

a. Code output:

Sum: 3180.03
Mean: 6.28463
Median: 6.2085
Range: 5.219

Sum: 11401.6
Mean: 22.5328
Median: 21.2
Range: 45

Covariance: 4.49345
Correlation: 0.69536

- b. Personally, manually coding the statistical functions myself really drove home the importance of having a dedicated language like R. Obviously, the coding itself wasn't difficult at all and took very minimal amounts of time and energy. That being said, the convenience of R was simply wonderful. I think at some point in data analysis you don't quite care about the formula for correlation as much as the correlation itself, which is exactly where R shines. Overall, in a sentiment that I'm sure is shared by many, for simple statistical measures R is not *essential*, but certainly nice to have.
- c. The mean of a dataset is its unweighted average, found by summing all data and dividing by the total amount of data. The median of a dataset is the value in the exact middle of a dataset. Finally, the range is the difference between the highest and lowest values in a dataset.

The mean of a dataset is useful in giving a concise "typical value" in the data, but is especially sensitive to outliers. This may be helpful or harmful in various cases, but in the cases where it is not beneficial, the median is preferred. For example, the mean net worth between Elon Musk and the entire UTD student body is \$6 million. As a less silly example, the mean household income in the US is \$121,000, but the median is a much lower \$70,000.

Therefore, the median is sometimes used to better account for outliers. Together, both measure the "typical value" of data and allow for a limited understanding of how to evaluate and understand an otherwise large and complicated dataset.

The range of a dataset is a simple measure that indicates how spread out a dataset may be, as well as detection for possible outliers. In a dataset where values are typically in the thousands, a range in the ten-thousands would signal an outlier; a range in the tens would signal a very clumped or small dataset.

- d. Covariance is essentially the direction of relationship between variables, meaning that a positive covariance signals that the two are high at the same time and low at the same time, and vice versa for negative covariance. Correlation is a scaled covariance, which therefore measures both the direction and strength of the relationship. In the context of machine learning, the two measures are broadly applicable to many areas, particularly in early data exploration and verification of results or predictions.