

PROBABILISTIC INFERENCE AND LEARNING

LECTURE 04

PARAMETRIC GAUSSIAN REGRESSION

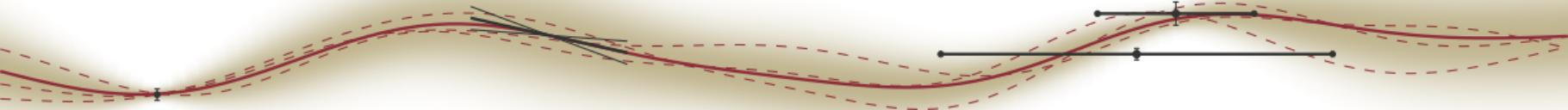
Philipp Hennig

29 October 2018

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

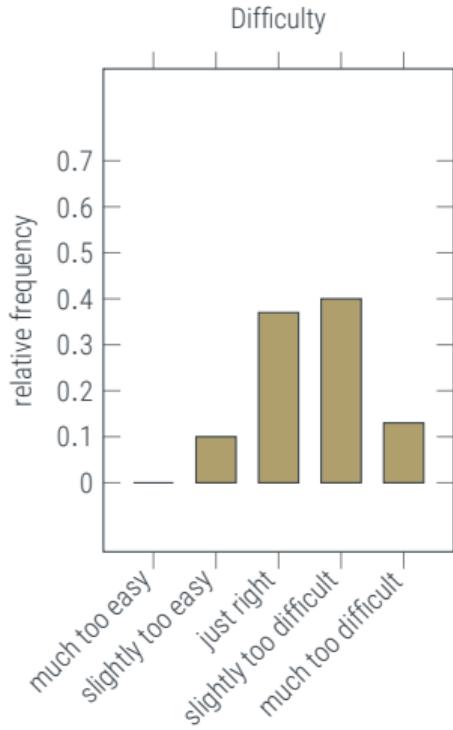
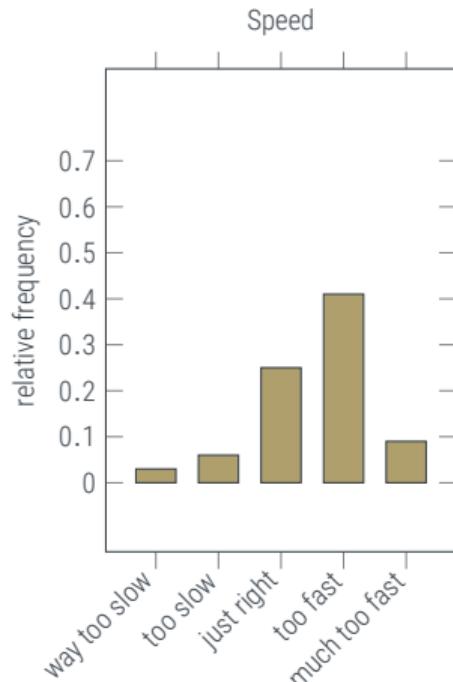
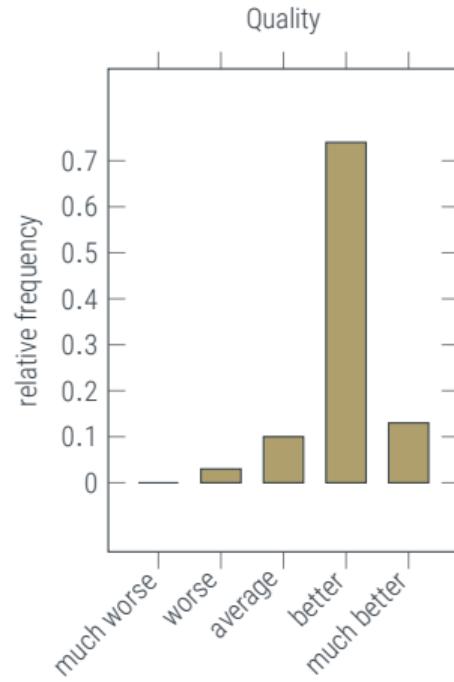


FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING



Last Lecture: Debrief

Feedback dashboard





Last Lecture: Debrief

Detailed Feedback

Things you did not like:

- ♦ use of blackboard !
- ♦ examples
- ♦ plots
- ♦ too fast, more proofs required
- ♦ too easy, we know this already
- ♦ too many equations
- ♦ the pace is too high to take notes on a computer (equations are hard)

Things you did not understand:

- ♦ map from ν to x in Gaussian examples
- ♦ why is $p(x_1 | x_2)$ unnormalizable in x_2 , yet still a Gaussian?
- ♦ proofs
- ♦ what does "explaining away" mean?
- ♦ why is axis aligned \Leftrightarrow independence?

Things you enjoyed:

- ♦ plots / visualizations
- ♦ use of blackboard
- ♦ properties of Gaussians
- ♦ historical notes
- ♦ break
- ♦ examples



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
3. Gaussian Probability Distributions
 - Gaussians map probabilistic inference to **linear algebra**

Today:

- using Gaussian distributions to **learn/infer functions** – the base case of **supervised learning**

Recap: Gaussian Distributions

Gaussian distributions provide the linear algebra of inference

- 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

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

- (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

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

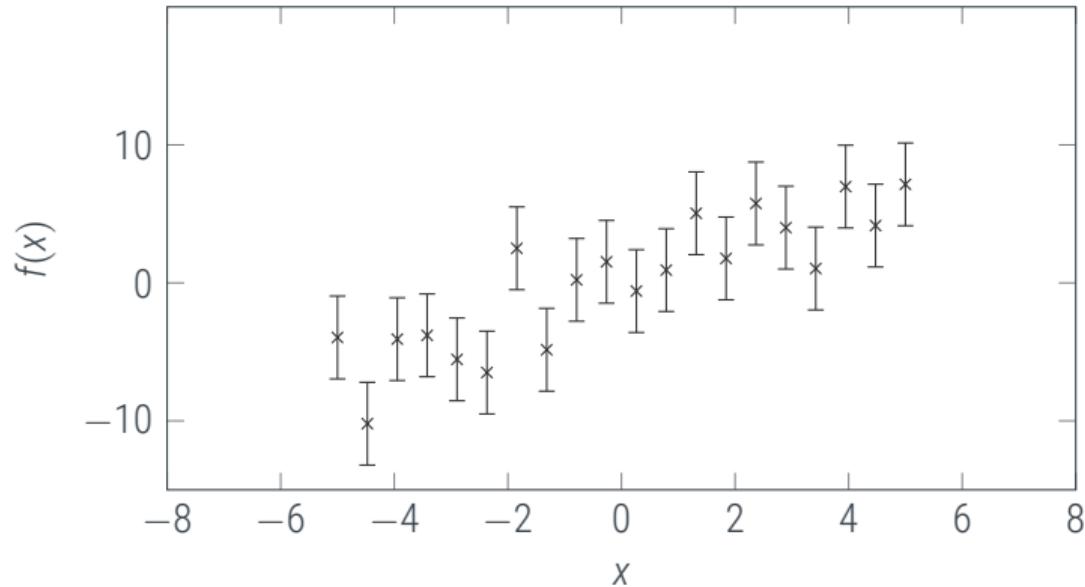
$$p(B^\top x + c | y) = \mathcal{N}[B^\top x + c; B^\top \mu + c + B^\top \Sigma A (A^\top \Sigma A + \Lambda)^{-1} (y - A^\top \mu - b), B^\top \Sigma B - B^\top \Sigma A (A^\top \Sigma A + \Lambda)^{-1} A^\top \Sigma B]$$



Supervised Regression

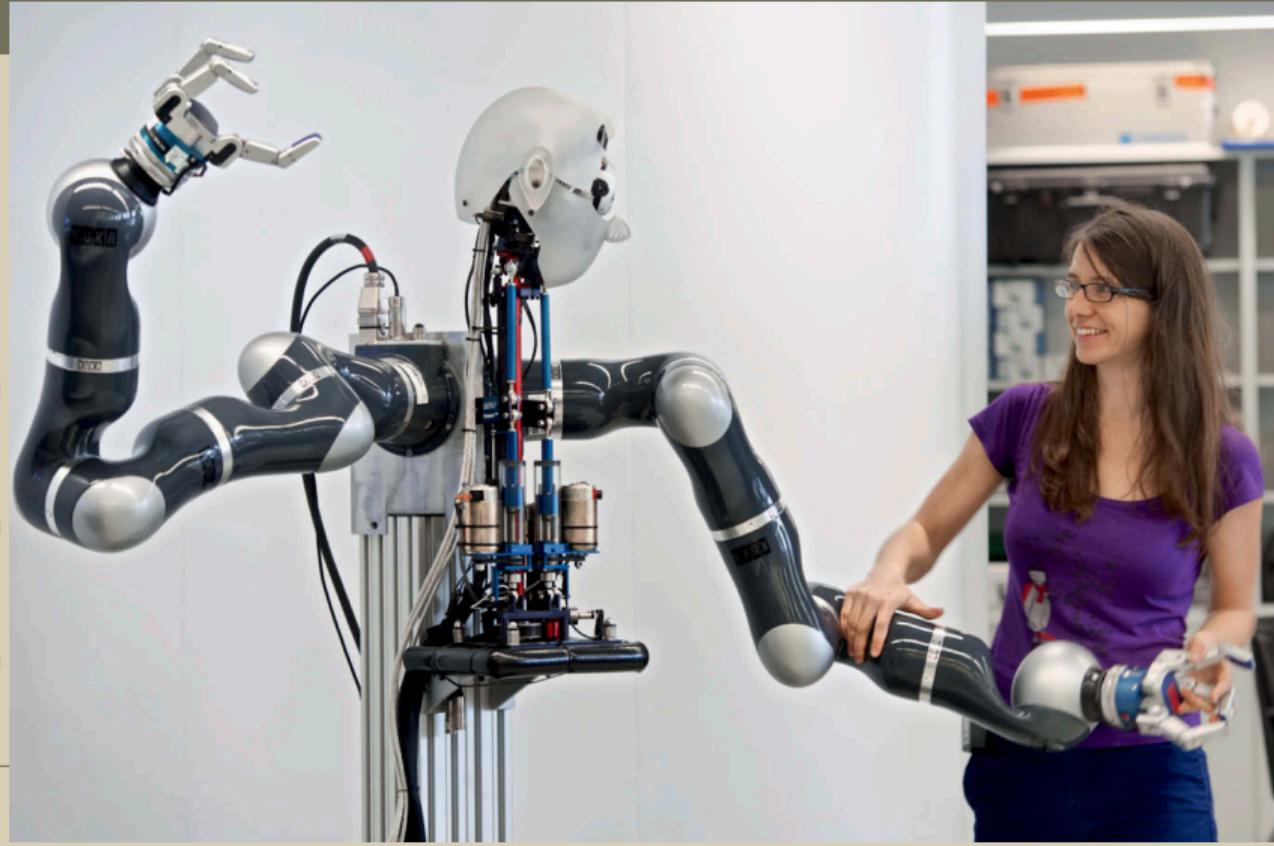
A data set

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

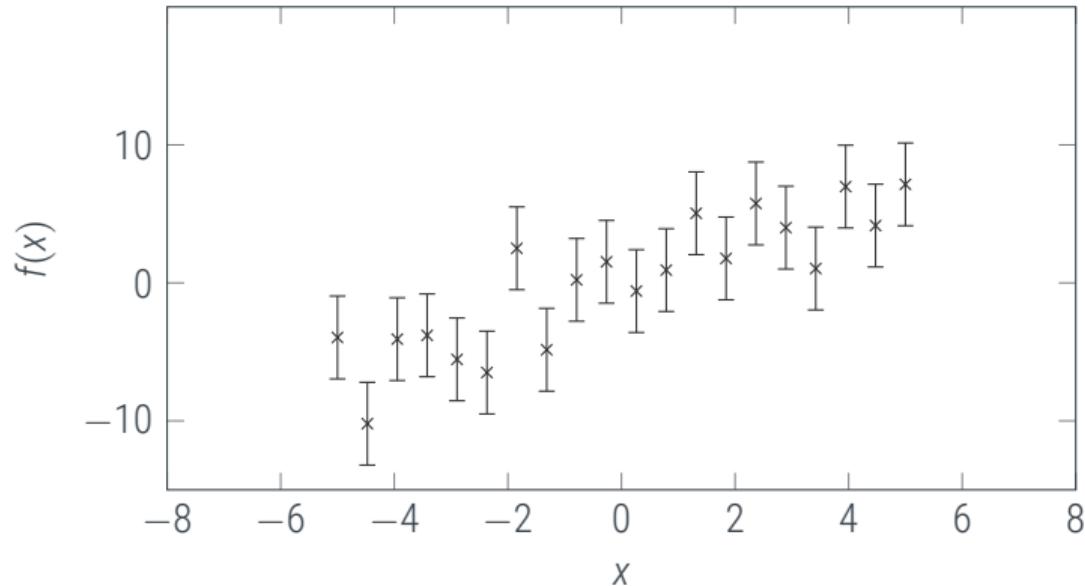




Supervised Regression

A data set

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

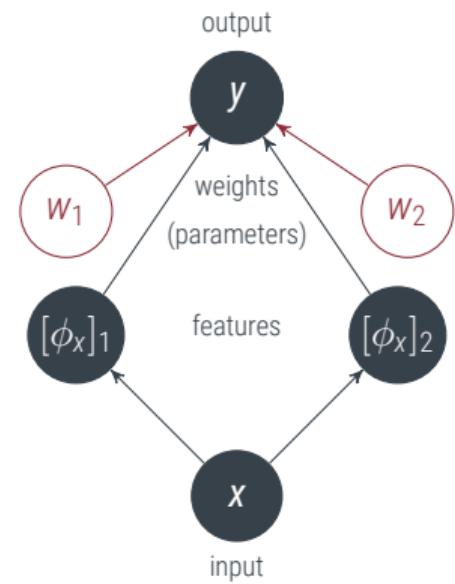
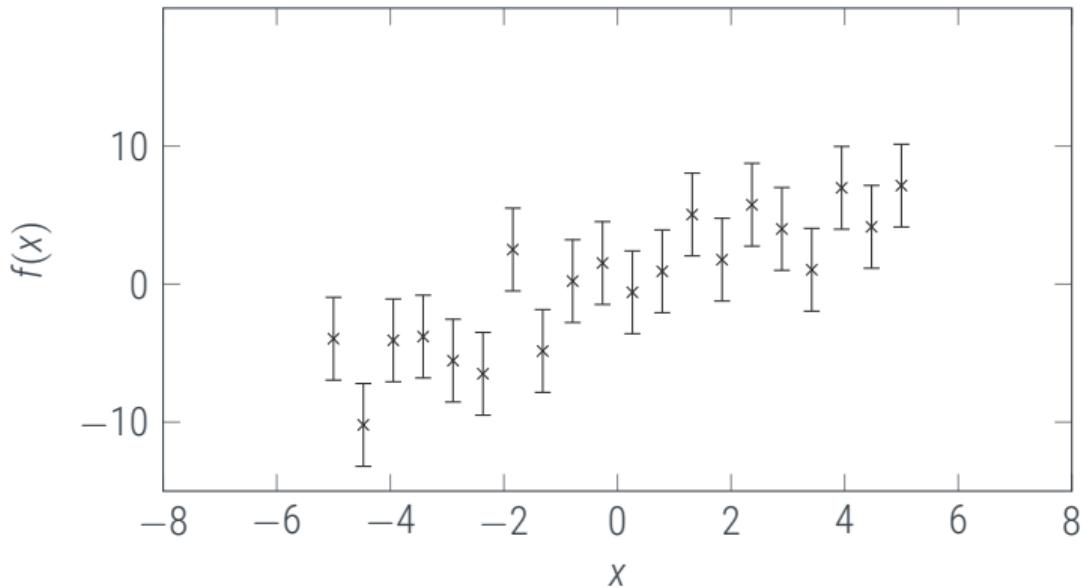




A Linear Model

linear regression

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



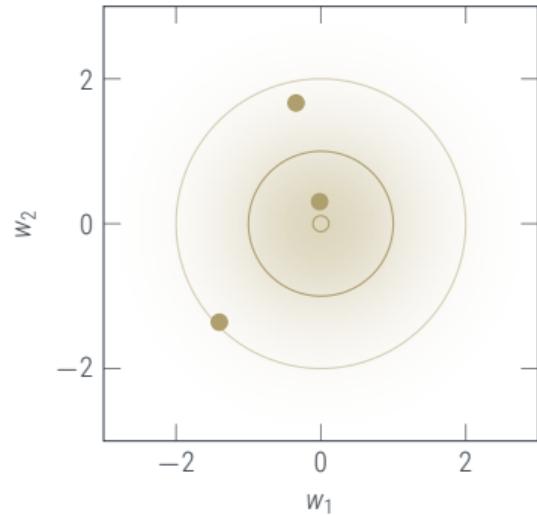


A linear generative model

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

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

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



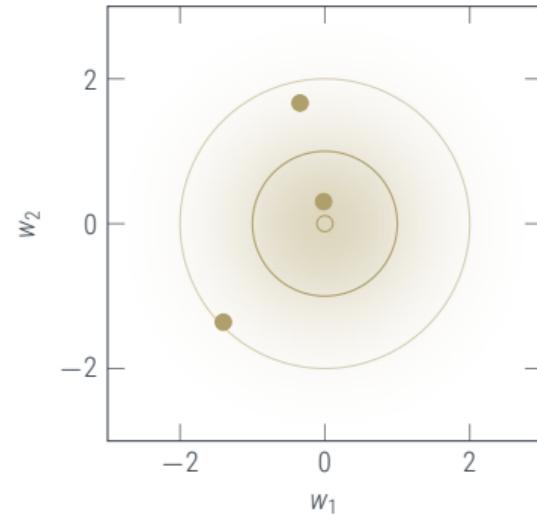


A linear generative model

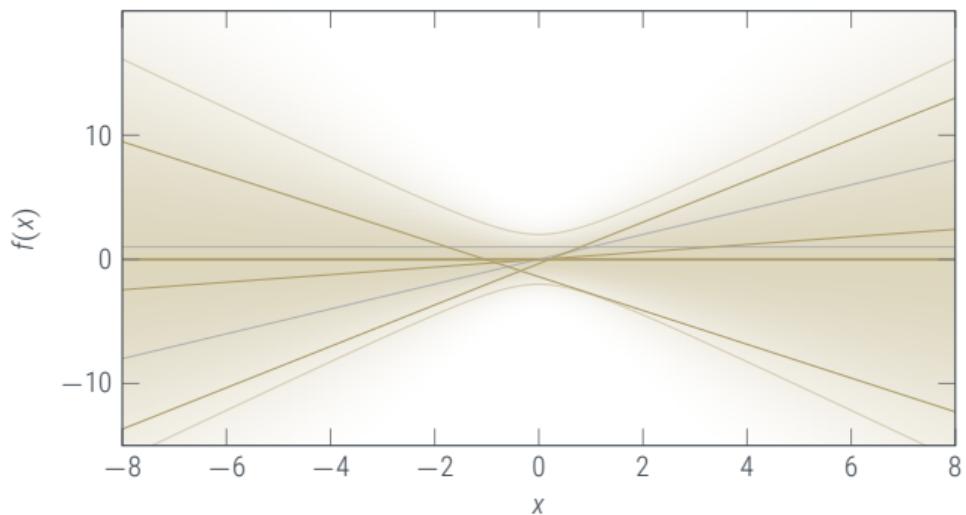
if every variable is Gaussian and every relationship is linear, all marginals and conditionals are also Gaussian

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

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



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





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

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

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

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

$$\begin{aligned} \text{prior} \quad p(w) &= \mathcal{N}(w; \mu, \Sigma) \quad \Rightarrow \quad p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x) \\ \text{likelihood} \quad p(y | w, \phi_x) &= \mathcal{N}(y; \phi_x^\top w, \sigma^2 I) = \mathcal{N}(y; f_x, \sigma^2 I) \end{aligned}$$

Gaussian Inference on a linear function

weight space / function space

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

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

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



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

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

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

Aside: Sampling from a Gaussian

Part I: Sampling from a Standard Gaussian – The Box-Muller Transform



[Box & Muller, 1958; also Paley & Wiener, 1934]

Illustration: (u/Cmglee, CC BY-SA 3.0)

https://upload.wikimedia.org/wikipedia/commons/1/1f/Box-Muller_transform_visualisation.svg

```
1 procedure RANDN(n)
2      $U_0 \leftarrow \text{RAND}$                                 // draw first unit random
3     for i=1,...,n do
4          $U_i \leftarrow \text{RAND}$                             // draw unit random
5          $r \leftarrow \sqrt{-2 \ln U_{i-1}}$                   // compute radius
6          $\theta \leftarrow 2\pi U_i$                           // compute angle
7         if  $i \% 2 == 0$  then
8              $Z_i \leftarrow r \cos \theta$                     // transform to Euclidean
9         else
10             $Z_i \leftarrow r \sin \theta$                     // transform to Euclidean
11        end if
12    end for
13    return  $z \sim \mathcal{N}(0, I_n)$ 
14 end procedure
```

Improvements

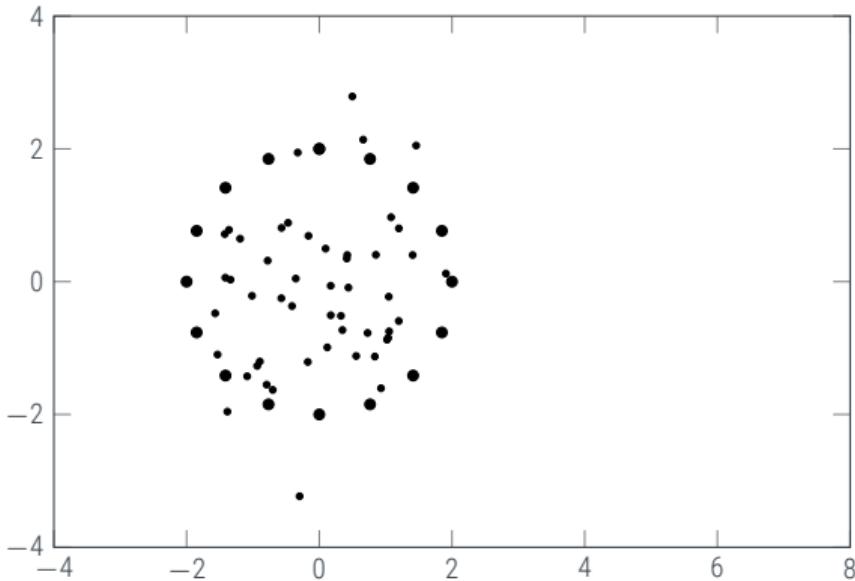
- Marsaglia's polar form: draw from unit circle instead of square (using rejection sampling). This avoids sin and cos, but requires rejection. Typically faster
- Marsaglia's Ziggurat algorithm: iteratively refined rejection sampling. Even faster, but tough to understand intuitively.

Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
U,D = eig(Sigma)
```

```
x = dot(dot(U,diag(sqrt(D))),randn(n,3)) + mu
```



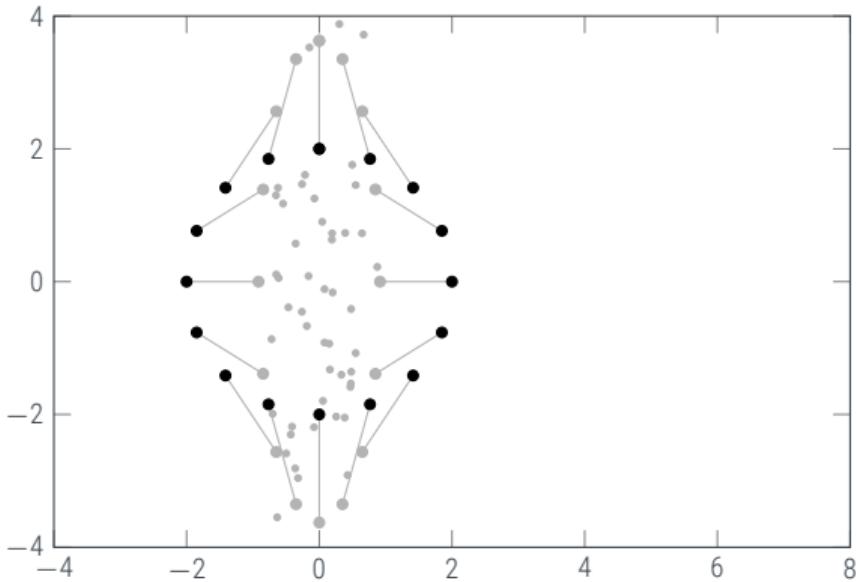
$$\Sigma = UDU^\top \quad x = u$$

Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
U,D = eig(Sigma)
```

```
x = dot(dot(U,diag(sqrt(D))),randn(n,3)) + mu
```



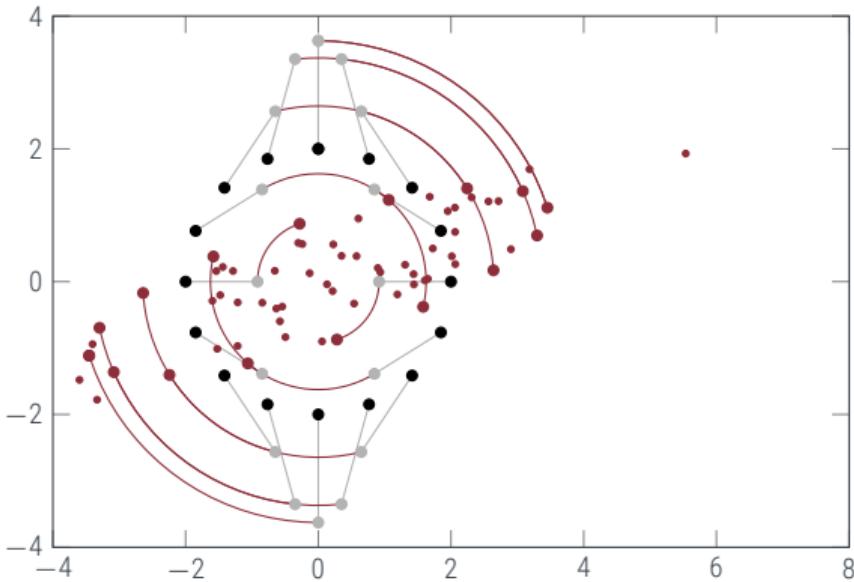
$$\Sigma = UDU^\top \quad x = D^{1/2}u$$

Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
U,D = eig(Sigma)
```

```
x = dot(dot(U,diag(sqrt(D))),randn(n,3)) + mu
```



$$\Sigma = UDU^\top \quad x = UD^{1/2}u$$

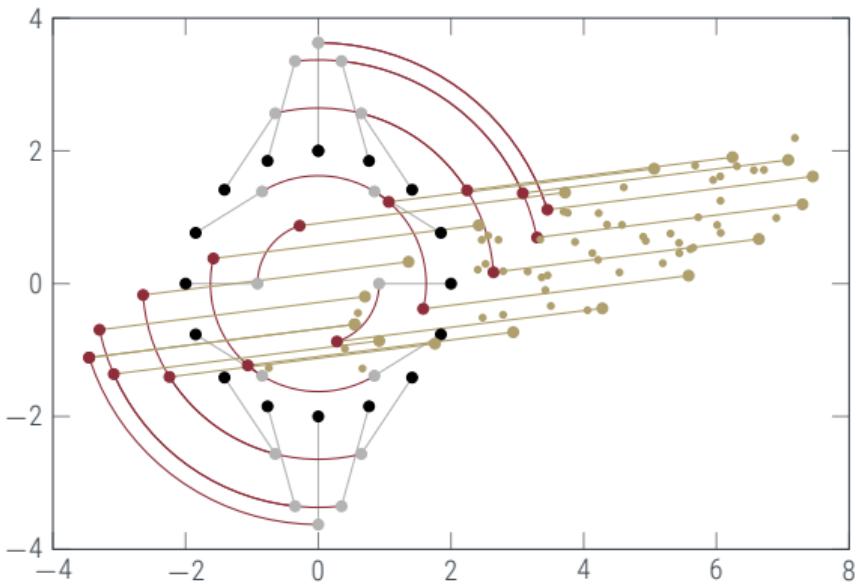


Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
U,D = eig(Sigma)
```

```
x = dot(dot(U,diag(sqrt(D))),randn(n,3)) + mu
```



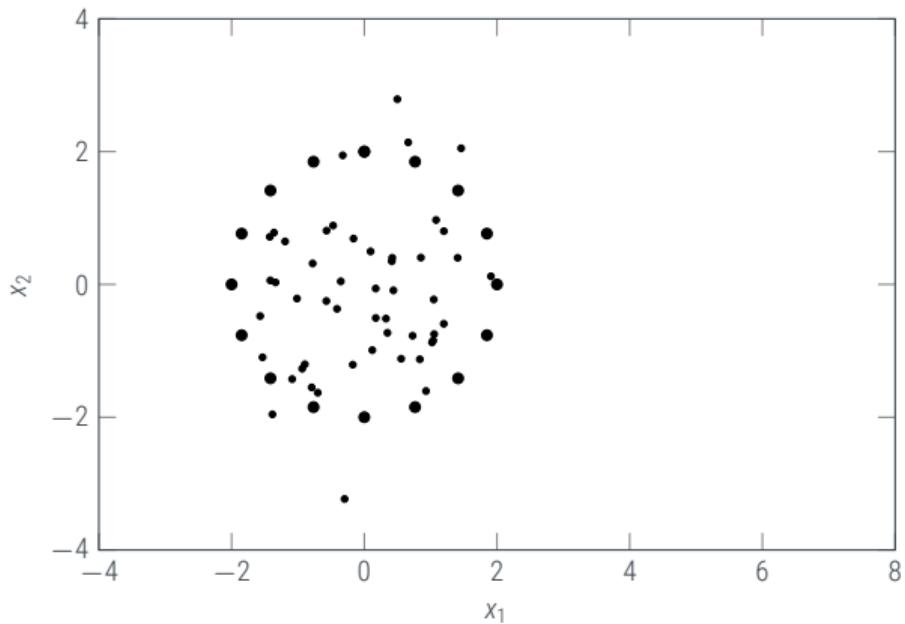
$$\Sigma = UDU^\top \quad x = UD^{1/2}u + \mu$$



Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
x = dot(cholesky(Sigma), randn(n,3)) + mu
```



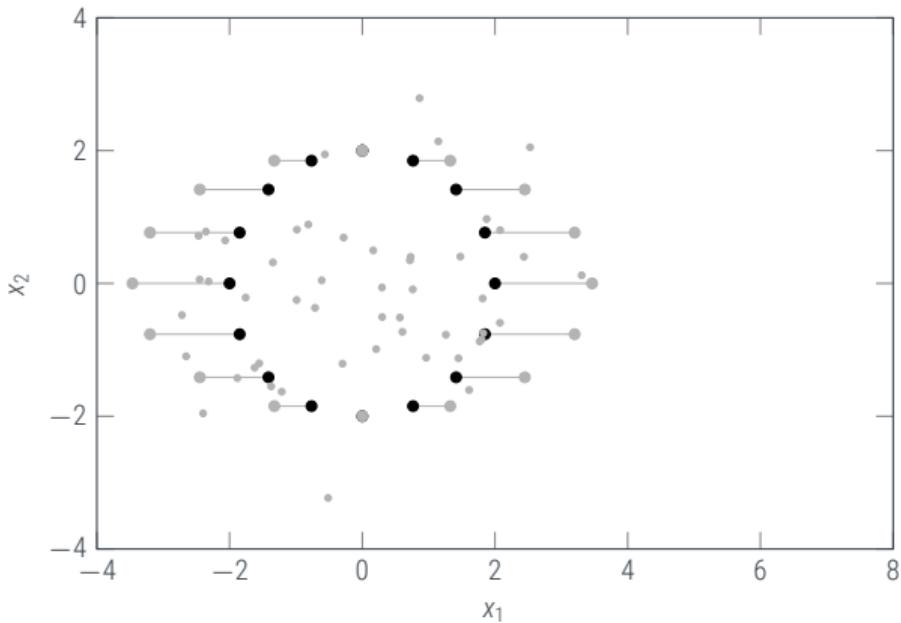
$$x = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$



Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
x = dot(cholesky(Sigma), randn(n,3)) + mu
```

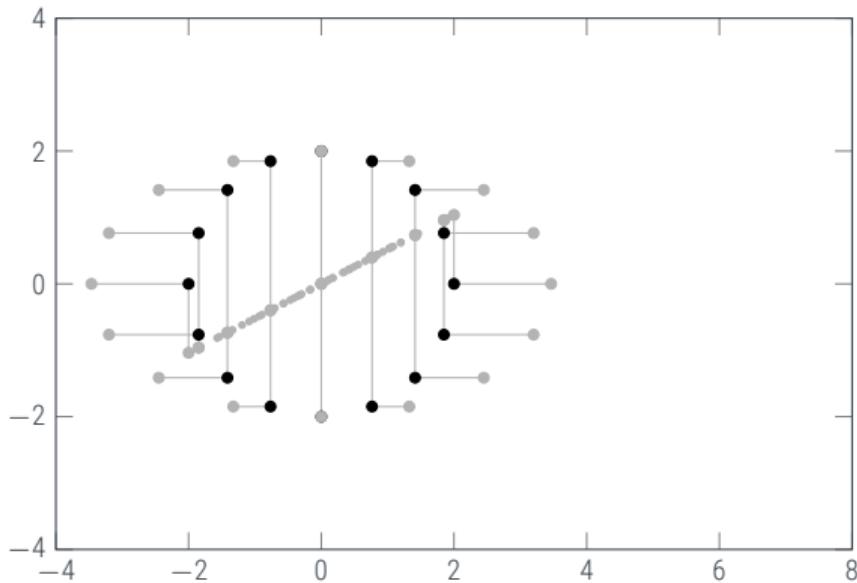


$$x = \begin{pmatrix} \sqrt{\Sigma_{11}} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
x = dot(cholesky(Sigma), randn(n,3)) + mu
```

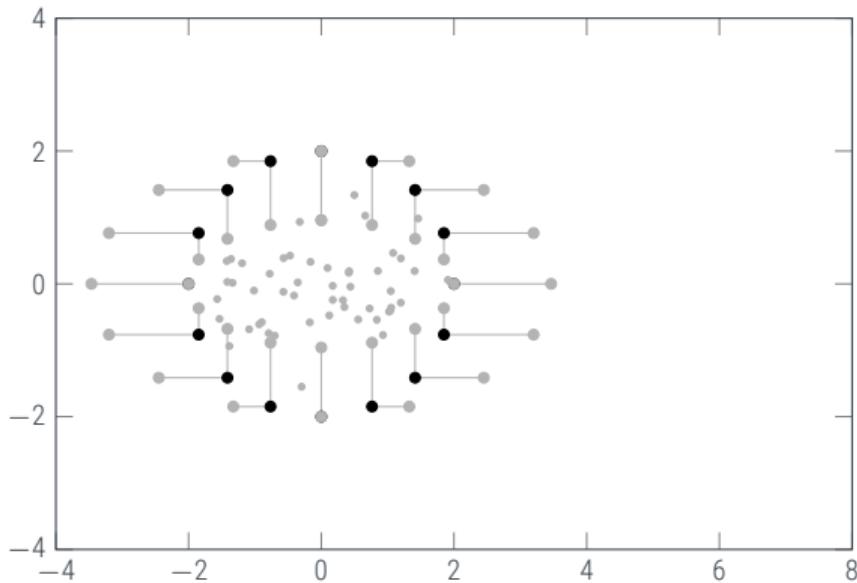


$$x = \begin{pmatrix} \sqrt{\Sigma_{11}} & 0 \\ \Sigma_{11}/\sqrt{\Sigma_{12}} & 0 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
x = dot(cholesky(Sigma), randn(n,3)) + mu
```



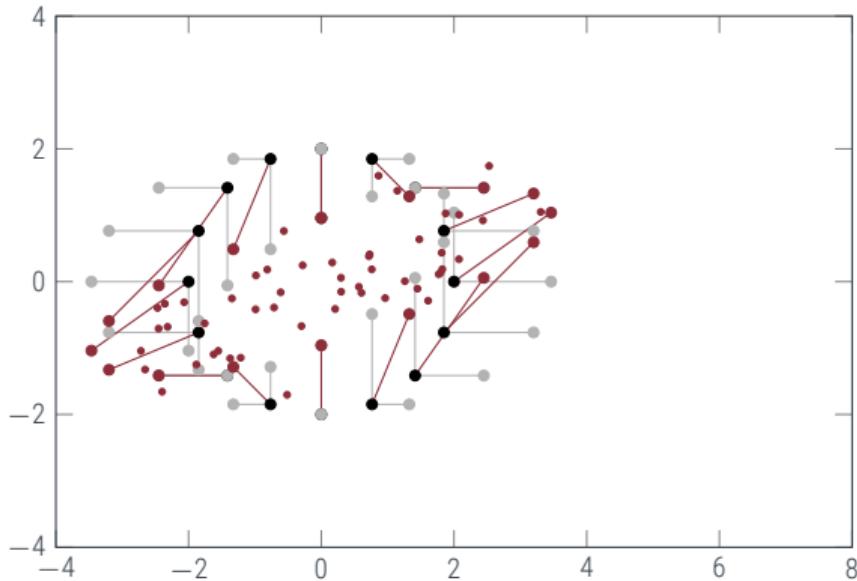
$$x = \begin{pmatrix} \sqrt{\Sigma_{11}} & 0 \\ 0 & \sqrt{\Sigma_{22} - \Sigma_{12}/\Sigma_{22}} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$



Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
x = dot(cholesky(Sigma), randn(n,3)) + mu
```

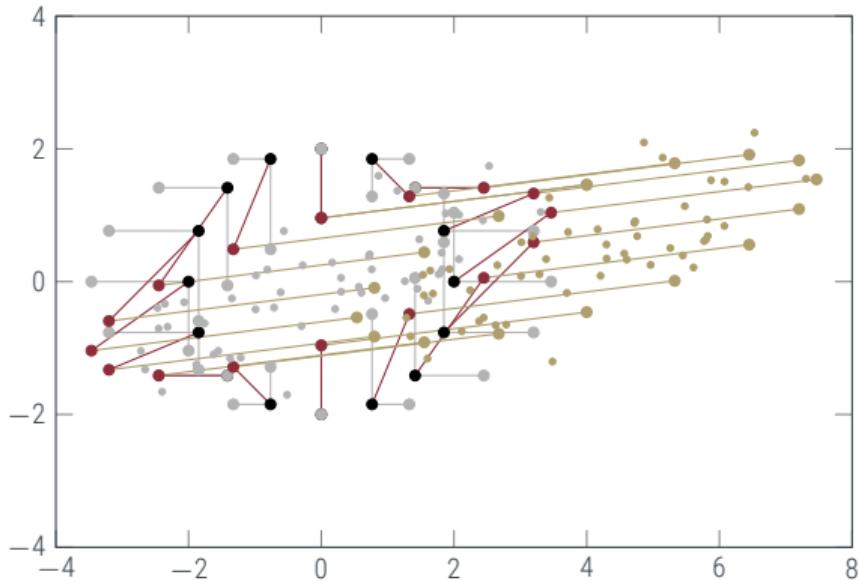


$$x = \begin{pmatrix} \sqrt{\Sigma_{11}} & 0 \\ \Sigma_{11}/\sqrt{\Sigma_{12}} & \sqrt{\Sigma_{22} - \Sigma_{12}/\Sigma_{11}} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

Aside: Sampling from a Gaussian

Part II: Mapping to non-standard mean and covariance

```
x = dot(cholesky(Sigma), randn(n,3)) + mu
```



$$x = \begin{pmatrix} \sqrt{\Sigma_{11}} & 0 \\ \Sigma_{11}/\sqrt{\Sigma_{12}} & \sqrt{\Sigma_{22} - \Sigma_{12}/\Sigma_{22}} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} + \mu$$

- cost: $n^3/3$ FLOPS, about 1/4 of Eigenvalue decomposition
- error bounded by $\|\Sigma\|_2$ (sqrt of largest eigenvalue) Golub & van Loan, 2013, 4.2.6

That's it for the Algebra

Gaussian Inference on a linear function

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

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

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

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



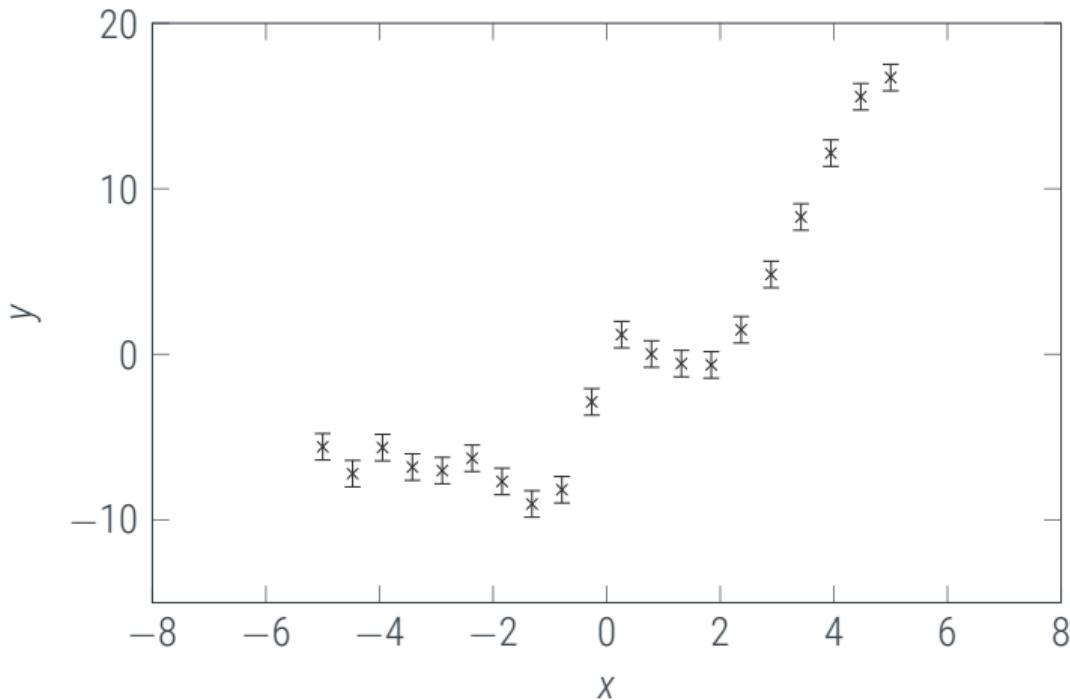
Code

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

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



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



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

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



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



Laguerre Regression

$$\phi(x) = \left[\frac{1}{n!} \left(\frac{d}{dx} - 1 \right)^n x^n \right]_{n=1,\dots,8}^{\mathsf{T}}$$



Laguerre Regression

$$\phi(x) = \left[\frac{1}{n!} \left(\frac{d}{dx} - 1 \right)^n x^n \right]_{n=1,\dots,8}^{\mathsf{T}}$$



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



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



Bell Curve Regression

$$\phi(x) = \begin{bmatrix} e^{-\frac{1}{2}(x-8)^2} & e^{-\frac{1}{2}(x-7)^2} & e^{-\frac{1}{2}(x-6)^2} & \dots \end{bmatrix}^\top$$



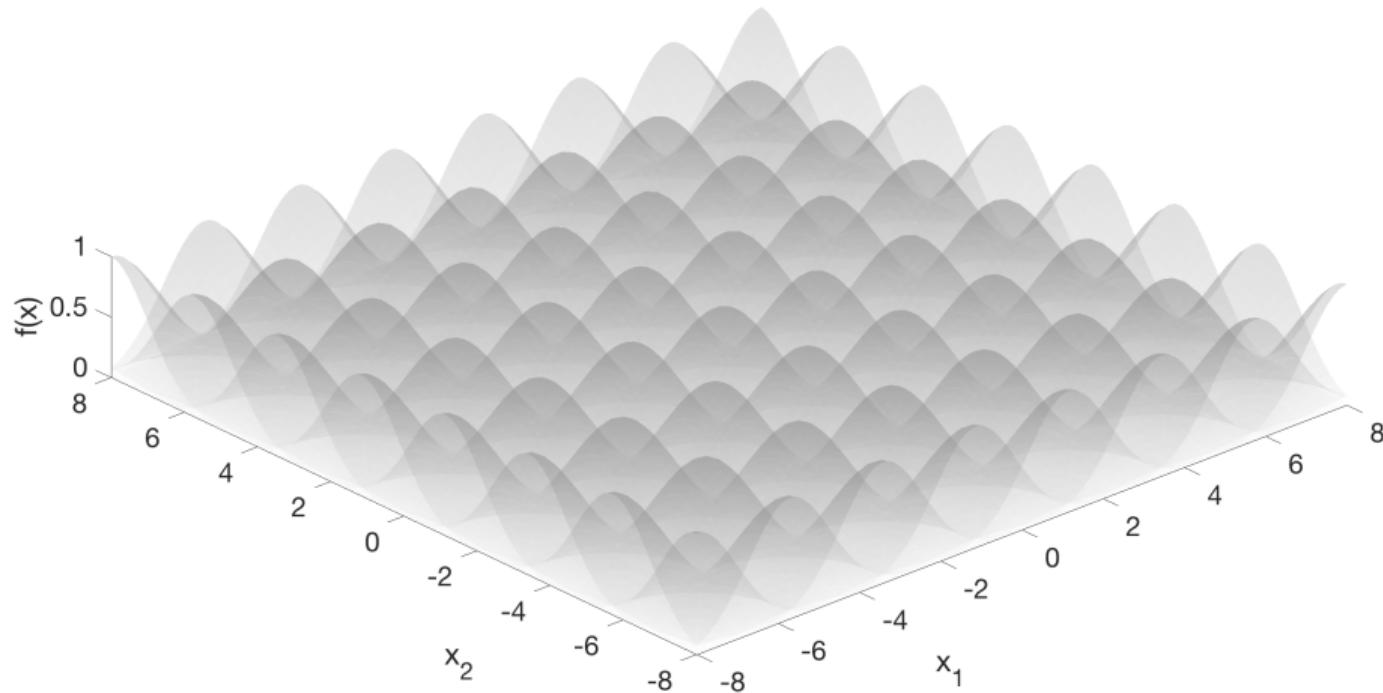
Bell Curve Regression

$$\phi(x) = \begin{bmatrix} e^{-\frac{1}{2}(x-8)^2} & e^{-\frac{1}{2}(x-7)^2} & e^{-\frac{1}{2}(x-6)^2} & \dots \end{bmatrix}^\top$$



Multiple Inputs

Input domain \mathbb{X} can be anything



Multiple Inputs

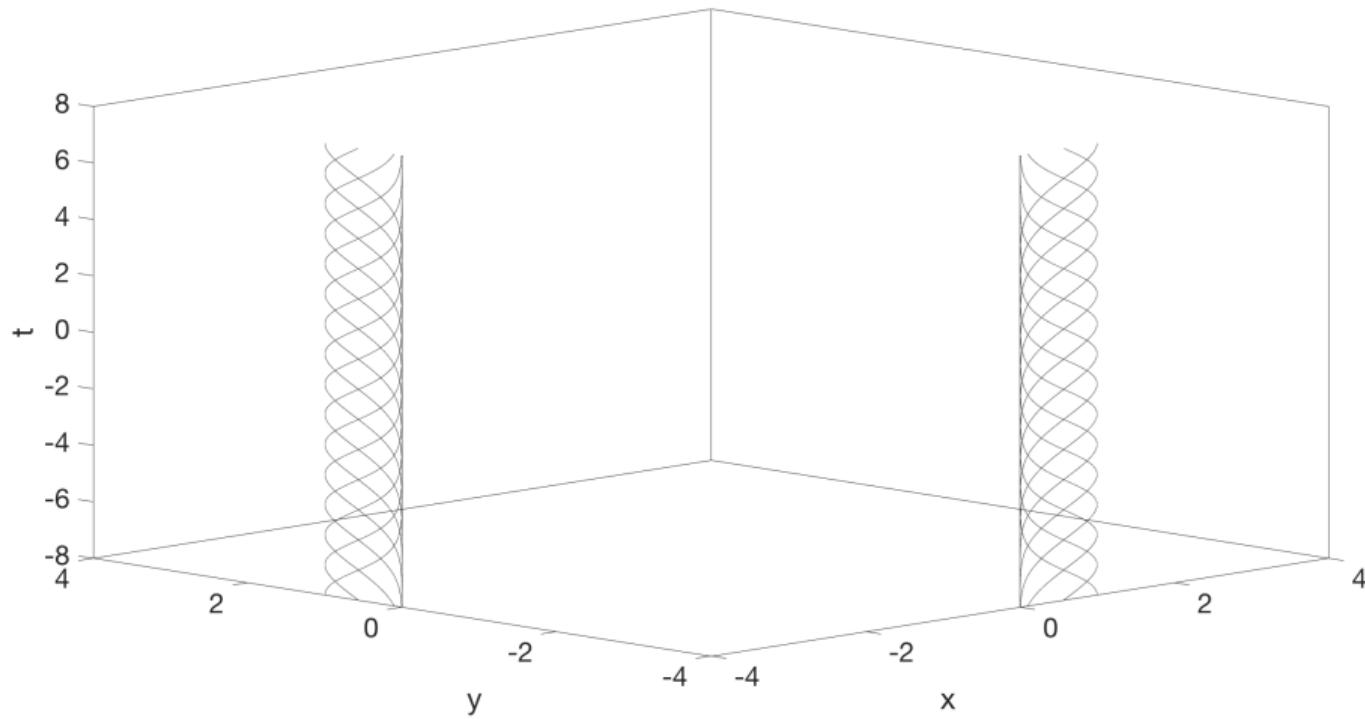
Input domain \mathbb{X} can be anything





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