

# Exploring the relationship between sensory and instrumental data with component-based methods



**KoSfoST International Symposium  
and Annual Meeting 2025**

*Pioneering Future Connection in FoodTech*  
Gwangju, Korea · 2-4 July 2025

**John Castura**

Dr. Philos., M.Sc.  
Research Fellow

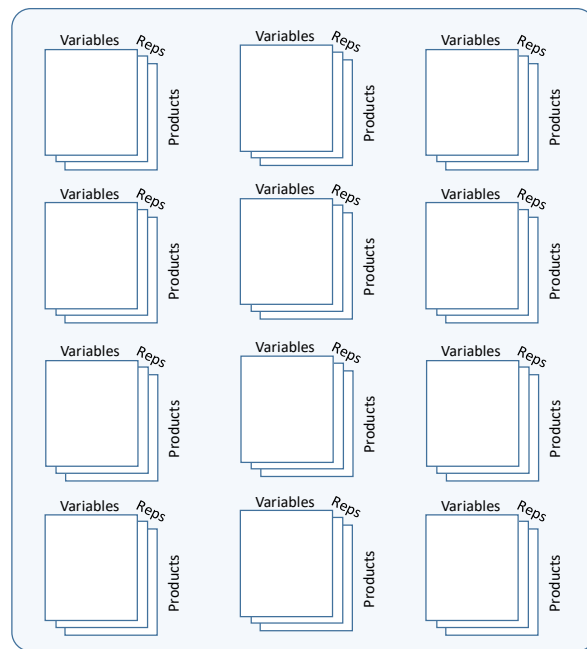
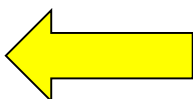
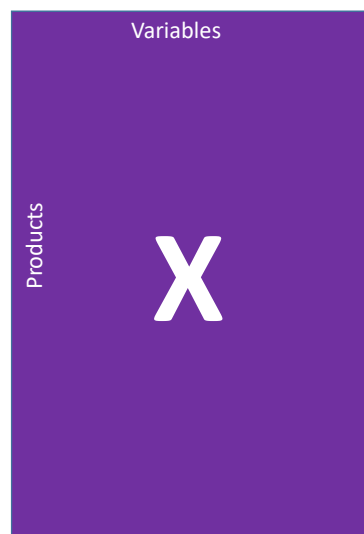


**Compusense®**

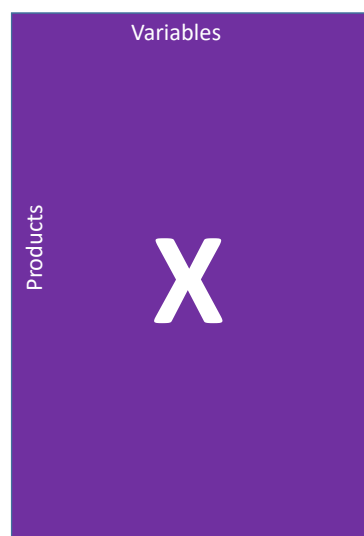




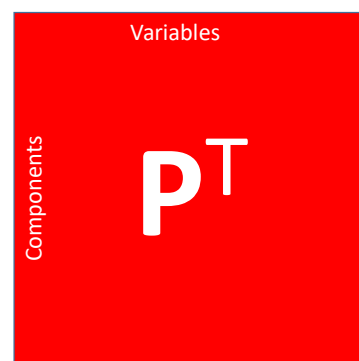
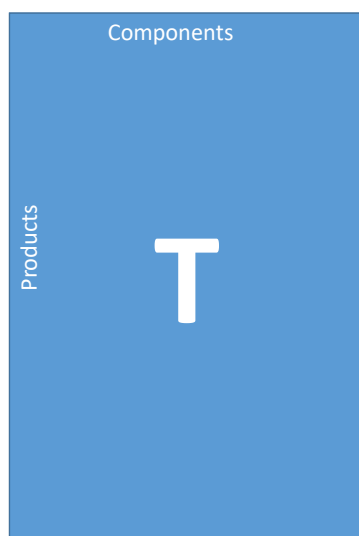
## Aggregate into data matrix



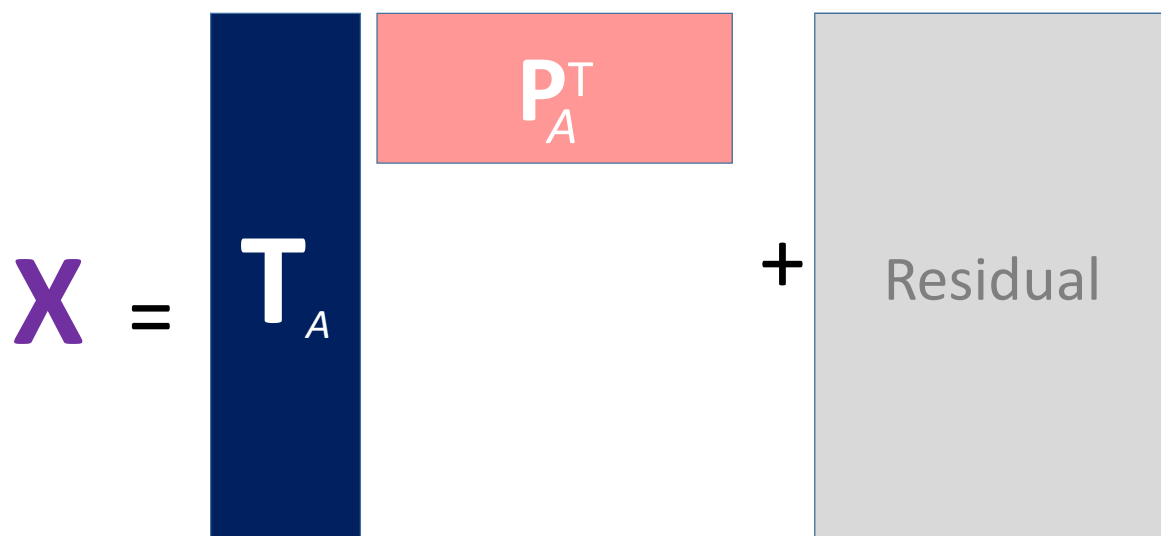
## Principal component analysis



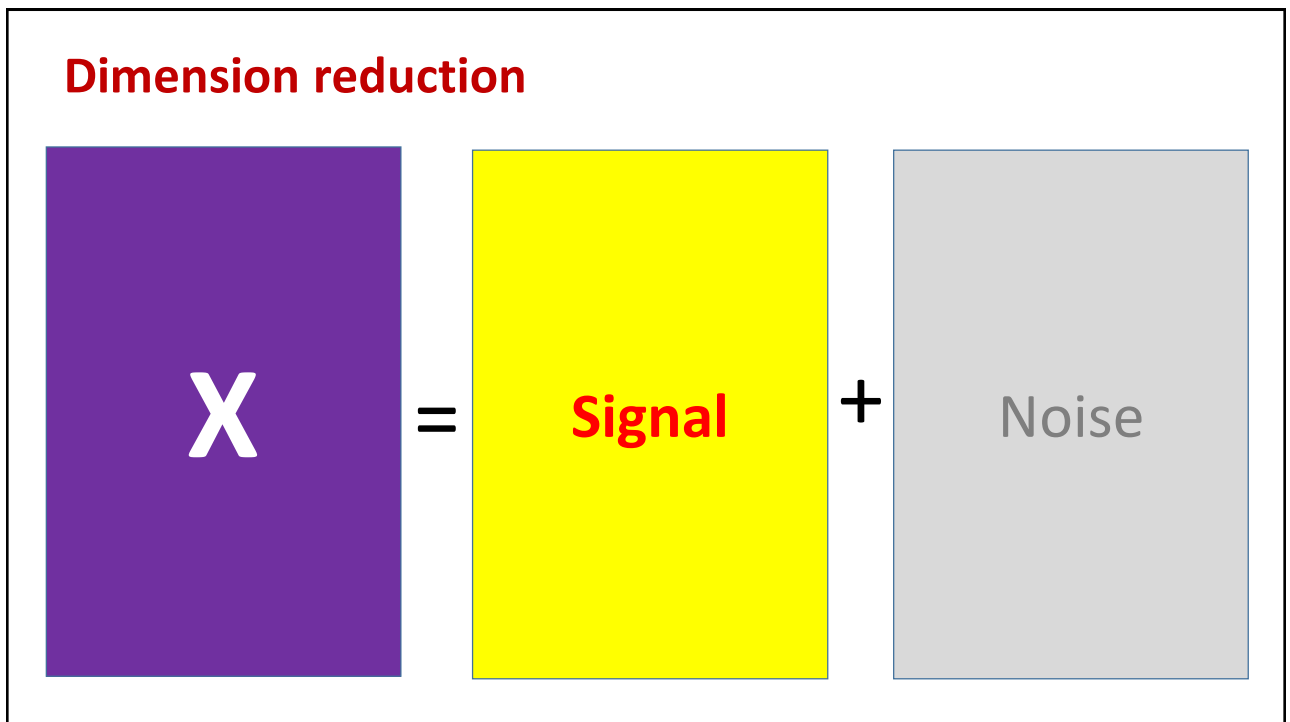
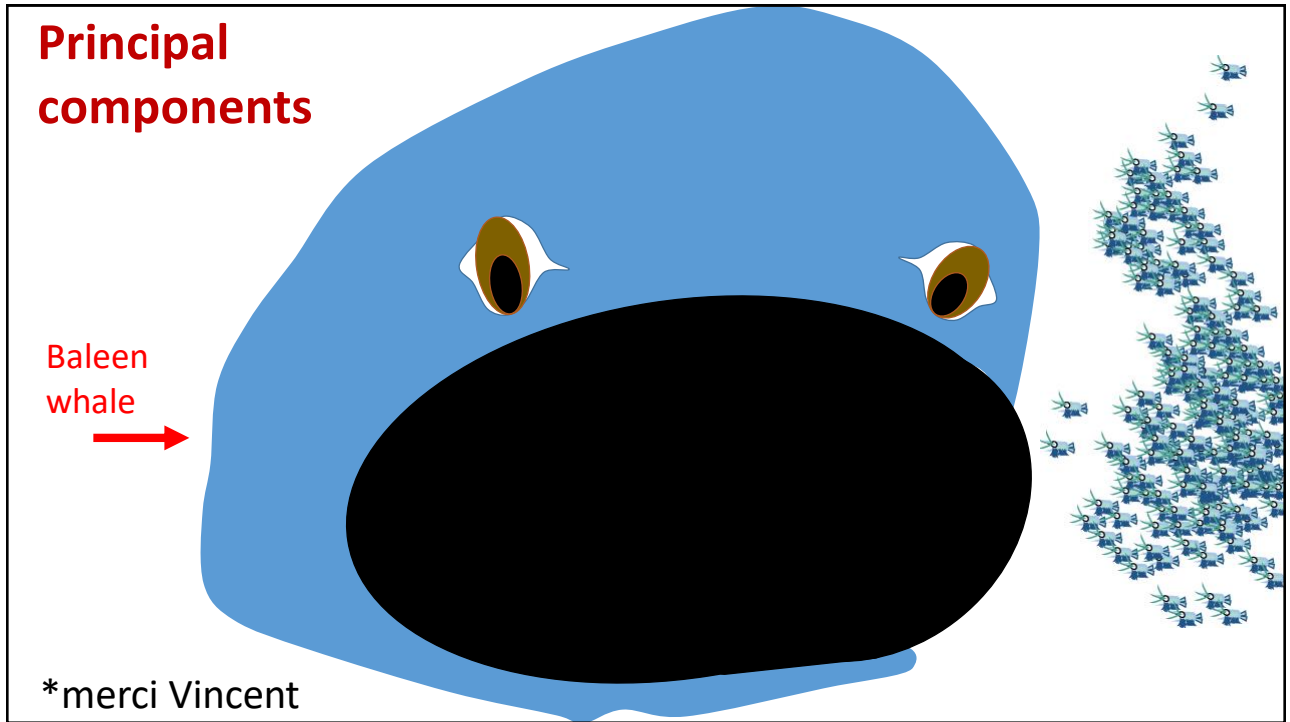
=



## Dimension reduction to $A$ components

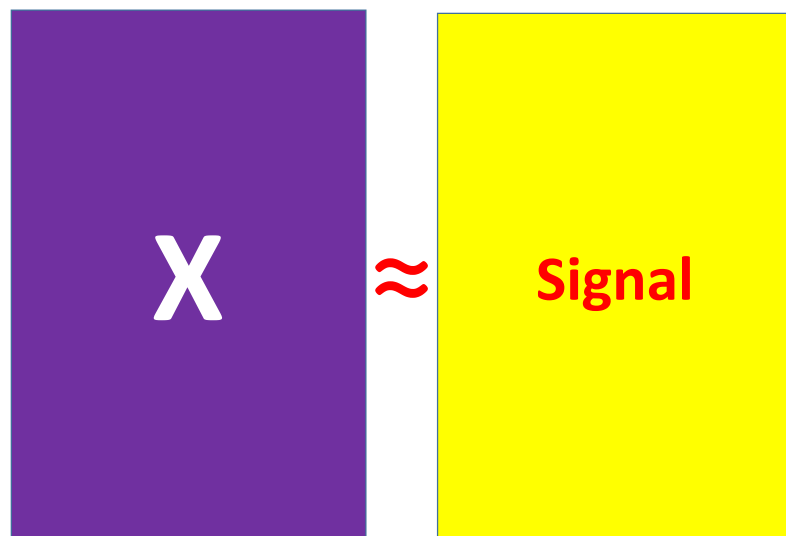
$$\mathbf{X} = \mathbf{T}_A \mathbf{P}_A^T + \text{Residual}$$




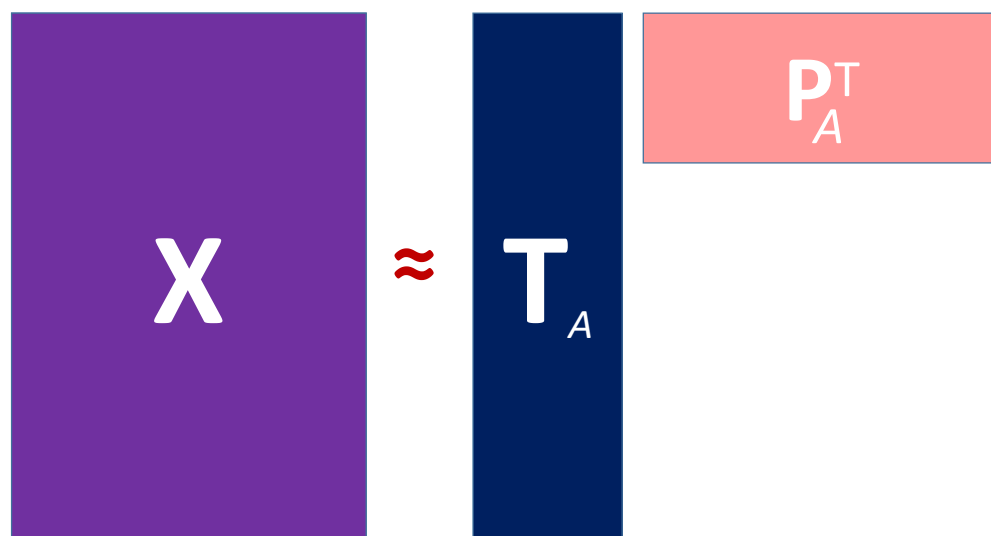




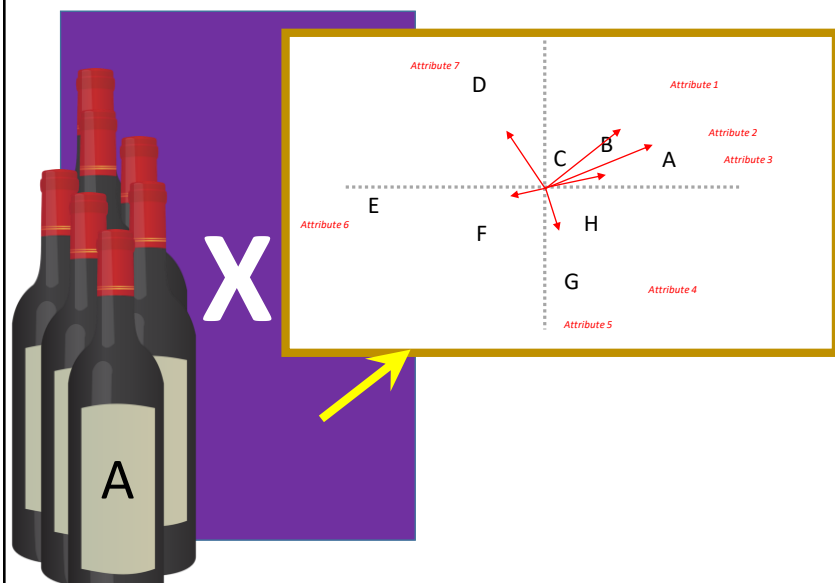
## Dimension reduction



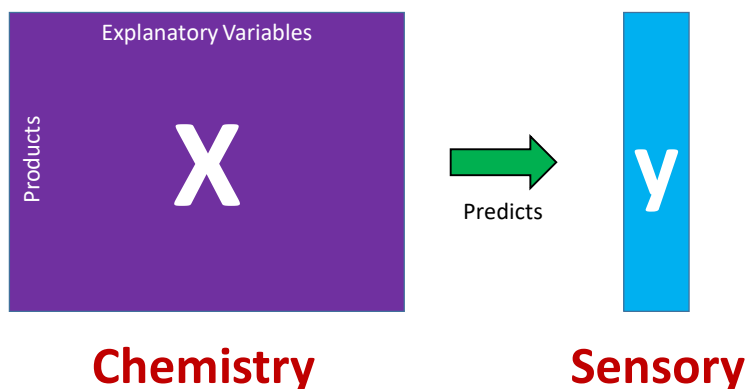
## Dimension reduction to $A$ components



## Visualizing PCA results



## Investigating data relationships



## Why component-based methods?

- Summarize and visualize complicated data
- Often relatively few underlying sources of variation
- Useful approximation of data
- Separation of signal/noise
- Outlier detection
- Confirm hypotheses
- Generate hypotheses

Næs, Varela, Castura, Bro & Tomic (2023). Why use component-based methods in sensory science? *Food Quality and Preference*, 112, 105028. <https://doi.org/10.1016/j.foodqual.2023.105028>

## Data for French Pinot Noir Wines



## Volatile organic compounds (VOCs)

Headspace measurements of VOCs obtained from  
headspace—solid phase micro-extraction—gas chromatography—mass spectrometry (HS-SPME-GC-MS)

**X**

1-hexanol	acetaldehyde	ethyl acetate	isoamyl acetate
1-octanol	acetic acid	ethyl butyrate	isoamyl propionate
1-phenoxy-2-propanol	alpha-ionone	ethyl caproate	isovaleric acid
2,3-butanedione	beta-ionone	ethyl isobutyrate	methional
2-ethylhexan-1-ol	butyl acetate	ethyl isovalerate	methionol
2-isobutyl-3-methoxypyrazine	butyric acid	ethyl lactate	pentyl propionate
2-methyl-1-butanol	damascenone	ethyl octanoate	phenol
2-methylbutyl acetate	dimethyl sulfide	ethyl propionate	phenylacetaldehyde
2-phenylethanol	ethyl 2-methylbutyrate	furaneol	phenylacetic acid
3-methyl-1-butanol	ethyl 3-hydroxybutyrate	hexyl acetate	propionic acid
4-ethyl-2-methoxyphenol	ethyl 6-hydroxyhexanoate	homofuraneol	trans-3-hexen-1-ol
4-ethylphenol			

Villière et al. (2019)

## Pinot noir wines

**T1** 2010 Bourgogne PDO

**T5** 2009 Savigny-lès-Beaune PDO

**T2** 2009 Bourgogne PDO

**T6** 2010 Maranges PDO

**T3** 2009 Bourgogne PDO

**C** 2009 Côte de Nuits-villages PDO

**T4** 2009 Bourgogne Hautes Côtes de Beaune PDO

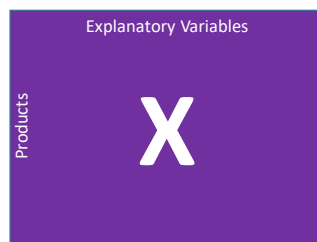
**T7** 2009 Ladoix PDO

Villière et al., 2019, <https://doi.org/10.1016/j.dib.2019.103725>

## Sensory Variables

y

Artichoke	Cherry fresh	Hay	Smoky
Bell pepper	Cherry stone	Leather	Strawberry cooked
Blackberry fresh	Clove	Musk	Strawberry fresh
Blackcurrant bud	Cut grass	Pepper	Toasty
Blackcurrant fresh	Elderflower	Plum cooked	Undergrowth
Blueberry fresh	Ethanol	Plum fresh	Vanilla
Brioche	Firestone	Prune	Violet
Butter	Geranium	Raspberry fresh	Woody
Cherry cooked			



## Principal component regression

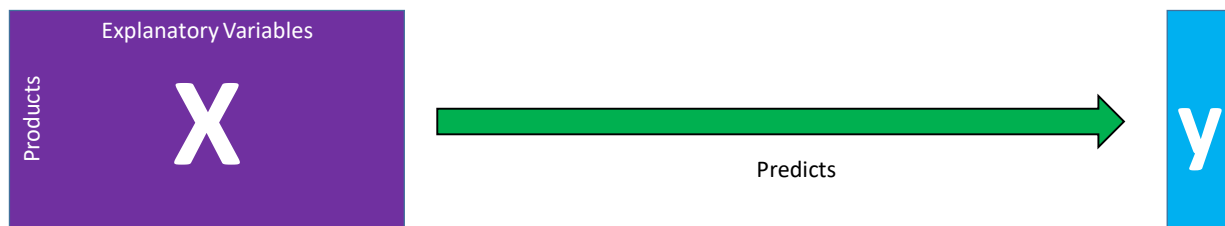
Unsupervised & Supervised

We want to ***predict*** the response **y** from multivariate **X**.

Response **y** is regressed on **principal components** of **X**.

## Investigating data relationships

with **multiple linear regression**



## Investigating data relationships

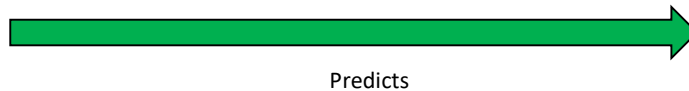
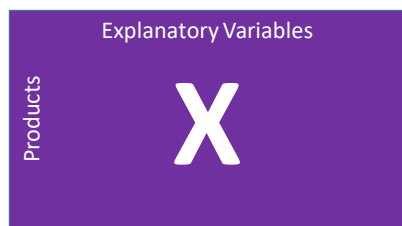
with multiple linear regression

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{f}$$

Regression coefficients

## Investigating data relationships

with multiple linear regression

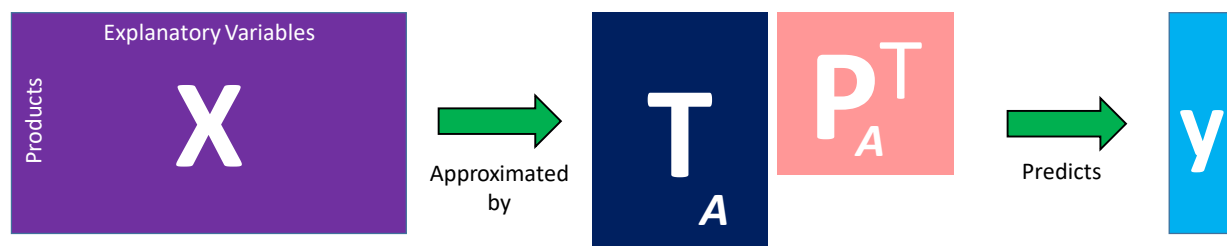


$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{f}$$
$$\mathbf{b} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

Multicollinearity?  
More variables than objects?

## Investigating data relationships

with principal component regression (PCR)



## Investigating data relationships

with principal component regression (PCR)

The diagram shows the mathematical model for PCR:

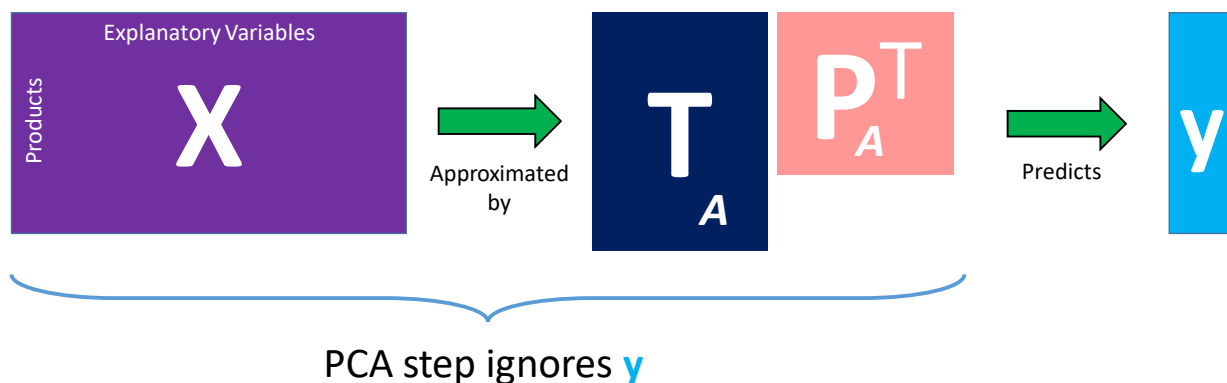
$$y = Xb + f$$

Where:

- $y$  is the response variable (blue box).
- $X$  is the matrix of explanatory variables (purple box).
- $b$  is the vector of **Regression coefficients from PCR** (red box, indicated by a red arrow).
- $f$  is the error term (gray box).

## Investigating data relationships

with **principal component regression (PCR)**

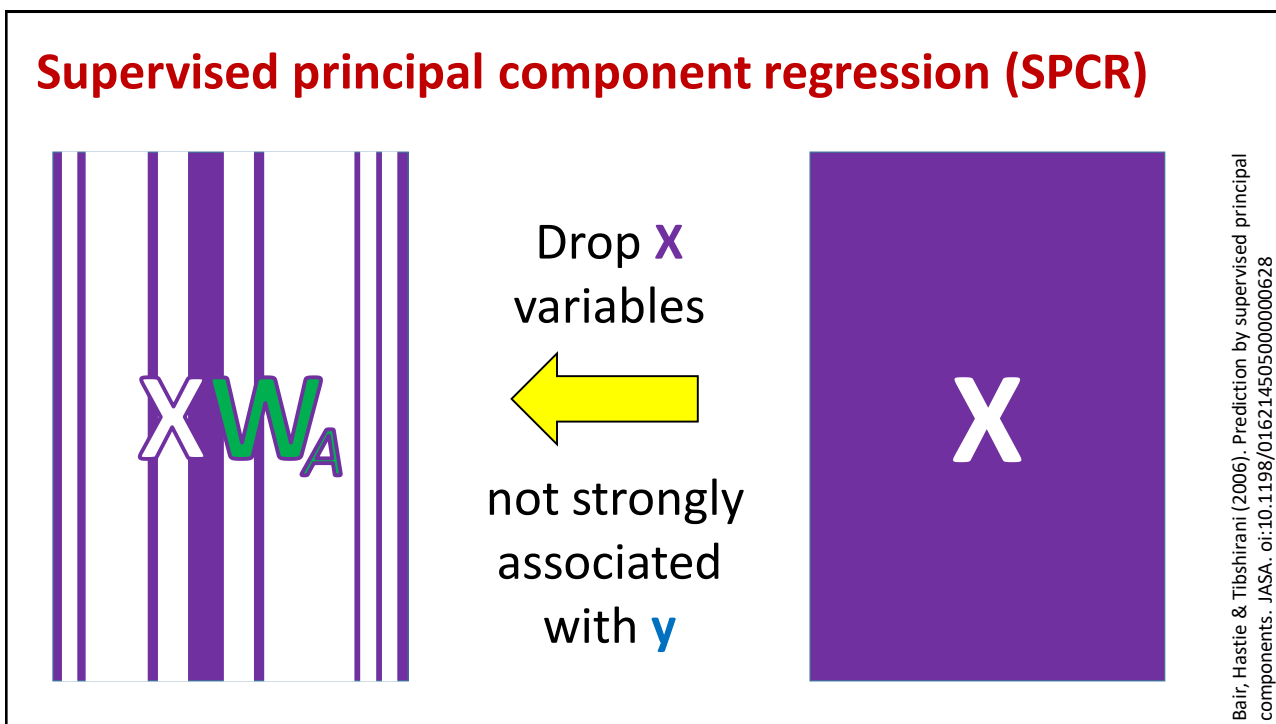
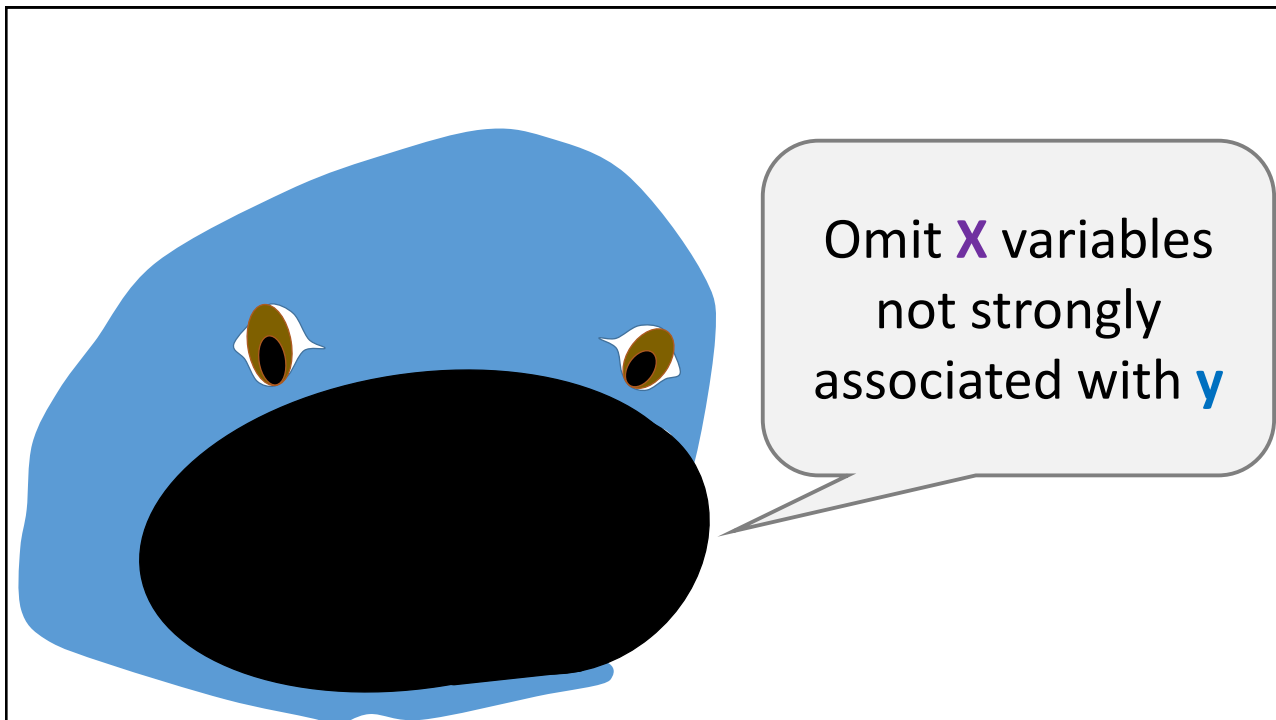


## Supervised principal component regression (SPCR)

We want to ***predict*** the response  $y$  from the multivariate  $X$ .

Response  $y$  is regressed on principal components of  $X$ .

**SUPERVISED**



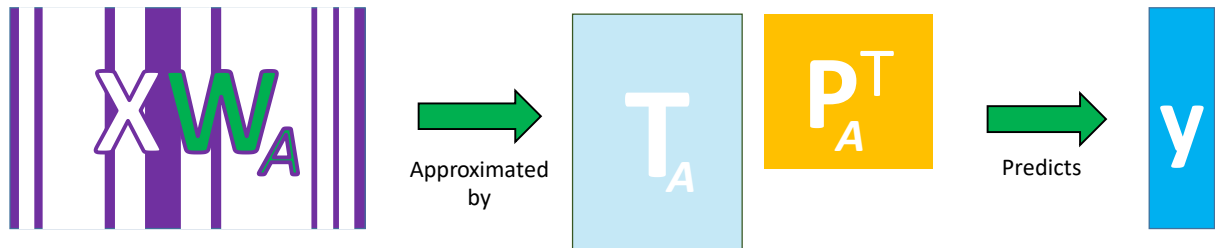


## Dimension reduction in SPCR



Bair, Hastie & Tibshirani (2006)

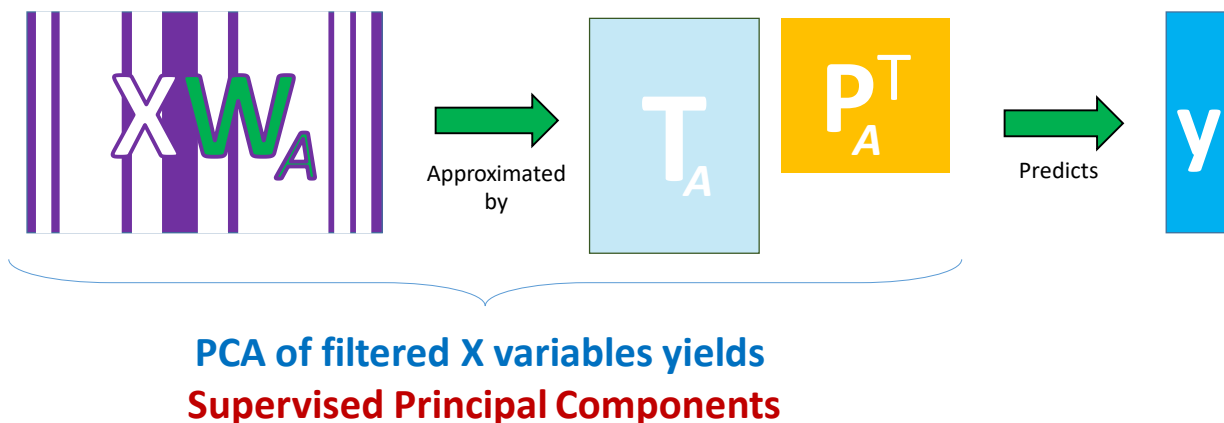
## Supervised principal component regression (SPCR)



PCA of filtered  $X$  variables yields  
Supervised Principal Components

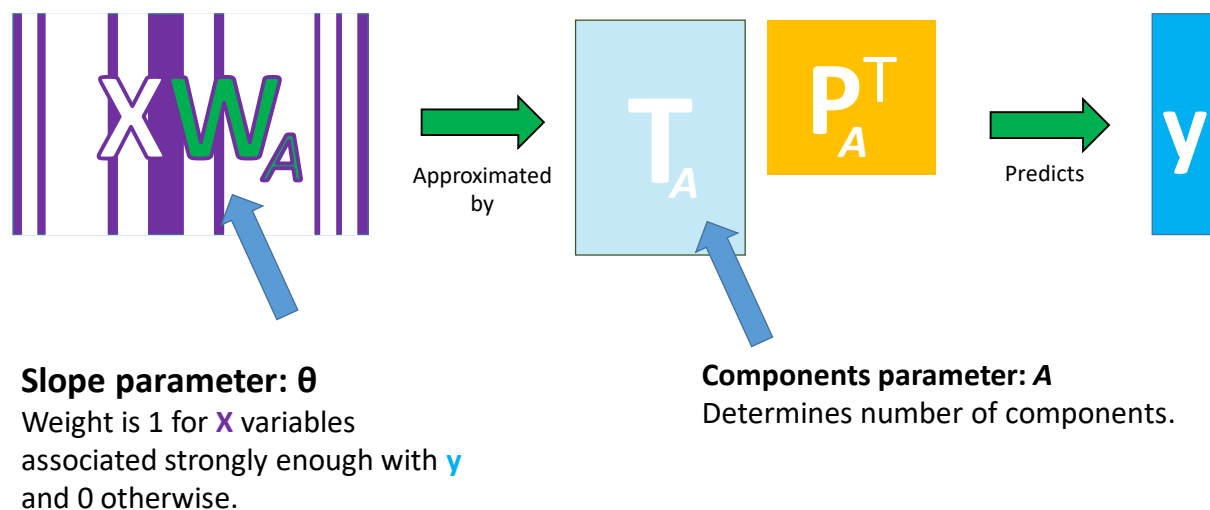
Bair, Hastie & Tibshirani (2006)

## Supervised principal component regression (SPCR)



Bair, Hastie & Tibshirani (2006)

## Supervised principal component regression (SPCR)



Bair, Hastie & Tibshirani (2006)

Many candidate values for parameters  $\theta$  and  $A$

### *How to choose?*

**Strategy:** Choose model that gives the most accuracy predictions for data withheld from model training

*Bair et al.:* 2-fold cross-validation (2-fold cv)  
*Us:* leave-one-out cross-validation (loocv)

### Supervised principal component regression (SPCR)

$$y = Xb + f$$

Regression coefficients

Bair, Hastie & Tibshirani (2006)

Usually, when interpret results, we draw conclusions about...

**relationships between variables,**

**relationships between objects, and**

**differences between objects.**

Usually, when interpret results, we draw conclusions about...

**relationships between variables,**

**relationships between objects, and**

**differences between objects.**

# Investigating paired comparisons after principal component analysis

Castura, J.C., Varela, P, & Næs, T. (2023).  
Investigating paired comparisons after principal component  
analysis. *Food Quality and Preference*, 106, 104814.  
<https://doi.org/10.1016/j.foodqual.2023.104814>

## Investigating paired comparisons after PCA

### “Crossdiff-unfolding”

$X$

$X$  is a column-centered ( $J \times M$ ) matrix

Every row is subtracted from  
every row

$X \ominus X$

$X \ominus X$  is a column-centered ( $J^2 \times M$ ) matrix

## Investigating paired comparisons after PCA

### “Crossdiff-unfolding”

$X$

The covariance matrix of  $X$  and the covariance matrix of  $X \ominus X$  are identical except for a multiplier.

$X \ominus X$

Next, consider PCA of  $X$  and PCA of  $X \ominus X$ .

## Investigating paired comparisons after PCA

PCA of  $X$

$$X = \text{[blue rectangle]} P^T$$

PCA of  $X \ominus X$

$$X \ominus X = \text{[blue rectangle]} P^T$$

## Investigating paired comparisons after PCA

PCA of  $X$

$$X = \text{[blue bar]} \circledast P^T$$

*Key result #1:*

Loading matrices obtained from these two PCA solutions are identical.

PCA of  $X \ominus X$

$$X \ominus X = \text{[blue bar]} \circledast P^T$$

## Investigating paired comparisons after PCA

PCA of  $X$

$$X = T \circledast P^T$$

*Key result #2:*

If we crossdiff-unfold scores from PCA of  $X$ , we get scores from PCA of  $X \ominus X$ .

PCA of  $X \ominus X$

$$X \ominus X = T \ominus T \circledast P^T$$



## Paired comparisons

This demonstrates that objects and all their paired comparisons are optimally investigated **in the same principal components**.

## Paired comparisons

Therefore, we can just do PCA of  $\mathbf{X}$  and know the PCA of  $\mathbf{X} \ominus \mathbf{X}$  without actually doing this PCA.

This lays necessary theoretical groundwork to justify a strategy for doing paired comparisons after PCA.

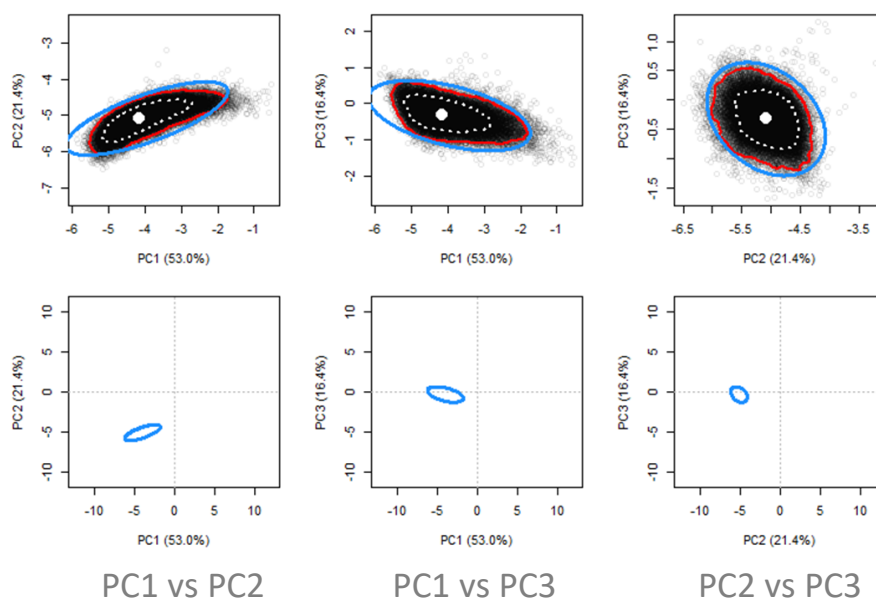
We need to **account for mutual dependencies in the data** when investigating paired comparisons.

## Evaluation of complementary numerical and visual approaches for investigating pairwise comparisons after principal component analysis

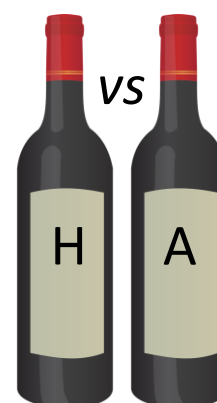
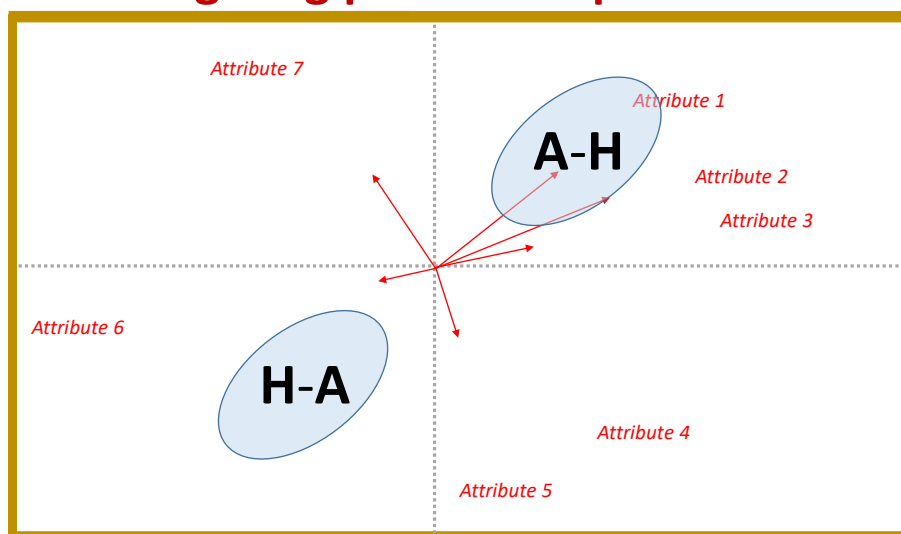
Castura, J.C., Varela, P. & Næs, T. (2023).  
Evaluation of complementary numerical and visual approaches for investigating pairwise comparisons after principal component analysis. *Food Quality and Preference*, 107, 104843. <https://doi.org/10.1016/j.foodqual.2023.104843>

## ...complementary numerical and visual approaches...

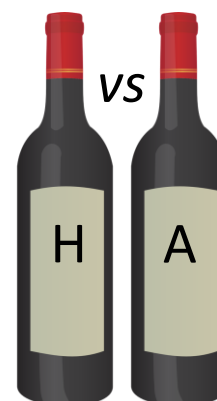
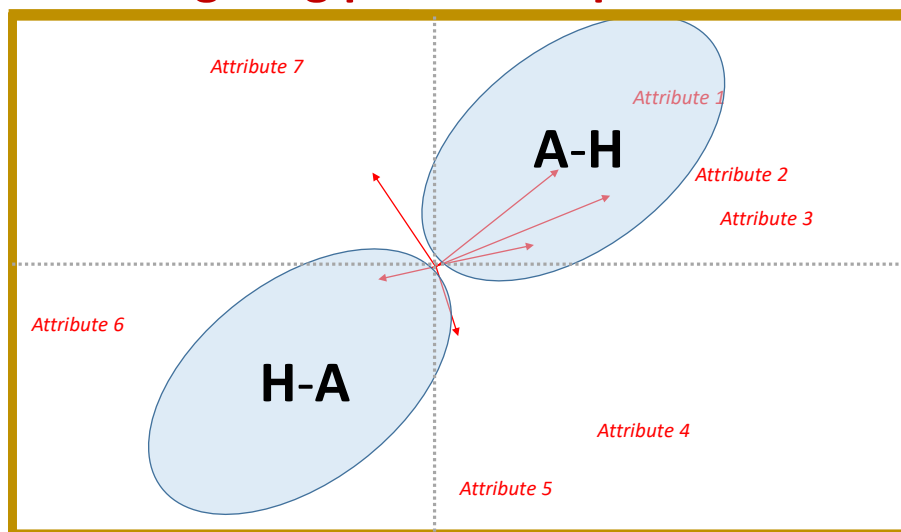
Suppl. Fig. 1a from Castura, Varela & Naes (2023a)



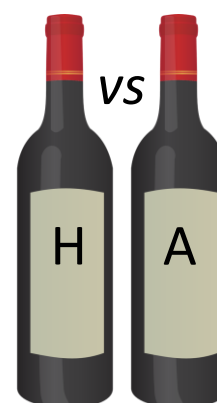
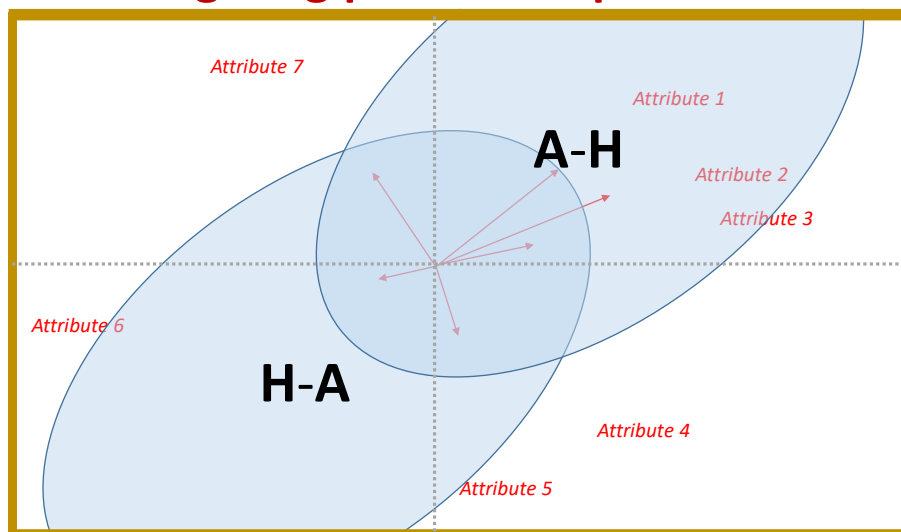
## Investigating paired comparisons



## Investigating paired comparisons



## Investigating paired comparisons



In some studies, only a subset of  
paired comparisons is of primary  
interest...

## Investigating only a subset of paired comparisons after principal component analysis

Castura, J.C., Varela, P., & Næs, T. (2023).  
Investigating only a subset of paired comparisons after principal  
component analysis. *Food Quality and Preference*, 110, 104941.  
<https://doi.org/10.1016/j.foodqual.2023.104941>

## Investigating a subset of paired comparisons after PCA

Example where only a subset of paired comparisons are of primary interest:

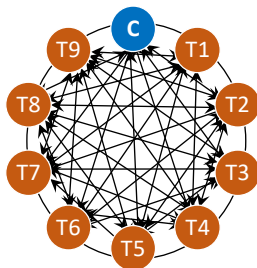
### Many Test Products vs One Control

Focus on Test-Control pairs ...*not Test-Test pairs*

## 1 Control vs 9 Test Products

$X \ominus X$

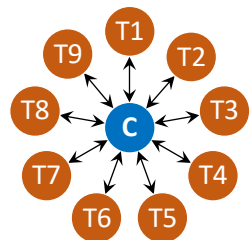
All Pairs



45 paired comparisons  
90 paired differences

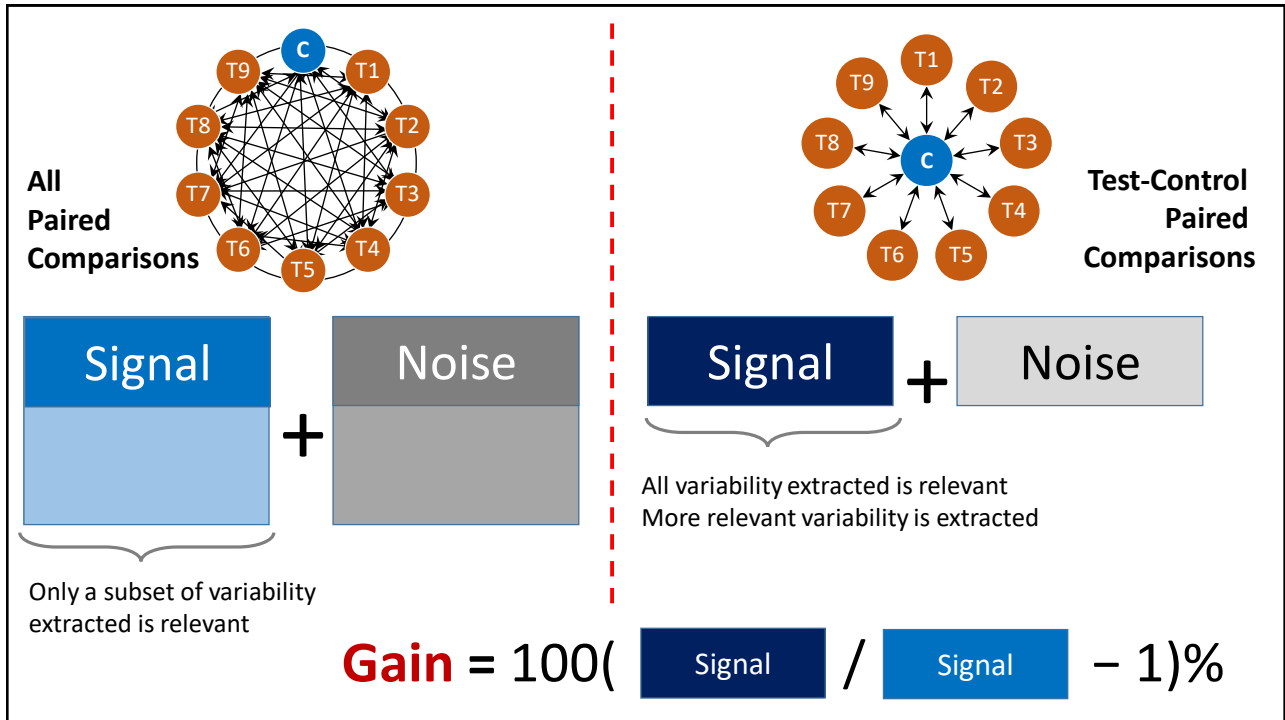
$\Delta^*$

Test-Control Pairs



9 paired comparisons  
18 paired differences

Castura, Cariou & Næs (2024)



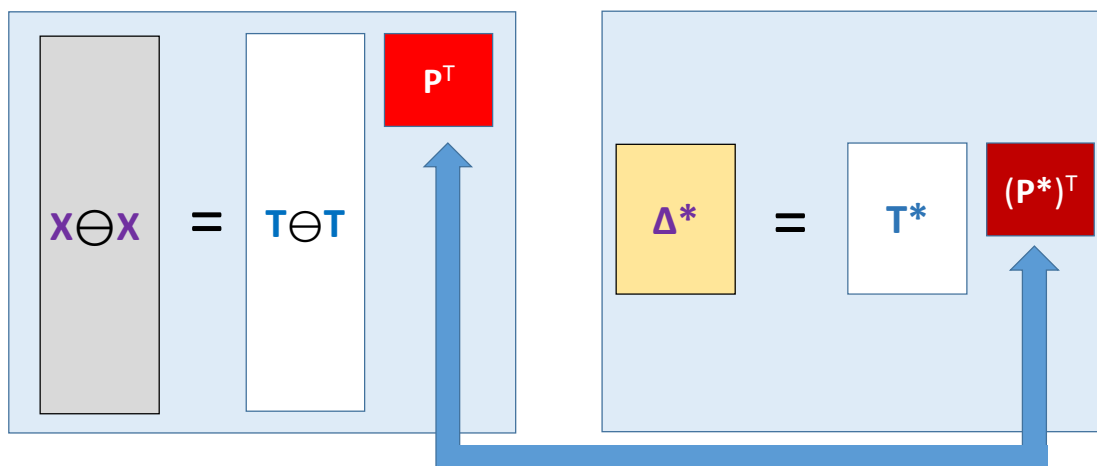
## Investigating only a subset of paired comparisons

“...the interrelationships between the variables might be different for the subset of paired comparisons than it is for all paired comparisons. So the covariance matrix for a matrix of all paired comparisons and the covariance matrix of selected paired comparisons will differ depending on the data. ”

Castura, J.C., Varela, P., & Næs, T. (2023). Investigating only a subset of paired comparisons after principal component analysis. *Food Quality and Preference*, 110, 104941.



## Investigating a subset of paired comparisons after PCA

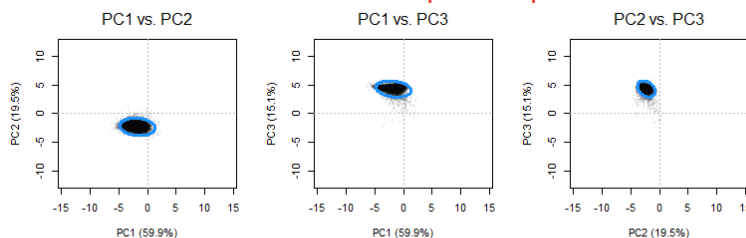


$\Delta^*$  contains a subset of the rows in  $X \ominus X$

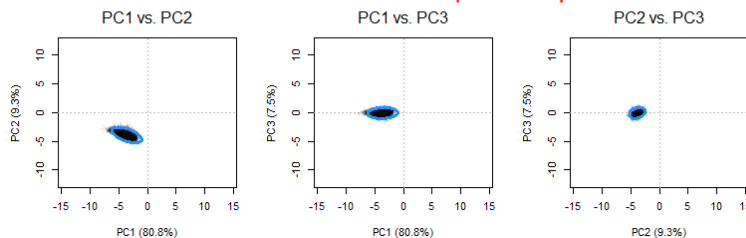
Loading matrices differ

## Investigating a subset of paired comparisons after PCA

T3-C based on PCA of all paired comparisons



T3-C based on PCA of selected paired comparisons



Gain:

1 PC:  
15%

2 PCs:  
14%

3 PCs:  
1%

Castura, Varela, & Næs (2023) [eComponent]  
doi:10.1016/j.foodqual.2023.104941

## Investigating a subset of paired comparisons after PCA

Another example where only a subset of paired comparisons are of primary interest:

### Temporal sensory data

Focus:

Paired comparisons *within* timepoints, **not** *across* timepoints

## Investigating a subset of paired comparisons after PCA

### All Pairs

- 8 yogurts  $\times$  56 timepoints
- 448 combinations
- All pairs = 100,028

$\mathbf{X} \ominus \mathbf{X}$  has dimension  
 $100028 \times 10$

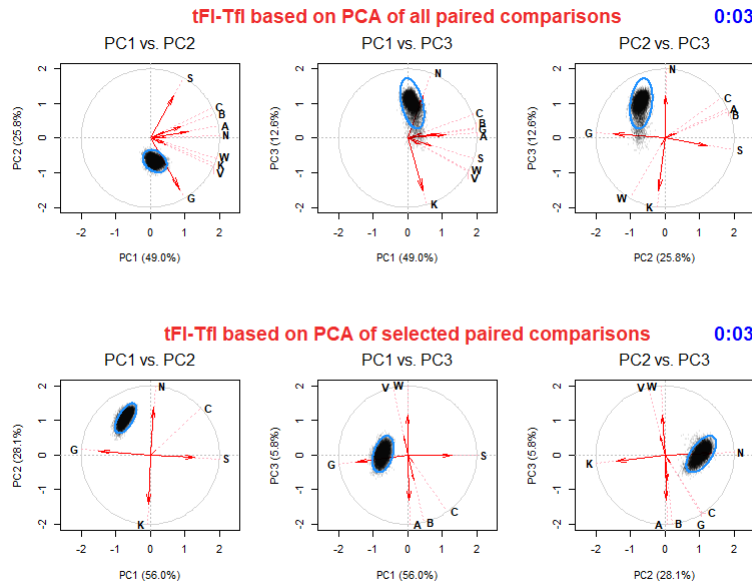
### Within-timepoint Pairs

- 28 within-timepoint pairs
- 56 timepoints
- $C = 28 \times 56 = 1568$

$\Delta^*$  matrix has dimension  
 $3136 \times 10$

## Investigating a subset of paired comparisons after PCA

Castura, Varela, & Næs (2023) [eComponent]  
doi:10.1016/j.foodqual.2023.104941



**Gain:**

1 PC:  
**>3500%**

2 PCs:  
**52%**

3 PCs:  
**1%**

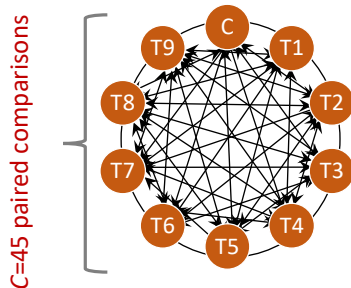
## Investigating control-centred results after uncentred principal component analysis

Castura, J.C., Cariou, V., & Næs, T. (2025). Investigating control-centred results after uncentred principal component analysis. *Zenodo*. Preprint.  
[Manuscript under review. Preprint not peer reviewed.]

***This preprint to be updated very soon!***  
<https://doi.org/10.5281/zenodo.15073361>

# 1 Control & 9 Test Products

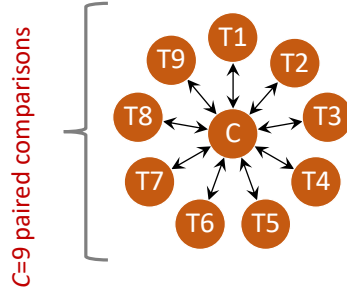
All Pairs



$$X \ominus X$$

$$J^2=100 \text{ rows}$$

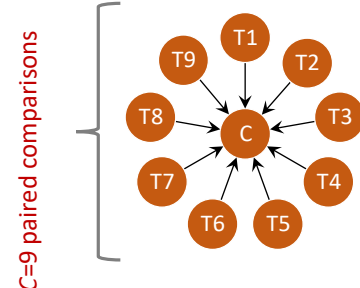
Test-Control Pairs



$$\Delta^*$$

$$2C=18 \text{ rows}$$

Test-Control Differences

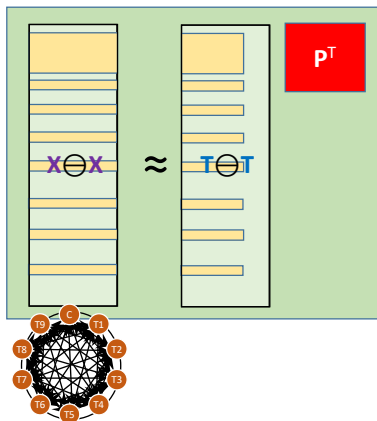


$$X^C$$

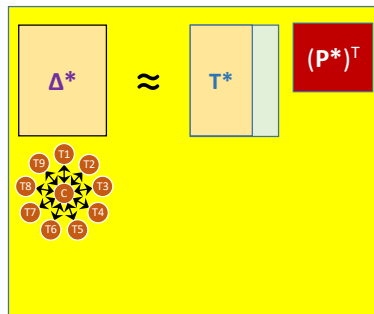
$$C+1=10 \text{ rows}$$

Castura, Cariou & Næs (2025)

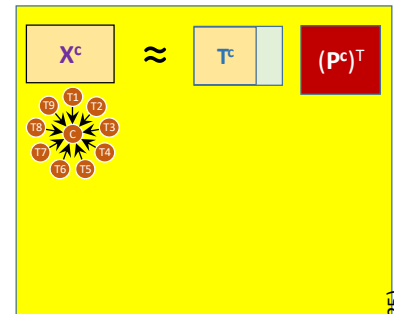
Centred PCA of  $X \ominus X$



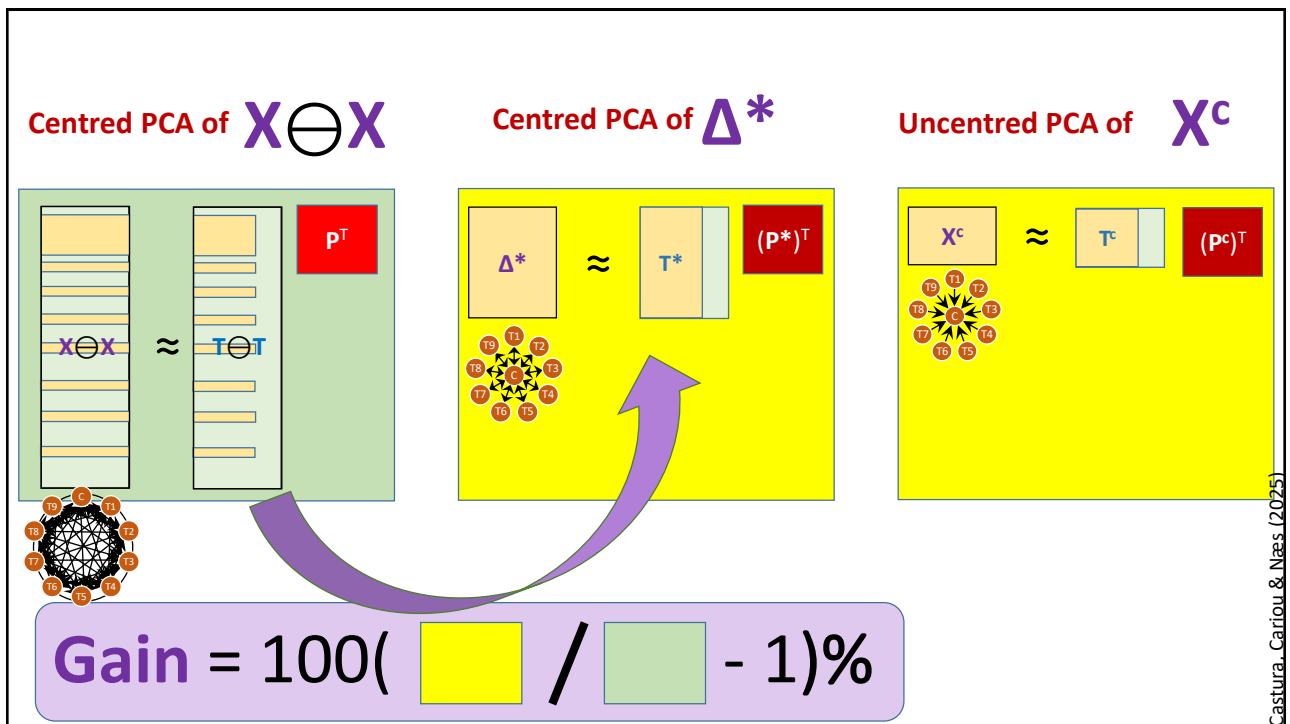
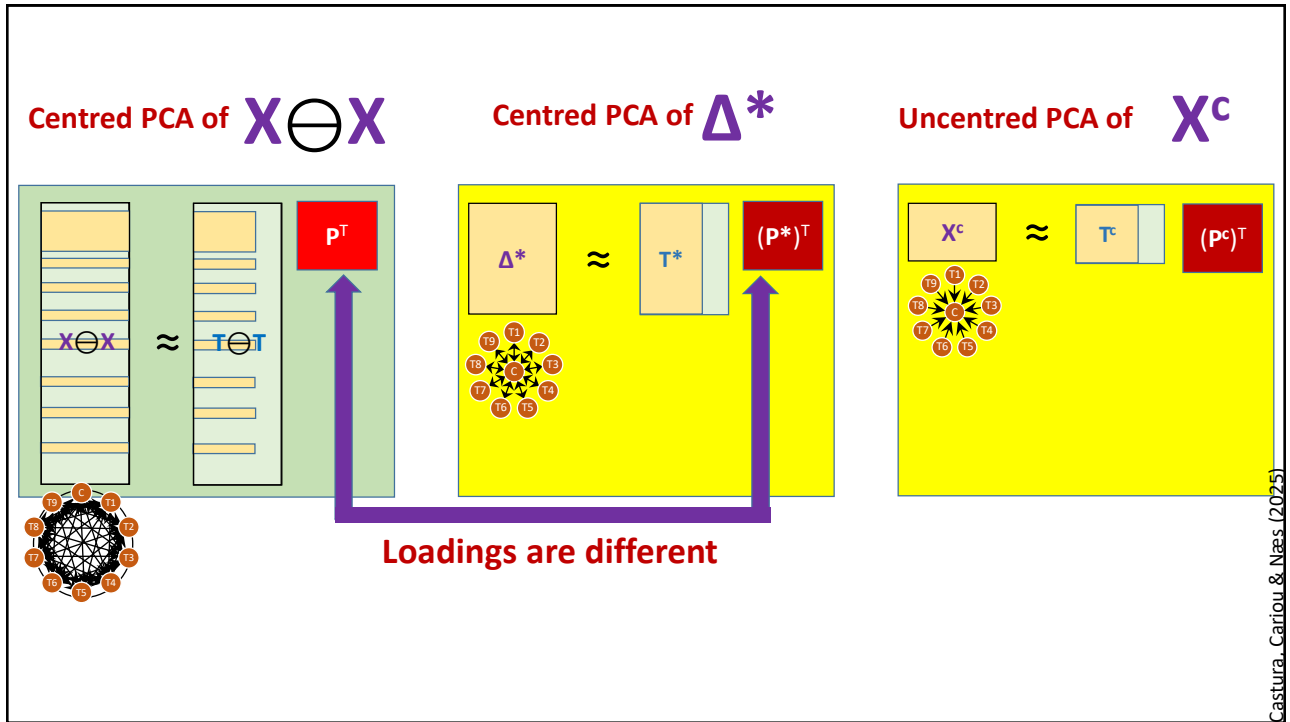
Centred PCA of  $\Delta^*$

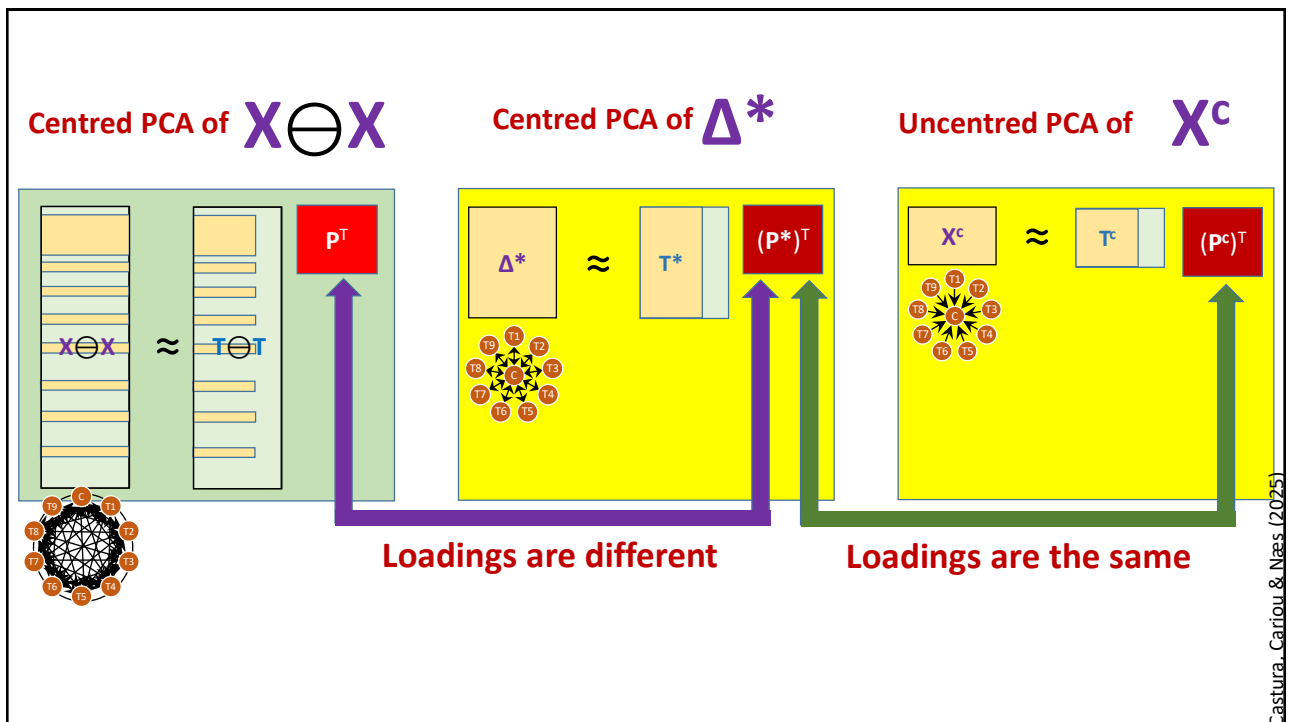
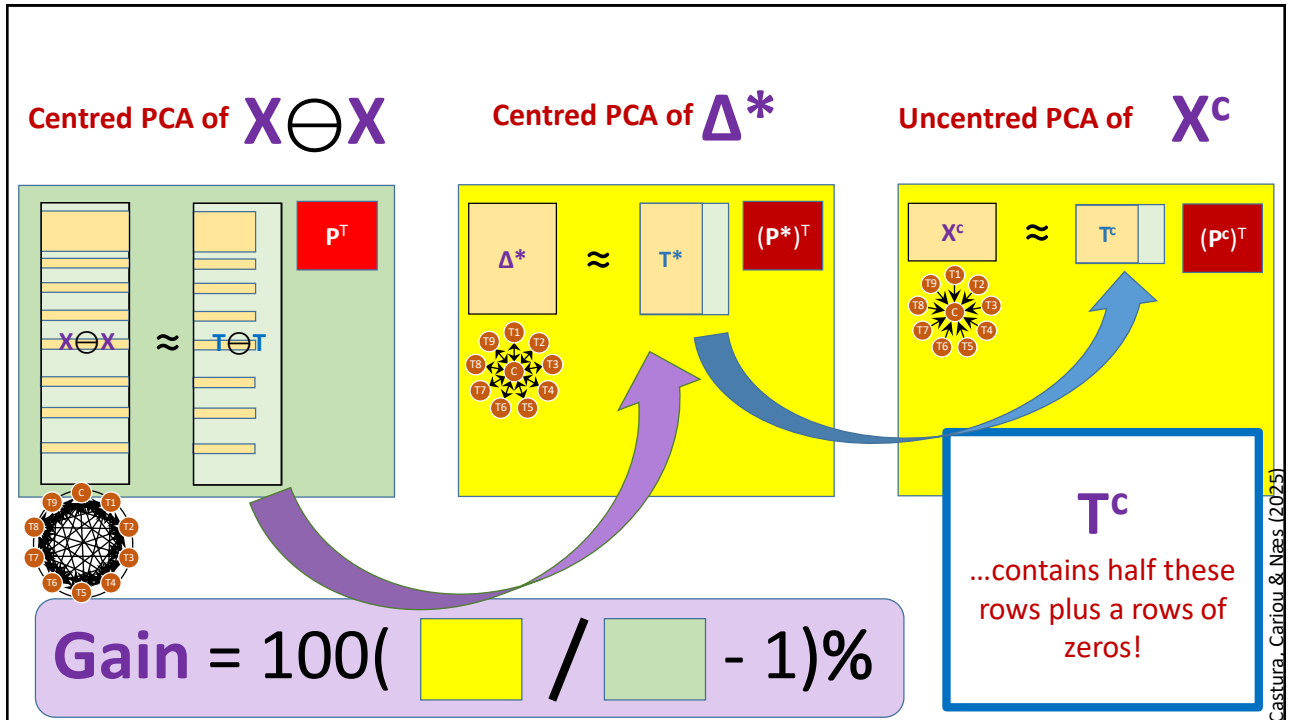


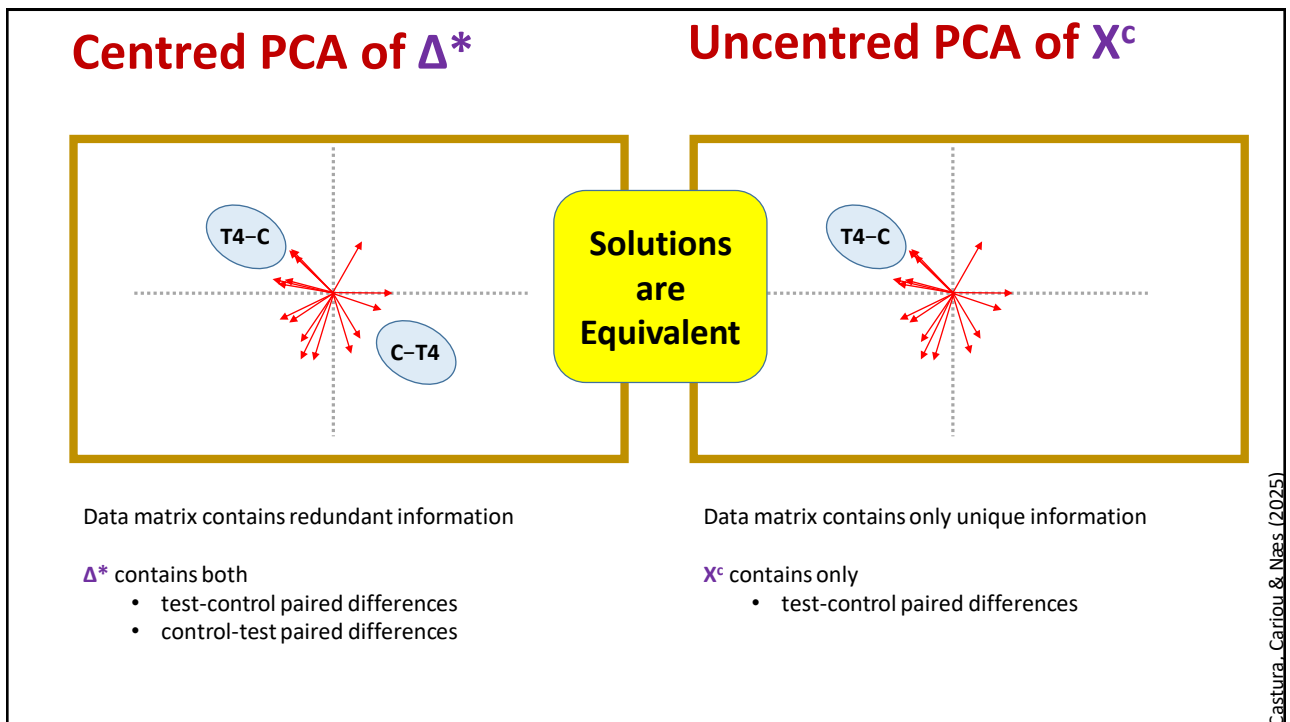
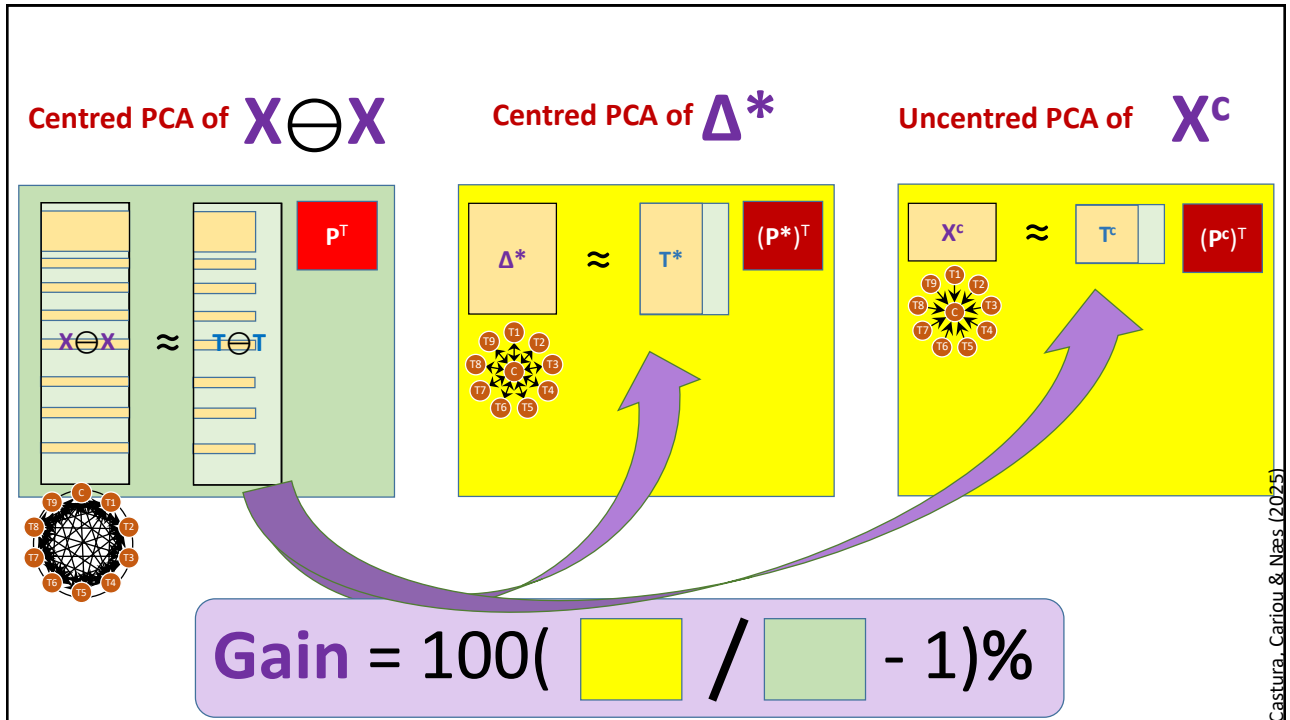
Uncentred PCA of  $X^C$



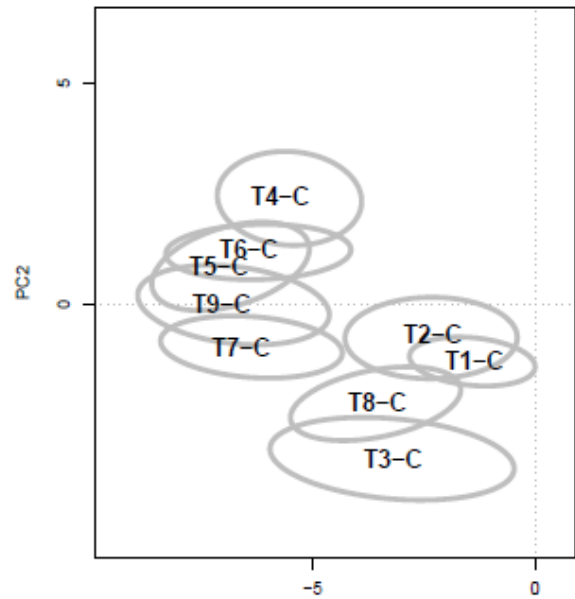
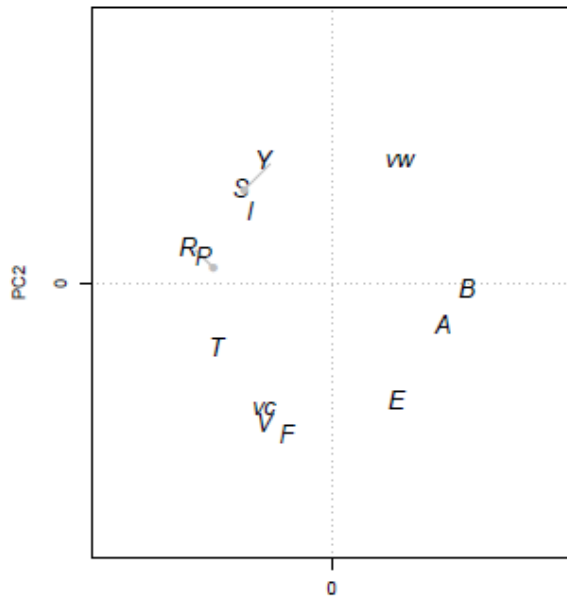
Castura, Cariou & Næs (2025)





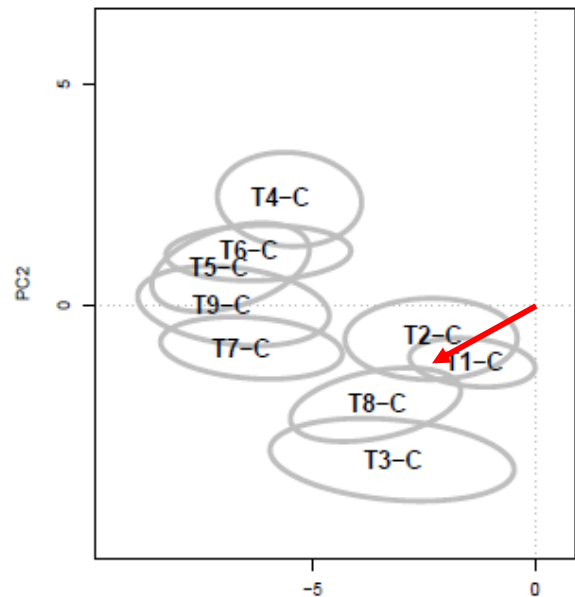
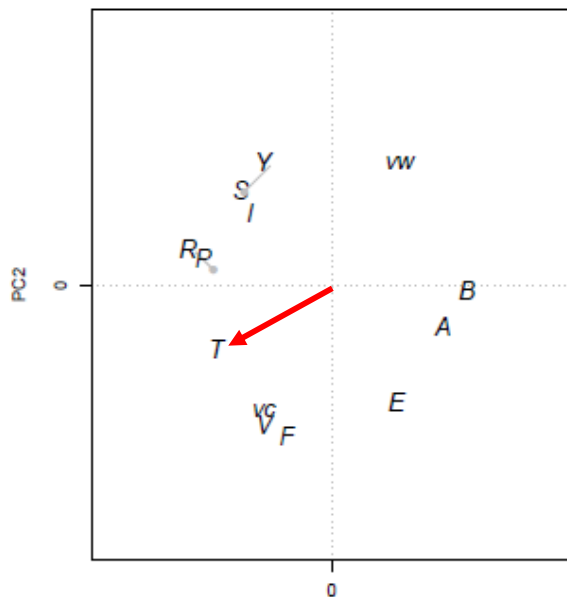


## Uncentred PCA of $X^c$



Castura, Cariou & Næs (2025)

## Uncentred PCA of $X^c$



Castura, Cariou & Næs (2025)



## For investigating Test-Control paired comparisons...

### Advantages of uncentred PCA of $X^c$

- “best” subspace for investigating the relevant test-control pairs
- components ordered by importance
- these particular uncentred PCA results have a conventional “variance interpretation”
- origin interpretation: “no difference from control”

Castura, Cariou & Næs (2025)

## Supervised principal component regression of selected paired comparisons

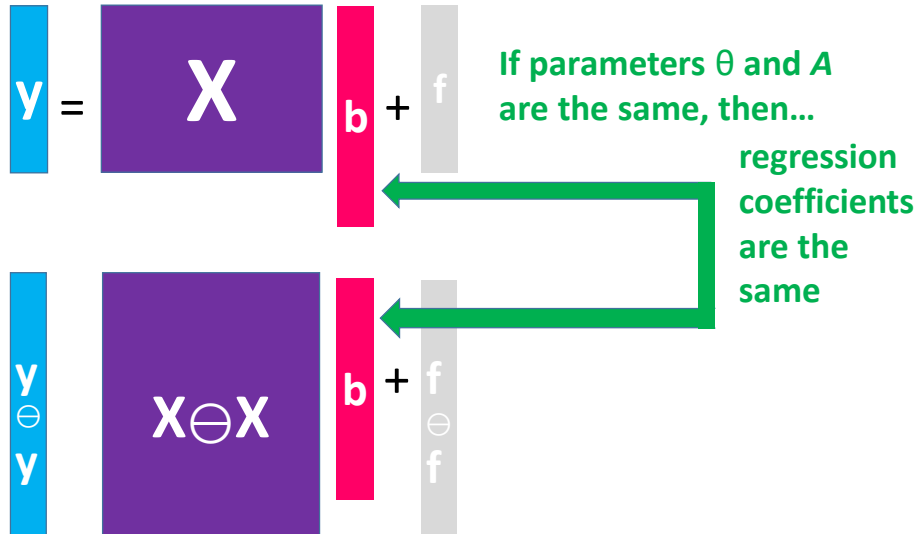
**Castura, J.C., & Tomic, O. (2024).** Supervised principal component regression of select paired comparisons. *Zenodo*.  
[Manuscript under review. Preprint not peer reviewed].  
<https://doi.org/10.5281/zenodo.11663995>

***This preprint to be updated very soon!***

**Castura, J.C., & Tomic, O. (2025).** Supervised principal component regression of select paired comparisons.  
[To be uploaded soon to the *Zenodo* preprint server!]  
[Manuscript under review. Preprint not peer reviewed].

## Investigating data relationships

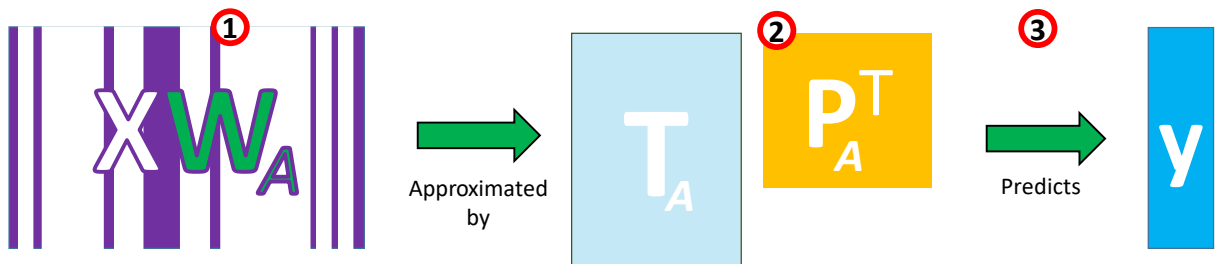
with supervised principal component regression (SPCR)



Castura & Tomic (2024)

## Investigating data relationships

with supervised principal component regression (SPCR)



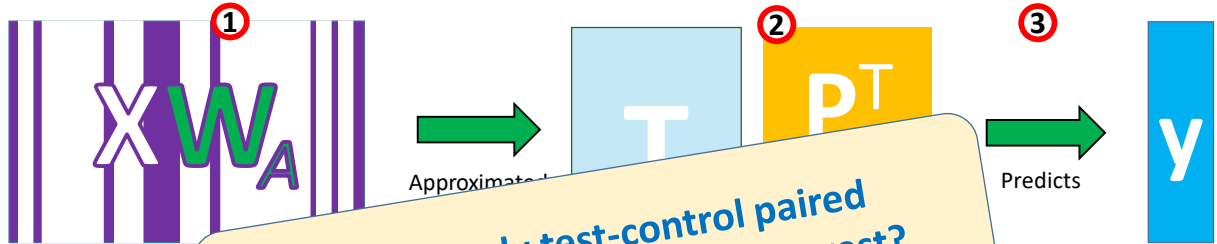
SPCR of  $X$  and  $y$  is equivalent to SPCR of  $X \ominus X$  and  $y \ominus y$  because regression coefficients are the same, which occurs the following are the same:

- ①  $X$  weights
- ②  $X$  loadings
- ③  $Y$  loadings

Castura & Tomic (2024)

## Investigating data relationships

with supervised principal component regression (SPCR)



What if only test-control paired differences are of primary interest?

SPCR of  $X$  and  $y$  occurs the following are the same:

- ①  $X$  weights
- ②  $X$  loadings
- ③  $y$  loadings

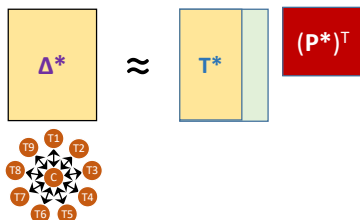
Castura & Tomic (2024)

## Investigating data relationships

with supervised principal component regression (SPCR)

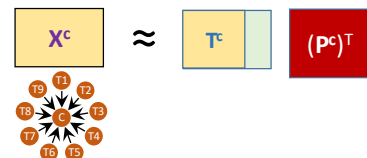
We know we can focus on paired comparisons or paired differences since...

Centred PCA of  $\Delta^*$



...is equivalent to...

Uncentred PCA of  $X^c$

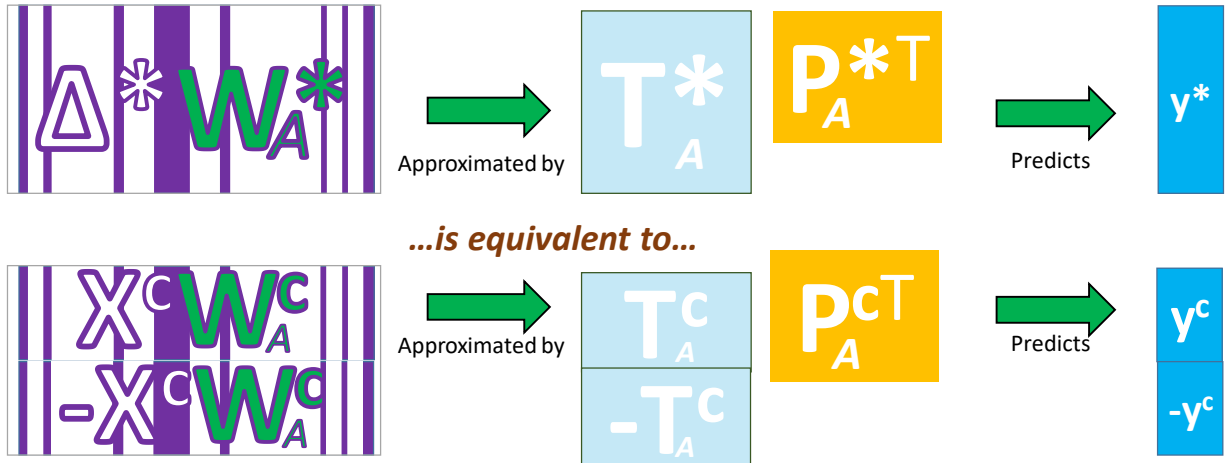


...however, using  $\Delta^*$  introduces a subtle problem...

Castura & Tomic (2025)

## Investigating data relationships

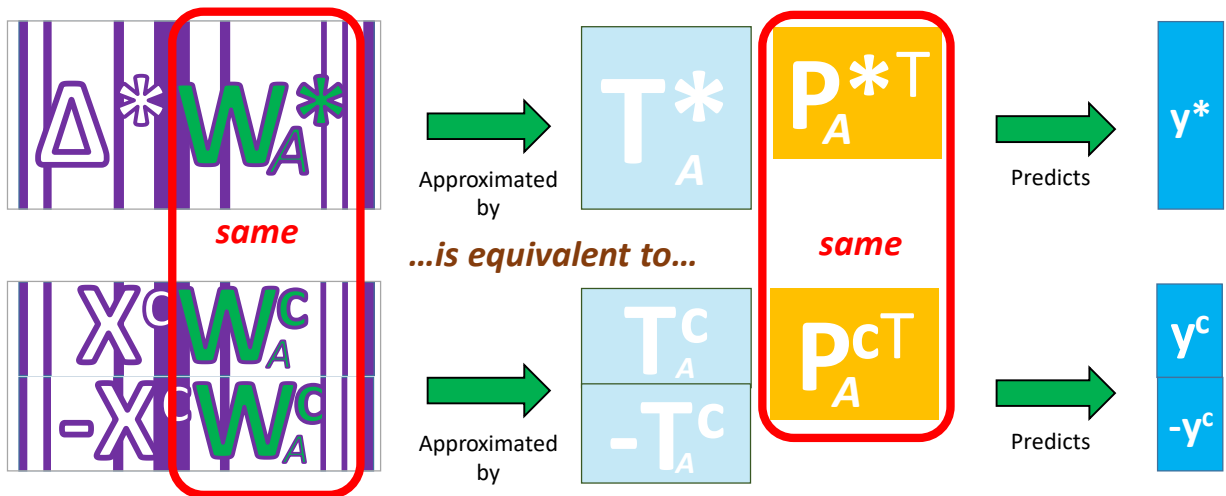
with supervised principal component regression (SPCR)



Castura & Tomic (2025)

## Investigating data relationships

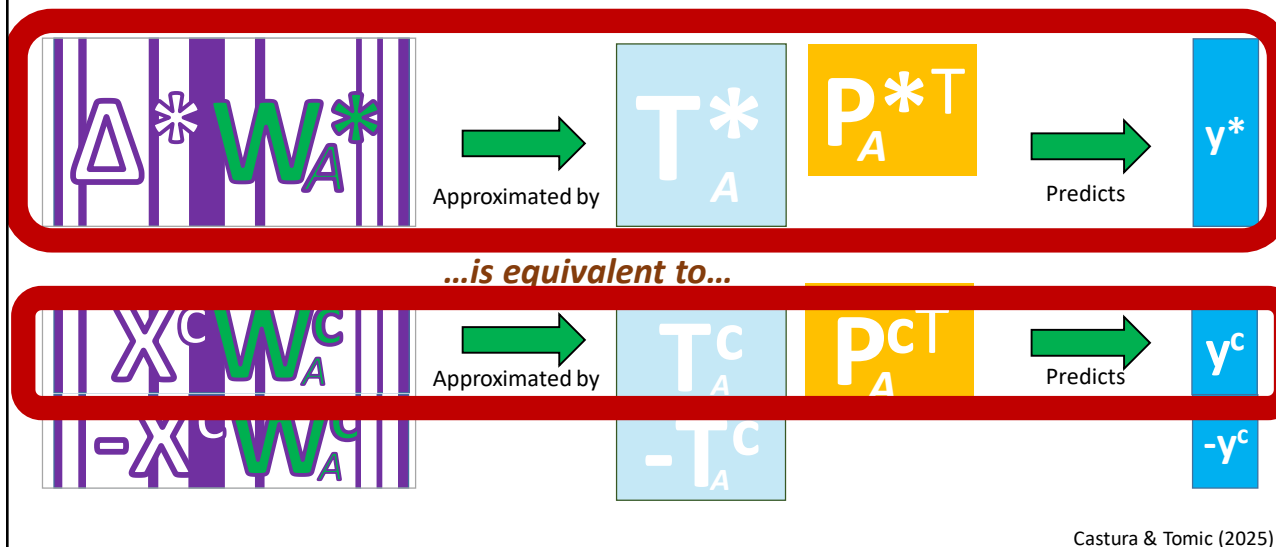
with supervised principal component regression (SPCR)



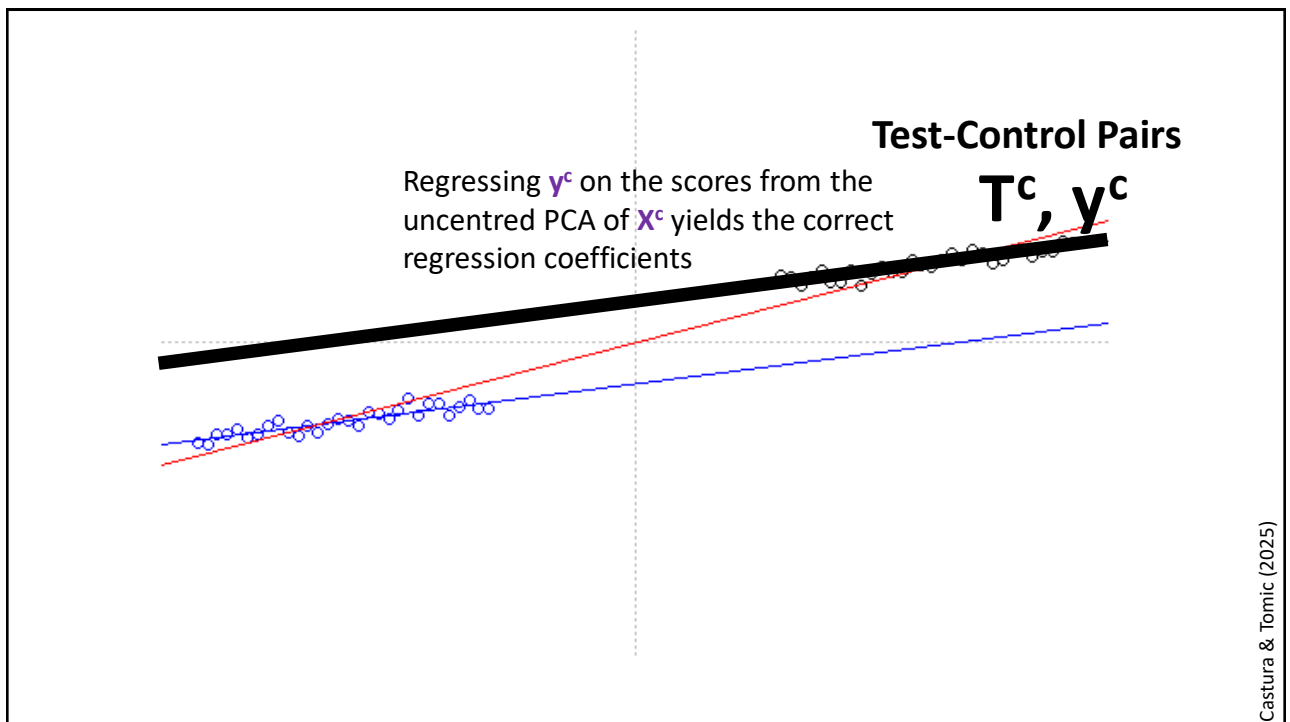
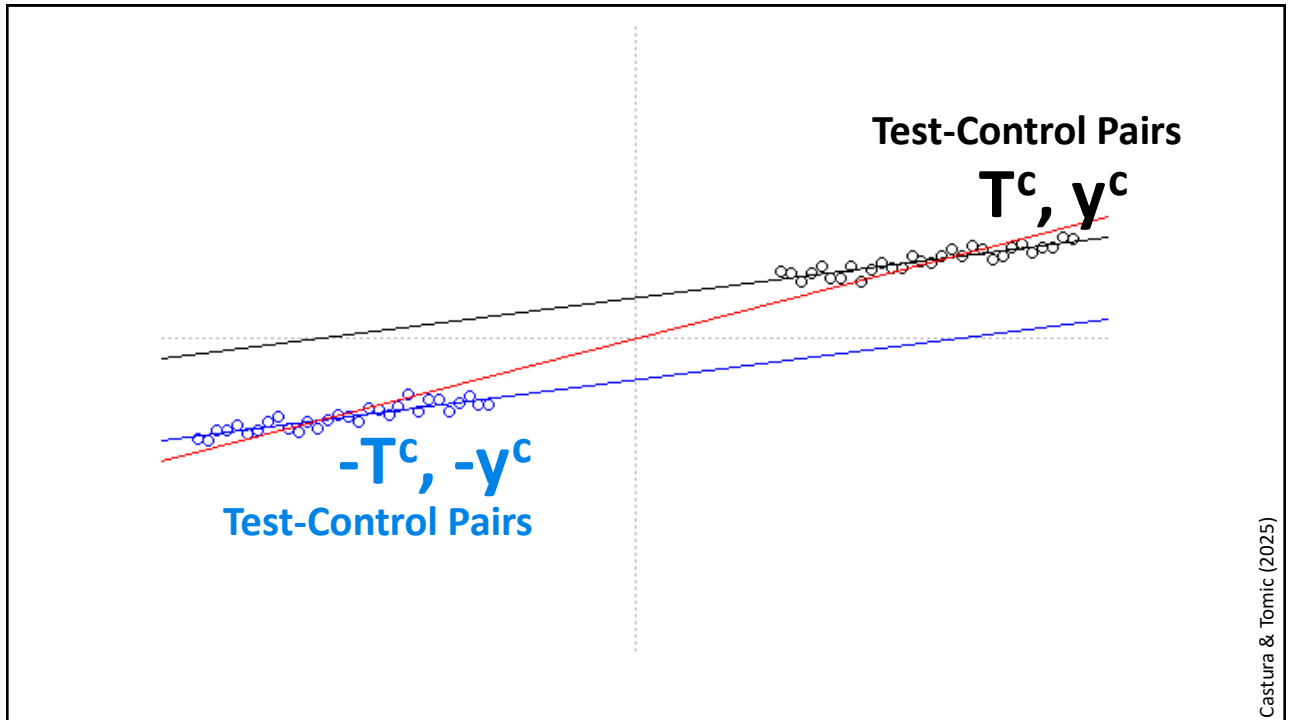
Castura & Tomic (2025)

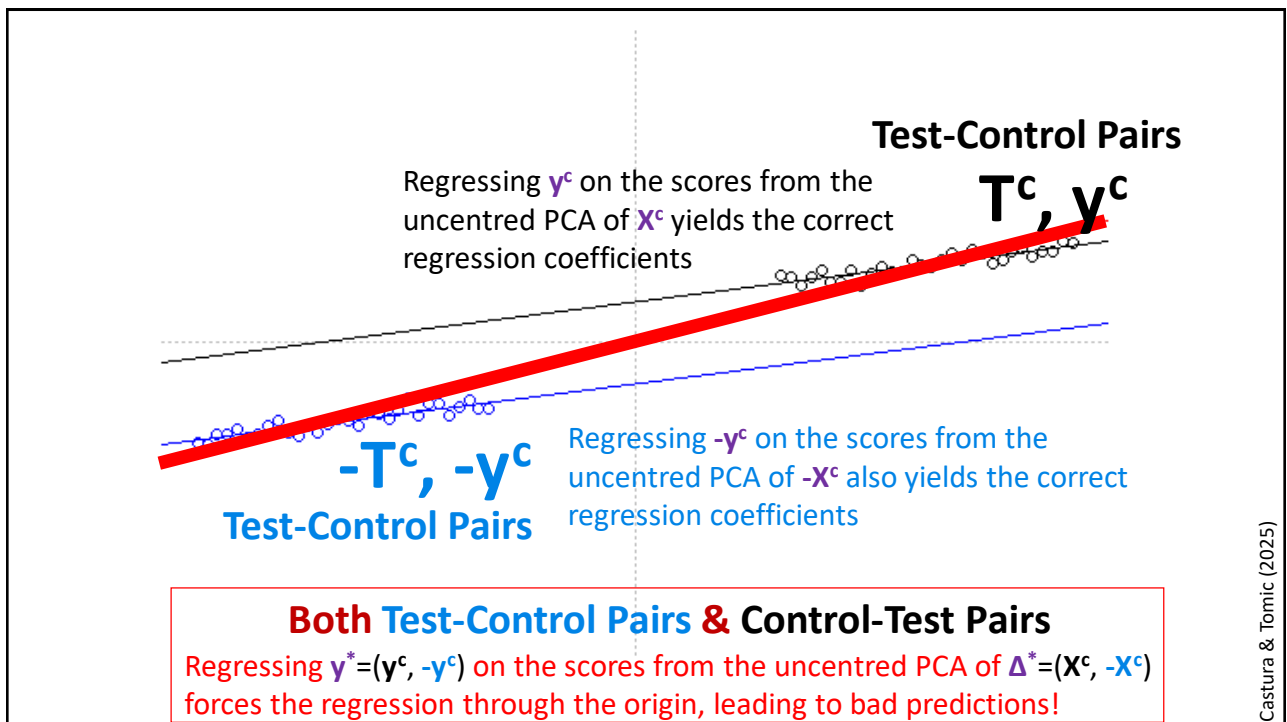
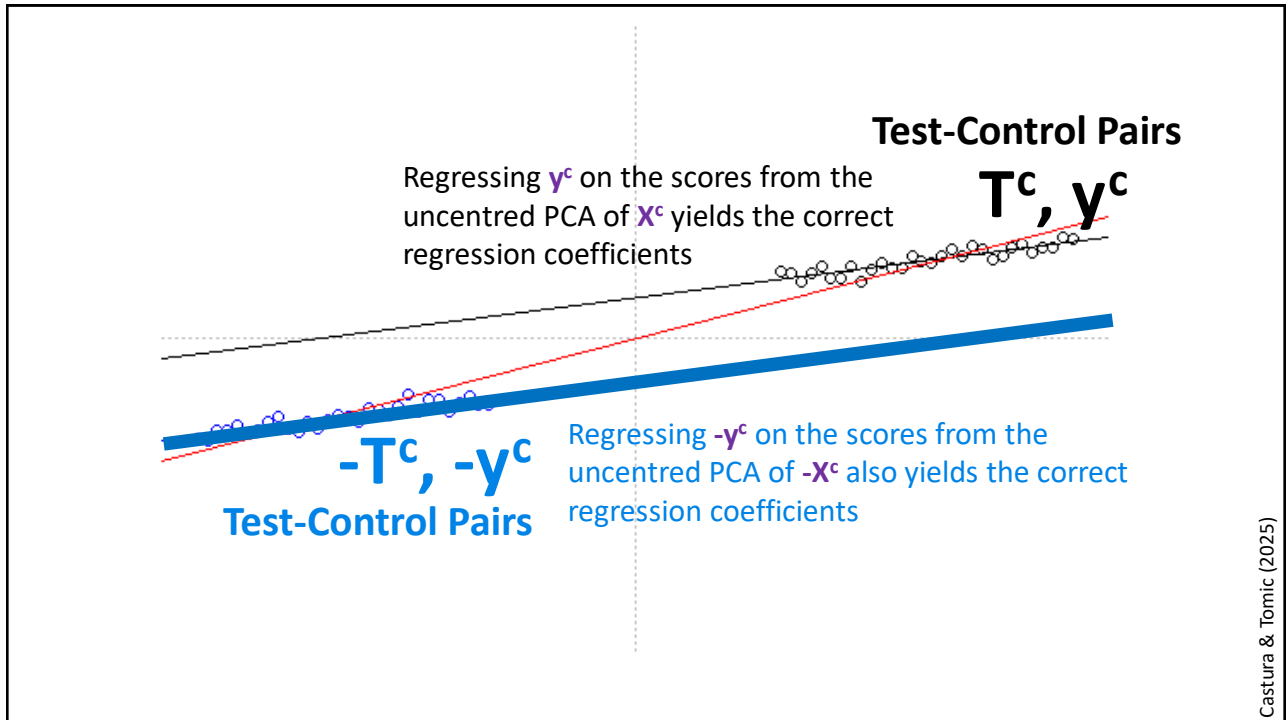
## Investigating data relationships

with supervised principal component regression (SPCR)



...however, using  $\Delta^*$  introduces a subtle problem...





Many candidate values for  
parameters  $\theta$  and  $A$

Since we are focused on paired comparisons,  
we **adapted leave-one-out cross-validation**  
to handle *dependent data*.

Castura & Tomic (2025)

### Conventional loocv

One row left out per fold.

Each row is an object.

No information from the  
left-out object leaks into  
the model during training.

Software for conventional  
loocv widely available.

### Adapted loocv

One object left out per fold.

Multiple rows are left out.

No information from the  
left-out object can leak into  
the model during training.

Functions for adapted loocv  
must be coded.

Castura & Tomic (2025)



## Predicted error sum of squares (PRESS)

Models evaluated by calculating PRESS statistic calculated for left-out test-control paired differences.

$$\sum (predicted - observed)^2$$

Castura & Tomic (2025)

## Investigating data relationships

with supervised principal component regression (SPCR)

$$y = Xb + f$$

If using loocv...

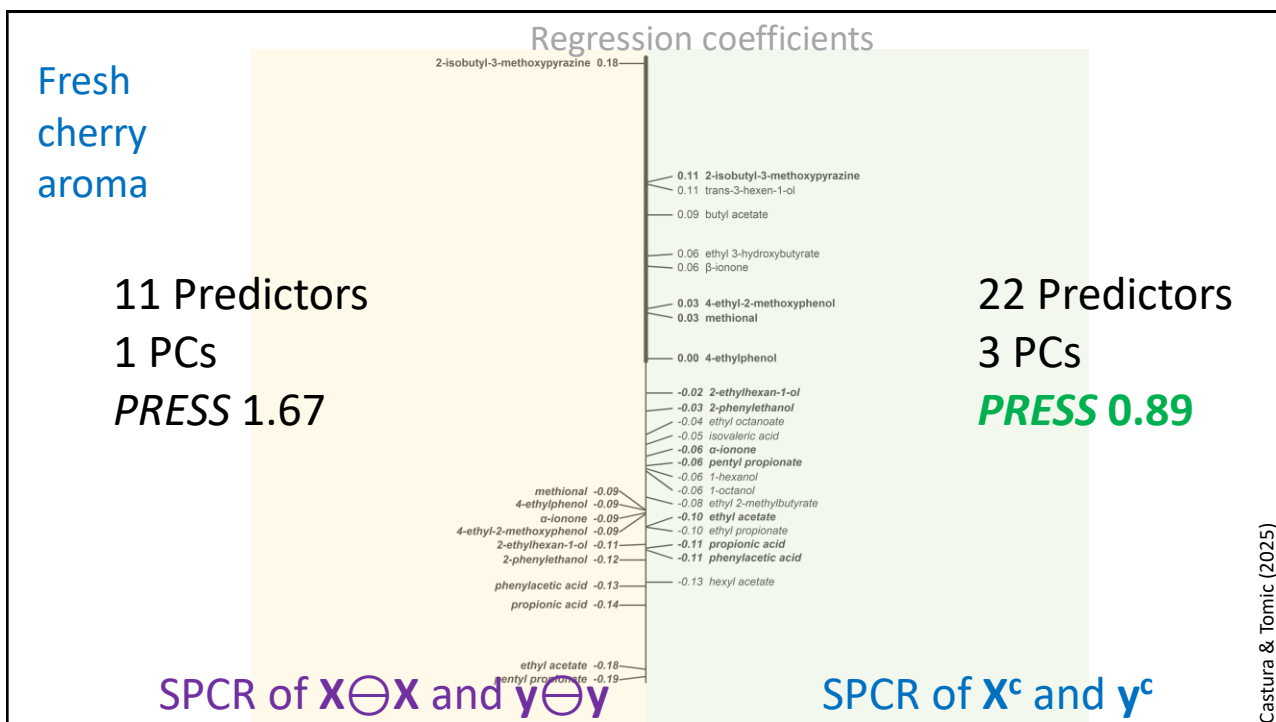
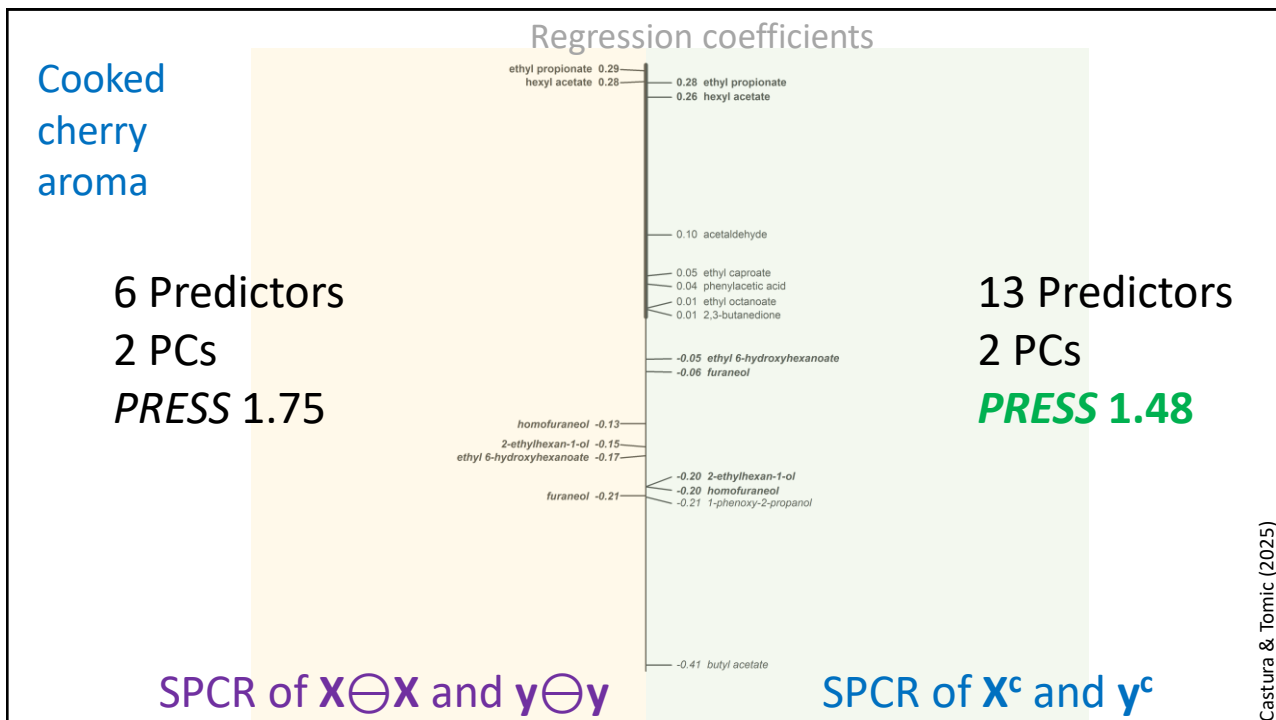
...SPCR of test-control paired differences has different regression coefficients

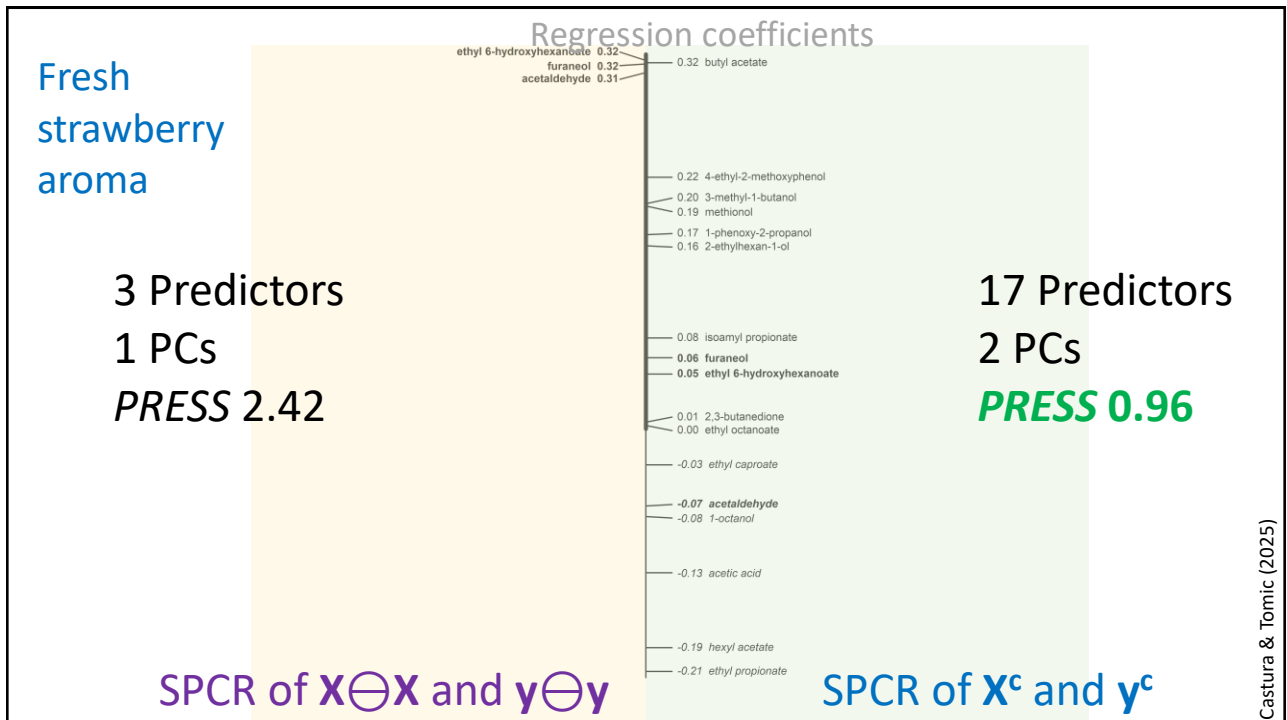
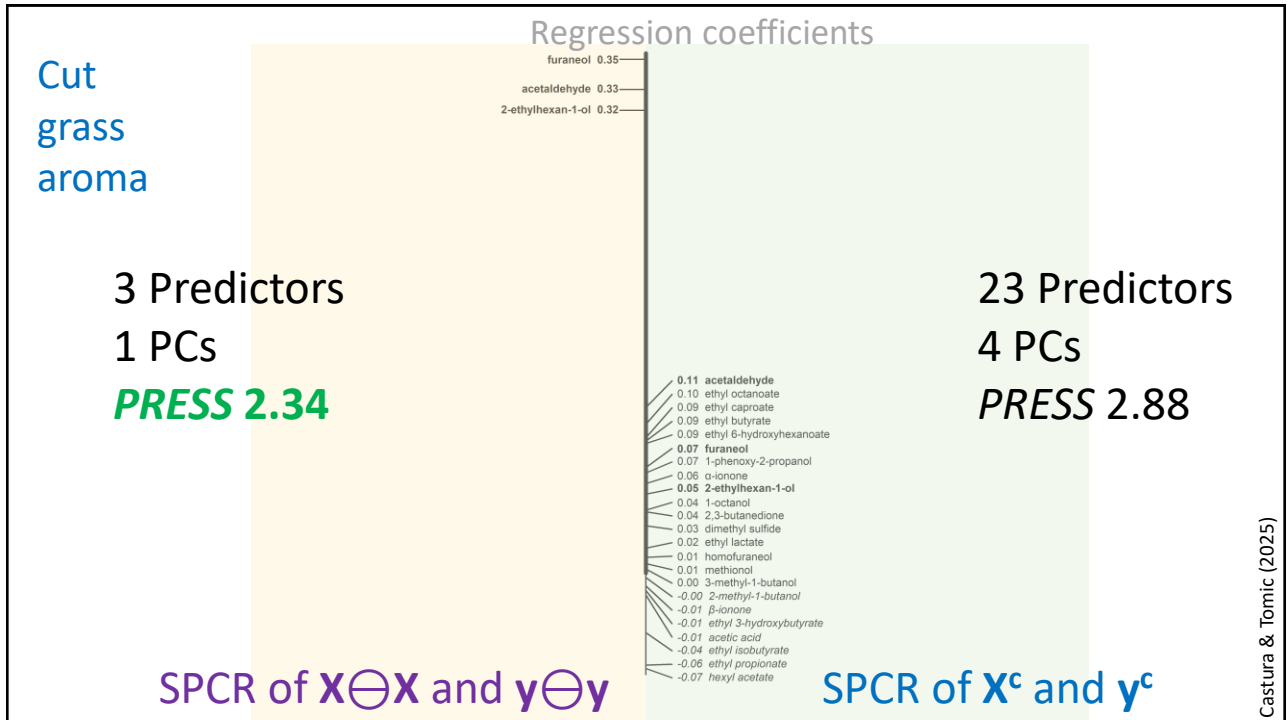
$$y \ominus y = X \ominus X b + f \ominus f$$

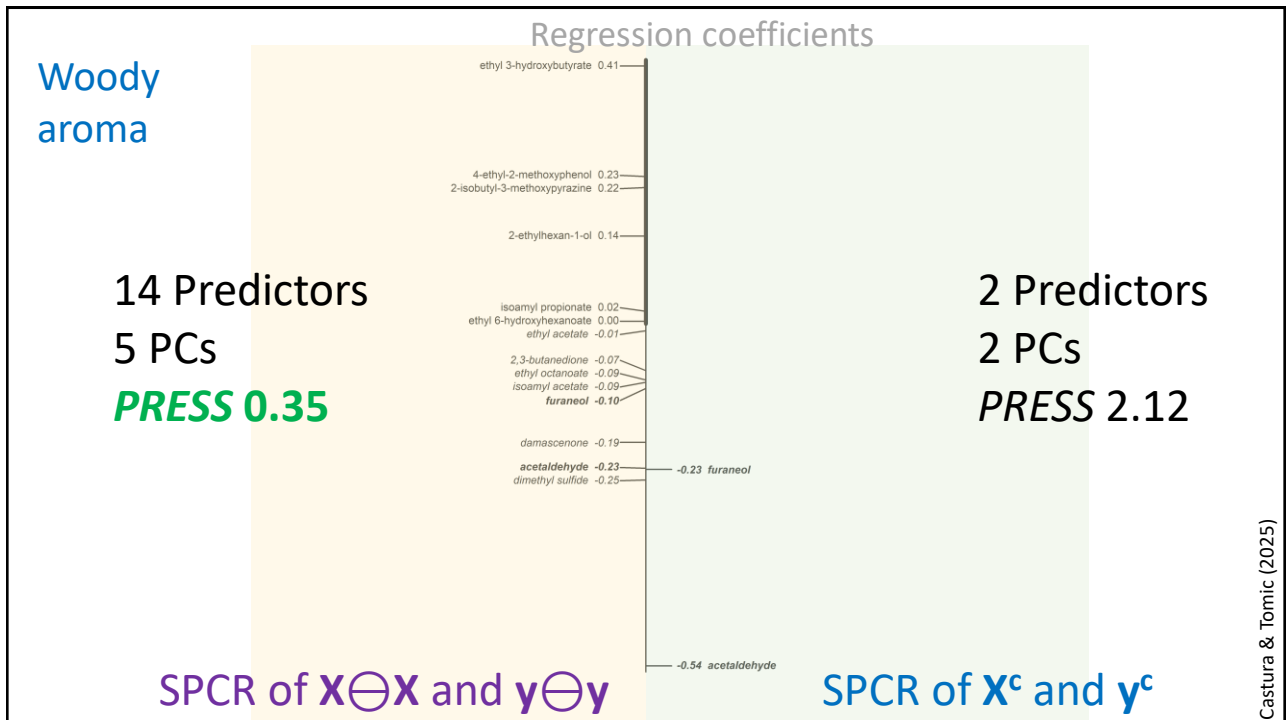
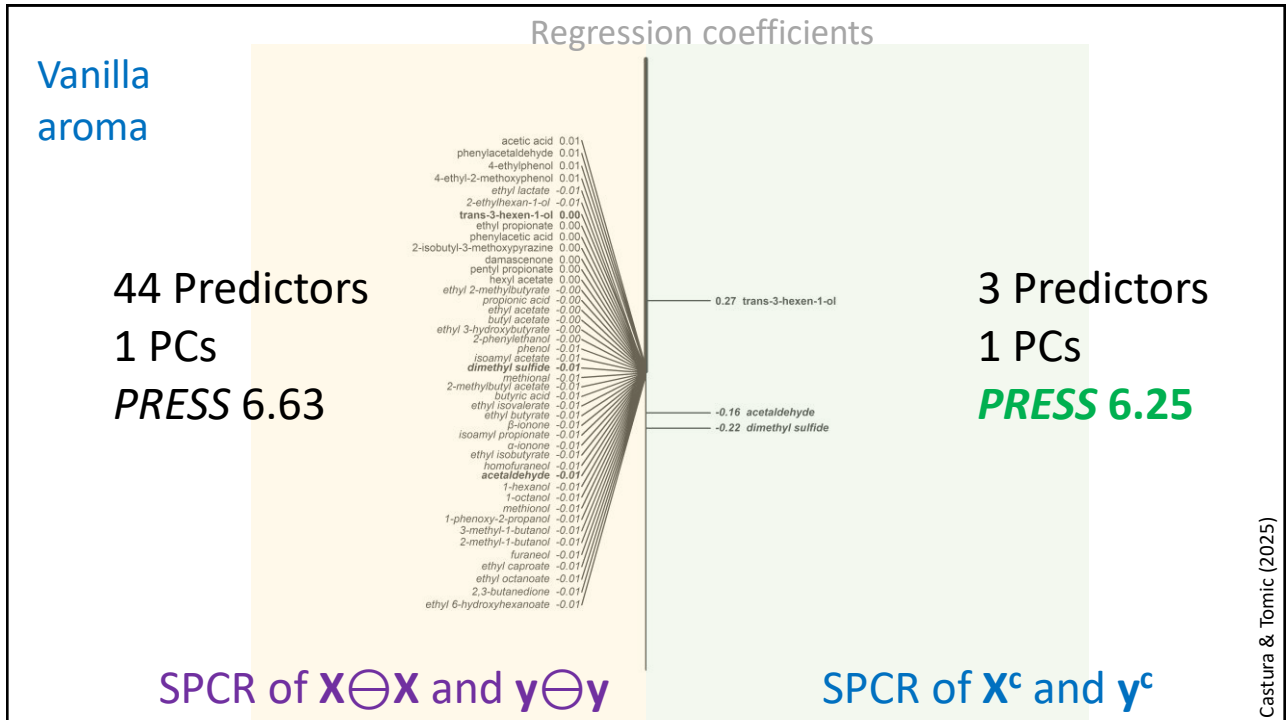
...SPCR of all objects has the same regression coefficients as SPCR of all paired differences

$$y^c = X^c b^c + f^c$$

Castura & Tomic (2025)





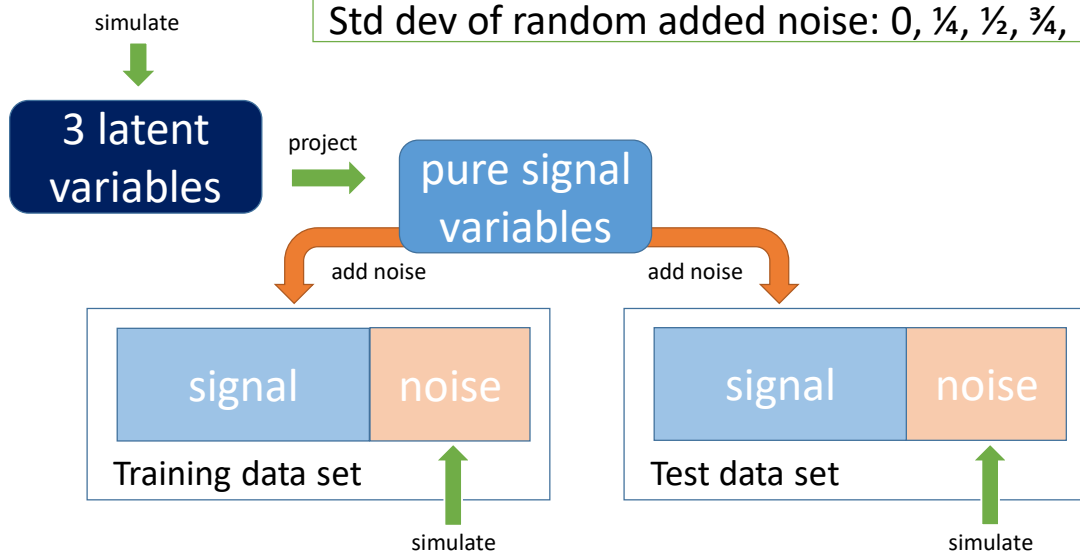


## Simulation study

Each condition was a combination of...

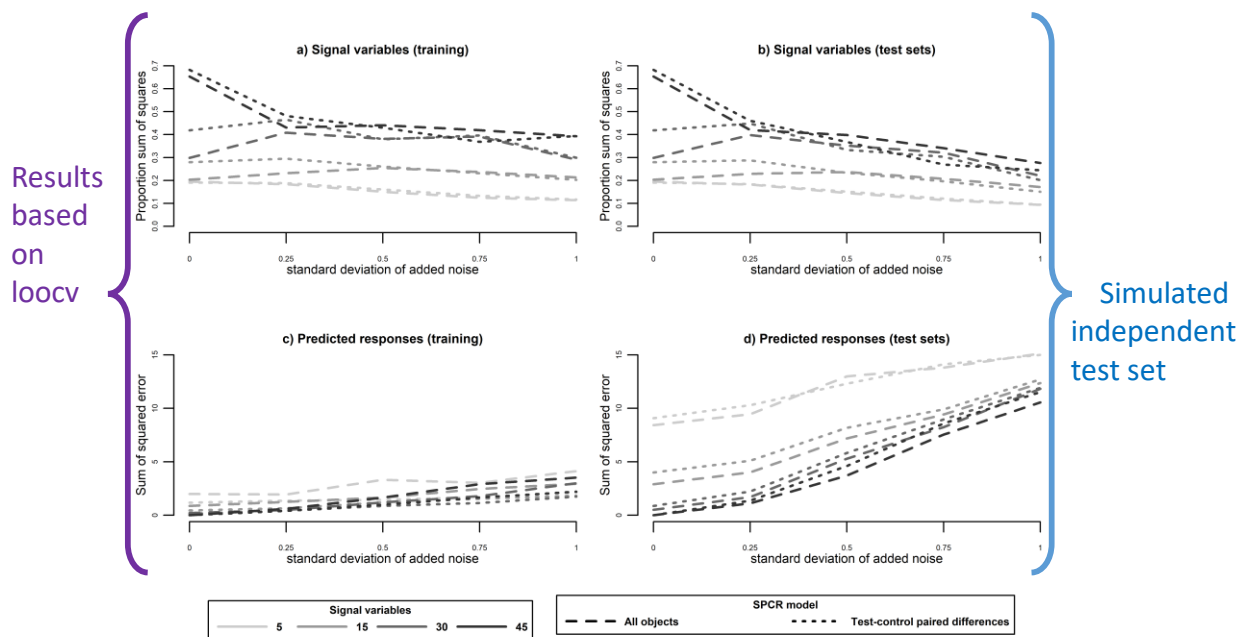
Number of “signal variables”: 5, 15, 30, 45.

Std dev of random added noise: 0,  $\frac{1}{4}$ ,  $\frac{1}{2}$ ,  $\frac{3}{4}$ , 1.



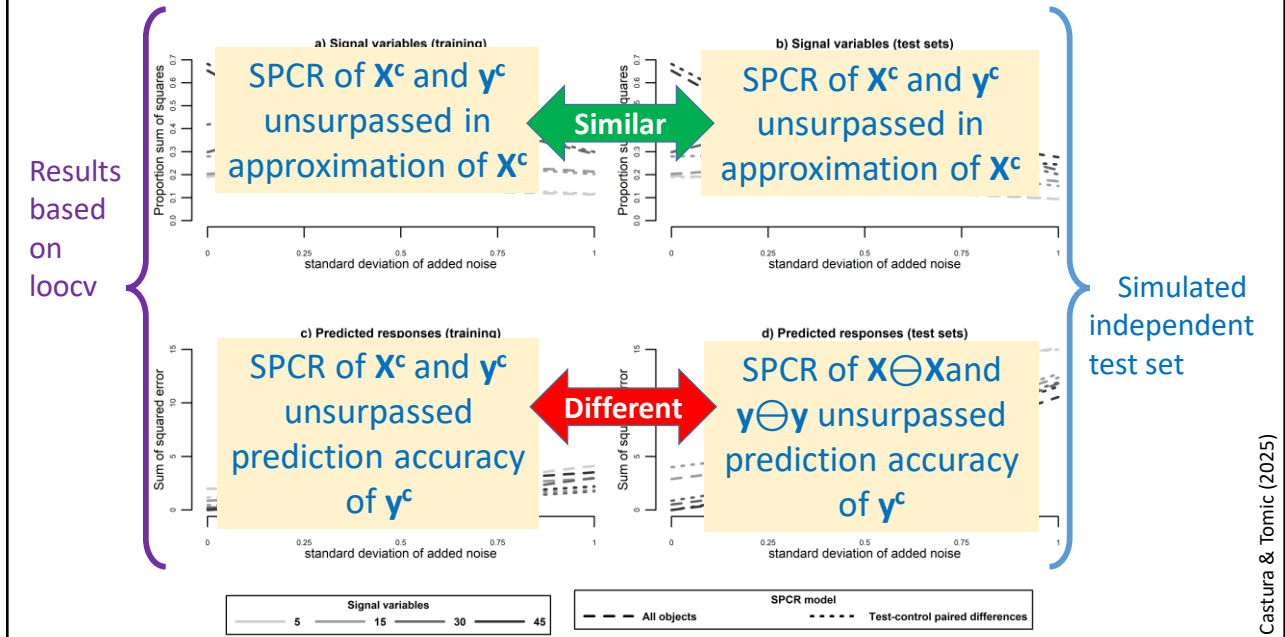
Castura & Tomic (2025)

## Simulation study



Castura & Tomic (2025)

## Simulation study

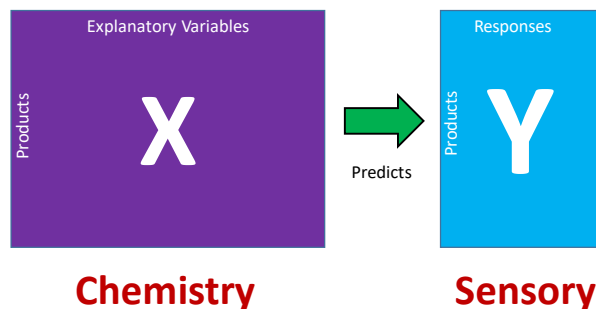


The simulation study shows the superiority of the SPCR of test-control paired differences model for the training set does not generalize to test sets simulated from the same data generating process.

**We are only interested in the test-control paired differences, not the test-test paired differences, but since both types of pairs carry information about the association between variables, the relationship between predictor and response variables is modelled better by the SPCR of all objects model than by the SPCR of test-control paired differences model.**

This finding shows **overfitting appears in many guises**. For this reason, the SPCR of all objects model is recommended for most routine analyses.

Castura & Tomic (2025)



## Partial least squares regression (PLSR)

### Partial least squares regression (PLSR)

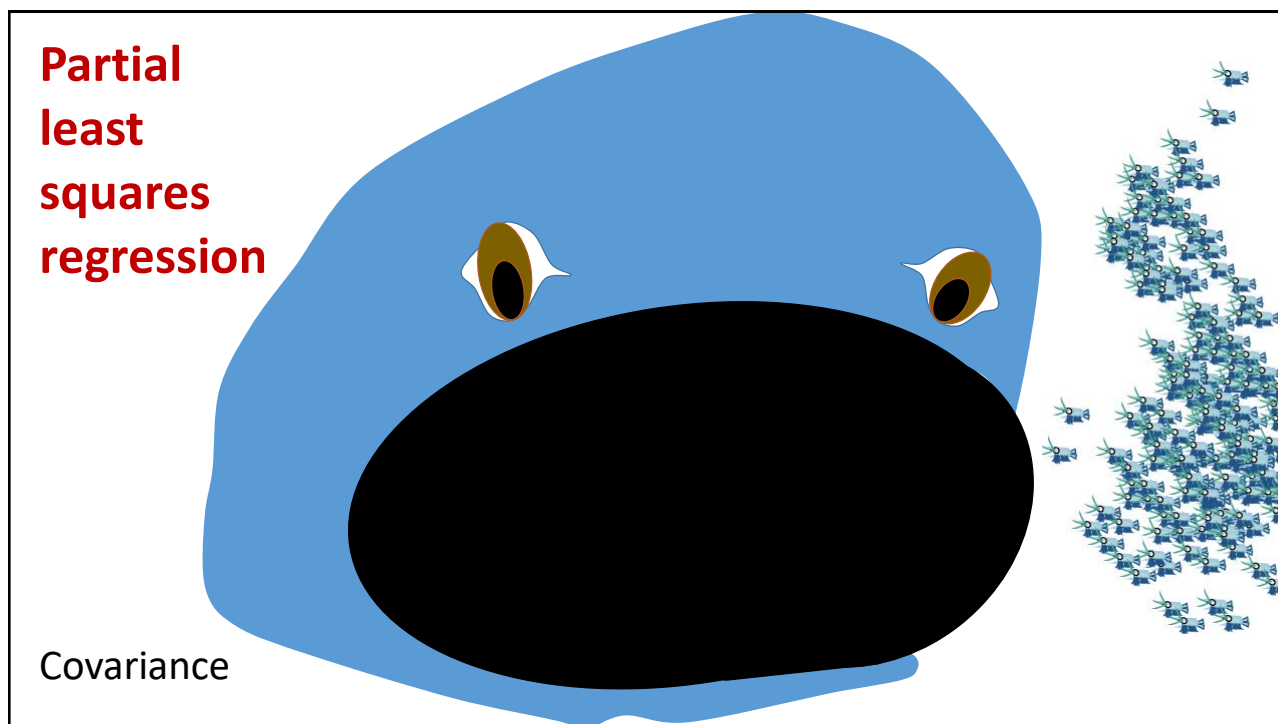
We want to ***predict*** multivariate **Y** from the multivariate **X**.

Successive PLS components extract ***covariation*** between **X** and **Y** maximally.

## Partial least squares regression (PLSR)

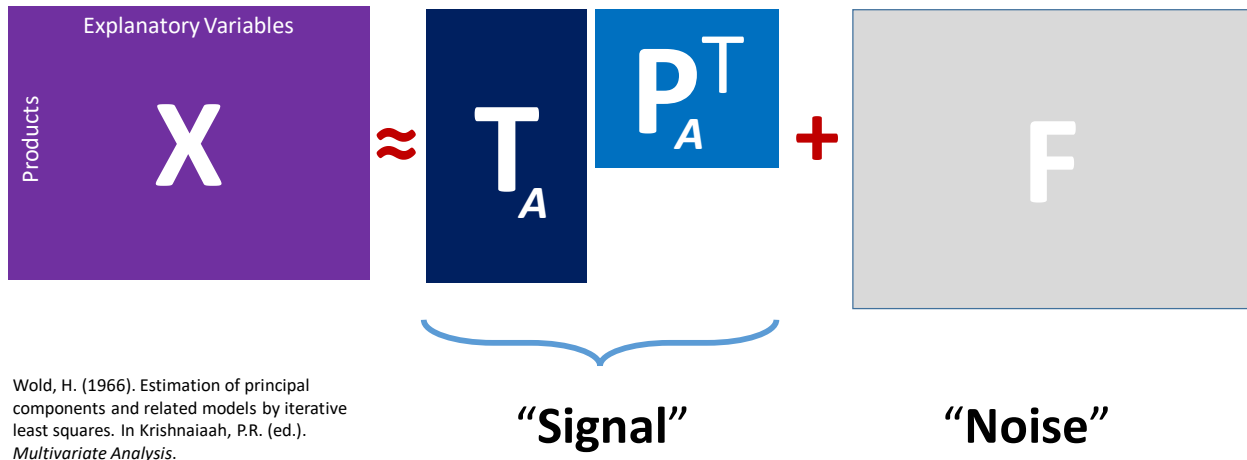
We want to ***predict*** multivariate **Y** from the multivariate **X**.

Successive **PLS components** extract ***covariation*** between **X** and **Y** maximally.

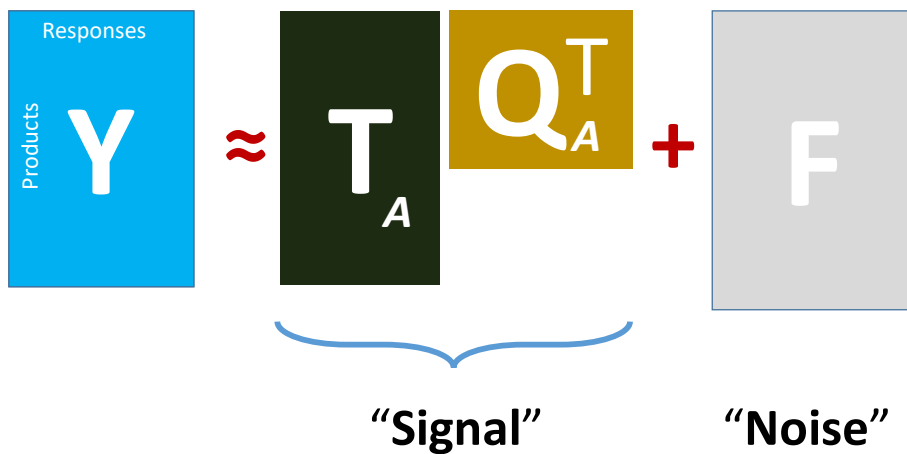




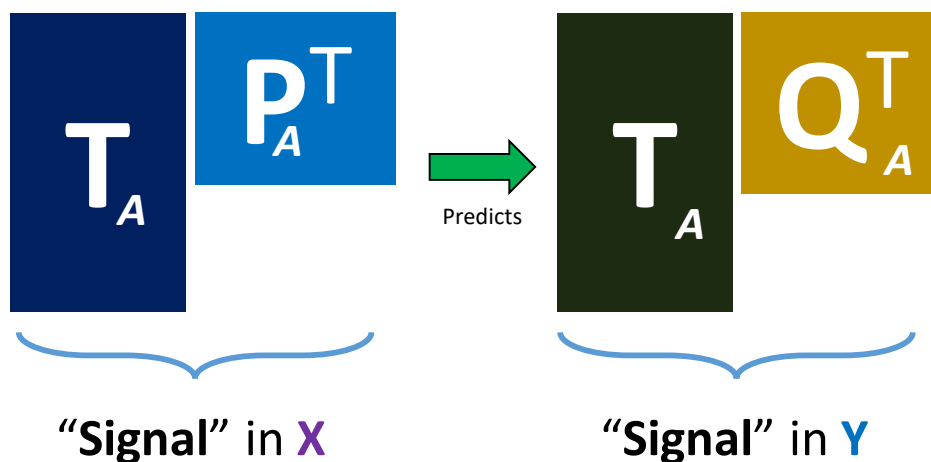
## Explaining and predicting data relationships



## Explaining and predicting data relationships



## Partial least squares regression (PLSR)



Wold (1966)

## PLS components

Linear  
combination  
of variables in

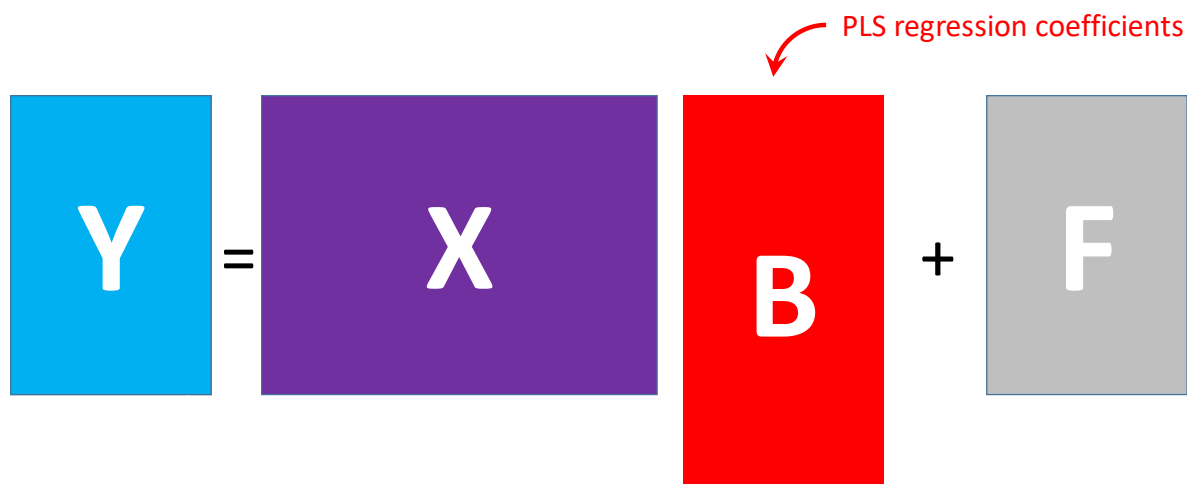
$X$

&

Linear  
combination  
of variables in

$Y$

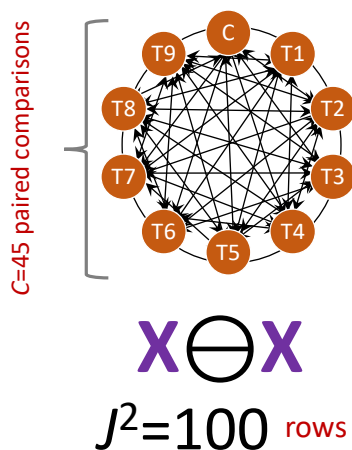
## Investigating data relationships



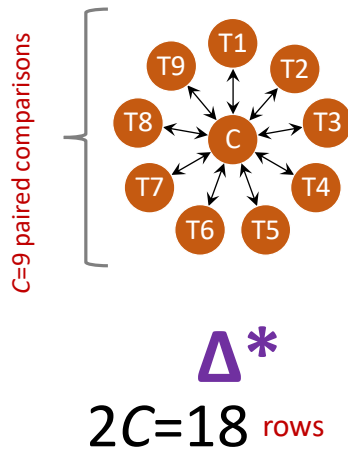
Wold (1966)

## Paired comparisons after PLSR

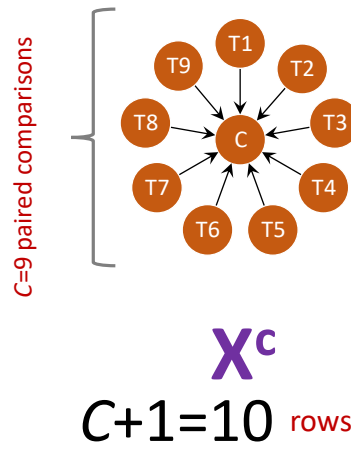
All Pairs



Test-Control Pairs



Test-Control Differences



Castura, Tomic & Næs (unpublished)

For further information, contact  
[jcastura@compusense.com](mailto:jcastura@compusense.com)



John Castura



## Acknowledgements



Oliver Tomic



Paula Varela



Tormod Næs



Véronique Cariou



# Exploring the relationship between sensory and instrumental data with component-based methods



## KoSFoST International Symposium and Annual Meeting 2025

*Pioneering Future Connection in FoodTech*  
Gwangju, Korea · 2-4 July 2025

**John Castura**

Dr. Philos., M.Sc.  
Research Fellow

