Using any self-reported data source comes with inherent data quality limitations due to the lack of quality control and standardization enforcement. Text fields where the user may type their own responses are prone to spelling errors, abbreviations, grammar mistakes, and inconsistent capitalization. They are also not guaranteed to be factually true, since users may input whatever information they want at their discretion. Other user-fillable fields are prone to similar errors; any user on TheGradCafe may embellish or mistype their GRA and GRE scores, for example. Aside from errors, the lack of standardization makes the data inherently more challenging to analyze since values that mean the same thing to a human are treated as distinct values when read by a computer. This can be mitigated by using fields where the user must select from pre-filled options, but whether this is the case is completely up to the creator of the website.

Another limitation of using self-reported data sources for any analytical purpose is the inherent selection bias in the dataset. For example, graduate school applicants who are accepted into their program of choice may be more motivated to share their results with the world than those who were rejected. These kinds of biases will be present in any dataset reliant on users choosing to report data for themselves and will necessarily lead to results based on analysis from a different distribution than the 'true' distribution (were the dataset to be populated without any biases). There may be other confounding factors as well, and some may be unique to different data sources. If one is trying to create a regression model that predicts the chance of acceptance based on GRE scores using only the latest one year of data, then the model's future predictive ability could certainly diminish in years where the GRE is 'harder' or 'easier' to most students, since the exam itself changes annually. Public policy may also influence these datasets, since a change in student visa policies can lead to more/fewer international student applications which also skews the distribution of GRE and GPA scores.