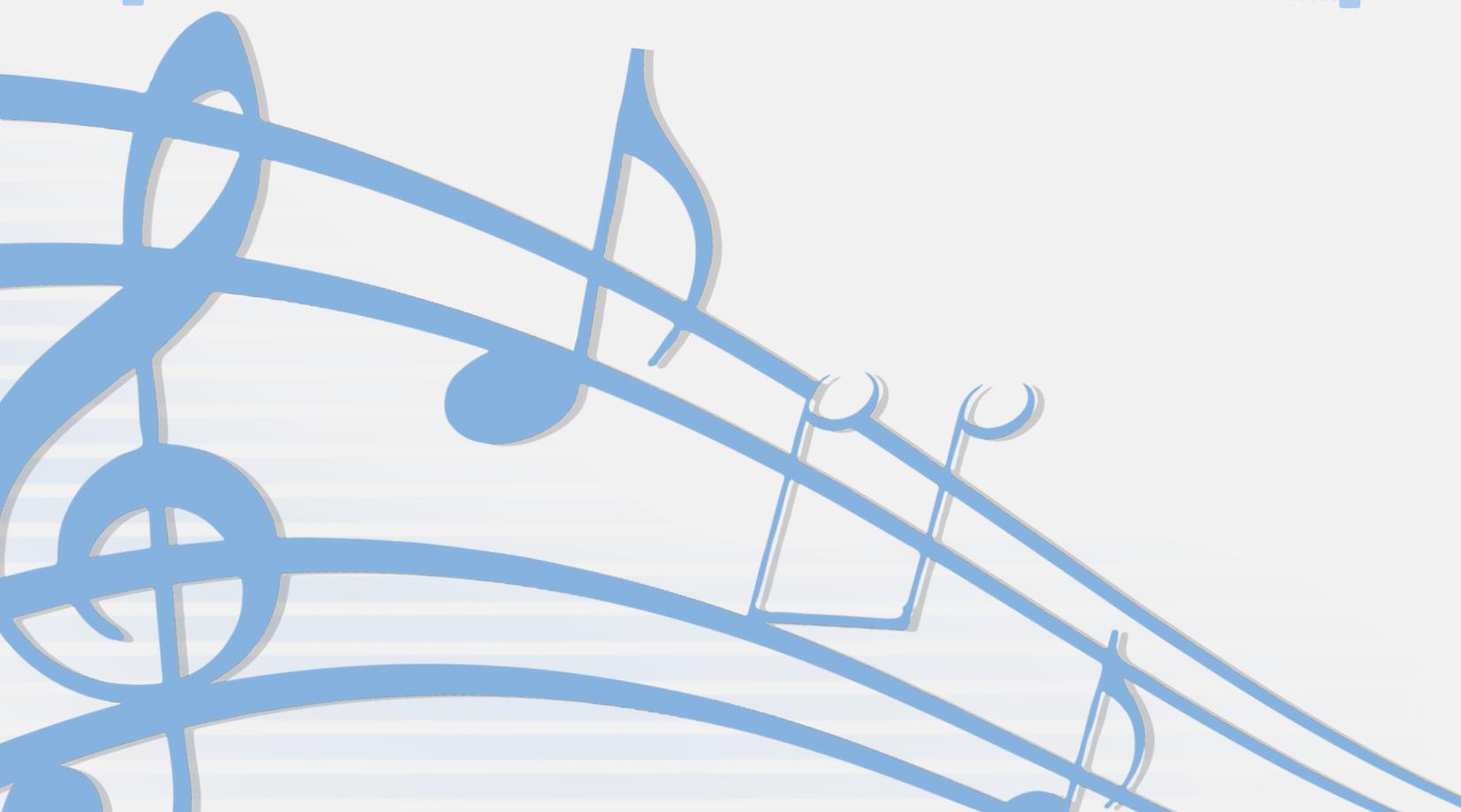


# **Collaboration vs Solo Songs - Popularity Analysis**

**By Zixu Xing**

# Research Question:

- Are collaboration songs more popular than solo songs?



# Measures of popularity:

- Average log(Streams)
- Probability of charting (charted = 1 if in\_spotify\_charts > 0)

# Data Cleaning and Preparation:

◆ Goal: create an analysis-ready dataset (df2)

◆ Main steps:

1. Made a working copy: df2 <- df
2. Converted streams from character to numeric: streams\_num
3. Removed rows with invalid numeric conversion (1 row removed)
4. Created group labels:
  - a. solo if artist\_count == 1
  - b. collab if artist\_count > 1
5. Created log\_streams = log(streams\_num)
6. Created charted:
  - a. 1 if in\_spotify\_charts > 0
  - b. 0 otherwise

◆ Final sample sizes:

- collab: 366
- solo: 586
- charted = 1: 548
- charted = 0: 404

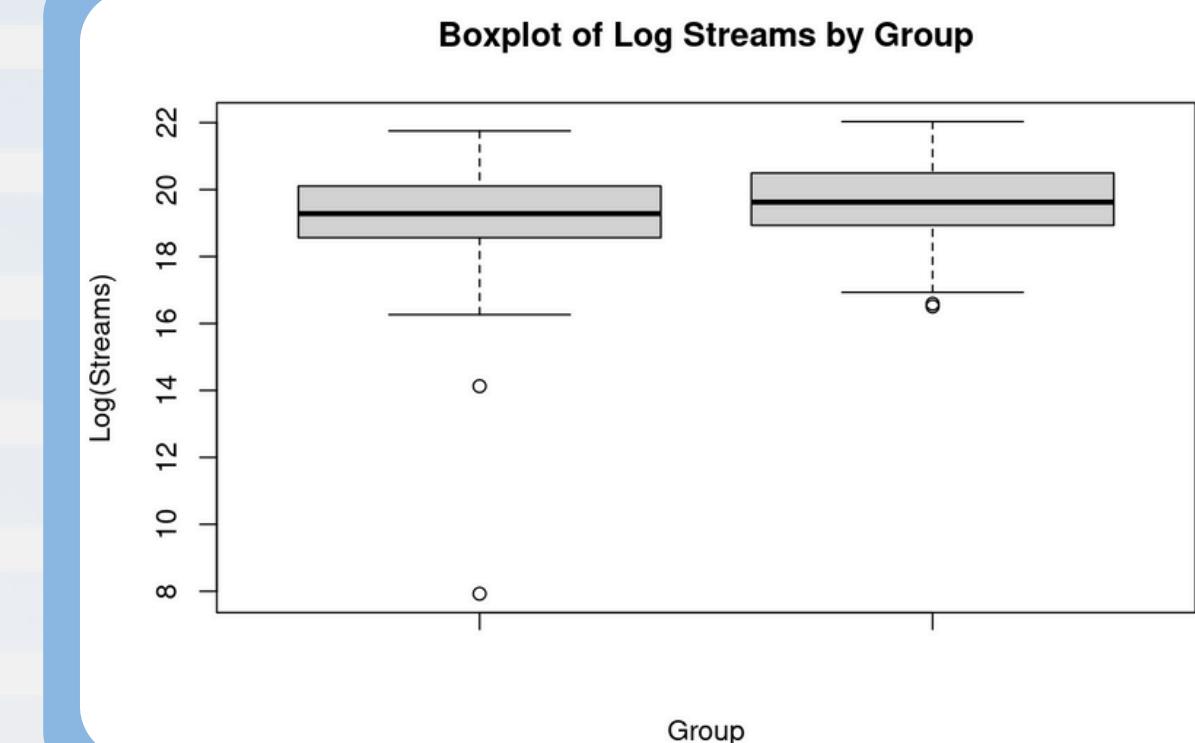
# Mean Popularity(t-test on log-streams)

- **Method:** Welch two-sample t-test
- $H_0: \mu_{\text{collab}} - \mu_{\text{solo}} \leq 0$
- $H_1: \mu_{\text{collab}} - \mu_{\text{solo}} > 0$
- **Key results:**
- Mean log-streams:
  - ◆ collab: 19.29
  - ◆ solo: 19.64
- p-value: 1
- **Conclusion:** Fail to reject  $H_0$   
(no evidence collab songs are higher)

```
> x <- df2$log_streams[df2$group == 'collab']
> y <- df2$log_streams[df2$group == 'solo']
> t.test(x, y, alternative = 'greater')
```

```
Welch Two Sample t-test

data: x and y
t = -4.4329, df = 696.78, p-value = 1
alternative hypothesis: true difference in means is greater than 0
95 percent confidence interval:
-0.4741073      Inf
sample estimates:
mean of x mean of y
19.29387 19.63954
```



# Charting Probability (2-proportion test)

◆ **Question:** Are collab songs more likely to chart than solo songs?

◆ **Method: 2-sample test for proportions**

- H<sub>0</sub>: p<sub>collab</sub> ≤ p<sub>solo</sub>
- H<sub>1</sub>: p<sub>collab</sub> > p<sub>solo</sub>

```
> prop.test(c(218, 330), c(366, 586), alternative = 'greater')
```

2-sample test for equality of proportions with continuity correction

◆ **Counts (charted = 1):**

- collab: 218 / 366 = 0.5956
- solo: 330 / 586 = 0.5631

```
data: c out of c218 out of 366330 out of 586  
X-squared = 0.84499, df = 1, p-value = 0.179  
alternative hypothesis: greater  
95 percent confidence interval:  
-0.02373351 1.00000000  
sample estimates:  
prop 1 prop 2  
0.5956284 0.5631399
```

◆ **Key results:**

- p-value: 0.179
- Conclusion: Fail to reject H<sub>0</sub> (no significant evidence collab chart more)

(Optional mini-table to show the 2×2 counts)

# Power + Overall Conclusion

◆ **Power goal:** detect a true 10% difference in charting rates with 80% detection chance at  $\alpha = 0.05$

◆ **Assumption:**  $p_1 = 0.60$  vs  $p_2 = 0.50$  (difference = 0.10)

**Power result (from R):**

- Required n per group: 388

```
> power.prop.test(p1 = 0.6, p2 = 0.5, power = 0.80, sig.level = 0.05)
```

Two-sample comparison of proportions power calculation

```
n = 387.3385  
p1 = 0.6  
p2 = 0.5  
sig.level = 0.05  
power = 0.8  
alternative = two.sided
```

◆ **Compare to dataset:**

- solo: 586 (enough)
- collab: 366 (slightly below required)

NOTE: n is number in \*each\* group

◆ **Overall conclusion:**

- Collab songs did not show higher mean log-streams than solo songs
- Collab songs had slightly higher chart rate, but not statistically significant
- Dataset may be slightly underpowered to detect a 10% chart-rate difference for collabs

◆ **Limitations (1–2 bullets):**

- Observational data, not causal
- Popular-song dataset may not represent all Spotify songs

