

Kevin Wang

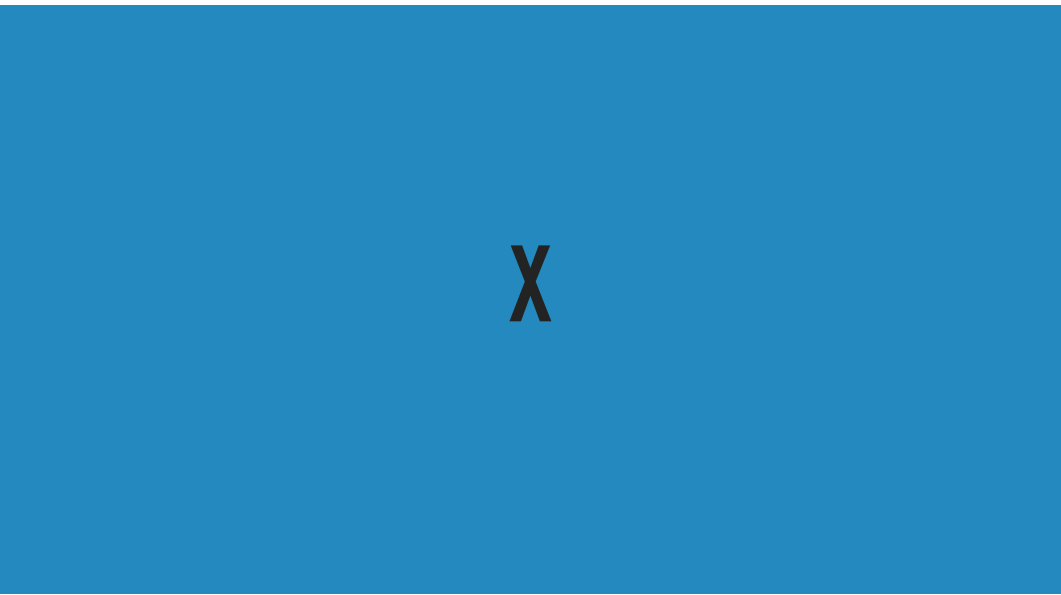
Dimensional Reduction

Acknowledgement

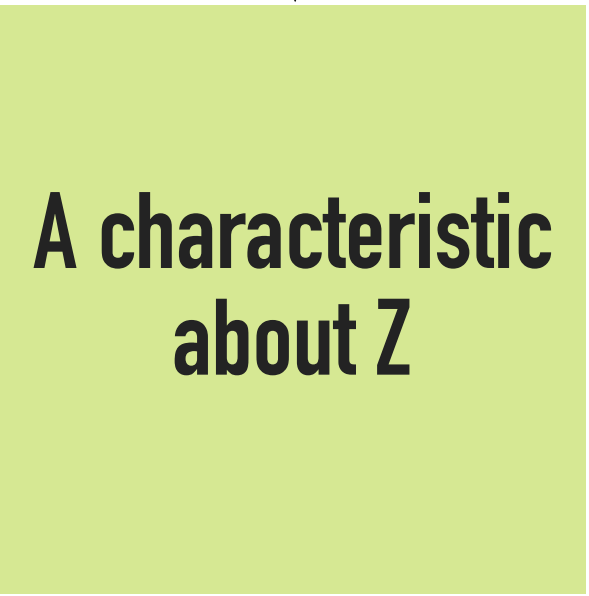
- ▶ A large proportion of this material was adapted from the Honours thesis of Nelson Ma, formerly at the School of Mathematics and Statistics, the University of Sydney.

Dimensional reduction

- ▶ High dimensional data are tricky:
 - ▶ Correlation between variables could contain redundant information
 - ▶ Humans eyes are not great beyond 3 dimensions
 - ▶ Humans brains are not great at handling non-linear relationships
- ▶ Reduce the dimension of our data, while **preserving** one key characteristic

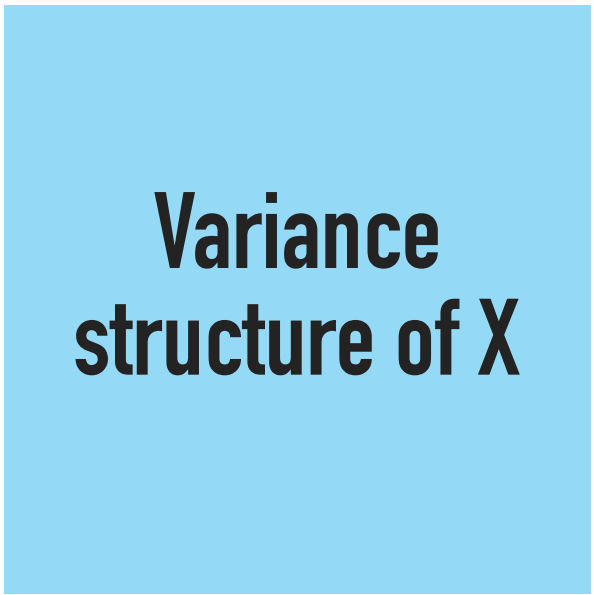
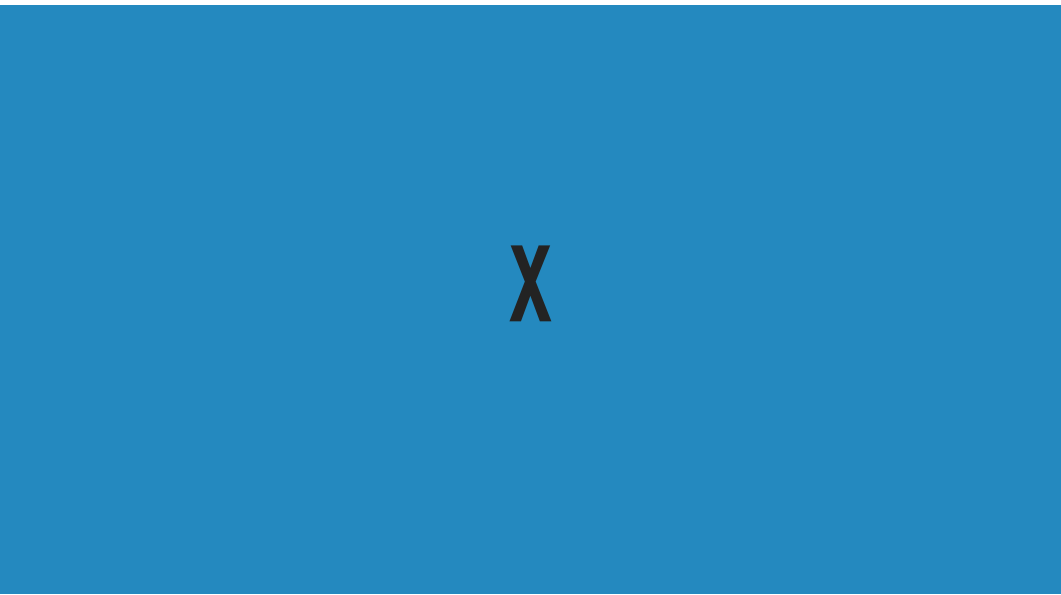


Preserve/minimise

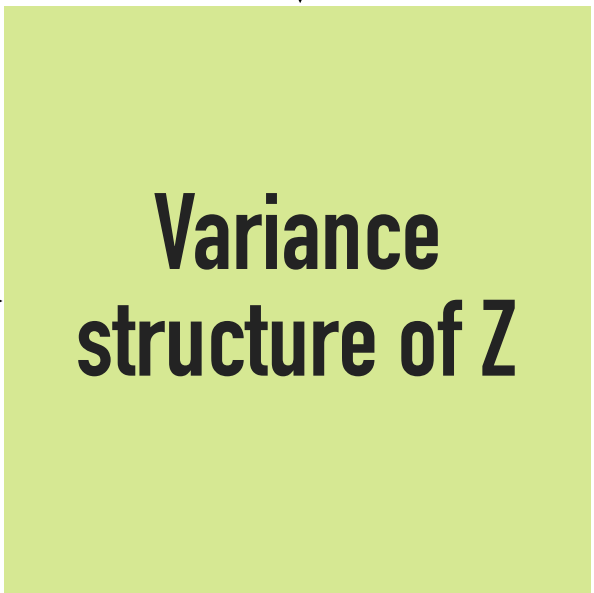


PCA

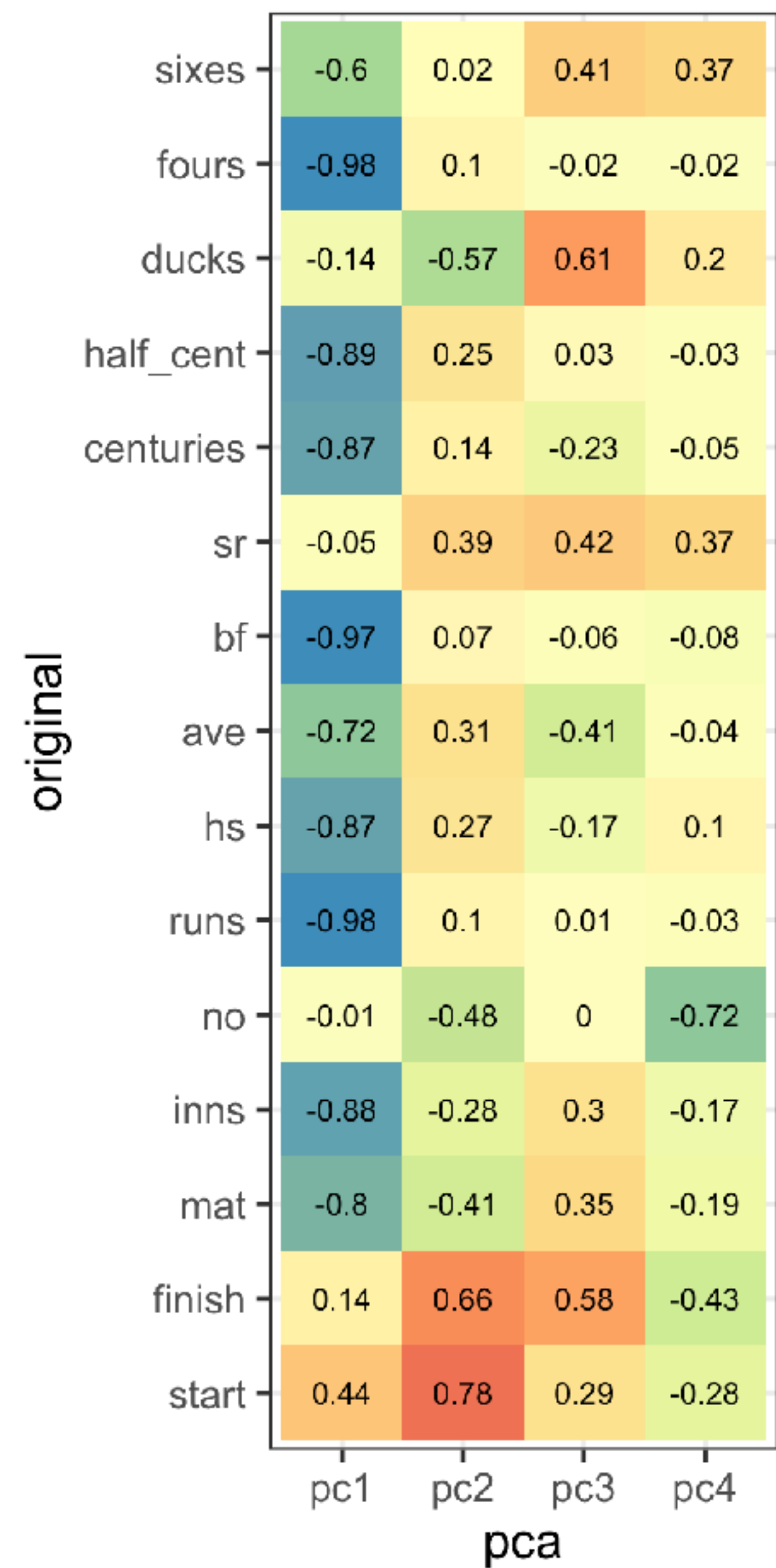
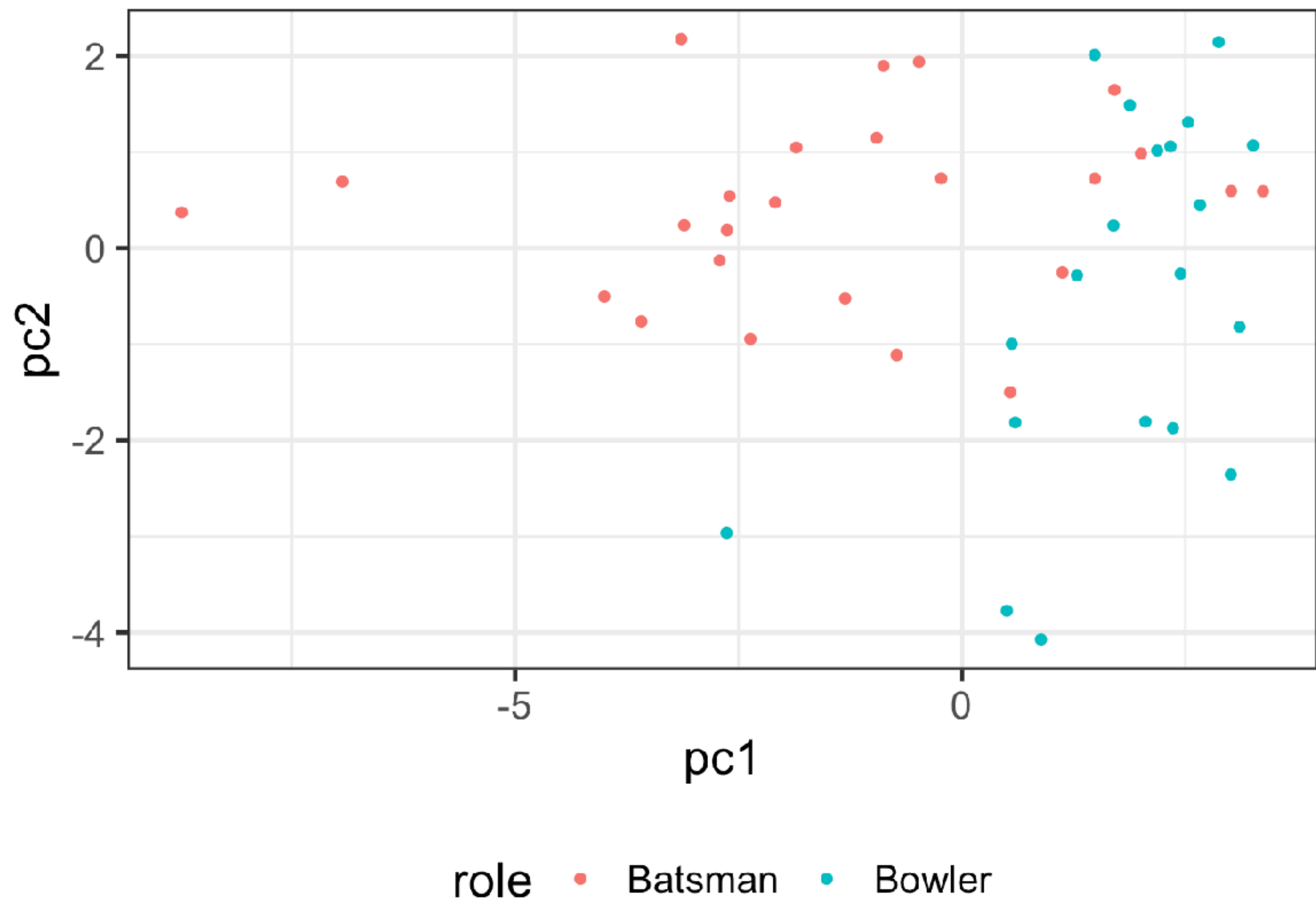
- ▶ Decompose the correlation matrix $\Sigma = U\Lambda U^\top$
- ▶ Create a score matrix: $Z = XU$
- ▶ The score matrix has the **same amount of variance** as the original data matrix
- ▶ Columns of score matrix **successively** inherit the maximum possible variance from X
- ▶ This is why the first few columns of the score matrix can be used for visualisation: they already captured a large amount of variation in the original data.



Preserved



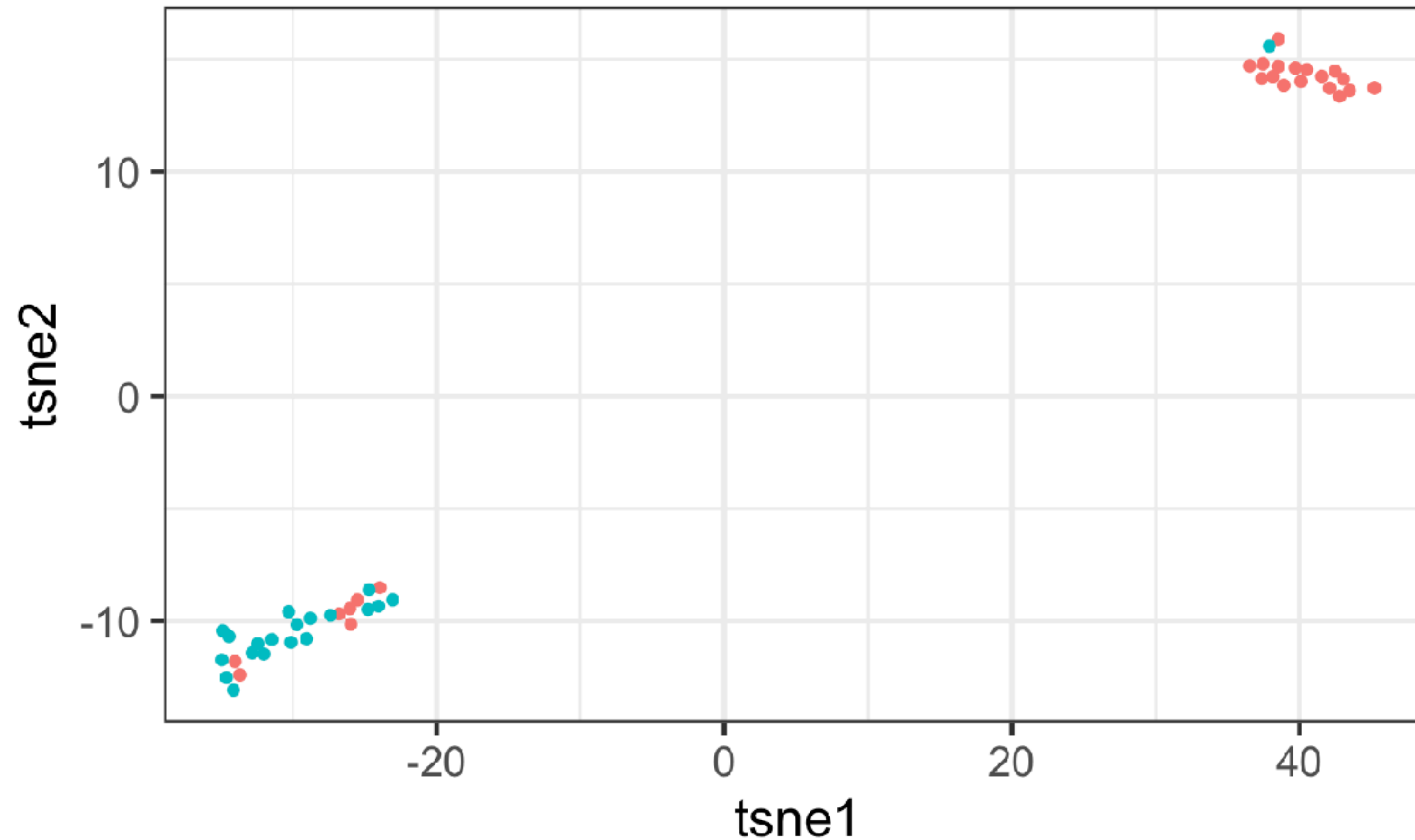
PCA visualisation



tSNE: t-distributed stochastic neighbor embedding

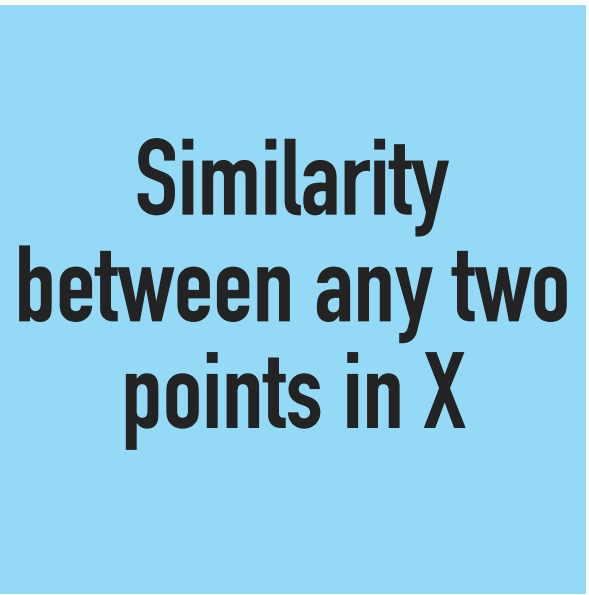
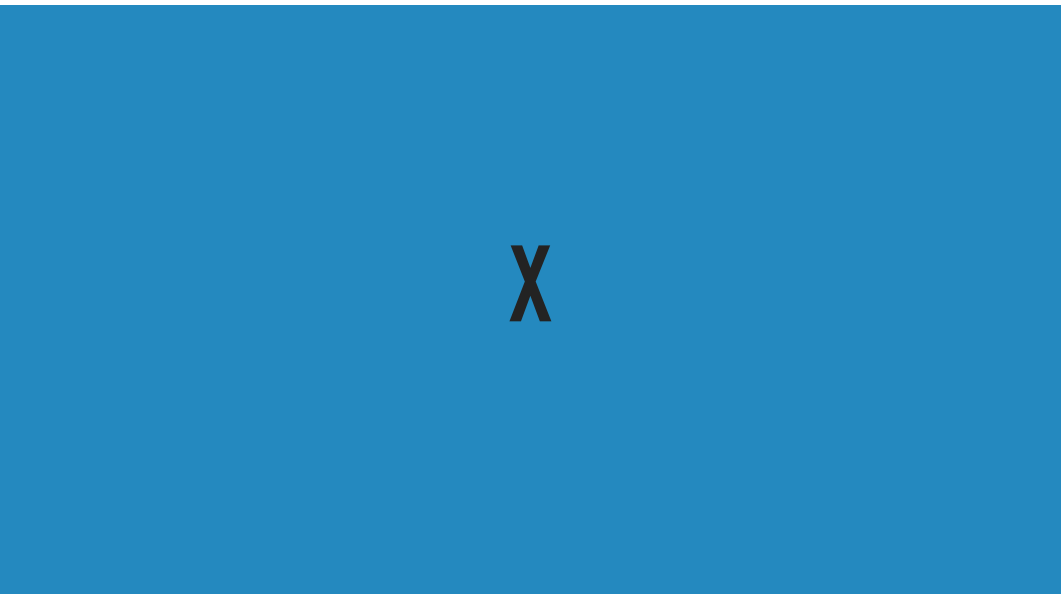
- ▶ tSNE was invented in 2008 as a non-linear alternative of PCA
- ▶ Unlike PCA, the output matrix of tSNE does not have an interpretation, but its major advantage is in the visualisation
- ▶ (Speaking from personal experience) For complex data in my research, tSNE tends to produce more separation of clusters

tSNE visualisation

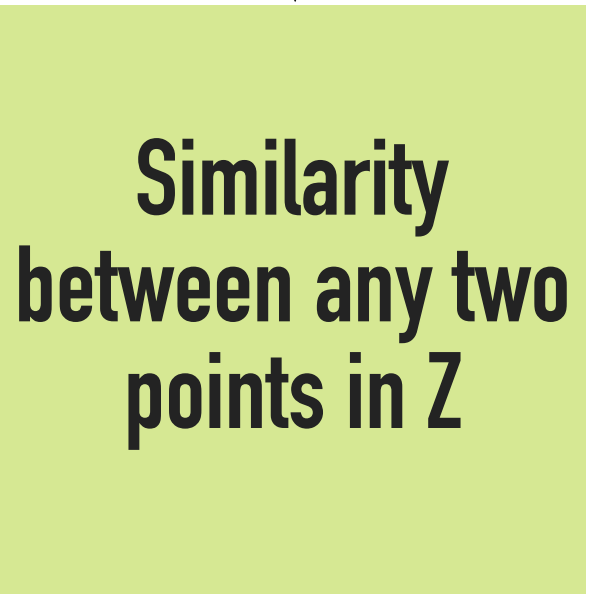


role ● Batsman ● Bowler

Points that are close to each other in the plot are also close in the original dimension



minimise



Summary

	PCA	tSNE
Relationship captured	linear	non-linear
What is preserved/ minimised between X and Z	variance	similarity between points
Interpretation of output numerical matrix	yes	no