

02477 – Bayesian Machine Learning: Lecture 5

Michael Riis Andersen

Technical University of Denmark,
DTU Compute, Department of Applied Math and Computer Science

Outline

- ➊ Towards prior distributions for function spaces
- ➋ A visual approach towards Gaussian process regression
- ➌ Gaussian process regression
- ➍ Covariance functions
- ➎ Hyperparameters and the marginal likelihood

Multitude of Gaussian processes applications

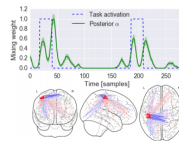
■ Regression (supervised learning)

- Time series analysis
- EEG brain imaging
- Survival analysis for cancer data
- Predicting rainfall
- Robot dynamics
- ...

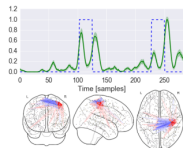
■ Classification (supervised learning)

- Recognizing human movements
- Brain decoding
- ...

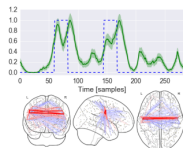
- Used as building block in more complex models
- Dimensionality reduction (unsupervised learning)
- Optimization of black functions (Bayesian optimization)
- Numerical integration (Bayesian quadrature)
- Solving differential equations (probabilistic numerics)



(a) Right hand tapping



(c) Left hand tapping



(b) Tongue wagging

Towards prior distributions for function spaces

Parametric models

- In week 3, we studied linear models of the form

$$y_n = f_n + e_n = \phi(\mathbf{x}_n)^T \mathbf{w} + e_n$$

- In week 4, we studied Bayesian logistic regression

$$y_n | f_n \sim \text{Ber}(\sigma(f_n))$$

$$f_n = \phi(\mathbf{x}_n)^T \mathbf{w}$$

- Typical workflow

1. Specify prior $p(\mathbf{w})$ and likelihood $p(\mathbf{y}|\mathbf{w})$
2. Calculate posterior distribution $p(\mathbf{w}|\mathbf{y})$
3. Make predictions based on the predictive distribution $p(y^*|\mathbf{y}, \mathbf{x}^*)$

- All we care about is the parameters \mathbf{w} - Once we have calculated the posterior distribution $p(\mathbf{w}|\mathbf{y})$, we don't need the training data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ anymore

- Linear and logistic regression are both *parametric models*: probability distributions indexed by finite dimensional parameters

From parameters to functions I

- Our linear model

$$y_n = f_n + e_n = \phi(\mathbf{x}_n)^T \mathbf{w} + e_n$$

- Our goal was to learn the latent function

$$f_n = \phi(\mathbf{x}_n)^T \mathbf{w}$$

- We focussed on \mathbf{w} and the joint distribution via the product rule

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$$

- Let \mathbf{f} denote the function values, i.e. $\mathbf{f} = [f(\phi(\mathbf{x}_1)) \quad f(\phi(\mathbf{x}_2)) \quad \dots \quad f(\phi(\mathbf{x}_N))]$

$$p(\mathbf{y}, \mathbf{f}, \mathbf{w}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{w})p(\mathbf{w})$$

- The model is the same - we can recover the old model formulation via the sum rule

$$p(\mathbf{y}, \mathbf{w}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{w})p(\mathbf{w})d\mathbf{f}$$

From parameters to functions II

- The augmented model

$$p(\mathbf{y}, \mathbf{f}, \mathbf{w}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{w})p(\mathbf{w})$$

- What if we integrate out the parameters instead?

$$p(\mathbf{y}, \mathbf{f}) = p(\mathbf{y}|\mathbf{f}) \int p(\mathbf{f}|\mathbf{w})p(\mathbf{w})d\mathbf{w} = p(\mathbf{y}|\mathbf{f})p(\mathbf{f})$$

$$\text{where } p(\mathbf{f}) = \int p(\mathbf{f}|\mathbf{w})p(\mathbf{w})d\mathbf{w}$$

- Let's study the distribution of $\mathbf{f} = \Phi\mathbf{w}$ for our Gaussian prior on $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$

$$p(\mathbf{f}) = \int p(\mathbf{f}|\mathbf{w})p(\mathbf{w})d\mathbf{w} = \int p(\mathbf{f}|\mathbf{w})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})d\mathbf{w}$$

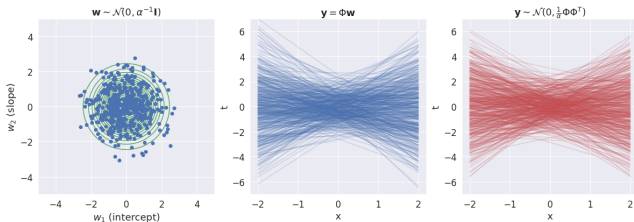
- We could do the integral directly, let's use this result instead

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}) \quad \Rightarrow \quad \mathbf{a} + \mathbf{B}\mathbf{x} \sim \mathcal{N}(\mathbf{a} + \mathbf{B}\mathbf{m}, \mathbf{B}\mathbf{V}\mathbf{B}^T)$$

- What is the distribution of \mathbf{f} ?

Changing perspective from weight space to function space

$$p(\mathbf{f}) = \int p(\mathbf{f}|\mathbf{w})p(\mathbf{w})d\mathbf{w} = \int p(\mathbf{f}|\mathbf{w})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})d\mathbf{w} = \mathcal{N}(\mathbf{f}|\mathbf{0}, \alpha^{-1}\Phi\Phi^T)$$



■ Two ways to generate samples of $\mathbf{f} \sim p(\mathbf{f})$

■ Weight space-perspective

Step 1: Generate a sample $\mathbf{w}^{(i)} \sim p(\mathbf{w})$

Step 2: Compute $\mathbf{f}^{(i)} = \Phi\mathbf{w}^{(i)}$

■ Function space-perspective

Step 1: Generate a sample $\mathbf{f}^{(i)} \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\Phi\Phi^T)$

A closer look at the covariance

- A prior on linear functions: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$, where $\mathbf{K} = \frac{1}{\alpha} \Phi \Phi^T$
- A closer look on the covariance between f_i and f_j

$$\begin{aligned} K_{ij} &= \text{cov}(y_i, y_j) = \text{cov}(f(\mathbf{x}_i), f(\mathbf{x}_j)) \\ &= \text{cov}(\phi(\mathbf{x}_i)^T \mathbf{w}, \phi(\mathbf{x}_j)^T \mathbf{w}) \\ &= \phi(\mathbf{x}_i)^T \text{cov}(\mathbf{w}, \mathbf{w}) \phi(\mathbf{x}_j) \\ &= \phi(\mathbf{x}_i)^T \mathbb{V}(\mathbf{w}) \phi(\mathbf{x}_j) \\ &= \phi(\mathbf{x}_i)^T \frac{1}{\alpha} \mathbf{I} \phi(\mathbf{x}_j) \\ &= \frac{1}{\alpha} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ &\equiv k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

- What happens if we change the *covariance function* k ?

Covariance functions

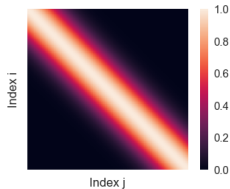
Linear

$$k(x_i, x_j) = \frac{1}{\alpha} \phi(x_i)^T \phi(x_j)$$



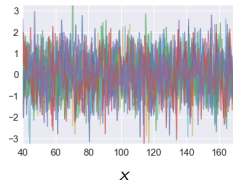
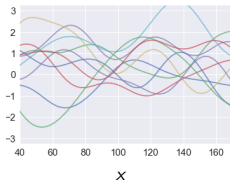
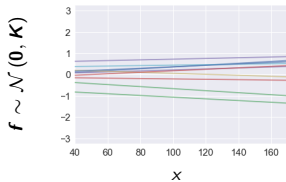
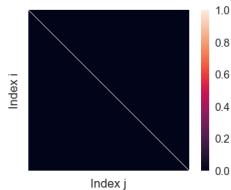
Squared exponential

$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$



White noise

$$k(x_i, x_j) = \delta(x_i - x_j)$$



The form of the covariance function determines the characteristics of the functions

The big picture: Summary so far

1. We started with a Bayesian linear model

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$$

2. We introduced \mathbf{f} into the model and marginalized over the weights \mathbf{w}

$$p(\mathbf{y}, \mathbf{f}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{w})p(\mathbf{w})d\mathbf{w} = p(\mathbf{y}|\mathbf{f})p(\mathbf{f})$$

3. This gave us a prior for linear functions in function space $p(\mathbf{f})$, where the covariance function for \mathbf{f} was given by

$$k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\alpha} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

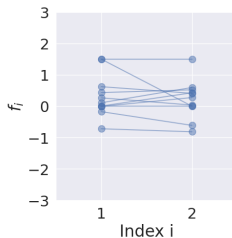
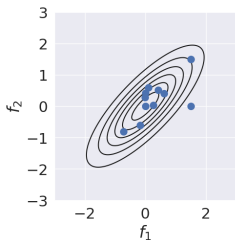
4. By changing the form of the covariance function $k(\mathbf{x}_i, \mathbf{x}_j)$, we can model much more interesting functions

A visual approach towards Gaussian process regression

Visualizing samples in higher dimensions

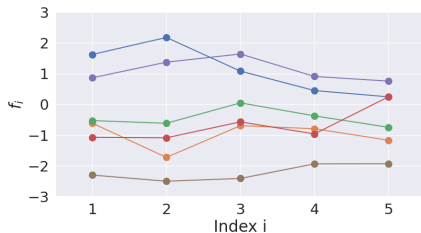
- How can a multivariate normal distribution represent functions?
- Visualizations in 2D

$$\mathbf{K} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



Visualizing samples in higher dimensions

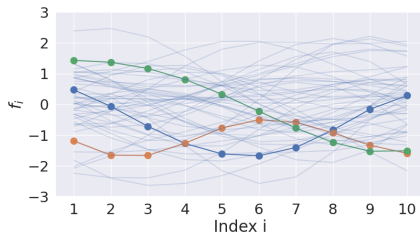
■ Visualizations in 5D



$$K = \begin{bmatrix} 1 & 0.8^1 & 0.8^2 & 0.8^3 & 0.8^4 \\ 0.8^1 & 1 & 0.8^1 & 0.8^2 & 0.8^3 \\ 0.8^2 & 0.8^1 & 1 & 0.8^1 & 0.8^2 \\ 0.8^3 & 0.8^2 & 0.8^1 & 1 & 0.8^1 \\ 0.8^4 & 0.8^3 & 0.8^2 & 0.8^1 & 1 \end{bmatrix}$$

Visualizing samples in higher dimensions

■ Visualizations in 10D



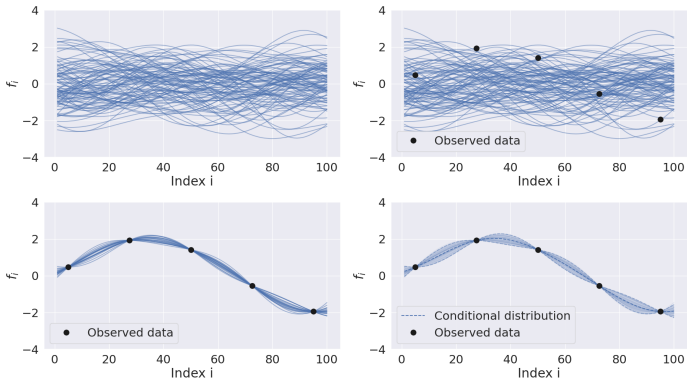
$$K = \begin{bmatrix} 1 & 0.8^1 & 0.8^2 & \dots & 0.8^9 \\ 0.8^1 & 1 & 0.8^1 & & \vdots \\ 0.8^2 & 0.8^1 & 1 & & \vdots \\ \vdots & & & \ddots & \vdots \\ 0.8^9 & \dots & \dots & \dots & 1 \end{bmatrix}$$

Conditioning

- So far, we have seen samples from the distribution $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- We can also write $p(\mathbf{f}) = p(f_1, \mathbf{f}_{2:10})$
- We now observe $f_1 = 0$
- Let's sample from the conditional distribution $p(\mathbf{f}_{2:10} | f_1 = 0)$

Conditioning II

- Let's now consider a case with $\mathbf{f} \in \mathbb{R}^{100}$ dimensions with 5 observations



- Informally: We can think functions as vectors with infinite dimensions
- Using conditioning in multivariate Gaussian distributions, we can do non-linear regression!

Formal definitions

Definition of the multivariate Gaussian distribution

A random vector $\mathbf{x} = [x_1, x_2, \dots, x_D]$ is said to have the **multivariate Gaussian distribution** if all linear combinations of \mathbf{x} are (univariate) Gaussian distributed:

$$f = a_1 x_1 + a_2 x_2 + \dots + a_D x_D \sim \mathcal{N}(m, v)$$

for all $\mathbf{a} \in \mathbb{R}^D$

Definition of Gaussian process

A **Gaussian process** (GP) is a collection of random variables, any finite number of which have a joint Gaussian distribution.

Notation and characterization

- We'll use the notation

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

- A Gaussian process can be considered as a prior distribution over functions $f : \mathcal{X} \rightarrow \mathbb{R}$ (the domain \mathcal{X} is typically \mathbb{R}^D)
- A Gaussian process is completely characterized by its mean function $m(\mathbf{x})$ and its covariance function $k(\mathbf{x}, \mathbf{x}')$.

$$\begin{aligned} m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})] \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \end{aligned}$$

- This means that $f(\mathbf{x})$ and $f(\mathbf{x}')$ are jointly Gaussian distributed with covariance $k(\mathbf{x}, \mathbf{x}')$
- Not all functions are valid covariance functions - more on that later

Gaussian process regression

Recall: Linear Gaussian-systems in general (see Section 3.3 in Murphy1)

- For *linear* systems: the Gaussian distribution is *conjugate* to itself
- The *posterior* for a *linear* Gaussian model with Gaussian prior is also *Gaussian*

$$p(\mathbf{y}|\mathbf{z}) = \mathcal{N}(\mathbf{y}|\mathbf{W}\mathbf{z} + \mathbf{b}, \Sigma_y) \qquad p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_z, \Sigma_z)$$

- The *joint* distribution $p(\mathbf{z}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{z} \\ \mathbf{y} \end{bmatrix} \middle| \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_z \\ \mathbf{W}\boldsymbol{\mu}_z + \mathbf{b} \end{bmatrix} \qquad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_z & \Sigma_z \mathbf{W}^T \\ \mathbf{W}\Sigma_z & \Sigma_y + \mathbf{W}\Sigma_z \mathbf{W}^T \end{bmatrix}$$

- The *posterior* distribution of \mathbf{z} given \mathbf{y}

$$\begin{aligned} p(\mathbf{z}|\mathbf{y}) &= \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{z|y}, \Sigma_{z|y}) \\ \Sigma_{z|y}^{-1} &= \Sigma_z^{-1} + \mathbf{W}^T \Sigma_y \mathbf{W} \\ \boldsymbol{\mu}_{z|y} &= \Sigma_{z|y} \left[\mathbf{W}^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{b}) + \Sigma_z^{-1} \boldsymbol{\mu}_z \right] \end{aligned}$$

- The *marginal* distribution \mathbf{y}

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{W}\boldsymbol{\mu}_z + \mathbf{b}, \Sigma_y + \mathbf{W}\Sigma_z \mathbf{W}^T)$$

Conditioning for multivariate Gaussians (Murphy1 Section 3.2.3)

- Suppose $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ is jointly Gaussian $\mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with mean and covariance

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$$

- The precision matrix $\boldsymbol{\Lambda}$

$$\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} = \begin{bmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{bmatrix}$$

- The marginals are given by

$$p(\mathbf{y}_1) = \mathcal{N}(\mathbf{y}_1|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

$$p(\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_2|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$$

- The conditional

$$p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1|\boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

$$\boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{y}_2 - \boldsymbol{\mu}_2)$$

$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1}$$

Gaussian process regression I

- Our model

$$y_n = f(\mathbf{x}_n) + e_n$$

- Likelihood for all datapoints (assuming homoscedastic noise)

$$p(\mathbf{y}|\mathbf{f}) = \prod_{n=1}^N \mathcal{N}(y_n|f_n, \beta^{-1}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \beta^{-1}\mathbf{I})$$

- We impose a prior directly on the *function values* $\mathbf{f} = [f(\mathbf{x}_1) \quad f(\mathbf{x}_2) \quad \dots \quad f(\mathbf{x}_N)]$

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}) \quad \text{for} \quad (\mathbf{K})_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

- Goal: compute predictive distribution for $y^* = y(\mathbf{x}^*)$ given data \mathbf{y} , i.e. $p(y^*|\mathbf{y}, \mathbf{x}^*)$

- Two-step strategy

1. Calculate the joint Gaussian distribution $p(\mathbf{y}, y^*|\mathbf{x}^*)$
2. Use rule for conditioning in Gaussian distributions to compute $p(y^*|\mathbf{y}, \mathbf{x}^*)$

Gaussian process regression II

- Recall: General Linear Gaussian systems (Murphy1 page 86-77).

If

$$\begin{aligned}p(\mathbf{z}) &= \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) \\ p(\mathbf{y} | \mathbf{z}) &= \mathcal{N}(\mathbf{y} | \mathbf{W}\mathbf{z} + \mathbf{b}, \boldsymbol{\Sigma}_y)\end{aligned}$$

then

$$p(\mathbf{y}) = \int p(\mathbf{y} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} = \mathcal{N}(\mathbf{y} | \mathbf{W}\boldsymbol{\mu}_z + \mathbf{b}, \boldsymbol{\Sigma}_y + \mathbf{W}\boldsymbol{\Sigma}_z\mathbf{W}^T)$$

- We can compute the marginal distribution of \mathbf{y} using the *sum rule*

$$\begin{aligned}p(\mathbf{y}) &= \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}) d\mathbf{f} \\ &= \int \mathcal{N}(\mathbf{y} | \mathbf{f}, \beta^{-1} \mathbf{I}) \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K}) d\mathbf{f} \\ &= \mathcal{N}(\mathbf{y} | ?, ?)\end{aligned}$$

- Spend 5 minutes calculating the mean and variance of $p(\mathbf{y})$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | ?, ?)$$

Gaussian process regression III

- The distribution of $\mathbf{y} \in \mathbb{R}^N$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}) \quad \text{for} \quad \mathbf{C} = \mathbf{K} + \beta^{-1}\mathbf{I}$$

- Let $\tilde{\mathbf{y}} = [y(\mathbf{x}^*) \quad \mathbf{y}]^T \in \mathbb{R}^{N+1}$, then

$$p(\tilde{\mathbf{y}}) = \mathcal{N}(\tilde{\mathbf{y}}|\mathbf{0}, \tilde{\mathbf{C}}) \quad \text{for} \quad \tilde{\mathbf{C}} = \begin{bmatrix} c & \mathbf{k} \\ \mathbf{k}^T & \mathbf{C} \end{bmatrix}$$

where

$$c = k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1}$$

$$\mathbf{k} = [k(\mathbf{x}^*, \mathbf{x}_1) \quad k(\mathbf{x}^*, \mathbf{x}_2) \quad \dots \quad k(\mathbf{x}^*, \mathbf{x}_N)]$$

- What is the mean and variance for the following distribution?

$$p(y^*|\mathbf{y}) = \mathcal{N}(y^* | \mu_{y^*|\mathbf{y}}, \sigma_{y^*|\mathbf{y}}^2)$$

Suppose $\mathbf{y} = (y_1, y_2)$ is jointly Gaussian $\mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$p(y_1|y_2) = \mathcal{N}(y_1|\mu_{1|2}, \Sigma_{1|2})$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(y_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

Example

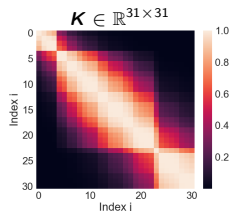
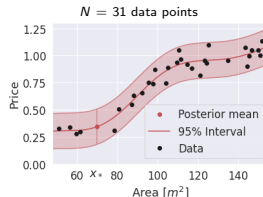
Key equations for Gaussian process regression

$$p(y^* | \mathbf{y}) = \mathcal{N}(y^* | \mu_{y^* | \mathbf{y}}, \sigma_{y^* | \mathbf{y}}^2)$$

$$\mu_{y^* | \mathbf{y}} = \mathbf{k} (\mathbf{K} + \beta^{-1} \mathbf{I})^{-1} \mathbf{y}$$

$$\sigma_{y^* | \mathbf{y}}^2 = c - \mathbf{k} (\mathbf{K} + \beta^{-1} \mathbf{I})^{-1} \mathbf{k}^T$$

- Predict $y^* = f(x_*) + e^*$ for test input $x_* = 70$
- Observation vector $\mathbf{y} = [y_1, y_2, \dots, y_{31}]^T \in \mathbb{R}^{31 \times 1}$
- $k(x, x') = \text{cov}(f(x), f(x')) = \exp \left[-\frac{(x-x')^2}{2 \cdot 20^2} \right]$
- Covariance matrix for training data: $[\mathbf{K}]_{ij} = k(x_i, x_j)$
- Cov. between test and training $[\mathbf{k}]_j = k(x_*, x_j)$
- Covariance of test point $y^* = y(x_*)$: $c = k(x_*, x_*) + \beta^{-1}$
- Now we have all the ingredients for the key equations



Gaussian process intuition

- Gaussian process implements the assumption

$$\mathbf{x} \approx \mathbf{x}' \Rightarrow f(\mathbf{x}) \approx f(\mathbf{x}')$$

- In words: If the inputs are similar, the outputs should be similar as well.

- Using the squared exponential covariance function as example

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$$

- Then covariance between $f(\mathbf{x})$ and $f(\mathbf{x})'$ is given by

$$\text{cov}[f(\mathbf{x}), f(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$$

- Note: the covariance between outputs are given in terms of the inputs

True or false?

Key equations for Gaussian process regression

$$\begin{aligned}p(y^*|\mathbf{y}) &= \mathcal{N}\left(y^*|\mu_{y^*|\mathbf{y}}, \sigma_{y^*|\mathbf{y}}^2\right) \\ \mu_{y^*|\mathbf{y}} &= \mathbf{k}(\mathbf{K} + \beta^{-1}\mathbf{I})^{-1}\mathbf{y} \\ \sigma_{y^*|\mathbf{y}}^2 &= c - \mathbf{k}(\mathbf{K} + \beta^{-1}\mathbf{I})^{-1}\mathbf{k}^T\end{aligned}$$

True or false?

- Spend 5 minutes on the DTU Learn quiz: “Lecture 5: Key equations for GP Regression.”

Covariance functions

Covariance functions

- A covariance function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ maps a pair of inputs $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ from some input space \mathcal{X} to the real line \mathbb{R}

- Recall: the covariance / kernel matrix is given by

$$\mathbf{K}_{ij} = \text{cov}(f(\mathbf{x}_i), f(\mathbf{x}_j)) = k(\mathbf{x}_i, \mathbf{x}_j)$$

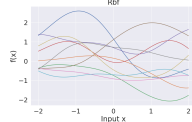
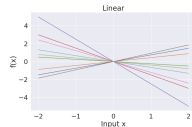
- Covariance functions must be symmetric & Positive Semi-Definite such that

$$\text{(Symmetric)} \quad \mathbf{K} = \mathbf{K}^T$$

$$\text{(PSD)} \quad \forall \mathbf{x} \neq 0 : \mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$$

- Must hold for all possible data sets $\{\mathbf{x}_n\}_{n=1}^N \subset \mathcal{X}$ in the input space \mathcal{X}

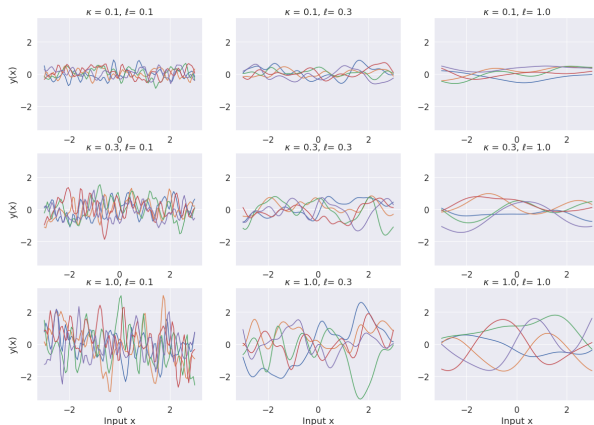
- Covariance functions as prior information



The squared exponential kernel - prior samples

$$k(\mathbf{x}, \mathbf{x}') = \kappa^2 \exp \left[-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2} \right]$$

Parameter ℓ is called the *lengthscale* and parameter κ is called the *magnitude*



Constructing new kernels from old ones

Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

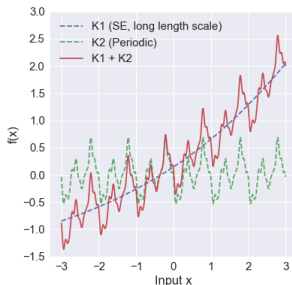
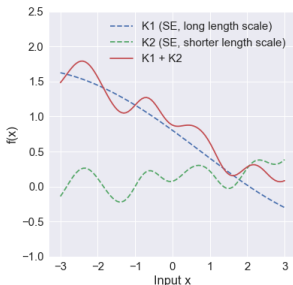
where $c > 0$ is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(\mathbf{x})$ is a function from \mathbf{x} to \mathbb{R}^M , $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M , \mathbf{A} is a symmetric positive semidefinite matrix, \mathbf{x}_a and \mathbf{x}_b are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and k_a and k_b are valid kernel functions over their respective spaces.

Additive kernels

- Adding two SEs kernels to model long term trends (long length scale) and short term fluctuations (short length scale)

$$k(\mathbf{x}, \mathbf{x}') = \kappa_1^2 \exp \left[-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell_1^2} \right] + \kappa_2^2 \exp \left[-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell_2^2} \right]$$

- Adding SE and period kernels to model long term trends (long length scale) and periodic fluctuations



Hyperparameters and the marginal likelihood

The marginal likelihood I

- Let θ denote all hyperparameters, then marginal likelihood for Gaussian likelihood

$$\begin{aligned} p(\mathbf{y}|\theta) &= \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta_K)d\mathbf{f} \\ &= \int \mathcal{N}(\mathbf{y}|\mathbf{f}, \beta^{-1}\mathbf{I}) \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}) d\mathbf{f} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \mathbf{K}) \end{aligned}$$

- We can tune the hyperparameters of the model by optimizing the marginal likelihood as we did for linear regression
 1. Hyperparameters of the likelihood (e.g. β or σ)
 2. Hyperparameters of the kernel θ_K (e.g. lengthscales and magnitudes)
- In practice, we compute the gradient of $p(\mathbf{y}|\theta)$ wrt. θ and use numerical optimization

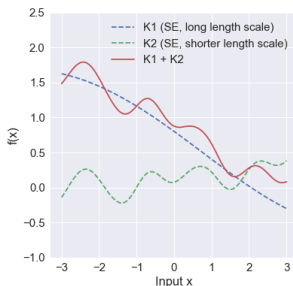
The marginal likelihood II

- Suppose we have 5 hyperparameters in total

$$\theta = \{\sigma, \kappa_1, \ell_1, \kappa_2, \ell_2\}$$

- Suppose we want to estimate those using 10-fold cross-validation and test out 10 values for each hyperparameter. How many times do we need to train the model?

$$10 \cdot 10^5 = 10^6$$



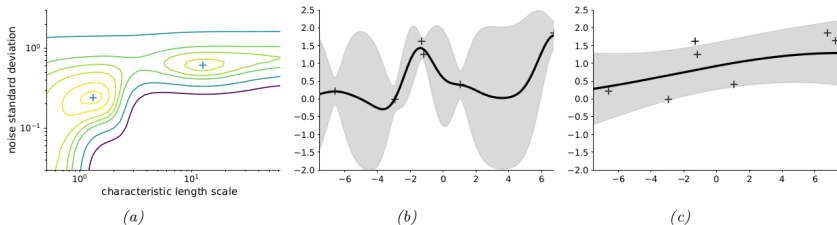
The marginal likelihood III

- The gradients of the marginal likelihood wrt. hyperparameters are given by

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^T - \mathbf{K}^{-1}) \frac{\partial \mathbf{K}}{\partial \theta_j} \right),$$

where $\boldsymbol{\alpha} = \mathbf{K}^{-1} \mathbf{y}$ and $\frac{\partial \mathbf{K}}{\partial \theta_j}$ depends on the specific choice of kernel.

- $\log p(\mathbf{y}|\boldsymbol{\theta})$ is also multimodal wrt. $\boldsymbol{\theta}$



From the Murphy1 book (p. 578)

The marginal likelihood: numerics

- In practice, we should avoid computing determinants and inverses!

$$\ln p(\mathbf{y}|\boldsymbol{\theta}) = -\frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |\beta^{-1}\mathbf{I} + \mathbf{K}| - \frac{1}{2} \mathbf{y}^T (\beta^{-1}\mathbf{I} + \mathbf{K})^{-1} \mathbf{y}$$

- In numpy: $|0.1\mathbf{I}_{400 \times 400}| = 0.0$, but $\ln |0.1\mathbf{I}_{400 \times 400}| = -2302.58$ and $\exp(-2302.58) > 0$
- Step 1: Compute Cholesky factorization of $\mathbf{C} = \beta^{-1}\mathbf{I} + \mathbf{K}$ such that $\mathbf{C} = \mathbf{L}\mathbf{L}^T$
- Step 2: Compute the log determinant term as follows

$$\ln |\mathbf{C}| = \ln |\mathbf{L}\mathbf{L}^T| = \ln |\mathbf{L}| \cdot |\mathbf{L}^T| = \ln |\mathbf{L}|^2 = 2 \ln |\mathbf{L}| = 2 \ln \prod_{n=1}^N L_{nn} = 2 \sum_{n=1}^N \ln L_{nn}$$

- Step 3: Compute quadratic term as follows

$$\mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} = \mathbf{y}^T (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{y} = \mathbf{y}^T \mathbf{L}^{-T} \mathbf{L}^{-1} \mathbf{y} = (\mathbf{L}^{-1} \mathbf{y})^T \underbrace{(\mathbf{L}^{-1} \mathbf{y})}_{=\mathbf{v}} = \mathbf{v}^T \mathbf{v}$$

- Step 4: Sum components

$$\ln p(\mathbf{y}|\boldsymbol{\theta}) = -\frac{N}{2} \ln(2\pi) - \frac{1}{2} 2 \sum_{n=1}^N \ln L_{nn} - \frac{1}{2} \mathbf{v}^T \mathbf{v}$$

- Note that we never compute the determinant or the inverse of \mathbf{C} directly!

Computational complexity of Gaussian Processes

Key equations for Gaussian process regression

$$\begin{aligned}p(y^*|\mathbf{y}) &= \mathcal{N}\left(y^*|\mu_{y^*|\mathbf{y}}, \sigma_{y^*|\mathbf{y}}^2\right) \\ \mu_{y^*|\mathbf{y}} &= \mathbf{k} (\mathbf{K} + \beta^{-1}\mathbf{I})^{-1} \mathbf{y} \\ \sigma_{y^*|\mathbf{y}}^2 &= c - \mathbf{k} (\mathbf{K} + \sigma^2\mathbf{I})^{-1} \mathbf{k}^T\end{aligned}$$

- Gaussian processes are *non-parametric models*
- Recall: If $\mathbf{A} \in \mathbb{R}^{N \times M}$ and $\mathbf{b} \in \mathbb{R}^M$, then the cost of computing \mathbf{Ab} is $\mathcal{O}(NM)$
- Recall: If $\mathbf{C} \in \mathbb{R}^{N \times N}$, then the cost of computing \mathbf{C}^{-1} is $\mathcal{O}(N^3)$
- What is computational complexity for computing the posterior distribution for 1 test point based on a data set with N observations? $\mathcal{O}(N^3)$
- What about the memory footprint? $\mathcal{O}(N^2)$