## 7.    APPLICATION 5: DZNMF

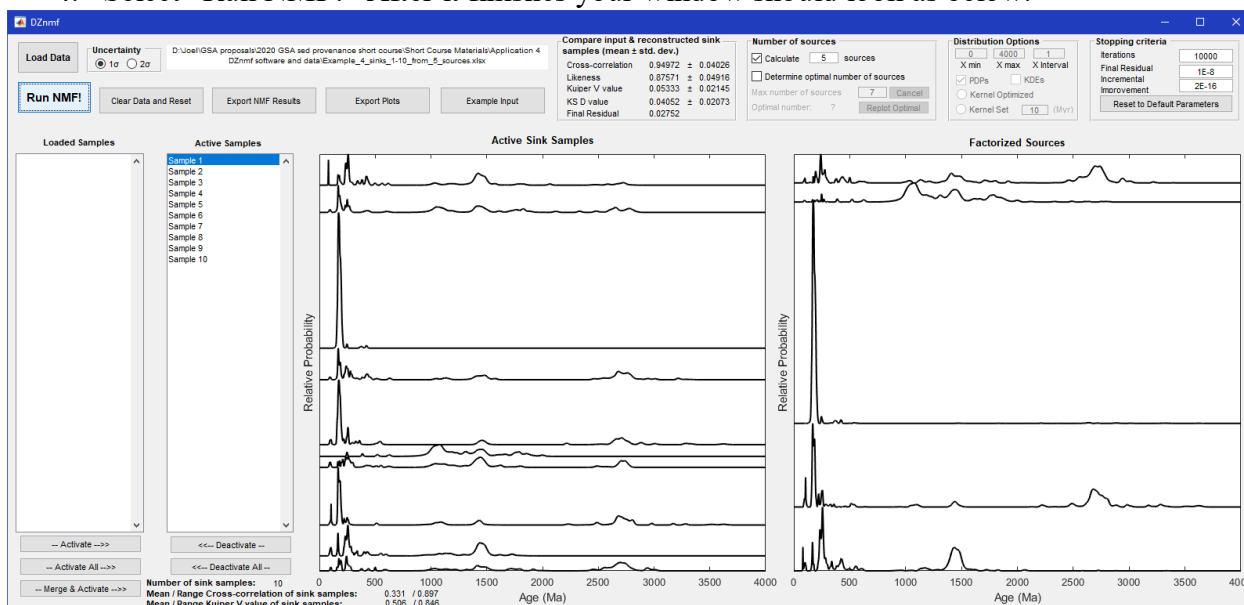Learning goals:

-How to determine the optimum number of sources for factorization using breakpoint analysis.

-How to factorize a sink sample set using DZnmf

-Recognize the effect of the number of sinks on the quality of the source factorization

-Recognize the non-uniqueness of factorization and understand methods to mitigate the likelihood of identifying a local (rather than global) solution
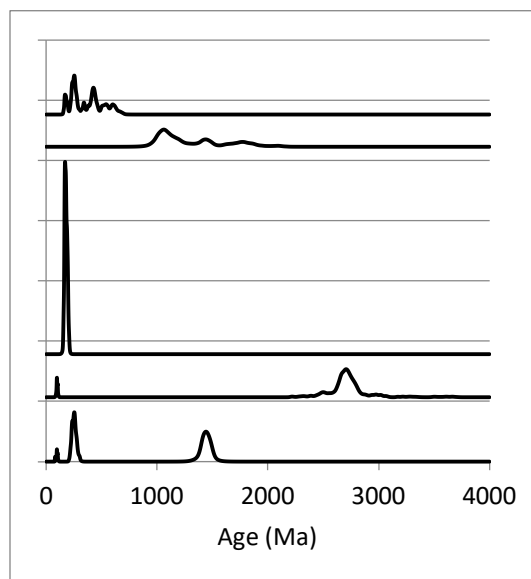
### 7.1.    Factorizing a synthetic data set

We will be using a couple of data sets that are mixtures of five source distributions.

1. Open DZnmf by double clicking the icon in in the folder where you saved the application.
2. Select "Load Data" and navigate to the file, "Application_5_sinks_1-10_from_5_sources.xlsx".  Leave the options at their defaults.
3. Select "Activate All" at the bottom left.
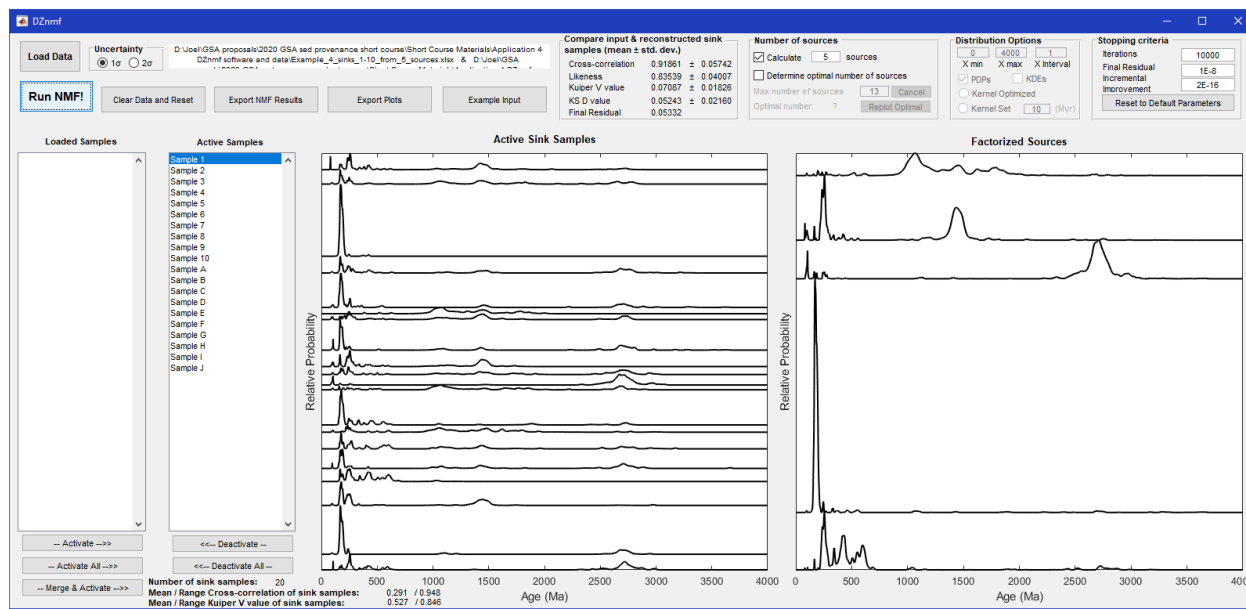4. Select "Run NMF!" After it finishes your window should look as below.

Because we created these samples we know what the original sources looked like (right). Clearly the algorithm is doing a poor job of identifying these sources. The topmost and second from bottom sources are particularly poorly separated.



### 7.2. Impact of the number of samples on factorization

1. Let's increase the number of samples that we are including in the analysis. Without clearing the data from the previous step, select "Load Data" and navigate to the file, "Application_5_sinks_A-J_from_5_sources.xlsx". Leave the options at their defaults.
2. Select "Activate All" at the bottom left.
3. Note that we now have 20 active samples (samples 1–10 and A–J). DZnmf can merge data sets on the fly to maximize the ability to explore data sets. You can also select multiple samples in the "Loaded Samples" list and merge them into one sample in the "Active Samples" list by highlighting them, and selecting "Merge & Activate" at the bottom left.
4. Select "Run NMF!" After it finishes your window should look as below.

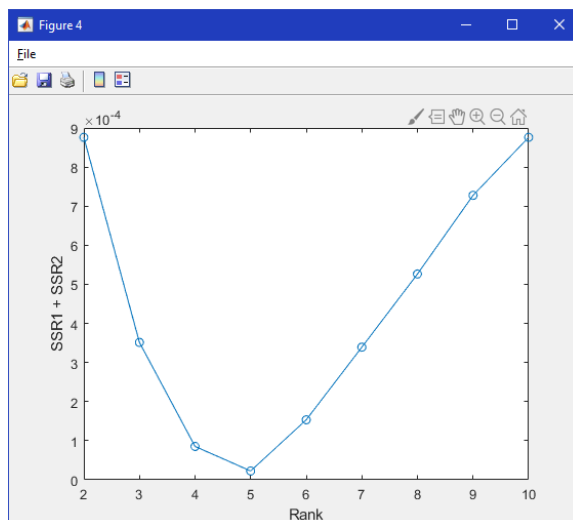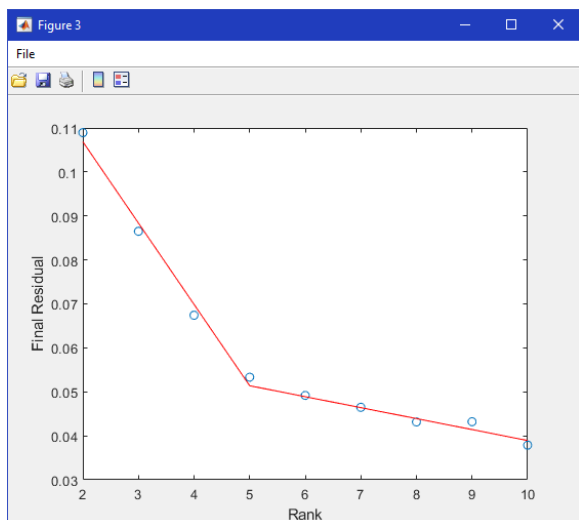**Quantitative analysis, visualization, and modelling of detrital geochronology data**



5. Comparing the new results to the known age distributions above shows two things. First, the **order** of the results does not matter. Secondly, and more importantly, adding more samples resulted in much more discriminatory power and a closer match between the known source distributions and the factorized distributions.

## 7.3. Determining the optimum number of sources
This would usually be the first step in an NMF analysis.
1. Without clearing the data, select "Determine optimal number of sources".
2. Set the "Max number of sources" to 10.
3. Select "Run NMF!" again. The algorithm will cycle through a factorization of all ranks from 2–10 and calculate the final residual. It should take about 10 minutes.
4. When it finishes, two windows will open, as shown below. The plot on the left shows the Final Residual for each rank between 2 and 10 as blue data points and the optimized segmented linear regression as a red line. The plot on the right shows the sum of summed squared residuals (SSSR) based on the segmented linear regression. The optimum rank is the one that minimizes the SSSR. Note that the calculated optimum rank (5) coincides with the known number of sources.
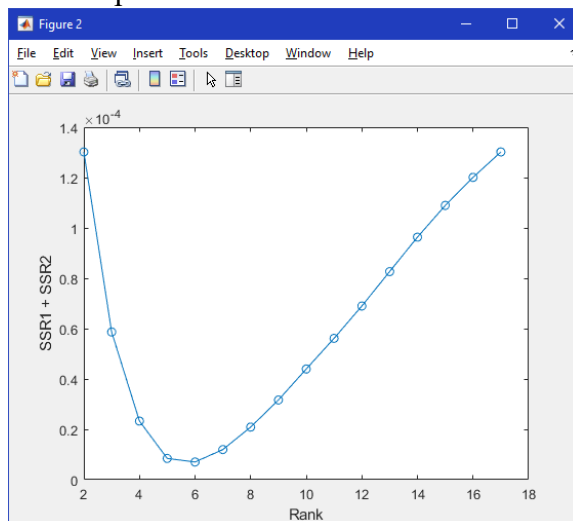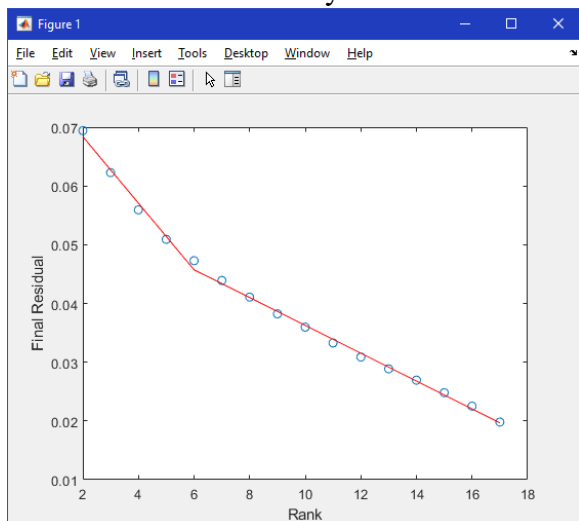
**Quantitative analysis, visualization, and modelling of detrital geochronology data**



## 7.4.    Non-negative matrix factorization of an empirical data set

We will be using a couple of empirical data sets from Neoproterozoic–Triassic strata from western Laurentia (Gehrels et al., 2011; Gehrels and Pecha, 2014) to explore some of the complexities of using DZnmf.
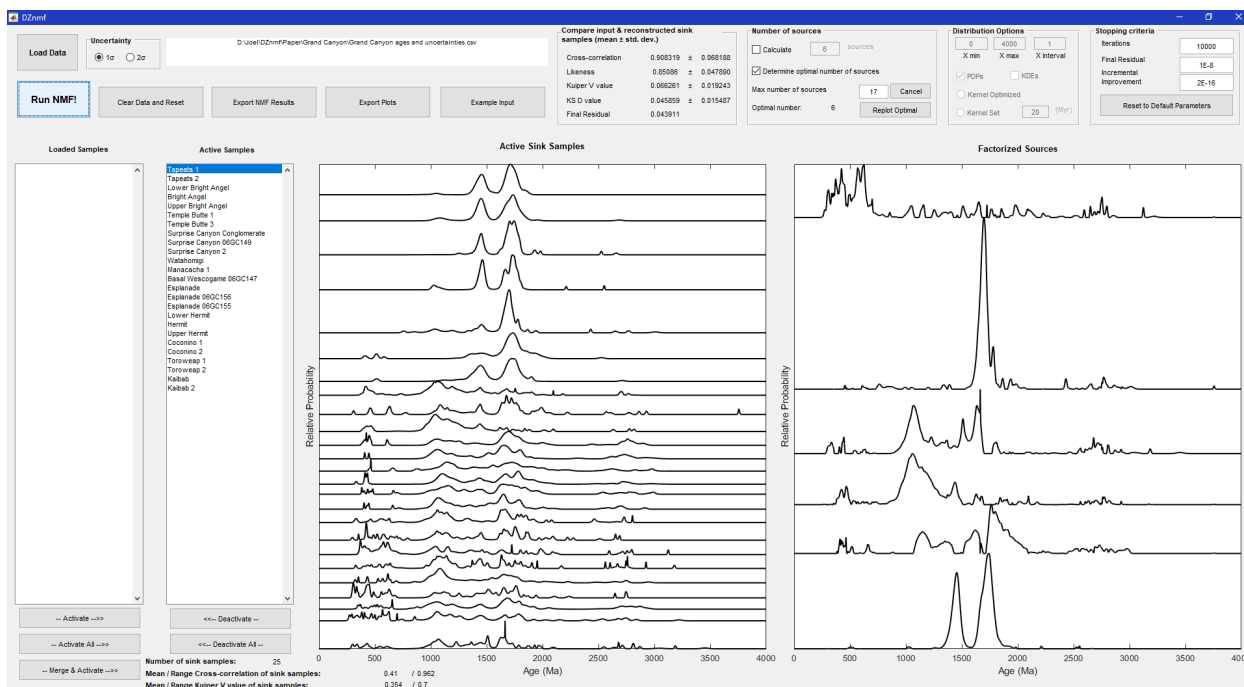
1. Select "Load Data" and navigate to the file, "Application_5_Gehrels_et_al_2011_Grand_Canyon_U-Pb.xlsx".  Leave the options at their defaults. Uncertainty is reported as "1sigma", and we will visualize the samples as PDPs.
2. Select "Activate All" at the bottom left.
3. Usually the first step would be to factorize the data set to a wide range of ranks to determine the optimum rank. Because this process takes about 20 minutes, I have done it already and show the results below. The optimum rank is 6.
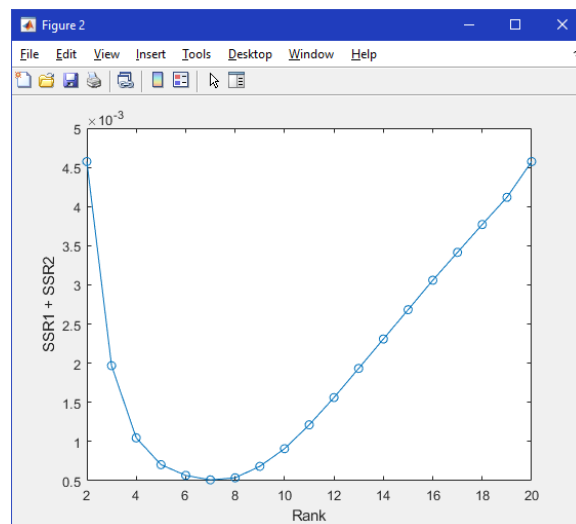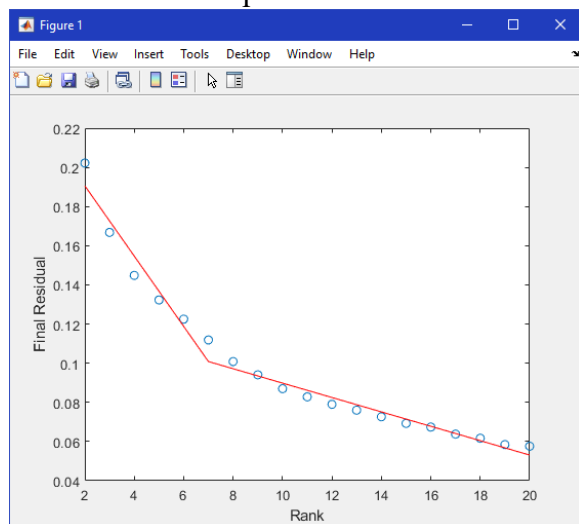


4. Under "Number of sources" select the checkbox for "Calculate" and input "6" in the sources box.
5. Select "Run NMF!"

6. Once the factorization is complete, your DZnmf window should be as below. Note that the order of the sources may not match this figure. Compare these factorized sources to the empirical sources reported by Saylor et al. (2019).
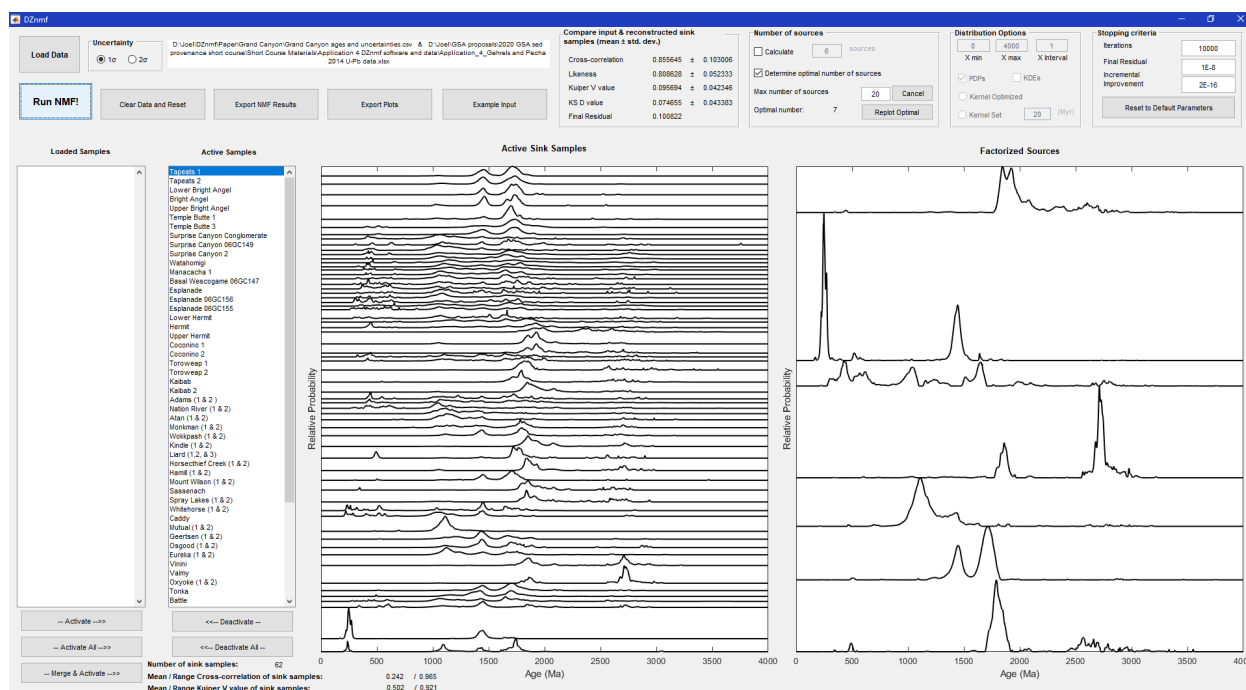


7. We will take one more step: adding in additional data and repeating the factorization.
8. Select "Deactivate All"
9. Select "Load Data" and navigate to the file, "Application_5_Gehrels_and_Pecha_2014_W_Laurentia_U-Pb.xlsx". Leave the options at their defaults. Uncertainty is reported as "1sigma", and we will visualize the samples as PDPs.



10. As above, we would usually factorize the data set to a wide range of ranks to determine the optimum rank. I have done it already and show the results above. The optimum rank is 7.
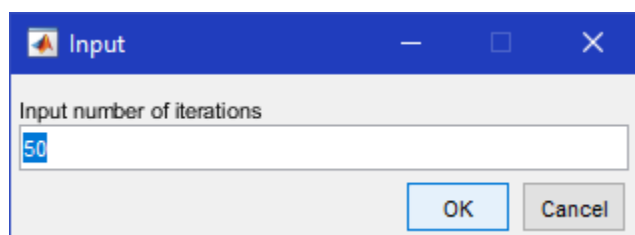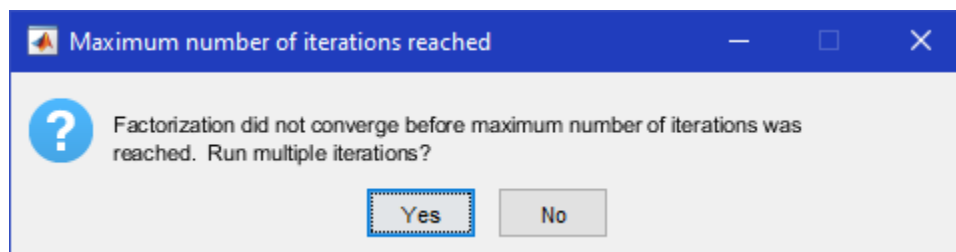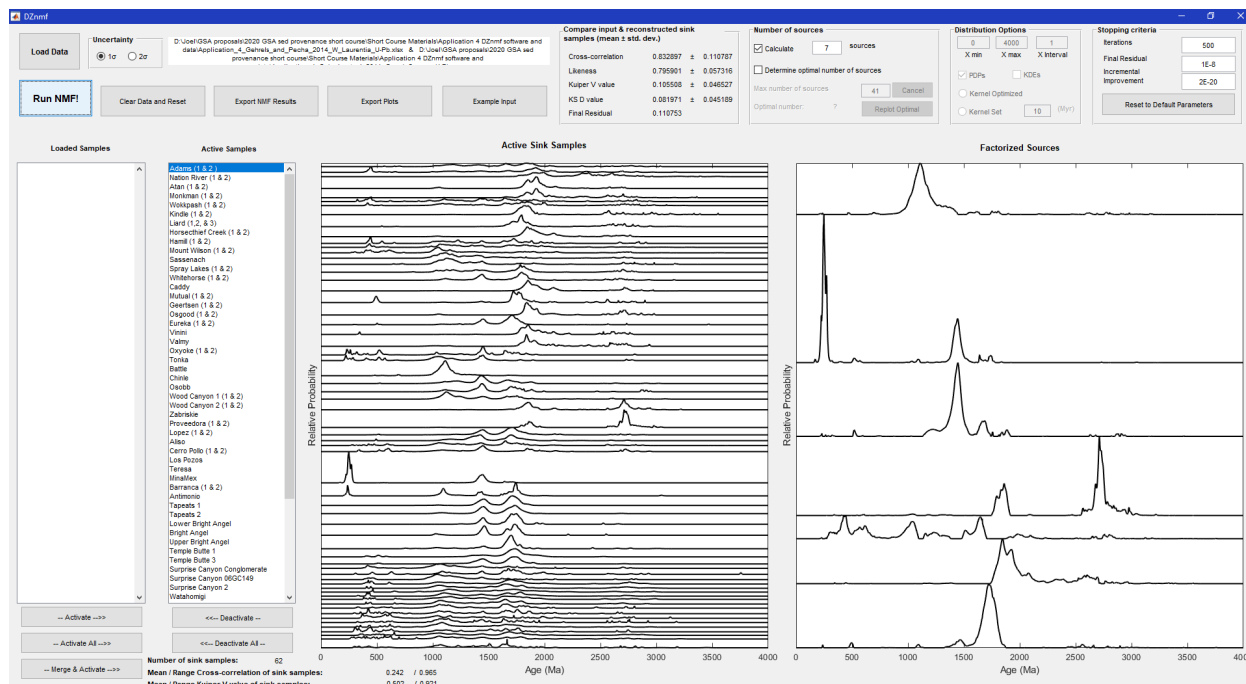
11. Under "Number of sources" select the checkbox for "Calculate" and input "7" in the sources box.
12. Select "Run NMF!"
13. Once the factorization is complete, your DZnmf window should be as below. Note that the order of the sources may not match this figure. Compare these results to the results from step 6. From bottom (1) to top (7):
    a. Step 13 Source 1 => Step 6 Source 5
    b. Step 13 Source 2 => Step 6 Source 1
    c. Step 13 Source 3 => Step 6 Source 3
    d. Step 13 Source 4 => Novel
    e. Step 13 Source 5 => Step 6 Sources 4 & 6
    f. Step 13 Source 6 => Novel
    g. Step 13 Source 7 => Step 6 Source 2?



14. One final point. I ran the same data set through a factorization with only 500 iterations and got the result below. Notice that the age distributions are different from above. Notice also that the factorization below has a higher final residual and lower mean cross-correlation coefficient. NMF is non-unique. We have tried to account for this by including the possibility of running multiple trials to maximize the probability of finding a global minimum in final residual.
15. In order to use this feature, set both the final residual criteria and incremental improvement criteria to extremely low values (i.e., 1E-20). The algorithm will factorize the data set but if it reaches the target iterations before it reaches the ending

**Quantitative analysis, visualization, and modelling of detrital geochronology data**

criteria, it will ask if you want to repeat the factorization (see below). After running the prescribed number of factorizations it will select the one with the lowest final residual. We recommend repeating the factorization at least 50 times prior to accepting the results.

**Quantitative analysis, visualization, and modelling of detrital geochronology data**

## Works Cited

Gehrels, G., Pecha, M., 2014. Detrital zircon U-Pb geochronology and Hf isotope geochemistry of Paleozoic and Triassic passive margin strata of western North America. Geosphere 10, 49–65.

Gehrels, G.E., Blakey, R., Karlstrom, K.E., Timmons, J.M., Dickinson, B., Pecha, M., 2011. Detrital zircon U-Pb geochronology of Paleozoic strata in the Grand Canyon, Arizona. Lithosphere 3, 183–200.

Kuiper, N.H., 1960. Tests concerning random points on a circle. Proceedings of the Koninklijke Nederlandse Akademie van Wetenschappen 63, 38–47.

Licht, A., Pullen, A., Kapp, P., Abell, J., Giesler, N., 2016. Eolian cannibalism: Reworked loess and fluvial sediment as the main sources of the Chinese Loess Plateau. Geol. Soc. Am. Bull. 128, 944–956.

Massey Jr, F.J., 1951. The Kolmogorov-Smirnov test for goodness of fit. Journal of the American statistical Association 46, 68-78.

Saylor, J.E., Knowles, J.N., Horton, B.K., Nie, J.S., Mora, A., 2013. Mixing of Source Populations Recorded in Detrital Zircon U-Pb Age Spectra of Modern River Sands. J. Geol. 121, 17–33.

Saylor, J.E., Stockli, D.F., Horton, B.K., Nie, J., Mora, A., 2012. Discriminating rapid exhumation from syndepositional volcanism using detrital zircon double dating: Implications for the tectonic history of the Eastern Cordillera, Colombia. Geol. Soc. Am. Bull. 124, 762–779.

Saylor, J.E., Sundell, K.E., 2016. Quantifying comparison of large detrital geochronology data sets. Geosphere 12, 203–220.

Saylor, J.E., Sundell, K.E., Sharman, G.R., 2019. Characterizing sediment sources by non-negative matrix factorization of detrital geochronological data. Earth Planet. Sci. Lett. 512, 46-58.

Sharman, G.R., Sharman, J.P., Sylvester, Z., 2018. detritalPy: A Python-based toolset for visualizing and analysing detrital geo-thermochronologic data. The Depositional Record 4, 202-215.

Sundell, K.E., Saylor, J.E., 2017. Unmixing detrital geochronology age distributions. Geophysics, Geochemistry, Geosystems 18.