sample size increases (e.g., $N = 1000, 10000, 100000, 1000000$), the histogram becomes smoother and more closely matches the shape of the theoretical PDF.

2 Maximum Likelihood Estimation

Problem 1. Machines in Hamtaro's factory have their lifetime modelled by exponential distribution with an unknown parameter. Hamtaro found out that his machines failed after x_1, x_2, \dots, x_n years. Estimate the unknown parameter.

Solution. The estimated parameter λ is calculated as the ratio of the number of observations n to the sum of the observed lifetimes $\sum_{i=1}^n x_i$.

$$\hat{\lambda} = \frac{n}{\sum_{i=1}^n x_i}$$

Using Python code, we can implement the MLE for the exponential distribution parameter λ .

```
1 import numpy as np
2
3 machine_failed_time = [2, 3, 1, 3, 4]    # In class
4 # machine_failed_time = sample_exponential(5, 0.3) #
5 # Sampled from exponential distribution
6
7 def prob1_mle(X):
8     return len(X)/np.sum(X)
9
10 print("The estimated parameter is: {}".format(prob1_mle(
11     machine_failed_time)))
```

For example, if the observed lifetimes of the machines are $[2, 3, 1, 3, 4]$ years, the estimated parameter λ would be:

$$\hat{\lambda} = \frac{5}{2 + 3 + 1 + 3 + 4} = \frac{5}{13} \approx 0.3846$$

The MLE $\hat{\lambda}$ is: 0.3846

Problem 2. Cappy is learning to perfectly replicate a new hat design. The number of attempts he needs follows a Geometric distribution with unknown parameter. For n different designs, he failed x_1, x_2, \dots, x_n times before succeed. Find the MLE of the parameter.

Solution. The estimated parameter p is calculated as the ratio of the number of successes n to the total number of trials $\sum_{i=1}^n x_i$.

$$\hat{p} = \frac{n}{\sum_{i=1}^n x_i}$$

Using Python code, we can implement the MLE for the geometric distribution parameter p .

```
1 import numpy as np
2
3 X = [0, 0, 2] # In Class Example
4 # X = sample_geometric(10, 0.9) # Sampled from actual
   geometric distribution
5
6 def prob2_mle(X):
7     # [YOUR CODE HERE]
8     return len(X)/np.sum(X)
9
10 print("The MLE is {}".format(prob2_mle(X)))
```

For example, if the number of failures before success for different designs are $[0, 0, 2]$, the estimated parameter p would be:

$$\hat{p} = \frac{3}{0+0+2} = \frac{3}{2} = 1.5$$

The MLE \hat{p} is: 1.5

Problem 3. Suppose our data x_1, x_2, \dots, x_n is randomly drawn from uniform distribution $U(a, b)$. Find MLE of a and b .

Solution. The estimated parameters \hat{a} and \hat{b} can be found using the following formulas:

$$\hat{a} = \min(x_1, x_2, \dots, x_n)$$

$$\hat{b} = \max(x_1, x_2, \dots, x_n)$$

Using Python code, we can implement the MLE for the uniform distribution parameters a and b .

```

1 import numpy as np
2
3 X = sample_uniform(100, 60, 78)
4
5 def prob3_mle(X):
6     # [YOUR CODE HERE]
7     a = np.min(X)
8     b = np.max(X) - a
9     return a, b
10
11 a, b = prob3_mle(X)
12 print("The MLE is ({}, {})".format(a, b))

```

By the code above, we can get the MLE of a and b .

$$\hat{a} = \min(x_1, x_2, \dots, x_n) = 61.2829$$

$$\hat{b} = \max(x_1, x_2, \dots, x_n) = 75.0727$$

The MLE (\hat{a}, \hat{b}) is: (61.2829, 75.0727)

TO SUBMIT

Problem 4. Dexter tracks the growth of his prized sunflower over three days (day 0, 1, and 2). He believes the sunflower's height at the end of each day y_{t+1} is its height from the previous day y_t multiplied by a secret growth factor, α , plus some random daily noise.

For a two-day period, we can observe the following Markov process:

$$P(y_2, y_1, y_0 | \alpha) = P(y_2 | y_1) P(y_1 | y_0) P(y_0 | \alpha)$$

where $y_2 \sim \mathcal{N}(\alpha y_1, \sigma^2)$, $y_1 \sim \mathcal{N}(\alpha y_0, \sigma^2)$, $y_0 \sim \mathcal{N}(0, \lambda)$

Find the Maximum Likelihood Estimate (MLE) for the secret growth factor, α , given the observed heights at the end of each day y_2, y_1, y_0 . In other words, compute for the value of α that maximizes $P(y_2, y_1, y_0 | \alpha)$.

Solution. The probability density functions for the normal distributions are given by:

$$P(y_2|y_1) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_2 - \alpha y_1)^2}{2\sigma^2}}$$

$$P(y_1|y_0) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_1 - \alpha y_0)^2}{2\sigma^2}}$$

$$P(y_0|\alpha) = \frac{1}{\sqrt{2\pi\lambda}} e^{-\frac{y_0^2}{2\lambda}}$$

Thus, the joint probability is:

$$\begin{aligned} P(y_2, y_1, y_0|\alpha) &= P(y_2|y_1)P(y_1|y_0)P(y_0|\alpha) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_2 - \alpha y_1)^2}{2\sigma^2}} \times \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_1 - \alpha y_0)^2}{2\sigma^2}} \times \frac{1}{\sqrt{2\pi\lambda}} e^{-\frac{y_0^2}{2\lambda}} \end{aligned}$$

Using **log likelihood**, we have:

$$\ln(P(y_2|y_1)) = \ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \left(-\frac{(y_2 - \alpha y_1)^2}{2\sigma^2}\right)$$

$$\ln(P(y_1|y_0)) = \ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \left(-\frac{(y_1 - \alpha y_0)^2}{2\sigma^2}\right)$$

$$\ln(P(y_0|\alpha)) = \ln\left(\frac{1}{\sqrt{2\pi\lambda}}\right) + \left(-\frac{y_0^2}{2\lambda}\right)$$

Therefore, the log-likelihood function is:

$$\begin{aligned} \ln(P(y_2, y_1, y_0|\alpha)) &= \ln(P(y_2|y_1)) + \ln(P(y_1|y_0)) + \ln(P(y_0|\alpha)) \\ \Rightarrow \left[\ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \ln\left(\frac{1}{\sqrt{2\pi\lambda}}\right) \right] &- \frac{(y_2 - \alpha y_1)^2 + (y_1 - \alpha y_0)^2}{2\sigma^2} - \frac{y_0^2}{2\lambda} \end{aligned}$$

To find the MLE of α , we take the derivative of the log-likelihood with respect to α and set it to zero:

$$\begin{aligned} \frac{d}{d\alpha} \ln(P(y_2, y_1, y_0|\alpha)) &= -\frac{-2y_1(y_2 - \alpha y_1) - 2y_0(y_1 - \alpha y_0)}{2\sigma^2} \\ &= \frac{y_1(y_2 - \alpha y_1) + y_0(y_1 - \alpha y_0)}{\sigma^2} \end{aligned}$$

$$0 = y_1(y_2 - \alpha y_1) + y_0(y_1 - \alpha y_0)$$

$$\Rightarrow y_1(y_2 - y_1\hat{\alpha}) = -y_0(y_1 - y_0\hat{\alpha})$$

$$y_1y_2 - y_1^2\hat{\alpha} = -y_0y_1 + y_0^2\hat{\alpha}$$

$$\hat{\alpha}(y_1^2 + y_0^2) = y_0y_1 + y_1y_2$$

Therefore: $\Rightarrow \hat{\alpha} = \frac{y_0y_1 + y_1y_2}{y_1^2 + y_0^2}$

3 Maximum A Posteriori Estimation

Problem 5. Hamtaro is trying to find acorns hidden by Boss. Boss has three favourite types of hiding spots:

- Type A - $P[\text{Acorn}] = c_a$
- Type B - $P[\text{Acorn}] = c_b$
- Type C - $P[\text{Acorn}] = c_c$

Boss has a habit to use Type A p_a of the time, Type B p_b of the time, and Type C p_c of the time. Find the MAP estimate for finding acorns at a new spot. Using the following data:

```
1 num_spot = 3
2 spot_acorn_prob = [0.8, 0.5, 0.4] # From slide
3 spot_select_prob = [0.4, 0.4, 0.2] # From slide
4 n = 5
5 h = 2
```

Solution. Firstly, implement the function to calculate the posterior probability of finding acorns at each type of hiding spot.

```
1 def spot_posterior(n, h, head_prob, select_prob):
2     # [YOUR CODE HERE]
3     return math.comb(n, h) * ((head_prob) ** h) * ((1 -
4         head_prob) ** (n - h)) * select_prob
```

Using the provided data, we can calculate the posterior probabilities for each type of hiding spot.

```
1 p_map = 0
2 p_map_val = 0
3 for i in range(num_spot):
4     posterior = spot_posterior(n, h, spot_acorn_prob[i],
5         spot_select_prob[i])
6     print("Spot {} has posterior of {}".format(i, posterior))
7     if posterior > p_map_val:
8         p_map_val = posterior
9         p_map = spot_acorn_prob[i]
10 print("The estimated parameter is {}".format(p_map))
```

The result of this code will give us the posterior probabilities for each type of hiding spot and the MAP estimate for finding acorns at a new spot.

```
Spot 0 has posterior of 0.020479999999999999
Spot 1 has posterior of 0.125
Spot 2 has posterior of 0.06912
The estimated parameter is 0.5
```

The MAP estimate for finding acorns at a new spot is: 0.5

Problem 6. From <https://xkcd.com/1132/>. Assume that chance of the sun actually explode is 10^{-6} . What are the chance that the machine said the sun exploded when it actually isn't?

Solution. To solve this problem, we can simulate the scenario using Python code.

```

1 def calculate_posterior_prob_of_lie(sun_prior, lie_prob):
2     """
3     Calculates the posterior probability that the machine
4     is lying (the event
5     did not happen) given that the machine reported 'YES'.
6     """
7     truth_prob = 1 - lie_prob
8
9     # P(A) - Prior probability of the event happening
10    p_A = sun_prior
11    # P(~A) - Prior probability of the event NOT happening
12    p_not_A = 1 - p_A
13
14    # P(B|A) - Likelihood of machine saying 'YES' if event
15    # happened (it tells the truth)
16    p_B_given_A = 1 - lie_prob
17    # P(B|~A) - Likelihood of machine saying 'YES' if event
18    # did not happen (it lies)
19    p_B_given_not_A = lie_prob
20
21    # P(B) - Total probability of the machine saying 'YES'
22    # (the evidence)
23    # This is the sum of true positives and false positives
24    .
25    p_B = p_B_given_A * p_A + p_B_given_not_A * p_not_A
26
27    # P(~A|B) - The posterior probability we want to find.
28    # This is the probability of a false positive, given a
29    # positive result.
30    p_not_A_given_B = p_B_given_not_A * p_not_A / p_B
31    return p_not_A_given_B
32
33 # Parameters from the problem
34 sun_prior = 1e-6
35 lie_prob = 1/36
36
37 # Calculate and print the result
38 chance_of_lie = calculate_posterior_prob_of_lie(
39     sun_prior, lie_prob)
40 print("Given the machine said 'YES'")
41 print(f"The probability it's a false alarm is: {
42     chance_of_lie:.8f}")

```

The result of this code will give us the posterior probability that the machine is lying given that it reported 'YES'.

Given the machine said 'YES'.

The probability it's a false alarm is: 0.999965

Problem 7. Go back to problem 1 - 2, and try to play with input size and parameter. Observe the change in result. Explain in detail.

7.1 Problem 1: Machines in Hamtaro's factory have their lifetime modelled by exponential distribution with an unknown parameter.

Solution. This problem involves estimating the parameter λ of an exponential distribution for the different sample sizes.

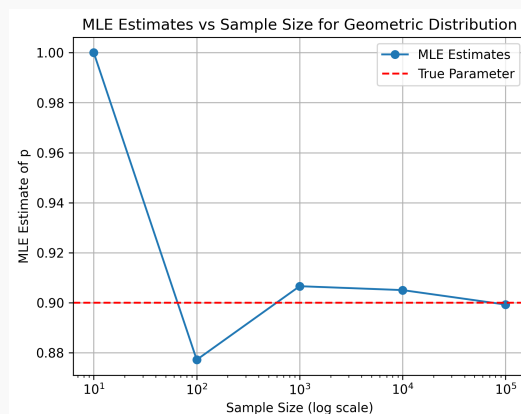
```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 sample_sizes = [10, 100, 1000, 10000, 100000]
5 true_param = 0.9
6
7 mles = np.ndarray(len(sample_sizes))
8
9 np.random.seed(1)
10
11 for i, sample_size in enumerate(sample_sizes):
12     X = sample_geometric(sample_size, true_param)
13     mle = prob2_mle(X)
14     mles[i] = mle
15
16 plt.plot(sample_sizes, mles, marker="o", label="MLE
    Estimates")
17 plt.axhline(y=true_param, color="r", linestyle="--",
    label="True Parameter")
18 plt.xscale("log")
19 plt.xlabel("Sample Size (log scale)")
20 plt.ylabel("MLE Estimate of p")
21 plt.title("MLE Estimates vs Sample Size for Geometric
    Distribution")
22 plt.legend()
23 plt.grid(True)
24 plt.show()

```

This experiment shows how the MLE estimate of the parameter p converges to the true parameter value as the sample size increases.

The following figure shows the MLE estimates for different sample sizes.



7.2 Problem 2: Cappy is learning to perfectly replicate a new hat design.

Solution. This problem involves estimating the parameter a, b of a uniform distribution for the different sample sizes.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 sample_sizes = [10, 100, 1000, 10000, 100000]
5 true_param_a = 60
6 true_param_b = 78
7
8 mles = np.ndarray((len(sample_sizes), 2))
9
10 np.random.seed(1)
11
12 for i, sample_size in enumerate(sample_sizes):
13     X = sample_uniform(sample_size, true_param_a,
14                        true_param_b)
15     mle = prob3_mle(X)
16     mles[i] = mle
17
18 fig, axes = plt.subplots(1, 2, figsize=(12, 5))
19 axes = axes.flatten()
20
21 for i, p in enumerate(["A", "B"]):
22     axes[i].plot(sample_sizes, mles[:, i], marker="o",
23                label=f"MLE Estimate of {p}")
24     axes[i].axhline(y=[true_param_a, true_param_b][i],
25                    color="r", linestyle="--", label=f"True Parameter {p}")
26
27 axes[i].set_xscale("log")
28 axes[i].set_xlabel("Sample Size (log scale)")
29 axes[i].set_ylabel(f"MLE Estimate of {p}")
30 axes[i].set_title("MLE Estimates vs Sample Size for Uniform Distribution")
31 axes[i].legend()
32 axes[i].grid(True)
33
34 plt.show()

```

Like the previous problem, this experiment shows how the MLE estimates of the parameters a and b converge to the true parameter values as the sample size increases.

The following figure shows the MLE estimates for different sample sizes.

