

Finding an Optimal Balance Between Agreement and Performance in an Online Reciprocal Peer Evaluation System

The research paper investigates the trade-off between evaluation reliability and performance improvement in reciprocal peer evaluation (RPE) systems.

Key Findings and Results:

1. Reliability of Peer Evaluations:

- The study measured **effective reliability** using a formula based on Spearman-Brown's equation.
- It was found that **at least four peer raters** are necessary to achieve acceptable reliability (0.6) (*Rosenthal & Rosnow, 1991*), ensuring consistent assessments.
- More raters improve reliability but with **diminishing returns** after a certain point.

2. Outliers in Peer Evaluations:

- Outliers (ratings significantly different from the rest) were examined.
- The occurrence of **false outliers decreased** as the number of raters increased.
- With **six or more raters**, the number of false outliers was minimal, closely aligning with the true outlier rate.

3. Impact on Performance Improvement:

- Students revised their documents based on peer feedback.
- The study found an **inverted U-shaped relationship** between the number of raters and performance improvement.
- Performance improved up to **six raters**, but beyond this point, additional feedback led to **cognitive overload** and diminished benefits.
- The optimal number of peer raters for **maximizing student performance was found to be around six**.

4. Quality and Usefulness of Feedback:

- More raters provided more feedback, but students were selective in the feedback they considered useful.
- The study **rejected the assumption that more feedback automatically improves learning**.
- Students tended to ignore excess feedback when they received too many comments from multiple raters.

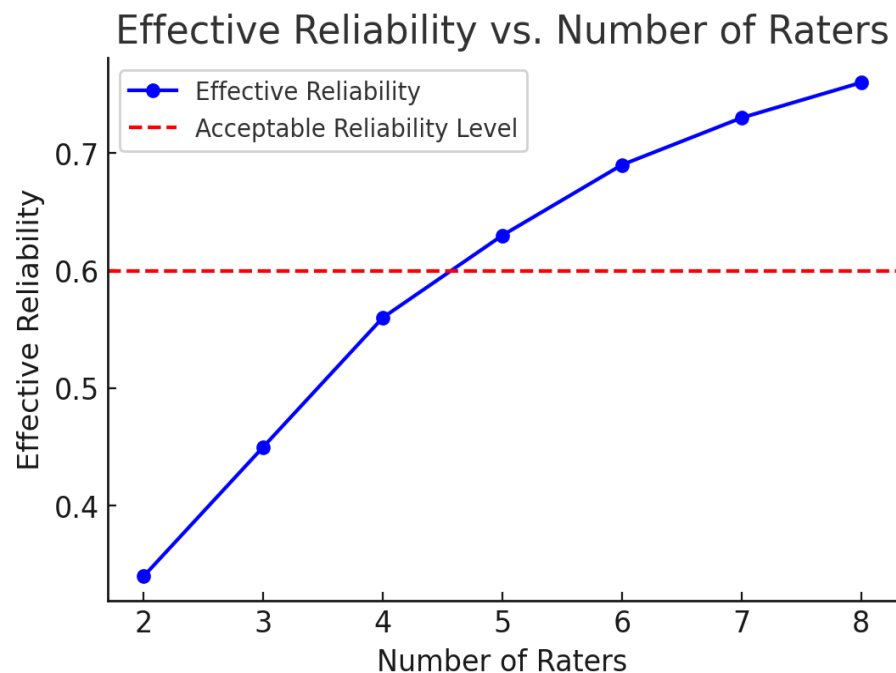
5. Time Investment in Evaluation:

- Students spent about **one hour per document** (reading + providing evaluation).
- With six raters, **12 hours of total reviewing time was required**.
- Beyond six raters, the additional time commitment yielded little benefit in improving reliability or performance.

Conclusion:

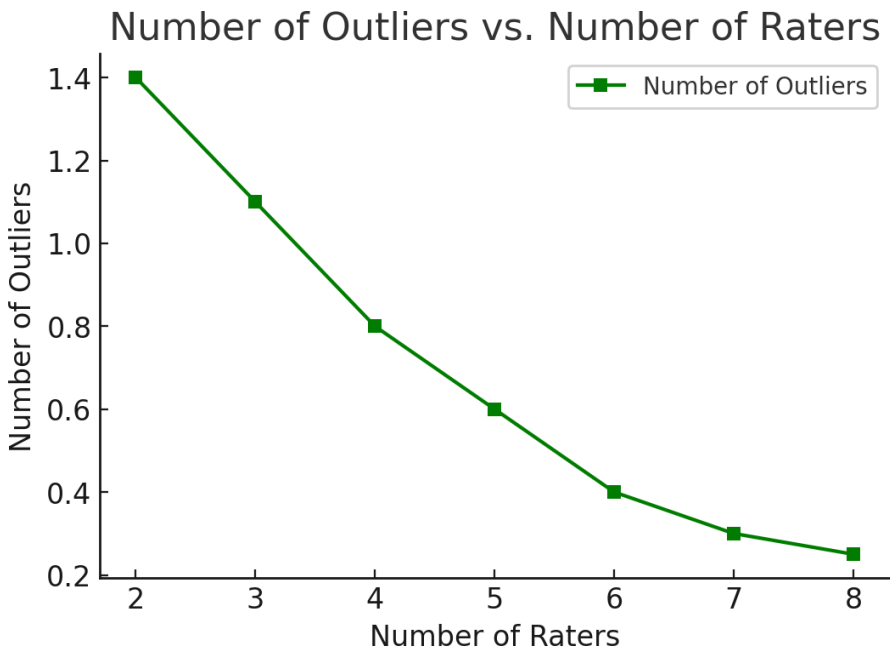
- The **maxima strategy** (more raters improve assessment quality) **holds true for reliability but not for performance**.
- **Five to six peer raters** is the **optimal balance** between reliability and student performance.
- Increasing raters beyond this number **does not enhance learning outcomes** and can even negatively impact performance due to **cognitive overload**.
- Peer evaluation systems should **balance the number of raters** to avoid unnecessary time investment and information overload.

Results in the form of Graphs:



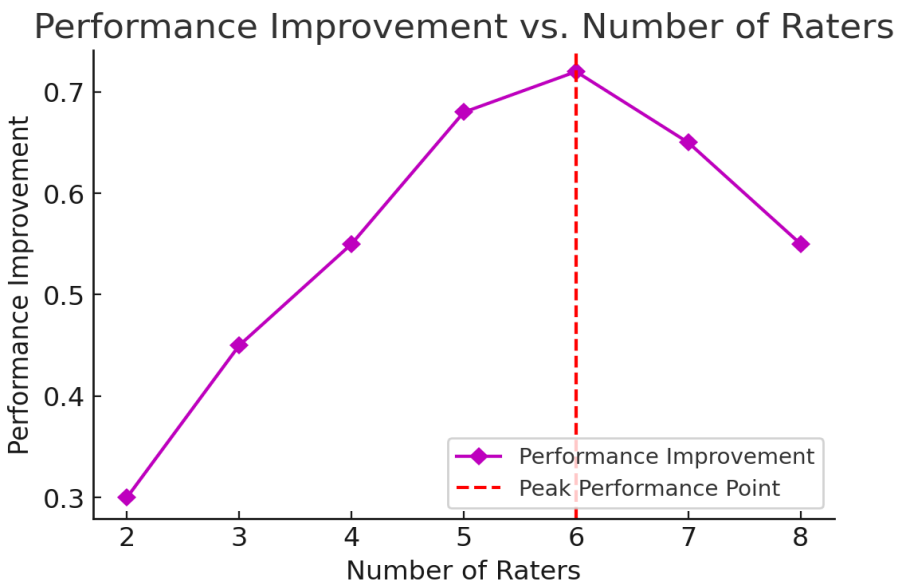
Effective Reliability vs. Number of Raters

- Shows that reliability increases with more raters but plateaus after 4-6 raters.
- The red dashed line indicates the **acceptable reliability level (0.6)**.



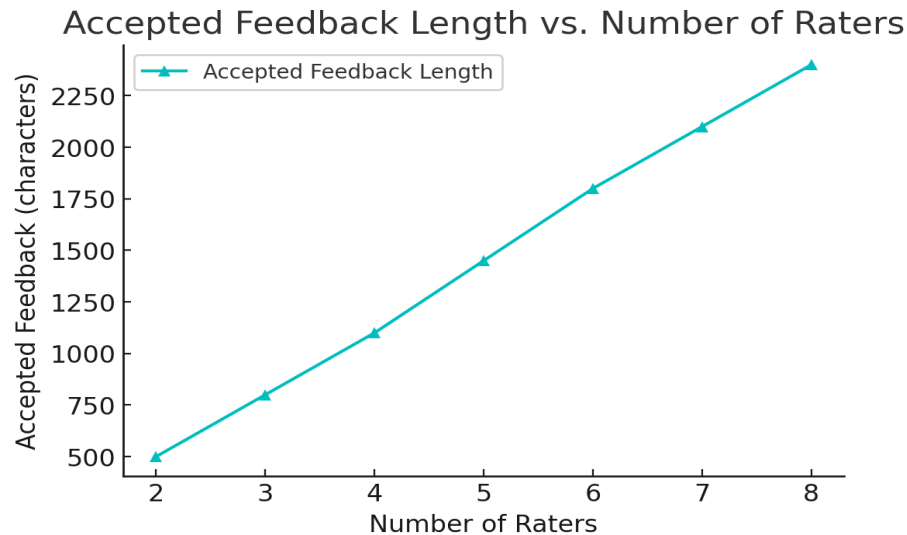
Number of Outliers vs. Number of Raters

- Demonstrates that false outliers **decrease exponentially** as the number of raters increases.
- At **6+ raters**, the number of false outliers becomes minimal.



Performance Improvement vs. Number of Raters

- Follows an **inverted U-shape**: Performance improves up to **6 raters** and then declines due to cognitive overload.
- The red dashed line marks the **peak performance point at 6 raters**.



Accepted Feedback Length vs. Number of Raters

- Indicates that **students continue to accept more feedback** as raters increase, contradicting the assumption that they ignore extra feedback.

Key Insights from the Graphs

- **Reliability increases with more raters** but plateaus at **4-6 raters**.
- **False outliers decrease** as raters increase, stabilizing at **6+ raters**.
- **Performance follows an inverted U-shape**, peaking at **6 raters** before declining.
- **Feedback acceptance increases** but does not explain performance drop beyond 6 raters.
- **6 peer raters appear to be the optimal balance** for maximizing reliability while improving student performance.