Image source: alamy.com

| observation | Var 1 | Var 2 |
|---|---|---|
| 1 | Rep 1 | Rep 1 |
| 1 | Rep 2 | Rep 2 |
| 1 | Rep 3 | Rep 3 |
| 2 | Rep 2 | Rep 2 |

| Observation | Var 1 | Var 2 | Var 3 |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |

Design of experiment (DoE)

➤ Biological repeats

➤ Methodological repeats

➤ Single phenomenon

Multiple measurements

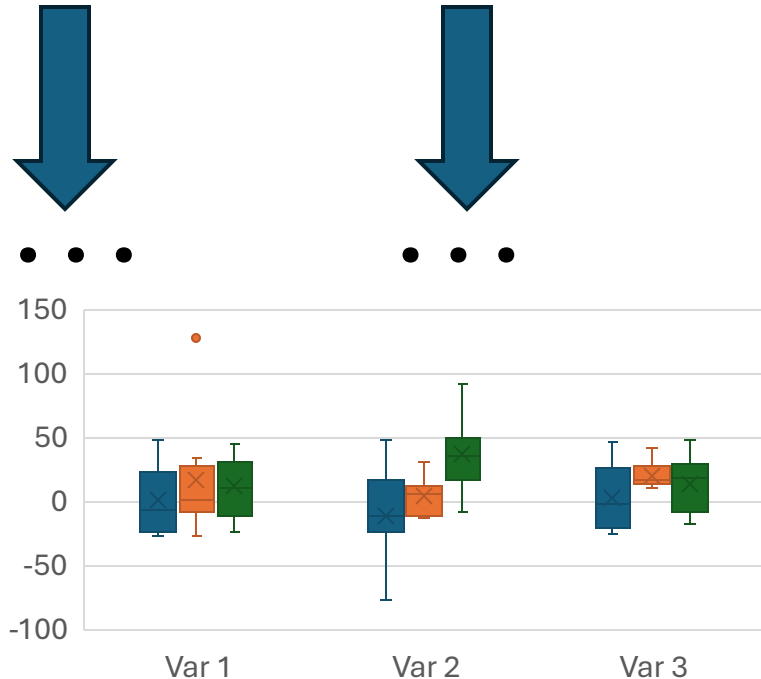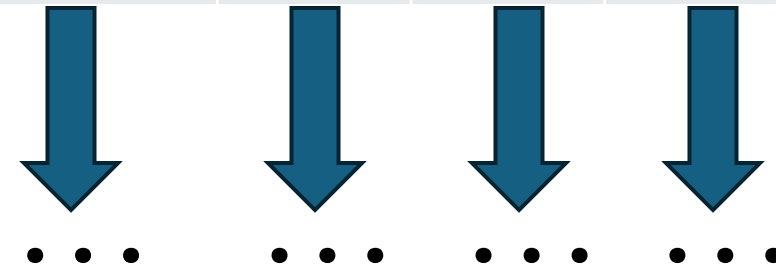➤ Precision methodology – repeats not needed

➤ Multiple phenomena

# Univariate to multivariate analysis

| observation | Var 1 | Var 2 |
|---|---|---|
| 1 | Rep 1 | Rep 1 |
| 1 | Rep 2 | Rep 2 |
| 1 | Rep 3 | Rep 3 |
| 2 | Rep 2 | Rep 2 |

| Observation | Var 1 | Var 2 | Var 3 |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |

| | | | |
|---|---|---|---|
| 1.9763 | 2 | 10.4 | 10.3 |
| 1.576 | 1 | 12 | 11.8 |
| 1.7067 | 2 | 10.1 | 10.8 |
| 1.7401 | 2 | 11.1 | 11.2 |
| 0.9642 | 2 | 10.6 | 10.8 |
| 2.0098 | 2 | 8.9 | 9.1 |
| 1.2985 | 1 | 9.1 | 9.2 |
| 1.4708 | 1 | 11.5 | 11.1 |
| 1.3306 | 1 | 11.2 | 11.1 |
| 1.3314 | 1 | 16.55 | 15.78 |
| 1.6034 | 2 | 15.11 | 15.13 |
| 2.7886 | 2 | 15.5 | 18.45 |

➢ Aggregation method: sum, counts, etc.

➢ Treating NaNs – empty vs zeros (consequences?)
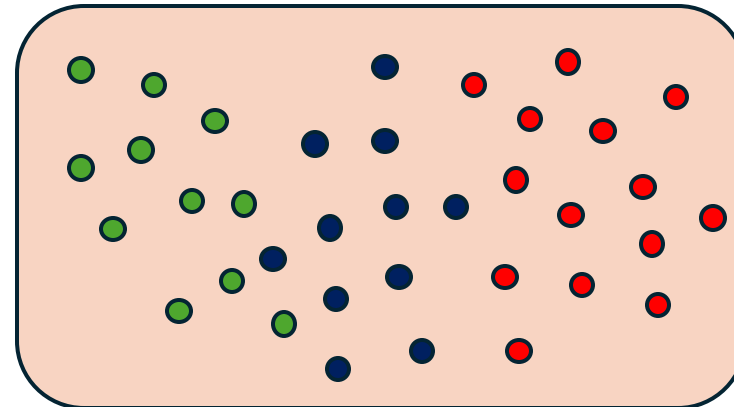
➢ Grouping

Wrangling libraries and techniques:
R - tidyverse
Python – pandas, numpy
SQL - count

**Grouping observations**
The variance across observational groups of **samples** is greater than between each sample

**Grouping variables**
The variance across observational groups of **measurements** is greater than between each measurements

➢ Aggregation method: sum, counts, etc.

➢ Treating NaNs, NULLs: empty vs zeros
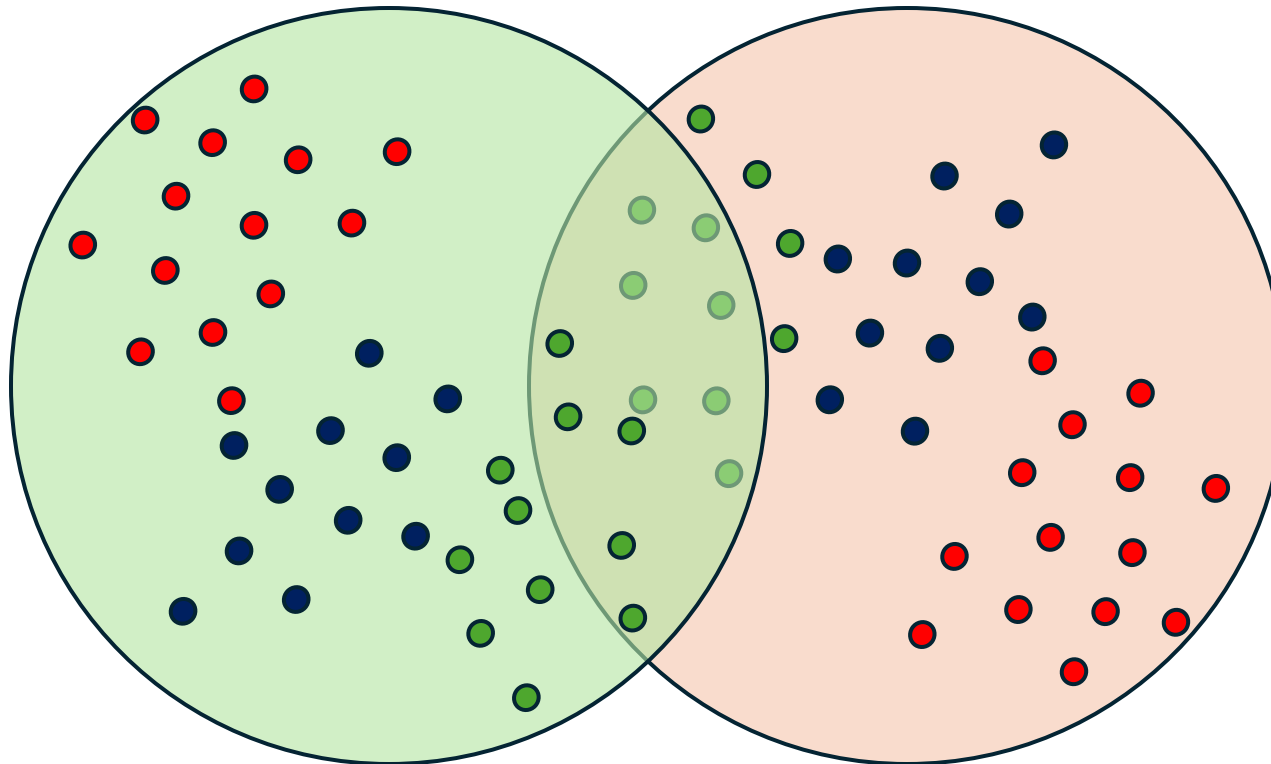
(consequences?)

➢ Grouping

**Combined effects**

Wrangling libraries and techniques:
R - tidyverse
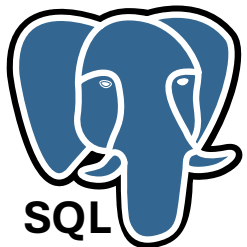Python – pandas, numpy
SQL - count

Automate data wrangling and analysis

- Captured format is different from analysis format

- Capture data in the simplest format

- Try to change data in the processing stage instead of in

  spreadsheets (for reproducibility)

# Data wrangling & Caveats in MVDA

Base libraries:

"readxl"
"writexl"

**Dimension reduction**

Some come with a common visualization, some do not. Ultimately, they are just **mathematical calculations**, applied according to the type of data (e.g., Categorical, continuous, discreet, nominal, sparse). You get to represent the results in a communicative way that elucidates insight.

**1** | Orthogonal decomposition

PCA

MDS

**Dimension reduction**

Some come with a common visualization, some do not. Ultimately, they are just **mathematical calculations**, applied according to the type of data (e.g., Categorical, continuous, discreet, nominal, sparse). You get to represent the results in a communicative way that elucidates insight.

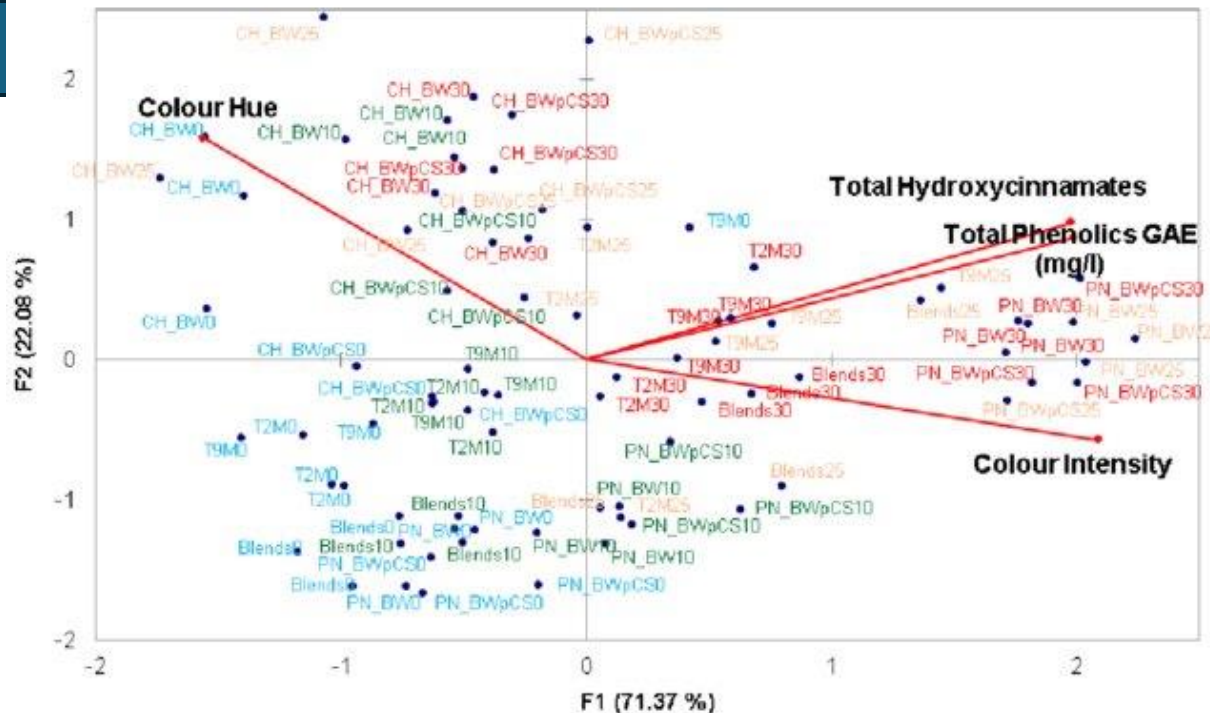**1**     Orthogonal decomposition

**CA**



**MCA**

## Dimension reduction

Some come with a common visualization, some do not. Ultimately, they are just **mathematical calculations**, applied according to the type of data (e.g., Categorical, continuous, discreet, nominal, sparse). You get to represent the results in a communicative way that elucidates insight.
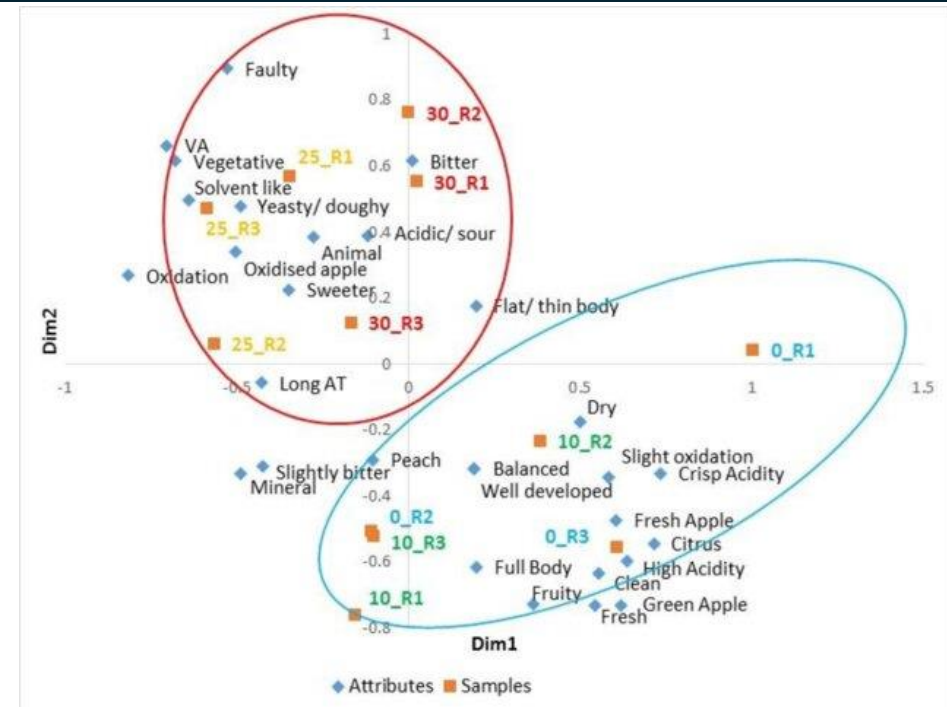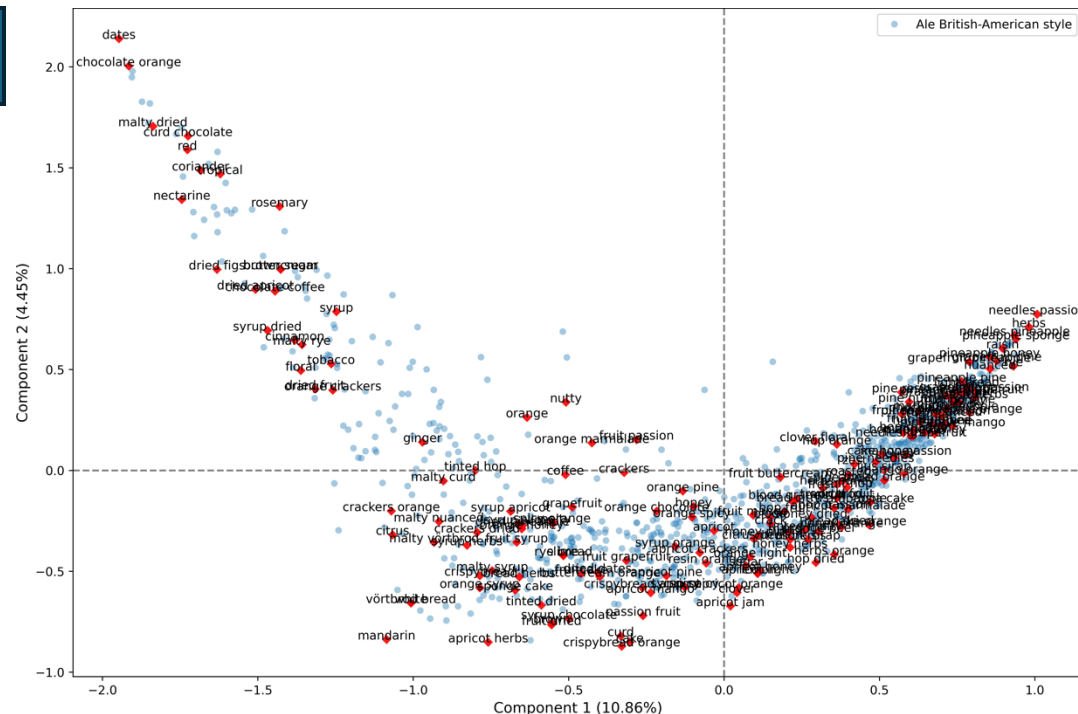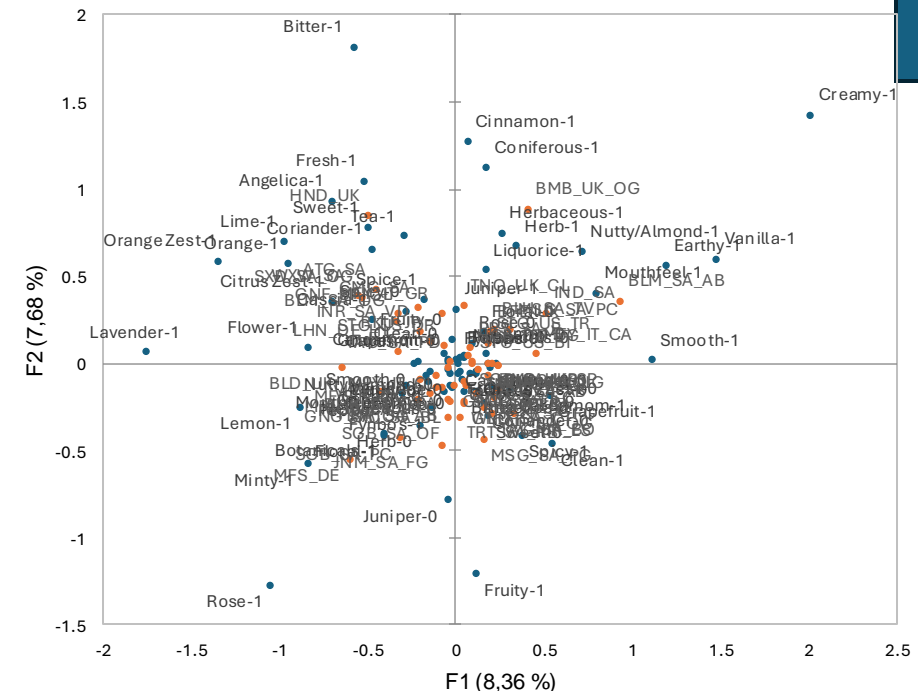
**1** Orthogonal decomposition

PCA

| Sample ID | Class | Vine age (yrs) | 3991.9861 | 3987.8664 | 3983.7467 | 3979.6270 | 3975.5073 | 3 |
|---|---|---|---|---|---|---|---|---|
| 751 | Young | 30 | -0.000378 | -0.000475 | -0.000567 | -0.000628 | -0.000635 | - |
| 752 | Young | | -0.000876 | -0.00099 | -0.001077 | -0.001102 | -0.001063 | - |
| 753 | Young | 5 | -0.001342 | -0.00143 | -0.001513 | -0.001569 | -0.00158 | - |
| 754 | Young | 30 | -0.002598 | -0.002718 | -0.002826 | -0.002887 | -0.002888 | - |
| 755 | Young | 30 | -0.003074 | -0.003125 | -0.00318 | -0.003223 | -0.003229 | - |
| 756 | Young | 20 | -0.003272 | -0.003334 | -0.003394 | -0.003429 | -0.003418 | - |
| 757 | Young | 30 | -0.003459 | -0.00354 | -0.003611 | -0.00366 | -0.003683 | - |
| 758 | Old | 35 | -0.003599 | -0.003673 | -0.00375 | -0.003806 | -0.003816 | - |
| 759 | Young | 29 | -0.000259 | -0.000324 | -0.000393 | -0.000443 | -0.00046 | - |
| 760 | Old | 40 | -0.000462 | -0.000488 | -0.000517 | -0.000546 | -0.00057 | - |

CA

| Primary ID | Old | Teenager | Young | Textured | Structured | Robust | Rich | Ripe | Nutty | Wood |
|---|---|---|---|---|---|---|---|---|---|---|
| 751 | 13 | 1 | 12 | 3 | 0 | 1 | 1 | 3 | 1 | 0 |
| 752 | 11 | 2 | 17 | 1 | 1 | 1 | 1 | 2 | 0 | 0 |
| 753 | 17 | 2 | 10 | 5 | 0 | 1 | 4 | 5 | 2 | 0 |
| 754 | 19 | 1 | 11 | 4 | 0 | 2 | 3 | 6 | 2 | 0 |
| 755 | 18 | 1 | 11 | 5 | 3 | 1 | 5 | 5 | 3 | 0 |
| 756 | 17 | 1 | 11 | 1 | 1 | 2 | 5 | 5 | 3 | 0 |
| 757 | 18 | 2 | 9 | 2 | 0 | 2 | 4 | 2 | 3 | 0 |
| 758 | 11 | 2 | 16 | 4 | 0 | 0 | 5 | 3 | 3 | 1 |
| 759 | 8 | 3 | 18 | 3 | 0 | 0 | 1 | 4 | 0 | 0 |
| 760 | 24 | 1 | 5 | 5 | 2 | 2 | 8 | 7 | 4 | 0 |
| 761 | 13 | 4 | 12 | 2 | 1 | 1 | 2 | 5 | 1 | 1 |
| 762 | 20 | 1 | 10 | 3 | 3 | 2 | 4 | 7 | 1 | 0 |
| 763 | 16 | 4 | 11 | 4 | 1 | 0 | 2 | 1 | 3 | 0 |
| 764 | 22 | 0 | 8 | 5 | 2 | 0 | 5 | 4 | 3 | 0 |
| 765 | 1 | 4 | 18 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

MDS

| Primary ID | 751 | 752 | 753 | 754 | 755 | 756 | 757 | 758 | 759 |
|---|---|---|---|---|---|---|---|---|---|
| 751 | 32 | 13 | 13 | 13 | 14 | 9 | 13 | 14 | 14 |
| 752 | 13 | 32 | 11 | 16 | 11 | 14 | 16 | 15 | 15 |
| 753 | 13 | 11 | 32 | 16 | 12 | 12 | 14 | 12 | 14 |
| 754 | 13 | 16 | 16 | 32 | 12 | 16 | 20 | 13 | 13 |
| 755 | 14 | 11 | 12 | 12 | 32 | 13 | 8 | 11 | 16 |
| 756 | 9 | 14 | 12 | 16 | 13 | 32 | 19 | 13 | 7 |
| 757 | 13 | 16 | 14 | 20 | 8 | 19 | 32 | 13 | 7 |
| 758 | 14 | 15 | 12 | 13 | 11 | 13 | 13 | 32 | 16 |
| 759 | 14 | 15 | 14 | 13 | 16 | 7 | 7 | 16 | 32 |

MCA

| A Typical Old Vine CB Word Association | B Judge 1 | C Judge 2 | D Judge 3 | E Judge 4 | F Judge 5 | G Judge 6 | H Judge 7 | I Judge 8 | J Judge 9 | K Judge 10 | L Judg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Concentration | 1 | | 1 | | | | | | | | |
| Rich | | | | | | | | 1 | 1 | | 1 |
| Balanced | 1 | | 1 | | | 1 | | | | | |
| Complexity | | | | | | 1 | 1 | 1 | | | 1 |
| Long AT | | | | | 1 | 1 | | | | | 1 |
| Full bodied | | | | | | | | 1 | 1 | 1 | |
| Minerality | | | | 1 | | 1 | | | | | |
| Tropical | | | | | | 1 | | 1 | | | |
| Structure | | 1 | | | | | 1 | | | 1 | |
| Good mouthfeel | | | | | | 1 | | 1 | | | |
| Fruity | 1 | | | | | | 1 | | | | |
| Depth | | | 1 | | | | | | | | |
| Stone fruit | | | | | | | | | | | |
| Oily | | 1 | | | | | | | | | |
| Marketing | | | | | | | | | | | |
| Potential | | 1 | | | | | | | | | 1 |
| Yellow fruit | | | | | | | | | | | |

➢ Aggregation method: sum, counts, etc.
➢ Treating NaNs – empty vs zeros (consequences?)
➢ Grouping

Wrangling libraries and techniques:
R - tidyverse
Python – pandas, numpy
SQL - count

**Grouping observations**
The variance across observational groups of **samples** is greater than between each sample

**Grouping variables**
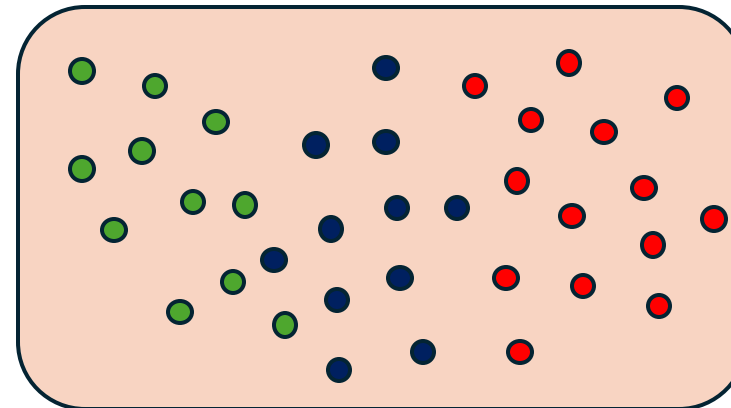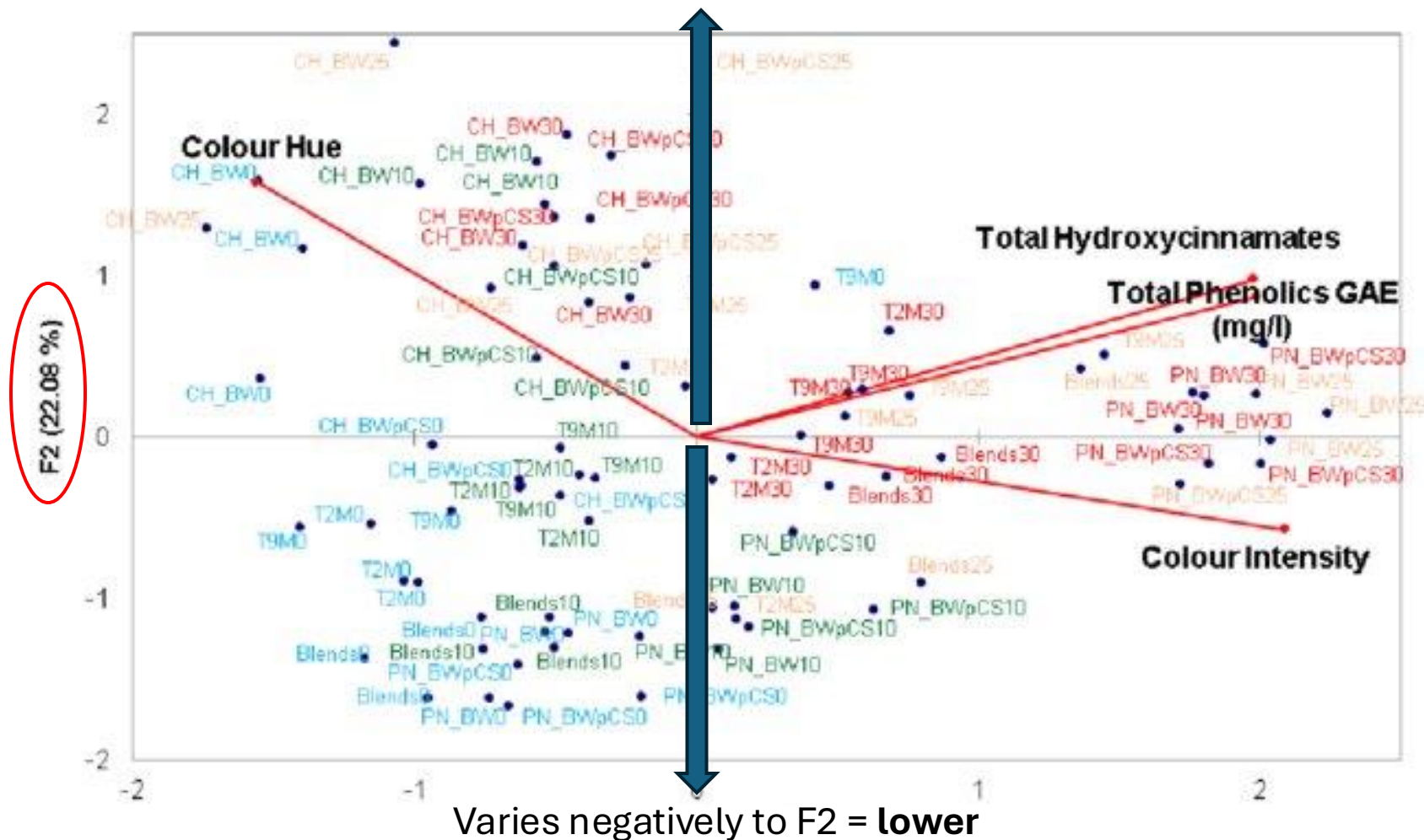The variance across observational groups of **measurements** is greater than between each measurements

# Common MVDA in wine sciences

**Reading biplots/cartesian covariance representations**

## Cluster analysis

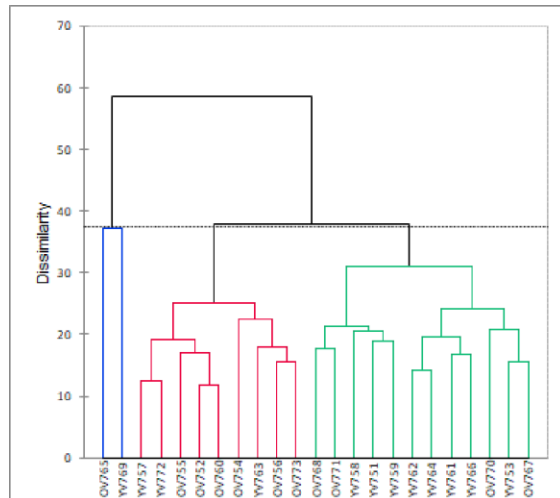Often calculated on results from dimension reduction (recommended!).
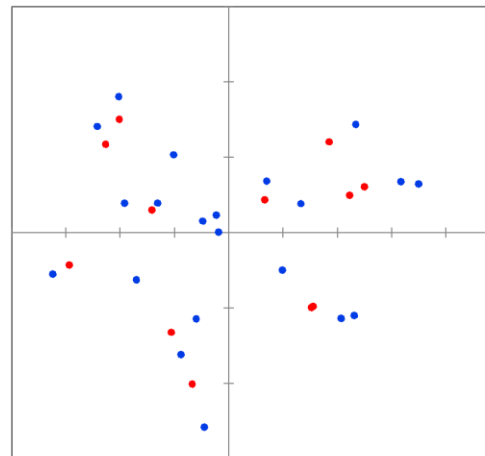
*e.g.,* Hierarchical clustering (HCA),
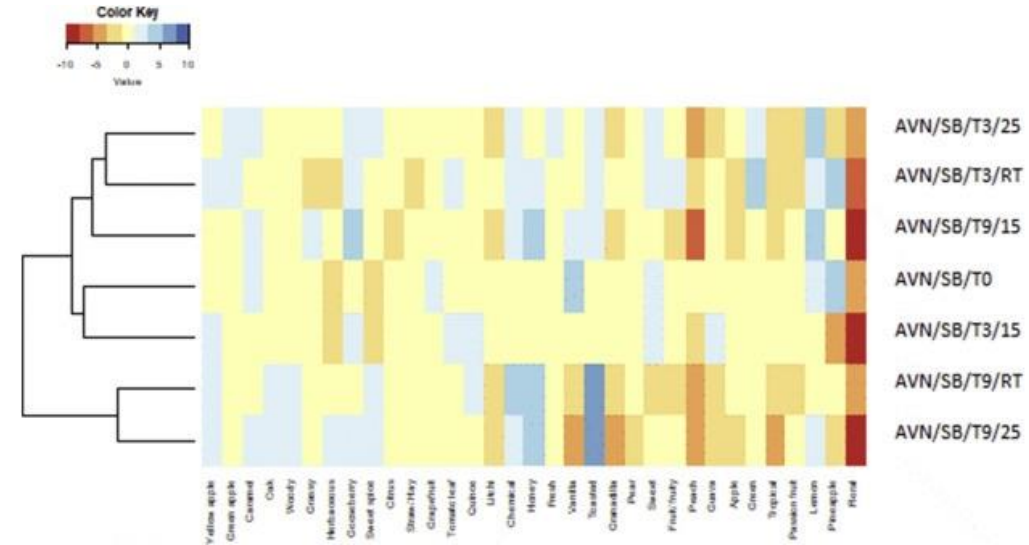
**2** | Cluster analysis

Number of clusters: 9

Cluster: 1 755
Cluster: 1 756
Cluster: 1 761
Cluster: 1 762
Cluster: 1 764
Cluster: 1 766
Cluster: 1 767
Cluster: 1 770
Cluster: 2 752
Cluster: 2 754
Cluster: 2 757
Cluster: 2 758
Cluster: 2 759
Cluster: 2 760
Cluster: 3 763
Cluster: 3 771
Cluster: 3 773
Cluster: 4 765
Cluster: 5 769
Cluster: 6 753
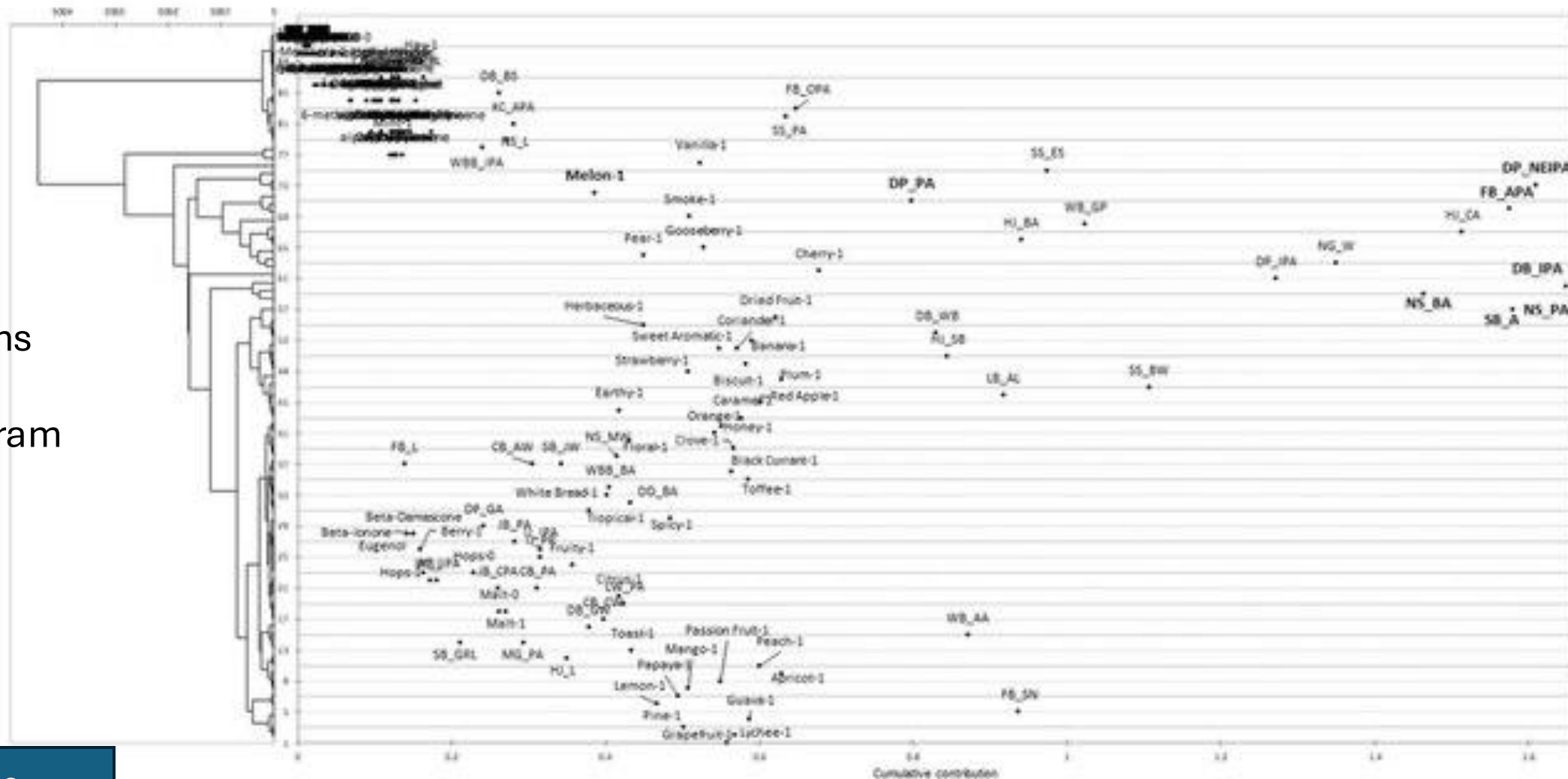Cluster: 7 768
Cluster: 8 772
Cluster: 9 751



Dendrogram



scatterplot



heatmap

**Common MVDA in wine sciences**

PCA correlations
+
Cluster dendrogram

Just a scatterplot, so . . .

# Thank you!

https://github.com/mpho-mafata

https://www.linkedin.com/in/mafatampho

https://orcid.org/0000-0002-6468-7193

**A chemometric approach to investigating South African wine behaviour using chemical and sensory markers**

by

**Mpho Mafata**

Dissertation presented for the degree of
**Doctor of Philosophy (Agricultural Sciences)**

at
**Stellenbosch University**
Department of viticulture and Oenology, Faculty of AgriSciences

*Supervisor:* Dr. Astrid Buica
*Co-supervisors:* Dr. Jeanne Brand and Prof. Andrei V. Medvedovici

March 2021

https://doi.org/10.13140/RG.2.2.17899.00804