# Multivariate Analysis Module I: A Bird's Eye View

Dr. Mark Williamson

DaCCoTA

University of North Dakota

# Introduction

## What is Multivariate analysis?

- "multiple dependent variables resulting in one outcome" [1]
- "statistical models that have two or more dependent or outcome variables" [2]
- "matrices" [3]

## Confusions

- Sometimes uses interchangeably with multivariable [2]
- Other classifications ANOVA, logistic regression, etc.
- 'Multivariate regression' also thrown around wildly

**Multivariate:**
multiple dependent variables or other more complicated structures such as ordination or non-linearity

$$(1)\ y = \alpha + x\beta + \varepsilon$$

$$(2)\ y = \alpha + x_1\beta_1 + x_2\beta_2 + \ldots + x_k\beta_k + \varepsilon$$

$$(3)\ Y_{n \times p} = X_{n \times (k+1)}\ \beta_{(k+1) \times p} + \varepsilon$$

- **Single response and single predictor (simple regression)**
- **Single response and multiple predictors (multiple regression)**
- **Multiple responses and predictors (multivariate regression)**
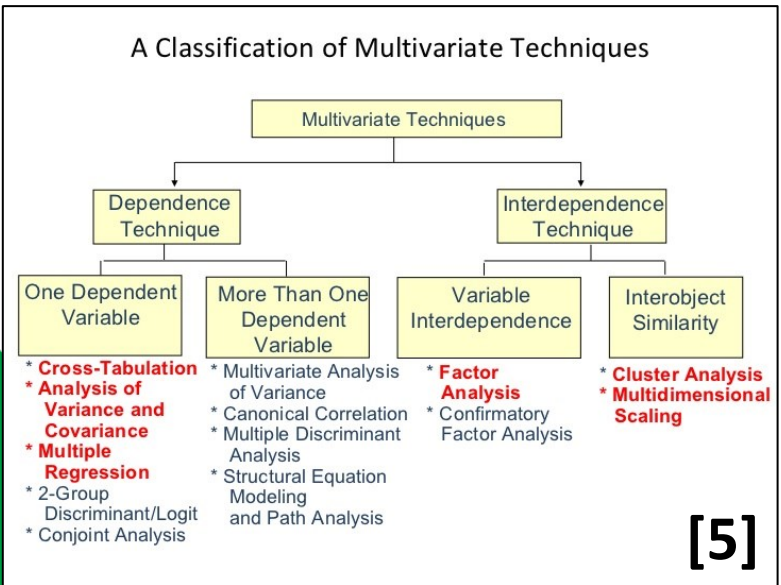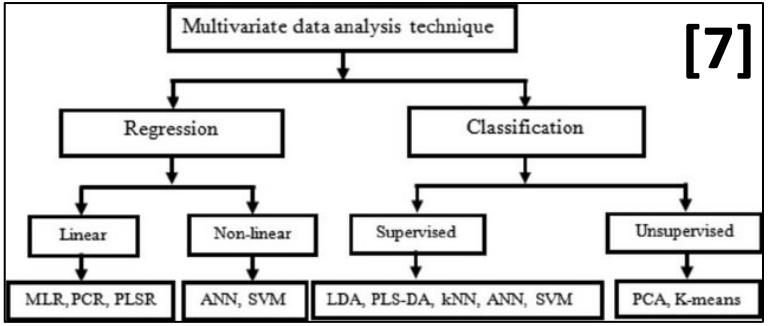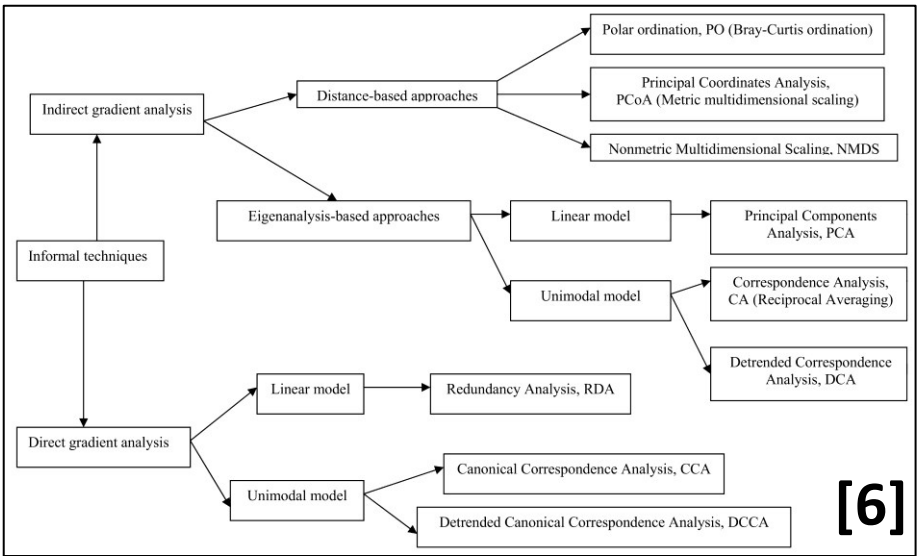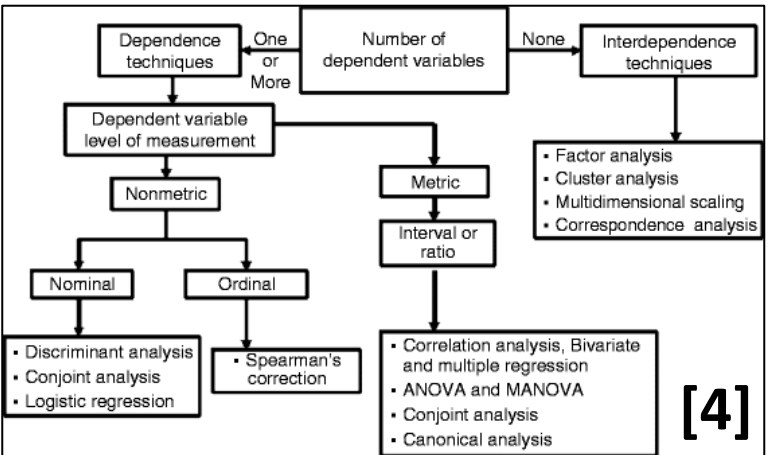
# Classification Problems

# Landscape

**Multivariate Analysis**

**Regression & ANOVA extensions**
- ❖ MANOVA
- ❖ MANCOVA
- ❖ Multivariate Linear Regression
- ❖ Artificial Neural Networks
- ❖ Support Vector Machines
- ❖ Conjoint Analysis

**Grouping Methods**
- ❖ Cluster Analysis
- ❖ Recursive Partitioning

**Canonical Analysis**
- ❖ Canonical Correlation Analysis
- ❖ Linear Discriminant Analysis
- ❖ Redundancy Analysis
- ❖ Canonical Correspondence Analysis

**Simultaneous & Structural Equations**
- ❖ Simultaneous Equation Modeling
- ❖ Structural Equation Modeling

**Dimensional Reduction**
- ❖ Factor Analysis
- ❖ Principal Components Analysis
- ❖ Principal Coordinates Analysis
- ❖ Correspondence Analysis

**Time Based**
- ❖ Vector Autoregression
- ❖ Principal Response Curves

# Structures and Uses

## Regression & ANOVA extensions

**MANOVA:**
◊ Multivariate Analysis of Variance [10]
◊ Extension of ANOVA with multiple Y variables
◊ Tests for the difference in two or more vectors

**MANCOVA:**
◊ Multivariate Analysis of Co-Variance [10]
◊ Extension of ANCOVA with multiple Y variables
◊ Covariate accounted for to reduce noise

**Multivariate Linear Regression:**
◊ Extension of Linear Regression with multiple Y variables

**Support Vector Machines:**
◊ Extends regression into non-linearity [12]
◊ Machine learning algorithms
◊ Supervised learning model
◊ Fast and less requirements to train

**Artificial Neural Networks:**
◊ Same base characteristics as SVM
◊ Simulates how brains process (collection of connected nodes)

**Conjoint Analysis:**
◊ Extension of dummy variable regression [15]
◊ Survey-based analysis method
◊ Helps extract consumer preferences
◊ Sub-techniques include Choice-based conjoint, Adaptive conjoint, MaxDiff conjoint, etc.

[11]


TWO-WAY MANOVA EXAMPLE


MANCOVA EXAMPLE

[13]


[14]


[15]


[12]

A linear decision boundary


A non-linear decision boundary

# Structures and Uses

**Canonical Analysis**

**Canonical Correlation Analysis:**
◊ Symmetrical [16] and linear [17]
◊ Method where two sets of variables have maximized correlation

**Linear Discriminant Analysis:**
◊ Non-symmetric, Y is class of objects
◊ Essentially an extension of multinomial regression

**Redundancy Analysis:**
◊ Non-symmetric, linear model, Euclidean distances
◊ Related to multiple linear regression, ordination added on [20]

**Canonical Correspondence Analysis:**
◊ Non-symmetric, unimodal model, Chi-squared distances
◊ Related to multiple linear regression, ordination added on



(a) Simple ordination of matrix **Y**: principal comp. analysis (PCA) correspondence analysis (CA)

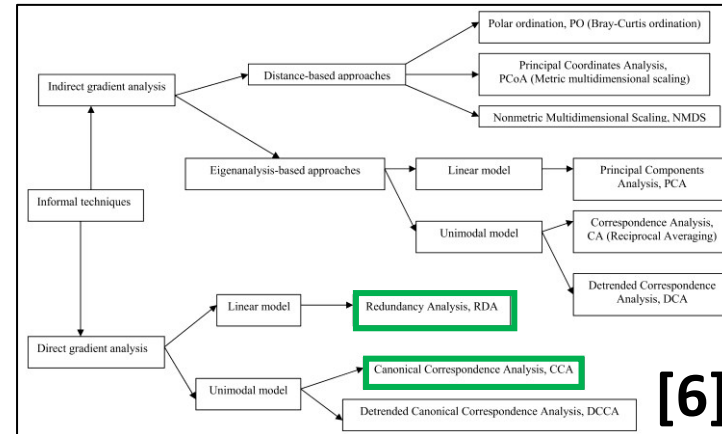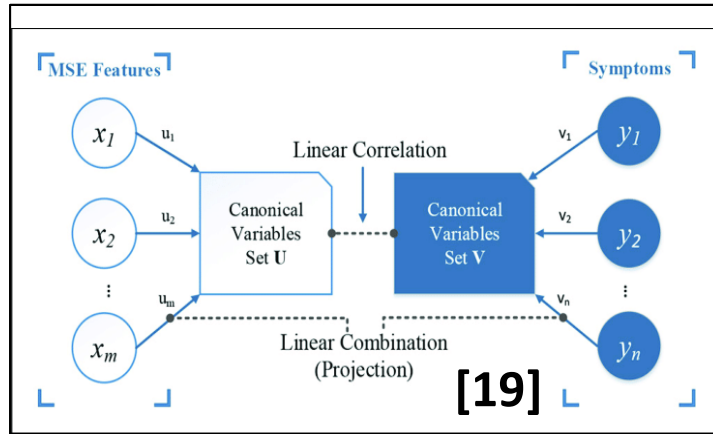(b) Ordination of **y** (single axis) under constraint of **X**: multiple regression

Model: $\hat{y} = b_0 + b_1x_1 + \ldots + b_mx_m$

(c) Ordination of **Y** under constraint of **X**: redundancy analysis (RDA) canonical correspondence analysis (CCA)

[16]

[19]

[6]

Decision tree:
- Symmetric?
  - Yes → Canonical Correlation Analysis
  - No → Y-variable categorical?
    - Yes → Linear Discriminant Analysis
    - No → Distances?
      - Euclidean → Redundancy Analysis
      - Chi-squared → Canonical Correspondence Analysis

# Structures and Uses

**Dimensional Reduction**

**Factor Analysis:**
◊ Interdependence method that used ordination [21]
◊ Reducing observable variables into latent factors
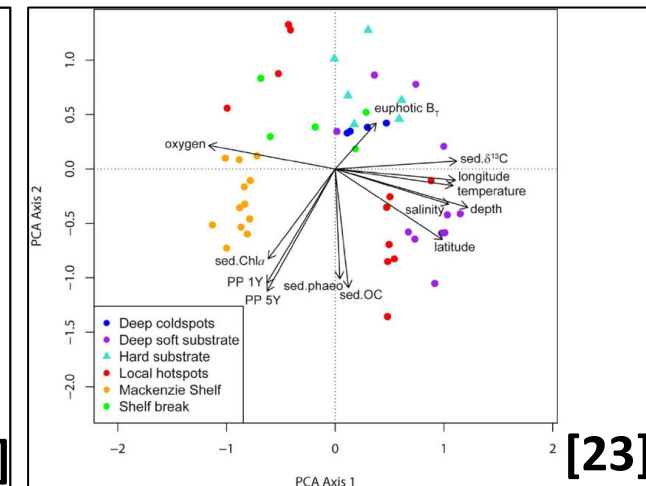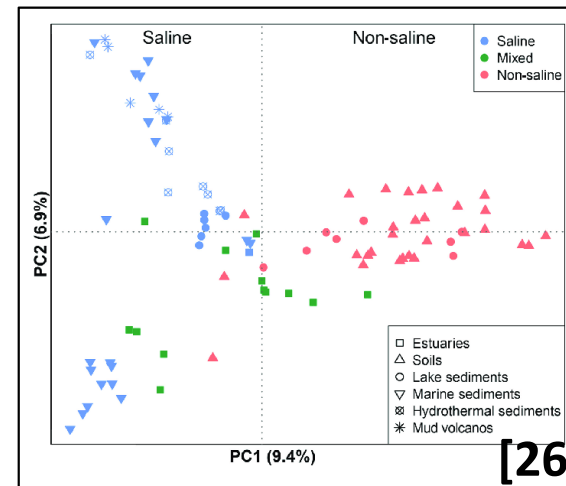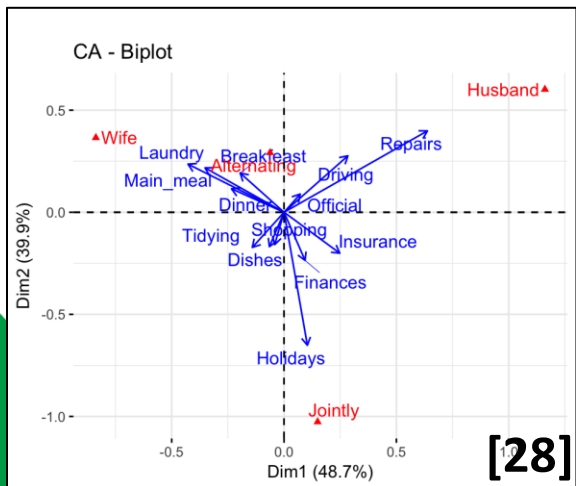◊ Both exploratory and confirmatory

**Principal Components Analysis (PCA):**
◊ Ordination technique
◊ Arranges data along gradients to simply
◊ Reduces dimensionality and still preserves most info
◊ Only highly correlated variables are together


[22]


[22]

**Principal Coordinates Analysis (PCoA):**
◊ Ordination technique
◊ Arranges data along gradients to simply
◊ Represents distances between samples in low dimensional space [24]
◊ Better for missing data or fewer individuals than characteristics [25]

**Correspondence Analysis:**
◊ Like PCA but applies categorical rather than continuous data
◊ Applies contingency tables [27]


[28]


[26]


[23]

# Structures and Uses
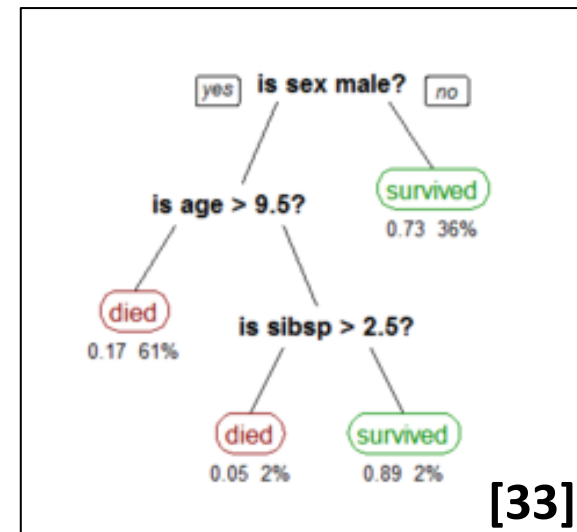
**Cluster Analysis:**
- ◊ Grouping set of objects so that each group is more similar to members than non-members [29]
- ◊ Difficult computational problem
- ◊ Specific types include hierarchical, k-means, distribution-based, density-based, etc. [30]

**Recursive Partitioning:**
- ◊ Inverse of clustering [32]
- ◊ Creates a decision tree to attempt to correctly classify members of the population based on several dichotomous dependent variables (yes/no)
- ◊ Part of the more general technique of decision trees
- ◊ Intuitive models that can be tweaked for sensitivity or specificity but can overfit data [33]

**Grouping Methods**



[31]



[33]

# Structures and Uses

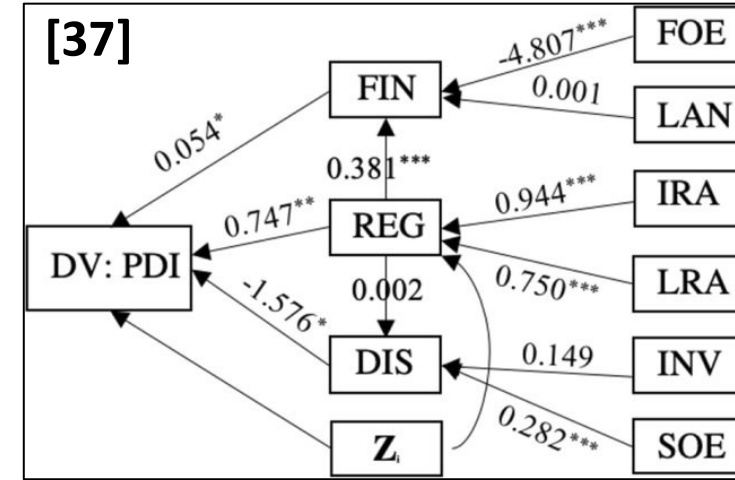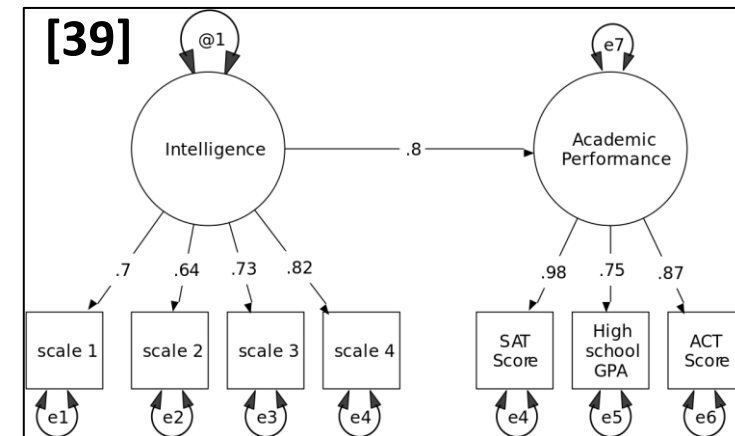**Simultaneous Equation Modeling:**
- ◊ System of two or more linear simultaneous equations [34]
- ◊ Variables used are classified as endogenous (jointly determined, dependent) and exogenous (predetermined, independent) [35]
- ◊ Endogenous influenced by exogenous by not the other way round
- ◊ Complete model is when the number of endogenous variables equals the number of equations and is a Structural Equation Model
- ◊ Solution is determined by equilibrium among opposing forces [36]
- ◊ Structural form includes multiple endogenous variables; reduced form only has one; reduced can be estimated by least squares
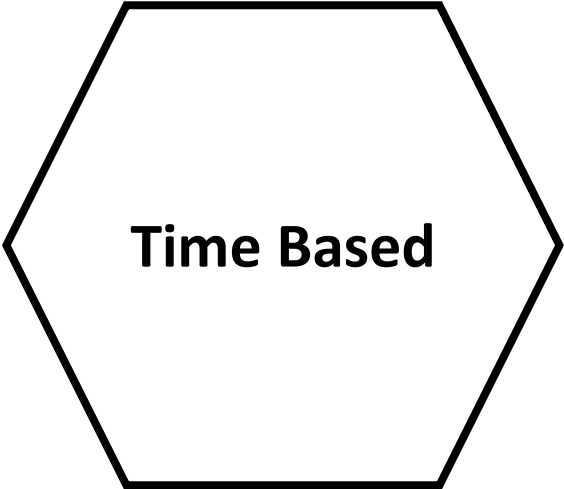
**Structural Equation Modeling:**
- ◊ Combination of factor analysis and regression [38]
- ◊ Interest is usually on latent factors that underlie observable variables
- ◊ Can be used to impute relationship between those latent factors from the observable variables
- ◊ Includes confirmatory factor analysis, confirmatory composite analysis, path analysis, partial least squares modeling, and latent growth modeling [29]
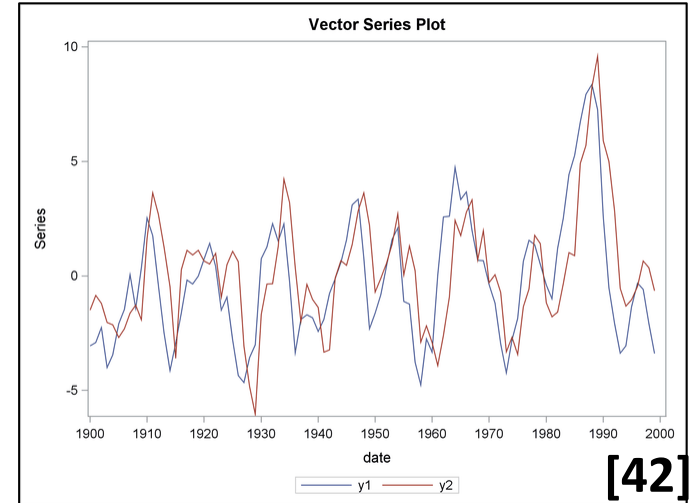
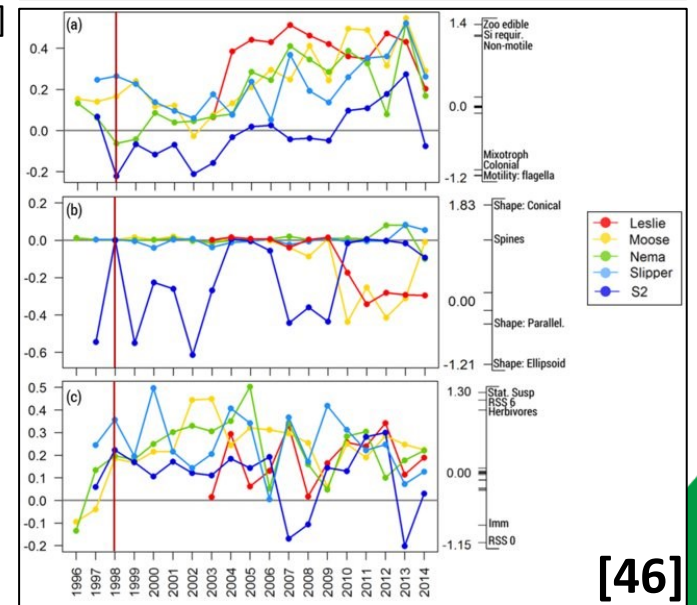**Simultaneous & Structural Equations**

# Structures and Uses

**Vector Autoregression:**
◊   Extension of univariate autoregressive model
◊   Stochastic process model that tries to understand the change in multiple quantities over time [40]
◊   Used when two or more time series influence each other [41]
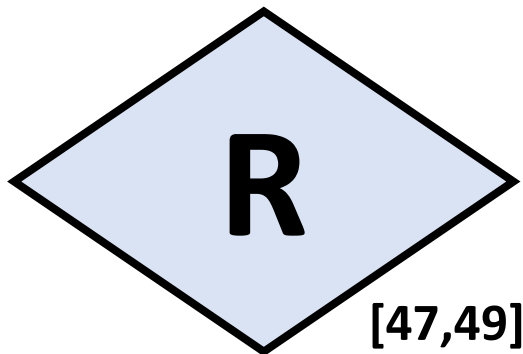◊   Each variable has an equation modeling its change over time, including past (lagged) values
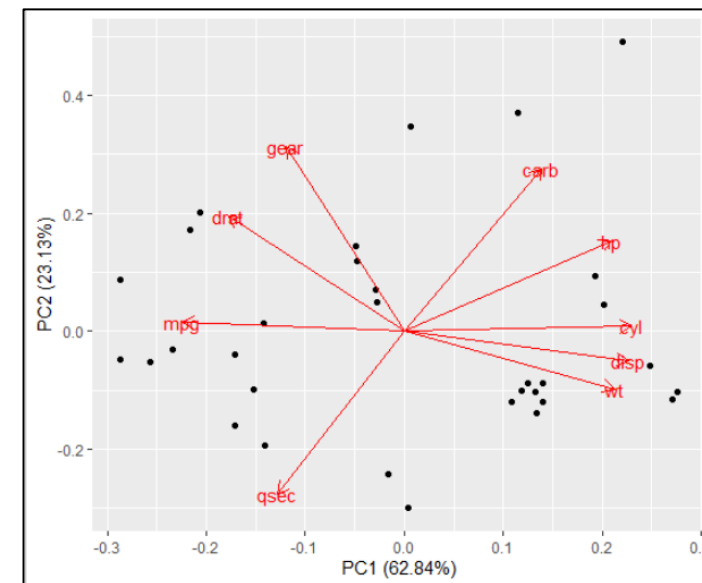


[42]

**Principal Response Curves:**
◊   Used to analysis of treatment effects in repeated measures [43]
◊   Corrects for temporal trends in control treatments
◊   Special kind of redundancy analysis [44]
◊   Allows for summarizing and plotting of the results that is much more interpretable than a bi-plot [45]

**Time Based**



[46]

# Examples 📊



```
library(ggfortify)

head(mtcars)
mtcars2 <-mtcars[,c(1:7,10,11)] #remove cat vars (vs & am)

PCA1 <-prcomp(mtcars2, center=TRUE, scale.=TRUE)

summary(PCA1) #PC1 is 63% of var, PC2 is 23%

autoplot(PCA1, loadings=TRUE, loadings.label=TRUE)
```
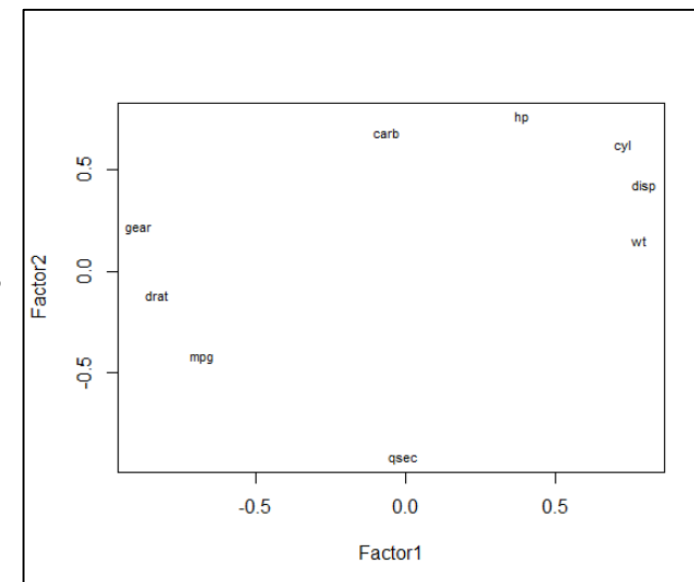
**R**

**[47,49]**

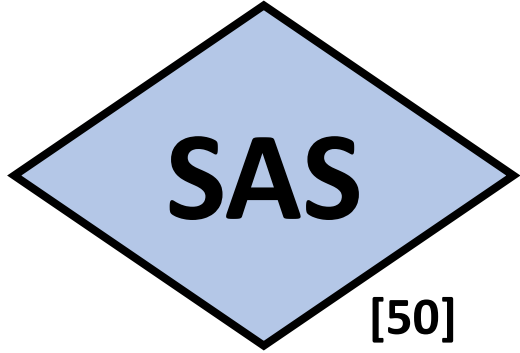**Importance of components:**

|  | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 |
|---|---|---|---|---|---|---|---|---|---|
| **Standard deviation** | 2.3782 | 1.4429 | 0.71008 | 0.51481 | 0.42797 | 0.35184 | 0.32413 | 0.2419 | 0.14896 |
| **Proportion of Variance** | 0.6284 | 0.2313 | 0.05602 | 0.02945 | 0.02035 | 0.01375 | 0.01167 | 0.0065 | 0.00247 |
| **Cumulative Proportion** | 0.6284 | 0.8598 | 0.91581 | 0.94525 | 0.96560 | 0.97936 | 0.99103 | 0.9975 | 1.00000 |

```
FA1 <-factanal(mtcars2, 3, rotation="varimax")
print(FA1, digits=2, cutoff=.3, sort=TRUE)

FA1_load <-FA1$loadings[,1:2]
plot(FA1_load, type='n')
text(FA1_load, labels=names(mtcars2), cex=0.7)
```

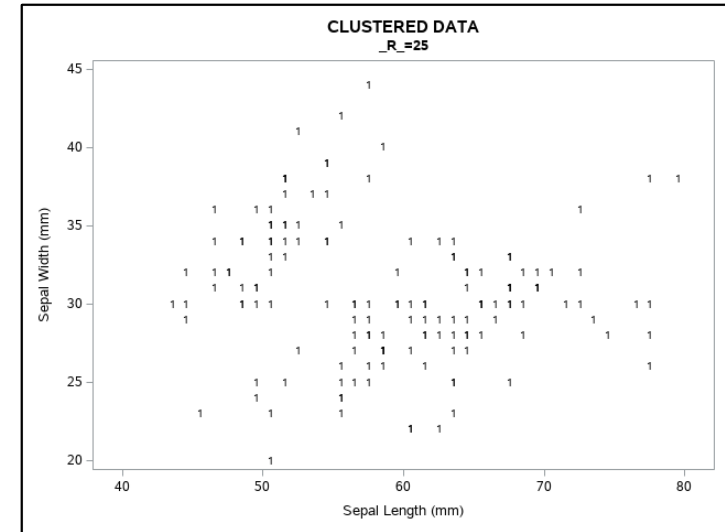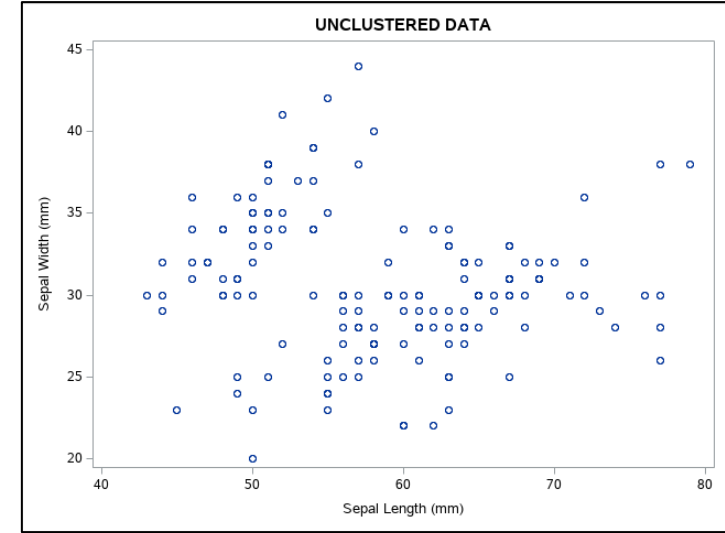|  | Factor1 | Factor2 | Factor3 |
|---|---|---|---|
| **SS loadings** | 3.85 | 2.74 | 1.28 |
| **Proportion Var** | 0.43 | 0.30 | 0.14 |
| **Cumulative Var** | 0.43 | 0.73 | 0.87 |

# Examples

## SAS [50]

```
PROC MODECLUS data=SASHELP.IRIS
    method=1 r=(5 10 25) out=out;
PROC SGPLOT;
    scatter y=SEPALWIDTH x=SEPALLENGTH;
    title'UNCLUSTERED DATA';
PROC SGPLOT data=out;
    scatter y=SEPALWIDTH x=SEPALLENGTH/ markerchar=cluster;
    by _R_;
    title'CLUSTERED DATA';
```
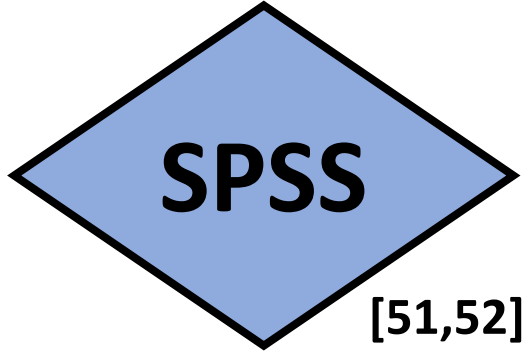
### The MODECLUS Procedure — R=5 METHOD=1

**Cluster Statistics**

| Cluster | Frequency | Maximum Estimated Density | Boundary Frequency | Estimated Saddle Density |
|---|---|---|---|---|
| 1 | 49 | 0.00007565 | 0 | . |
| 2 | 22 | 0.00003458 | 9 | 0.00003026 |
| 3 | 35 | 0.00003242 | 10 | 0.00002162 |
| 4 | 27 | 0.00003026 | 3 | 0.00002162 |
| 5 | 4 | 8.64607E-6 | 0 | . |
| 6 | 3 | 6.48456E-6 | 0 | . |
| 7 | 2 | 4.32304E-6 | 0 | . |
| 8 | 2 | 4.32304E-6 | 0 | . |
| 9 | 1 | 2.16152E-6 | 0 | . |
| 10 | 1 | 2.16152E-6 | 0 | . |
| 11 | 1 | 2.16152E-6 | 0 | . |
| 12 | 1 | 2.16152E-6 | 0 | . |
| 13 | 1 | 2.16152E-6 | 0 | . |
| 14 | 1 | 2.16152E-6 | 0 | . |

### The MODECLUS Procedure — R=10 METHOD=1

**Cluster Statistics**

| Cluster | Frequency | Maximum Estimated Density | Boundary Frequency | Estimated Saddle Density |
|---|---|---|---|---|
| 1 | 100 | 7.8355E-6 | 0 | . |
| 2 | 50 | 6.48456E-6 | 0 | . |

### The MODECLUS Procedure — R=25 METHOD=1

**Cluster Statistics**

| Cluster | Frequency | Maximum Estimated Density | Boundary Frequency | Estimated Saddle Density |
|---|---|---|---|---|
| 1 | 150 | 3.66594E-7 | 0 | . |

### The MODECLUS Procedure

**Cluster Summary**

| R | Number of Clusters | Frequency of Unclassified Objects |
|---|---|---|
| 5 | 14 | 0 |
| 10 | 2 | 0 |
| 25 | 1 | 0 |



UNCLUSTERED DATA



CLUSTERED DATA _R_=25

# Examples

# Summary and Conclusion

- Multivariate analysis is a series of advanced methods that typically feature multiple predictor (Y) variables

- Features such as ordination, non-linearity, or repeated measures are also common

- Uses include data simplification, data exploration, and hypothesis testing

- Tune in next time for a more detailed look at multivariate analysis in Multivariate Analysis Module II: Leaves and Trees

# Assessment & Acknowledgements

- Please take the 5-question assessment at:

    https://und.qualtrics.com/jfe/form/SV_eLhvOwQNlPAYNPE

- References cited in this presentation are available here:

    multivariate_analysis_module_1_refs

- The DaCCoTA is supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729.

- For the labs that use the Biostatistics, Epidemiology, and Research Design Core in any way, including this Module, please acknowledge us for publications. ***"Research reported in this publication was supported by DaCCoTA (the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729)".***