

# PRINCIPAL COMPONENT ANALYSIS

A STORY WITH WINES  
GUEST STAR: THE EIGEN FAIRY

HERVÉ ABDI & JU-CHI YU

# TODAY'S

- The core of it all:
- Principal Component Analysis (PCA)

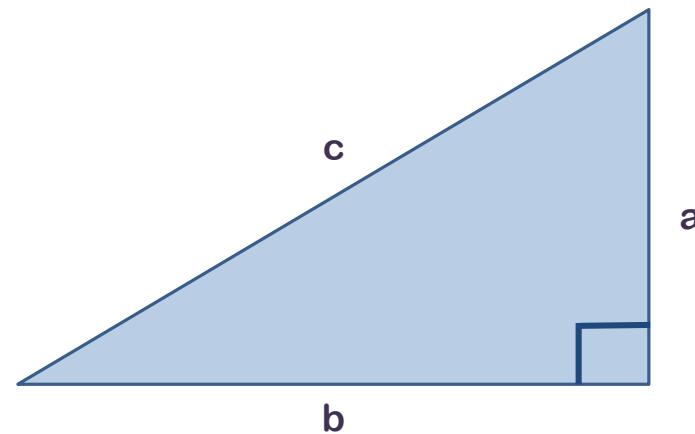
# WHAT DO YOU REALLY NEED TO KNOW

→ **The THEOREM:**

→ **Yes, this one!**

→ **The Pythagorean Theorem**

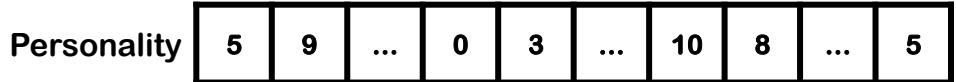
$$a^2 + b^2 = c^2$$



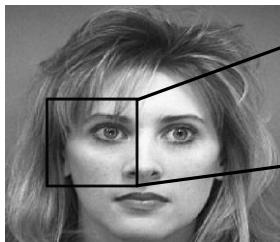
## WHAT DO WE MEASURE?

# WHAT ARE MULTIVARIATE DATA?

→ Surveys



→ Images



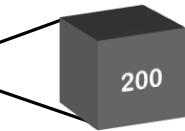
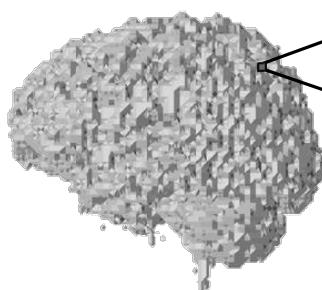
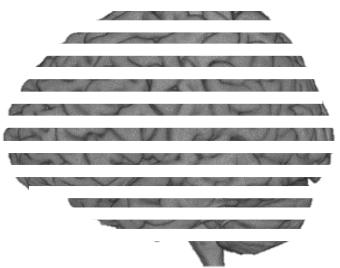
90	73
42	26

A 2x2 grid of numbers, likely representing pixel values or features extracted from the face image.

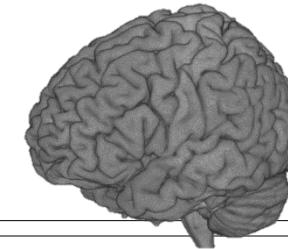
Face

A horizontal bar with a gradient from dark to light gray, representing a feature vector for the face.

→ Brains



Brain

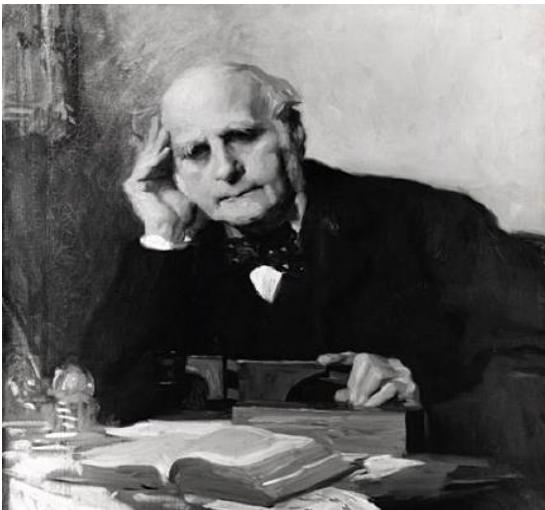
A horizontal bar with a gradient from dark to light gray, representing a feature vector for the brain.

# PRINCIPAL COMPONENT ANALYSIS

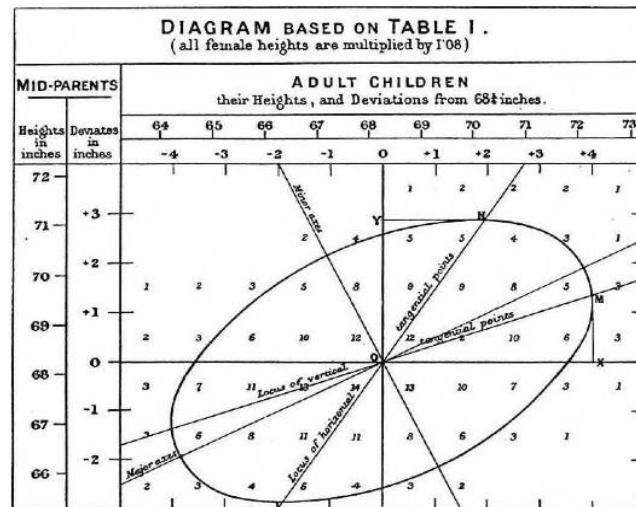
- You already did it today!
- Huh?
- Example 1.
- Example 1. Google (original page rank)
- Example 2
- Example 2. Music (.mp3)
- Example 3
- Example 3. Netflix
- Example 4, 5, ....

# PRINCIPAL COMPONENT ANALYSIS

- Oldest multivariate technique
    - Idea: Cauchy (1829); Galton (1859); Pearson (1902)
    - AKA: Eigen-analysis; singular value decomposition, Hotelling transform, Karhunen-Loève decompositions, ...



# Sir Francis Galton (1822-1911)



# (A Kind of) PCA Comparing Mid-Parental to Adult Children Height

# PRINCIPAL COMPONENT ANALYSIS

- **Oldest multivariate technique**
  - Idea: Cauchy (1829); Galton (1859); Pearson (1902)
  - AKA: Eigen-analysis; singular value decomposition, Hotelling transform, Karhunen-Loève decompositions, ...
- **But really doable only with computers**
- **Goals of PCA:**
  - Extracts important information from a data table
  - Represents the information as nicer variables
  - Replaces “nasty numbers” by “nice maps”
    - For both observations and variables

# WORKING WITH MULTIVARIATE DATA

## → Common characteristics

- Often more variables than observations ( $N \ll P$ )
  - Traditional analyses cannot handle this
- Inference: usual assumptions may not hold

## → Goal of multivariate analysis techniques

- Extract most important information
- Present it well

# **PCA EXAMPLE WITH 20 WINES AND 2 VARIABLES**

### Meet the Expert



Our Expert tasted 20 (+1) wines

Rated them on 2 scales (0 to 20):

*Sugar (Sweet)*

*Astringency*

10 Wines were from the US (Zinfandel)

10 Wines were from France (Cabernet + Merlot)

Later on, rated the wines on *Bitter* and *Acidic*

Hypothesis:

Commented:  
**Sugar & Astringent inversely perceived**  
Some wines were *Fruity*, some were *Woody*

# 20 WINES



(Long) Name

Short Name	1. Bordeaux	2. Stone	3. Listrac	4. Canyon Creek	5. Côtes de Bourg	6. Poor Hill	7. Hollow	8. St Estèphe	9. Wooden Hill	10. Blaye	11. Côtes de Blaye	11. Sun Set	12. Black Bird	13. Médoc	14. St Julien	14. St Julien	15. Pauillac	16. Rush	17. Oak Ville	18. Gravé	19. St Emillion	20. Temecula	
F: Sugar	12	6	2	6	2	9	4	5	9	4	11	7	4	11	5	4	9	—	4	10	4	15	10
W: Astringency	14	7	11	9	9	4	6	11	5	8	8	2	2	4	12	9	9	U	4	10	13	15	6
Origin	F	U	F	U	F	U	U	F	U	F	F	U	U	F	F	F	F	U	U	F	F	U	



## 20 WINES: ALL DATA

(Long) Name	1. Bordeaux	2. Black Stone	3. Listrac	4. Canyon Creek	5. Côtes de Bourg	6. Post Hill	7. Hollow	8. St. Estèphe	9. Wooden Hill	11. Côtes de Blaye	11. Sun Set	12. Black Bird	13. Médoc	14. St. Julien	15. Pauillac	16. Gold Rush	17. Oak Vale	18. Grave	19. Emil	20. Tornillo
Short Name	1. Bordeaux	2. Stone	3. Listrac	4. Canyon Creek	5. Bourg	6. Hill	7. Hollow	8. Estephe	9. Wooden	10. Blaye	11. Sun	12. Bird	13. Medoc	14. Julien	15. Pailiac	16. Bush	17. Oak	18. Grave	19. St. Brillion	20. Tornillo
F: Sugar	3	6	2	6	2	9	6	5	9	4	7	11	8	11	5	8	4	5	10	10
W: Astringency	14	7	11	9	9	4	8	11	5	5	8	8	12	12	9	8	1	13	15	6
Origin	F	U	F	U	F	U	U	F	U	F	U	U	F	F	F	U	U	F	U	
Fruity			F		F			F	F				F	F	F					
Woody													W	W			W	W	W	
Acid	2	5	1	1	9	1	2	2	1	1	2	1	2	1	2	9	14	2	1	2
Bitter	8	3	16	3	11	1	1	1	9	1	8	3	2	9	12	10	2	1	10	7



## 20 WINES

(Long) Name

	1. Bordeaux	2. Stone	3. Listrac	4. Canyon Creek	5. Côtes de Bourg	6. Poor Hill	7. Hollow	8. St Estèphe	9. Wooden Hill	11. Côtes de Blaye	11. Sun Set	12. Black Bird	13. Médoc	14. St Julien	15. Pauillac	16. Rush	17. Oak Ville	18. Gravé	19. St Emillion	20. Temasello
Short Name	1. Bord	2. Stone	3. Listrac	4. Canyon	5. Bourg	6. Hill	7. Hollow	8. Estephe	9. Wooden	10. Blaye	11. Sun	12. Bird	13. Medoc	14. Julien	15. Pauillac	16. Rush	17. Oak	18. Gravé	19. Emil	20. Temas
F: Sugar	12	6	2	6	2	9	4	5	9	4	7	11	5	5	9	9	10	4	15	10
W: Astringency	14	7	11	9	9	4	8	11	5	8	2	11	4	12	9	1	4	13	6	6
Origin	F	U	F	U	F	U	U	F	U	F	U	U	F	F	F	U	U	F	F	U

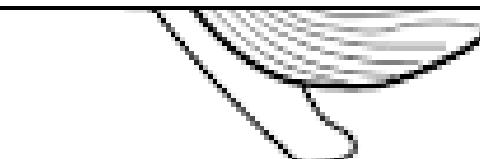
Hypothesis:

Sugar & Astringent inversely perceived

# WE ARE LOST WITH NUMBERS!

(Long) Name	1. Bordeaux	2. Black Stone	3. Lismac	4. Canyon Creek	5. Côtes de Bourg	6. Rock Hill	7. Harbor	8. St. Esprithe	9. Wooden Hill	11. Côtes de Blaye	11. Sun Sea	12. Black Bird	13. Médoc	14. St. Julian	15. Paillac	16. Rush	17. Oak Ville	18. Grave	19. St. Emilion	20. Tomassello
Short Name	1. Bordour	2. Stone	3. Listrac	4. Canyon	5. Bourg	6. Hill	7. Harbor	8. Esprithe	9. Wooden	10. Blaye	11. Sun	12. Bird	13. Medoc	14. Julian	15. Paillac	16. Rush	17. Oak	18. Grave	19. Emil	20. Tommas
F: Sugar	13	6	2	6	2	9	6	5	9	4	7	11	5	12	5	6	9	10	4	10
W: Astringency	14	7	11	9	9	8	11	5	8	8	2	4	12	9	6	1	4	13	15	6

Numbers

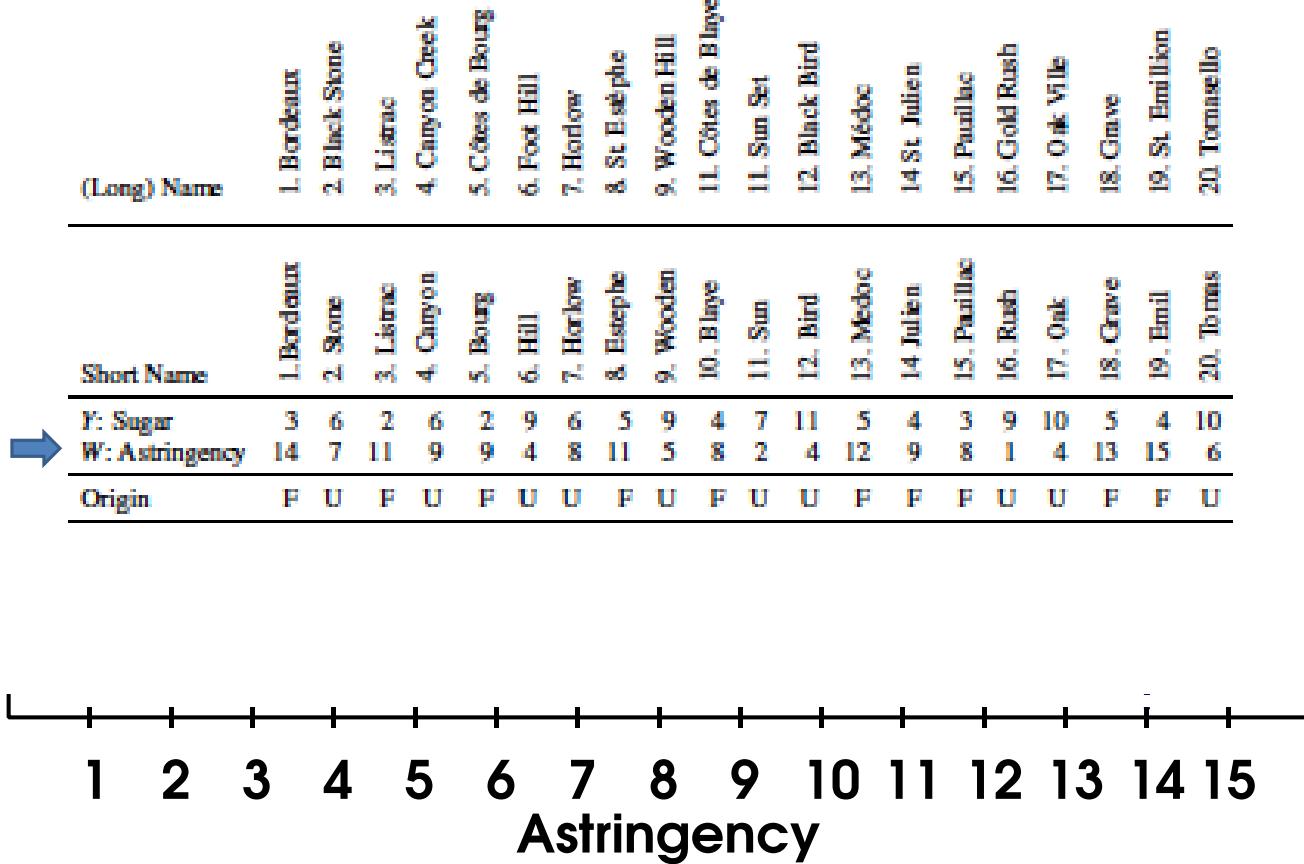


# So: MAKE A PICTURE!

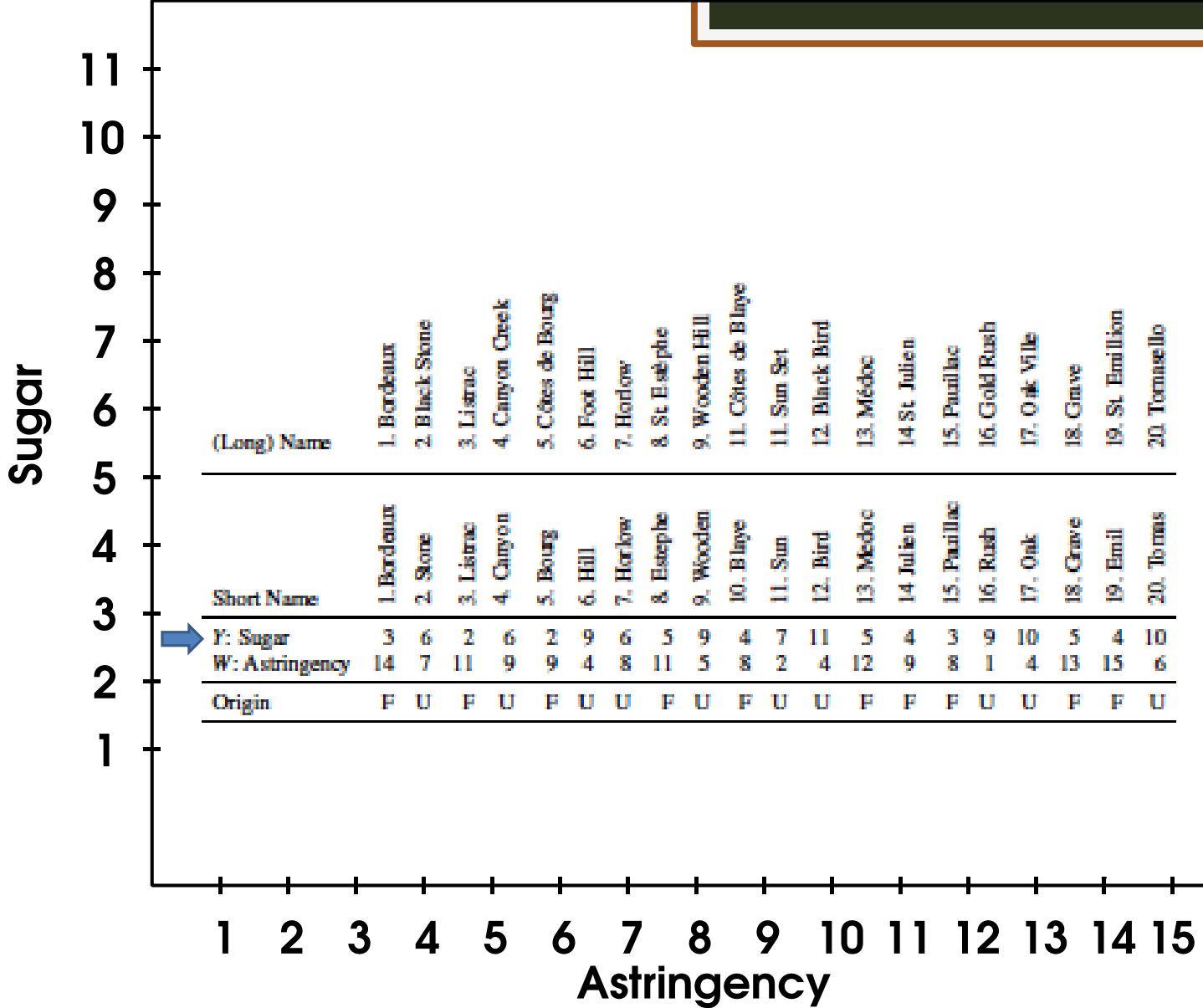
(Long) Name	1. Bordeaux	2. Black Stone	3. Listrac	4. Canyon Creek	5. Cotes de Bourg	6. Foot Hill	7. Hollow	8. St Estephe	9. Wooden Hill	11. Cotes de Blaye	11. Sun Set	12. Black Bird	13. Medoc	14. St Julien	15. Pauillac	16. Rush	17. Oak Ville	18. Grave	19. St Emilion	20. Tornasello
Short Name	3	6	2	4. Canyon	5. Bourg	6. Hill	7. Harkow	8. Estephe	9. Wooden	10. Blaye	11. Sun	12. Bird	13. Medoc	14. Julien	15. Pauillac	16. Rush	17. Oak Ville	18. Grave	19. St Emilion	20. Tornas
F: Sugar	14	7	11	9	9	4	8	11	5	8	2	11	12	13	9	8	1	4	15	16
W: Astringency	P	U	F	U	F	U	U	F	U	F	U	F	U	F	F	U	U	F	F	U
Origin	F	U	F	U	F	U	U	F	U	F	U	F	U	F	F	U	U	F	F	U

(Long) Name	1. Bordeaux	1. Bordeaux
Short Name	2. Stone	2. Black Stone
F: Sugar	3. Listrac	3. Listrac
W: Astringency	4. Canyon	4. Canyon Creek
Origin	5. Bourg	5. Côtes de Bourg
F	6. Hill	6. Foot Hill
U	7. Harkow	7. Hollow
F	8. Estephe	8. St Estèphe
U	9. Wooden	9. Wooden Hill
U	10. Blaye	11. Côtes de Blaye
F	11. Sun	11. Sun Set
U	12. Bird	12. Black Bird
U	13. Medoc	13. Médoc
F	14. Julien	14. St Julien
F	15. Pauillac	15. Pauillac
F	16. Rush	16. Gold Rush
U	17. Oak	17. Oak Ville
U	18. Grave	18. Grave
F	19. Emil	19. St Emillion
F	20. Tornas	20. Tornasello
U		

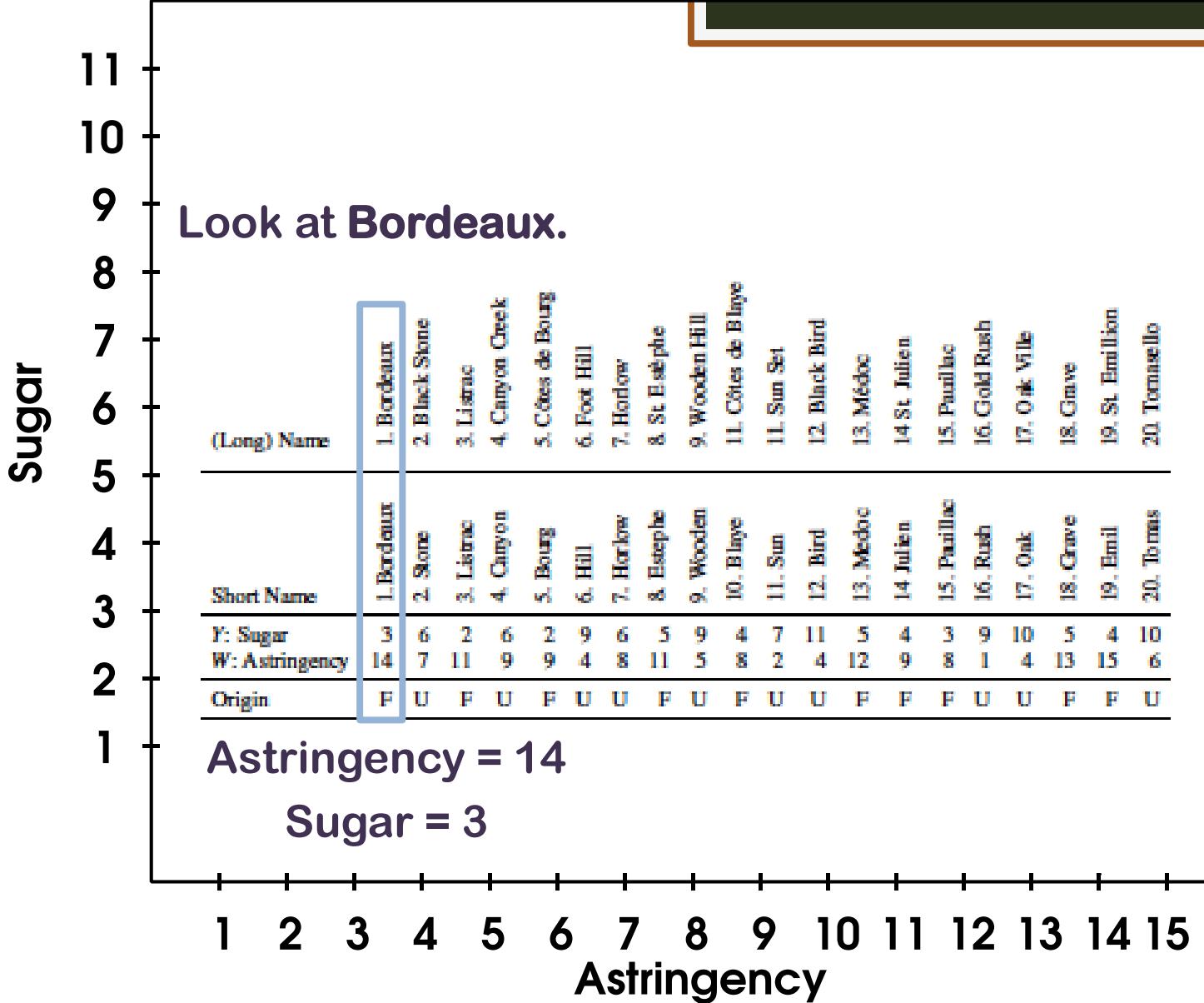
# A PICTURE IS WORTH ...



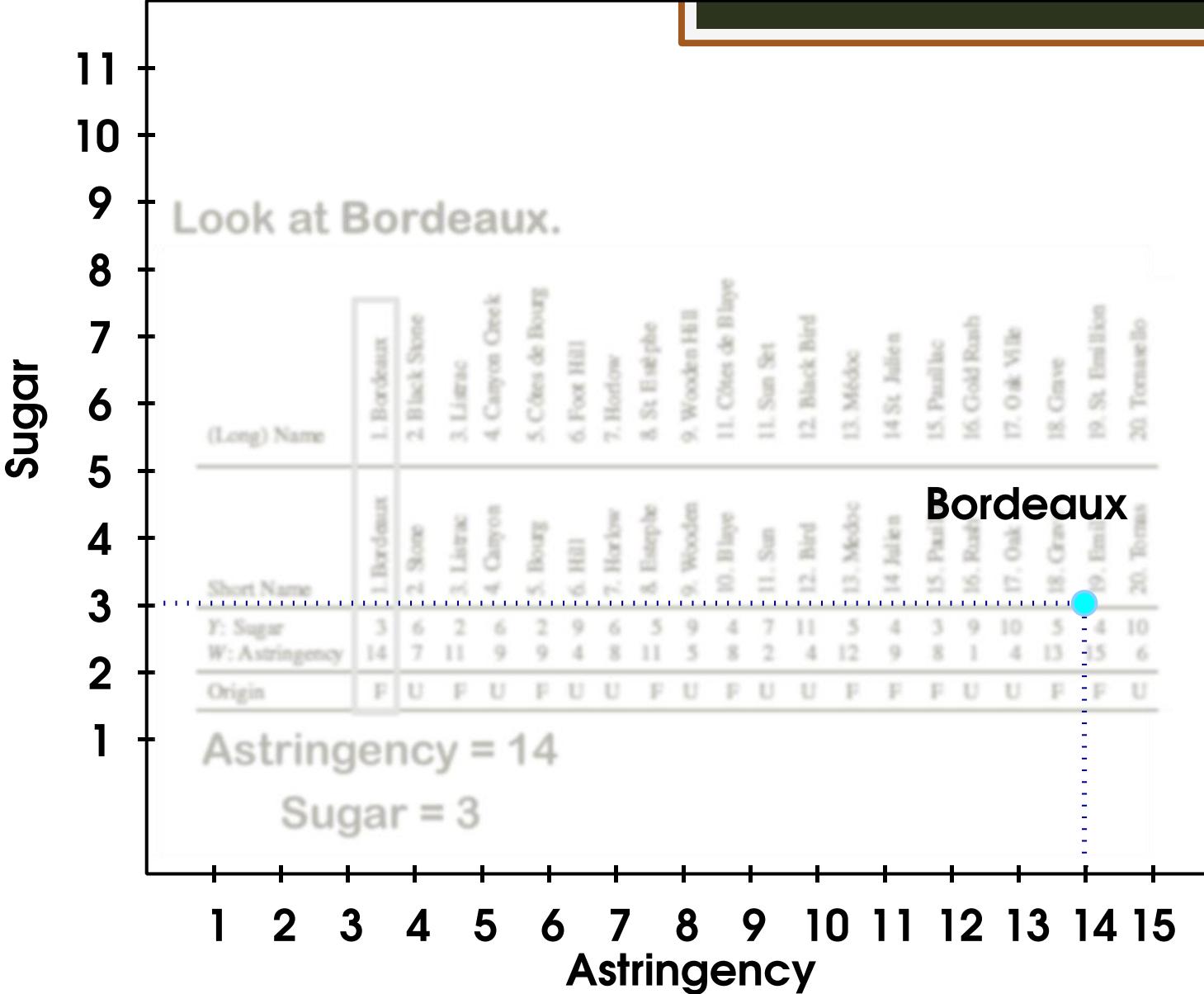
## A PICTURE IS WORTH ...



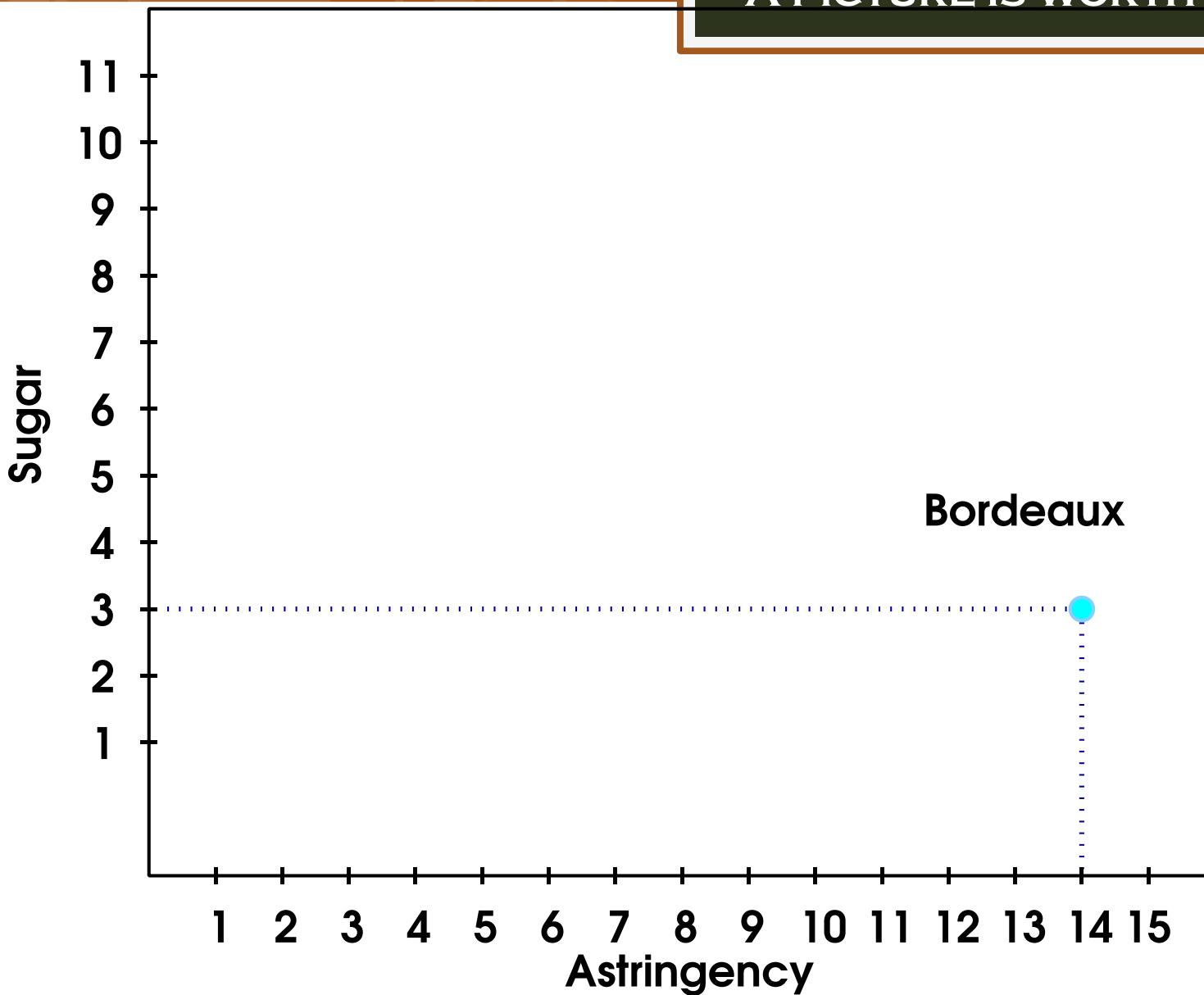
A PICTURE IS WORTH ...



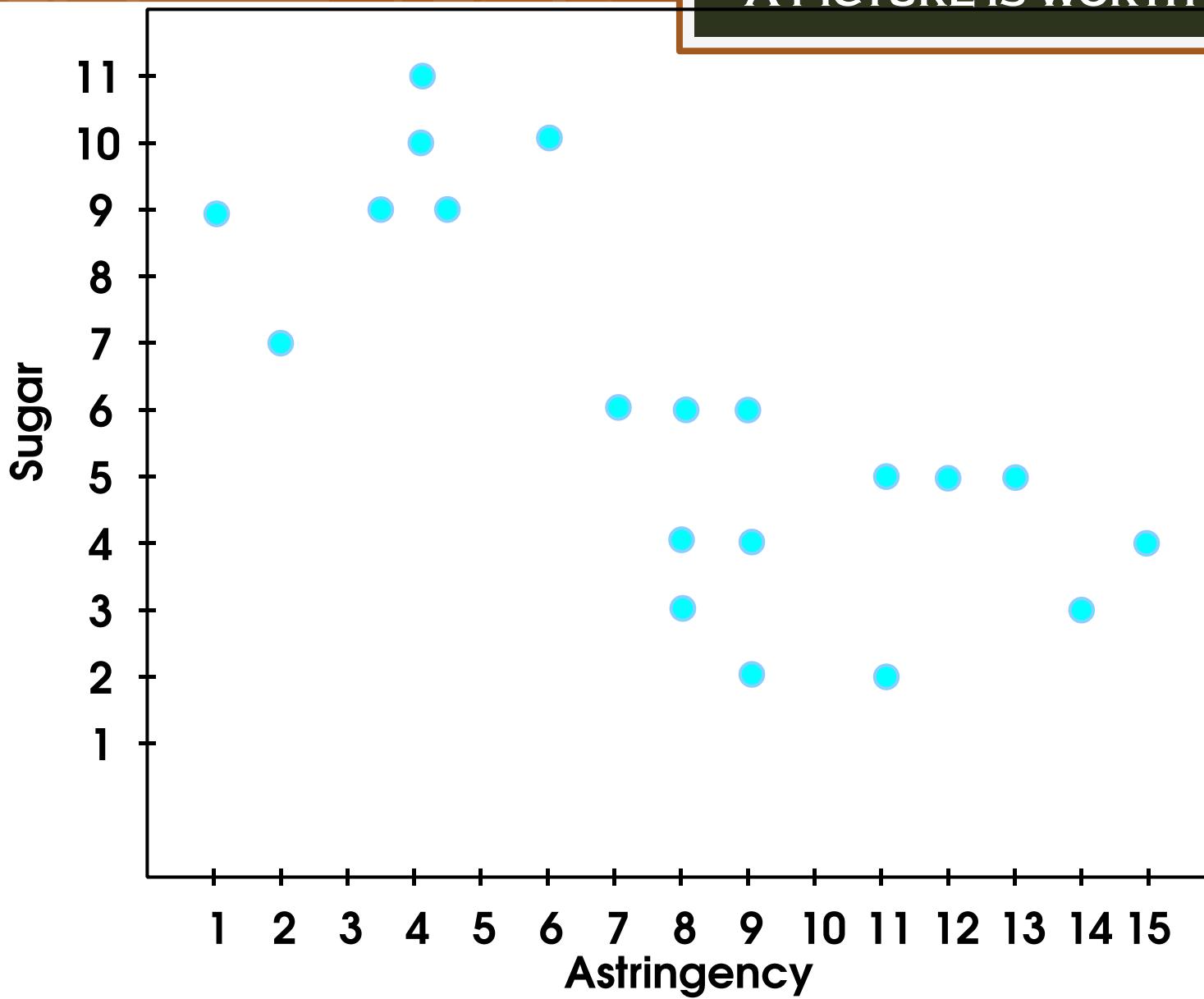
A PICTURE IS WORTH ...

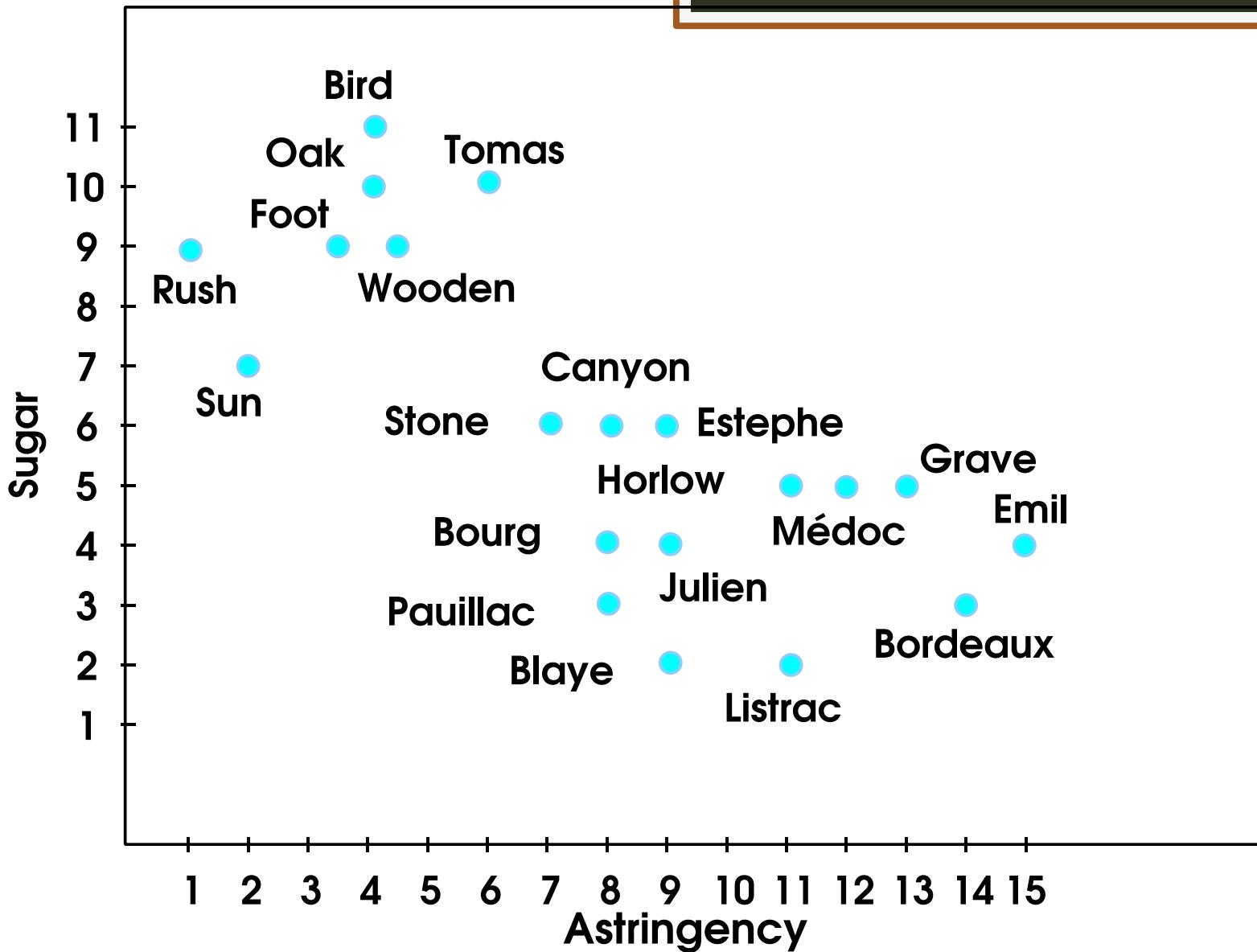


A PICTURE IS WORTH ...

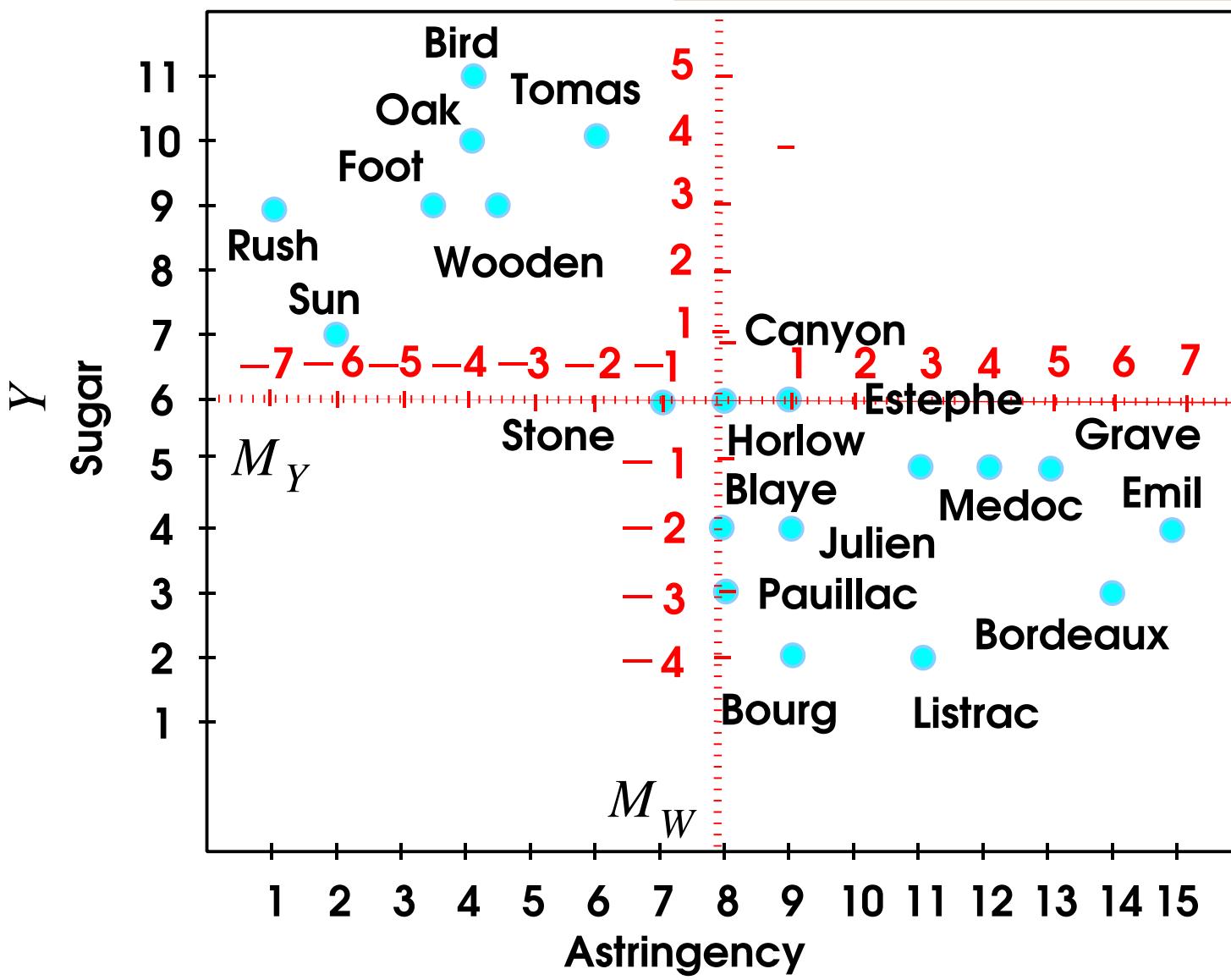


A PICTURE IS WORTH ...



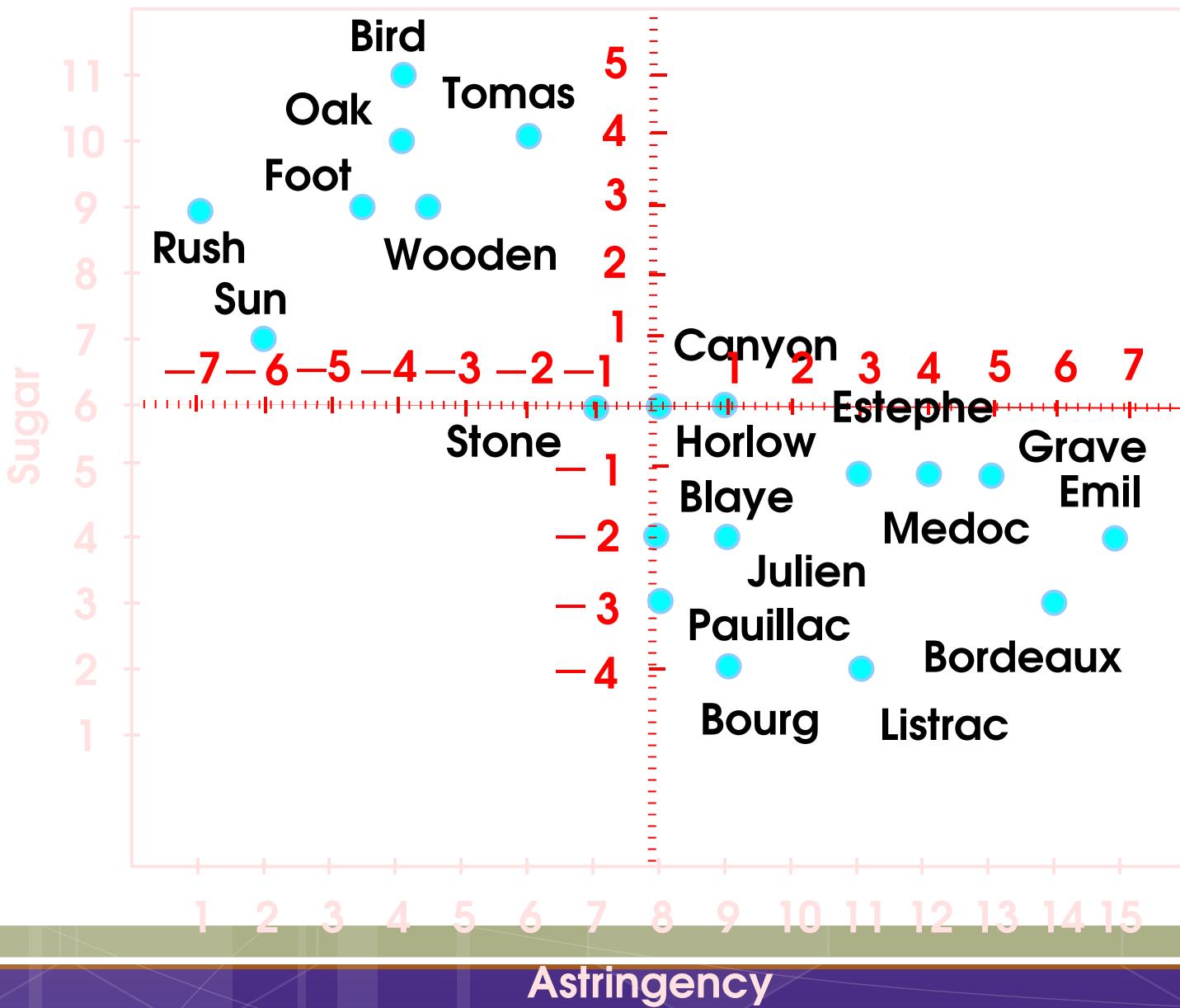
**A PICTURE ...**

A PICTURE ...

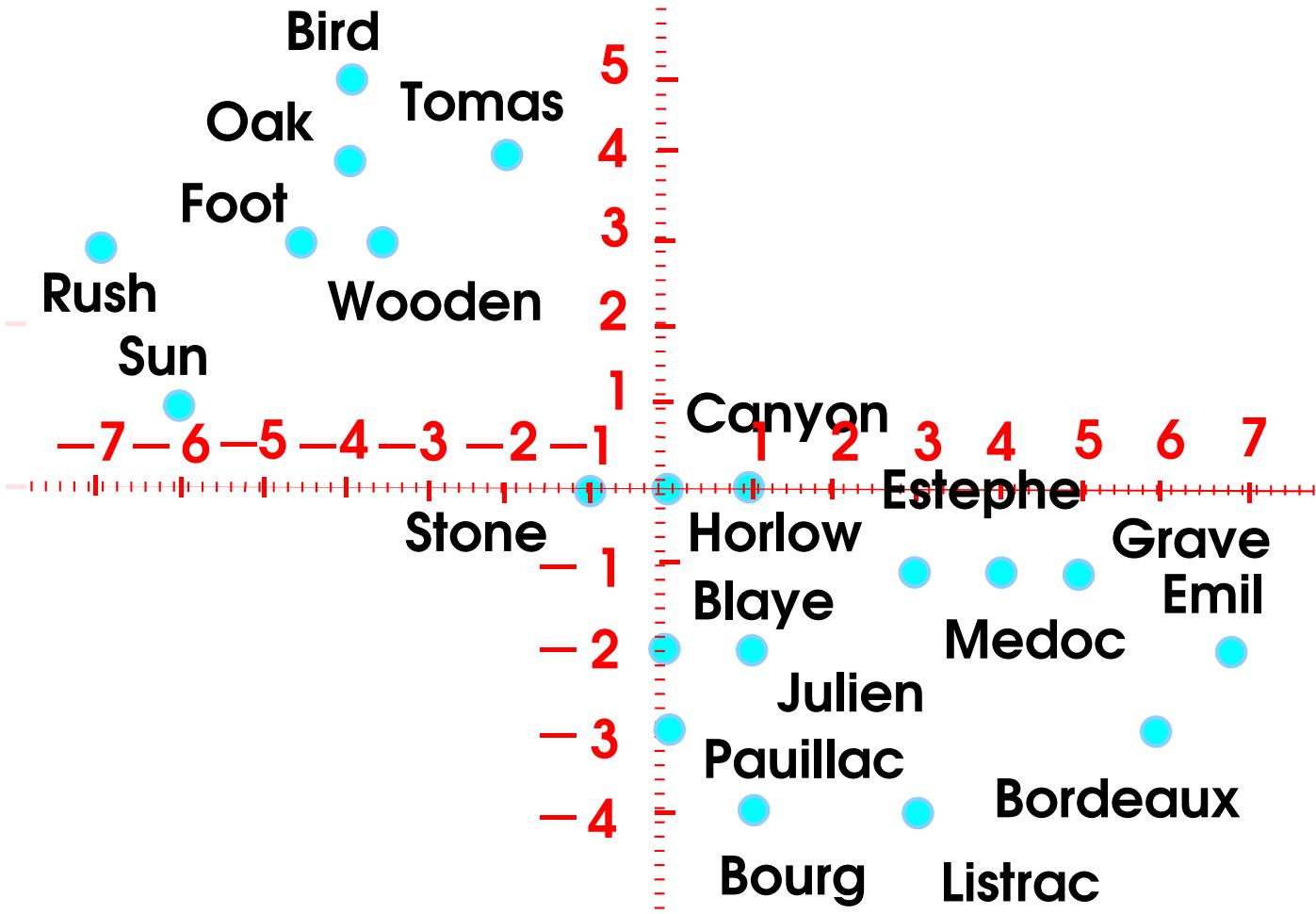


W

## A PICTURE ...

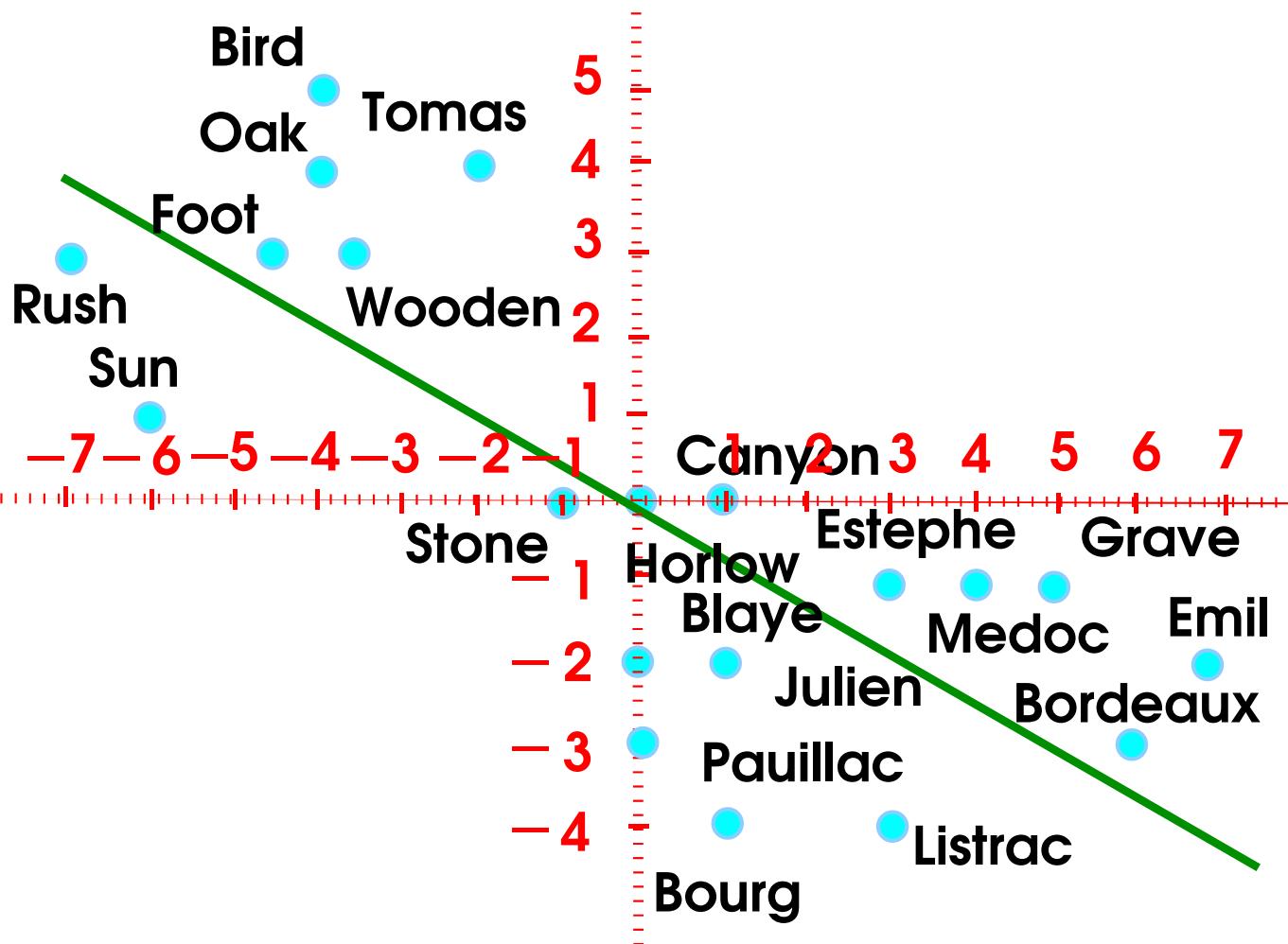


## A PICTURE ... CENTERING



## A PICTURE ... CENTERING

Could we get a line passing in the middle of this cloud?

**A PICTURE ... PC 1**

## A PICTURE ... CENTERING

How to do it?

EYEBALLING!

## How to do it?

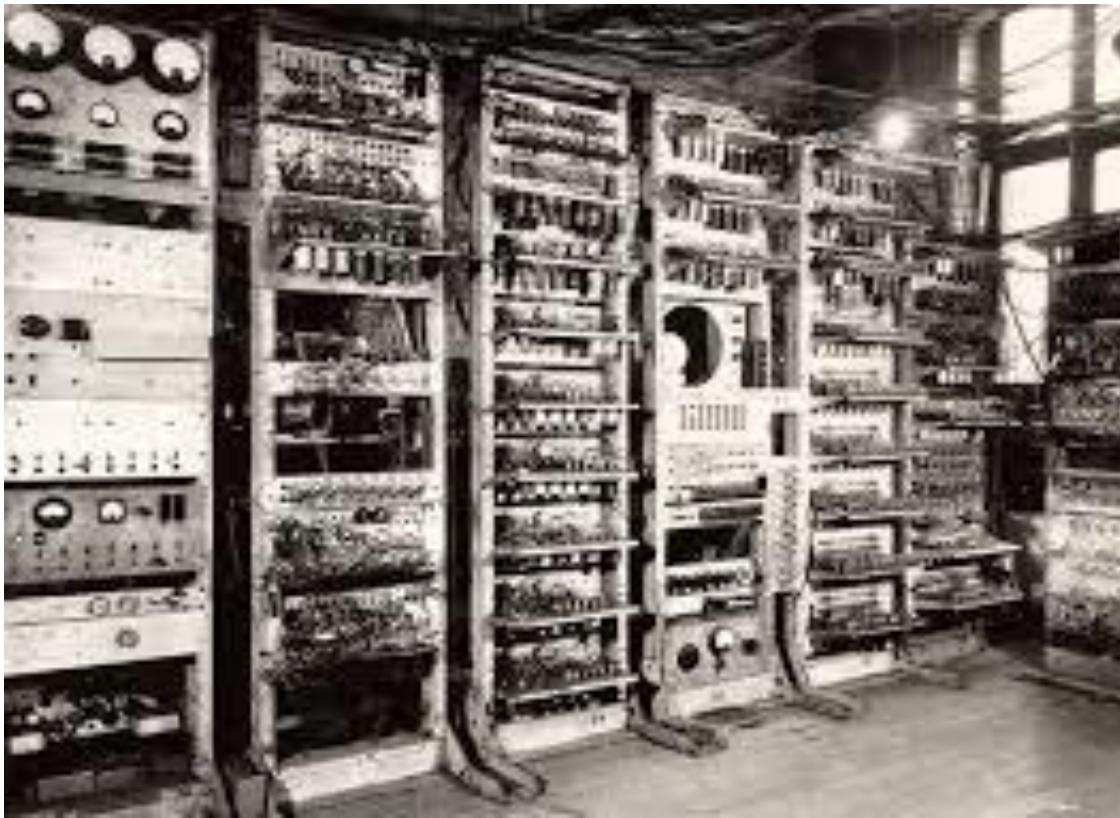


FOR THE HANDY ONES ...

## How to do it? The old way ...



## SOME COMPUTING ...



## THE EIGEN FAIRY

# OR EVEN MATH MAGIC! THE EIGEN FAIRY



### WHAT TO ASK FOR?

EIGEN FAIRY, PLEASE  
FIND THE ROTATION  
WITH **SMALLEST SUM OF  
SQUARED DISTANCES**

$$D^2 = A^2 + B^2$$

# PYTHAGOREAN MAGIC:



## SIMILARLY, SMALLEST SUM OF SQUARED DISTANCES

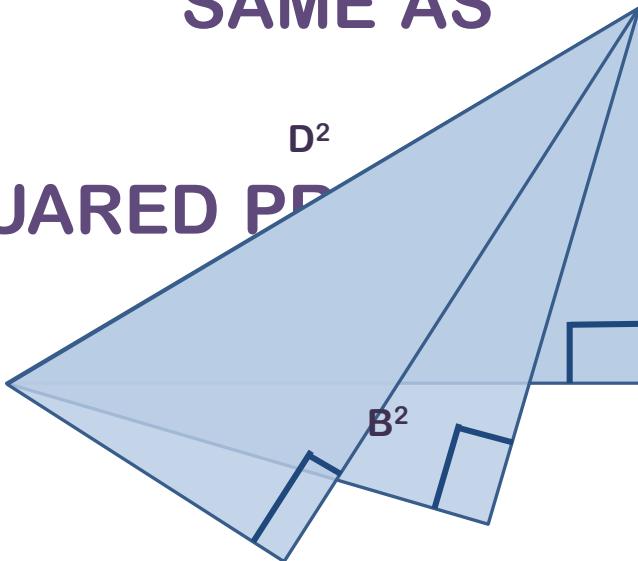
SAME AS

LARGEST SQUARED PROJECTIONS (VARIANCE)

$$D^2 = A^2 + B^2$$

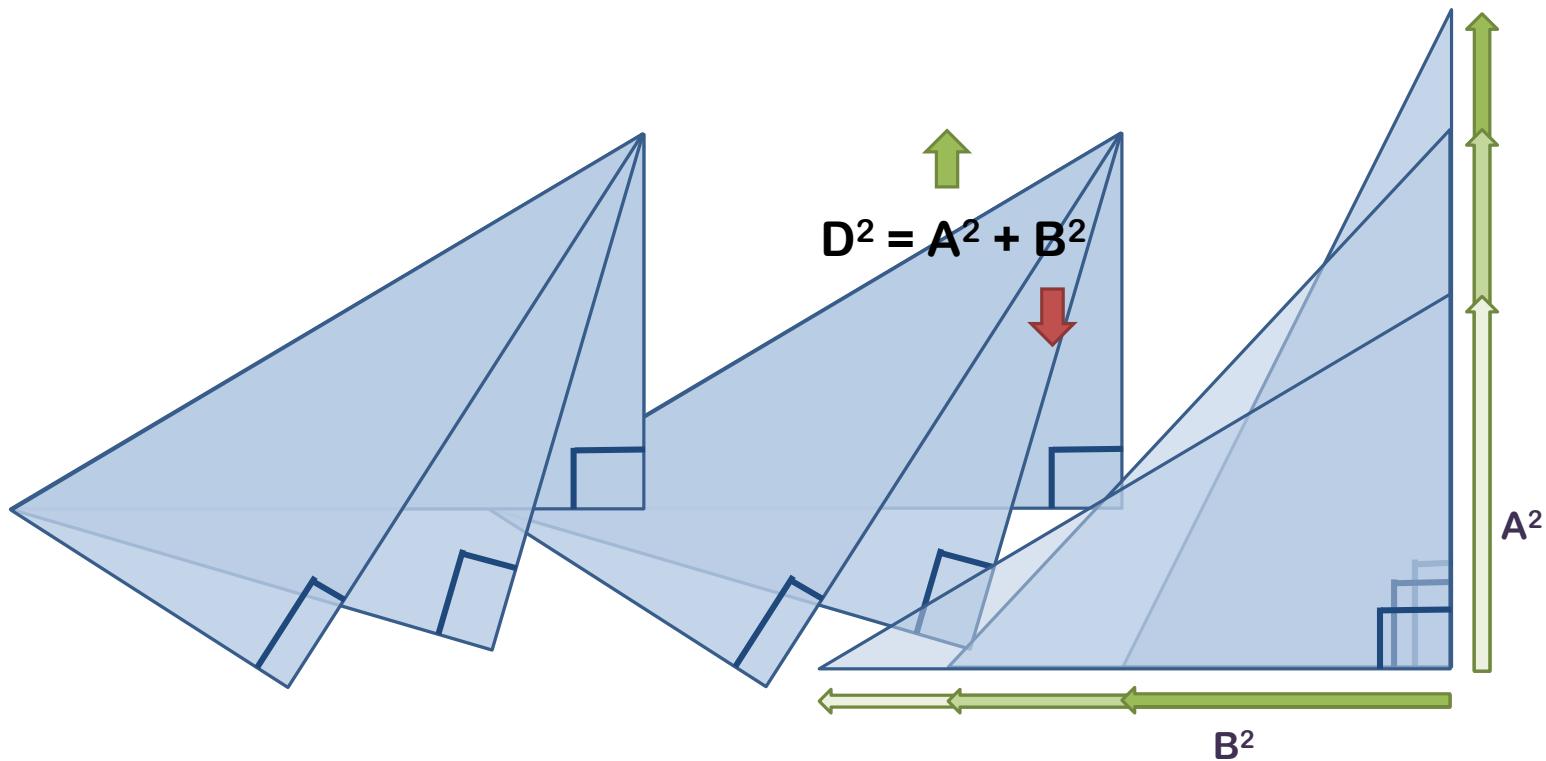
# PYTHAGOREAN MAGIC: SMALLEST SUM OF SQUARED DISTANCES

SAME AS  
LARGEST SQUARED PERP DISTANCE (VARIANCE)

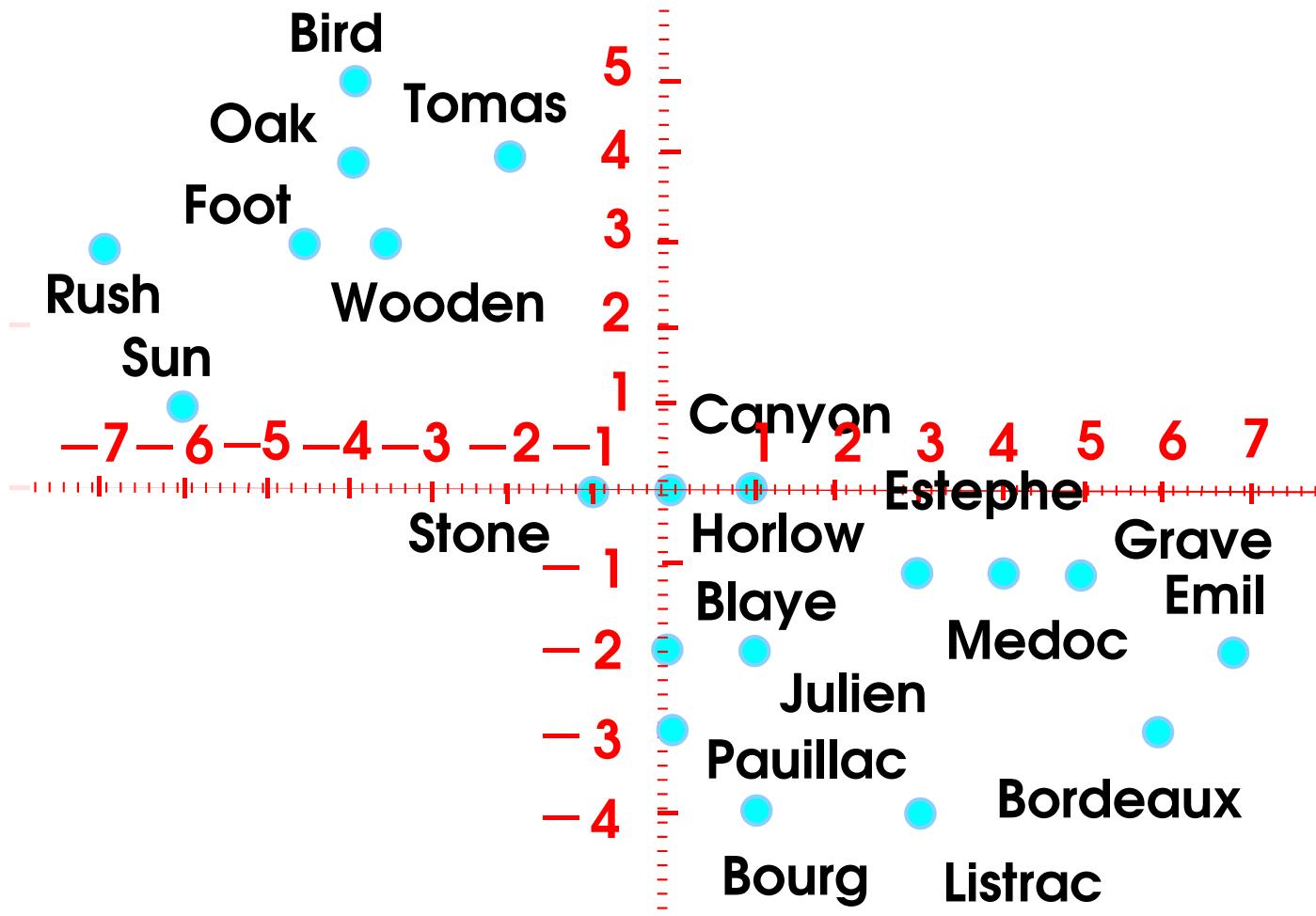


$$D^2 = A^2 + B^2$$

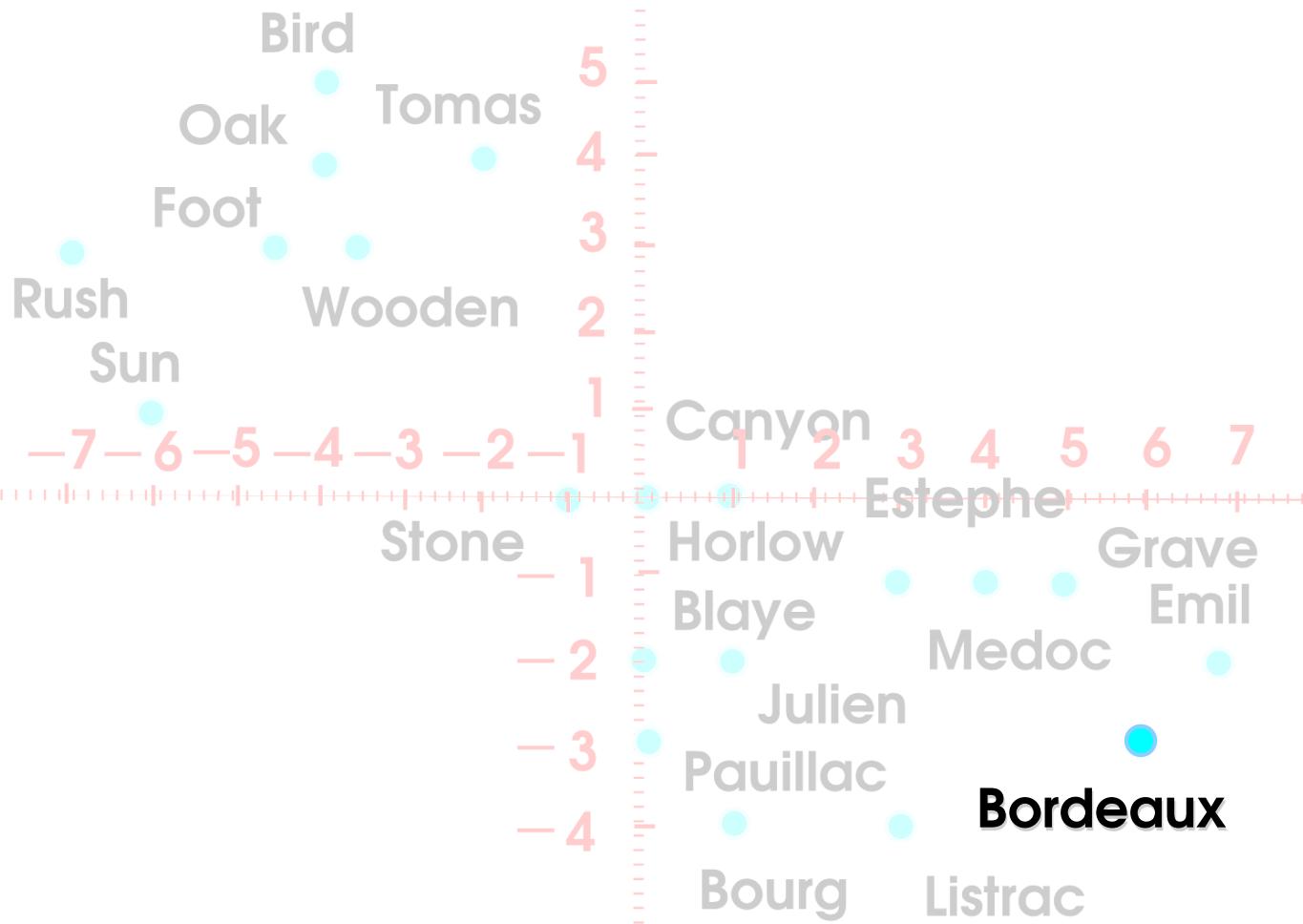
# PYTHAGOREAN MAGIC: SMALLEST SUM OF SQUARED DISTANCES SAME AS LARGEST SQUARED PROJECTIONS (VARIANCE)



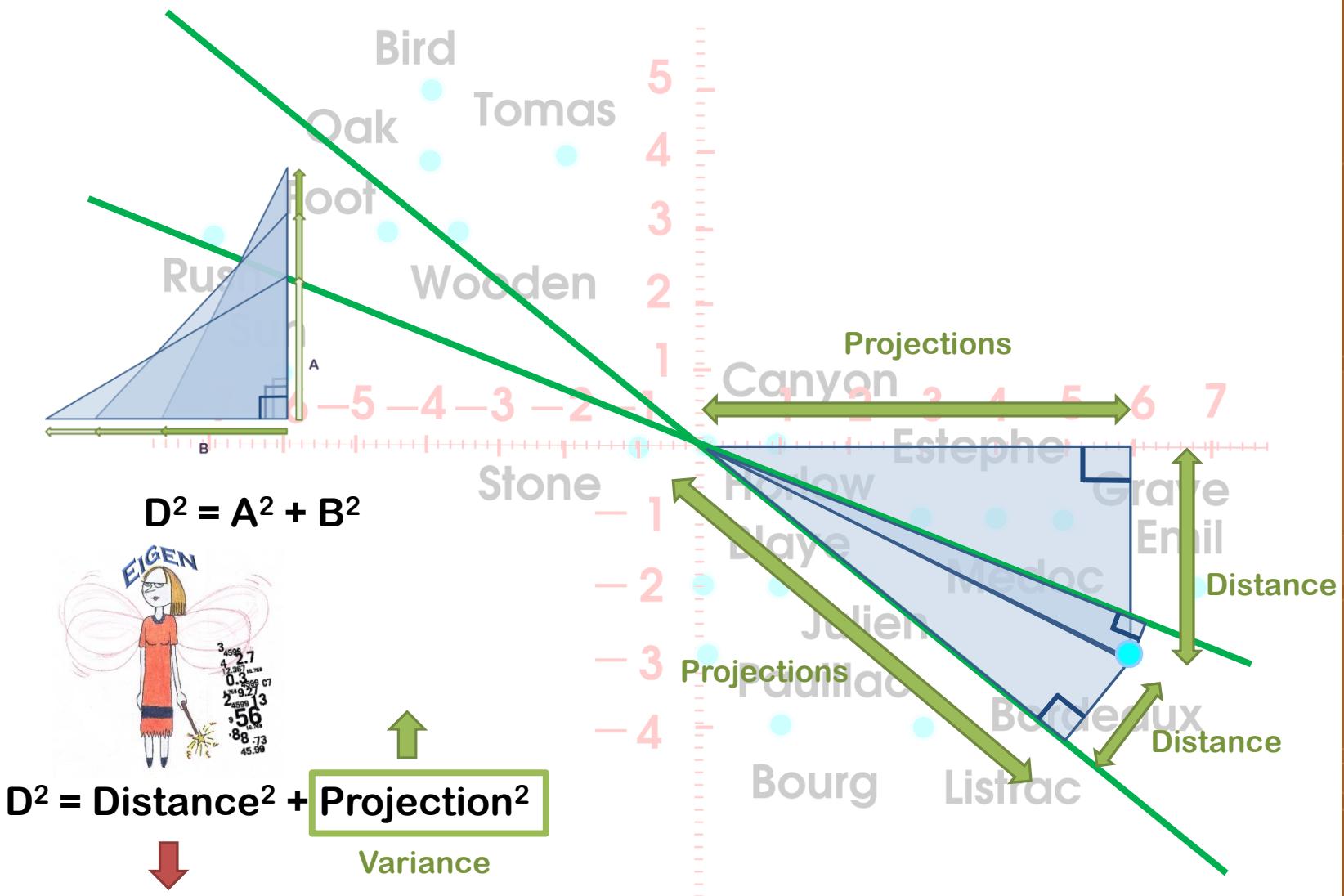
## BACK TO THE PICTURE



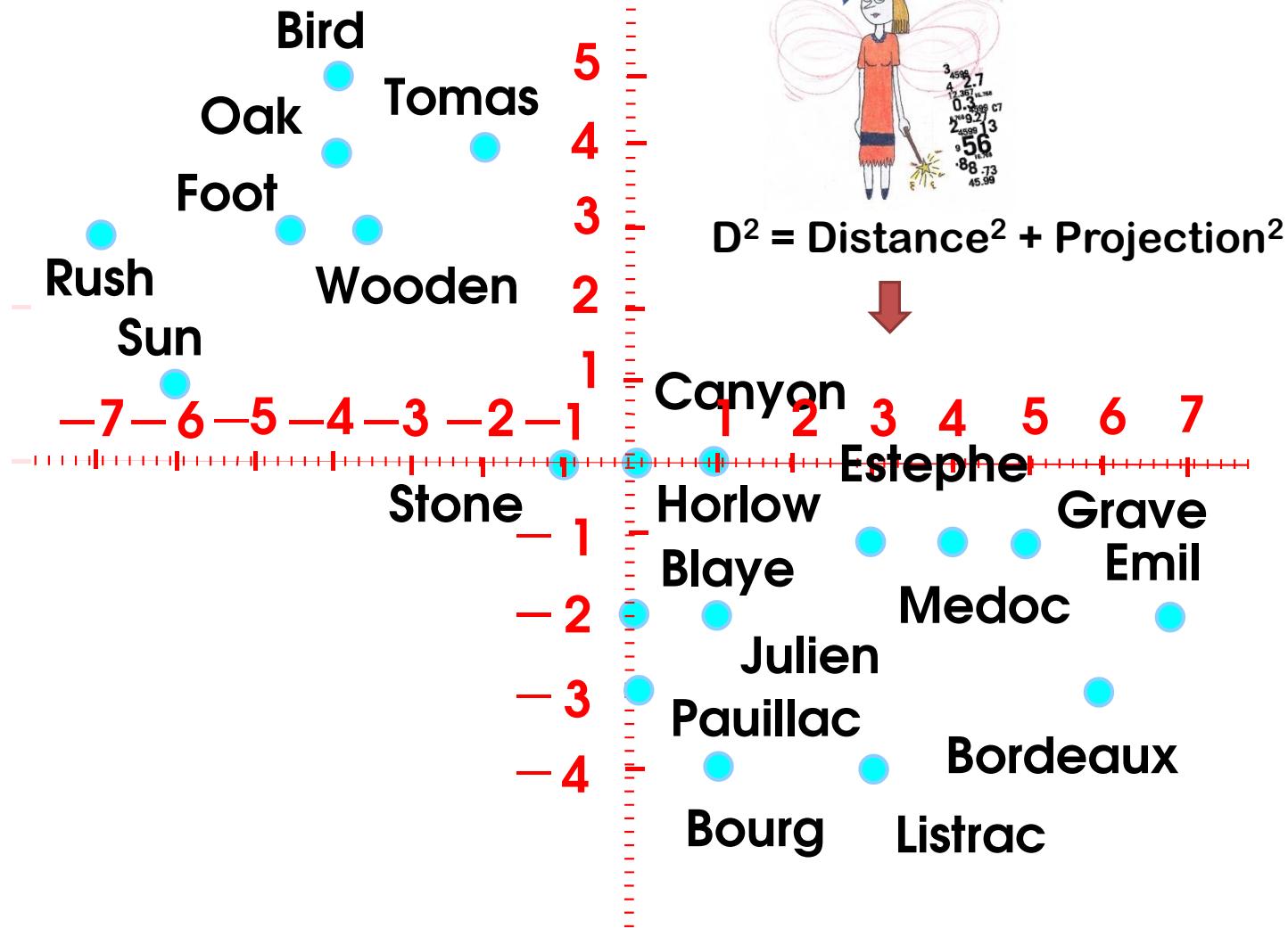
## BACK TO THE PICTURE



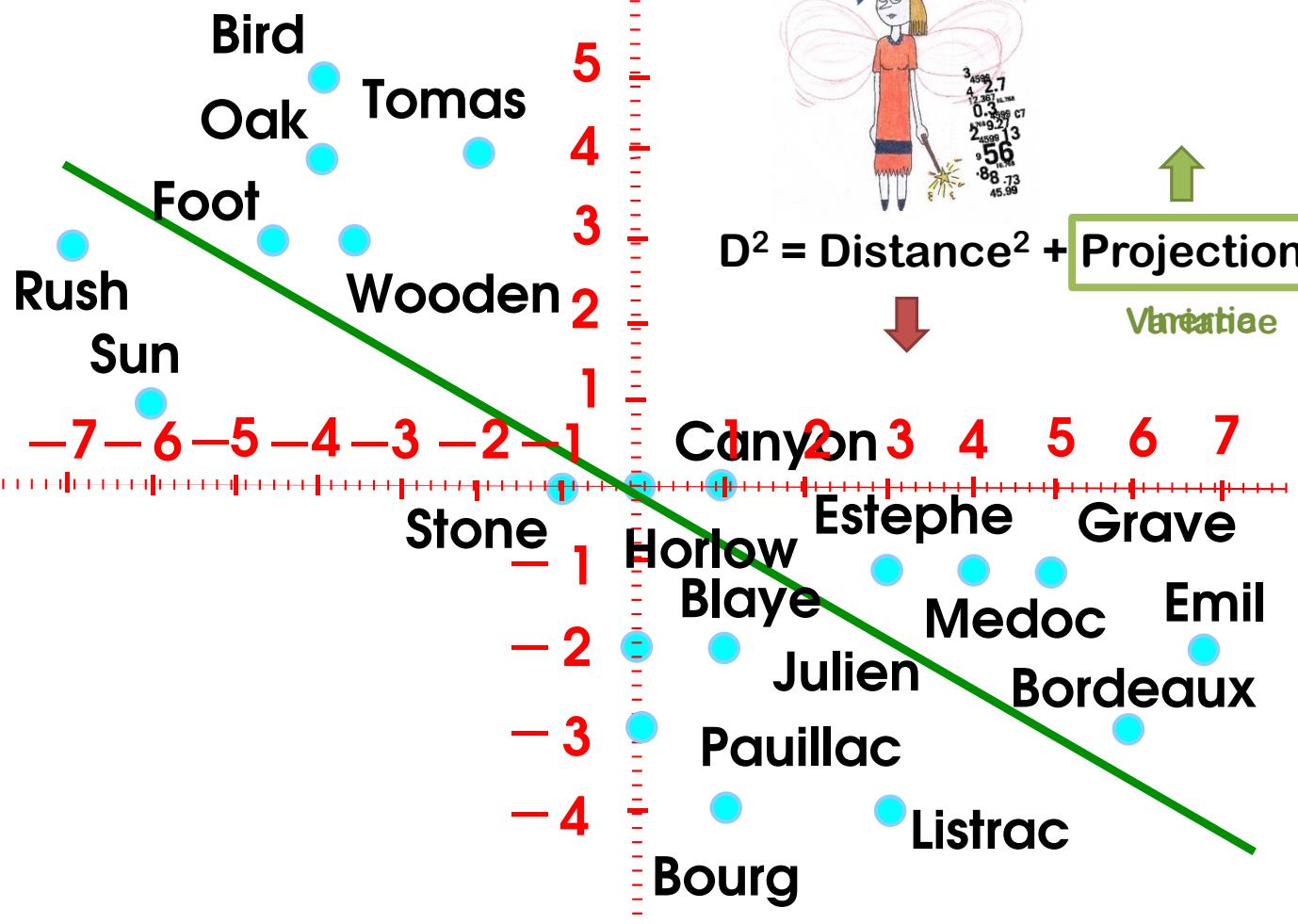
## BACK TO THE PICTURE

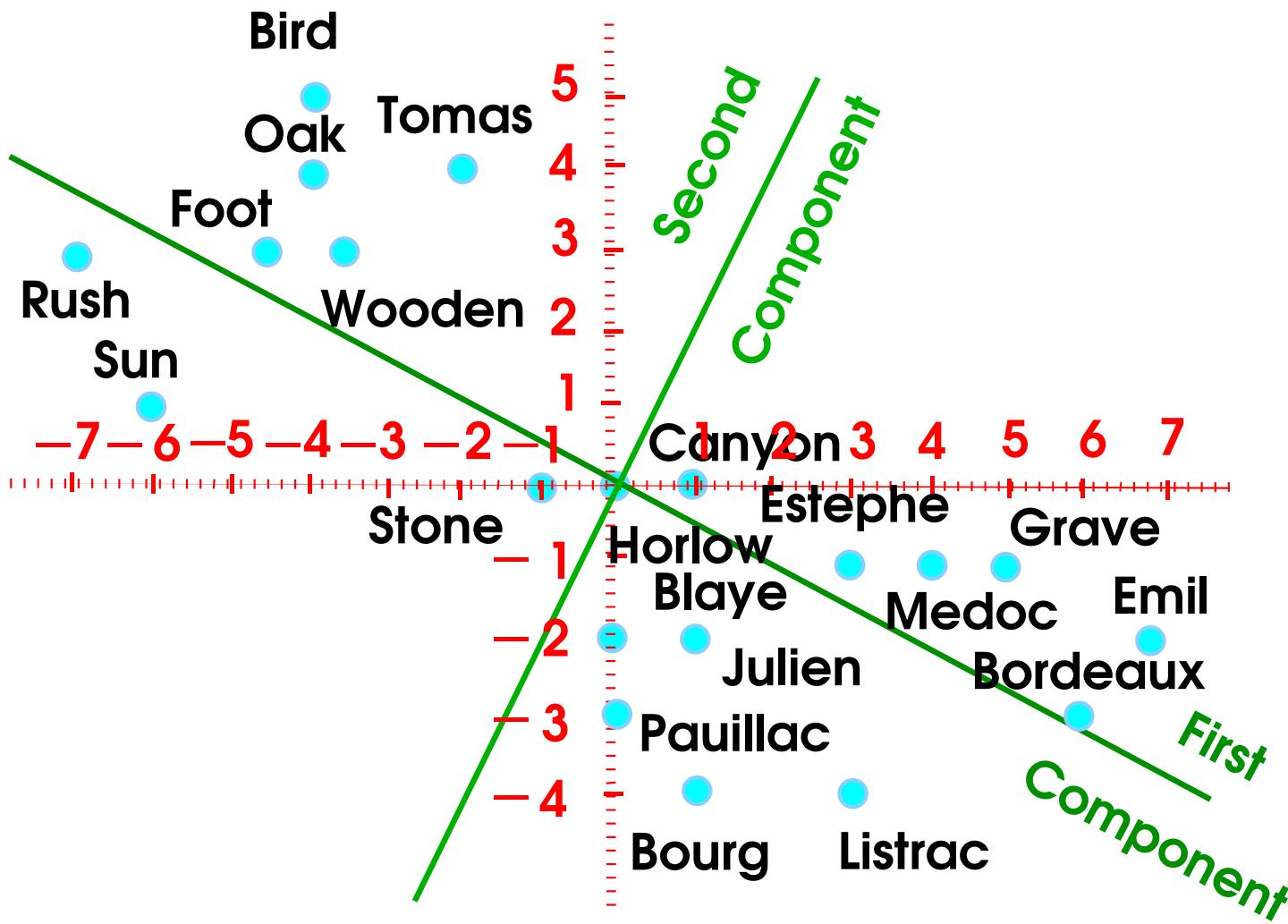


# A PICTURE ... PC 1

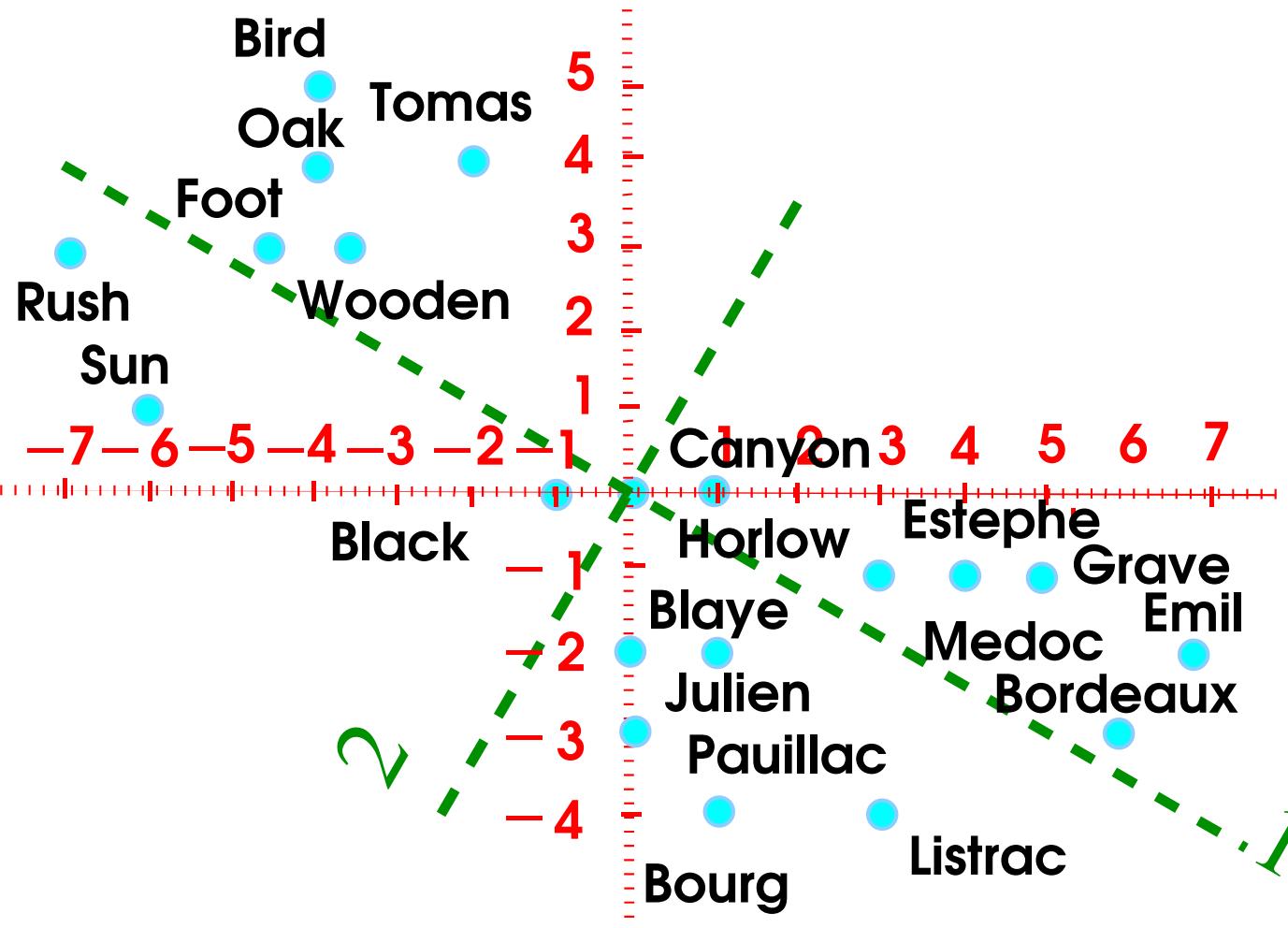


# A PICTURE ... PC 1

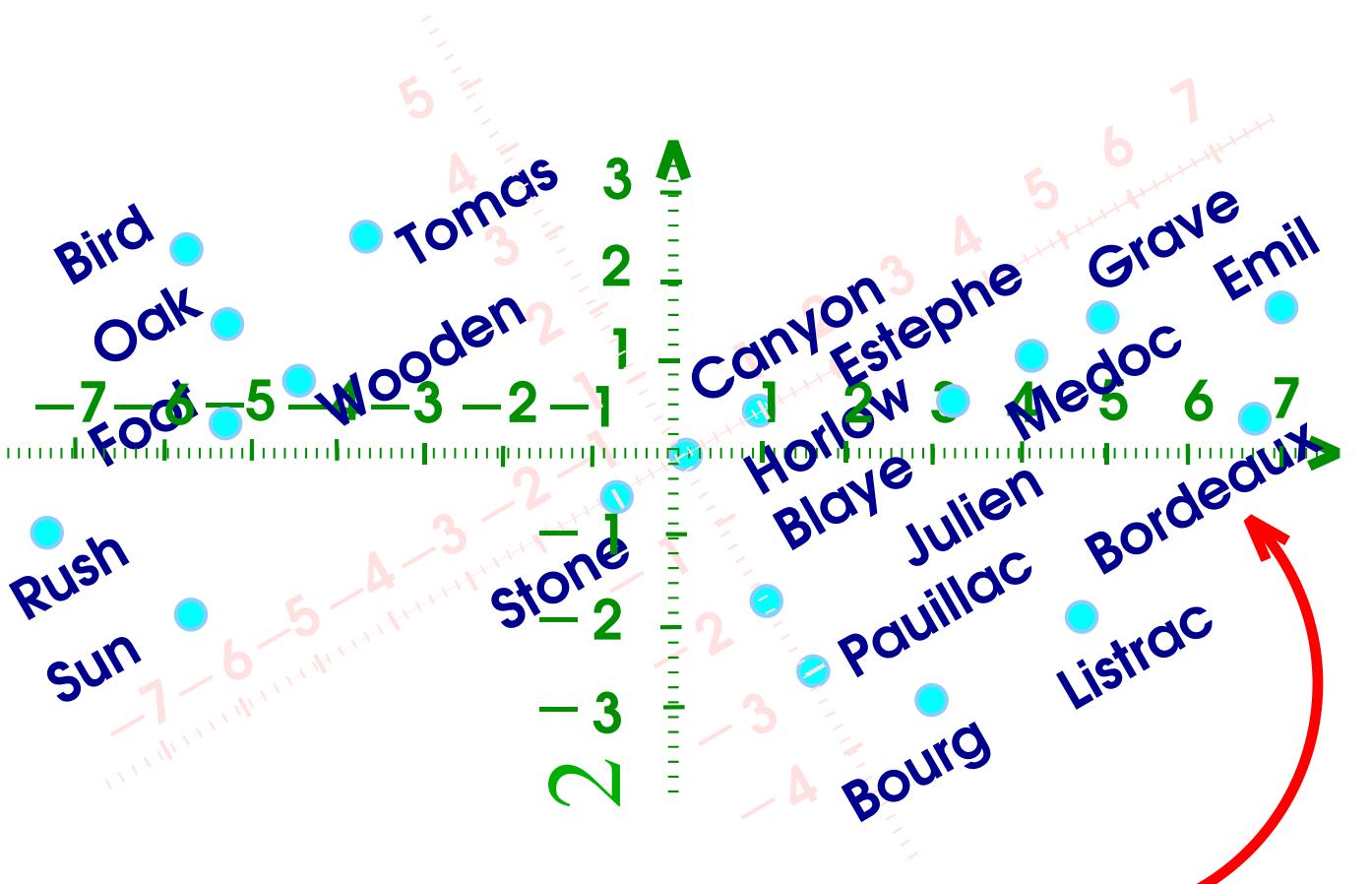


**A PICTURE ... PC 1**

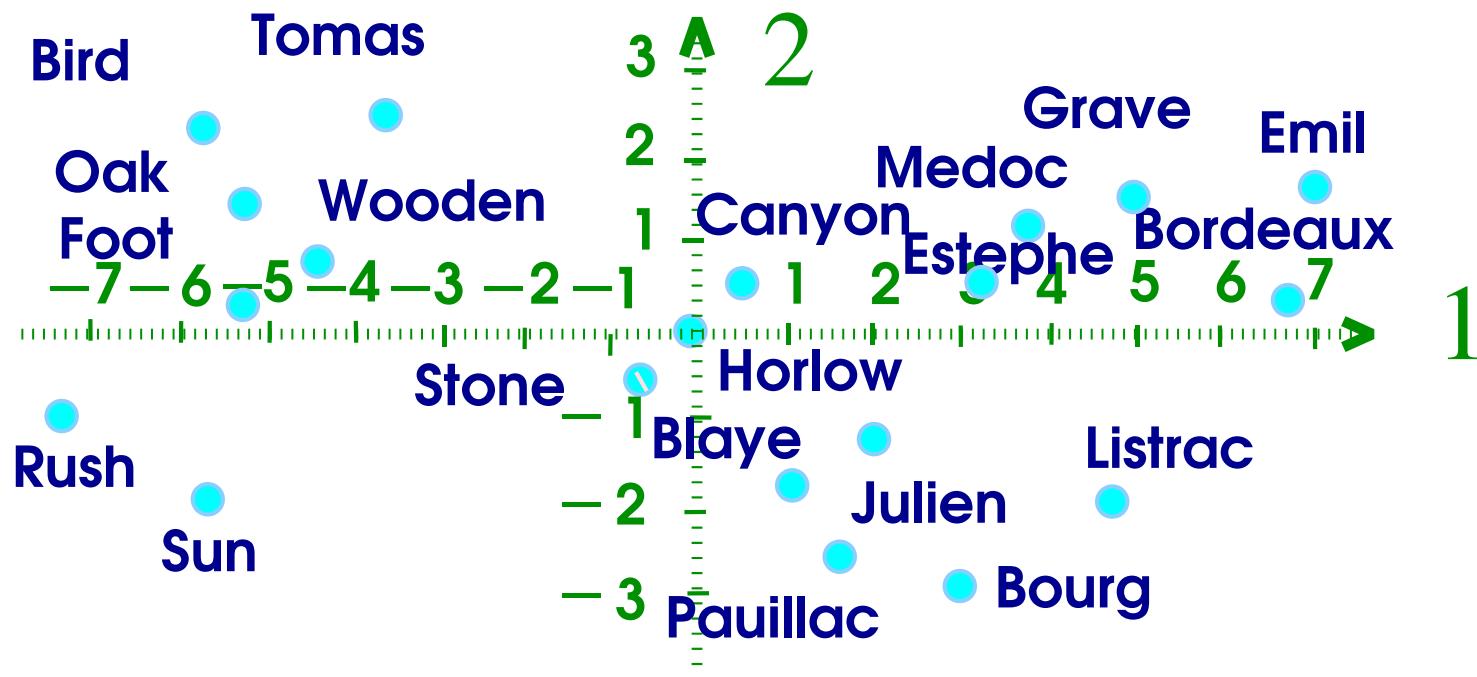
## A PICTURE ... PC1&amp;2



## A PICTURE ... PC 1



## A PICTURE ... PC 1 &amp; 2





**THIS IS ALL THERE IS TO PCA!**

# HOW TO INTERPRET THE DIMENSIONS?

## 20 WINES: ALL DATA



Here are the data

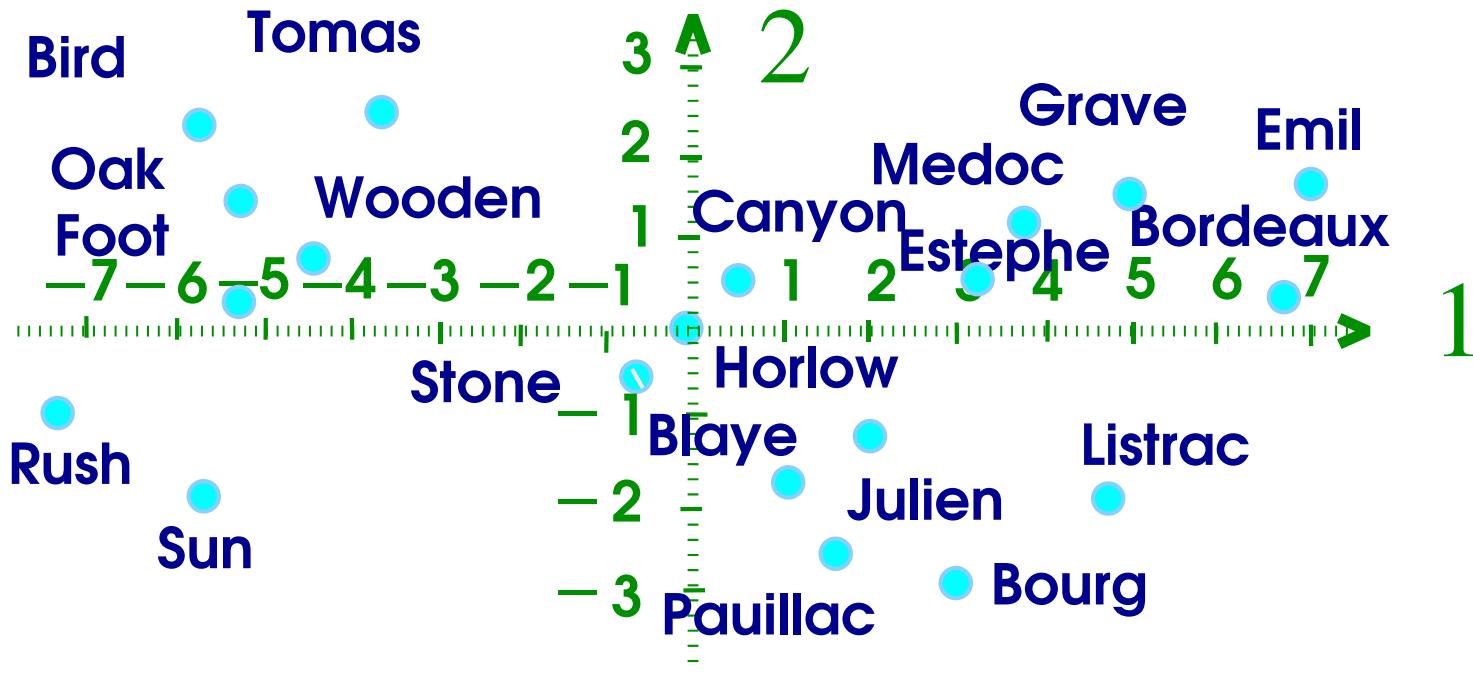
In case you had forgotten. USA vs France

(Long) Name	1. Bordeaux	2. Stone	3. Listrac	4. Canyon Creek	5. Bourg	6. Hill	7. Hollow	8. Estephe	9. Wooden Hill	10. Blaye	11. Sun Set	12. Black Bird	13. Medoc	14. St. Julien	15. Pauillac	16. Rush	17. Oak	18. Grave	19. Emil	20. Tomassello
Short Name	1. Bordeaux	2. Stone	3. Listrac	4. Canyon	5. Bourg	6. Hill	7. Hollow	8. Estephe	9. Wooden	10. Blaye	11. Sun Set	12. Black Bird	13. Medoc	14. St. Julien	15. Pauillac	16. Rush	17. Oak	18. Grave	19. Emil	20. Tomassello
F: Sugar	9	6	2	6	2	9	9	4	8	11	5	9	4	12	9	1	4	13	15	10
W: Astringency	14	7	11	9	9	4	8	11	5	8	4	11	5	12	9	1	4	13	15	6
Origin	F	U	F	U	F	U	U	F	U	F	U	U	F	F	F	U	U	F	F	U

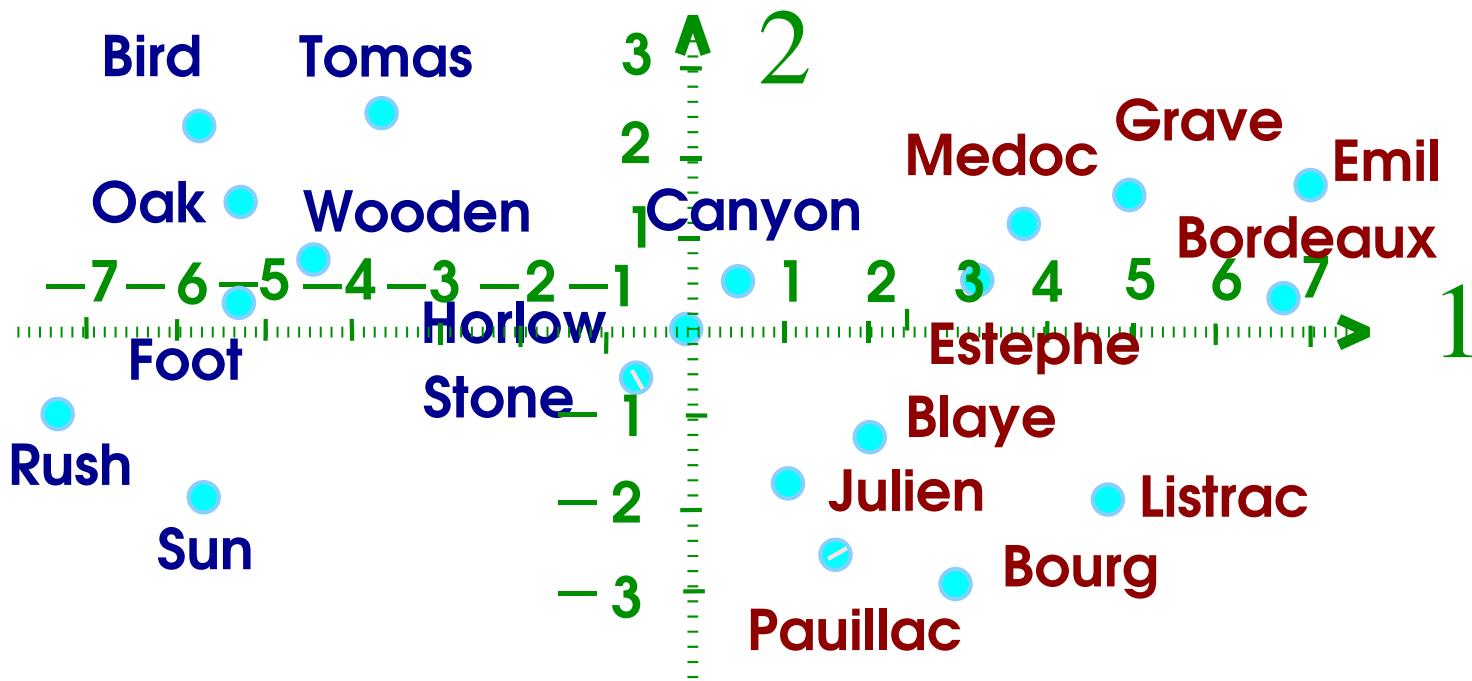
## A PICTURE ... PC 1 & 2

How to interpret the dimensions?

Use **COLORS**

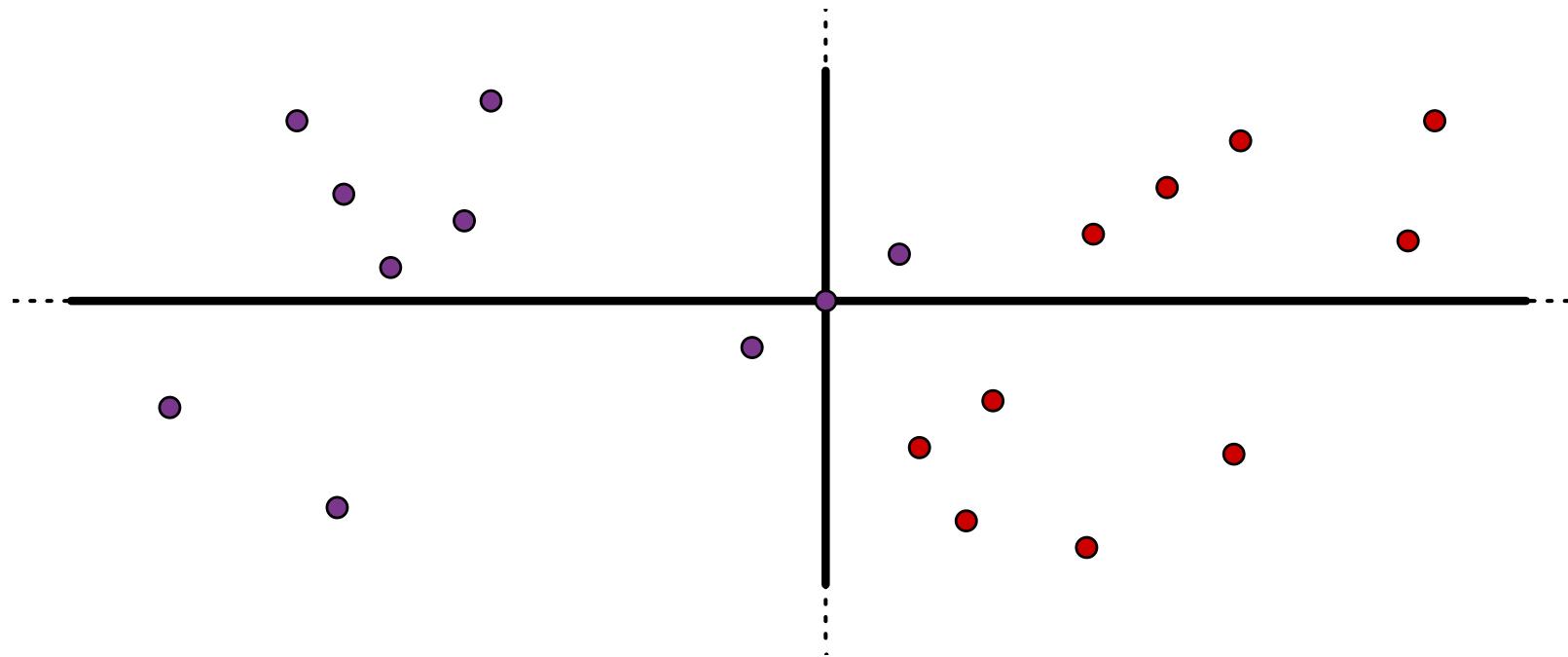


## USA versus France

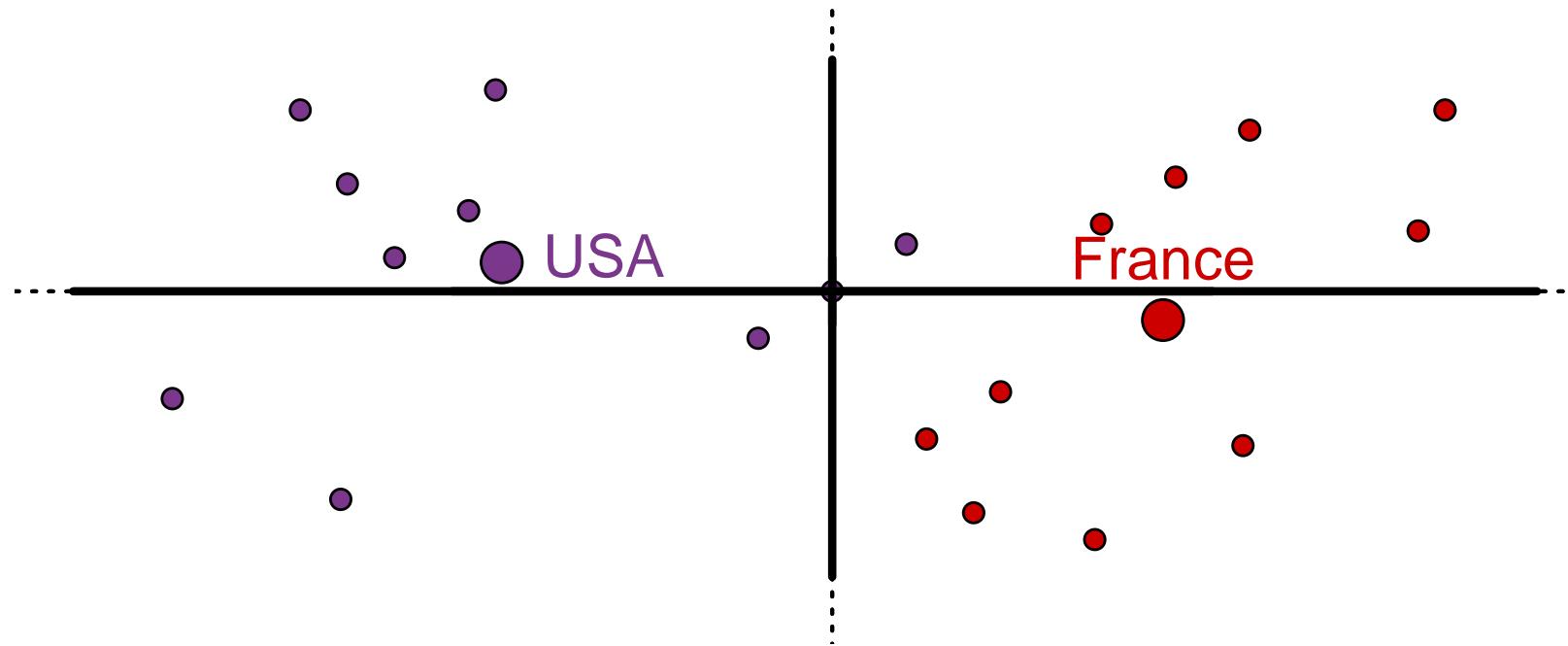


## GROUPS AS MEANS

# WITH GROUP MEANS AS LABELS



# WITH GROUP MEANS AS LABELS



## 20 WINES: ALL DATA



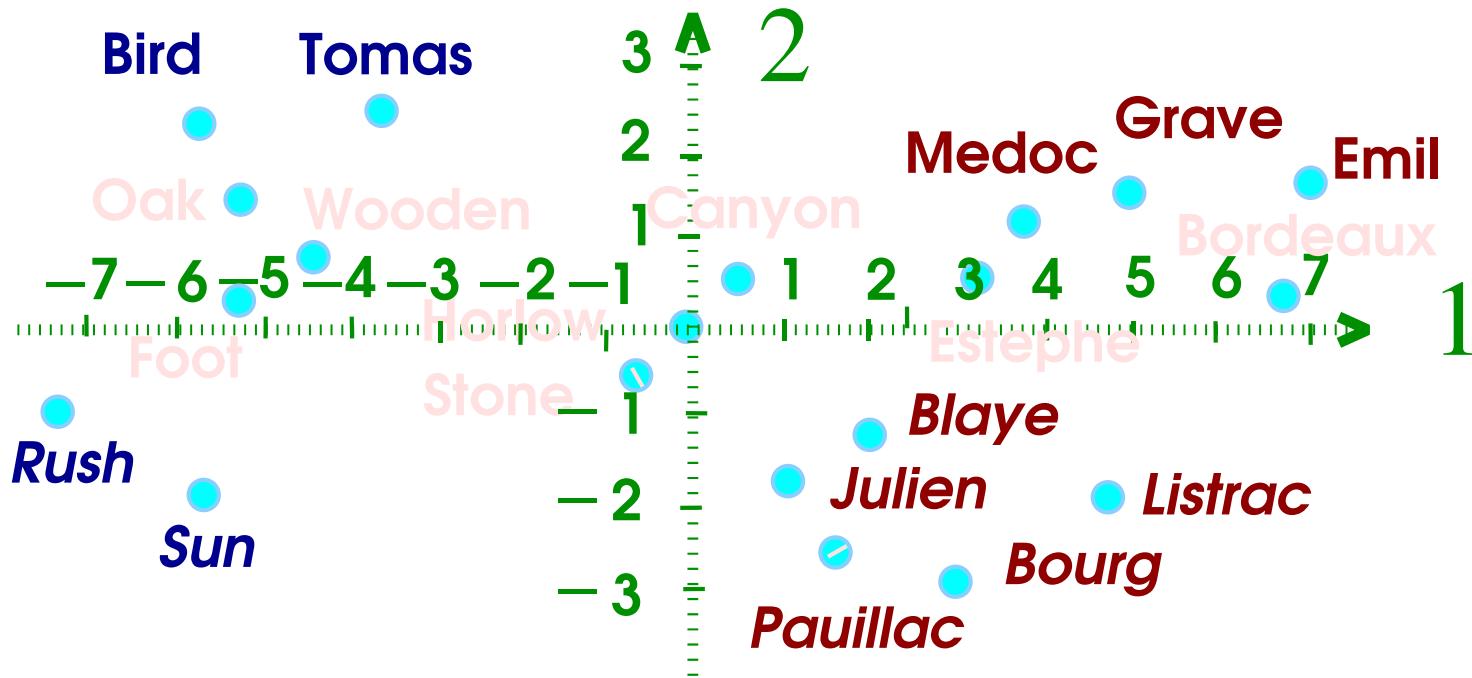
Here are the data  
In case you had forgotten.

(Long) Name	1. Bordeaux	2. Black Stone	3. Listrac	4. Canyon Creek	5. Clos de Boug	6. Fox Hill	7. Hordow	8. St Etéphe	9. Wooden Hill	10. Blaye	11. Côtes de Blaye	11. Sun Set	12. Black Bird	13. Medoc	14. St. Julien	15. Pauillac	16. Rush	17. Oak	18. Gravé	19. Emil	20. Thomas
Short Name	L.Bordeaux	S.Stone	L.Listrac	C.Canyon	B.Boug	H.Hill	H.Hordow	E.Stéphe	W.Wooden	B.Blaye	C.Côtes	S.Sun	B.Black	M.Medoc	J.StJ	P.Pauillac	R.Rush	O.Oak	G.Gravé	E.Emil	T.Thomas
F: Sugar	3	6	2	6	2	9	6	5	9	4	7	11	4	12	5	9	3	10	5	4	10
W: Astringency	14	7	11	9	9	4	8	11	5	8	2	2	4	12	9	8	1	4	13	15	6
Origin	F	U	F	U	F	U	U	F	U	F	U	U	F	F	F	F	U	U	F	F	U
Fruity		F		F				F	F	F	F		F	F	F	F					
Woody													W	W				W	W	W	
Acid	2	5	1	1	9	1	2	2	1	1	2	1	2	1	9	14	2	1	2	1	1
Bitter	8	3	16	3	11	1	1	9	1	8	3	2	2	9	12	10	2	1	10	7	3

# SOME WINES ARE FRUITY, SOME ARE WOODY

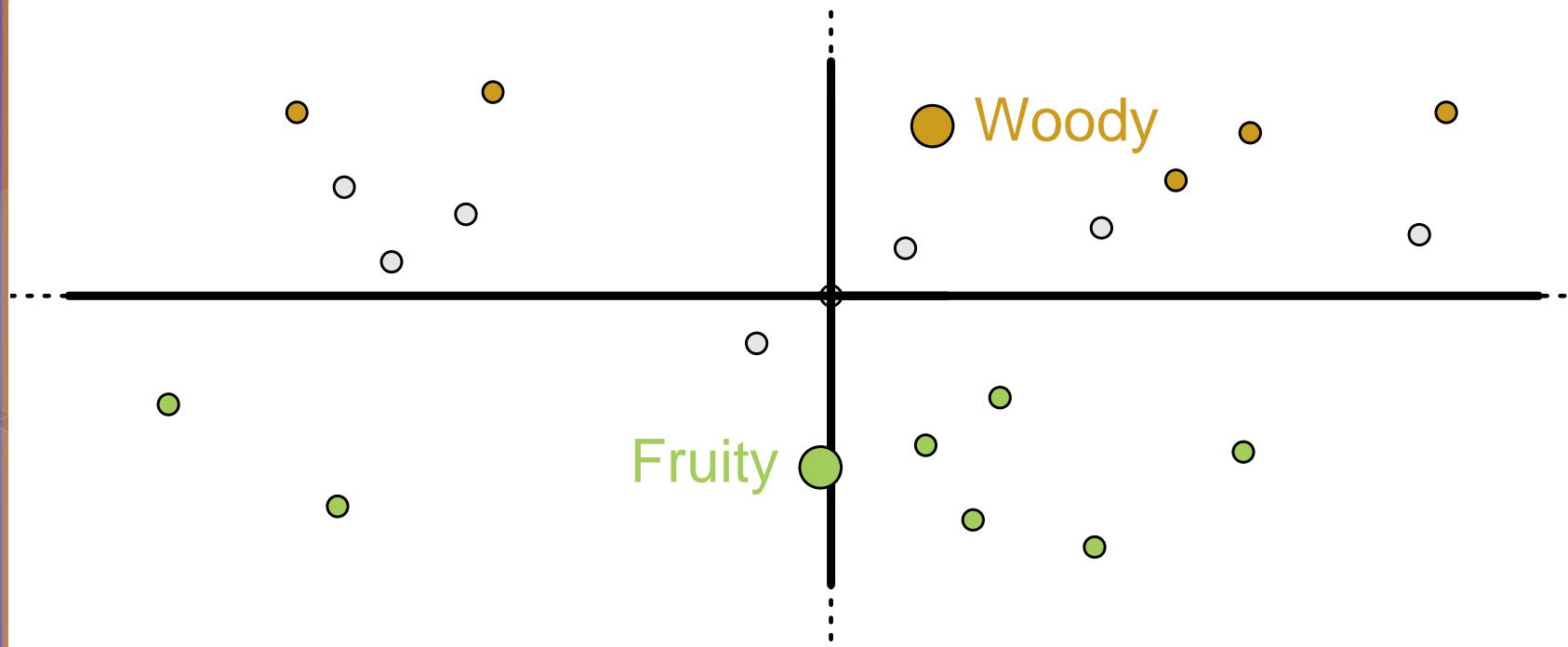
# SOME WINES ARE FRUITY, SOME ARE WOODY

## Woody versus *Fruity*



## WOODY-FRUITY OH RUTTY

SOME WINES ARE FRUITY, SOME ARE WOODY

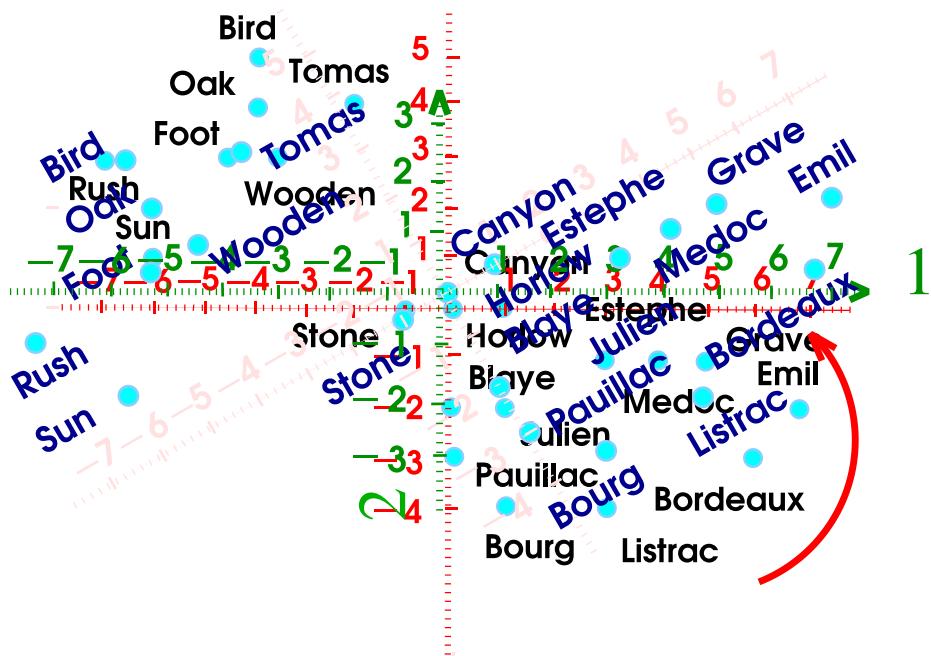


# A BIT MORE FORMAL

# VOCABULARY AND SOME GREEK

## FACTOR SCORES

	<i>Y</i>	<i>W</i>	<i>y</i>	<i>w</i>
Bordeaux	3	14	-3	6
Black Stone	6	7	0	-1
Listrac	2	11	-4	3
Canyon Creek	6	9	0	1
Côtes de Bourg	2	9	-4	1
Foot Hill	9	4	3	-4
Hollow	6	8	0	0
St. Estphe	5	11	-1	3
Wooden Hill	9	5	3	-3
Côtes de Blaye	4	8	-2	0
Sun Set	7	2	1	-6
Black Bird	11	4	5	-4
Médoc	5	12	-1	4
St Julien	4	9	-2	1
Pauillac	3	8	-3	0
Gold Rush	9	1	3	-7
Oak Ville	10	4	4	-4
Grave	5	13	-1	5
St Emilion	4	15	-2	7
Tomasello	10	6	4	-2
<b><math>\Sigma</math></b>	<b>120</b>	<b>160</b>	<b>0</b>	<b>0</b>

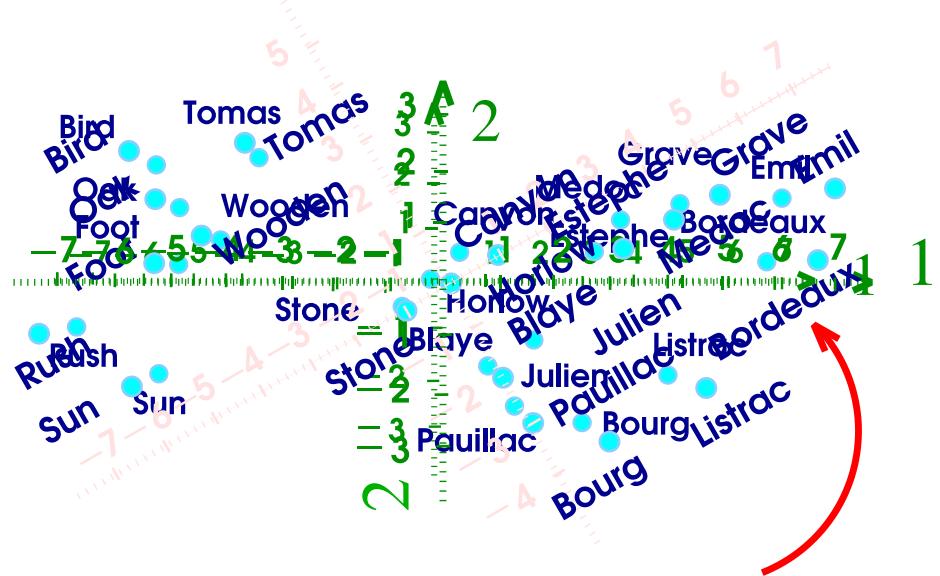


## FACTOR SCORES

# CALL

# THE COORDINATES: FACTOR SCORES

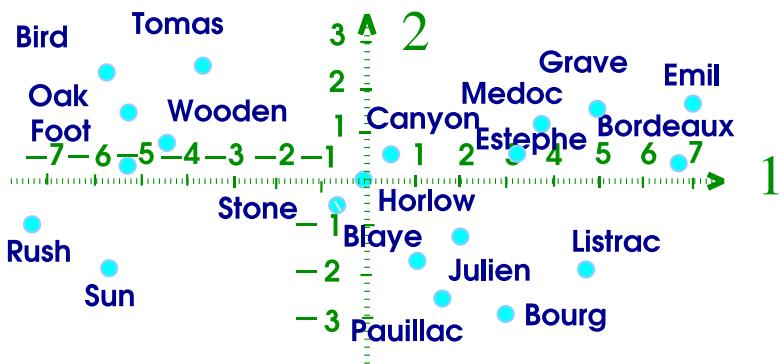
	<i>Y</i>	<i>W</i>	<i>y</i>	<i>w</i>	<i>F</i> <sub>1</sub>	<i>F</i> <sub>2</sub>
Bordeaux	3	14	-3	6	6.67	0.69
Black Stone	6	7	0	-1	-0.84	-0.54
Listrac	2	11	-4	3	4.68	-1.76
Canyon Creek	6	9	0	1	0.84	0.54
Côtes de Bourg	2	9	-4	1	2.99	-2.84
Foot Hill	9	4	3	-4	-4.99	0.38
Horlow	6	8	0	0	0.00	0.00
St. Estphe	5	11	-1	3	3.07	0.77
Wooden Hill	9	5	3	-3	-4.14	0.92
Côtes de Blaye	4	8	-2	0	1.07	-1.69
Sun Set	7	2	1	-6	-5.60	-2.38
Black Bird	11	4	5	-4	-6.06	2.07
Médoc	5	12	-1	4	3.91	1.30
St Julien	4	9	-2	1	1.92	-1.15
Pauillac	3	8	-3	0	1.61	-2.53
Gold Rush	9	1	3	-7	-7.52	-1.23
Oak Ville	10	4	4	-4	-5.52	1.23
Grave	5	13	-1	5	4.76	1.84
St Emilion	4	15	-2	7	6.98	2.07
Tomasello	10	6	4	-2	-3.83	2.30
<b><math>\Sigma</math></b>	<b>120</b>	<b>160</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>



## FACTOR SCORES

# CALL SUM OF SQUARED FACTOR SCORES: EIGENVALUE

	<i>F</i>	<i>W</i>	<i>y</i>	<i>w</i>	<i>F<sub>1</sub></i>	<i>F<sub>2</sub></i>	<i>F<sub>1</sub><sup>2</sup></i>	<i>F<sub>2</sub><sup>2</sup></i>
Bordeaux	3	14	-3	6	6.67	0.69	44.52	0.48
Black Stone	6	7	0	-1	-0.84	-0.54	0.71	0.29
Listrac	2	11	-4	3	4.68	-1.76	21.89	3.11
Canyon Creek	6	9	0	1	0.84	0.54	0.71	0.29
Côtes de Bourg	2	9	-4	1	2.99	-2.84	8.95	8.05
Foot Hill	9	4	3	-4	-4.99	0.38	24.85	0.15
Horlow	6	8	0	0	0.00	0.00	0	0.00
St Estiphe	5	11	-1	3	3.07	0.77	9.41	0.59
Wooden Hill	9	5	3	-3	-4.14	0.92	17.15	0.85
Côtes de Blaye	4	8	-2	0	1.07	-1.69	1.15	2.85
Sun Set	7	2	1	-6	-5.60	-2.38	31.35	5.65
Black Bird	11	4	5	-4	-6.06	2.07	36.71	4.29
Médoc	5	12	-1	4	3.91	1.30	15.30	1.70
St Julien	4	9	-2	1	1.92	-1.15	3.68	1.32
Pauillac	3	8	-3	0	1.61	-2.53	2.59	6.41
Gold Rush	9	1	3	-7	-7.52	-1.23	56.49	1.51
Oak Ville	10	4	4	-4	-5.52	1.23	30.49	1.51
Grave	5	13	-1	5	4.76	1.84	22.61	3.39
St Emilion	4	15	-2	7	6.98	2.07	48.71	4.29
Tomasello	10	6	4	-2	-3.83	2.30	14.71	5.29
$\Sigma$	120	160	0	0	0	0	392	52
							$\lambda_1$	$\lambda_2$





**EXPLAINED VARIANCE?**

**THE REVENGE OF PYTHAGORAS**

**THE RETURN OF THE EIGEN-FAIRY**

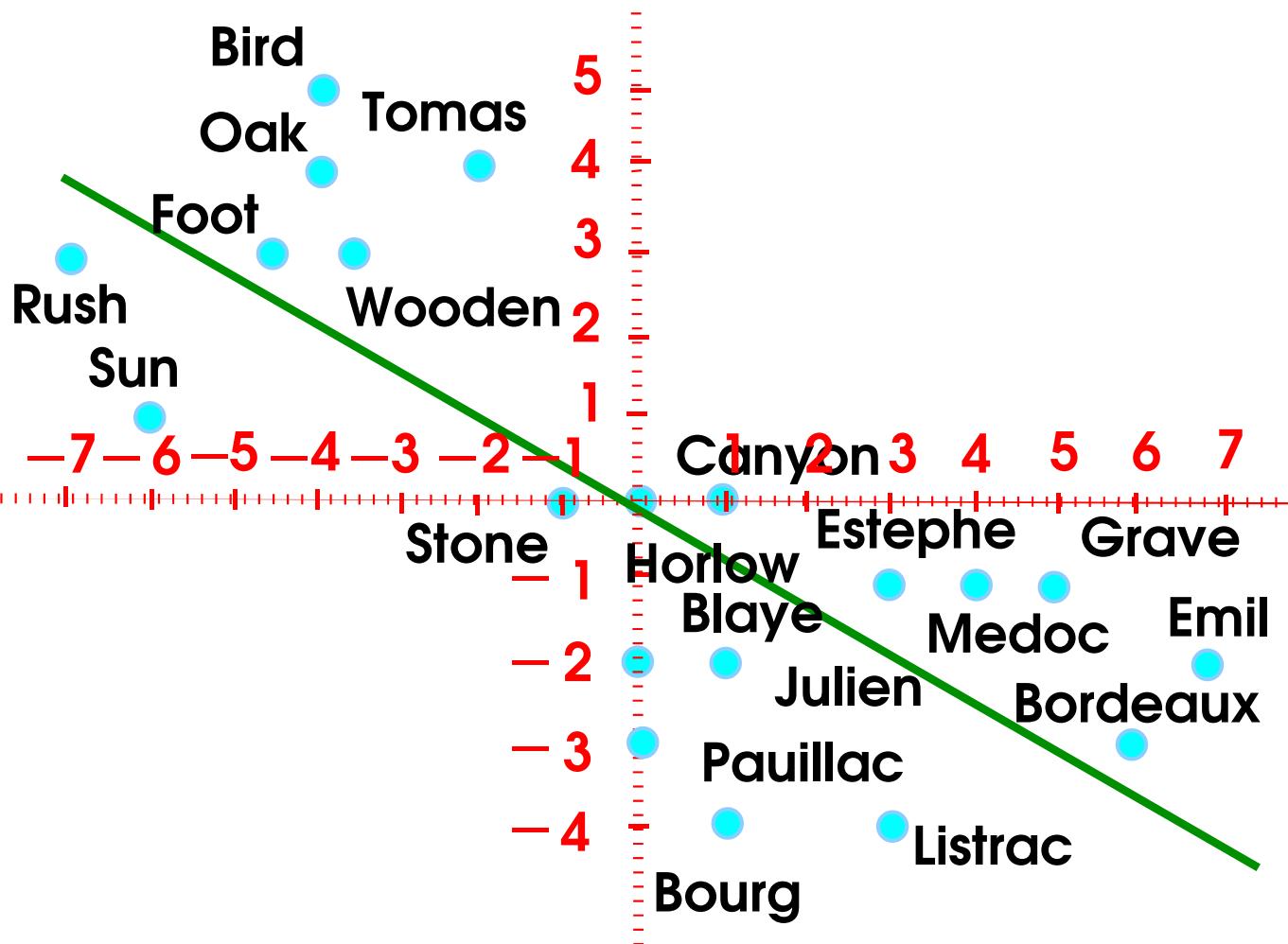
## THE EIGEN FAIRY

# BACK TO THE EIGEN FAIRY

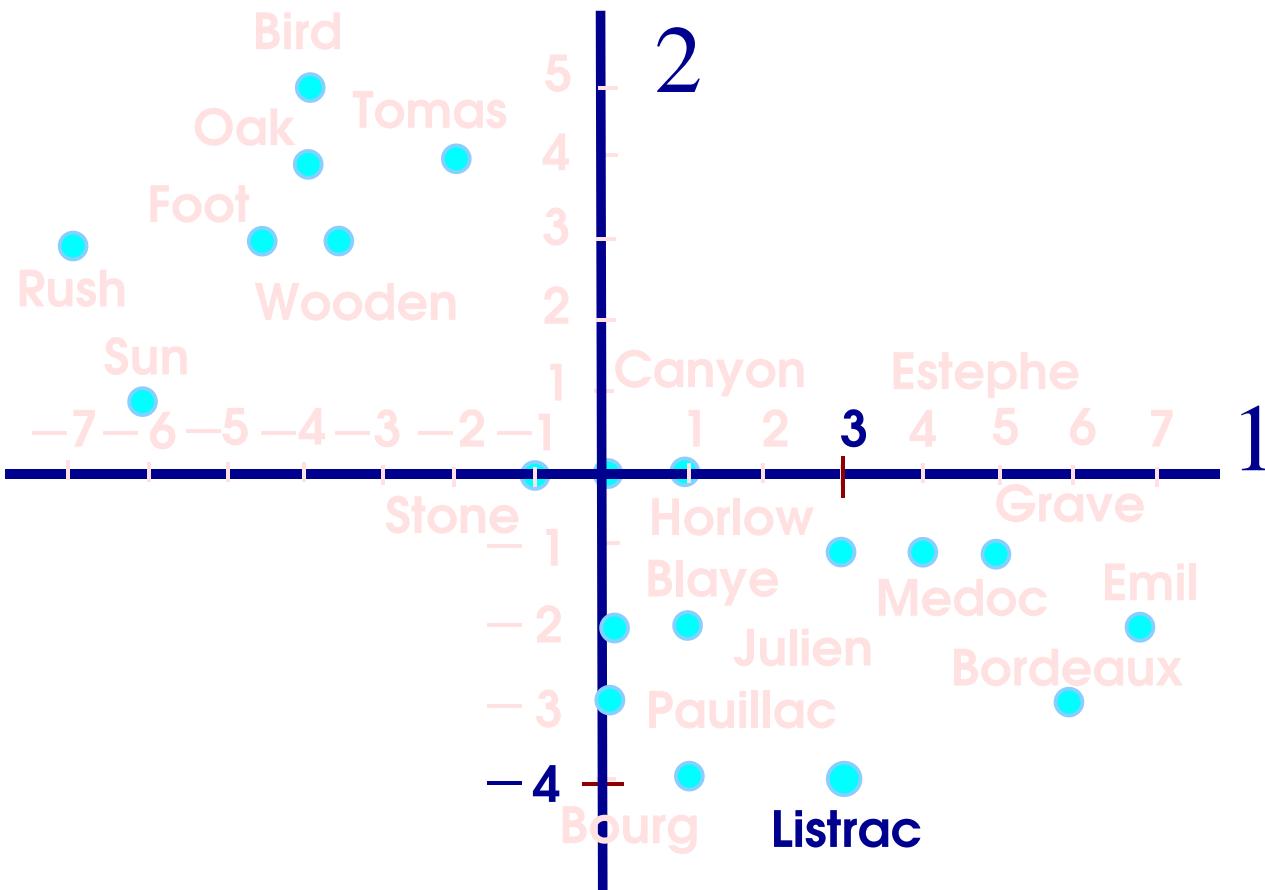


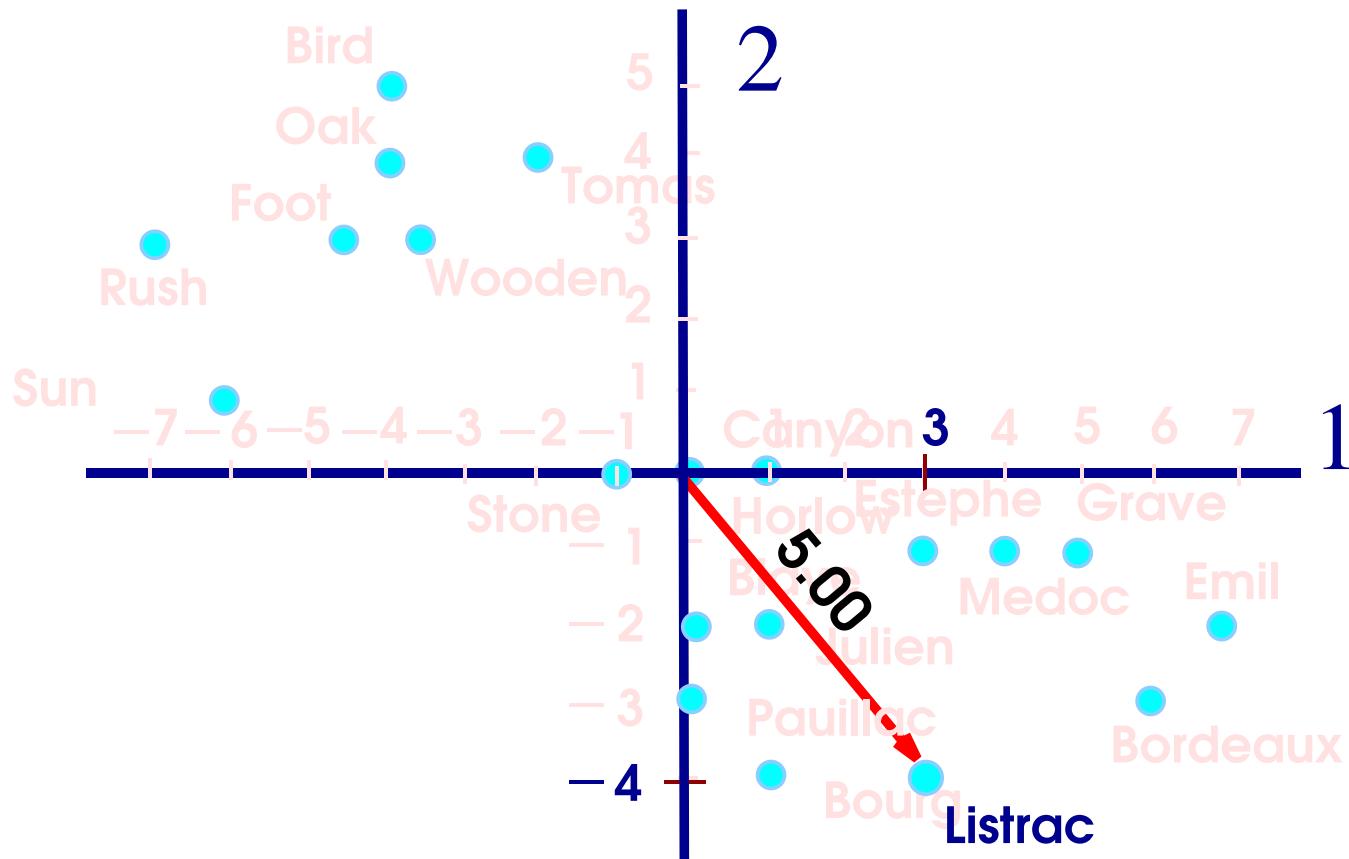
WHAT DID WE ASK FOR?

EIGEN FAIRY, PLEEEEZ  
FIND THE ROTATION  
WITH SMALLEST  
SQUARED DISTANCES

**A PICTURE ... PC 1**

# THE SQUARE DISTANCES





**REMEMBER**

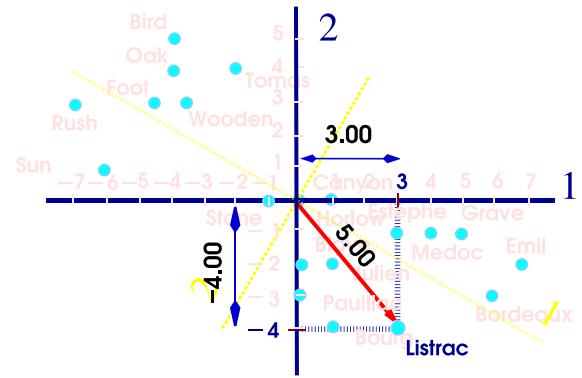
# **THE THEOREM!**

# SLICES OF INERTIA

# **INERTIA: WHAT DO WE DECOMPOSE IN PCA?**

$$d_{\text{Listrac,g}}^2 = w_{\text{Listrac,g}}^2 + y_{\text{Listrac,g}}^2$$

$$3^2 + (-4)^2 = 9 + 16 = 25$$



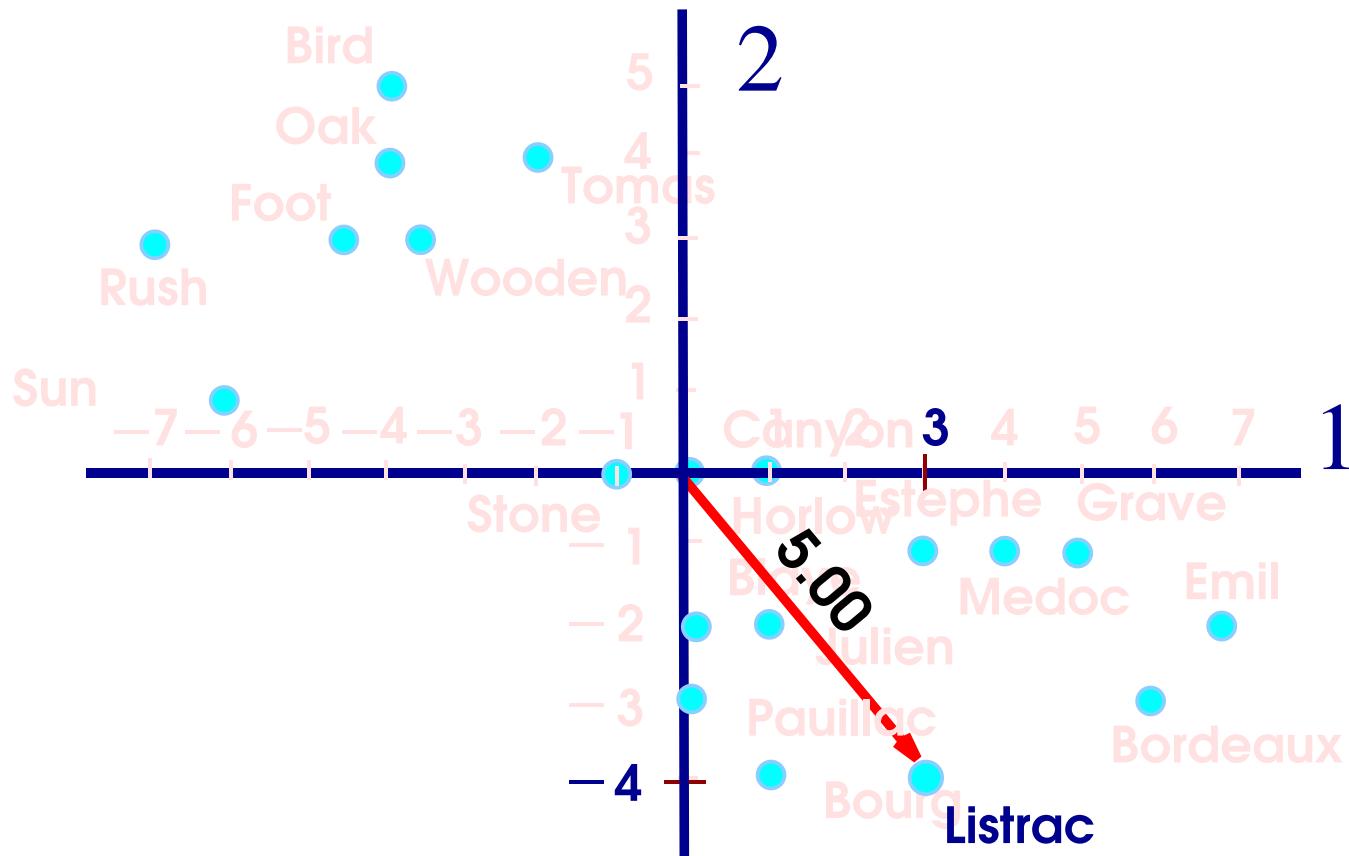
# **INERTIA: SUM OF SQUARED DISTANCES**

$$\mathcal{I} = \sum_i^I d_{i,\mathbf{g}}^2$$

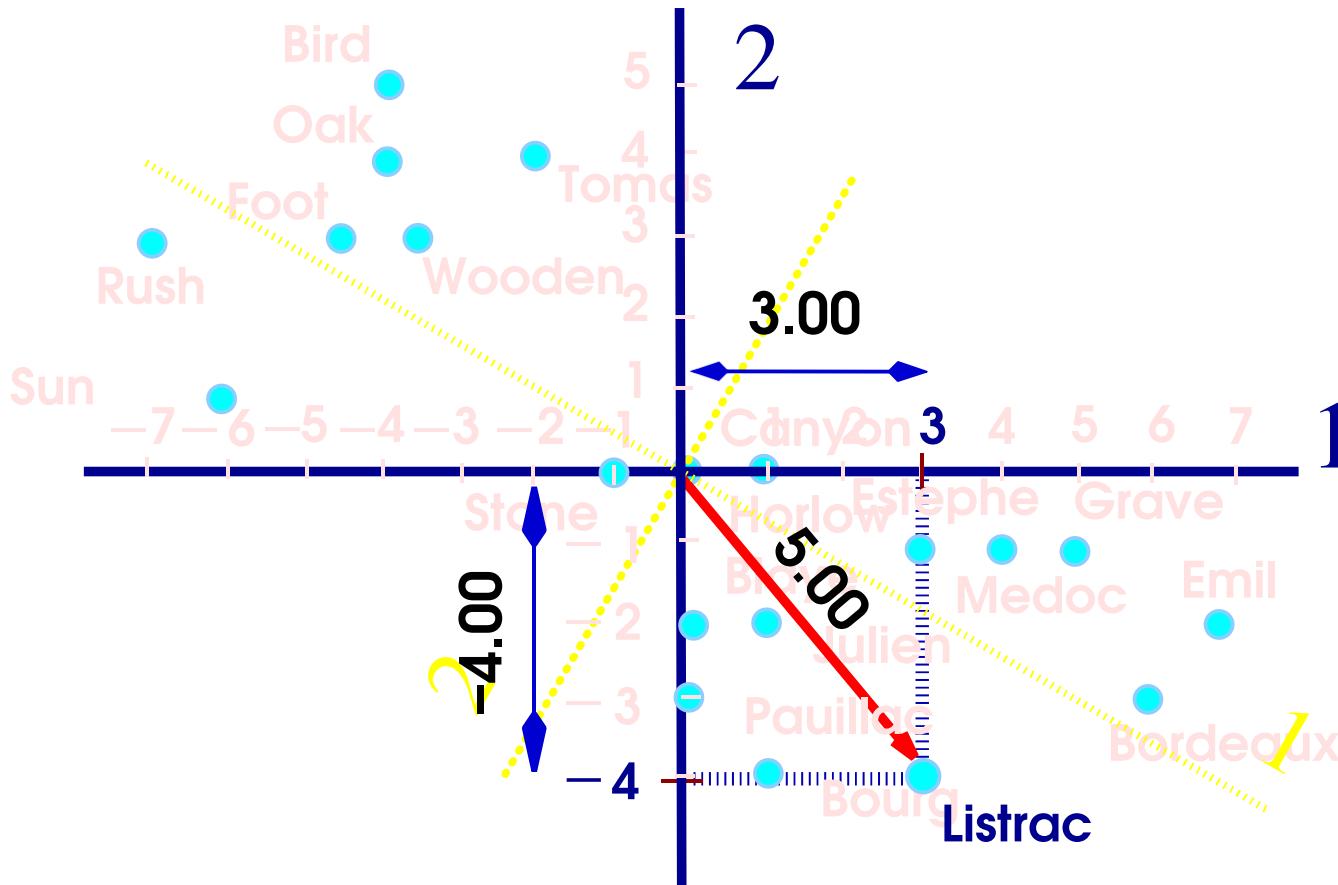
## I FROM VARIABLES

# INERTIA FROM THE VARIABLES

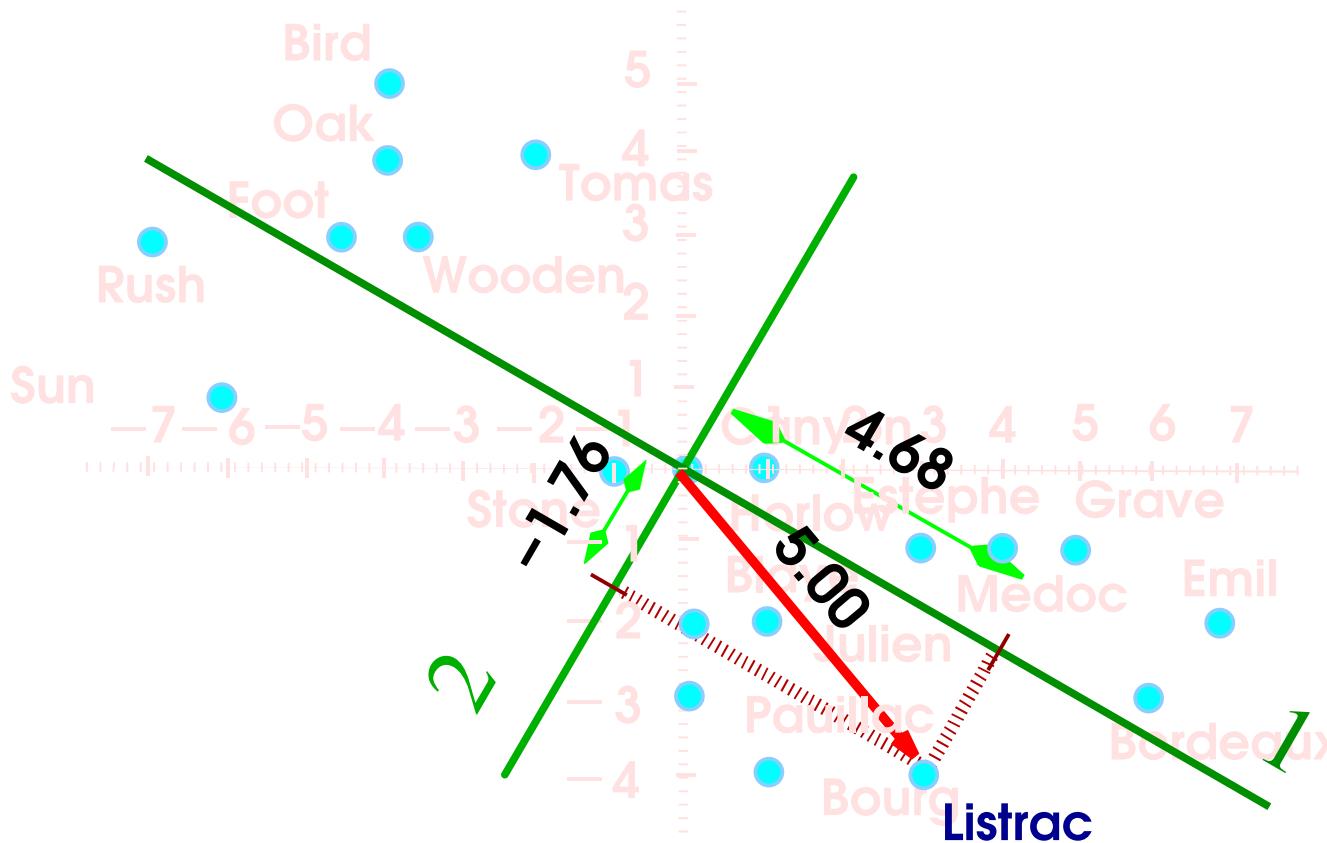
$$\mathcal{I} = \sum_i^I d_{i,\mathbf{g}}^2 = \sum_i^I (w_{i,\mathbf{g}}^2 + y_{i,\mathbf{g}}^2)$$



# DISTANCE FROM VARIABLES



# DISTANCE FROM COMPONENTS

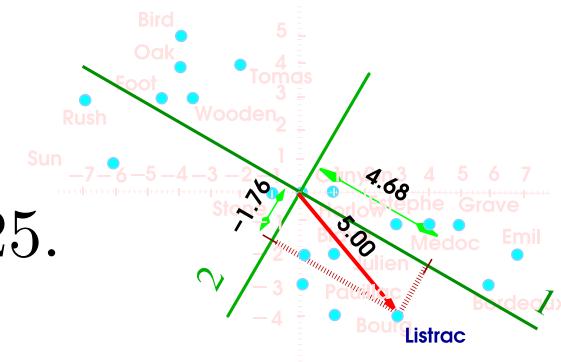


## I FROM COMPONENT

# INERTIA FROM THE COMPONENTS

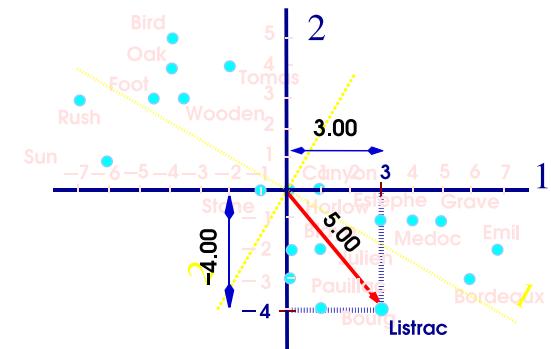
Squared distance from “Listrac” to center of the space  
From components

$$4.68^2 + (-1.76)^2 = 21.89 + 3.11 = 25.$$



Squared distance from “Listrac” to center of the space  
From variable

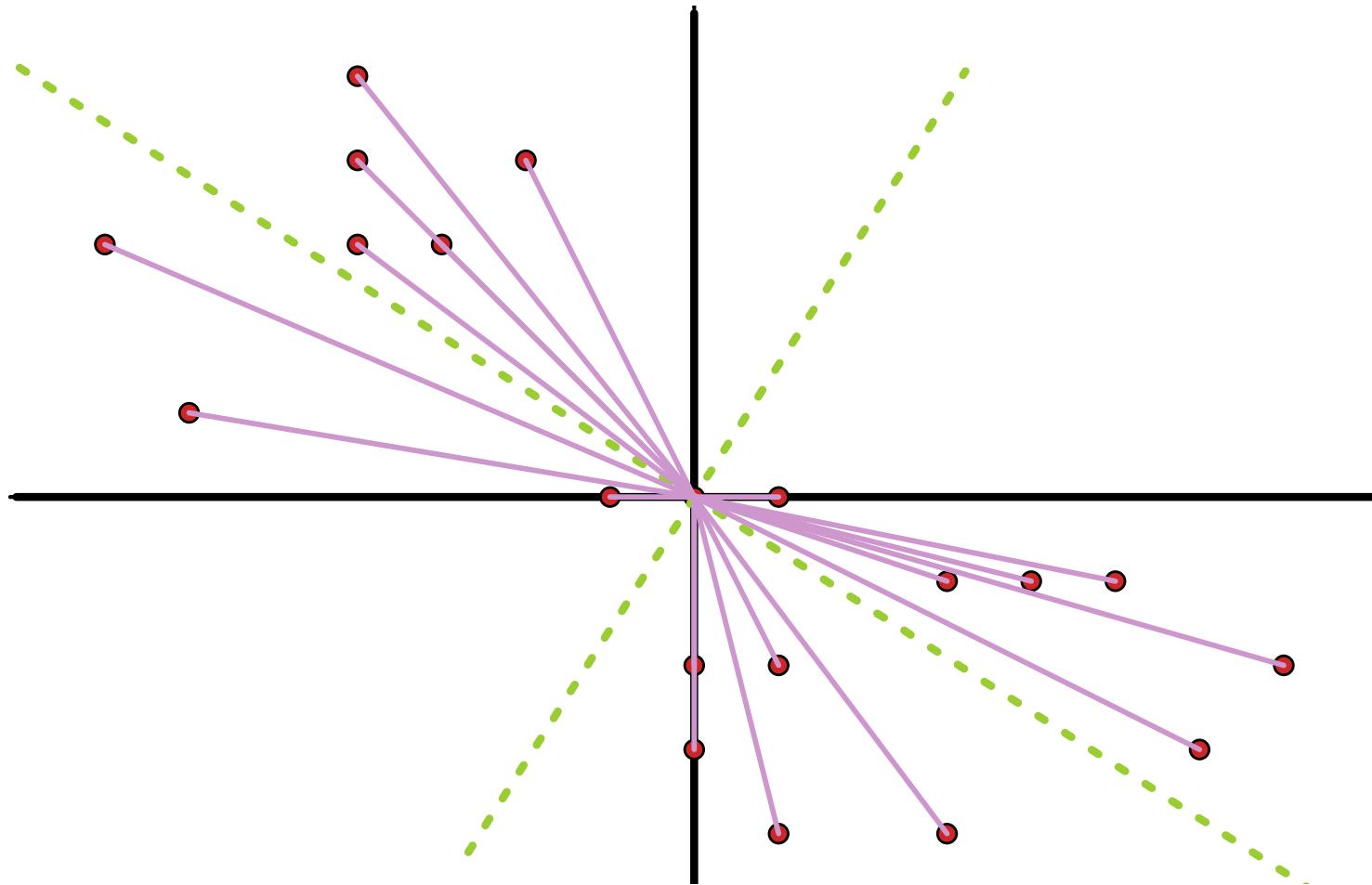
$$3^2 + (-4)^2 = 9 + 16 = 25$$



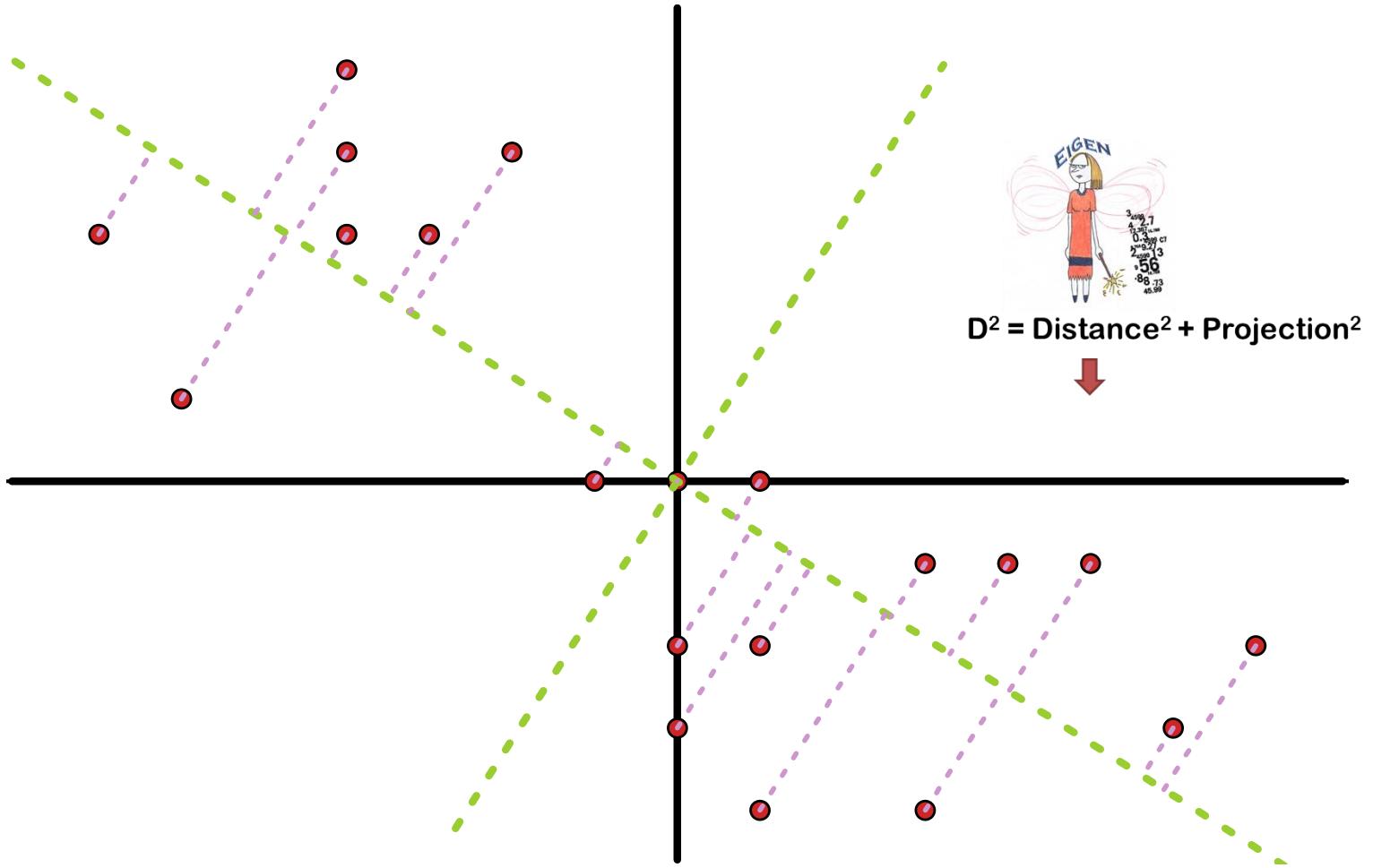
## INERTIA

# ALL OBSERVATIONS: TOTAL INERTIA

# TOTAL VARIANCE (AKA INERTIA)



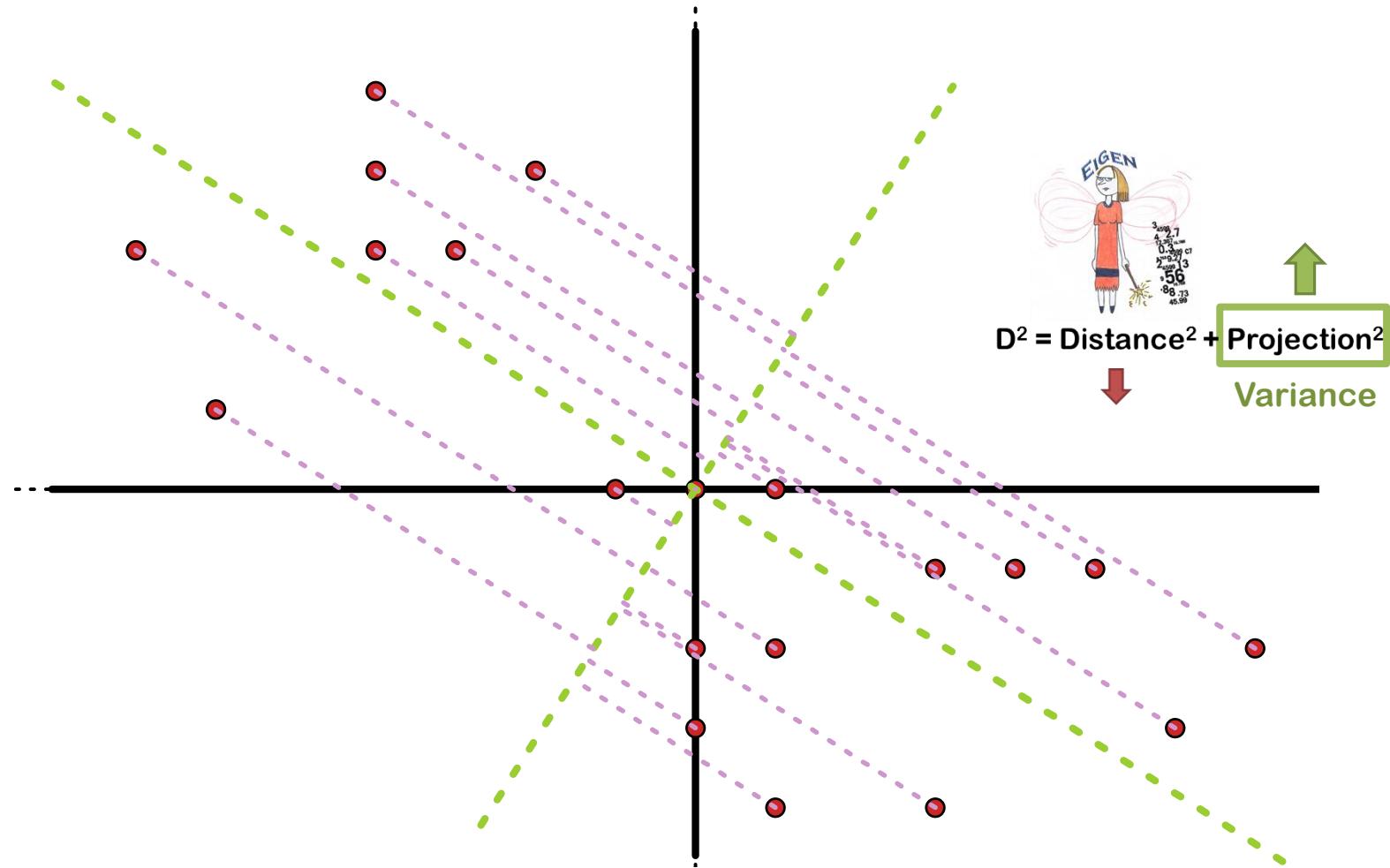
# MIN ERROR VARIANCE (AKA INERTIA)



# PYTHAGOREAN MAGIC

## INERTIA/VARIANCE

# MAX VARIANCE (AKA INERTIA)



$$D^2 = A^2 + B^2$$

## PYTHAGOREAN MAGIC:

SIMILARLY,  
SMALLEST SUM OF SQUARED DISTANCES

SAME AS

LARGEST SQUARED PROJECTIONS (VARIANCE)



$$D^2 = \text{Distance}^2 + \text{Projection}^2$$



# MIN ERROR = MAX VARIANCE

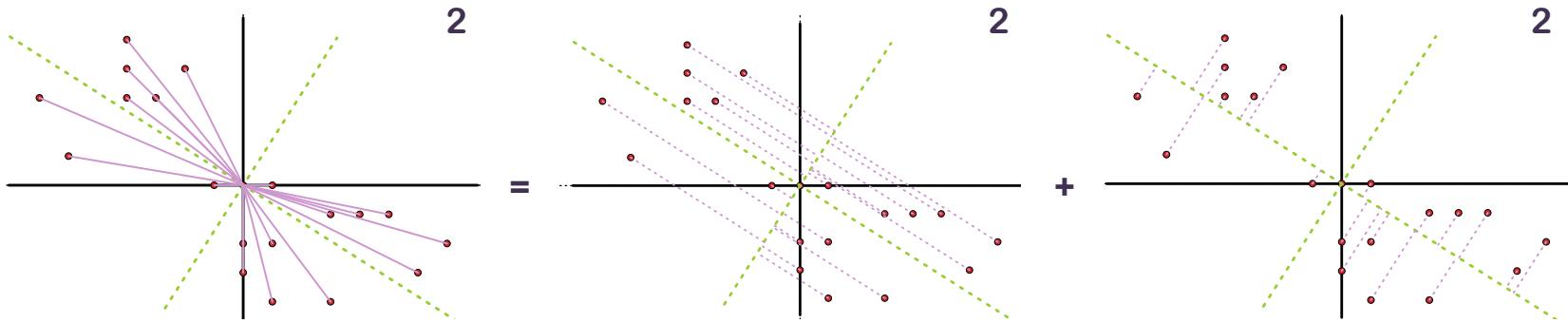
VARIANCE IS CALLED:  
AN EIGEN-VALUE  
A CHARACTERISTIC ROOT  
LATENT ROOT  
PROPER ROOT  
FRENCH: “VALEUR PROPRE”



$$D^2 = A^2 + B^2$$

# PYTHAGOREAN MAGIC:

TOTAL INERTIA = SUM OF EIGENVALUES



## How GOOD?

# GOOD COMPONENTS EXPLAIN A LOT

Component	$\lambda_i$ (eigenvalue)	Cumulated (eigenvalues)	Percent of of inertia	Cumulated (percentage)
1	392	392	83.29	83.29
2	52	444	11.71	100.00

COMPARED TO

Variable	$SS_i$	Cumulative $SS_i$	Percent of of inertia	Cumulative (percentage)
Astringent	294	294	66.22	66.22
Sugar	150	444	33.78	100.00

EIGEN-FAIRY: YOU CANNOT BETTER PCA!

## VARIABLES?

# HOW TO REPRESENT THE VARIABLES?

## PART 1. CORRELATIONS AKA LOADINGS

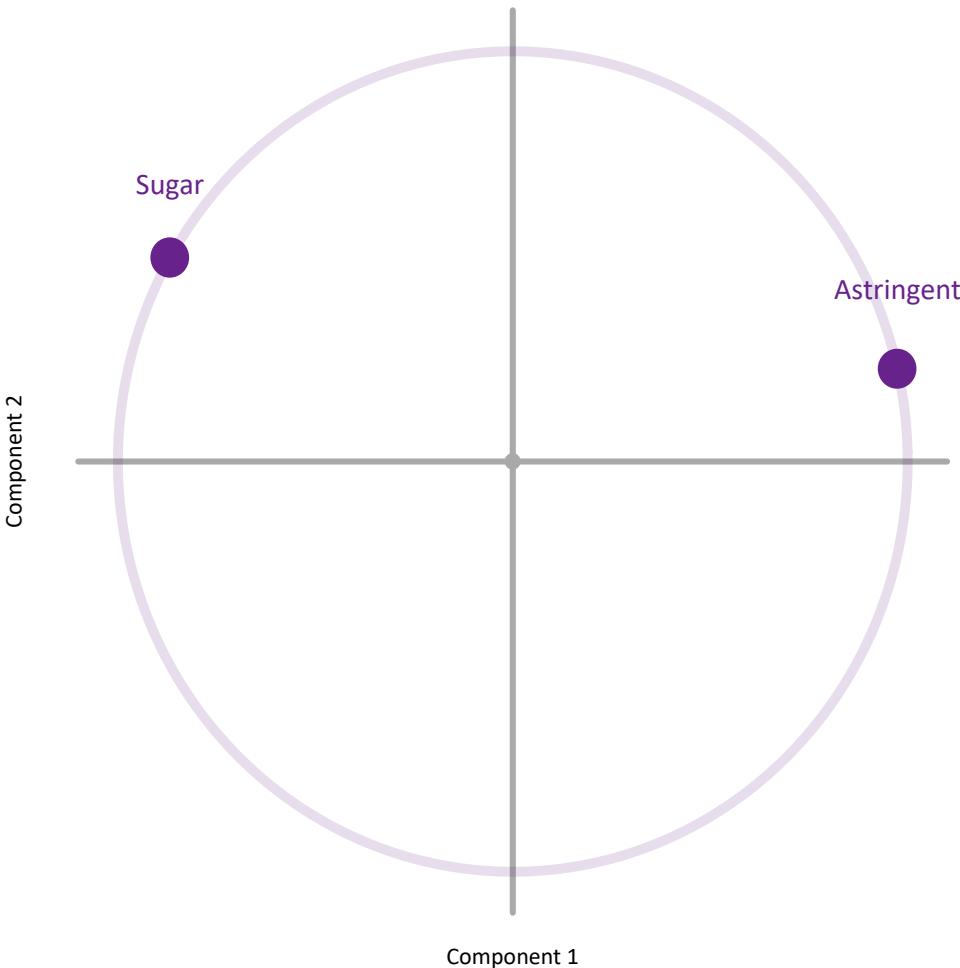
CORRELATION BETWEEN  
VARIABLES AND FACTOR SCORES  
(TRÈS FRENCH, BY THE WAY)

# FIRST COMPUTE THE CORRELATIONS

$Y$ : Astringent.  $W$ : Sugar

Component	Loadings	
	$Y$	$W$
1	.9742	-.8670
2	.2258	.4967
3		

# CIRCLE OF CORRELATION



## DESCARTES'S CIRCLE

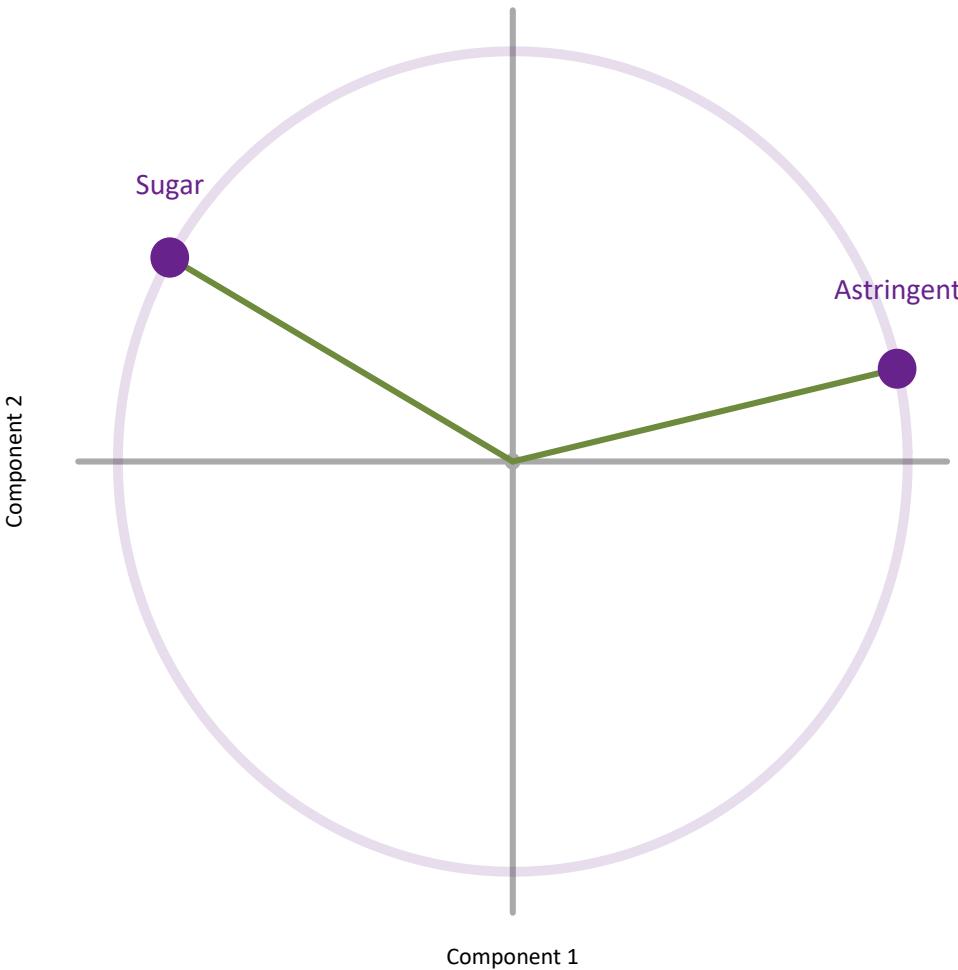
# WHY A CIRCLE OF RADIUS 1?



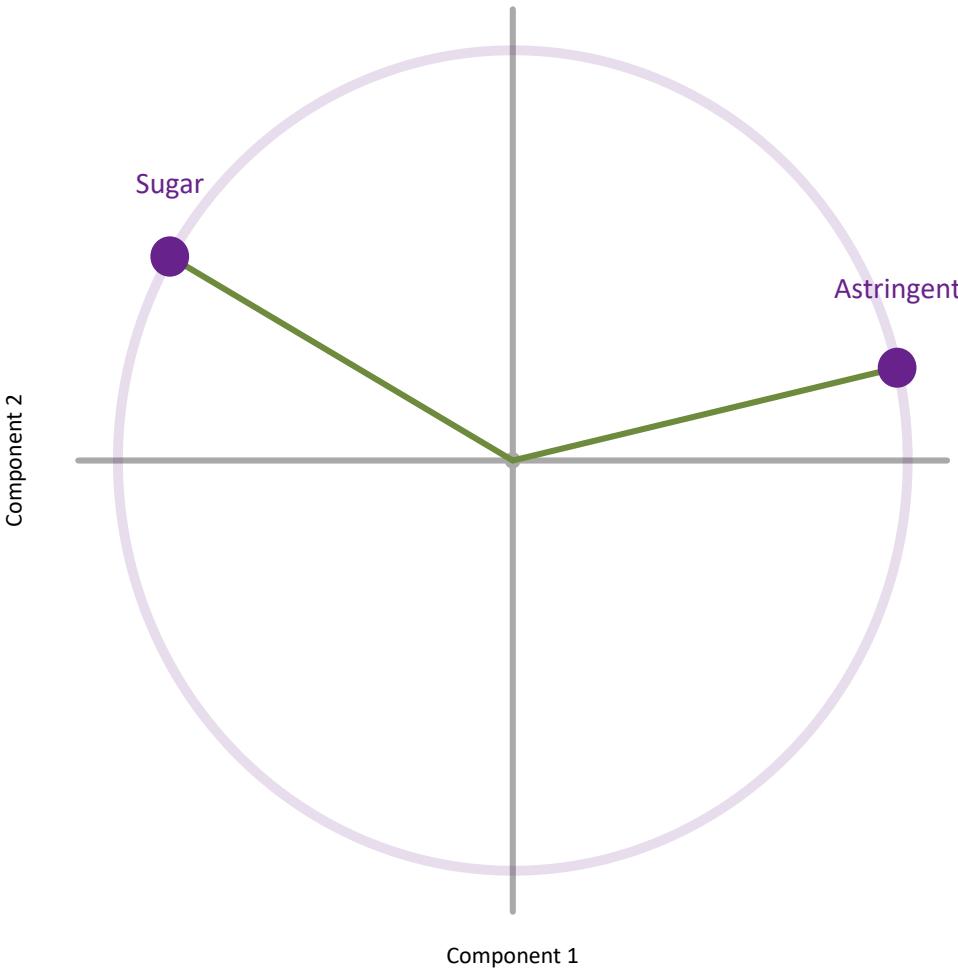
THIS IS DESCARTES'S FAULT

A CIRCLE IS DEFINED AS  $x^2 + y^2 = \text{RADIUS}$

# CIRCLE OF CORRELATION

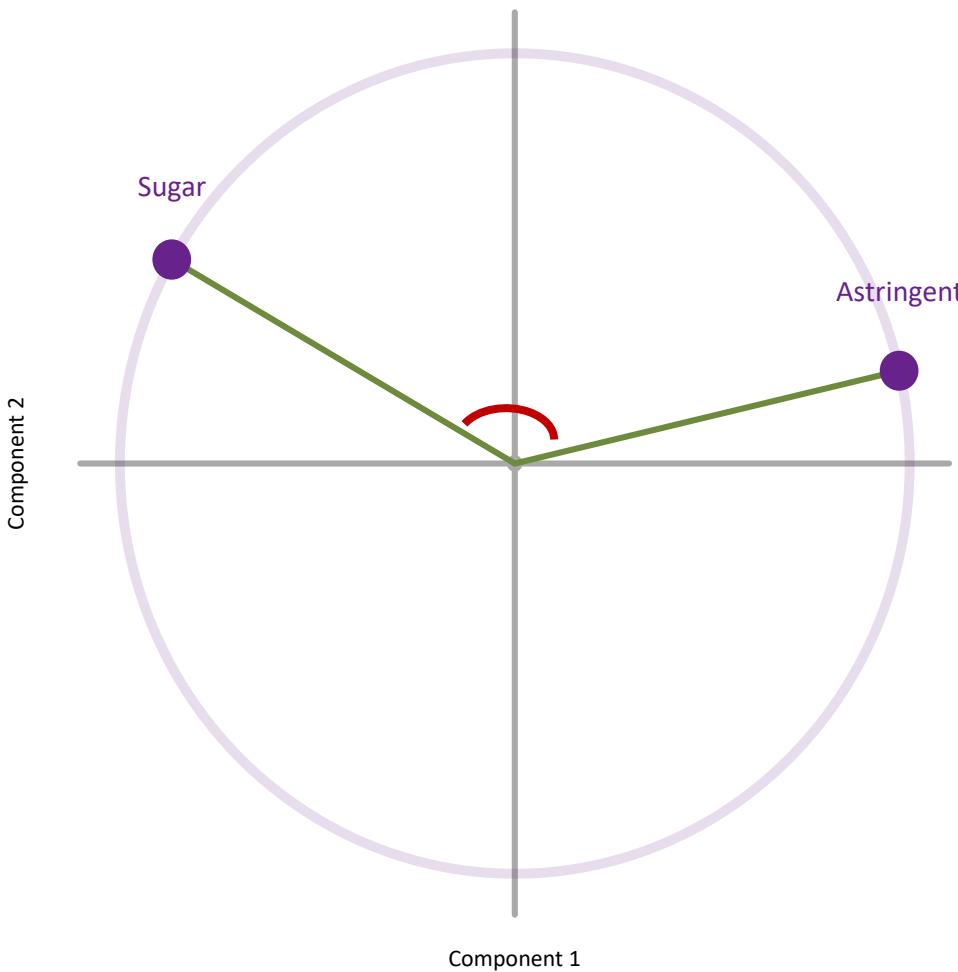


# CIRCLE OF CORRELATION



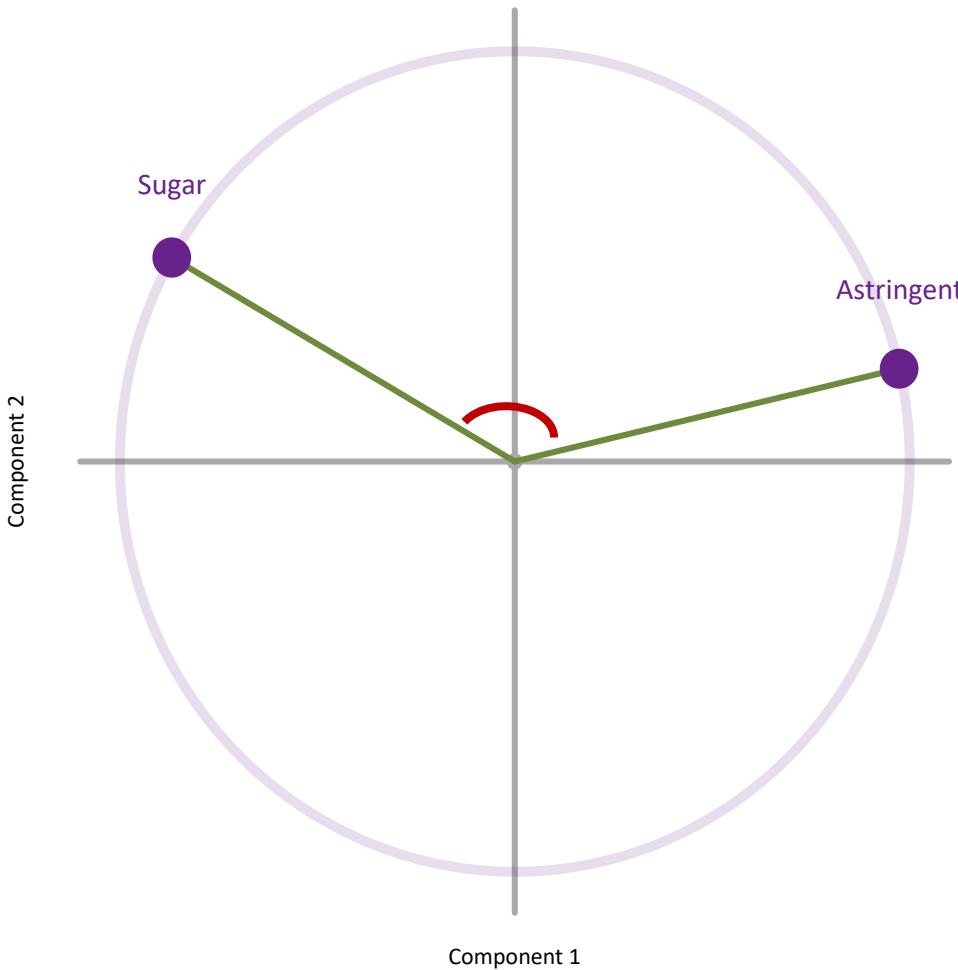
**HOW TO INTERPRET:**

# COSINES APPROXIMATE CORRELATION



**HOW TO INTERPRET:**

# COSINES APPROXIMATE CORRELATION



HOW TO INTERPRET: ANGLE =  $137^\circ$ ,  $R = -.73$

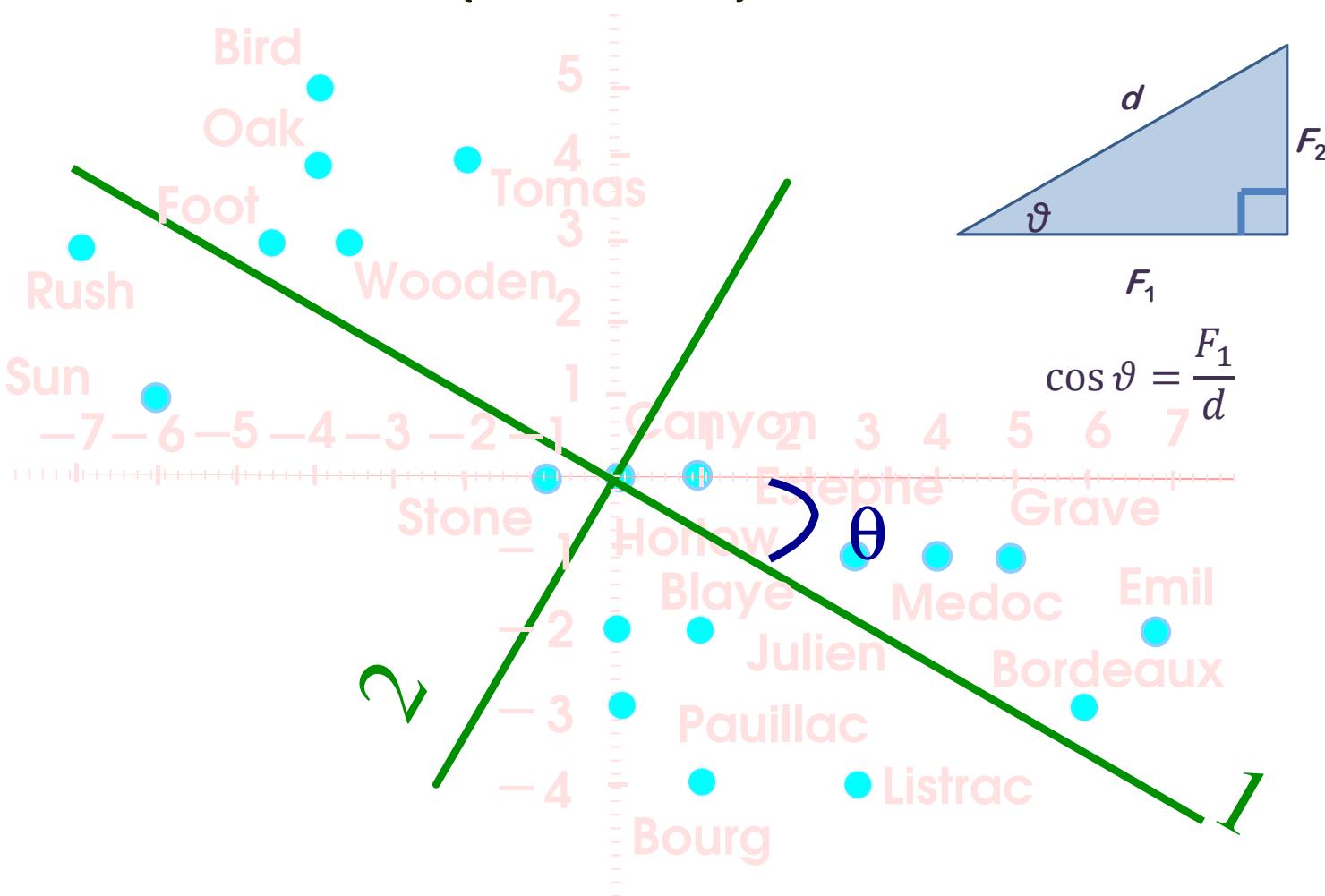
## LOADINGS: AS COSINES

# HOW TO REPRESENT THE VARIABLES?

## PART 2. COSINES (ANGLES) AS LOADINGS

## VARIABLES?

# PART 2. COSINES (ANGLES) AS LOADINGS



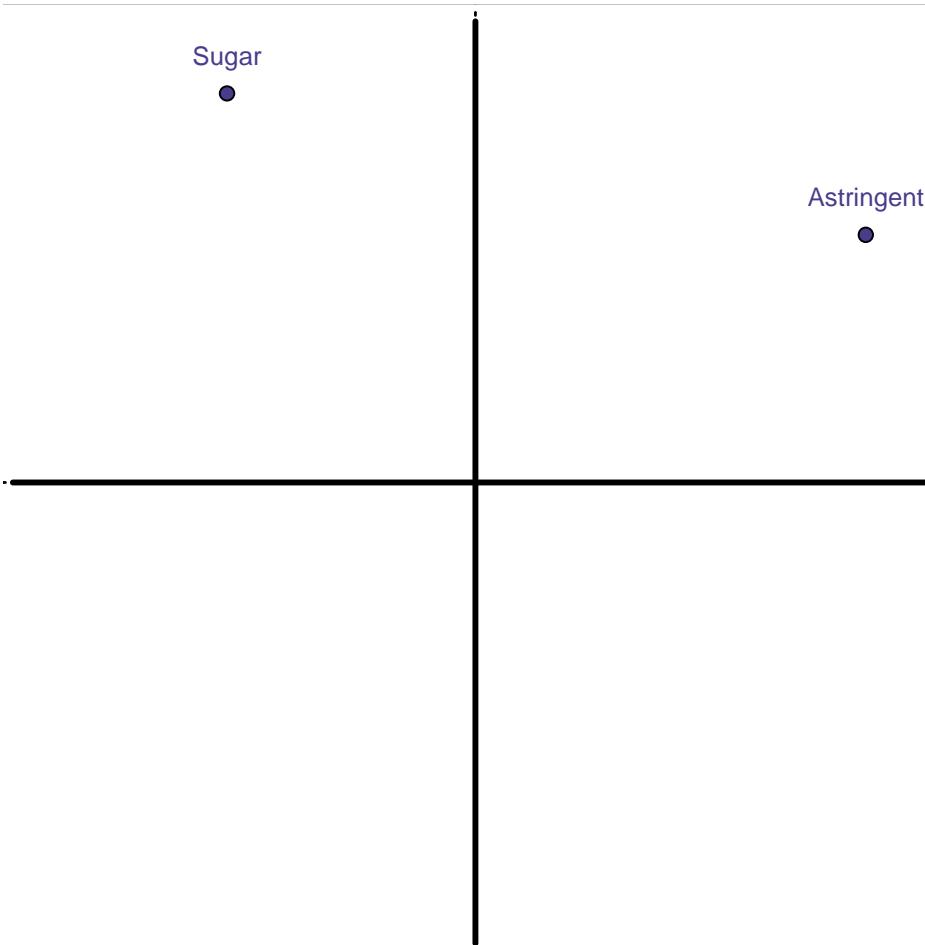
## COSINES

## LOADINGS AS COSINES

$Y$ : Astringent.  $W$ : Sugar

Component	Loadings		Squared Loadings $\Sigma$		
	$Y$	$W$	$Y$	$W$	
1	.8437	-.5369	.7118	.2882	1
2	-.5369	.8437	.2882	.7118	1

# PLOT THEM!



# LOADINGS AS “SLICES OF INERTIA”

$Y$ : Astringent.  $W$ : Sugar

*Sum of row: Component Inertia (eigenvalue)*

Component	Loadings		Squared Loadings		$\Sigma$
	$Y$	$W$	$Y$	$W$	
1	16.70	-10.63	279	113	392
2	3.87	6.08	15	37	52
	294	150			444

*Sum of columns: Variable Inertia*

Exposition: `resPCA$ExPosition.Data$fj`

## A NEW WINE?

# WHAT TO DO WITH A NEW WINE? (TRIUS, A WINE FROM ONTARIO)



RATED:

SUGAR: 3/20

ASTRINGENT: 12/20

CENTERED RATING:

SUGAR:  $3 - 6 = -3$

ASTRINGENT:  $12 - 8 = 4$

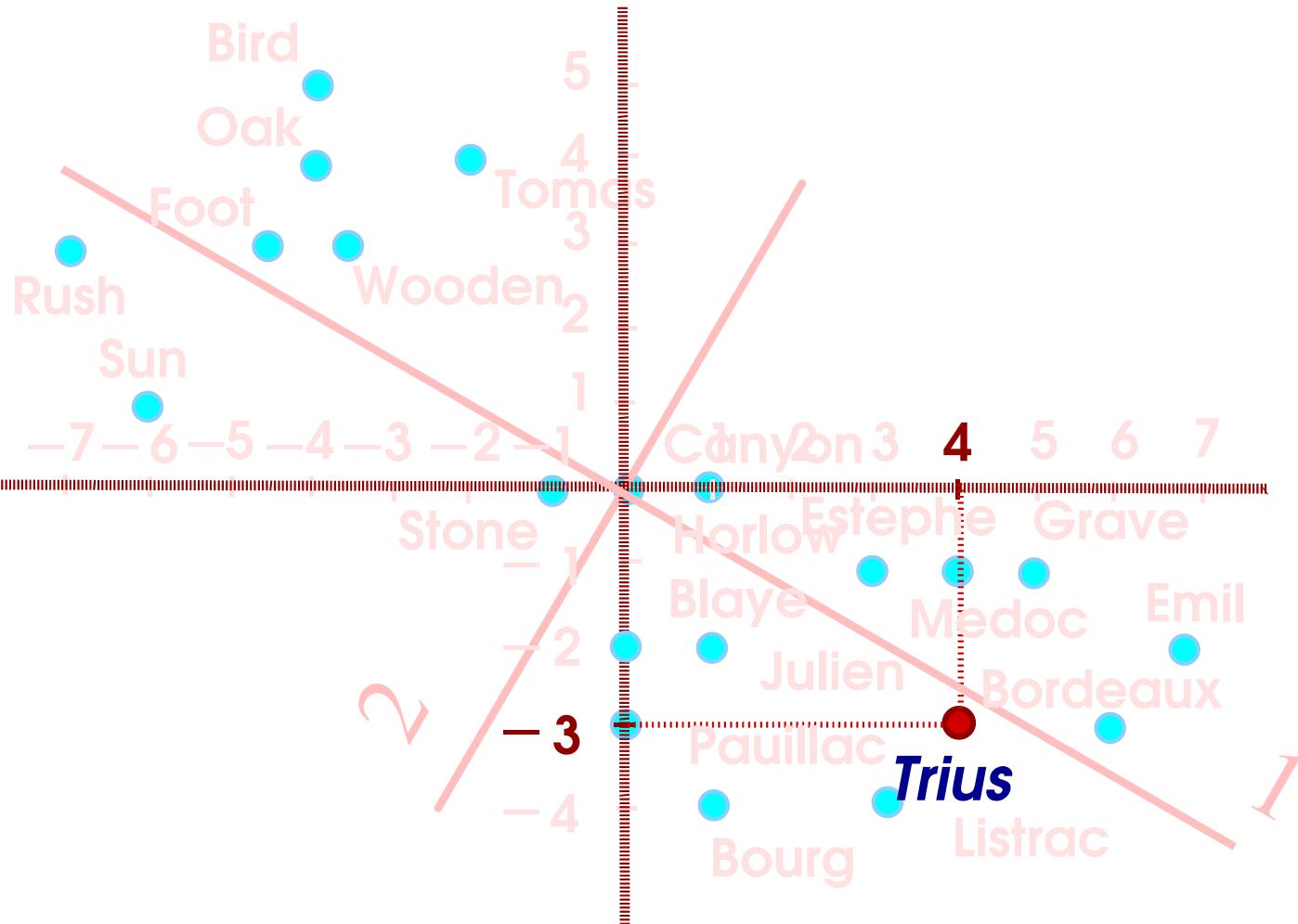
**THIS IS CALLED:**

**A SUPPLEMENTARY OBSERVATION**

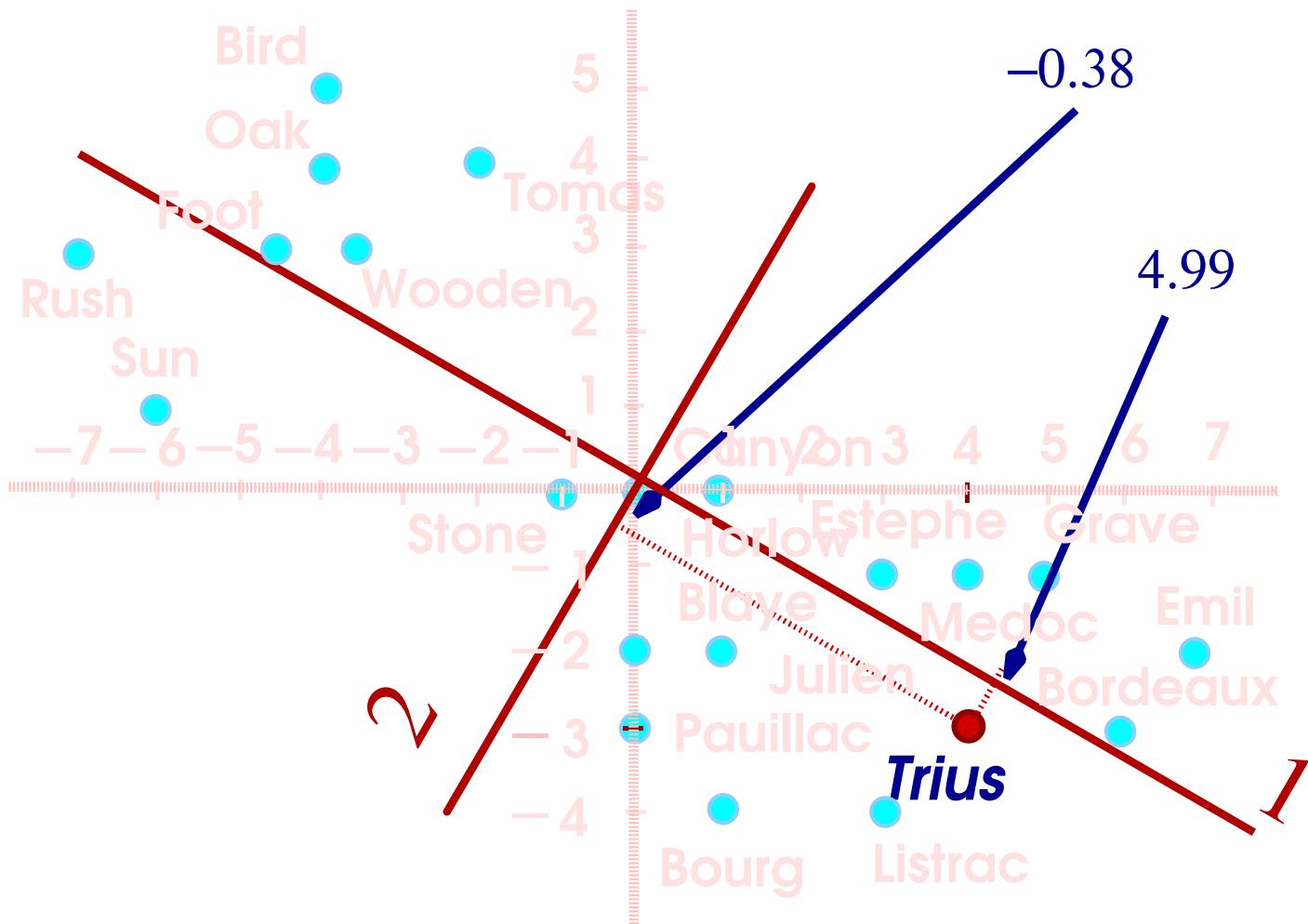
**AN ILLUSTRATIVE OBSERVATION**

**AN “OUT OF SAMPLE” OBSERVATION**

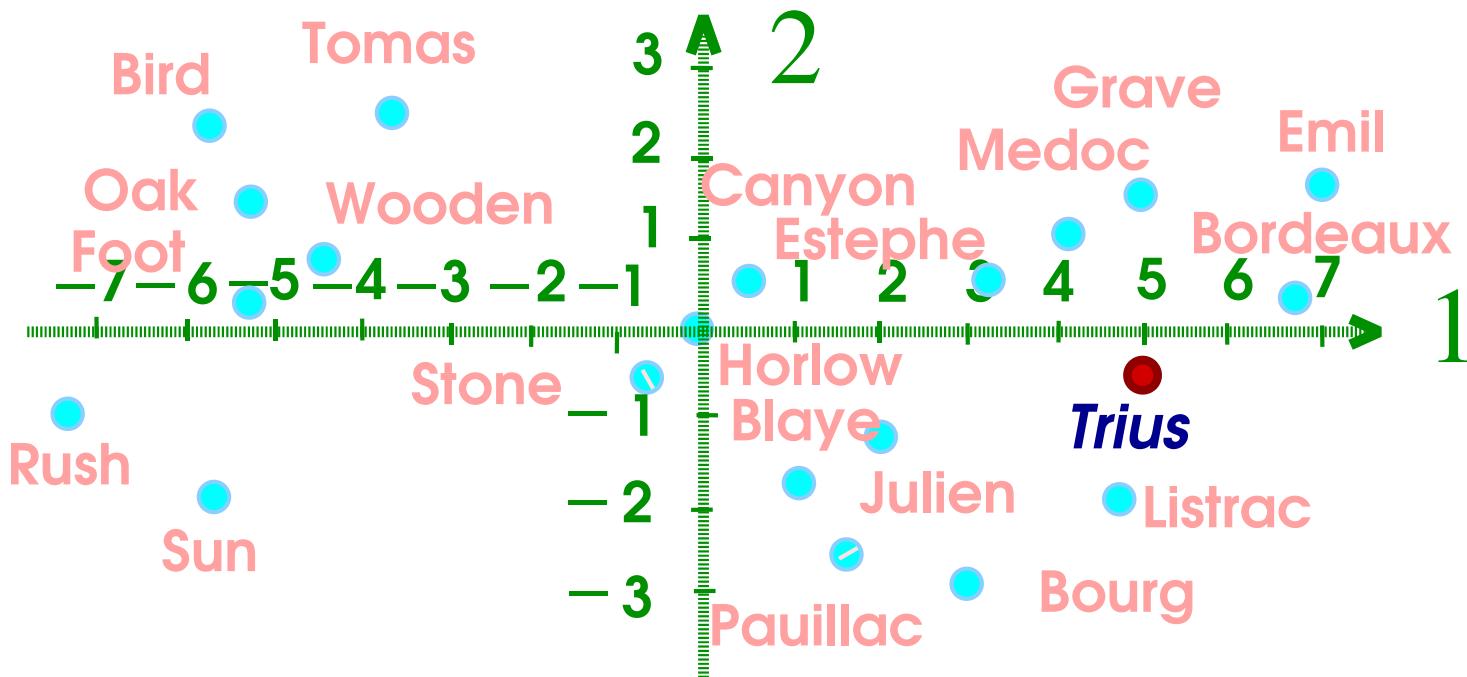
# JUST PLOT IT IN THE SPACE



# BACK-PROJECT ON COMPONENTS



## “HORIZONTALIZE” COMPONENT 1



LOOKS LIKE TRIUS IS CABERNET + MERLOT  
(AND IT IS!)

# WHAT ABOUT NEW VARIABLES?

## 20 WINES: ALL DATA



Here are the data  
In case you had forgotten.

(Long) Name	1. <a href="#">Bordeaux</a>	2. <a href="#">Black Stone</a>	3. <a href="#">Listrac</a>	4. <a href="#">Canyon Creek</a>	5. <a href="#">Côtes de Bourg</a>	6. <a href="#">Foot Hill</a>	7. <a href="#">Horizon</a>	8. <a href="#">St Esprithe</a>	9. <a href="#">Wooden Hill</a>	10. <a href="#">Blaye</a>	11. <a href="#">Côtes de Blaye</a>	11. <a href="#">Sun Set</a>	12. <a href="#">Black Bird</a>	13. <a href="#">Medoc</a>	14. <a href="#">St Julien</a>	15. <a href="#">Pauillac</a>	16. <a href="#">Rush</a>	17. <a href="#">Oak</a>	18. <a href="#">Grave</a>	19. <a href="#">Emil</a>	20. <a href="#">Tomas</a>	
Short Name	<a href="#">L</a>	<a href="#">S</a>	<a href="#">U</a>	<a href="#">F</a>	<a href="#">U</a>	<a href="#">U</a>	<a href="#">F</a>	<a href="#">U</a>	<a href="#">W</a>	<a href="#">B</a>	<a href="#">W</a>	<a href="#">S</a>	<a href="#">B</a>	<a href="#">M</a>	<a href="#">J</a>	<a href="#">P</a>	<a href="#">R</a>	<a href="#">O</a>	<a href="#">G</a>	<a href="#">E</a>	<a href="#">T</a>	
R: Sugar	3	6	2	6	2	9	6	5	9	4	7	7	11	5	4	8	2	4	12	5	10	
W: Astringency	14	7	11	9	9	4	8	11	5	8	2	4	11	2	1	12	9	8	1	13	15	6
Origin	F	U	F	U	F	U	U	F	U	F	U	U	F	F	F	P	U	U	F	F	U	
Fruity		F		F					F	F				F	F	F	F					
Woody																						
Acid	2	5	1	1	9	1	2	2	1	1	2	1	2	1	2	9	14	2	1	2	1	1
Bitter	8	3	16	3	11	1	1	9	1	8	3	2	9	12	10	2	1	10	7	3		

**WHAT ABOUT NEW VARIABLES?**

**THIS IS CALLED:**

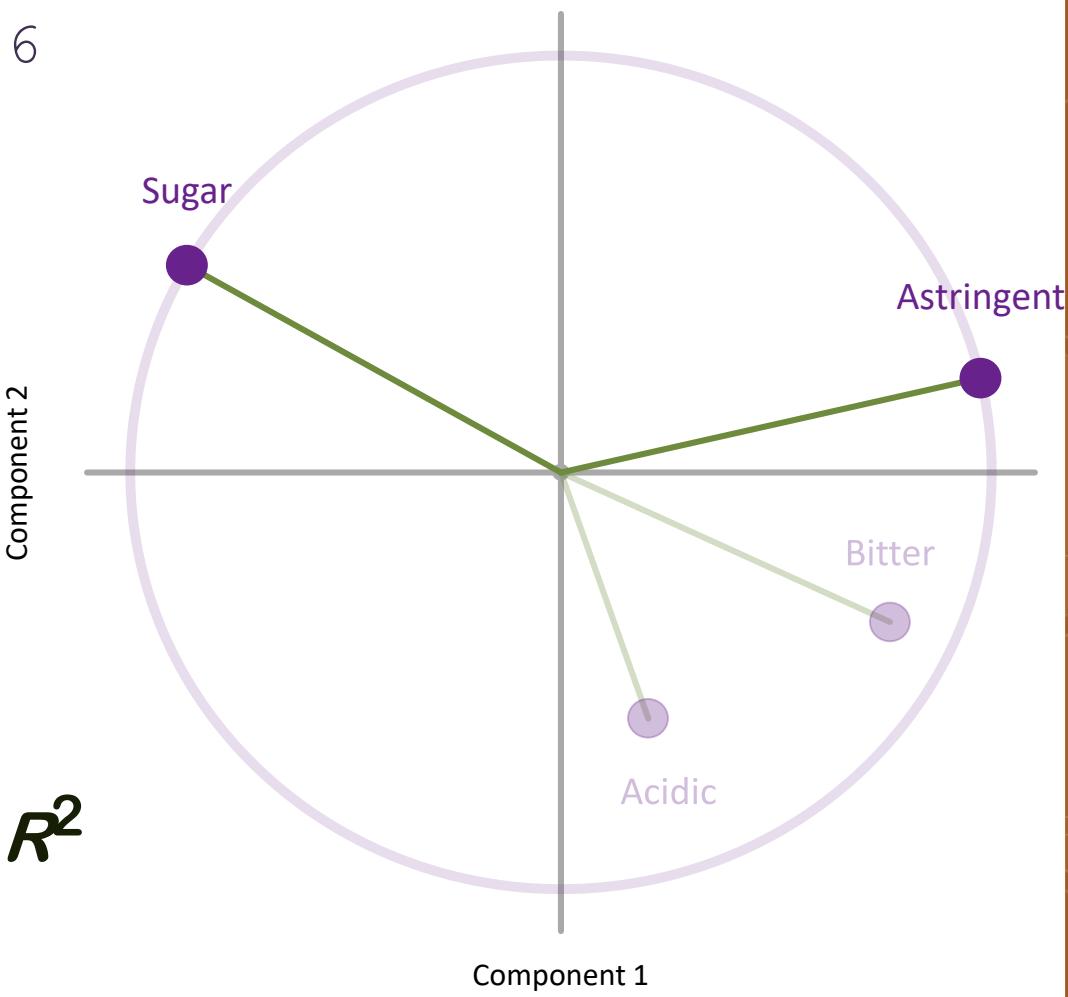
**A SUPPLEMENTARY VARIABLE**

**AN ILLUSTRATIVE VARIABLE**

**How: COMPUTE CORRELATION → LOADINGS**

**PLOT THEM**

	PC1	PC2
<i>Acidic</i>	.20	-.59
<i>Bitter</i>	.76	-.36



**CORRELATION:**  
***SUPPLEMENTARY R<sup>2</sup>***  
***DO NOT SUM TO 1***

## MID-TALK WRAP-UP

**MAPS ARE BEST!**

**NEARBY OBSERVATIONS ARE ALIKE**



**COSINES ARE FOR VARIABLES:  $[-1, 1]$**

**FACTOR SCORES ARE BEST MIXTURE**

**GOOD FACTORS EXPLAIN A LOT**

**“MERCI” TO THE EIGEN-FAIRY**

# WHAT ARE THE IMPORTANT COMPONENTS

(DÉJÀ VU) HOW GOOD?

# GOOD COMPONENTS EXPLAIN A LOT

Component	$\lambda_i$ (eigenvalue)	Cumulated (eigenvalues)	Percent of of inertia	Cumulated (percentage)
1	392	392	83.29	83.29
2	52	444	11.71	100.00

## COMPARED TO

Variable	$SS_i$	Cumulative $SS_i$	Percent of of inertia	Cumulative (percentage)
Astringent	294	294	66.22	66.22
Sugar	150	444	33.78	100.00

EIGEN-FAIRY: YOU CANNOT BETTER PCA!

# WHAT ARE THE IMPORTANT COMPONENTS

**RULE OF THUMB:**

**A GOOD COMPONENT EXPLAINS MORE  
THAN THE AVERAGE INERTIA.**

# IMPORTANT OBSERVATIONS

FIRST A REFRESHER

## FACTOR SCORES

# EIGENVALUES

# SUM OF SQUARED PROJECTIONS

	<i>F</i>	<i>W</i>	<i>y</i>	<i>w</i>	<i>F<sub>1</sub></i>	<i>F<sub>2</sub></i>	<i>F<sub>1</sub><sup>2</sup></i>	<i>F<sub>2</sub><sup>2</sup></i>
Bordeaux	3	14	-3	6	6.67	0.69	44.52	0.48
Black Stone	6	7	0	-1	-0.84	-0.54	0.71	0.29
Listrac	2	11	-4	3	4.68	-1.76	21.89	3.11
Canyon Creek	6	9	0	1	0.84	0.54	0.71	0.29
Côtes de Bourg	2	9	-4	1	2.99	-2.84	8.95	8.05
Foot Hill	9	4	3	-4	-4.99	0.38	24.85	0.15
Horlow	6	8	0	0	0.00	0.00	0	0.00
St. Estphe	5	11	-1	3	3.07	0.77	9.41	0.59
Wooden Hill	9	5	3	-3	-4.14	0.92	17.15	0.85
Côtes de Blaye	4	8	-2	0	1.07	-1.69	1.15	2.85
Sun Set	7	2	1	-6	-5.60	-2.38	31.35	5.65
Black Bird	11	4	5	-4	-6.06	2.07	36.71	4.29
Médoc	5	12	-1	4	3.91	1.30	15.30	1.70
St Julien	4	9	-2	1	1.92	-1.15	3.68	1.32
Pauillac	3	8	-3	0	1.61	-2.53	2.59	6.41
Gold Rush	9	1	3	-7	-7.52	-1.23	56.49	1.51
Oak Ville	10	4	4	-4	-5.52	1.23	30.49	1.51
Grave	5	13	-1	5	4.76	1.84	22.61	3.39
St Emilion	4	15	-2	7	6.98	2.07	48.71	4.29
Tomasello	10	6	4	-2	-3.83	2.30	14.71	5.29
<b><math>\Sigma</math></b>	<b>120</b>	<b>160</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>392</b>	<b>52</b>
							$\lambda_1$	$\lambda_2$

## CONTRIBUTION

# IMPORTANT OBSERVATIONS:

**Rule of thumb:**

**A Good observation “explains”  
more than the average Inertia.**

**How to find the “explained inertia?”**

# FIRST MAKE IT A PROPORTION

$F_1^2$	$F_2^2$
44.52	0.48
0.71	0.29
21.89	3.11
0.71	0.29
8.95	8.05
24.85	0.15
0.00	0.00
9.41	0.59
17.15	0.85
1.15	2.85
31.35	5.65
36.71	4.29
15.30	1.70
3.68	1.32
2.59	6.41
56.49	1.51
30.49	1.51
22.61	3.39
48.71	4.29
14.71	5.29
<hr/>	
$\lambda$	392 52
<hr/>	

$$44.52 / 392 = .114$$

## IMPORTANT VARIABLES

**CALL THIS RATIO: THE CONTRIBUTION**

**IMPORTANT OBSERVATIONS:  
CONTRIBUTION >**

## FACTOR SCORES

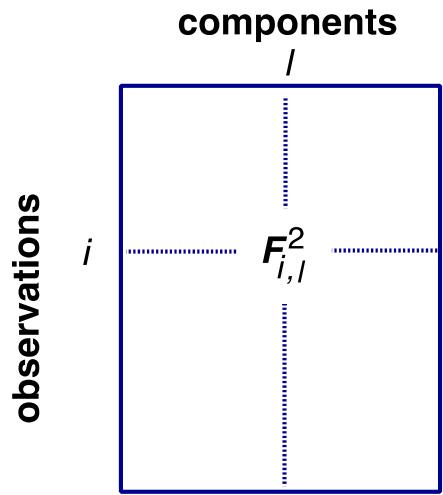
# EIGENVALUES

## SUM OF SQUARED PROJECTIONS

$$(1/20) = .05$$

	<i>Y</i>	<i>W</i>	<i>y</i>	<i>w</i>	<i>F<sub>1</sub></i>	<i>F<sub>2</sub></i>	<i>ctr<sub>1</sub></i>	<i>ctr<sub>2</sub></i>	<i>F<sub>1</sub><sup>2</sup></i>	<i>F<sub>2</sub><sup>2</sup></i>	<i>d<sup>2</sup></i>
							<i>x</i>	<i>x</i>	100	100	
Bordeaux	3	14	-3	6	6.67	0.69	11	1	44.52	0.48	45
Black Stone	6	7	0	-1	-0.84	-0.54	0	1	0.71	0.29	1
Listrac	2	11	-4	3	4.68	-1.76	6	6	21.89	3.11	25
Canyon Creek	6	9	0	1	0.84	0.54	0	1	0.71	0.29	1
Côtes de Bourg	2	9	-4	1	2.99	-2.84	2	15	8.95	8.05	17
Foot Hill	9	4	3	-4	-4.99	0.38	6	0	24.85	0.15	25
Horlow	6	8	0	0	0.00	0.00	0	0	0.00	0.00	0
St. Estphe	5	11	-1	3	3.07	0.77	3	1	9.41	0.59	10
Wooden Hill	9	5	3	-3	-4.14	0.92	5	2	17.15	0.85	18
Côtes de Blaye	4	8	-2	0	1.07	-1.69	0	5	1.15	2.85	4
Sun Set	7	2	1	-6	-5.60	-2.38	8	11	31.35	5.65	37
Black Bird	11	4	5	-4	-6.06	2.07	9	8	36.71	4.29	41
Médoc	5	12	-1	4	3.91	1.30	4	3	15.30	1.70	17
St. Julien	4	9	-2	1	1.92	-1.15	1	3	3.68	1.32	5
Pauillac	3	8	-3	0	1.61	-2.53	1	12	2.59	6.41	9
Gold Rush	9	1	3	-7	-7.52	-1.23	14	3	56.49	1.51	58
Oak Ville	10	4	4	-4	-5.52	1.23	8	3	30.49	1.51	32
Grave	5	13	-1	5	4.76	1.84	6	7	22.61	3.39	26
St Emilion	4	15	-2	7	6.98	2.07	12	8	48.71	4.29	53
Tomasello	10	6	4	-2	-3.83	2.30	4	10	14.71	5.29	20
$\Sigma$		120	160	0	0	0	100	100	392	52	444
									$\lambda_1$	$\lambda_2$	$\delta$

# CONTRIBUTIONS IN A NUTSHELL



## SQUARED COSINE

**CONTRIBUTIONS: IMPORTANT  
OBSERVATIONS FOR A COMPONENT**

**WHAT ARE THE IMPORTANT COMPONENTS  
FOR AN OBSERVATION?  
GOOD COMPONENTS HAVE LARGE COSINES<sup>2</sup>**

## FACTOR SCORES

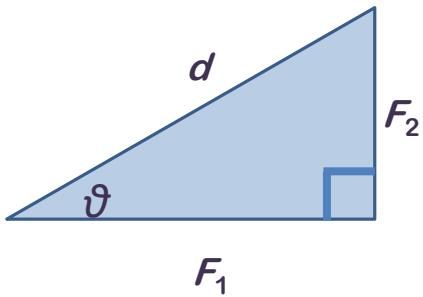
# DISTANCE<sup>2</sup>

**SUM OF SQUARED PROJECTIONS ( $F_1^2 + F_2^2 = D^2$ )**

$$\frac{44.52}{45} = .99$$

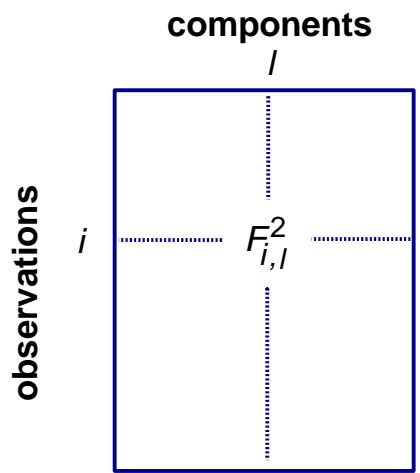
This is a cosine<sup>2</sup>

$$\cos^2 \vartheta = \frac{F_1^2}{d^2}$$



	<i>F</i>	<i>W</i>	<i>y</i>	<i>w</i>	<i>F<sub>1</sub></i>	<i>F<sub>2</sub></i>	<i>F<sub>1</sub><sup>2</sup></i>	<i>F<sub>2</sub><sup>2</sup></i>	<i>d<sup>2</sup></i>	<i>cos<sup>2</sup> θ</i>	<i>×</i>	<i>100</i>	<i>100</i>
Bordeaux	3	14	-3	6	6.67	0.69	44.52	0.48	45	99	1		
Black Stone	6	7	0	-1	-0.84	-0.54	0.71	0.29	1	71	29		
Listrac	2	11	-4	3	4.68	-1.76	21.89	3.11	25	88	12		
Canyon Creek	6	9	0	1	0.84	0.54	0.71	0.29	1	71	29		
Côtes de Bourg	2	9	-4	1	2.99	-2.84	8.95	8.05	17	53	47		
Foot Hill	9	4	3	-4	-4.99	0.38	24.85	0.15	25	99	1		
Horlow	6	8	0	0	0.00	0.00	0	0.00	0	0	0		
St. Estphe	5	11	-1	3	3.07	0.77	9.41	0.59	10	94	6		
Wooden Hill	9	5	3	-3	-4.14	0.92	17.15	0.85	18	95	5		
Côtes de Blaye	4	8	-2	0	1.07	-1.69	1.15	2.85	4	29	71		
Sun Set	7	2	1	-6	-5.60	-2.38	31.35	5.65	37	85	15		
Black Bird	11	4	5	-4	-6.06	2.07	36.71	4.29	41	90	10		
Médoc	5	12	-1	4	3.91	1.30	15.30	1.70	17	90	10		
St Julien	4	9	-2	1	1.92	-1.15	3.68	1.32	5	74	26		
Pauillac	3	8	-3	0	1.61	-2.53	2.59	6.41	9	29	71		
Gold Rush	9	1	3	-7	-7.52	-1.23	56.49	1.51	58	97	3		
Oak Ville	10	4	4	-4	-5.52	1.23	30.49	1.51	32	95	5		
Grave	5	13	-1	5	4.76	1.84	22.61	3.39	26	87	13		
St Emilion	4	15	-2	7	6.98	2.07	48.71	4.29	53	92	8		
Tomasello	10	6	4	-2	-3.83	2.30	14.71	5.29	20	74	26		
<b><i>Σ</i></b>	<b>120</b>	<b>160</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>392</b>	<b>52</b>	<b>444</b>				
							<b><i>λ<sub>1</sub></i></b>	<b><i>λ<sub>2</sub></i></b>	<b><i>λ</i></b>				

# COSINES<sup>2</sup> IN A NUTSHELL



## VARIANTS OF PCA

AND NOW FOR SOMETHING COMPLETELY  
NOT DIFFERENTLY ALIKE ...

# THE UNIT PROBLEM: NORMALIZATION.



**CHOUX ET CAROTTES**

**APPLES AND ORANGES**



# How?

# WHEN *NOT* TO NORMALIZE

# Ax EXAMPLE: PCA WITH COVARIANCE

## The story

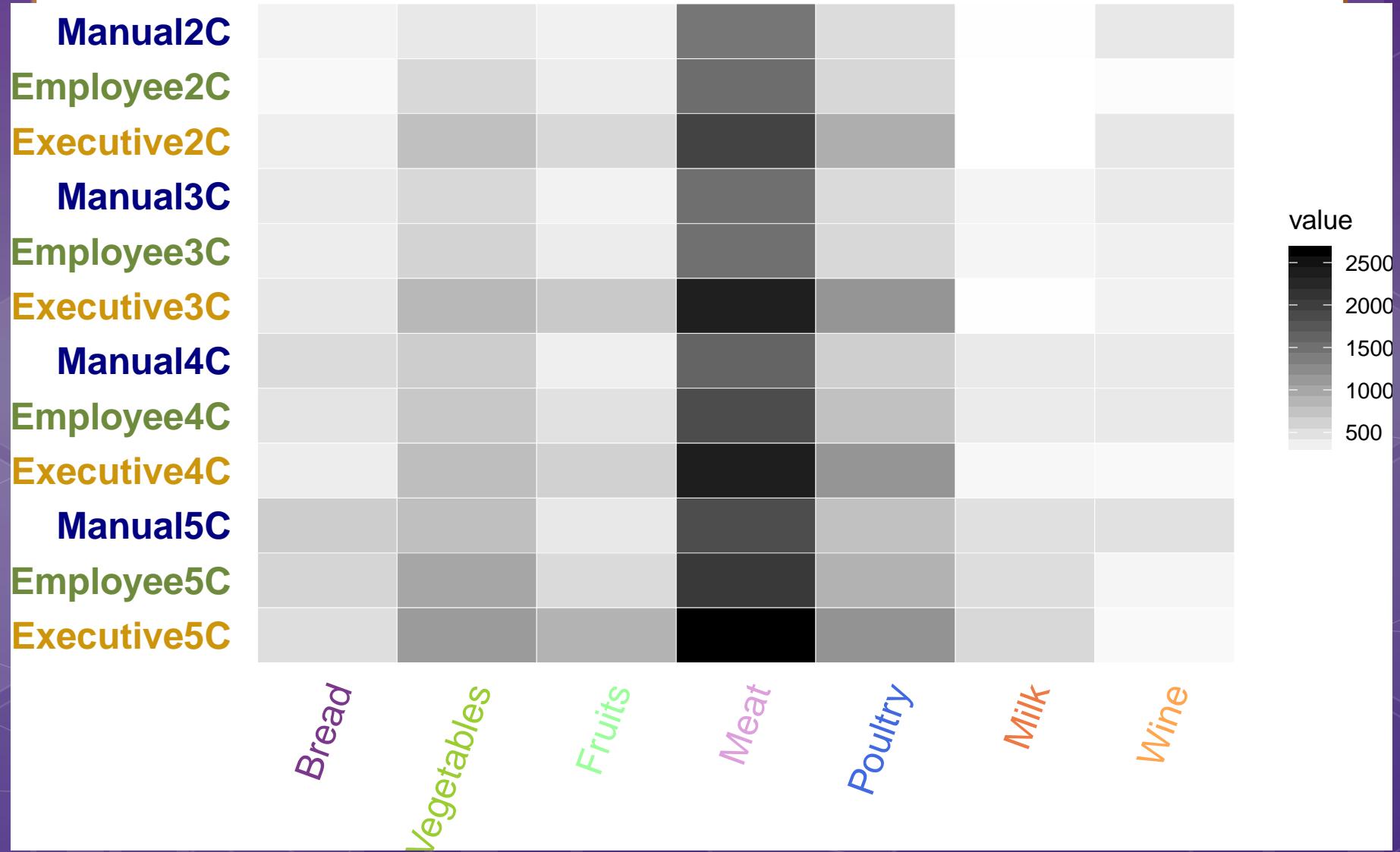
These data give the consumption in Francs of different "types of food" according to social class and number of children. The observations correspond to the average amount of money spent per month on a given type of food for a given social class and a given number of children. Because a franc spent on one item has the same value as one franc spent on another one, we want to keep the same unit of measurement for the complete space. Therefore we will analyze only the centered data, equivalently, we will analyze the covariance matrix.

## THE DATA

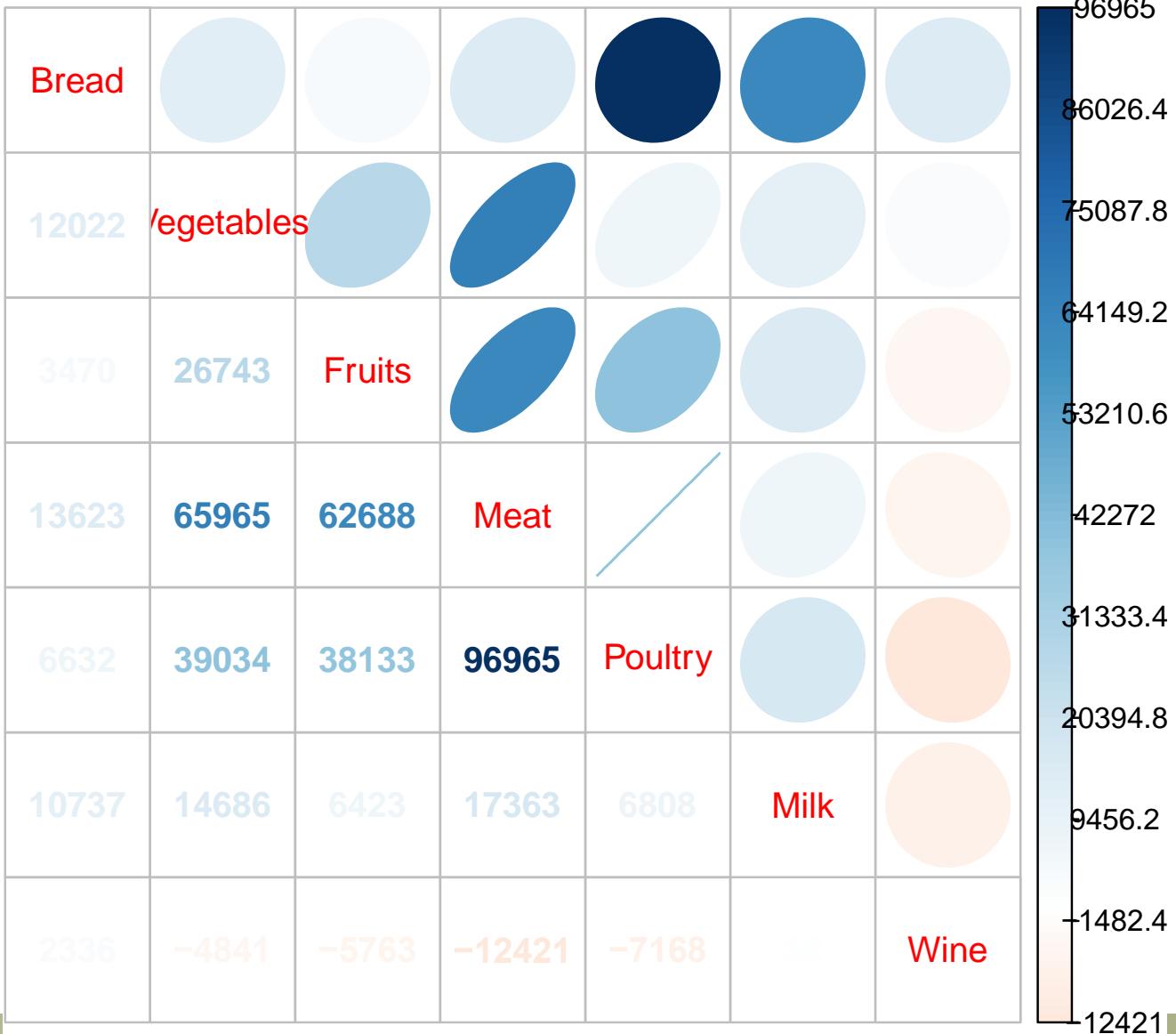
Category	Bread	Vegies	Fruits	Meat	Pltr	Milk	Vine
Blue Clr, 2C.	332	428	354	1437	526	247	427
White Clr, 2C.	293	559	388	1527	567	239	258
Upper Cls, 2C.	372	767	562	1948	927	235	433
Blue Clr, 3C.	406	563	341	1507	544	324	407
White Clr, 3C.	386	608	396	1501	558	319	363
Upper Cls, 3C.	438	843	689	2345	1148	243	341
Blue Clr, 4C.	534	660	367	1620	638	414	407
White Clr, 4C.	460	699	484	1856	762	400	416
Upper Cls, 4C.	385	789	621	2366	1149	304	282
Blue Clr, 5C.	655	776	423	1848	759	495	486
White Clr, 5C.	584	995	548	2056	893	518	319
Upper Cls, 5C.	515	1097	887	2630	1167	561	284
Mean	447	732	505	1887	803	358	369
$\hat{S}$	107	189	165	396	250	117	72

Table 13: The Data

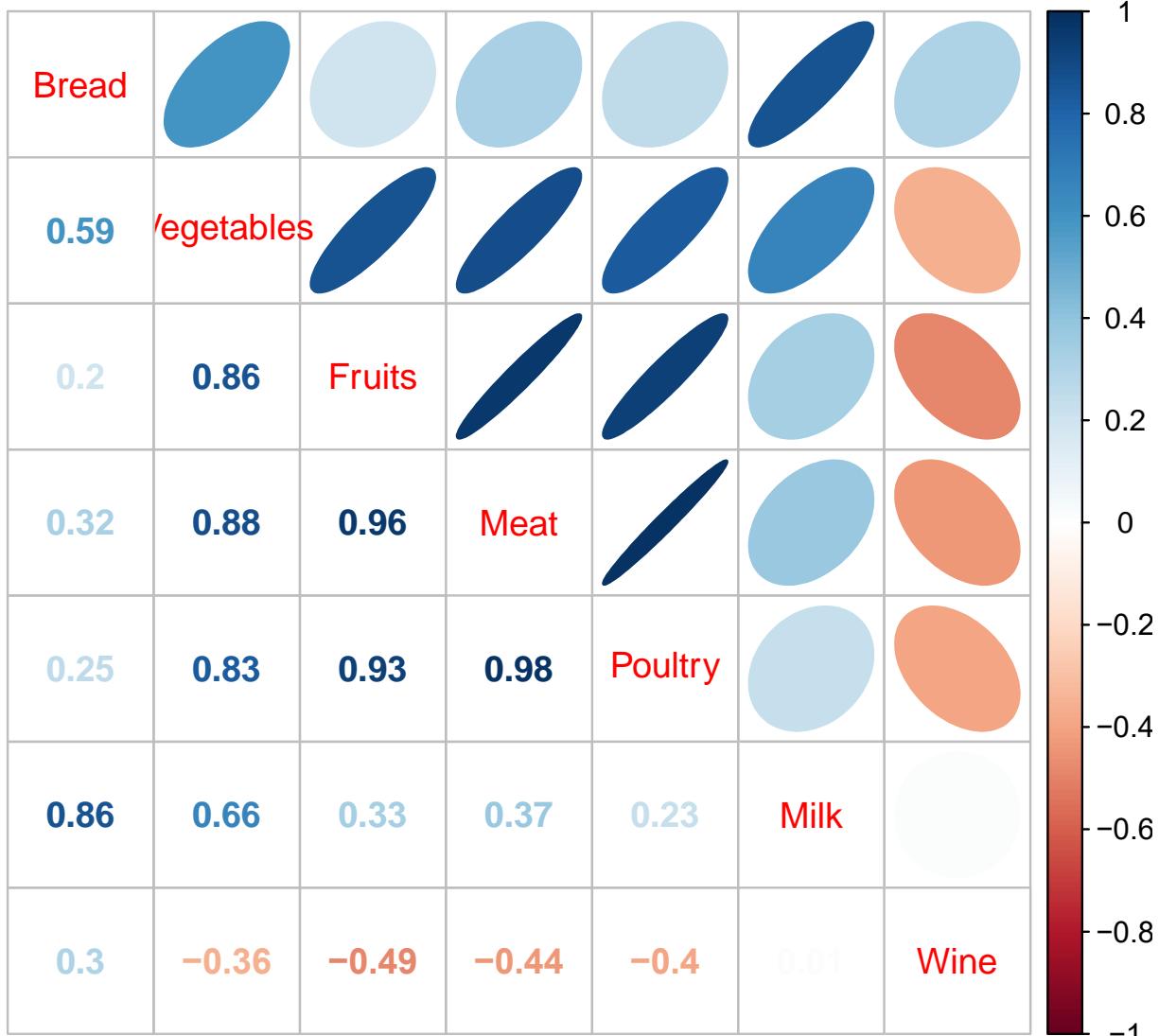
# DATA AS PICTURE



## BETTER: COVARIANCE



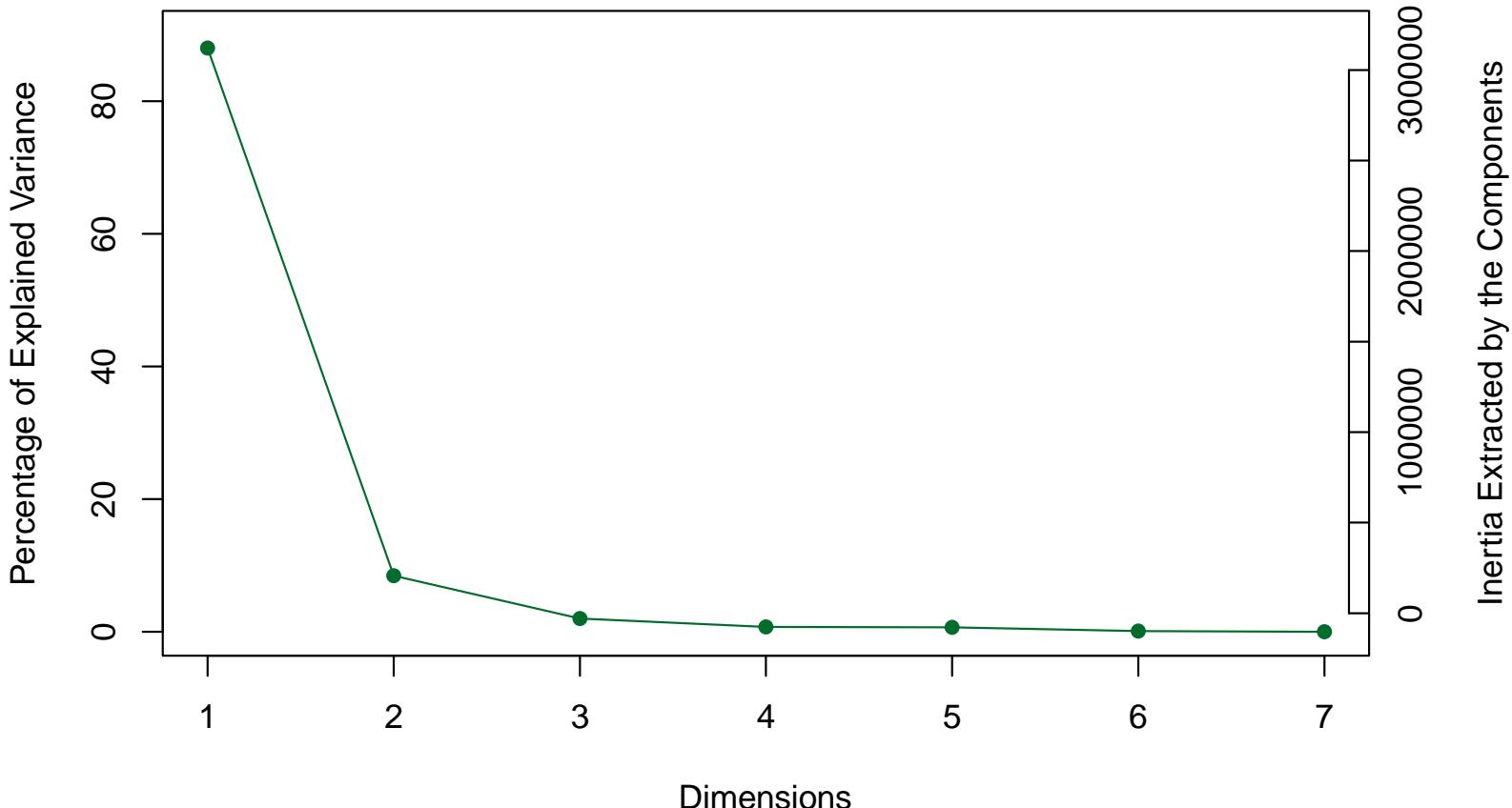
## CORRELATION



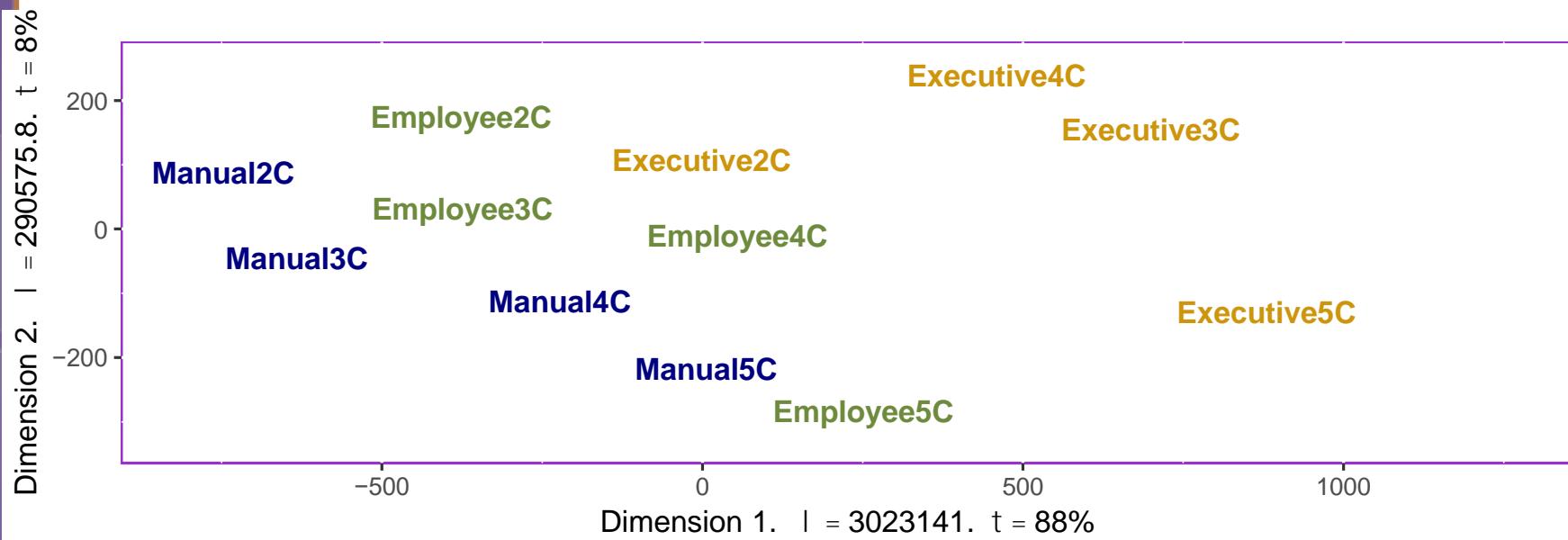
## PCA

**RESPCA = EPPCA(X, SCALE = FALSE)**

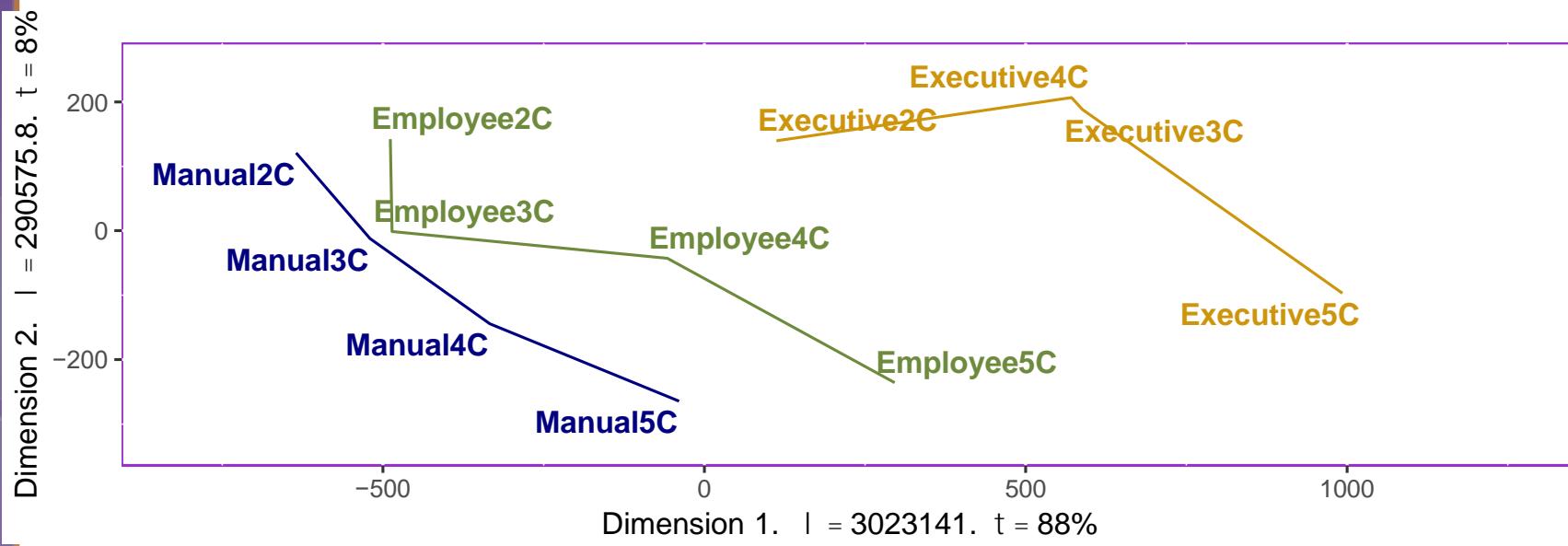
Explained Variance per Dimension



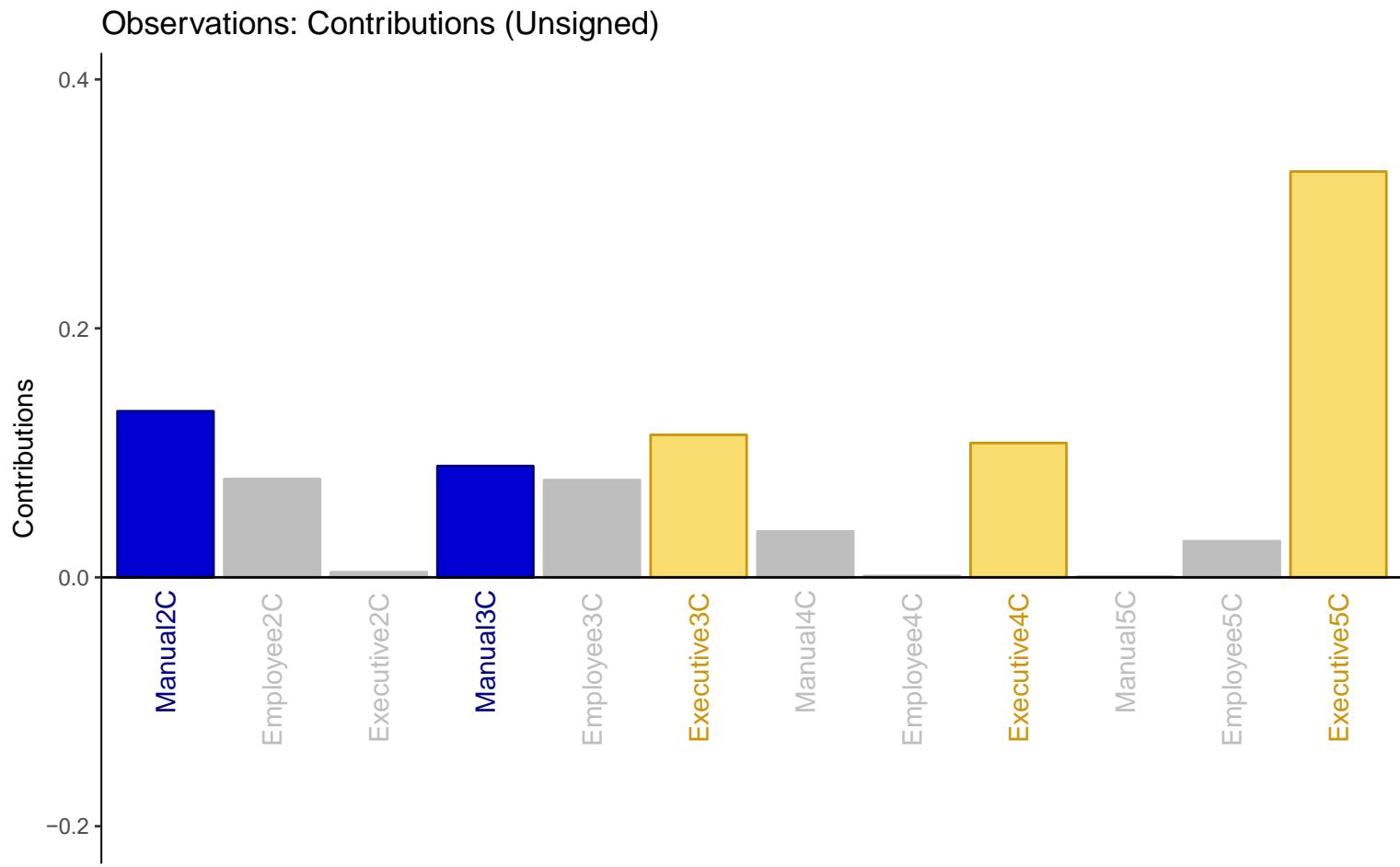
## FACTOR SCORES



## WITH LINES

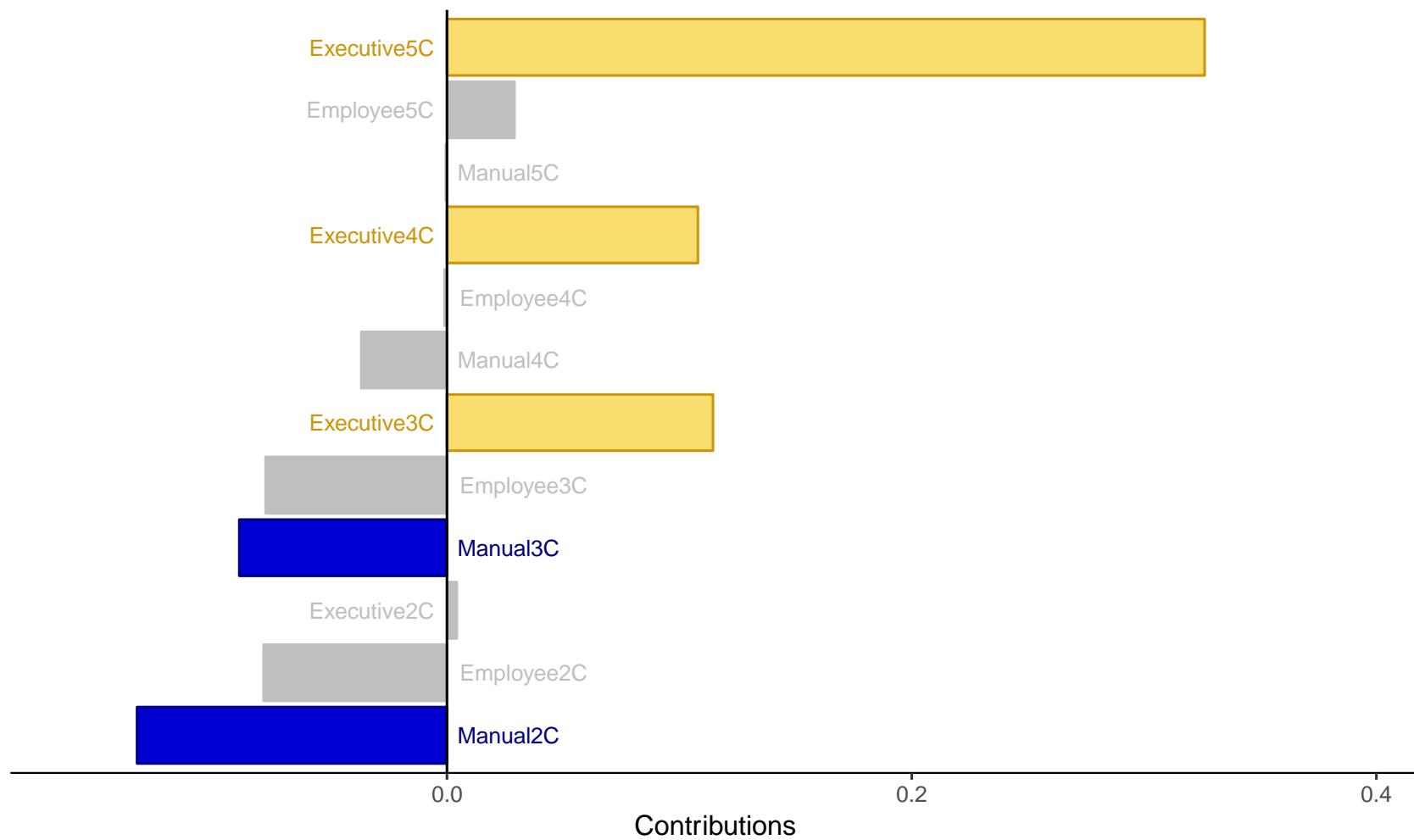


# CONTRIBUTIONS

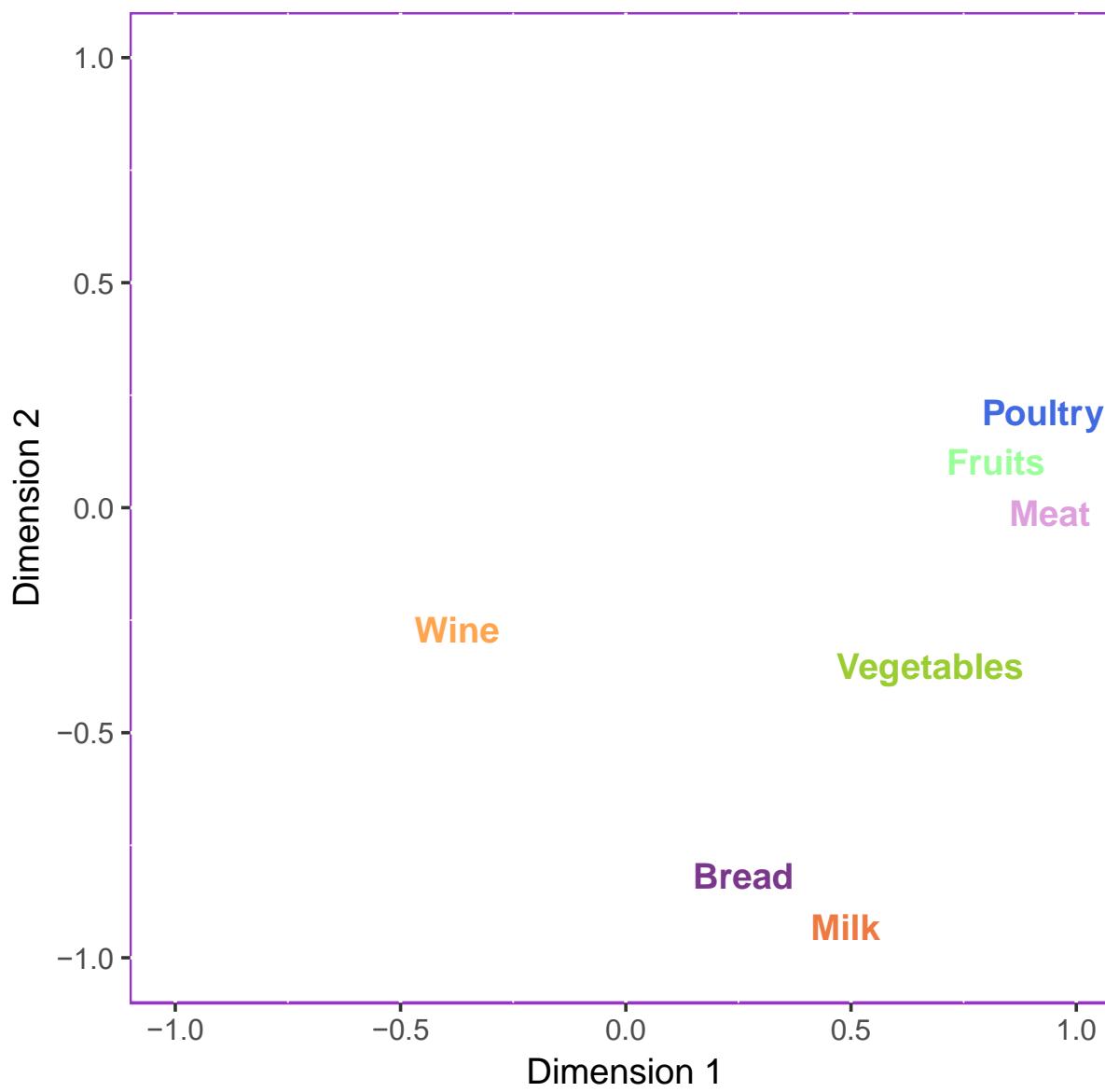


**SIGNED**

### Observations: Contributions (Signed)



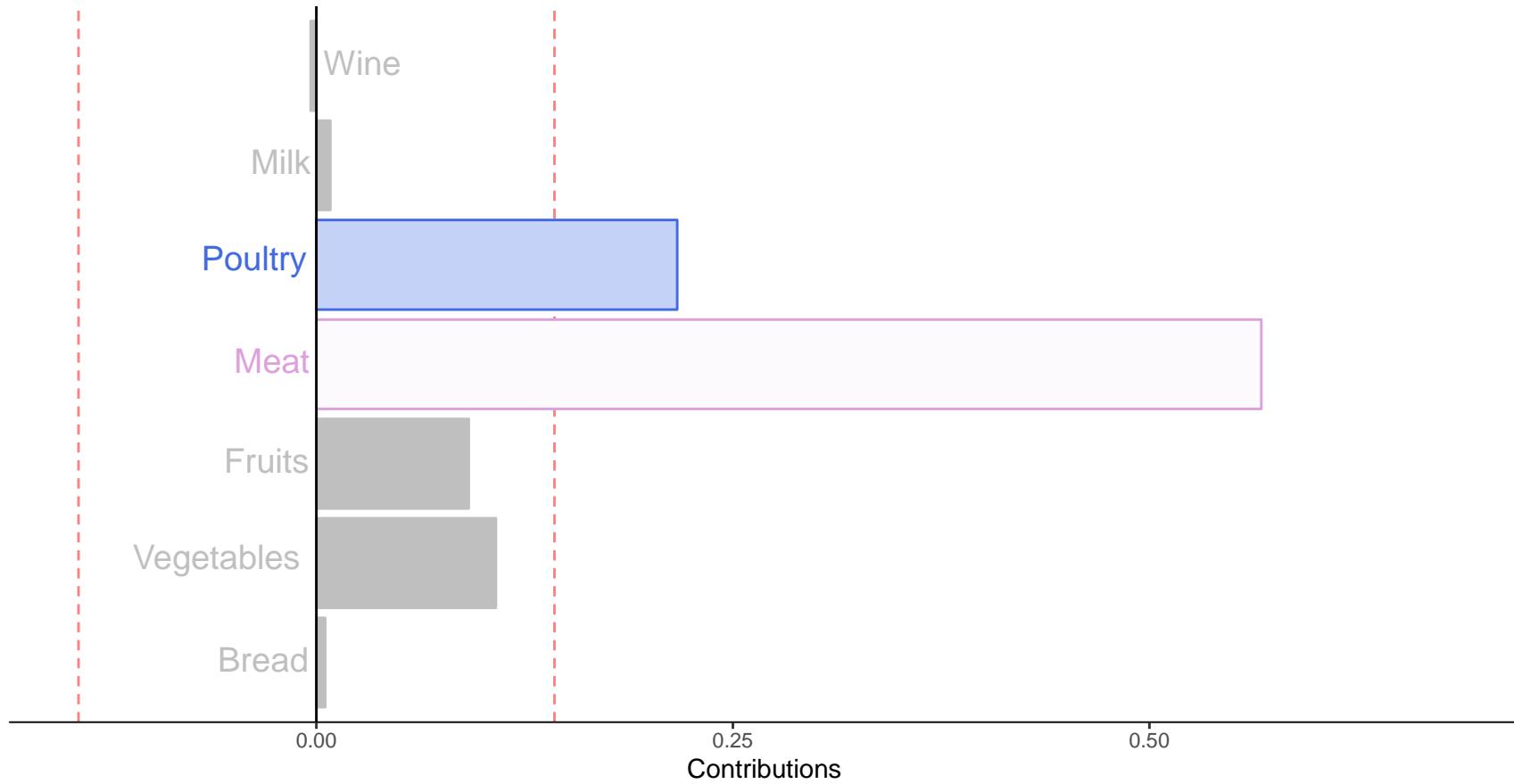
## THE J-SET



# WHAT ABOUT THE CONTRIBUTIONS

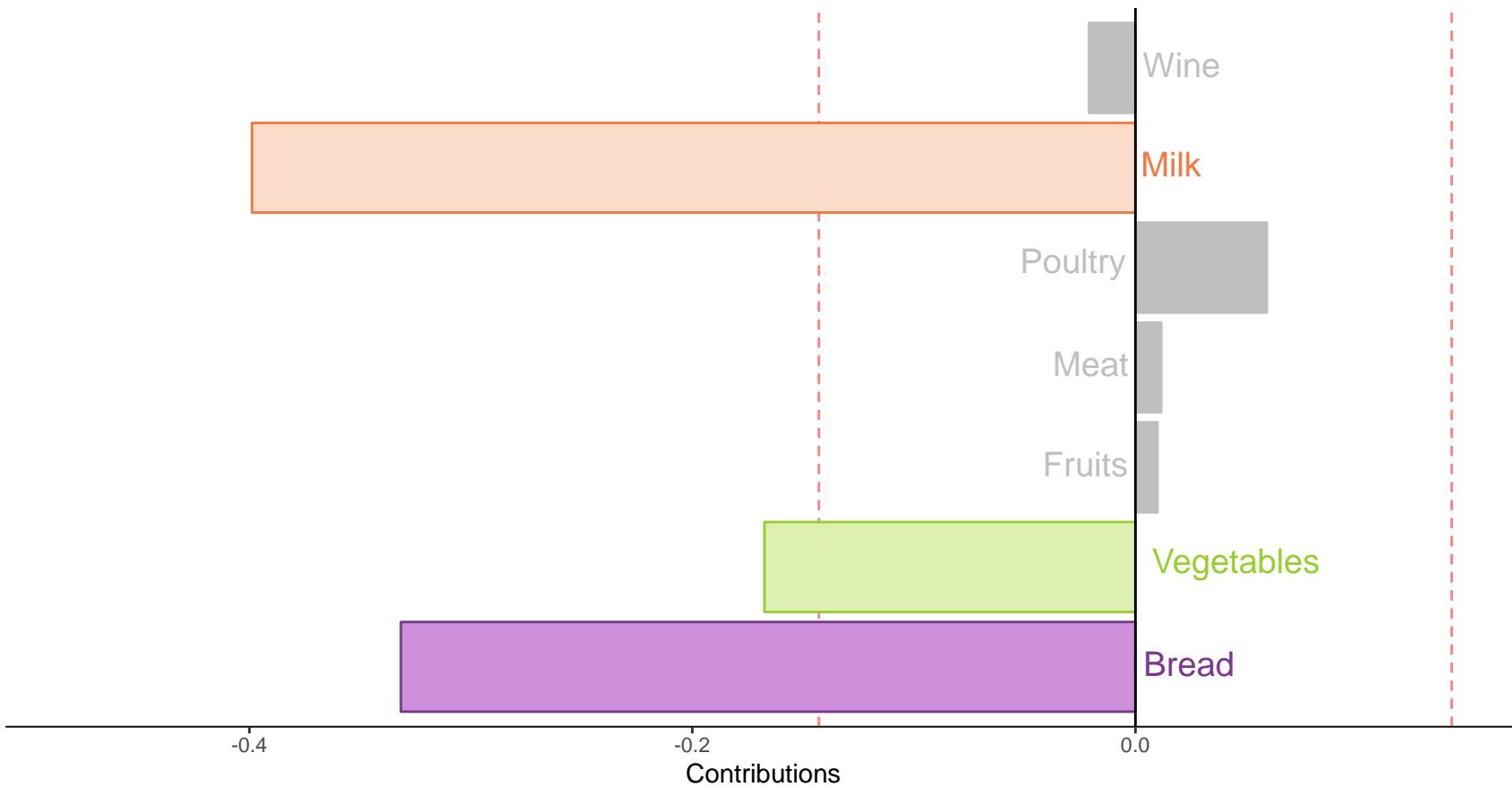
# DIMENSION 1.

Signed Contributions. Dimension 1



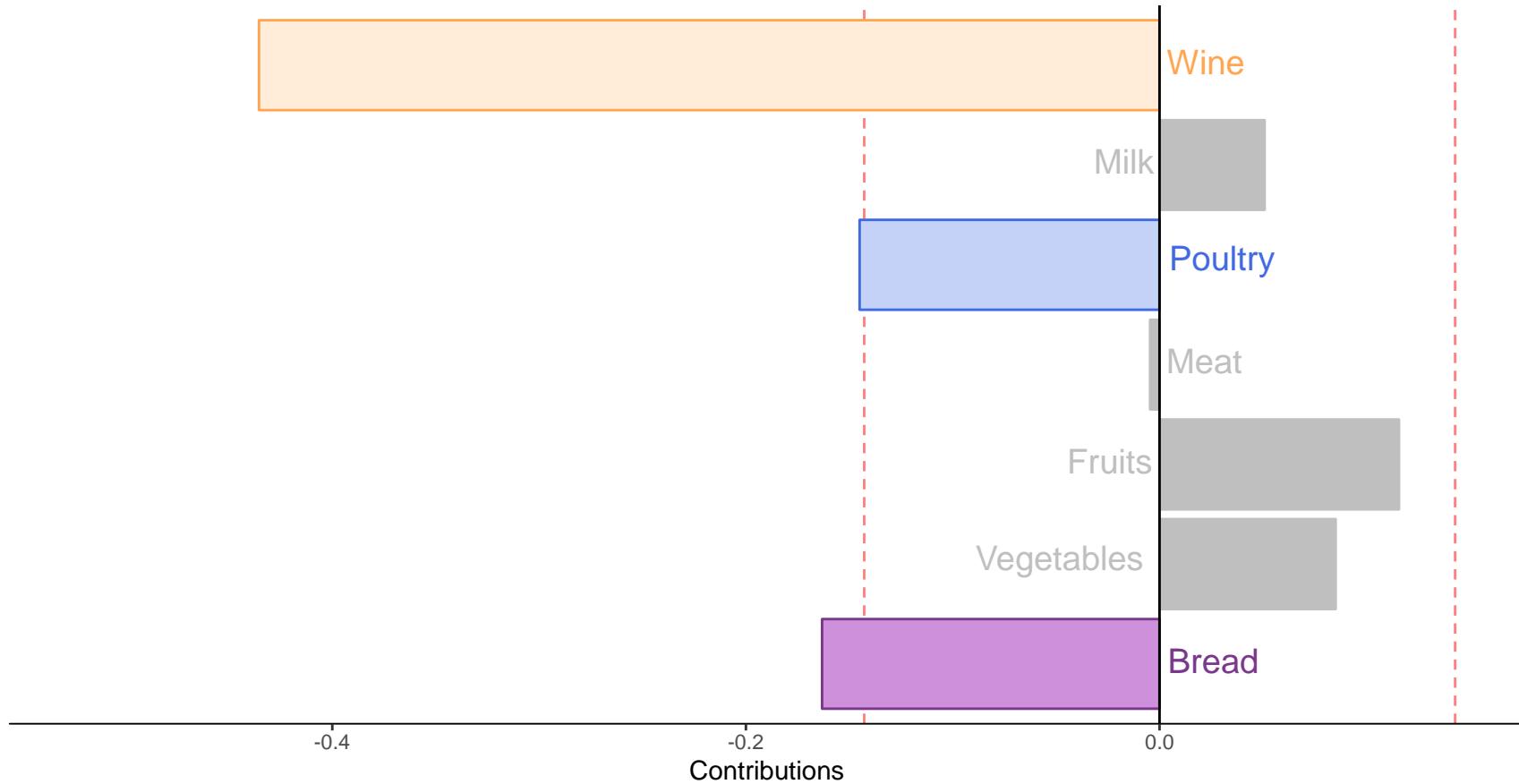
# DIMENSION 2.

Signed Contributions. Dimension 2



# DIMENSION 3

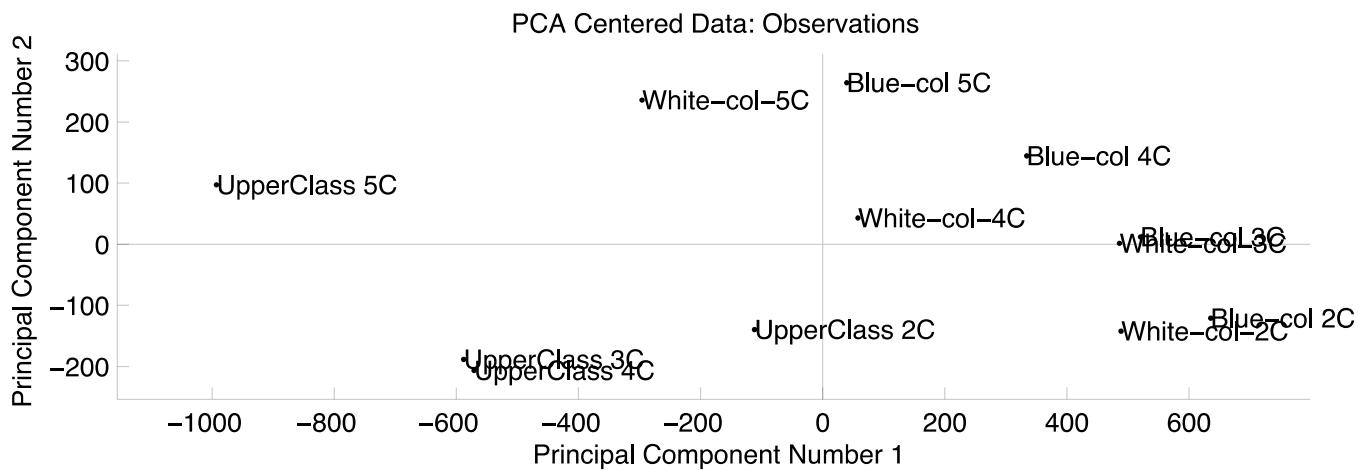
Signed Contributions. Dimension 3

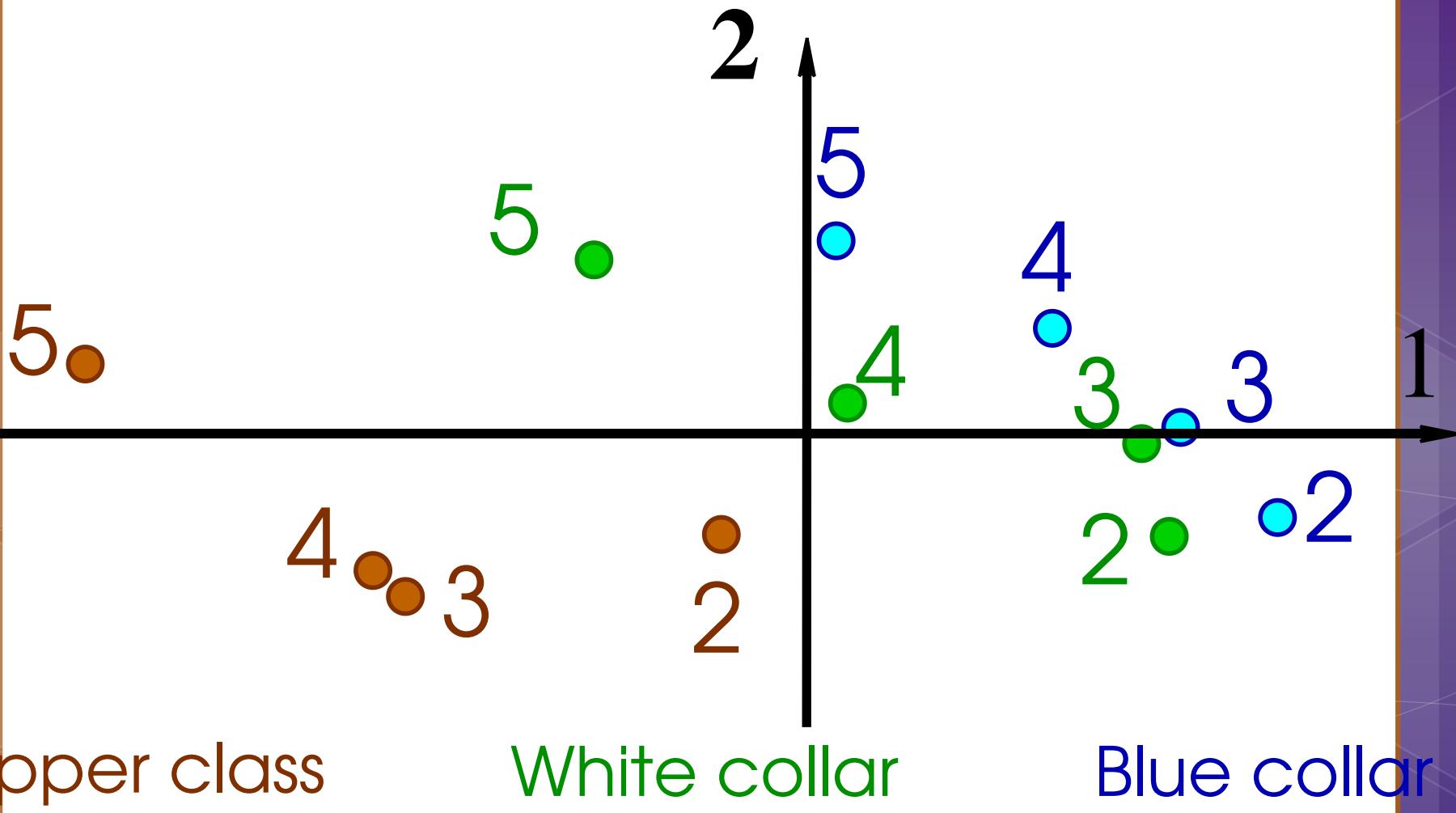


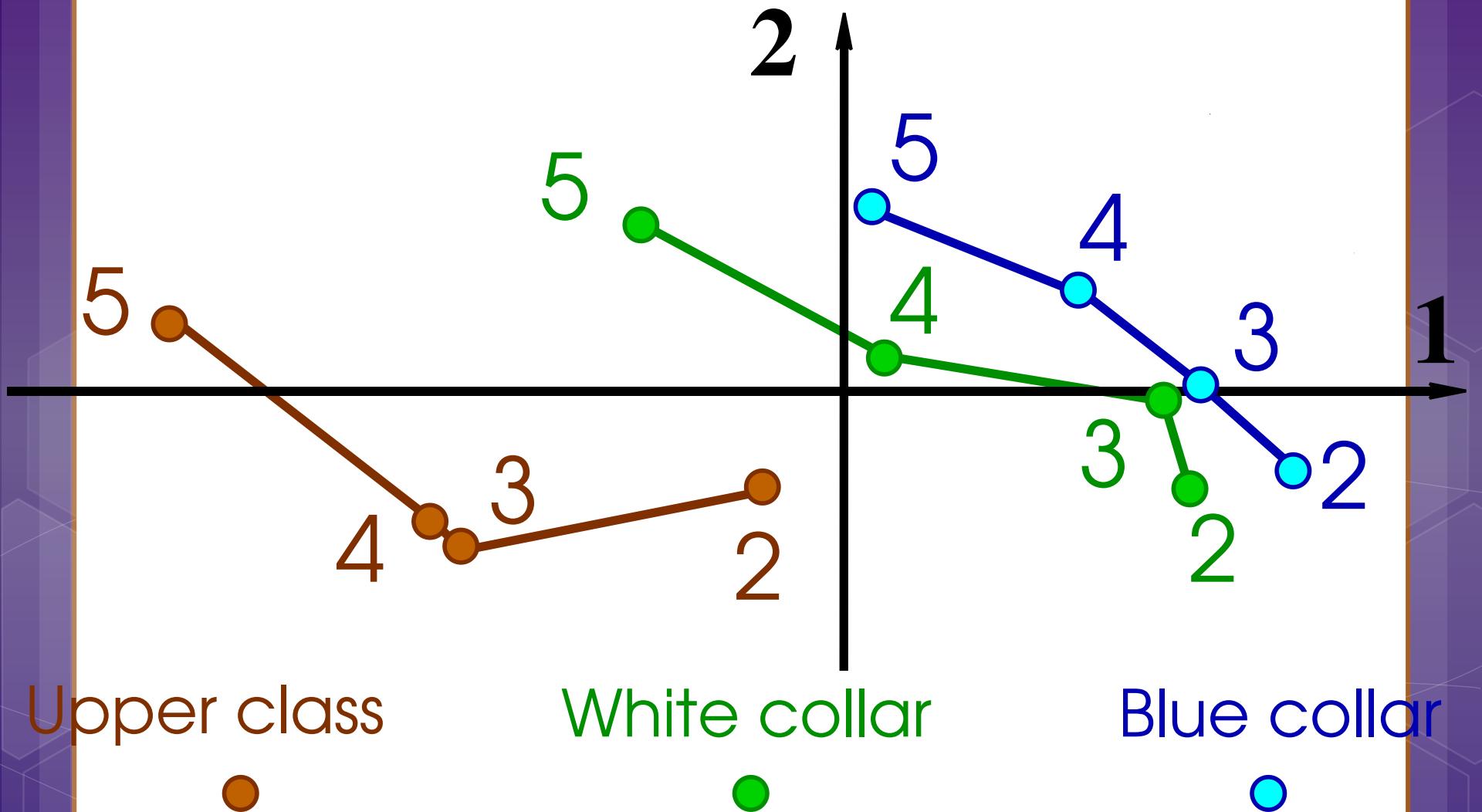
## CONCLUSIONS FOOD

# PCA. ANALYSIS WITH MATLAB +

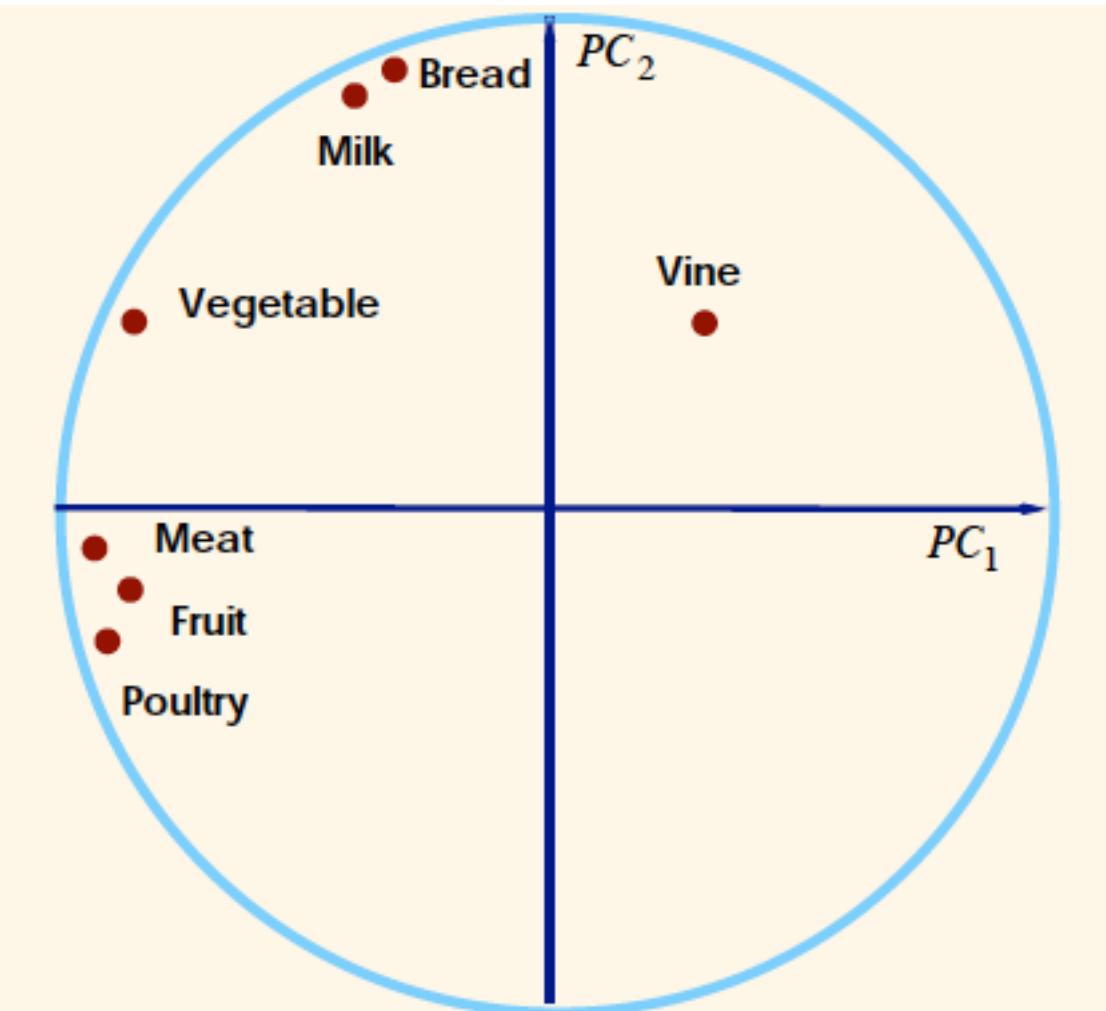
# NEVER DO THAT: THE HORRIBLE WAY





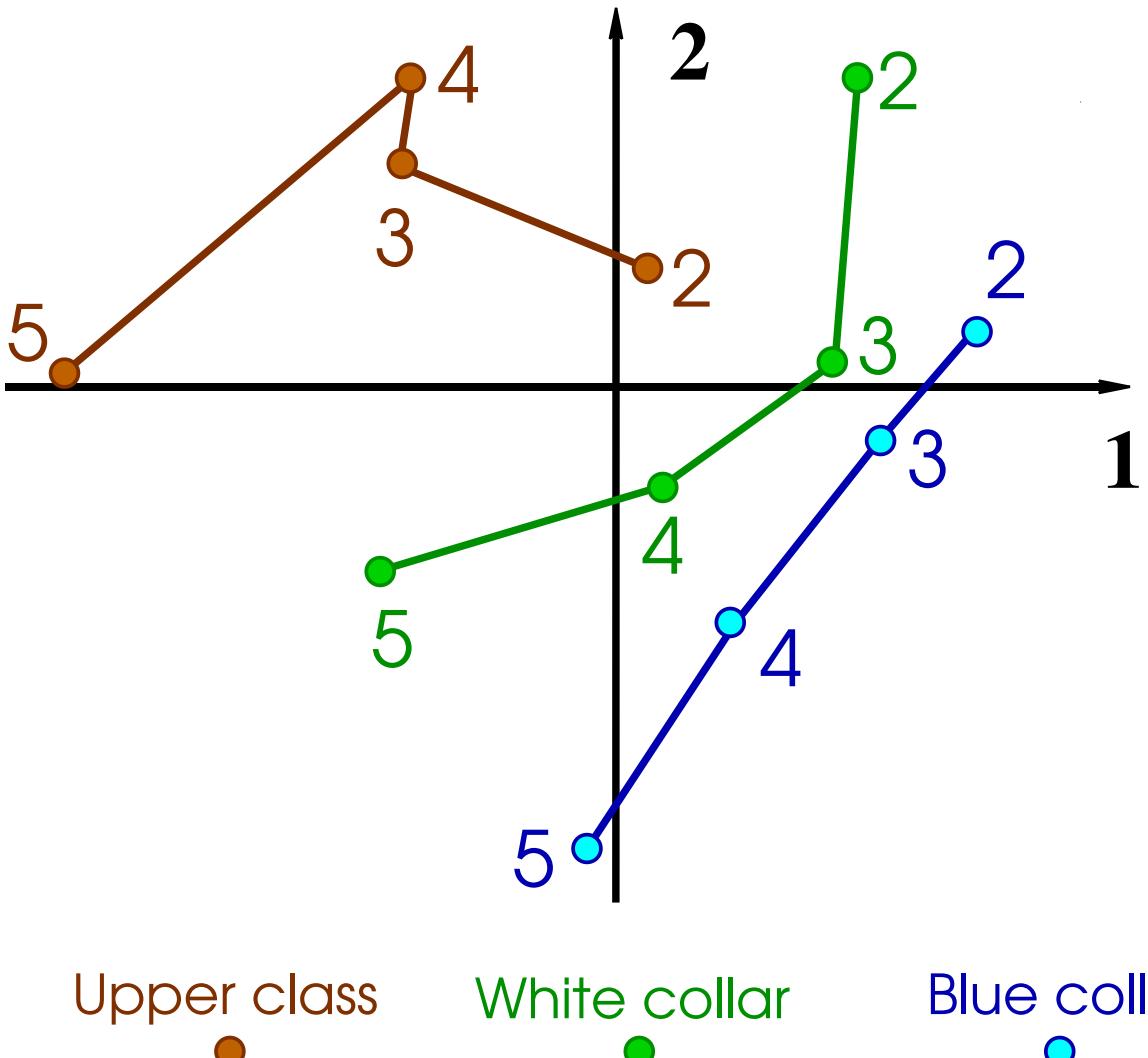


# THE VARIABLES

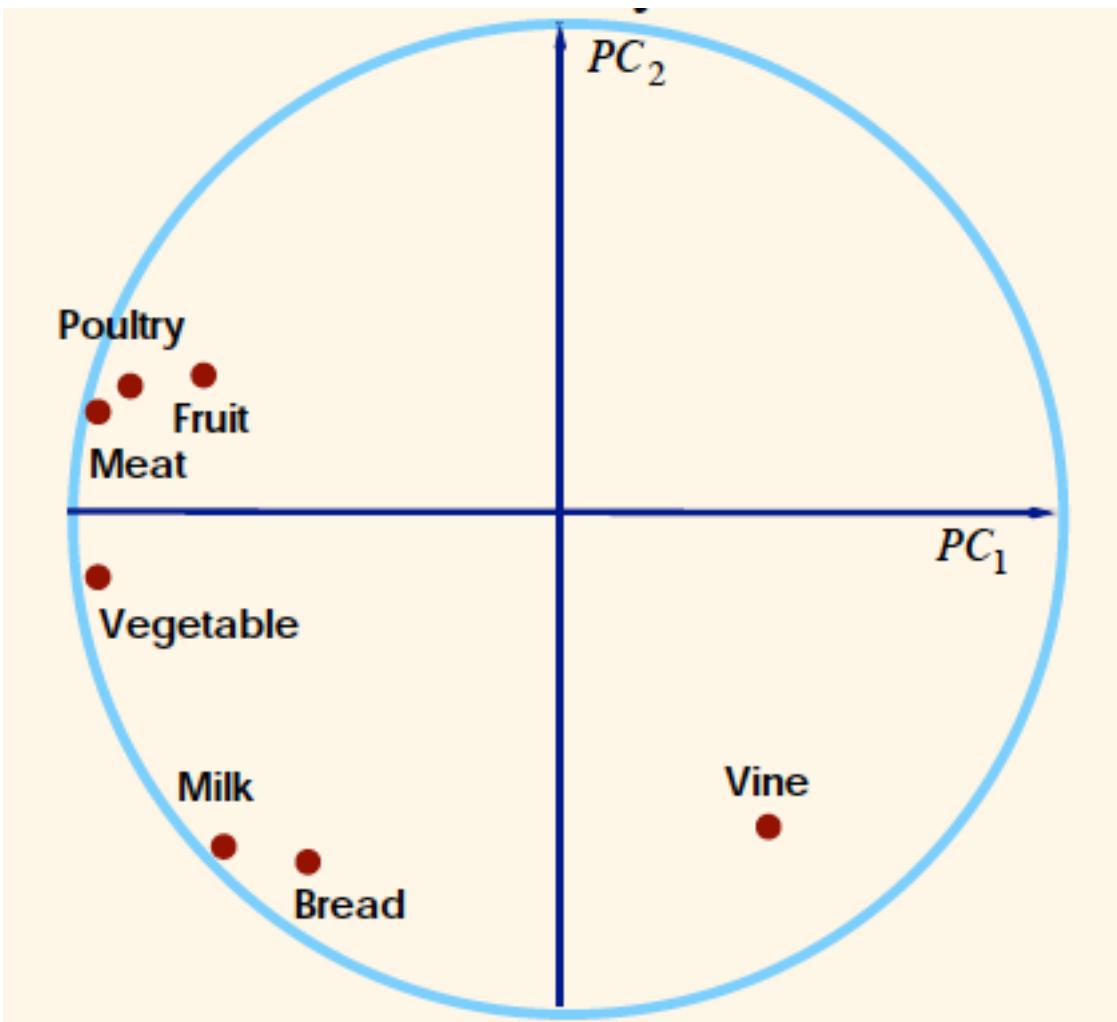


# THE Z STORY

## THE Z-STORY



# THE Z-STORY



## ROTATION

# IMPORTANT VARIABLES: TAKE II

# THE GREAT DIVIDE: ROTATION

Why? Psychometric Tradition

- Advocates: Thurstone (1947) and Cattell (1978)
- simplifies the factor structure
- facilitates interpretation
- increases reliability / replicability (with *different* data samples).

# SIMPLE STRUCTURE: THURSTONE'S 5 CRITERIA

1. Each row contains at least one zero
2. For each column, there are at least as many zeros as there are columns (i.e., number of factors kept)
3. For any pair of factors, there are some variables with zero loadings on one factor and large loadings on the other factor;
4. For any pair of factors, there is a sizable proportion of zero loadings
5. For any pair of factors, there is only a small number of large loadings

# HOW TO DO

**First: Choose the number of Factors to Keep**

**Second: Rotate to match a criterion**

# ORTHOGONAL ROTATIONS

- Most of Rotations is Varimax
  - maximizes the variance of the loadings
- Quartimax
  - minimizes number of factors per variable
- Equimax
  - combine varimax and quartimax

# OBLIQUE ROTATIONS

- Promax: a 2 step procedure.

**Goal: create binary factors**

1. Create ?target matrix:? Start with Varimax and raise to power 3 or 4 (or more).  
This creates binary factors.
2. Fit the original factors to the target with procrustes analysis (What is that?).

- Newcomer:

- Independent Component Analysis (ICA)

# NEW IDEAS: SPARSIFICATION AND LASSO

# WHEN TO ROTATE

- French School: Never!
- American school: Always!
- So, the question is:
  - When does rotating make sense?
  - When the Psychometric model is valid!

# SILLY EXAMPLE FOR ROTATION

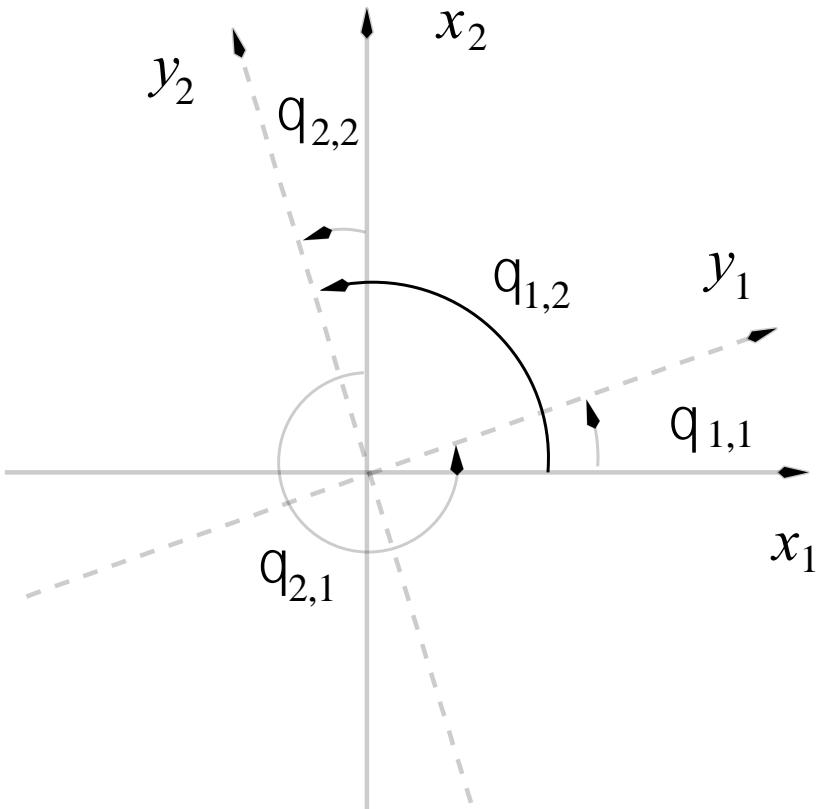


Table 10: An (artificial) example for PCA and rotation. Five wines are described by seven variables.

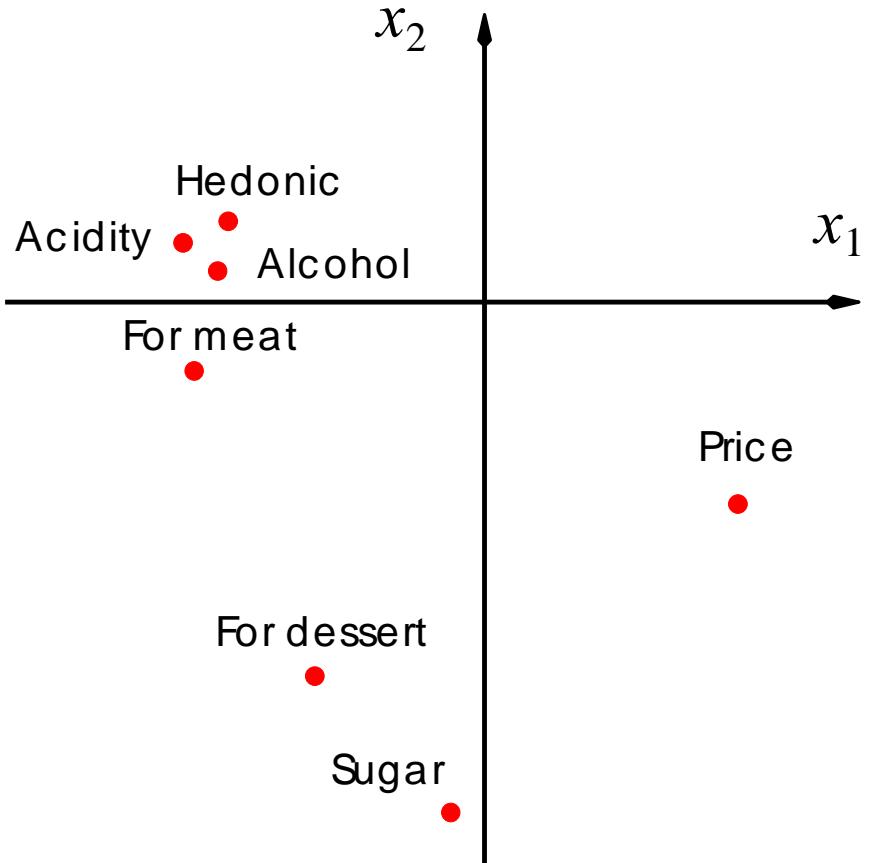
	Hedonic	For meat	For dessert	Price	Sugar	Alcohol	Acidity
Wine 1	14	7	8	7	7	13	7
Wine 2	10	7	6	4	3	14	7
Wine 3	8	5	5	10	5	12	5
Wine 4	2	4	7	16	7	11	3
Wine 5	6	2	4	13	3	10	3

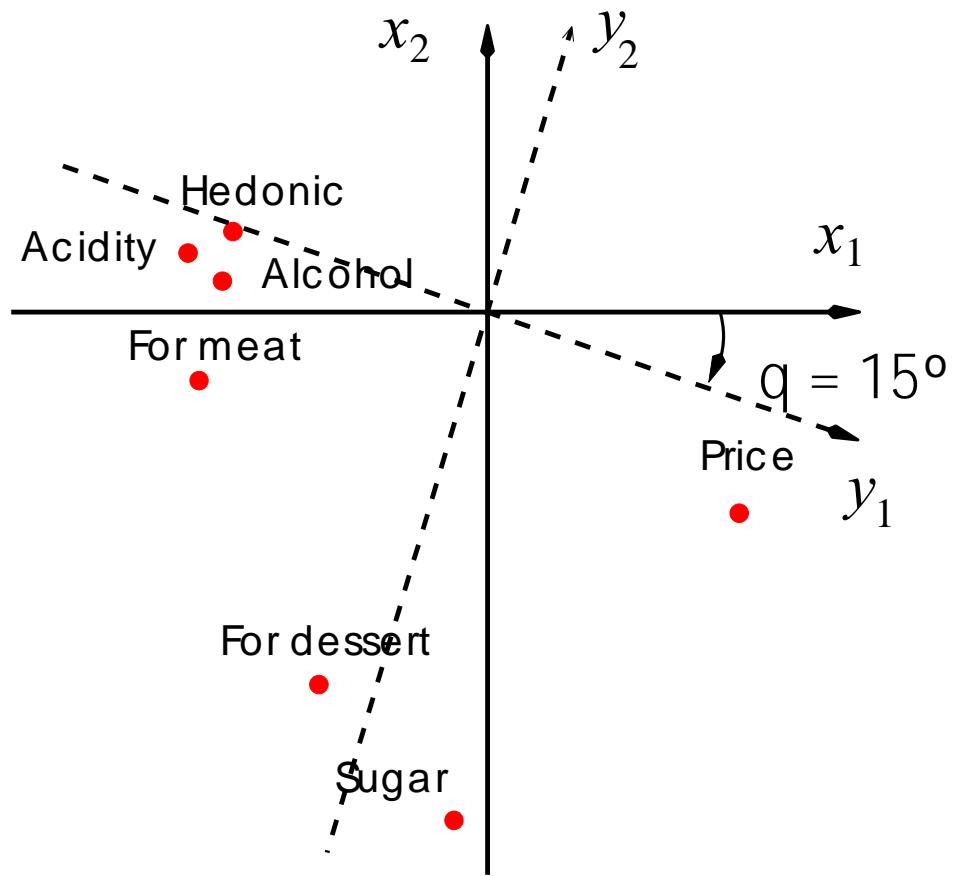
4 Factors. (Eigenvalues of 4.7627, 1.8101, 0.3527, 0.0744, respectively), 2 factors: 94% of the variance.

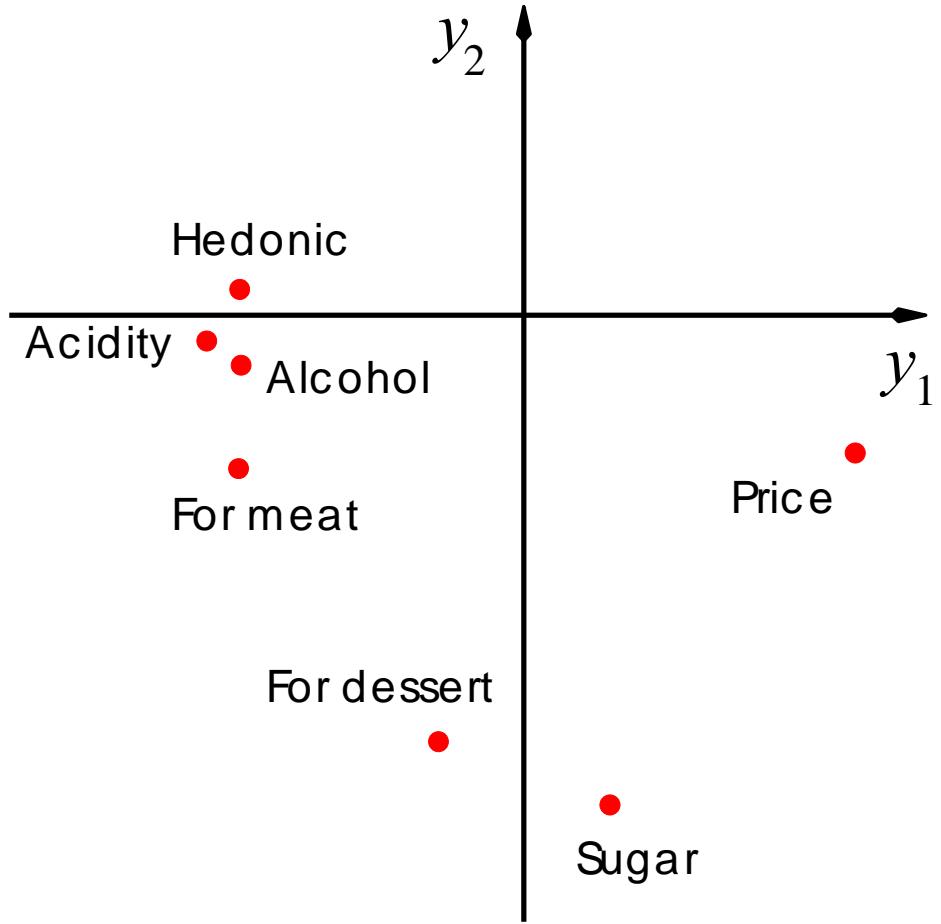
Table 11: Wine example: Original loadings of the seven variables on the first two components.

	Hedonic	For meat	For dessert	Price	Sugar	Alcohol	Acidity
Factor 1	-0.3965	-0.4454	-0.2646	0.4160	-0.0485	-0.4385	-0.4547
Factor 2	0.1149	-0.1090	-0.5854	-0.3111	-0.7245	0.0555	0.0865

# ORIGINAL LOADINGS







# ROTATION: A BETTER EXAMPLE

# ROTATIONS?

Spatial 1, 9, 13, 14, 17, 22, 25, 28, 30

I was very good in 3-D geometry as a student.

I can easily imagine and mentally rotate 3-dimensional geometric figures.

I can easily sketch a blueprint for a building that I am familiar with.

I am a good Tetris player.

I have excellent abilities in technical graphics.

I find it difficult to imagine how a 3-dimensional geometric figure

My graphic abilities would make a career in architecture relatively easy for me.

Object 4, 12, 17, 22, 25, 28, 30,

My images are very colourful and bright.

My images are very vivid and photographic.

When I imagine the face of a friend, I have a perfectly clear and bright image.

Sometimes my images are so vivid and persistent that it is difficult to ignore them.

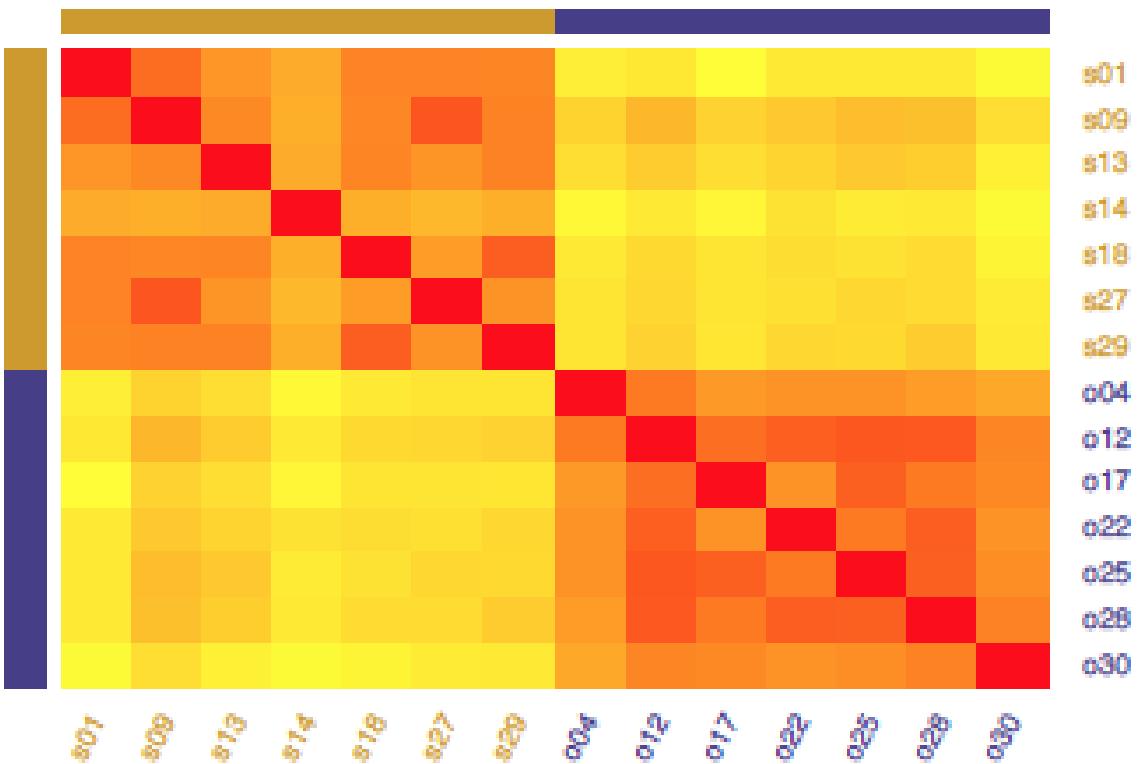
I can close my eyes and easily picture a scene that I have experienced.

My visual images are in my head all the time. They are just right there.

When I hear a radio announcer or a DJ I've never actually seen, I usually find myself picturing what he or she might look like.

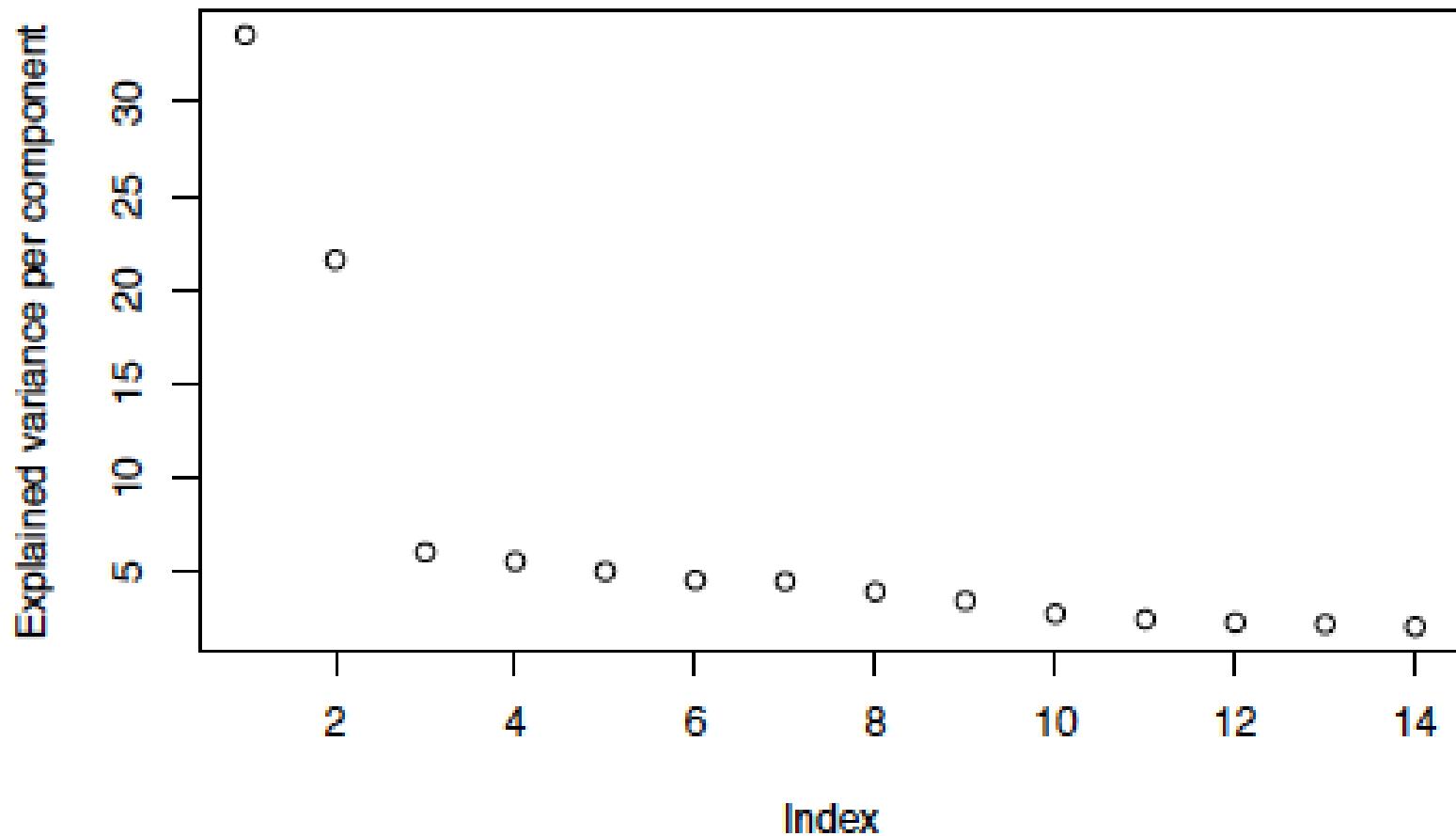
**MINI-OSIQ 14Q**

# A HEAT MAP

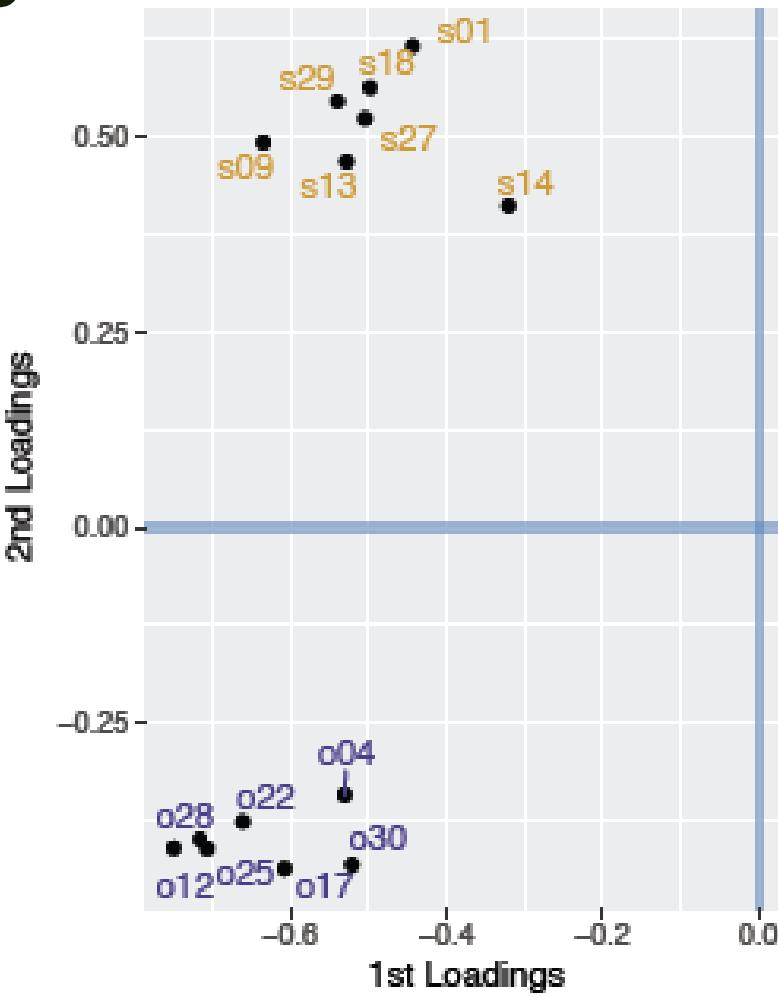
**OSIQ**

# SCREE

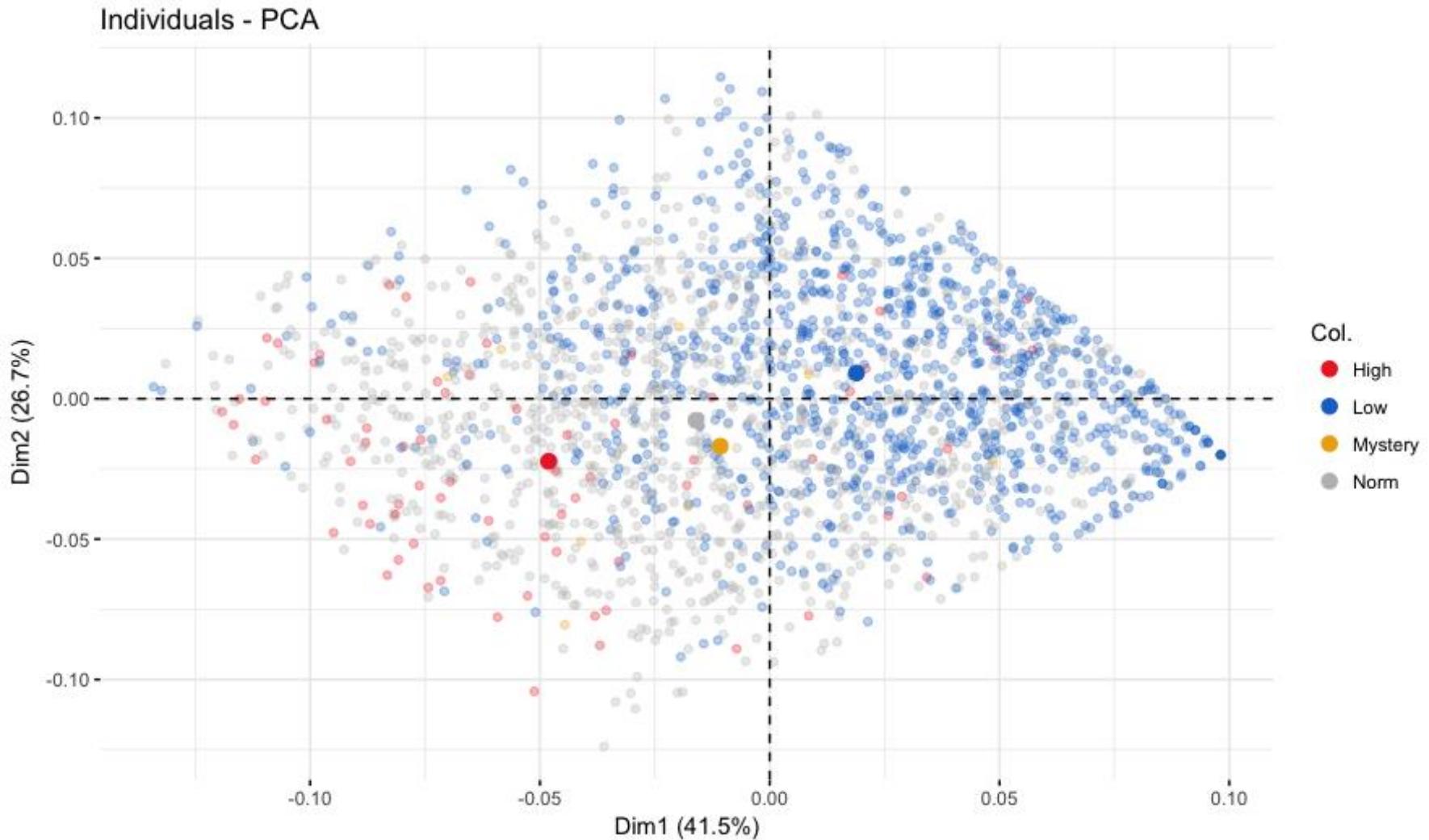
Scree plot – OSIQ

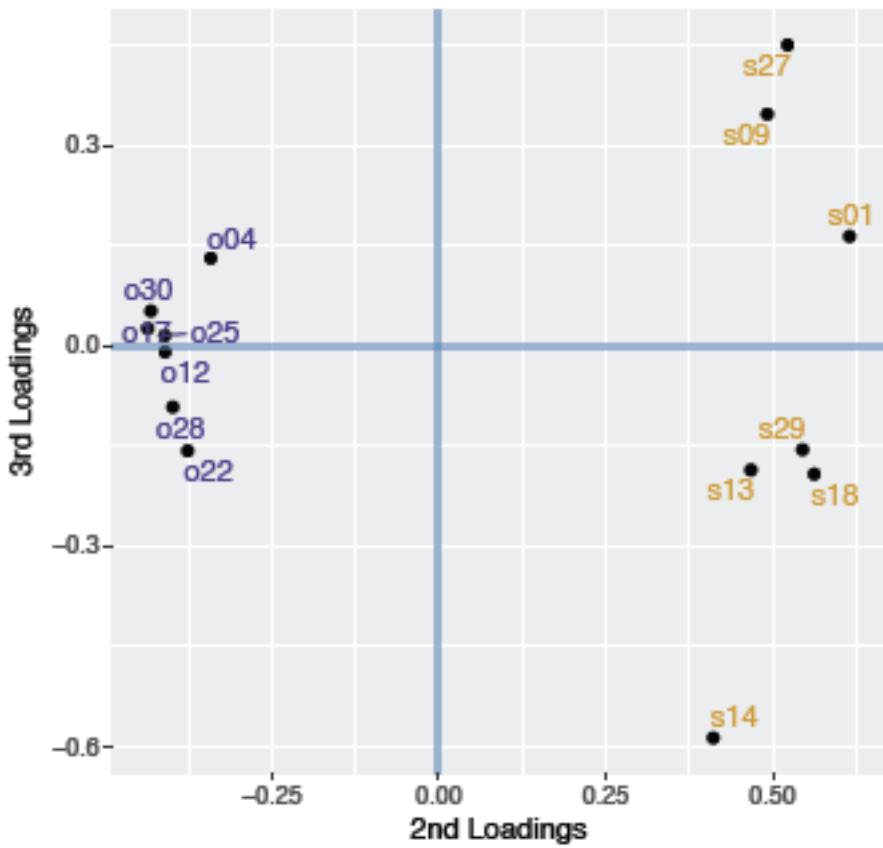


# LOADINGS



# UNROTATED I-SET





## VARIMAX HERE

```
dim.rot <- 2 # how many factors do we want to keep for rotation  
Y_varimax = varimax(ReaPCA_Y$ExPosition.Data$fj[,1:dim.rot],  
                      normalize = TRUE, eps = 1e-5)
```

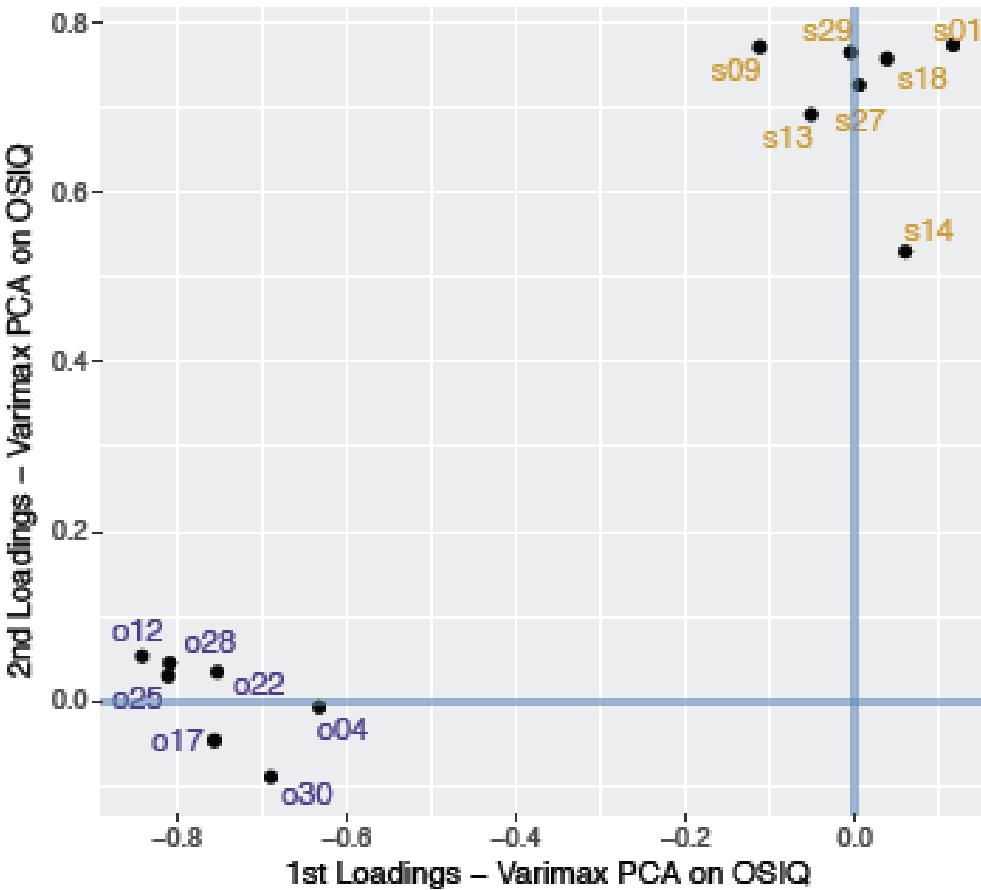
Rotation matrix from varimax:

	[,1]	[,2]
[1,]	0.7868091	-0.6171964
[2,]	0.6171964	0.7868091

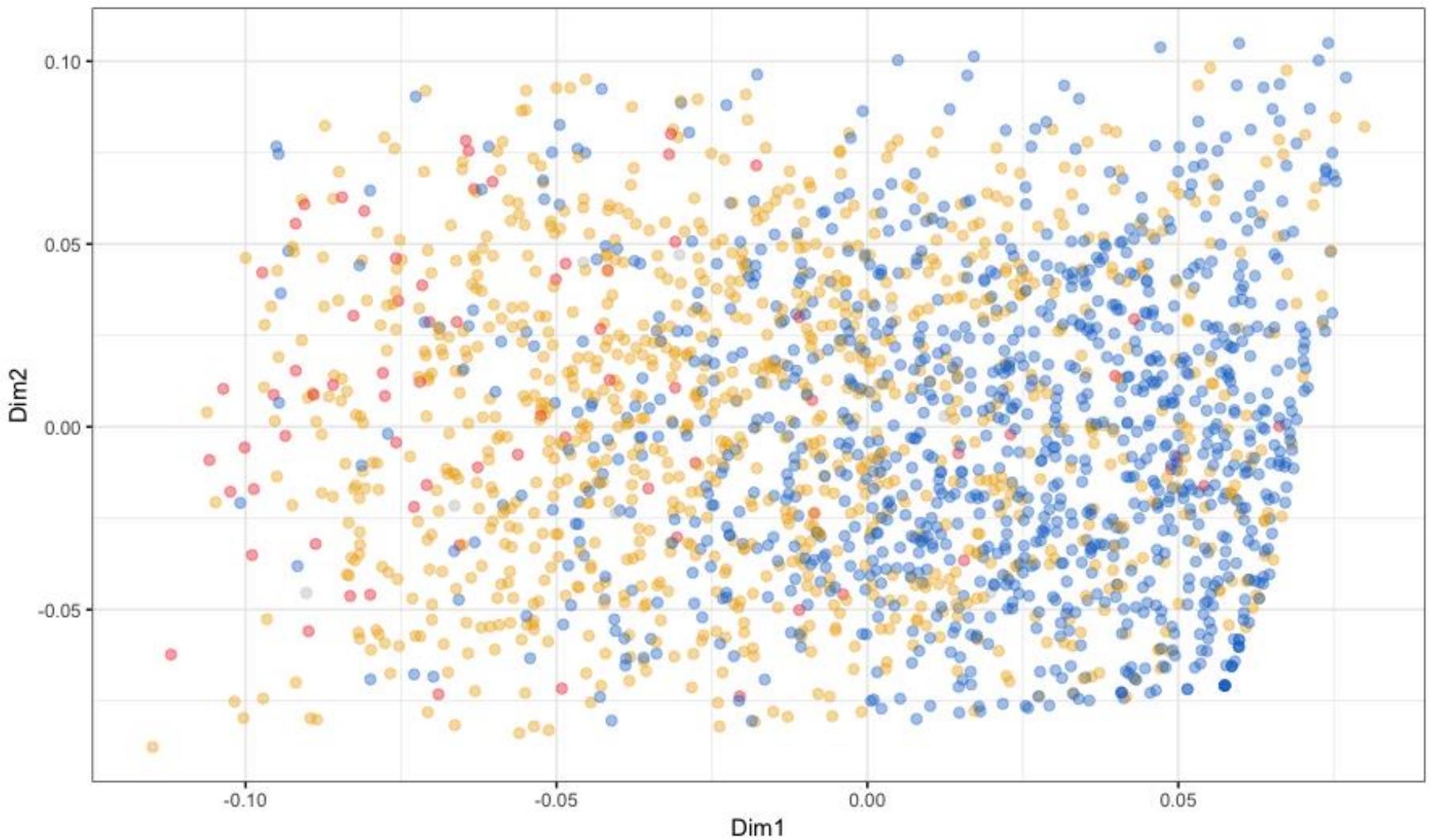
Rotation:  $\arccos(.617) = 52$  degrees

**VARIMAX**

# VARIMAX ON LOADINGS



# ROTATED I-SET



## SPARSIFICATION?

# IMPORTANT VARIABLES: TAKE III

## INFERENCES?

**ALL THAT IS DESCRIPTIVE.**

## THE EARTH IS ROUND

**THE EARTH IS ROUND:**

$P < .05$  (CF. JACOB COHEN, 1994)

**ARE THE RESULTS TO BE TRUSTED?**

**QUESTION 1. IS THERE ANYTHING IN THE DATA?**

**QUESTION 2. HOW ROBUST ARE THE RESULTS?**

## QUESTION 1.

**IS THERE ANYTHING IN THE DATA**

**COMPARE RESULTS TO CHANCE (NOTHING)**

## WHOM TO BLAME?

# PERMUTATION TEST

William S. Gosset *aka* Student  
(1876-1937)



Ronald A. Fisher  
(1890-1962)

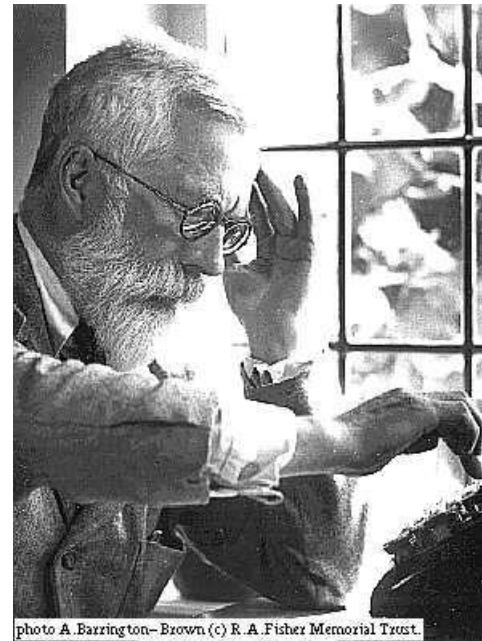


photo A Barrington-Brown (c) R.A.Fisher Memorial Trust.



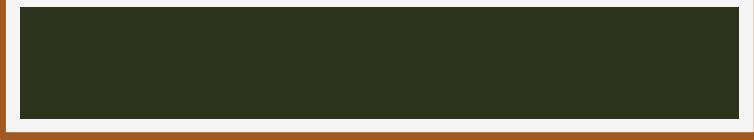
**HOW TO GET CHANCE?**

**PUT THE DATA FOR VARIABLE 1 IN A HAT**

HOW TO GET CHANCE?

PUT THE DATA FOR VARIABLE 1 IN A HAT





**HOW TO GET CHANCE?**

**PUT THE DATA FOR VARIABLE 1 IN A HAT**



# A SHAKE!



CHANCE

# HOW TO GET CHANCE?



GET THEM OUT!

## HOW TO GET CHANCE?

PUT THE DATA FOR VARIABLE 2 IN A HAT



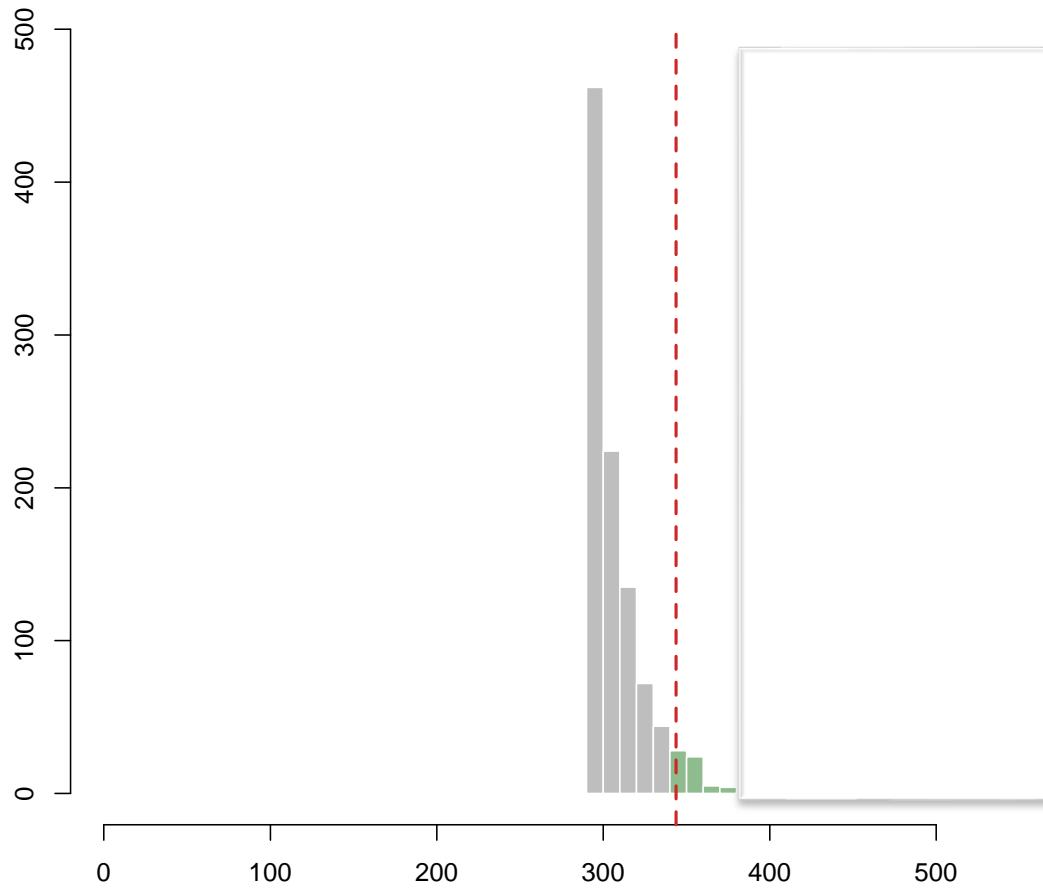
GET THEM OUT!

**SO THE LINK BETWEEN VARIABLES IS BROKEN!**

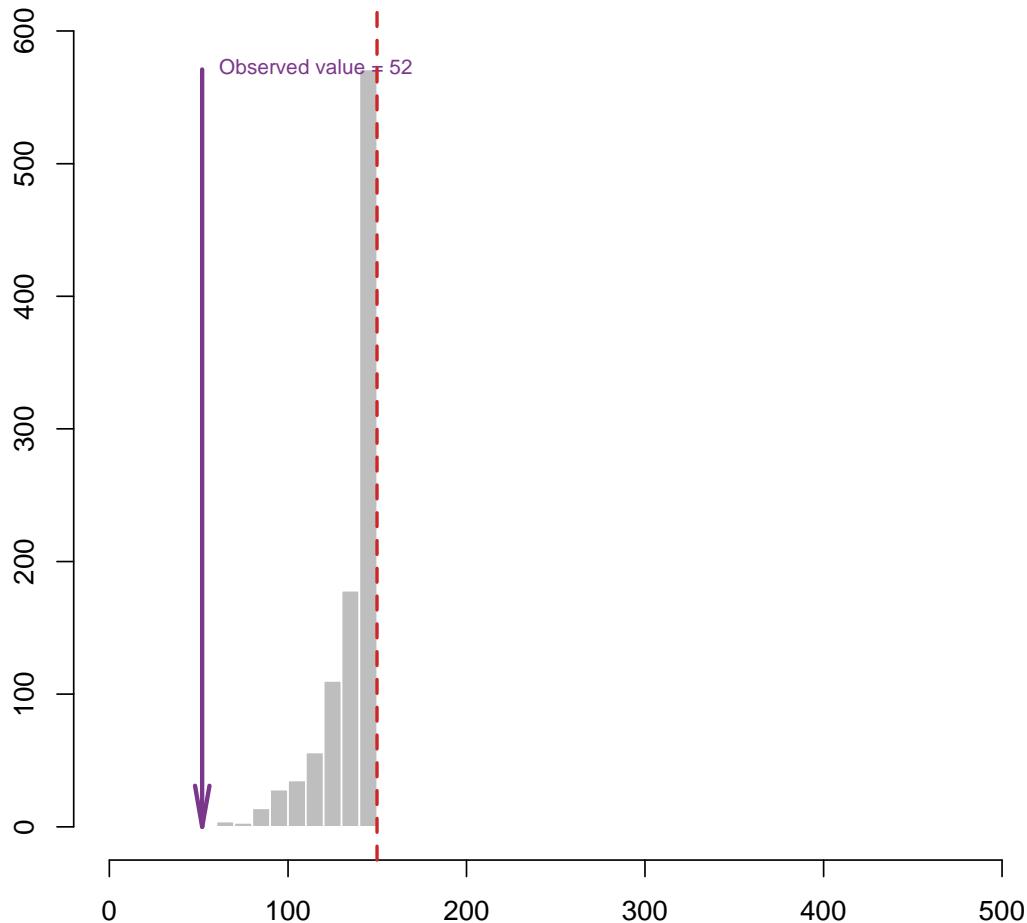
**REPEAT ALL THAT A LOT (1,000 TIMES)**

**GET THE DISTRIBUTION**

# FIRST COMPONENT (EXPLAINS 392 OUT OF 444). $P < .001$



# SECOND COMPONENT (EXPLAINS 52 OUT OF 444). VERY N.S.



# HOW MANY COMPONENTS?

SCREE TEST

EIGEN-VALUES LARGER THAN 1

EIGEN-VALUES LARGER THAN AVERAGE

PARALLEL TEST

ETC.

## BOOTSTRAP

**IMPORTANT VARIABLES FOR A COMPONENT**

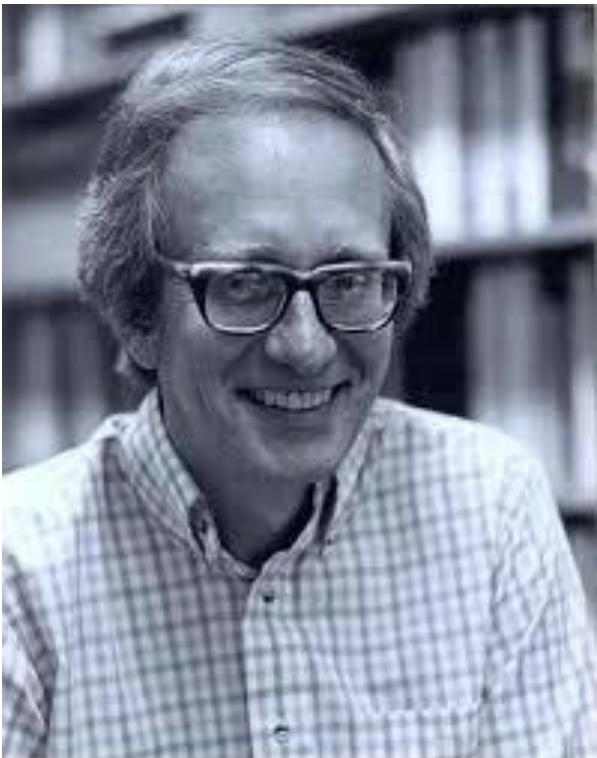
**CREATE NEW SAMPLES, AND SEE WHAT STAYS**

**THIS IS CALLED THE “BOOTSTRAP”**

WHOM TO BLAME?

# BOOTSTRAP TEST

**Bradley Efron**



- **Ideally (Population):**
  - Look for “large loading” variables
- **How to find them:**
  - Sample repeatedly from this *infinite* Population
  - Infinity big property: probabilities do *not* change

- **No Population:**
  - What to do then?
- Replace *population* by *sample*
  - Sample repeatedly from *finite* sample
  - Make sure that probabilities do not change
    - This means

**sample *with* replacement!**



# **BOOTSTRAP SAMPLING: TAKE THE 20 WINES**



**BOOTSTRAP SAMPLING:  
TAKE THE 20 WINES.  
I MEAN THE NUMBERS 1 TO 20 ...**



1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20



# AND THE BOOTSTRAP HAT



Sampling with Replacement:  
 $p$  does not change!



**Sampling with Replacement:  
 $p$  does not change!**



**Sampling with Replacement:  
 $p$  does not change!**



A top hat containing 20 numbered balls (1 through 20) arranged in four rows. A red arrow points from the number 12 in the second row to the brim of the hat.

1	5	7	8	10	11	13	16
2	3	4	9	12	15	17	19
				14	18		20

Sampling with Replacement:  
 $p$  does not change!

A SHAKE!





10



10



10



10

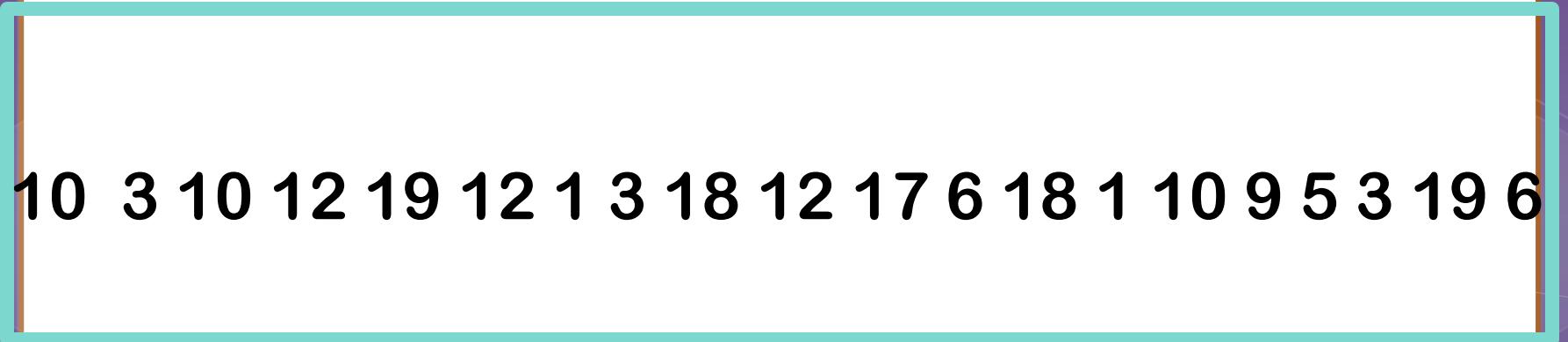


10



**AND SO ON ...**

**TILL WE HAVE 20 OBSERVATIONS**



**10 3 10 12 19 12 1 3 18 12 17 6 18 1 1 0 9 5 3 19 6**

**SAMPLING WITH REPLACEMENT  
SOME ARE REPEATED, SOME ARE LOST**

DO THAT AGAIN

RUN A PCA ON NEW SAMPLE

REPEAT THAT A LOT (1,000 TIMES)

COMPUTE LOADINGS FOR VARIABLES

GET CONFIDENCE INTERVALS

COMPUTE MEANS AND STANDARD DEVIATIONS

BOOTRAP RATIOS:  $BR = M / S$

*BR* ARE LIKE STUDENT'S  $T$

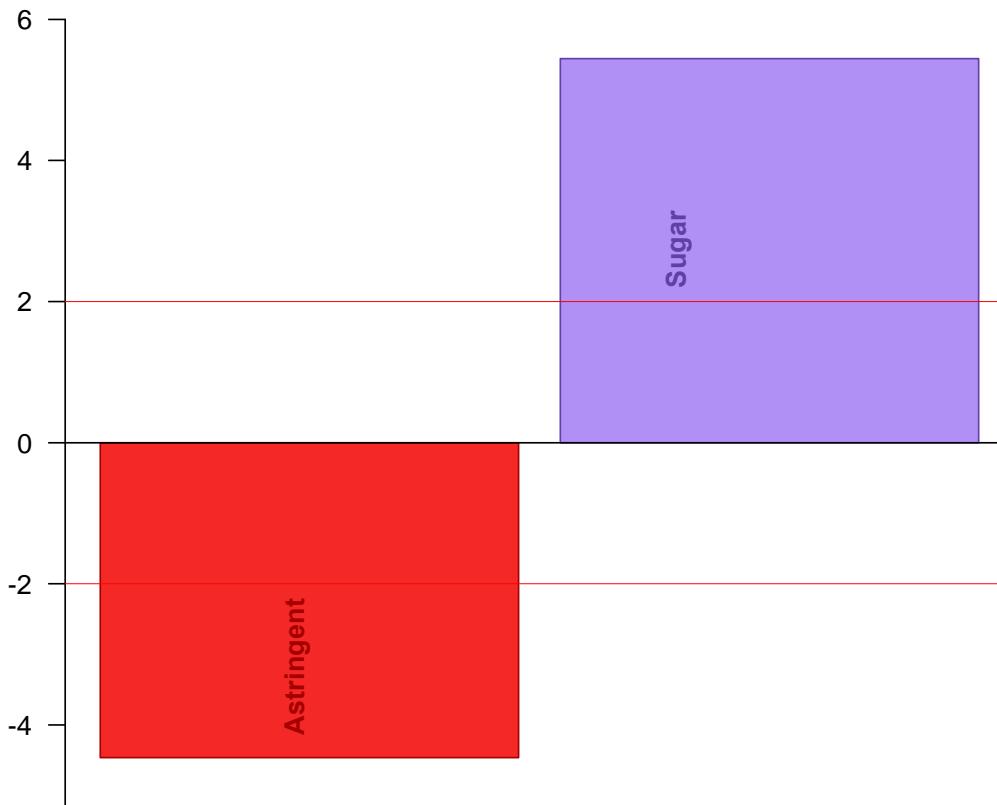
SO  $BR > 2$  SAYS  $P < .05$

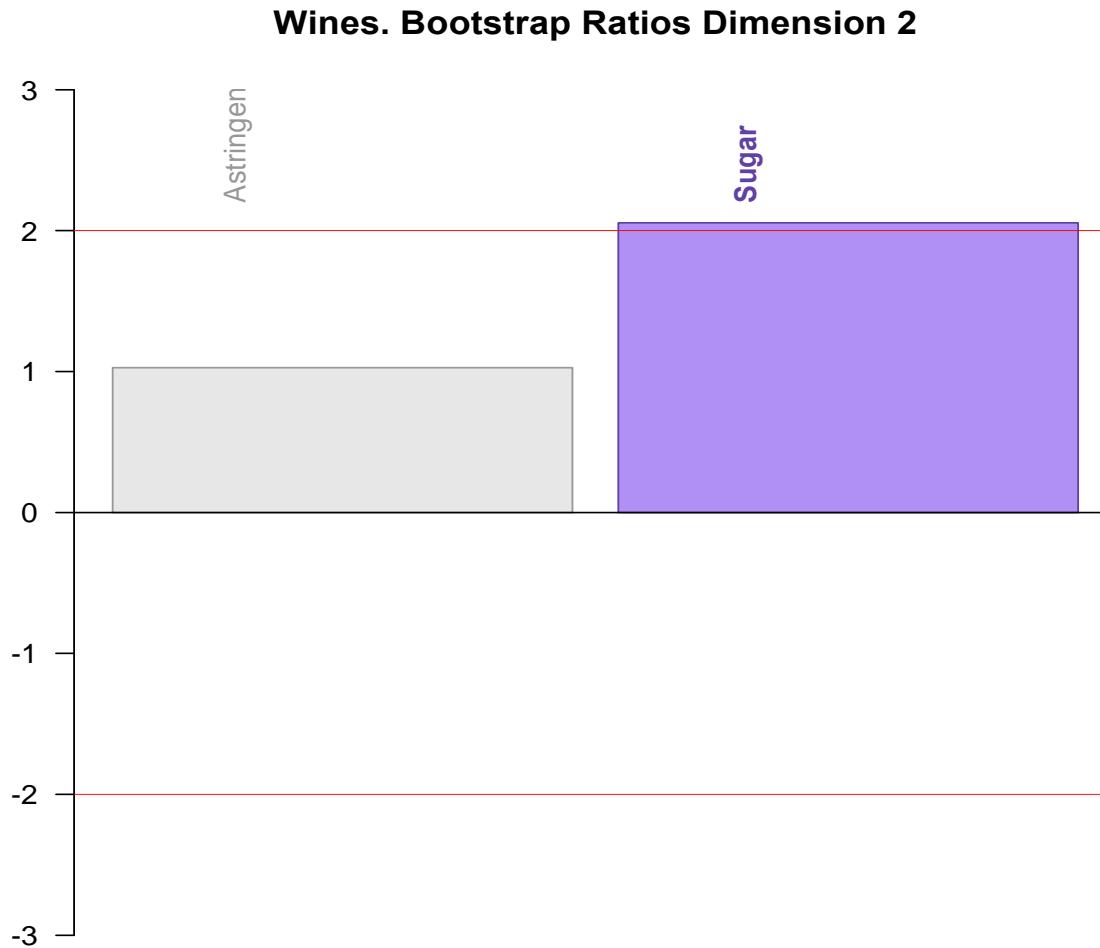
... AND, OF COURSE, THE EARTH IS ROUND!

## DIMENSION 1

# BOOTSTRAP RATIOS. PC1. BR > 2, P < .05

Wines. Bootstrap Ratios Dimension 1



**DIMENSION 2****BOOTSTRAP RATIOS. PC2. BR > 2, P < .05**

## THE STORY SO FAR

**WHAT HAVE LEARN SO FAR?**

**WE CAN DESCRIBE MULTIVARIATE DATA  
(WITH PICTURES)**

**INFERENCE 1:  
WHAT IS REAL (PERMUTATION)**

**INFERENCE 2:  
WHAT IS IMPORTANT (BOOTSTRAP)**

THIS IS THE END, MY  
FRIEND!...

**THANK YOU ALL!**

**MERCI ...**

## YOUR TURN

**QUESTIONS?**

**COMMENTS?**

## ADDITIONAL SLIDES:

WHAT IS BEHIND  
THE EIGEN-FAIRY,  
THE SVD FAIRY  
AND  
THE MAGIC LAGRANGIAN ...

## A MAXIMIZATION PROBLEM

$\mathbf{X}$  is an  $I$  by  $J$  matrix

We want vectors:

$I \times 1$   $\mathbf{p}$  and  $J \times 1$   $\mathbf{q}$

such that:

$$\arg \max_{\mathbf{p}, \mathbf{q}} \delta = \mathbf{p}^\top \mathbf{X} \mathbf{q} \text{ with } \mathbf{p}^\top \mathbf{p} = \mathbf{q}^\top \mathbf{q} = 1$$

Equivalent to

$$\mathbf{X}\mathbf{q} = \delta\mathbf{p}$$

$$\mathbf{X}^\top \mathbf{p} = \delta\mathbf{q}.$$

...

$$\mathbf{X}\mathbf{q} = \delta\mathbf{p}$$

$$\mathbf{X}^\top \mathbf{p} = \delta\mathbf{q}.$$

and also equivalent to

$$\mathbf{q}^\top \mathbf{X}^\top \mathbf{X}\mathbf{q} = \delta^2 \mathbf{p}^\top \mathbf{p} = \delta^2$$

$$\mathbf{p}^\top \mathbf{X}\mathbf{X}^\top \mathbf{p} = \delta^2 \mathbf{q}^\top \mathbf{q} = \delta^2.$$

# THE LAGRANGIAN!

$$\mathcal{L} = \mathbf{q}^\top \mathbf{X}^\top \mathbf{X} \mathbf{q} + \mathbf{p}^\top \mathbf{X} \mathbf{X}^\top \mathbf{p}$$

$$- 2\delta^2$$

$$- \alpha (\mathbf{q}^\top \mathbf{q} - 1)$$

$$- \beta (\mathbf{p}^\top \mathbf{p} - 1)$$

Partial derivatives of  $\mathcal{L}$  for  $\mathbf{q}$  and  $\mathbf{p}$

## PARTIAL DERIVATIVES

Partial derivatives of  $\mathcal{L}$  for  $\mathbf{q}$  and  $\mathbf{p}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{q}} = 2\mathbf{X}^\top \mathbf{X}\mathbf{q} - 2\alpha\mathbf{q}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{p}} = 2\mathbf{X}\mathbf{X}^\top \mathbf{p} - 2\beta\mathbf{p} .$$

... SVD ...

## SET THE PARTIAL DERIVATIVES TO 0

$$\mathbf{X}^\top \mathbf{X} \mathbf{q} = \alpha \mathbf{q}$$

$$\mathbf{X} \mathbf{X}^\top \mathbf{p} = \beta \mathbf{p} .$$

LOOKS LIKE AN EIGEN-PROBLEM!

EIGEN?

$$\mathbf{X}^\top \mathbf{X} \mathbf{q} = \alpha \mathbf{q} \quad \mathbf{X} \mathbf{X}^\top \mathbf{p} = \beta \mathbf{p}$$