| EXPT NO:6 DATE: 09.02.2026 | OVER-PLOTTING REDUCTION TECHNIQUES |
|---|---|

**PRE-LAB QUESTIONS**

1.  **Why is over-plotting common in big data visualization?**
    Over-plotting is common because big data contains thousands or millions of points, and many points fall in the same positions on a graph. This causes points to overlap and hide the actual distribution.

2.  **How does data density affect perception?**
    When data density is high, the plot looks like a dark clustered region. This makes it difficult to identify important patterns like clusters, trends, and outliers, reducing the readability of the visualization.

3.  **What trade-offs exist between detail and clarity?**
    *   Showing every data point gives high detail, but it creates clutter and reduces clarity.
    *   Using aggregation/binning improves clarity, but it reduces the visibility of individual data points.
        So the trade-off is between full detail vs clear understanding.

4.  **How do AI datasets increase visualization complexity?**
    AI datasets are usually large-scale, high-dimensional, and continuous, with many features and classes. This increases overlap in plots and makes it harder to visualize relationships clearly using normal plotting methods.

5.  **Why is over-plotting a serious analytical risk?**
    Over-plotting can hide important information such as:
    *   minority patterns
    *   outliers
    *   true clusters
    *   real correlations
    This can lead to wrong conclusions and wrong decisions, especially in AI and business analytics.

**OBJECTIVE** : To apply techniques that reduce visual clutter in large-scale datasets.

**SCENARIO A social media analytics company visualizes millions of user interactions to study engagement patterns.**

**IN-LAB TASKS (Using R Language) • Apply alpha blending • Implement jittering techniques • Use aggregation and binning**

**CODE:**

```r
1 # ================================================================
2 # EXPT NO: 7- OVER-PLOTTING REDUCTION TECHNIQUES
3 # ROLL NO: 23BAD122
4 # ================================================================
5
6 if (!require(ggplot2)) install.packages("ggplot2", dependencies = TRUE)
7 if (!require(readr)) install.packages("readr", dependencies = TRUE)
8
9 library(ggplot2)
10 library(readr)
11
12 df <- read_csv("7.social_media_interactions.csv")
13
14 print(head(df))
15
16 x_col <- "Likes"
17 y_col <- "Comments"
18
19 df_clean <- df[!is.na(df[[x_col]]) & !is.na(df[[y_col]]), ]
20
21 cat("✅ Rows in dataset:", nrow(df_clean), "\n")
22
23 # ================================================================
24 # TASK 1: ALPHA BLENDING
25 # ================================================================
26 p_alpha <- ggplot(df_clean, aes(x = Likes, y = Comments)) +
27   geom_point(alpha = 0.2) +
28   ggtitle("Alpha Blending (Likes vs Comments)") +
29   xlab("Likes") + ylab("Comments")
30
31 print(p_alpha)
32
33 # Save output
34 ggsave("Alpha_Blending_Output.png", p_alpha, width = 7, height = 5)
35
36 # ================================================================
37 # TASK 2: JITTERING
38 # ================================================================
39 p_jitter <- ggplot(df_clean, aes(x = Likes, y = Comments)) +
40   geom_jitter(width = 30, height = 10, alpha = 0.3) +
```

**OUTPUT:**

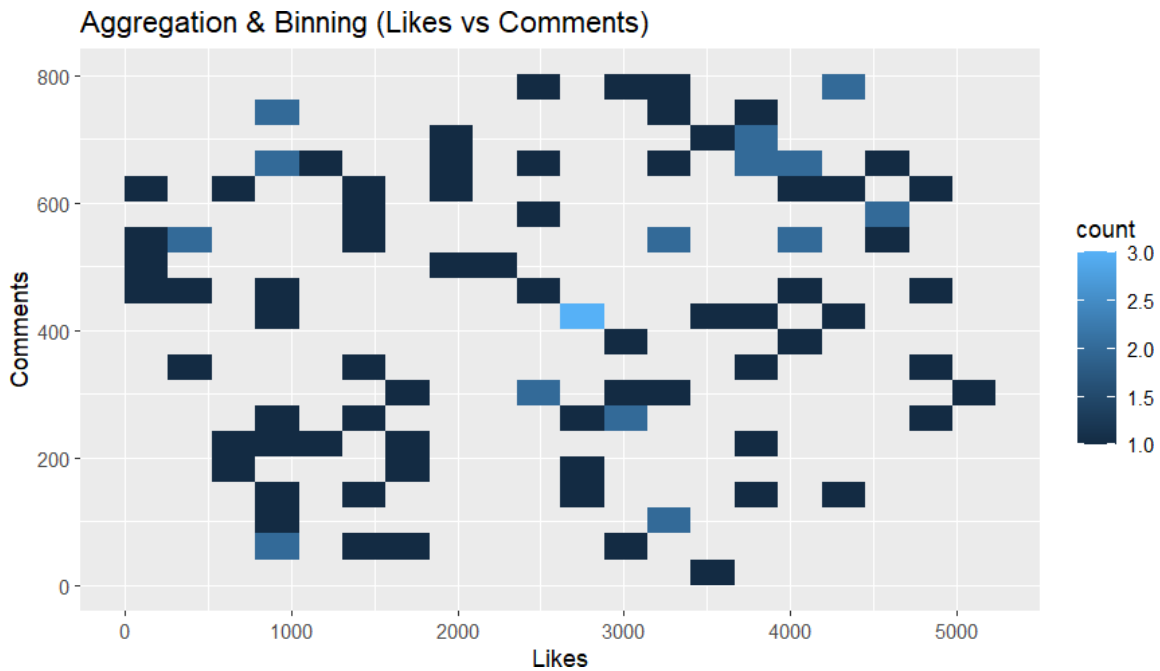## Alpha Blending (Likes vs Comments)



## Jittering (Likes vs Comments)

## Aggregation & Binning (Likes vs Comments)



**POST-LAB QUESTIONS**

1. **Which technique provided the best clarity and why?**
   Aggregation and binning (2D bin / hexbin) provided the best clarity because it groups overlapping points into bins and shows data density clearly. This makes patterns visible even when the dataset is large.

2. **How does over-plotting distort analytical conclusions?**
   Over-plotting hides important patterns such as:
   - clusters
   - trends
   - outliers
   - rare behaviors

   This can make the data look misleading and can cause incorrect conclusions about relationships between variables.

3. **When should aggregation be preferred over raw plotting?**
   Aggregation should be preferred when:
   - the dataset is very large
   - points overlap heavily
   - the goal is to understand overall patterns (macro trends)
   - raw scatter plots become unreadable

4. **How do these techniques support scalable AI analytics?**
   These techniques help scalable AI analytics by:
   - revealing real patterns in large datasets
   - improving readability and interpretability
   - supporting better feature analysis
   - helping detect bias, imbalance, and outliers
   - making large AI datasets easier to validate visually

5. **Explain real-world consequences of ignoring over-plotting.**
   Ignoring over-plotting can cause:
   - wrong business decisions (marketing, targeting, engagement)
   - missing fraud/spam patterns
   - incorrect trend detection
   - poor AI model training insights
   - failure to identify rare but important user behaviors

**LEARNING OUTCOME: Students master over-plotting reduction for big data visual analytics.**

## ASSESSMENT

| Description | Max Marks | Marks Awarded |
|---|---|---|
| Pre Lab Exercise | 5 | |
| In Lab Exercise | 10 | |
| Post Lab Exercise | 5 | |
| Viva | 10 | |
| Total | 30 | |
| Faculty Signature | | |