

# Simple & Multiple Correspondence Analyses

Contingency, categorical, ordinal, continuous and mixed data

Derek Beaton

Rotman Research Institute

October 29, 2019

Before we get started

# Our new best friends

## CONTINUOUS

measured data, can have  $\infty$  values within possible range.



I AM 3.1" TALL  
I WEIGH 34.16 grams

## DISCRETE

OBSERVATIONS CAN ONLY EXIST AT LIMITED VALUES, OFTEN COUNTS.



I HAVE 8 LEGS  
and  
4 SPOTS!

@allison\_horst

via @allison\_horst

## NOMINAL

UNORDERED DESCRIPTIONS



## ORDINAL

ORDERED DESCRIPTIONS



## BINARY

ONLY 2 MUTUALLY EXCLUSIVE OUTCOMES



@allison\_horst

via @allison\_horst

## CONTINUOUS

measured data, can have  $\infty$  values within possible range.



I AM 3.1" TALL  
I WEIGH 34.16 grams

## DISCRETE

OBSERVATIONS CAN ONLY EXIST AT LIMITED VALUES, OFTEN COUNTS:



I HAVE 8 LEGS  
and 4 SPOTS!

## NOMINAL

UNORDERED DESCRIPTIONS



I'M A TURTLE!  
I'M A SNAIL!  
I'M A BUTTERFLY!

## ORDINAL

ORDERED DESCRIPTIONS



I AM UNHAPPY  
I AM OK.  
I AM AWESOME!!!

## BINARY

ONLY 2 MUTUALLY EXCLUSIVE OUTCOMES



@alissa\_horch

- ▶ What do we do with all of these in a PCA like way?

## CONTINUOUS

measured data, can have  $\infty$  values within possible range.



I AM 3.1" TALL  
I WEIGH 34.16 grams

## DISCRETE

OBSERVATIONS CAN ONLY EXIST AT LIMITED VALUES, OFTEN COUNTS:



I HAVE 8 LEGS  
and 4 SPOTS!

## NOMINAL

UNORDERED DESCRIPTIONS



I'M A TURTLE!  
I'M A SNAIL!  
I'M A BUTTERFLY!

## ORDINAL

ORDERED DESCRIPTIONS



## BINARY

ONLY 2 MUTUALLY EXCLUSIVE OUTCOMES



@alissa\_horch

- ▶ What do we do with all of these in a PCA like way?
- ▶ Some are very difficult and effectively ignored

## CONTINUOUS

measured data, can have  $\infty$  values within possible range.



I AM 3.1" TALL  
I WEIGH 34.16 grams

## DISCRETE

OBSERVATIONS CAN ONLY EXIST AT LIMITED VALUES, OFTEN COUNTS:



I HAVE 8 LEGS  
and 4 SPOTS!

## NOMINAL

UNORDERED DESCRIPTIONS



i'm a  
snail! —  
i'm a  
butterfly!

## ORDINAL

ORDERED DESCRIPTIONS



-I am unhappy  
-I am OK.  
-I am Awesome!!!

## BINARY

ONLY 2 MUTUALLY EXCLUSIVE OUTCOMES



i am  
extinct! —  
HA.

@alissa\_horch

- ▶ What do we do with all of these in a PCA like way?
- ▶ Some are very difficult and effectively ignored
  - ▶ We won't do that!

## CONTINUOUS

measured data, can have  $\infty$  values within possible range.



I AM 3.1" TALL  
I WEIGH 34.16 grams

## DISCRETE

OBSERVATIONS CAN ONLY EXIST AT LIMITED VALUES, OFTEN COUNTS:



I HAVE 8 LEGS  
and 4 SPOTS!

## NOMINAL

UNORDERED DESCRIPTIONS



## ORDINAL

ORDERED DESCRIPTIONS



## BINARY

ONLY 2 MUTUALLY EXCLUSIVE OUTCOMES



@alissa\_horch

@alissa\_horch

- ▶ What do we do with all of these in a PCA like way?
- ▶ Some are very difficult and effectively ignored
  - ▶ We won't do that!
- ▶ See SS Steven's typology:  
[https://en.wikipedia.org/wiki/Level\\_of\\_measurement](https://en.wikipedia.org/wiki/Level_of_measurement)

## Motivation & Objectives

- ▶ Not everything is a number

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!
- ▶ We need to recognize when this happens

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!
- ▶ We need to recognize when this happens
  - ▶ And know what to do

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!
- ▶ We need to recognize when this happens
  - ▶ And know what to do
- ▶ Introduce CA, MCA, and tricks

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!
- ▶ We need to recognize when this happens
  - ▶ And know what to do
- ▶ Introduce CA, MCA, and tricks
  - ▶ Leave you overwhelmed, but knowing that

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!
- ▶ We need to recognize when this happens
  - ▶ And know what to do
- ▶ Introduce CA, MCA, and tricks
  - ▶ Leave you overwhelmed, but knowing that
  - ▶ PCA is sometimes the most wrong approach

## Motivation & Objectives

- ▶ Not everything is a number
  - ▶ Sometimes numbers aren't numbers!
- ▶ We need to recognize when this happens
  - ▶ And know what to do
- ▶ Introduce CA, MCA, and tricks
  - ▶ Leave you overwhelmed, but knowing that
  - ▶ PCA is sometimes the most wrong approach
  - ▶ CA & MCA are suitably less wrong

## Where to find everything

- ▶ Generally: <https://github.com/derekbeaton/workshops>

## Where to find everything

- ▶ Generally: <https://github.com/derekbeaton/workshops>
- ▶ Today: [https://github.com/derekbeaton/Workshops/tree/master/Misc/CA\\_MCA](https://github.com/derekbeaton/Workshops/tree/master/Misc/CA_MCA)

# Overview

- ▶ Revisit PCA

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data

## Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis
  - ▶ and many of its connections

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis
  - ▶ and many of its connections
- ▶ Multiple correspondence analysis

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis
  - ▶ and many of its connections
- ▶ Multiple correspondence analysis
  - ▶ generalizes CA (amongst many other things)

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis
  - ▶ and many of its connections
- ▶ Multiple correspondence analysis
  - ▶ generalizes CA (amongst many other things)
  - ▶ and how to handle various data types

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis
  - ▶ and many of its connections
- ▶ Multiple correspondence analysis
  - ▶ generalizes CA (amongst many other things)
  - ▶ and how to handle various data types
- ▶ A whole bunch of bonuses

# Overview

- ▶ Revisit PCA
- ▶ Looking at some data
- ▶ Simple correspondence analysis
  - ▶ and many of its connections
- ▶ Multiple correspondence analysis
  - ▶ generalizes CA (amongst many other things)
  - ▶ and how to handle various data types
- ▶ A whole bunch of bonuses
  - ▶ Robustness, PLS, Software

## Revisting PCA

## What is PCA for?

- ▶ When we can compute a covariance or correlation matrix

## What is PCA for?

- ▶ When we can compute a covariance or correlation matrix
- ▶ Break data into components

## What is PCA for?

- ▶ When we can compute a covariance or correlation matrix
- ▶ Break data into components
  - ▶ Orthogonal

## What is PCA for?

- ▶ When we can compute a covariance or correlation matrix
- ▶ Break data into components
  - ▶ Orthogonal
  - ▶ Rank ordered

## What is PCA for?

- ▶ When we can compute a covariance or correlation matrix
- ▶ Break data into components
  - ▶ Orthogonal
  - ▶ Rank ordered
  - ▶ Made of bits & pieces of original measures

Some data

## Diagnosis and education

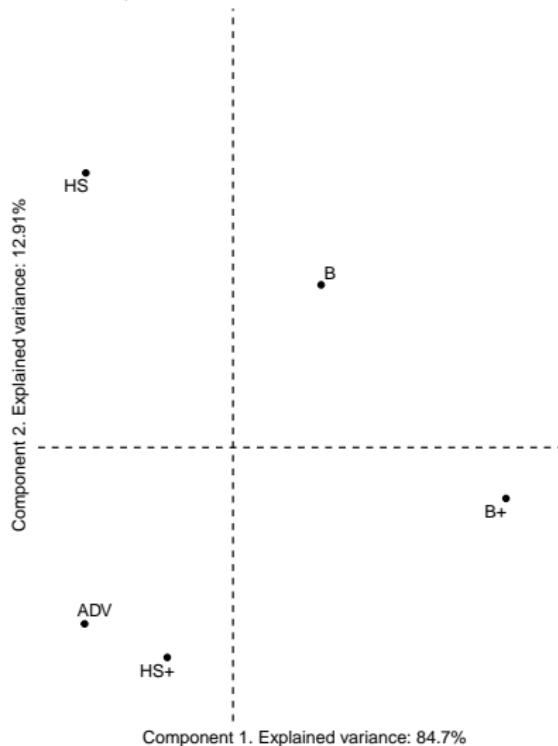
	CN	Dementia	MCI
<i>ADV</i>	39	7	54
<i>B</i>	57	17	75
<i>B+</i>	75	19	113
<i>HS</i>	25	13	46
<i>HS+</i>	39	9	77

- ▶ Given a table, and asked for a multivariate analysis

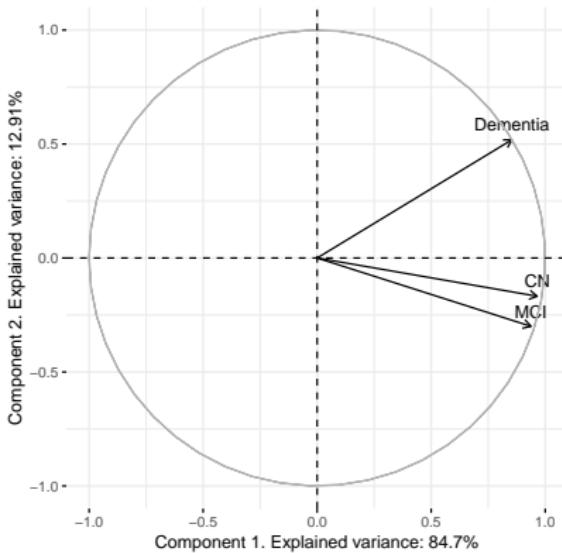
- ▶ Given a table, and asked for a multivariate analysis
- ▶ We do what we know: PCA



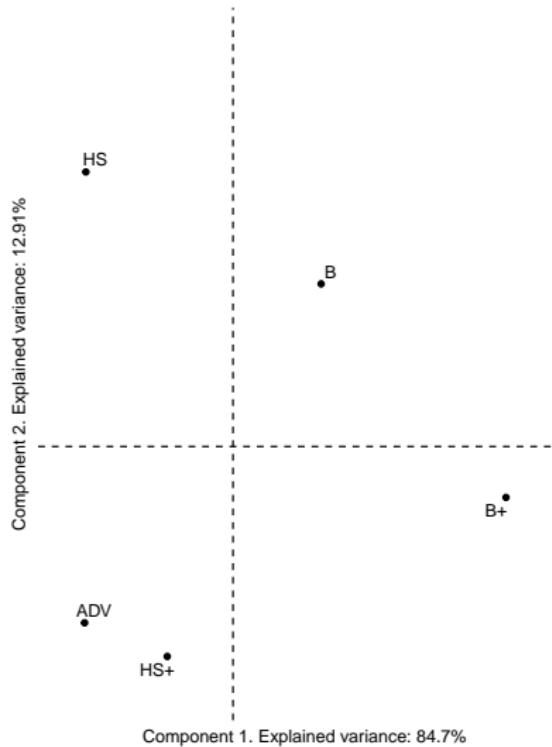
PCA:  
Row component scores



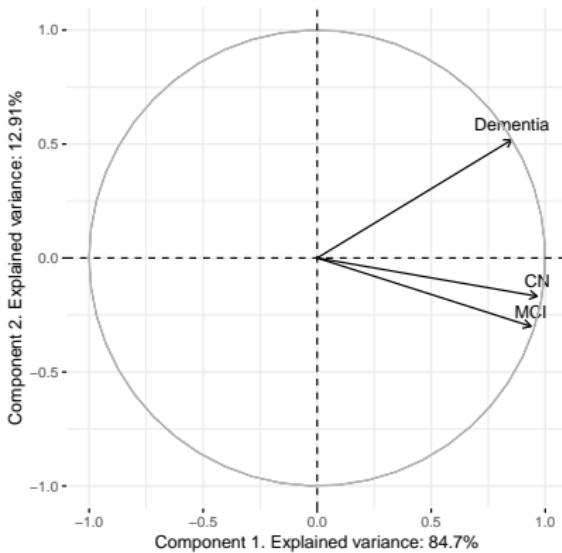
PCA:  
Variable–Component Correlations



PCA:  
Row component scores



PCA:  
Variable–Component Correlations



## What did we analyze?

	CN	Dementia	MCI
CN	1.000	0.730	0.921
Dementia	0.730	1.000	0.652
MCI	0.921	0.652	1.000

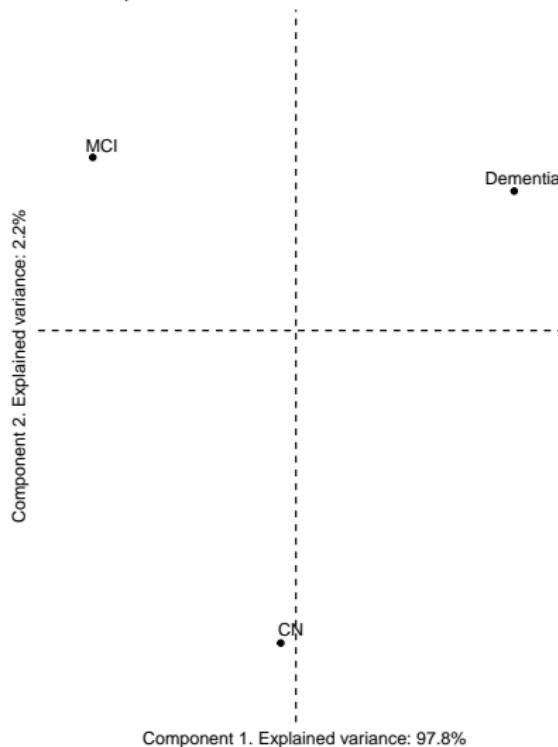
## What did PCA detect?

	CN	Dementia	MCI	<b><i>Row sums</i></b>
<i>ADV</i>	39	7	54	<b>100</b>
<i>B</i>	57	17	75	<b>149</b>
<i>B+</i>	75	19	113	<b>207</b>
<i>HS</i>	25	13	46	<b>84</b>
<i>HS+</i>	39	9	77	<b>125</b>

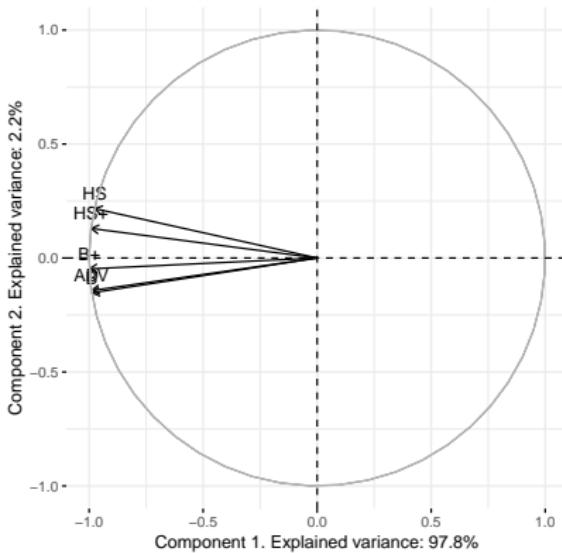
Let's try something different!

	ADV	B	B+	HS	HS+
<i>CN</i>	39	57	75	25	39
<i>Dementia</i>	7	17	19	13	9
<i>MCI</i>	54	75	113	46	77

PCA:  
Row component scores



PCA:  
Variable–Component Correlations



## What did PCA analyze?

	ADV	B	B+	HS	HS+
ADV	1.000	1.000	0.995	0.935	0.963
B	1.000	1.000	0.994	0.932	0.960
B+	0.995	0.994	1.000	0.965	0.984
HS	0.935	0.932	0.965	1.000	0.996
HS+	0.963	0.960	0.984	0.996	1.000

## What did PCA detect?

	ADV	B	B+	HS	HS+	<b><i>Row sums</i></b>
<i>CN</i>	39	57	75	25	39	<b>235</b>
<i>Dementia</i>	7	17	19	13	9	<b>65</b>
<i>MCI</i>	54	75	113	46	77	<b>365</b>

## What is PCA for?

- ▶ When we can compute a *meaningful* covariance or correlation matrix

Let's take another look

	CN	Dementia	MCI	<b><i>Row sums</i></b>
<b><i>ADV</i></b>	39	7	54	<b><i>100</i></b>
<b><i>B</i></b>	57	17	75	<b><i>149</i></b>
<b><i>B+</i></b>	75	19	113	<b><i>207</i></b>
<b><i>HS</i></b>	25	13	46	<b><i>84</i></b>
<b><i>HS+</i></b>	39	9	77	<b><i>125</i></b>
<b><i>Column sums</i></b>	<b><i>235</i></b>	<b><i>65</i></b>	<b><i>365</i></b>	

- ▶ Tell me things about this matrix

Let's take another look

	CN	Dementia	MCI	<b><i>Row sums</i></b>
<b><i>ADV</i></b>	39	7	54	<b><i>100</i></b>
<b><i>B</i></b>	57	17	75	<b><i>149</i></b>
<b><i>B+</i></b>	75	19	113	<b><i>207</i></b>
<b><i>HS</i></b>	25	13	46	<b><i>84</i></b>
<b><i>HS+</i></b>	39	9	77	<b><i>125</i></b>
<b><i>Column sums</i></b>	<b><i>235</i></b>	<b><i>65</i></b>	<b><i>365</i></b>	

- ▶ Tell me things about this matrix
- ▶ What kind of problem does this look like?

## Simple correspondence analysis

## What is CA?

- ▶ Initially: *visualize contingency tables* (**a la PCA**, factor analyses)

## What is CA?

- ▶ Initially: *visualize contingency tables* (**a la PCA**, factor analyses)
  - ▶ Text (corpus) of philosophy, biblical passages, literature

## What is CA?

- ▶ Initially: *visualize contingency tables* (**a la PCA**, factor analyses)
  - ▶ Text (corpus) of philosophy, biblical passages, literature
  - ▶ From Benzecri (1964) & Escofier (1965)

## What is CA?

- ▶ Initially: *visualize contingency tables* (**a la PCA**, factor analyses)
  - ▶ Text (corpus) of philosophy, biblical passages, literature
  - ▶ From Benzecri (1964) & Escofier (1965)
  - ▶ Fully developed by Escofier (1969)

## What is CA?

- ▶ Initially: *visualize contingency tables* (**a la PCA**, factor analyses)
  - ▶ Text (corpus) of philosophy, biblical passages, literature
  - ▶ From Benzecri (1964) & Escofier (1965)
  - ▶ Fully developed by Escofier (1969)
- ▶ Explosion of the technique in France

## What is CA?

- ▶ Initially: *visualize contingency tables* (**a la PCA**, factor analyses)
  - ▶ Text (corpus) of philosophy, biblical passages, literature
  - ▶ From Benzecri (1964) & Escofier (1965)
  - ▶ Fully developed by Escofier (1969)
- ▶ Explosion of the technique in France
  - ▶ Across virtually every field (except psychology and neuroscience)

## History

- ▶ Hotelling (1933) & Thurstone (1933)

## History

- ▶ Hotelling (1933) & Thurstone (1933)
- ▶ Hirschfeld (1935) & Horst (1935)

## History

- ▶ Hotelling (1933) & Thurstone (1933)
- ▶ Hirschfeld (1935) & Horst (1935)
- ▶ Guttman (1941)

## History

- ▶ Hotelling (1933) & Thurstone (1933)
- ▶ Hirschfeld (1935) & Horst (1935)
- ▶ Guttman (1941)
- ▶ Burt (1950)

## History

- ▶ Hotelling (1933) & Thurstone (1933)
- ▶ Hirschfeld (1935) & Horst (1935)
- ▶ Guttman (1941)
- ▶ Burt (1950)
- ▶ And then Benzecri (1964) & Escofier (1965)

## History

- ▶ Hotelling (1933) & Thurstone (1933)
- ▶ Hirschfeld (1935) & Horst (1935)
- ▶ Guttman (1941)
- ▶ Burt (1950)
- ▶ And then Benzecri (1964) & Escofier (1965)
- ▶ Many more very important characters to re-discover CA

- ▶ See Lebart's History & Prehistory of CA

- ▶ See Lebart's History & Prehistory of CA
  - ▶ [http://www.dtmvic.com/doc/About\\_the\\_History\\_of\\_CA.pdf](http://www.dtmvic.com/doc/About_the_History_of_CA.pdf)

- ▶ See Lebart's History & Prehistory of CA
  - ▶ [http://www.dtmvic.com/doc/About\\_the\\_History\\_of\\_CA.pdf](http://www.dtmvic.com/doc/About_the_History_of_CA.pdf)
- ▶ And Beh & Lombardo's series

- ▶ See Lebart's History & Prehistory of CA
  - ▶ [http://www.dtmvic.com/doc/About\\_the\\_History\\_of\\_CA.pdf](http://www.dtmvic.com/doc/About_the_History_of_CA.pdf)
- ▶ And Beh & Lombardo's series
  - ▶ A genealogy of CA:  
<https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-842X.2012.00676.x>

- ▶ See Lebart's History & Prehistory of CA
  - ▶ [http://www.dtmvic.com/doc/About\\_the\\_History\\_of\\_CA.pdf](http://www.dtmvic.com/doc/About_the_History_of_CA.pdf)
- ▶ And Beh & Lombardo's series
  - ▶ A genealogy of CA:  
<https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-842X.2012.00676.x>
  - ▶ A genealogy of CA 2: <http://siba-ese.unisalento.it/index.php/ejasa/article/view/19785>

## We're diving in

	CN	Dementia	MCI
<i>ADV</i>	39	7	54
<i>B</i>	57	17	75
<i>B+</i>	75	19	113
<i>HS</i>	25	13	46
<i>HS+</i>	39	9	77

CA:  
Row component scores

Component 2. Explained variance: 37.01%

HS+

ADV

Component 1. Explained variance: 62.99%

B+

B

HS

CA:  
Column component scores

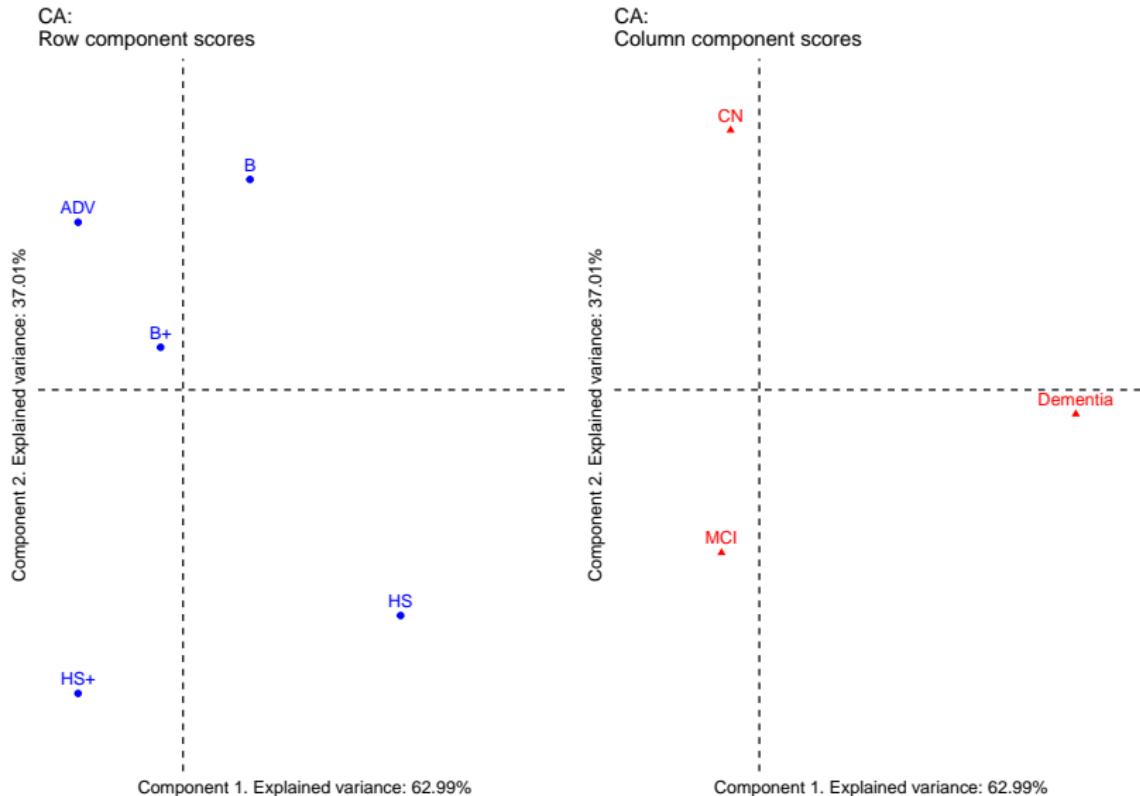
Component 2. Explained variance: 37.01%

CN

Dementia

MCI

Component 1. Explained variance: 62.99%

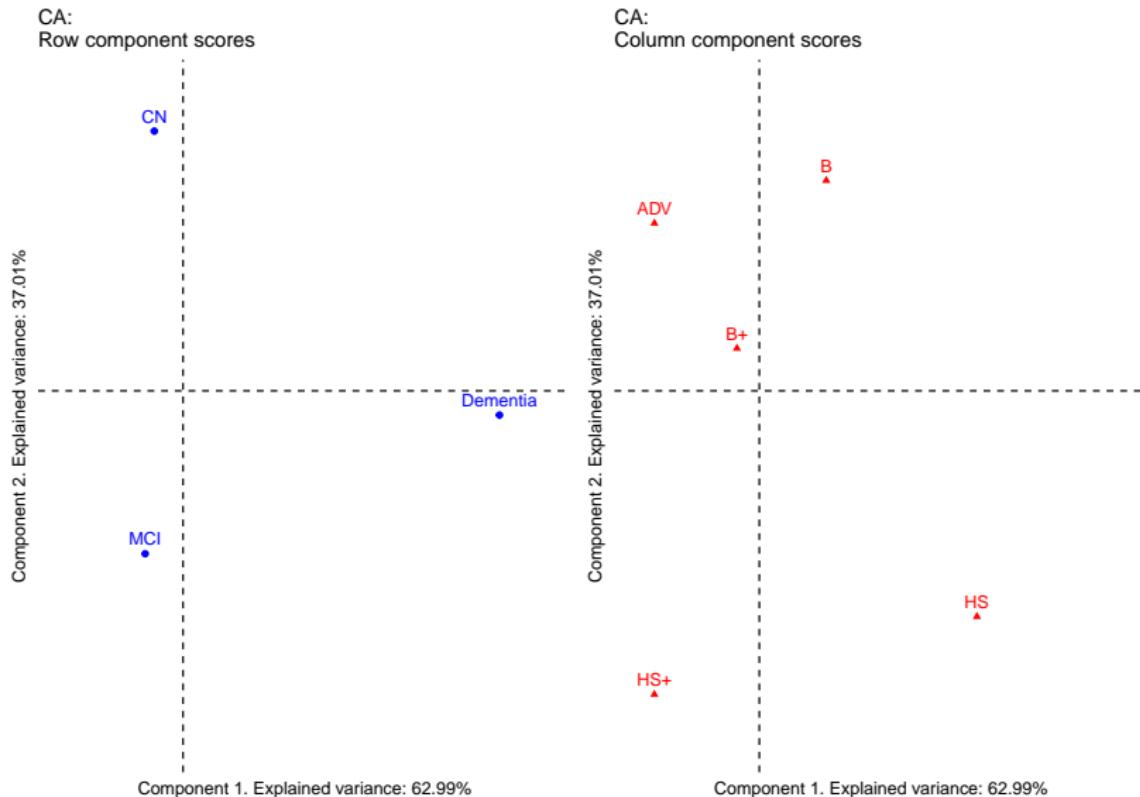


Want to see a cool trick?

	CN	Dementia	MCI
<i>ADV</i>	39	7	54
<i>B</i>	57	17	75
<i>B+</i>	75	19	113
<i>HS</i>	25	13	46
<i>HS+</i>	39	9	77

	ADV	B	B+	HS	HS+
<i>CN</i>	39	57	75	25	39
<i>Dementia</i>	7	17	19	13	9
<i>MCI</i>	54	75	113	46	77

What if we perform CA on this?



# How did that happen?

Table 1: Data

	CN	Dementia	MCI
<i>ADV</i>	39	7	54
<i>B</i>	57	17	75
<i>B+</i>	75	19	113
<i>HS</i>	25	13	46
<i>HS+</i>	39	9	77

Table 2: Observed probabilities

	CN	Dementia	MCI
<i>ADV</i>	0.059	0.011	0.081
<i>B</i>	0.086	0.026	0.113
<i>B+</i>	0.113	0.029	0.170
<i>HS</i>	0.038	0.020	0.069
<i>HS+</i>	0.059	0.014	0.116

Table 3: Observed probabilities and margins

	CN	Dementia	MCI	<b>Row sums</b>
<i>ADV</i>	0.059	0.011	0.081	<b>0.150</b>
<i>B</i>	0.086	0.026	0.113	<b>0.224</b>
<i>B+</i>	0.113	0.029	0.170	<b>0.311</b>
<i>HS</i>	0.038	0.020	0.069	<b>0.126</b>
<i>HS+</i>	0.059	0.014	0.116	<b>0.188</b>
<b>Column sums</b>	<b>0.353</b>	<b>0.098</b>	<b>0.549</b>	

Table 4: Expected probabilities and margins

	CN	Dementia	MCI	<b>Row sums</b>
<i>ADV</i>	0.053	0.015	0.083	<b>0.150</b>
<i>B</i>	0.079	0.022	0.123	<b>0.224</b>
<i>B+</i>	0.110	0.030	0.171	<b>0.311</b>
<i>HS</i>	0.045	0.012	0.069	<b>0.126</b>
<i>HS+</i>	0.066	0.018	0.103	<b>0.188</b>
<b>Column sums</b>	<b>0.353</b>	<b>0.098</b>	<b>0.549</b>	

Table 5: Deviations: Observed - Expected

	CN	Dementia	MCI
<i>ADV</i>	0.006	-0.004	-0.001
<i>B</i>	0.007	0.004	-0.010
<i>B+</i>	0.003	-0.002	-0.001
<i>HS</i>	-0.007	0.007	0.000
<i>HS+</i>	-0.008	-0.005	0.013

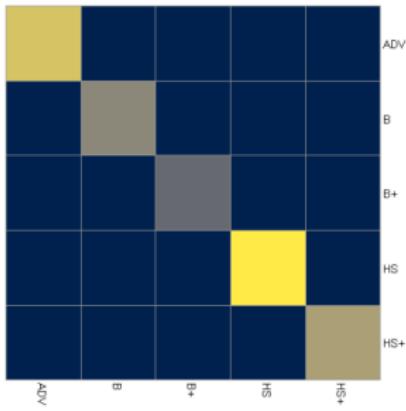
Table 6: Row constraints (inverse row margins)

	ADV	B	B+	HS	HS+
<i>ADV</i>	6.65	0.000	0.000	0.000	0.00
<i>B</i>	0.00	4.463	0.000	0.000	0.00
<i>B+</i>	0.00	0.000	3.213	0.000	0.00
<i>HS</i>	0.00	0.000	0.000	7.917	0.00
<i>HS+</i>	0.00	0.000	0.000	0.000	5.32

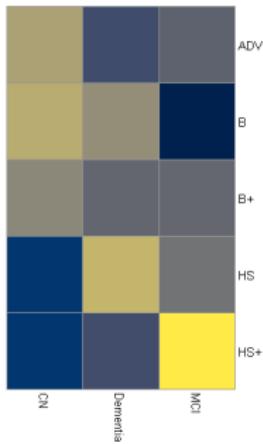
Table 7: Column constraints (inverse column margins)

	CN	Dementia	MCI
<i>CN</i>	2.83	0.000	0.000
<i>Dementia</i>	0.00	10.231	0.000
<i>MCI</i>	0.00	0.000	1.822

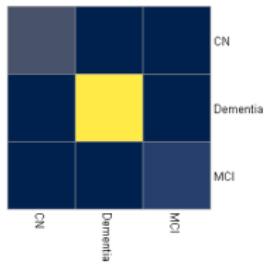
# What CA needs



**R:** Row constraints  
(inverse row probabilities)



**Z:** Deviations



**C:** Column constraints  
(inverse column probabilities)

- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )

- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )
- ▶ Uses but generalizes the SVD

- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )
- ▶ Uses but generalizes the SVD
  - ▶ Uses row & column weights (constraints)

- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )
- ▶ Uses but generalizes the SVD
  - ▶ Uses row & column weights (constraints)
- ▶ Gives back

- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )
- ▶ Uses but generalizes the SVD
  - ▶ Uses row & column weights (constraints)
- ▶ Gives back
  - ▶ Component (factor) scores

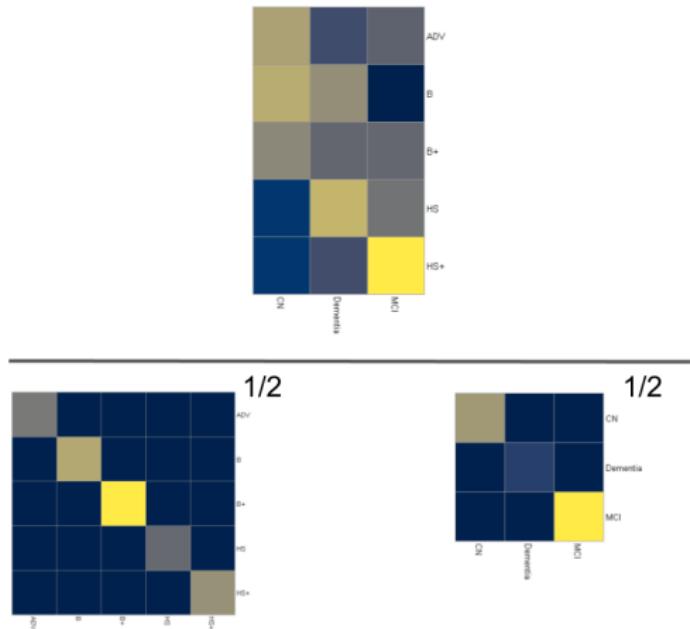
- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )
- ▶ Uses but generalizes the SVD
  - ▶ Uses row & column weights (constraints)
- ▶ Gives back
  - ▶ Component (factor) scores
  - ▶ Eigenvalues, singular values, & singular vectors

- ▶ GSVD( $\mathbf{R}$ ,  $\mathbf{X}$ ,  $\mathbf{C}$ )
- ▶ Uses but generalizes the SVD
  - ▶ Uses row & column weights (constraints)
- ▶ Gives back
  - ▶ Component (factor) scores
  - ▶ Eigenvalues, singular values, & singular vectors
  - ▶ *Generalized* singular vectors

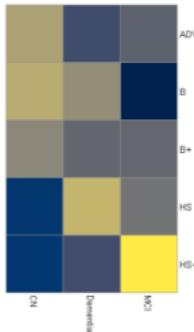
# What we really decompose



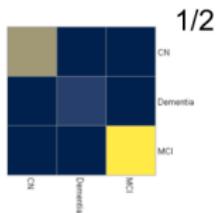
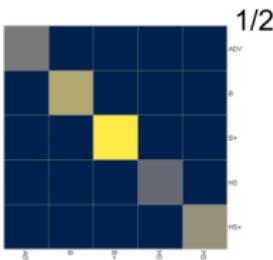
=



- A rectangle
- Deviations: Observed - Expected
  - Expected from Observed's margins



- Two squares
- Row margins and column margins



**Z**

$\overline{\mathbf{R}^{\frac{1}{2}}\mathbf{C}^{\frac{1}{2}}}$

$$\frac{(O - E)}{E^{\frac{1}{2}}}$$

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

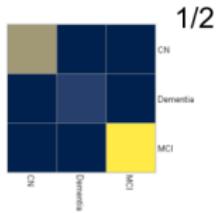
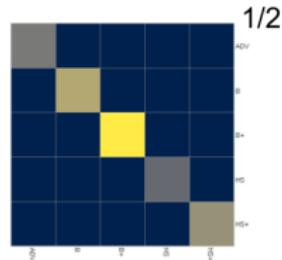
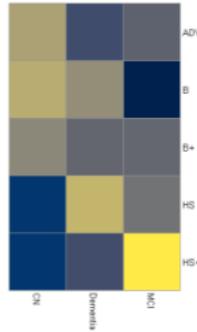
## CA's secrets

```
EDU <- amerge_subset$PTEDUCAT
DX <- amerge_subset$DX
edu_dx_table <- table(EDU, DX)

chisq.test(edu_dx_table)

##
## Pearson's Chi-squared test
##
## data: edu_dx_table
## X-squared = 8.648, df = 8, p-value = 0.3729
edu_dx_ca <- epCA(edu_dx_table, graphs = F)
sum(edu_dx_ca$ExPosition.Data$eigs) * sum(edu_dx_table)

## [1] 8.647979
```



Besides  $\chi^2$  this looks really familiar. What else are rectangles over squares?

$$r = \frac{cov(\mathbf{x}, \mathbf{y})}{\sigma_{\mathbf{x}} \times \sigma_{\mathbf{y}}}$$

## More of CA's secrets

- ▶ CA generalizes canonical correlation analysis (CCA)

## More of CA's secrets

- ▶ CA generalizes canonical correlation analysis (CCA)
- ▶ CA is the CCA between two *nominal* tables

## More of CA's secrets

- ▶ CA generalizes canonical correlation analysis (CCA)
- ▶ CA is the CCA between two *nominal* tables
- ▶ How do we create a contingency table?

## Nominal data

# NOMINAL

## UNORDERED DESCRIPTIONS



<b>EDU</b>	<b>DX</b>
B	Dementia
B	MCI
B+	Dementia
HS	Dementia
B+	CN

<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>
1	0	0	0	0
1	0	0	0	0
0	1	0	0	0
0	0	0	0	1
0	1	0	0	0

<b>MCI</b>	<b>CN</b>	<b>Dementia</b>
0	0	1
1	0	0
0	0	1
0	0	1
0	1	0

	<b>CN</b>	<b>Dementia</b>	<b>MCI</b>		<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>		<b>MCI</b>	<b>CN</b>	<b>Dementia</b>
<i>ADV</i>	39	7	54		1	0	0	0	0		0	0	1
<i>B</i>	57	17	75		1	0	0	0	0		1	0	0
<i>B+</i>	75	19	113		0	1	0	0	0		0	0	1
<i>HS</i>	25	13	46		0	0	0	0	1		0	0	1
<i>HS+</i>	39	9	77		0	1	0	0	0		0	1	0

<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>
1	0	0	0	0
1	0	0	0	0
0	1	0	0	0
0	0	0	0	1
0	1	0	0	0

<b>MCI</b>	<b>CN</b>	<b>Dementia</b>
0	0	1
1	0	0
0	0	1
0	0	1
0	1	0

## How to analyze nominal data?

- ▶ “coding categorical variables with the indicator matrix of dummy variables and considering them as Gaussian, for instance, is almost a crime.” in *Jan de Leeuw and the French School of Data Analysis* (Husson, Josse, Saporta)

## How to analyze nominal data?

- ▶ “coding categorical variables with the indicator matrix of dummy variables and considering them as Gaussian, for instance, is almost a crime.” in *Jan de Leeuw and the French School of Data Analysis* (Husson, Josse, Saporta)
- ▶ We *could* perform PCA on nominal data, but what would we get?

	<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>	<b>MCI</b>	<b>CN</b>	<b>Dementia</b>
<i>B</i>	1	-0.361	-0.226	-0.259	-0.204	-0.049	0.033	0.03
<i>B+</i>	-0.361	1	-0.283	-0.323	-0.256	-0.004	0.013	-0.013
<i>ADV</i>	-0.226	-0.283	1	-0.202	-0.16	-0.008	0.032	-0.039
<i>HS+</i>	-0.259	-0.323	-0.202	1	-0.183	0.065	-0.042	-0.042
<i>HS</i>	-0.204	-0.256	-0.16	-0.183	1	-0.001	-0.044	0.073
<i>MCI</i>	-0.049	-0.004	-0.008	0.065	-0.001	1	-0.815	-0.363
<i>CN</i>	0.033	0.013	0.032	-0.042	-0.044	-0.815	1	-0.243
<i>Dementia</i>	0.03	-0.013	-0.039	-0.042	0.073	-0.363	-0.243	1

	<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>	<b>MCI</b>	<b>CN</b>	<b>Dementia</b>
<i>B</i>	149	0	0	0	0	75	57	17
<i>B+</i>	0	207	0	0	0	113	75	19
<i>ADV</i>	0	0	100	0	0	54	39	7
<i>HS+</i>	0	0	0	125	0	77	39	9
<i>HS</i>	0	0	0	0	84	46	25	13
<i>MCI</i>	75	113	54	77	46	365	0	0
<i>CN</i>	57	75	39	39	25	0	235	0
<i>Dementia</i>	17	19	7	9	13	0	0	65

## Multiple correspondence analysis

## Multiple correspondence analysis

- ▶ Two perspectives:

## Multiple correspondence analysis

- ▶ Two perspectives:
  - ▶ *Weighted PCA* for nominal data

## Multiple correspondence analysis

- ▶ Two perspectives:
  - ▶ *Weighted PCA* for nominal data
  - ▶ Generalized CA for N-way contingency tables

## Multiple correspondence analysis

- ▶ Two perspectives:
  - ▶ *Weighted PCA* for nominal data
  - ▶ Generalized CA for N-way contingency tables
- ▶ So much more than nominal

## We're diving in

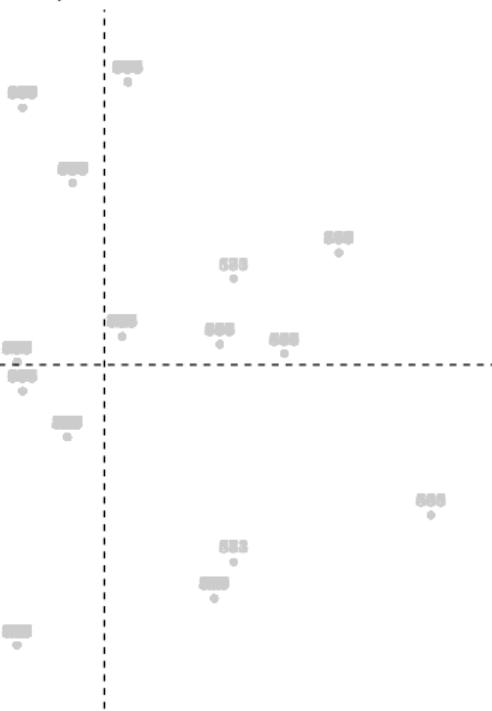
B	B+	ADV	HS+	HS	MCI	CN	Dementia
1	0	0	0	0	0	0	1
1	0	0	0	0	1	0	0
0	1	0	0	0	0	0	1
0	0	0	0	1	0	0	1
0	1	0	0	0	0	1	0

This is the kind of table we're analyzing. It has  $N = 665$ .

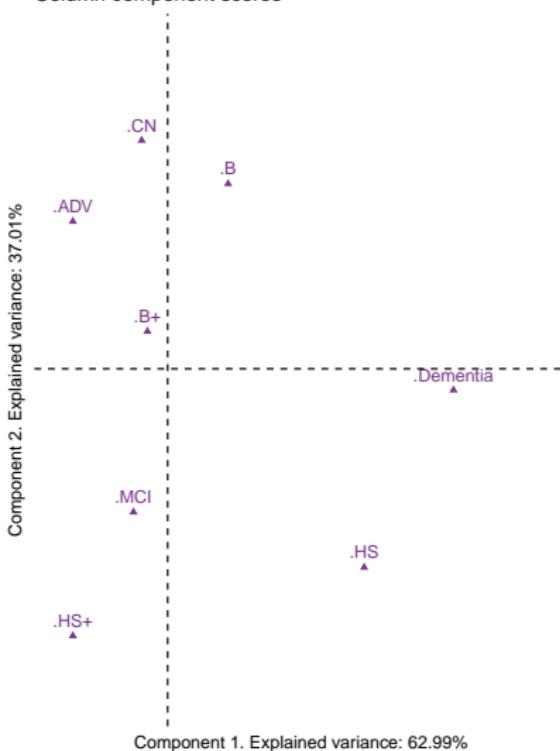
MCA:  
Row component scores

Component 2. Explained variance: 37.01%

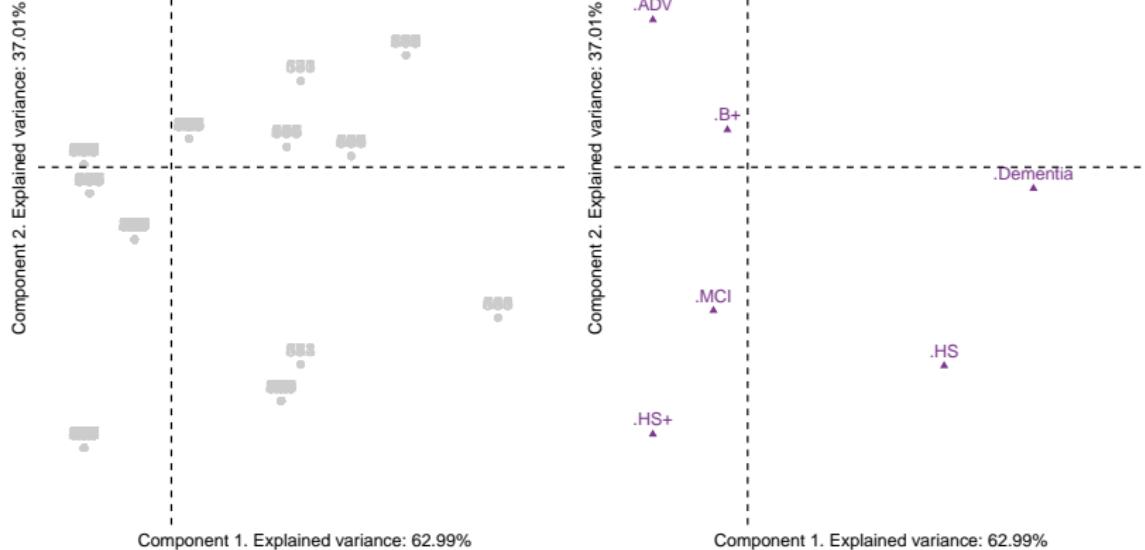
Component 1. Explained variance: 62.99%



MCA:  
Column component scores

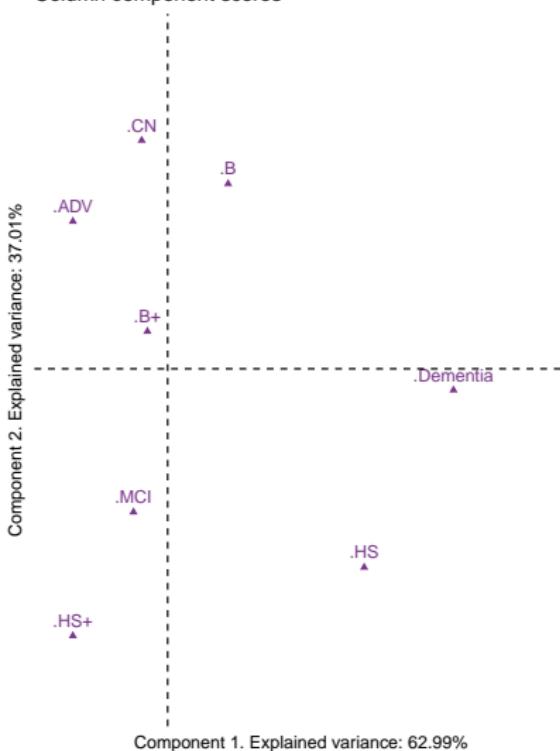


MCA:  
Row component scores



Does any of this look familiar?

MCA:  
Column component scores



CA:  
Row component scores

Component 2. Explained variance: 37.01%

CN

MCI

Dementia

Component 1. Explained variance: 62.99%

MCA:  
Column component scores

Component 2. Explained variance: 37.01%

.CN

.B+

.MCI

.HS+

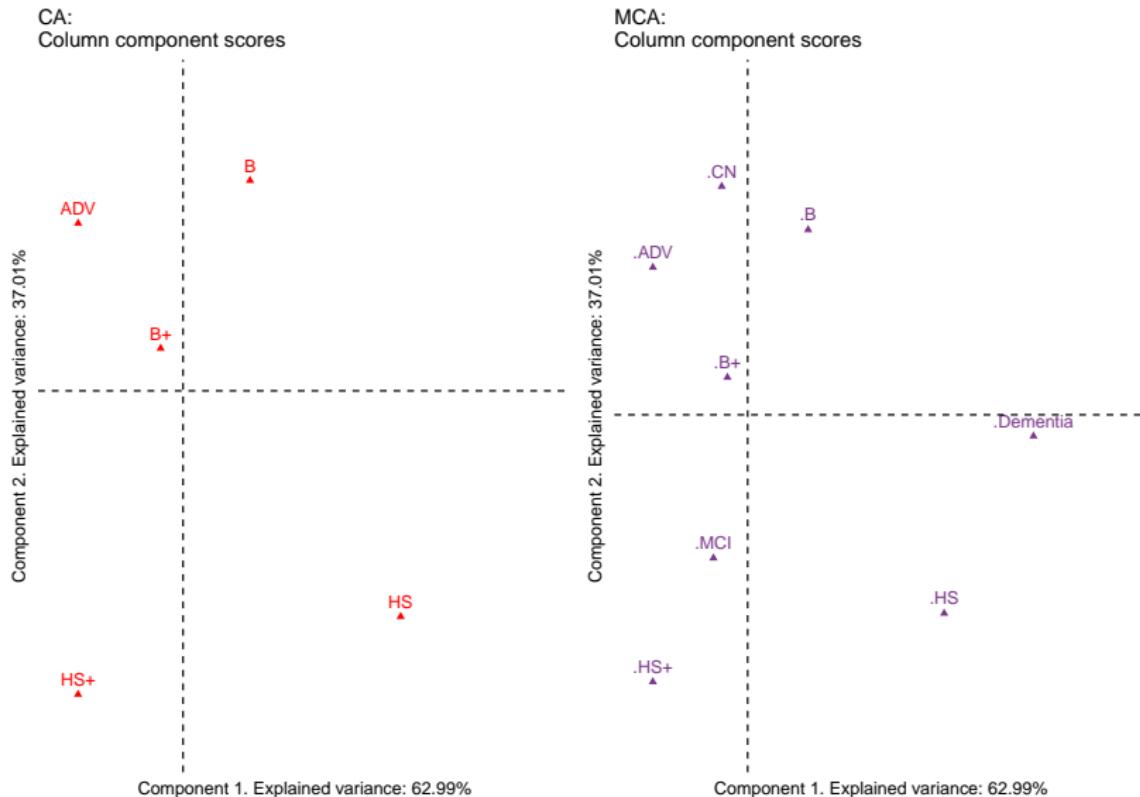
.ADV

.B

.Dementia

.HS

Component 1. Explained variance: 62.99%



# CA & MCA Magic!

	CN	Dementia	MCI
ADV	39	7	54
B	57	17	75
B+	75	19	113
HS	25	13	46
HS+	39	9	77

B	B+	ADV	HS+	HS	MCI	CN	Dementia
1	0	0	0	0	0	0	1
1	0	0	0	0	1	0	0
0	1	0	0	0	0	0	1
0	0	0	0	1	0	0	1
0	1	0	0	0	0	1	0

Same technique on two *different* tables: same result

## Scaling up

- ▶ Let's bring in ApoE

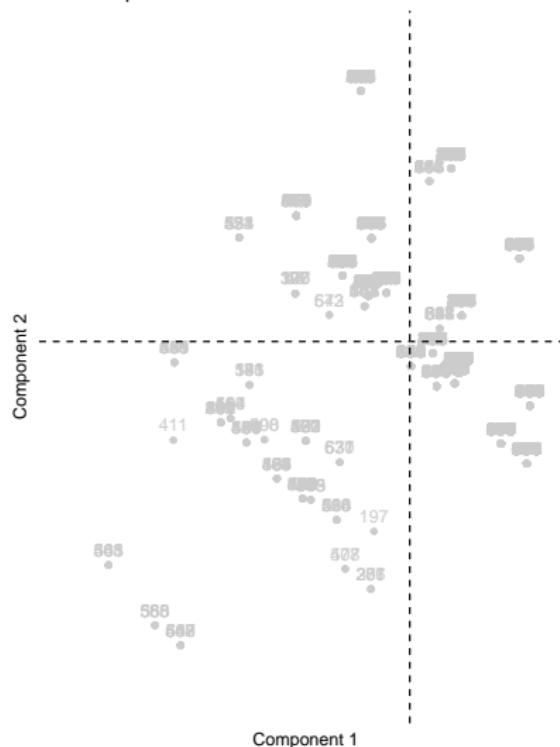
## Scaling up

- ▶ Let's bring in ApoE
- ▶ It has 3 levels: 0 copy, 1 copy, 2 copies

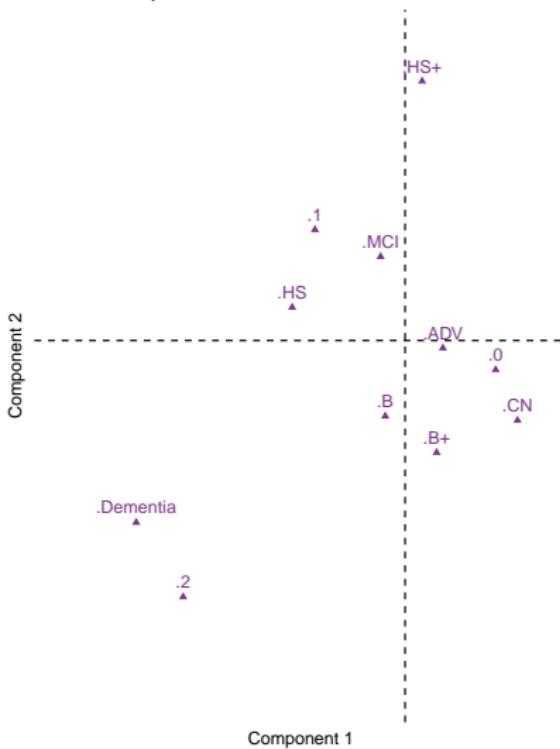
<b>EDU</b>	<b>DX</b>	<b>APOE</b>
B	Dementia	2
B	MCI	0
B+	Dementia	2
HS	Dementia	2
B+	CN	0

<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>	<b>MCI</b>	<b>CN</b>	<b>Dementia</b>	<b>0</b>	<b>1</b>	<b>2</b>
1	0	0	0	0	0	0	1	0	0	1
1	0	0	0	0	1	0	0	1	0	0
0	1	0	0	0	0	0	1	0	0	1
0	0	0	0	1	0	0	1	0	0	1
0	1	0	0	0	0	1	0	1	0	0

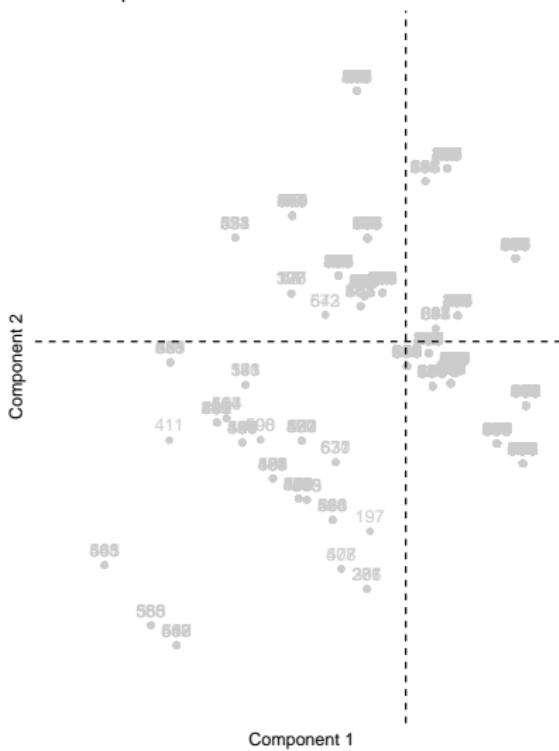
MCA:  
Row component scores



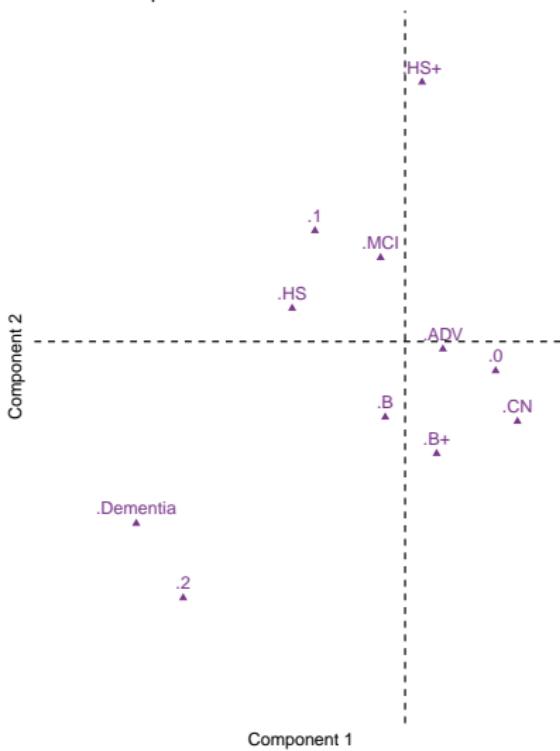
MCA:  
Column component scores



MCA:  
Row component scores



MCA:  
Column component scores



# Crisp vs. fuzzy coding

EDU	DX	APOE
B	Dementia	2
B	MCI	0
B+	Dementia	2
HS	Dementia	2
B+	CN	0

B	B+	ADV	HS+	HS	MCI	CN	Dementia	0	1	2
1	0	0	0	0	0	0	1	0	0	1
1	0	0	0	0	1	0	0	1	0	0
0	1	0	0	0	0	0	1	0	0	1
0	0	0	0	1	0	0	1	0	0	1
0	1	0	0	0	0	1	0	1	0	0

<b>EDU</b>	<b>DX</b>	<b>APOE</b>
B	Dementia	2
B	MCI	0
B+	Dementia	2
HS	Dementia	2
B+	CN	0

<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>	<b>MCI</b>	<b>CN</b>	<b>Dementia</b>	<b>0</b>	<b>1</b>	<b>2</b>
1	0	0	0	0	0	0	1	0	0	1
1	0	0	0	0	0.5	0	0.5	1	0	0
0	1	0	0	0	0	0	1	0	0	1
0	0	0	0	1	0	0	1	0	0	1
0	1	0	0	0	0	1	0	1	0	0

# Our first fuzzy friend

## ORDINAL

ORDERED DESCRIPTIONS



- ▶ Modified Hachinski

- ▶ Modified Hachinski
  - ▶ 0, 1, 2, 3 (in these data)

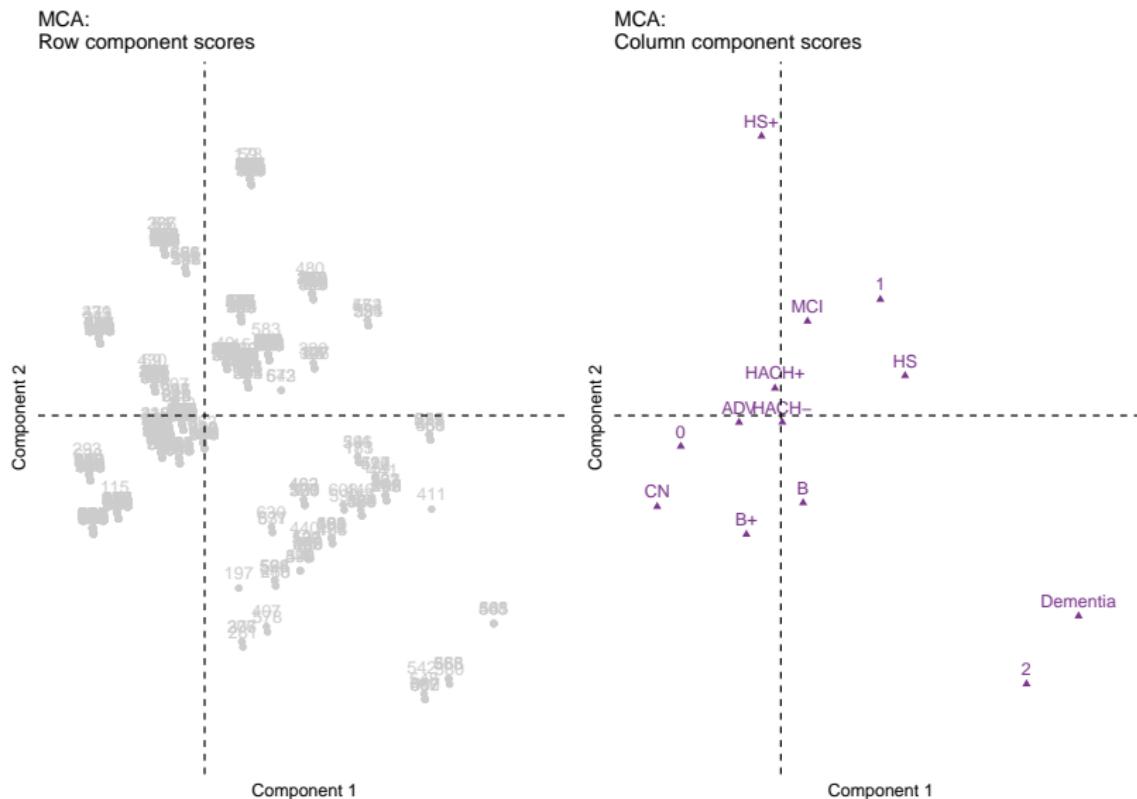
- ▶ Modified Hachinski
  - ▶ 0, 1, 2, 3 (in these data)
- ▶ Specific form of fuzzy coding: “bipolar”

HACH
0
3
1
2
1

HACH-	HACH+
1	0
0	1
0.667	0.333
0.333	0.667
0.667	0.333

<b>EDU</b>	<b>DX</b>	<b>APOE</b>	<b>HACH</b>
B	Dementia	2	1
B	MCI	0	1
B+	Dementia	2	0
HS	Dementia	2	0
B+	CN	0	0

<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>	<b>MCI</b>	<b>CN</b>	<b>Dementia</b>	<b>0</b>	<b>1</b>	<b>2</b>	<b>HACH-</b>	<b>HACH+</b>
1	0	0	0	0	0	0	1	0	0	1	0.667	0.333
1	0	0	0	0	1	0	0	1	0	0	0.667	0.333
0	1	0	0	0	0	0	1	0	0	1	1	0
0	0	0	0	1	0	0	1	0	0	1	1	0
0	1	0	0	0	0	1	0	1	0	0	1	0



Our second fuzzy friend

# CONTINUOUS

measured data, can have  $\infty$  values within possible range.



I AM 3.1" TALL

- ▶ Age: 55.00 - 89.60

- ▶ Age: 55.00 - 89.60
  - ▶ But we need to scale it (Z-score)

- ▶ Age: 55.00 - 89.60
  - ▶ But we need to scale it (Z-score)
- ▶ We use two columns again:

- ▶ Age: 55.00 - 89.60
  - ▶ But we need to scale it (Z-score)
- ▶ We use two columns again:
  - ▶  $\frac{(1-x)}{2}$  &  $+\frac{(1+x)}{2}$

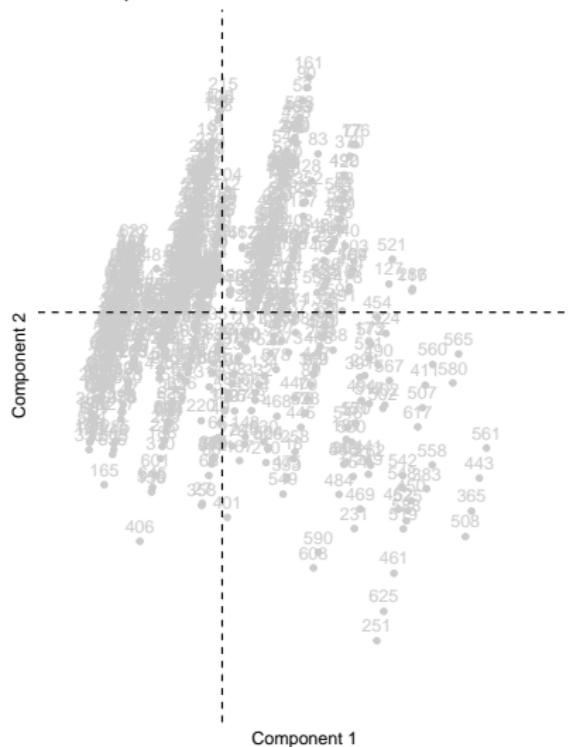
<b>AGE</b>	<b>AGE (Z)</b>
76.3	0.637
76.5	0.666
64.4	-1.095
62.9	-1.314
63.9	-1.168

<b>AGE-</b>	<b>AGE+</b>
0.181	0.819
0.167	0.833
1.048	-0.048
1.157	-0.157
1.084	-0.084

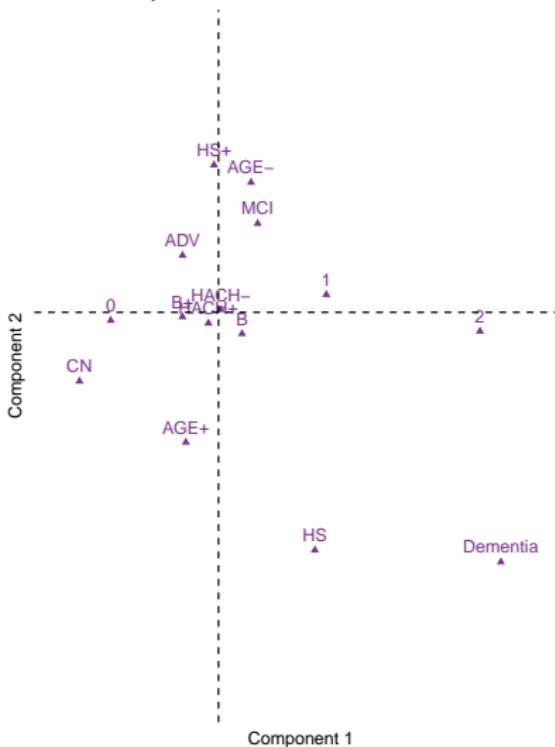
<b>EDU</b>	<b>DX</b>	<b>APOE</b>	<b>HACH</b>	<b>AGE</b>
B	Dementia	2	1	76.3
B	MCI	0	1	76.5
B+	Dementia	2	0	64.4
HS	Dementia	2	0	62.9
B+	CN	0	0	63.9

<b>B</b>	<b>B+</b>	<b>ADV</b>	<b>HS+</b>	<b>HS</b>	<b>MCI</b>	<b>CN</b>	<b>Dementia</b>	<b>0</b>	<b>1</b>	<b>2</b>	<b>HACH-</b>	<b>HACH+</b>	<b>AGE-</b>	<b>AGE+</b>
1	0	0	0	0	0	0	1	0	0	1	0.667	0.333	0.181	0.819
1	0	0	0	0	1	0	0	1	0	0	0.667	0.333	0.167	0.833
0	1	0	0	0	0	0	1	0	0	1	1	0	1.048	-0.048
0	0	0	0	1	0	0	1	0	0	1	1	0	1.157	-0.157
0	1	0	0	0	0	1	0	1	0	0	1	0	1.084	-0.084

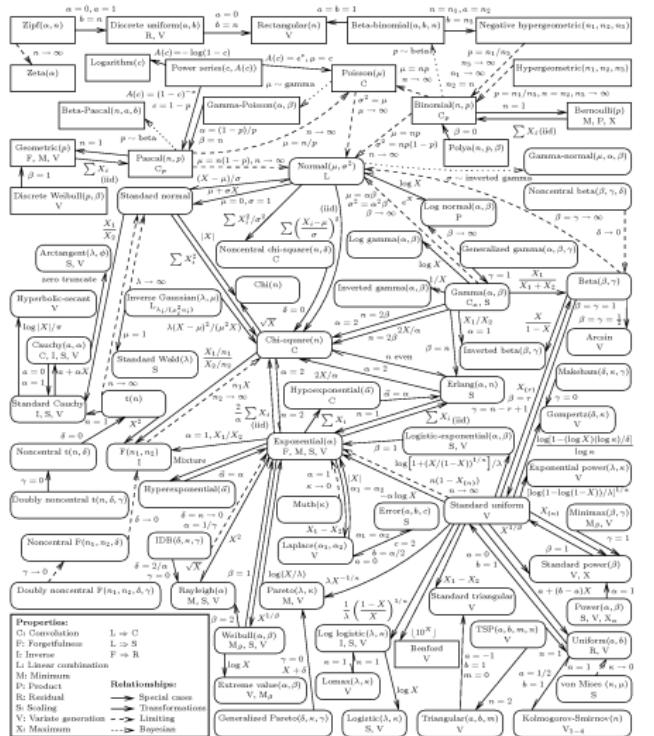
MCA:  
Row component scores



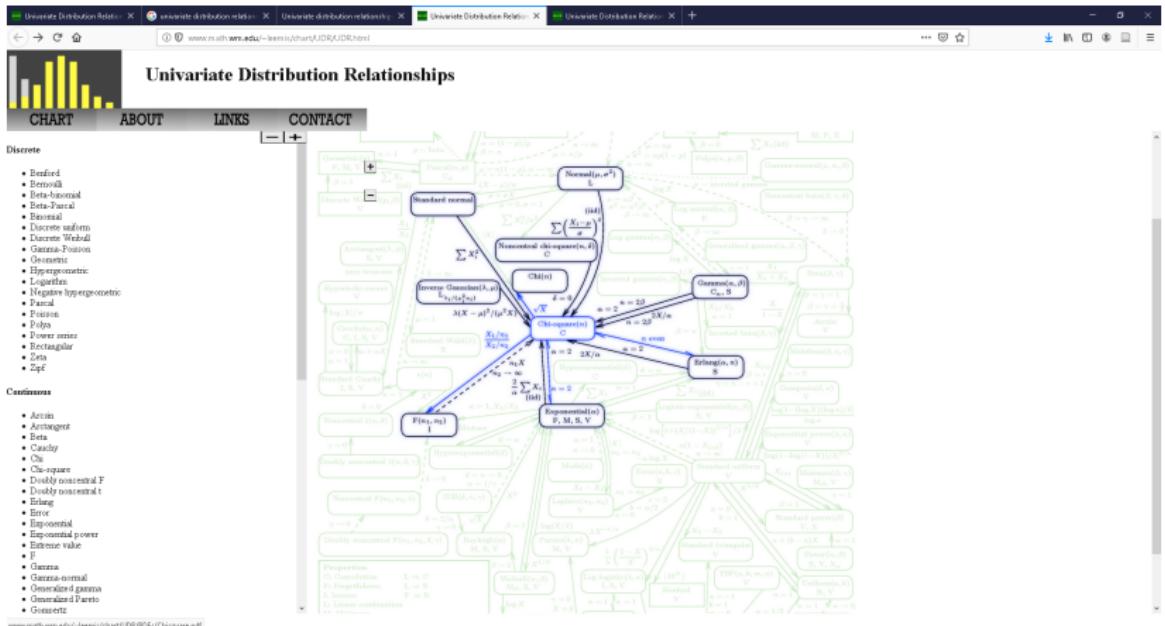
MCA:  
Column component scores



How is this magic possible?!



See here



[www.mathservice.edu/~leemis/chart/UDR/UDR.html](http://www.mathservice.edu/~leemis/chart/UDR/UDR.html)

See here

## Conclusions

- ▶ There's so much I'm not telling you

## Conclusions

- ▶ There's so much I'm not telling you
  - ▶ I wish I could!

## Conclusions

- ▶ There's so much I'm not telling you
  - ▶ I wish I could!
- ▶ But that's kind of the point

## Conclusions

- ▶ There's so much I'm not telling you
  - ▶ I wish I could!
- ▶ But that's kind of the point
  - ▶ CA is a *massive* world

## Conclusions

- ▶ There's so much I'm not telling you
  - ▶ I wish I could!
- ▶ But that's kind of the point
  - ▶ CA is a *massive* world
  - ▶ Solves many problems we typically ignore

# Conclusions

- ▶ There's so much I'm not telling you
  - ▶ I wish I could!
- ▶ But that's kind of the point
  - ▶ CA is a *massive* world
  - ▶ Solves many problems we typically ignore
- ▶ Learn to recognize data types

# Conclusions

- ▶ There's so much I'm not telling you
  - ▶ I wish I could!
- ▶ But that's kind of the point
  - ▶ CA is a *massive* world
  - ▶ Solves many problems we typically ignore
- ▶ Learn to recognize data types
  - ▶ Learn what to do with them

Some many bonuses!

## A workshop

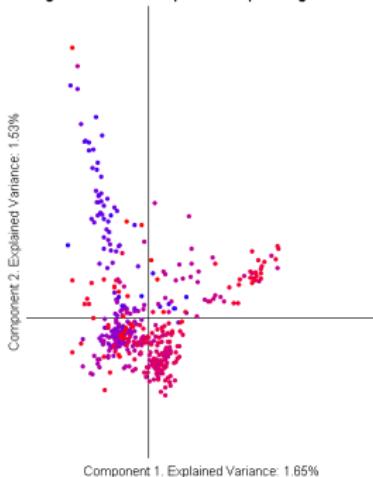
- ▶ [https://github.com/derekbeaton/Workshops/tree/master/R\\_TC/PCA\\_MCA\\_Resampling](https://github.com/derekbeaton/Workshops/tree/master/R_TC/PCA_MCA_Resampling)

## A workshop

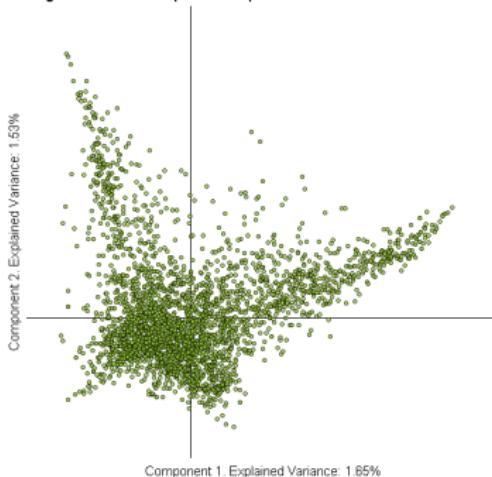
- ▶ [https://github.com/derekbeaton/Workshops/tree/master/R\\_TC/PCA\\_MCA\\_Resampling](https://github.com/derekbeaton/Workshops/tree/master/R_TC/PCA_MCA_Resampling)
  - ▶ Extraordinary detail on all of this

# The Mueller report

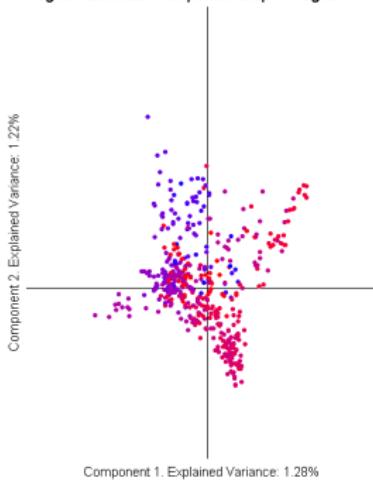
Pages x Words. Component Map of Pages.



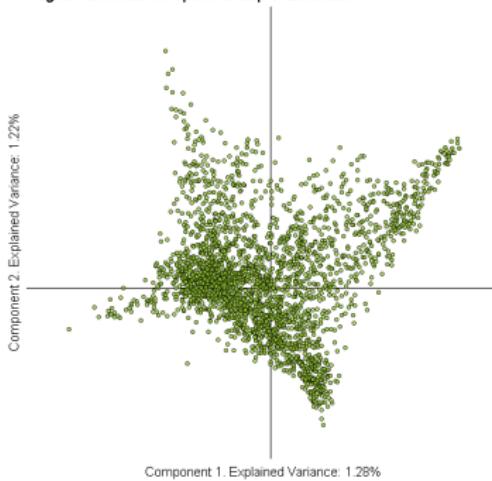
Pages x Words. Component Map of Words.



Pages x Lemmas. Component Map of Pages



Pages x Lemmas Component Map of Lemmas.



# The Marvel Cinematic Universe

- ▶ Actually super cool

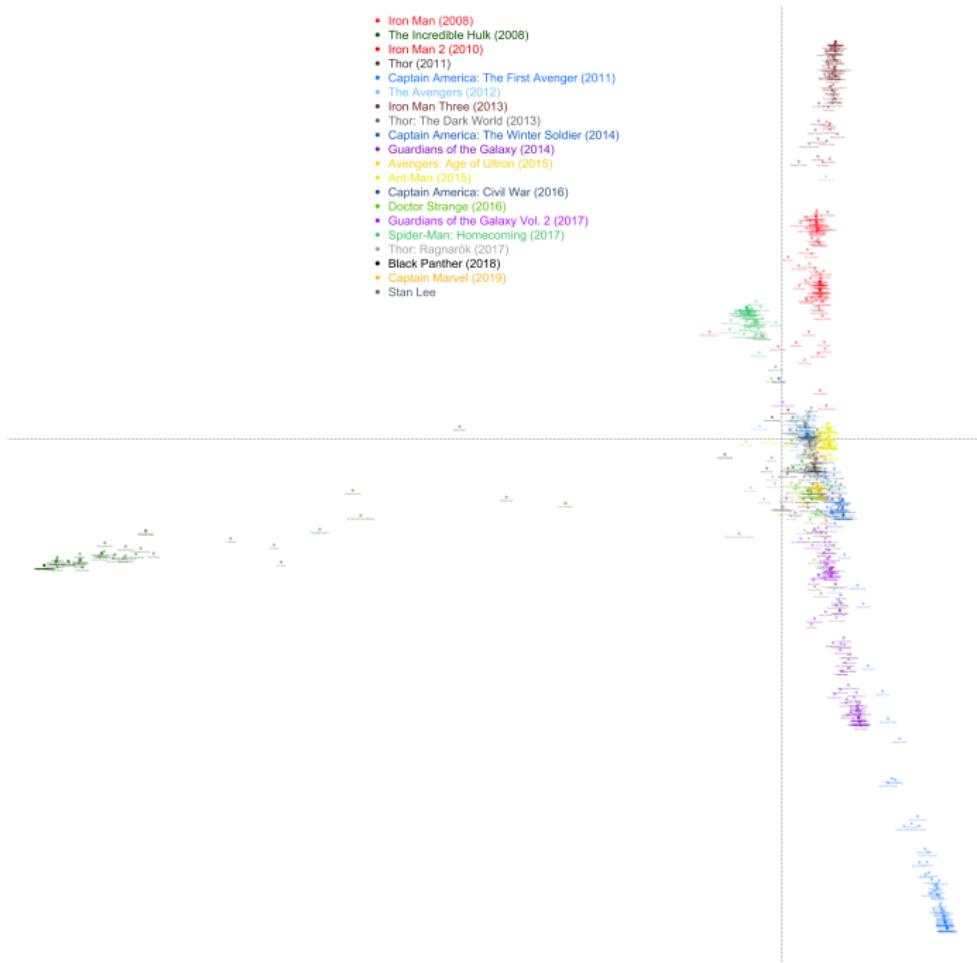
## The Marvel Cinematic Universe

- ▶ Actually super cool
- ▶ CA naturally handles networks

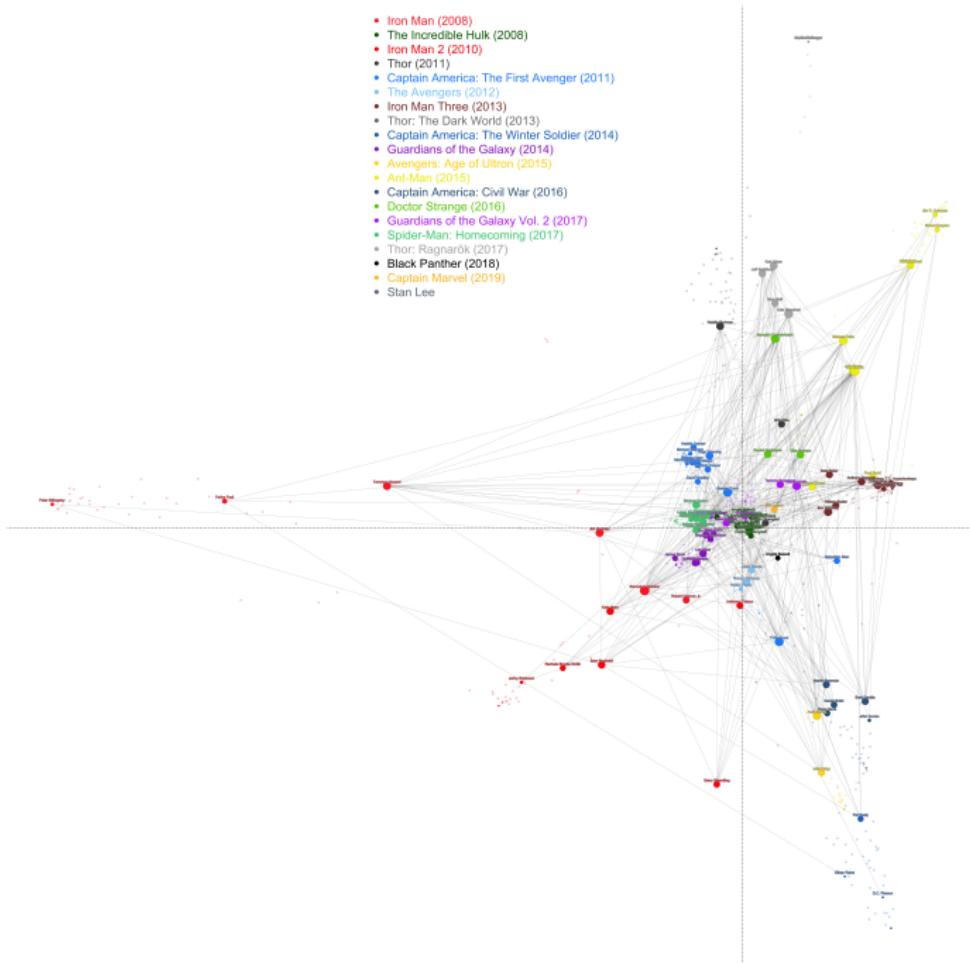
## The Marvel Cinematic Universe

- ▶ Actually super cool
- ▶ CA naturally handles networks
- ▶ [https://github.com/derekbeaton/Marvel-Cinematic-Universe\\_Network/](https://github.com/derekbeaton/Marvel-Cinematic-Universe_Network/)

- Iron Man (2008)
- The Incredibile Hulk (2008)
- Iron Man 2 (2010)
- Thor (2011)
- Captain America: The First Avenger (2011)
- The Avengers (2012)
- Iron Man Three (2013)
- Thor: The Dark World (2013)
- Captain America: The Winter Soldier (2014)
- Guardians of the Galaxy (2014)
- Avengers: Age of Ultron (2015)
- Ant-Man (2015)
- Captain America: Civil War (2016)
- Doctor Strange (2016)
- Guardians of the Galaxy Vol. 2 (2017)
- Spider-Man: Homecoming (2017)
- Thor: Ragnarök (2017)
- Black Panther (2018)
- Captain Marvel (2019)
- Stan Lee



- Iron Man (2008)
- The Incredible Hulk (2008)
- Iron Man 2 (2010)
- Thor (2011)
- Captain America: The First Avenger (2011)
- The Avengers (2012)
- Iron Man Three (2013)
- Thor: The Dark World (2013)
- Captain America: The Winter Soldier (2014)
- Guardians of the Galaxy (2014)
- Avengers: Age of Ultron (2015)
- Ant-Man (2015)
- Captain America: Civil War (2016)
- Doctor Strange (2016)
- Guardians of the Galaxy Vol. 2 (2017)
- Spider-Man: Homecoming (2017)
- Thor: Ragnarök (2017)
- Black Panther (2018)
- Captain Marvel (2019)
- Stan Lee



## GMCD & ours

- ▶ Generalized minimum covariance determinant

## GMCD & ours

- ▶ Generalized minimum covariance determinant
- ▶ ours (another R package)

## GMCD & ours

- ▶ Generalized minimum covariance determinant
- ▶ ours (another R package)
  - ▶ New package for outliers

## GMCD & ours

- ▶ Generalized minimum covariance determinant
- ▶ ours (another R package)
  - ▶ New package for outliers
  - ▶ Has some important bells-and-whistles

## GMCD & ours

- ▶ Generalized minimum covariance determinant
- ▶ ours (another R package)
  - ▶ New package for outliers
  - ▶ Has some important bells-and-whistles
  - ▶ <https://github.com/derekbeaton/ours>

# GPLS

- ▶ First: a PLS for mixed data types

# GPLS

- ▶ First: a PLS for mixed data types
  - ▶ Including those not discussed here

# GPLS

- ▶ First: a PLS for mixed data types
  - ▶ Including those not discussed here
- ▶ Second: unify the “two-table” techniques

# GPLS

- ▶ First: a PLS for mixed data types
  - ▶ Including those not discussed here
- ▶ Second: unify the “two-table” techniques
  - ▶ PLS, CCA, RRR/RDA

# GPLS

- ▶ First: a PLS for mixed data types
  - ▶ Including those not discussed here
- ▶ Second: unify the “two-table” techniques
  - ▶ PLS, CCA, RRR/RDA
- ▶ Package & preprint

# GPLS

- ▶ First: a PLS for mixed data types
  - ▶ Including those not discussed here
- ▶ Second: unify the “two-table” techniques
  - ▶ PLS, CCA, RRR/RDA
- ▶ Package & preprint
  - ▶ <https://github.com/derekbeaton/gpls>

# GPLS

- ▶ First: a PLS for mixed data types
  - ▶ Including those not discussed here
- ▶ Second: unify the “two-table” techniques
  - ▶ PLS, CCA, RRR/RDA
- ▶ Package & preprint
  - ▶ <https://github.com/derekbeaton/gpls>
  - ▶ Github issues where I routinely call myself a “dummy”

# ExPosition

- ▶ ExPosition

# ExPosition

- ▶ ExPosition
  - ▶ Family of packages

## ExPosition

- ▶ ExPosition
  - ▶ Family of packages
  - ▶ Includes resampling

# ExPosition

- ▶ ExPosition
  - ▶ Family of packages
  - ▶ Includes resampling
  - ▶ Lots of PCA & CA techniques

## Some alternatives

- ▶ FactoMineR

## Some alternatives

- ▶ FactoMineR
- ▶ ade4

## Some alternatives

- ▶ FactoMineR
- ▶ ade4
- ▶ ca

## Some alternatives

- ▶ FactoMineR
- ▶ ade4
- ▶ ca
- ▶ MASS

## Some alternatives

- ▶ FactoMineR
- ▶ ade4
- ▶ ca
- ▶ MASS
- ▶ psych

## Some alternatives

- ▶ FactoMineR
- ▶ ade4
- ▶ ca
- ▶ MASS
- ▶ psych
- ▶ factoextra (visualization)

## Some alternatives

- ▶ FactoMineR
- ▶ ade4
- ▶ ca
- ▶ MASS
- ▶ psych
- ▶ factoextra (visualization)
- ▶ So many others

## (Some) References

## Expansions & data details

- ▶ Beaton, D., Saporta, G., Abdi, H., & Alzheimer's Disease Neuroimaging Initiative. (2019). A generalization of partial least squares regression and correspondence analysis for categorical and mixed data: An application with the ADNI data. bioRxiv, 598888.

## Expansions & data details

- ▶ Beaton, D., Saporta, G., Abdi, H., & Alzheimer's Disease Neuroimaging Initiative. (2019). A generalization of partial least squares regression and correspondence analysis for categorical and mixed data: An application with the ADNI data. bioRxiv, 598888.
- ▶ Beaton, D., Sunderland, K. M., Levine, B., Mandzia, J., Masellis, M., Swartz, R. H., ... & Strother, S. C. (2019). Generalization of the minimum covariance determinant algorithm for categorical and mixed data types. bioRxiv, 333005.

And these

- ▶ Beaton, D., Fatt, C. R. C., & Abdi, H. (2014). An ExPosition of multivariate analysis with the singular value decomposition in R. *Computational Statistics & Data Analysis*, 72, 176-189.

## And these

- ▶ Beaton, D., Fatt, C. R. C., & Abdi, H. (2014). An ExPosition of multivariate analysis with the singular value decomposition in R. *Computational Statistics & Data Analysis*, 72, 176-189.
- ▶ Beaton, D., Dunlop, J., & Abdi, H. (2016). Partial least squares correspondence analysis: A framework to simultaneously analyze behavioral and genetic data. *Psychological methods*, 21(4), 621.

## Techniques

- ▶ Greenacre, M. (2017). Correspondence analysis in practice. CRC press.

## Techniques

- ▶ Greenacre, M. (2017). Correspondence analysis in practice. CRC press.
- ▶ Greenacre, M. J. (1984). Theory and Applications of Correspondence Analysis. Retrieved from <http://books.google.com/books?id=LsPaAAAAMAAJ>

## Techniques

- ▶ Greenacre, M. J. (2010). Correspondence analysis. Wiley Interdisciplinary Reviews: Computational Statistics, 2(5), 613–619. <https://doi.org/10.1002/wics.114>

## Techniques

- ▶ Greenacre, M. J. (2010). Correspondence analysis. Wiley Interdisciplinary Reviews: Computational Statistics, 2(5), 613–619. <https://doi.org/10.1002/wics.114>
- ▶ Lebart, L., Morineau, A., & Warwick, K. M. (1984). Multivariate descriptive statistical analysis: correspondence analysis and related techniques for large matrices. Wiley.

## Techniques

- ▶ Greenacre, M. J. (2010). Correspondence analysis. Wiley Interdisciplinary Reviews: Computational Statistics, 2(5), 613–619. <https://doi.org/10.1002/wics.114>
- ▶ Lebart, L., Morineau, A., & Warwick, K. M. (1984). Multivariate descriptive statistical analysis: correspondence analysis and related techniques for large matrices. Wiley.
- ▶ Nguyen, L. H., & Holmes, S. (2019). Ten quick tips for effective dimensionality reduction. PLOS Computational Biology, 15(6), e1006907.

## Data

- ▶ Escofier, B. (1978). Analyse factorielle et distances répondant au principe d'équivalence distributionnelle. *Revue de Statistique Appliquée*, 26(4), 29–37.

## Data

- ▶ Escofier, B. (1978). Analyse factorielle et distances répondant au principe d'équivalence distributionnelle. *Revue de Statistique Appliquée*, 26(4), 29–37.
- ▶ Escofier, B. (1979). Traitement simultané de variables qualitatives et quantitatives en analyse factorielle. *Cahiers de l'Analyse Des Données*, 4(2), 137–146.

## Data

- ▶ Escofier, B. (1978). Analyse factorielle et distances répondant au principe d'équivalence distributionnelle. *Revue de Statistique Appliquée*, 26(4), 29–37.
- ▶ Escofier, B. (1979). Traitement simultané de variables qualitatives et quantitatives en analyse factorielle. *Cahiers de l'Analyse Des Données*, 4(2), 137–146.
- ▶ Greenacre, M. (2014). Data Doubling and Fuzzy Coding. In J. Blasius & M. Greenacre (Eds.), *Visualization and Verbalization of Data* (pp. 239–253). Philadelphia, PA, USA: CRC Press.

## History

- ▶ Holmes S, Josse J. Discussion of “50 Years of Data Science”. Journal of Computational and Graphical Statistics. 2017, V26(4) 768-769. <https://www.tandfonline.com/doi/full/10.1080/10618600.2017.1385471>