## IERG 3050 Assignment 3        Due: 11 April 2025

- Submit a single .pdf file containing all your answers to the Blackboard before the due date.

- Answer all questions.

- Type or write your work neatly.

---

1. We want to predict the weight $(y)$ of whitefishes from their heights $(h)$ and widths $(w)$. Here are the data obtained from a fish market:

| $y$ | $h$ | $w$ |
|------|---------|--------|
| 270 | 8.3804 | 4.2476 |
| 270 | 8.1454 | 4.2485 |
| 306 | 8.778 | 4.6816 |
| 540 | 10.744 | 6.562 |
| 800 | 11.7612 | 6.5736 |
| 1000 | 12.354 | 6.525 |

Find the plane $y = a_0 + a_1 h + a_2 w$ by multiple linear regression. What is the predicted weight of a whitefish of height 10 and width 5.5?

2. Apply Monte Carlo simulation with importance sampling to compute $\int_0^2 \frac{dx}{\sqrt{x}}$.

3. Generate 300 samples of $U(0, 1)$ by any programming language or software you like. Perform the uniformity test of these samples by chi-square goodness-of-fit test with 5 equal-length bins and at 5% significance level. Also show the histogram of the bins.

4. Generate 5 samples from a Student's $t$-distribution with 30 degrees of freedom, and then standardize them. Sort the values in ascending order, and let the values be the 0.1-, 0.3-, 0.5-, 0.7- and 0.9-quantiles. Plot the Q-Q plot of these values against the standard normal distribution. What is the conclusion of the Q-Q plot?

1. We want to predict the weight ($y$) of whitefishes from their heights ($h$) and widths ($w$). Here are the data obtained from a fish market:

| $y$ | $h$ | $w$ |
|---|---|---|
| 270 | 8.3804 | 4.2476 |
| 270 | 8.1454 | 4.2485 |
| 306 | 8.778 | 4.6816 |
| 540 | 10.744 | 6.562 |
| 800 | 11.7612 | 6.5736 |
| 1000 | 12.354 | 6.525 |

Find the plane $y = a_0 + a_1 h + a_2 w$ by multiple linear regression. What is the predicted weight of a whitefish of height 10 and width 5.5?

```
q1.py > ...
1    import pandas as pd
2
3    data = {
4        'y': [270, 270, 306, 540, 800, 1000],
5        'h': [8.3804, 8.1454, 8.778, 10.744, 11.7612, 12.354],
6        'w': [4.2476, 4.2485, 4.6816, 6.562, 6.5736, 6.525]
7    }
8    df = pd.DataFrame(data)
9    from sklearn.linear_model import LinearRegression
10
11   X = df[['h', 'w']]
12   Y = df['y']
13
14   regr = LinearRegression()
15   regr.fit(X, Y)
16
17   print('Intercept:', regr.intercept_)
18   print('Coefficients:', regr.coef_)
19
20   new_h = 10
21   new_w = 5.5
22   predicted_y = regr.predict([[new_h, new_w]])
23   print('Predicted Weight:', predicted_y[0])
```

```
(.venv) lifehater@LifedeMacBook-Pro temp %
win-arm64/bundled/libs/debugpy/adapter/..//..
Intercept: -1263.789283808669
Coefficients: [ 275.77358715 -177.31218183]
Predicted Weight: 518.7295875821146
```

2. Apply Monte Carlo simulation with importance sampling to compute $\int_0^2 \frac{dx}{\sqrt{x}}$.

For $I = \int_0^2 \frac{1}{\sqrt{x}} dx$

Exact analytical solution:

$$I = \int_0^2 x^{-\frac{1}{2}} dx = \left[ 2x^{\frac{1}{2}} \right]_0^2 = 2\sqrt{2}$$

For $f(x) = \frac{1}{\sqrt{x}}$ ⟹ we select $p(x) = \frac{1}{2\sqrt{x}}$ for $x \in (0,2]$
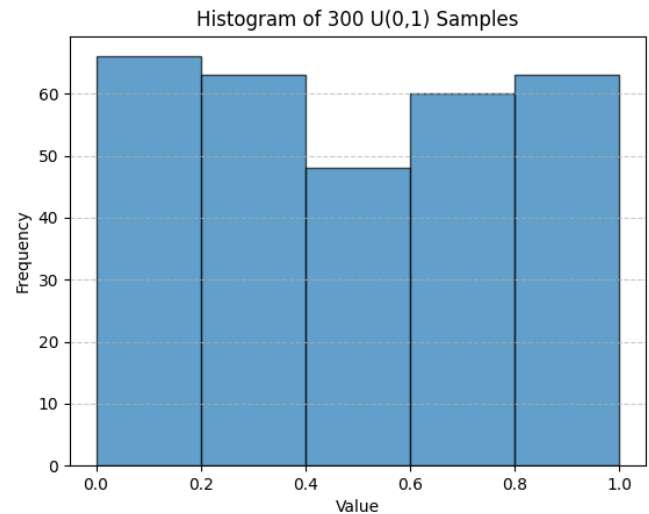
Importance sampling estimate: $I \approx \frac{1}{N} \sum_{i=1}^{N} \frac{f(x_i)}{p(x_i)} = \frac{1}{N} \sum_{i=1}^{N} \frac{\frac{1}{\sqrt{x_i}}}{\frac{1}{2\sqrt{x_i}}} = \frac{1}{N} \sum_{i=1}^{N} 2 = 2$ (✗)

Correct normalized distribution: $p(x) = \frac{1}{2\sqrt{2}x}$

$$I \approx \frac{1}{N} \sum_{i=1}^{N} \frac{\frac{1}{\sqrt{x_i}}}{\frac{1}{2\sqrt{2}\sqrt{x_i}}} = \frac{1}{N} \sum_{i=1}^{N} 2\sqrt{2} = 2\sqrt{2} \quad \text{//}$$

3. Generate 300 samples of $U(0,1)$ by any programming language or soft-ware you like. Perform the uniformity test of these samples by chi-square goodness-of-fit test with 5 equal-length bins and at 5% significance level. Also show the histogram of the bins.
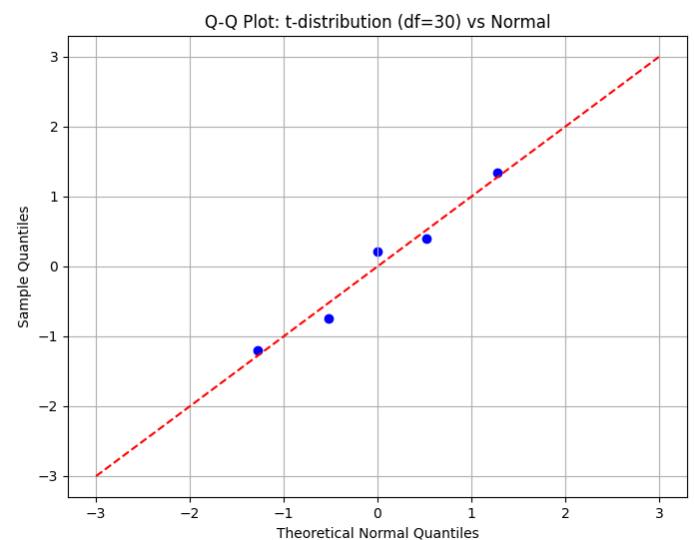
```python
q3.py > ...
1    import numpy as np
2    import matplotlib.pyplot as plt
3    from scipy.stats import chisquare
4
5    # 1. Generate 300 samples from U(0, 1)
6    np.random.seed(42)  # For reproducibility
7    samples = np.random.uniform(0, 1, 300)
8
9    # 2. Define 5 equal-length bins between 0 and 1
10   bin_edges = np.linspace(0, 1, 6)  # 6 edges for 5 bins
11   observed_counts, _ = np.histogram(samples, bins=bin_edges)
12
13   # 3. Expected counts for each bin (uniform): 300 / 5 = 60
14   expected_counts = np.array([60] * 5)
15
16   # 4. Perform chi-square goodness-of-fit test
17   chi2_stat, p_value = chisquare(f_obs=observed_counts, f_exp=expected_counts)
18
19   # 5. Output results
20   print("Observed counts:", observed_counts)
21   print("Expected counts:", expected_counts)
22   print(f"Chi-square statistic: {chi2_stat:.4f}")
23   print(f"P-value: {p_value:.4f}")
24
25   if p_value < 0.05:
26       print("Reject the null hypothesis: Data is not uniformly distributed.")
27   else:
28       print("Fail to reject the null hypothesis: Data may be uniformly distributed.")
29
30   # 6. Plot histogram
31   plt.hist(samples, bins=bin_edges, edgecolor='black', alpha=0.7)
32   plt.title("Histogram of 300 U(0,1) Samples")
33   plt.xlabel("Value")
34   plt.ylabel("Frequency")
35   plt.xticks(bin_edges)
36   plt.grid(axis='y', linestyle='--', alpha=0.7)
37   plt.show()
```



Histogram of 300 U(0,1) Samples

```
○ (.venv) lifehater@LifedeMacBook-Pro temp %
win-arm64/bundled/libs/debugpy/adapter/../.
Observed counts: [66 63 48 60 63]
Expected counts: [60 60 60 60 60]
Chi-square statistic: 3.3000
P-value: 0.5089
```

4. Generate 5 samples from a Student's $t$-distribution with 30 degrees of free-dom, and then standardize them. Sort the values in ascending order, and let the values be the 0.1-, 0.3-, 0.5-, 0.7- and 0.9-quantiles. Plot the Q-Q plot of these values against the standard normal distribution. What is the conclusion of the Q-Q plot?

```python
q4.py > ...
1    import numpy as np
2    import scipy.stats as stats
3    import matplotlib.pyplot as plt
4
5    # Set random seed for reproducibility
6    np.random.seed(42)
7
8    # 1. Generate 5 samples from t-distribution with 30 df
9    t_samples = stats.t.rvs(df=30, size=5)
10
11   # 2. Standardize the samples (though t with high df is already ~N(0,1))
12   standardized = (t_samples - np.mean(t_samples)) / np.std(t_samples, ddof=1)
13
14   # 3. Sort the values
15   sorted_samples = np.sort(standardized)
16
17   # 4. Define desired quantiles
18   quantiles = np.array([0.1, 0.3, 0.5, 0.7, 0.9])
19
20   # 5. Get theoretical normal quantiles
21   normal_quantiles = stats.norm.ppf(quantiles)
22
23   # 6. Create Q-Q plot
24   plt.figure(figsize=(8, 6))
25   plt.scatter(normal_quantiles, sorted_samples, color='blue')
26   plt.plot([-3, 3], [-3, 3], 'r--')  # y=x reference line
27   plt.title('Q-Q Plot: t-distribution (df=30) vs Normal')
28   plt.xlabel('Theoretical Normal Quantiles')
29   plt.ylabel('Sample Quantiles')
30   plt.grid(True)
31   plt.show()
32
33   # 7. Calculate correlation for additional insight
34   correlation = np.corrcoef(normal_quantiles, sorted_samples)[0, 1]
35   print(f"Correlation between theoretical and sample quantiles: {correlation:.4f}")
```



Q-Q Plot: t-distribution (df=30) vs Normal

Correlation between theoretical and sample quantiles: 0.9850