

# PROBABILISTIC INFERENCE AND LEARNING

## LECTURE 03

### GAUSSIAN PROBABILITY DISTRIBUTIONS

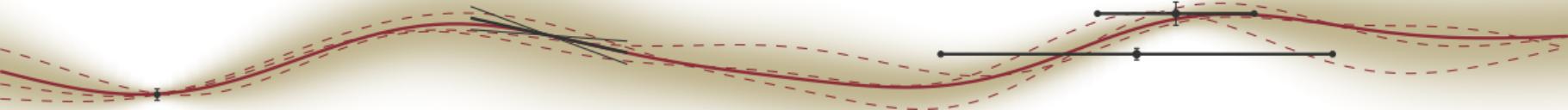
Philipp Hennig

24 October 2018

EBERHARD KARLS  
**UNIVERSITÄT**  
TÜBINGEN

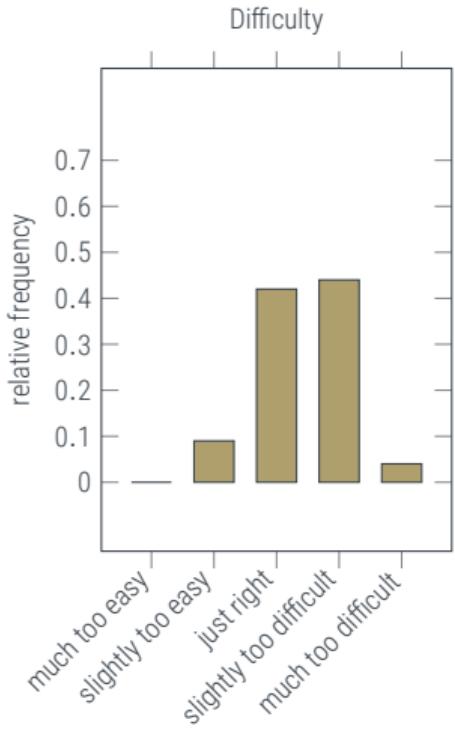
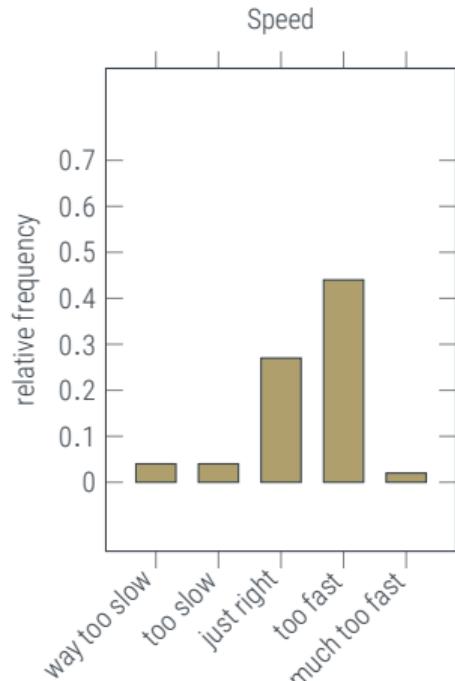
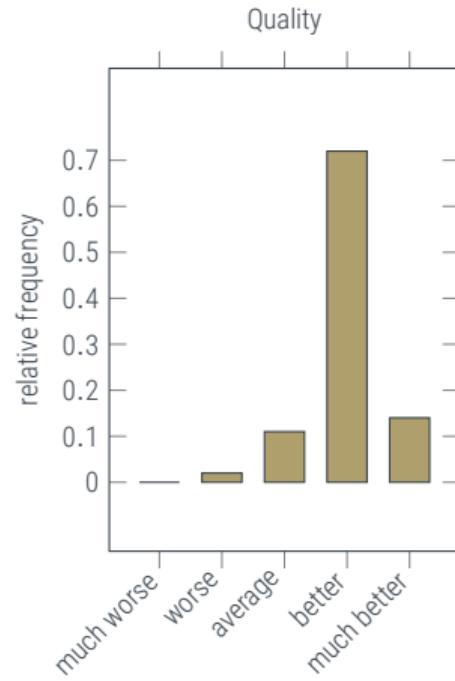


FACULTY OF SCIENCE  
DEPARTMENT OF COMPUTER SCIENCE  
CHAIR FOR THE METHODS OF MACHINE LEARNING



# Last Lecture: Debrief

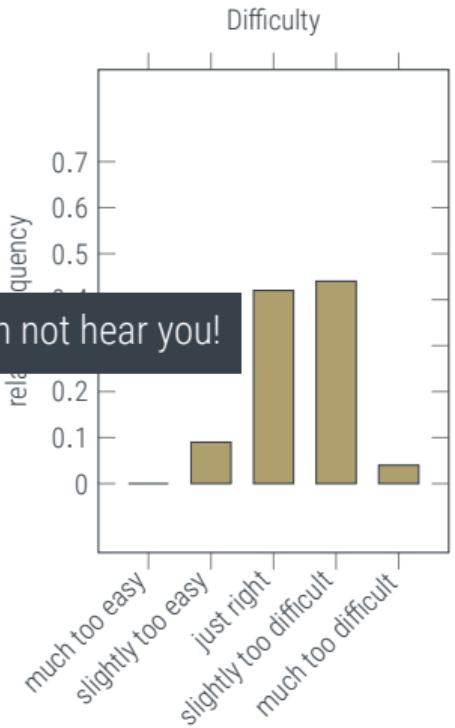
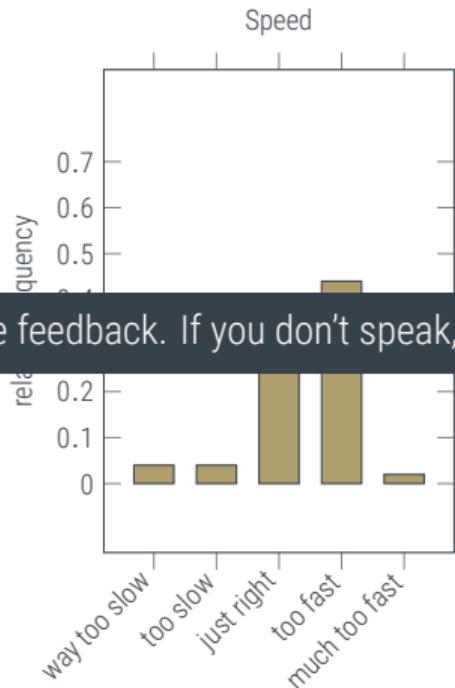
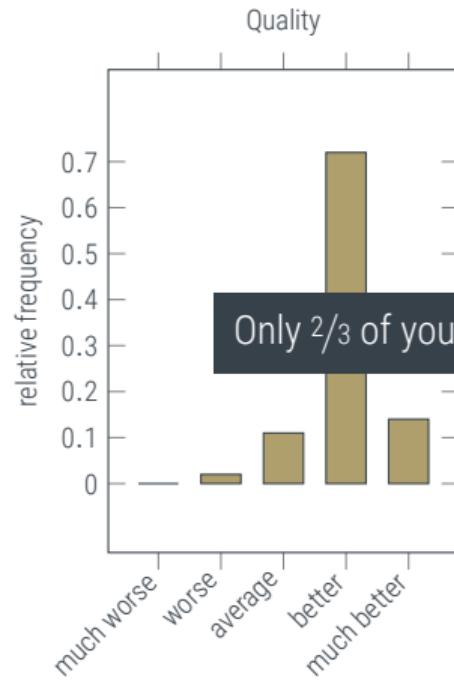
Feedback dashboard





# Last Lecture: Debrief

## Feedback dashboard



Only 2/3 of you gave feedback. If you don't speak, I can not hear you!



# Last Lecture: Debrief

## Detailed Feedback

### Things you did not like:

- ♦ Please write on the blackboard again (!)
- ♦ I do not like the lengthy quotes
- ♦ "the plots were not helpful"
- ♦ the definitions went by too fast
- ♦ the definitions took too long
- ♦ Labelling combinatorics as uninteresting
- ♦ please keep the slides printer-friendly, upload before lecture

### Things you did not understand:

- ♦ change of variable
- ♦ the 3D plot of densities
- ♦ Why are there some black slides?
- ♦ why is the Beta-dist. better than a uniform prior? It introduces new parameters to simulate evidence.  
Doesn't that introduce bias?

### Things you enjoyed:

- ♦ the "glasses" example (!!!!!!!)
- ♦ mathematical definitions
- ♦ explanation on priors
- ♦ citing Laplace
- ♦ interesting side-notes / "meta comments"
- ♦ the 3D plot of densities



## Overview of Lectures so far:

0. Introduction to Reasoning under Uncertainty
  - Probabilities are the mathematical formalization of uncertainty
1. Probabilistic Reasoning
  - Probabilities extend deductive to plausible reasoning. Conditional independence affects complexity
2. Probabilities over Continuous Variables
  - Probability **densities** distribute probability over continuous domains
  - for continuous variables, sum & product rule apply to the PDF, not the CDF
  - densities transform nontrivially under a change of basis
  - It is possible to **infer** probabilities from observations

## Today:

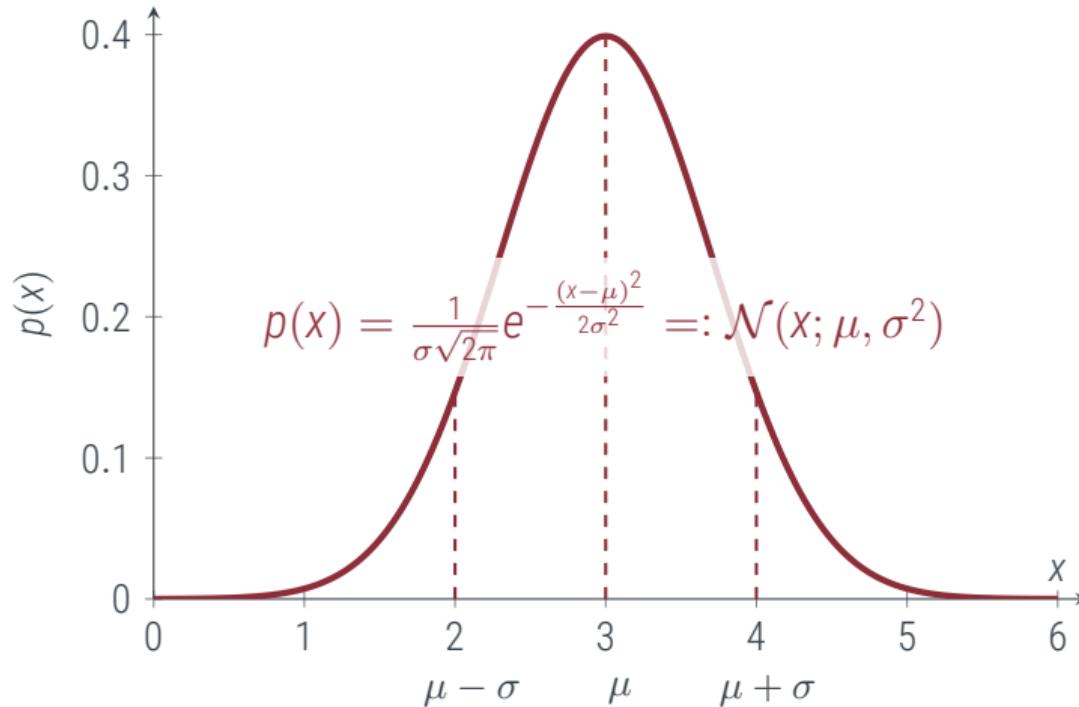
- Gaussian distributions provide the **linear algebra of inference**

00S8



# The (univariate) Gaussian distribution

an exponentiated square



$\mu$  the mean of  $x$   
 $\sigma^2$  the variance of  $x$   
 $\sigma$  the standard deviation of  $x$

# Univariate Gaussians

some observations and notations, conventions

## Definition

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

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

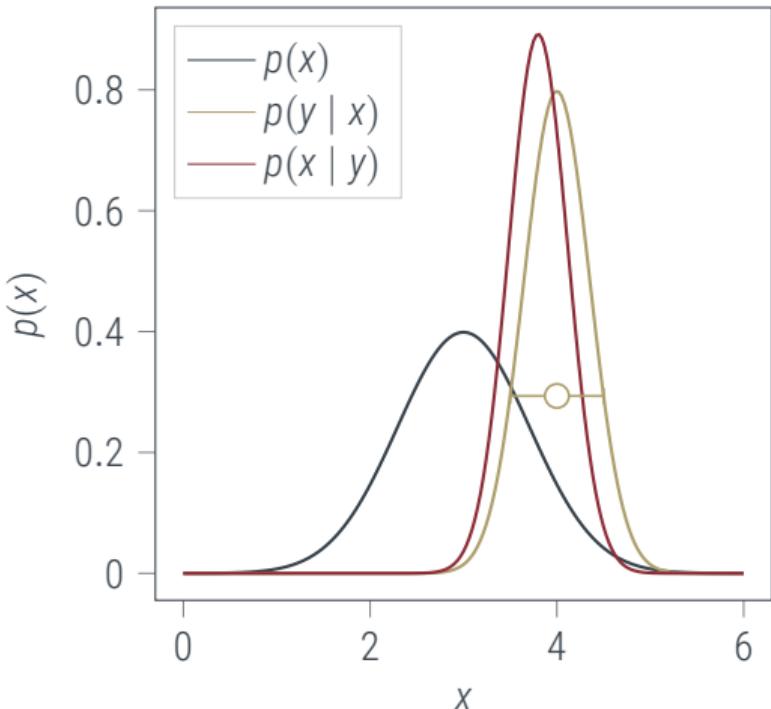
- $\int \mathcal{N}(x; \mu, \sigma^2) dx = 1$  and  $\mathcal{N}(x; \mu, \sigma^2) > 0 \forall x \in \mathbb{R}$ . So  $\mathcal{N}$  is the density of a probability measure.
- Symmetry in  $x$  and  $\mu$ :  $\mathcal{N}(x; \mu, \sigma^2) = \mathcal{N}(\mu, x, \sigma^2)$
- An **exponential of a quadratic polynomial** of the **natural parameters**  $(a, \eta, \tau)$ :

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



# Gaussian Inference

The Gaussian is its own conjugate prior.



Let

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

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

Then

$$\begin{aligned} p(x | y) &= \frac{p(x)p(y | x)}{\int p(x)p(y | x) dx} \\ &= \mathcal{N}(x; m, s^2), \text{ with} \end{aligned}$$

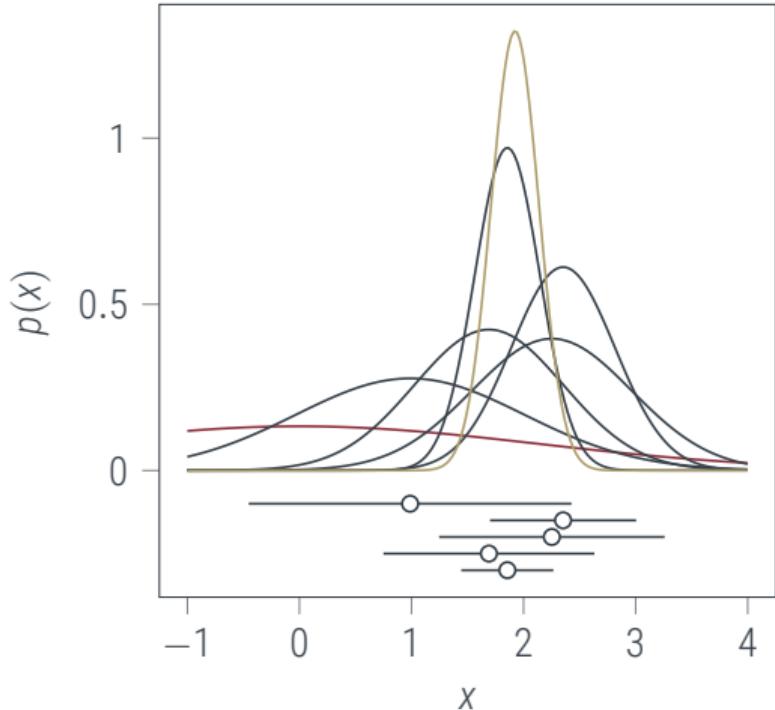
$$s^2 := \frac{1}{\sigma^{-2} + \nu^{-2}}$$

$$m := \frac{\sigma^{-2}\mu + \nu^{-2}y}{\sigma^{-2} + \nu^{-2}}$$



# Gaussian Inference

## Least-Squares Estimation



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

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

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

$$\begin{aligned} p(x | y) &= \frac{p(x)p(y | x)}{\int p(x)p(y | x) dx} \\ &= \mathcal{N}(x; m, s^2), \text{ with} \end{aligned}$$

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

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



# 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]

so wird allgemein sein müssen  $\varphi'(M-p) + \varphi'(M'-p) + \varphi'(M''-p) + \text{etc.} = 0$ , wenn für  $p$  der Werth  $\frac{1}{\mu}(M+M'+M''+\text{etc.})$  substituirt wird, welches positive Ganze nun auch durch  $\mu$  ausdrückt sein mag. Setzt man daher voraus  $M = M' = \text{etc.} = M - \mu N$ , so wird allgemein, d. h. für jeden ganzen positiven Werth für  $\mu$ , sein  $\varphi'(\mu-1)N = (1-\mu)\varphi'(-N)$ , woraus man leicht sieht, dass allgemein  $\frac{\varphi(A)}{A}$  eine constante Grösse sein müsse, welche ich mit  $k$  bezeichnen will. Hieraus wird  $\log \varphi A = \frac{1}{2}kAA + \text{Const.}$ , oder wenn man die Basis der hyperbolischen Logarithmen mit  $e$  bezeichnet und die Constante  $= \log z$  setzt,

$$\varphi A = ze^{kAA}.$$

Ferner sieht man leicht ein, dass  $k$  nothwendig negativ sein müsse, damit  $\Omega$  in der That ein Grösster werden könnte, weshalb wir setzen  $\frac{1}{2}k = -\frac{h}{\pi}$ ; und da vermittelst des eleganten, zuerst von Laplace\*) gefundenen Theorems das Integral  $\int e^{-\frac{hAA}{\pi}} dA$ , von  $A = -\infty$  bis zu  $A = +\infty$ , wird  $= \frac{V\pi}{h}$  (wobei  $\pi$  den halben Kreisumfang für den Radius = 1 bezeichnet), so wird unsere Function werden:

$$\varphi A = \frac{h}{V\pi} e^{-\frac{hAA}{\pi}}.$$

178.

Die so eben ermittelte Function kann zwar nicht in aller Strenge die Wahrscheinlichkeiten der Fehler ausdrücken; denn da die möglichen Fehler stets in gewisse Grenzen eingewängt sind, so müsste die Wahrscheinlichkeit grösserer Fehler immer = 0 herauskommen, während unsere Formel stets einen begrenzten Werth darstellt. Dennoch aber ist dieser Mangel, an welchem jede analytische Function ihrer Natur nach laboriren muss, für jeden praktischen

\*) In v. Zach „Monatliche Correspondenz“ Band 21, S. 280 äussert Gauss: „Dass Euler schon das Theorem gefunden hat, woraus der schöne, von mir Laplace beigelegte Lehrsatz sehr leicht abgeleitet werden kann, fiel mir selbst schon früher ein, als über die Stelle S. 212 schon abgedruckt war; ich wollte es aber nicht unter die Erstes setzen, weil Laplace wenigstens das obige Theorem doch erst in der dort gebrauchten Form aufgestellt hat.“  
*Anmerkung des Übersetzers.*





# The Multivariate Gaussian distribution

An exponentiated quadratic form

Definition (multivariate Gaussian distribution)

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

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



# The Multivariate Gaussian distribution

An exponentiated quadratic form

Definition (multivariate Gaussian distribution)

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

$\Sigma$  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  $A = A^\top$ , and

$$v^\top A v \geq 0 \quad \forall v \in \mathbb{R}^n.$$

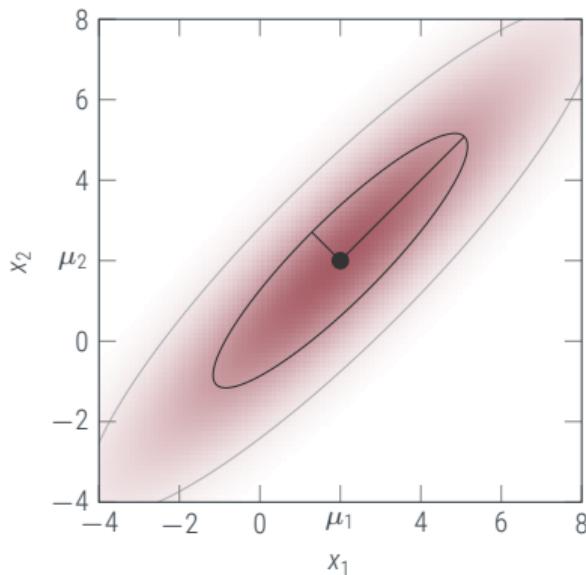
Equivalent statement: All eigenvalues of the symmetric matrix  $A$  are non-negative.



# The Multivariate Gaussian distribution

Equiprobability lines are ellipsoids

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



- $\int \mathcal{N}(x; \mu, \Sigma) = 1$  and  $\mathcal{N}(x; \mu, \Sigma) > 0 \forall x \in \mathbb{R}^n$ .
- Symmetry in  $x$  and  $\mu$ :  $\mathcal{N}(x; \mu, \Sigma) = \mathcal{N}(\mu; x, \Sigma)$
- An exponential of a quadratic polynomial:

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

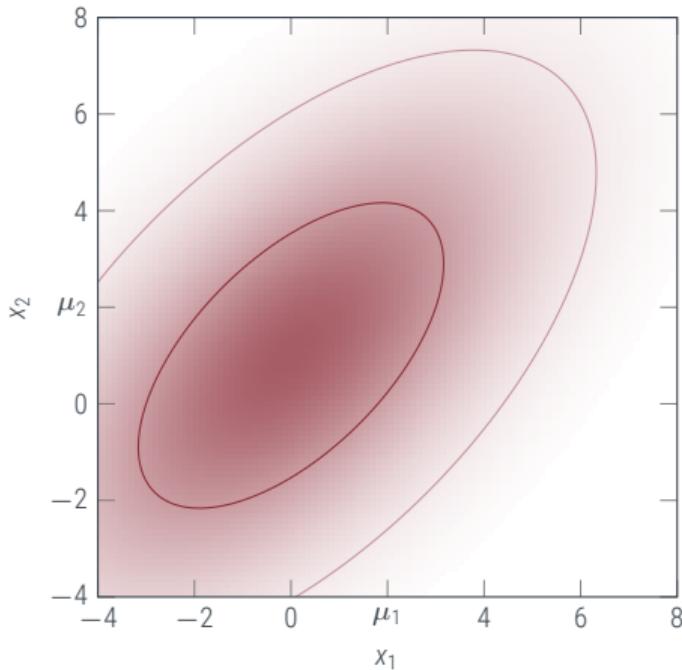
$$= \exp\left(a + \eta^\top x - \frac{1}{2} \text{tr}(xx^\top \Lambda)\right) \quad (2)$$

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



# Products of Gaussians are Gaussians

Closure under Multiplication



To multiply Gaussians, add the natural parameters

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

$$C = (A^{-1} + B^{-1})^{-1}$$

$$c = C(A^{-1}a + B^{-1}b)$$

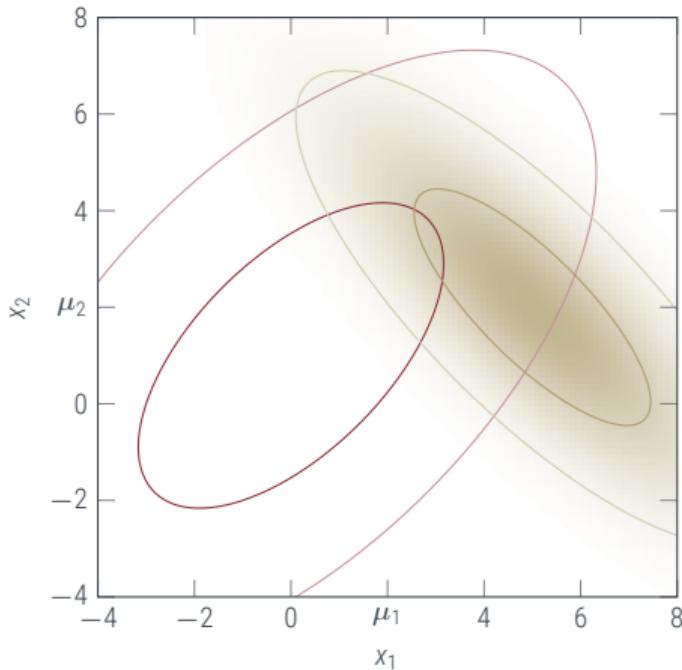
$$Z = \mathcal{N}(a; b, A + B)$$

Note similarity to univariate case.



# Products of Gaussians are Gaussians

Closure under Multiplication



To multiply Gaussians, add the natural parameters

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

$$C = (A^{-1} + B^{-1})^{-1}$$

$$c = C(A^{-1}a + B^{-1}b)$$

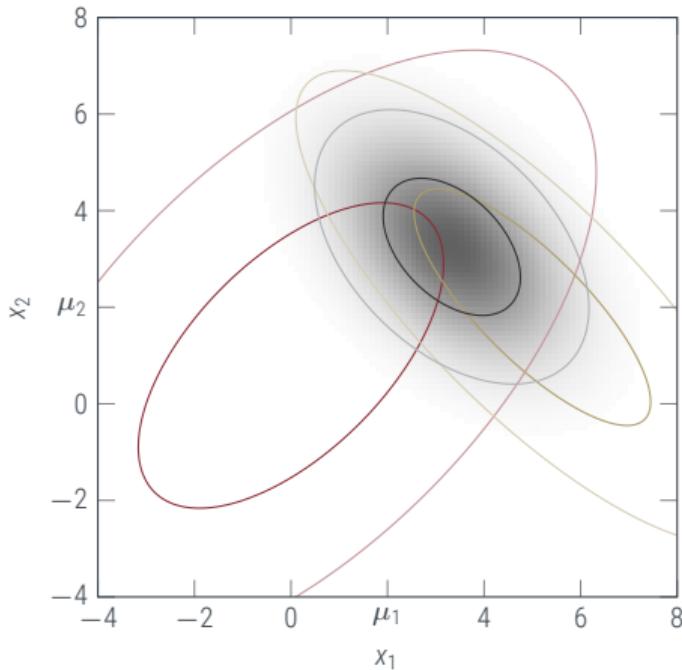
$$Z = \mathcal{N}(a; b, A + B)$$

Note similarity to univariate case.



# Products of Gaussians are Gaussians

Closure under Multiplication



To multiply Gaussians, add the natural parameters

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

$$C = (A^{-1} + B^{-1})^{-1}$$

$$c = C(A^{-1}a + B^{-1}b)$$

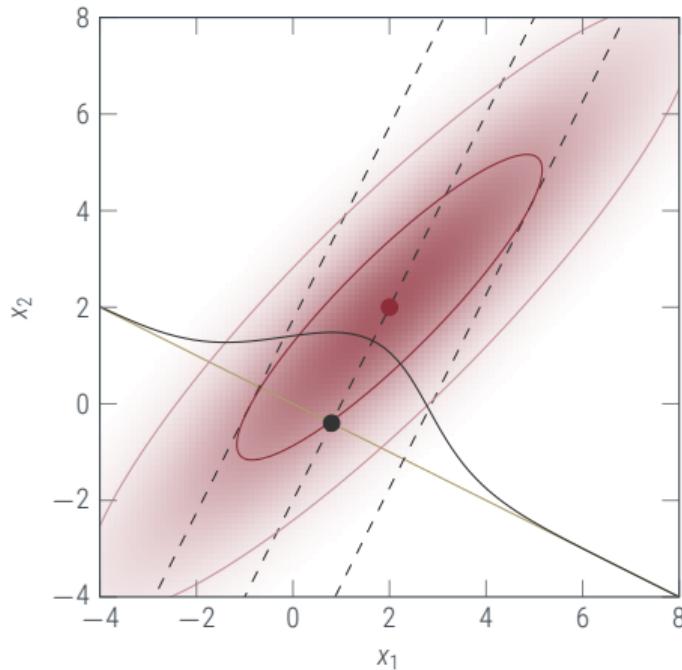
$$Z = \mathcal{N}(a; b, A + B)$$

Note similarity to univariate case.



# Linear Projections of Gaussians are Gaussians

Closure under linear maps

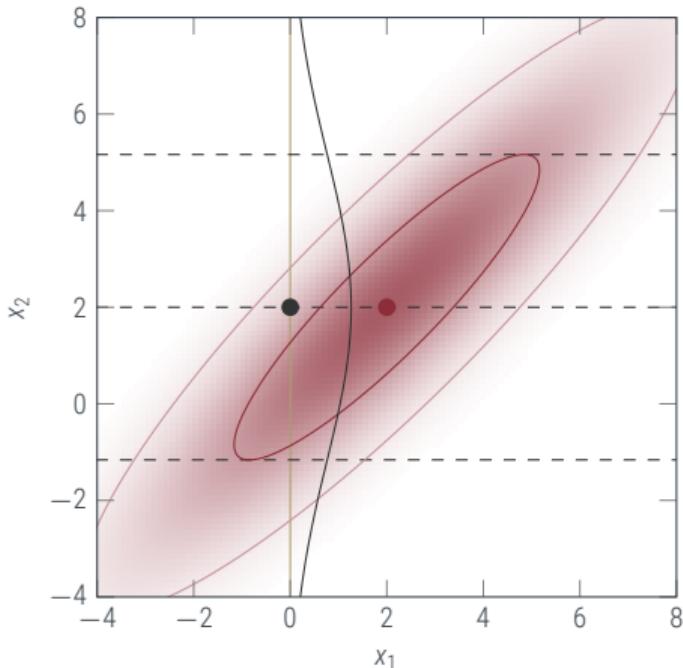


To linearly project a Gaussian variable,  
project the parameters

$$\begin{aligned} p(z) &= \mathcal{N}(z; \mu, \Sigma) \\ \Rightarrow p(Az) &= \mathcal{N}(Az, A\mu, A\Sigma A^\top) \end{aligned}$$

# Marginals of Gaussians are Gaussians

Closure under marginalization



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

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

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

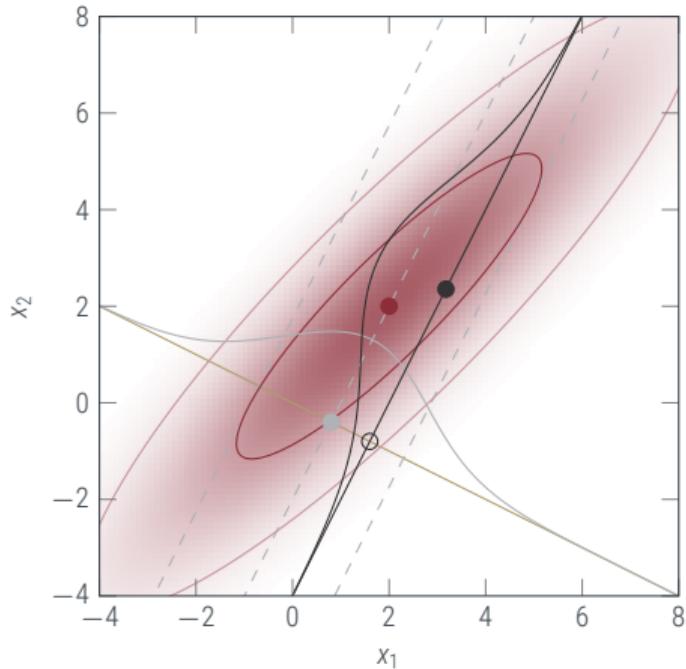
- this is the **sum rule**

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

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

# Cuts through Gaussians are Gaussians

Closure under conditioning



$$\begin{aligned}
 p(x | Ax = y) &= \frac{p(x, y)}{p(y)} \\
 &= \mathcal{N}(x; \mu + \Sigma A^T (A\Sigma A^T)^{-1} (y - A\mu), \\
 &\quad \Sigma - \Sigma A^T (A\Sigma A^T)^{-1} A\Sigma)
 \end{aligned}$$

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

# Inference with Gaussians

Since conditioning and marginalization are mapped to linear algebra, so is Bayes' Theorem

## Theorem

If  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$

and  $p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y}; A\mathbf{x} + \mathbf{b}, \Lambda)$ ,

then  $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; A\boldsymbol{\mu} + \mathbf{b}, \Lambda + A\Sigma A^\top)$

and  $p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} + \underbrace{\Sigma A^\top (A\Sigma A^\top + \Lambda)^{-1}}_{\text{gain}} \underbrace{(\mathbf{y} - (A\boldsymbol{\mu} + \mathbf{b}))}_{\text{residual}}, \Sigma - \underbrace{\Sigma A^\top (A\Sigma A^\top + \Lambda)^{-1} A \Sigma}_{\text{Gram matrix}})$

$$= \mathcal{N}(\mathbf{x}; (\underbrace{\Sigma^{-1} + A^\top \Lambda^{-1} A}_{\text{precision matrix}})^{-1} (A^\top \Lambda^{-1} (\mathbf{y} - \mathbf{b}) + \Sigma^{-1} \boldsymbol{\mu}), (\underbrace{\Sigma^{-1} + A^\top \Lambda^{-1} A}_{\text{precision matrix}})^{-1})$$

# Gaussians provide the linear algebra of inference

if all joints are Gaussian and all observations are linear, all posteriors are Gaussian

- products of Gaussians are Gaussians

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

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

- 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

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

- linear projections of Gaussians are Gaussians

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

Bayesian inference becomes linear algebra

$$p(x) = \mathcal{N}(x; \mu, \Sigma) \quad p(y | x) = \mathcal{N}(y; Ax + b, \Lambda)$$

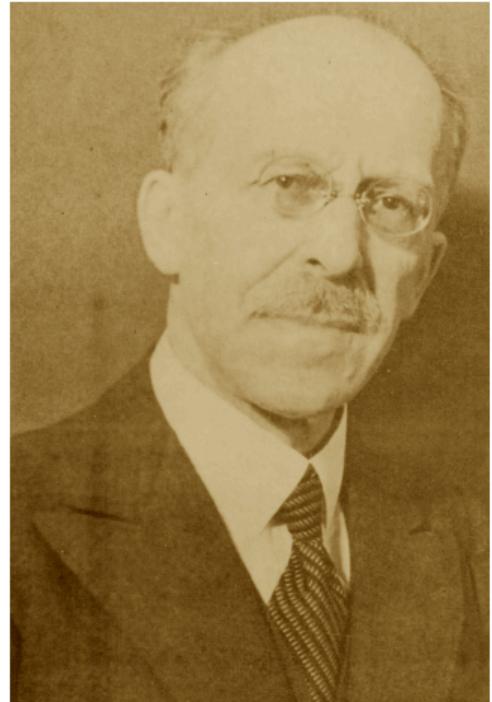
$$p(B^\top x + c | y) = \mathcal{N}[B^\top x + c; B^\top \mu + c + B^\top \Sigma A (A^\top \Sigma A + \Lambda)^{-1} (y - A^\top \mu - b), B^\top \Sigma B - B^\top \Sigma A (A^\top \Sigma A + \Lambda)^{-1} A^\top \Sigma B]$$

# The Core Insight for All of This

Gaussian inference is linear algebra at its core



$$A = \begin{bmatrix} P & Q \\ R & S \end{bmatrix} \quad M := (S - RP^{-1}Q)^{-1}$$
$$A^{-1} = \begin{bmatrix} P^{-1} + P^{-1}QMRP^{-1} & -P^{-1}QM \\ -MRP^{-1} & M \end{bmatrix}$$
$$(Z + UWV^\top)^{-1} = Z^{-1} - Z^{-1}U(W^{-1} + V^\top Z^{-1}U)^{-1}V^\top Z^{-1}$$
$$|Z + UWV^\top| = |Z| \cdot |W| \cdot |W^{-1} + V^\top Z^{-1}U|$$



Issai Schur (1875–1941)

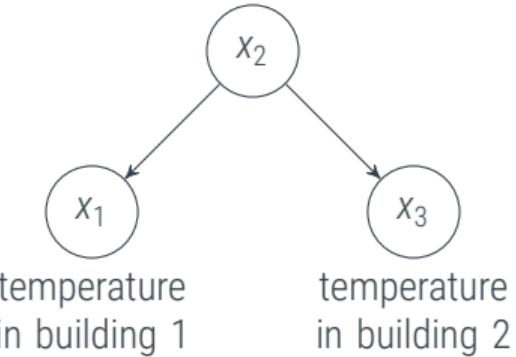
# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians



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

temperature outside



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

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

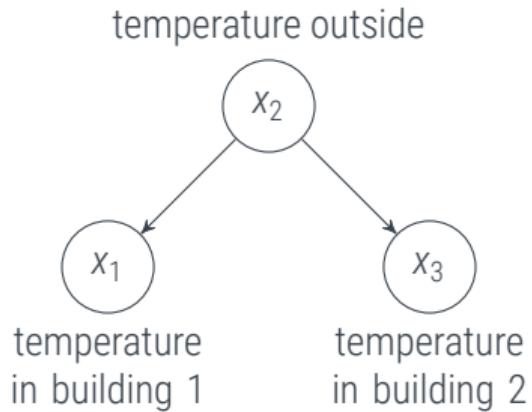
$$x_3 = w_3 x_2 + \nu_3 \quad p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians



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



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

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

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

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

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

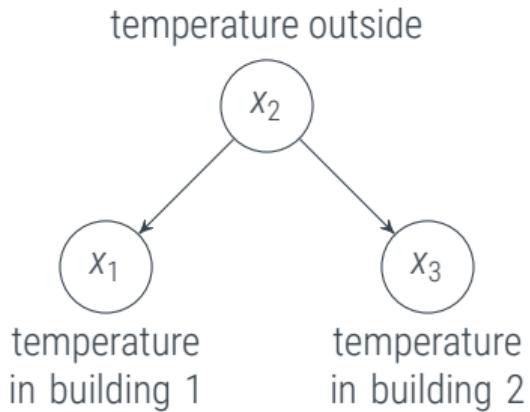
$$p(\mathbf{x} = A\boldsymbol{\nu}) = \mathcal{N}\left(\mathbf{x}; \underbrace{A\boldsymbol{\mu}}_{=:m}, \underbrace{\begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1\sigma_2^2 & w_1w_3\sigma_2^2 \\ \sigma_2^2 & w_3\sigma_2^2 & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix}}_{=: \Sigma}\right)$$

# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians



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



$$x_2 = \nu_2$$

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

$$x_1 = w_1 x_2 + \nu_1$$

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

$$x_3 = w_3 x_2 + \nu_3$$

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

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

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

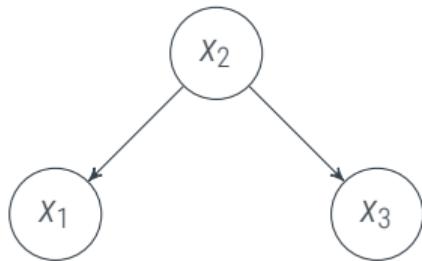
$$p(\mathbf{x} = A\boldsymbol{\nu}) = \mathcal{N}\left(\mathbf{x}; \underbrace{A\boldsymbol{\mu}}_{=:m}, \underbrace{\begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1\sigma_2^2 & w_1w_3\sigma_2^2 \\ w_1\sigma_2^2 & \sigma_2^2 & w_3\sigma_2^2 \\ w_1w_3\sigma_2^2 & w_3\sigma_2^2 & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix}}_{=: \Sigma}\right)$$

From graph:  $x_1 \perp\!\!\!\perp x_3 | x_2$ . 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]



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

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

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

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

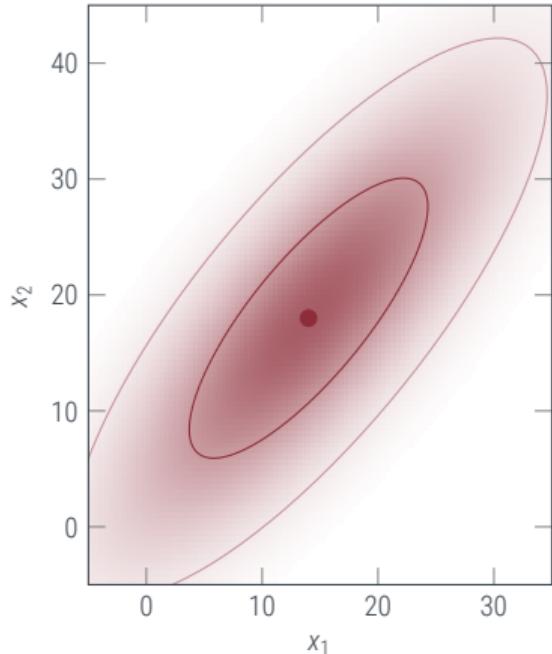
$$= \frac{1}{Z_1 Z_2 Z_3} \exp \left( -\frac{1}{2} \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} \frac{1}{\sigma_1^2} & -\frac{w_1}{\sigma_1^2} & 0 \\ -\frac{w_1}{\sigma_1^2} & \left( \frac{1}{\sigma_2^2} + \frac{w_1^2}{\sigma_1^2} + \frac{w_3^2}{\sigma_3^2} \right) & -\frac{w_3}{\sigma_3^2}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

# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians



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



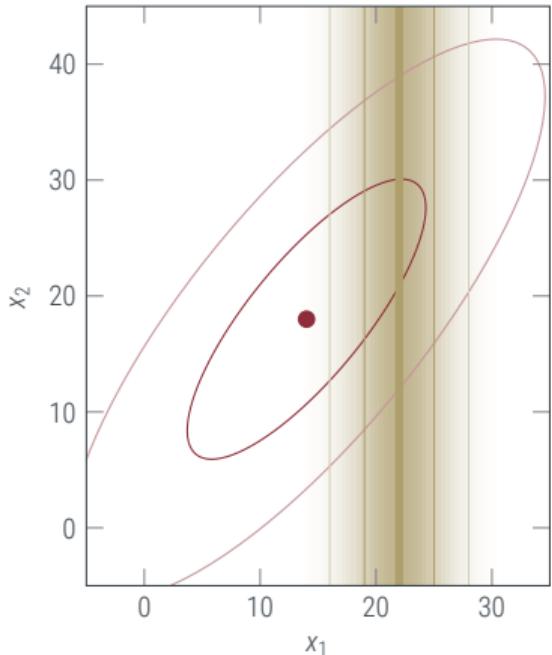
$$\begin{aligned} p \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix} \right) &= \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} w_1\mu_2 + \mu_1 \\ w_3\mu_2 + \mu_3 \end{bmatrix}, \begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1w_3\sigma_2^2 \\ w_1w_3\sigma_2^2 & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix} \right) \\ &= \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} m_1 \\ m_3 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{13} \\ \Sigma_{31} & \Sigma_{33} \end{bmatrix} \right) \end{aligned}$$

# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians



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



$$p \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix} \right) = \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} w_1\mu_2 + \mu_1 \\ w_3\mu_2 + \mu_3 \end{bmatrix}, \begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1w_3\sigma_2^2 \\ w_1w_3\sigma_2^2 & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix} \right)$$

$$= \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} m_1 \\ m_3 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{13} \\ \Sigma_{31} & \Sigma_{33} \end{bmatrix} \right)$$

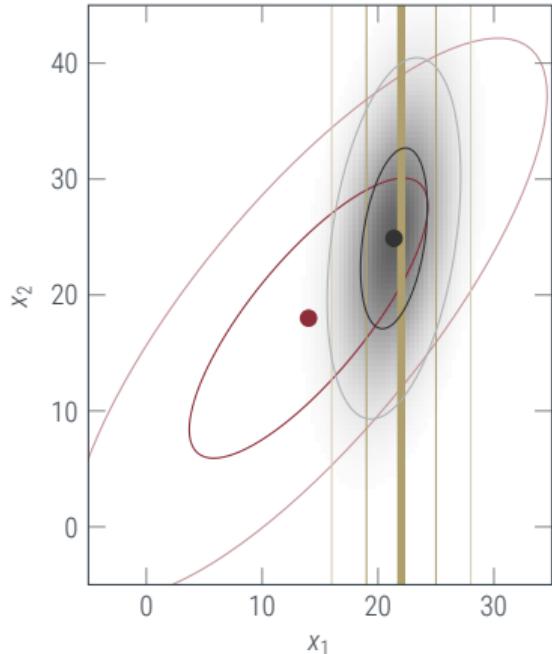
$$p(y | x_1, x_3) = \mathcal{N}(y; x_1, \xi^2) = \mathcal{N} \left( y; \underbrace{\begin{bmatrix} 1 & 0 \end{bmatrix}}_{=:B} \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}, \xi^2 \right)$$

# Example 1: Conditional Independence, Marginal Correlation

Bayesian Inference with Gaussians



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



$$\begin{aligned} p\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}\right) &= \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} w_1\mu_2 + \mu_1 \\ w_3\mu_2 + \mu_3 \end{bmatrix}, \begin{bmatrix} w_1\sigma_2^2 + \sigma_1^2 & w_1w_3\sigma_2^2 \\ w_1w_3\sigma_2^2 & w_3^2\sigma_2^2 + \sigma_3^2 \end{bmatrix}\right) \\ &= \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} m_1 \\ m_3 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{13} \\ \Sigma_{31} & \Sigma_{33} \end{bmatrix}\right) \end{aligned}$$

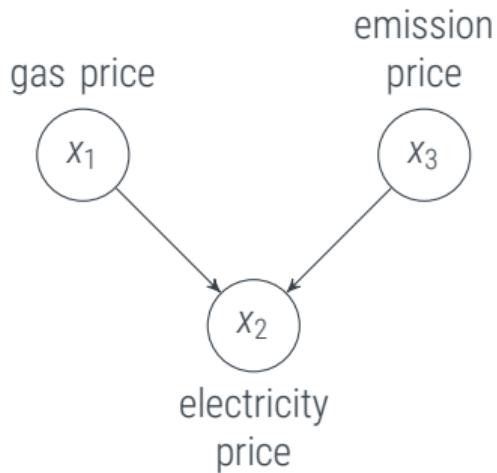
$$p(y | x_1, x_3) = \mathcal{N}(y; x_1, \xi^2) = \mathcal{N}\left(y; \underbrace{\begin{bmatrix} 1 & 0 \end{bmatrix}}_{=:B} \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}, \xi^2\right)$$

$$p(x_1, x_3 | y) = \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; m + \Sigma B^\top \frac{(y - Bm)}{B\Sigma B^\top + \xi^2}, \Sigma - \frac{1}{B\Sigma B^\top + \xi^2} \Sigma B^\top B\right)$$

# Example 2: Explaining away

Bayesian Inference with Gaussians

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



$$x_1 = \nu_1$$

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

$$x_3 = \nu_3$$

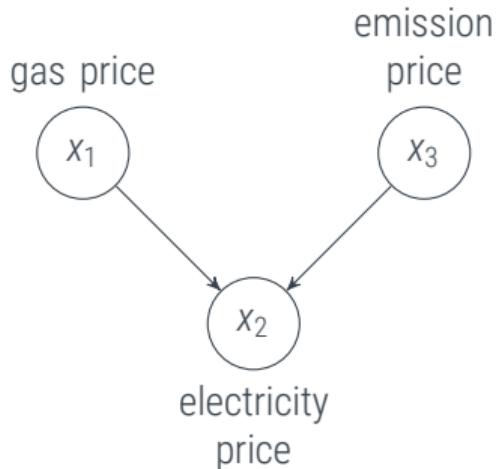
$$p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2)$$

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

# Example 2: Explaining away

Bayesian Inference with Gaussians

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



$$x_1 = \nu_1$$

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

$$x_3 = \nu_3$$

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

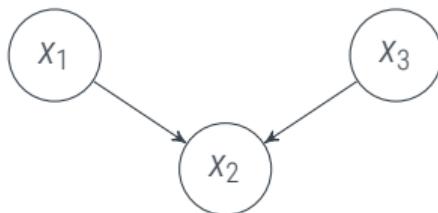
$$x_2 = w_1 x_1 + w_3 x_3 + \nu_2 \quad p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2)$$

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

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

## Example 2: Explaining away

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



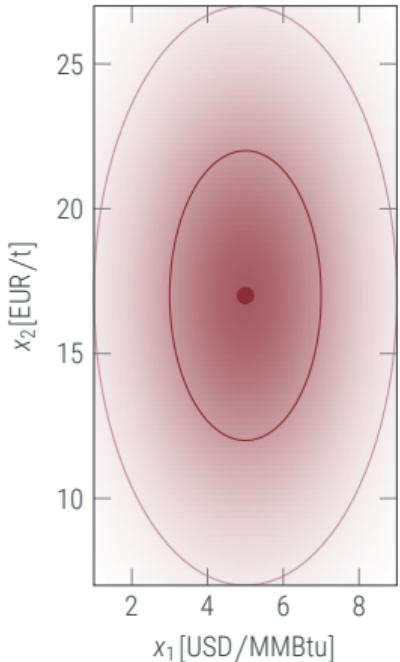
$$\begin{array}{lll} x_1 = \nu_1 & p(\nu_1) = \mathcal{N}(\nu_1; \mu_1, \sigma_1^2) \\ x_3 = \nu_3 & p(\nu_3) = \mathcal{N}(\nu_3; \mu_3, \sigma_3^2) \\ x_2 = w_1 x_1 + w_3 x_3 + \nu_2 & p(\nu_2) = \mathcal{N}(\nu_2; \mu_2, \sigma_2^2) \end{array}$$

$$\begin{aligned} p(x_1, x_2, x_3) &= p(x_1) \cdot p(x_3) \cdot p(x_2 | x_1, x_3) \\ &= \frac{1}{Z_1 \cdot Z_2 \cdot Z_3} \exp \left( -\frac{1}{2} \left( \frac{x_1}{\sigma_1^2} + \frac{x_3}{\sigma_3^2} + \frac{x_2 - w_1 x_1 - w_3 x_3}{\sigma_2^2} \right) \right) \\ &= \frac{1}{Z_1 \cdot Z_2 \cdot Z_3} \exp \left( -\frac{1}{2} \left( x_1^2 \left( \frac{1}{\sigma_1^2} + \frac{w_1^2}{\sigma_2^2} \right) + x_2^2 \frac{1}{\sigma_2^2} - 2x_1 x_2 \frac{w_1}{\sigma_1^2} + x_3^2 \left( \frac{1}{\sigma_3^2} + \frac{w_3^2}{\sigma_2^2} \right) - 2x_2 x_3 \frac{w_3}{\sigma_2^2} - 2x_3 x_1 \frac{w_1 w_3}{\sigma_2^2} \right) \right) \\ &= \frac{1}{Z_1 \cdot Z_2 \cdot Z_3} \exp \left( -\frac{1}{2} \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} \left( \frac{1}{2\sigma_1^2} + \frac{w_1^2}{\sigma_2^2} \right) & -\frac{w_2}{\sigma_2^2} & \frac{w_1 w_3}{\sigma_2^2} \\ -\frac{w_2}{\sigma_2^2} & \frac{1}{\sigma_2^2} & -\frac{w_3}{\sigma_2^2} \\ \frac{w_1 w_3}{\sigma_2^2} & -\frac{w_3}{\sigma_2^2} & \left( \frac{1}{2\sigma_3^2} + \frac{w_3^2}{\sigma_2^2} \right) \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \end{aligned}$$

# Example 2: Explaining away

Bayesian Inference with Gaussians

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

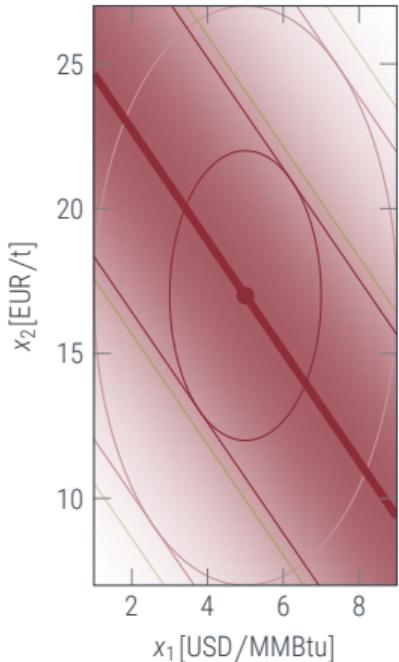


$$p(x_1, x_3) = \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix} \right)$$

# Example 2: Explaining away

Bayesian Inference with Gaussians

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



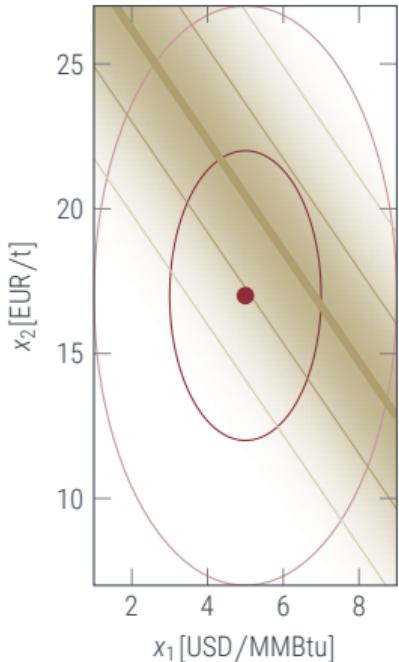
$$p(x_1, x_3) = \mathcal{N} \left( \begin{bmatrix} x_1 \\ x_3 \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_3^2 \end{bmatrix} \right)$$

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

# Example 2: Explaining away

Bayesian Inference with Gaussians

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



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

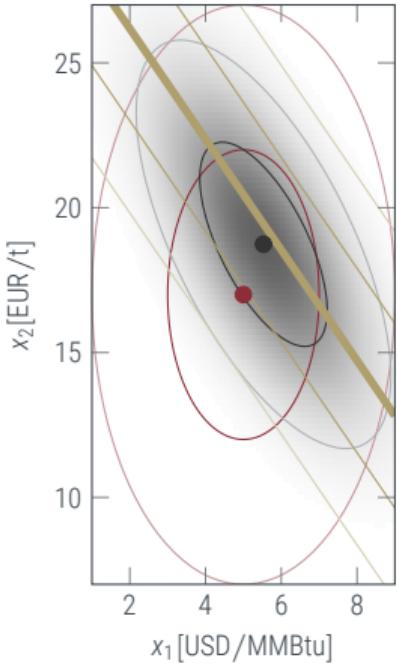
$$p(x_2) = \mathcal{N} \left( x_2; w_1\mu_1 + w_3\mu_3 + \mu_2, \sigma_2^2 + w_1^2\sigma_1^2 + w_3^2\sigma_3^2 \right)$$

$$p(x_2 | x_1, x_3) = \mathcal{N}(x_2; w_1x_1 + w_3x_3 + \mu_2, \sigma_2^2)$$

# Example 2: Explaining away

Bayesian Inference with Gaussians

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



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

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

$$p(x_2 | x_1, x_3) = \mathcal{N}(x_2; w_1x_1 + w_3x_3 + \mu_2, \sigma_2^2)$$

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

$$\left. \Sigma_{1,3} - \Sigma_{1,3} W^\top \frac{1}{W\Sigma_{1,3}W^\top + \sigma_2^2} W\Sigma_{1,3} \right)$$

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

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

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

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

$$\begin{aligned} [\Sigma]_{ij} = 0 \quad &\Rightarrow \quad p(x_i, x_j) = \mathcal{N}(x_i; [\boldsymbol{\mu}]_i, [\Sigma]_{ii}) \cdot \mathcal{N}(x_j; [\boldsymbol{\mu}]_j, [\Sigma]_{jj}) \\ [\Sigma^{-1}]_{ij} = 0 \quad &\Rightarrow \quad p(x_i, x_j \mid x_{\neq i,j}) = \mathcal{N}(x_i \mid x_{\neq i,j}) \cdot \mathcal{N}(x_j \mid x_{\neq i,j}) \end{aligned}$$