

02477 Bayesian Machine Learning 2024: Assignment 3

This is the last assignment out of three in the Bayesian machine learning course 2024. The assignment is a group work of 3-5 students (please use the same groups as in assignment 1,2 if possible) and hand in via DTU Learn). The assignment is **mandatory**. The deadline is **5th of May 23:59**.

Part 1: Fully Bayesian inference for Gaussian process regression

In this part, we will extend our Gaussian process regression analysis of the bike sharing dataset (from week 5) to fully Bayesian inference on hyperparameter level. Use the code below (Fig. 1) to load and preprocess the data.

```
# load data from disk
data = np.load('./data_exercise5b.npz')
X = data['day']
y = np.log(data['bike_count'])

# remove mean and scale to unit variance
ym, ys = np.mean(y), np.std(y)
y = (y-ym)/ys
```

Figure 1: Code for loading and preprocessing the bike sharing dataset.

The Gaussian process regression model for the dataset $\mathcal{D} = \{x_n, y_n\}_{n=1}^N$ is given below:

$$y_n = f(x_n) + \epsilon_n, \quad (1)$$

where $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$ and $f(x) \sim \mathcal{GP}(0, k(x, x'))$ for $k(x, x') = \kappa^2 \exp\left(-\frac{\|x-x'\|_2^2}{2\ell^2}\right)$.

We will impose the following prior distributions on the hyperparameters:

$$\begin{aligned} \kappa &\sim \mathcal{N}_+(0, 1) \\ \ell &\sim \mathcal{N}_+(0, v) \\ \sigma &\sim \mathcal{N}_+(0, 1), \end{aligned}$$

where $v > 0$ is a positive constant (which you will determine in the next task) and $\mathcal{N}_+(m, v)$ is the half-normal distribution. These assumptions lead to the following joint distribution

$$\begin{aligned} p(\mathbf{y}, \mathbf{f}, \sigma, \kappa, \ell) &= p(\mathbf{y}|\mathbf{f}, \sigma^2)p(\mathbf{f}|\kappa, \ell)p(\kappa)p(\ell)p(\sigma) \\ &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})p(\mathbf{f}|\mathbf{0}, \mathbf{K})\mathcal{N}_+(\kappa|0, 1)\mathcal{N}_+(\ell|0, v)\mathcal{N}_+(\sigma|0, 1). \end{aligned}$$

In this exercise, we want to impose a prior that prevent the lengthscale from becoming too large.

Task 1.2: Choose a value for v such that the prior probability of observing a lengthscale larger than 100 is approximately 1%, i.e. $p(\ell > 100) \approx 0.01$.

Hints: You can do this in several ways, e.g. numerically or analytically

Task 1.3: Compute the distribution $p(\mathbf{y}, \sigma, \kappa, \ell)$ by marginalizing out \mathbf{f}

The next goal is to approximate the posterior distribution over the hyperparameters, i.e. $p(\kappa, \ell, \sigma|\mathbf{y})$, using a Metropolis-sampler. Define $\theta = \{\kappa, \ell, \sigma^2\}$ to be the set of hyperparameters of the model. Since the scale of the hyperparameters are quite different, we will use an anisotropic proposal distribution:

$$q(\theta^*|\theta^{k-1}) = \mathcal{N}(\theta^*|\theta^{k-1}, \Sigma) \quad \text{for} \quad \Sigma = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 100 & 0 \\ 0 & 0 & 0.01 \end{bmatrix}. \quad (2)$$

That is, the proposed step-size will generally be larger for the lengthscale dimension and so on.

Task 1.4: Implement a Metropolis sampler using the proposal distribution in eq. (2) for generating samples from the posterior $p(\kappa, \ell, \sigma | \mathbf{y})$. Run 4 chains for 10000 iterations each. Discard the warm-up samples and report the number of samples discarded.

Task 1.5: Plot the trace for each parameter and report the convergence diagnostics \hat{R} and S_{eff} for each parameter.

Task 1.6: Estimate and report the posterior mean for each hyperparameter. Report the MCSE for each estimate.

Task 1.7: Estimate a 95% posterior credibility interval for each hyperparameter.

Part 2: Variational inference, KL-divergences and entropy

Variational inference is a flexible tool approximating posterior distributions $p(\mathbf{w} | \mathbf{y})$, where $\mathbf{w} \in \mathbb{R}^D$ is a set of parameters to be estimated given some data \mathbf{y} . The central object in variational inference is the evidence lower bound (ELBO) given by

$$\mathcal{L}[q] \equiv \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{w})] - \mathbb{E}_q[\ln q(\mathbf{w})]. \quad (3)$$

The first term of the ELBO depends both on the model $p(\mathbf{w}, \mathbf{y}) = p(\mathbf{y} | \mathbf{w})p(\mathbf{w})$ and the variational approximation $q(\mathbf{w})$. The second term only depends on the variational approximation q and is equivalent to the *entropy* of $q(\mathbf{w})$, i.e. $\mathcal{H}[q] \equiv -\mathbb{E}_q[\ln q(\mathbf{w})]$. In this assignment, we will study the significance and interpretation of the entropy term in the context of mean-field Gaussian variational families

$$q(\mathbf{w}) = \prod_{i=1}^D q(w_i) = \prod_{i=1}^D \mathcal{N}(w_i | \mu_i, \sigma_i^2). \quad (4)$$

That is, all members in our variational family \mathcal{Q} can be written as in eq. (4), where $\{\mu_i, \sigma_i^2\}_{i=1}^D$ are the *variational parameters*.

For many variational families, including mean-field Gaussians, the entropy $\mathcal{H}[q(\mathbf{w})]$ can be computed analytically, and it is the purpose of the next few tasks to derive this quantity analytically.

Task 2.1: Assuming a mean-field Gaussian variational family, show that the entropy of the variational approximation $q(\mathbf{w})$ is equal to the sum of the marginal entropies, i.e. $\mathcal{H}[q(\mathbf{w})] = \sum_{i=1}^D \mathcal{H}[q(w_i)]$.

Hints: Use the definition of entropy and insert the mean-field family from eq. (4). Use the fact that the distribution $q(\mathbf{w})$ factorizes to simplify the expression and obtain the desired result.

Task 2.2: Assuming a mean-field Gaussian variational family, derive the analytical expression for the entropy of the marginal distribution $q(w_i)$ is $\mathcal{H}[q(w_i)] = \frac{1}{2} \log(2\pi\sigma_i^2) + \frac{1}{2}$.

Hints: Marginal distributions are easy to obtain for mean-field families, i.e. the approximate marginal distribution of w_j is simply $q(w_j) = \mathcal{N}(w_j | \mu_j, \sigma_j^2)$.

Task 2.3: Assuming a mean-field Gaussian variational family, show that $\mathcal{H}[q(\mathbf{w})] = \frac{1}{2} \sum_{i=1}^D \ln(2\pi e \sigma_i^2)$.

Hints: Combine the results from the previous two tasks. The constant e is the base of the natural logarithm.

This result implies that the ELBO objective for mean-field Gaussian families can be written as follows

$$\mathcal{L}[q] \equiv \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{w})] + \frac{1}{2} \sum_{i=1}^D \ln(2\pi e \sigma_i^2). \quad (5)$$

When fitting our variational approximation $q(\mathbf{w})$ by maximizing eq. (5), the two terms above will each ‘favor’ different solutions. Next, we will investigate which types of distributions $q(\mathbf{w})$ that are ‘encouraged’ or favored by the entropy term.

Task 2.4: Compute the partial derivative of the entropy term with respect to σ_j^2 , i.e. $\frac{\partial}{\partial \sigma_j^2} \mathcal{H}[q(\mathbf{w})]$, and argue that the gradient is always strictly positive.

Task 2.5: Which types of mean-field Gaussian distributions $q(\mathbf{w})$ are encouraged by the entropy-term?

Hint: When maximizing the ELBO using gradient-based methods, each of the two terms of the ELBO will contribute to the gradient. What would happen to q if the first term in the ELBO was not present?

Next, we will rewrite the ELBO to enable a different interpretation of the objective. To do this, we will focus on a broad class of models for supervised learning

$$p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{w})p(\mathbf{w}), \quad (6)$$

where $p(y_n|\mathbf{w})$ is the likelihood for the n ’th data point and $p(\mathbf{w})$ is the prior distribution. This model family includes regression, classification, neural networks, Gaussian processes etc.

Task 2.6: Assume a probabilistic model of the form in eq. (6). Show that the ELBO can be written as in eq. (7), where the KL-divergence is now between the approximation $q(\mathbf{w})$ and the prior $p(\mathbf{w})$.

$$\mathcal{L}[q] = \sum_{n=1}^N \mathbb{E}_q[\ln p(y_n|\mathbf{w})] - \text{KL}[q(\mathbf{w})||p(\mathbf{w})] \quad (7)$$

Task 2.7: Which types of approximate distributions $q(\mathbf{w})$ are encouraged by the KL-term in eq. (7)?

Hint: When maximizing the ELBO using gradient-based methods, each of the two terms of the ELBO will contribute to the gradient. What would happen to q if the first term (the entire sum) in the ELBO was not present? When is the KL-divergence minimized?

Task 2.8: Which mean-field Gaussian distributions $q(\mathbf{w})$ are encouraged by the first term, i.e. the expected log likelihood? Argue why if you can.

Task 2.9: Derive the analytical expression (i.e solve the integral) for the KL-divergence $\text{KL}[q_1||q_2]$ between univariate Gaussian distributions $q_1(w) = \mathcal{N}(w|\mu_1, \sigma_1^2)$ and $q_2(w) = \mathcal{N}(w|\mu_2, \sigma_2^2)$.

Part 3: Regression modelling using mixture of experts

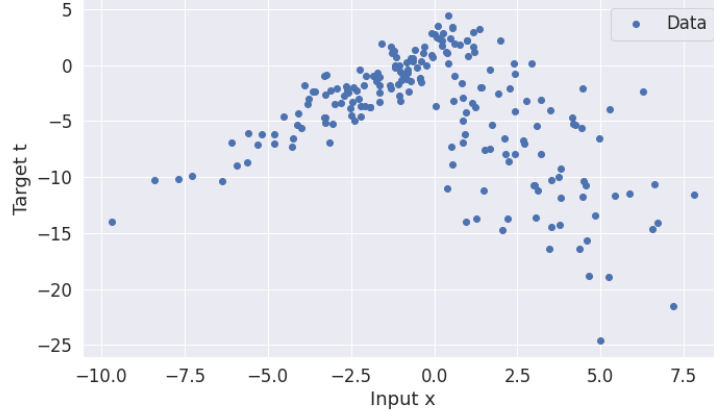


Figure 2: Dataset for regression

Consider the dataset for regression given in Figure 2. It is evident from the plot that there is a strong non-linear dependency between the target variables y_n and the input variables x_n . We also observe that the noise variance seems to depend on the input x . Consequently, assuming $y_n = y(x_n) + e_n$, where e_n is independent and identically distributed (i.i.d) noise would not be appropriate.

In this part, we will work with a so-called *Mixture of experts* (MoE) model for regression. The core idea is to model the dataset using several submodels, which are ‘experts’ in their own region of the input space. For example, the dataset in Figure 2 could be well-represented using two linear models: a linear model with positive slope and relative low noise variance for the ‘left half’ of the dataset and another linear model with negative slope and larger noise variance on the ‘right half’ of the dataset. Our goal is to simultaneously learn these two linear models as well as when to use which model.

To construct this model, we will consider two different linear models and then introduce a latent binary variable for each data point, $z_n \in \{0, 1\}$, to control which of the two linear models that should explain the data point y_n . That is, if $z_n = 0$ we assume y_n should be explained a linear model with parameters \mathbf{w}_0 and if $z_n = 1$, then y_n should be explained by a linear model with weights \mathbf{w}_1 .

We model each z_n using Bernoulli distributions as follows

$$p(z_n|\pi_n) = \text{Ber}(z_n|\pi_n), \quad (8)$$

where $\pi_n = \sigma(\mathbf{v}^T \mathbf{x}_n)$ is the probability of $z_n = 1$, $\sigma(\cdot)$ is the logistic sigmoid function and \mathbf{v} is a parameter vector to be estimated. That is, we basically model the latent z_n variable using a logistic regression model. We can set up a conditional likelihood as follows

$$p(y_n|x_n, z_n, \mathbf{w}_0, \mathbf{w}_1, \sigma_0^2, \sigma_1^2) = \begin{cases} \mathcal{N}(y_n|\mathbf{w}_1^T \mathbf{x}_n, \sigma_1^2) & \text{if } z_n = 1, \\ \mathcal{N}(y_n|\mathbf{w}_0^T \mathbf{x}_n, \sigma_0^2) & \text{if } z_n = 0, \end{cases} = \mathcal{N}(y_n|\mathbf{w}_{z_n}^T \mathbf{x}_n, \sigma_{z_n}^2) \quad (9)$$

where we also allow the noise variance to depend on z_n . We complete the model with generic priors

$$\tau, \sigma_0^2, \sigma_1^2 \sim \mathcal{N}_+(0, 1) \quad (10)$$

$$\mathbf{w}_0, \mathbf{w}_1, \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}) \quad (11)$$

$$z_n|\mathbf{v} \sim \text{Ber}(\sigma(\mathbf{v}^T \mathbf{x}_n)) \quad (12)$$

$$y_n|z_n \sim \mathcal{N}(\mathbf{w}_{z_n}^T \mathbf{x}_n, \sigma_{z_n}^2), \quad (13)$$

where $\mathcal{N}_+(0, 1)$ is the half-normal distribution. This leads to a joint model of the form

$$p(\mathbf{y}, \mathbf{w}_1, \mathbf{w}_0, \mathbf{v}, \mathbf{z}, \tau, \sigma_0, \sigma_1) = \left[\prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}_{z_n}^T \mathbf{x}_n, \sigma_{z_n}^2) \text{Ber}(z_n | \sigma(\mathbf{v}^T \mathbf{x}_n)) \right] \mathcal{N}(\mathbf{w}_0 | \mathbf{0}, \tau^2 \mathbf{I}) \mathcal{N}(\mathbf{w}_1 | \mathbf{0}, \tau^2 \mathbf{I}) \mathcal{N}(\mathbf{v} | \mathbf{0}, \tau^2 \mathbf{I}) \mathcal{N}_+(\tau^2 | 0, 1) \mathcal{N}_+(\sigma_0 | 0, 1) \mathcal{N}_+(\sigma_1 | 0, 1). \quad (14)$$

The purpose is now to fit this model to the dataset in Figure 2 using MCMC. The dataset can be found on DTU Learn as loaded as follows:

```
data = np.load('./data_assignment3.npz')
x, y = data['x'], data['t']
```

To avoid handling the intercepts in the linear models explicitly, we will use the convention $\mathbf{x}_n = [x_n, 1]$.

Task 3.1: Marginalize out each z_n from to joint model in eq. (14) to obtain a joint distribution, where the likelihood for each observation is a mixture of two Gaussian distributions.

Hints: Use the sum rule to marginalize out z_n .

Task 3.2: Implement a Python function to evaluate the marginalized log joint distribution

Hints: Let θ denote all the parameters of the model, i.e. $\theta = \{\mathbf{w}_0, \mathbf{w}_1, \mathbf{v}, \tau, \sigma_0^2, \sigma_1^2\}$, then implement a function that takes θ and returns $\ln p(\mathbf{y}, \mathbf{w}_1, \mathbf{w}_0, \mathbf{v}, \tau, \sigma_0, \sigma_1)$. The half-normal distribution can be implemented as follows

```
from scipy.stats import norm
def log_halfnormal(x):
    if x < 0: # negative values not supported
        return -np.Inf
    else:
        return np.log(2) + norm.logpdf(x, 0, 1)
```

Task 3.3: Run a Metropolis-Hasting sampler to infer all parameters. Explain the settings you used (number of iterations, proposal distribution etc).

Hints: Feel free to use the Metropolis-Hasting code from the exercises. Plot the trace for all parameters to assess convergence (convergence diagnostics might be tricky due to the multimodality of the model). Compute the posterior mean of the weights $\mathbf{w}_0, \mathbf{w}_1, \mathbf{v}$, and plot the resulting function on top of the data as a sanity check. Make sure to initialize the sampler using a valid parameter vector, i.e. positive scale parameters $\sigma_0, \sigma_1, \tau > 0$.

Task 3.4: Report the posterior mean and 95% credibility intervals for all parameters.

Task 3.5: For $x \in [-12, 12]$, plot posterior predictive distribution for $p(\pi^* | \mathbf{y}, \mathbf{x}^*)$, $p(y^* | \mathbf{y}, \mathbf{x}^*, z^* = 0)$ and $p(y^* | \mathbf{y}, \mathbf{x}^*, z^* = 1)$ on top of the data.

Hints: For inspiration on how to plot these distributions using samples, see the exercise on Gibbs sampling for change point detection.

Task 3.6: For $x \in [-12, 12]$, plot posterior predictive distribution for $p(y^* | \mathbf{y}, \mathbf{x}^*)$

Task 3.7: Suppose we were analyzing another dataset requiring three experts rather two. Describe how you would adapt the model and write the corresponding probabilistic model in the same format as in eq. (10)–(13).

Hints: A softmax function is required somewhere.