# PCA in R

OR: SOPHISTICATED SIMPLICITY

DR ERIN I WALSH  |  21/08/2019

# The plan for today

- Why this is useful

- How to think in $n$ dimensions

- From raw score to principal component

- PCA in general

- PCA in R

Image source: dead___artist on unsplash.com

# Why this is useful

**Abstract**

**Background and purpose**

Inflammation and oxidative stress (OS) have been clearly linked to neurodegeneration. **However, studies investigating the associations between peripheral markers of inflammation and cognitive decline have produced mixed results. This is possibly due to the fact that markers are typically tested individually despite the fact that biologically they function interactively.** Thus, the aim of this study was to investigate the association between a combination of OS/inflammation markers and outcomes including mild cognitive impairment (MCI) diagnosis, cognitive decline and hippocampal atrophy.

**Methods**

Oxidative stress/inflammation status was assessed in 380 older community-living individuals. Thirteen blood markers were assayed. **Principal component analysis (PCA) of all markers was conducted to identify the more salient inflammatory components.** Associations between significant principal components, MCI diagnosis, previous change in Mini-Mental State Examination (MMSE) score and hippocampal atrophy were investigated through logistic and linear multiple regression.

**Results**

**Two factors (PC1 and PC2) reflecting predominantly broad pro-inflammatory activity and two factors (PC3 and PC4) reflecting predominantly OS activity were identified by PCA analysis.** PC3 and PC4 were predictive of MCI. PC3 was also predictive of prior MMSE change. PC1, PC2 and PC3 were significantly associated with hippocampal atrophy.

**Conclusions**

Combined analysis of complex and interacting biomarkers revealed a protective association between antioxidant activity and MCI that is consistent with lifestyle factors shown to reduce risk of cognitive decline. **OS and broad systemic inflammation were also found to be associated with hippocampal atrophy further highlighting the benefits of the PCA methodology applied in this study.**

# Why this is useful

Walsh, E. I., Shaw, M. E., Oyarce, D. A. E., Fraser, M., & Cherbuin, N. (2019). Assumption-Free Assessment of Corpus Callosum Shape: Benchmarking and Application. *Concepts in Magnetic Resonance Part A*, 2019.

# How to think in *n* dimensions



▶ Faces have many possibly salient features that vary

    ▶ e.g. eyes

# How to think in *n* dimensions

▶ Each of these features can vary in some way

  ▶ e.g. eyes can vary in position. This is a *dimension*: eye position

# How to think in *n* dimensions

▶ Each of these features can vary in some way

  ▶ e.g. eyes can vary in position. This is a *dimension*: eye position



▶ In any sample, there is an average on every dimension

  ▶ e.g. mean eye position is at the **centre** of the dimension

    ▶ Not necessarily a real individual, could be somewhere in between

# How to think in *n* dimensions

- Each of these features can vary in some way
  - e.g. eyes can vary in position. This is a *dimension*: eye position



- In any sample, there is an average on every dimension
  - e.g. mean eye position is at the **centre** of the dimension
    - Not necessarily a real individual, could be somewhere in between
- Every feature of an individual in the population can be characterised by their distance from this average
  - This is the position on the dimension
  - Has distribution and variability like any other variable

# How to think in *n* dimensions

▶ You can physically arrange them along one salient dimension at a time

Eye height



Eye separation

- Or arrange across two dimensions at once

Eye height →

Eye separation →

▶ Or arrange across two dimensions at once



Eye height →

↑ Eye separation ↓

Low and close

Average

High and wide

Or arrange across three…

Eye height

Eye separation

Mouth height

Or arrange across three…

Eye height

Eye separation

Mouth height

► But physical arrangements fall apart across more dimensions

▶ This was a constrained example (thanks, Skyrim character generator!)

► But face space is real (and very complex)!



Projecting the images in the 2-D eigenspace

# How to think in *n* dimensions

Faces are a good start, but there are many *n* dimensional constructs: everyone come up with a construct with more than 4 dimensions and share!



Image source: https://independent.co.uk/

# Drawing the line: Raw score to PC

Conceptually, a **component** is a variable summarising the degree to which individuals in a sample vary on a target dimension of meaning.

**Principal** components are new variables that are constructed as linear combinations or mixtures of the initial variables.

# Drawing the line: Raw score to PC

Raw scales could be considered a "component", but they…

▶ May not relate directly to the dimension of interest

# Drawing the line: Raw score to PC

Raw scales could be considered a "component", but they…

▶ May not relate directly to the dimension of interest

▶ Have different degrees of variability (on different scales)

  ▶ A principal component maximises the amount of variation captured.



High on both

Eye width

Low on both

Eye height

# Drawing the line: Raw score to PC

Raw scales could be considered a "component", but they…

▶ May not relate directly to the dimension of interest

▶ Have different degrees of variability (on different scales)

  ▶ A principal component maximises the amount of variation captured.

High on both

Admixture of eye width and height at each corner

Component 2

Component 1

Eye width

Low on both

Eye height

# Drawing the line: Raw score to PC

Raw scales could be considered a "component", but they…

▶ May not relate directly to the dimension of interest

▶ Have different degrees of variability (on different scales)

▶ A principal component maximises the amount of variation captured.



Eye width

Eye height

Component 2

Component 1

# Drawing the line: Raw score to PC

Raw scales could be considered a "component", but they…

▶ May not relate directly to the dimension of interest

▶ Have different degrees of variability (on different scales)

   ▶ A principal component maximises the amount of variation captured.

   ▶ This is conceptually the same as finding the line of best fit in a linear model.

# Drawing the line: Raw score to PC

Raw scales could be considered a "component", but they...

▶ May not relate directly to the dimension of interest

▶ Have different degrees of variability (on different scales)

　▶ A principal component maximises the amount of variation captured.

　▶ This is conceptually the same as finding the line of best fit in a linear model.

　▶ Standardize all variables first

# Drawing the line: Raw score to PC
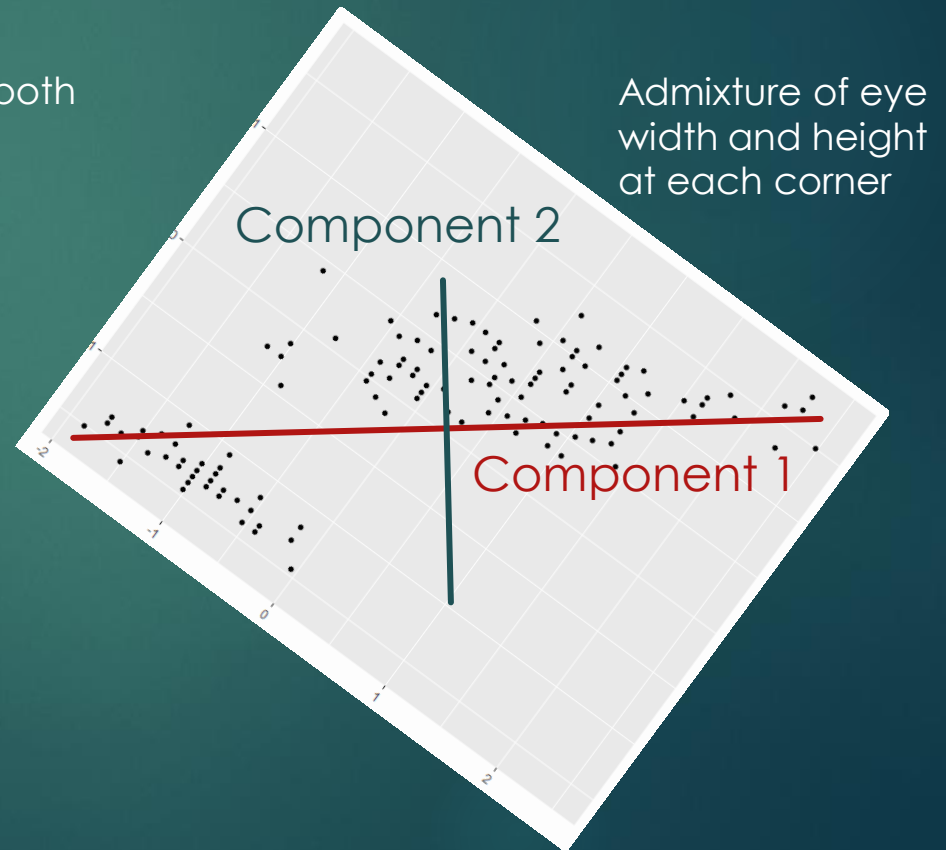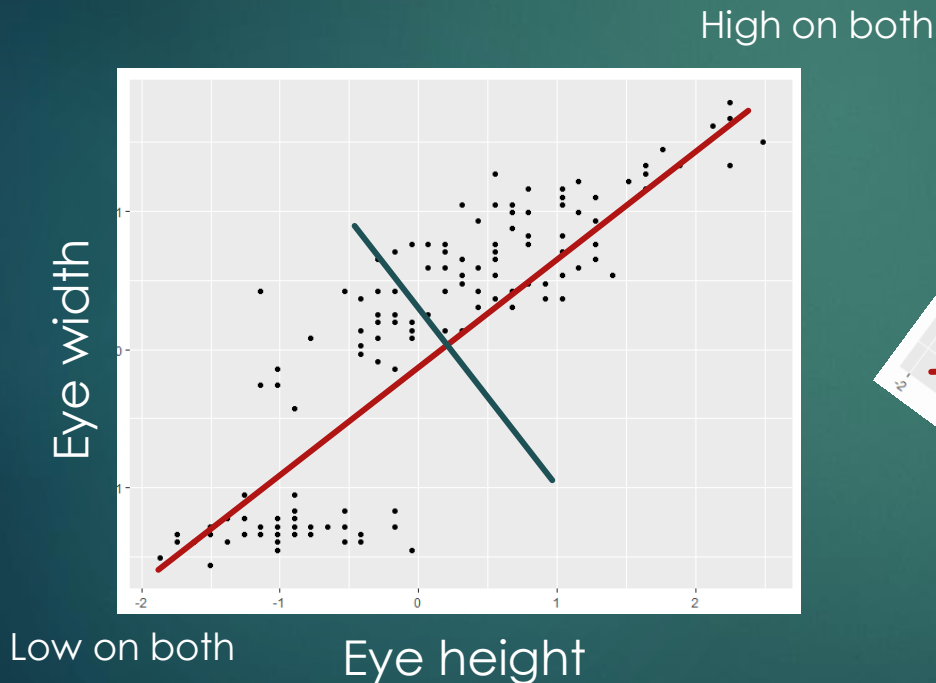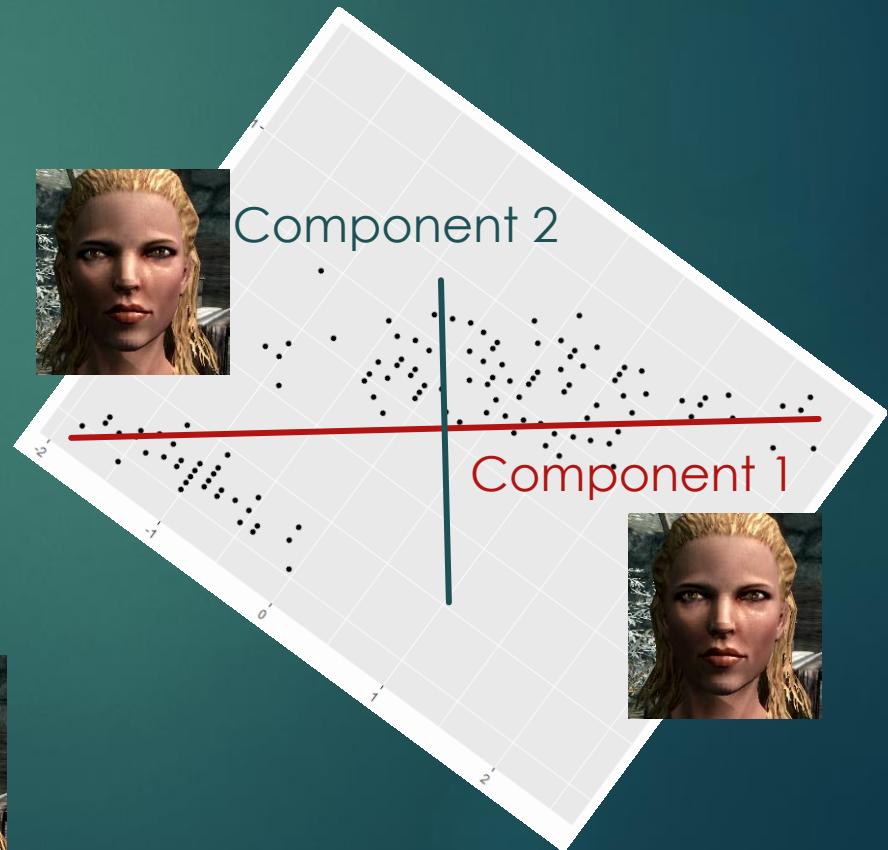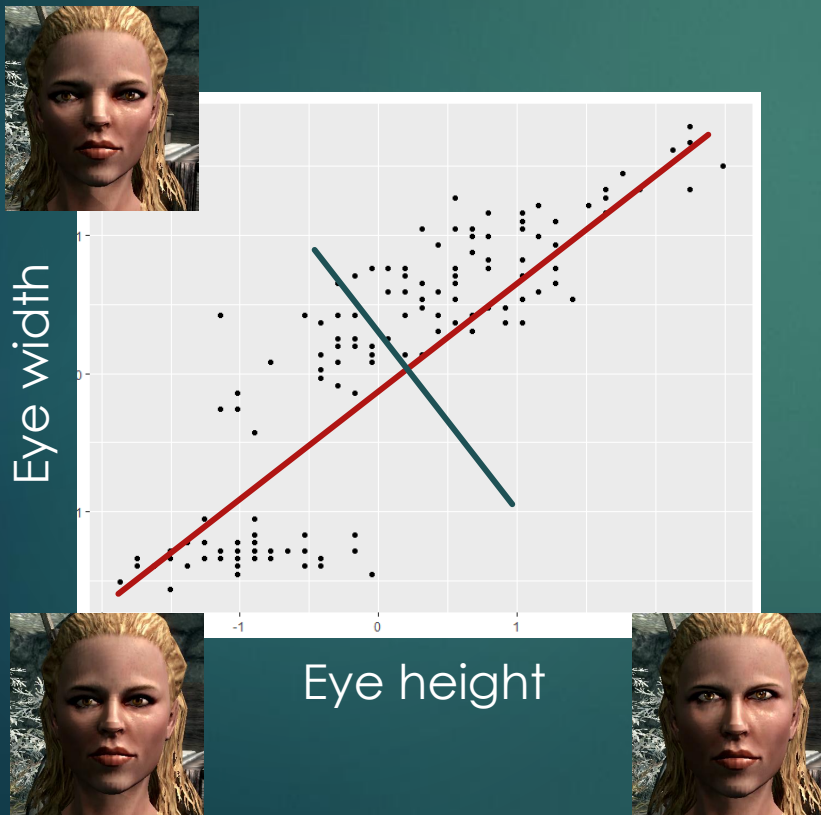
Raw scales could be considered a "component", but they…

▶ May not relate directly to the dimension of interest

▶ Have different degrees of variability (on different scales)

▶ Have differing relationships to one another (covariance)

   ▶ Typically Principal Components just take one of two highly correlated items

   ▶ Side note: rotation (technically no longer PCA if you rotate)

PCA: find linear combination of best fit for all points

Varimax = orthogonal, does not allow them to be correlated

Oblimin = oblique, allows them to be correlated

# Drawing the line: Raw score to PC
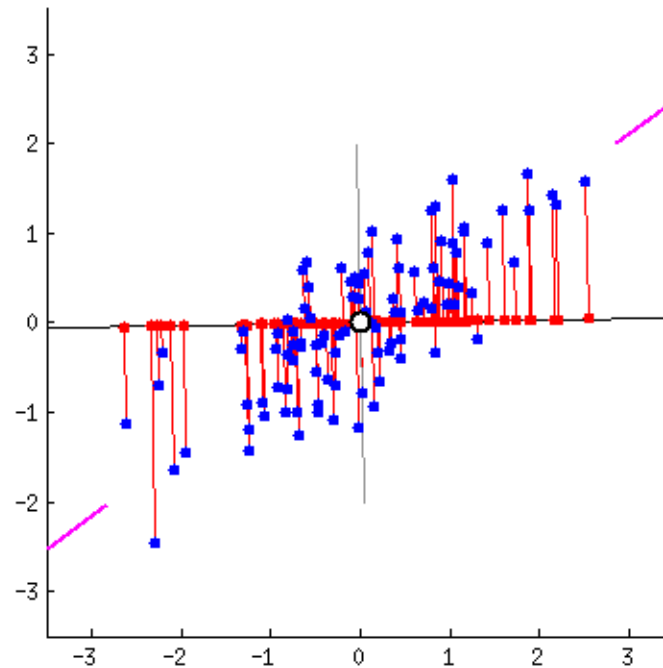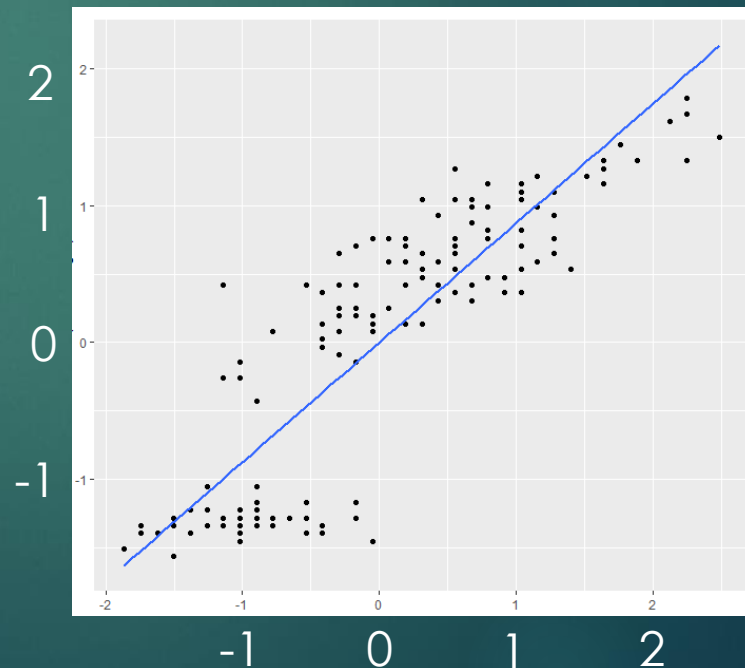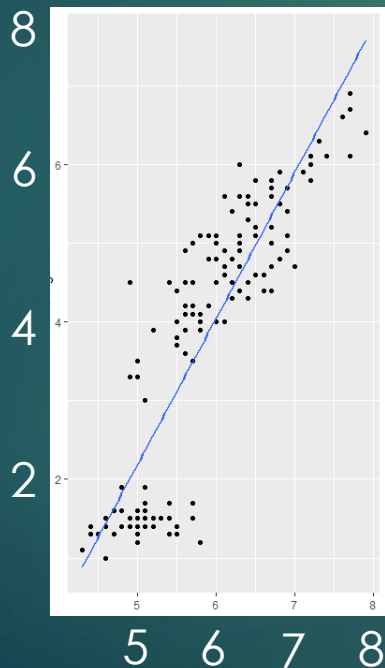
Raw scales could be considered a "component", but they…

▶   May not relate directly to the dimension of interest

▶   Have different degrees of variability (on different scales)

▶   Have differing relationships to one another (covariance)

▶   Give no guidance as to what is or is not a useful.

Principal components are fit iteratively, so first combinations explain the most variability. Less explained with each new component, so you can pick just a few based on a cut off (e.g. "elbow" of a scree plot)



% of variance explained from original variables

Number of components

# Some terms…

- Dimension: any axis of possibly meaningful variation in your data

- Component: a variable summarising the degree to which individuals in a sample vary on a dimension

- Principal component: variables constructed from linear combinations or mixtures of the initial variables.

- Principal Component Analysis: dimension reduction method that transforming a large set of variables into a smaller set (principal components) that still contains most of the information (variability) in the original set.

THE ELUSIVE LIKENESS

Image source: https://www.lovelif edrawing.com/

# PCA in general

1. Standardization (conversion to z scores):
   - ▸ Homogenizes range so they contribute equally to regression line of best fit
2. Covariance matrix computation
   - ▸ Computes associations between input variables, to detect what is redundant/can be dropped
3. Find the eigenvectors
   - ▸ The directions of the axes where there is the most variance. These are our principal components!

   … and the eigenvalues
   - ▸ The amount of variance each principal component explains. This is the values on the scree plot
4. Find the feature vector
   - ▸ Arrange our eigenvectors by eigenvalues to see which to keep
5. Recast the data to this new shape space
   - ▸ Gives you the information you're after
     - ▸ Eigenvectors: how much each of the original variables contributes to each principal component
       - ▸ Similar to concept of "loading" in FA, though here 'loadings' are orthonormal Eigenvectors
     - ▸ Position: where each individual falls on each principal component
       - ▸ This is what you'll use in later analyses if you were using PCA for dimension reduction

# PCA in R

▶ Standardization

▶ Covariance matrix computation

▶ Fitting the PCA

▶ Evaluating the PCA

    ▶ Eigenvalues

    ▶ Eigenvectors

    ▶ Position on components

We'll run through an example together, then you try!

Image source: Daniel Gonzalez on unsplash.com