

Math 611 HW6

Jeff Gould

10/8/2020

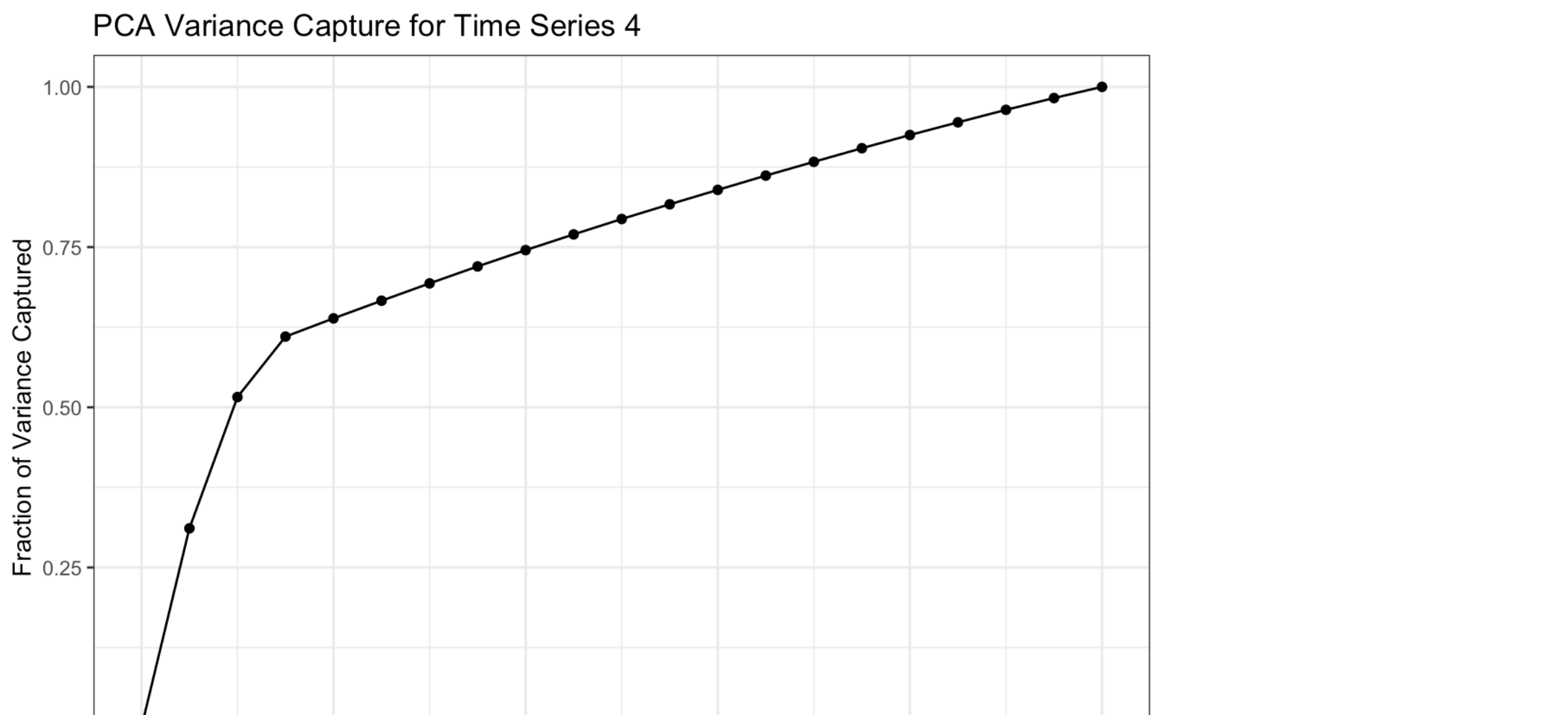
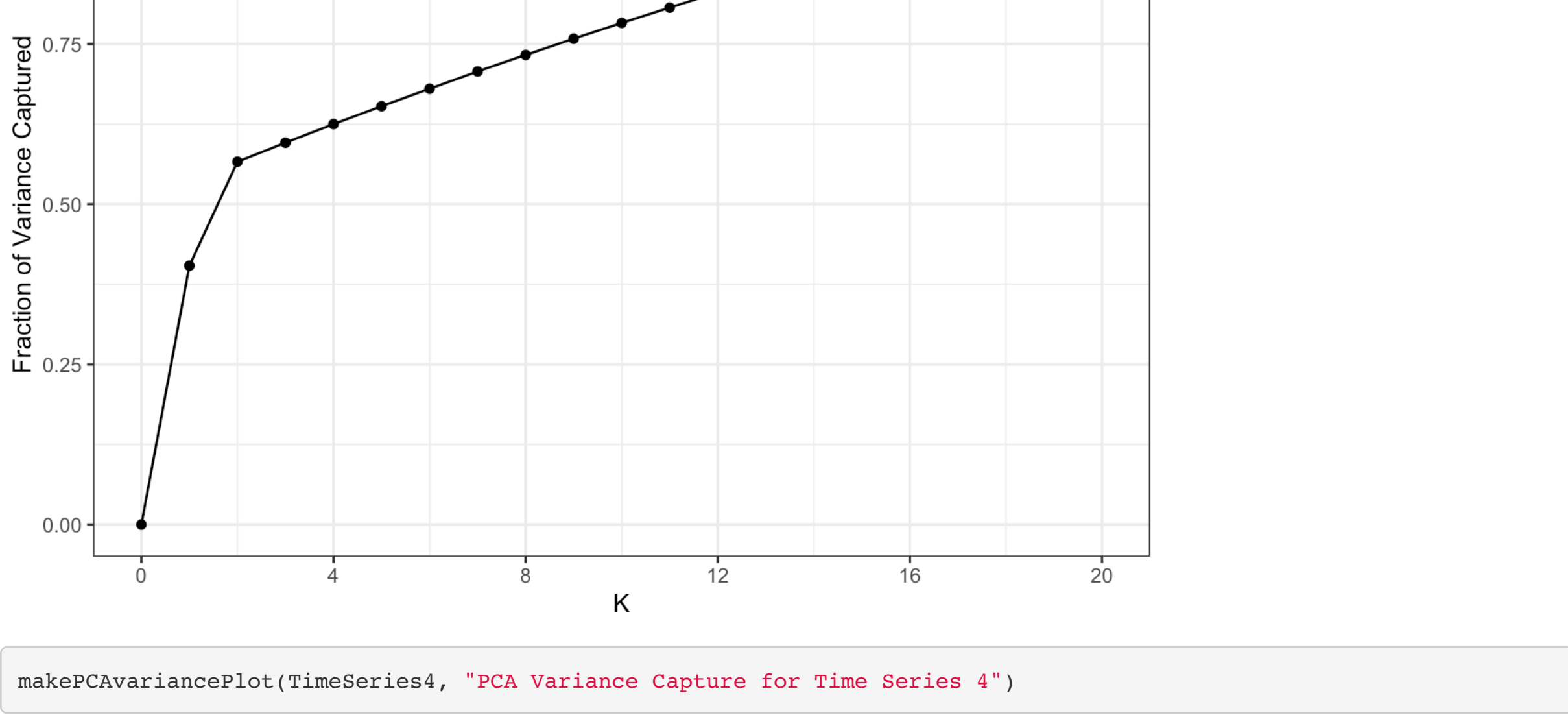
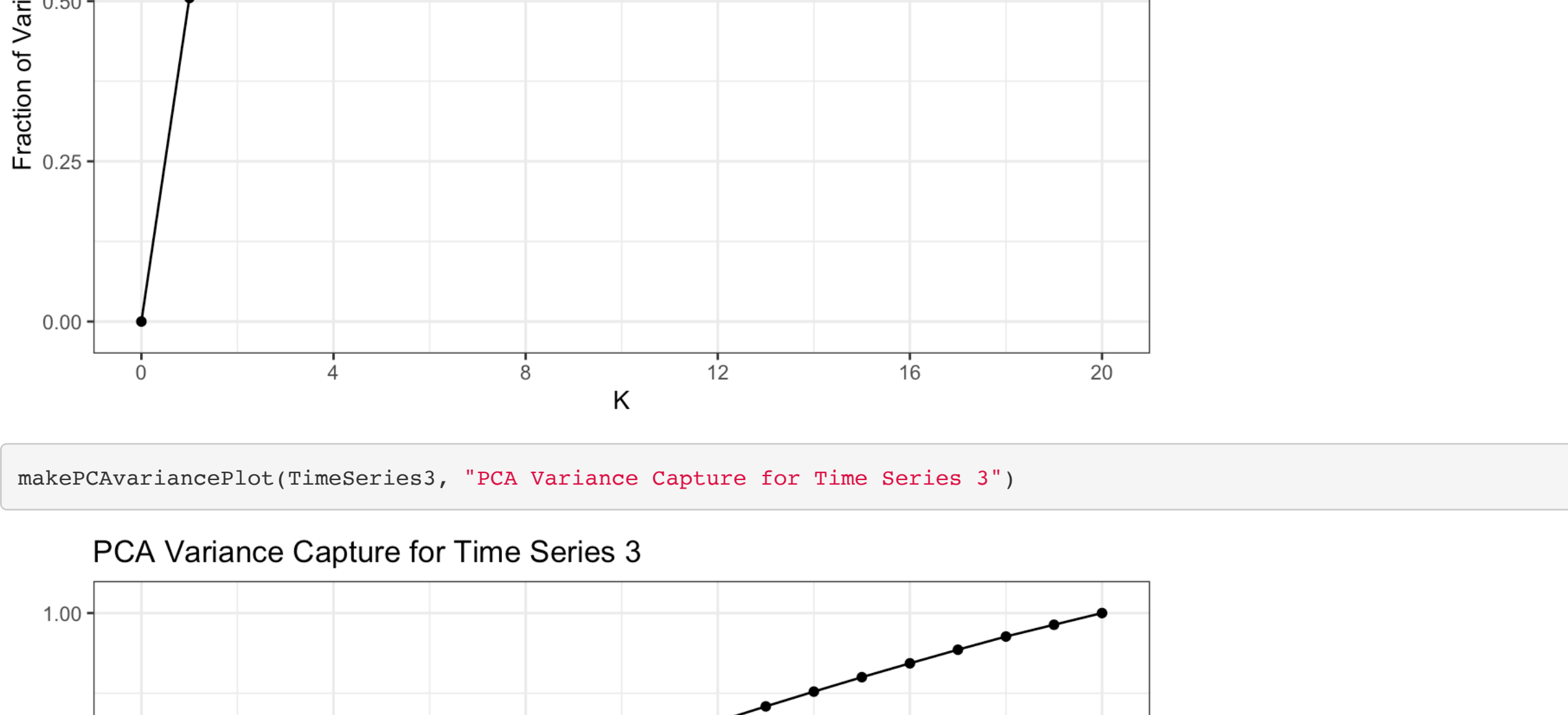
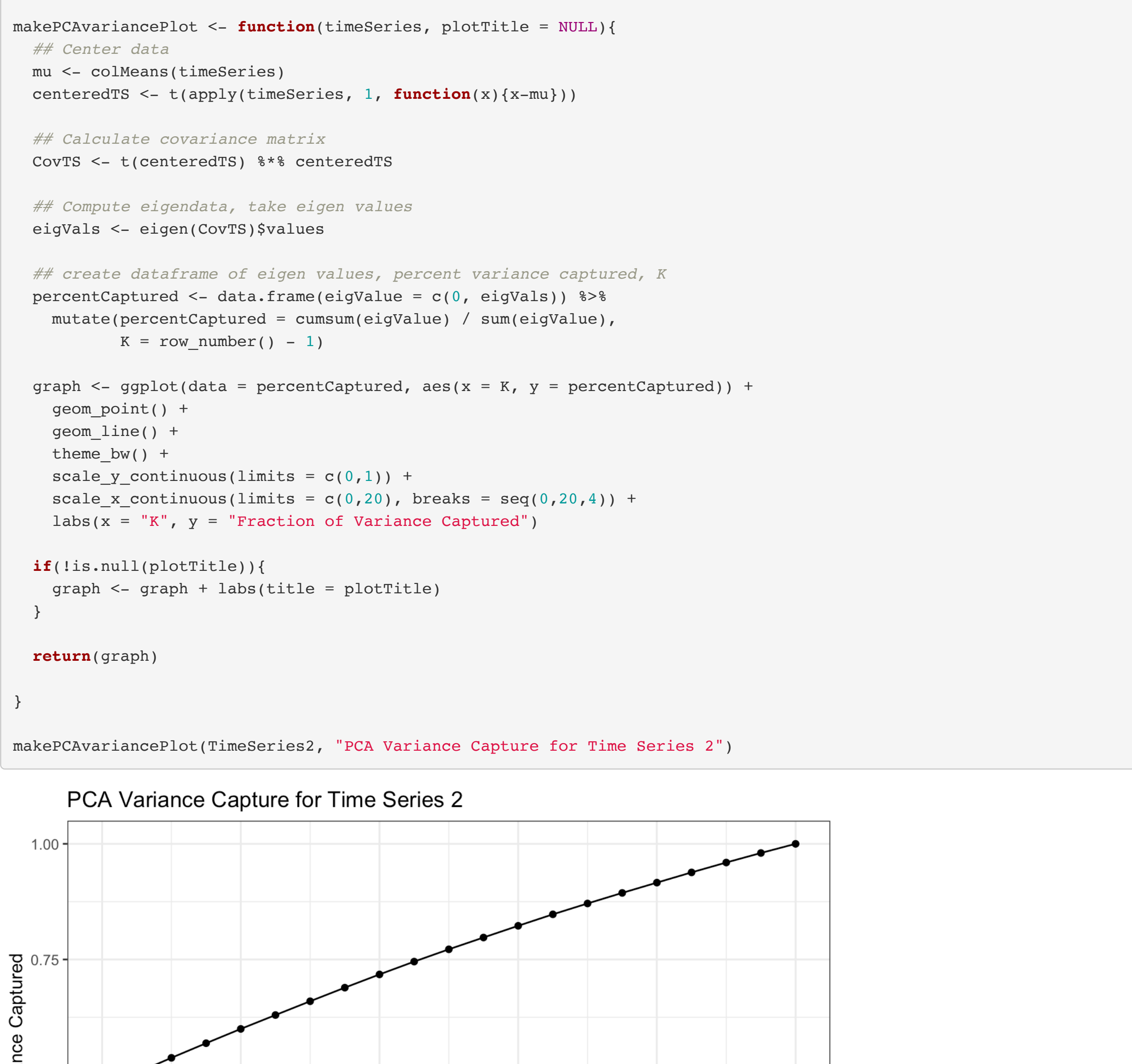
1. Time Series

a)

As we have more variables, we will need to use more dimensions of PCA in order to capture the increase in variance. For `TimeSeries2`, since there are just two base series that we need to differentiate the variation between, we should be able to use K to either 1 or 2 and capture most of the variance. As we increase the number of base series to differentiate, we will likely need to increase K by a similar amount. So if we find $K = 1$ to be sufficient for `TimeSeries2`, then $K = 2$ should be sufficient for `TimeSeries3`, and $K = 3$ for `TimeSeries4`.

b)

i

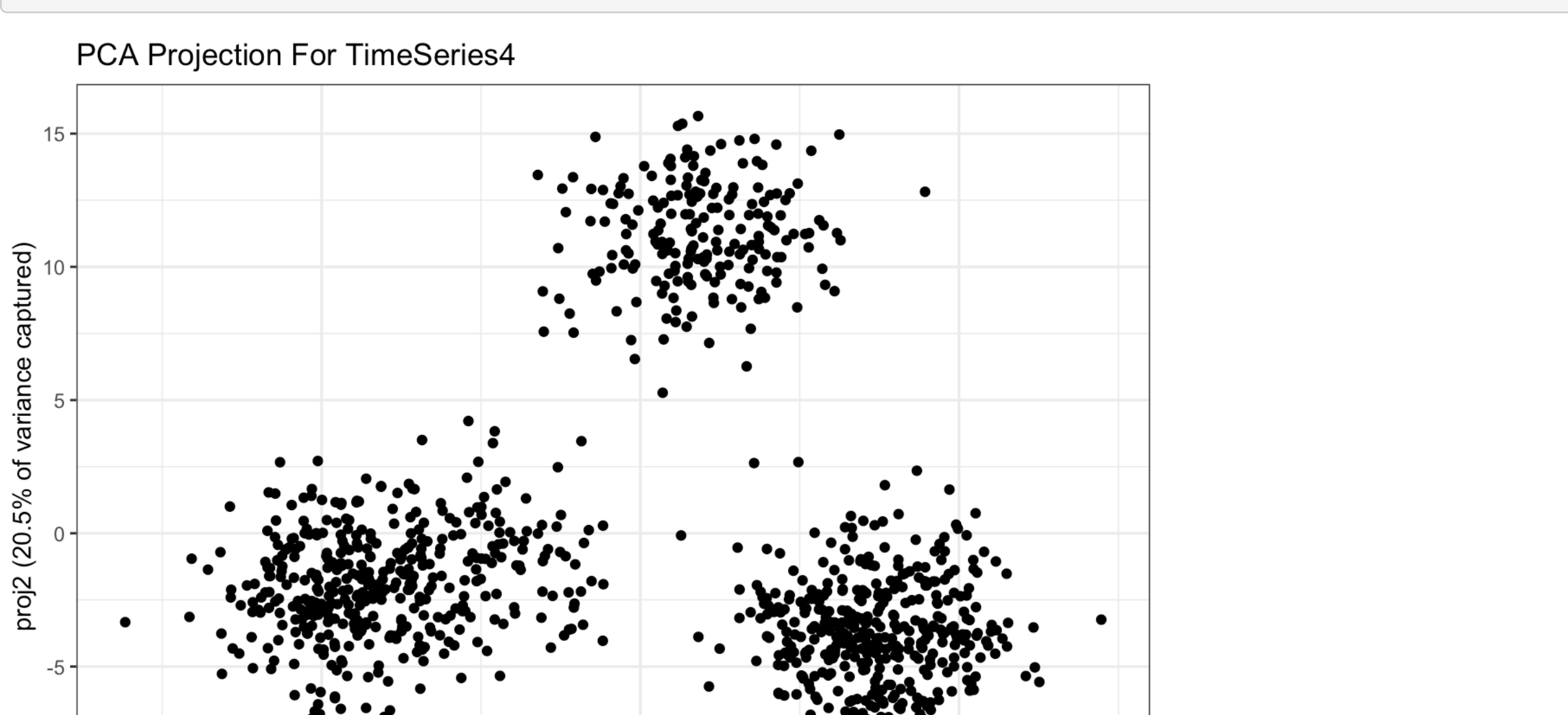
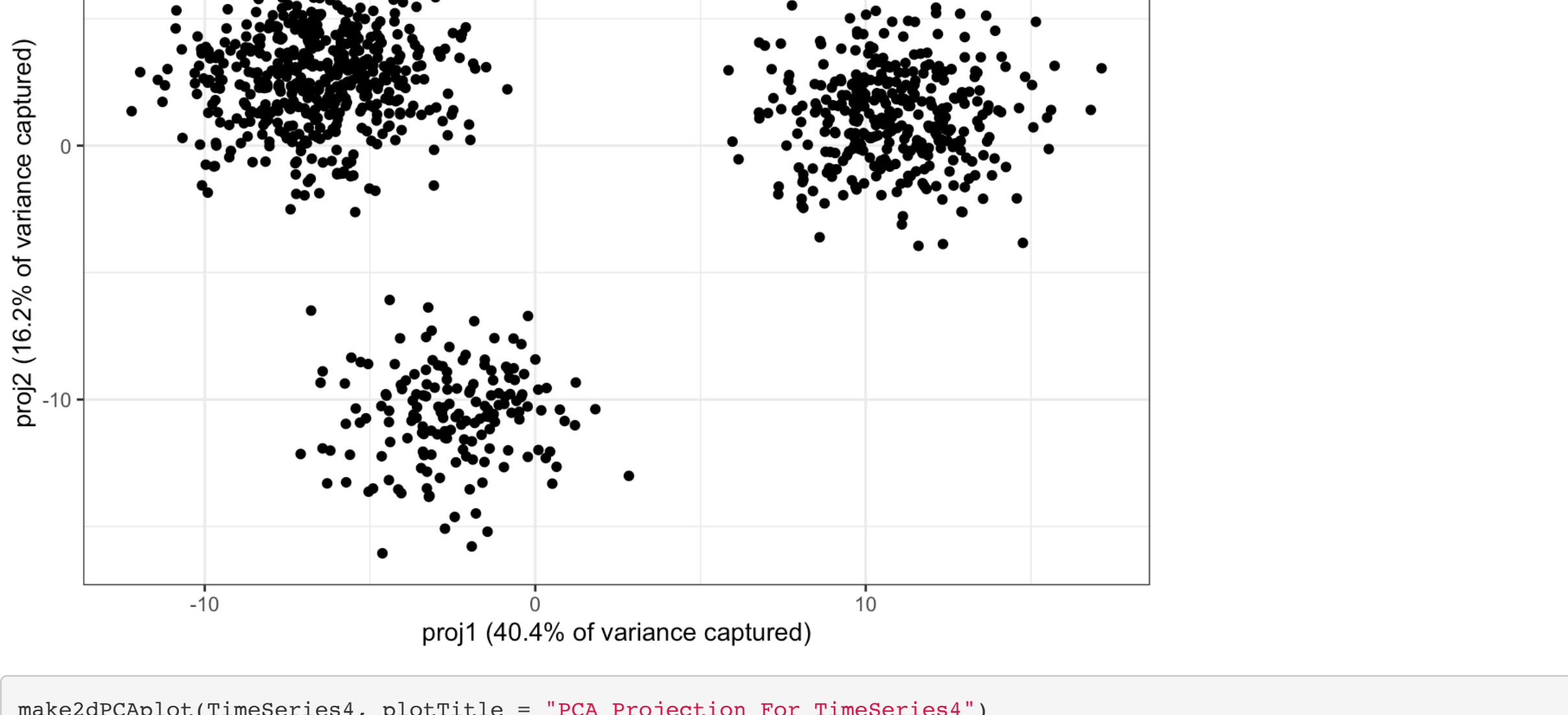
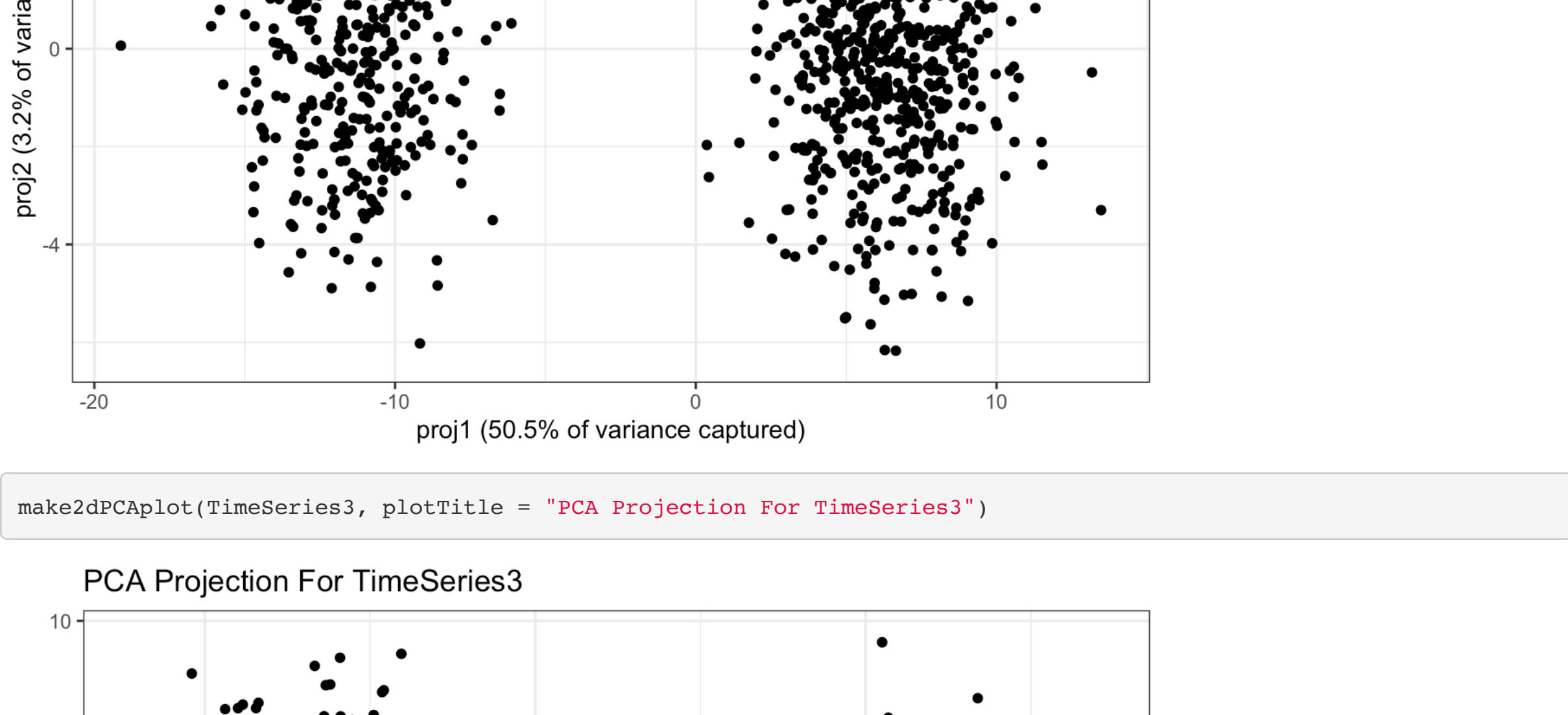


For `TimeSeries2`, we find that we are able to capture just over 50% of the variance with the first dimension of the PCA, and then gain around 2-3% for the rest of the dimensions, monotonically decreasing.

For `TimeSeries3`, we capture about 40% of the variance with the first dimension, around 57% after capturing the first two dimensions, and then around 2-3% for each additional dimension, monotonically decreasing.

For `TimeSeries4` we capture 31% of the variance after the first dimension, 52% after the second dimension, and 61% after three dimensions. Again, we then gain about 2-3% per additional dimension, monotonically decreasing.

ii

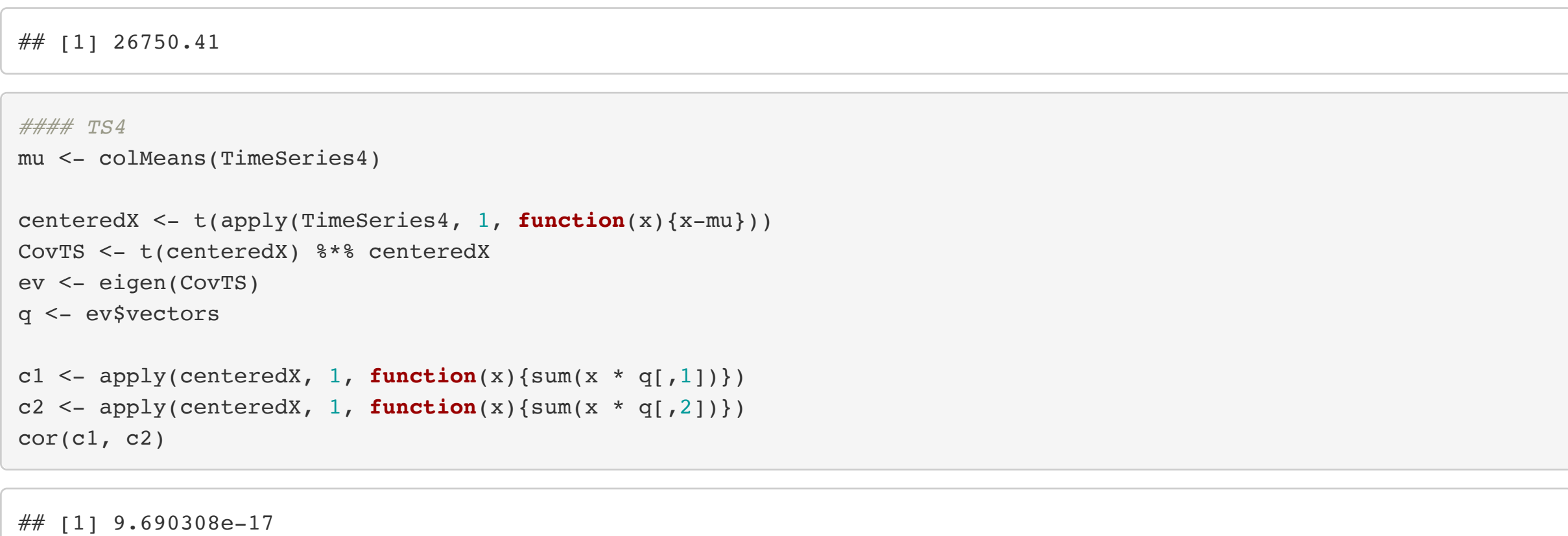
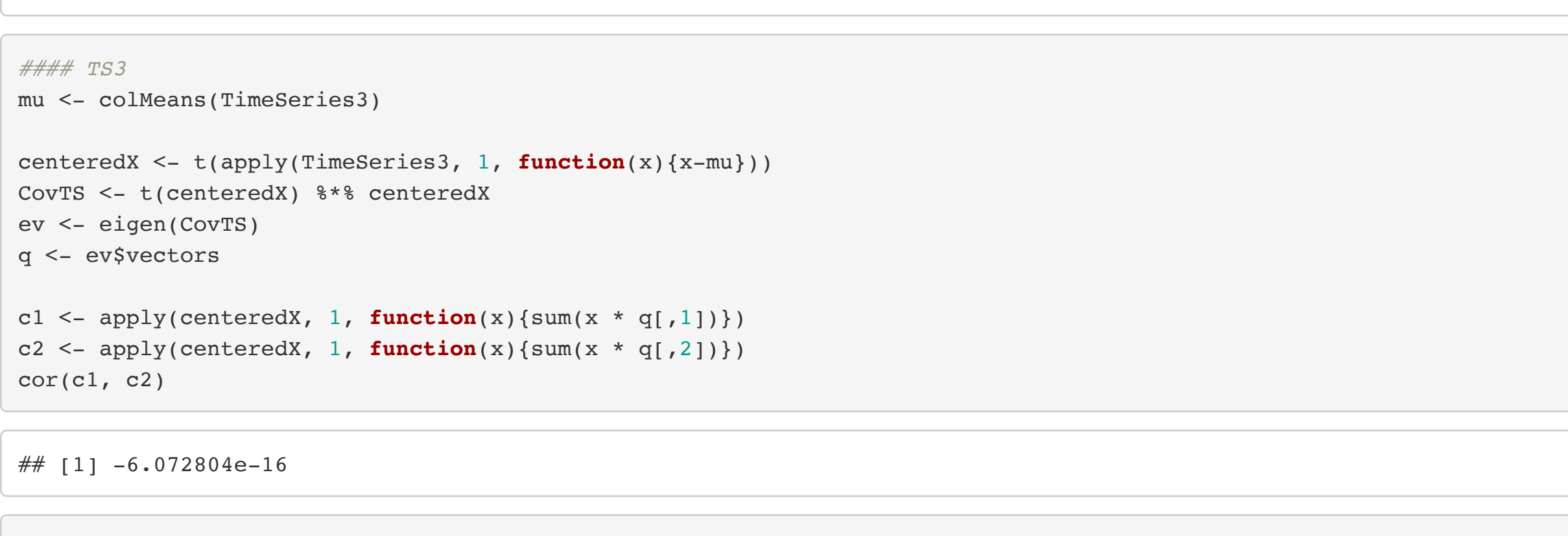
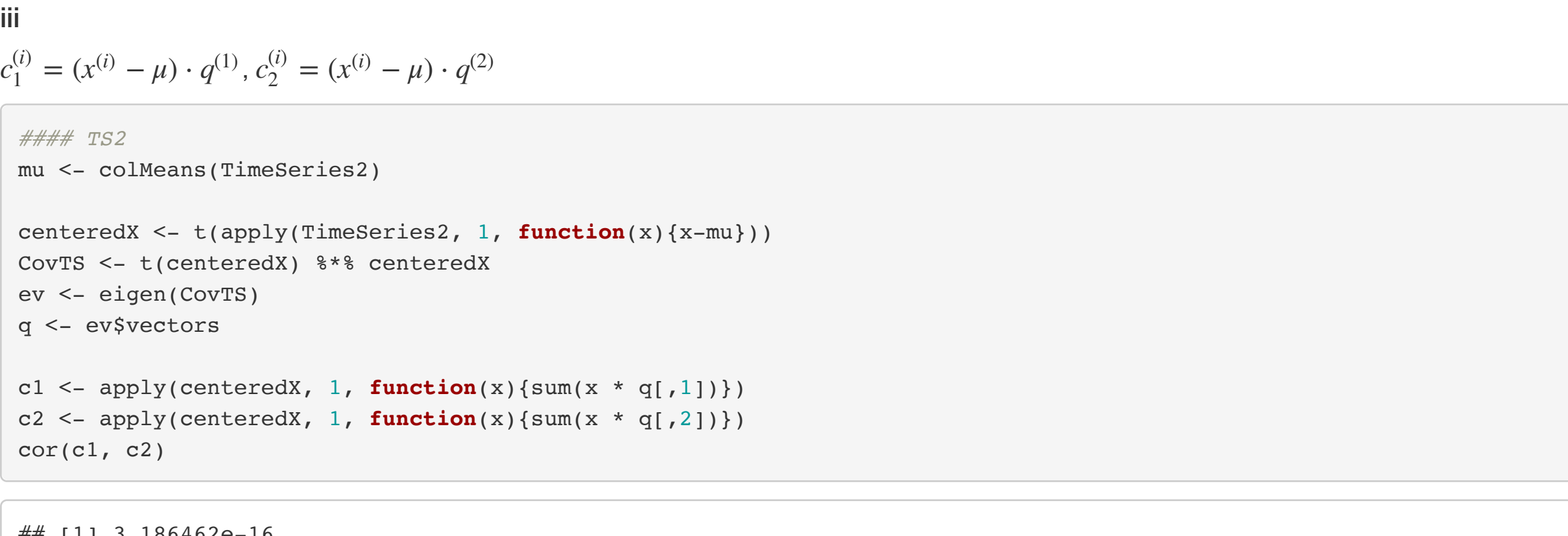


With `TimeSeries2` we see the data is clearly split into two groups, so we should easily be able to classify the data to one of the series. We also see that the differentiation is all in the first component. Adding the second dimension is not helpful in separating the series. This also makes sense, as the first dimension captures 50% of the variance, and the second component only captures 3.2%.

For `TimeSeries3`, the data is also grouped into three distinct clusters. We could also group the series into three clusters pretty easily.

In `TimeSeries4`, we only get three clusters instead of 4, and the space between the clusters is not as distinct as in `TimeSeries2` and `TimeSeries3`. It looks as though the top cluster might be one series, the bottom right cluster may be one series, and the bottom left cluster may be a blend of the two other series, but it is impossible to say for certain. We would need to expand to three dimensions in order to better separate the clusters.

iii



For each `TimeSeries`, there is no correlation between the $c_1^{(i)}, c_2^{(2)}$ s. The variance of $c_k^{(i)}$ for each time series corresponds to q^k , just scaled by 2 magnitudes. This is due to the normalization of the eigen data, so we can say that variance of $c_k^{(i)} = q^k$, before normalizing.

2. Roll 100 die



The average number of 6's rolled when the die sum to 450 is `mean(E_6) : 33.7254`

The average number of 1's rolled when the die sum to 450 is `mean(E_1) : 4.9176`

The probability that we rolled fewer than 30 1's is `mean(E_1 < 30) : 1`

Note that there is one result where we could end up with exactly 30 1's, and that would be if we also rolled 70 6's. However, this never happened in our simulation. The probability of this event is $\frac{1}{6^{100}}$ (not conditioned on the die summing to 450)