

# PROBABILISTIC MACHINE LEARNING

## LECTURE 08

### GAUSSIAN PROCESSES

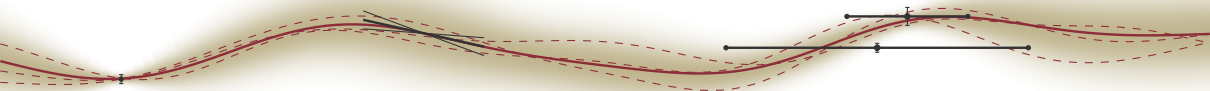
Philipp Hennig

15 May 2023

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



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



## Summary:

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

$$f(x) = \phi(x)^{\top} \mathbf{w} = \phi_x^{\top} \mathbf{w}$$

- ▶ The choice of features  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  is essentially unconstrained



$$f(x) := \phi_x^\top \mathbf{w} \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad p(\mathbf{y} | f) = \mathcal{N}(\mathbf{y}; f_X, \sigma^2 I)$$

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

$$= \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \quad \text{where}$$

$$\underbrace{\boldsymbol{\mu}_y}_{\mathbb{E}(\mathbf{w}|\mathbf{y})} := \underbrace{\boldsymbol{\mu}}_{\mathbb{E}(\mathbf{w})} + \underbrace{\boldsymbol{\Sigma} \phi_X}_{\text{COV}(\mathbf{w}, \mathbf{y})} \underbrace{(\phi_X^\top \boldsymbol{\Sigma} \phi_X + \sigma^2)^{-1}}_{\text{COV}(\mathbf{y}, \mathbf{y})^{-1}} \underbrace{(\mathbf{y} - \phi_X^\top \boldsymbol{\mu})}_{\mathbf{y} - \mathbb{E}(\mathbf{y})} \quad \text{and}$$

$$\underbrace{\boldsymbol{\Sigma}_y}_{\text{COV}(\mathbf{w}, \mathbf{w}|\mathbf{y})} := \underbrace{\boldsymbol{\Sigma}}_{\text{COV}(\mathbf{w}, \mathbf{w})} - \underbrace{\boldsymbol{\Sigma} \phi_X}_{\text{COV}(\mathbf{w}, \mathbf{y})} \underbrace{(\phi_X^\top \boldsymbol{\Sigma} \phi_X + \sigma^2)^{-1}}_{\text{COV}(\mathbf{y}, \mathbf{y})^{-1}} \underbrace{\phi_X^\top \boldsymbol{\Sigma}}_{\text{COV}(\mathbf{y}, \mathbf{w})}$$

$$p(f_x | \mathbf{y}, \phi_X) = \mathcal{N}(f_x; \phi_X^\top \boldsymbol{\mu}_y, \phi_X^\top \boldsymbol{\Sigma}_y \phi_X)$$

$$= \mathcal{N}\left(f_x; \phi_X^\top \boldsymbol{\mu} + \phi_X^\top \boldsymbol{\Sigma} \phi_X (\phi_X^\top \boldsymbol{\Sigma} \phi_X + \sigma^2 I)^{-1} (\mathbf{y} - \phi_X^\top \boldsymbol{\mu}),\right.$$

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

# Some Things to Note

## Probabilistic Parametric Regression

$$\begin{aligned}
 f(x) &:= \phi_x^\top w & p(w) &= \mathcal{N}(w; \mu, \Sigma) & p(y | f) &= \mathcal{N}(y; f_x, \sigma^2 l) \\
 p(w | y, \phi_x) &= \mathcal{N}\left(w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} (y - \phi_x^\top \mu), \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} \phi_x^\top \Sigma\right) \\
 &= \mathcal{N}(w; \mu_y, \Sigma_y) & \text{where} \\
 \mu_y &:= \mu + \underbrace{\Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1}}_{\nabla_y \mu_y} (y - \phi_x^\top \mu) & \text{and} \\
 \Sigma_y &:= \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} \phi_x^\top \Sigma. & \text{Thus,} \\
 &= \Sigma - \Sigma \phi_x (\nabla_y \mu_y)^\top
 \end{aligned}$$

Computing the posterior covariance is **cheap**! It measures the *remaining capacity* in the model after seeing the data. **If you can auto-diff, you can be uncertain!**

# Shallow Learning

Gaussian linear regressors are single-layer neural networks with quadratic loss

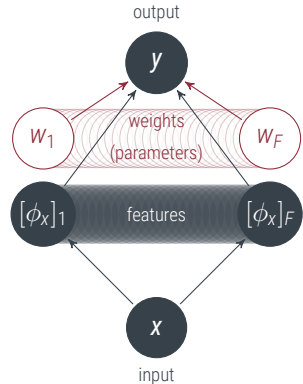
$$p(\mathbf{w} \mid \mathbf{y}, X) = \frac{1}{p(\mathbf{y})} \mathcal{N}(\mathbf{y} \mid \phi_X^\top \mathbf{w}, \sigma I) \mathcal{N}(\mathbf{w}; \mu, \Sigma)$$

$$-\log p(\mathbf{w} \mid \mathbf{y}, X) = -\log \mathcal{N}(\mathbf{y} \mid \phi_X^\top \mathbf{w}, \sigma I) - \log \mathcal{N}(\mathbf{w}; \mu, \Sigma) + \text{const.}$$

$$\log p(\mathbf{w} \mid \mathbf{y}, X) \triangleq \frac{1}{2\sigma^2} \sum_i \|y_i - \phi(x_i)^\top \mathbf{w}\|^2 + \frac{1}{2} (\mathbf{w} - \mu)^\top \Sigma^{-1} (\mathbf{w} - \mu)$$

$$\begin{aligned} \nabla \log p(\mathbf{w} \mid \mathbf{y}, X) &= \frac{1}{\sigma^2} \phi_X (\mathbf{y} - \phi_X^\top \mathbf{w}) - \Sigma^{-1} (\mathbf{w} - \mu) \\ &= (\sigma^{-2} \phi_X \phi_X^\top + \Sigma^{-1}) \mathbf{w} - \sigma^{-2} \phi_X \mathbf{y} - \Sigma^{-1} \mu \stackrel{!}{=} 0 \\ \mathbf{w}_* &= (\sigma^{-2} \phi_X \phi_X^\top + \Sigma^{-1})^{-1} (\sigma^{-2} \phi_X \mathbf{y} + \Sigma^{-1} \mu) \end{aligned}$$

$$\nabla \nabla^\top \log p(\mathbf{w} \mid \mathbf{y}, X) = \sigma^{-2} \phi_X \phi_X^\top + \Sigma^{-1}$$



$$\begin{aligned}
 p(\mathbf{w} \mid \mathbf{y}, X) &= \frac{1}{p(\mathbf{y})} \mathcal{N}(\mathbf{y} \mid \phi_X^\top \mathbf{w}, \sigma^2) \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
 &= \mathcal{N}(\mathbf{w}; \mathbf{w}_* := \arg \min \mathcal{L}(\mathbf{y}, X), \nabla \nabla^\top \mathcal{L}(\mathbf{y}, X))
 \end{aligned}$$

- ▶ Gaussian linear regression is training a *single-layer* neural network with quadratic loss
- ▶ The minimum-loss weights of the neural network are the *posterior mean* (MAP estimate) of the Gaussian linear regression
- ▶ The curvature of the loss function is the *posterior covariance* of the Gaussian linear regression

We now have two options:

1. Add more features to the single layer
2. Add more layers / train the layers



# Can we learn *functions* abstractly?

An observation from last lecture

We currently have

---

```

1 @dataclasses.dataclass
2 class Gaussian:
3     # Gaussian distribution with mean mu and covariance Sigma
4     mu: jnp.ndarray # shape (D,)
5     Sigma: jnp.ndarray # shape (D,D)

```

---

But we're trying to learn **functions**  $y \approx f(x)$  from observations  $(X, y) \in (\mathbb{X} \times \mathbb{R})^n$

$$p(y \mid f(X)) = \mathcal{N}(y; f_X, \sigma^2 I) = \prod_{i=1}^n \mathcal{N}(y_i; f(x_i), \sigma^2)$$

So can we do something like this?

---

```

1 @dataclasses.dataclass
2 class GaussianFunction:
3     # mean function
4     m: Callable[[jnp.ndarray], jnp.ndarray] # maps |X|^n -> |R|^n
5     # covariance function
6     k: Callable[[jnp.ndarray, jnp.ndarray], jnp.ndarray] # maps (|X| x |X|)^n -> |R|^{n x n}

```

---

Note that “in practice”, functions are only defined through their values at finitely many concrete points

---

```
1 @dataclasses.dataclass
2 class GaussianProcess:
3     # mean function
4     m: Callable[[jnp.ndarray], jnp.ndarray]
5     # covariance function
6     k: Callable[[jnp.ndarray, jnp.ndarray], jnp.ndarray]
7
8     def __call__(self, x):
9         return Gaussian(mu=self.m(x), Sigma=self.k(x[:, None, :], x[None, :, :]))
```

---

What is  $k$ ?





# 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]



J. Carl F. Gauss, 1777 Brunswick-1855 Göttingen

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 Ganzes 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  $\log \varphi A = \frac{1}{2} k A A + \text{Const.}$ , oder wenn man die Basis der hyperbolischen Logarithmen mit  $e$  bezeichnet und die Constante  $= \log x$  setzt,

$$\varphi A = x e^{\frac{1}{2} k A A}.$$

Ferner sieht man leicht ein, dass  $k$  nothwendig negativ sein müsse, damit  $\Omega$  in der That ein Grösstes werden könne, weshalb wir setzen  $\frac{1}{2} k = -h h$ ; und da vermittelt des eleganten, zuerst von Laplace\*) gefundenen Theorems das Integral  $\int e^{-h h A A} dA$ , 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:

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

178.

Die so eben ermittelte Function kann zwar nicht in aller Strenge die Wahrscheinlichkeiten der Fehler ausdrücken; denn da die möglichen Fehler (213) stets in gewisse Grenzen eingezwängt 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 küsset Gauss: „Das Euler schon das Theorem gefunden hat, wovon der schön, von mir Laplace beigelegte Lehrsatz sehr leicht abgeleitet werden kann, den mir selbst schon früher ein, als aber die Stelle S. 212 schon abgedruckt war; ich wollte es aber nicht unter die Errata setzen, weil Laplace wenigstens das obige Theorem doch erst in der dort gebrauchten Form aufgestellt hat.“ *Anmerkung des Uebersetzers.*

## Definition (kernel)

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a (Mercer / positive definite) **kernel** if, for any finite collection  $X = [x_1, \dots, x_N]$ , the matrix  $k_{XX} \in \mathbb{R}^{N \times N}$  with  $[k_{XX}]_{ij} = k(x_i, x_j)$  is **positive semidefinite**.

## Definition (positive definite matrix)

A symmetric matrix  $A \in \mathbb{R}^{N \times N}$  is called **positive (semi-) definite** if

$$v^T A v \geq 0 \quad \forall v \in \mathbb{R}^N.$$

Equivalently:

- ▶ All eigenvalues of the symmetric matrix  $A$  are non-negative
- ▶  $A$  is a Gram matrix – the outer product of  $N$  vectors  $[\phi_i]_{i=1, \dots, N}$

# Gaussian processes

are programs that always produce valid Gaussians

## Definition

Let  $\mu : \mathcal{X} \rightarrow \mathbb{R}$  be any function,  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a Mercer kernel. A **Gaussian process**  $p(f) = \mathcal{GP}(f; \mu, k)$  is a probability distribution over the function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , such that every finite restriction to function values  $f_X := [f_{x_1}, \dots, f_{x_N}]$  is a Gaussian distribution  $p(f_X) = \mathcal{N}(f_X; \mu_X, k_{XX})$ .

---

```

1 @dataclasses.dataclass
2 class GaussianProcess:
3     # mean function
4     m: Callable[[jnp.ndarray], jnp.ndarray]
5     # covariance function
6     k: Callable[[jnp.ndarray, jnp.ndarray], jnp.ndarray]
7
8     def __call__(self, x):
9         return Gaussian(mu=self.m(x), Sigma=self.k(x[:, None, :], x[None, :, :]))

```

---

# How do we build a kernel?

let's look at that algebra again

$$\begin{aligned}
 p(f(\bullet) \mid w) &= \mathcal{N}(f(\bullet); \phi(\bullet)^\top w, \sigma I) \\
 p(f(\bullet)) &= \int p(f(\bullet) \mid w) p(w) dw \\
 &= \mathcal{N}(f(\bullet); \phi(\bullet)^\top \mu, \phi(\bullet)^\top \Sigma \phi(\bullet) + \sigma I)
 \end{aligned}$$

using the abstraction / encapsulation

$m(\bullet) := \phi(\bullet)^\top \mu$	$m : \mathbb{X} \rightarrow \mathbb{R}$	mean function
$k(\bullet, \circ) := \phi(\bullet)^\top \Sigma \phi(\circ)$	$k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$	covariance function, aka. <b>kernel</b>

# What are we actually doing with those features?

let's look at that algebra again

$$\begin{aligned}
 p(f_{\bullet} \mid y, \phi_X) &= \mathcal{N}(f_{\bullet}; \phi_{\bullet}^{\top} \mu + \phi_{\bullet}^{\top} \Sigma \phi_X (\phi_X^{\top} \Sigma \phi_X + \sigma^2 I)^{-1} (y - \phi_X^{\top} \mu), \\
 &\quad \phi_{\bullet}^{\top} \Sigma \phi_o - \phi_{\bullet}^{\top} \Sigma \phi_X (\phi_X^{\top} \Sigma \phi_X + \sigma^2 I)^{-1} \phi_X^{\top} \Sigma \phi_o) \\
 &= \mathcal{N}(f_{\bullet}; m_{\bullet} + k_{\bullet X} (k_{XX} + \sigma^2 I)^{-1} (y - m_X), \\
 &\quad k_{\bullet o} - k_{\bullet X} (k_{XX} + \sigma^2 I)^{-1} k_{Xo})
 \end{aligned}$$

using the abstraction / encapsulation

$$m_{\bullet} := \phi_{\bullet}^{\top} \mu$$

$$m : \mathbb{X} \rightarrow \mathbb{R}$$

mean function

$$k_{\bullet o} := \phi_{\bullet}^{\top} \Sigma \phi_o$$

$$k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$$

covariance function, aka. **kernel**

# Inner products are sums

an observation

- In feature-based (i.e. parametric) regression, we have

$$k_{\bullet\circ} := \phi_{\bullet}^{\top} \Sigma \phi_{\circ} = \sum_{i=1}^F \sum_{j=1}^F \phi(\bullet)_i \Sigma_{ij} \phi(\circ)_j$$

- For simplicity, let's fix  $\Sigma = \frac{\sigma^2(c_{\max} - c_{\min})}{F} I$ , i.e.  $\Sigma_{ij} = \frac{\sigma^2(c_{\max} - c_{\min})}{F} \delta_{ij}$

$$\text{thus: } \phi(x_i)^{\top} \Sigma \phi(x_j) = \frac{\sigma^2(c_{\max} - c_{\min})}{F} \sum_{i=1}^F \sum_{j=1}^F \phi(x_i)^2 \delta_{ij} \phi(x_j)^2 = \frac{\sigma^2(c_{\max} - c_{\min})}{F} \sum_{\ell=1}^F \phi_{\ell}(x_i) \phi_{\ell}(x_j)$$

- Sometimes sums can be finite even if they contain infinitely many terms...

$$\sum_{i=0}^{\infty} x^i = \frac{1}{1-x} \quad \text{for } |x| < 1$$

# Sometimes, **more** features make things **cheaper**

an example



[DJC MacKay, 1998]

- For simplicity, let's fix  $\Sigma = \frac{\sigma^2(c_{\max} - c_{\min})}{F}I$ , i.e.  $\Sigma_{ij} = \frac{\sigma^2(c_{\max} - c_{\min})}{F}\delta_{ij}$

$$\text{thus: } \phi(x_i)^\top \Sigma \phi(x_j) = \frac{\sigma^2(c_{\max} - c_{\min})}{F} \sum_{i=1}^F \sum_{j=1}^F \phi(x_i)^2 \delta_{ij} \phi(x_j)^2 = \frac{\sigma^2(c_{\max} - c_{\min})}{F} \sum_{\ell=1}^F \phi_\ell(x_i) \phi_\ell(x_j)$$

- especially, for  $\phi_\ell(x) = \exp\left(-\frac{(x - c_\ell)^2}{2\lambda^2}\right)$

$$\begin{aligned} \phi(x_i)^\top \Sigma \phi(x_j) &= \frac{\sigma^2(c_{\max} - c_{\min})}{F} \sum_{\ell=1}^F \exp\left(-\frac{(x_i - c_\ell)^2}{2\lambda^2}\right) \exp\left(-\frac{(x_j - c_\ell)^2}{2\lambda^2}\right) \\ &= \frac{\sigma^2(c_{\max} - c_{\min})}{F} \exp\left(-\frac{(x_i - x_j)^2}{4\lambda^2}\right) \sum_{\ell} \exp\left(-\frac{(c_\ell - \frac{1}{2}(x_i + x_j))^2}{\lambda^2}\right) \end{aligned}$$

# Sometimes, **more** features make things **cheaper**

an example



[DJC MacKay, 1998]

$$\phi(x_i)^\top \Sigma \phi(x_j) = \frac{\sigma^2(c_{\max} - c_{\min})}{F} \exp\left(-\frac{(x_i - x_j)^2}{4\lambda^2}\right) \sum_{\ell}^F \exp\left(-\frac{(c_{\ell} - \frac{1}{2}(x_i + x_j))^2}{\lambda^2}\right)$$

now increase  $F$  so # of features in  $\delta c$  approaches  $\frac{F \cdot \delta c}{(c_{\max} - c_{\min})}$

$$\phi(x_i)^\top \Sigma \phi(x_j) \rightarrow \sigma^2 \exp\left(-\frac{(x_i - x_j)^2}{4\lambda^2}\right) \int_{c_{\min}}^{c_{\max}} \exp\left(-\frac{(c - \frac{1}{2}(x_i + x_j))^2}{\lambda^2}\right) dc$$

let  $c_{\min} \rightarrow -\infty$ ,  $c_{\max} \rightarrow \infty$

$$k(x_i, x_j) := \phi(x_i)^\top \Sigma \phi(x_j) \rightarrow \sqrt{\pi} \lambda \sigma^2 \exp\left(-\frac{(x_i - x_j)^2}{4\lambda^2}\right)$$

$$k_{x_i x_j} = \sqrt{\pi} \lambda \sigma^2 \exp\left(-\frac{(x_i - x_j)^2}{4\lambda^2}\right)$$





# Demo



# Our construction is enough

It is fine to think of them as inner products

**Lemma:**  $k$  is a Mercer kernel if it can be written as (assuming  $\mathcal{L}$  is positive set for  $\nu$ )

$$k(x, x') = \sum_{\ell \in \mathcal{L}} \phi_{\ell}(x) \phi_{\ell}(x') \quad \text{or} \quad = \int_{\mathcal{L}} \phi_{\ell}(x) \phi_{\ell}(x') d\nu(\ell)$$

Proof:  $\forall X \in \mathbb{X}^N, v \in \mathbb{R}^N : v^{\top} k_{XX} v = \int \sum_i^N v_i \phi_{\ell}(x_i) \sum_j^N v_j \phi_{\ell}(x_j) d\nu(\ell) = \int \left[ \sum_i v_i \phi_{\ell}(x_i) \right]^2 d\nu(\ell) \geq 0 \quad \square$

**Theorem (Mercer, 1909):** (precise form in Lecture 11). *If  $k$  is a Mercer kernel, then it can be written as above (with countably many terms).*

All kernels admit a (potentially infinite) feature expansion

## Summary: Gaussian Processes

- ▶ Sometimes it is possible to consider **infinitely many** features at once, by extending from a sum to an integral. This requires some regularity assumption about the features' locations, shape, etc.
- ▶ The resulting **nonparametric model** is known as a **Gaussian process**
- ▶ Inference in GPs is tractable (though at polynomial cost  $\mathcal{O}(N^3)$  in the number  $N$  of datapoints)

Please cite this course, as

```
@techreport{Tuebingen_ProbML23,
  title =
    {Probabilistic Machine Learning},
  author = {Hennig, Philipp},
  series = {Lecture Notes
            in Machine Learning},
  year = {2023},
  institution = {Tübingen AI Center}}
```

