

Applied Statistics Computational Project 2

For the dataset, We have taken the Video game sales Data.

Conclusions from Gamma Distribution Fitting of Global_Sales Data

```
1 import pandas as pd
2 import numpy as np
3 from scipy.stats import gamma
4
5 # Load the Excel file
6 df = pd.read_excel("newgame.xlsx")
7
8 # Extract Global_Sales and filter positive values
9 global_sales = df['Global_Sales']
10 global_sales = global_sales[global_sales > 0].dropna()
11
12 # Method of Moments Estimation
13 sample_mean = global_sales.mean()
14 sample_var = global_sales.var()
15
16 # MoM formulas for Gamma(a, b)
17 b_mom = sample_var / sample_mean
18 a_mom = sample_mean / b_mom
19
20 print("Method of Moments Estimates:")
21 print(f"a (shape) = {a_mom:.4f}")
22 print(f"b (scale) = {b_mom:.4f}")
23
24 # Maximum Likelihood Estimation (fit Gamma with loc=0)
25 a_mle, loc_mle, b_mle = gamma.fit(global_sales, floc=0)
26
27 print("\nMaximum Likelihood Estimates:")
28 print(f"a (shape) = {a_mle:.4f}")
29 print(f"b (scale) = {b_mle:.4f}")
```

The objective of the analysis is to model the distribution of `Global_Sales` using the Gamma distribution. Two estimation methods were used:

- Method of Moments (MoM)
- Maximum Likelihood Estimation (MLE)

Estimated Parameters

Method	Shape (a)	Scale (b)
Method of Moments	2.1766	6.7504
Maximum Likelihood Estimation	3.9330	3.7359

Interpretations and Conclusions

1. Both estimation methods suggest that the Gamma distribution is a reasonable model for the `Global_Sales` data. This is justified since the sales data is non-negative and skewed, which aligns with the properties of the Gamma distribution.
2. The two methods yield different parameter values:
 - The MoM estimates indicate a lower shape parameter and a higher scale parameter, suggesting a more right-skewed distribution.
 - The MLE estimates yield a higher shape parameter and a lower scale parameter, suggesting less skewness and more concentration near the mean.
3. MLE is generally more accurate and statistically efficient than MoM, especially with larger datasets. While MoM provides a quick and intuitive estimation based on sample moments, MLE makes use of the full likelihood function.
4. For robust statistical inference or simulation purposes, the MLE-based Gamma model is preferable.

Confidence Interval for the Variance of Global_Sales

```
1 import pandas as pd
2 from scipy.stats import chi2
3
4 # Sample statistics
5 n = len(global_sales)
6 sample_var = global_sales.var(ddof=1) # unbiased estimator
7
8 # Confidence level
9 alpha = 0.05
10
11 # Chi-squared critical values (corrected)
12 chi2_lower = chi2.ppf(alpha / 2, df=n - 1)
13 chi2_upper = chi2.ppf(1 - alpha / 2, df=n - 1)
14
15 # Confidence interval for the variance
16 lower_bound = (n - 1) * sample_var / chi2_upper
17 upper_bound = (n - 1) * sample_var / chi2_lower
18
19 print(f"95% Confidence Interval for the Variance: ({lower_bound:.4f}
    }, {upper_bound:.4f})")
```

Using the sample of `Global_Sales` (filtered to include only positive values), we compute a 95% confidence interval for the population variance based on the Chi-squared distribution.

Method

Let s^2 denote the sample variance computed from n independent observations. The confidence interval for the true variance σ^2 of a normal population is given by:

$$\left(\frac{(n-1)s^2}{\chi_{n-1,1-\alpha/2}^2}, \frac{(n-1)s^2}{\chi_{n-1,\alpha/2}^2} \right)$$

where:

- n is the sample size,
- s^2 is the unbiased sample variance,
- $\chi_{n-1,\alpha/2}^2$ and $\chi_{n-1,1-\alpha/2}^2$ are the critical values from the Chi-squared distribution with $n-1$ degrees of freedom,
- $\alpha = 0.05$ corresponds to a 95% confidence level.

Results

For the given data:

95% Confidence Interval for the Variance: (76.3680, 134.0711)

Interpretation

We are 95% confident that the true variance of global sales lies between approximately 76.37 and 134.07. This interval quantifies the uncertainty around the sample variance estimate due to sampling variability.

Confidence Interval for Difference in Mean Global Sales: Sports vs Platform

```
1 import pandas as pd
2 import numpy as np
3 from scipy import stats
4
5 # Load Excel file
6 df = pd.read_excel('newgame.xlsx')
7
8
9 # Confidence Interval for Difference in Means
10
11
12 # Choose two non-overlapping genres
13 genre1 = 'Sports'
14 genre2 = 'Platform'
15
16 # Filter data for each genre
17 sales1 = df[df['Genre'] == genre1]['Global_Sales'].dropna()
18 sales2 = df[df['Genre'] == genre2]['Global_Sales'].dropna()
19
20 # Calculate means, standard deviations, and sample sizes
21 n1, n2 = len(sales1), len(sales2)
22 mean1, std1 = sales1.mean(), sales1.std(ddof=1)
23 mean2, std2 = sales2.mean(), sales2.std(ddof=1)
24
25 # Compute standard error and margin of error
26 se_diff = np.sqrt(std1**2 / n1 + std2**2 / n2)
27 z_critical = stats.norm.ppf(0.975) # 95% CI
28 margin_of_error = z_critical * se_diff
29
30 # Confidence Interval
31 ci_lower = (mean1 - mean2) - margin_of_error
32 ci_upper = (mean1 - mean2) + margin_of_error
33
34 # Output for Part 1
35 print("=== 95% Confidence Interval for Difference in Means ===")
36 print(f"Genre 1: {genre1} | Mean: {mean1:.3f}")
37 print(f"Genre 2: {genre2} | Mean: {mean2:.3f}")
38 print(f"95% CI for (Mean of {genre1} - Mean of {genre2}): ({ci_lower:.3f}, {ci_upper:.3f})\n")
```

Objective

To determine whether there is a statistically significant difference in the average global sales between the **Sports** and **Platform** (Non - Overlapping populations) genres, we compute a 95% confidence interval for the difference in their means.

Method

Let:

- \bar{x}_1 , s_1 , and n_1 be the sample mean, standard deviation, and size for **Sports**,
- \bar{x}_2 , s_2 , and n_2 be the same for **Platform**.

The confidence interval for the difference in means ($\mu_1 - \mu_2$) is given by:

$$(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

with $z_{\alpha/2} = 1.96$ for a 95% confidence level.

Results

- Mean Global Sales for **Sports**: $\bar{x}_1 = 29.502$
- Mean Global Sales for **Platform**: $\bar{x}_2 = 15.745$
- 95% Confidence Interval for ($\mu_{\text{Sports}} - \mu_{\text{Platform}}$):

$$(-8.942, 36.456)$$

Interpretation

The confidence interval includes 0, indicating that the difference in mean global sales between **Sports** and **Platform** genres is not statistically significant at the 95% confidence level. That is, we do not have sufficient evidence to conclude that one genre significantly outsells the other on average.

Hypothesis Test for Proportion of High-Selling Games

```
1
2 # Hypothesis Test for Binary Response (Bernoulli)
3
4
5 # Define binary response variable: success if Global_Sales > 10
  million
6 df['Success'] = (df['Global_Sales'] > 10).astype(int)
7
8 # Sample proportion and size
9 p_hat = df['Success'].mean()
10 n = len(df['Success'])
11 p0 = 0.5 # Null hypothesis value
12
13 # Z-test statistic and p-value
14 z_stat = (p_hat - p0) / np.sqrt(p0 * (1 - p0) / n)
15 p_value = 1 - stats.norm.cdf(z_stat)
16
17 # Output for Part 2
18 print("=== Hypothesis Test for Proportion (Bernoulli) ===")
19 print("H0: p <= 0.5 vs H1: p > 0.5")
20 print(f"Sample proportion (p ): {p_hat:.3f}")
21 print(f"Z-statistic: {z_stat:.3f}")
22 print(f"P-value: {p_value:.4f}")
23 if p_value < 0.05:
24     print("Conclusion: Reject H0 at 5% significance level. Evidence
      suggests p > 0.5.")
25 else:
26     print("Conclusion: Fail to reject H0. No strong evidence that p
      > 0.5.")
```

Objective

To test whether the proportion of games with global sales greater than 10 million units exceeds 0.5.

Definitions

- Define a binary variable `Success` such that:

$$\text{Success} = \begin{cases} 1, & \text{if Global_Sales} > 10\text{million} \\ 0, & \text{otherwise} \end{cases}$$

- Let \hat{p} denote the sample proportion of successful games.
- Let n be the total number of observations.

Hypotheses

$H_0 : p \leq 0.5$ (Proportion of successful games is 50% or less)

$H_1 : p > 0.5$ (Proportion of successful games exceeds 50%)

Test Statistic

The test statistic is computed using the formula for a one-sample proportion Z-test:

$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

where:

- $\hat{p} = 0.626$ (sample proportion),
- $p_0 = 0.5$ (null hypothesis proportion),
- n is the total sample size.

Results

- Z-statistic: $Z = 2.513$
- P-value: 0.0060

Conclusion

Since the p-value is less than the significance level $\alpha = 0.05$, we reject the null hypothesis H_0 .

Conclusion: There is statistically significant evidence at the 5% level to suggest that the proportion of high-selling games (over 10 million units) exceeds 50%.

Done By :

- **Puli Dinesh - AI23BTECH11019**
- **Rathod sai dhanush - AI23BTECH11021**
- **Jatavath ajay- AI23BTECH11011**