

UPPSALA UNIVERSITET

FÖRELÄSNINGSKOMMENTARER

# Multivariate Methods

*Rami Abou Zahra*

Inlämningsdatum  
January 16, 2024

## CONTENTS

1. Introduction	4
1.1. MANOVA	4
1.2. Regression analysis	4
2. Sample & Random Matrices	5
2.1. Slide 3 - Expectation	5
2.2. Slide 4 - Covariance Matrix	5
2.3. Slide 5 - Covariance Matrix	5
2.4. Slide 6 - Linear Combination	5
2.5. Slide 7 - Linear Combination	6
2.6. Slide 9 - Independence	6
2.7. Slide 10 - Random Sample	6
2.8. Slide 12 - Some Notes on Sample Covariance Matrix	7
2.9. Slide 17 - Sample Covariance Matrix	7
3. Multivariate Normal Distribution	8
3.1. Slide 4-5 - From Univariate to Multivariate Normal	8
3.2. Slide 6 - Special Case: Bivariate Normal	8
3.3. Slide 7 - Contour of Bivariate Normal Density	8
3.4. Slide 8 - Linear Combinations	8
3.5. Slide 10 - Normal and Chi-Square	8
3.6. Slide 11 - Subset of Variables	9
3.7. Slide 12 - Example: Subset of Variables	9
3.8. Slide 13 - Subset of Variables	9
3.9. Slide 15 - Marginal Normal and Joint Distribution	9
3.10. Slide 23 - Likelihood of Normal Random Sample	10
3.11. Slide 32 - Limit of MLE	10
4. Inference for Several Sample	11
4.1. Slide 3 - Paired Data	11
4.2. Slide 9 - Two Populations	11
4.3. Slide 10 - Pooled Sample Covariance Matrix	11
4.4. Slide 20 - MANOVA Model	11
4.5. Slide 28 - Multivariate Two-Way Fixed Effects Model with Interaction	12
4.6. Slide 33 - Test of No Interaction	12
5. Regression	13
5.1. Slide 6 - Classic Linear Regression	13
5.2. Slide 7 - Matrix Notation	13
5.3. Slide 9 - ANOVA With $g = 2$	13
5.4. Slide 10 - Anova With $g = 2$ and $b = 2$	13
5.5. Slide 11 - Ordinary Least Squares	13
5.6. Slide 12 - OLS Estimator	13
5.7. Slide 15 - Sampling Properties of OLS Estimators	14
5.8. Slide 17 - Distribution of Regression Coefficients	14
5.9. Slide 18 - Confidence Region	14
5.10. Slide 19 - Confidence interval	14
5.11. Slide 20 - More Than One Responses	14
5.12. Slide 22 - Assumptions	14
5.13. Slide 23 - Least Squares	15
5.14. Slide 27 - Regression Coefficients With Zero Constraints	15
5.15. Slide 31 - LRT when $m = 1$	15
5.16. Slide 32 - Prediction of Regression Function	15
5.17. Slide 34 - Forecast New Response	15
6. Principal Component Analysis	16
6.1. Slide 3 - Motivation	16
6.2. Slide 4 - Task of Principal Component Analysis (PCA)	16
6.3. Slide 5 - Restriction	16
6.4. Slide 6/7 - Principal Components and Two useful Lemmas	16
6.5. Slide 7 - Two useful Lemmas	16
6.6. Slide 9 - Principal Components	16

6.7.	Slide 10 - Total Variation Explained by Principal Components	17
6.8.	Slide 12 - Principal Components From Correlation Matrix	17
6.9.	Slide 14 - Sample Principal Components	17
7.	Factor Analysis	18
7.1.	Slide 3 - Latent Variable Modelling	18
7.2.	Slide 4 - The Model	18
7.3.	Slide 5 - Scale Indeterminacy	18
7.4.	Slide 7 - Model Implied Covariance Matrix	18
7.5.	Slide 8 - Existence of Decomposition	18
7.6.	Slide 9 - Indeterminacy	19
7.7.	Slide 10 - Scale Invariant	19
7.8.	Slide 11 - Popular Estimation Methods	19
7.9.	Slide 12 - Spectral Decomposition	19
7.10.	Slide 14 - Determine Number of Factors	19
7.11.	Inbetween lectures	19
7.12.	Slide 23 - Versus PCA	19
7.13.	Slide 24 - Orthogonal Rotation	19
7.14.	Slide 26 - Oblique Rotation	19
7.15.	Slide 28 - Bartlett Score	19
8.	Canonical Correlation Analysis	20
8.1.	Slide 3 - Motivation	20
8.2.	Slide 4 - Task	20
8.3.	Slide 5 - Correlation Coefficient	20
8.4.	Slide 7 - Canonical Variates	20
8.5.	Slide 9 - Find Canonical Variates	20
8.6.	Slide 12 - Scale Invariant: Coefficient Vector	21
8.7.	Slide 14 - Proportion of Explained Variance	21
8.8.	Slide 15 - Sample Canonical Variate Pair	21
8.9.	Slide 19 - Special case: $p = 1$	21
8.10.	Slide 22 - Almost same thing	22
8.11.	Slide 24 - Maximize Covariance	22
8.12.	Slide 26 - PLS Regression	22
9.	Discriminant Analysis and Classification	23
9.1.	Slide 3 - Motivation	23
9.2.	Slide 4 - Two-class problem	23
9.3.	Slide 6 - Classification Table	23
9.4.	Slide 7 - F-Score of Binary Classification	23
9.5.	Slide 8 - Cost of Misclassification	24
9.6.	Slide 9 - Minimizing ECM	24
9.7.	Slide 10 - An Example Using ECM	24
9.8.	Slide 13 - Special Case III: Highest Posterior Probability	24
9.9.	Slide 14 - Naive Bayes: Gaussian Discriminant Analysis	24
9.10.	Slide 15 - Gaussian Discriminant Analysis: Classification	25
9.11.	Slide 16 - Gaussian Discriminant Analysis: Decision Boundary	25
9.12.	Slide 17 - Fishers Linear Discriminant Analysis	25
9.13.	Slide 18 - MANOVA-Like Idea	25
9.14.	Slide 19 - Within Versus Between Variation	25
9.15.	Slide 20 - Fishers LDA	26
9.16.	Slide 22 - Case I: $\Sigma_1 = \Sigma_2 = \Sigma$	26
9.17.	Slide 24 - Connection to Fishers LDA	26
9.18.	Slide 26 - Case II: $\Sigma_1 \neq \Sigma_2$	26
9.19.	Slide 27 - Logistic Model For Two Populations	26
9.20.	Slide 28 - Maximum Likelihood Estimator	26
9.21.	Slide 29 - Penalized Logistic Regression	26
9.22.	Slide 30 - Limitation	26
9.23.	Slide 31 - Euclidean Inner Product	27
9.24.	Slide 32 - New Features	27
9.25.	Slide 33 - Kernel Function and Kernel Matrix	27
9.26.	Slide 38 - Motivation: Margin	27

9.27.	Slide 40 - Hinge Loss: Brief Intro	27
9.28.	Slide 41 - Hinge Loss Vs Log-Likelihood Loss	27
9.29.	Slide 45 - Best Partition	27
9.30.	Slide 47 - A Tree Versus A Forest	27
9.31.	Slide 56 - Linear Discriminants	27
10.	Clustering	28
10.1.	Slide 3 - Applications of Multivariate Analysis: Clustering	28
10.2.	Slide 5 - Distance for Items/Observations	28
10.3.	Slide 6 - Clustering Algorithms	28
10.4.	Slide 7 - Linkage Methods	28
10.5.	Slide 9 - An example	28
10.6.	Slide 16 - Dendrogram	28
10.7.	Slide 18 - Wards Hierarchical Clustering Method	28
10.8.	Slide 20 - Pros and Cons	28
10.9.	Slide 27 - CH Index and Gap Statistic	28
10.10.	Slide 29 - Parametric Approach: Mixing Distribution	28
10.11.	Slide 30 - EM Algorithm	28
10.12.	Slide 31 - E step	29
11.	Old Exam 2022-01-05	30
12.	Old exam - training exam	33
13.	Old exam 2021-01-11	34
14.	Tips/rumors for the exam	36

## 1. INTRODUCTION

Analysis dealing with simultaneous measurements on many variables.

We may want to do some statistical analysis on not only salary, but factor in things such as gender, whether or not one has been to uni etc.

One should always strive to use as much information as possible, you want to remove any chance to miss a pattern.

In general, if you arrive to a conclusion, think of why/what caused this and factor everything in your data and analysis.

### 1.1. MANOVA.

MANOVA is a method to measure if a data-set shares a similar mean. For example, with different flower types we may want to check if "does sweden has a similar income as norwegian citizens", comparing the sample from sweden to norwegian. We will get different numbers but that is something that we take into analysis.

### 1.2. Regressionanalysis.

Allows us to predict a variable  $y$  from an observation  $x$ .  $x = \text{bmi}$ , while  $y$  is your blood pressure.

## 2. SAMPLE &amp; RANDOM MATRICES

## 2.1. Slide 3 - Expectation.

For a discrete random variable we use summation, for a continuous random variable we use integrals. What do we use for vectors/matrices?

⇒ We perform the operations elementwise in the matrix. Take  $\mathbb{E}(X_{ij})$

## 2.2. Slide 4 - Covariance Matrix.

Recall

$$\text{Cov}(X, Y) = \mathbb{E}(X - \mathbb{E}(X))(Y - \mathbb{E}(Y)) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) \quad (1)$$

for scalars.

What about  $\text{Cov} \left( \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}, \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \right)$ ?

We can pick any pair  $(X_i, Y_j)$  and compute  $\text{Cov}(X_i, Y_j)$  leading to the same as (1) but with  $X_i, Y_j$  instead.

In the case  $\begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}, \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$ , we get a  $3 \times 2$  matrix where the  $i, j$ th elements corresponds to  $\text{Cov}(X_i, Y_j)$ .

Think of it like

$$XY^T = \begin{pmatrix} X_1Y_1 & X_1Y_2 \\ X_2Y_1 & X_2Y_2 \\ X_3Y_1 & X_3Y_2 \end{pmatrix} \quad (2)$$

Now look at  $\mathbb{E}(XY^T)$ , same as (2) but  $\mathbb{E}(X_iY_j)$ .

Then we can easily see that  $\text{Cov}(X, Y) = \mathbb{E}(XY^T) - \mu_X \mu_Y^T$

*What if  $X$  is continuous and  $Y$  discrete?*

*What if  $Y = X$ ?*

$$\text{Cov}(X_i, X_i) = \mathbb{E}(X_i^2) - (\mathbb{E}(X_i))^2 = \text{Var}(X_i)$$

## 2.3. Slide 5 - Covariance Matrix.

Since in the scalar case  $\text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i)$ , then  $\text{Cov}(X, Y) = \sum = \text{symmetric \& positive definite}$ .

**Definition/Sats 2.1: Positive & Semi-definite**

Definite matrix  $A$ :

$$A > 0 \Leftrightarrow x^T A x > 0$$

Semi-definite matrix  $A$ :

$$A \geq 0 \Leftrightarrow x^T A x \geq 0$$

## 2.4. Slide 6 - Linear Combination.

You can view the vector  $c$  as regression values for example

## 2.5. Slide 7 - Linear Combination.

**Example:**

$$\text{Var}(X_1 + 2X_2 + 4X_3) \sim \text{Var} \left( \begin{pmatrix} 1 & 2 & 4 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \right)$$

A tip for remembering where to put  $c^T$ , think of it like matching dimensions of left hand side and right hand side.

We only want to compute expectation for the random stuff, so we can chuck coefficients and constants out.

## 2.6. Slide 9 - Independence.

For simplicity, we define independence in the continuous case as  $f(X, Y) = f(X)f(Y)$  and in the discrete case as  $P(X, Y) = P(X)P(Y)$

**Anmärkning:** Jist because  $\text{Cov}(X, Y) = 0$  does not imply independence. Take the unit circle and the contour as pairs over  $(X, Y)$ . It is clear that  $(X, Y)$  are dependant but their covariance is 0 since for every point on the circle you can reflect the  $X, Y$  and therefore, by  $\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$ , you would be adding a bunch of 0. Same goes for any function that can be reflected.

## 2.7. Slide 10 - Random Sample.

**Example** (Scalar case):

Let  $\mathbf{x} \sim x_1 x_2 x_3 \dots$  be a random sample from  $N(\mu, \sigma^2)$

We look at what it means for scalar random variables to be independent:

$$\begin{aligned} F(X, Y) &= F(X)F(Y) \\ f(x, y) &= f(x)f(y) \\ p(x, y) &= p(x)p(y) \end{aligned}$$

The same principle goes for random vectors, eg:

$$X_{n \times p} = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix}$$

Think of each row as a sample from a different place  $\Rightarrow$  independence in row  $\Rightarrow$  random sample.

**Non-example:** Looking at the pulse of 1 person is not an independent response since it is only about 1 person. Even if you sampled a bunch of values from the same person into a matrix, that would still be a non-independent sample since we only sample from 1 person.

**Non-example:** Let us assume there is a competition between Uppsala and Lund in Multivariate Analysis. Everyone in the class at Uppsala has had the same teacher, so the values collected from that class are not independent.

## 2.8. Slide 12 - Some Notes on Sample Covariance Matrix.

Unbiased becomes biased during non-linear & non-affine transformations.

Even for large  $n$ , sometimes you cannot ignore the difference between  $S_n$  and  $S$  (eg. determining exact distributions)

## 2.9. Slide 17 - Sample Covariance Matrix.

$$\begin{bmatrix} \mathbf{x}'_1 - \bar{\mathbf{x}}' \\ \mathbf{x}'_2 - \bar{\mathbf{x}}' \\ \vdots \\ \mathbf{x}'_n - \bar{\mathbf{x}}' \end{bmatrix} = \underbrace{\mathbf{X}}_{n \times p} - \underbrace{\mathbf{1}}_{n \times 1} \underbrace{\bar{\mathbf{x}}'}_{1 \times p}$$

$$X - \frac{1}{n} \mathbf{1} \mathbf{1}^T X = (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) X$$

So for  $(X - \frac{1}{n} \mathbf{1} \mathbf{1}^T X)^T (X - \frac{1}{n} \mathbf{1} \mathbf{1}^T X)$ :

$$X^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T)^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) X = X^T \left[ I - \frac{1}{n} \mathbf{1} \mathbf{1}^T - \frac{1}{n} \mathbf{1} \mathbf{1}^T - \frac{1}{n} \mathbf{1} \mathbf{1}^T + \frac{1}{n} \mathbf{1} \underbrace{\mathbf{1}^T \mathbf{1}}_{=n} \mathbf{1}^T \right] X$$

$$X^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T - \frac{1}{n} \mathbf{1} \mathbf{1}^T + \frac{1}{n} \mathbf{1} \mathbf{1}^T) X = X^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) X$$

$$X^T X - X^T \mathbf{1} \mathbf{1}^T X \Rightarrow S_n = \frac{1}{n} \underbrace{X^T X}_{\text{Data matrix}} - (\frac{1}{n} X^T \mathbf{1}) (\frac{1}{n} \mathbf{1}^T X)$$

$$\text{Cov}(X) = \mathbb{E}(X X^T) - \mathbb{E}(X) \mathbb{E}(X)^T$$

### Anmärkning:

$\mathbf{1}$  is an  $n \times 1$  vector of ones.



### 3. MULTIVARIATE NORMAL DISTRIBUTION

#### 3.1. Slide 4-5 - From Univariate to Multivariate Normal.

Recall that in the univariate case we had:

$$(x - \mu) \frac{1}{\sigma^2}$$

In the multivariate case, we swap  $x$  and  $\mu$  for vectors instead.

Since variance matrix is expressed by  $(x - \mu)^T \Sigma^{-1} (x - \mu)$ , instead of  $\sigma^2$  we have have

$$\frac{1}{\sigma\sqrt{2\pi}} \sim \rightarrow \frac{1}{(2\pi)^{p/2} \sqrt{\det(\Sigma)}}$$

#### Anmärkning:

Covariance matrix must be positive definite! Not semi.

There is no requirement for slide 4 with  $\Sigma$

The  $(2\pi)^{p/2}$  comes from multiplying  $z_1 z_2 \cdots z_p$   $p$ -times.

#### 3.2. Slide 6 - Special Case: Bivariate Normal.

#### Anmärkning:

$\rho$  denotes the correlation coefficient

$\sigma_{11}$  &  $\sigma_{22}$  correspond to our variance

$\sigma_{12}$  &  $\sigma_{21}$  correspond to our covariance

$$\text{Corr}(x_1, x_2) = \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}}$$

#### 3.3. Slide 7 - Contour of Bivariate Normal Density.

We change the correlation to see what happens.

#### 3.4. Slide 8 - Linear Combinations.

For the univariate case, we had that if we scaled  $X \sim N(\mu, \sigma^2)$  with an affine transformation, we got  $aX + b \sim N(a\mu, a^2\sigma^2)$ .

One thing that is good to keep in the back of the head is that the linear combination/affine transformation of normally distributed random variables will remain normal.

Let us look at what happens when we look at the multivariate case:

$$\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \quad Y_1 \sim N \quad Y_2 \sim N$$

$$\Rightarrow \begin{pmatrix} x_1 + x_2 \\ x_1 - x_2 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix} \sim N_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, A \Sigma A^T \right)$$

From result 4.2, we can get the result of multi-linear combinations

#### 3.5. Slide 10 - Normal and Chi-Square.

If  $X$  has a linear combination will it still be  $p$ -degrees of freedom? Answer is surprisingly yes!

$$\Sigma^{-1} = \Sigma^{-1/2} \Sigma^{-1/2} \quad X \sim N_p(\mu, \Sigma)$$

$$\Rightarrow Z = \Sigma^{-1/2}(x - \mu) = \underbrace{\Sigma^{-1/2}x}_A \underbrace{- \Sigma^{-1/2}\mu}_d \sim N_p(0, \Sigma^{-1/2} \Sigma \Sigma^{-1/2})$$

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = Z^t Z = \sum_{j=1}^p Z_j^2$$

### 3.6. Slide 11 - Subset of Variables.

Using result 4.4, we can choose subsets however we want, it will stay normal.

### 3.7. Slide 12 - Example: Subset of Variables.

From the slide we have the following:

Suppose that:

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \sim N_2 \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{bmatrix} \right)$$

Find the distribution of  $\begin{bmatrix} X_1 \\ X_3 \end{bmatrix}$  as well as the distribution of

$$\begin{bmatrix} X_1 & X_3 \end{bmatrix} \begin{bmatrix} \sigma_{11} & \sigma_{13} \\ \sigma_{31} & \sigma_{33} \end{bmatrix}^{-1} \begin{bmatrix} X_1 \\ X_3 \end{bmatrix}$$

In the first one, what we really essentially are looking for is the following:

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \sim N_3 \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{bmatrix} \right)$$

If we want  $\begin{bmatrix} X_1 \\ X_3 \end{bmatrix}$ , then:

$$\begin{bmatrix} X_1 \\ X_3 \end{bmatrix} \sim N_2 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{11} & \sigma_{13} \\ \sigma_{31} & \sigma_{33} \end{bmatrix} \right)$$

So:

$$\begin{bmatrix} X_1 & X_3 \end{bmatrix} \begin{bmatrix} \sigma_{11} & \sigma_{13} \\ \sigma_{31} & \sigma_{33} \end{bmatrix}^{-1} \begin{bmatrix} X_1 \\ X_3 \end{bmatrix} \sim \chi^2_2$$

It is really important to remember that linear combinations of normal variables, are still normal variables. Since linear combinations can be regarded as linear/affine transformations, the "crossing out the  $X_2$ " part of the computation is really just matrix-multiplication, since:

$$\begin{bmatrix} X_1 \\ X_3 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_A \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$$

In order to find the distribution of the second question, we see that it is really just the multivariate  $\chi^2$ :

$$\begin{bmatrix} X_1 & X_3 \end{bmatrix} \Sigma^{-1} \begin{bmatrix} X_1 \\ X_3 \end{bmatrix} \Rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu) \text{ where } \mu = 0 \\ \Rightarrow \chi^2_2$$

### 3.8. Slide 13 - Subset of Variables.

#### Anmärkning:

Since what we really care about is what happens during the transpose, sometimes we write  $\Sigma_{12}$  for  $\Sigma_{12} = \Sigma_{21} = 0$

### 3.9. Slide 15 - Marginal Normal and Joint Distribution.

Usually, if they are independent, they are normal.

### 3.10. Slide 23 - Likelihood of Normal Random Sample.

$$a^T B a = \text{tr}(a^T B a) = \text{tr}(B a a^T)$$

Of course, in order to maximize the likelihood we sometimes need to find the derivative of the matrix/vector.

**Example:**

$$\underbrace{\begin{bmatrix} x_1 & x_2 \end{bmatrix}}_{x^T} \underbrace{\begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix}}_B \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_x \Rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \begin{bmatrix} b_{11}x_1 + b_{12}x_2 \\ b_{21}x_1 + b_{22}x_2 \end{bmatrix}$$

$$\Rightarrow b_{11}x_1^2 + b_{12}x_1x_2 + b_{21}x_1x_2 + b_{22}x_2^2 = f(x_1, x_2)$$

Now we can just collect the partials in a vector (or a matrix if we end up with a matrix):

$$\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2b_{11}x_1 + b_{12}x_2 + b_{21}x_2 \\ 2b_{22}x_2 + b_{12}x_1 + b_{21}x_1 \end{bmatrix} = \begin{bmatrix} 2b_{11} & b_{12} + b_{21} \\ b_{12} + b_{21} & 2b_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

### 3.11. Slide 32 - Limit of MLE.

$$\underbrace{\frac{n}{n-1}}_{\substack{n \rightarrow \infty \\ \rightarrow 1}} \underbrace{(\mu_1 - \hat{X}_i)}_{\rightarrow 0} \underbrace{(\hat{X}_k - \mu_k)}_{\rightarrow 0} \rightarrow \frac{1}{n-1} \sum \approx \sigma_{ik}$$

#### 4. INFERENCE FOR SEVERAL SAMPLE

##### 4.1. Slide 3 - Paired Data.

Here, *paired* means 2 tests/observations from the **same** subject  $x_{j_1}$  and  $x_{j_2}$  are always correlated since they are about the same person.

##### 4.2. Slide 9 - Two Populations.

Different people, but 2 populations (different countries, people, etc).

$X_{ij}$ , where  $j$  could be the  $j$ :th person in the  $i$ :th "country"/group

But different countries may have different amounts in population, what happens to  $D_i$ ? Well, we will allow  $t$  and define our own  $\mathbb{E}$  and  $\Sigma$

##### 4.3. Slide 10 - Pooled Sample Covariance Matrix.

$$X_{11}, \dots, X_{1n} \sim N(\mu_1, \Sigma) \quad \text{Estimate of } \Sigma: \hat{\Sigma}_1 = \frac{1}{n_1 - 1} \sum_{j=1}^{n_1} (X_{1j} - \bar{X}_1)^2$$

$$X_{21}, \dots, X_{2n} \sim N(\mu_2, \Sigma) \quad \text{Estimate of } \Sigma: \hat{\Sigma}_2 = \frac{1}{n_2 - 1} \sum_{j=1}^{n_2} (X_{2j} - \bar{X}_2)^2$$

Here we are not using all our possible data to get a good approximation/estimate. Sure,  $\hat{\Sigma}_1$  may be unbiased, but it can be better:

$$\left. \begin{aligned} \mathbb{E}(\hat{\Sigma}_1) &= (n_1 - 1)\Sigma \\ \mathbb{E}(\hat{\Sigma}_2) &= (n_2 - 1)\Sigma \end{aligned} \right\} \Rightarrow (n_1 + n_2 - 2)\Sigma = S_{\text{pooled}}$$

If  $\mu_1 = \mu_2$ , then we can estimate  $\Sigma$  using:

$$\left. \frac{\sum_{j=1}^{n_1} (X_{1j} - \bar{Z})^2 + \sum_{j=1}^{n_2} (X_{2j} - \bar{Z})^2}{n_1 + n_2} \right\} \quad \text{Where } \bar{Z} = \frac{X_1 + \dots + X_{1n} + X_{21} + \dots + X_{2n}}{n_1 + n_2}$$

##### 4.4. Slide 20 - MANOVA Model.

$\tau_i$  denotes population  $i$  where  $\tau_i$  is how much that population deviates from the mean. This can be useful, since we can look at some statistic over nordic countries and let  $\mu$  be the mean over all nordic countries (by adding all statistics from every country and dividing by the nordic population, not by taking the mean of the mean in every country)

A one way MANOVA indicates that we are looking at one category of population, ie nationality. We can of course include things like nationality, race, gender, etc. but then it will be two-way/more MANOVA.

Since we have variations either above or below the average (per definition of the average), some  $n_i \tau_i$  will be negative while others might be positive. That is why we set  $\sum n_i \tau_i = 0$ .

If we do not do this, we might as well write:

$$\begin{aligned} \mu + \tau_{\text{SWE}} &= \mu + c + \tau_{\text{SWE}} - c \\ \mu + \tau_{\text{NOR}} &= \mu + c + \tau_{\text{NOR}} - c \end{aligned}$$

#### 4.5. Slide 28 - Multivariate Two-Way Fixed Effects Model with Interaction.

$$\mu + \underbrace{\tau_l}_{\text{property 1 in nordic}} + \underbrace{\beta_k}_{\text{property 2 in nordic}} + \underbrace{\gamma_{lk}}_{\text{property 1} \wedge \text{property 2 in nordic}} + e_{lkr}$$

A *marginalising parameter* is setting it as a summation index, eg:  $\sum_j \gamma_{jk} \rightarrow j$  is marginalised

$\tau, \beta = \text{main effect}$ , while  $\gamma$  is called the *interaction term*

#### 4.6. Slide 33 - Test of No Interaction.

It makes no sense (often) to test  $\tau_i$  since even if  $\sum \tau_i = 0$ , it may/will have effect on  $\gamma_{lk}$ . This is called *principal of marginality*.

## 5. REGRESSION

## 5.1. Slide 6 - Classic Linear Regression.

$$Y = Z^T \beta + e \rightarrow \mathbb{E}(e|Z) = 0$$

$$\mathbb{E}(Y|Z) = \mathbb{E}(Z^T \beta + e|Z) = \underbrace{\mathbb{E}(Z^T \beta|Z)}_{\mathbb{E}(Z^T \beta)} + \underbrace{\mathbb{E}(e|Z)}_{=0}$$

Why is it then called linear when we do not always approximate using linear functions but curves? Well,  $Y = Z^T \beta + e = \beta_1 z_1 + \dots + \beta_r z_r$ , this is just a *linear* combination of our regression-coefficients.

An example,  $Y = \beta_1 z_1 + \beta_2 z_2^2$  is still linear regression, since it is linear in  $\beta$ , what happens with  $Z$  is not what we care about.

However,  $Y = e^{\beta_1 z_1} / \sin(\beta_2 z_2)$  is not a linear regression.

## 5.2. Slide 7 - Matrix Notation.

*Heteroscedasticity* = every observation variance depends on observation. Can also be dependant on  $Z$ , so  $\sigma_i^2$

Estimation methods still valid for heteroscedastic variances, although maybe not optimal.

5.3. Slide 9 - ANOVA With  $g = 2$ .

Note that we only need 2 columns to find the last rank  $\rightarrow 1$  restriction:

$$\sum n_l \tau_l = 0 \Rightarrow \tau_l = 0$$

5.4. Slide 10 - Anova With  $g = 2$  and  $b = 2$ .

Instead of restriction, construct a submatrix with the bad (linearly dependant) columns deleted. Estimation depends on rank.

## 5.5. Slide 11 - Ordinary Least Squares.

$$-2Z^T(y - Z\beta) = 0 \Leftrightarrow Z^T y = Z^T Z \beta = \hat{\beta}_{\text{OLS}} = (Z^T Z)^{-1} Z^T y$$

## 5.6. Slide 12 - OLS Estimator.

$$Y = Z\beta + e \quad \mathbb{E}(Y) = Z\beta \quad \hat{Y} = Z\hat{\beta}$$

*Residual* is given by  $\hat{e} = y - Z\hat{\beta} = y - Hy = (I - H)y$

Interesting things:

$$Z^T \hat{e} = Z^T (I - H)y = (Z^T - \underbrace{Z^T Z (Z^T Z)^{-1} Z^T}_H)y = 0$$

We note that the residual is perpendicular to observed values! This makes sense.

$$\hat{y}^T \hat{e} = y^T H (I - H)y = y^T (H - H^2)y = 0$$

$$H^2 = ZZ^T (Z^T Z)^{-1} Z^T Z (Z^T Z)^{-1} Z^T = H \quad (\text{idempotent})$$

Predicted value is perpendicular to  $\hat{e}$

### 5.7. Slide 15 - Sampling Properties of OLS Estimators.

$\mathbb{E}(\hat{e}^T \hat{e}) = (n - r)\sigma^2$  if  $e$  has some distribution of  $\mu = 0$  and  $\Sigma = \sigma^2 I$

$\frac{1}{n-1}$  comes from  $\frac{\hat{e}^T \hat{e}}{n-r}$ , since  $\underbrace{Z}_{n \times r} \underbrace{\beta}_{r \times 1}$ , but for constants/1D we have  $r = 1$

$$\begin{aligned} \text{Cov}(\hat{\beta}, \hat{e}) &= \text{Cov}\left(\underbrace{(Z^T Z)^{-1} Z^T}_A y, \underbrace{(I - H)}_B y\right) = (Z^T Z)^{-1} Z^T \sigma^2 I (I - H)^T \\ &\Rightarrow \sigma^2 (Z^T Z)^{-1} Z [I - Z(Z^T Z)^{-1} Z] = \sigma^2 ((Z^T Z)^{-1} Z^T - (Z^T Z)^{-1} Z (Z^T Z)^{-1} Z^T) = 0 \\ &\Rightarrow \text{unbiased} \end{aligned}$$

### 5.8. Slide 17 - Distribution of Regression Coefficients.

By assuming  $e \sim N$  distributed, we could do inference on  $\beta$

#### Anmärkning:

- Normal distribution  $\Rightarrow$  every marginal distribution is normal
- Sum of squares of normal random variables  $\sim \chi^2$
- Standard normal ( $N(0, 1)$ ) divided by  $\chi^2$  divided by degrees of freedom  $\sim t_{n-r} \rightarrow$  degrees of freedom
- Joint statistics = how many things you chuck in the conf. intern.

### 5.9. Slide 18 - Confidence Region.

If  $\hat{\beta} - \beta \ll 1$ , then we are close. We capture this in our test.

### 5.10. Slide 19 - Confidence interval.

You will get some  $F$  distribution (**CHECK**)

#### Anmärkning:

Some nomenclature:

- *Multiple regression*  $r \geq 2, 3, \dots$
- *Multivariate regression*  $Y$  is a matrix

### 5.11. Slide 20 - More Than One Responses.

$\underbrace{Y}_{m \times 1}$  here is for one subject, where  $m$  is the amount of responses. If we have  $n$  subjects, we get what is on slide 21.

### 5.12. Slide 22 - Assumptions.

In the second point  $e_{(i)} = i$ th thing to compare, eg price/time and not subject such as apartment.  $\text{Cov}(e_{(i)}, e_{(k)})$  compares price and time simultaneously.

5.13. **Slide 23 - Least Squares.**

$$\underbrace{(Y_Z \beta)^T}_{m \times n} \underbrace{(YZ \beta)}_{n \times m} \sim m \times m$$

**Anmärkning:**

- $Y_{(i)}$  =  $i$ th column  $Y_i$  =  $i$ th row
- Wishart = Generalisation of  $\chi^2$  in multivariate case

5.14. **Slide 27 - Regression Coefficients With Zero Constraints.**

When we reduce to  $Y_{n \times m} = Z_1 \beta_1 + E$ , we can go back to multiple regression by letting  $Z = Z_1$ ,  $\beta = \beta_1$

What happens to  $E$ ? It never changes, we just use the one that corresponds with the column we test in the multiple regression model.

5.15. **Slide 31 - LRT when  $m = 1$ .**

Let  $w$  = numerator =  $(Y - Z\hat{\beta})^T(Y - Z\hat{\beta})$

Let  $w_1$  = denominator =  $(Y - Z_1\hat{\beta}_1)^T(Y - Z_1\hat{\beta}_1)$

Result 7.6 says  $\frac{w_1 - w}{w}$  but if  $\frac{w}{w_1}$  small, then  $\frac{w_1}{w}$  must be big  $\Rightarrow \frac{w_1}{w} - 1$  is still big.

$F$  test tests if every  $\beta_i$  is 0 except the intercept.

5.16. **Slide 32 - Prediction of Regression Function.**

$\beta_{(i)}$  has dimension  $1 \times r$   $\hat{\beta}_{(i)}^T$  has dimension  $r \times 1$ , dimension of  $z_0$   $1 \times r$ . This gives us that  $\hat{\beta}_{(i)} z_0$  is a scalar  $\Rightarrow N_1(\mu, \sigma^2)$

5.17. **Slide 34 - Forecast New Response.**

$\mathbb{E}(Y)$  is predicting mean, but we want to find/predict  $Y$



## 6. PRINCIPAL COMPONENT ANALYSIS

## 6.1. Slide 3 - Motivation.

PCA is mostly used in pre-processing these days, instead of being the actual analysis.

## 6.2. Slide 4 - Task of Principal Component Analysis (PCA).

$\mathbf{a}_3$  maximizes  $\text{Var}(\mathbf{a}_3^T \mathbf{X})$  and  $\text{Cov}(\mathbf{a}_3^T \mathbf{X}, \mathbf{a}_j^T \mathbf{X}) = 0$ . In the covariance term, we look at all  $j < 3$  and not just  $j = 1$ . That is, our requirement is that  $\text{Cov}(\mathbf{a}_3^T \mathbf{X}, \mathbf{a}_1^T \mathbf{X}) = 0 \wedge \text{Cov}(\mathbf{a}_1^T \mathbf{X}, \mathbf{a}_2^T \mathbf{X}) = 0$

Big variation is good since it covers more cases. Think of it like salary analysis, with low variance you may only have asked the CEO/higher ups and you will not get as great of a picture as if you used the whole wide company.

## 6.3. Slide 5 - Restriction.

$$\text{Cov}(\mathbf{a}_i^T \mathbf{X}, \mathbf{a}_k^T \mathbf{X}) = \mathbf{a}_i^T \underbrace{\text{Cov}(\mathbf{X}, \mathbf{X})}_{=\Sigma} \mathbf{a}_k \Rightarrow \mathbf{a}_i^T \Sigma \mathbf{a}_k$$

## 6.4. Slide 6/7 - Principal Components and Two useful Lemmas.

Maximize  $\text{Var}(\mathbf{a}_1^T \mathbf{X})$  such that  $\mathbf{a}_1^T \mathbf{a}_1 = 1 \Leftrightarrow \text{maximize } f(\mathbf{a}_1) = \text{Var}(\mathbf{a}_1^T \mathbf{X}) - \underbrace{\lambda}_{\text{Lagrange multiplier}} (\mathbf{a}_1^T \mathbf{a}_1 - 1)$

This uses the Lagrange multiplier method.

Adding more constraints, you add more Lagrange multipliers (*KKT condition*)

In order to maximise, we want  $\frac{df}{da_1} = 0 \wedge \frac{df}{d\lambda} = 0$

Note that:

$$\begin{aligned} \frac{df}{d\lambda} &= -(a_1^T a_1 - 1) = 0 \wedge \frac{df}{da_1} = 1 \\ &\Rightarrow 2\Sigma a_1 - 2\lambda a_1 = 0 \end{aligned}$$

Zero only when  $\Sigma a_1 = \lambda a_1$

## 6.5. Slide 7 - Two useful Lemmas.

Reason we want to use the largest eigenvalue is because we want to maximise variance:

$$Y_1 = a_1^T X \quad (\Sigma a_1 = \lambda a_1) \rightarrow \text{Var}((Y_1)) = a_1^T \Sigma a_1 = \lambda \underbrace{a_1^T a_1}_{=1} = \lambda$$

First thing (maximise variance) is done, second step:

$$\begin{aligned} \max(\text{Var}(a_2^T X)) &= a_2^T \Sigma a_2 \text{ s.t. } a_2^T a_1 = 1 \quad \underbrace{\text{Cov}(a_2^T X, a_1^T X)}_{\substack{= a_2^T \Sigma a_1 = a_2^T \lambda a_1 = \lambda a_2^T a_1 \\ a_2 \notin \text{span}\{a_1\} \Leftarrow a_2^T a_1 = 0}} = 0 \\ &\Rightarrow \max(f(a_2)) = a_2^T \Sigma a_2 - \lambda(a_2^T a_2 - 1) \end{aligned}$$

The whole text under the covariance can be boiled down to implying that  $a_2$  has to be a span of other eigenvectors. Then they will be orthogonal to each other.

For the last row, to  $f(a_2)$ , we use the second largest eigenvector.

## 6.6. Slide 9 - Principal Components.

Even if we have eigenvalues with duplicate values this holds.

### 6.7. Slide 10 - Total Variation Explained by Principal Components.

By having orthogonal  $Y_i$ :s (due to eigenvectors), we have reduced dependency from all  $Y_i$ :s. Any non-orthogonality yields some correlation between some  $Y_i$  and  $Y_k$ , and we have now removed that.

Using  $\frac{\lambda_k}{\sum \lambda_i}$  gives us the contribution from  $\lambda_k$ , but we can look for say  $\frac{\sum_k^j \lambda_k}{\sum \lambda_i}$  until we get a % we are satisfied with.

### 6.8. Slide 12 - Principal Components From Correlation Matrix.

Reason we standardized is to be able to compare with other data of different scale

$$V = \begin{bmatrix} \sigma_{11} & & \\ & \sigma_{22} & \\ & & \ddots \end{bmatrix}$$

### 6.9. Slide 14 - Sample Principal Components.

Let  $\Sigma$  be our sample covariance matrix instead. Then we carry out as usual.

#### **Anmärkning:**

Centered = mean is 0. Taking away some of the data would yield an almost 0 mean (numerically 0)

## 7. FACTOR ANALYSIS

## 7.1. Slide 3 - Latent Variable Modelling.

LVM = factor analysis. Find values such as personality using a proxy. Personality is the factor/latent variable.

In PCA, we had a bunch of values and we wanted to simplify them and keep them as concise as possible while still retaining as much of the information as possible. In factor analysis, we go the other direction, we have some "simplified" data set and we want to draw more conclusions from this.

## 7.2. Slide 4 - The Model.

$i$  =  $i$ th question in questionnaire

$$\underbrace{X_i}_{\text{Math-score/what you know}} = \mu_i + \underbrace{\ell_{i1}}_{\text{Algebra ability/what you want to predict}} + \cdots + \varepsilon \rightarrow \text{math ability still there, but } \varepsilon \text{ may be latent}$$

Depending on application, sometimes we need to find  $\ell$ .

Tasks is usually lower than factors.

In multiple regression we had  $X = X\beta + \varepsilon$ , but this is for all people/subjects, while  $X = \mu + \ell F + \varepsilon$  is for 1 person/subject, like in multivariate multiple regression model, difference is  $Y_i = \beta^T \underbrace{X_i}_{\text{observed}} + \varepsilon$

A regressor is values you do not really observe such as IQ, but you still want to build a model using observations.

## 7.3. Slide 5 - Scale Indeterminacy.

We can redefine scale, continuing with the IQ example, there really is nothing stopping us from saying that the IQ scale should lay inbetween  $[-1,1]$  by just compressing the Gaussian.

## 7.4. Slide 7 - Model Implied Covariance Matrix.

$$\text{Cov}(LF + e) = \text{Cov}(LF, LF) + \underbrace{\text{Cov}(LF, e)}_{L\text{Var}(F)L^T + \text{Var}(e) = LL^T + \psi = \text{Cov}(X)} + \underbrace{\text{Cov}(e, LF)}_{=0} + \text{Cov}(e, e)$$

The reason for the name model implied covariance matrix, is because it is implied by the setup of the model.

**Example:**

$$\underbrace{\begin{bmatrix} l_{11} & l_{12} \\ l_{21} & l_{22} \\ l_{31} & l_{32} \end{bmatrix}}_L \underbrace{\begin{bmatrix} l_{11} & l_{21} & l_{31} \\ l_{12} & l_{22} & l_{32} \end{bmatrix}}_{L^T} + \begin{bmatrix} \psi_1 & & \\ & \psi_2 & \\ & & \psi_3 \end{bmatrix}$$

$$= \begin{bmatrix} l_{11}^2 + l_{12}^2 + \psi_1 & 0 & 0 \\ 0 & l_{21}^2 + l_{22}^2 + \psi_2 & 0 \\ 0 & 0 & l_{31}^2 + l_{32}^2 + \psi_3 \end{bmatrix}$$

Here  $l_{11}^2 + l_{12}^2 + \psi_1 = \text{Var}(X_1)$ ,  $l_{21}^2 + l_{22}^2 + \psi_2 = \text{Var}(X_2)$ ,  $l_{31}^2 + l_{32}^2 + \psi_3 = \text{Var}(X_3)$  Want communality over uniqueness

## 7.5. Slide 8 - Existence of Decomposition.

Note that the invalid decomposition is invalid due to negative  $\psi$ ,  $\text{Var}(e)$  not as assumption.

**Anmärkning:** PCA always doable, not the same with factor analysis.

### 7.6. Slide 9 - Indeterminacy.

Rotation (multiplication by diagonal matrix  $T$ ) is invariant  $\Rightarrow$  indeterminate

### 7.7. Slide 10 - Scale Invariant.

**Curiosity:** What happens if  $c\mu$  or  $cz$ ?

### 7.8. Slide 11 - Popular Estimation Methods.

**Curiosity:** What if only one person in our sample has  $p + 1 \times 1$ ?

### 7.9. Slide 12 - Spectral Decomposition.

Recurring small eigenvalues yields a possibility to find  $\psi$  such that  $\psi$  is diagonal.

### 7.10. Slide 14 - Determine Number of Factors.

A factor/factors = things we want to test, while  $m$  is the number of things we want to test. Say for example we use the multivariate analysis exam as an example. It might test our knowledge in multivariate analysis ( $m = 1$ ), but might also in addition test our R knowledge ( $m = 2$ )

Kaiser criterion sucks.

### 7.11. Inbetween lectures.

We want the RMSEA index  $< 0.1$ . Or rank of the two sample test  $> 3 = \text{bad?}$

### 7.12. Slide 23 - Versus PCA.

Notice that they are very similar. I can do PCA to solve FA.

### 7.13. Slide 24 - Orthogonal Rotation.

Multiplication will not change anything.

### 7.14. Slide 26 - Oblique Rotation.

All factors cannot always be orthogonal. For example, say we test English knowledge, then spelling will have some correlation to writing  $\rightarrow$  oblique.

If we have the following correlation matrix  $\begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$ , then 0.9 (since it is too close to 1), might be the same so we could reduce 1 factor.

### 7.15. Slide 28 - Bartlett Score.

A factor score is the estimate of  $F$

## 8. CANONICAL CORRELATION ANALYSIS

**Anmärkning:**

Tasks is coumns/abilites, such as listening etc

Demean = remove mean, that is if something hsa value  $\mu + x$ , after demeaning it only has value  $x$

*Scores* is tranformed data, ie.  $x$ -scores =  $a^T x$  (values of linear combination of  $U$ )

The scores (latent variables) are what we want for further analysis

**8.1. Slide 3 - Motivation.**

In PCA, we had a lot of variables and we simplified them to less variables, and in factor analysis we did the opposite and were able to draw conclusions from our data.

In CCA, we have 2 sets of variables, and simplify them to have some lower dimension (think PCA on 2 sets of data).

Note that just like FA was the "opposite" of PCA, there is an "opposite" to CCA (not covered) called structural equation modelling.

The intuition is, say  $X^{(2)}$  has  $q$  entries, and  $X^{(1)}$  has  $p$  entries. That is  $p \cdot q$  scatter plots to read, we want to minimize this number.

In PCA, we maximize  $\text{Var} \left( a^T \begin{bmatrix} X^{(1)} \\ X^{(2)} \end{bmatrix} \right)$

**8.2. Slide 4 - Task.**

When  $\text{Cor}(U, V)$  is maxed, we capture as much information from lower dimensions as we do with higher dimensions (ie, adding more dimensions will not yield more data). This can be seen as regression, want  $U$  to describe  $V$  as good as possible.

**8.3. Slide 5 - Correlation Coefficient.**

Even if we let  $a, b$  blow up, we normalise it in he correlation through the denominator

**8.4. Slide 7 - Canonical Variates.**

$(U_2, V_2)$  will be orthogonal to  $(U_1, V_1)$ , same with  $(U_k, V_k)$ , it will be uncorrelated to all the previous ones.

**8.5. Slide 9 - Find Canonical Variates.**

In PCA we used the Lagrange multipliers in order to maximize, we will do the same but with 2 constraints ( $a^T \Sigma_{11} a = 1$  and  $b^T \Sigma_{22} b = 1$ ):

$$\max f = a^T \Sigma_{12} b - \lambda_1 (a^T \Sigma_{11} a - 1) - \lambda_2 (b^T \Sigma_{22} b - 1)$$

Note that usually it will be  $(a^T \Sigma_{12} b)$ , but we only care about positive values, so we do  $(a^T \Sigma_{12} b)$

Maximizing:

$$\begin{aligned} \frac{\partial f}{\partial a} = 0 &= \Sigma_{12} b - \lambda_1 (\Sigma_{11} + \Sigma_{11}^T) a - 2 \Sigma_{11} & \frac{\partial f}{\partial b} = 0 &= \Sigma_{21} a - \lambda_2 (\Sigma_{22} + \Sigma_{22}^T) b - 2 \Sigma_{22} \\ \frac{\partial x^T \Sigma x}{\partial x} &= (\Sigma + \Sigma^T) x & \frac{\partial b^T x}{\partial x} &= b \\ a^T \Sigma_{12} b &= 2 \lambda_1 \underbrace{a^T \Sigma_{11} a}_{=1} \\ b^T \Sigma_{21} a &= a^T \Sigma_{12} b \Rightarrow b^T \Sigma_{21} a = 2 \lambda_2 \underbrace{b^T \Sigma_{22} b}_{=1} \end{aligned} \left. \vphantom{\begin{aligned} \frac{\partial f}{\partial a} = 0 \\ \frac{\partial f}{\partial b} = 0 \\ \frac{\partial x^T \Sigma x}{\partial x} = (\Sigma + \Sigma^T) x \\ \frac{\partial b^T x}{\partial x} = b \\ a^T \Sigma_{12} b = 2 \lambda_1 \underbrace{a^T \Sigma_{11} a}_{=1} \\ b^T \Sigma_{21} a = 2 \lambda_2 \underbrace{b^T \Sigma_{22} b}_{=1} \end{aligned}} \right\} \Rightarrow \lambda_2 = \lambda_1 = \lambda$$

We will see that for  $\Sigma_{12}b$  and  $\Sigma_{12}a$ :

$$\Sigma_{12}b = 2\lambda\Sigma_{11}a$$

$$\Sigma_{21}a = 2\lambda\Sigma_{22}b \Rightarrow b = \frac{1}{2\lambda}\Sigma_{22}^{-1}\Sigma_{21}a \quad (\text{assuming } \Sigma_{22} \text{ is invertible, which it is since it is the covariance matrix})$$

Plugging this definition of  $b$  in the first equation yields:

$$\Rightarrow \Sigma_{12} \left( \frac{1}{2\lambda}\Sigma_{22}^{-1}\Sigma_{21}a \right) = 2\lambda\Sigma_{22}b$$

$$\Rightarrow \Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}a = 4\lambda^2a \Rightarrow a \quad \text{is an eigenvector to } \Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \quad \text{with eigenvalue } (2\lambda)^2$$

By the useful lemma on slide 8:

$$\begin{aligned} \Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} &= \Sigma_{11}^{-1/2}\Sigma_{22}^{-1}\Sigma_{12}^T\Sigma_{11}^{-1/2} \quad \text{since:} \\ &\Sigma_{11}^{1/2}\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2}\Sigma_{11}^{1/2} \\ &= \Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2}(\Sigma_{11}^{1/2}a) = (2\lambda)^2(\Sigma_{11}^{1/2}a) \\ &\Rightarrow \Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2}e = (\rho^*)^2e \quad \rho^* = 2\lambda \\ &e = \Sigma_{11}^{1/2}a \Rightarrow a = \Sigma_{11}^{-1/2}e \\ &b = \frac{1}{2\lambda}\Sigma_{22}^{-1}\Sigma_{21}a \end{aligned}$$

Using  $a^T\Sigma_{11}a$ :

$$e^T\Sigma_{11}^{-1/2}\Sigma_{11}\Sigma_{11}^{-1/2}e = e^Te = 1$$

Similarly:

$$\begin{aligned} b &= \frac{1}{(2\lambda)^2}a^T\Sigma_{12}\underbrace{\Sigma_{22}^{-1}\Sigma_{22}}_{=1}\Sigma_{22}^{-1}\Sigma_{21}a \\ &\Rightarrow \frac{1}{(2\lambda)^2}a^T\underbrace{\Sigma_{11}\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}}_{\substack{\text{add this} \\ =(2\lambda)^2a}}a \end{aligned}$$

Here, we are using  $a^T\Sigma_{11}a = 1$  (linear constraint is satisfied and variance is 1)

Now:

$$a^T\Sigma_{12}b = \frac{1}{2\lambda}a^T\Sigma_{11}\Sigma_{11}^{-1}\Sigma_{12}^T\Sigma_{12}^{-1}\Sigma_{12}a \Rightarrow 2\lambda\underbrace{a^T\Sigma_{11}a}_{=1} = 2\lambda$$

Correlation is therefore given by the square root of eigenvalues. Once we know  $a$ , we can find  $b$  (that is the meaning behind proportionality).

## 8.6. Slide 12 - Scale Invariant: Coefficient Vector.

Even though it is scale invariant, it is good practice to scale/normalize.

## 8.7. Slide 14 - Proportion of Explained Variance.

$r^2$  = how much of the variance that is explained. Sample covariance matrix *or* sample correlation matrix work in slide 9.

## 8.8. Slide 15 - Sample Canonical Variate Pair.

Sample variance should be one.

Since we use maximum correlation, we can use one set to describe the other

## 8.9. Slide 19 - Special case: $p = 1$ .

$\alpha$  = proportional

## 8.10. Slide 22 - Almost same thing.

- Demean  $\hat{Y} = 0$   $\bar{x}_1 = \bar{x}_2 = \bar{x}_3 = 0$
- $\underbrace{S_{11}^{-1/2} S_{12} S_{22}^{-1} S_{21} S_{11}^{-1/2}}_{\substack{Y \\ \text{mean} = 0}} = \alpha + \sum_i \underbrace{\beta_i x_i}_{\substack{\beta_i x_i \\ \text{mean} = 0}} \Rightarrow$  forces  $\alpha$  to be 0

Multiple  $r$ -squared = "how much variation/variance in  $Y$  is explained by our model"

## 8.11. Slide 24 - Maximize Covariance.

Multicollinearity is a problem with numerical 0:es, such when trying to invert the following matrix:

$$\begin{bmatrix} 1 & 1e16 \\ 1e16 & 1 \end{bmatrix}$$

Very similar to CCA, difference is how we manage restriction:

$$\max \text{Corr}(a^T X, b^T Y) \stackrel{\text{Var}=1}{=} \text{Cov}(a^T X, b^T Y)$$

**Anmärkning:**

Restriction can be  $a^T \sum_{xx} a = 1$  and  $b^T \sum_{yy} b = 1$ . In CCA we downgrade dimension of both  $X, Y$  in PLS we only downgrade  $X$

## 8.12. Slide 26 - PLS Regression.

Since  $t$  is a vector,  $t^T t$  is a scalar, so  $(t^T t)^{-1}$  is just the reciprocal.

## 9. DISCRIMINANT ANALYSIS AND CLASSIFICATION

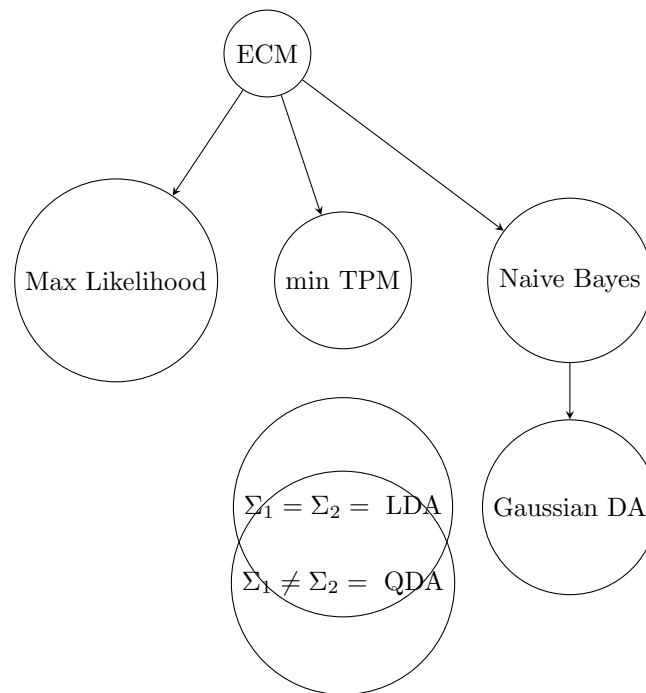


FIGURE 1.

## 9.1. Slide 3 - Motivation.

Discrimination here means to separate subjects.

*Use labeled observations to build a classification rule*, means using previous data.

Classes are well defined, that is they are fixed in some sense. Having two classes such as cat and trying to classify a picture of a horse will only lead to the picture of the horse being classified into either cat or dog. It will not "create a new class" just because the data does not match. It will find which class it matches the best in, and sort into there.

## 9.2. Slide 4 - Two-class problem.

Dividing  $\Omega$  into two disjoint sets  $R_1$  and  $R_2$  is the discrimination step.

*Probabilistic* = vague group belonging

*Deterministic* = you are either in group  $A$  or group  $B$

## 9.3. Slide 6 - Classification Table.

$m_{ij}$  where  $j$  = observed, and  $i$  = predicted

There is another thing we can compute, accuracy =  $\frac{m_{11} + m_{22}}{m_{11} + m_{12} + m_{21} + m_{22}}$

Note that using absolute rates, ie  $\frac{m_{12} + m_{21}}{m_{11} + m_{22}}$  does not take into account the cost of misclassification, it cares more about maximizing  $m_{12} + m_{21}$

## 9.4. Slide 7 - F-Score of Binary Classification.

If F-score  $> 0.5$ , then we are good (DIY rule of thumb)



### 9.5. Slide 8 - Cost of Misclassification.

$$\text{ECM} = \underbrace{c(2|1)}_{\text{cost of misclas.}} \underbrace{P(2|1)}_{\text{prob. of misclas.}} p_1 + \dots$$

Bayesian means "I have some knowledge, I get new data, I update my previous knowledge"

### 9.6. Slide 9 - Minimizing ECM.

Intuitive idea: If it costs a lot to misclass into class 1, then we "want" to shift our mistake-making into the other class (ie shift such that whenever we misclass it is a higher probability that we misclass into the "less costly" class).

### 9.7. Slide 10 - An Example Using ECM.

$$\begin{array}{lll} \pi_1 : X \sim N(0, \Sigma) & p_1 = 0.8 & c(2|1) = 5 \\ \pi_2 : X \sim N(\mu, \Sigma) & p_2 = 0.2 & c(1|2) = 10 \end{array}$$

We now have everything we need to compute:

$$R_1 : \frac{f_1(x)}{f_2(x)} = \frac{\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}x^T \Sigma^{-1}x\right)}{\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)} > \frac{1}{2}$$

The cost is subject to your own judgement, can be biased to you/the one who decides what the cost of misclassification is.

### 9.8. Slide 13 - Special Case III: Highest Posterior Probability.

Assigning class depending on posterior probability.

- Prior:  $P(\pi_i)$
- Posterior:  $P(\pi_i \underbrace{x}_{\text{data}})$

After reading mail (spam not spam mail problem), I may update my knowledge of the probabilities of belonging to each class.

Naive Bayes does this, after updating probabilities it adds into the class with highest probability.

### 9.9. Slide 14 - Naive Bayes: Gaussian Discriminant Analysis.

If we do not know something, either estimate or use your own judgement.

Log-likelihood given by:

$$\begin{aligned} \ell = \sum_j^n Z_j & \left[ \ln(\phi) - \frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma|) - \frac{1}{2} (x_j - \mu_1)^T \Sigma^{-1} (x_j - \mu_1) \right] + \\ & (1 - Z_j) \left[ \ln(1 - \phi) - \frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma|) - \frac{1}{2} (x_j - \mu_2)^T \Sigma^{-1} (x_j - \mu_2) \right] \end{aligned}$$

$$\text{Where } \hat{\phi} = \frac{\text{numbers of 1}}{n} = \hat{p}_1 \Rightarrow \hat{p}_2 = 1 - \phi = \frac{\text{number of 2}}{n}$$

MLE: In order to optimize (find  $\widehat{\mu}_1$ ), we can use the optimization lemma from chapter 4 (slide 24):

$$\left. \begin{aligned} \widehat{\mu}_1 &= \frac{\Sigma}{\Sigma Z_j = \frac{1}{n} \text{ part}} \overbrace{x_j}^{\substack{\text{if in group 2, } Z=0, \text{ so no contrib.} \\ Z_j}} \end{aligned} \right\} \text{average of group 1}$$

$$\widehat{\mu}_2 = \frac{\Sigma(1 - Z_j)x_j}{\Sigma(1 - Z_j)} \left. \right\} \text{average of group 2}$$

Only thing we do not know is  $\Sigma$ , but we use the same lemma from chapter 4.

#### Anmärkning:

Quadratics are the same trace!

$$\begin{aligned} (x_j - \mu_1)^T \Sigma^{-1} (x_j - \mu_1) &= \text{tr}(\Sigma^{-1} (x_j - \mu_1)(x_j - \mu_1)^T) \\ \Rightarrow -\frac{1}{2} [\Sigma Z_j + \Sigma(1 - Z_j)] \ln(|\Sigma|) - \frac{1}{2} \Sigma \text{tr} &\underbrace{(\Sigma^{-1} Z_j (x_j - \mu_1)(x_j - \mu_1)^T + (1 - Z_j)(x_2 - \mu_2)(x_2 - \mu_2)^T)}_{=A} \\ &= -\frac{1}{2} \text{tr}(\Sigma \Sigma^{-1} (Z_j (x_j - \mu_1) \cdots)) \\ q &= \Sigma Z_j + \Sigma(1 - Z_j) \end{aligned}$$

$$\begin{aligned} P(Z = 1) &= \widehat{p}_1 & X|Z = 1 &\sim N(\widehat{\mu}_1, \widehat{\Sigma}) & P(Z = 0) &= \widehat{p}_2 & X|Z = 0 &\sim N(\widehat{\mu}_2, \widehat{\Sigma}) \\ \Rightarrow P(Z|x_0) &= \frac{P(x_0|Z)P(Z)}{\sum_Z P(x_0|Z)P(Z)} \end{aligned}$$

Given you have cancer, you will observe  $X$ . Logistic regression is opposite, given  $X$  determine if the patient has cancer.

#### 9.10. Slide 15 - Gaussian Discriminant Analysis: Classification.

The equivalence is same as slide 16, by taking the logarithm

#### Anmärkning:

Gaussian  $\Rightarrow$  only for normally distributed

#### 9.11. Slide 16 - Gaussian Discriminant Analysis: Decision Boundary.

Changing  $\geq$  to  $=$  yields your *decision boundary*

#### 9.12. Slide 17 - Fishers Linear Discriminant Analysis.

Assuming  $\Sigma_1 = \Sigma_2 = \Sigma$  yields linearity. Gaussian discriminant analysis is a type of ECM. Fisher ECM works through projections instead.

#### 9.13. Slide 18 - MANOVA-Like Idea.

$W$  = within group variation

$B$  = between group variation.

Want things in the "within" set to be as close to each other, ie small  $a^T W a$  and big  $a^T B a$

#### 9.14. Slide 19 - Within Versus Between Variation.

Note that we have not made any assumptions on the distribution.

### 9.15. Slide 20 - Fishers LDA.

Idea is the same as Guassiam, instead of  $\Sigma$ , we use  $W$  and we have "mean - mean" $\Sigma$  (or  $W$ )

### 9.16. Slide 22 - Case I: $\Sigma_1 = \Sigma_2 = \Sigma$ .

The  $-\frac{1}{2}(\mu_1 - \mu_2)^T$  looks similar to  $\frac{1}{2}(\bar{x}_1 - \bar{x}_2)^T$

### 9.17. Slide 24 - Connection to Fishers LDA.

Normality yields robust results, but it can be from other distributions and not necessarily normal.

### 9.18. Slide 26 - Case II: $\Sigma_1 \neq \Sigma_2$ .

QDA allows  $\Sigma$  to be different, contrary to LDA where it has to be the same.

### 9.19. Slide 27 - Logistic Model For Two Populations.

Probabilistic model of  $X$  to model probability

### 9.20. Slide 28 - Maximum Likelihood Estimator.

For Bernoulli! No close method for finding  $\alpha$  and  $\beta$ . We can either guess or use some numerical method to find it/approximate it.

### 9.21. Slide 29 - Penalized Logistic Regression.

Think of  $\alpha$  like the intercept and  $\beta$  like the slope coefficients

Maximizing likelihood function is the same as minimizing the likelihood function after you multiply it by -1.

SPLINE approach (yields continuous and differential function)

Recall noise in data, what can happen to complicated models? Well they would end up modelling the noise in their attempt to fit the curve to the data-points.

Doupple descent phenomenon

#### Example:

$$-\ell(\alpha, \beta) + \lambda[|\beta_1| + |\beta_2|] \Leftrightarrow \max(\ell(\alpha, \beta)) \text{ such that } |\beta_1| + |\beta_2| = t$$

Here  $t$  is some number. Tuning  $t$  tunes  $\lambda$ . The above example is an example of LASSO. Below is an example of Ridge:

$$\beta_1^2 + \beta_2^2$$

One way to tune  $\lambda$  is by cross validation. 3-fold cross validation works in the following way:

- Find a sequence of lambdas beforehand, testing will yield which  $\lambda$  to pick  $(\lambda = 0, \dots, \lambda_{\max})$   
somewhere in here
- Split data into 3 disjoint sets
- Use the first 2 parts to predict 3
- Tune accordingly
- Repeat previous 2 steps  $\binom{3}{2}$  times

#### Anmärkning:

10-fold method is the most popular.

### 9.22. Slide 30 - Limitation.

Kind of like how Newtons method works to find 0:es of a function

### 9.23. Slide 31 - Euclidean Inner Product.

We replace the estimate  $\hat{\beta}$  with  $\Sigma d_j x_j$

### 9.24. Slide 32 - New Features.

Here *features* = covariates.  $\delta(x)$  includes "old" features:

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \xrightarrow{\delta} \delta \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2 x_1 \\ x_2^2 \\ x_1^3 \end{bmatrix}$$

### 9.25. Slide 33 - Kernel Function and Kernel Matrix.

The output of the function  $\kappa$  is a scalar.

$\delta^T(x)\delta(z)$  = Euclidean inner product of  $\delta(x)$  and  $\delta(z)$

### 9.26. Slide 38 - Motivation: Margin.

Idea of a *support vector machine*; if we have multiple ways to pick our decision line, this will give us the best lines.

$Y$  here is some binary set.

### 9.27. Slide 40 - Hinge Loss: Brief Intro.

Term after hinge loss is the Ridge penalty

### 9.28. Slide 41 - Hinge Loss Vs Log-Likelihood Loss.

You always lose some with logistic regression. With Hinge loss we will lose more because of how the curve is.

#### Anmärkning:

In the R code,  $\gamma$  is the inverse of  $\Sigma$

### 9.29. Slide 45 - Best Partition.

$\hat{p}_{mk}$  is the proportion of say  $x_1$  in that region.

*Purity* is about how red is in the black area etc. A very pure area is homogenous and has a small Gini index.

### 9.30. Slide 47 - A Tree Versus A Forest.

$B$  = the number of trees in the forest.

We cannot mimick the population due to biases, but we can see the data as our population instead and pick out a subset of our data and pretend we are picking out a subset of our population

Line 3 is the bootstrapping method.

### 9.31. Slide 56 - Linear Discriminants.

$x_0$  is our new variable  $\Rightarrow$  still linear

## 10. CLUSTERING

Using 1 model to describe your data is not good, clustering and modelling per cluster = much better

**10.1. Slide 3 - Applications of Multivariate Analysis: Clustering.**

We don't know which group it comes from, as opposed to classification.

The first step is to guess how many groups we have to classify our data into.

**10.2. Slide 5 - Distance for Items/Observations.**

Need to define a metric of how "different" something is.

**10.3. Slide 6 - Clustering Algorithms.**

(1) Forward search (start chaos (lots of groups)  $\rightarrow$  purify)

(2) Backwards (start with 1 group  $\rightarrow$  chop into pieces)

**10.4. Slide 7 - Linkage Methods.**

With forward search we will always end up with 1 cluster. You will be needing to trim the search to your liking. It is exploratory.

**10.5. Slide 9 - An example.**

Note that the matrix is symmetric since metrics are symmetric.

Distance between cluster (12) and (3) and (4) is smallest.

**Anmärkning:**

For average linkage, we only look for distance between 2 clusters.

**10.6. Slide 16 - Dendrogram.**

Height tells us ordering of joining clusters

**10.7. Slide 18 - Wards Hierarchical Clustering Method.**

At every step you lose something. This tries to account for sum of squares.

Large sum of squares = distance is large

**10.8. Slide 20 - Pros and Cons.**

Monotone in the sense  $x \mapsto x^2, x \mapsto x^3$

**10.9. Slide 27 - CH Index and Gap Statistic.**

$W(K)$ , take euclidean distance to center and sum. Minimized when each point in its own cluster.

**Idea: Difference of clusters**

Collect distances in vector with closest first, second closest second etc...

Make  $B$  large, make  $W$  small.

Gap Statistic: If random cluster is similar, you're not doing well. Our model should beat a purely random model.

$a_i$  small,  $b_i$  large (= sparse points)

**10.10. Slide 29 - Parametric Approach: Mixing Distribution.**

Complete is best scenario

**10.11. Slide 30 - EM Algorithm.**

$\theta = \{p_k, \mu_k, \Sigma_k, \forall k\}$

### 10.12. Slide 31 - E step.

Need to know expected value of  $\log(f(\underbrace{X}_{\text{complete}}, \underbrace{Z}_{\text{distribution}}, \theta) | X; \underbrace{\hat{\theta}}_{\text{guessed}})$

Expectation is for  $Z|X$

$$\mathbb{E}(\log(\dots)) = \sum_{k=1}^K \log(f(X, Z; \theta)) P(Z = k | X; \hat{\theta}^{(t)})$$

By Bayes:

$$P(Z = k | X, \hat{\theta}^{(t)}) = \frac{f(X | Z = k, \hat{\theta}^{(t)}) \hat{P}_k^{(t)}}{\sum_{k=1}^K f(X | Z = k, \hat{\theta}^{(t)}) \hat{P}_k^{(t)}}$$

With normal assumption:

$$\Rightarrow \sum_{k=1}^K \left[ \underbrace{\log(p_k) - \frac{q}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}_{X|Z} \right] P(Z = k | X; \hat{\theta}^{(t)})$$

Because we have a sample of  $N$  observations, we get:

$$\sum_{j=1}^N \sum_{k=1}^K \log(f(x_j, Z; \theta)) P(Z_j = k | x_j; \hat{\theta}^{(t)}) = Q$$

In order to do M step (maximize Q):

$$\begin{aligned} \frac{\partial Q}{\partial p_k} &= \sum_{j=1}^N \left[ \frac{1}{p_k} P(Z_j = k | x_j, \hat{\theta}^{(t)}) - \frac{1}{p_k} P(Z_j = k | x_j; \hat{\theta}^{(t)}) \right] = 0 \\ &\Rightarrow p_k = \frac{\sum_{j=1}^N P(Z_j = k | x_j; \hat{\theta}^{(t)}) P_k}{\sum_{j=1}^N P(Z_j = K | x_j; \hat{\theta}^{(t)})} \\ &\Rightarrow \frac{\sum_{k=1}^{K-1} \sum_{j=1}^N P(Z_j = k | x_j; \hat{\theta}^{(t)}) P_k}{\sum_{j=1}^N P(Z_j = K | x_j, \hat{\theta}^{(t)})} = 1 \\ &\Rightarrow \hat{p}_k^{(t+1)} = \frac{\sum_{j=1}^N P(Z_j = k | x_j; \hat{\theta}^{(t)})}{\sum_{k=1}^K \sum_{j=1}^N P(Z_j = k | x_j; \hat{\theta}^{(t)})} \leftarrow \text{updated value at } p_k \end{aligned}$$

Next parameter is  $\mu_k$ , but that depends on the term  $-\frac{1}{2} (x_j - \mu_k)^T \Sigma_k^{-1} (x_j - \mu_k) \Rightarrow$  essentially want to maximize the following:

$$\sum_{j=1}^N (x_j - \mu_k)^T \Sigma_k^{-1} (x_j - \mu_k) P(Z_j = k | x; \hat{\theta}^{(t)})$$

Similar to Naive Bayes. We can use the lemma from Chapter 4. This yields to something complicated, but there exists a closed form expression of  $\mu_k^{(t+1)}$

Next parameter is  $\Sigma_k$ , we want to maximize:

$$-\frac{N}{2} \log(|\Sigma_k|) - \frac{1}{2} \sum_{j=1}^N (x_j - \mu_k)^T \Sigma_k^{-1} (x_j - \mu_k) P(Z_j = k | x_j; \hat{\theta}^{(t)})$$

to get  $\Sigma_k$

### 10.13. Slide 32 - M step.

Need to determine which cluster its from, but that is already done! Since the  $P(Z = k \dots)$  part gives probability of being in cluster  $k$

### 10.14. Slide 33 - EM Estimator.

We update using this  $P(Z = k \dots)$  probability. Like slide 15 in Gaussian DA.

**10.15. Slide 34 - Gaussian Mixture.**

Assume  $Z$  is multinomial:  $P(Z_j = k) = p_k$

$x_j | Z_j = k \sim N_q(\mu_k, \Sigma_k) \leftarrow$  conditional only even though marginals are not normal.

## 11. OLD EXAM 2022-01-05

1. a) Remember than any linear combination of normally distributed random variables is normal:

$$AX \sim N(A\mu, A\Sigma A^T)$$

So let  $A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$ :

$$A \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_2(A\mu, A\Sigma A^T)$$

- b)  $X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$ , then  $X_1|X_2 = a \sim N(\mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(a - \mu_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$

In our case,

$$\begin{aligned} X_1 &= \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} & \mu_1 &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ X_2 &= \begin{bmatrix} X_3 \\ X_4 \end{bmatrix} & \mu_2 &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \Sigma_{12} &= \begin{bmatrix} 1/4 & 0 \\ 0 & 0 \end{bmatrix} & \Sigma_{22} &= \begin{bmatrix} 2 & 1/2 \\ 1/2 & 2 \end{bmatrix} \\ \Sigma_{11} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} & \Sigma_{21} = \Sigma_{12}^T &= \begin{bmatrix} 1/4 & 0 \\ 0 & 0 \end{bmatrix} \end{aligned}$$

**Anmärkning:** Normally distributed variables are independent if and only if the covariance matrix is diagonal

2. a) Two-way MANOVA (CHapter 6).

The measured variables are  $\begin{bmatrix} O \\ V \\ C \\ B \end{bmatrix} = \underbrace{X}_{4 \times 1}$ , depends on brands of coffee and milk type, so we get:

$$X_{ijk} = \mu + \underbrace{\alpha_i}_{\text{coffee brand}} + \underbrace{\beta_j}_{\text{milk}} + \underbrace{r_{ij}}_{\text{interaction between brand \& milk}} + \underbrace{\varepsilon_{ijk}}_{\text{error}}$$

$k$  = repeated measures. Our model assumptions are assumptions of  $\alpha, \beta, r$

we need  $\sum \alpha_i = 0 \sum \beta_i = \sum r_{ij} = 0$

Error term should be normailly distributed  $N_4(0, \Sigma)$

- b) Testing wether interaction is 0 yields  $H_0 : r_{ij} = 0$  for any  $i, j$

In that case, the model becomes additive model (how the milk influences coffee is same regardless of brand) i.e  $X_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$

If we want to test wether brand has effect then we test  $\alpha_i = 0$

3. a) **Anmärkning:** How do you determine if it is paired data or uncorrolated?

In our case, think of it like this, if we take the sample from the same cup and add oat milk to 1/2 of the sample and cream to the other then it is paired since the coffee comes from the same sample.

In the case of paired data, we can use Hotellings  $T^2$ .

In our case however, we let  $X(4 \times 1)$  be the sample tested with cream, and  $Y(4 \times 1)$  be tested with oat milk. We have a 2-sample problem.



We test  $H_0 : \mu_X = \mu_Y$ , in order to do so we need some assumptions about the distributions as well as i.i.d:

- $X \sim N(\mu_X, \Sigma)$  (cream)
- $Y \sim N(\mu_Y, \Sigma)$  (oat milk)

Note that they have the same variance.

- b) For Hotellings  $T^2$  we need our hamburger 

Normal	Wishart	Normal
--------	---------	--------

:

$$(\bar{X} - \bar{Y})^T S^{-1} (\bar{X} - \bar{Y})$$

- c) In b), we had  $S^{-1}$  which comes from  $S_{\text{pooled}}$ , which can be estimated by (we can do this since  $\Sigma_X = \Sigma_Y$ ):

$$\frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2} \quad (3)$$

In order to determine the distribution of our estimation (3), we need to find the distribution of the terms:

- $(n_1 - 1)S_1 = \sum_{i=1}^{n_1} (X_i - \bar{X})(X_i - \bar{X})^T \sim W_4(\Sigma, n_1 - 1)$
- $(n_2 - 1)S_2 = \sum_{i=1}^{n_2} (X_i - \bar{X})(X_i - \bar{X})^T \sim W_4(\Sigma, n_2 - 1)$

We also know that Wishart + Wishart = Wishart if they are independent which we have assumed, so we get

$$(n_1 - 1)S_1 + (n_2 - 1)S_2 \sim W_4(\Sigma, n_1 + n_2 - 2)$$

With a scaling of  $\frac{1}{n_1 + n_2 - 2}$  in front.

- d) To estimate  $S_1$ , we need to ask ourselves if we are looking for a biased or unbiased estimate. We arrive at the following:

- *Biased*:  $\frac{1}{n_1} \sum_{i=1}^{n_1} (X_i - \bar{X})(X_i - \bar{X})^T$
- *Unbiased*:  $\frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_i - \bar{X})(X_i - \bar{X})^T$

4. a) Recall that in CCA, the main task is to find linear combinations of  $\left. \begin{matrix} U = a^T X \\ V = b^T Y \end{matrix} \right\}$  such that their

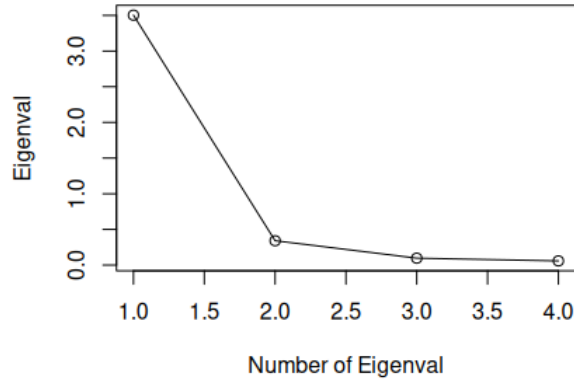
correlation is maximized and  $\left. \begin{matrix} \text{Var}(U) = 1 \\ \text{Var}(V) = 1 \end{matrix} \right\}$

- b) In order to find the canonical variates, we can use Result 10.1 but we would need the eigenvalues of  $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$  which we can compute using R:

```
eigen() decomposition
$values
[1] 0.5458401879 0.0008598121
```

- c) Notice that the second eigenvalue is numerically 0 (i.e it does not contribute that much).

5. a) A scree plot is essentially a plot of  $(x, y) = (\text{number of eigenvalue}, \text{eigenvalue})$ . In our case, we can plot it by hand or in R to get the following:



- b) There are numerous methods to decide how many components to include in the study. We can use the elbow-method in the scree plot or we can use the explained variance proportion which is given by:

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{i=1}^n \lambda_i} = \text{variance explained by the } k\text{th eigenvalues}$$

Personally I would have used 2 components, since this seems to be the elbow in the scree plot. Computing the variance explained yields:

$$\frac{3.50293864 + 0.34142136}{3.50293864 + 0.34142136 + 0.09706136 + 0.05857864} \approx 96\% \text{ explained}$$

- c) By result 8.3, the first component is given by  $[, 1]X$ , and second component by  $[, 2]X$
6. a) This is a matter of classification since we know the classes.  
The ECM is given by:

$$p_1 P(2|1) c(2|1) + p_2 P(1|2) c(1|2) = \frac{1}{2} \int_{R_2} f_1(x) dx c + \frac{1}{2} \int_{R_1} f_2(x) dx 2c$$

- b) & c)

$$R_1 = \left\{ x : \frac{f_1(x)}{f_2(x)} \geq \frac{c(1|2)p_2}{c(2|1)p_1} \right\} \Leftrightarrow \left\{ x : \frac{f_1(x)}{f_2(x)} \geq 2 \right\}$$

8. We have multivariate regression, our assumption is independence.

## 12. OLD EXAM - TRAINING EXAM

**Anmärkning:** There are certain questions that have come in the homeworks, see respective homework solution.

4. a) The test problem that has been considered here is if the mean is  $\begin{bmatrix} 60 \\ 29 \\ 60 \\ 1000 \end{bmatrix}$ . This is a one sample  $T^2$  test where we assume our data is i.i.d and normal.

- b) We know that we can express the  $T^2$  distribution using the  $F$  distribution through the following relationship:

$$T^2(p, n-1) \sim \frac{(n-1)p}{n-p} F_{p, n-p} \quad \left. \begin{array}{l} n-p=199 \\ p=4 \end{array} \right\} n=403$$

- c) Yes they can if we look individually for each  $X_i$ , no if they are joint.
- d) We are not testing any interaction between the data points, only if they have the same value. The MANOVA model is therefore specified through:

$$X_{ij} = \mu + \alpha_i + \varepsilon_{ij} \quad H_0 : \alpha_i = 0 \quad \forall i$$

- e) Our assumptions correspond to the normal MANOVA assumptions, that is  $\varepsilon_{ij}$  is multivariate normal, data is independent across stations.
- f) We can use 2 methods, either the elbow method (subjective) or we fix an acceptable percentage of explained variance such as 95%.

- g) We sort our  $n$  eigenvalues, and then compute the  $k$ th component by:  $\frac{\sum_{j=1}^k \lambda_j}{\sum_{i=1}^n \lambda_i}$

In our case, this yields  $\frac{3}{3+0.6+0.3+0.1} = 75\%$

- h)

$$X = \mu + LF + \varepsilon \quad \text{Cov}(X) = LL^T + \Psi$$

Under rotation, this becomes:

$$L \mapsto LT \Rightarrow (LT)(LT)^T + \Psi = L^T T^T L^T + \Psi$$

We can see here that the matrix  $T$  did not disappear, so the communalities after rotations are not the same as the communalities before rotation they need to be rotated as well.

- i) PCA is dimension reduction while FA is environmental/causal stuff.
5. a) Y-axis groups distance, so if three clusters are desired then we pick those top 3 highest  $Y$  values  $\Rightarrow 29, 6, 12$
- b) Using the max CH-index, we get 3 clusters.
- c) In  $k$ -means, the data needs to be standardized and the  $nstart$  is too low otherwise it would be unstable.

## 13. OLD EXAM 2021-01-11

2. a) A useful model could be  $\underbrace{X_{ij}}_{4 \times 1} = \mu + \underbrace{\alpha_i}_{\text{brand}} + \varepsilon_{ij}$ , here we test if the mean is the same across all  $\alpha_i = 0 \quad \forall i$

- b) Mathematically, the testing problem can be expressed as the following:

$$H_0 : \alpha_i = 0 \quad \forall i = 1, 2, 3, 4$$

Identification restriction is given by  $\Sigma \alpha_i = 0$

- c)

$$\sum_{i=1}^4 \sum_j (X_{ij} - \bar{X})(X_{ij} - \bar{X})^T = B + W$$

Here the  $B$  is the *between* sum of squares (compares observations between groups)

Here the  $W$  is the *within* sum of squares (compares observations between each group)

- d)  $H = : \mu = \mu_0 \quad X_1, \dots, X_n \sim N(\mu, \Sigma)$ . Testing normal mean

3. If  $m \neq n$ , then we can exclude the "paired data test" directly since they need to have the same dimensions

- a) We are therefore testing  $\mu_X = \mu_Y$

- b) Only  $S_{\text{pooled}}$  is for 2-sample test (which is what we have for this case)

What we do is assume  $\Sigma_1 = \Sigma_2$  and estimate by Chapter 6 p.10:

$$S_{\text{pooled}} = \frac{\sum_{j=1}^{n_1} (X_{1j} - \bar{X}_1)(X_{1j} - \bar{X}_1)^T + \sum_{j=1}^{n_2} (X_{2j} - \bar{X}_2)(X_{2j} - \bar{X}_2)^T}{n_1 + n_2 - 2} = \frac{(n_1 - 1)S_1}{n_1 + n_2 - 2} + \frac{(n_2 - 1)S_2}{(n_1 + n_2 - 2)}$$

Before we proceed with showing which distribution, it will be easier if we first do c) and then use the results here

- c)

$$S_1 = \frac{1}{n-1} \sum (X_{ij} - \bar{X})(X_{ij} - \bar{X})^T \quad \text{for unbiased}$$

$$S_1 = \frac{1}{n} \sum (X_{ij} - \bar{X})(X_{ij} - \bar{X})^T \quad \text{for MLE}$$

$$\Rightarrow \sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^T \sim W_p(\Sigma, n-1)$$

$$\Rightarrow S_1 \text{ is Wishart} \Leftrightarrow (n_1 - 1)S_1 \text{ is Wishart}$$

$$\Rightarrow S_{\text{pool}} \sim W_p(\Sigma, n_1 - 1) + W_p(\Sigma, n_2 - 1) \Rightarrow (n_1 + n_2 - 2)S_{\text{pool}} = W_p(\Sigma, n_1 + n_2 - 2)$$

Due to the independance, Wishart + Wishart is still Wishart.

If we test the mean of the first sample, can we get Hotellings  $T^2$ ?

$$\bar{X} \sim N(\mu, \Sigma) \Rightarrow n(\bar{X} - \mu_X) \sim N(0, \Sigma)$$

- d)  $\Sigma$  is common for both samples. We shall use some results from  $S_{\text{pool}}$ :

$$X_1, \dots, X_{n_1} \sim N\left(\mu_X, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$Y_1, \dots, Y_{n_2} \sim N\left(\mu_Y, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

If we only care about  $\Sigma_{11}$  then for every observation we take out the first 2 since:

$$\begin{aligned} \forall X_i \quad \begin{bmatrix} X_{i1} \\ X_{i2} \\ X_{i3} \\ X_{i4} \end{bmatrix} &\rightarrow \begin{bmatrix} X_{i1} \\ X_{i2} \end{bmatrix} \\ \forall Y_i \quad \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \\ Y_{i4} \end{bmatrix} &\rightarrow \begin{bmatrix} Y_{i1} \\ Y_{i2} \end{bmatrix} \end{aligned}$$

$\Rightarrow$  the estimator has distribution  $W_2(\Sigma_{11}, n_1 + n_2 - 2)$

- e) In order to use  $ny^T M^{-1}y \sim T^2(\Sigma, n)$ , we need  $\bar{x} \sim N(0, \Sigma)$  and  $M \sim W(\Sigma, n)$  (where they share  $\Sigma$ ):

$$\begin{aligned} \bar{x} &\sim N(\mu, \Sigma_{11}/n) \Rightarrow \sqrt{n}(\bar{x} - \mu_1) \sim N(0, \Sigma_{11}) = y \\ (n_1 + n_2 - 2)S_{\text{pool}} &= M \end{aligned}$$

For two-sample testing:

$$\begin{aligned} (\bar{x} - \bar{y}) &\sim N\left([\mu_{X_1} - \mu_{Y_1}], \left(\frac{1}{n_1} + \frac{1}{n_2}\right) \Sigma_{11}\right) \\ \Rightarrow \frac{1}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} [\bar{x} - \bar{y} - (\mu_{X_1} - \mu_{Y_1})] &\sim N(0, \Sigma_{11}) \end{aligned}$$

## 14. TIPS/RUMORS FOR THE EXAM

- Chapter 5, applying  $nx^T S^{-1}x \sim T^2$  will come on the exam
- Be able to state assumptions of models such as factor analysis
- Scale invariancy for different things is good to know
- For distributions, be able to tell what happens during scalar multiplication and addition of scalars
- (canonical correlation)<sup>2</sup> = eigenvectors/values  $\rho^2$
- Chapter 11; two class problems will come on exam
- Know the basic idea of LDA, QDA, logistic regression, RKHS, SVM, RF (when to use them)
- Derive ECM (we have done Gaussian)
- Interpret dendrograms and output of k-means
- EM-algorithm is a *must* for a higher grade
  - Get expression of  $Q$
  - Explore how EM can be used if distribution is normal/exponential
  - You should be able to explain how to get  $Q$  and maximize
- Find distribution or find normal
- R test (communals, output, interpretation)
- 6-8 points