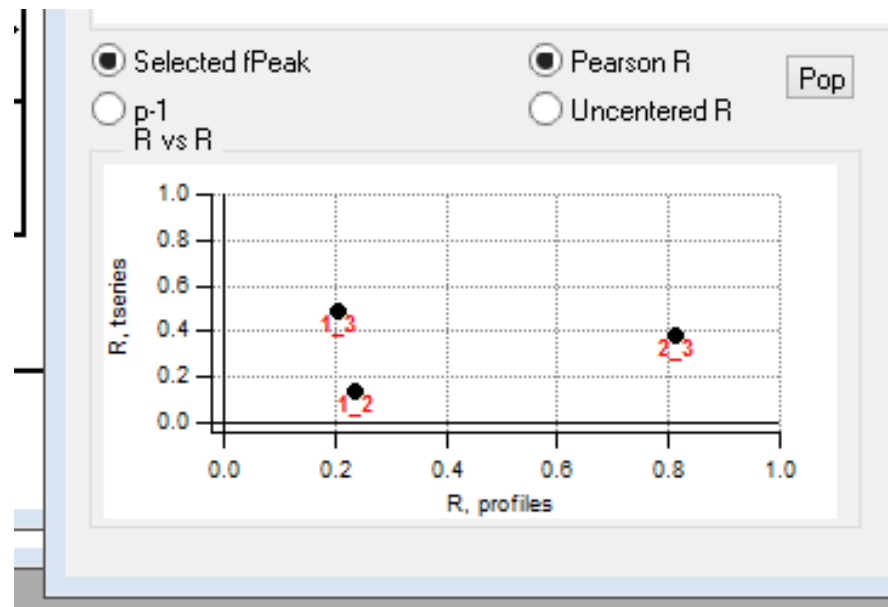


Centred vs Uncentred r

James Allan

Motivation

- Frequently need to compare mass spectra, e.g. PMF output with library spectra
- Need a metric that is a good discriminator



Generic unweighted closeness metric

- Data vectors A and B of length n
- Generate closeness metric S

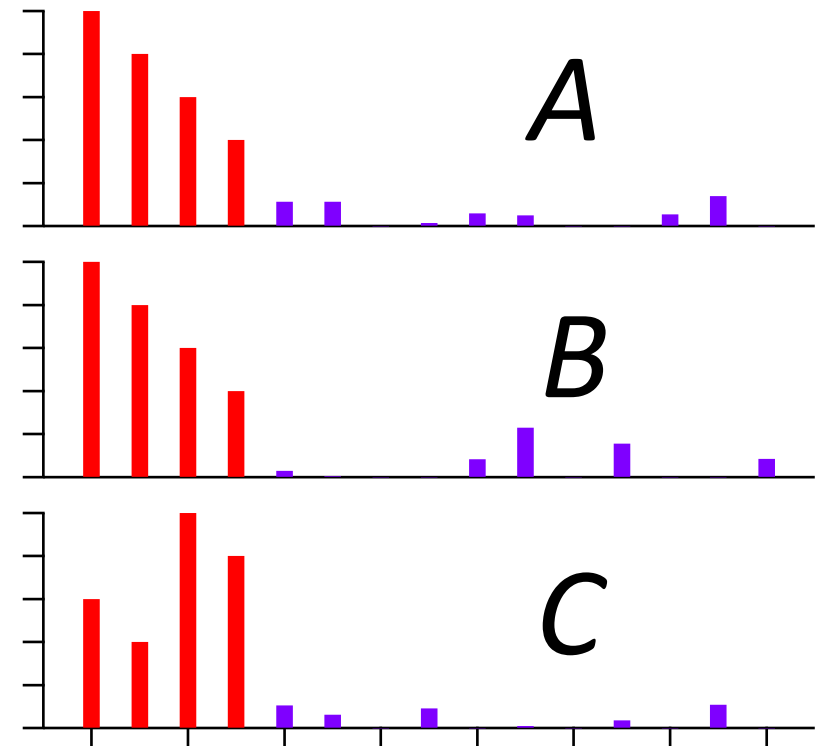
$$S(A, B) = \frac{\sum_{i=1}^n (A_i - A')(B_i - B')}{\sqrt{\sum_{i=1}^n (A_i - A')^2 \sum_{i=1}^n (B_i - B')^2}}$$

- For r : $A' = \bar{A} = n^{-1}$ (if normalised)
- For uncentred r : $A' = 0$
 - AKA normalised dot product:

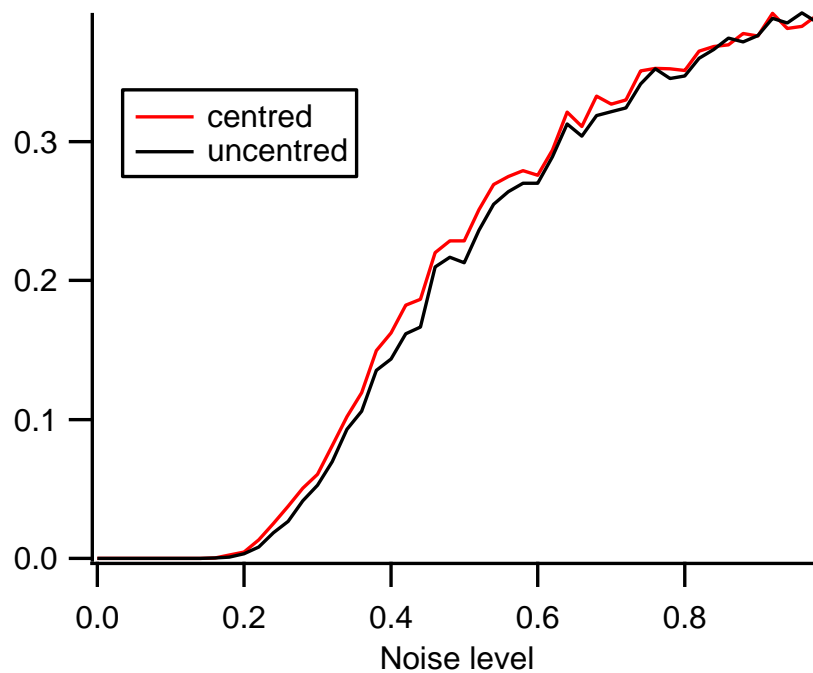
$$S(A, B) = \frac{A \cdot B}{\sqrt{(A \cdot A)(B \cdot B)}}$$

Testing Metrics

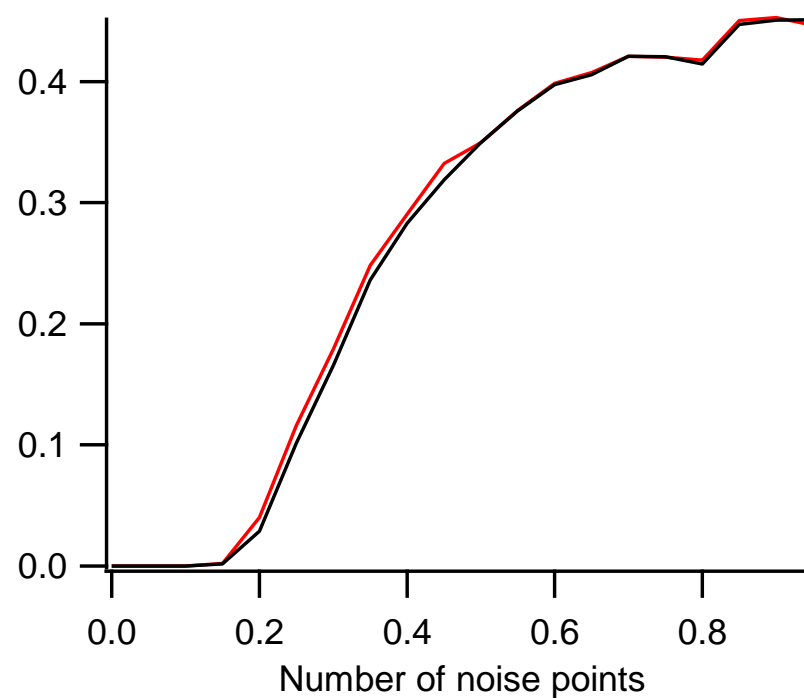
- Generate three vectors of synthetic data
 - Part 'real', part 'noise'
- $S(A,B)$ should always be greater than $S(A,C)$
- Repeat multiple times with different random data, record fraction of misclassified results



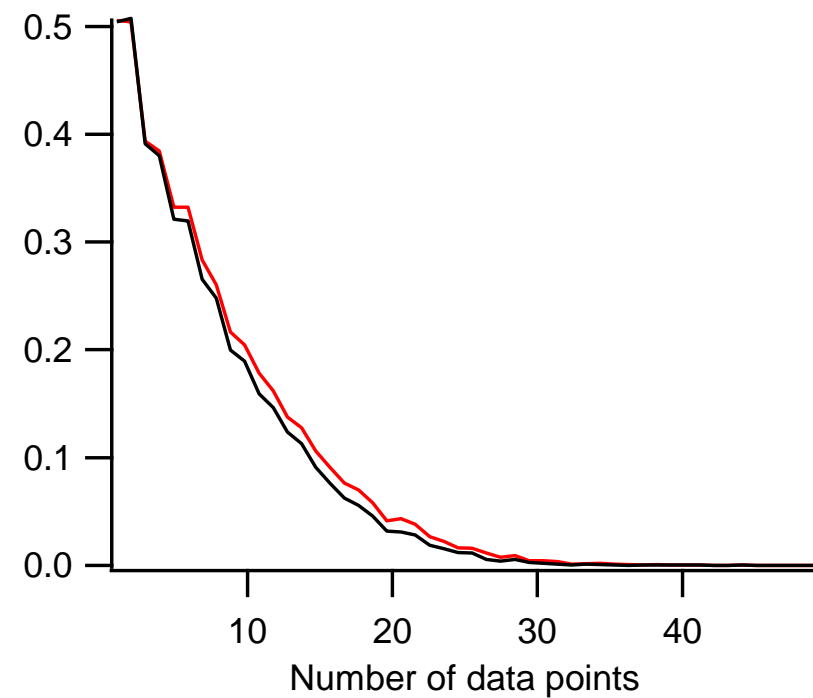
Misclassified fraction



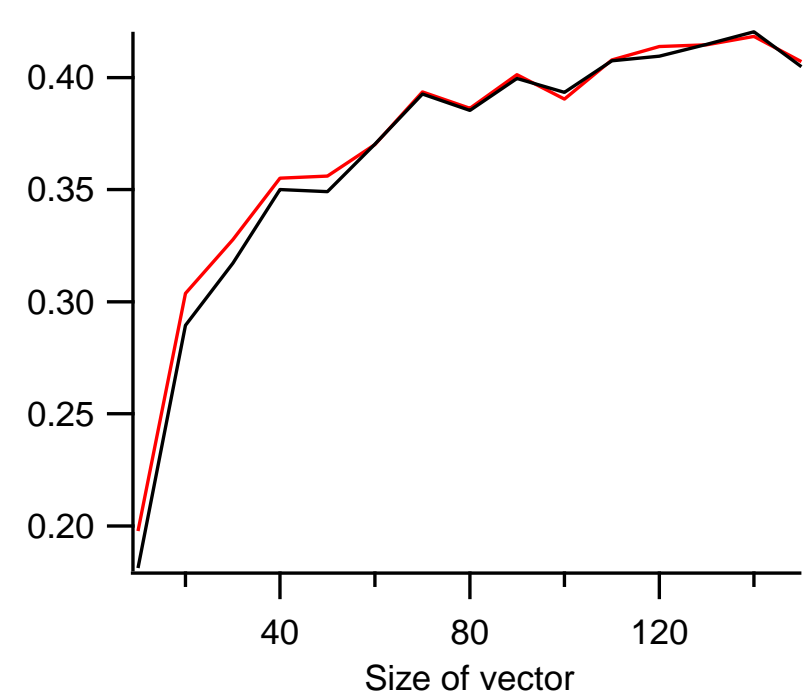
Misclassified fraction



Misclassified fraction



Misclassified fraction



Conclusions

- Uncentred r is more mathematically sound and performs better than centred r for comparing mass spectra
 - But not by much
- Centred r remains better for time series as the mean of a time series is meaningful

