

# PROBABILISTIC INFERENCE AND LEARNING

## LECTURE 07

### PARAMETRIC REGRESSION

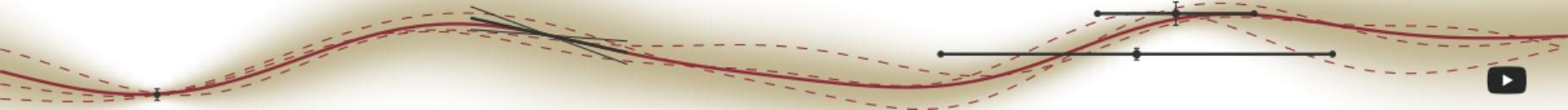
Philipp Hennig

11 May 2020

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
2	21.04.	Reasoning under Uncertainty		15	15.06.	Exponential Families	8
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.	<b>Parametric Regression</b>	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	



# Recap: Gaussian Distributions

Gaussian distributions provide the linear algebra of inference

- ▶ products of Gaussians are Gaussians

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

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

- ▶ marginals of Gaussians are Gaussians

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

- ▶ (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

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

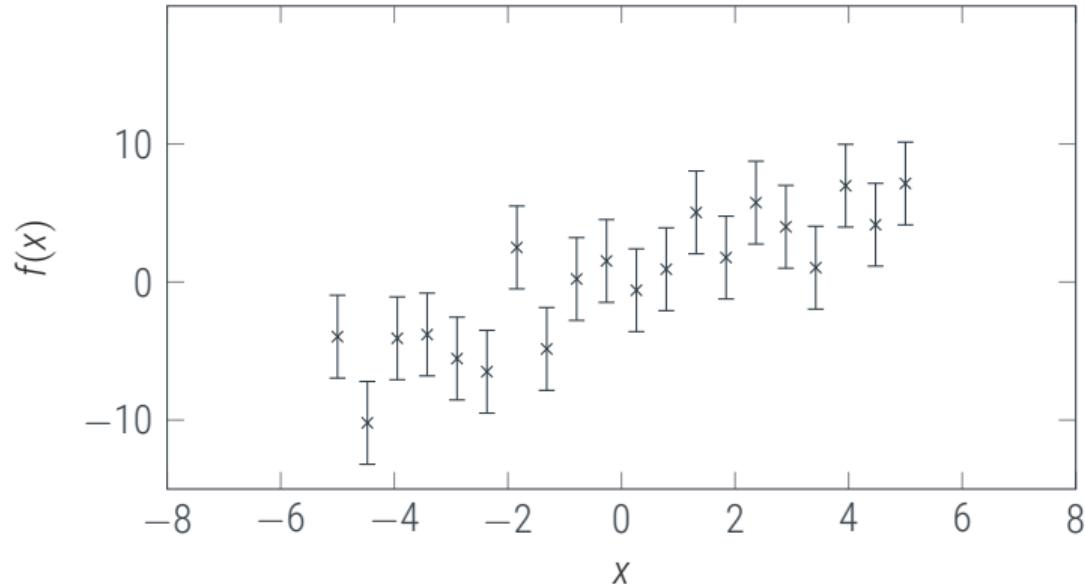
$$p(B^\top x + c | y) = \mathcal{N}[B^\top x + c; B^\top \mu + c + B^\top \Sigma A (A^\top \Sigma A + \Lambda)^{-1} (y - A^\top \mu - b), B^\top \Sigma B - B^\top \Sigma A (A^\top \Sigma A + \Lambda)^{-1} A^\top \Sigma B]$$



# Supervised Regression

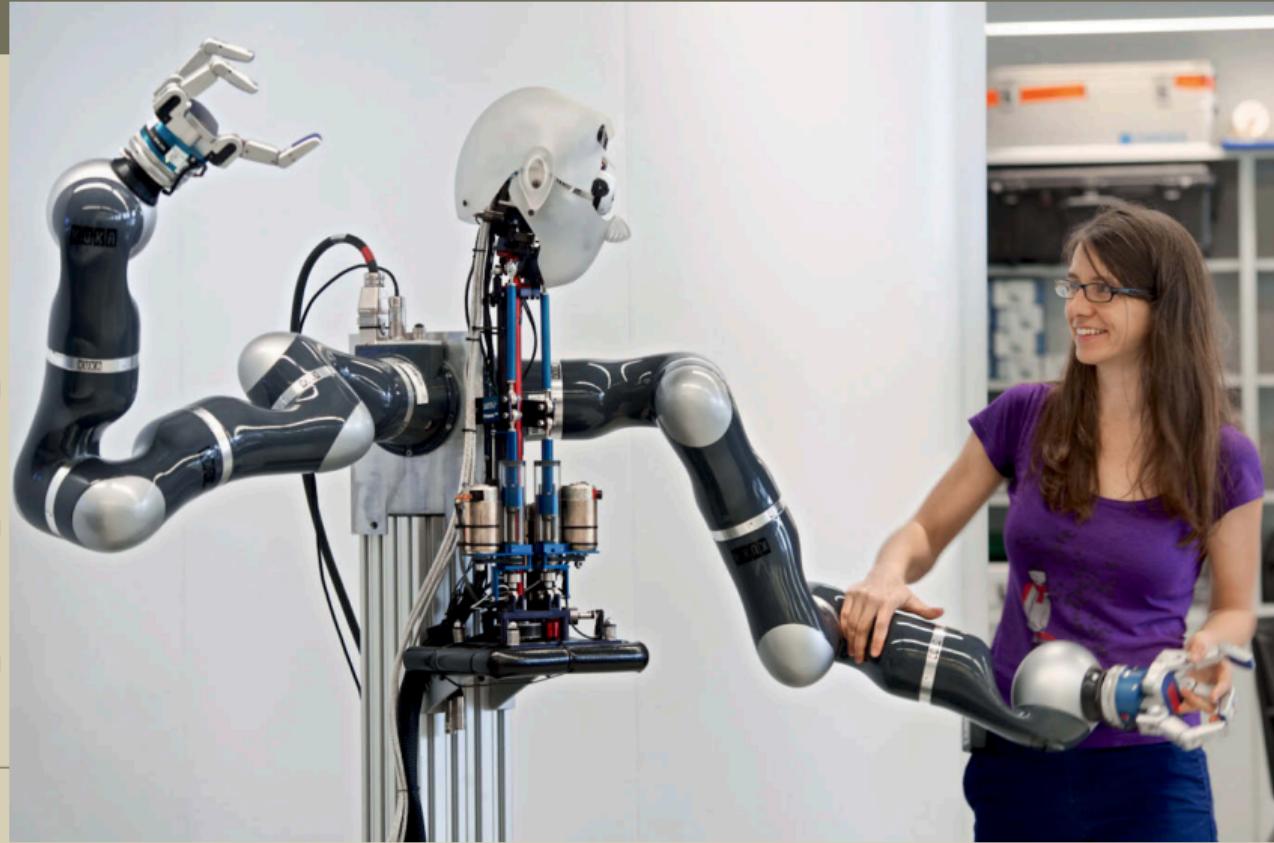
A data set

given:  $y \in \mathbb{R}^N$ ,  $p(y | f) = \mathcal{N}(y; f(x), \sigma^2 I_N)$ . What is  $f$ ?



# Supervised Regression

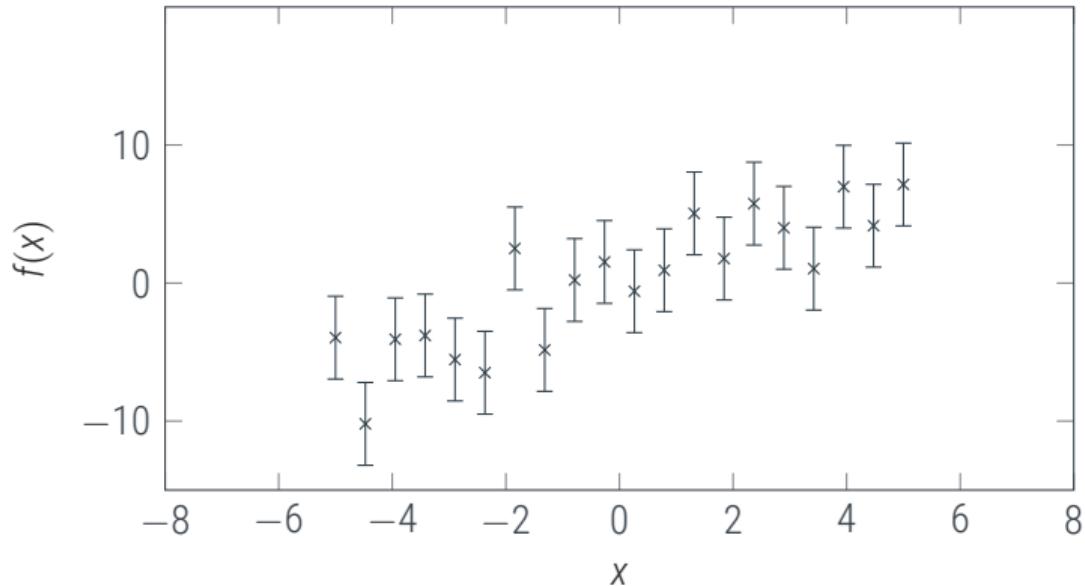
A data set



# Supervised Regression

A data set

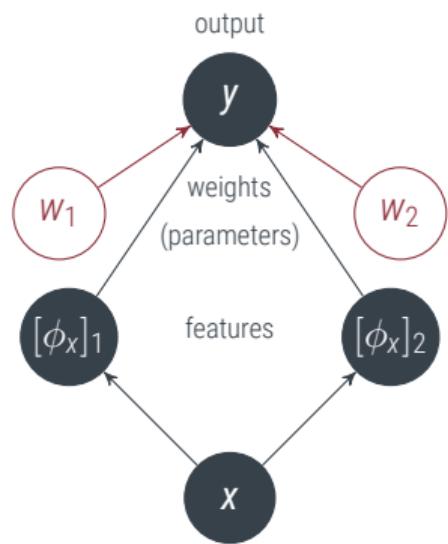
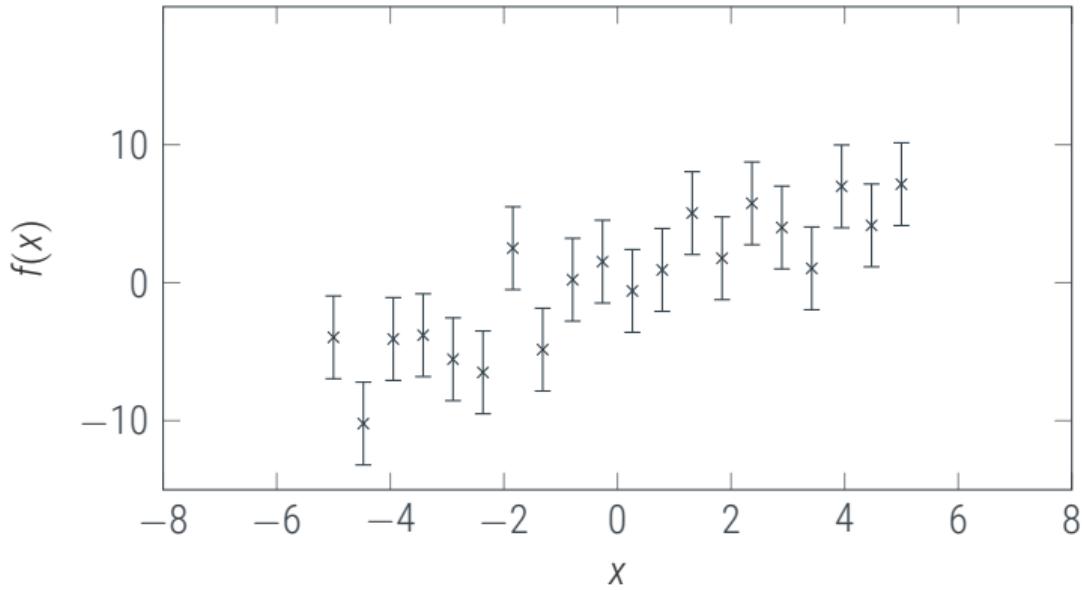
given:  $y \in \mathbb{R}^N$ ,  $p(y | f) = \mathcal{N}(y; f(x), \sigma^2 I_N)$ . What is  $f$ ?



# A Linear Model

linear regression

Assume linear function  $f(x) = w_1 + w_2 x = \phi_x^T w$  with features  $\phi(x) = \begin{bmatrix} 1 \\ x \end{bmatrix} =: \phi_x$



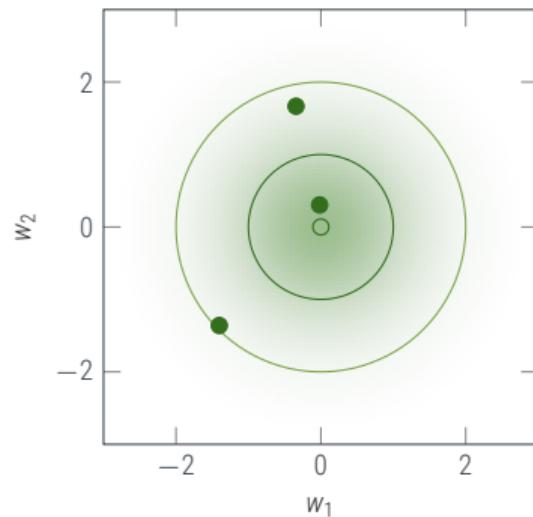


# A linear generative model

if every variable is Gaussian and every relationship is linear, all marginals and conditionals are also Gaussian

$$f(x) = w_1 + w_2 x = \phi_x^T w$$

$$p(w) = \mathcal{N}(w; \mu, \Sigma)$$

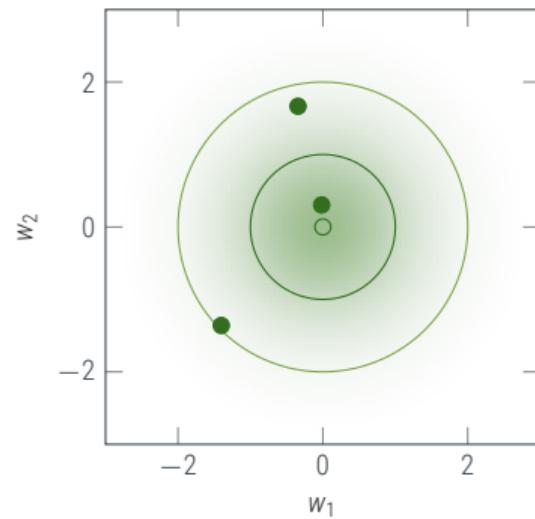


# A linear generative model

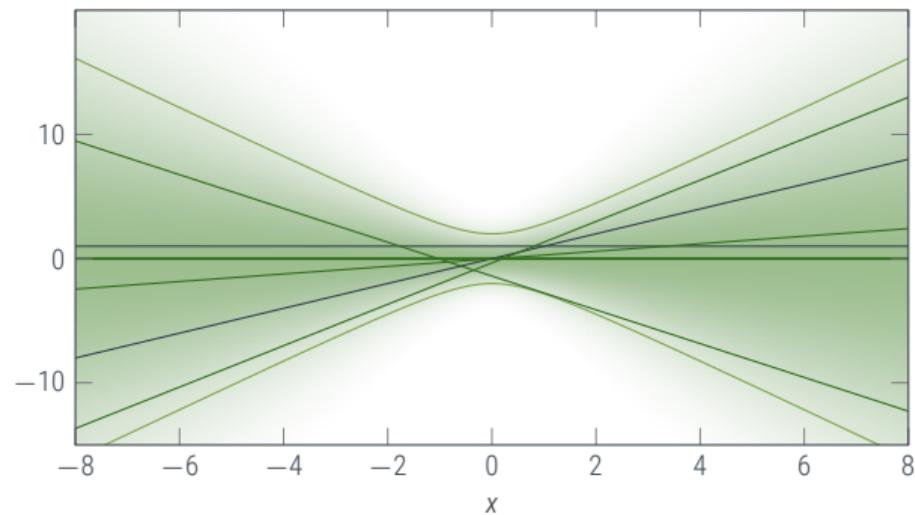
if every variable is Gaussian and every relationship is linear, all marginals and conditionals are also Gaussian

$$f(x) = w_1 + w_2 x = \phi_x^T w$$

$$p(w) = \mathcal{N}(w; \mu, \Sigma)$$



$$p(f) = \mathcal{N}(f; \phi_x^T \mu, \phi_x^T \Sigma \phi_x)$$





# A linear generative model

if every variable is Gaussian and every relationship is linear, all marginals and conditionals are also Gaussian

$$f(x) = w_1 + w_2 x = \phi_x^T w$$

$$p(w) = \mathcal{N}(w; \mu, \Sigma)$$

$$p(f) = \mathcal{N}(f; \phi_x^T \mu, \phi_x^T \Sigma \phi_x)$$



# Notation

this will become exceedingly helpful later on

Dataset:  $X := \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{X}^N$ ,  $y := \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$ . We will use the following very sloppy notation, sloppily

$$\phi_x := \phi(x) = \begin{bmatrix} 1 \\ x \end{bmatrix} \in \mathbb{R}^F \quad \phi_X := [\phi(x_1) \quad \phi(x_2) \quad \cdots \quad \phi(x_N)] = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \end{bmatrix} \in \mathbb{R}^{F \times N}$$

$$f_x := f(x) \in \mathbb{R} \quad f_X := \phi_X^\top w = \begin{bmatrix} \phi_1(x_1) & \phi_2(x_1) \\ \phi_1(x_2) & \phi_2(x_2) \\ \vdots & \vdots \\ \phi_1(x_N) & \phi_2(x_N) \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \phi_{x_1}^\top w \\ \phi_{x_2}^\top w \\ \vdots \\ \phi_{x_N}^\top w \end{bmatrix} = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_N) \end{bmatrix} \in \mathbb{R}^N$$

Think of  $f$  as an infinitely long vector, indexed by  $x$ :

$$v \in \mathbb{R}^N, l \in \mathbb{N}^d \Rightarrow v_l := [v_{l_1}, \dots, v_{l_d}] \in \mathbb{R} \iff f \in \mathbb{R}^\infty, X \in \mathbb{R}^N \Rightarrow f_X := [f_{x_1}, \dots, f_{x_N}] \in \mathbb{R}^N.$$



# Gaussian Inference on a linear function

weight space / function space

$$\text{prior } p(w) = \mathcal{N}(w; \mu, \Sigma) \Rightarrow p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x)$$

$$\text{likelihood } p(y | w, \phi_x) = \mathcal{N}(y; \phi_x^\top w, \sigma^2 I) = \mathcal{N}(y; f_x, \sigma^2 I)$$



# Gaussian Inference on a linear function

weight space / function space

$$\text{prior } p(w) = \mathcal{N}(w; \mu, \Sigma) \Rightarrow p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x)$$

$$\text{likelihood } p(y | w, \phi_x) = \mathcal{N}(y; \phi_x^\top w, \sigma^2 I) = \mathcal{N}(y; f_x, \sigma^2 I)$$

$$\text{posterior on } w \quad p(w | y, \phi_x) = \mathcal{N}(w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu),$$

$$\Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma)$$

$$= \mathcal{N}\left(w; (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right),\right.$$

$$\left.(\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1}\right)$$

# Gaussian Inference on a linear function

weight space / function space

$$\text{prior } p(w) = \mathcal{N}(w; \mu, \Sigma) \Rightarrow p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x)$$

$$\text{likelihood } p(y | w, \phi_x) = \mathcal{N}(y; \phi_x^\top w, \sigma^2 I) = \mathcal{N}(y; f_x, \sigma^2 I)$$

$$\text{posterior on } w \quad p(w | y, \phi_x) = \mathcal{N}(w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu),$$

$$\Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma)$$

$$= \mathcal{N}\left(w; (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right), (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1}\right)$$

$$\text{posterior on } f \quad p(f_x | y, \phi_x) = \mathcal{N}(f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu),$$

$$\phi_x^\top \Sigma \phi_x - \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \phi_x)$$

$$\mathcal{N}\left(f_x; \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right), \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \phi_x^\top\right)$$



# Graphical View

inference in weight- and function-space

$$p(w | y, \phi_x) = \mathcal{N} \left( w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right)$$

$$p(f_x | y, \phi_x) = \mathcal{N} \left( f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \phi_x^\top \left( \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right) \phi_x \right)$$





# Graphical View

inference in weight- and function-space

$$p(w | y, \phi_x) = \mathcal{N} \left( w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right)$$

$$p(f_x | y, \phi_x) = \mathcal{N} \left( f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \phi_x^\top \left( \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right) \phi_x \right)$$





# Graphical View

inference in weight- and function-space

$$p(w | y, \phi_x) = \mathcal{N} \left( w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right)$$

$$p(f_x | y, \phi_x) = \mathcal{N} \left( f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \phi_x^\top \left( \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right) \phi_x \right)$$





# Graphical View

inference in weight- and function-space

$$p(w | y, \phi_x) = \mathcal{N} \left( w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right)$$

$$p(f_x | y, \phi_x) = \mathcal{N} \left( f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \phi_x^\top \left( \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right) \phi_x \right)$$





# Graphical View

inference in weight- and function-space

$$p(w | y, \phi_x) = \mathcal{N} \left( w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right)$$

$$p(f_x | y, \phi_x) = \mathcal{N} \left( f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \phi_x^\top \left( \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \right) \phi_x \right)$$



# That's it for the Algebra

Gaussian Inference on a linear function

$$\text{prior } p(w) = \mathcal{N}(w; \mu, \Sigma) \Rightarrow p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x)$$

$$\text{likelihood } p(y | w, \phi_x) = \mathcal{N}(y; \phi_x^\top w, \sigma^2 I) = \mathcal{N}(y; f_x, \sigma^2 I)$$

$$\begin{aligned} \text{posterior on } w \quad p(w | y, \phi_x) &= \mathcal{N}(w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \\ &\quad \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma) \\ &= \mathcal{N}\left(w; (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right), \right. \\ &\quad \left. (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1}\right) \end{aligned}$$

$$\begin{aligned} \text{posterior on } f \quad p(f_x | y, \phi_x) &= \mathcal{N}(f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} (y - \phi_x^\top \mu), \\ &\quad \phi_x^\top \Sigma \phi_x - \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 I)^{-1} \phi_x^\top \Sigma \phi_x) \\ &= \mathcal{N}\left(f_x; \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \left(\Sigma^{-1} \mu + \sigma^{-2} \phi_x y\right), \right. \\ &\quad \left. \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x^\top \phi_x)^{-1} \phi_x^\top\right) \end{aligned}$$



# Code

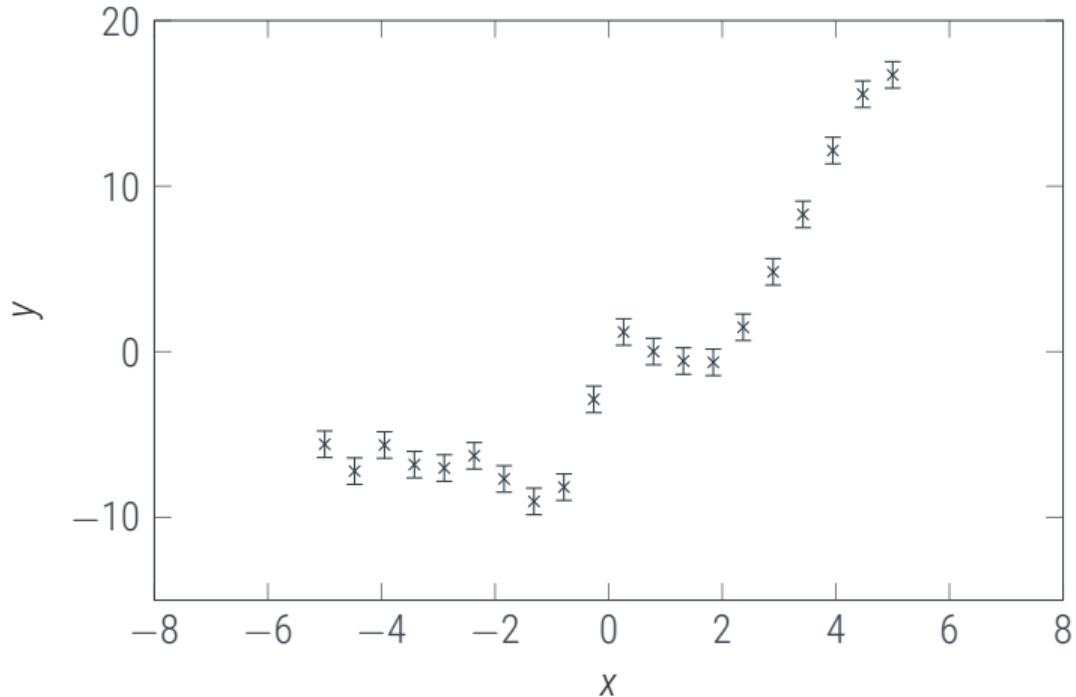
Gaussian\_Linear\_Regression.ipynb





# A more Realistic Dataset

General linear regression





$$f(x) = w_1 + w_2 x = \phi_x^T w$$

$$\phi_x := \begin{bmatrix} 1 \\ x \end{bmatrix}$$



# Cubic Regression



$$f(x) = \phi(x)^T w \quad \phi(x) = [1 \quad x \quad x^2 \quad x^3]^T$$

# Cubic Regression



$$f(x) = \phi(x)^T w \quad \phi(x) = [1 \quad x \quad x^2 \quad x^3]^T$$



# Septic Regression (?)

$$f(x) = \phi(x)^T w \quad \phi(x) = [1 \quad x \quad x^2 \quad \dots \quad x^7]^T$$



# Septic Regression (?)

$$f(x) = \phi(x)^T w \quad \phi(x) = [1 \quad x \quad x^2 \quad \dots \quad x^7]^T$$

# Fourier Regression



$$f(x) = \phi(x)^T w \quad \phi(x) = [\cos(x) \quad \cos(2x) \quad \cos(3x) \quad \dots \quad \sin(x) \quad \sin(2x) \quad \dots]^T$$

# Fourier Regression



$$f(x) = \phi(x)^T w \quad \phi(x) = [\cos(x) \quad \cos(2x) \quad \cos(3x) \quad \dots \quad \sin(x) \quad \sin(2x) \quad \dots]^T$$

# Pixel Regression



$$\phi(x) = -1 + 2 \begin{bmatrix} \theta(x - 8) & \theta(8 - x) & \theta(x - 7) & \theta(7 - x) & \dots \end{bmatrix}^\top$$

# Pixel Regression



$$\phi(x) = -1 + 2 \begin{bmatrix} \theta(x - 8) & \theta(8 - x) & \theta(x - 7) & \theta(7 - x) & \dots \end{bmatrix}^\top$$

# Switch Regression



$$\phi(x) = [\theta(x - 8) \quad \theta(8 - x) \quad \theta(x - 7) \quad \theta(7 - x) \quad \dots]^T$$

# Switch Regression



$$\phi(x) = [\theta(x - 8) \quad \theta(8 - x) \quad \theta(x - 7) \quad \theta(7 - x) \quad \dots]^T$$



# V Regression

$$\phi(x) = [|x - 8| - 8 \quad |x - 7| - 7 \quad |x - 6| - 6 \quad \dots]^T$$



# V Regression

$$\phi(x) = [|x - 8| - 8 \quad |x - 7| - 7 \quad |x - 6| - 6 \quad \dots]^T$$

# Legendre Regression



$$\phi(x) = [b^0 P_0(x), b^1 P_1(x), \dots, b^{13} P_{13}(x)]^\top \quad P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n$$

# Legendre Regression



$$\phi(x) = [b^0 P_0(x), b^1 P_1(x), \dots, b^{13} P_{13}(x)]^\top \quad P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n$$

# Eiffel Tower Regression



$$\phi(x) = [e^{-|x-8|} \quad e^{-|x-7|} \quad e^{-|x-6|} \quad \dots]^T$$

# Eiffel Tower Regression



$$\phi(x) = [e^{-|x-8|} \quad e^{-|x-7|} \quad e^{-|x-6|} \quad \dots]^T$$



# Bell Curve Regression

$$\phi(x) = \left[ e^{-\frac{1}{2}(x-8)^2} \quad e^{-\frac{1}{2}(x-7)^2} \quad e^{-\frac{1}{2}(x-6)^2} \quad \dots \right]^T$$



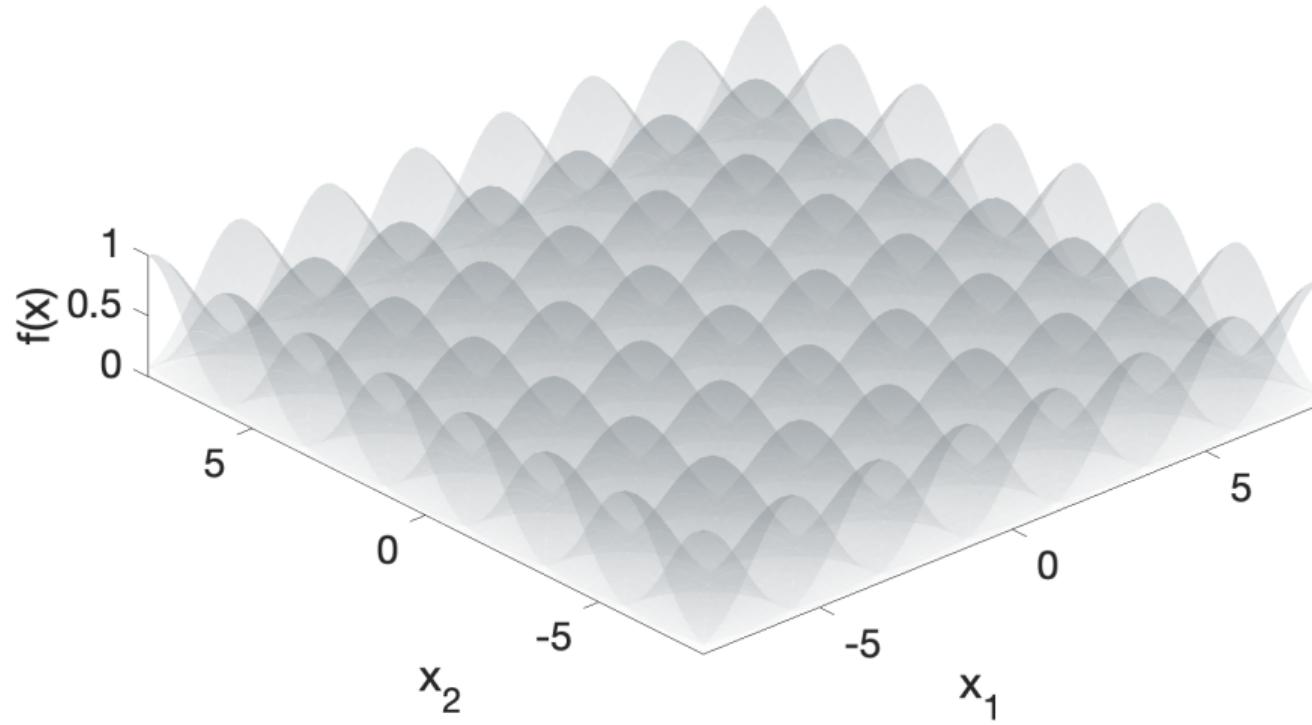
# Bell Curve Regression

$$\phi(x) = \left[ e^{-\frac{1}{2}(x-8)^2} \quad e^{-\frac{1}{2}(x-7)^2} \quad e^{-\frac{1}{2}(x-6)^2} \quad \dots \right]^T$$



# Multiple Inputs

Input domain  $\mathbb{X}$  can be anything





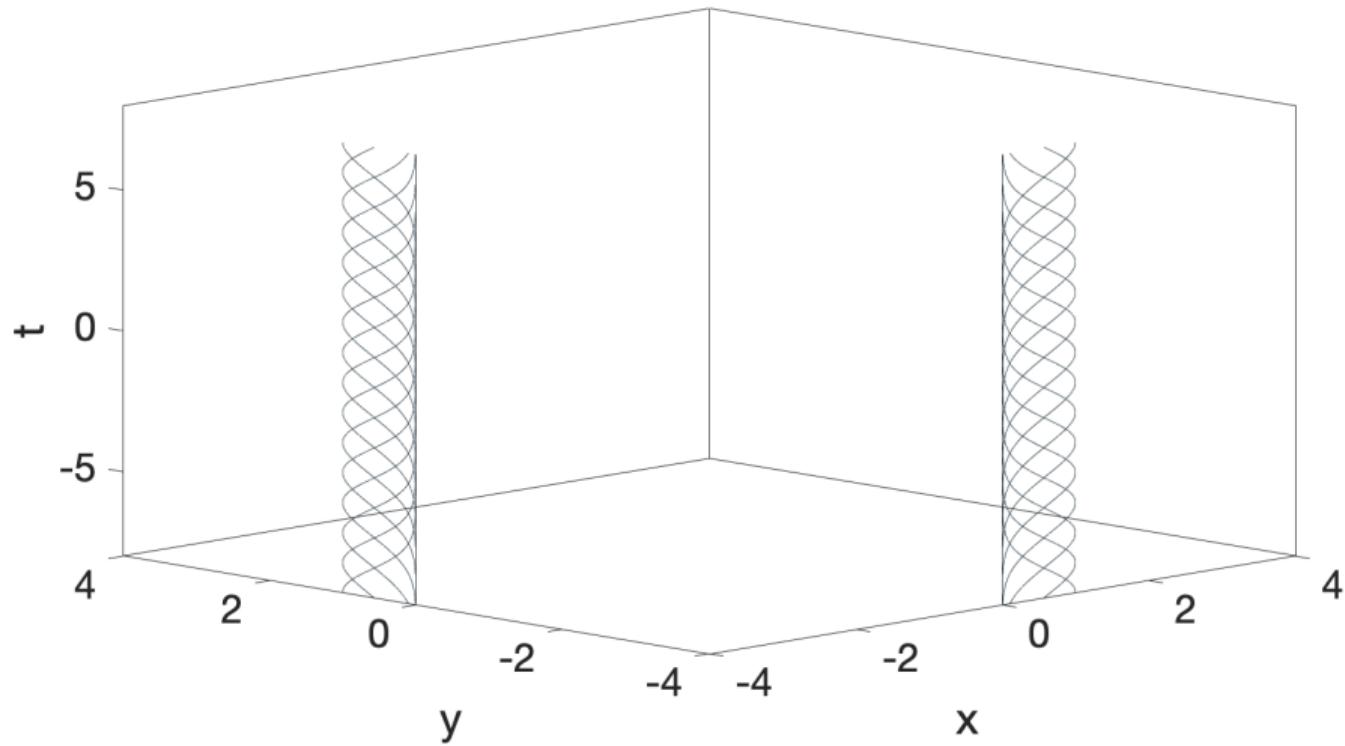
# Multiple Inputs

Input domain  $\mathbb{X}$  can be anything



# Multiple Outputs

The output domain can be anything isomorphic to  $\mathbb{R}^N$





# Multiple Outputs

The output domain can be anything isomorphic to  $\mathbb{R}^N$



# Multiple Outputs

The output domain can be anything isomorphic to  $\mathbb{R}^N$



## Summary:

- ▶ Gaussian distributions can be used to **learn functions**
- ▶ Analytical inference is possible using **general linear models**

$$f(x) = \phi(x)^T w = \phi_x^T w$$

- ▶ Then the posterior on both  $w$  and  $f$  is Gaussian
- ▶ The choice of features  $\phi : \mathbb{X} \rightarrow \mathbb{R}$  is essentially unconstrained

