

BAYESIAN MACHINE LEARNING

Exercise 4: Gaussian Processes

Prof. Yair Weiss

TA: Roy Friedman

Deadline: January 5, 2023

1 Theoretical

In this section we will look at the effects of the bandwidth parameter β on the prior induced by the Gaussian process. Consider the following Gaussian process priors:

- $p_1(f)$: $f \sim \mathcal{GP}(0, \exp[-\beta\|x - x'\|^2])$ where $\beta = \frac{1}{\alpha}$
 - $p_2(f)$: $f \sim \mathcal{GP}(0, \exp[-\beta\|x - x'\|^2])$ where $\beta = \alpha$
1. Let x_1 and x_2 be some vectors such that $\|x_1 - x_2\|^2 = \Delta$. Define $f_i = f(x_i)$ and $\mathbf{f} \triangleq (f_1, f_2)^T$. What is the correlation¹ between f_1 and f_2 under the two priors p_1 and p_2 when $\alpha \rightarrow \infty$?
 2. Suppose we observe data $\{y_1, \dots, y_N\}$ where we assume that $y_i = f(x_i)$ under prior p_2 . In addition, we notice that for every pair of neighbors x_1 and x_2 , the average distance is (again) $\|x_1 - x_2\|^2 = \Delta$ and the correlation between y_1 and y_2 is around $\text{corr}(y_1, y_2) \approx \frac{1}{4}$. Which value of α best reflects such a function under the prior p_2 ?
 3. Consider some set of vectors $\{x_i\}_{i=1}^M$ sampled uniformly from the 20 dimensional unit sphere and define $\mathbf{f} \triangleq (f_1, \dots, f_M)^T$. Consider the 2^M possible vectors \mathbf{f} such that $\forall i, f_i \in \{-1, 1\}$ and let $\alpha \rightarrow \infty$:
 - (a) Show that under p_1 there are two vectors that are equally likely, while all others have vanishing probability
 - (b) Show that under p_2 all vectors have the same probability

The Gram matrix defined by p_1 at the limit $\alpha \rightarrow \infty$ is a singular matrix (it is the all-ones matrix). This makes it difficult to talk about the Gaussian process at the limit $\alpha \rightarrow \infty$. To mitigate this, we will now introduce the Gaussian distribution when $\Sigma = \mathbf{1}\mathbf{1}^T$:

Rank Degenerate Gaussian let $p(x) = \mathcal{N}(x|0, \Sigma)$ be a Gaussian distribution over $x \in \mathbb{R}^d$ with the rank degenerate covariance matrix $\Sigma = \mathbf{1}\mathbf{1}^T$. Then, we will define the PDF of this distribution to be equal to²:

$$p(x) = p(x_1, \dots, x_d) = \mathcal{N}(x_1|0, 1) \prod_{i=2}^d \delta(x_1 - x_i) \quad (1.1)$$

$$\text{where } \delta(x_1 - x_i) = \begin{cases} 0 & x_1 \neq x_i \\ \infty & x_1 = x_i \end{cases}.$$

4. Find an analytical form for the evidence of a vector \mathbf{f} under p_1 and p_2 at the limit $\alpha \rightarrow \infty$.
 - (a) Let \mathbf{f}_1 be a vector of all ones and suppose we are using Bayesian model selection to choose between p_1 and p_2 . Which model will be selected given \mathbf{f}_1 ?
 - (b) Let \mathbf{f}_2 be a vector whose even components are 1 and odd components are -1. Which model will be selected given \mathbf{f}_2 ?

¹As a reminder, the correlation between two random variables x and y is given by: $\text{corr}(x, y) = \frac{\text{cov}[x, y]}{\sqrt{\text{var}[x]}\sqrt{\text{var}[y]}}$

²This definition is similar to the one in [Wikipedia for degenerate covariance matrices](#).

2 Practical

In this part of the assignment we will look at the effects different kernels have on Gaussian processes. Consider the following kernels:

- Laplacian kernel: $k_1(x, x') = \alpha \cdot e^{-\beta \|x - x'\|_1}$ with $\alpha, \beta > 0$
- RBF kernel: $k_2(x, x') = \alpha \cdot e^{-\beta \|x - x'\|^2}$ with $\alpha, \beta > 0$
- Gibbs' kernel: $k_3(x, x') = \sqrt{\frac{2\ell(x)\ell(x')}{\ell^2(x) + \ell^2(x')}} e^{-\frac{\|x - x'\|^2}{\ell^2(x) + \ell^2(x')}}}$ with $\ell(x) = \alpha \cdot e^{-\beta \|x - \delta\|^2} + \gamma$ for $\alpha, \beta, \gamma > 0$ and δ is any vector
- Neural network kernel³: $k_4(x, x') = \alpha \cdot \frac{2}{\pi} \sin^{-1} \left(\frac{2\beta(x^T x' + 1)}{\sqrt{(1 + 2\beta(1 + x^T x))(1 + 2\beta(1 + x'^T x'))}} \right)$ for $\alpha, \beta > 0$

1. Implement a Gaussian process model that receives as its input a user-specified kernel function and the sample noise⁴
2. For each of the kernels described above, choose 3 different parameter settings while keeping the sample noise at $\sigma^2 = 0.05$. For each of these settings, plot the mean and confidence interval of the prior in the interval $x \in [-5, 5]$. Sample 5 functions from the prior and plot them together with the confidence interval of the prior⁵
3. Consider the following 5 data points:

x	-2	-1	0	1	2
y	-2.1	-4.3	0.7	1.2	3.9

Calculate the posterior $\mathbf{f} | \{(x_1, y_1), \dots, (x_5, y_5)\}$ and plot the posterior mean and confidence intervals in the range $x \in [-5, 5]$ for each parameter setting chosen in the previous question. Sample 5 functions from the posterior and plot them together with the posterior mean

The purpose of the above two questions is to make you play around with the kernels a bit. As you probably noticed, the last two kernels are a bit weird - try to play around with them and get them to act like you want them to. You'll see that they are pretty expressive and can form rather weird functions. The Gibbs kernel is a type of generalized RBF kernel (if you set α to be pretty low and γ to be pretty high, then γ will act like the typical length-scale of the RBF) while the NN kernel is what you would get if you try to make predictions using a (fully connected) neural network with infinitely-wide layers and an error-function (sigmoid) activation function (if you're interested, read more about it in [Rasmussen and Williams](#)). Notice that in the supplied skeleton for your solution, we already did (most) of the heavy lifting for the plotting, so all you have to do is set the kernel parameters you want to see and run the code. We *really* recommend that you play around with the kernels, this will give you more of an intuition than anything we can show in class.

4. For the data points defined above, plot the evidence function of the RBF kernel with 101 β s evenly spaced in the range $\beta \in [1, 15]$ and noise $\sigma^2 = 0.15$ as a function of β . Which value of β has the highest evidence score? Plot the points, posterior mean and confidence interval of the best, worst and median β s according to the evidence

³See section 4.2.3 of [Rasmussen and Williams](#) for more information

⁴See Algorithm 2.1 on page 37 of [Rasmussen and Williams](#) for pseudo-code of a numerically stable version

⁵You can use the [visualization tool from Moodle](#) as a basic guideline for how these plots should look, although most of the plotting work is already there in the skeleton

3 Submission Guidelines

Submit a single zip file named “ex4_<YOUR ID>.zip”. This file should contain your code, along with an “ex4.pdf” file in which you should write your answers to the theoretical part and add the figures/text for the practical part. Please write readable code, as the code will also be checked manually (and you may find it useful in the following exercises). In the submitted code, please make sure that you write a basic main function in a file named “ex4.py” that will run (without errors) and produce all of the results that you showed in the pdf of answers that you submitted. The only packages you should use are `numpy`, `scipy` and `matplotlib`. You may also reuse code from your previous exercise in order to answer the questions in this exercise, if needed.

In general, it is better if you type your homework, but if you prefer handwriting your answers, please make sure that the text is readable when you scan it.

Part of your assignment will be graded by submitting your answers through Moodle, at [this link](#). In each of the questions, write the answer to the corresponding question for grading. These answers will be graded automatically, so write only numeric values where needed.

4 Supplementary Code

We have supplied an outline code which you can use to get started in `ex4.py`. You don’t have to use the format we outlined, but your code must run without errors and you must submit the plots required in the exercise description.

Good luck!