

PROBABILISTIC MACHINE LEARNING

LECTURE 10

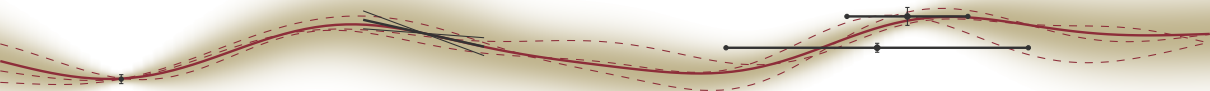
GAUSSIAN PROCESS REGRESSION: AN EXTENSIVE EXAMPLE

Philipp Hennig
25 May 2023

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



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

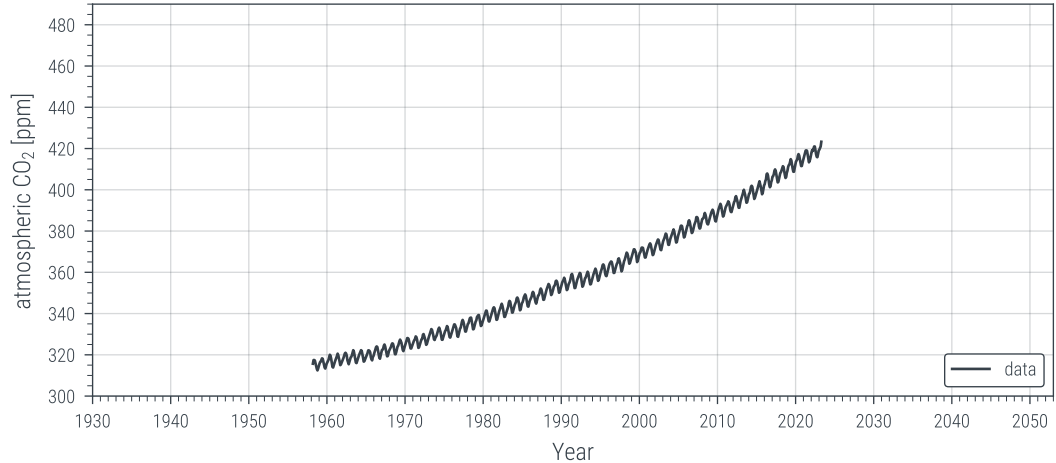


A Dataset

NOAA Mauna Loa CO₂ data



adapted from Rasmussen & Williams, 2007



Additive Kernels and Multi-Output GPs

a sum in the kernel implies the data is emerges from a sum of functions

- Our data represents an (approximately linear) combination of multiple functions:

$$p(f_1) = \mathcal{GP}(f_1; 0, k_1)$$

$$p(f_2) = \mathcal{GP}(f_2; 0, k_2)$$

$$p(f_1, f_2) = \mathcal{GP} \left(\begin{bmatrix} f_1 \\ f_2 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k_1 & 0 \\ 0 & k_2 \end{bmatrix} \right)$$

$$\begin{aligned} p(f) &= \mathcal{GP} \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}; \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} k_1 & 0 \\ 0 & k_2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) \\ &= \mathcal{GP}(f; 0, k_1 + k_2) \end{aligned}$$



Can we **learn** the kernel?

Hierarchical Bayesian Inference

Bayesian model adaptation

$$p(f | y, x, \theta) = \frac{p(y | f, x, \theta)p(f |, \theta)}{\int p(y | f, x, \theta)p(f |, \theta) df} = \frac{p(y | f, x, \theta)p(f |, \theta)}{p(y | x, \theta)}$$

► **Model parameters** like θ are also known as **hyper-parameters**.

► This is largely a computational, practical distinction:

data are observed

→ **condition**

variables are the things we care about

→ **full probabilistic treatment**

parameters are the things we have to deal with to get the model right

→ **integrate out**

hyper-parameters are the top-level, too expensive to properly infer

→ **fit**

The **model evidence** in Bayes' Theorem is the (marginal) **likelihood** for the model. So we would like

$$p(\theta | y) = \frac{p(y | \theta)p(\theta)}{\int p(y | \theta')p(\theta') d\theta'}$$



Can We Learn the Kernel?

Bayesian Hierarchical Inference

$$p(f | \theta) = \mathcal{GP}(f; m_\theta, k_\theta) \quad \text{e.g.} \quad m_\theta(\bullet) = \phi(\bullet)^\top \theta, \text{ or } k_\theta(\bullet, \circ) = \theta_1 \exp \left(-\frac{(\bullet - \circ)^2}{2\theta_2^2} \right).$$

- The **evidence** in Bayes' theorem is the **marginal likelihood for the model**

$$p(f | y, x, \theta) = \frac{p(y | f, x, \theta) p(f | \theta)}{\int p(y | f, x, \theta) p(f | \theta) df} = \frac{p(y | f, x, \theta) p(f | \theta)}{p(y | x, \theta)}$$

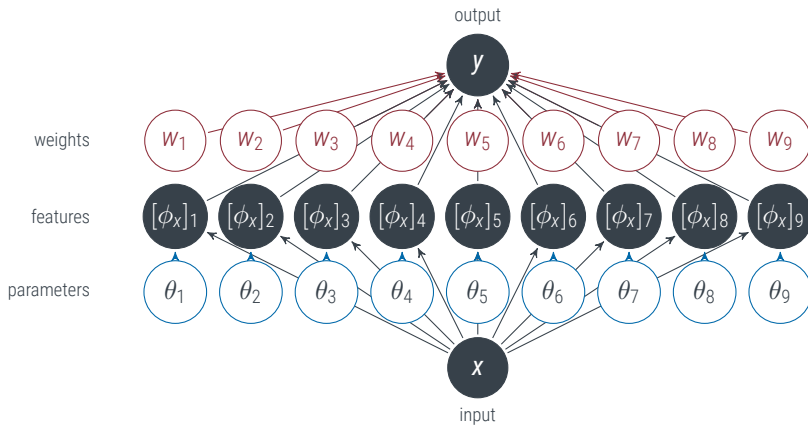
- For Gaussians and Gaussian processes, the evidence has **analytic form**:

$$\underbrace{\mathcal{N}(y; \phi_X^{\theta^\top} w + b, \Lambda)}_{p(y|f,x,\theta)} \cdot \underbrace{\mathcal{N}(w, \mu, \Sigma)}_{p(f)} = \underbrace{\mathcal{N}(w; m_{\text{post}}^\theta, V_{\text{post}}^\theta)}_{p(f|y,x,\theta)} \cdot \underbrace{\mathcal{N}(y; \phi_X^{\theta^\top} \mu + b, \phi_X^{\theta^\top} \Sigma \phi_X^\theta + \Lambda)}_{p(y|\theta,x)}$$

$$\mathcal{N}(y; f^\theta(X), \Lambda^\theta) \cdot \mathcal{GP}(f, \mu^\theta, k^\theta) = \mathcal{GP}(f; m_{\text{post}}^\theta, V_{\text{post}}^\theta) \cdot \mathcal{N}(y; \mu^\theta(X), \Lambda^\theta + k^\theta(X, X))$$

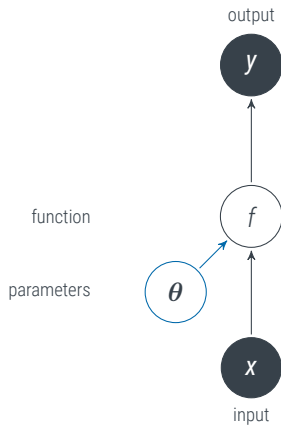
Kernel Learning is Fitting an Entire Population of Features

The Graphical Model View



Kernel Learning is Fitting an Entire Population of Features

The Graphical Model View



$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} p(\mathbf{y} \mid \mathbf{x}, \theta) = \arg \max_{\theta} \int \mathcal{N}(\mathbf{y}; f(\mathbf{X}), \Lambda_{\theta}) \cdot \mathcal{N}(f_{\mathbf{X}}, \mu_{\mathbf{X}}, k_{\mathbf{X}\mathbf{X}}) d\mathbf{f}_{\mathbf{X}} \\&= \arg \max_{\theta} \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}; k_{\mathbf{X}\mathbf{X}} + \Lambda_{\theta}) \int \mathcal{N}(f_{\mathbf{X}}; \mu_{\mathbf{y}, \mathbf{X}}, v_{\mathbf{y}, \mathbf{X}\mathbf{X}}) d\mathbf{f}_{\mathbf{X}} \\&= \arg \max_{\theta} \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}^{\theta}, k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}) \\&= \arg \max_{\theta} \log \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}^{\theta}, k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}) \\&= \arg \min_{\theta} -\log \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}^{\theta}, k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}) \\&= \arg \min_{\theta} \frac{1}{2} \left(\underbrace{(\mathbf{y} - \mu_{\mathbf{X}}^{\theta})^{\top} (k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta})^{-1} (\mathbf{y} - \mu_{\mathbf{X}}^{\theta})}_{\text{square error}} + \underbrace{\log |k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}|}_{\text{model complexity / Occam factor}} \right) + \frac{N}{2} \log 2\pi\end{aligned}$$

Entities should not be
posited without necessity

***Numquam ponenda est pluralitas sine necessitate.
Plurality must never be posited without necessity.***

William of Occam
(1285 (Occam, Surrey)–1349 (Munich, Bavaria))
stained-glass window by Lawrence Lee

What does the Optimizer need from us?

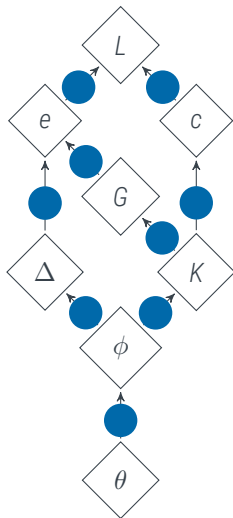
A bit of algorithmic wizardry

$$L(\boldsymbol{\theta}) = \frac{1}{2} \left(\underbrace{\underbrace{(y - \phi_X^{\boldsymbol{\theta}^\top} \mu)^\top}_{=:K} \underbrace{\left(\underbrace{\phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}}}_{=:G} + \Lambda \right)^{-1}}_{=:e}}_{=:e} \underbrace{(y - \phi_X^{\boldsymbol{\theta}^\top} \mu)}_{=: \Delta} + \underbrace{\log \left| \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda \right|}_{=:c} \right)$$



What does the Optimizer need from us?

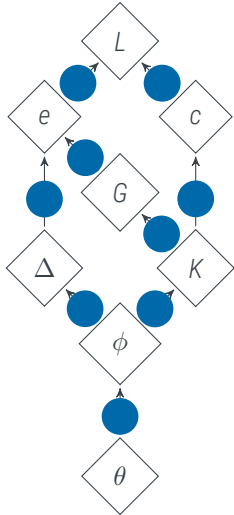
Automatic Differentiation



$$\begin{aligned}
 L(\theta) &= \frac{1}{2} \underbrace{\left((y - \underbrace{\phi_X^{\theta^\top} \mu}_{=:K})^\top \underbrace{\left(\underbrace{\phi_X^{\theta^\top} \Sigma \phi_X}_{=:G} + \Lambda \right)^{-1}}_{=:e} \right)}_{=:e} \underbrace{\left(\underbrace{(y - \phi_X^{\theta^\top} \mu)}_{=: \Delta} + \log \underbrace{\left| \phi_X^{\theta^\top} \Sigma \phi_X + \Lambda \right|}_{=:c} \right)}_{=:c} \\
 &= m_9 + m_8 = (m_6^\top m_5 m_6) + \log |m_7 + \Lambda| \\
 &= \dots
 \end{aligned}$$

What does the Optimizer need from us?

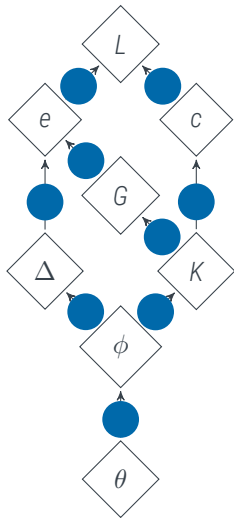
Automatic Differentiation – Forward Mode



$$L(\theta) = \frac{1}{2} \underbrace{\left((y - \underbrace{\phi_X^{\theta^T} \mu}_{=:K})^T \underbrace{\left(\underbrace{\phi_X^{\theta^T} \Sigma \phi_X^{\theta}}_{=:G} + \Lambda \right)^{-1}}_{=:e} \right)}_{=:c} \underbrace{\left(\underbrace{(y - \phi_X^{\theta^T} \mu)}_{=: \Delta} + \log \left| \underbrace{\phi_X^{\theta^T} \Sigma \phi_X^{\theta} + \Lambda}_{=:c} \right| \right)}$$

$$\begin{aligned}
 \frac{\partial L}{\partial \theta} &= \frac{\partial L}{\partial e} \frac{\partial e}{\partial \theta} + \frac{\partial L}{\partial c} \frac{\partial c}{\partial \theta} = \dot{m}_9 \frac{\partial e}{\partial \theta} + \dot{m}_8 \frac{\partial c}{\partial \theta} = \dot{m}_9 \left(\frac{\partial e}{\partial \Delta} \frac{\partial \Delta}{\partial \theta} + \frac{\partial e}{\partial G} \frac{\partial G}{\partial \theta} \right) + \dot{m}_8 \frac{\partial c}{\partial K} \frac{\partial K}{\partial \theta} \\
 &= \dot{m}_9 \left(\dot{m}_6 \frac{\partial \Delta}{\partial \theta} + \dot{m}_5 \frac{\partial G}{\partial \theta} \right) + \dot{m}_8 \dot{m}_7 \frac{\partial K}{\partial \theta} = \dot{m}_9 \left(\dot{m}_6 \frac{\partial \Delta}{\partial \phi} \frac{\partial \phi}{\partial \theta} + \dot{m}_5 \frac{\partial G}{\partial K} \frac{\partial K}{\partial \theta} \right) + \dot{m}_8 \dot{m}_7 \frac{\partial K}{\partial \theta} \\
 &= \dot{m}_9 \dot{m}_6 \dot{m}_2 \frac{\partial \phi}{\partial \theta} + (\dot{m}_9 \dot{m}_5 \dot{m}_4 + \dot{m}_8 \dot{m}_7) \frac{\partial K}{\partial \theta} = \left(\dot{m}_9 \dot{m}_6 \dot{m}_2 + (\dot{m}_9 \dot{m}_5 \dot{m}_4 + \dot{m}_8 \dot{m}_7) \right) \frac{\partial \phi}{\partial \theta} \frac{\partial \phi}{\partial \theta} \\
 &= (\dot{m}_9 \dot{m}_6 \dot{m}_2 + (\dot{m}_9 \dot{m}_5 \dot{m}_4 + \dot{m}_8 \dot{m}_7) \dot{m}_3) \frac{\partial \phi}{\partial \theta} \frac{\partial \theta}{\partial \theta} \\
 &= (\dot{m}_9 \dot{m}_6 \dot{m}_2 + (\dot{m}_9 \dot{m}_5 \dot{m}_4 + \dot{m}_8 \dot{m}_7) \dot{m}_3) \dot{m}_1 1
 \end{aligned}$$

Automatic Differentiation – Forward Mode



$$L(\theta) = \underbrace{\frac{1}{2} \left((y - \phi_X^{\theta^\top} \mu)^\top \underbrace{\left(\underbrace{\phi_X^{\theta^\top} \Sigma \phi_X^\theta}_{=:K} + \Lambda \right)}_{=:G} \right)^{-1}}_{=:e} \underbrace{(y - \phi_X^{\theta^\top} \mu)}_{=: \Delta} + \underbrace{\log \left| \phi_X^{\theta^\top} \Sigma \phi_X^\theta + \Lambda \right|}_{=:c}$$

$$\dot{m}_9 = \frac{\partial L}{\partial e} = 1/2 \quad \dot{m}_8 = \frac{\partial L}{\partial c} = 1/2 \quad [\dot{m}_7]_{ij} = \frac{\partial c}{\partial K_{ij}} = K_{ij}^{-1}$$

$$[\dot{m}_6]_i = \frac{\partial e}{\partial \Delta_i} = 2[G\Delta]_i \quad [\dot{m}_5]_{jj} = \frac{\partial e}{\partial G_{jj}} = \Delta_i \Delta_j \quad [\dot{m}_4]_{ij,kl} = \frac{\partial G_{ij}}{\partial K_{kl}} = -G_{ik}G_{jl}$$

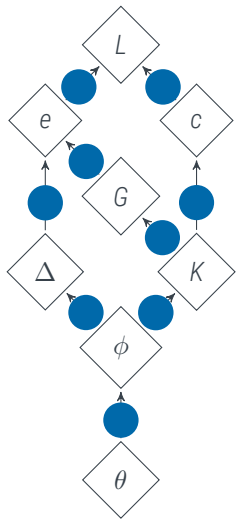
$$[\dot{m}_3]_{ij,ab} = \frac{\partial K_{ij}}{\partial \phi_{ab}} = \delta_{ia} [\Sigma \phi]_{bj} + \delta_{jb} [\Sigma \phi]_{ki}$$

$$[\dot{m}_2]_{i,ab} = \frac{\partial \Delta_i}{\partial \phi_{ab}} = -\delta_{ia} \mu_b \quad [\dot{m}_1]_{ab,\ell} = \frac{\partial \phi_{ab}}{\partial \theta_\ell} = \text{your choice!}$$

What does the Optimizer need from us?

Automatic Differentiation – Backward Mode

[Seppo Linnainmaa, 1970]



$$L(\theta) = \frac{1}{2} \left(\underbrace{(y - \phi_X^{\theta^T} \mu)^T}_{=:G} \underbrace{\left(\underbrace{\phi_X^{\theta^T} \Sigma \phi_X^{\theta}}_{=:K} + \Lambda \right)}_{=:e}^{-1} \underbrace{(y - \phi_X^{\theta^T} \mu)}_{=:c} + \log \underbrace{|\phi_X^{\theta^T} \Sigma \phi_X^{\theta} + \Lambda|}_{=:c} \right)$$

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial \phi} \frac{\partial \phi}{\partial \theta} =: \bar{m}_1 = \left(\frac{\partial L}{\partial \Delta} \frac{\partial \Delta}{\partial \phi} + \frac{\partial L}{\partial K} \frac{\partial K}{\partial \phi} \right) \frac{\partial \phi}{\partial \theta} =: (\bar{m}_2 + \bar{m}_3) \frac{\partial \phi}{\partial \theta}$$

$$\bar{m}_2 = \frac{\partial L}{\partial e} \frac{\partial e}{\partial \Delta} \frac{\partial \Delta}{\partial \phi} =: \bar{m}_6 \frac{\partial \Delta}{\partial \phi} \quad \bar{m}_3 = \left(\frac{\partial L}{\partial G} \frac{\partial G}{\partial K} + \frac{\partial L}{\partial c} \frac{\partial c}{\partial K} \right) \frac{\partial K}{\partial \phi} =: (\bar{m}_4 + \bar{m}_7) \frac{\partial K}{\partial \phi}$$

$$\bar{m}_4 = \frac{\partial L}{\partial e} \frac{\partial e}{\partial G} \frac{\partial G}{\partial K} =: \bar{m}_5 \frac{\partial G}{\partial K} \quad \bar{m}_5 = \frac{\partial L}{\partial e} \frac{\partial e}{\partial G} =: \bar{m}_9 \frac{\partial e}{\partial G} \quad \bar{m}_6 = \frac{\partial L}{\partial e} \frac{\partial e}{\partial \Delta} =: \bar{m}_9 \frac{\partial e}{\partial \Delta}$$

$$\bar{m}_7 = \frac{\partial L}{\partial c} \frac{\partial c}{\partial K} =: \bar{m}_8 \frac{\partial c}{\partial K} \quad \bar{m}_8 = \bar{m}_9 = 1/2$$

$\bar{w}_i = \frac{\partial L}{\partial \text{subgraph}_i}$ are known as *adjoints*. Traverse graph backward to collect the derivative. This is faster than forward-mode for single-output-many-input functions, but requires storing the above structure (known as a *Wengert list*). (cf. "Backpropagation")

Source Separation

extracting individual signals from data in which they are linearly mixed

From Lecture 6:

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

In our case $y = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$ and $p(f^1, f^2) = \mathcal{GP} \left(\begin{bmatrix} f^1 \\ f^2 \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k^1 & 0 \\ 0 & k^2 \end{bmatrix} \right)$ Thus we can *extract* the individual signals f^1 and f^2 from the data y by setting $B = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ and get, in particular

$$p(f_{\bullet}^1 | y) = \mathcal{GP}(f_{\bullet}^1; k_{\bullet, X}^1 (k_{XX}^1 + k_{XX}^2 + \Lambda)^{-1} y, \quad k_{\bullet, \circ}^1 - k_{\bullet, X}^1 (k_{XX}^1 + k_{XX}^2 + \Lambda)^{-1} k_{X, \circ}^2)$$

Summary:

- ▶ An unstructured kernel regression model can only do so much. **Extrapolation** and extracting **structural knowledge** require prior knowledge about the **causal structure**.
- ▶ Linear models with elaborate features can be quite expressive, while remaining interpretable
- ▶ Complicated processes require complicated (and questionable!) prior assumptions

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}}
```

The ability to build structured predictive models is a **key skill**. Everyone can run a TensorFlow script! Masters of structured probabilistic inference are highly sought after.

