

# INTRODUCTION TO PARTIAL LEAST SQUARES REGRESSION

HERVÉ ABDI, PH.D.  
THE UNIVERSITY OF TEXAS AT DALLAS

# PARTIAL LEAST SQUARES METHODS

## ➔ Partial Least Squares

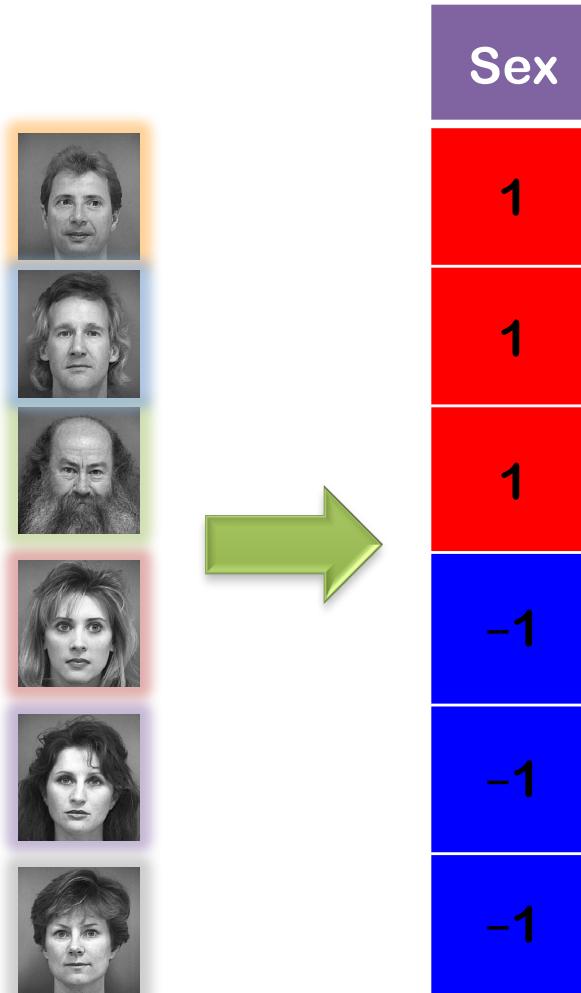
- ➔ PLS Regression – Predict performance
- ➔ PLS Correlation – Extract commonalities

## ➔ Partial Least Squares Regression

- ➔ Econometrics; Chemometrics
- ➔ H. Wold (1982), M. Martens, H. Martens, S. Wold (1983)

# PARTIAL LEAST SQUARES REGRESSION

- Goal of PLSR
  - Predict one table from another
- Components
  - Latent variables
- 6 Faces
  - Boys +1
  - Girls -1
- Constraint
  - Best prediction of Y

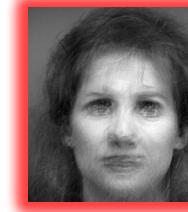
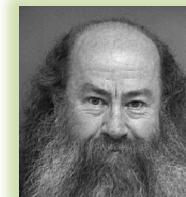
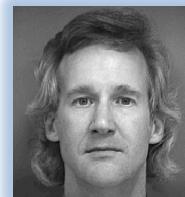


# What in these faces separate boys from girls

- To find out:
  - Mix the faces up
    - with the constraint that this mixture gives the best prediction of  $y$
- This gives the first *latent face*:
  - from  $X$  contributing to the prediction of  $y$

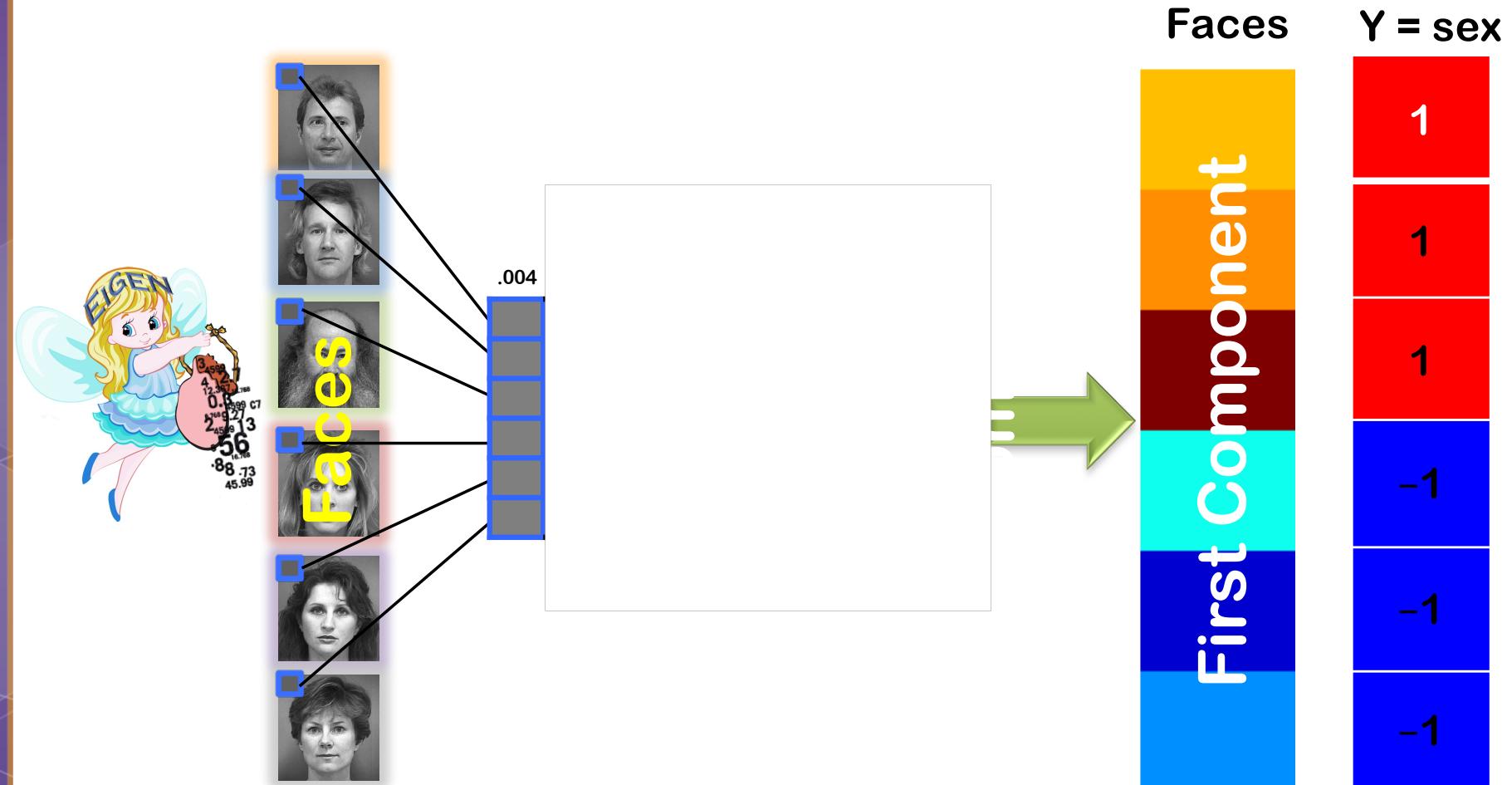
CENTER & NORMALIZE

# REMOVE WHAT IS COMMON



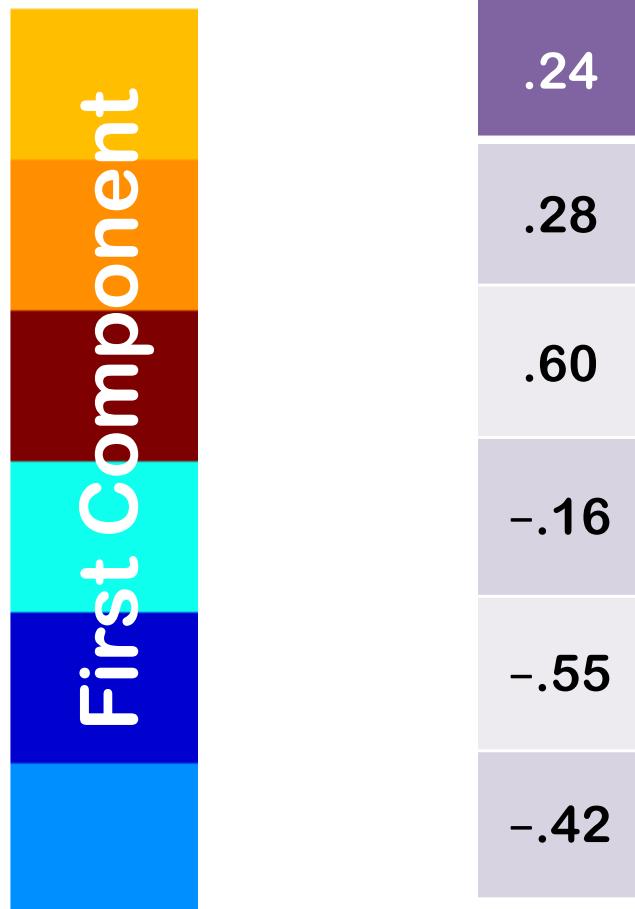
## WEIGHTS

# MIXING PIXELS TO GET SEX ...



# WITH NUMBERS: LATENT PIXELS

Faces



# HOW DOES IT LOOK? THE PIXEL COEFFICIENTS

**ANSWER** The answer is 1000.

# 55,200 of these ...

Looks  
like this

1

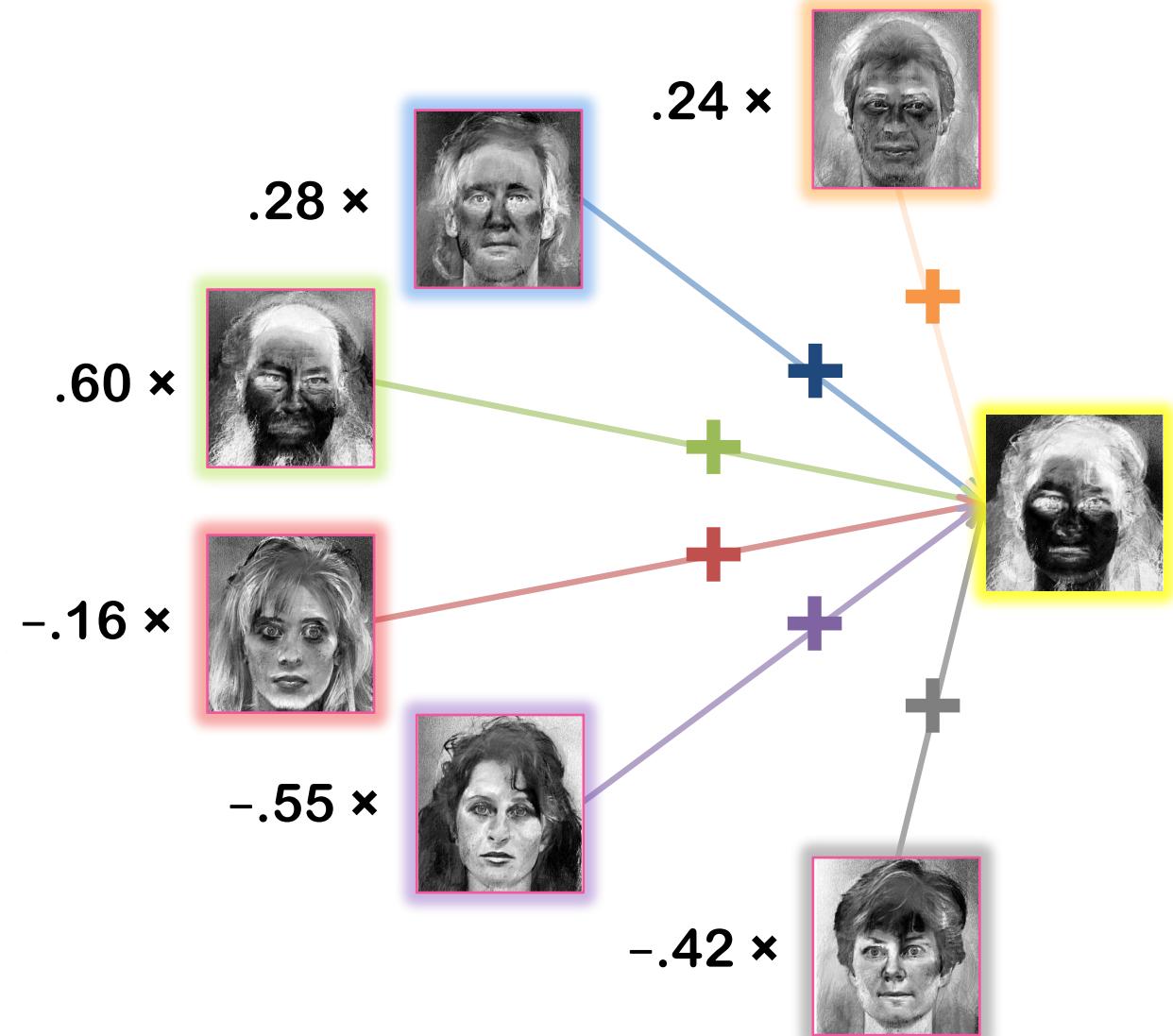


# Call it a latent face

• • •

## FIRST LATENT FACE

# DOUBLE DUTY FOR T: FACES TO LATENT-FACE



## SO THE FIRST LATENT FACE

Is a linear combination of the faces

with *positive* and *negative* weights

# BACK TO THE LATENT PIXEL OR X-LATENT VARIABLE



**QUESTION****WHAT DO WE WANT TO KNOW?**

$t_1$
.24
.28
.60
-.16
-.55
-.42



Sex
1
1
1
-1
-1
-1

... PREDICT Y FROM T

**SOLUTION**

IT'S A SIMPLE (LINEAR) REGRESSION PROBLEM!

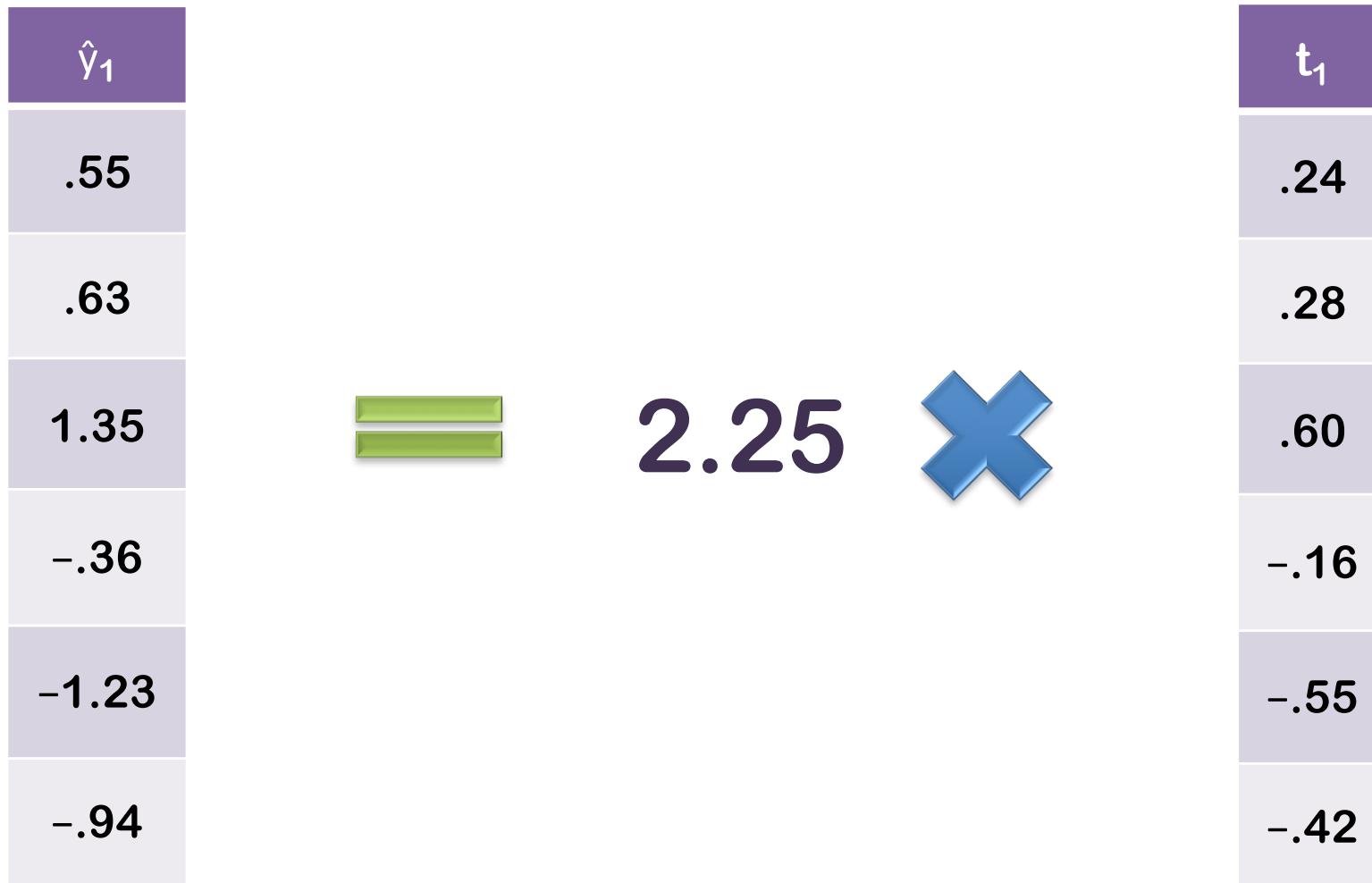
$$\hat{y}_1 = b_1$$

$t_1$
.24
.28
.60
-.16
-.55
-.42

WHY NO A?

## PREDICTING Y

### PREDICT Y FROM X



## PARTIALLING OUT

# WHAT'S LEFT BEHIND IN Y?

y
1
1
1
-1
-1
-1

$\hat{y}_1$
.55
.63
1.35
-.36
-1.23
-.94

$y - \hat{y}_1$
.45
.37
-.35
-.64
.23
-.06

## PARTIALLING OUT

# WHAT IS LEFT BEHIND IN X?



-



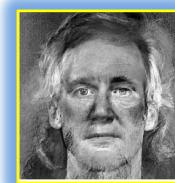
=



-



=



-



=



|-



-



=



-



=



-



=



**SUBTRACT = TO PARTIAL OUT**

$$Y_2 = Y_1 - \hat{Y}_1$$

$$X_2 = X_1 - \hat{X}_1$$

**SO WE PARTIALED OUT THE EFFECT OF  
THE LATENT VARIABLE FROM X AND Y**

## WHAT'S IN A NAME

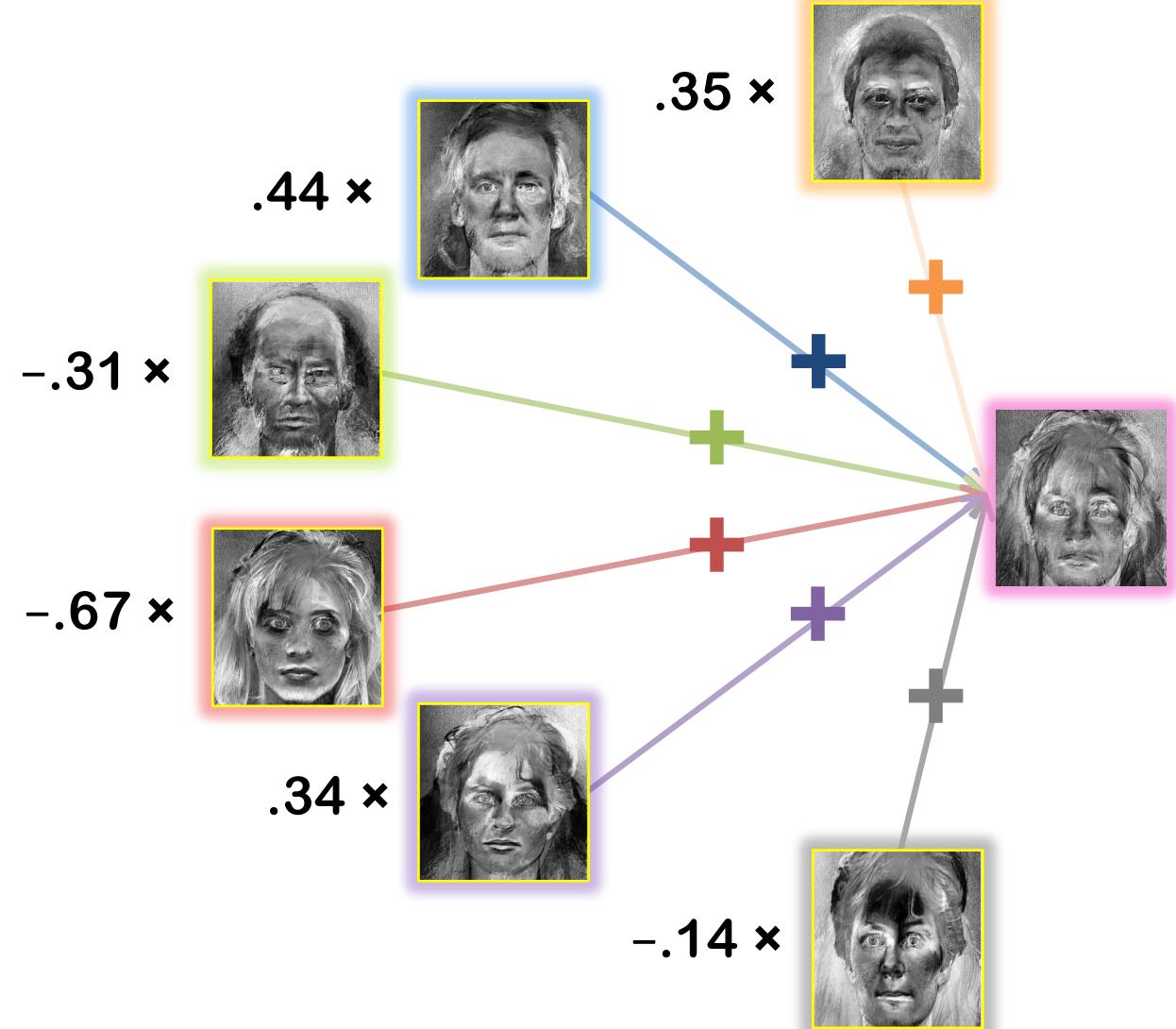
AND THEN WE KEEP ON DOING THE SAME STUFF

HENCE: *PARTIAL-LEAST SQUARE REGRESSION*

NOTE: WE NEED AN SVD FOR EVERY-ITERATION

## SECOND LATENT FACE

# WHAT'S THE NEXT BEST PREDICTION?



**SOLUTION**

$$B_2 = .95$$

$$\hat{y}_2 \quad \equiv \quad b_2$$

$t_2$
.35
.44
-.31
-.67
.34
-.14

## PARTIALLING OUT

### PREDICT $\hat{Y}_2$

$\hat{Y}_2$
.33
.42
-.30
-.64
.32
-.14



.95



$t_2$
.35
.44
-.31
-.67
.34
-.14

## PARTIALLING OUT

# WHAT IS LEFT BEHIND IN $\hat{Y}_1$ ?

$y_1$	$\hat{y}_2$	$y_1 - \hat{y}_2$
.45	.33	.12
.37	.42	-.05
-.35	-.30	-.06
-.64	-.64	-.00
.23	.32	-.09
-.06	-.14	.08

## PARTIALLING OUT

# WHAT IS LEFT BEHIND IN X?



-



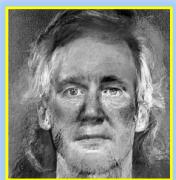
=



-



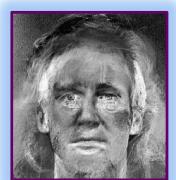
=



-



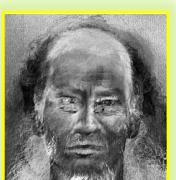
=



-



=



-



=



-



=



AND SO ON

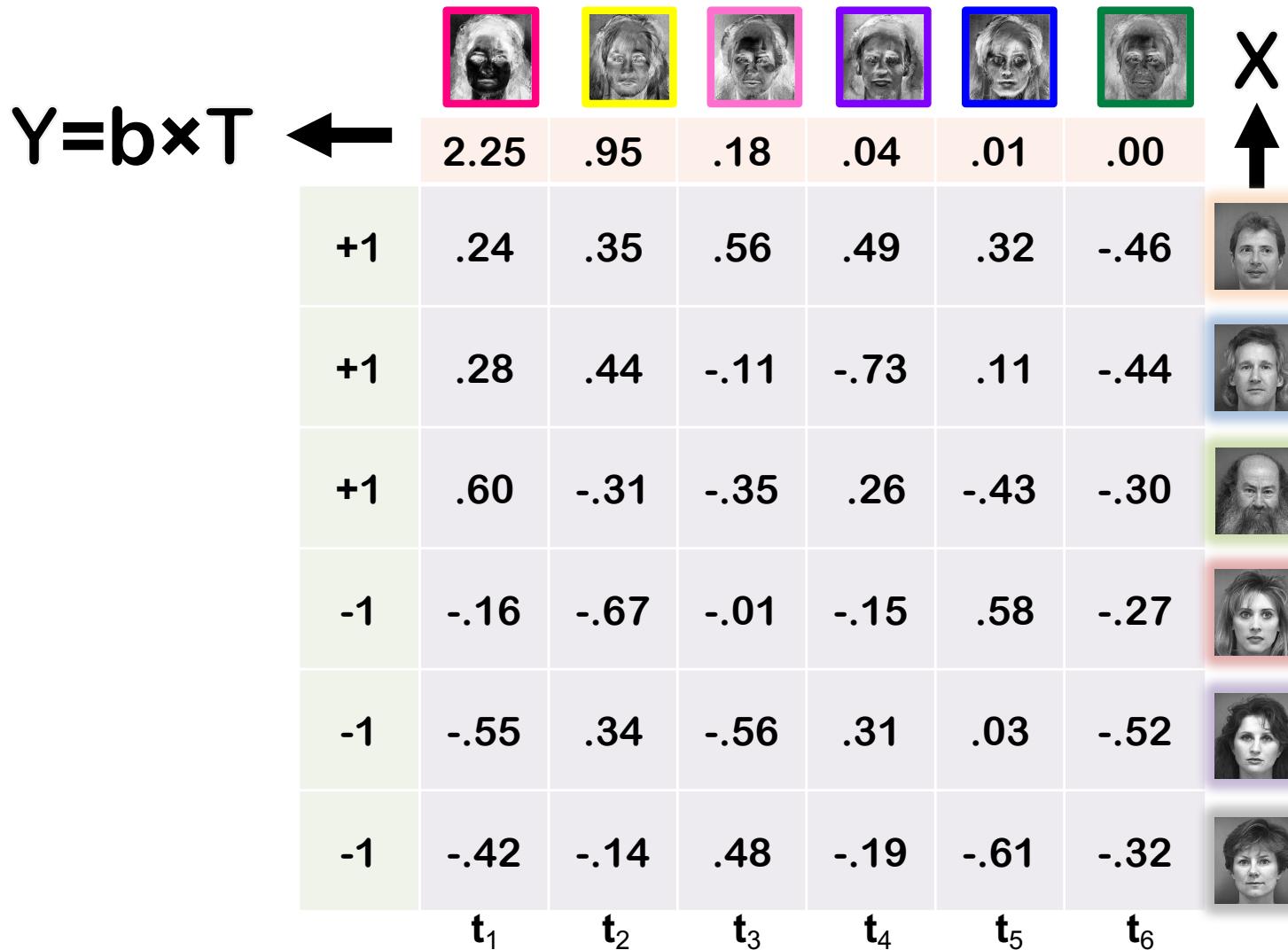
**KEEP ON, KEEP ON!**

**SUBTRACT, PREDICT, PARTIAL OUT ...**

**... TILL NOTHING IS LEFT ...**

## LATENT VARIABLES

# WHEN Do WE STOP?



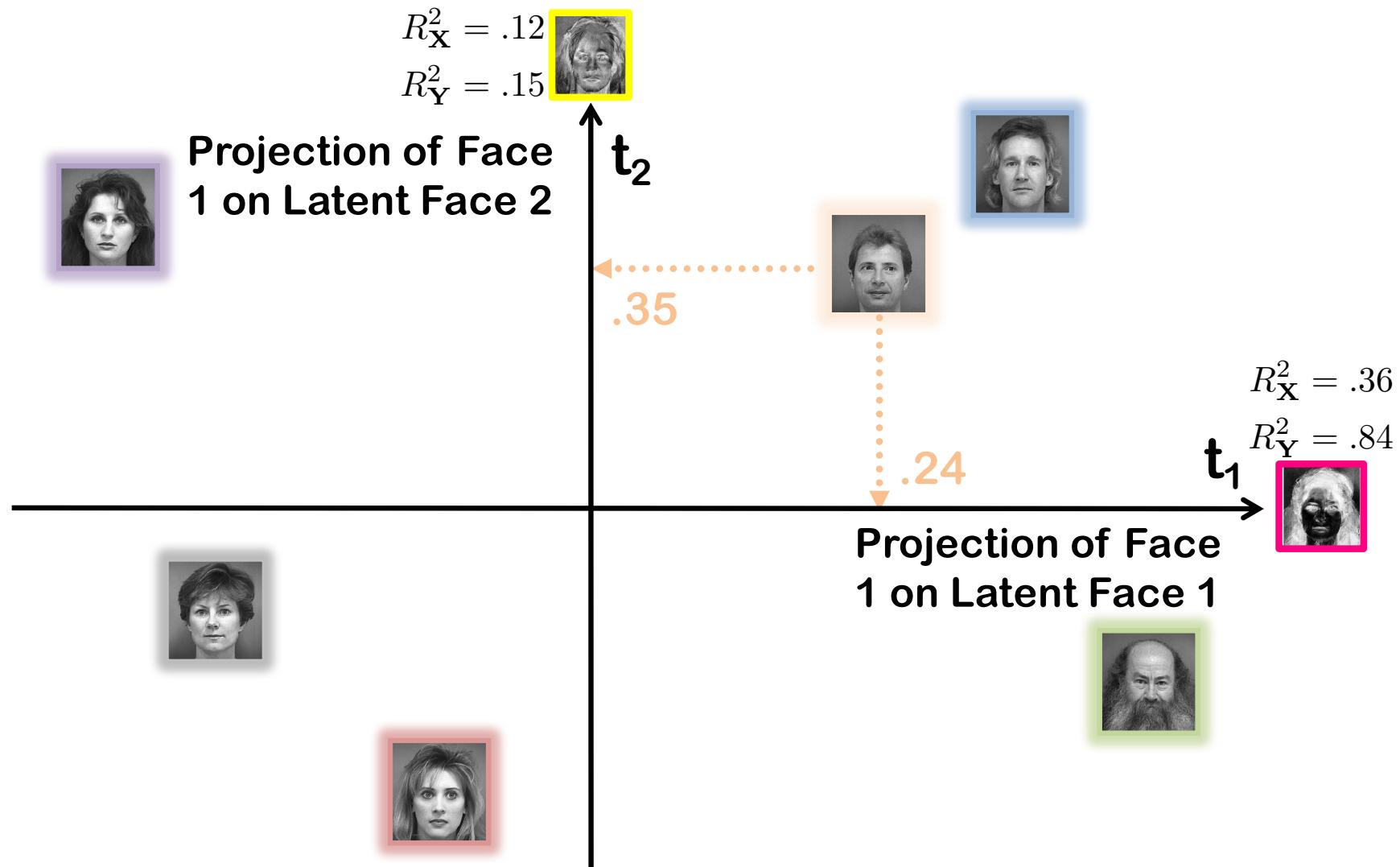
## VARIANCE EXPLAINED

# How Good Is THE PREDICTION?

	2.25	.95	.18	.04	.01	.00
						
Y % Variance	84	15	1	0	0	0
Sum Y % Variance	84	99	100	100	100	100
X % Variance	36	12	21	16	15	0
Sum X % Variance	36	48	69	85	100	100
Sum X eigenvalues	94	96	97	98	99	100
	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$

## PLOTTING

# WHAT TO DO WITH LATENT VARIABLES?



# How Do WE Do It?



## FORMAL PLS

**GET FORMAL!**

**PLS IS A 2-STEP PROCEDURE**

## KIND OF DOUBLE SVD (NOT A REAL ONE)

$$\mathbf{X} = \mathbf{T}\mathbf{P}^T \text{ with } \mathbf{T}^T\mathbf{T} = \mathbf{I}$$

$$\hat{\mathbf{Y}} = \mathbf{T}\mathbf{B}\mathbf{C}^T$$

## AS A REGRESSION PROBLEM (PSEUDO-INVVERSE)

$$\hat{\mathbf{Y}} = \mathbf{T}\mathbf{B}\mathbf{C}^T = \mathbf{X}\mathbf{B}_{PLS}$$

$$\mathbf{B}_{PLS} = \mathbf{X}^+ \mathbf{Y} = \mathbf{P}^{T+} \mathbf{B}\mathbf{C}^T$$

# HOW TO PREDICT Y?

- From the latent vectors, use:

$$\hat{Y} = TBC^T$$

- From X, with Z-scores (for X and Y), use:

$$\hat{Z}_Y = Z_X B_{PLS}$$

- From X, use:

$$\hat{Y} = XB_{PLS}^*$$

# NIPALS: THE ORIGINAL ALGORITHM

- Get a linear combination from X
- With maximal covariance with Y

Start with E = X, and F = Y. Initiate a random u

- Step 1.  $w \propto E^T u$  (estimate X weights)
- Step 2.  $t \propto Ew$  (estimate X factor scores)
- Step 3.  $c \propto F^T t$  (estimate Y weights)
- Step 4.  $u = Fc$  (estimate Y scores)

Iterate till t converges, then compute

$$\beta = t^T u, \text{ and } p = E^T t$$

Now partial out t from both E and F as:

$$E = E - tp^T, \text{ and } F = F - \beta tc^T.$$

# How GOOD IS THE PREDICTION

- ➔ The problem:
  - ➔ The prediction is sample-dependent
- ➔ The question:
  - ➔ How would we do on a *new* sample?
- ➔ One answer:
  - ➔ Cross-validation, e.g. *Jackknife*
  - ➔ A.K.A Leave-One-Out (LOO)

WHOM TO BLAME?

# THE JACKKNIFE

John W. Tukey  
(1915-2000)

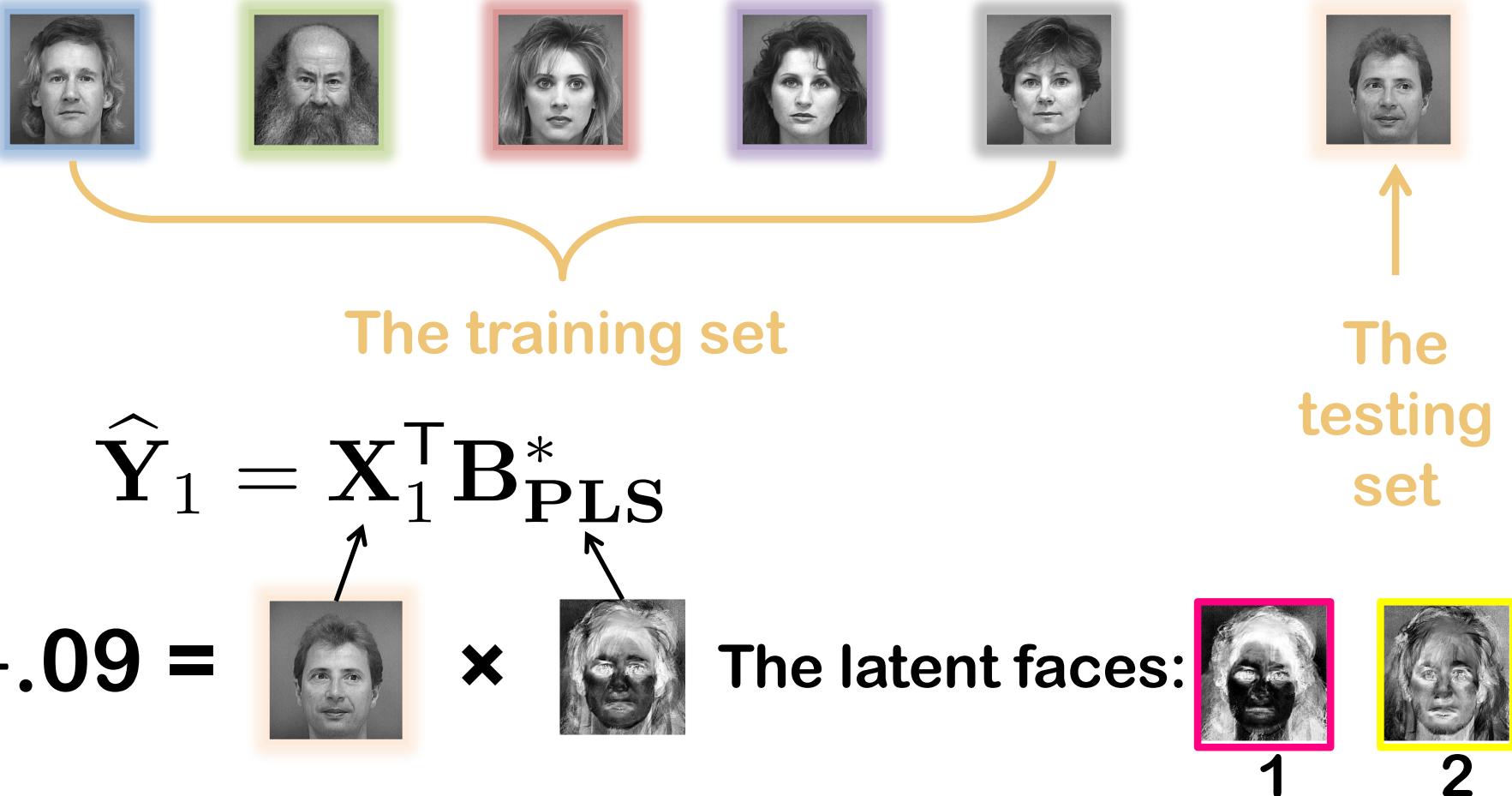


A Jackknife



## FIRST ITERATION

# JACKKNIFE



## SECOND ITERATION

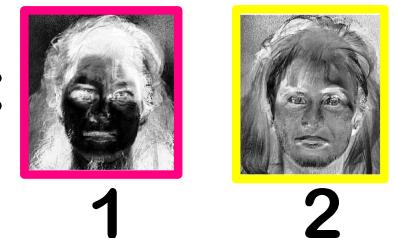
# JACKKNIFE



$$\hat{Y}_2 = X_2^T B_{PLS}^*$$

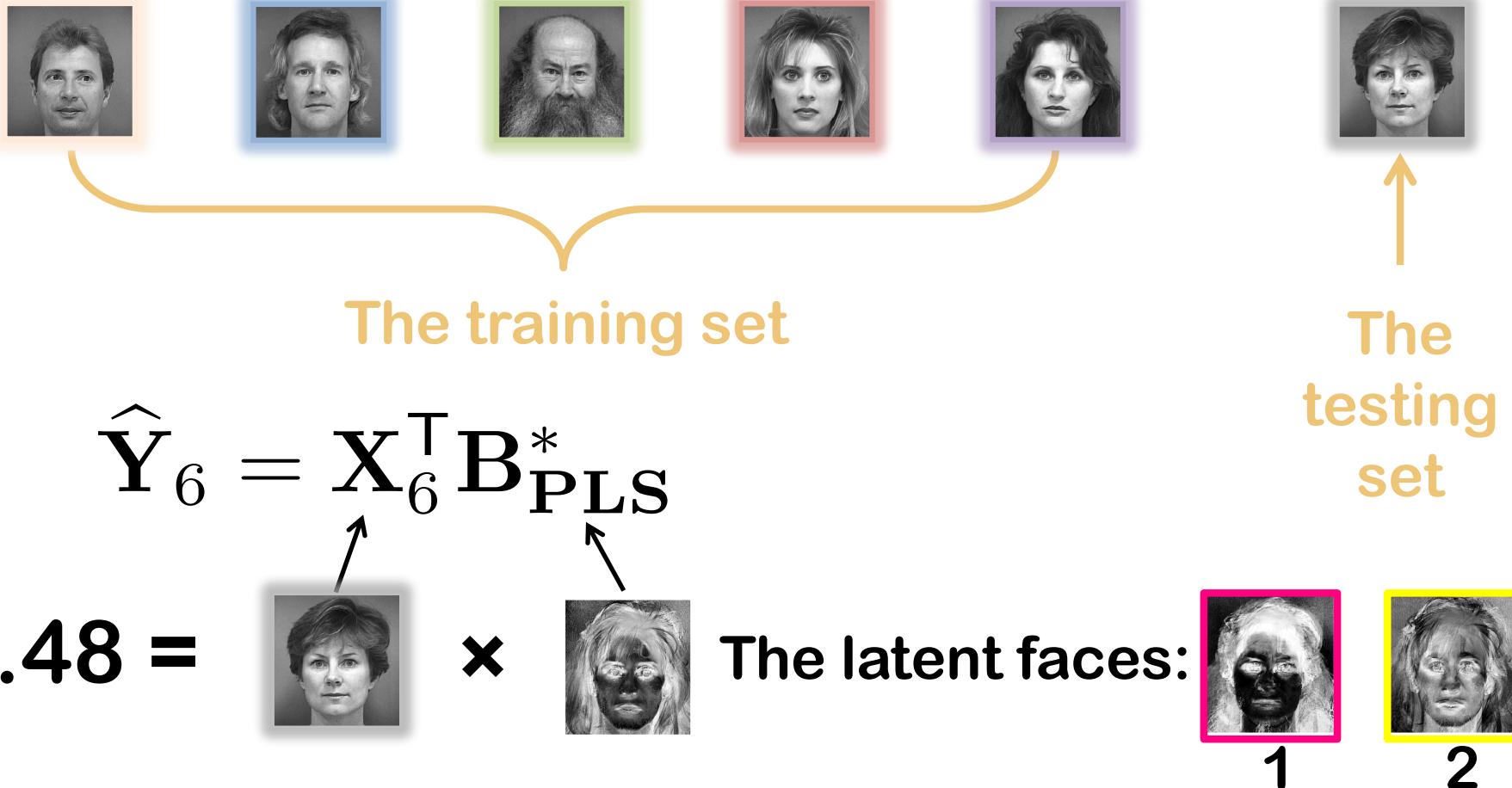
-.05 =   $\times$  

**The latent faces:**



## LAST ITERATION

# JACKKNIFE



**ALL ITERATIONS**

# JACKKNIFE SUMMARY

**The faces**

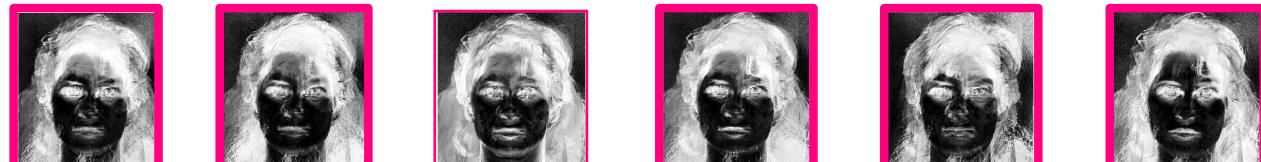


**The predictors:**

$B_{PLS}$



**1<sup>st</sup> latent faces**



**2<sup>nd</sup> latent faces**



## PCR Vs. PLSR

# MORE ABOUT SEX

134 Faces

83 Women

51 Men

PCR

	2	4	3	6	39
$R^2_Y$	.29	.13	.09	.05	.03
Sum	.29	.42	.51	.52	.59
$R^2_X = \lambda$	.02	.01	.01	.00+	.00+

PLSR

	2	4	3	6	39
$R^2_Y$	.50	.18	.06	.07	.05
Sum	.50	.69	.75	.82	.87
$R^2_X$	.16	.09	.10	.06	.04

# JACKKNIFE FOR PCR AND PLSR

## Jackknife (2 latent vector solution)

$\hat{Y}$	-.12	-.08	2.08	.40	-.50	-.60
$Y$	1	1	1	-1	-1	-1

$$r_{Y,\hat{Y}} = .48$$

$$r^2_{Y,\hat{Y}} = .23$$

## Sample (2 latent vector solution)

$\hat{Y}$	.79	.96	.97	-.91	-.83	-.98
$Y$	1	1	1	-1	-1	-1

$$r_{Y,\hat{Y}} = .99$$

$$r^2_{Y,\hat{Y}} = .98$$

# QUALITY OF CLASSIFICATION

## PCR

Fixed Effect		Jackknife	
R <sup>2</sup>	Misclassifications	R <sup>2</sup>	Misclassifications
.77	17/134	.74	20/134

## PLSR

Fixed Effect		Jackknife	
R <sup>2</sup>	Misclassifications	R <sup>2</sup>	Misclassifications
.93	1/134	.79	14/134

## HOW MANY LV'S

**ADDITIONAL PROBLEM:**

**THE CLASSIFICATION IS OPTIMAL FOR 2 LV'S**

**BUT WE USED THE DATA TO FIND THIS NUMBER**

**THIS IS ANOTHER CASE OF “DOUBLE DIPPING”**

THREE SAMPLES ...

WE NEED (AT LEAST) THREE INDEPENDENT SAMPLES

THE CLASSIFICATION IS OPTIMAL FOR 2 LVs

BUT WE USED THE DATA TO FIND THIS NUMBER

THIS IS ANOTHER CASE OF “DOUBLE DIPPING”

# A (FAKED) WINE EXAMPLE

# WHAT Do WE WANT To Do?

- ➔ Predict subjective evaluation of 5 wines
- ➔ The dependent variables for each wine:
  - ➔ Likeability,
  - ➔ How well it goes with meat
  - ➔ How well it goes with dessert
- ➔ The predictors are:
  - ➔ Price
  - ➔ Sugar
  - ➔ Alcohol
  - ➔ Acidity

**MATRIX X**

# THE PREDICTORS

Wine	Price	Sugar	Alcohol	Acidity
1	7	7	13	7
2	4	3	14	7
3	10	5	12	5
4	16	7	11	3
5	13	3	10	3

**MATRIX Y**

# THE DEPENDENT VARIABLES

Wine	Hedonic	Goes with Meat	Goes with Dessert
1	14	7	8
2	10	7	6
3	8	5	5
4	2	4	7
5	6	2	4

# FIRST STEP: GET Z-SCORES

MATRIX  $Z_X$ 

## THE PREDICTORS Z-SCORES

Wine	Price	Sugar	Alcohol	Acidity
1	-0.63	1	0.63	1
2	-1.26	-1	1.26	1
3	0	0	0	0
4	1.26	1	-0.63	-1
5	.0.63	-1	-1.26	-1

## THE DEPENDENT VARIABLES Z-SCORED

Wine	Hedonic	Goes with Meat	Goes with Dessert
1	1.34	0.94	1.26
2	0.45	0.94	0
3	0	0	-0.63
4	-1.34	-0.47	0.63
5	-0.45	-1.41	-1.26

## MATRIX R

### THE “R” MATRIX

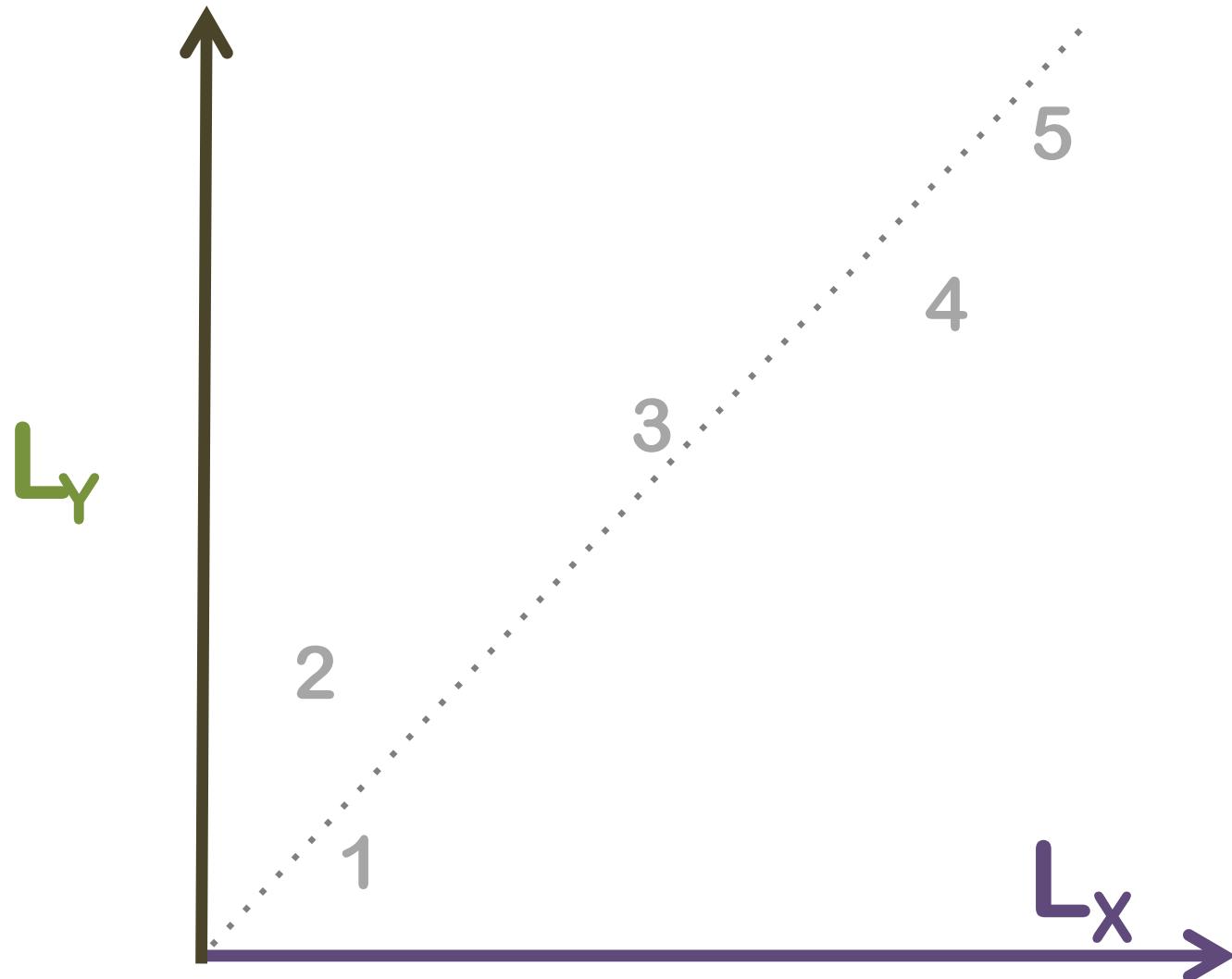
		Y		
		Hedonic	Meat	Dessert
X	Price	-3.39	-3.28	-0.80
	Sugar	0.00	0.94	3.16
	Alcohol	2.82	3.88	2.00
	Acidity	3.58	3.77	1.90

$$\mathbf{R} = \mathbf{Z}_X^\top \mathbf{Z}_Y = \mathbf{R}^* \times (I - 1)$$

(with  $\mathbf{R}^*$  being correlation matrix)

PLOT THEM

$L_X$  VS  $L_Y$



## NEW MATRIX $R_1$

### THE NEW “R” MATRIX

		$Y$		
		Hedonic	Meat	Dessert
$\times$	Price	-0.25	-0.34	1.10
	Sugar	-0.08	0.85	3.11
	Alcohol	-0.45	0.09	0.02
	Acidity	0.24	-0.08	-0.12

$$R_1 = E_1^T F_1$$

# KEEP ON DOING THAT TILL X IS GONE

**MATRIX T**

# THE LATENT VARIABLES

Wine	$t_1$	$t_2$	$t_3$
1	0.4538	-0.4662	0.5716
2	0.5399	0.4940	-0.4631
3	0	0	0
4	-0.4304	-0.5327	-0.5301
5	-0.5633	-0.5049	0.4217

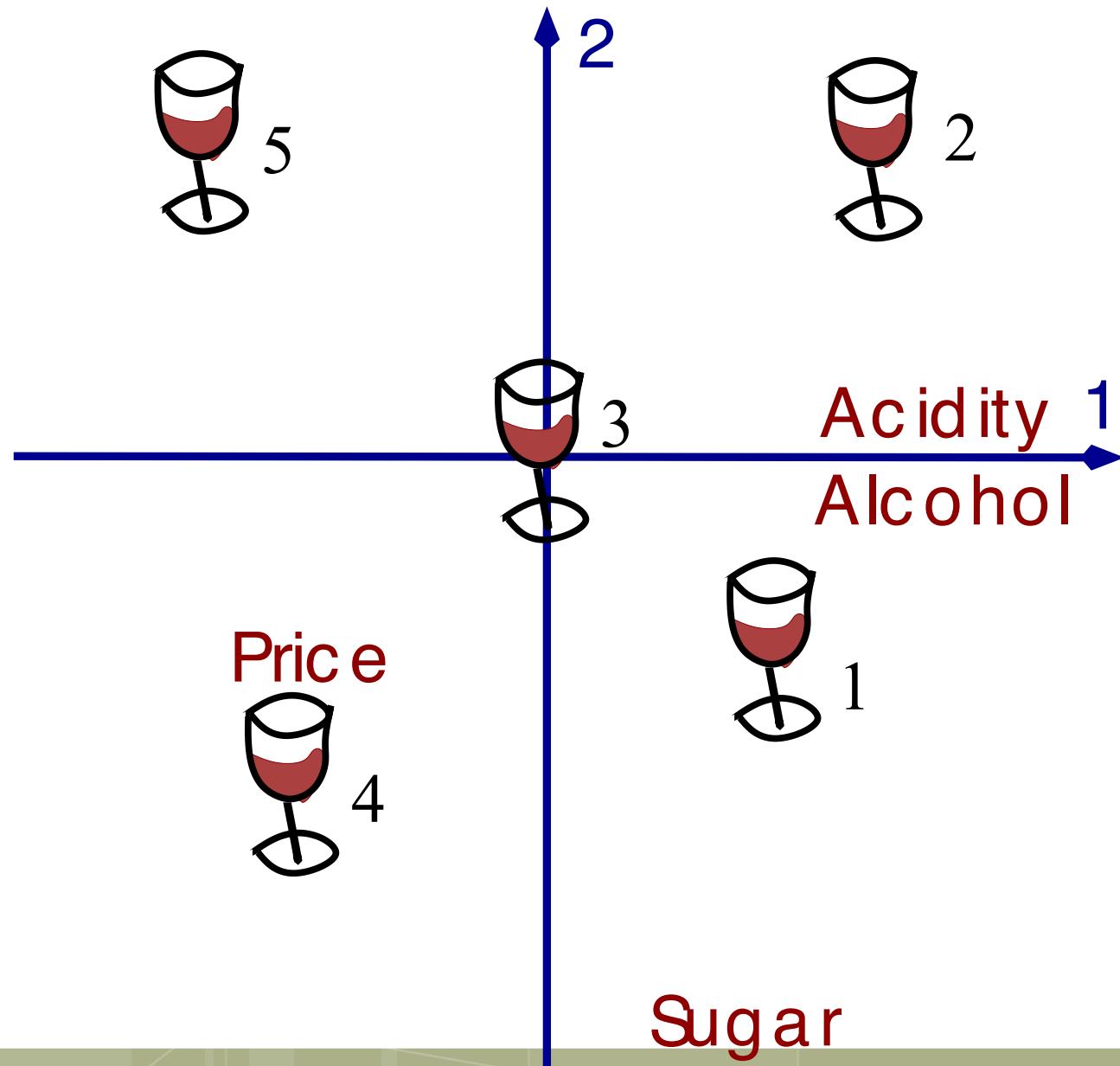
**MATRIX W****WEIGHTS FOR X VARIABLES**

Wine	w <sub>1</sub>	w <sub>2</sub>	w <sub>3</sub>
Price	-.51	-.34	-.35
Sugar	.20	-.94	.16
Alcohol	.57	-.02	-.82
Acidity	.61	.04	.42

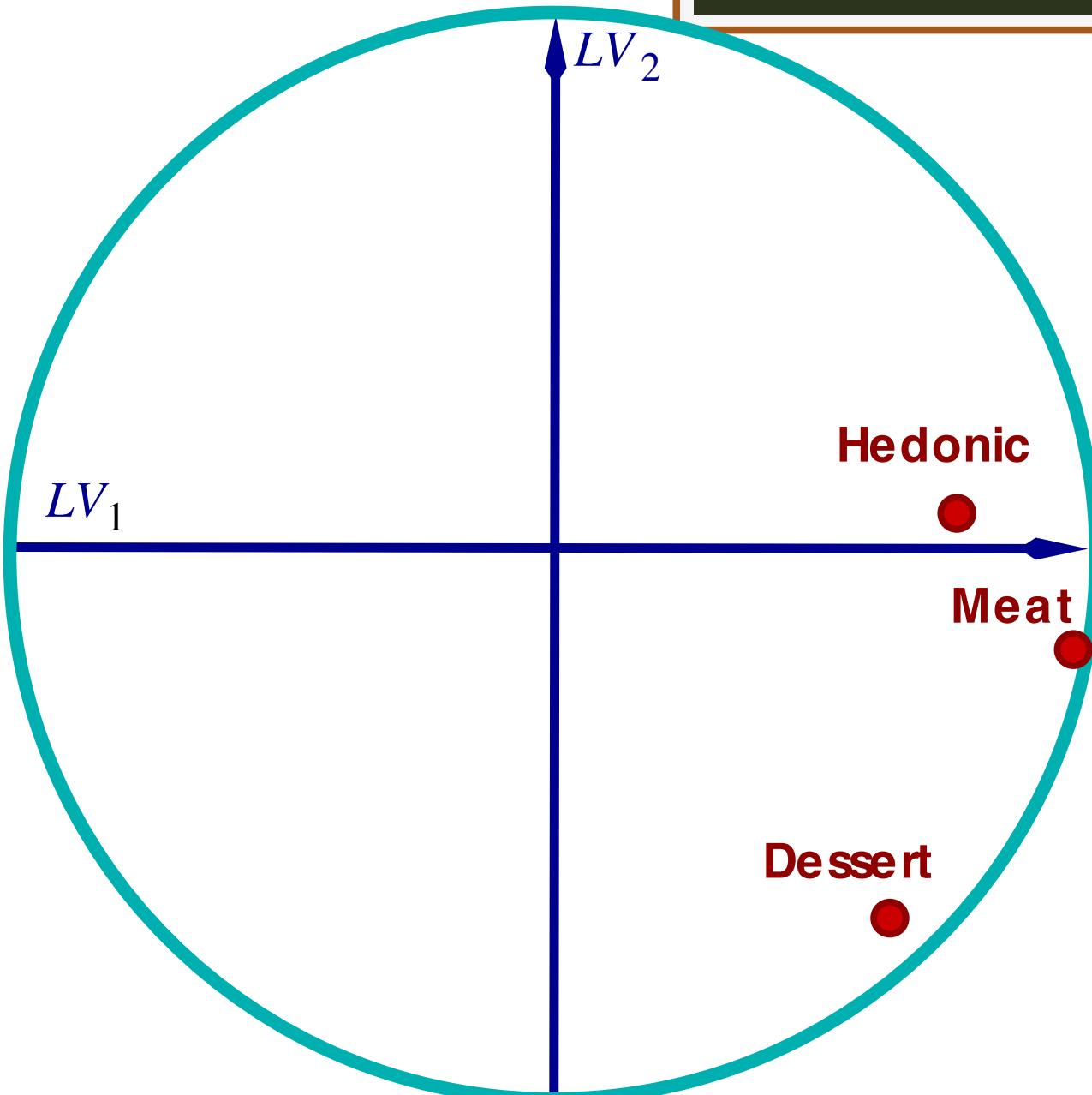
**MATRIX C****WEIGHTS FOR Y VARIABLES**

Wine	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>
Hedonic	.61	.05	.97
Meat	.70	-.27	-.22
Dessert	.37	-.96	-.13

# WHAT ABOUT SOME PICTURES

**PLOT T AND W BI PLOT**

## Y CORRELATION CIRCLE



# THE PREDICTION WITH 3 LATENT VARIABLES

		Y		
		Hedonic	Meat	Dessert
X	Intercept	48.50	-8.92	-3.85
	Price	-1.00	-0.03	0.04
	Sugar	0.75	0.28	0.59
	Alcohol	-4.00	1.00	0.50
	Acidity	2.75	0.18	0.09

A Wine costing \$10, with 5 g of sugar per litre, 10 degree of alcohol and an acidity of 5 (whatever) will have a hedonic rating of:

$$\hat{y}_1 = 48.50 - 1 \times 10 + .75 \times 5 - 4 \times 10 + 2.75 \times 5 = 16$$

# THE PREDICTION WITH 2 LATENT VARIABLES

		Y		
		Hedonic	Meat	Dessert
X	Intercept	-3.28	-3.38	-1.39
	Price	-0.26	-0.11	0.01
	Sugar	0.14	0.34	0.62
	Alcohol	-0.81	0.49	0.27
	Acidity	0.69	0.40	0.19

A Wine costing \$10, with 5 g of sugar per litre, 10 degree of alcohol and an acidity of 5 (whatever) will have a hedonic rating of:

$$\hat{y}_1 = -2.28 - .26 \times 10 + .14 \times 5 - .81 \times 10 + 0.69 \times 5 = 6.38$$

## HOW MANY

**BIG QUESTION:**

**HOW MANY LATENT VARIABLES?**

**WHAT'S GIVE THE BEST PREDICTION?**

**RISK OF OVER-FITTING**

**So 3 SETS. LEARNING, TUNING, TESTING**

MORE

## ANOTHER EXAMPLE: CRUNCH CRUCNH

**OBJECTIVE:** To study the relationships between the 3 texture attributes: *croustillant* (crispy), *craquant* (crackly), *croquant* (crunchy) and 9 attributes describing the biting sounds and other textural characteristics.

# Descriptive Analysis

Panel: 8 trained assessors

## Attributes

crispy  
crunchy  
crackly

+  
hard  
crumbly  
puffy  
brittle  
sticky  
duration  
loudness  
high-pitched  
low-pitched

eating noises {

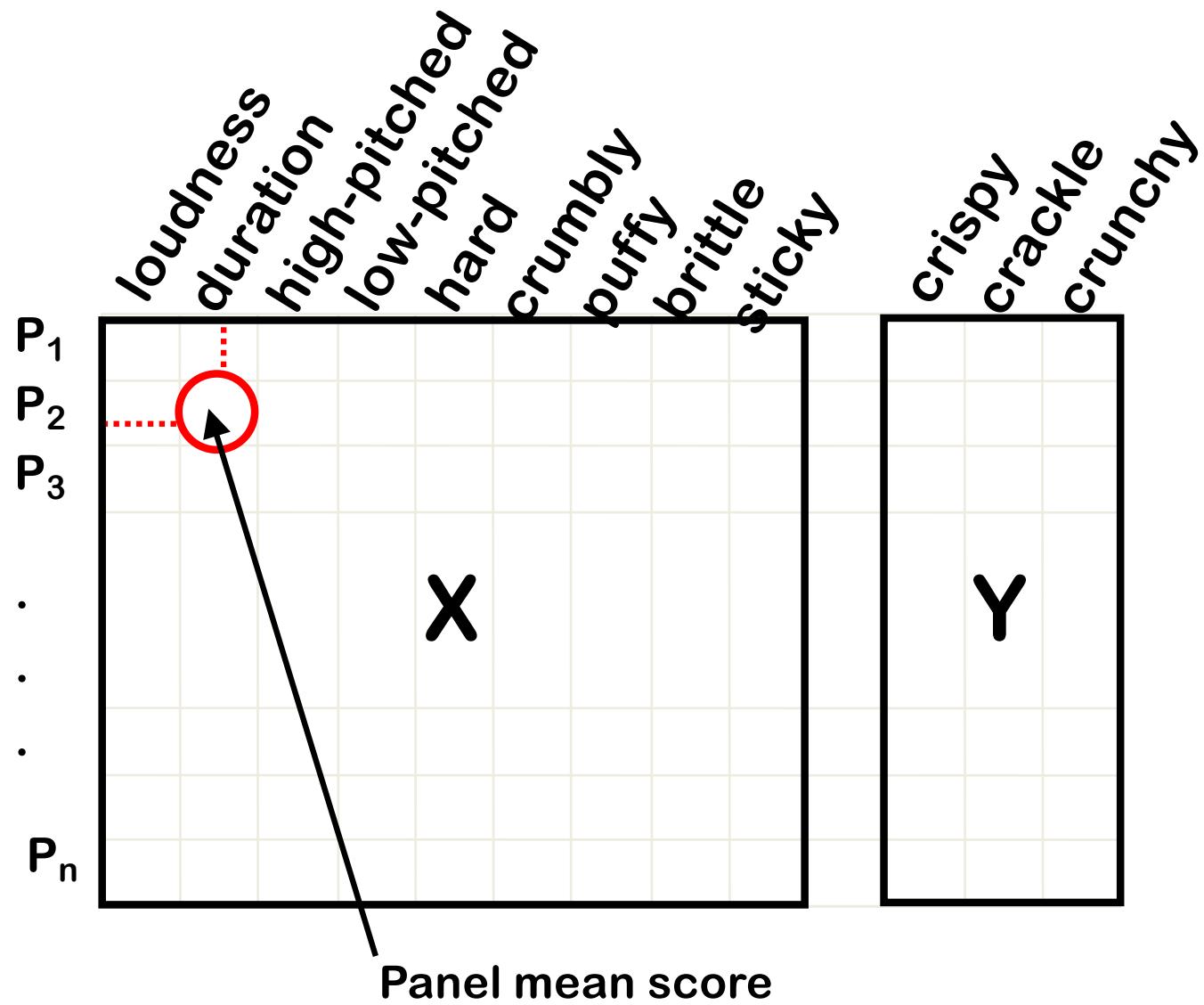
## 35 Products

*extruded bread & snacks  
dry biscuits, cold cereals  
chips, biscotte, chocolate  
raw vegetables & salad  
dry fruits & seeds*

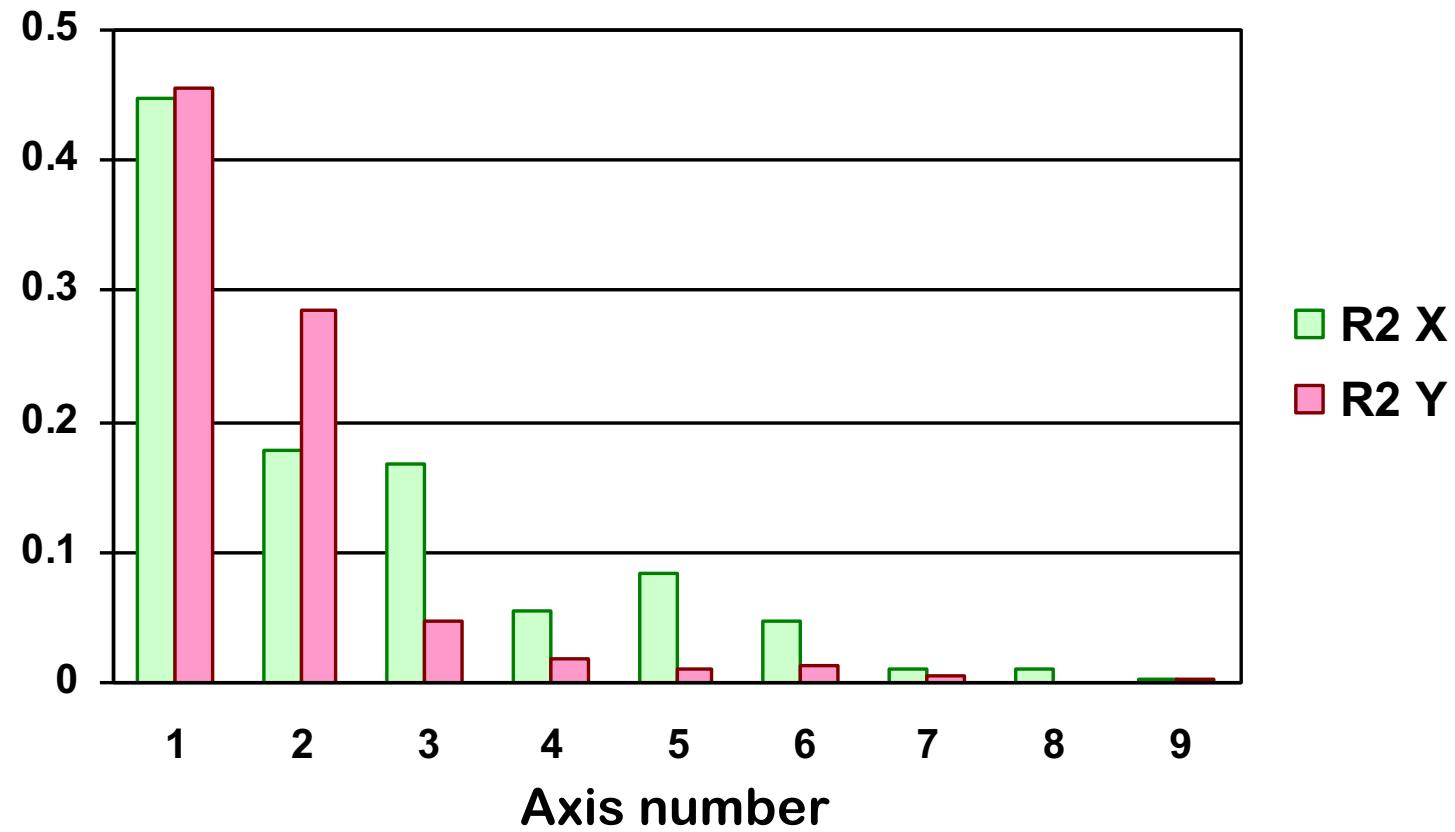
## Evaluation

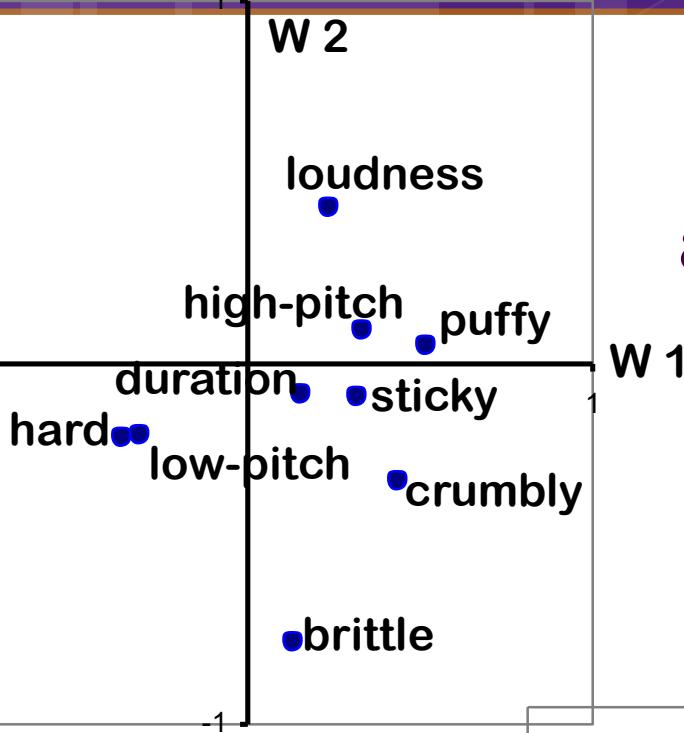
- ♦ in duplicate
- ♦ linear scale
- ♦ scores : 0 → 10

# Data

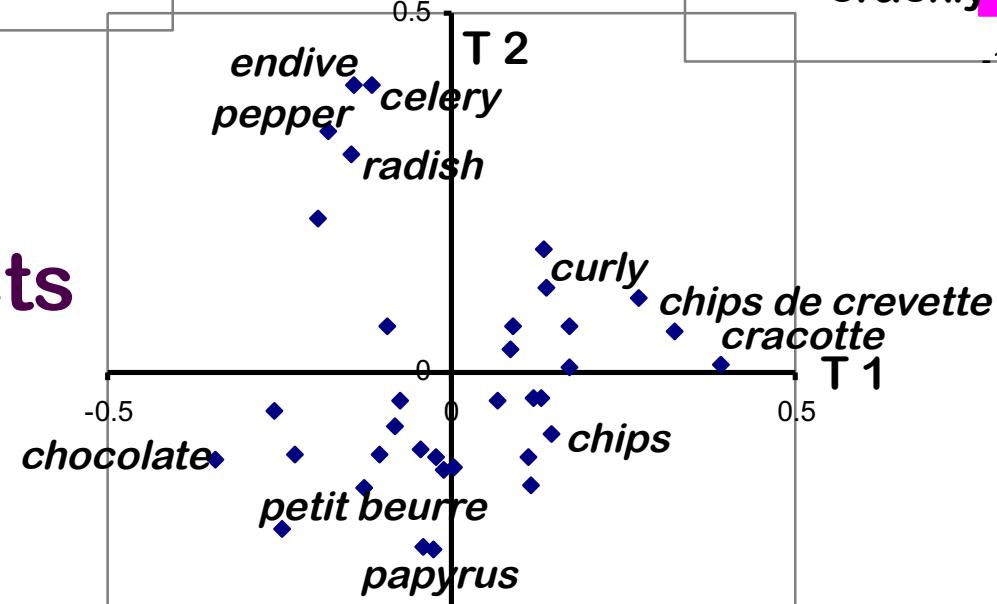
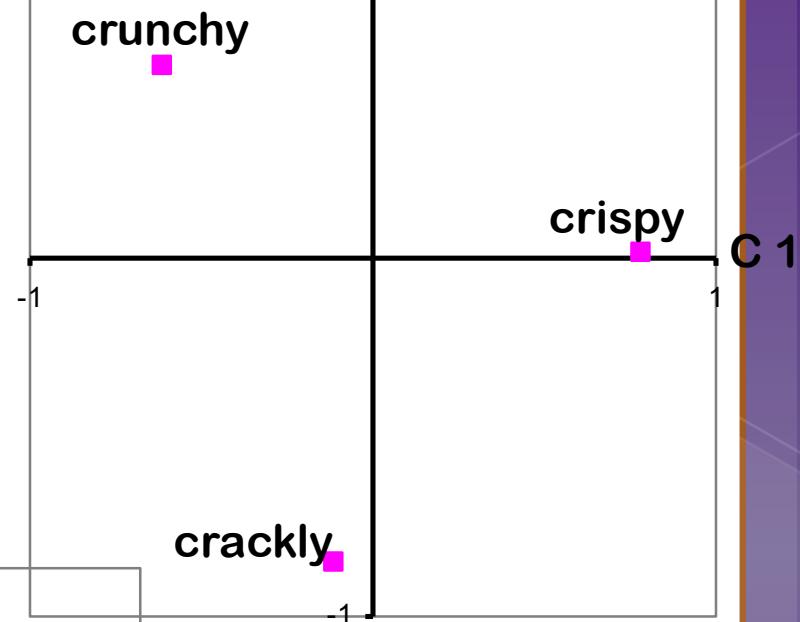


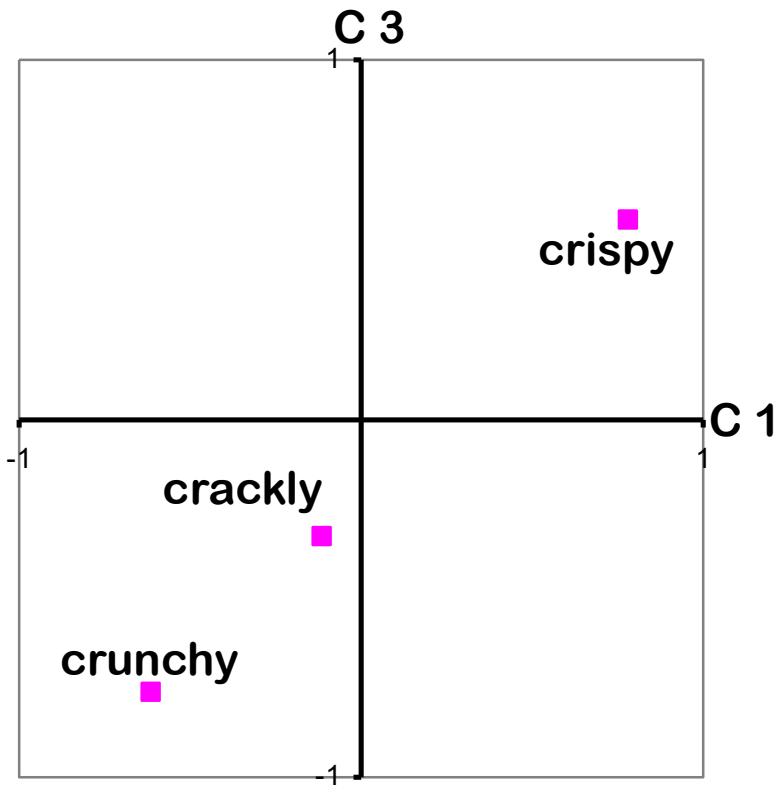
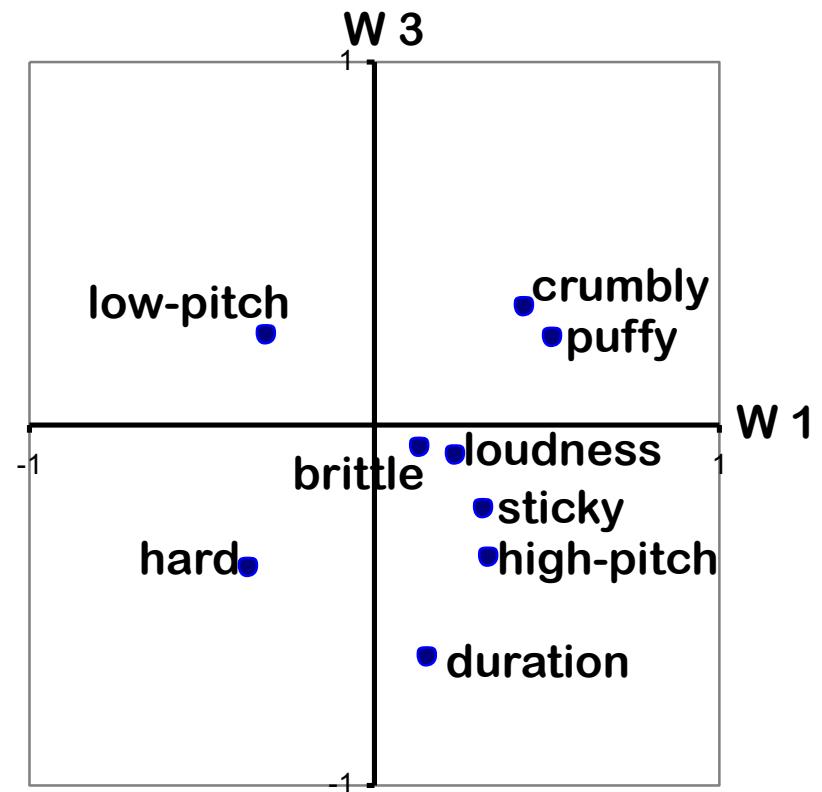
# Error

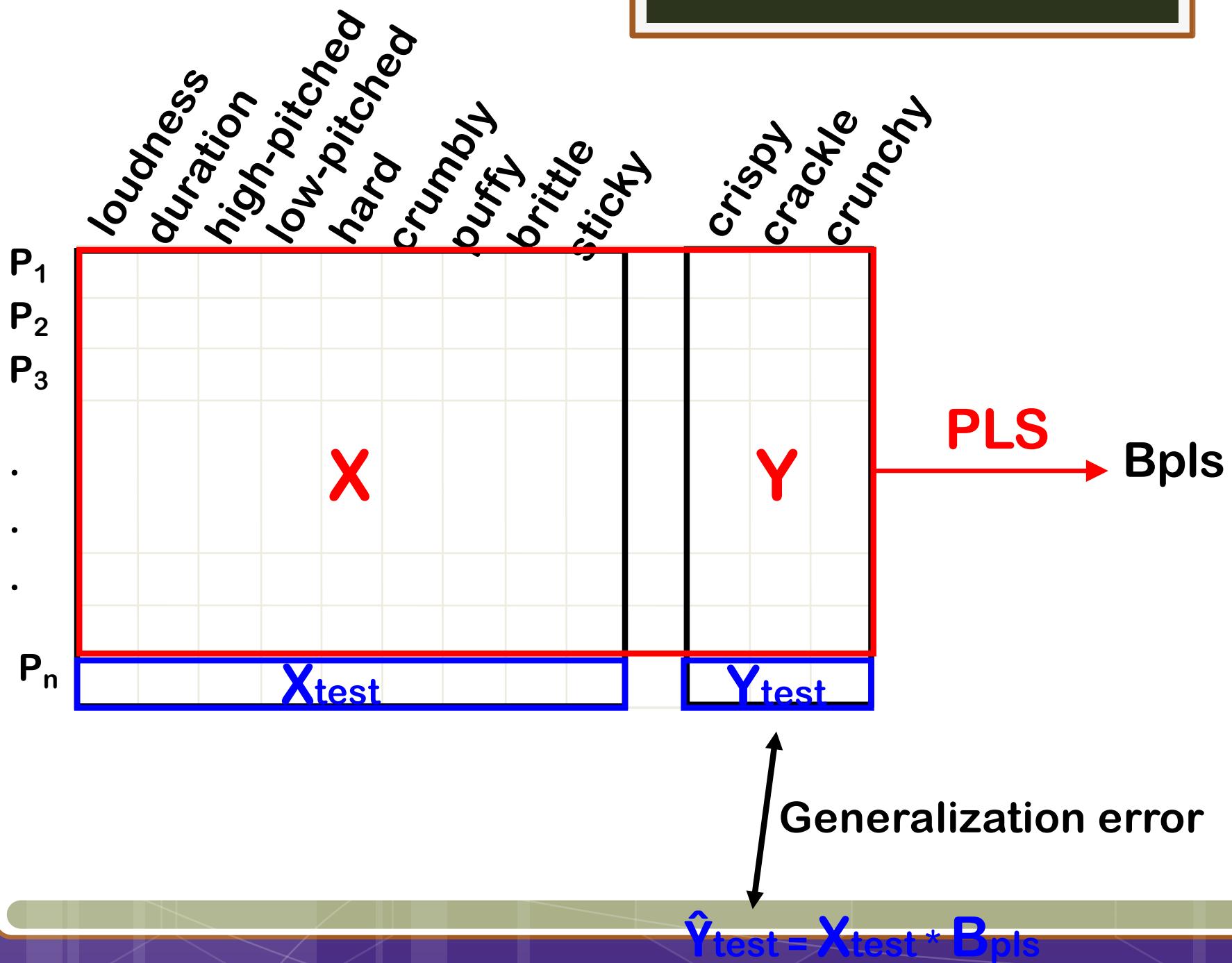




## attributes

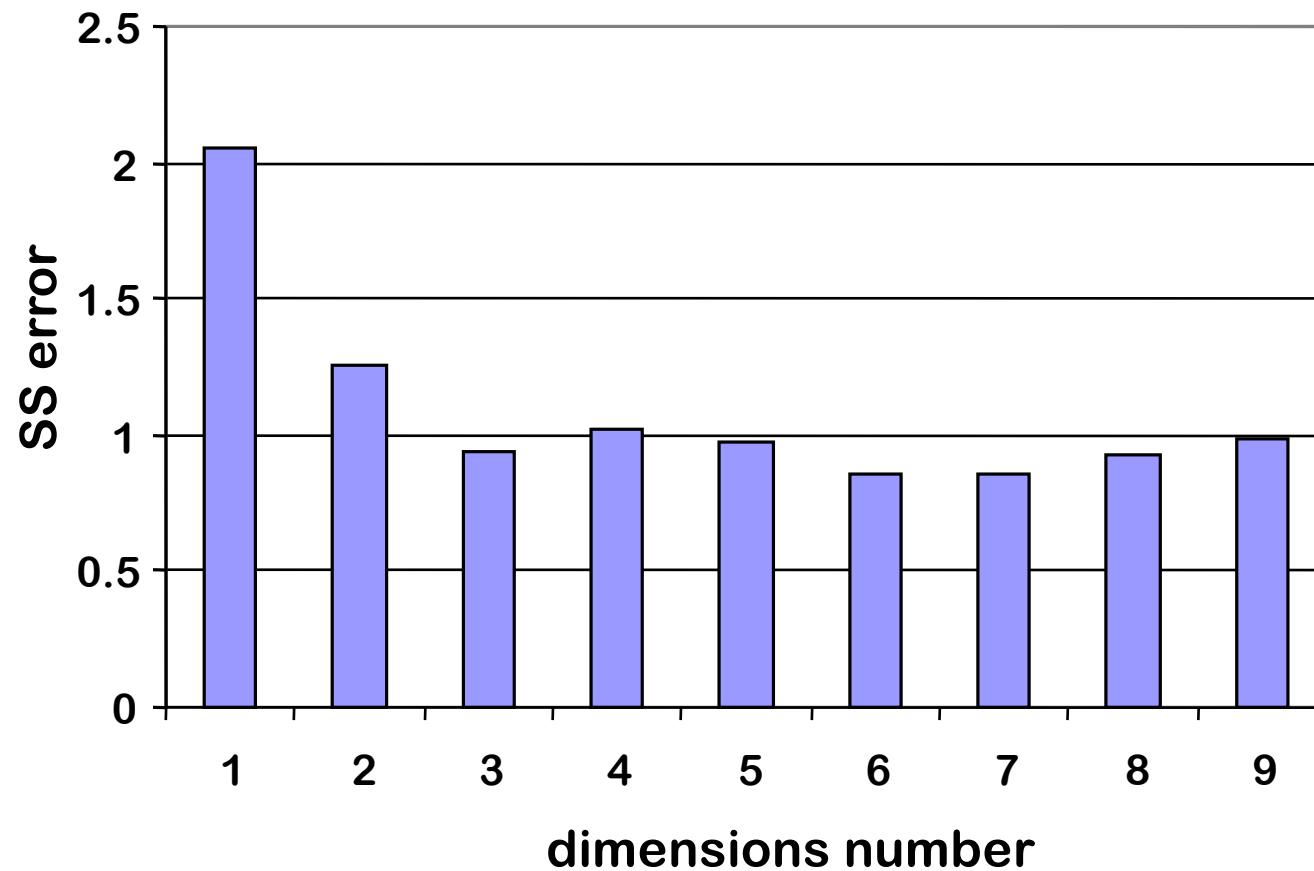




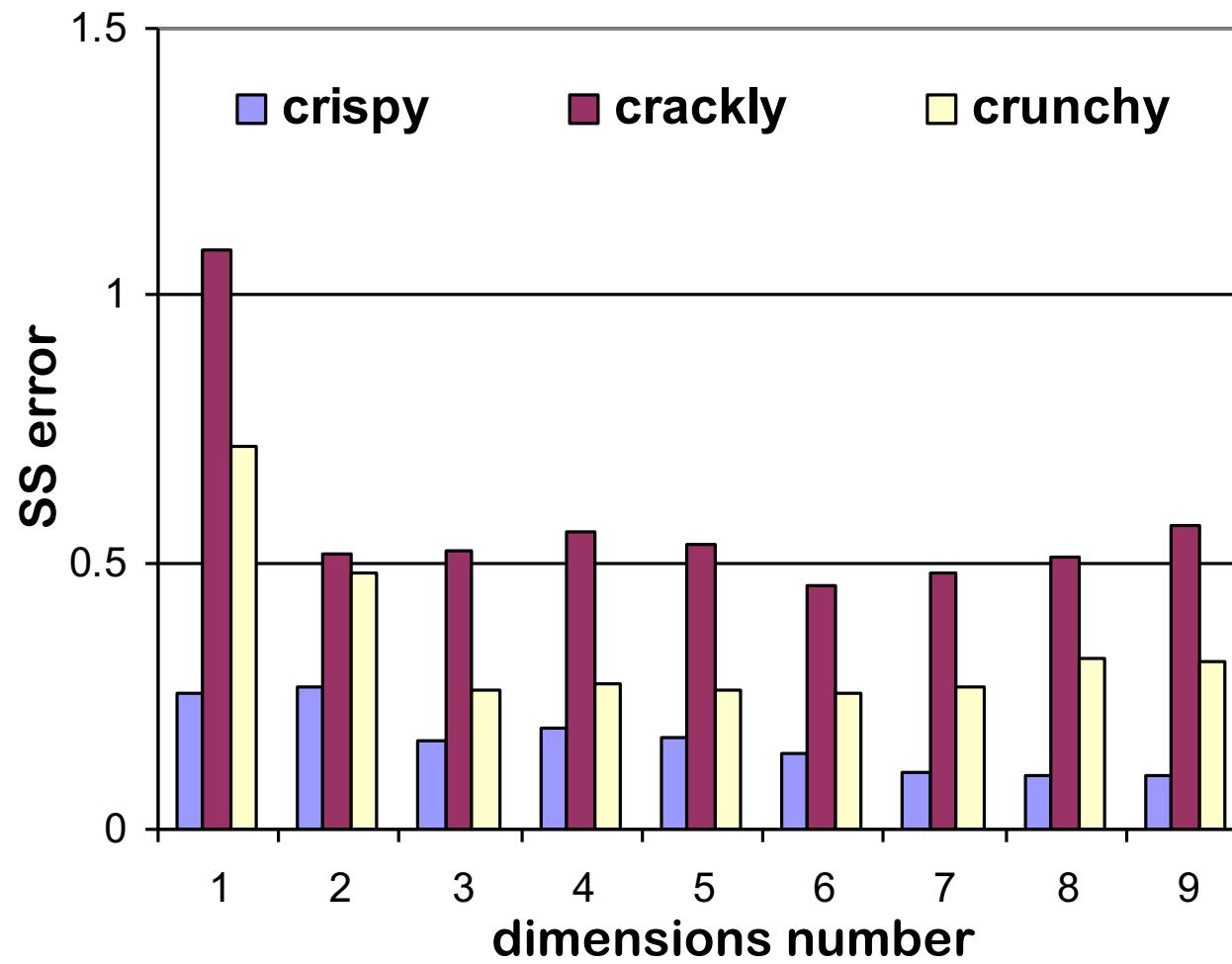


# Generalisation Error

Learning set: 34 products ; test set: 1 product  
35 runs



# Generalisation Error



# DEVELOPING PLS

- ➔ **PLS for Distances**
  - ➔ DISPLS Correlation
  - ➔ DISPLS Regression
- ➔ **PLS for Qualitative Data**
  - ➔ PLS Correspondence Analysis
- ➔ **PLS for multiple data tables**
  - ➔ Multi-Block PLS, DISPLS



**THANK YOU**