# Probabilistic Inference and Learning
## Lecture 06
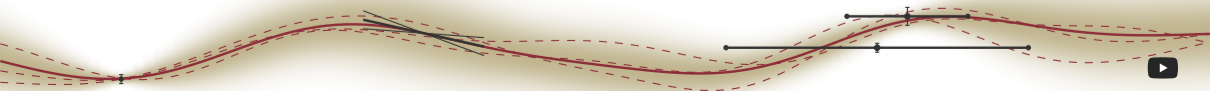## Gaussian Probability Distributions

Philipp Hennig

04 May 2021

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING

| # | date | content | Ex | # | date | content | Ex |
|---|------|---------|-----|---|------|---------|-----|
| 1 | 20.04. | Introduction | 1 | 14 | 09.06. | Logistic Regression | 8 |
| 2 | 21.04. | Reasoning under Uncertainty | | 15 | 15.06. | Exponential Families | |
| 3 | 27.04. | Continuous Variables | 2 | 16 | 16.06. | Graphical Models | 9 |
| 4 | 28.04. | Monte Carlo | | 17 | 22.06. | Factor Graphs | |
| 5 | 04.05. | Markov Chain Monte Carlo | 3 | 18 | 23.06. | The Sum-Product Algorithm | 10 |
| 6 | 05.05. | **Gaussian Distributions** | | 19 | 29.06. | Example: Topic Models | |
| 7 | 11.05. | Parametric Regression | 4 | 20 | 30.06. | Mixture Models | 11 |
| 8 | 12.05. | Understanding Deep Learning | | 21 | 06.07. | EM | |
| 9 | 18.05. | Gaussian Processes | 5 | 22 | 07.07. | Variational Inference | 12 |
| 10 | 19.05. | An Example for GP Regression | | 23 | 13.07. | Example: Topic Models | |
| 11 | 25.05. | Understanding Kernels | 6 | 24 | 14.07. | Example: Inferring Topics | 13 |
| 12 | 26.05. | Gauss-Markov Models | | 25 | 20.07. | Example: Kernel Topic Models | |
| 13 | 08.06. | GP Classification | 7 | 26 | 21.07. | Revision | |

DEUTSCHE BUNDESBANK

10

ZEHN DEUTSCHE MARK

1777 – 1855 Carl Friedr. Gauß

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}}\, e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

10

GN4480100S8

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} =: \mathcal{N}(x; \mu, \sigma^2)$$

$\mu$  the **mean** of $x$

$\sigma^2$  the **variance** of $x$

$\sigma$  the **standard deviation** of $x$

### Definition

$$\mathcal{N}(x; \mu, \sigma^2) =: \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \qquad \text{with } \mu, \sigma \in \mathbb{R}$$

will be called the **Gaussian** or **normal distribution** of $x$. We call $x$ the **argument** or **variable**, $\mu, \sigma^2$ the **parameters**. We write $x \sim \mathcal{N}(\mu, \sigma^2)$ to say that the variable $x$ is distributed with pdf $\mathcal{N}(x; \mu, \sigma^2)$.

▶ $\int \mathcal{N}(x; \mu, \sigma^2) \, dx = 1$ and $\mathcal{N}(x; \mu, \sigma^2) > 0 \, \forall x \in \mathbb{R}$. So $\mathcal{N}$ is the density of a probability measure.

▶ Symmetry in $x$ and $\mu$: $\mathcal{N}(x; \mu, \sigma^2) = \mathcal{N}(\mu; x, \sigma^2)$

▶ An **exponential of a quadratic polynomial** of the **natural parameters** $(a, \eta, \tau)$:

$$\mathcal{N}(x; \mu, \sigma^2) = \exp\left(a + \eta x - \frac{1}{2}\tau x^2\right) \quad \text{with} \quad \tau = \sigma^{-2} \text{ ("precision")}, \eta = \sigma^{-2}\mu$$

$$a = -\frac{1}{2}\left(\log(2\pi) - \log \lambda^2 + \lambda^2 \eta^2\right)$$

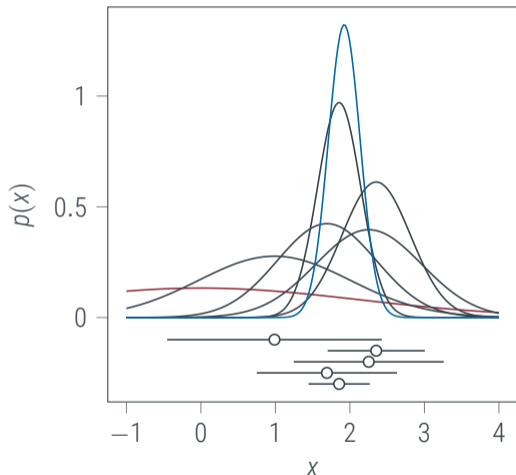The Gaussian is its own **conjugate prior**.



Let

$$p(x) = \mathcal{N}(x; \mu, \sigma^2)$$
$$p(y \mid x) = \mathcal{N}(y; x, \nu^2)$$

Then

$$p(x \mid y) = \frac{p(x)p(y \mid x)}{\int p(x)p(y \mid x)\, dx}$$
$$= \mathcal{N}(x; m, s^2), \text{ with}$$
$$s^2 := \frac{1}{\sigma^{-2} + \nu^{-2}}$$
$$m := \frac{\sigma^{-2}\mu + \nu^{-2}y}{\sigma^{-2} + \nu^{-2}}$$

$$p(x) = \mathcal{N}(x; \mu, \sigma^2)$$

$$p(\boldsymbol{y} \mid x) = \prod_{i=1}^{N} \mathcal{N}(y_i; x, \nu_i^2)$$

$$p(x \mid y) = \frac{p(x)p(\boldsymbol{y} \mid x)}{\int p(x)p(\boldsymbol{y} \mid x)\, dx}$$

$$= \mathcal{N}(x; m, s^2), \text{ with}$$

$$s^{-2} := \sigma^{-2} + \sum_{i=1}^{N} \nu_i^{-2}$$

$$s^{-2}m := \sigma^{-2}\mu + \sum_{i=1}^{N} \nu_i^{-2} y_i$$

If $\sigma^{-2} \to 0$, $\nu_i = \nu \, \forall i$, then $m$ is the **arithmetic mean**.

# The Method of Least Squares
The Gaussian distribution is the unique choice yielding a mean that is the mean of measurements.

[image: C.A. Jensen, 1840]

### Definition (multivariate Gaussian distribution)

$$\mathcal{N}(x; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\boldsymbol{\mu})^\mathsf{T}\Sigma^{-1}(x-\boldsymbol{\mu})\right) \quad x, \boldsymbol{\mu} \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n\times n}, \text{spd}.$$

$\Sigma$ must be **symmetric positive definite**.

Definition (multivariate Gaussian distribution)

$$\mathcal{N}(x; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \boldsymbol{\mu})^\mathsf{T}\Sigma^{-1}(x - \boldsymbol{\mu})\right) \quad x, \boldsymbol{\mu} \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}, \text{spd}.$$

$\Sigma$ must be **symmetric positive definite**.

Definition (symmetric positive definite matrix)

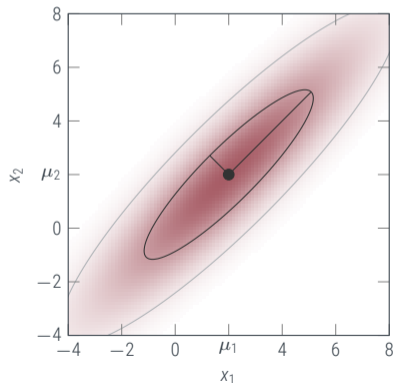A matrix $A \in \mathbb{R}^{n \times n}$ is called **symmetric positive (semi-) definite** if $A = A^\mathsf{T}$, and

$$v^\mathsf{T} A v \geq 0 \ \forall v \in \mathbb{R}^n.$$

Equivalent statement: All eigenvalues of the symmetric matrix $A$ are non-negative.

$$\mathcal{N}(x; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\boldsymbol{\mu})^{\mathsf{T}}\Sigma^{-1}(x-\boldsymbol{\mu})\right) \quad x, \boldsymbol{\mu} \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n\times n}, \text{spd.}$$



▶ $\int \mathcal{N}(x; \boldsymbol{\mu}, \Sigma) = 1$ and $\mathcal{N}(x; \boldsymbol{\mu}, \Sigma) > 0 \ \forall x \in \mathbb{R}^n$.

▶ Symmetry in $x$ and $\boldsymbol{\mu}$: $\mathcal{N}(x; \boldsymbol{\mu}, \Sigma) = \mathcal{N}(\boldsymbol{\mu}; x, \Sigma)$

▶ An exponential of a quadratic polynomial:

$$\mathcal{N}(x; \mu, \Sigma) = \exp\left(a + \eta^{\mathsf{T}}x - \frac{1}{2}x^{\mathsf{T}}\Lambda x\right) \quad (1)$$

$$= \exp\left(a + \eta^{\mathsf{T}}x - \frac{1}{2}\operatorname{tr}(xx^{\mathsf{T}}\Lambda)\right) \quad (2)$$

with the **natural parameters** $\Lambda = \Sigma^{-1}$ (precision matrix), $\eta = \Lambda\mu$, and the **sufficient statistics** $x, xx^{\mathsf{T}}$.

Closure under Multiplication



### To multiply Gaussians, add the natural parameters

$$\mathcal{N}(x; a, A)\mathcal{N}(x; b, B) = \mathcal{N}(x; c, C)Z$$
$$C = (A^{-1} + B^{-1})^{-1}$$
$$c = C(A^{-1}a + B^{-1}b)$$
$$Z = \mathcal{N}(a; b, A + B)$$

Note similarity to univariate case.

### To multiply Gaussians, add the natural parameters

$$\mathcal{N}(x; a, A)\mathcal{N}(x; b, B) = \mathcal{N}(x; c, C)Z$$
$$C = (A^{-1} + B^{-1})^{-1}$$
$$c = C(A^{-1}a + B^{-1}b)$$
$$Z = \mathcal{N}(a; b, A + B)$$

Note similarity to univariate case.

**To multiply Gaussians, add the natural parameters**
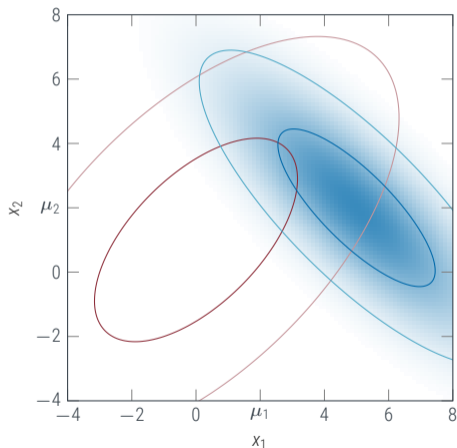
$$\mathcal{N}(x; a, A)\mathcal{N}(x; b, B) = \mathcal{N}(x; c, C)Z$$
$$C = (A^{-1} + B^{-1})^{-1}$$
$$c = C(A^{-1}a + B^{-1}b)$$
$$Z = \mathcal{N}(a; b, A + B)$$

Note similarity to univariate case.

To linearly project a Gaussian variable,
project the parameters

$$p(z) = \mathcal{N}(z; \mu, \Sigma)$$
$$\Rightarrow \quad p(Az) = \mathcal{N}(Az, A\mu, A\Sigma A^\intercal)$$

$$p(z) = \mathcal{N}(z; \mu, \Sigma) \quad \Rightarrow \quad p(Az) = \mathcal{N}(Az, A\mu, A\Sigma A^{\mathsf{T}})$$

$$\text{choose } A = \begin{pmatrix} 1 & 0 \end{pmatrix}$$

$$\int \mathcal{N}\left[ \begin{pmatrix} x \\ y \end{pmatrix}; \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right] dy = \mathcal{N}(x; \mu_x, \Sigma_{xx})$$

▶ this is the sum rule

$$\int p(x, y)\, dy = \int p(y \mid x) p(x)\, dy = p(x)$$

▶ so every finite-dim Gaussian is a marginal of infinitely many more

$$p(x \mid Ax = y) = \frac{p(x, y)}{p(y)}$$
$$= \mathcal{N}\big(x; \mu + \Sigma A^{\mathsf{T}}(A\Sigma A^{\mathsf{T}})^{-1}(y - A\mu),$$
$$\Sigma - \Sigma A^{\mathsf{T}}(A\Sigma A^{\mathsf{T}})^{-1}A\Sigma\big)$$

▶ this is the product rule
▶ so Gaussians are closed under the rules of probability

## Theorem

$$\text{If } p(x) = \mathcal{N}(x; \mu, \Sigma)$$
$$\text{and } p(y \mid x) = \mathcal{N}(y; Ax + b, \Lambda),$$
$$\text{then } p(y) = \mathcal{N}(y; A\mu + b, \Lambda + A\Sigma A^\mathsf{T})$$
$$\text{and } p(x \mid y) = \mathcal{N}(x; \mu + \underbrace{\Sigma A^\mathsf{T}(A\Sigma A^\mathsf{T} + \Lambda)^{-1}}_{gain} \underbrace{(y - (A\mu + b))}_{residual}, \Sigma - \underbrace{\Sigma A^\mathsf{T}(A\Sigma A^\mathsf{T} + \Lambda)^{-1}A\Sigma}_{Gram\ matrix})$$
$$= \mathcal{N}(x; \underbrace{(\Sigma^{-1} + A^\mathsf{T}\Lambda^{-1}A)^{-1}}_{precision\ matrix}(A^\mathsf{T}\Lambda^{-1}(y - b) + \Sigma^{-1}\mu), \underbrace{(\Sigma^{-1} + A^\mathsf{T}\Lambda^{-1}A)^{-1}}_{precision\ matrix})$$

$$A = \begin{bmatrix} P & Q \\ R & S \end{bmatrix} \quad M := (S - RP^{-1}Q)^{-1}$$

$$A^{-1} = \begin{bmatrix} P^{-1} + P^{-1}QMRP^{-1} & -P^{-1}QM \\ -MRP^{-1} & M \end{bmatrix}$$

$$(Z + UWV^{\mathsf{T}})^{-1} = Z^{-1} - Z^{-1}U(W^{-1} + V^{\mathsf{T}}Z^{-1}U)^{-1}V^{\mathsf{T}}Z^{-1}$$

$$|Z + UWV^{\mathsf{T}}| = |Z| \cdot |W| \cdot |W^{-1} + V^{\mathsf{T}}Z^{-1}U|$$



Issai Schur (1875–1941)

▶ products of Gaussians are Gaussians

$$\mathcal{N}(x; a, A)\mathcal{N}(x; b, B)$$
$$= \mathcal{N}(x; c, C)\mathcal{N}(a; b, A + B)$$
$$C := (A^{-1} + B^{-1})^{-1} \quad c := C(A^{-1}a + B^{-1}b)$$

▶ linear projections of Gaussians are Gaussians

$$p(z) = \mathcal{N}(z; \mu, \Sigma)$$
$$\Rightarrow \quad p(Az) = \mathcal{N}(Az, A\mu, A\Sigma A^{\mathsf{T}})$$

▶ marginals of Gaussians are Gaussians

$$\int \mathcal{N}\left[\begin{pmatrix} x \\ y \end{pmatrix}; \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}\right] dy = \mathcal{N}(x; \mu_x, \Sigma_{xx})$$
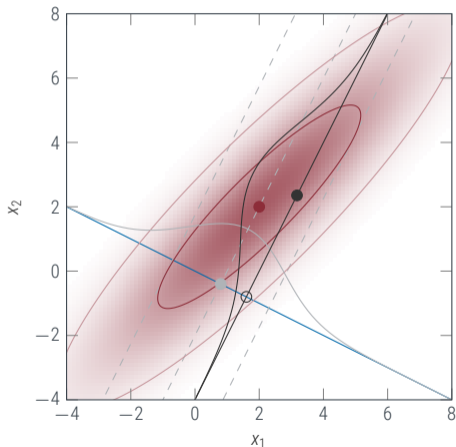
▶ (linear) conditionals of Gaussians are Gaussians

$$p(x \mid y) = \frac{p(x, y)}{p(y)}$$
$$= \mathcal{N}\left(x; \mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}\right)$$

### Bayesian inference becomes linear algebra

$$\text{If } p(x) = \mathcal{N}(x; \mu, \Sigma) \qquad \text{and} \qquad p(y \mid x) = \mathcal{N}(y; A^{\mathsf{T}}x + b, \Lambda), \text{ then}$$
$$p(B^{\mathsf{T}}x + c \mid y) = \mathcal{N}[B^{\mathsf{T}}x + c; B^{\mathsf{T}}\mu + c + B^{\mathsf{T}}\Sigma A(A^{\mathsf{T}}\Sigma A + \Lambda)^{-1}(y - A^{\mathsf{T}}\mu - b), B^{\mathsf{T}}\Sigma B - B^{\mathsf{T}}\Sigma A(A^{\mathsf{T}}\Sigma A + \Lambda)^{-1}A^{\mathsf{T}}\Sigma B]$$

# Example 1: Conditional Independence, Marginal Correlation

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Bayesian Inference with Gaussians

[DJC MacKay, *The humble Gaussian distribution*, 2006]

temperature outside

$x_2$

$x_1$

$x_3$

temperature
in building 1

temperature
in building 2

$$x_2 = \nu_2 \qquad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$

$$x_1 = w_1 x_2 + \nu_1 \quad p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$$

$$x_3 = w_3 x_2 + \nu_3 \quad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[DJC MacKay, *The humble Gaussian distribution*, 2006]

temperature outside



$x_2$

$x_1$            $x_3$

temperature
in building 1

temperature
in building 2

$$x_2 = \nu_2 \qquad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$
$$x_1 = w_1 x_2 + \nu_1 \quad p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$$
$$x_3 = w_3 x_2 + \nu_3 \quad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

$$p(\boldsymbol{\nu}) = \mathcal{N}(\boldsymbol{\nu}; \boldsymbol{\mu}, \mathrm{diag}(\boldsymbol{\sigma}^2))$$

$$A = \begin{bmatrix} 1 & w_1 & 0 \\ 0 & 1 & 0 \\ 0 & w_3 & 1 \end{bmatrix} \qquad \Longrightarrow$$

$$p(\boldsymbol{x} = A\nu) = \mathcal{N}\left(\boldsymbol{x}; \underbrace{A\boldsymbol{\mu}}_{=:m}, \underbrace{\begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1\sigma_2^2 & w_1 w_3 \sigma_2^2 \\ & \sigma_2^2 & w_3\sigma_2^2 \\ & & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix}}_{=:\Sigma}\right)$$

# Example 1: Conditional Independence, Marginal Correlation

UNIVERSITÄT
TÜBINGEN
EBERHARD KARLS

Bayesian Inference with Gaussians
[DJC MacKay, *The humble Gaussian distribution*, 2006]

temperature outside



$x_2$

$x_1$                    $x_3$

temperature          temperature
in building 1        in building 2
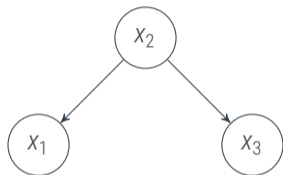
$$x_2 = \nu_2 \qquad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$
$$x_1 = w_1 x_2 + \nu_1 \quad p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$$
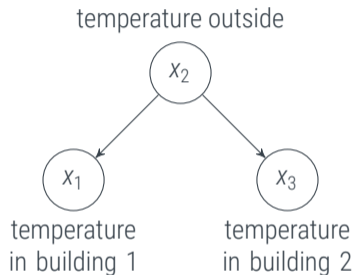$$x_3 = w_3 x_2 + \nu_3 \quad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

$$p(\boldsymbol{\nu}) = \mathcal{N}(\boldsymbol{\nu}; \boldsymbol{\mu}, \mathrm{diag}(\boldsymbol{\sigma}^2))$$

$$A = \begin{bmatrix} 1 & w_1 & 0 \\ 0 & 1 & 0 \\ 0 & w_3 & 1 \end{bmatrix} \quad \implies$$

$$p(\boldsymbol{x} = A\nu) = \mathcal{N}\left(\boldsymbol{x}; \underbrace{A\boldsymbol{\mu}}_{=:m}, \underbrace{\begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1\sigma_2^2 & w_1 w_3 \sigma_2^2 \\ \sigma_2^2 & w_3\sigma_2^2 \\ & & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix}}_{=:\Sigma}\right)$$

From graph: $x_1 \perp\!\!\!\perp x_3 \mid x_2$. Where can we see this in the pdf?

# Example 1: Conditional Independence, Marginal Correlation

A zero in the precision matrix means independence conditional on everything else    [DJC MacKay, *The humble Gaussian distribution*, 2006]

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

$$x_2 = \nu_2 \qquad\qquad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$
$$x_1 = w_1 x_2 + \nu_1 \qquad p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$$
$$x_3 = w_3 x_2 + \nu_3 \qquad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

to simplify exposition, set $\boldsymbol{\mu} = 0$.

$$
\begin{aligned}
p(x_1, x_2, x_3) &= p(x_2) \cdot p(x_1 \mid x_2) \cdot p(x_3 \mid x_2) \\
&= \frac{1}{Z_1 Z_2 Z_3} \exp\left(-\frac{1}{2}\left(\frac{x_2^2}{\sigma_2^2} + \frac{(x_1 - w_1 x_2)^2}{\sigma_1^2} + \frac{(x_3 - w_3 x_2)^2}{\sigma_3^2}\right)\right) \\
&= \frac{1}{Z_1 Z_2 Z_3} \exp\left(-\frac{1}{2}\left(x_2^2\left(\frac{1}{\sigma_2^2} + \frac{w_1^2}{\sigma_1^2} + \frac{w_3^2}{\sigma_3^2}\right) + x_1^2 \frac{1}{\sigma_1^2} - 2x_1 x_2 \frac{w_1}{\sigma_1^2} + x_3^2 \frac{1}{\sigma_3^2} - 2x_3 x_2 \frac{w_3}{\sigma_3^2}\right)\right) \\
&= \frac{1}{Z_1 Z_2 Z_3} \exp\left(-\frac{1}{2}\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}\begin{bmatrix} \frac{1}{\sigma_1^2} & -\frac{w_1}{\sigma_1^2} & 0 \\ -\frac{w_1}{\sigma_1^2} & \left(\frac{1}{\sigma_2^2} + \frac{w_1^2}{\sigma_1^2} + \frac{w_3^2}{\sigma_3^2}\right) & -\frac{w_3}{\sigma_3^2} \\ 0 & -\frac{w_3}{\sigma_3^2} & \frac{1}{\sigma_3^2} \end{bmatrix}\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right)
\end{aligned}
$$

# Example 2: Explaining away

Bayesian Inference with Gaussians

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[DJC MacKay, *The humble Gaussian distribution*, 2006]

gas price

emission
price

$x_1$

$x_3$

$x_2$

electricity
price

$$x_1 = \nu_1 \qquad\qquad\quad p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$$

$$x_3 = \nu_3 \qquad\qquad\quad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

$$x_2 = w_1 x_1 + w_3 x_3 + \nu_2 \quad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$

## Example 2: Explaining away

Bayesian Inference with Gaussians

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[DJC MacKay, *The humble Gaussian distribution*, 2006]

gas price

emission
price

$x_1$

$x_3$

$x_2$

electricity
price

$$p(x) = \mathcal{N}\left(x; m, \underbrace{\begin{bmatrix} \sigma_1^2 & w_1\sigma_1^2 & 0 \\ & \sigma_2^2 + w_1^2\sigma_1^2 + w_3^2\sigma_3^2 & w_3\sigma_3^2 \\ & & \sigma_3^2 \end{bmatrix}}_{\Sigma}\right)$$

$x_1 = \nu_1$      $p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$    $p(x_1, x_3) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix}\right)$

$x_3 = \nu_3$      $p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$

$x_2 = w_1 x_1 + w_3 x_3 + \nu_2$    $p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$

# Example 2: Explaining away

a $\pm$ value in the precision matrix implies $\mp$ correlation conditional on everything else [DJC MacKay, *The humble Gaussian distribution*, 2006]

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

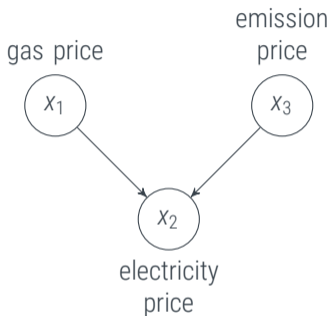$$x_1 = \nu_1 \qquad\qquad p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2)$$
$$x_3 = \nu_3 \qquad\qquad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$
$$x_2 = w_1 x_1 + w_3 x_3 + \nu_2 \qquad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$

$$p(x_1, x_2, x_3) = p(x_1) \cdot p(x_3) \cdot p(x_2 \mid x_1, x_3)$$

$$= \frac{1}{Z_1 \cdot Z_2 \cdot Z_3} \exp\left(-\frac{1}{2}\left(\frac{x_1}{\sigma_1^2} + \frac{x_3}{\sigma_3^2} + \frac{x_2 - w_1 x_1 - w_3 x_3}{\sigma_2^2}\right)\right)$$

$$= \frac{1}{Z_1 \cdot Z_2 \cdot Z_3} \exp\left(-\frac{1}{2}\left(x_1^2\left(\frac{1}{\sigma_1^2} + \frac{w_1^2}{\sigma_2^2}\right) + x_2^2 \frac{1}{\sigma_2^2} - 2x_1 x_2 \frac{w_1}{\sigma_2^2} + x_3^2\left(\frac{1}{\sigma_3^2} + \frac{w_3^2}{\sigma_2^2}\right) - 2x_2 x_3 \frac{w_3}{\sigma_2^2} + 2x_3 x_1 \frac{w_1 w_3}{\sigma_2^2}\right)\right)$$
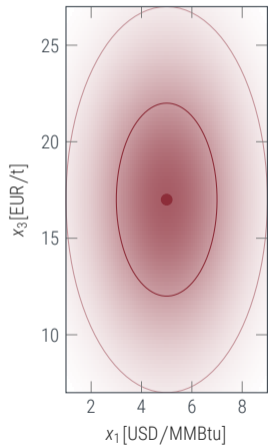
$$= \frac{1}{Z_1 \cdot Z_2 \cdot Z_3} \exp\left(-\frac{1}{2} \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} \left(\frac{1}{2\sigma_1^2} + \frac{w_1^2}{\sigma_2^2}\right) & -\frac{w_1}{\sigma_2^2} & \frac{w_1 w_3}{\sigma_2^2} \\ -\frac{w_1}{\sigma_2^2} & \frac{1}{\sigma_2^2} & -\frac{w_3}{\sigma_2^2} \\ \frac{w_1 w_3}{\sigma_2^2} & -\frac{w_3}{\sigma_2^2} & \left(\frac{1}{2\sigma_3^2} + \frac{w_3^2}{\sigma_2^2}\right) \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right)$$

# Example 2: Explaining away
Bayesian Inference with Gaussians

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[DJC MacKay, *The humble Gaussian distribution*, 2006]

$$p(x_1, x_3) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix}\right)$$

# Example 2: Explaining away

Bayesian Inference with Gaussians

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN
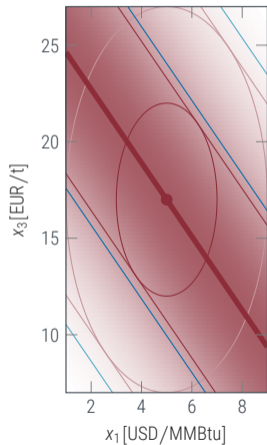
[DJC MacKay, *The humble Gaussian distribution*, 2006]

$$p(x_1, x_3) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix}\right)$$

$$p(x_2) = \mathcal{N}\left(x_2; w_1\mu_1 + w_3\mu_3 + \mu_2, \sigma_2^2 + w_1^2\sigma_1^2 + w_3^2\sigma_3^2\right)$$

# Example 2: Explaining away

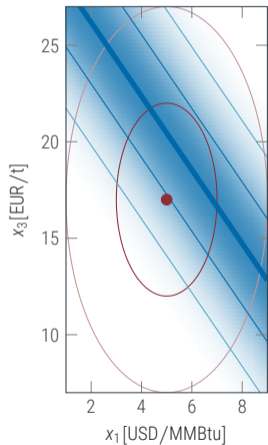Bayesian Inference with Gaussians

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

[DJC MacKay, *The humble Gaussian distribution*, 2006]

$$p(x_1, x_3) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix}\right)$$

$$p(x_2) = \mathcal{N}\left(x_2; w_1\mu_1 + w_3\mu_3 + \mu_2, \sigma_2^2 + w_1^2\sigma_1^2 + w_3^2\sigma_3^2\right)$$

$$p(x_2 \mid x_1, x_3) = \mathcal{N}(x_2; w_1 x_1 + w_3 x_3 + \mu_2, \sigma_2^2)$$

$$p(x_1, x_3) = \mathcal{N}\left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix} ; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix} \right)$$

$$p(x_2) = \mathcal{N}\left( x_2; w_1\mu_1 + w_3\mu_3 + \mu_2, \sigma_2^2 + w_1^2\sigma_1^2 + w_3^2\sigma_3^2 \right)$$

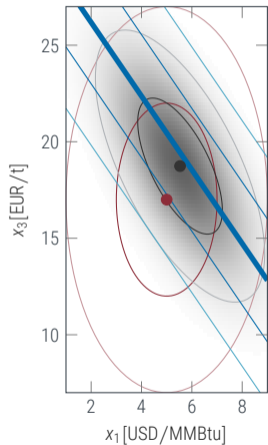$$p(x_2 \mid x_1, x_3) = \mathcal{N}(x_2; w_1 x_1 + w_3 x_3 + \mu_2, \sigma_2^2)$$

$$p(x_1, x_3 \mid x_2) = \mathcal{N}\left( x_{1,3}; \boldsymbol{\mu}_{1,3} - \Sigma_{1,3} w^\mathsf{T} \frac{x_2 - w\boldsymbol{\mu}_{1,3} - \mu_2}{w\Sigma_{1,3}w^\mathsf{T} + \sigma_2^2}, \right.$$

$$\left. \Sigma_{1,3} - \Sigma_{1,3}w^\mathsf{T} \frac{1}{w\Sigma_{1,3}w^\mathsf{T} + \sigma_2^2} w\Sigma_{1,3} \right)$$

$$= \mathcal{N}\left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix} ; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix} - \begin{bmatrix} w_1\sigma_1^2 \\ w_3\sigma_3^2 \end{bmatrix} \frac{x_2 - w_1\mu_1 - w_3\mu_3 - \mu_2}{w_1^2\sigma_1^2 + w_3^2\sigma_3^2 + \sigma_2^2}, \right.$$

$$\left. \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix} - \begin{bmatrix} w_1\sigma_1^2 \\ w_3\sigma_3^2 \end{bmatrix} \frac{1}{w_1^2\sigma_1^2 + w_3^2\sigma_3^2 + \sigma_2^2} \begin{bmatrix} w_1\sigma_1^2 & w_3\sigma_3^2 \end{bmatrix} \right)$$

$$\mathcal{N}(x; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \boldsymbol{\mu})^\mathsf{T} \Sigma^{-1}(x - \boldsymbol{\mu})\right)$$

Today:

▶ Gaussian distributions provide the linear algebra of inference.
  ▶ products of Gaussians are Gaussians
  ▶ linear maps of Gaussian variables are Gaussian variables
  ▶ marginals of Gaussians are Gaussians
  ▶ linear conditionals of Gaussians are Gaussians

If all variables in a generative model are linearly related, and the distributions of the parent variables are Gaussian, then all conditionals, joints and marginals are Gaussian, with means and covariances computable by linear algebra operations.

▶ A zero off-diagonal element in the covariance matrix implies independence if all other variables are integrated out
▶ A zero off-diagonal element in the precision matrix implies independence conditional on all other variables

$$[\Sigma]_{ij} = 0 \qquad \Rightarrow \qquad p(x_i, x_j) = \mathcal{N}(x_i; [\boldsymbol{\mu}]_i, [\Sigma]_{ii}) \cdot \mathcal{N}(x_j; [\boldsymbol{\mu}]_j, [\Sigma]_{jj})$$

$$[\Sigma^{-1}]_{ij} = 0 \qquad \Rightarrow \qquad p(x_i, x_j \mid x_{\neq i,j}) = \mathcal{N}(x_i \mid x_{\neq i,j}) \cdot \mathcal{N}(x_j \mid x_{\neq i,j})$$

The Toolbox

Framework:

$$\int p(x_1, x_2)\, dx_2 = p(x_1) \qquad p(x_1, x_2) = p(x_1 \mid x_2) p(x_2) \qquad p(x \mid y) = \frac{p(y \mid x) p(x)}{p(y)}$$

Modelling:
- ▶ Directed Graphical Models
- ▶ Gaussian Distributions
- ▶
- ▶
- ▶
- ▶

Computation:
- ▶ Monte Carlo
- ▶ Linear algebra / Gaussian inference
- ▶
- ▶
- ▶