**Abstract**: Microsimulations and synthetic data are often high-dimensional, requiring extensive validation and exploration. This article illustrates data visualization techniques for comparing a generated sample or population against a known sample or population. For implementation ease, we also outline an iterative workflow built with R Markdown that can be shared publicly on GitHub or privately with Amazon AWS s3. The lessons learned from this work apply to any analysis that compares data sets, deals with high-dimensional data, and/or involves summarizing several iterations of analyses.    We organize the paper as follows. Section 2 discusses data visualization techniques for directly comparing observations from two high-dimensional data sets. To deal with many observations, we use transparency, highlighting of important trends or observations, and visual summary techniques for univariate and multi-variate comparisons (e.g. one-dimensional and multi-dimensional binning, density estimation, and LOESS). To deal with many variables we use techniques for showing multiple relationships (e.g. color, shape, and faceting) and iteration. Section 3 discusses data visualization techniques for analytical comparisons of high-dimensional data sets. We visualize high-dimensional correlation matrices, many counts, many of the first four moments of variables common between data sets, Kolmogorov-Smirnov tests, and regression confidence intervals estimated on both data sets. Section 4 outlines an iterative workflow with R Markdown, GitHub pages, and AWS s3 to iterate and disseminate analysis results. All of the data visualizations are implemented in the R package ggplot2 (Wickham, 2010), all examples will be reproducible, and all code will be available.