



# Statistics Example Reports



# Key Statistical Concepts and Terminology

## Population vs. Sample

- **Population:** The entire group that you want to draw conclusions about.
- **Sample:** A subset of the population, selected for analysis.

## Parameter vs. Statistic

- **Parameter:** A numerical value summarizing a characteristic of the population (e.g., population mean).
- **Statistic:** A numerical value summarizing a characteristic of the sample (e.g., sample mean).

## Types of Data

- **Quantitative:** Numerical values (e.g., height, weight, sales).
  - **Continuous:** Can take any value within a range (e.g., height).
  - **Discrete:** Can take specific values (e.g., number of customers).
- **Qualitative:** Categorical data (e.g., gender, type of car).

# Practical Application: Analyzing Data with Artificial Intelligence

1. **Data Collection:** Gather data through surveys, experiments, or secondary data sources.
2. **Data Cleaning:** Handle missing values, remove duplicates, and correct errors.
3. **Descriptive Analysis:** Summarize the main features of the data.
4. **Inferential Analysis:** Draw conclusions about the population based on sample data.
5. **Reporting:** Communicate findings using visualizations and statistical summaries.

## Practice Exercise – Run your data in Artificial Intelligence like Chat GPT to get the results below:

1. Collect a sample dataset relevant to your field.
2. Calculate the mean, median, mode, range, variance, and standard deviation.
3. Create a histogram and box plot to visualize the distribution.
4. Conduct a simple hypothesis test (e.g., t-test) on the dataset.

Perform a linear regression analysis and interpret the results.

**Visuals:** Box plots, scatter plots, grouped bar charts for age and gender.

# Practical Application: Analyzing Data

Synthetic Dataset

<input checked="" type="checkbox"/>	Age	Experience	Salary
1	50	1	26740
2	36	35	51028
3	29	37	34065
4	42	14	57103
5	40	3	36208
6	44	1	73249

I have generated a synthetic dataset based on the requirements you provided. It includes columns for Age, Experience, and Salary for a sample size of 100 individuals.

# Next steps in the exercise

Next steps based on the exercise:

**1. Descriptive Analysis:**

- I'll calculate the mean, median, mode, range, variance, and standard deviation for the dataset.

**2. Visualization:**

- I will create a histogram and box plot to visualize the distribution of one of the variables (e.g., salary).

**3. Hypothesis Test:**

- I can conduct a t-test on the salary variable to compare against a hypothetical population mean.

**4. Linear Regression:**

- I will perform a simple linear regression to see if there is a relationship between experience and salary.

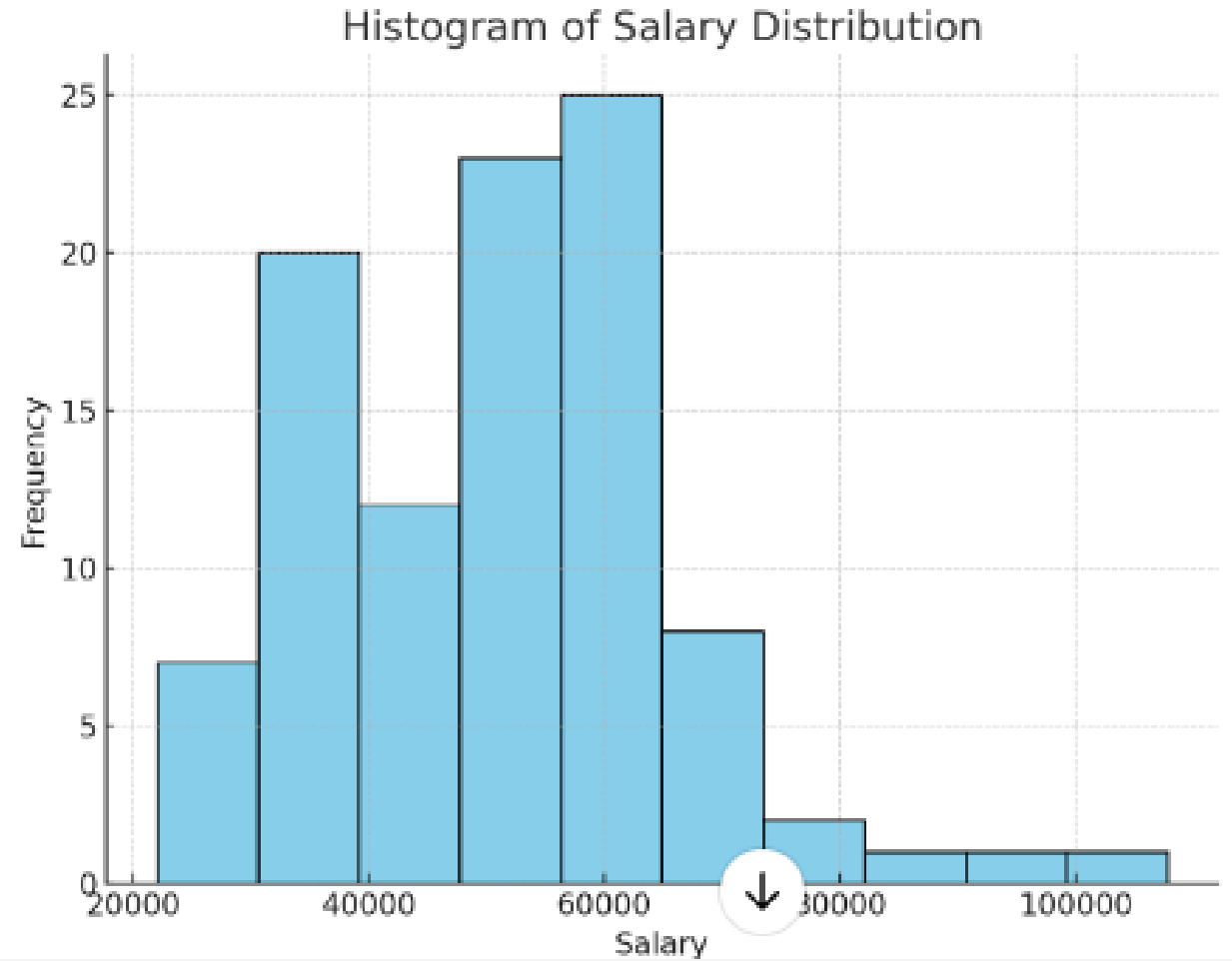
# Descriptive Statistics

Descriptive Statistics					↓	↗
<input checked="" type="checkbox"/>		Mean	Median	Mode		
1	Age	40.06	41.5	45		
2	Experience	20.48	22.5	32		
3	Salary	51231.69	53311.0	57725		

I have calculated the descriptive statistics for the dataset. The table includes the mean, median, mode, range, variance, and standard deviation for the variables Age, Experience, and Salary.

# Histogram of Salary Distribution

## Histogram Of Salary Distribution



Here is the histogram showing the distribution of salaries in the dataset. You can see that most salaries fall between \$30,000 and \$60,000, with fewer individuals earning over \$80,000.

3 rows

# T – Statistic and P - Value

The t-test results for the salary data are as follows:

- **T-Statistic:** 0.82
- **P-Value:** 0.41

Since the p-value is greater than the standard significance level (e.g., 0.05), we fail to reject the null hypothesis. This suggests that there is no significant difference between the sample mean salary and the hypothetical population mean of \$50,000.

The linear regression analysis results for Salary (dependent variable) and Experience (independent variable) are as follows:

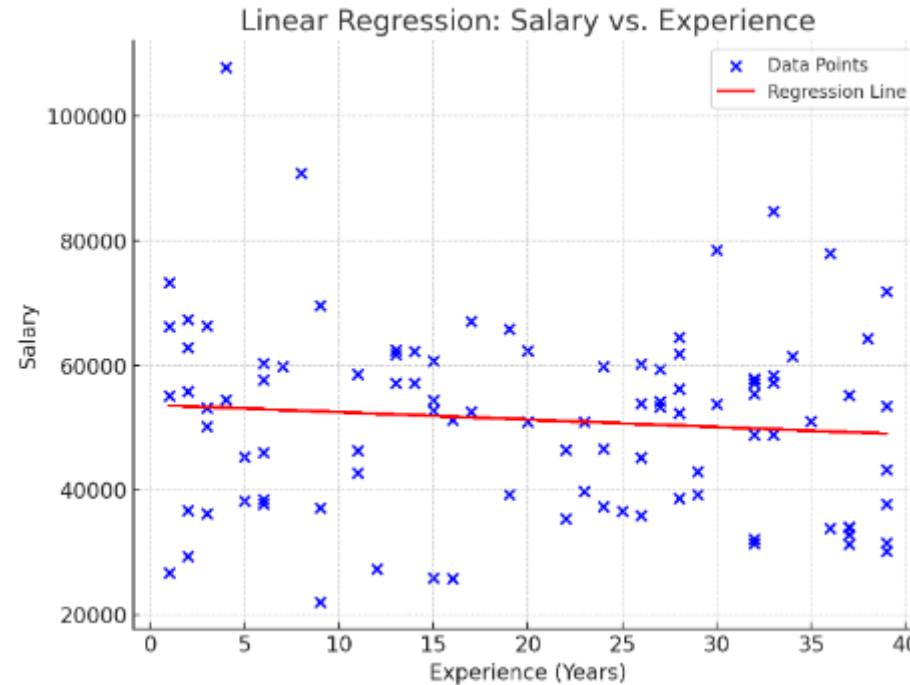
- **Slope (Coefficient):** -119.12 (suggesting that for each additional year of experience, the salary decreases by approximately \$119, though this is likely not significant)
- **Intercept:** 53,671.32 (the estimated salary when experience is 0)
- **R-Squared:** 0.0098 (indicating that only about 0.98% of the variation in salary is explained by years of experience)

This low R-squared value suggests that experience is not a strong predictor of salary in this dataset.



# Linear Regression

Linear Regression: Salary Vs. Experience



Here is the scatter plot of Salary versus Experience, with the regression line added in red. As you can see, the slope of the line is slightly negative, indicating a weak relationship between experience and salary. The scatter of data points suggests that experience alone doesn't explain much of the variation in salary.

# Summary of Exercise

The **p-value** from a hypothesis test helps us determine whether the observed relationship or result is statistically significant.

In the context of a **t-test** or **linear regression**:

**1. For the t-test** (one-sample t-test on salary):

- The p-value was **0.41**, which is much greater than the common significance level of 0.05.
- **Interpretation:** This means that we do not have sufficient evidence to reject the null hypothesis. In other words, the sample mean salary is not significantly different from the hypothetical population mean of \$50,000.

**2. For linear regression:**

- Although we didn't explicitly calculate the p-value for the regression slope, the low **R-squared** value and the near-zero slope indicate that the relationship between salary and experience is weak. If we were to calculate the p-value for the slope, it would likely be large (similar to the t-test result), suggesting that the relationship is not statistically significant.
- **Interpretation:** A large p-value would indicate that the observed negative slope could be due to random variation in the sample, and we cannot confidently claim that experience has a meaningful impact on salary.

In summary, a high p-value tells us that the observed result (whether a mean difference or a relationship in regression) is not statistically significant, and there's no strong evidence to support that the effect exists in the population.

# Gauss-Markov Theorem and Parameter Properties in AI: Unveiling Bias, Consistency, and Efficiency of Estimators

This course delves into the Gauss-Markov Theorem and its significance in artificial intelligence (AI) applications. The theorem's principles concerning parameter estimation, bias, consistency, and efficiency are explored, emphasizing their relevance in AI model development. Through comprehensive examples and analyses, this paper underscores how adhering to the Gauss-Markov Theorem can enhance the quality and reliability of AI estimators.

## 1. Introduction:

The Gauss-Markov Theorem, a cornerstone of statistical theory, plays a pivotal role in the field of parameter estimation. Originally formulated for linear regression, the theorem's concepts transcend into various AI applications. This paper investigates the theorem's key principles—bias, consistency, and efficiency—in the context of AI, showcasing their implications on the accuracy and robustness of estimators.

## 2. Gauss-Markov Theorem: An Overview:

The Gauss-Markov Theorem establishes that under certain assumptions, the ordinary least squares (OLS) estimator of regression coefficients possesses desirable properties. These properties include unbiasedness, consistency, and efficiency. These notions can be extended beyond linear regression, making the theorem applicable in AI domains.