

Special course on Gaussian processes: Session #1

Michael Riis Andersen

Aalto University

michael.riis@gmail.com

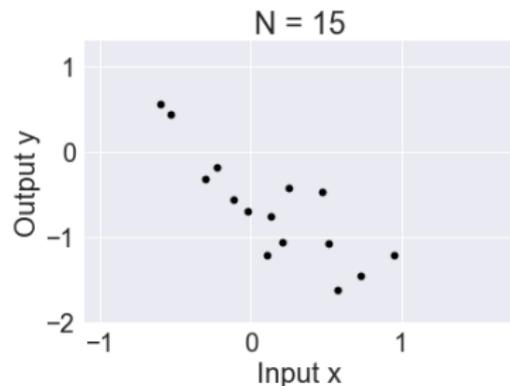
9/1-19

Agenda for today

- 1 Motivation for Gaussian processes
- 2 Course content, format, and evaluation
- 3 Warm up for Gaussian processes: Review of the multivariate Gaussian distribution
- 4 First assignment

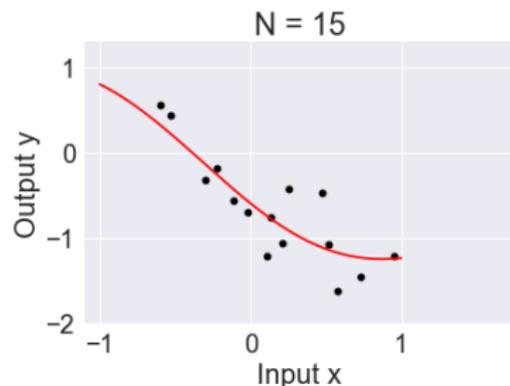
Gaussian processes in a nutshell

- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



Gaussian processes in a nutshell

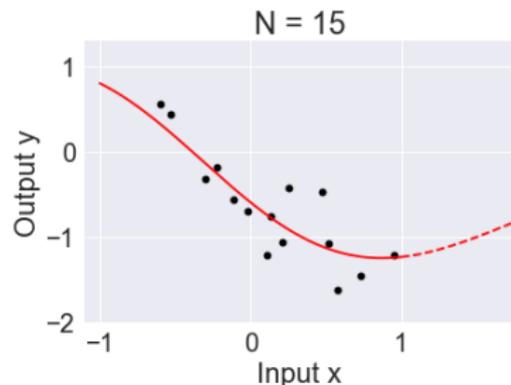
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data

Gaussian processes in a nutshell

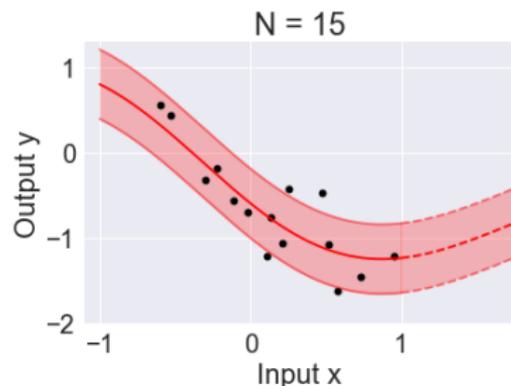
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs

Gaussian processes in a nutshell

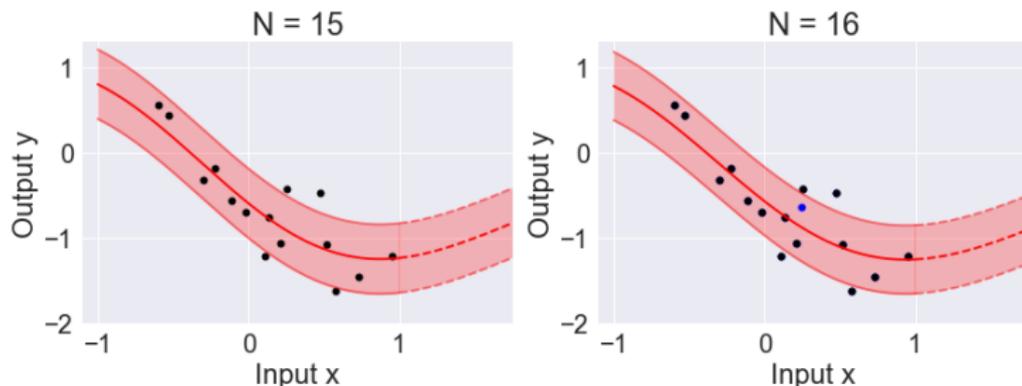
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties

Gaussian processes in a nutshell

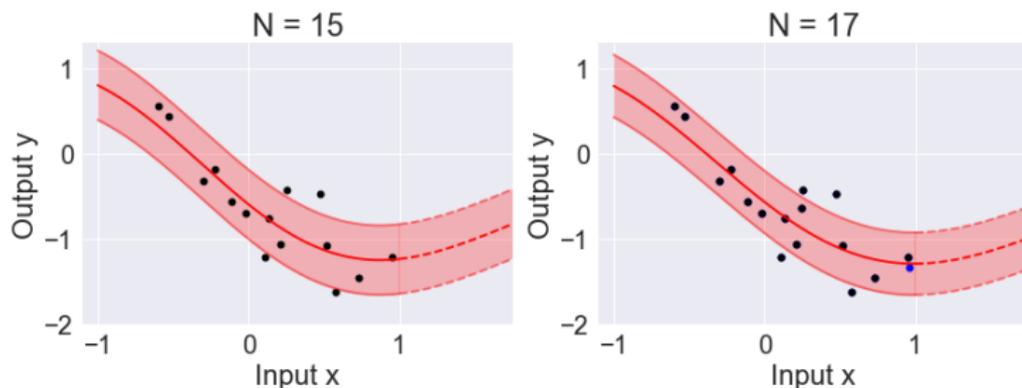
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

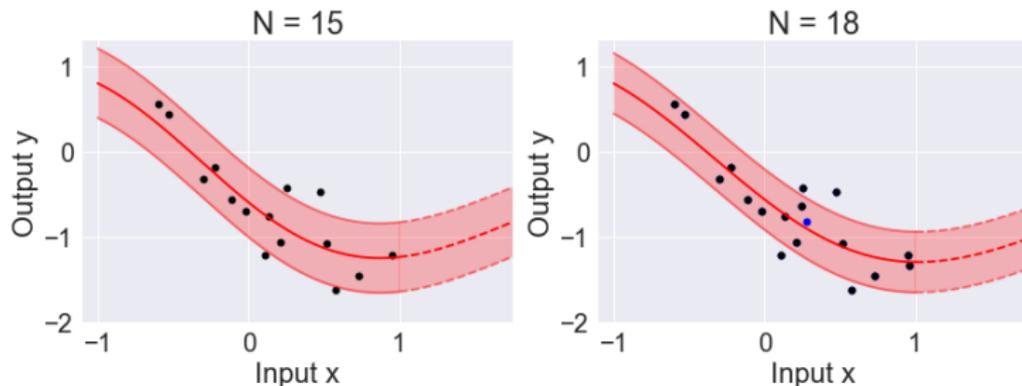
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

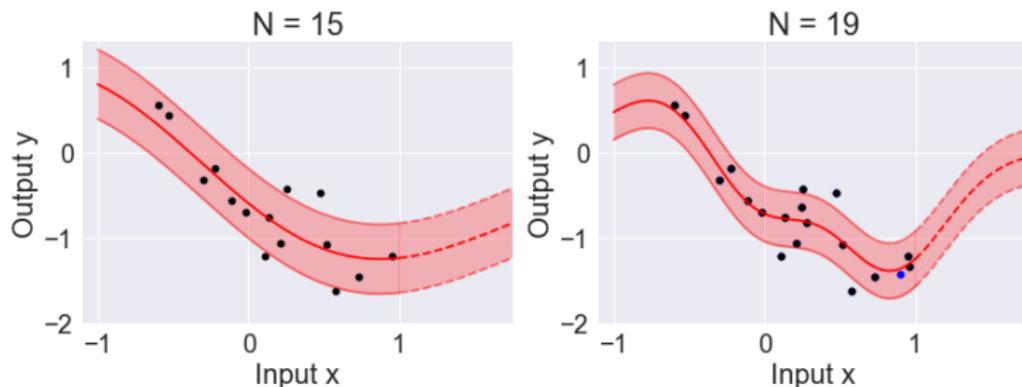
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

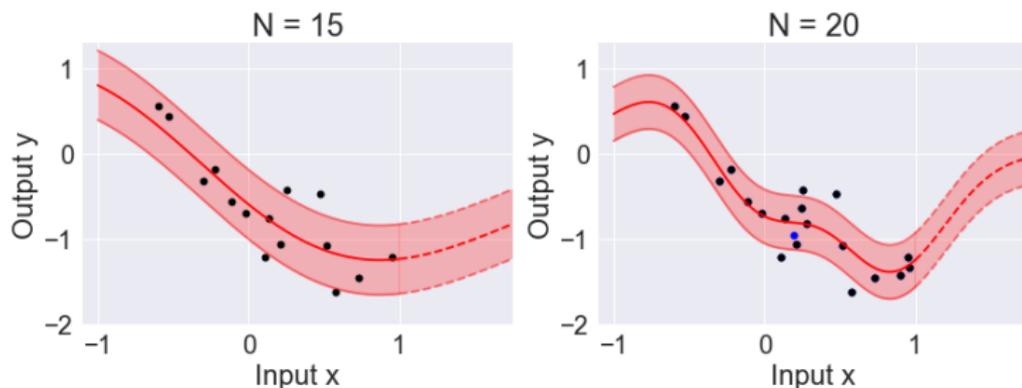
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

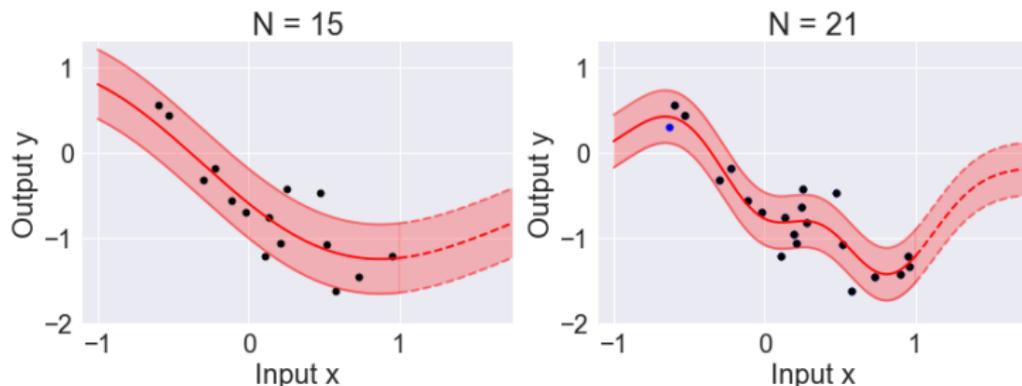
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

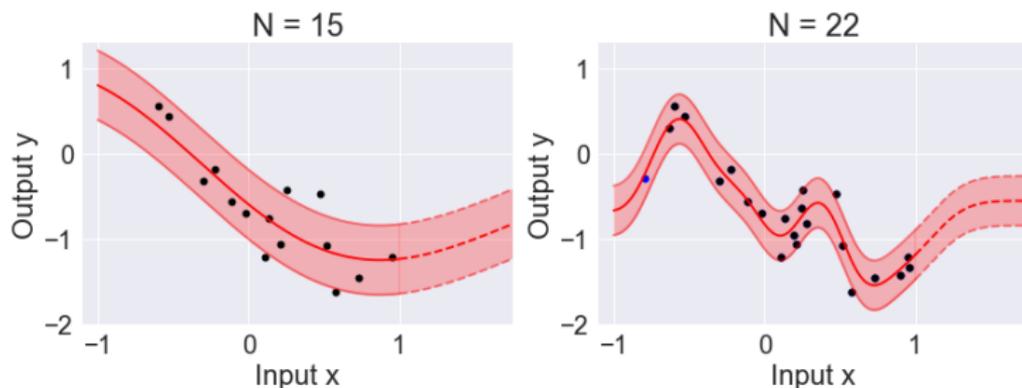
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

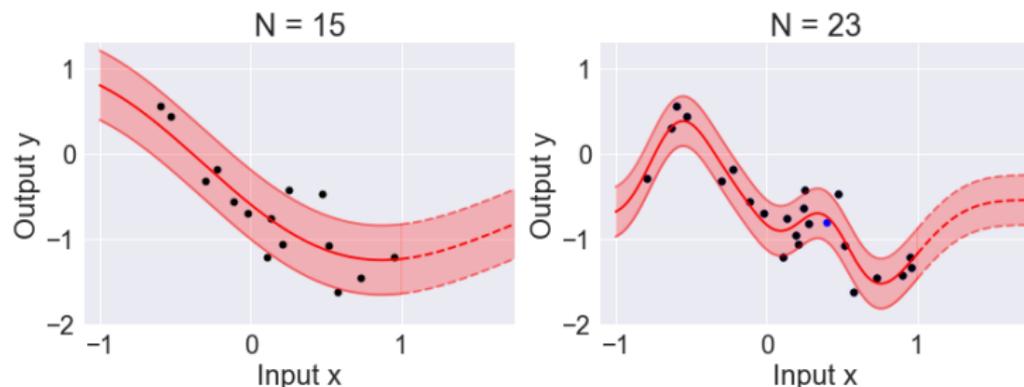
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

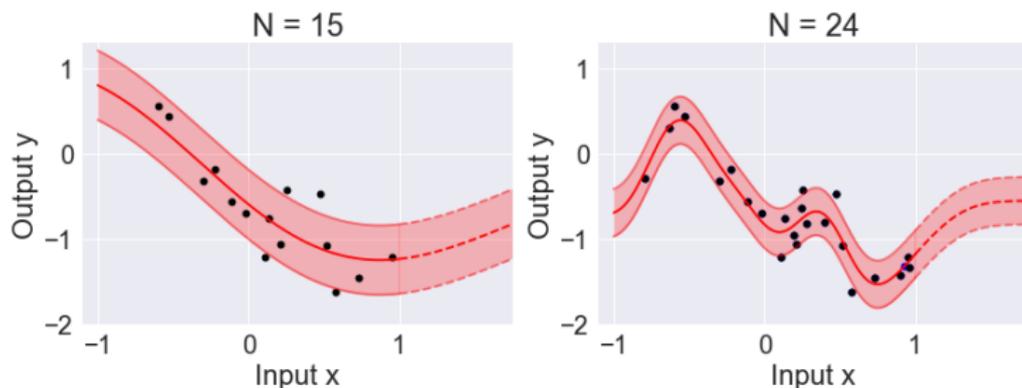
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

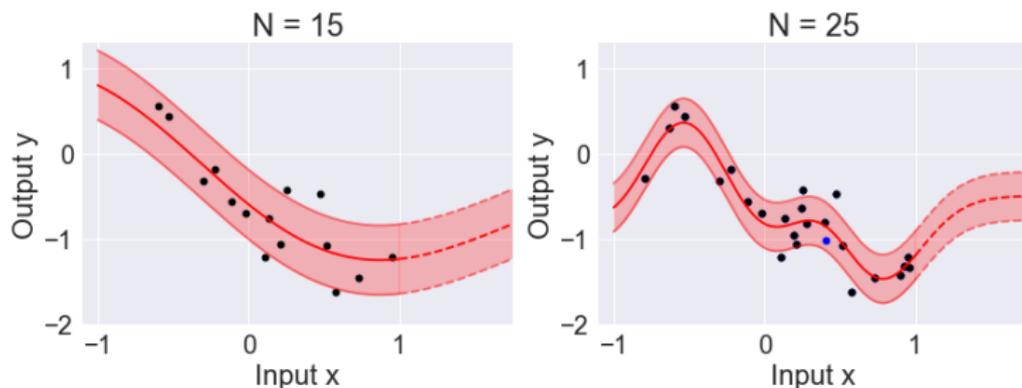
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

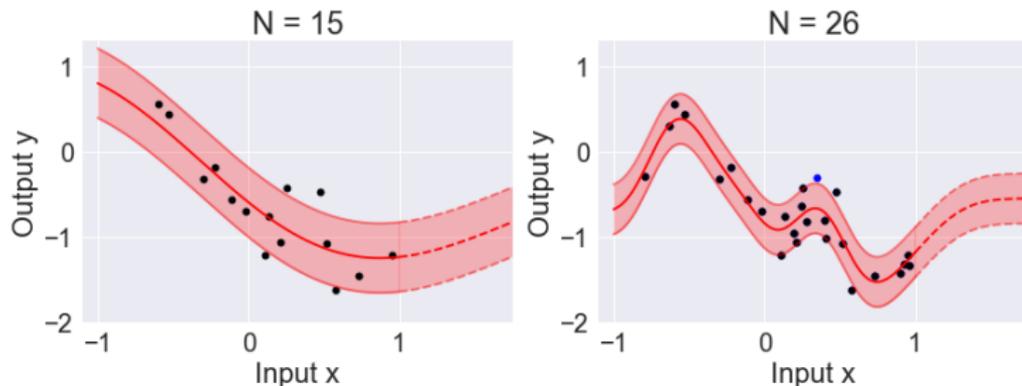
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

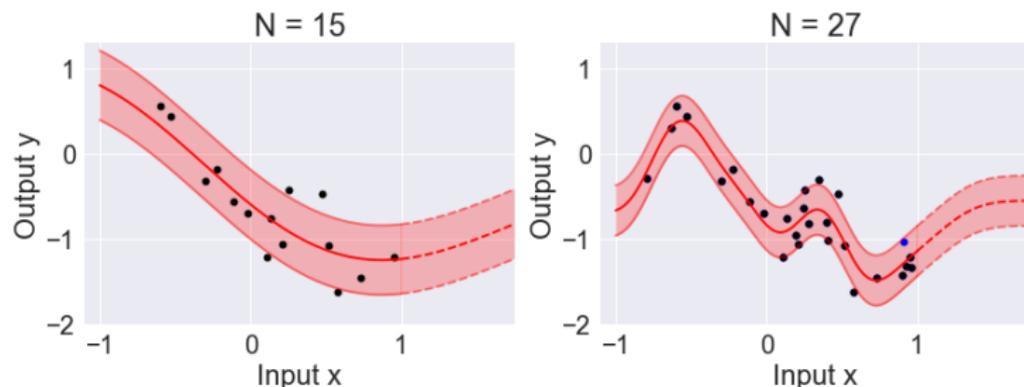
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

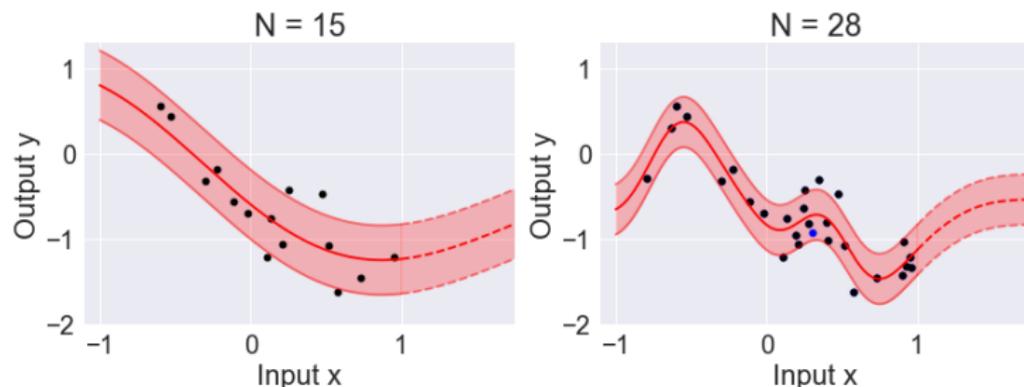
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

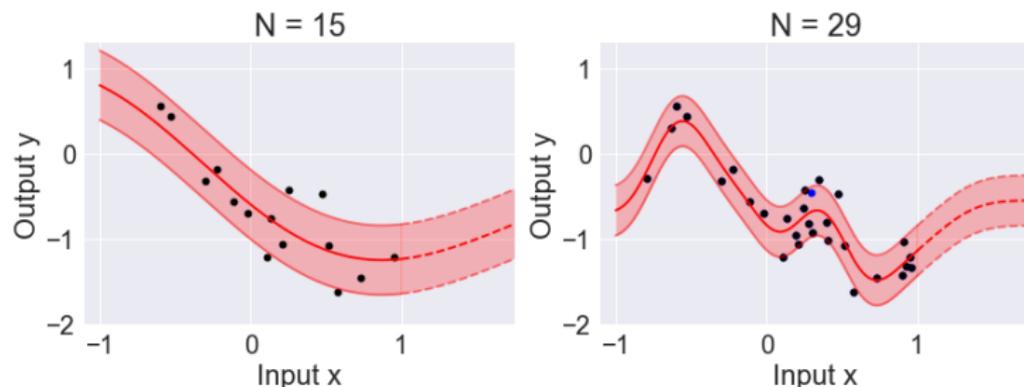
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Gaussian processes in a nutshell

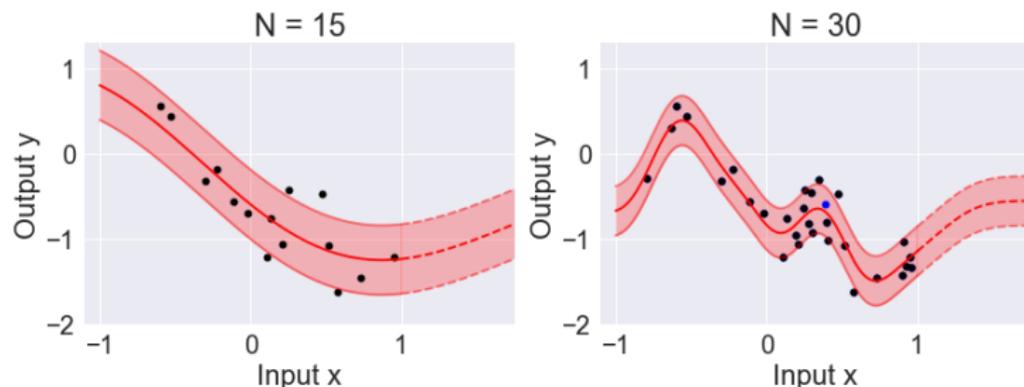
- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

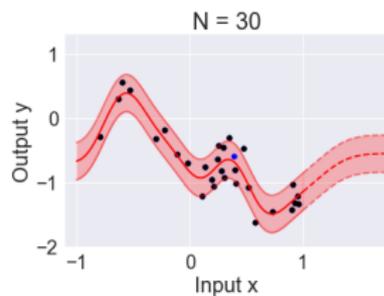
Gaussian processes in a nutshell

- It's all about learning functions from data
- Suppose we are given a data set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$



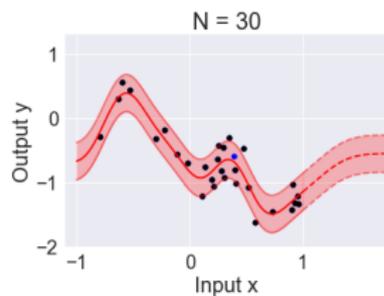
- Gaussian processes (GPs) can
 - ... fit non-linear functions to data
 - ... make predictions for new inputs
 - ... provide sensible uncertainties
 - ... adjust model complexity to data (nonparametric)

Functions with different domains

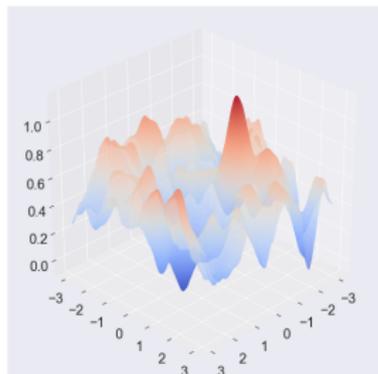


The real line

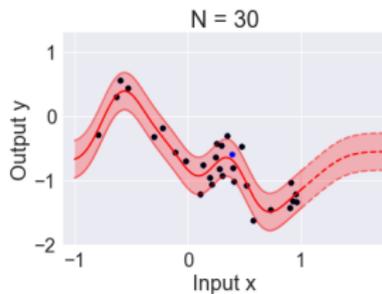
Functions with different domains



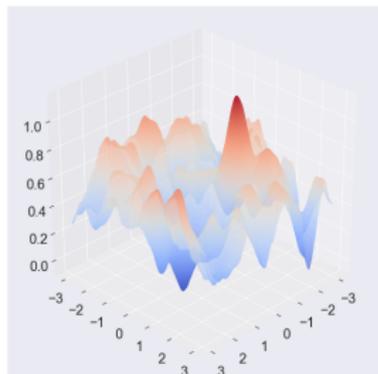
The real line



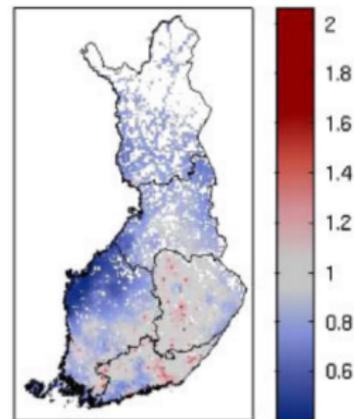
Functions with different domains



The real line

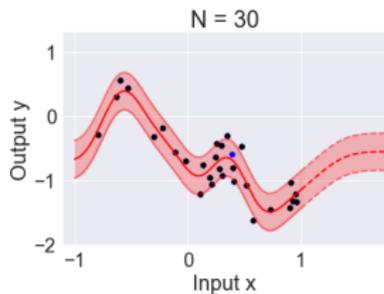


Higher dimensions

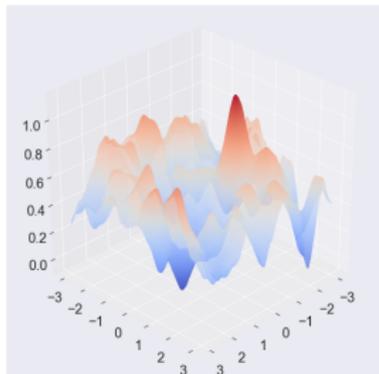


Finland

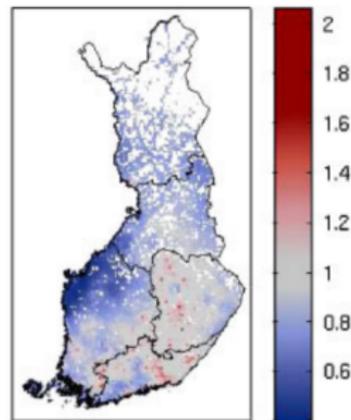
Functions with different domains



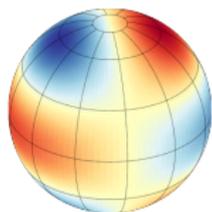
The real line



Higher dimensions

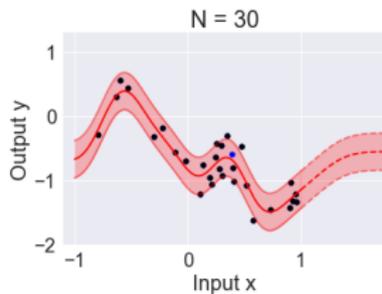


Finland

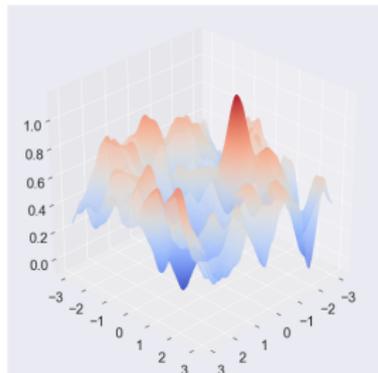


A sphere

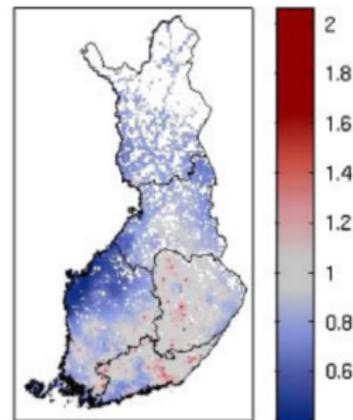
Functions with different domains



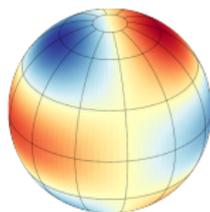
The real line



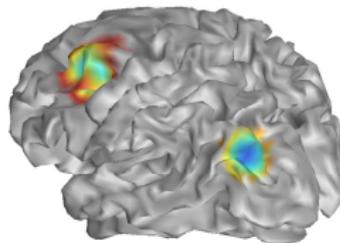
Higher dimensions



Finland



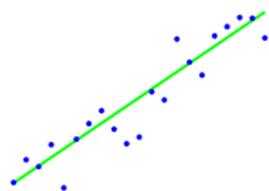
A sphere



A human brain

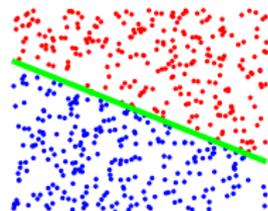
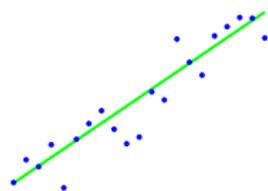
Multitude of Gaussian processes applications

- Regression (supervised learning)
 - Time series analysis
 - EEG brain imaging
 - Survival analysis for cancer data
 - Predicting rainfall
 - Robot dynamics
 - ...



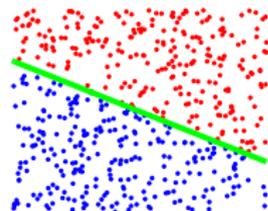
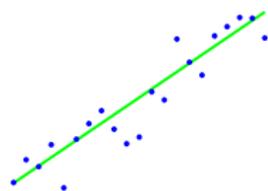
Multitude of Gaussian processes applications

- Regression (supervised learning)
 - Time series analysis
 - EEG brain imaging
 - Survival analysis for cancer data
 - Predicting rainfall
 - Robot dynamics
 - ...
- Classification (supervised learning)
 - Recognizing handwritten digits
 - Brain decoding
 - ...



Multitude of Gaussian processes applications

- Regression (supervised learning)
 - Time series analysis
 - EEG brain imaging
 - Survival analysis for cancer data
 - Predicting rainfall
 - Robot dynamics
 - ...
- Classification (supervised learning)
 - Recognizing handwritten digits
 - Brain decoding
 - ...
- Dimensionality reduction (unsupervised learning)
- Optimization of black box functions (Bayesian optimization)
- Numerical integration (Bayesian quadrature)
- Solving differential equations (probabilistic numerics)



Course content

- The goal of the course is to introduce you to Gaussian processes, applications and some recent research directions
- We will cover
 - ① ... Gaussian process regression & classification
 - ② ... model selection for Gaussian processes
 - ③ ... approximate inference & how to speed up GP inference
 - ④ ... spatio-temporal modelling
 - ⑤ ... some advanced topics based on your interests

Format of the course

- The course will be based on
 - shorts lectures
 - exercises (based on python notebooks)
 - project work + presentation in groups of 1-3 persons (optional)
- To pass the course, you need to
 - complete and hand in exercises
 - do project work (only for extra ECTS points)
 - 3 ECTS / 5 ECTS

Lectures

- Lecture 1: Warm up: Properties of the multivariate normal distribution
- Lecture 2: Linear Gaussian models and intro to Gaussian processes
- Lecture 3: Kernels and model selection
- Lecture 4: Inducing points method (.. or how to make GPs faster)
- Lecture 5: Spectral kernels (.. or how to make GPs more flexible)
- Lecture 6: Spatio-temporal models

Assignments

- Assignment #1 due 23rd of January (midnight)
- Assignment #2 due 6th of February (midnight)
- Assignment #3 due 20th of February (midnight)

The properties of the multivariate Gaussian distribution

The multivariate Gaussian distribution

- **Definition** A random vector $\mathbf{x} = [x_1, x_2, \dots, x_D]$ is said to have the multivariate Gaussian distribution if all linear combinations of \mathbf{x} are Gaussian distributed:

$$y = \mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \dots + a_D x_D \sim \mathcal{N}(m, v)$$

for all $\mathbf{a} \in \mathbb{R}^D$, where $\mathbf{a} \neq \mathbf{0}$

The multivariate Gaussian distribution

- **Definition** A random vector $\mathbf{x} = [x_1, x_2, \dots, x_D]$ is said to have the multivariate Gaussian distribution if all linear combinations of \mathbf{x} are Gaussian distributed:

$$y = \mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \dots + a_D x_D \sim \mathcal{N}(m, v)$$

for all $\mathbf{a} \in \mathbb{R}^D$, where $\mathbf{a} \neq \mathbf{0}$

- The multivariate Gaussian density for a variable $\mathbf{x} \in \mathbb{R}^D$:

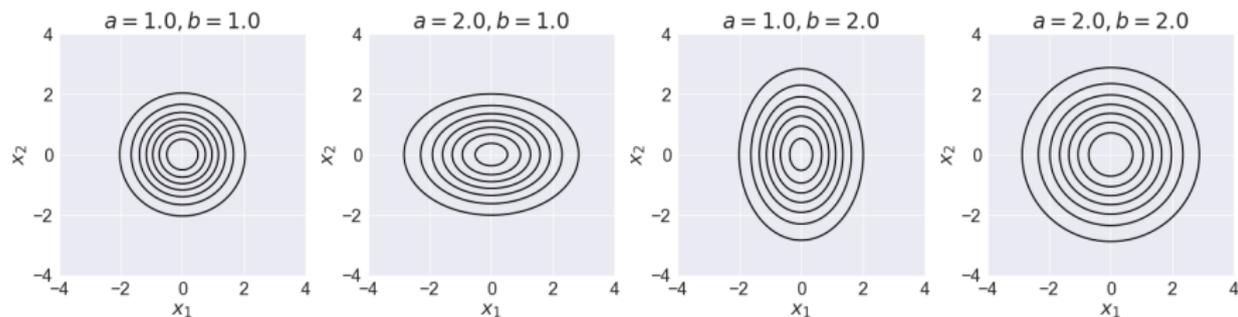
$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

- Completely described by its parameters:
 - $\boldsymbol{\mu} \in \mathbb{R}^D$ is the mean vector
 - $\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$ is the covariance matrix (positive definite)
- $(\boldsymbol{\Sigma})_{ij}$ is the covariance between the i 'th and j 'th elements in \mathbf{x}

Interpretation of the covariance matrix - 2D examples

The diagonal of the covariance controls the scaling/marginal variances

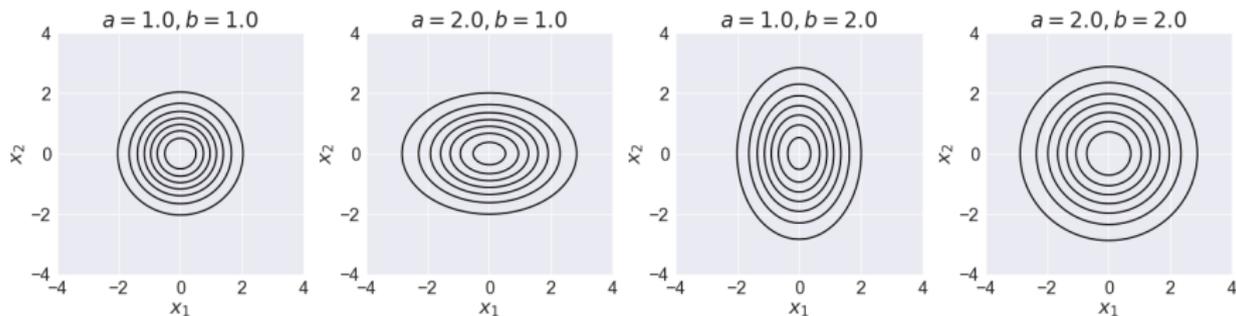
$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix} \quad (1)$$



Interpretation of the covariance matrix - 2D examples

The diagonal of the covariance controls the scaling/marginal variances

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix} \quad (1)$$



Questions to be discussed with your neighbor:

- 1 If $\boldsymbol{\Sigma}$ is diagonal, then x_1 and x_2 are uncorrelated? True or false?
- 2 If $\boldsymbol{\Sigma}$ is diagonal, then x_1 and x_2 are independent? True or false?
- 3 What is the volume (integral) of density?
- 4 Which of the four densities has the highest peak and why?

The density at the mode

- The density is given by

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- The mode (highest density value) is achieved at $\mathbf{x} = \boldsymbol{\mu}$

$$\mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}}$$

- The determinant of the covariance is

$$|\boldsymbol{\Sigma}| = \left| \begin{bmatrix} a & \rho \\ \rho & b \end{bmatrix} \right| = ab - \rho^2 \quad (2)$$

- Therefore

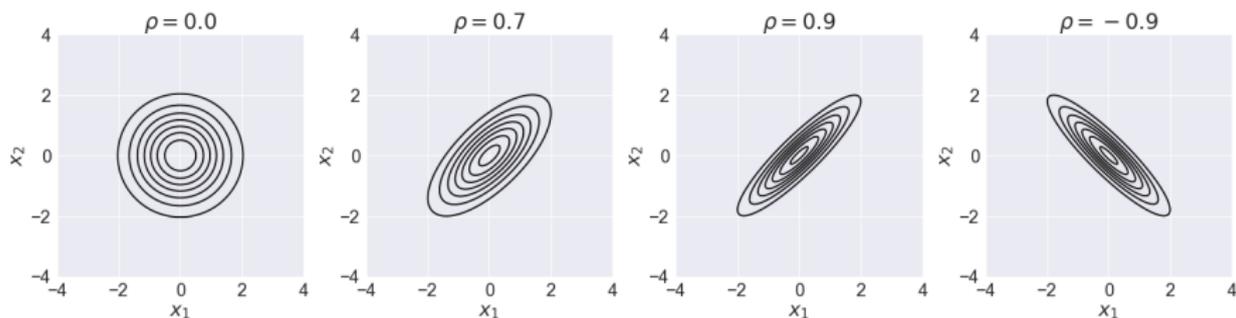
$$\mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} = (2\pi)^{-\frac{D}{2}} \frac{1}{\sqrt{ab - \rho^2}}$$

Interpretation of the covariance matrix

The off-diagonals control the covariances:

$$(\Sigma)_{ij} = \text{cov}(x_i, x_j) = \mathbb{E}[x_i x_j] - \mu_i \mu_j \quad (3)$$

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \quad (4)$$

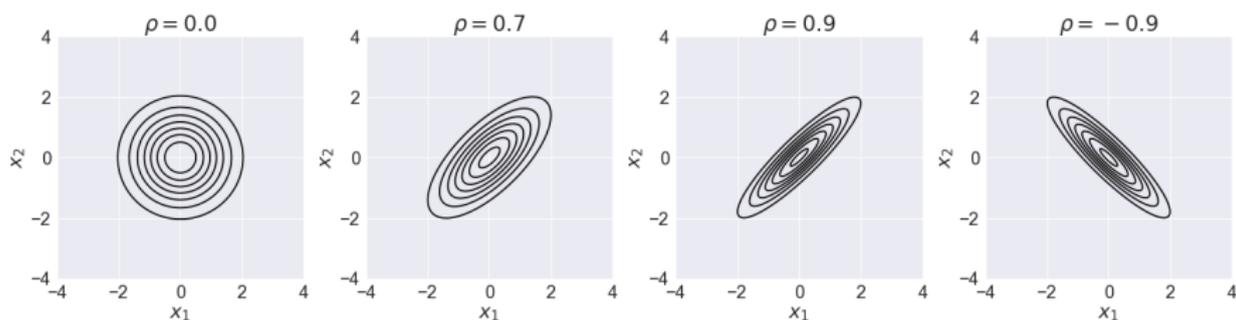


Interpretation of the covariance matrix

The off-diagonals control the covariances:

$$(\Sigma)_{ij} = \text{cov}(x_i, x_j) = \mathbb{E}[x_i x_j] - \mu_i \mu_j \quad (3)$$

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \quad (4)$$



Question:

- Which of the four densities has the highest peak and why?

Interpretation of the covariance matrix

Covariance matrices must be symmetric:

$$(\Sigma)_{ij} = \text{cov}(x_i, x_j) = \text{cov}(x_j, x_i) = (\Sigma)_{ji} \quad (5)$$

Consider the following set of covariance matrices:

$$\Sigma = \begin{bmatrix} a & c \\ c & b \end{bmatrix} \quad (6)$$

c is the covariance between x_1 and x_2 . Can c take any values?

Interpretation of the covariance matrix

Covariance matrices must be symmetric:

$$(\Sigma)_{ij} = \text{cov}(x_i, x_j) = \text{cov}(x_j, x_i) = (\Sigma)_{ji} \quad (5)$$

Consider the following set of covariance matrices:

$$\Sigma = \begin{bmatrix} a & c \\ c & b \end{bmatrix} \quad (6)$$

c is the covariance between x_1 and x_2 . Can c take any values?

$$|\rho| = \left| \frac{c}{\sqrt{a}\sqrt{b}} \right| \leq 1 \quad \Rightarrow \quad |c| \leq \sqrt{a}\sqrt{b} \quad (7)$$

Interpretation of the covariance matrix

Covariance matrices must be symmetric:

$$(\Sigma)_{ij} = \text{cov}(x_i, x_j) = \text{cov}(x_j, x_i) = (\Sigma)_{ji} \quad (5)$$

Consider the following set of covariance matrices:

$$\Sigma = \begin{bmatrix} a & c \\ c & b \end{bmatrix} \quad (6)$$

c is the covariance between x_1 and x_2 . Can c take any values?

$$|\rho| = \left| \frac{c}{\sqrt{a}\sqrt{b}} \right| \leq 1 \quad \Rightarrow \quad |c| \leq \sqrt{a}\sqrt{b} \quad (7)$$

Σ must be positive definite

Interpretation of the covariance matrix

Determine which of the following 5 matrices are valid covariance matrices and match them to the set of samples below.

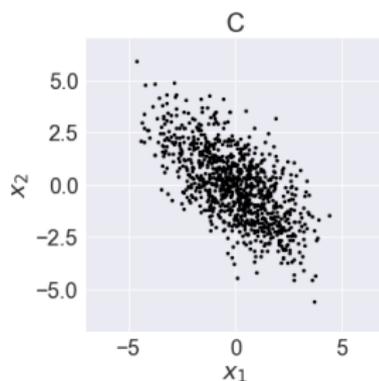
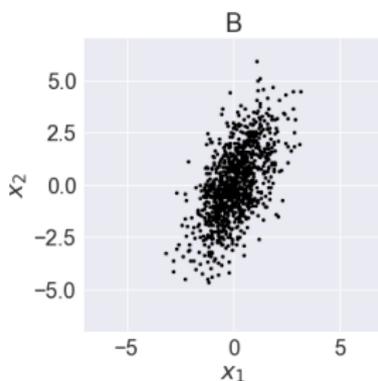
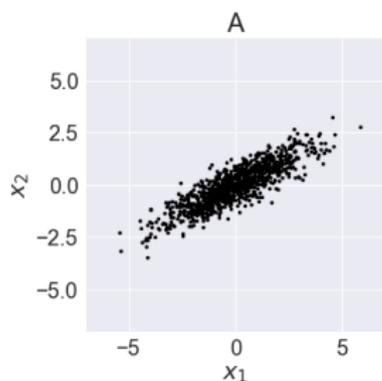
$$\Sigma_1 = \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix}$$

$$\Sigma_2 = \begin{bmatrix} 3 & 2 \\ 1.5 & 3 \end{bmatrix}$$

$$\Sigma_3 = \begin{bmatrix} 1 & 1 \\ 1 & 3 \end{bmatrix}$$

$$\Sigma_4 = \begin{bmatrix} 1 & -2 \\ -2 & 3 \end{bmatrix}$$

$$\Sigma_5 = \begin{bmatrix} 3 & 1.5 \\ 1.5 & 1 \end{bmatrix}$$



Discuss with your neighbor for 3 minutes

The multivariate Gaussian: Basic properties

- Gaussian distributions are closed under addition:

$$\mathbf{x}_1 \sim \mathcal{N}(\mathbf{m}_1, \mathbf{V}_1), \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_2, \mathbf{V}_2) \quad \Rightarrow \quad \mathbf{x}_1 + \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_1 + \mathbf{m}_2, \mathbf{V}_1 + \mathbf{V}_2)$$

The multivariate Gaussian: Basic properties

- Gaussian distributions are closed under addition:

$$\mathbf{x}_1 \sim \mathcal{N}(\mathbf{m}_1, \mathbf{V}_1), \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_2, \mathbf{V}_2) \quad \Rightarrow \quad \mathbf{x}_1 + \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_1 + \mathbf{m}_2, \mathbf{V}_1 + \mathbf{V}_2)$$

- For any finite number of independent variables:

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}_i, \mathbf{V}_i) \quad \Rightarrow \quad \sum_i \mathbf{x}_i \sim \mathcal{N}\left(\sum_i \mathbf{m}_i, \sum_i \mathbf{V}_i\right)$$

The multivariate Gaussian: Basic properties

- Gaussian distributions are closed under addition:

$$\mathbf{x}_1 \sim \mathcal{N}(\mathbf{m}_1, \mathbf{V}_1), \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_2, \mathbf{V}_2) \Rightarrow \mathbf{x}_1 + \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_1 + \mathbf{m}_2, \mathbf{V}_1 + \mathbf{V}_2)$$

- For any finite number of independent variables:

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}_i, \mathbf{V}_i) \Rightarrow \sum_i \mathbf{x}_i \sim \mathcal{N}\left(\sum_i \mathbf{m}_i, \sum_i \mathbf{V}_i\right)$$

- Gaussian distributions are closed under affine transformations:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}), \Rightarrow \mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{Am} + \mathbf{b}, \mathbf{AVA}^T)$$

The multivariate Gaussian: Basic properties

- Gaussian distributions are closed under addition:

$$\mathbf{x}_1 \sim \mathcal{N}(\mathbf{m}_1, \mathbf{V}_1), \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_2, \mathbf{V}_2) \Rightarrow \mathbf{x}_1 + \mathbf{x}_2 \sim \mathcal{N}(\mathbf{m}_1 + \mathbf{m}_2, \mathbf{V}_1 + \mathbf{V}_2)$$

- For any finite number of independent variables:

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}_i, \mathbf{V}_i) \Rightarrow \sum_i \mathbf{x}_i \sim \mathcal{N}\left(\sum_i \mathbf{m}_i, \sum_i \mathbf{V}_i\right)$$

- Gaussian distributions are closed under affine transformations:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}), \Rightarrow \mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{Am} + \mathbf{b}, \mathbf{AVA}^T)$$

- Hence, manipulating Gaussian distributions often boils down to linear algebra

Discuss with your neighbor...

... how to use the following two results

$$\begin{aligned} \mathbf{x}_i \sim \mathcal{N}(\mathbf{m}_i, \mathbf{V}_i) &\Rightarrow \sum_i \mathbf{x}_i \sim \mathcal{N}\left(\sum_i \mathbf{m}_i, \sum_i \mathbf{V}_i\right) \\ \mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) &\Rightarrow \mathbf{A}\mathbf{x} + \mathbf{b} \sim \mathcal{N}\left(\mathbf{A}\mathbf{m} + \mathbf{b}, \mathbf{A}\mathbf{V}\mathbf{A}^T\right), \end{aligned}$$

to calculate the distribution of \mathbf{Y} in the following linear model?

$$\mathbf{Y} = \boldsymbol{\mu} + \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon},$$

where

$$\mathbf{w} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Sampling from the multivariate Gaussian distribution

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad \Rightarrow \quad \mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{Am} + \mathbf{b}, \mathbf{AVA}^T)$$

- Suppose we know how to generate samples from a standardized univariate Gaussian distribution
- How can we use the above result to generate samples from an arbitrary multivariate Gaussian distribution $\mathbf{y} \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$?

Sampling from the multivariate Gaussian distribution

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad \Rightarrow \quad \mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{Am} + \mathbf{b}, \mathbf{AVA}^T)$$

- Suppose we know how to generate samples from a standardized univariate Gaussian distribution
- How can we use the above result to generate samples from an arbitrary multivariate Gaussian distribution $\mathbf{y} \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$?
 - 1 Compute the matrix square root of $\mathbf{V} = \mathbf{LL}^T$
 - 2 Generate a sample of \mathbf{x} such that $x_i \sim \mathcal{N}(0, 1)$, i.e. $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 3 Compute $\mathbf{y} = \mathbf{Lx} + \mathbf{m}$

Sampling from the multivariate Gaussian distribution

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad \Rightarrow \quad \mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{Am} + \mathbf{b}, \mathbf{AVA}^T)$$

- Suppose we know how to generate samples from a standardized univariate Gaussian distribution
- How can we use the above result to generate samples from an arbitrary multivariate Gaussian distribution $\mathbf{y} \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$?
 - 1 Compute the matrix square root of $\mathbf{V} = \mathbf{LL}^T$
 - 2 Generate a sample of \mathbf{x} such that $x_i \sim \mathcal{N}(0, 1)$, i.e. $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 3 Compute $\mathbf{y} = \mathbf{Lx} + \mathbf{m}$
- Why does it work?

$$\mathbf{y} = \mathbf{Lx} + \mathbf{m} \sim \mathcal{N}(\mathbf{L0} + \mathbf{m}, \mathbf{LIL}^T) = \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad (8)$$

The multivariate Gaussian: Marginalization

- Gaussian densities are closed on marginalization
- Let \mathbf{x}_1 and \mathbf{x}_2 be a partitioning of $\mathbf{x} = \mathbf{x}_1 \cup \mathbf{x}_2$, then

$$p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right) \quad (9)$$

The multivariate Gaussian: Marginalization

- Gaussian densities are closed on marginalization
- Let \mathbf{x}_1 and \mathbf{x}_2 be a partitioning of $\mathbf{x} = \mathbf{x}_1 \cup \mathbf{x}_2$, then

$$p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right) \quad (9)$$

then

$$p(\mathbf{x}_1) = \int p(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = \mathcal{N}(\mathbf{x}_1 \mid \mathbf{m}_1, \Sigma_{11}) \quad (10)$$

The multivariate Gaussian: Marginalization

- Gaussian densities are closed on marginalization
- Let \mathbf{x}_1 and \mathbf{x}_2 be a partitioning of $\mathbf{x} = \mathbf{x}_1 \cup \mathbf{x}_2$, then

$$p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right) \quad (9)$$

then

$$p(\mathbf{x}_1) = \int p(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = \mathcal{N}(\mathbf{x}_1 \mid \mathbf{m}_1, \Sigma_{11}) \quad (10)$$

and

$$p(\mathbf{x}_2) = \int p(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 = \mathcal{N}(\mathbf{x}_2 \mid \mathbf{m}_2, \Sigma_{22}) \quad (11)$$

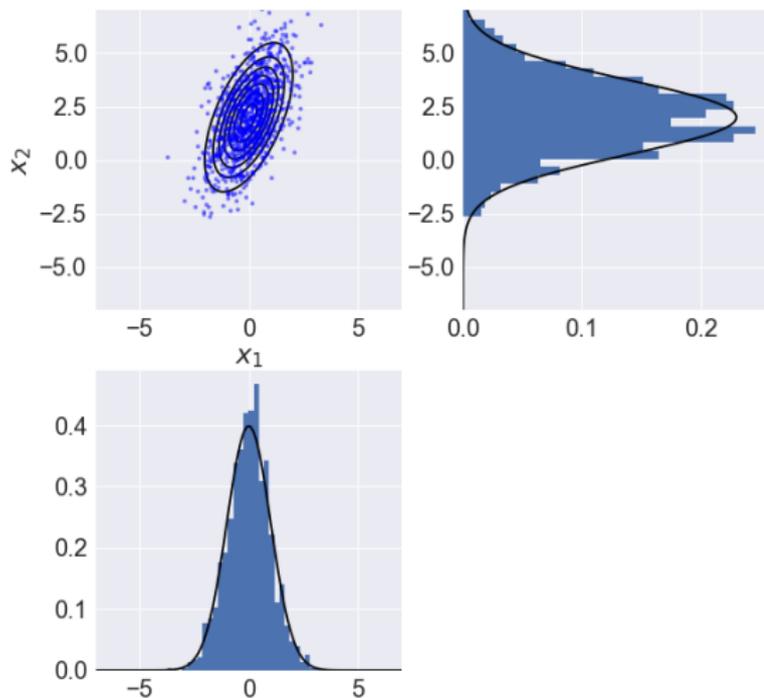
- The same is true for any partitioning

Marginalization example in 2D

$$\mathbf{x} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 1 & 3 \end{bmatrix}\right)$$

$$x_1 \sim \mathcal{N}(0, 1)$$

$$x_2 \sim \mathcal{N}(2, 3)$$



- Gaussian densities are closed under conditioning!
- Recall the definition of conditioning:

$$p(A|B) = \frac{p(A \cap B)}{p(B)}$$

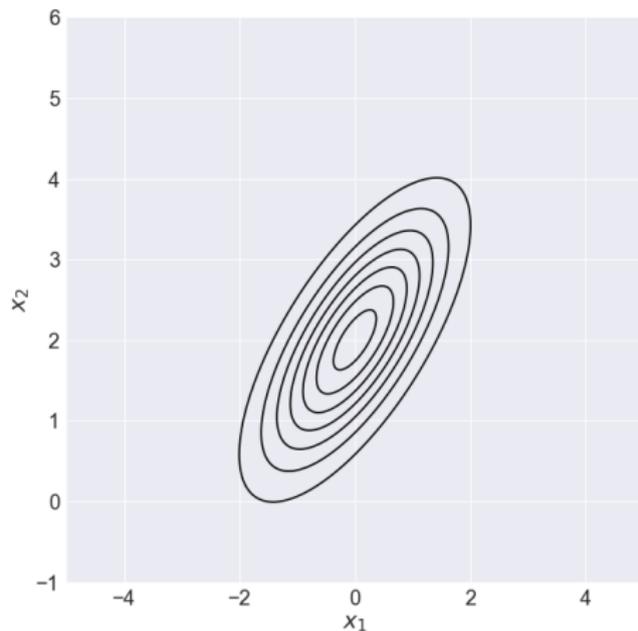
- Let \mathbf{x}_1 and \mathbf{x}_2 be a partitioning of $\mathbf{x} = \mathbf{x}_1 \cup \mathbf{x}_2$, then

$$p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

- The conditional of \mathbf{x}_1 is given \mathbf{x}_2 by:

$$p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N} \left(\mathbf{x}_1 \mid \Sigma_{12} \Sigma_{22}^{-1} [\mathbf{x}_2 - \boldsymbol{\mu}_2] + \mathbf{m}_1, \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right)$$

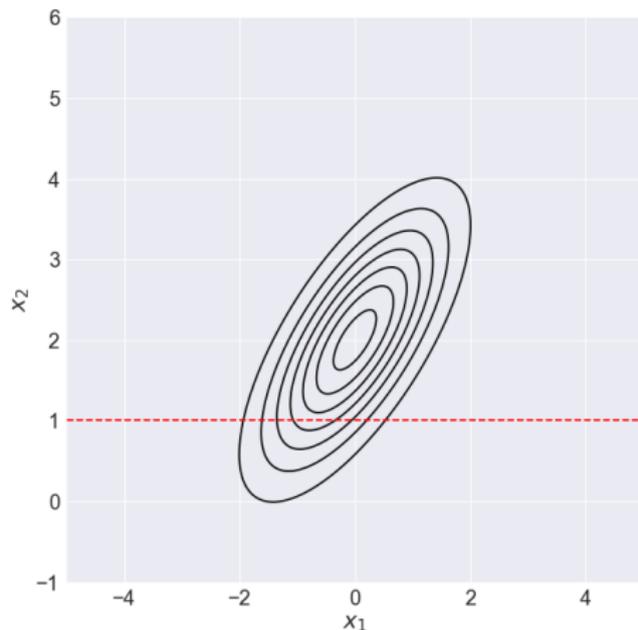
Conditioning example in 2D



- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

Conditioning example in 2D

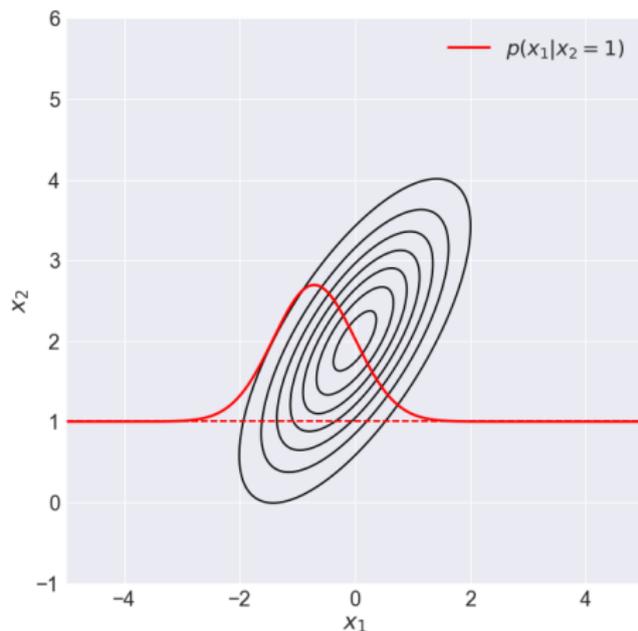


- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

- Assume we observe $x_2 = 1$

Conditioning example in 2D



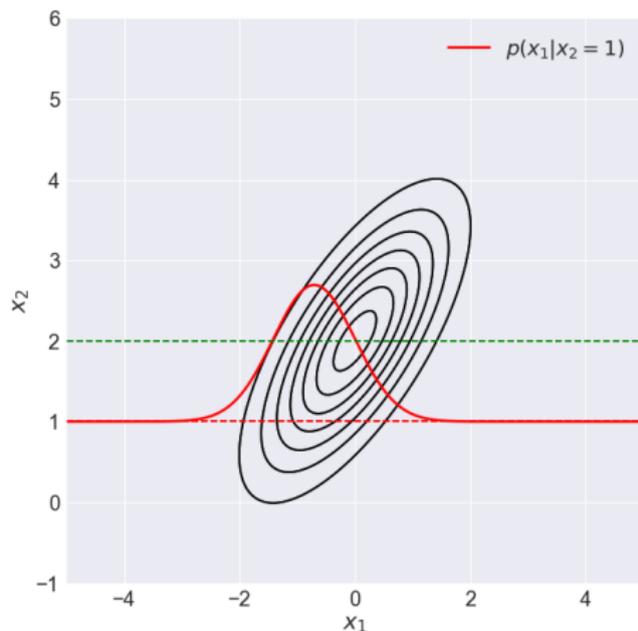
- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

- Assume we observe $x_2 = 1$
- The conditional distribution

$$p(x_1|x_2) = \mathcal{N}\left(x_1 \mid -\frac{\sqrt{2}}{2}, \frac{1}{2}\right)$$

Conditioning example in 2D

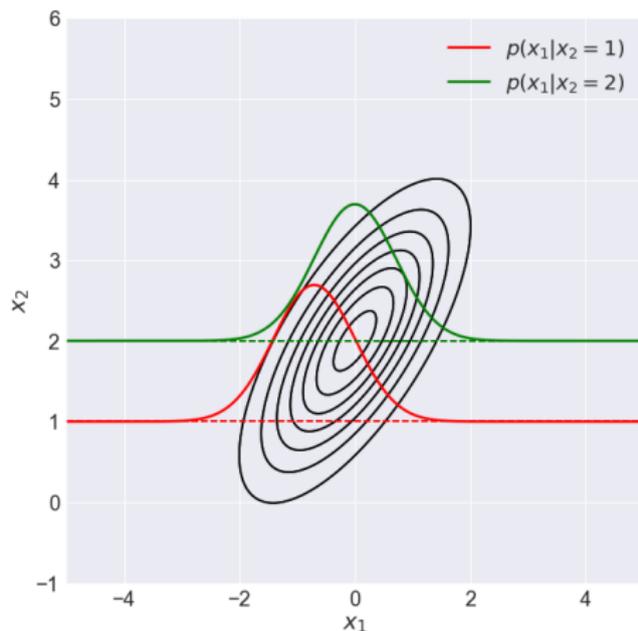


- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

- Assume we observe $x_2 = 2$

Conditioning example in 2D



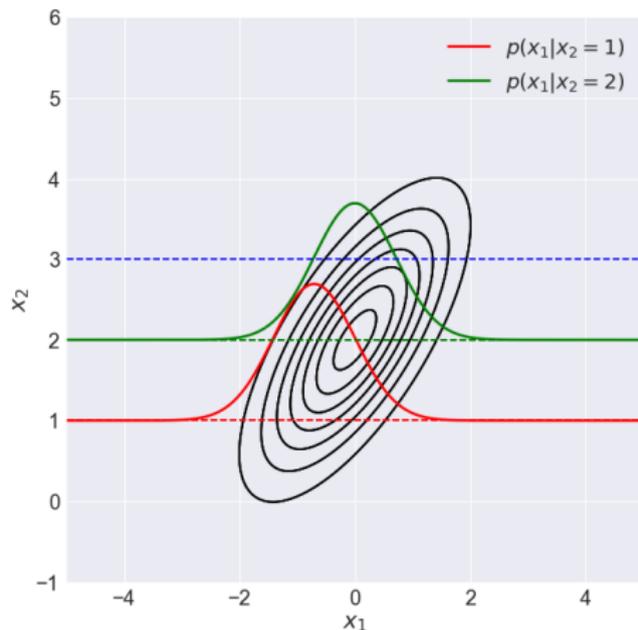
- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

- Assume we observe $x_2 = 2$
- The conditional distribution

$$p(x_1|x_2) = \mathcal{N}\left(x_1|0, \frac{1}{2}\right)$$

Conditioning example in 2D

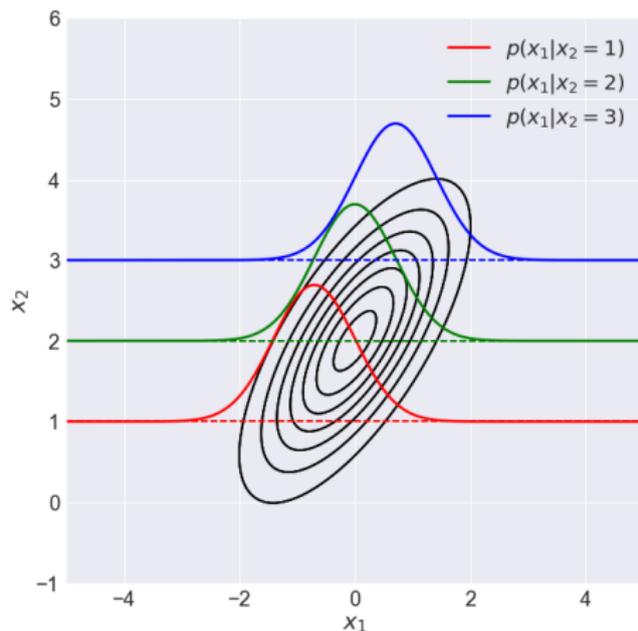


- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

- Assume we observe $x_2 = 3$

Conditioning example in 2D



- 2D example

$$\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

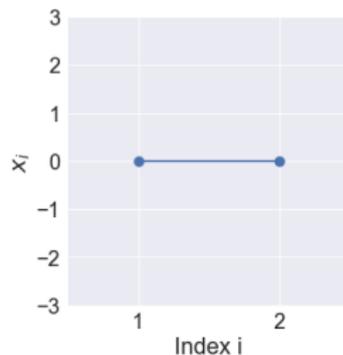
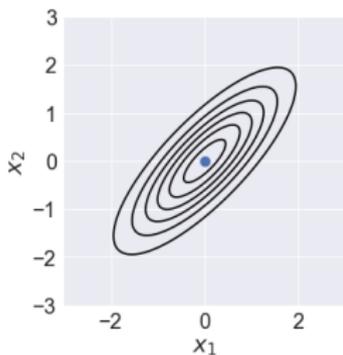
- Assume we observe $x_2 = 3$
- The conditional distribution

$$p(x_1|x_2) = \mathcal{N}\left(x_1 \mid \frac{\sqrt{2}}{2}, \frac{1}{2}\right)$$

Visualizing samples in higher dimensions

- Visualizations in 2D

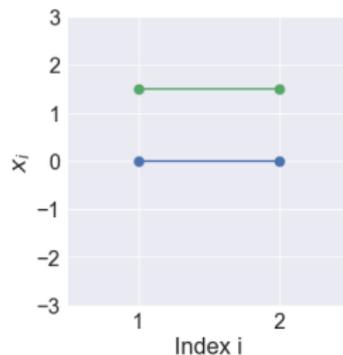
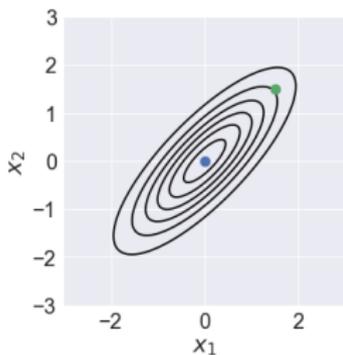
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

- Visualizations in 2D

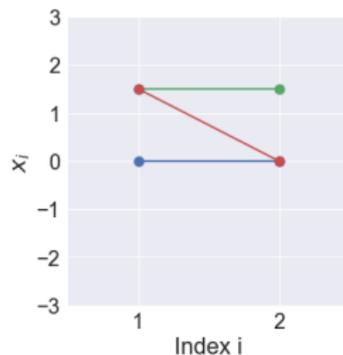
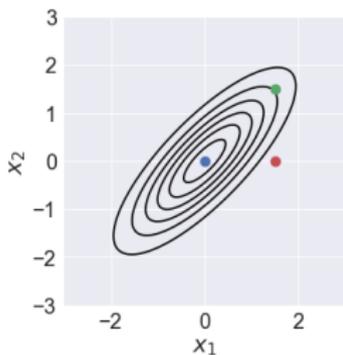
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

- Visualizations in 2D

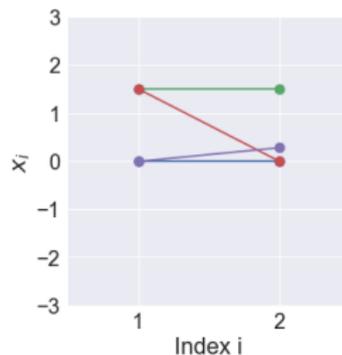
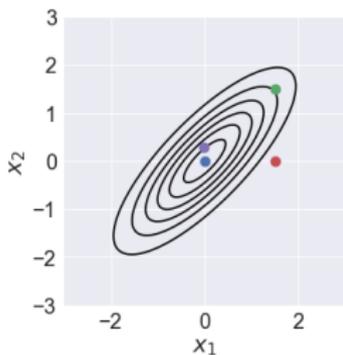
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

- Visualizations in 2D

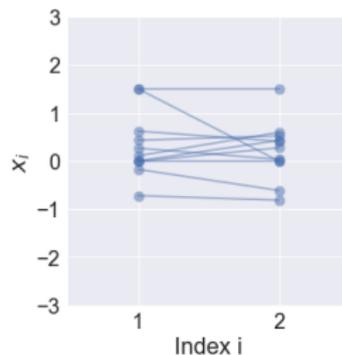
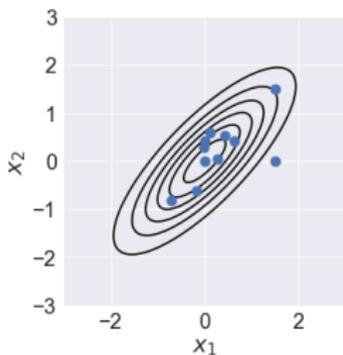
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

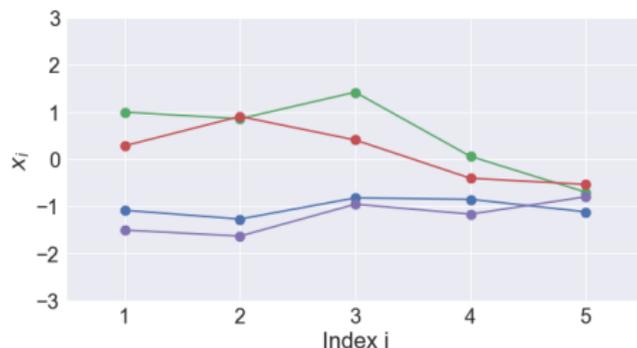
- Visualizations in 2D

$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

- Visualizations in 5D

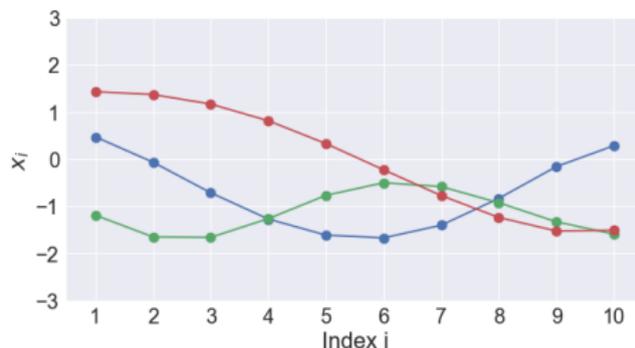


$$\Sigma = \begin{bmatrix} 1 & 0.8^1 & 0.8^2 & 0.8^3 & 0.8^4 \\ 0.8^1 & 1 & 0.8^1 & 0.8^2 & 0.8^3 \\ 0.8^2 & 0.8^1 & 1 & 0.8^1 & 0.8^2 \\ 0.8^3 & 0.8^2 & 0.8^1 & 1 & 0.8^1 \\ 0.8^4 & 0.8^3 & 0.8^2 & 0.8^1 & 1 \end{bmatrix}$$

Visualizing samples in higher dimensions

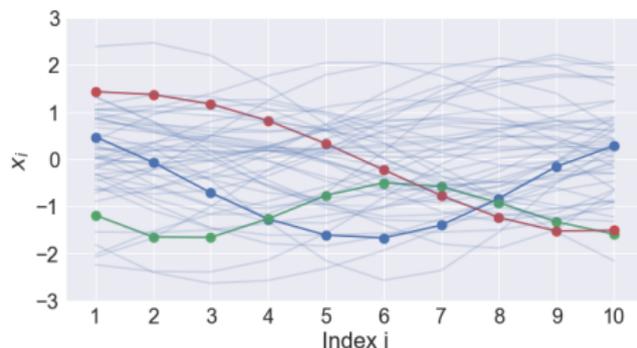
- Visualizations in 10D

$$\Sigma = \begin{bmatrix} 1 & 0.8^1 & 0.8^2 & \dots & 0.8^9 \\ 0.8^1 & 1 & 0.8^1 & & \vdots \\ 0.8^2 & 0.8^1 & 1 & & \vdots \\ \vdots & & & \ddots & \vdots \\ 0.8^9 & \dots & \dots & \dots & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

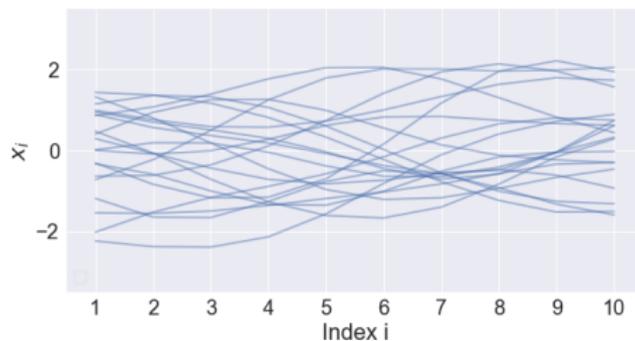
- Visualizations in 10D



$$\Sigma = \begin{bmatrix} 1 & 0.8^1 & 0.8^2 & \dots & 0.8^9 \\ 0.8^1 & 1 & 0.8^1 & & \vdots \\ 0.8^2 & 0.8^1 & 1 & & \vdots \\ \vdots & & & \ddots & \vdots \\ 0.8^9 & \dots & \dots & \dots & 1 \end{bmatrix}$$

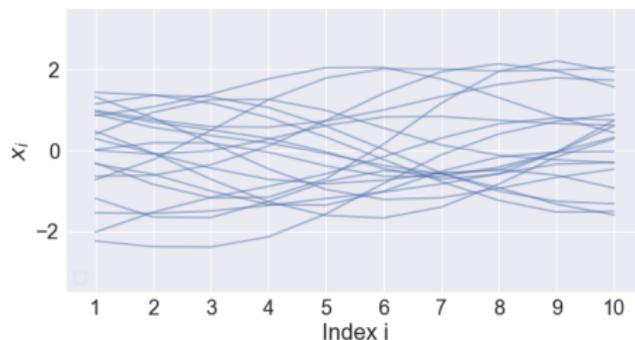
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$



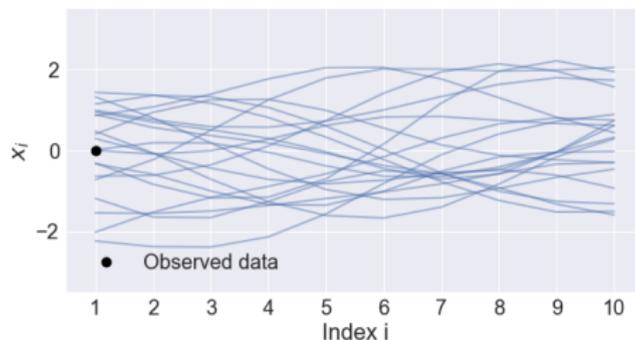
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$



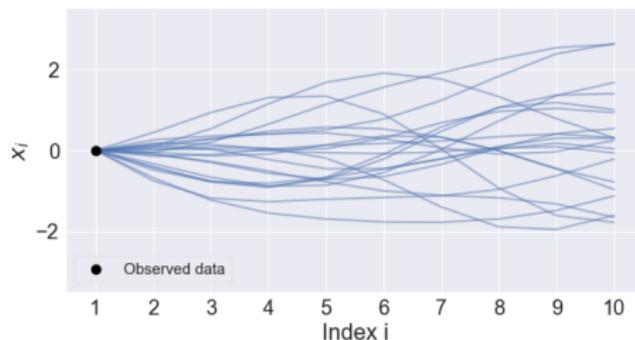
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$
- We now observe $x_1 = 0$



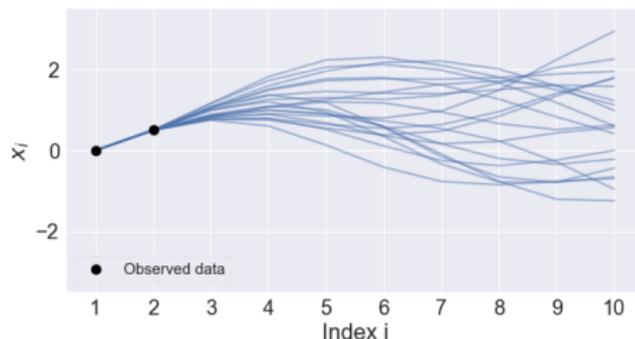
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$
- We now observe $x_1 = 0$
- Let's sample from the conditional distribution $p(\mathbf{x}_{2:10} | x_1 = 0)$



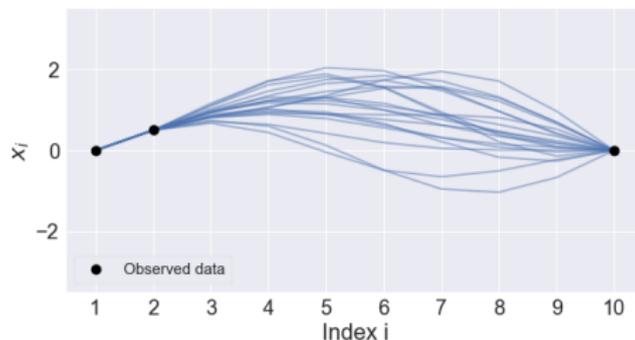
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$
- We now observe $x_1 = 0$
- Let's sample from the conditional distribution $p(\mathbf{x}_{2:10} | x_1 = 0)$



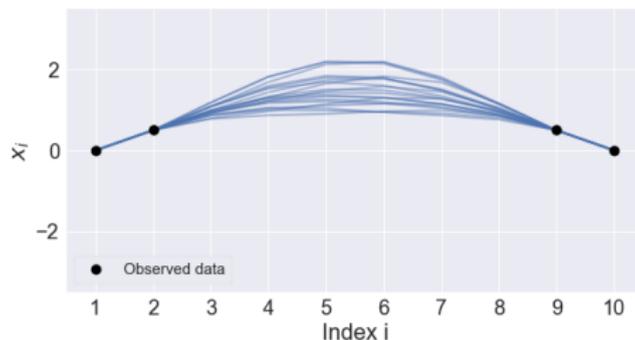
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$
- We now observe $x_1 = 0$
- Let's sample from the conditional distribution $p(\mathbf{x}_{2:10}|x_1 = 0)$



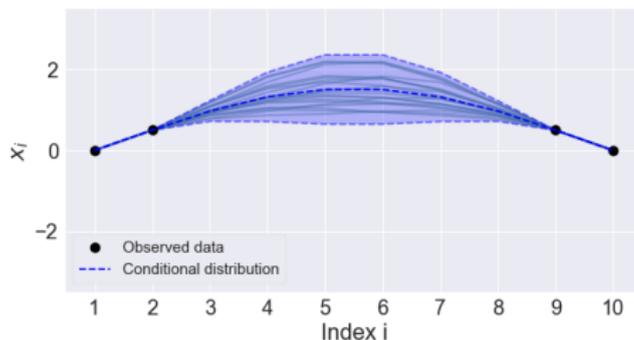
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$
- We now observe $x_1 = 0$
- Let's sample from the conditional distribution $p(\mathbf{x}_{2:10} | x_1 = 0)$



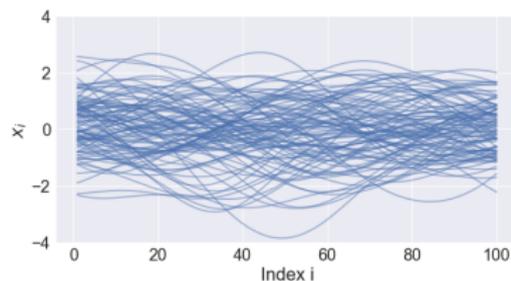
Back to conditioning

- So far, we have seen samples from the distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$
- We can also write $p(\mathbf{x}) = p(x_1, \mathbf{x}_{2:10})$
- We now observe $x_1 = 0$
- Let's sample from the conditional distribution $p(\mathbf{x}_{2:10}|x_1 = 0)$



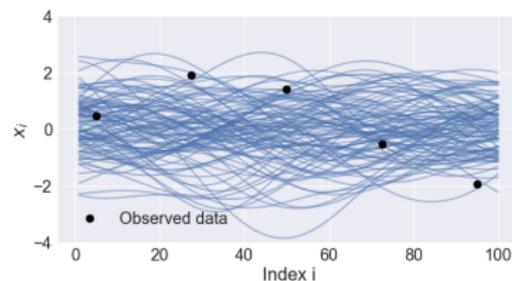
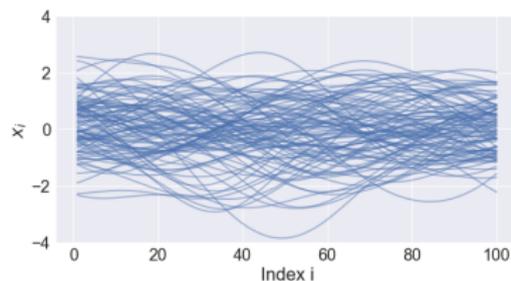
Back to conditioning II

- Let's now consider a case with $\mathbf{x} \in \mathbb{R}^{100}$ dimensions with 5 observations



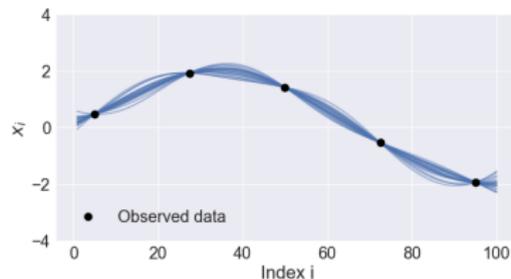
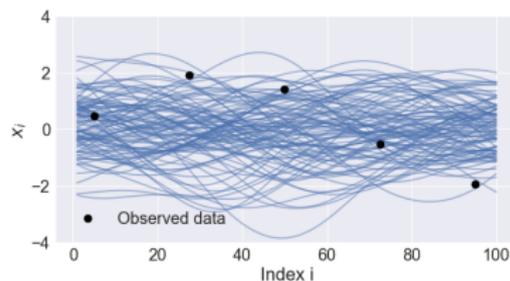
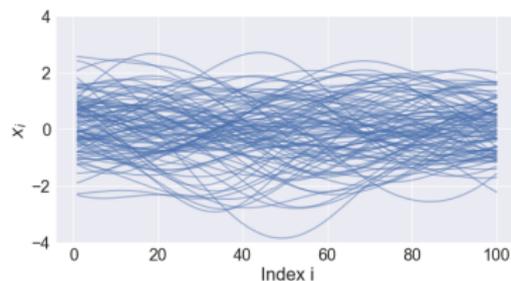
Back to conditioning II

- Let's now consider a case with $\mathbf{x} \in \mathbb{R}^{100}$ dimensions with 5 observations



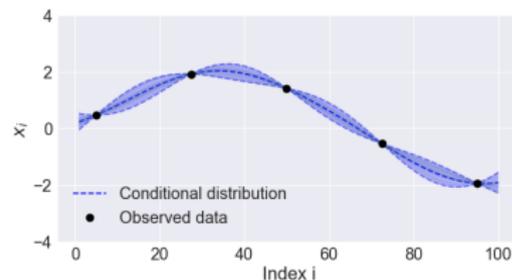
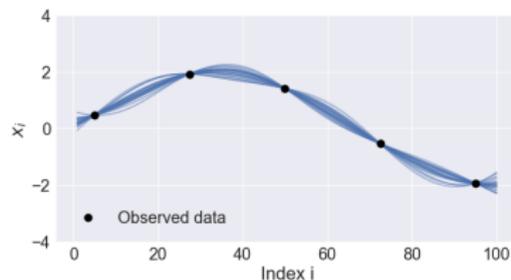
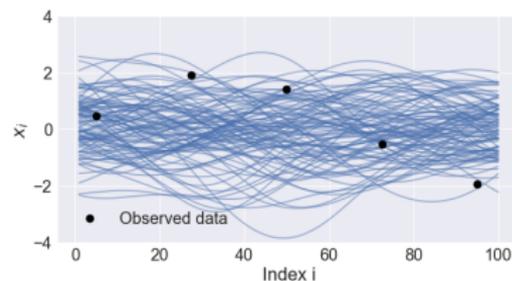
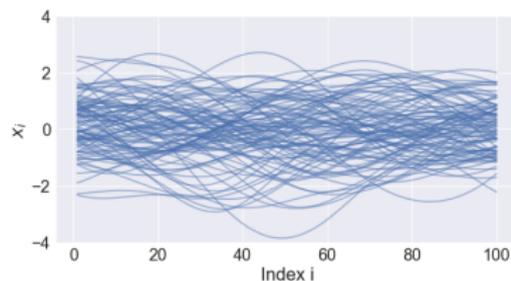
Back to conditioning II

- Let's now consider a case with $\mathbf{x} \in \mathbb{R}^{100}$ dimensions with 5 observations



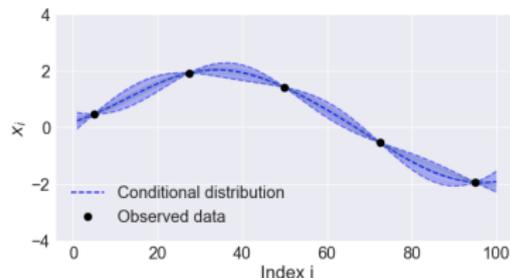
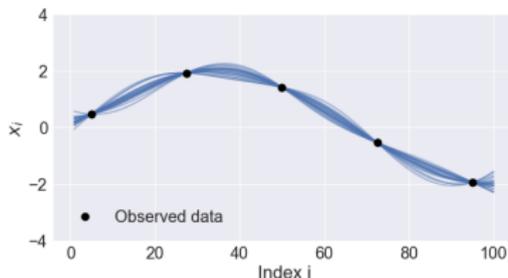
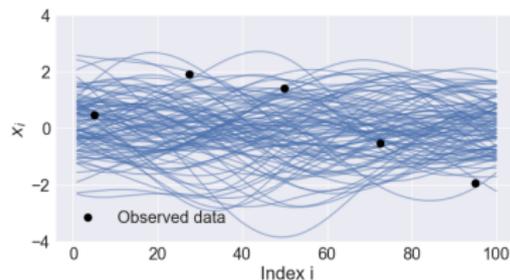
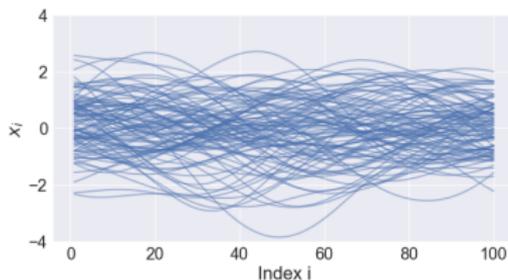
Back to conditioning II

- Let's now consider a case with $\mathbf{x} \in \mathbb{R}^{100}$ dimensions with 5 observations



Back to conditioning II

- Let's now consider a case with $\mathbf{x} \in \mathbb{R}^{100}$ dimensions with 5 observations



- Informally: We can think functions as vectors with infinite dimensions
- Using conditioning in Gaussian distributions, we can do non-linear regression!

The end of today's lecture

- Next time
 - We will introduce Gaussian processes more formally
 - Read Chapter 1 & 2 in Gaussian processes for Machine Learning by Carl Rasmussen (<http://www.gaussianprocess.org/gpml>)
- First assignment
 - Warm up for Gaussian processes
 - Reviews the basics of Bayesian inference
 - Reviews the multivariate Gaussian density
 - Must be handed in through MyCourses