

Simple & Multiple Correspondence Analyses

Contingency, categorical, ordinal, continuous and mixed data

Derek Beaton

Rotman Research Institute

October 28, 2019

Before we get started

Our new best friends

CONTINUOUS

measured data, can have ∞ 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



CONTINUOUS

measured data, can have ∞ values within possible range.



| AM 3.1" TALL
| WEIGH 34.16 grams

DISCRETE

observations can only exist at limited values, often counts.



I HAVE 8 LEGS
and
4 SPOTS!

@alison_horst

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



I am
EXTINCT!



HA

@alison_horst

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

CONTINUOUS

measured data, can have ∞ values within possible range.



| AM 3.1" TALL
| WEIGH 34.16 grams

DISCRETE

observations can only exist at limited values, often counts.



I HAVE 8 LEGS
and
4 SPOTS!

@alison_horst

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



I am
EXTINCT!



HA

@alison_horst

- ▶ 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 ∞ values within possible range.



| AM 3.1" TALL
| WEIGH 34.16 grams

DISCRETE

observations can only exist at limited values, often counts.



I HAVE 8 LEGS
and
4 SPOTS!

@alison_horst

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



I am
EXTINCT!



-HA

@alison_horst

- ▶ 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 oo values within possible range.



| AM 3.1" TALL
| WEIGH 34.16 grams

DISCRETE

observations can only exist at limited values, often counts.



I HAVE 8 LEGS
and
4 SPOTS!

@alison_hart

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



I am
EXTINCT!-



-HA

@alison_hart

- ▶ 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

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

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, Networks, Software

Revisiting 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

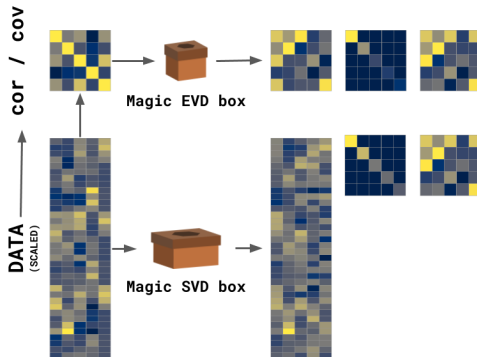
What is PCA for?

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

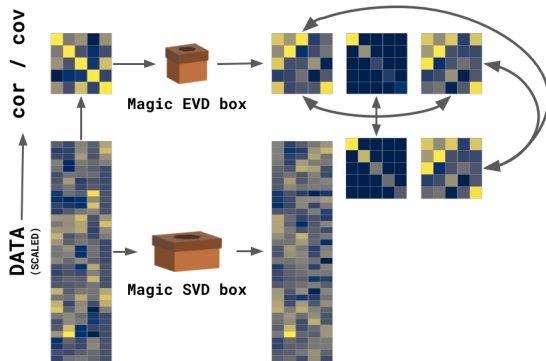
What is PCA for?

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

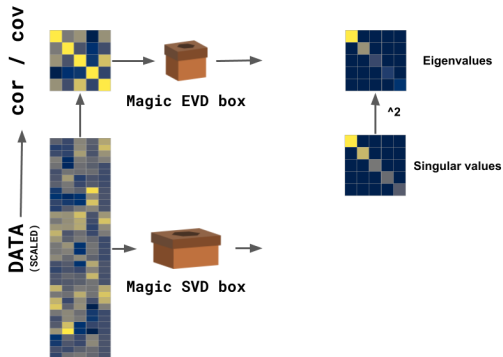
Eigen- and singular value decompositions



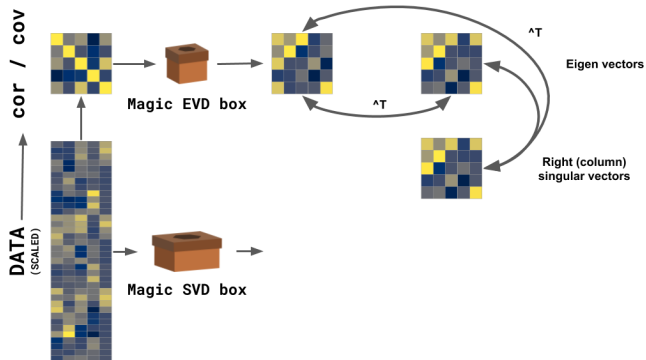
Eigen- and singular value decompositions



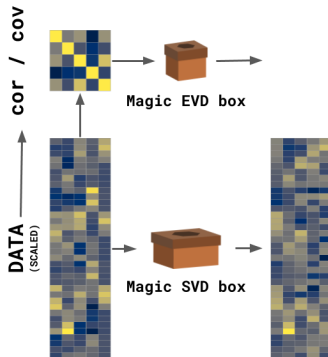
Eigen- and singular value decompositions



Eigen- and singular value decompositions



Eigen- and singular value decompositions



Left (row) singular
vectors

Some data

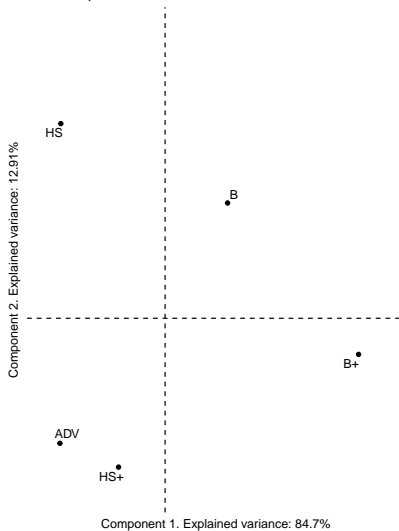
Diagnosis and education

	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

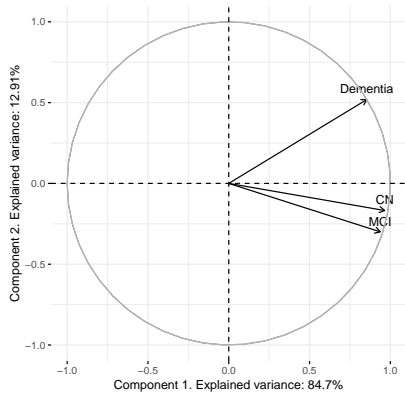
- ▶ Given a table, and asked for a multivariate analysis

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

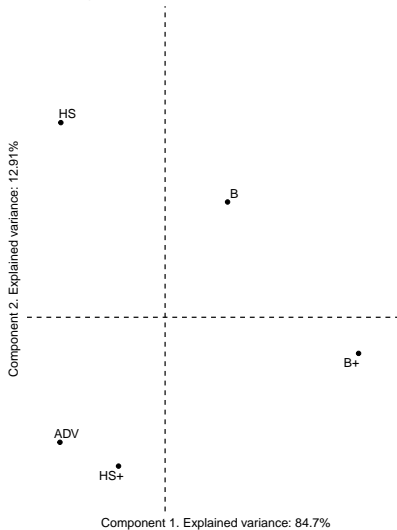
PCA:
Row component scores



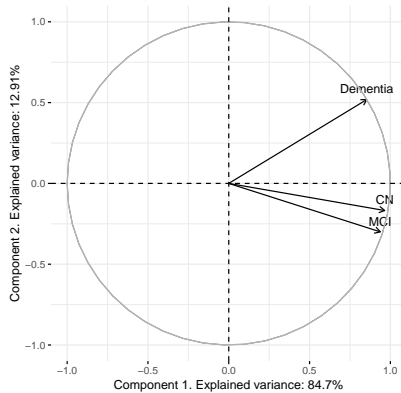
PCA:
Variable-Component Correlations



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

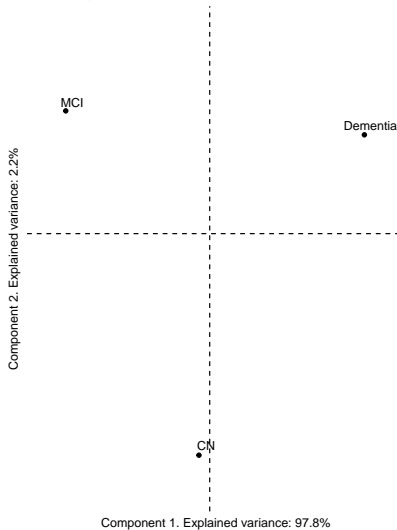
What did PCA detect?

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

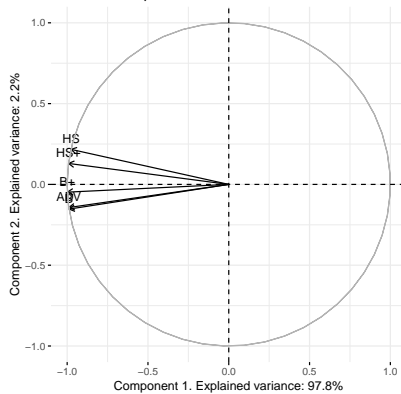
Let's try something different!

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

What is PCA for?

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

Let's take another look

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

- Tell me things about this matrix

Let's take another look

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

- ▶ 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

- ▶ See Lebart's History & Prehistory of CA
 - ▶ http://www.dtmvic.com/doc/About_the_History_of_CA.pdf
- ▶ And Beh & Lombardo's series

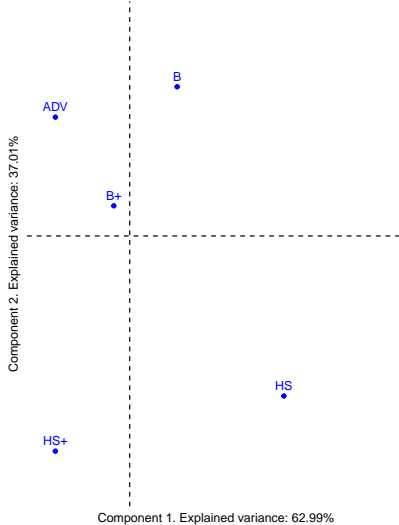
- ▶ See Lebart's History & Prehistory of CA
 - ▶ http://www.dtmvic.com/doc/About_the_History_of_CA.pdf
- ▶ And Beh & Lombardo's series
 - ▶ A geneaology of CA:
<https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1467-842X.2012.00676.x>

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

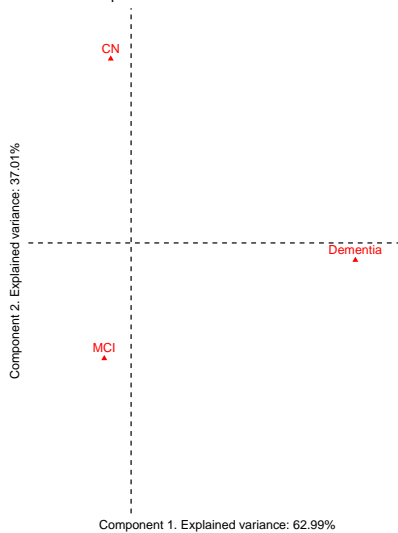
We're diving in

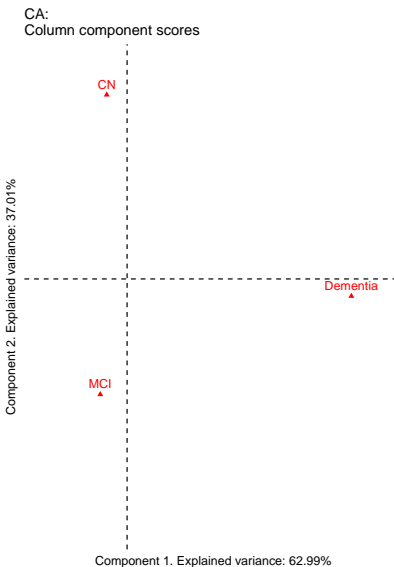
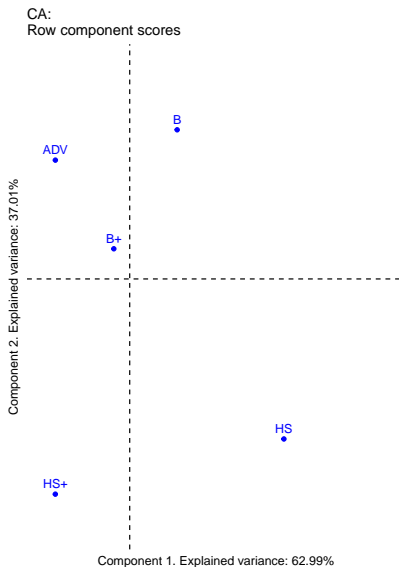
	CN	Dementia	MCI
<i>ADV</i>	39	7	54
<i>B</i>	57	17	75
<i>B+</i>	75	19	113
<i>HS</i>	25	13	46
<i>HS+</i>	39	9	77

CA:
Row component scores



CA:
Column component scores



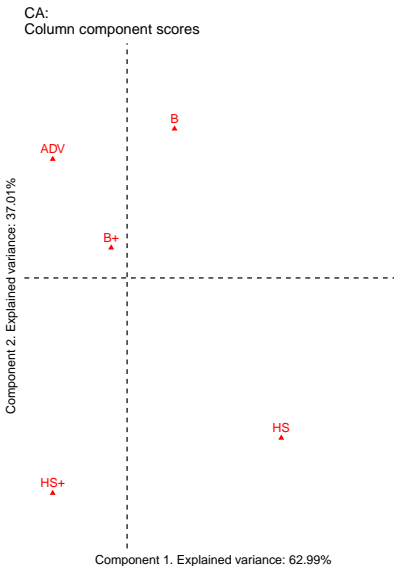
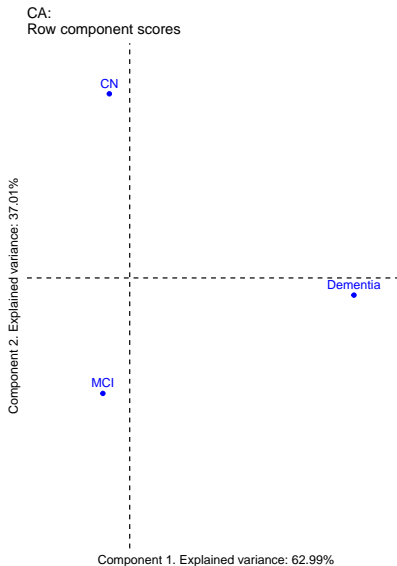


Want to see a cool trick?

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

Table 4: Expected probabilities and margins

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

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

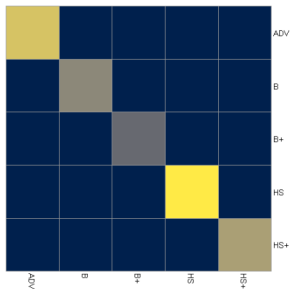
Table 6: Row constraints (inverse row margins)

	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

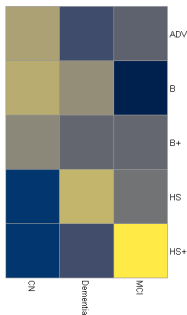
Table 7: Column constraints (inverse column margins)

	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

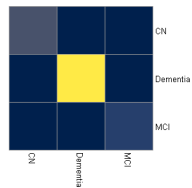
What CA needs



R: Row constraints
(inverse row probabilities)



Z: Deviations



C: Column constraints
(inverse column probabilities)

► GSVD(**R**, **X**, **C**)

- ▶ GSVD(**R**, **X**, **C**)
- ▶ Uses but generalizes the SVD

- ▶ GSVD(**R**, **X**, **C**)
- ▶ Uses but generalizes the SVD
 - ▶ Uses row & column weights (constraints)

- ▶ GSVD(**R**, **X**, **C**)
- ▶ Uses but generalizes the SVD
 - ▶ Uses row & column weights (constraints)
- ▶ Gives back

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

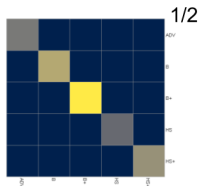
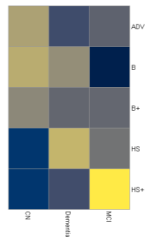
- ▶ GSVD(**R**, **X**, **C**)
- ▶ Uses but generalizes the SVD
 - ▶ Uses row & column weights (constraints)
- ▶ Gives back
 - ▶ Component (factor) scores
 - ▶ Eigenvalues, singular values, & singular vectors

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

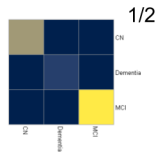
What we really decompose



=

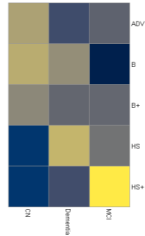


1/2

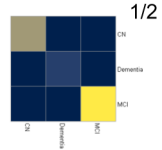
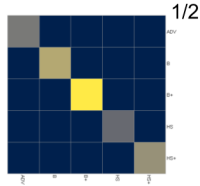


1/2

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



- Two squares
- Row margins and column margins



$$\frac{Z}{R^{\frac{1}{2}}C^{\frac{1}{2}}}$$

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

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

CA's secrets

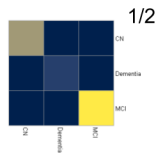
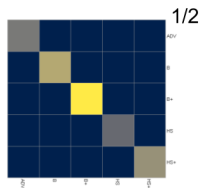
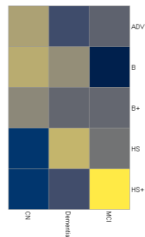
```
EDU <- amerge_subset$PTEDUCAT
DX <- amerge_subset$DX
edu_dx_table <- table(EDU, DX)
```

```
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 χ^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



EDU	DX
B	Dementia
B	MCI
B+	Dementia
HS	Dementia
B+	CN

B	B+	ADV	HS+	HS
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

MCI	CN	Dementia
0	0	1
1	0	0
0	0	1
0	0	1
0	1	0

	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

B	B+	ADV	HS+	HS
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

MCI	CN	Dementia
0	0	1
1	0	0
0	0	1
0	0	1
0	1	0

B	B+	ADV	HS+	HS
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

MCI	CN	Dementia
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+	ADV	HS+	HS	MCI	CN	Dementia
<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+	ADV	HS+	HS	MCI	CN	Dementia
<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

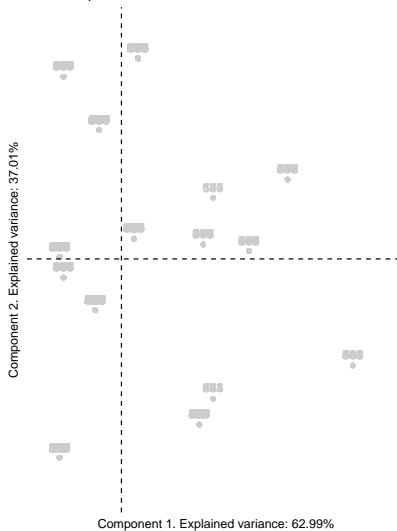
- ▶ Two perspectives:
 - ▶ *Weighted* PCA for nominal data
 - ▶ Generalized CA for N-way contingency tables
- ▶ So much more than nominal

We're diving in

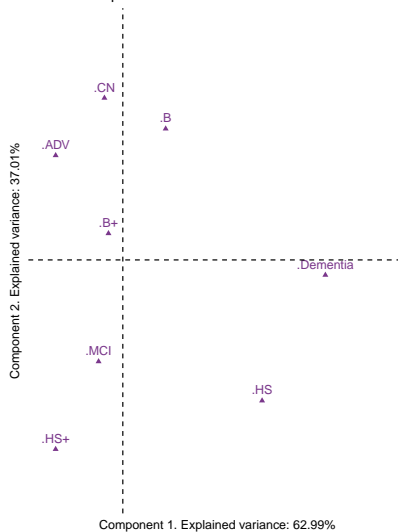
B	B+	ADV	HS+	HS	MCI	CN	Dementia
1	0	0	0	0	0	0	1
1	0	0	0	0	1	0	0
0	1	0	0	0	0	0	1
0	0	0	0	1	0	0	1
0	1	0	0	0	0	1	0

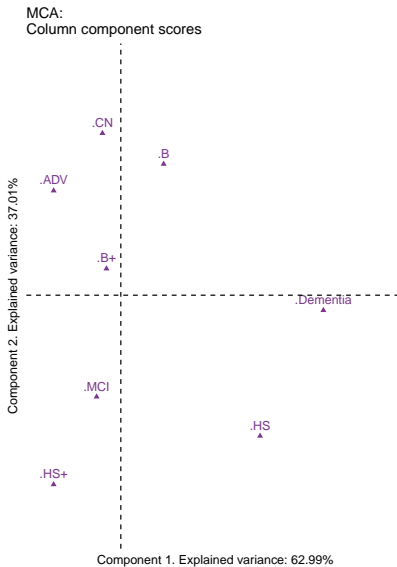
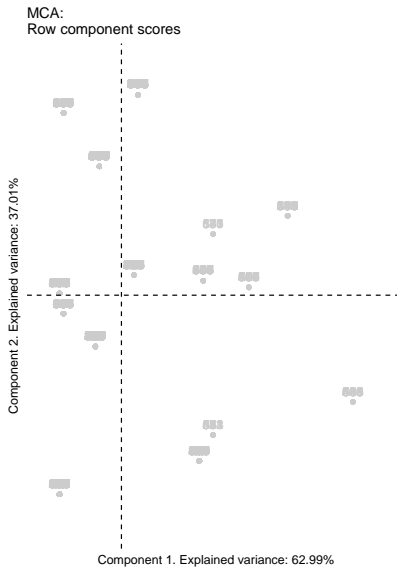
This is the kind of table we're analyzing. It has $N = 665$.

MCA:
Row component scores



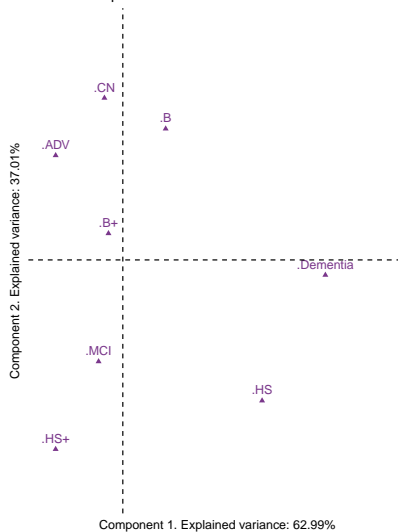
MCA:
Column component scores



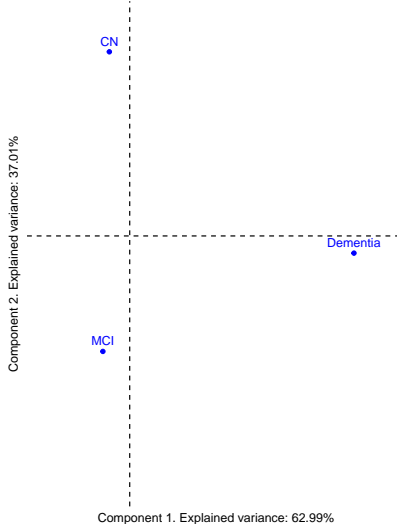


Does any of this look familiar?

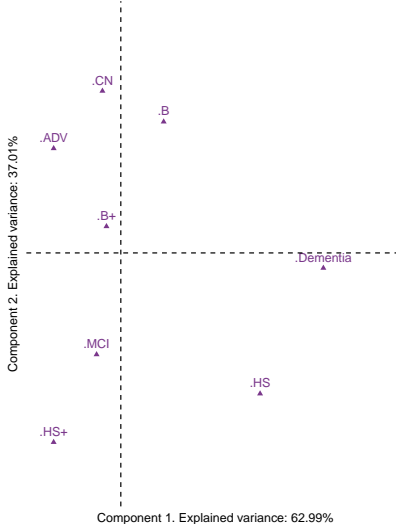
MCA:
Column component scores



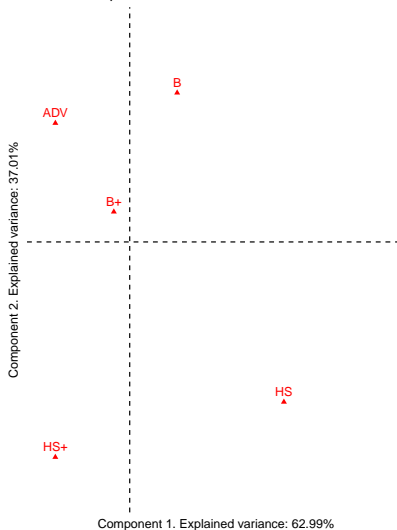
CA:
Row component scores



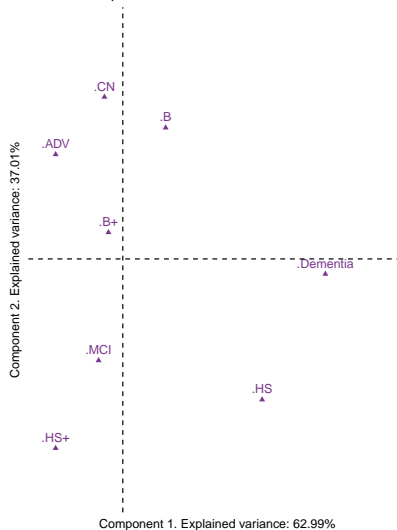
MCA:
Column component scores



CA:
Column component scores



MCA:
Column component scores



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

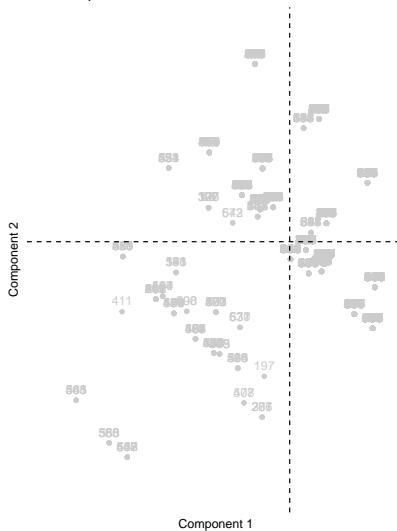
Scaling up

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

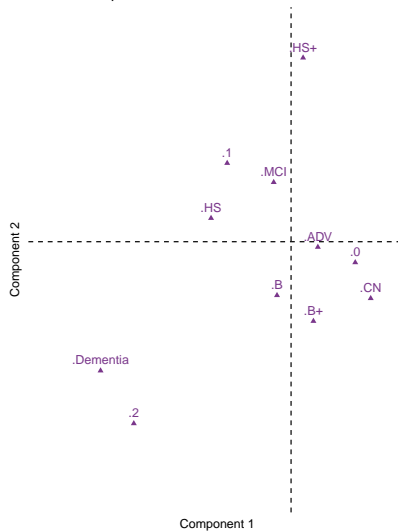
EDU	DX	APOE
B	Dementia	2
B	MCI	0
B+	Dementia	2
HS	Dementia	2
B+	CN	0

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

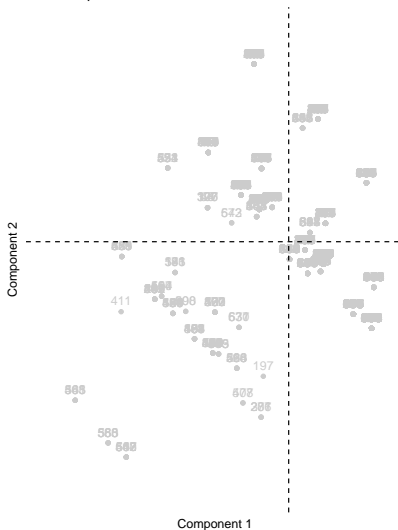
MCA:
Row component scores



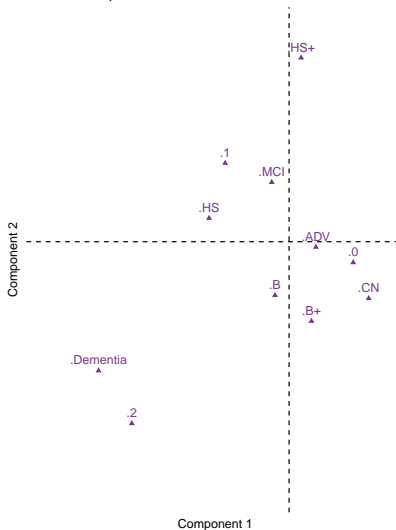
MCA:
Column component scores



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

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	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”

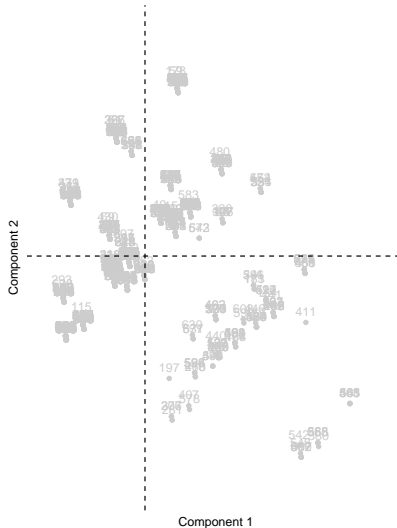
HACH
0
3
1
2
1

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

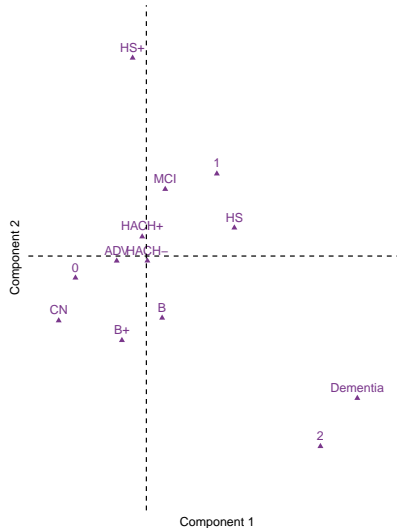
EDU	DX	APOE	HACH
B	Dementia	2	1
B	MCI	0	1
B+	Dementia	2	0
HS	Dementia	2	0
B+	CN	0	0

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

MCA:
Row component scores



MCA:
Column component scores



Our second fuzzy friend

CONTINUOUS

measured data, can have ∞
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:
 - ▶ $\frac{(1-x)}{2}$ & $+\frac{(1+x)}{2}$

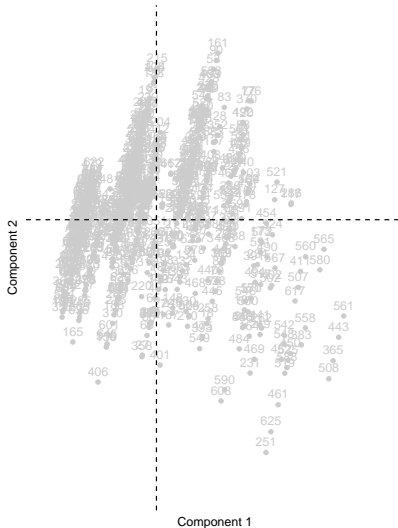
AGE	AGE (Z)
76.3	0.637
76.5	0.666
64.4	-1.095
62.9	-1.314
63.9	-1.168

AGE-	AGE+
0.181	0.819
0.167	0.833
1.048	-0.048
1.157	-0.157
1.084	-0.084

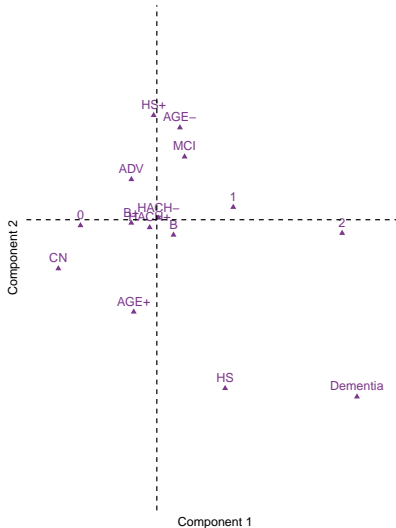
EDU	DX	APOE	HACH	AGE
B	Dementia	2	1	76.3
B	MCI	0	1	76.5
B+	Dementia	2	0	64.4
HS	Dementia	2	0	62.9
B+	CN	0	0	63.9

B	B+	ADV	HS+	HS	MCI	CN	Dementia	0	1	2	HACH-	HACH+	AGE-	AGE+
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

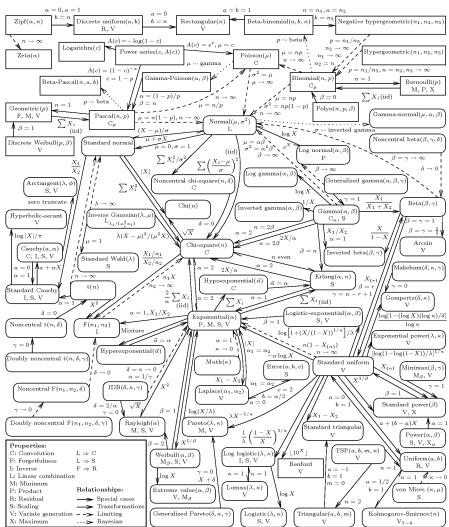
MCA:
Row component scores



MCA:
Column component scores

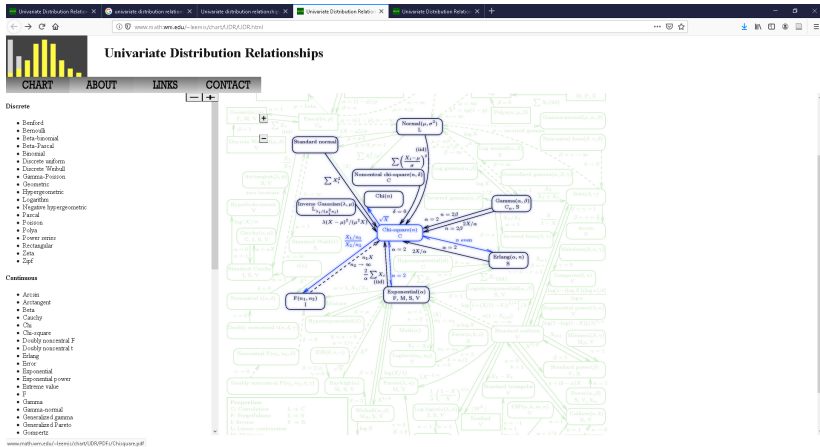


Chi-squared



See [here](#)

Chi-squared



See [here](#)

Some many bonuses!

(Some) References

See the reference sections of these

- ▶ 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.

See the reference sections of these

- ▶ 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

- ▶ Abdi, H., Guillemot, V., Eslami, A., & Beaton, D. (2017). Canonical correlation analysis. Encyclopedia of Social Network Analysis and Mining, 1-16.

And these

- ▶ Abdi, H., Guillemot, V., Eslami, A., & Beaton, D. (2017). Canonical correlation analysis. Encyclopedia of Social Network Analysis and Mining, 1-16.
- ▶ 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>