

732A96/TDDE15 Advanced Machine Learning

Gaussian Process Regression and Classification

Jose M. Peña
IDA, Linköping University, Sweden

Lectures 11: Kernels, Hyperparameter Learning and More

Contents

- ▶ Three Common Covariance Functions
- ▶ Learning the Hyperparameters of the Covariance Function
- ▶ More on Covariance Functions
- ▶ Gaussian Process Networks
- ▶ Bayesian Optimization

- ▶ Main source

- ▶ Rasmussen, C. E. and Williams, K. I. *Gaussian Processes for Machine Learning*. MIT Press, 2006. Chapters 2.3, 5.1-5.4.1.
- ▶ Friedman, N. and Nachman, I. Gaussian Process Networks. In *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence*, 211-219, 2000.
- ▶ Brochu, E., Cora, V. M. and de Freitas, N. *A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning*.
<https://arxiv.org/pdf/1012.2599>.

- ▶ Additional source

- ▶ Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006. Chapters 6.4.3-6.4.4.

Three Common Covariance Functions

- ▶ Let $r = \|\mathbf{x} - \mathbf{x}'\|$.
- ▶ Squared exponential (SE):

$$k_{SE}(r) = \sigma_f^2 \exp \left\{ -\frac{r^2}{2\ell^2} \right\}$$

where $\sigma_f^2 > 0, \ell > 0$. Very smooth.

- ▶ Rational quadratic (RQ):

$$k_{RQ}(r) = \sigma_f^2 \left(1 + \frac{r^2}{2\alpha\ell^2} \right)^{-\alpha}$$

$\sigma_f^2 > 0, \ell > 0, \alpha > 0$. k_{RQ} is an infinite sum of k_{SE} with different ℓ . As $\alpha \rightarrow \infty$, $k_{RQ}(r) \rightarrow k_{SE}(r)$.

- ▶ Matérn:

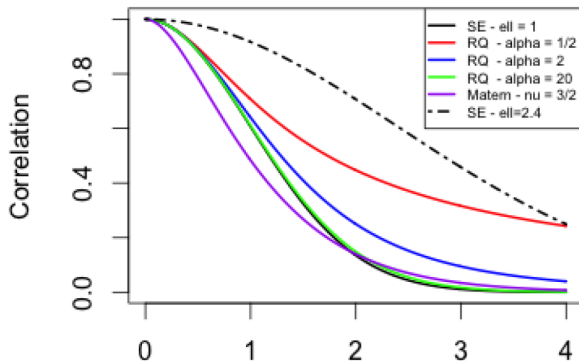
$$k_{Matern} = \sigma_f^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\ell} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}r}{\ell} \right)$$

where $\sigma_f^2 > 0, \ell > 0, \nu > 0$, and K_ν is the modified Bessel function. As $\nu \rightarrow \infty$, $k_{Matern}(r) \rightarrow k_{SE}(r)$.

- ▶ Demo of `GaussianProcesses.R` and `KernLabDemo.R`.

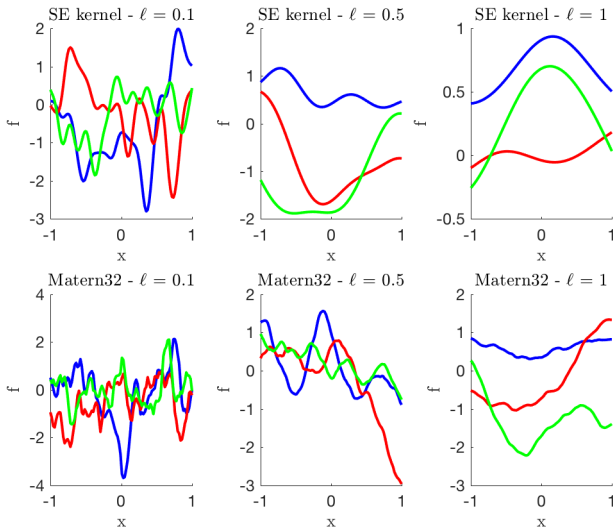
Three Common Covariance Functions

Correlation functions



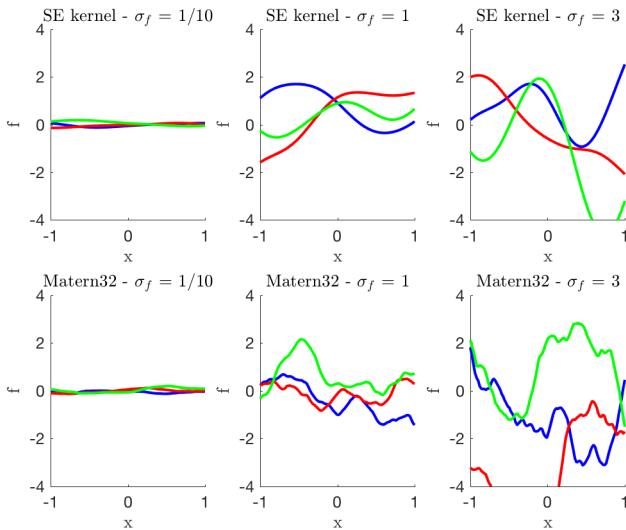
Three Common Covariance Functions

- The length scale ℓ determines the smoothness.



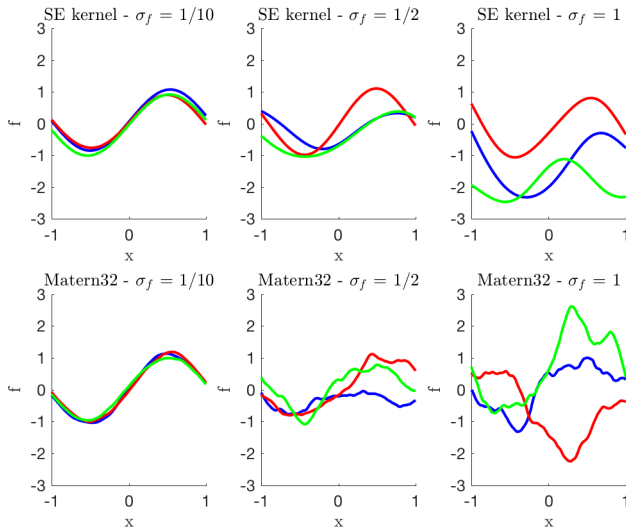
Three Common Covariance Functions

- The scale factor σ_f determines the variance.



Three Common Covariance Functions

- The mean can be arbitrary, e.g. $\sin(3x)$.



Learning the Hyperparameters of the Covariance Function

- ▶ Let θ denote the hyperparameters of the covariance function, i.e. $\theta = (\sigma_f, \ell)$ for k_{SE} , $\theta = (\sigma_f, \ell, \alpha)$ for k_{RQ} , and $\theta = (\sigma_f, \ell, \nu)$ for k_{Matern} .
- ▶ Choose the hyperparameters that maximize the marginal likelihood:

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

where $\mathbf{f}|\mathbf{X}, \theta \sim \mathcal{N}(0, K(\mathbf{X}, \mathbf{X}))$ and $\mathbf{y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma_n^2 I)$, which implies

$$\log p(\mathbf{y}|\mathbf{X}, \theta) = -\frac{1}{2} \mathbf{y}^T (K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I| - \frac{n}{2} \log 2\pi$$

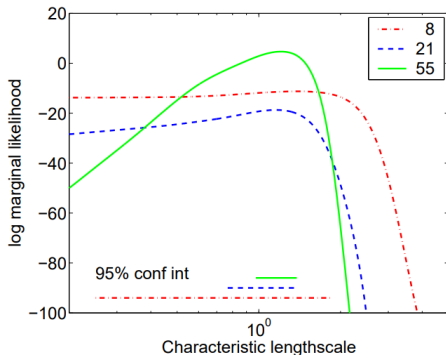
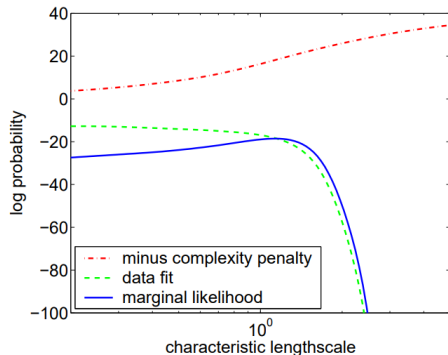
which alternatively can be obtained directly from

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 I & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right).$$

- ▶ In general, this is a non-convex optimization problem, and gradient methods are typically used. For most common covariance functions, the derivative of $K(\mathbf{X}, \mathbf{X})$ wrt θ is easy to compute.
- ▶ For a Bayesian approach, choose the hyperparameters that maximize the posterior distribution $p(\theta|\mathbf{y}, \mathbf{X}) \propto p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)$. It typically requires MCMC sampling or Laplace approximation.
- ▶ The methods above can also be used to select among covariance functions, i.e. simply include them as hyperparameters. Cross-validation is also an option.

Learning the Hyperparameters of the Covariance Function

$$\begin{aligned}\log p(\mathbf{y}|\mathbf{X}, \theta) &= -\frac{1}{2}\mathbf{y}^T (K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}| - \frac{n}{2} \log 2\pi \\ &= \text{data fit} - \text{model complexity} - \text{normalization constant}.\end{aligned}$$



More on Covariance Functions

- ▶ Anisotropic version of isotropic covariance function (i.e., it depends only on r) by setting $r^2 = (\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')$ where \mathbf{M} is positive definite.
- ▶ Automatic Relevance Determination: $\mathbf{M} = \text{diag}(\ell_1^{-2}, \dots, \ell_D^{-2})$, i.e. different length scales for different dimensions. In other words, ARM performs variable selection since a large ℓ_j means that the j -th dimension is essentially irrelevant for $f(\mathbf{x})$.
- ▶ Linear dimensionality reduction: $\mathbf{M} = \Lambda \Lambda^T + \Psi$ where Λ is a $D \times d$ matrix ($d < D$) whose columns define d directions of high relevance, and Ψ is a diagonal matrix capturing the axis aligned relevances.
- ▶ Periodic kernel with period d : $k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left\{ - \frac{2 \sin^2(\pi |\mathbf{x} - \mathbf{x}'|/d)}{\ell^2} \right\}$.
- ▶ The sum and product of two kernels is a kernel. For instance:
 - ▶ $k_{ARD}(\mathbf{x}, \mathbf{x}') = \prod_{d=1}^D k_{SE, \ell_d}(x_d, x'_d)$.
 - ▶ $k_{Periodic}(\mathbf{x}, \mathbf{x}') \times k_{SE}(\mathbf{x}, \mathbf{x}')$: Close peaks more dependent than distant ones.

More on Covariance Functions

- Assume x_1 is continuous (mg/week) and x_2 is binary (0=male, 1=female). Then,

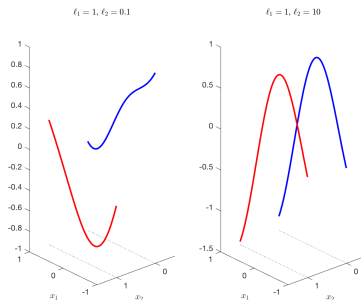
$$k_{ARD}((x_1, x_2), (x'_1, x'_2)) = \exp \left\{ - \frac{(x_1 - x'_1)^2}{2\ell_1^2} \right\} \exp \left\{ - \frac{(x_2 - x'_2)^2}{2\ell_2^2} \right\}$$

and thus

$$\text{cov}(f(x_1, 0), f(x_1, 1)) = \exp \left\{ - \frac{1}{2\ell_2^2} \right\}$$

which determines the similarity between the male and female profiles wrt x_1 , i.e. large (resp. small) ℓ_2 implies similar (resp. potentially different) profiles.

- For more than two categories, use one-hot encoding.



Gaussian Process Networks

- ▶ GPNs = continuous BNs + GPs. In other words,

$$p(\mathbf{y}) = \prod_i q(y_i | pa_G(y_i)) = \prod_i q(y_i | \mathbf{x}_i) \text{ where } y_i | \mathbf{x}_i \sim \mathcal{GP}(m(y_i), k(\mathbf{x}_i, \mathbf{x}'_i)).$$

- ▶ Learning means choosing the most likely DAG G given the data $d[1:N]$:

$$p(G | d[1:N]) = p(d[1:N] | G) p(G) / P(d[1:N]) \propto p(d[1:N] | G) p(G)$$

where

$$p(d[1:N] | G) = \prod_i p(y_i[1:N] | \mathbf{x}_i[1:N])$$

where

$$p(y_i[1:N] | \mathbf{x}_i[1:N]) = \mathcal{N}(m(y_i[1:N]), K(\mathbf{x}_i[1:N], \mathbf{x}_i[1:N])).$$

- ▶ Note that the last term involves hyperparameter learning, i.e. not really closed-form expression.
- ▶ Problem: The Lauritzen-Spiegelhalter algorithm does not work (?).
Solution: Approximate inference.

Bayesian Optimization

- ▶ Technique to find extreme values of an objective function. Particularly useful when the objective function is expensive to evaluate, no closed-form expression of it is available, its derivatives are not available, or it is non-convex. However, we should be able to evaluate the objective function at chosen values.
- ▶ Typical application: Hyperparameter learning.
- ▶ GPs are used to represent the posterior distribution over the objective function, which is used to choose the next point to evaluate. Specifically, the posterior mean of the GP (a.k.a. the surrogate function) can be seen as an estimate of the objective function. The next point to evaluate is chosen so as to maximize the so-called acquisition function.
- ▶ BO in a nutshell:
 - ▶ Assign GP prior over the unknown objective function f .
 - ▶ Evaluate f at some values $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - ▶ Update GP to the posterior $f|\mathbf{x}_1, \dots, \mathbf{x}_n, f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)$.
 - ▶ Use the GP posterior over f to find a new evaluation point \mathbf{x}_{n+1} by maximizing some acquisition function, which combines exploration with exploitation.
 - ▶ Iterate until the change in optimum is lower than some tolerance.

Bayesian Optimization

- ▶ Assume a maximization problem. Some acquisition functions follow.
- ▶ Upper confidence bound:

$$\mu(\mathbf{x}) + \kappa\sigma(\mathbf{x})$$

where κ is a user-defined parameter. Trade-off between exploration and exploitation.

- ▶ Probability of improvement:

$$p(f(\mathbf{x}) \geq f(\mathbf{x}^+)) = \Phi(z)$$

where $\Phi(\cdot)$ is the standard normal cumulative distribution function, $z = (\mu(\mathbf{x}) - f(\mathbf{x}^+))/\sigma(\mathbf{x})$, \mathbf{x}^+ is the current maximum, and $\mu(\mathbf{x})$ and $\sigma(\mathbf{x})$ are the GP posterior's mean and variance. Pure exploitation (?).

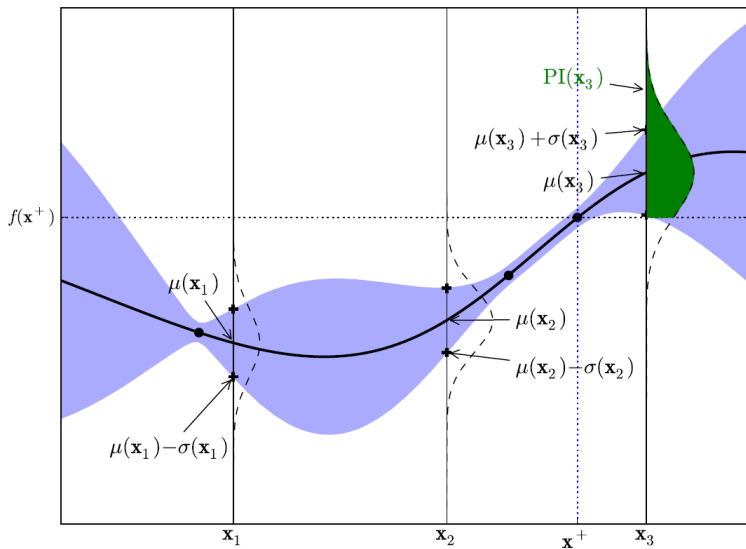
- ▶ Expected improvement:

$$(\mu(\mathbf{x}) - f(\mathbf{x}^+))\Phi(z) + \sigma(\mathbf{x})\mathcal{N}(z|0, 1)$$

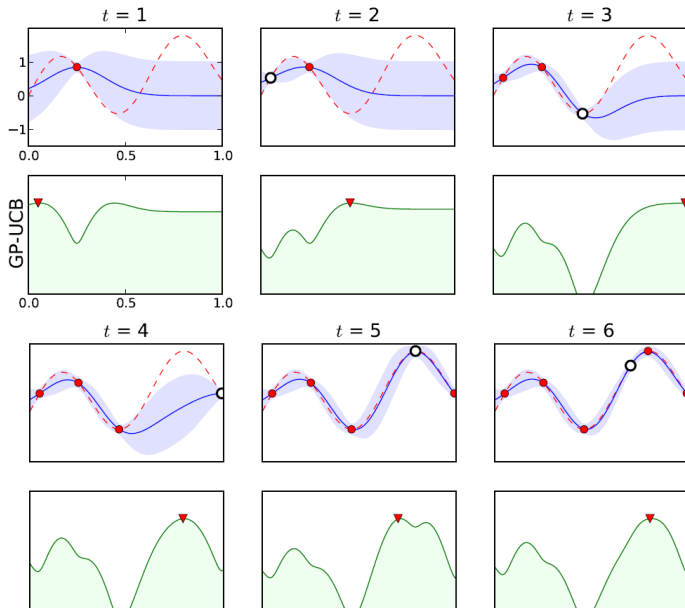
which combines exploration (choose points with high surrogate variance) with exploitation (choose points with high surrogate mean).

- ▶ Note that we still need to maximize the acquisition function. Despite being non-convex, this function is cheaper and simpler to maximize than the original one.

Bayesian Optimization



Bayesian Optimization



Contents

- ▶ Three Common Covariance Functions
- ▶ Learning the Hyperparameters of the Covariance Function
- ▶ More on Covariance Functions
- ▶ Gaussian Process Networks
- ▶ Bayesian Optimization

Thank you