

CSE 515T (Spring 2017) Midterm

- There are two ways to hand in this midterm. **Late submissions will not be accepted!** I do not recommend cutting it too close.
 - Physically in class. The due date for this option is **4:00 PM, Thursday, 9 March**.
 - Electronically on Piazza as a private message to the instructors. The due date for this option is **23:59 CST Friday, 10 March** (the “midnight” that occurs after the full day of Friday passes in the central time zone).
- Please do not discuss the questions with other members of the class.
- Please post any questions as a *private message to the instructors* on Piazza.
- Any corrections will be posted by the instructors on Piazza. This document will also be kept up-to-date on the course webpage and in GitHub.

We will consider a series of questions relating to an application of Bayesian inference to numerical analysis, specifically quadrature.¹

We are going to consider the function

$$f(x) = \exp(-x^2)$$

and its definite integral

$$Z = \int_{-\infty}^{\infty} f(x) dx.$$

The function f has no elementary antiderivative, so the calculation of Z is not straightforward. There is a famous method for computing Z with the trick of considering Z^2 instead, rewriting the resulting $2d$ integral in polar coordinates, and making a convenient substitution.² If you haven't seen this, it's beautiful and worth checking out. The result is

$$Z = \sqrt{\pi}.$$

We will consider modeling f with a Gaussian process prior distribution:

$$p(f) = \mathcal{GP}(f; \mu, K),$$

and conditioning on the following set of data $\mathcal{D} = (\mathbf{x}, \mathbf{y})$:

$$\mathbf{x} = [-3, -2, -1, 0, 1, 2, 3]^\top;$$

$$\mathbf{y} = \exp(-\mathbf{x}^2)$$

$$= [1.2341 \times 10^{-4}, 1.8316 \times 10^{-2}, 0.36788, 1, 0.36788, 1.8316 \times 10^{-2}, 1.2341 \times 10^{-4}]^\top.$$

We will fix the prior mean function μ to be identically zero; $\mu(x) = 0$.

1. First, let us consider the question of model, specifically kernel, selection.

Consider the following four choices for the covariance function K :

$$K_1(x, x') = \exp(-|x - x'|^2)$$

$$K_2(x, x') = \exp(-|x - x'|)$$

$$K_3(x, x') = (1 + \sqrt{3}|x - x'|) \exp(-\sqrt{3}|x - x'|)$$

$$K_4(x, x') = \delta(x - x'),$$

where δ is the Kronecker delta function. Note that I am not parameterizing any of these kernels; please consider them to be fixed as given.

Each kernel defines a Gaussian process model for the data in a natural way:

$$p(f | \mathcal{M}_i) = \mathcal{GP}(f; \mu, K_i).$$

Consider a uniform prior distribution over these models:

$$\Pr(\mathcal{M}_i) = 1/4 \quad i = 1, 2, 3, 4.$$

¹This is my attempt to reinforce some fascinating concepts to make up for a somewhat confusing lecture. I hope these exercises help clarify things!

²This is commonly credited to Gauss, but the idea goes back at least to Poisson.

- (a) Compute the log model evidence for each model given the data \mathcal{D} above.
 - (b) Compute the model posterior $\Pr(\mathcal{M} \mid \mathcal{D})$.
 - (c) Can you find a kernel with higher model evidence given the data above? (I will award an extra credit point to the person who provides the kernel with the highest evidence.)
2. Now let's turn to prediction.
- (a) For each kernel above, plot the predictive distribution over the interval $x^* \in [-5, 5]$. For each model \mathcal{M}_i , please plot, in a separate figure, the predictive mean $p(y^* \mid x^*, \mathcal{D}, \mathcal{M}_i)$ and a 95% credible interval. These plots should be the result of a computer program. Please add legends and axes labels, etc., and plot the true function on the same interval for reference.
 - (b) In addition, please write out the predictive mean and standard deviation at $x^* = 2.5$ for each of the kernels, $p(y^* \mid x^* = 2.5, \mathcal{D}, \mathcal{M}_i)$.
 - (c) What is the model-marginal predictive distribution, $p(y^* \mid x^*, \mathcal{D})$? Write this in terms of the model-conditional predictive distribution and the model posterior.
 - (d) Assume that the model posterior is uniform; $\Pr(\mathcal{M}_i \mid \mathcal{D}) = 1/4$ for all models i (this is *not* the case, if you are worried about your answer to 1(b)). Plot the model-marginal predictive mean function $\mathbb{E}[y^* \mid \mathcal{D}]$ over the interval $x^* \in [-5, 5]$.
3. Let us consider conditioning a Gaussian process on the observation of a derivative.
- (a) Write down the joint distribution implied by an arbitrary Gaussian process distribution on a function $f: \mathbb{R} \rightarrow \mathbb{R}$, $\mathcal{GP}(f; \mu, K)$, between the function value at an arbitrary point x , $f(x)$, and the value of a derivative at another arbitrary point x' , $f'(x) = \frac{df}{dx}|_{x'}$.
 - (b) Fix the zero mean function $\mu(x) = 0$ and the squared exponential covariance $K_1(x, x')$ from above. Write down the joint distribution from part (a), evaluating any derivatives or integrals you may encounter. (I do not expect the expression K_1 to appear in your answer, but maybe things like $\exp(\dots)$ instead. That is, I want explicit formulas.)
 - (c) Show that for this model, the function value at a point x and the value of the derivative at that point are uncorrelated *a priori*.
 - (d) Consider the Gaussian process from part 1 with covariance function K_1 . Condition the model on the data \mathcal{D} given previously as well as the observation that the derivative is 0 at $x = 0$, finding $p(f \mid \mathcal{D}, f'(0) = 0, \mathcal{M}_1)$. Prepare a plot of the predictive distribution on the interval $x^* \in [-5, 5]$ as before.
4. Now we will consider integration.
- Perform Bayesian quadrature to estimate the definite integral $\int_{-5}^5 f(x) dx$, using the model \mathcal{M}_1 from question 1. What is the predictive mean and standard deviation, $p(Z \mid \mathcal{D}, \mathcal{M}_1)$? (You may give a numeric answer.) How does this compare with the true answer?
5. Finally, we will consider a decision problem. Suppose we have already made some observations \mathcal{D} . How can we select the *most-informative* next observation $(x^*, f(x^*))$ to make? This is a decision problem where the action space (parametrizing the next observation location) is the domain, $x^* \in \mathcal{X}$.

Suppose that we are to estimate Z with a point estimate \hat{Z} , and that we have selected the squared loss

$$\ell(Z, \hat{Z}) = (Z - \hat{Z})^2.$$

- (a) Given a set of observations \mathcal{D} , what is the Bayesian optimal action? What is the expected loss of that action?
- (b) Compute the expected loss of the Bayesian optimal action after adding a new observation to \mathcal{D} located at a point x^* . Plot this result as a function of $x^* \in [-5, 5]$. What is the optimal location to measure the function next? (By symmetry there may be multiple equivalent answers.)
- (c) Condition the function on an observation of the function at the chosen location and plot the predictive distribution as in part 2(a). Recompute the predictive distribution for Z . Did our estimate improve?
- (d) How does the answer to the above part relate to the observation $f'(0) = 0$ we considered above? Which observation is preferable?