

Statistics and Data Analysis

Unit 02 – Lecture 06 Notes

In-Class Activity: Summarization + Interpretation

Tofik Ali

February 9, 2026

Purpose of This Activity

This activity is designed to help you practice descriptive statistics end-to-end. You will not only compute numbers, but also interpret what they mean and write a short conclusion.

Repository. <https://github.com/tali7c/Statistics-and-Data-Analysis>

Learning Outcomes

After this activity, you should be able to:

1. compute central tendency (mean, median, mode),
2. compute dispersion (range, IQR, sample variance, sample std),
3. compute Pearson correlation and interpret its sign and magnitude,
4. compare groups using grouped summaries,
5. write clear insights and limitations from statistical summaries.

1. Dataset Description

File: `data/activity_student_dataset.csv`

1.1 What the columns mean

- `program`: CSE / ECE / AIML (categorical).
- `attendance_pct`: attendance percentage (numeric).
- `study_hours`: study hours (numeric).
- `social_media_hours`: social media hours (numeric).
- `final_score`: final score (numeric).

2. Tasks (What You Must Compute)

Task 1: Central tendency

For `final_score`, compute:

- mean: $\bar{x} = \frac{1}{n} \sum x_i$
- median (middle value after sorting)
- mode (most frequent value)

Interpretation tip. Compare mean vs median:

- mean $>$ median often hints right skew,
- mean $<$ median often hints left skew.

This is only a clue; confirm with a plot.

Task 2: Dispersion

For `final_score`, compute:

- range = max – min
- quartiles Q_1, Q_3 and IQR = $Q_3 - Q_1$
- sample variance and sample standard deviation:

$$s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2, \quad s = \sqrt{s^2}$$

Robustness. IQR is more robust to outliers than standard deviation.

Task 3: Correlation

Compute Pearson correlation:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

for:

- (`study_hours, final_score`)
- (`social_media_hours, final_score`)

Very important. Correlation measures **linear association**. It does **not** prove causation.

Task 4: Grouped summaries

Group by `program` and compute:

- mean and median of `final_score`
- mean of `attendance_pct`

Task 5: Write-up

Write:

- 3 insights (what the dataset suggests)
- 2 limitations (why your conclusions might not generalize)

3. Expected Results (From the Provided Solution)

If you run the provided solution script (next section), you should obtain:

3.1 Overall results

- $\text{mean}(\text{final_score}) = 65.50$
- $\text{median}(\text{final_score}) = 65.50$
- $\text{mode}(\text{final_score}) = 60$
- $Q_1 = 60, Q_3 = 72, IQR = 12$
- sample std(final_score) ≈ 14.11

3.2 Correlations

- $\text{corr}(\text{study_hours}, \text{final_score}) \approx 0.5190$ (moderate positive linear association)
- $\text{corr}(\text{social_media_hours}, \text{final_score}) \approx -0.9771$ (strong negative linear association)

Interpretation note. Even a strong correlation does not prove that one variable causes the other. Other variables (motivation, prior knowledge, teaching quality) could be involved.

3.3 By program (mean final score)

- AIML: mean ≈ 75.83
- CSE: mean ≈ 59.17
- ECE: mean ≈ 61.50

4. Provided Solution Script (Mini Demo)

After you attempt the activity yourself, run:

```
python demo/activity_solution.py
```

It will save:

- `data/overall_results.csv`
- `data/summary_by_program.csv`
- plots in `images/` (scatter plots, bar chart, histogram)

5. Common Mistakes

- **Mixing up correlation and causation.** Correlation is not proof.
- **Using only the mean.** Always look at median and IQR too.
- **Ignoring groups.** A global average can hide group differences.
- **Not writing limitations.** Every conclusion must mention what could be wrong.
- **No plots.** Tables alone can hide skewness and outliers.

6. Extension Questions (Optional)

If you finish early, try any two:

1. Identify the lowest and highest scoring students and comment on their study/social media hours.
2. Compare correlation within each program separately (CSE vs ECE vs AIML).
3. Replace mean with median for program comparison and see if ranking changes.
4. Use IQR fences to flag potential outliers in `final_score`.

References

- Montgomery, D. C., & Runger, G. C. *Applied Statistics and Probability for Engineers*, Wiley, 7th ed., 2020.
- Freedman, D., Pisani, R., & Purves, R. *Statistics*, W. W. Norton, 4th ed., 2007.
- McKinney, W. *Python for Data Analysis*, O'Reilly, 2022.