TMA4267 Linear Statistical Models V2017 [L1]

Introduction to the course
Part 1: Multivariate random variables,
and the multivariate normal distribution

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 10, 2017

1/29

2/29

What is Statistics?

TMA4267 Linear Statistical Models

- Statistics, linear statistical models and movie recommender systems.
- Learning outcome.
- ► TMA4267 core and parts.
- ▶ Background knowledge in probability and statistical inference.
- ► TMA4267 course information.
- Voting and questionnaire.
- ► Part 1: Multivariate RVs and the multivariate normal distribution.

1/29

What is Statistics?

- ► The true foundation of theology is to ascertain the character of God.
- ▶ It is by the aid of *Statistics* that law in the social sphere can be ascertained and codified.
- ▶ and, certain aspects of the character of God hereby revealed.
- ▶ The study of statistics is thus a *religious service*.

What is Statistics?

- The true foundation of theology is to ascertain the character of God.
- ▶ It is by the aid of *Statistics* that law in the social sphere can be ascertained and codified.
- ▶ and, certain aspects of the character of God hereby revealed.
- ▶ The study of statistics is thus a *religious service*.

Florence Nightingale (1820-1910). Quotation from "Games, Gods and Gambling: A History of Probability and Statistical Ideas" by F. N. David.

2/29

Word cloud: Probability



What is Statistics?

- ▶ The goal of Statistics is to expand our knowledge based on collection and analysis of empirical data.
- ► Two branches:
 - Probability: the mathematical study of the probability of random events.
 - Statistical Inference: models and methods for collecting, describing, analysing and interpreting numerical data.



Drawing taken from http://www.nearingzero.net - now at http://www.lab-initio.com/

3 / 29

Word cloud: Statistical Inference



Linear Statistical Models

Simple linear regression (height of child explained by mid-parent height):

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

6 / 29

Linear Statistical Models

Simple linear regression (height of child explained by mid-parent height):

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

Multiple linear regression (also include other explanatory variables):

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

The multiple linear regression model is our linear statistical model! So, why is this course not called "Regression"?

Linear Statistical Models

Simple linear regression (height of child explained by mid-parent height):

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

Multiple linear regression (also include other explanatory variables):

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

The multiple linear regression model is our linear statistical model!

6 / 29

Linear Statistical Models

Simple linear regression (height of child explained by mid-parent height):

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

Multiple linear regression (also include other explanatory variables):

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

The multiple linear regression model is our linear statistical model! So, why is this course not called "Regression"? We include theory that focus on mathematical understanding: multivariate random variables, the multivariate normal distribution, projections, idempotent matrices, hypothesis tests, design of

experiments,

Recommender systems

▶ Recommender systems are a subclass of information filtering system that seek to predict the 'rating' or 'preference' that a user would give to an item.

7 / 29

Recommender systems

- ▶ Recommender systems are a subclass of information filtering system that seek to predict the 'rating' or 'preference' that a user would give to an item.
- ▶ Recommender systems have become extremely common in recent years, and are applied in a variety of applications. The most popular ones are probably movies, music, news, books, research articles, search queries, social tags, and products in general. However, there are also recommender systems for experts, collaborators, jokes, restaurants, financial services, life insurance, persons (online dating), and Twitter followers.

Recommender systems

- Recommender systems are a subclass of information filtering system that seek to predict the 'rating' or 'preference' that a user would give to an item.
- ▶ Recommender systems have become extremely common in recent years, and are applied in a variety of applications. The most popular ones are probably movies, music, news, books, research articles, search queries, social tags, and products in general. However, there are also recommender systems for experts, collaborators, jokes, restaurants, financial services, life insurance, persons (online dating), and Twitter followers.

7 / 29

Recommender systems

- ▶ Recommender systems are a subclass of information filtering system that seek to predict the 'rating' or 'preference' that a user would give to an item.
- ▶ Recommender systems have become extremely common in recent years, and are applied in a variety of applications. The most popular ones are probably movies, music, news, books, research articles, search queries, social tags, and products in general. However, there are also recommender systems for experts, collaborators, jokes, restaurants, financial services, life insurance, persons (online dating), and Twitter followers.

Source: Wikipedia: Recommender systems

The Netflix Price: 2006

Text from http:\www.netflixprice.com

▶ To help customers find those movies, we've developed our world-class movie recommendation system: Cinematch.

8 / 29

The Netflix Price: 2006

Text from http:\www.netflixprice.com.

- ► To help customers find those movies, we've developed our world-class movie recommendation system: Cinematch.
- ▶ Its job is to predict whether someone will enjoy a movie based on how much they liked or disliked other movies. We use those predictions to make personal movie recommendations based on each customer's unique tastes. And while Cinematch is doing pretty well, it can always be made better.
- ▶ We provide you with a lot of anonymous rating data, and a prediction accuracy bar that is 10% better than what Cinematch can do on the same training data set. (Accuracy is a measurement of how closely predicted ratings of movies match subsequent actual ratings.)

The Netflix Price: 2006

Text from http:\www.netflixprice.com.

- ► To help customers find those movies, we've developed our world-class movie recommendation system: Cinematch.
- ▶ Its job is to predict whether someone will enjoy a movie based on how much they liked or disliked other movies. We use those predictions to make personal movie recommendations based on each customer's unique tastes. And while Cinematch is doing pretty well, it can always be made better.

8 / 29

The Netflix Price: 2006

Text from http:\www.netflixprice.com.

- ► To help customers find those movies, we've developed our world-class movie recommendation system: Cinematch.
- ▶ Its job is to predict whether someone will enjoy a movie based on how much they liked or disliked other movies. We use those predictions to make personal movie recommendations based on each customer's unique tastes. And while Cinematch is doing pretty well, it can always be made better.
- ▶ We provide you with a lot of anonymous rating data, and a prediction accuracy bar that is 10% better than what Cinematch can do on the same training data set. (Accuracy is a measurement of how closely predicted ratings of movies match subsequent actual ratings.)
- Remark: At this point in time DVDs were sent to customers by mail - this was before the age of online streaming.

The Netflix Price: 2006

Text from http:\www.netflixprice.com

- ► To help customers find those movies, we've developed our world-class movie recommendation system: Cinematch.
- ▶ Its job is to predict whether someone will enjoy a movie based on how much they liked or disliked other movies. We use those predictions to make personal movie recommendations based on each customer's unique tastes. And while Cinematch is doing pretty well, it can always be made better.
- ▶ We provide you with a lot of anonymous rating data, and a prediction accuracy bar that is 10% better than what Cinematch can do on the same training data set. (Accuracy is a measurement of how closely predicted ratings of movies match subsequent actual ratings.)
- Remark: At this point in time DVDs were sent to customers by mail - this was before the age of online streaming.

8 / 29

9 / 29

Cinematch

The Cinematch recommender system: use statistical linear models with a lot of data conditioning. I have not found any other information on the algorithm online.

"Simple" linear suggestion:
(predicted score on movie for person)=
(some overall score for this movie)+
(some overall score used by this person)+
(similarity of this movie with other movie this person has seen)*
(how much this person liked that movie)+
the same for all the movies this person has rated+
error term.

The Netflix Price: 2006

Text from http:\www.netflixprice.com.

- ► To help customers find those movies, we've developed our world-class movie recommendation system: Cinematch.
- ▶ Its job is to predict whether someone will enjoy a movie based on how much they liked or disliked other movies. We use those predictions to make personal movie recommendations based on each customer's unique tastes. And while Cinematch is doing pretty well, it can always be made better.
- ▶ We provide you with a lot of anonymous rating data, and a prediction accuracy bar that is 10% better than what Cinematch can do on the same training data set. (Accuracy is a measurement of how closely predicted ratings of movies match subsequent actual ratings.)
- Remark: At this point in time DVDs were sent to customers by mail - this was before the age of online streaming.

The prize was awarded the team *BellKor's Pragmatic Chaos* in 2009.

8 / 29

10 / 29

The Netflix Price: Training data

- ► The training data set consists of more than 100 million ratings from over 480 thousand randomly-chosen, anonymous customers on nearly 18 thousand movie titles.
- The ratings are on a scale from 1 to 5 (integral) stars. The date of each rating and the title and year of release for each movie are provided.
- No other customer or movie information is provided. No other data were employed to compute Cinematch's accuracy values used in this Contest.

Text from http:\www.netflixprice.com.

The Netflix Price: Test data

- A qualifying test set is provided containing over 2.8 million customer/movie id pairs with rating dates but with the ratings withheld.
- ► Eligible algorithms must provide predictions for all the withheld ratings for each customer/movie id pair in the qualifying set.
- ► The qualifying set is divided into two disjoint subsets containing randomly selected pairs from the qualifying set. The assignment of pairs to these subsets is not disclosed.
 - ► The Site will score each subset by computing the square root of the averaged squared difference between each prediction and the actual rating (the root mean squared error or "RMSE") in the subset, rounded to the nearest .0001.
 - The RMSE for the first "quiz" subset will be reported publicly on the Site,
 - the RMSE for the second "test" subset will not be reported publicly but will be employed to qualify a submission as described below

Text from http:\www.netflixprice.com.

11 / 29

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

► Most entries into the competition looked at the problem as a set of algorithms – focus on prediction rather than on understanding what *drives* the preditions.

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

12 / 29

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

- ► Most entries into the competition looked at the problem as a set of algorithms focus on prediction rather than on understanding what *drives* the preditions.
- ► Complex models are prone to over fitting or matching small details rather than the big picture, especially where data are scarce (importance of cross-validation).

12 / 29

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

- ► Most entries into the competition looked at the problem as a set of algorithms focus on prediction rather than on understanding what *drives* the preditions.
- ► Complex models are prone to over fitting or matching small details rather than the big picture, especially where data are scarce (importance of cross-validation).
- ➤ The final model is an *ensemble model* combining many different prediction models (at least more than 100), including nearest neighbour methods, latent factor models, neural networks, weighting determined by ridge regression.

12 / 29

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

- ► Most entries into the competition looked at the problem as a set of algorithms focus on prediction rather than on understanding what *drives* the preditions.
- ► Complex models are prone to over fitting or matching small details rather than the big picture, especially where data are scarce (importance of cross-validation).
- ► The final model is an *ensemble model* combining many different prediction models (at least more than 100), including nearest neighbour methods, latent factor models, neural networks, weighting determined by ridge regression.
- ► The winning model was never implemented by Netflix, partly due to implementation issues but also due to the increase of available data after "sending DVDs by mail" was replaced by online streaming.

12 / 29

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

- ► Most entries into the competition looked at the problem as a set of algorithms focus on prediction rather than on understanding what *drives* the preditions.
- ► Complex models are prone to over fitting or matching small details rather than the big picture, especially where data are scarce (importance of cross-validation).
- ➤ The final model is an *ensemble model* combining many different prediction models (at least more than 100), including nearest neighbour methods, latent factor models, neural networks, weighting determined by ridge regression.
- ► The winning model was never implemented by Netflix, partly due to implementation issues but also due to the increase of available data after "sending DVDs by mail" was replaced by online streaming.

12 / 29

The winning algorithm: lessons to learn

Bell, Koren and Volinsky (2010): All Together Now: A Perspective on the Netflix Price, Chance, 23, p. 24-29.

- ▶ Most entries into the competition looked at the problem as a set of algorithms focus on prediction rather than on understanding what *drives* the preditions.
- ► Complex models are prone to over fitting or matching small details rather than the big picture, especially where data are scarce (importance of cross-validation).
- ► The final model is an *ensemble model* combining many different prediction models (at least more than 100), including nearest neighbour methods, latent factor models, neural networks, weighting determined by ridge regression.
- ► The winning model was never implemented by Netflix, partly due to implementation issues but also due to the increase of available data after "sending DVDs by mail" was replaced by online streaming.

Read more: Link to talk with interesting points raised.

TMA4267 Linear statistical methods Learning outcome, Knowledge

► The student has strong theoretical knowledge about the most popular statistical models and methods that are used in science and technology, with emphasis on regression-type statistical models.

13 / 29

TMA4267 Linear statistical methods Learning outcome, Skills

▶ The student knows how to design an experiment and

TMA4267 Linear statistical methods Learning outcome, Knowledge

- ► The student has strong theoretical knowledge about the most popular statistical models and methods that are used in science and technology, with emphasis on regression-type statistical models.
- ➤ The statistical properties of the multivariate normal distribution are well known to the student, and the student is familiar with the role of the multivariate normal distribution within linear statistical models.

13 / 29

TMA4267 Linear statistical methods Learning outcome, Skills

- ▶ The student knows how to design an experiment and
- how to collect informative data of high quality to study a phenomenon of interest.

TMA4267 Linear statistical methods Learning outcome, Skills

- ▶ The student knows how to design an experiment and
- how to collect informative data of high quality to study a phenomenon of interest.
- Subsequently, the student is able to choose a suitable statistical model,

14 / 29

TMA4267 Linear statistical methods Learning outcome, Skills

- ▶ The student knows how to design an experiment and
- how to collect informative data of high quality to study a phenomenon of interest.
- Subsequently, the student is able to choose a suitable statistical model.
- ▶ apply sound statistical methods, and
- perform the analyses using statistical software.

TMA4267 Linear statistical methods Learning outcome, Skills

- ▶ The student knows how to design an experiment and
- how to collect informative data of high quality to study a phenomenon of interest.
- ► Subsequently, the student is able to choose a suitable statistical model.
- ▶ apply sound statistical methods, and

14 / 29

TMA4267 Linear statistical methods Learning outcome, Skills

- ▶ The student knows how to design an experiment and
- how to collect informative data of high quality to study a phenomenon of interest.
- Subsequently, the student is able to choose a suitable statistical model.
- ▶ apply sound statistical methods, and
- ▶ perform the analyses using statistical software.
- ► The student knows how to present the results from the statistical analyses, and how to draw conclusions about the phenomenon under study.

TMA4267: Parts

- ► Part 1: Multivariate RVs and the multivariate normal distribution [week 2-5].
 - Data consists of simultaneous measurements on many variables: we work with random vectors and random matrices.
 - ► There is a strong connection between the *multivariate normal* distribution and the classical linear model.
- ▶ Part 2: The classical linear model [week 6-9]
 - ► We want to understand the relationship between many variables: with focus on linear relationships through the classical linear model (multiple linear regression).
- ▶ Part 3: Hypothesis tests and analysis of variance [week 9-11]
 - ► Is there and association between a response and an explanatory variable? Does a response vary between treatment groups?
- ▶ Part 4: Design of Experiments [week 12-13+project]
 - ▶ If we want to collect data, we need to do know how to design an experiment.

15 / 29

Do you know this?

Recommended background: TMA4240/TMA4245 Statistics.

- ▶ Probability: (continuous) random variables (RV), probability distribution function (pdf), cumulative distribution function (cdf), mean E, variance Var, covariance Cov, correlation Corr, transformation formula, momentgenerating function (MFG), normal, chi-square and t-distributions.
- ▶ Inference: population and sample philosophy, parameter estimation, confidence interval, hypothesis test, *p*-value, simple linear regression.

Do you know this?

Recommended background: TMA4240/TMA4245 Statistics.

Probability: (continuous) random variables (RV), probability distribution function (pdf), cumulative distribution function (cdf), mean E, variance Var, covariance Cov, correlation Corr, transformation formula, momentgenerating function (MFG), normal, chi-square and t-distributions.

16 / 29

Do you know this?

Recommended background: TMA4240/TMA4245 Statistics.

- ▶ Probability: (continuous) random variables (RV), probability distribution function (pdf), cumulative distribution function (cdf), mean E, variance Var, covariance Cov, correlation Corr, transformation formula, momentgenerating function (MFG), normal, chi-square and t-distributions.
- Inference: population and sample philosophy, parameter estimation, confidence interval, hypothesis test, p-value, simple linear regression.
- ► Linear methods: vector and matrix algebra (trace, determinant, eigenvalues/vectors), real vector spaces, orthogonality, spectral decomposition.

TMA4267 Linear Statistical Models Course information

https://innsida.ntnu.no/bb

- ► Course information.
- Course material.
- ► Lectures (and handouts).
- ► Statistical software.
- ► Exercises (6 recommended and 4 compulsory).
- ► Exam (80% of portfolio assessment).

17 / 29

Is this the correct course for you?

Are you afraid that this course have a too strong focus on theory and to little on the practical aspects of statistics? You may also look at at the following similar courses (that is, a second course in statistics, with focus on inference)

- ► ST2304 Statistical modelling for biology/biotechnology: https://wiki.math.ntnu.no/st2304/
- ► TMA4255 Applied statistics, for all siv.ing. studiprograms (except IndMat): https://wiki.math.ntnu.no/tma4255/
- ► KLMED Medical statistics II: https://www.ntnu.no/studier/emner/KLMED8005

Is this the correct course for you?

Are you afraid that this course have a too strong focus on theory and to little on the practical aspects of statistics?

18 / 29

Electronic voting

- more than an anonymous show of hands?

For student: check that topics are understood, compare to class,

focus on the question asked, while preserving

anonymity.

For lecturer: collect data to design sessions that are more

contingent.

Software: clicker. math.ntnu.no (single questions), Kahoot!

(end-of-lecture sum-up), quiz in Blackboard.

Future studies?

What is your current plan of topic for future studies?

- ► A: Statistics
- ► B: Mathematics
- ► C: Numerics
- D: Other
- ► E: Don't know

20 / 29

Part 1: Multivariate random vectors and the multivariate normal distribution

- Härdle and Simar (2015): Applied Multivariate Statistical Analysis. Springer.
 - Chapter 2 (p. 53-76): A Short Excursion into Matrix Algebra (partly lectured, manly assumed known).
 - ► Chapter 3.3 (p. 89-93): Summary statistics.
 - ► Chapter 4.1-4.5 (p. 117-149): Multivariate Distributions.
 - Chapter 5.1 (p. 183-190): Elementary Properties of the Multinormal.
- ► Fahrmair, Kneib, Lang and Marx (2013): Regression. Springer.
 - Appendix B: Def B.11 (chis q), B.13 (t), B14 (F), Theorem B.2 and B3.3 (distribution of quadratic forms).

A merged pdf named TMA4267Part1.pdf is available from Bb. Both eBooks and can be downloaded without charge for NTNU students.

22 / 29

Electronic voting

Use your smart phone, or other devise with internet access and go to http://clicker.math.ntnu.no/, and then select TMA4267 as classroom

Answers

- A: Statistics
- ▶ B: Mathematics
- ► C: Numerics
- ▶ D. Other
- ► E: Don't know

Start voting now!

21 / 29

The Cork deposit data

- ► Classical data set from Rao (1948).
- ▶ Weigth of bark deposits of n = 28 cork trees in p = 4 directions (N, E, S, W).

Tree	N	N E		W
1	72	66	76	77
2	60	53	66	63
3	56	57	64	58
Ė	:	:	:	:
28	48	54	57	43

How may we define a random vector in connection to the cork deposit data set?

Hands-on

Let $X_{(2\times 1)}$ have joint pdf (see the 3D-printed figure)

$$f(x_1, x_2) = \frac{1}{2\pi} e^{-\frac{1}{2}(x_1^2 + x_2^2)}$$
 for $-\infty < x_1, x_2 < \infty$

Find:

- 1. the marginal distributions $f_1(x_1)$ and $f_2(x_2)$,
- 2. the conditional distributions $f(x_1 \mid x_2)$ and $f(x_2 \mid x_1)$.
- 3. What about $F(x_1, x_2)$?

Hint (why?):

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx = 1$$

24 / 29

Sklar's Theorem [H4.1, p121-122]

Let F be a joint (cumulative) distribution function with marginal distribution functions F_1 and F_2 . Then a copula C exists with

$$F(x_1, x_2) = C(F_1(x_1), F_2(x_2))$$

for every $x_1, x_2 \in \mathbb{R}$. If F_1 and F_2 are continuous, then C is unique. On the other hand, if C is a copula and F_1 and F_2 are (cumulative) distribution functions, then the function F defined above, is a joint distribution function with marginals F_1 and F_2 .

Remark: this is not part of the core of the course (not suitable as an exam question), but it is a nice concept and you should have heard about it.

26 / 29

Copula [H4.1, p120]

A two-dimensional copula is a function $C:[0,1]^2 \to [0,1]$ with the following properties:

- ▶ For every $u \in [0,1]$: C(0,u) = C(u,0) = 0.
- ▶ For every $u \in [0,1]$: C(u,1) = u and C(1,u) = u.
- ▶ For every $(u_1, u_2), (v_1, v_2) \in [0, 1] \times [0, 1]$ with $u_1 \leq v_1$ and $u_2 \leq v_2$:

$$C(v_1, v_2) - C(v_1, u_2) - C(u_1, v_2) + C(u_1, u_2) \ge 0$$

(The last property is called "2-increasing".)

Remark: this is not part of the core of the course (not suitable as an exam question), but it is a nice concept and you should have heard about it.

25 / 29

Bivariate Copulas

Farlie-Gumbel-Morgenstern family

$$C(u, v) = uv + \theta uv(1-u)(1-v), \ \theta \in [-1, 1]$$

The only copulas that are polynomial quadratic in \boldsymbol{u} and \boldsymbol{v} , symmetric.

Normal (Gaussian) copulas [H.p141]

$$f(x_1, x_2) = \frac{1}{2\pi} \frac{1}{\sigma_1 \sigma_2 \sqrt{1 - \rho^2}} e^{-\frac{1}{2}Q(x_1, x_2)}$$

$$Q(x_1, x_2) = \frac{1}{1 - \rho^2} \left[\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1} \right) \left(\frac{x_2 - \mu_2}{\sigma_2} \right) \right]$$

$$C(u, v) = \int_{-\infty}^{\Phi_1^{-1}(u)} \int_{-\infty}^{\Phi_2^{-1}(v)} f(x_1, x_2) dx_1 dx_2$$

Read more? Properties and applications of copulas: A brief survey, Roger B. Nelsen (And same remark as before.)

More on copulas?

We will later in Part 1 (recommended exercise 2) look at data and contour plots from Gaussian copulas, and other copulas poplar in finance - using R and the copula library in R.

28 / 29

What have we worked with today?

- ► A random vector is ... a vector of random variables.
- Joint distribution function.

What have we worked with today?

► A random vector is ... a vector of random variables.

29 / 29

What have we worked with today?

- ► A random vector is ... a vector of random variables.
- ▶ Joint distribution function.
- From joint distribution function to marginal and conditional distributions.

29 / 29

What have we worked with today?

- ▶ A random vector is ... a vector of random variables.
- ▶ Joint distribution function.
- From joint distribution function to marginal and conditional distributions.
- Cumulative distribution.

29 / 29

What have we worked with today?

- ▶ A random vector is ... a vector of random variables.
- ▶ Joint distribution function.
- From joint distribution function to marginal and conditional distributions.
- Cumulative distribution.
- ► Independence.
- From marginal cumulative distribution functions to joint using copula.

What have we worked with today?

- ▶ A random vector is ... a vector of random variables.
- ▶ Joint distribution function.
- From joint distribution function to marginal and conditional distributions
- Cumulative distribution.
- ► Independence.

29 / 29

What have we worked with today?

- ▶ A random vector is ... a vector of random variables.
- ▶ Joint distribution function.
- From joint distribution function to marginal and conditional distributions.
- ► Cumulative distribution.
- ► Independence.
- From marginal cumulative distribution functions to joint using copula.

What have we worked with today?

- ▶ A random vector is ... a vector of random variables.
- Joint distribution function.
- From joint distribution function to marginal and conditional distributions
- Cumulative distribution.
- ► Independence.
- From marginal cumulative distribution functions to joint using copula.

Next lecture: Mean vector and covariance matrix. You may want to look into how to define a positive definite matrix, how to define eigenvalues/vectors and results for symmetric matrices.

29 / 29

Lat
$$X = \begin{bmatrix} X_A \\ X_B \end{bmatrix}$$
 $X_A \in \mathbb{R}^k$
 $X_B \in \mathbb{R}^k$
 $X_B \in \mathbb{R}^{k-1}$

3) X_A has margoral distribution functor

 $f_A(x_A) = \iint_{\mathbb{R}^2} \int_{\mathbb{R}^2} f(x) dx_B$

"Integrating out the vorable rot of intent"

4) and conditional distribution of $X_B = X_A$
 $f(x_B \mid x_A) = \int_{\mathbb{R}^2} f(x_A, x_B) f(x_A) > 0$
 $Ex: \int_{\mathbb{R}^2} f(x_B, x_B) = \frac{1}{2\pi\pi} e^{-\frac{1}{2}(x_B^2 + x_B^2)}$
 $f(x_B \mid x_A) = \int_{\mathbb{R}^2} e^{-\frac{1}{2}(x_B^2 + x_B^2)} dx_B$

1) find $f_1(x_A) = \int_{\mathbb{R}^2} f(x_A, x_B) dx_B$
 $f(x_B \mid x_A) = \int_{\mathbb{R}^2} f(x_A, x_B) dx_B$

Port 1: multivariate random vectors and the multivariate normal distribution [14] looks Rondom vector, fix), fix)

Rendom vector = vector of rendom variables (RU)

We will focus on continuous random vonables.

- 1) It has a joint distribution (density) function (pdf) $f(x) = f(x_1, x_2, ..., x_q)$ Ramamber $f(x_1) \ge 0$ $\forall x$ $\iint \int_{0}^{\infty} f(x_1, x_2, ..., x_q) dx_1 \cdots dx_q = 1$
- 2) I has a joint cumulative distribution function (cdf) $F(x): P(X \in x) = P(X_1 \in X_1, X_2 \in x_2, X_2 \in x_2)$



$$= \int_{0}^{\infty} \frac{1}{\ln x} e^{-\frac{1}{2}x_{1}^{2}} dx_{2}$$

$$= \int_{0}^{\infty} \frac{1}{\ln x} e^{-\frac{1}{2}x_{1}^{2}} dx_{2}$$

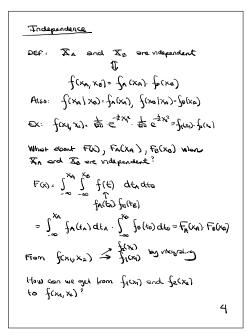
$$= \int_{0}^{\infty} \frac{1}{\ln x} e^{-\frac{1}{2}x_{1}^{2}} dx_{1}$$

$$= \frac{1}{\ln e^{-\frac{1}{2}x_{k}}} = \int_{2}^{x_{k}} (x_{k})$$

3)
$$F(x) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int (t_0, t_0) dt_1 dt_2 = \cdots$$

$$F(x_0, x_0) \qquad \text{not on closed from, under Resolutions}$$

3



TMA4267 Linear Statistical Models V2017 [L2]

Part 1: Multivariate RVs, and the multivariate normal distribution Moments: mean and covariance [H:4.2]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 13, 2017

fich) and fector of f(x1, x2). fich) fich)

XI and X2 independent

But what if XI and X2 are not independent?

Solution is Capula > see share!

Last lecture

1/15

- A random vector X_(p×1) is ... a p-dimensional vector of random variables.
 - Weight of cork deposits in p = 4 directions (N, E, S, W).
 - Rent index in Munich: rent, area, year of construction, location, bath condition, kitchen condition, central heating, district
- ▶ Joint distribution function: f(x).
- From joint distribution function to marginal (and conditional distributions).

$$f_1(x_1) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, x_2, \dots, x_p) dx_2 \cdots dx_p$$

- Cumulative distribution (definite integrals!) used to calculate probabilites.
- ▶ Independence: $f(x_1, x_2) = f_1(x_1) \cdot f(x_2)$ and $f(x_1 | x_2) = f_1(x_1)$.
- From marginal cumulative distribution functions to joint using copula.

Word cloud: Probability



2/15

The Cork deposit data

- ► Classical data set from Rao (1948).
- ▶ Weigth of bark deposits of n = 28 cork trees in p = 4 directions (N, E, S, W).

Tree	Ν	Ε	S	W
1	72	66	76	77
2	60	53	66	63
3	56	57	64	58
÷	:	÷	:	÷
28	48	54	57	43

How may we define a random vectors and random matrices for cork trees?

Today

- ► Moments: important properties about the distribution of **X**.
- ► E: Mean of random vector and random matrices.
- ► Cov: Covariance matrix.
- ► Corr: Correlation matrix.
- ▶ E and Cov of multiple linear combinations.

3 / 15

The Cork deposit data

Draw a random sample of size n = 28 from the population of cork treed and observe a p = 4 dimensional random vector for each tree.

$$\boldsymbol{X}_{(28\times4)} = \left[\begin{array}{ccccc} X_{11} & X_{12} & X_{13} & X_{14} \\ X_{21} & X_{22} & X_{23} & X_{24} \\ X_{31} & X_{32} & X_{33} & X_{34} \\ \vdots & \vdots & \ddots & \vdots \\ X_{28,1} & X_{28,2} & X_{28,3} & X_{28,4} \end{array} \right]$$

and
$$E(\boldsymbol{X}) = \{E(X_{ij})\}.$$

Random vectors and matrices: rules for means

▶ Random vector $\boldsymbol{X}_{(p\times 1)}$ with mean vector $\boldsymbol{\mu}_{(p\times 1)}$:

$$m{X}_{(
ho imes 1)} = \left[egin{array}{c} X_1 \ X_2 \ dots \ X_{
ho} \end{array}
ight], \qquad m{\mu}_{(
ho imes 1)} = \mathrm{E}(m{X}) = \left[egin{array}{c} \mathrm{E}(X_1) \ \mathrm{E}(X_2) \ dots \ \mathrm{E}(X_{
ho}) \end{array}
ight]$$

▶ 1) Random matrix $\boldsymbol{X}_{(n \times p)}$ and random matrix $\boldsymbol{Y}_{(n \times p)}$:

$$E(\boldsymbol{X} + \boldsymbol{Y}) = E(\boldsymbol{X}) + E(\boldsymbol{Y})$$

2) Random matrix X_(n×p) and conformable constant matrices A and B:

$$E(AXB) = AE(X)B$$

6 / 15

Hands-on

Let $\boldsymbol{X}_{4 \times 1}$ have variance-covariance matrix

$$\mathbf{\Sigma} = \left[egin{array}{cccc} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 1 \\ 0 & 0 & 2 & 1 \\ 0 & 1 & 1 & 2 \end{array}
ight].$$

Explain to your neighbour what this means.

Variance-covariance matrix

▶ Random vector $\boldsymbol{X}_{(p\times 1)}$ with mean vector $\boldsymbol{\mu}_{(p\times 1)}$:

$$m{X}_{(m{
ho} imes 1)} = \left[egin{array}{c} X_1 \ X_2 \ dots \ X_{m{
ho}} \end{array}
ight], \quad m{\mu}_{(m{
ho} imes 1)} = \left[egin{array}{c} \mathrm{E}(X_1) \ \mathrm{E}(X_2) \ dots \ \mathrm{E}(X_{m{
ho}}) \end{array}
ight] = \left[egin{array}{c} \mu_1 \ \mu_2 \ dots \ \mu_{m{
ho}} \end{array}
ight]$$

Variance-covariance matrix Σ (real and symmetric)

$$\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \mathrm{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{12} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1p} & \sigma_{2p} & \cdots & \sigma_{pp} \end{bmatrix}$$

7 / 15

Correlation matrix

Correlation matrix ρ (real and symmetric)

$$\boldsymbol{\rho} = \begin{bmatrix} \frac{\sigma_{11}}{\sqrt{\sigma_{11}\sigma_{11}}} & \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}} & \cdots & \frac{\sigma_{1p}}{\sqrt{\sigma_{11}\sigma_{pp}}} \\ \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}} & \frac{\sigma_{22}}{\sqrt{\sigma_{22}\sigma_{22}}} & \cdots & \frac{\sigma_{2p}}{\sqrt{\sigma_{22}\sigma_{pp}}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_{1p}}{\sqrt{\sigma_{11}\sigma_{pp}}} & \frac{\sigma_{2p}}{\sqrt{\sigma_{22}\sigma_{pp}}} & \cdots & \frac{\sigma_{pp}}{\sqrt{\sigma_{pp}\sigma_{pp}}} \end{bmatrix} = \begin{bmatrix} 1 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{12} & 1 & \cdots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{1p} & \rho_{2p} & \cdots & 1 \end{bmatrix}$$

$$oldsymbol{
ho} = (oldsymbol{V}^{rac{1}{2}})^{-1} oldsymbol{\Sigma} (oldsymbol{V}^{rac{1}{2}})^{-1}, ext{ where } oldsymbol{V}^{rac{1}{2}} = \left[egin{array}{cccc} \sqrt{\sigma_{11}} & 0 & \cdots & 0 \ 0 & \sqrt{\sigma_{22}} & \cdots & 0 \ dots & dots & \ddots & dots \ 0 & 0 & \cdots & \sqrt{\sigma_{pp}} \end{array}
ight]$$

Hands-on

Let $X_{4\times 1}$ have variance-covariance matrix

$$\mathbf{\Sigma} = \left[\begin{array}{cccc} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 1 \\ 0 & 0 & 2 & 1 \\ 0 & 1 & 1 & 2 \end{array} \right].$$

Find the correlation matrix.

10 / 15

12 / 15

Hands-on: Focus on **C**?

$$\boldsymbol{X} = \begin{bmatrix} X_N \\ X_E \\ X_S \\ X_W \end{bmatrix}, \ \boldsymbol{\mu} = \begin{bmatrix} \mu_N \\ \mu_E \\ \mu_S \\ \mu_W \end{bmatrix}, \ \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{NN} & \sigma_{NE} & \sigma_{NS} & \sigma_{NW} \\ \sigma_{NE} & \sigma_{EE} & \sigma_{ES} & \sigma_{EW} \\ \sigma_{NS} & \sigma_{EE} & \sigma_{SS} & \sigma_{SW} \\ \sigma_{NW} & \sigma_{EW} & \sigma_{SW} & \sigma_{WW} \end{bmatrix}$$

- ▶ Scientists would like to compare the following three contrasts: N-S, E-W and (E+W)-(N+S),
- ▶ and define a new random vector $Y_{(3\times1)} = C_{(3\times4)}X_{(4\times1)}$ giving the three contrasts.
- ▶ Write down *C*.
- ▶ Use the formulas we just developed and explain how to find $E(Y_1)$ and $Cov(Y_1, Y_3)$.

Linear combinations

- lacktriangleright Random vector $oldsymbol{\mathcal{X}}_{(p imes1)}$ with mean vector $oldsymbol{\mu_{oldsymbol{X}}}=\mathrm{E}(oldsymbol{X})$ and variance-covariance matrix $\Sigma_{\mathbf{X}} = \operatorname{Cov}(\mathbf{X})$.
- ightharpoonup The linear combinations Z = CX have

$$\mu_{Z} = E(Z) = E(CX) = C\mu_{X}$$

 $\Sigma_{Z} = Cov(Z) = Cov(CX) = C\Sigma_{X}C^{T}$

11 / 15

Exam V2014: Problem 1a

Let
$$\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$$
 be a random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\boldsymbol{X}) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ and covariance matrix $\boldsymbol{\Sigma} = \mathrm{Cov}(\boldsymbol{X}) = \boldsymbol{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$. Further, let

and covariance matrix
$$\mathbf{\Sigma} = \operatorname{Cov}(\mathbf{X}) = \mathbf{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$
. Further, let

$$\mathbf{A} = \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{pmatrix} \text{ be a matrix of constants.}$$

Define
$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \end{pmatrix} = \mathbf{AX}$$
.

Find $E(\mathbf{Y})$ and $Cov(\mathbf{Y})$.

Are X_1 and X_2 independent?

Are Y_1 and Y_2 independent? Justify your answers.

Find the mean of $\mathbf{X}^T \mathbf{A} \mathbf{X}$.

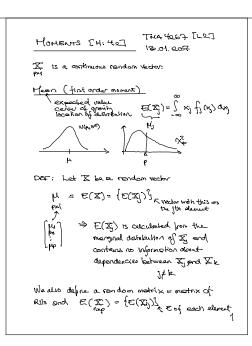
The covariance matrix

Random vector $\pmb{X}_{(\rho \times 1)}$ with mean vector $\pmb{\mu}_{(\rho \times 1)}$ and covariance matrix

$$\mathbf{\Sigma} = \operatorname{Cov}(\mathbf{X}) = \operatorname{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{12} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1p} & \sigma_{2p} & \cdots & \sigma_{pp} \end{bmatrix}$$

The covariance matrix is by construction symmetric, and we would only consider covariance matrices that are positive definite (PD). Why would we only consider PD matrices? Homework for next lecture: Read H.Chapter 2.1-2.2 to remind yourself of spectral decomposition (diagonalization), positive definite matrix, eigenvalues and eigenvectors.

14 / 15



What have we worked with today?

- ▶ Mean: $\mu_X = E(X) = E(X_i)$
- ► Covariance: $Cov(\boldsymbol{X}, \boldsymbol{Y}) = E((\boldsymbol{X} \boldsymbol{\mu}_{\boldsymbol{X}})(\boldsymbol{Y} \boldsymbol{\mu}_{\boldsymbol{Y}})^T)$.
- Variance-covariance: $\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \mathrm{E}((\mathbf{X} \mu_X)(\mathbf{X} \mu_X)^T)$, also sometimes denoted $\mathrm{Var}(\mathbf{X})$.
- ► Correlation: $Corr(X) = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}}$
- ▶ CX: $E(CX) = C\mu_X$ and $Cov(CX) = C\Sigma C^T$.

Next lecture: First work with the covariance matrix and positive definiteness, then start with the multivariate normal distribution (where we use moment generating functions and a multivariate version of the transformation formula).

15 / 15

RULES FOR MEAN

I and Y are random matrices (n hey bed)

and A and B are constant matrices

property

1)
$$E(X+Y) = E(X) + E(Y)$$

froof: Look at alament $Z:j = X_{ij} + Y_{ij}$ and see that $E(Zi_{j}) = E(Xi_{j} + Y_{i}) = E(Xi_{j}) + E(Yi_{j})$

2) $E(A X B) = A E(X) B$

Proof: look at alament (i,j) of AXX

 $E:j = \sum_{k=1}^{n} \alpha_{ik} \sum_{l=1}^{n} X_{kl} \text{ bij}$

and see that $E(e_{ij})$ is the alament (ij) of $A E(X)B$.

2

Coveriance metrix [44.2]

- (1) From THAYBYD/YBRSTALOY: X end Y ETE Sceler RUS with E(X)= phr and E(Y)= phr.
 - a) Cov (X,Y)= E((X-px)(Y-px))

 E(X:Y) pxpy

Ramanber: Car (35, 4) > 0: high \$7 and high? < 0: high x and low? = 0: no linear poten

- b) Var(8)= Car(8,8)= E(x2) 1/2
- d) (av (5,7)=0 tf X-and Y are independent, but not necessarily the opposite
- (2) Now: Pendom Vectors Ex Chemical procus

 a) DEF [H4.16] Fe guilt, specific for any section rendom vectors.

 yet and y section rendom vectors.

bxi dxi

E(X)= Mx and E(Y)= My

S

Example: Yx4 metrix

Observations: $Cor(X_1, X_2) = Cor(X_2, X_1)$ = $E[(X_1, \mu_1)(X_2, \mu_2)]$ = $\iint_{\mathbb{R}^n} (x_1, \mu_1)(x_2, \mu_2) f(x_1, \mu_2) dx_1 dx_2$

I is a symmetric motrix.

Cont

C)
$$Cov(X,Y) = E(XY^T) - \mu_x, \mu_x^T$$

$$eccEx4.82$$

$$Cov(X) = E(XX^+) - \mu_x \mu_x^T$$

$$ex8$$

$$ex8$$

$$ex8$$

$$ex8$$

$$ex8$$

d) As before if X_i end Y_i ere independent $\Rightarrow (av(X_i,Y_i)=0)$. And if all pairs of X_i and Y_i are yidospendent $\Rightarrow (av(X,Y_i)=0)$

$$Cr(X,Y) = E \left[(X-\mu_X)(Y-\mu_Y)^T \right]$$

$$= \sum_{\substack{x \in X \\ x \neq y \neq x = x \\ x \neq y \neq x = x \\ x \neq y \neq x = x }} erd Y era$$

$$erd Cov(X,Y) = Cr(X_v,Y_v) Cv(X_v,Y_v) Cv(X_v,Y_v)$$

$$freq Cov(X_v,Y_v) ev(X_v,Y_v) Cv(X_v,Y_v)$$

$$freq Cov(X_v,Y_v) = Cv(X_v,Y_v) ev(X_v,Y_v)$$

$$freq Cov(X_v,Y_v) = E \left[(X-\mu_V)(X-\mu_V)^T \right]$$

$$freq Cov(X_v,X_v) = E \left[(X-\mu_V)(X-\mu_V)^T \right]$$

Correlation

$$\mathbb{X}_1$$
 end \mathbb{X}_2 ; $\mathbb{C}_{\text{err}}(\mathbb{X}_1, \mathbb{X}_2) = \frac{\mathbb{C}_{\text{err}}(\mathbb{X}_1, \mathbb{X}_2)}{\sqrt{\text{Ver}(\mathbb{X}_1) \cdot \text{Ver}(\mathbb{X}_2)}}$

Rec Ext. Pli do this in R See metrix formule on slider

Ex: Given I find 9:

Q: con information i Z be condensed?

tr
$$(Z) = \sum_{i=1}^{p} Var(X_i) =$$
" total varience" trace = sum of diagonal elements

det (E) = "generalized vanence"

TMA4267 Linear Statistical Models V2017 [L3]

Part 1: Multivariate RVs and normal distribution (L3) Covariance and positive definiteness [H:2.2,2.3,3.3],
Principal components [H11.1-11.3]
Quiz with Kahoot!

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 17, 2017

Previously

- **X**_{$(p\times 1)$}: random vector, described by
- ▶ joint probability distribution function (pdf) f(x) and cumulative distribution function (cdf) F(x), or, as we will see (in L4), by the (multivariate) moment generating function (MGF) $M_X(t) = e^{t^T X}$.
- ▶ Important aspects: moments.
- ► E(X): mean of a random vector (or matrix) is found as the mean of each element.
- ► $Cov(\mathbf{X}) = E((\mathbf{X} \mu)(\mathbf{X} \mu)^T)$: $p \times p$ variance-covariance matrix, with variances on the diagonal and covariance off-diagonal, real, symmetric.
- ▶ Rules for vector of linear combinations CX: $E(CX) = C\mu$ and $Cov(CX) = C\Sigma C^T$.

Today!

- Requirements and properties of Σ = Cov(X): symmetric, positive definite (SPD), via spectral decomposition (eigenvalues/eigenvectors).
- ► The square root matrix.
- ► Linear combinations with maximal variability: principal components are linear combinations made from eigenvectors.
- PCA-plots.
- ► Kahoot! on what we have worked with so far.
- ▶ Next lecture: move on to multivariate normal data!

2/26

4 / 26

The variance-covariance matrix and positive definiteness

 $m{X}_{(p \times 1)}$ with symmetric $m{\Sigma} = \mathrm{Cov}(m{X})$

- ▶ Want variance of linear combination to be positive: $c^T \Sigma c$ for all $c \neq 0$.
- which means that Σ needs to positive definite. Write $\Sigma > 0$.
- This is true if all eigenvalues of Σ are positive (eigenvalues of symmetric matrix are real).
- Spectral decompositions (diagonalization): $Σ = PΛP^T$.
- Square root matrix defined from spectral decomposition: $\Sigma^{\frac{1}{2}} = P \Lambda^{\frac{1}{2}} P^{T}.$

Drinking habits data set

	Coffee	Too	Cocos	Liquer	Wine	Beer
••						
Norway	9.800000				6.4	
Danmark	10.400001	0.39	0.54	1.4	20.7	123.2
Finland	12.450000	0.17	0.03	3.1	5.4	79.0
Iceland	8.270001	0.23	0.00	2.2	5.2	23.7
Sweden	10.710000	0.32	0.16	1.8	12.3	57.4
France	5.490000	0.20	1.18	2.5	73.8	40.5
Ireland	0.550000	3.14	2.76	1.4	3.9	114.0
Italy	4.670000	0.08	0.97	1.0	67.0	22.7
Jugoslavia	3.100000	0.11	0.59	1.6	20.3	46.8
The Netherlands	10.970000	0.82	15.35	2.0	14.7	87.0
Poland	1.400000	0.54	0.56	4.3	7.6	30.8
Portugal	3.080000	0.03	0.02	0.8	51.5	60.7
Soviet Union	0.300000	0.98	0.56	1.9	6.6	19.3
Spain	4.250000	0.03	1.11	2.8	38.3	70.7
Schweitz	9.400000	0.25	3.06	1.9	49.6	69.3
Great Britain	2.060000	2.62	2.78	1.8	11.5	110.5
Chech Repl	2.200000	0.13	1.21	3.3	13.6	132.9
Germany	8.970000	0.22	3.94	2.0	26.0	143.0
Hungary	2.270000	0.07	0.85	4.6	21.5	103.9
Austria	10.220000	0.16	1.80	1.5	34.8	119.5
New Zealand	1.920000	1.46	0.03	1.4	14.6	114.2

3 / 26

5 / 26

Drinking habits

	Coffee	Tea	Cocoa	Liquer	Wine	Beer
Norway	9.800000	0.21	0.61	1.1	6.4	52.0
Danmark	10.400001	0.39	0.54	1.4	20.7	123.2
Finland	12.450000	0.17	0.03	3.1	5.4	79.0
Iceland	8.270001	0.23	0.00	2.2	5.2	23.7
Sweden	10.710000	0.32	0.16	1.8	12.3	57.4
France	5.490000	0.20	1.18	2.5	73.8	40.5
Ireland	0.550000	3.14	2.76	1.4	3.9	114.0
Italy	4.670000	0.08	0.97	1.0	67.0	22.7
Jugoslavia	3.100000	0.11	0.59	1.6	20.3	46.8
The Netherlands	10.970000	0.82	15.35	2.0	14.7	87.0
Poland	1.400000	0.54	0.56	4.3	7.6	30.8
Portugal	3.080000	0.03	0.02	0.8	51.5	60.7
Soviet Union	0.300000	0.98	0.56	1.9	6.6	19.3
Spain	4.250000	0.03	1.11	2.8	38.3	70.7
Schweitz	9.400000	0.25	3.06	1.9	49.6	69.3
Great Britain	2.060000	2.62	2.78	1.8	11.5	110.5
Chech Repl	2.200000	0.13	1.21	3.3	13.6	132.9
Germany	8.970000	0.22	3.94	2.0	26.0	143.0
Hungary	2.270000	0.07	0.85	4.6	21.5	103.9
Austria	10.220000					119.5
New Zealand	1.920000	1.46	0.03	1.4	14.6	114.2

Estimators for μ and Σ [H3.3]

$$\boldsymbol{X}_1, \boldsymbol{X}_2, \dots, \boldsymbol{X}_n$$
 i.i.d $\mathrm{E}(\boldsymbol{X}) = \boldsymbol{\mu}$ and $\mathrm{Cov}(\boldsymbol{X})$.

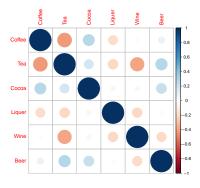
$$ar{m{X}} = rac{1}{n} \sum_{j=1}^n m{X}_j$$
 $m{S}^2 = rac{1}{n-1} \sum_{j=1}^n (m{X}_j - ar{m{X}}) (m{X}_j - ar{m{X}})^T$

are two commonly used estimators for the mean and covariance matrix.

RecEx1.P7: we may write S^2 using a centering matrix.

6/26

Correlation plot



8 / 26

Drinking habits data set

7 / 26

prcomp

```
drink <- read.csv("drikke.TXT",sep=",",header=TRUE)
drink <- na.omit(drink) # remove missing data
pca <- prcomp(drink,scale=TRUE) # scale: variables are scaled, and automatically center=TRUE
[1] "sdev"
             "rotation" "center" "scale"
> pca$rotation # the loadings
                        PC2
                                   PC3
                                                PC4
Coffee -0.26029733 0.66788815 -0.22475187 0.4132467433 0.07431918 0.5092751
      0.65540048 -0.09539757 0.36756357 -0.0002927055 -0.12503940 0.6407898
Cocoa 0 23510209 0 57754726 -0 06603093 -0 4200858712 -0 61199325 -0 2362164
Liquer 0.02190508 -0.32118904 -0.79997824 -0.3292322714 -0.12307455 0.3644878
Wine -0.50599685 0.06551597 0.37109534 -0.6765579799 0.15862233 0.3450672
      > s <- cor(drink) # cor, not cov, since covariates are scaled
> eigen(s) # same as pca$rotations - opposite sign of some vectors
$values
[1] 1.7204307 1.4295795 1.1408597 0.7731249 0.7354586 0.2005467
[1,] 0.26029733 0.66788815 -0.22475187 0.4132467433 0.07431918 -0.5092751
[2,] -0.65540048 -0.09539757 0.36756357 -0.0002927055 -0.12503940 -0.6407898
\hbox{\tt [4,]} \  \  \, \hbox{\tt -0.02190508} \  \  \, \hbox{\tt -0.32118904} \  \  \, \hbox{\tt -0.79997824} \  \  \, \hbox{\tt -0.3292322714} \  \  \, \hbox{\tt -0.12307455} \  \  \, \hbox{\tt -0.3644878}
[5,] 0.50599685 0.06551597 0.37109534 -0.6765579799 0.15862233 -0.3450672
```

Principal components

- Let Σ be the covariance matrix associated with the random vector $\boldsymbol{X}_{p\times 1}$. The covariance matrix has the eigenvalue-vector pairs $(\lambda_j, \boldsymbol{e}_j)$, where $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0$.
- ▶ The *m*th principal component is given by

$$Y_m = e_m^T X = e_{m1} X_1 + e_{m2} X_2 + \cdots + e_{mp} X_p$$

and has

$$Var(Y_m) = \mathbf{e}_m^T \mathbf{\Sigma} \mathbf{e}_m = \lambda_m, \quad i = 1, 2, ..., p$$
$$Cov(Y_i, Y_m) = \mathbf{e}_i^T \mathbf{\Sigma} \mathbf{e}_m = 0 \quad i \neq m$$

10 / 26

Principal components: idea

3. We choose principal component 3, $PC_1 = c_3^T X$, to have maximal variance and be uncorrelated with PC_1 and PC_2 .

$$\max_{m{c}_3
eq 0, m{c}_3^{ op} m{c}_3 = 1} ext{Var}(m{c}_3^{ op} m{X}) ext{ and } m{c}_i^{ op} m{\Sigma} m{c}_3 = 0$$

for i = 1, 2.

4. and so on.

It can be shown that choosing $c_i = e_i$ (ith eigenvector of Σ) fulfills these requirements.

Principal components: idea

1. We choose principal component 1, $PC_1 = c_1^T X$, to have maximal variance

$$\max_{\boldsymbol{c}_1 \neq 0, \boldsymbol{c}_1^T \boldsymbol{c}_1 = 1} \mathrm{Var}(\boldsymbol{c}_1^T \boldsymbol{X})$$

2. We choose principal component 2, $PC_2 = c_2^T X$, to have maximal variance and to be uncorrelated with PC_1 .

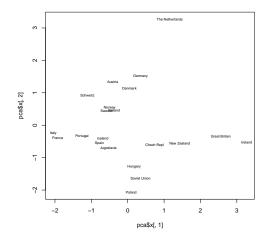
$$\max_{\boldsymbol{c}_2 \neq \boldsymbol{0}, \boldsymbol{c}_2^{\mathcal{T}} \boldsymbol{c}_2 = 1} \mathrm{Var}(\boldsymbol{c}_2^{\mathcal{T}} \boldsymbol{X}) \text{ and } \boldsymbol{c}_1^{\mathcal{T}} \boldsymbol{\Sigma} \boldsymbol{c}_2 = 0$$

11 / 26

Principal component scores

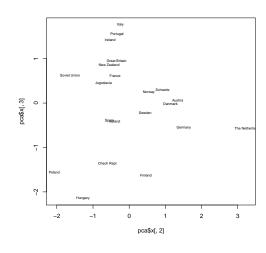
```
> pca$x
               -0.49755266 0.5423321 0.24730564 1.64969605 -0.141194597 -0.259224773
Norway
Danmark
                0.04736942 1.1407629 -0.01389387 0.62447776 1.286968641 0.032325551
Finland
               -0.68875794 -0.4028320 -0.40970218 1.46441365 -0.741785804
Iceland
                                                                        0.092078151
               -0.57111547 0.4409190 -0.21949104 1.33928137 0.005514957
Sweden
                                                                        0.299908368
France
               -1.92532968 -0.3909512 0.61668510 -1.42238114 -0.200223837
                3.27114811 -0.5256052 1.42206819 -0.04493512 -0.014828505
Ireland
Italy
               -2.04540738 -0.2394487 1.76095443 -0.63624103 -0.367324127 0.054512802
Jugoslavia
               -0.53258860 -0.7075896  0.45055250  0.39569929 -0.299213104 -0.741778021
The Netherlands 1.15707155 3.2747502 -0.57171870 -0.93460573 -2.324855354 -0.377956367
Poland
                0.09667998 \ -2.0690624 \ -1.55793948 \ -0.11990952 \ -1.118921154 \ \ 0.196370232
Portugal
               -1.25572594 -0.3355743 1.55429933 -0.39165369 0.418471226 -0.587655889
Soviet Union
               0.35234451 -1.6341096 0.62745644 0.66818676 -1.138865302 -0.446485485
               -0.77751355 -0.5475718 -0.39572230 -0.70847071 0.084179130 -0.055157063
Spain
               -1.10755090 0.9248648 0.29621598 -0.49777226 -0.047226891 0.484657020
Schweitz
Great Britain
               2.55842016 -0.3473694 0.94755863 -0.24512125 0.028322241 0.555177625
Chech Repl
                0.74010278 -0.6054459 -1.36137719 -0.75417409 0.950808503 -0.708681733
Germany
                0.35362098 1.5089269 -0.54743939 -0.46975244 0.999859480 -0.289909548
                0.17766893 -1.2757040 -2.13925810 -1.16454698 0.376532314 -0.009844322
Hungary
               -0.41110347 1.3402182 0.06024028 -0.01769055 1.108477110 0.029376612
New Zealand
               1.43034542 -0.5543289 0.86167320 0.09326849 0.850499977 -0.250026615
```

Scores PCA1 vs PCA2



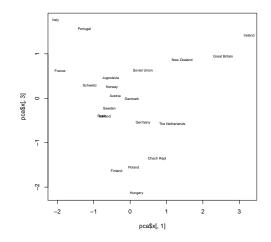
14 / 26

Scores PCA2 vs PCA3



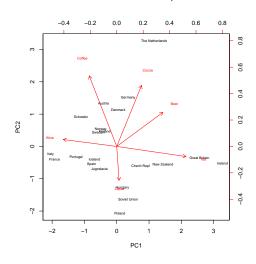
16 / 26

Scores PCA1 vs PCA3

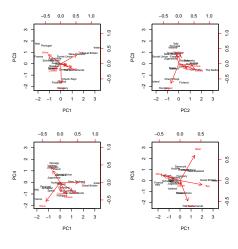


15 / 26

Scores PCA1 vs PCA2 with biplot



Biplots



18 / 26

How many PCs are needed?

Dependent on:

- ► The proportion of the total sample variance that we would like to explain. 80%? More?
- ► Look at the eigenvalues; small eigenvalues may be an evidence of collinearity problems.

Proportion of total population variance

► Total population variance:

$$\sum_{j=1}^{p} \operatorname{Var}(X_{j}) = \operatorname{tr} \mathbf{\Sigma} = \sum_{j=1}^{p} \lambda_{j} = \sum_{j=1}^{p} \operatorname{Var}(Z_{j}).$$

▶ Proportion of total population variance explained by PC *m*:

$$\frac{\lambda_m}{\sum_{j=1}^p \lambda_j}$$

▶ Proportion of total population variance explained by the first m PCs:

$$\frac{\sum_{j=1}^{m} \lambda_j}{\sum_{j=1}^{p} \lambda_j}$$

19 / 26

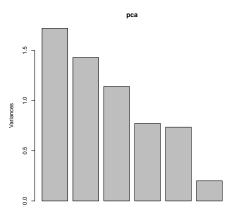
Importance of components

```
> summary(pca)
Importance of components:
```

PC1 PC2 PC3 PC4 PC5 1.3117 1.1957 1.0681 0.8793 0.8576 0.44782 Standard deviation Proportion of Variance 0.2867 0.2383 0.1901 0.1288 0.1226 0.03342 Cumulative Proportion 0.2867 0.5250 0.7151 0.8440 0.9666 1.00000

\$values [1] 1.7204307 1.4295795 1.1408597 0.7731249 0.7354586 0.2005467 > sqrt(eigen(s)\$values)
[1] 1.3116519 1.1956502 1.0681103 0.8792752 0.8575888 0.4478244





22 / 26

Principal components from singular value decomposition

The singular value decomposistion of a (data) matrix $\boldsymbol{X}_{n \times p}$ is given by:

$$\boldsymbol{X}_{n \times p} = \boldsymbol{U}_{n \times p} \boldsymbol{D}_{p \times p} \boldsymbol{V}_{p \times p}^{T}$$

where

- the columns of \boldsymbol{U} are the eigenvectors of $\boldsymbol{X}\boldsymbol{X}^T$
- ▶ D is a diagonal matrix with singular values on the diagonal, i.e. the square root of the eigenvalues of XX^T and X^TX (they have the same eigenvalues).
- the columns of V are the eigenvectors of $X^T X$.

And, the principal components (scores) of the data are defined as the columns of

$$Z = XV = UD$$

24 / 26

PC from standardized variables

X can be standardized to have mean **0** and unit variances.

$$oldsymbol{\mathcal{X}}^* = oldsymbol{V}^{-rac{1}{2}}(oldsymbol{\mathcal{X}} - oldsymbol{\mu})$$

- Principal components made from standardized variables will be based on the eigenvalues and eigenvectors of the correlation matrix $\rho = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}}$.
- Achilles heel: Since Σ and ρ do not have the same eigenvectors/eigenvalues, the principal components made from Σ and ρ will not be the same.
- ► Unless we have a good reason to compare the variances for the different *X_j*s we should make PCs from the standardized variables.
- For standardized variables $\sum_{i=1}^{p} \operatorname{Var}(X_{i}^{*}) = p$, and
- Proportion of total population variance explained by PC m: $\frac{\lambda_m}{p}$.

23 / 26

PCR: summary

- ▶ PCA finds linear combinations **Y** that "best" represents the **X**.
- ► The PCs are found in an unsupervised way. The "truth" is not known
- ► A plot of PC1 vs PC2 is often used to see if there is separation (subgroups in the data).
- ► The principal component loadings are often given interpretation (overall consumption,
- ▶ PCA can be combined with linear regression.

Quiz

with Kahoot! at kahoot.it. - based on what we have gone through so far!

26 / 26

So, what do we need?

this is the delimitar of a S P D metrix Symmetric portue definite

Useful result:

Proof via spectral decomposition:

Let PAPT= = wher P= [e, ez -- ep] is a matrix with the normalized ergenvectors of E as column rectors end 1 = diag (hi, ..., hp) is a diagonal metrix with the eigenvalues of E on the diagonal.

Remember: (hi, ei) soliste Zei = lie: and use $\lambda_1 > \lambda_2 > ... > \lambda_p$. And that a (real) symmetric metric has real enganvalues (end eigenvectors of distinct eigenvalues are orthogonal)

THAYSET LS The covanance motors [H22, H: 23] 17.01. 2017

E(X)- pm X rondom vector COV(X)= Z = E[(X-M(X-M+)

$$\begin{bmatrix} \operatorname{Vor}(X_1, X_2) & \cdots & \operatorname{Vor}(X_n, X_1) \\ \operatorname{Cor}(X_1, X_2) & \operatorname{Vor}(X_n) & \cdots \\ \end{array}$$

EX : Drinking hebits X1 = coffe, X2 = tca, x3 = cocor, Xy= higher, IT= Wine, Xx= beer 1.1.d n=21 countrie.

I is real and symmetric. Other requirements? We will need I's o we require that I is invertible, that, det(E) +0 (or we might use generalized

But, we boiled at CX and found GV(CX)=CZCT. If ked then CTX and Cov(CTX) = CT IC

ty pol

Int

Scales

We went Var(cTX) = cT Ec to be positive (because we don't want O ar negative vanences)

CTEC > 0 for all C+O CT PAPT C > 0 Pis investigation yt / y > 0 for all y +0 since
ARP PRO PTC = y

This is the when all

If $y=\begin{pmatrix} 1 \\ 6 \end{pmatrix} \rightarrow A_1$, some for all the eigenvalues

this also implies

Rectal.P4 aux (E)= #x:>0 4(E) = Ex:>0

Homework: Whip Is det (5) = TTA; and Tr(E) = Zh ?

What about E-1

If I is SPD than I'l is also SPD, and the eigenvalues of Z-1 se the invesce of the eigenvalues of E.

PROOF: $Z^{-1} = (PP)^{-1} = (PT)^{-1} N^{1}P^{-1}$ $(PT)^{-1} = P \quad \text{ond} \quad P^{-1} = P^{T} \text{ for attheoperal metric}$ $Z^{-1} = P N^{-1}P^{T} = P \begin{bmatrix} \frac{1}{2}N & \frac{1}{2} & 0 \\ \frac{1}{2}N & \frac{1}{2} & 0 \end{bmatrix} P^{T}$ exgreedors $exgreedors \quad exgradum of Z^{-1}$

Finally:
$$\mathbb{Z}^{\frac{1}{2}} = \bigcap_{i \in \mathbb{Z}} \Lambda^{\frac{1}{2}} \bigcap_{i \in \mathbb{Z}_p} \Lambda^{\frac{1}{2}}$$

RecExt. Pt: $Z^{\frac{1}{2}}Z^{\frac{1}{2}} = Z$ $Z^{\frac{1}{2}} = (Z^{\frac{1}{2}})^{-1} = P \bigwedge^{-\frac{1}{2}} P^{+}$ $Z^{\frac{1}{2}}Z^{-\frac{1}{2}} = Z \text{ (sum)} \qquad \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{pmatrix}$ The transformation $Y = Z^{-\frac{1}{2}}(X - \mu)$ Notheropis here from

Principal Component Analysis [HILL-11.3]

why is I important and what can CX be?

$$\frac{S^{2}}{\sum_{i=1}^{n}} = \frac{1}{\sum_{j=1}^{n}} \left(\sum_{j=1}^{n} \bar{\Sigma} \right) \left(\sum_{j=1}^{n} \bar{\Sigma} \right)^{T} \quad \text{prop matrix}$$
put $1 \times p$

We will work with scaled venzblo so Z= g.

Azzavalg: 6 Tiddlibom 9 Sut

Elutenfri=

Vegeler: 2+3

Q: When faced with a large sut of correlated vanebles - is it possible to define a set of linear Comfanations of the arroyand ranables that capture a large part of the vanability of the data?

- interpretability of Ains!
- A: Yes. Lot (\(\lambda_i, ei\) be engenvalue/vector

construct Yi = QT X

principal loadings or rotalions

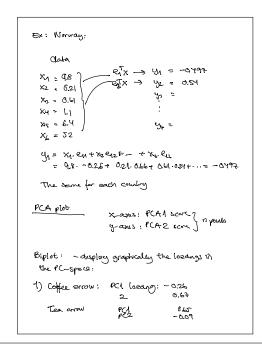
Components

Important property:

X, $Cr(X) = \Sigma$ and (A;e) eigenvalue (rector pairs of Σ , $1 = l_{r_1}p$)

Let $Y_i = e_i^T X$ and $Y = \begin{cases} Y_i \\ Y_i \end{cases} = p^T X$

> sum up next line ?



f(x)

Let X_1 and X_2 be two (continuous) RVs, and $f(x_1, x_2)$ be the joint pdf and $f_1(x_1)$ and $f_2(x_2)$ be the marginal pdfs, and C is a copula. What is true?

A
$$f(x_1, x_2) = f_1(x_1) \cdot f_2(x_2)$$

- **B** $f(x_1, x_2)$ is found from $f_1(x_1)$ and $f_2(x_2)$ alone
- C $f_1(x_1)$ is found from $f(x_1, x_2)$

$$D f(x_1, x_2) = C(f_1(x_1), f_2(x_2))$$

Mean of sum

 \boldsymbol{X} and \boldsymbol{Y} are two bivariate random vectors with $E(\boldsymbol{X}) = (1,2)^T$ and $E(\boldsymbol{Y}) = (2,0)^T$. What is $E(\boldsymbol{X}+\boldsymbol{Y})$?

- **A** $(1.5, 1)^T$
- **B** $(3,2)^T$
- $(-1,2)^T$
- \mathbf{D} $(1,-2)^T$

Mean of linear combination

 ${\pmb X}$ is a 2-dimensional random vector with ${\rm E}({\pmb X})=(2,5)^T$, and ${\pmb b}=(0.5,0.5)^T$ is a constant vector. What is ${\rm E}({\pmb b}^T{\pmb X})$?

- **A** 3.5
- **B** 7
- **D** 5

Covariance

 \boldsymbol{X} is a p-dimensional random vector with mean μ . Which of the following defines the covariance matrix?

- $\mathbf{A} \quad E[(\mathbf{X} \mathbf{\mu})^T (\mathbf{X} \mathbf{\mu})]$
- $\mathsf{B} \quad E[(\mathbf{X} \boldsymbol{\mu})(\mathbf{X} \boldsymbol{\mu})^T]$
- $E[(\mathbf{X} \mathbf{\mu})(\mathbf{X} \mathbf{\mu})]$
- $\mathbf{D} \quad E[(\mathbf{X} \mathbf{\mu})^T (\mathbf{X} \mathbf{\mu})^T]$

Mean of linear combinations

 \boldsymbol{X} is a p-dimensional random vector with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. \boldsymbol{C} is a constant matrix. What is then the mean of the k-dimensional random vector $\boldsymbol{Y} = \boldsymbol{C}\boldsymbol{X}$?

- A Cu
- \mathbf{B} \mathbf{C}
- $C C \mu C^T$
- $\mathbf{D} \quad \mathbf{C} \mathbf{\Sigma} \mathbf{C}^{\mathsf{T}}$

Covariance of linear combinations

 \boldsymbol{X} is a p-dimensional random vector with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. \boldsymbol{C} is a constant matrix. What is then the covariance of the k-dimensional random vector $\boldsymbol{Y} = \boldsymbol{C}\boldsymbol{X}$?

- **Α** *C* μ
- B $C\Sigma$
- $C C \mu C^T$
- $D C \Sigma C^T$

Symmetric positive definite matrix

Which of the following is not correct for a symmetric positive definite matrix?

- A The trace equals the rank of the matrix.
- B The determinant is positive.
- C The trace is the sum of the eigenvalues.
- D All the eigenvalues are positive.

Correlation

X is a 2-dimensional random vector with covariance matrix

$$\mathbf{\Sigma} = \left[\begin{array}{cc} 4 & 0.8 \\ 0.8 & 1 \end{array} \right]$$

Then the correlation between the two elements of \boldsymbol{X} are:

- **A** 0.10
- **B** 0.25
- C 0.40
- **D** 0.80

PCA interpretation

Data set: student's score on a Math test, a Physics test, a Reading comprehension test, and a Vocabulary test.

First PC represents overall academic ability, second PC represents a contrast between quantitative ability and verbal ability.

What loadings would be consistent with that interpretation?

- A (0.5,0.5,0.5,0.5) and (0.71,0.71,0,0)
- B (0.5,0.5,0.5,0.5) and (0.5,0.5,-0.5,-0.5)
- C (0.71,0.71,0,0) and (0,0,0.71,-0.71)
- D (0.71,0,-0.71,0) and (0,0.71,0,-0.71)

Correct?

Are you sure you want to read the correct answers? Maybe try first? The answers are explained on the next two slides.

Answers

- 5. A: $C\mu$ is the mean of Y = CX.
- 6. D: $\mathbf{C} \mathbf{\Sigma} \mathbf{C}^T$ is the covariance matrix of $\mathbf{Y} = \mathbf{C} \mathbf{X}$.
- 7. C: Correlation is 0.40 since covariance was 0.8 and variances 4 and 1.
- 8. A: NOT true for a symmetric positive definite matrix: the trace is in general not equal to the rank but it is for idempotent symmetric matrices.
- 9. B: average means equal weight for all values, difference between quantitative and verbal means opposite signs for quantitative (maths and physics) and verbal (reading and vocabular).

Answers

- 1. C: We go from joint to marginal distribution by integration. The product of marginals equal the joint only for independent variables. We need information on the dependency structure to construct a joint from marginals, and that is what is done with the copula but the formula is based on the cumulative distribution functions.
- 2. B: Mean of sum $(1,2)^T + (2,0)^T = (3,2)^T$.
- 3. A: Mean of linear combination $(0.5, 0.5)^T(2, 5) = 3.5$.
- 4. B: Covariance matrix defined as $E\{(\boldsymbol{X} \boldsymbol{\mu})(\boldsymbol{X} \boldsymbol{\mu})^T\}$. This was the only formula that gave a $p \times p$ matrix. A gave a scalar and C and D did not match in dimensions.

TMA4267 Linear Statistical Models V2017 [L4]

Part 1: Multivariate RVs, and the multivariate normal distribution The multivariate normal distribution (pdf and mgf) [H:4.2-4.4]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 20, 2017

What we know, and the plan for this lecture

1/8

What we know, and the plan for this lecture

- A random vector X can be described by the joint pdf f(x).
- ▶ Mean: $\mu = E(X) = \{E(X_j)\}$

What we know, and the plan for this lecture

▶ A random vector X can be described by the joint pdf f(x).

1/8

What we know, and the plan for this lecture

- A random vector \boldsymbol{X} can be described by the joint pdf $f(\boldsymbol{x})$.
- ▶ Mean: $\mu = E(X) = \{E(X_j)\}$
- ▶ Covariance matrix: $Cov(\mathbf{X}) = E((\mathbf{X} \boldsymbol{\mu})(\mathbf{X} \boldsymbol{\mu})^T)$, symmetric and we often require the matrix to be positive definite.

1/8

What we know, and the plan for this lecture

- A random vector \boldsymbol{X} can be described by the joint pdf $f(\boldsymbol{x})$.
- ▶ Mean: $\mu = E(X) = \{E(X_j)\}$
- ▶ Covariance matrix: $Cov(\mathbf{X}) = E((\mathbf{X} \boldsymbol{\mu})(\mathbf{X} \boldsymbol{\mu})^T)$, symmetric and we often require the matrix to be positive definite.
- Linear combinations CX: $E(CX) = C\mu_X$ and $Cov(CX) = C\Sigma C^T$.

1/8

Why is the mulitivariate normal distribution so important in statistics?

- ▶ Many natural phenomena may be modelled using this distribution (just as in the univariate case).
- Multivariate version of the central limit theorem- the sample mean will be approximately multivariate normal for large samples.
- ▶ Good interpretability of the covariance.
- Mathematically tractable.
- ▶ Building block in many models and methods.

What we know, and the plan for this lecture

- \blacktriangleright A random vector **X** can be described by the joint pdf f(x).
- ▶ Mean: $\mu = E(X) = \{E(X_i)\}$
- Covariance matrix: $Cov(\mathbf{X}) = E((\mathbf{X} \boldsymbol{\mu})(\mathbf{X} \boldsymbol{\mu})^T)$, symmetric and we often require the matrix to be positive definite.
- Linear combinations CX: $E(CX) = C\mu_X$ and $Cov(CX) = C\Sigma C^T$.
- ► Now: derive the joint pdf and the moment generating function for the multivariate normal distribution.

1/8

Cramer-Wold and moment generating functions

 $m{X}_{(p \times 1)}$ is a random vector. The distribution of $m{X}$ is completely determined by the set of all one-dimensional distributions of the linear combinations $Y = m{t}^T m{X} = \sum_{i=1}^p t_i X_i$ where $m{t}$ ranges over all fixed p-vectors.

- $Y = t^T X$ has MGF $M_Y(s) = \mathbb{E}(\exp(sY)) = \mathbb{E}(\exp(st^T X))$.
- ▶ If we choose s = 1 $M_Y(1) = E(\exp(t^T X)) = M_X(t)$, which is the MGF of X and thus determines the distribution of X.

Härdle and Simes (2015) use characteristic functions, $E(e^{it^TX})$ but we stick with moment generating functions $E(e^{t^TX})$. Why: we will only work with nice distributions and do not have problems with integrals not existing, and we know MGFs from previous course.

Multivariate transformation formula [H:4.3]

$$X = u(Y) \tag{4.43}$$

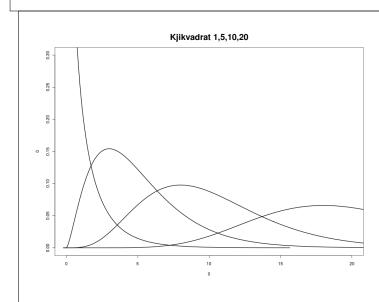
for a one-to-one transformation $u: \mathbb{R}^p \to \mathbb{R}^p$. Define the Jacobian of u as

$$\mathcal{J} = \left(\frac{\partial x_i}{\partial y_j}\right) = \left(\frac{\partial u_i(y)}{\partial y_j}\right)$$

and let $abs(|\mathcal{J}|)$ be the absolute value of the determinant of this Jacobian. The pdf of Y is given by

$$f_Y(y) = \operatorname{abs}(|\mathcal{J}|) \cdot f_X\{u(y)\}. \tag{4.44}$$

4/8



6/8

The Chi-square distribution

pdf χ_p^2 :

$$f(y) = \frac{1}{2^{p/2}\Gamma(p/2)} y^{p/2-1} e^{(-y/2)} \text{ for } y > 0$$

MGF χ_p^2 :

$$M_Y(t) = \frac{1}{(1-2t)^{p/2}}$$

Addition property:

Let $X_1\sim \chi_p^2$ and $X_2\sim \chi_q^2$, and let X_1 and X_2 be independent. Then $X_1+X_2\sim \chi_{p+q}^2$.

Subtraction property:

Let $X=X_1+X_2$ with $X_1\sim\chi_p^2$ and $X\sim\chi_{p+q}^2$. Assume that X_1 and X_2 are independent. Then $X_2\sim\chi_q^2$.

5/8

This lecture: derived the MGF and pdf of the multivariate normal distribution

- 1. $Z \sim N_1(0,1)$
 - MGF: $M_Z(t) = E(e^{tz}) = e^{\frac{1}{2}t^2}$
- 2. Z_1, Z_2, \ldots, Z_p iid $N_1(0,1) \rightarrow \boldsymbol{Z}_{p \times 1} \sim N_p(\boldsymbol{0}, \boldsymbol{I})$
 - MGF: $M_{\mathbf{Z}}(\mathbf{t}) = E(e^{\mathbf{t}^T \mathbf{z}}) = e^{\frac{1}{2}\mathbf{t}^T \mathbf{t}}$
- 3. $m{X} = m{A} m{Z} + m{\mu}, \ m{A} m{A}^T = m{\Sigma} \ ext{gives} \ m{X}_{p imes 1} \sim N_p(m{\mu}, m{\Sigma})$
 - $\blacktriangleright \mathsf{MGF} \colon M_{\boldsymbol{X}}(\boldsymbol{t}) = \mathrm{E}(\boldsymbol{e}^{\boldsymbol{t}^T\boldsymbol{x}}) = \boldsymbol{e}^{\boldsymbol{t}^T\boldsymbol{\mu} + \frac{1}{2}\boldsymbol{t}^T\boldsymbol{t}}$
 - pdf (invertible):

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{p}{2}} |\mathbf{\Sigma}|^{\frac{1}{2}}} \exp\{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\}$$

Properties of the mvN - plan for L5

Let $X_{(p\times 1)}$ be a random vector from $N_p(\mu, \Sigma)$.

- 1. Probability density function f(x) (both when Σ is invertible and not).
- 2. Moment generating function: $M_X(t) = \exp(t^T \mu + \frac{1}{2} t^T \Sigma t)$
- 3. Graphical display, contours (ellipsoids), and chisq-distributed $(\mathbf{X} \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{X} \boldsymbol{\mu})$.
- Linear combinations of components of X are (multivariate) normal.
- 5. All subsets of the components of \boldsymbol{X} are (multivariate) normal.
- 6. Zero covariance implies that the corresponding components are independently distributed.
- 7. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent.
- 8. The conditional distributions of the components are (multivariate) normal. $\boldsymbol{X}_2 \mid (\boldsymbol{X}_1 = \boldsymbol{x}_1) \sim N_{p2}(\mu_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\boldsymbol{x}_1 \mu_1), \boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})$

And then remains estimators for parameters and properties of quadratic forms in L6.

8/8

Characterizing
$$X$$
:

- far (pdf), Far (cdf)

- $M_X(t) = E(e^{t^TX}) = E(e^{t_tX_1t_tX_2t-\tau t_pX_p})$

And $P^{XI} = \int_{-\infty}^{\infty} e^{t^X} far dx_1...dx_p$

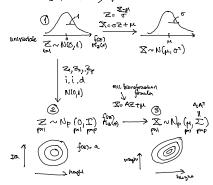
MAF: easy for proofs.

2

Multivanele normal distribution

THAY267 LY 20,01, 2017

PIEN: derive the pdf and moment generaling function (MEN $M_E(E)$: $E(e^{tx})$) of the multivariete normal.



1

$$E(8) = 0, \text{ Vol}(2) = 1$$

$$f(8) = 0, \text{ Vol}(2) = 1$$

$$f(8) = 0, \text{ for } 2 = 0$$

$$f(8) = 0, \text{ for } 2 = 0$$

$$f(8) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 2 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for } 3 = 0$$

$$f(9) = 0, \text{ for }$$

$$(1) \rightarrow (2)$$

$$Z_{1}, Z_{2}, ..., Z_{p} \text{ independent N(s,t) PUS}$$

$$f(z) = \prod_{i=1}^{p} \frac{1}{c_{ir}} e^{-\frac{i}{2}Z_{i}^{2}} = (\frac{1}{2\pi})^{\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{i=1}^{p}Z_{i}^{2}\right]$$

$$= (\frac{1}{2\pi})^{\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{i=1}^{p}Z_{i}^{2}\right]$$

$$= (\frac{1}{2\pi})^{\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{i=1}^{p}Z_{i}^{2}\right]$$

$$E(2) = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$
 and $C_{m}(2) = I_{php}$
 $C_{m}(2c, 3) = 0$ if

MGF: mulhrenzle version

$$H_{z}(t) = E(exp(t^{r}z))$$

$$t_{z_{1}}t_{z_{2}}t_{z_{3}}t_{1}t_{2}$$

$$E(X) = AE(2) + E(\mu) = A \cdot O + \mu = \mu$$

$$Con(X) = A (n)(2) A^{T} = AA^{T} = 2$$

S = TAA = TA (S) A = (R) since A has full rente AAT 18 possitive definite (can be proven wang) $(Ax)^{\tau}(Ax)$

What can A be? Z1 but also other possibilities.

6

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{t_1 e_1} \cdot e^{t_2 e_2} \cdot e^{t_2 e_3} \cdot f(e_3) \cdot f(e_3) \cdot f(e_4) de_4 \cdot de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_2} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_4 \cdot \dots \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty} e^{t_1 e_3} f(e_3) de_9 \cdot \dots de_9$$

$$= \int_{-\infty}^{\infty}$$

5

3>0 Frut Har of X

$$H_{X}(t) = H_{A2+\mu}(t) = E\left(e^{t^{T}(A2+\mu)}\right)$$

$$= E\left(e^{t^{T}A2} \cdot e^{t^{T}\mu}\right) = e^{t^{T}\mu} E\left(e^{t^{T}A2}\right)$$

2-3 fox), and the mr. transformation fermula

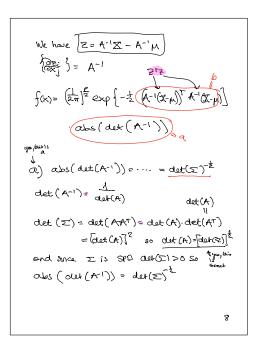
The mr. bronsf. formula

$$X = Az + \mu \iff z = A^{-1}(x - \mu)$$

$$f(x) = f_z(z(x)) \cdot abs(f)$$

$$f = det(f = x_i)_{(i,j)} \cdot denut((i,j))$$

$$f = another with that as for another with that as for another x with that as for an another x with that a second x with the x with that a second x with the x with th$$



Comment: when A does not have full renu We can use a singular veryon of the pdf.

10

(NON)
$$\Sigma A^{-1}(X^{-}\mu) = \frac{1}{2} \left(X^{-}\mu \right) = \frac{1}{2} \left(X^{-}\mu$$

TMA4267 Linear Statistical Models V2017 [L5]

Part 1: Multivariate RVs, and the multivariate normal distribution Properties of the multivariate normal distribution [H:2.6,4.4,5.1]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 24, 2017

Last lecture: derived the MGF and pdf of the multivariate normal distribution

- 1. $Z \sim N_1(0,1)$ • MGF: $M_Z(t) = E(e^{tz}) = e^{\frac{1}{2}t^2}$
- 2. Z_1, Z_2, \ldots, Z_p iid $N_1(0,1) \rightarrow \boldsymbol{Z}_{p \times 1} \sim N_p(\boldsymbol{0}, \boldsymbol{I})$
 - MGF: $M_{Z}(t) = E(e^{t^{T}z}) = e^{\frac{1}{2}t^{T}t}$
- 3. $\boldsymbol{X} = \boldsymbol{A}\boldsymbol{Z} + \boldsymbol{\mu}$, $\boldsymbol{A}\boldsymbol{A}^T = \boldsymbol{\Sigma}$ gives $\boldsymbol{X}_{p \times 1} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
 - $\blacktriangleright \mathsf{MGF} \colon M_{X}(t) = \mathrm{E}(e^{t^{T}x}) = e^{t^{T}\mu + \frac{1}{2}t^{T}t}$
 - ▶ pdf (∑ invertible):

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{\rho}{2}} |\mathbf{\Sigma}|^{\frac{1}{2}}} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\}$$

1/20

Today: six properties of the mvN

Let $\boldsymbol{X}_{(p\times 1)}$ be a random vector from $N_p(\mu, \boldsymbol{\Sigma})$.

1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).

Why is the mulitivariate normal distribution so important in statistics?

- ► Many natural phenomena may be modelled using this distribution (just as in the univariate case).
- ► Multivariate version of the central limit theorem- the sample mean will be approximately multivariate normal for large samples.
- ► Good interpretability of the covariance.
- ► Mathematically tractable.
- ▶ Building block in many models and methods.

2 / 20

Today: six properties of the mvN

Let $\boldsymbol{X}_{(p\times 1)}$ be a random vector from $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).

3 / 20

Today: six properties of the mvN

Let $X_{(p\times 1)}$ be a random vector from $N_p(\mu, \Sigma)$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- All subsets of the components of X are (multivariate) normal (special case of the above).

3 / 20

Today: six properties of the mvN

Let $X_{(p\times 1)}$ be a random vector from $N_p(\mu, \Sigma)$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- All subsets of the components of X are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent (will be very important in Part 2)

Today: six properties of the mvN

Let $\boldsymbol{X}_{(p\times 1)}$ be a random vector from $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- Linear combinations of components of X are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).

3 / 20

Today: six properties of the mvN

Let $X_{(p\times 1)}$ be a random vector from $N_p(\mu, \Sigma)$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- All subsets of the components of X are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $\mathbf{X}_2 \mid (\mathbf{X}_1 = \mathbf{x}_1) \sim N_{02}(\mu_2 + \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}(\mathbf{x}_1 \mu_1), \mathbf{\Sigma}_{22} \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12}).$

3 / 20

Diabetes data

We will study a data set on diabetes in Part 2. The data set has measurements on n=442 diabetes patients, and p=11 different measurements are taken for each patients. These measurements are:

- age
- sex
- ▶ body mass index (bmi)
- ▶ mean arterial blood pressure (map)
- ➤ six blood serum measurements: total cholesterol (tc), Idl cholesterol (Idl), Idl cholesterol (Idl), tch, Itg, glu.
- ► a quantitative measurement of disease progression one year after baseline (prog)

We will look at the four variables bmi, map, to and Idl. Can we assume that these follow a multivariate normal distribution?

4 / 20

Example: Slightly modified version of Exam K2014 1b

Let $\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ be a bivariate normal random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\boldsymbol{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix $\boldsymbol{\Sigma} = \mathrm{Cov}(\boldsymbol{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$.

You find the eigenvalues and eigenvectors of the covariance matrix ${f \Sigma}$ on the next slide.

Describe the graph of the equation $(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = b$ where b > 0 is a constant.

Make a drawing of the graph, for b = 1 found above.

What is the probability that a random sample from this distribution will be inside this graph?

Contours of multivariate normal distribution

 Contours of constant density for the p-dimensional normal distribution are ellipsoids defined by x such that

$$(\mathbf{x} - \mathbf{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \mathbf{\mu}) = b$$

where b > 0 is a constant.

These ellipsoids are centered at μ and have axes $\pm \sqrt{b\lambda_i} e_i$, where $\Sigma e_i = \lambda_i e_i$, for i = 1, ..., p.

- $\triangleright (\mathbf{x} \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} \boldsymbol{\mu})$ is distributed as χ_p^2
- ightharpoonup The volume inside the ellipsoid of x values satisfying

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq \chi_p^2(\alpha)$$

has probability $1 - \alpha$.

5 / 20

Example: Exam K2014 1b

```
> sigma <- matrix(c(1,0.5,0.5,2),ncol=2)
```

> eigen(sigma)

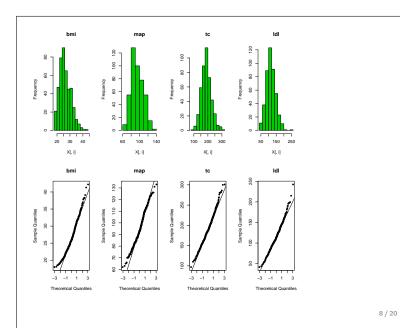
\$values

[1] 2.2071068 0.7928932

\$vectors

[1,] 0.3826834 -0.9238795

[2,] 0.9238795 0.3826834



Multivariate distributions - in 3D: task for the intermission!

Let
$$\mathbf{\Sigma} = \left[egin{array}{ccc} \sigma_{x}^{2} &
ho\sigma_{x}\sigma_{y} \
ho\sigma_{x}\sigma_{y} & \sigma_{y}^{2} \end{array}
ight]$$

The following four 3D-printed figures have been made:

- A: $\sigma_x = 1$, $\sigma_y = 2$, $\rho = 0.3$
- ▶ B: $\sigma_x = 1$, $\sigma_y = 1$, $\rho = 0$
- ► C: $\sigma_x = 1$, $\sigma_y = 1$, $\rho = 0.5$
- ▶ D: $\sigma_x = 1$, $\sigma_y = 2$, $\rho = 0$

The figures have the following colours:

- white
- purple
- ▶ red
- ▶ black

Task: match letter and colour by writing the correct letter after the name of the colour on the available sheets and take the sheet with you. We report on the solution after the intermission.

10 / 20

Today: six properties of the mvN

Let $\boldsymbol{X}_{(p\times 1)}$ be a random vector from $N_p(\mu, \boldsymbol{\Sigma})$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $\boldsymbol{X}_2 \mid (\boldsymbol{X}_1 = \boldsymbol{x}_1) \sim N_{p2}(\mu_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\boldsymbol{x}_1 \mu_1), \boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}).$

11 / 20

Example: Exam K2014 1a

Let
$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$$
 be a bivariate normal random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\mathbf{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix

$$\mathbf{\Sigma} = \operatorname{Cov}(\mathbf{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}.$$

Let
$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$$
, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$.

What is the distribution of \mathbf{Y} ?

What is the distribution of Y_1 ?

Let $Z = X_1 + aX_2$. How can you choose a so that Z and Y_2 are independent?

12 / 20

Example: Exam K2014 1a (slightly modified)

Let
$$m{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$$
 be a bivariate normal random vector with mean $m{\mu} = \mathrm{E}(m{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix $m{\Sigma} = \mathrm{Cov}(m{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$. Let $m{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$.

What is the distribution of Y_1 ?

What is the distribution of \mathbf{Y} ?

Let $Z=X_1+aX_2$. How can you choose a so that Z and Y_2 are independent?

14/20

Example: Exam K2014 1a (slightly modified)

Let
$$\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$$
 be a bivariate normal random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\boldsymbol{X}) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix $\boldsymbol{\Sigma} = \mathrm{Cov}(\boldsymbol{X}) = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$. Let $\boldsymbol{Y} = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$, where $Y_1 = 3X_1 - 2X_2$ and $Y_2 = X_1 + X_2$. What is the distribution of \boldsymbol{Y} ?

What is the distribution of Y_1 ?

Let $Z = X_1 + aX_2$. How can you choose a so that Z and Y_2 are independent?

13 / 20

Today: six properties of the mvN

Let $\boldsymbol{X}_{(p\times 1)}$ be a random vector from $N_p(\mu, \boldsymbol{\Sigma})$

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $\boldsymbol{X}_2 \mid (\boldsymbol{X}_1 = \boldsymbol{x}_1) \sim N_{p2}(\mu_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\boldsymbol{x}_1 \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}).$

Independent variables?

Let $\boldsymbol{X}_{p\times 1} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with

$$\mathbf{\Sigma} = \left[\begin{array}{cccc} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 1 \\ 0 & 0 & 2 & 1 \\ 0 & 1 & 1 & 2 \end{array} \right]$$

List the pairs of variables that are independent.

16 / 20

Today: six properties of the mvN

Let $X_{(p\times 1)}$ be a random vector from $N_p(\mu, \Sigma)$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition).
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF).
- All subsets of the components of X are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF).
- 5. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent (will be very important in Part 2)
- 6. The conditional distributions of the components are (multivariate) normal. $\boldsymbol{X}_2 \mid (\boldsymbol{X}_1 = \boldsymbol{x}_1) \sim N_{o2}(\mu_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\boldsymbol{x}_1 \mu_1), \boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}).$

18 / 20

Example: Exam K2014 1a - cont.

Let ${m X}=\begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ be a bivariate normal random vector with mean ${m \mu}={
m E}({m X})=\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and covariance matrix ${m \Sigma}={
m Cov}({m X})=\begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix}$.

Let
$$m{Y}=\left(egin{array}{c} Y_1 \\ Y_2 \end{array}
ight)$$
, where $Y_1=3X_1-2X_2$ and $Y_2=X_1+X_2.$

Let $Z = X_1 + aX_2$. How can you choose a so that Z and Y_2 are independent?

17 / 20

Example: Exam V2010, Problem 1

Let
$$\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \sim N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
 where $\boldsymbol{\mu} = \begin{pmatrix} 4 \\ -3 \\ 1 \end{pmatrix}$ and $\boldsymbol{\Sigma} = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 1 & -1.5 \\ 0 & -1.5 & 5 \end{pmatrix}$.

a) Find the distribution of $X_1 + X_2 + X_3$ and of X_2 given $X_1 = x_1$ and $X_3 = x_3$.

Help: for
$$\pmb{X} = \begin{pmatrix} \pmb{X}_1 \\ \pmb{X}_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \pmb{\mu}_1 \\ \pmb{\mu}_2 \end{pmatrix}, \begin{pmatrix} \pmb{\Sigma}_{11} & \pmb{\Sigma}_{12} \\ \pmb{\Sigma}_{21} & \pmb{\Sigma}_{22} \end{pmatrix} \right)$$
 we have $\pmb{X}_2 \mid (\pmb{X}_1 = \pmb{x}_1) \sim \mathcal{N}(\pmb{\mu}_2 + \pmb{\Sigma}_{21} \pmb{\Sigma}_{11}^{-1} (\pmb{x}_1 - \pmb{\mu}_1), \pmb{\Sigma}_{22} - \pmb{\Sigma}_{21} \pmb{\Sigma}_{11}^{-1} \pmb{\Sigma}_{12})$

Today: six properties of the mvN

Let $X_{(p\times 1)}$ be a random vector from $N_p(\mu, \Sigma)$.

- 1. The grapical contours of the mvN are ellipsoids (shown using spectral decomposition). [CompEx1.1b]
- 2. Linear combinations of components of **X** are (multivariate) normal (proof using MGF). [CompEx1.1a]
- 3. All subsets of the components of **X** are (multivariate) normal (special case of the above).
- 4. Zero covariance implies that the corresponding components are independently distributed (proof using MGF). [CompEx1.1a]
- 5. $\mathbf{A} \mathbf{\Sigma} \mathbf{B}^T = \mathbf{0} \Leftrightarrow \mathbf{A} \mathbf{X}$ and $\mathbf{B} \mathbf{X}$ are independent (will be very important in Part 2). [CompEx1.2b]
- 6. The conditional distributions of the components are (multivariate) normal. $\boldsymbol{X}_2 \mid (\boldsymbol{X}_1 = \boldsymbol{x}_1) \sim N_{p2}(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\boldsymbol{x}_1 \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}).$

20 / 20

Properties of the 7MAY269 LS 2401.2017.

multiveriese normal distribution (H.26,1445.1)

$$\begin{split} & \underset{\text{pM}}{\mathbb{X}} \sim \mathcal{N}_{p}(\mu, \Sigma) \\ & \text{for} : (e\pi)^{\frac{p}{2}} \text{ det}(\Sigma)^{-\frac{1}{2}} \exp\{-\frac{1}{2}(x-\mu)^{T} \Sigma^{-1}(x-\mu)^{T} \\ & \text{Hxlt}) : \mathcal{E}(e^{+TX}) : e^{-T\mu + \frac{p}{2}t^{2}} \text{ Zet} \end{split}$$

(1) Contous of the mull

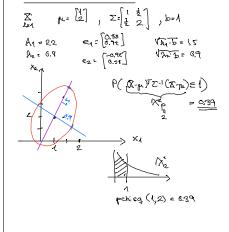
For some constant a (a>0) the solution to f(x)=a are called the contour of X.

 $f(x) = const. exp (-\frac{1}{2}(x-\mu)^T Z^{-1}(x-\mu)^2)$ rework with

(X-m) = b (6>0)

What is this graphical object?

Exem KROIY Ibrodilied



3

2 Linear combinetions

$$X \sim N_p(\mu, \Sigma)$$
 and B and CER² pm constants

PROOF: Mx(t) = exp(tTu+ 2tT I+)

TWS IS MIN with ECY) = Butc end GV(Y) = BEOT

Exan K2014 la

(Cover, one = 0 implies independence for my N

a) If
$$X_1$$
 and X_2 are independent.

Find

Cor $(X_1, X_2) = 0$

for general X 's.

But: Cou(X1, X2)=0 does not in general imply that I and I are independent, e.g.

How does Mx(t) look like when X, and X, ere

How does right) lose blue when
$$X_1$$
 and X_2 are independent when $X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_{\varphi}(\mu, \Sigma)$.

 $\begin{cases} p & p_1 \\ p_2 & p_3 \end{cases}$
 $\begin{cases} p_1 & p_2 \\ p_3 & p_4 \end{cases}$
 $\begin{cases} p_1 & p_3 \\ p_4 & p_4 \end{cases}$
 $\begin{cases} p_1 & p_3 \\ p_4 & p_4 \end{cases}$
 $\begin{cases} p_1 & p_3 \\ p_4 & p_4 \end{cases}$
 $\begin{cases} p_1 & p_4 \\ p_4 & p_4 \end{cases}$
 $\begin{cases} p_1 & p_4 \\ p_4 & p_4 \end{cases}$

6

X~ Np(µ, Z), partition into

$$\begin{array}{c}
X = \begin{pmatrix} X_{A} \\ P_{1} \times 1 \\ X_{2} \\ P_{2} \times A \end{pmatrix} \qquad \begin{array}{c}
P_{1} = \begin{pmatrix} P_{1} \\ P_{1} \\ P_{2} \end{pmatrix} / Z = \begin{pmatrix} Z_{1} & Z_{1} \\ Z_{2} & Z_{2} \\ Z_{2} & Z_{2} \end{pmatrix} \\
\begin{array}{c}
Z_{1} = G_{V}(X_{A}) \\
P_{1} \times P_{2} \\
Z_{2} = G_{V}(X_{A}, X_{2}) \\
P_{1} \times P_{2} \\
Z_{2} = Z_{2}
\end{array}$$

$$X_1$$
 and X_2 are independent iff $X_2 = Cov(X_1, X_2) = Over Cover C$

Why: Z12=0 (Zz1=0 and then

Mx(t): Mx1(ta) · Mx2(tz).

NB only if Np. E-multivanate normal.

Ex: Independent vonebles: (1,3) (1,4), (2,5)

Homework: K2014 la on sliae 17

(5) Independence of AX and BX

$$X \sim N_p(\mu, \Sigma)$$

AX and BX are independent of AIBT = 0.

unow: Hyle), exp(+++++=+)

which:

So If $A \subseteq O = ^T O \subseteq A \subseteq A$ and BR end until the second of the second

(This vorsion did not rely on result @ and @.)

6 Conditional distribution of min

a: What is the distribution of Iz | X = x1.

Observe E(X2 | X1=x1) is linear in X1
Cov (X2 (X1=x1) not dependent on X1

Added after the lecture: shorter version of "why" that builds on both @ and Q.

X ~ Np(M, Z) and Q, B
quipp quipp
Y = \big[Y_1] = \big[AX] \cdot [A] \times - QX.

From @ we know that Yn Ng(G,CCC).

From @ we know that Yn and Yz

one independent iff Car(Y_1, Y_2) = 0;

and Car(Y_1, Y_2) = Car(AB, BX)

= A Car(X, X) Bt = A \sum BT, so

Car(X)

A \sum Bt = 0 \implies AX and BX

independent.

9

TMA4267 Linear Statistical Models V2017 [L6]

Part 1: Multivariate RVs, and the multivariate normal distribution Estimators for mean and covariance Quadratic forms [H:3.3,4.5,5.1,F:AppB3]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: January 27, 2017

Plan for today

- estimators for mean and covariance
- quadratic forms and rules for quadratic forms
- ▶ idempotent matrices
- ▶ more rules for quadratic forms with idempotent matrices

1/17

Maximum likelihood estimators

Let X_1, X_2, \ldots, X_n be a random sample of size n from the multivariate normal distribution $N_p(\mu, \Sigma)$. The maximum likelihood estimators for are found by maximizing the likelihood:

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{j=1}^{n} f(\boldsymbol{x}_{j}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$= \prod_{j=1}^{n} (\frac{1}{2\pi})^{\frac{p}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\{-\frac{1}{2} (\boldsymbol{x}_{j} - \boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_{j} - \boldsymbol{\mu})\}$$

Maximum likelihood estimators

Let $X_1, X_2, ..., X_n$ be a random sample of size n from the multivariate normal distribution $N_n(\mu, \Sigma)$.

2/17

Maximum likelihood estimators

Let X_1, X_2, \ldots, X_n be a random sample of size n from the multivariate normal distribution $N_p(\mu, \Sigma)$. The maximum likelihood estimators for are found by maximizing the likelihood:

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{j=1}^{n} f(\boldsymbol{x}_{j}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$= \prod_{j=1}^{n} (\frac{1}{2\pi})^{\frac{p}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\{-\frac{1}{2} (\boldsymbol{x}_{j} - \boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_{j} - \boldsymbol{\mu})\}$$

Could take In and then partial derivatives, but easier to add and subtract the mean \bar{x} and rewrite (using trace-formulas)

Maximum likelihood estimators

Let X_1, X_2, \ldots, X_n be a random sample of size n from the multivariate normal distribution $N_p(\mu, \Sigma)$. The maximum likelihood estimators for are found by maximizing the likelihood:

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{j=1}^{n} f(\boldsymbol{x}_{j}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$= \prod_{i=1}^{n} (\frac{1}{2\pi})^{\frac{p}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\{-\frac{1}{2} (\boldsymbol{x}_{j} - \boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_{j} - \boldsymbol{\mu})\}$$

Could take In and then partial derivatives, but easier to add and subtract the mean \bar{x} and rewrite (using trace-formulas)

$$L(\mu, \mathbf{\Sigma}) = (\frac{1}{2\pi})^{\frac{np}{2}} \det(\mathbf{\Sigma})^{-\frac{n}{2}}$$

$$\exp\{-\frac{1}{2} \left[\operatorname{tr}(\mathbf{\Sigma}^{-1} \sum_{i=1}^{n} (\mathbf{x}_{j} - \bar{\mathbf{x}})(\mathbf{x}_{j} - \bar{\mathbf{x}})^{T}) + n(\bar{\mathbf{x}} - \mu)^{T} \mathbf{\Sigma}^{-1} (\bar{\mathbf{x}} - \mu) \right] \}$$

2/17

Maximum likelihood estimators: then for Σ

$$\begin{split} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (\frac{1}{2\pi})^{\frac{n\rho}{2}} \det(\boldsymbol{\Sigma})^{\frac{n}{2}} \\ &\exp\{-\frac{1}{2}[\operatorname{tr}(\boldsymbol{\Sigma}^{-1} \sum_{i=1}^{n} (\boldsymbol{x}_{j} - \bar{\boldsymbol{x}})(\boldsymbol{x}_{j} - \bar{\boldsymbol{x}})^{T}) + n(\bar{\boldsymbol{x}} - \boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1}(\bar{\boldsymbol{x}} - \boldsymbol{\mu})]\} \end{split}$$

A maximization theorem for matrices it used to find that the MLE for Σ is

$$\widehat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{j=1}^{n} (\boldsymbol{X}_{j} - \bar{\boldsymbol{X}}) (\boldsymbol{X}_{j} - \bar{\boldsymbol{X}})^{T}$$

Maximum likelihood estimators: first for μ

$$egin{aligned} \mathcal{L}(oldsymbol{\mu}, oldsymbol{\Sigma}) &= (rac{1}{2\pi})^{rac{np}{2}} \det(oldsymbol{\Sigma})^{-rac{n}{2}} \ \exp\{-rac{1}{2}[\mathrm{tr}(oldsymbol{\Sigma}^{-1}\sum_{j=1}^{n}(oldsymbol{x}_{j}-ar{oldsymbol{x}})(oldsymbol{x}_{j}-ar{oldsymbol{x}})^{T}) + n(ar{oldsymbol{x}}-oldsymbol{\mu})^{T}oldsymbol{\Sigma}^{-1}(ar{oldsymbol{x}}-oldsymbol{\mu})]\} \end{aligned}$$

and see directly for SPD $oldsymbol{\Sigma}$ that the maximum is achieved for $\mu=ar{oldsymbol{x}},$ so that the MLE for μ is

$$\bar{\boldsymbol{X}} = \frac{1}{n} \sum_{j=1}^{n} \boldsymbol{X}_{j}$$

3 / 17

Properties of the ML estimators

- $ightharpoonup ar{X}$ is distributed as $N_p(\mu, \frac{1}{n}\Sigma)$
- ▶ nS is distributed as a Wishart random matrix with n − 1 degrees of freedom.
- $ightharpoonup ar{X}$ and $nm{S}$ are independent.

The Wishart distribution is not on the reading list for TMA4267. General properties of maximum likelihood estimation is covered in detail in TMA4295 Statistical Inference.

Quadratic forms - first results [F:B3.3, Theorem B.2]

We stay with our random vector \boldsymbol{X} with $\boldsymbol{\mu}$ and covariance matrix Σ , and a symmetric constant matrix A.

- \blacktriangleright What is a quadratic form? X^TAX
- ▶ The "trace-formula": $E(X^TAX)$.

6 / 17

Useful facts about the trace [H:2.1] and [F:Theorem A.18]

Let A, B and C be conformable matrices

$$tr(A + B) = tr(A) + tr(B)$$

 $tr(AB) = tr(BA)$
 $tr(ABC) = tr(CAB) = tr(BCA)$

8/17

Exam V2014: Problem 1a

Let
$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$$
 be a random vector with mean $\boldsymbol{\mu} = \mathrm{E}(\mathbf{X}) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ and covariance matrix $\mathbf{\Sigma} = \mathrm{Cov}(\mathbf{X}) = \mathbf{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$. Further, let

and covariance matrix
$$\mathbf{\Sigma} = \operatorname{Cov}(\mathbf{X}) = \mathbf{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$
. Further, let

$$\mathbf{A} = \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{pmatrix} \text{ be a matrix of constants.}$$

Define
$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \end{pmatrix} = \mathbf{AX}$$
.

Find $E(\mathbf{Y})$ and $Cov(\mathbf{Y})$. Are X_1 and X_2 independent? Are Y_1 and Y_2 independent? Justify your answers.

Find the mean of X^TAX .

7 / 17

Quadratic forms - last results [F:B3.3, Theorem B.2]

Now: \boldsymbol{X} is multivariate normal with mean μ and covariance matrix I, and we also have a symmetric and idempotent matrix $R_{(p \times p)}$ with rank r.

- ▶ Properties of an idempotent matrix.
- ▶ Distribution of $\boldsymbol{X}^T \boldsymbol{R} \boldsymbol{X} \sim \chi_r^2$.
- ▶ Distribution of a ratio of two quadratic forms and the Fisher distribution.

Properties of symmetric idempotent matrices

A symmetric matrix \mathbf{A} is idempotent, $\mathbf{A}^2 = \mathbf{A}$, and has the following properties (to be proven in RecEx1.P7).

10 / 17

Properties of symmetric idempotent matrices

A symmetric matrix \boldsymbol{A} is idempotent, $\boldsymbol{A}^2 = \boldsymbol{A}$, and has the following properties (to be proven in RecEx1.P7).

- 1. The eigenvalues are 0 and 1.
- 2. The rank of a symmetric matrix (actually: a diagonalizable quadratic matrix) equals the number of nonero eigenvaluse of the matrix. Should be known from previous courses.

Properties of symmetric idempotent matrices

A symmetric matrix \boldsymbol{A} is idempotent, $\boldsymbol{A}^2 = \boldsymbol{A}$, and has the following properties (to be proven in RecEx1.P7).

1. The eigenvalues are 0 and 1.

10 / 17

Properties of symmetric idempotent matrices

A symmetric matrix \boldsymbol{A} is idempotent, $\boldsymbol{A}^2 = \boldsymbol{A}$, and has the following properties (to be proven in RecEx1.P7).

- 1. The eigenvalues are 0 and 1.
- 2. The rank of a symmetric matrix (actually: a diagonalizable quadratic matrix) equals the number of nonero eigenvaluse of the matrix. Should be known from previous courses.
- 3. (Combining 1+2). If a $(n \times n)$ symmetric idempotent matrix \boldsymbol{A} has rank r then r eigenvalues are 1 and n-r are 0.

10 / 17

Properties of symmetric idempotent matrices

A symmetric matrix \boldsymbol{A} is idempotent, $\boldsymbol{A}^2 = \boldsymbol{A}$, and has the following properties (to be proven in RecEx1.P7).

- 1. The eigenvalues are 0 and 1.
- 2. The rank of a symmetric matrix (actually: a diagonalizable quadratic matrix) equals the number of nonero eigenvaluse of the matrix. Should be known from previous courses.
- 3. (Combining 1+2). If a $(n \times n)$ symmetric idempotent matrix **A** has rank r then r eigenvalues are 1 and n-r are 0.
- 4. The trace and rank of a symmetric projection matrix are equal: $tr(\mathbf{A}) = rank(\mathbf{A}).$

10 / 17

The Chi-square distribution

pdf χ_p^2 :

$$f(y) = \frac{1}{2^{p/2}\Gamma(p/2)}y^{p/2-1}e^{(-y/2)}$$
 for $y > 0$

MGF χ_p^2 :

$$M_Y(t) = rac{1}{(1-2t)^{p/2}}$$

Addition property:

Let $X_1\sim \chi_p^2$ and $X_2\sim \chi_q^2$, and let X_1 and X_2 be independent. Then $X_1+X_2\sim \chi_{p+q}^2$.

Subtraction property:

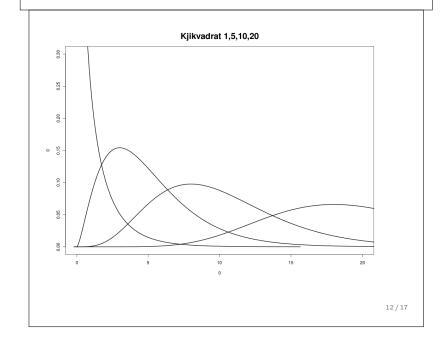
Let $X=X_1+X_2$ with $X_1\sim\chi_p^2$ and $X\sim\chi_{p+q}^2$. Assume that X_1 and X_2 are independent. Then $X_2\sim\chi_q^2$.

Properties of symmetric idempotent matrices

A symmetric matrix \boldsymbol{A} is idempotent, $\boldsymbol{A}^2 = \boldsymbol{A}$, and has the following properties (to be proven in RecEx1.P7).

- 1. The eigenvalues are 0 and 1.
- 2. The rank of a symmetric matrix (actually: a diagonalizable quadratic matrix) equals the number of nonero eigenvaluse of the matrix. Should be known from previous courses.
- 3. (Combining 1+2). If a $(n \times n)$ symmetric idempotent matrix **A** has rank r then r eigenvalues are 1 and n-r are 0.
- 4. The trace and rank of a symmetric projection matrix are equal: $tr(\mathbf{A}) = rank(\mathbf{A}).$
- 5. The matrix $\mathbf{I} \mathbf{A}$ is also idempotent, and $\mathbf{A}(\mathbf{I} \mathbf{A}) = 0$.

10 / 17



The Fisher distribution [F: B.1 Def 8.14], RecEx2.P5+6

"Tabeller og formeler i statistikk":

If Z_1 and Z_2 are independent and χ^2 -distributed with ν_1 and ν_2 degrees of freedom, then

$$F = \frac{Z_1/\nu_1}{Z_2/\nu_2}$$

is F(isher)-distributed with ν_1 and ν_2 degrees of freedom.

- ▶ The expected value of F is $E(F) = \frac{\nu_2}{\nu_2 2}$.
- ► The mode is at $\frac{\nu_1 2}{\nu_1} \frac{\nu_2}{\nu_2 + 2}$.
- Identity:

$$f_{1-\alpha,\nu_1,\nu_2} = \frac{1}{f_{\alpha,\nu_2,\nu_1}}$$

13 / 17

Quadratic forms [F:B3.3, Theorem B.2]

Random vector ${\pmb X}$ with mean ${\pmb \mu}$ and covariance matrix ${\pmb \Sigma}$, symmetric constant matrix ${\pmb A}$.

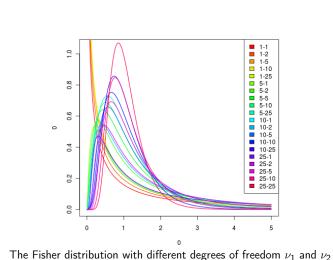
- ▶ Quadratic form: $X^T A X$.
- ► The "trace-formula": $E(X^TAX) = tr(AΣ) μ^TAμ$.

Then, let $m{X} \sim N_p(\mathbf{0}, m{I})$, and $m{R}$ is a symmetric and idempotent matrix with rank r.

$$\boldsymbol{X}^T \boldsymbol{R} \boldsymbol{X} \sim \chi_r^2$$

Now, also ${\bf S}$ is a symmetric and idempotent matrix with rank ${\bf s}$, and ${\bf R}{\bf S}={\bf 0}$.

$$\frac{s\boldsymbol{X}^T\boldsymbol{R}\boldsymbol{X}}{r\boldsymbol{X}^T\boldsymbol{S}\boldsymbol{X}}\sim F_{r,s}$$



The Fisher distribution with different degrees of freedom ν_1 and ν_2 (given in the legend).

14 / 17

Plan for the last week of Part 1

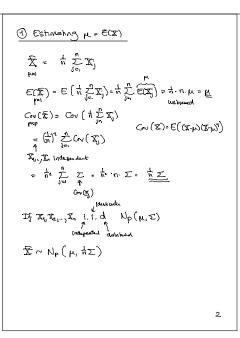
Supervision in lecture times.

See Blackboard: Part 1: dates and places for supervision.

Kahoot!

Summing up the last three lectures with a few multiple choice questions.

17 / 17



Foremeter schmetton: μ and Σ 27.01.207

[145.3, 1445]

X rendom vector equality of solven

PAT PM

THE PM

THE PM

THE PM

TO THE POPULATION

O VENIN FOR the population

3 Vonetions on Reform to the population

3 Vonetions on Reform to the population

3 Vonetions on Reform to the population

3 Vonetions on Refore the PM

Estimating \mathbb{Z} Plotivation: $\mathbb{Z} \cdot \mathbb{E} \left((X, \mu)(X, \mu)^T \right)$ $S = \frac{1}{n-1} \int_{\mathbb{R}^n}^{\Lambda} (X_j - \widehat{X})(X_j - \widehat{X})^T$ Prop $\mathbb{E}(S) = \dots = \mathbb{Z}$ unbiased

See proof, separate the

end I = Cov(X)?

3

Quedretic forms [F: 833, Theorem 8.2]

X Pandom rector, A constant matrix

$$X^T A X = \sum_{i=1}^{p} \sum_{j=1}^{p} a_{ij} X_i \cdot X_j$$
 16 a quadrete form.

 $\mu = G(Z), Z = Gv(Z)$

The "brece formula": Cross)

E(XTAX) = tr(AZ) + pTAp

Ex: 12014 Ma

4

Astribution of quadratic form

$$X \sim N_{P}(0, I)$$
PRI $X_{1...,X_{P}}$ se independent

R 15 symmetric and idempotent with santi(R)=r.

Coult: XTRX ~ Xx

Proof via R=PAPT

XTRX = sum of ~ NO.D2 ~ X2,

In comp & 1.12c: X~ Nn (pl) of I)

Nn(0, I)

B

Symmetric and idempotent matrices

A 15 symmetric; A=AT A 13 idempotent; AR=A

* The eigenvalues of A = 0 and 1

K The rent of a metrix 13 the number of linearly independent rows, and also given as the number of nonzero eigenvalues.

F Genual rule: tr(A)= \(\frac{1}{2} \) \(\hat{\chi} \) , since \(\hat{\chi} = 0 \) or \(\hat{\chi} \) must be the nucleu of nonzero exogeneous.

⇒ tr(A)=fenk(A) for symmetric and lampolent A.

Nov: example = the centering matrix. (RECEXI. PY, CompGN. 2a)

5

Finally: Ratio of quadrate forms (more in Part 2-3)

- * X ~ No(0, I)
- * R and S symmetric and idempotent with rank (R)=r, rank(S)=S
- * RS=0

Result: (a) $X^TKX \sim X^C / X^TSX \sim X^S$ End independent.

Because. RR and SX ere under if $R(x/X)S^T=0$ We have Car(X)=T and S=ST so RTST=RS=0 \Rightarrow yes.'

Than a function of R first and function of R g(8%) an also independent R See Not get R with part R with part R and R and R and R are R and R are

This result will be used alot in Pat 2+3.

7-

Here
$$f(RX) = \uparrow (RX) \uparrow RX = \uparrow X \uparrow R \uparrow R X$$

= $\uparrow X \uparrow R X$

and the same for $g(SX)$.

Trivariate normal pdf

What graphical form has the solution to $f(\mathbf{x}) = \text{constant}$?

A Circle

B Parabola

C Ellipsoid

D Bell shape

Multivariate normal pdf

The probability density function is $(\frac{1}{2\pi})^{\frac{p}{2}}\det(\mathbf{\Sigma})^{-\frac{1}{2}}\exp\{-\frac{1}{2}Q\}$ where Q is

A
$$(\mathbf{x} - \mathbf{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \mathbf{\mu})$$

B
$$(x - \mu)\Sigma(x - \mu)^T$$

$$\Sigma - \mu$$

Multivariate normal distribution

 $m{X}_p \sim N_p(\mu, m{\Sigma})$, and $m{C}$ is a $k \times p$ constant matrix. $m{Y} = m{C} m{X}$ is

- f A Chi-squared with k degrees of freedom
- B Multivariate normal with mean $k\mu$
- C Chi-squared with *p* degrees of freedom
- **D** Multivariate normal with mean $C\mu$

Independence

Let
$$\boldsymbol{X} \sim N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
, with $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 1 & 0 \\ 2 & 3 & 1 \\ 0 & 2 & 5 \end{bmatrix}$. Which

two variables are independent?

- **A** X_1 and X_2
- \mathbf{B} X_1 and X_3
- C X_2 and X_3
- D None but two are uncorrelated.

Conditional distribution: mean

 $\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ is a bivariate normal random vector. What is true for the conditional mean of X_2 given $X_1 = x_1$?

- A Not a function of x_1
- B A linear function of x_1
- C A quadratic function of x_1

Constructing independent variables?

Let $X \sim N_p(\mu, \Sigma)$. How can I construct a vector of independent standard normal variables from X?

- $\mathbf{C} \quad \mathbf{\Sigma}^{-\frac{1}{2}}(\mathbf{X} \mathbf{\mu})$
- $\mathbf{D} \quad \mathbf{\Sigma}^{\frac{1}{2}}(\mathbf{X} + \mathbf{\mu})$

Conditional distribution: variance

 $\boldsymbol{X} = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ is a bivariate normal random vector. What is true for the conditional variance of X_2 gi-

ven $X_1 = x_1$?

- A Not a function of x_1
- **B** A linear function of x_1
- C A quadratic function of x_1

Estimator for mean

 $\boldsymbol{X}_1, \boldsymbol{X}_2, \ldots, \boldsymbol{X}_n$ is a random sample from $N_p(\mu, \Sigma)$. What is the distribution of the estimator $\bar{\boldsymbol{X}}$ for the mean?

A $N_n(\mu, \Sigma)$

B $N_p(\mu, \frac{1}{n}\Sigma)$

C χ_p^2

 \mathbf{D} χ_n^2

Unbiased estimators

 $\boldsymbol{X}_1, \boldsymbol{X}_2, \ldots, \boldsymbol{X}_n$ is a random sample of size n of a p-dimensional random vector. An unbiased estimator for the covariance matrix $\boldsymbol{\Sigma}$ is.

$$\mathbf{A} \quad \frac{1}{n} \sum_{j=1}^{n} (\mathbf{X}_{j} - \bar{\mathbf{X}}) (\mathbf{X}_{j} - \bar{\mathbf{X}})^{T}$$

$$egin{array}{ll} egin{array}{ll} rac{1}{n-1} \sum_{j=1}^n (oldsymbol{X}_j - ar{oldsymbol{X}}) (oldsymbol{X}_j - ar{oldsymbol{X}})^T \end{array}$$

$$\mathbf{C}$$
 $\frac{1}{n}\sum_{j=1}^{n}(\mathbf{X}_{j}-\bar{\mathbf{X}})^{T}(\mathbf{X}_{j}-\bar{\mathbf{X}})$

$$\mathbf{D} \quad \frac{1}{n-1} \sum_{j=1}^{n} (\mathbf{X}_{j} - \bar{\mathbf{X}})^{T} (\mathbf{X}_{j} - \bar{\mathbf{X}})$$

Distribution of quadratic form

 $X \sim N_p(\mathbf{0}, \mathbf{I})$, and \mathbf{R} is a symmetric and idempotent matrix with rank r. What is the distribution of $X^T R X$?

A $N_p(\mu, rI)$

 \mathbf{B} $N_r(\mathbf{0}, \mathbf{I})$

C χ_r^2

 \mathbf{D} χ_p^2

Correct?

Are you sure you want to read the correct answers? Maybe try first? The answers are explained on the next two slides.

Answers

- 1. A: exponent quadratic form is $(\mathbf{x} \mathbf{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} \mathbf{\mu})$.
- 2. C: contours are ellipsoids in general. In two dimensions we have ellipses. For two dimensions and equal variance and correlation 0 we have circles.
- 3. D: linear combinations of mvN are also mvN.
- 4. B: $Cov(X_1, X_3) = 0$ and X_1 and X_3 are thus independent.
- 5. C: The Mahlanobis transform is $\Sigma^{-\frac{1}{2}}(\boldsymbol{X} \boldsymbol{\mu})$.

Answers

- 6. B: Conditional mean is linear in x_1 , which will be very useful when we start with multiple linear regression.
- 7. A: Conditional variance (covariance) is not a function of x_1 .
- 8. B: The mean is also mvN with mean μ and covariance $\frac{1}{n}\Sigma$.
- 9. B: $\frac{1}{n-1}\sum_{j=1}^n (\boldsymbol{X}_j \bar{\boldsymbol{X}})(\boldsymbol{X}_j \bar{\boldsymbol{X}})^T$ is the unbiased estimator for Σ . Observe the (n-1) and that the dimension is $p \times p$ (to place the transpose). Not a quadratic form.
- 10. C: Quadratic form is related to χ^2 .