

### **Limitations: Module 3**

Analytics performed on anonymously submitted datasets inherently contain several limitations. First, the completeness of the data depends entirely on individual users. Some users submit detailed information while others provide only partial entries, which can skew aggregate statistics. In this dataset, many students did not report GPA and/or GRE scores. The subset of users who choose to disclose academic metrics may represent a non-random population — potentially higher-performing students — introducing selection bias. For example, applicants with lower scores may choose not to report them, which would artificially inflate averages. Additionally, user-generated entries are not standardized. Some applicants placed start-term information in free-text comments rather than designated fields, requiring special parsing logic (implemented in the data-cleaning code) to interpret these edge cases. Variations in terminology (such as “Physics Phd” versus “Physics PhD”) further complicate consistent classification unless explicitly handled by cleaning logic. Additionally, automated data capture introduces the possibility of parsing or formatting errors when scraping large datasets. While such system-level inaccuracies can occur, their impact is likely negligible compared to behavioral and reporting biases, particularly given the large sample size.

The GRE statistics highlight another limitation of voluntary anonymous reporting. The average GRE quantitative score observed in this dataset is higher than national published averages. This discrepancy is likely due to selection bias: individuals who post to forums like Grad Café may be more confident in their academic performance and therefore more willing to disclose their scores. Conversely, lower-scoring applicants may be underrepresented because they are less likely to post scores — or to post at all. Participation bias may also occur, where applicants who receive interviews or acceptances are more motivated to share their outcomes, while less successful applicants remain silent. Additionally, because entries are self-reported and unverified, inaccuracies — whether intentional or accidental — cannot be ruled out. In contrast, official GRE and GPA statistics are collected under controlled, standardized conditions. Grad Café data is voluntary, unverified, and inconsistently structured, meaning it reflects a biased subset of applicants rather than a true population sample.

Therefore, analytics derived from anonymous datasets should be interpreted cautiously, as they may reflect behavioral patterns of contributors rather than true population characteristics.