

# TEK5040 Mandatory Assignment 3

## Bayesian Deep Learning

October 31, 2024

### 1 Introduction

This assignment deals with Bayesian deep learning with neural networks. You are going to investigate Monte Carlo Dropout method and Variational Inference methods in this exercise.

#### 1.1 Package Contents

<code>mc_train_test.py</code>	Tensorflow script for training and testing of Monte Carlo Dropout.
<code>vi_train_test.py</code>	Tensorflow script for training and testing for Variational Inference
<code>common.py</code>	Common routines and definitions, including data generation
<code>densevariational.py</code>	Tensorflow script defining a layer supporting Variational inference
<code>bayesian_exercise.pdf</code>	This document describing the task

In addition there is a `requirement.txt` file for a setup based on `tensorflow ver 2.18`.

### 2 Operation

We use a simple regression task for evaluating Bayesian deep learning techniques. The training data set is generated using `gen_data()` in `common.py`. Figure 1 shows the generated training data which follow a noisy sinusoidal. The task is to predict values of  $y$  axis given values of  $x$ -axis. We use Bayesian learning to generate the variance of the predicted values.

You can run `mc_train_test.py` in a Tensorflow environment to train and test the system based on Monte Carlo Drop out. i.e. `python3 mc_train_test.py`. Make sure that you have `tensorflow-probability` installed using for example `pip`.

```
pip3 install tensorflow-probability==x.x.x.
```

where `x.x.x` is the version which should be selected so that it matches your `tensorflow` version (see <https://github.com/tensorflow/probability/releases> for version compatibility).

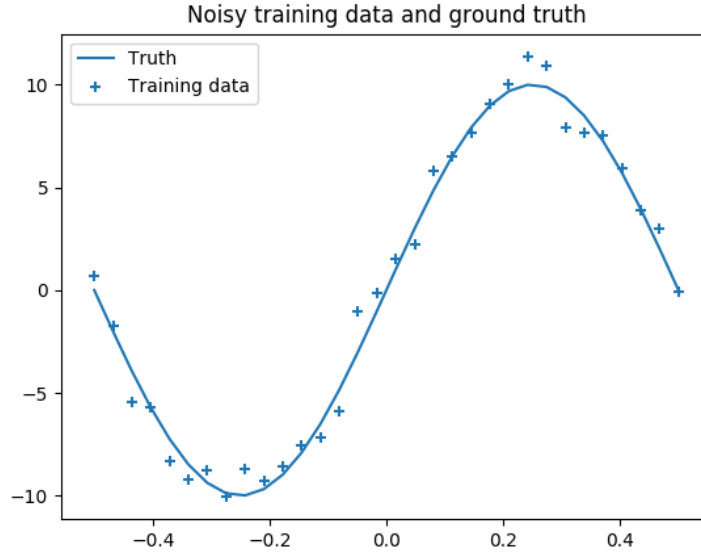


Figure 1: Training data

**NOTE:** `tensorflow_probability` depends on `keras` version 2. Newer `tensorflow` versions (i.e. versions greater than 2.16) use `keras` version 3, and therefore in such cases install `tensorflow_probability` using the command `pip3 install tensorflow-probability[tf]`. This will install the correct `keras` version, i.e. version 2, in addition.

You may also need to install `tqdm` and `matplotlib`.

Similarly you can run `python3 vi_train_test.py` for Variational Inference based training and testing.

If the scripts work correctly, you should see the typical Keras output of training losses and the progress bar of testing for 500 times. Further, it shows the predicted results in graphical form.

## 3 Theory and Implementation

### 3.1 Monte Carlo Dropout

This is a very straight-forward implementation where dropout layers in the architecture absorb the uncertainty. Dropout layers are used in testing as well to calculate the output variance.

### 3.2 Variational Inference

We use *Bayes by Backprop* approach with the *reparameterization trick*. `densevariational.py` implements a layer supporting these techniques. In the implementation, we assume the prior probability of the parameters (weights and biases) of the network

$$p(\mathbf{w}) = \prod_i p(w_i) = \prod_i [\pi_1 \mathcal{N}(w_i|0, \sigma_1) + (1 - \pi_1) \mathcal{N}(w_i|0, \sigma_2)] \quad (1)$$

where  $\mathcal{N}(w_i|0, \sigma)$  is a normal (Gaussian) distribution with zero mean and variance  $\sigma$  for each network parameter  $w_i$ . Further,  $\pi_1, \sigma_1$  and  $\sigma_2$  are fixed hyper-parameters.

We also assume that each network parameter has an associated variational distribution  $q(w_i)$  which is also assumed to be a normal (Gaussian) distribution  $\mathcal{N}(w_i|\mu_i, \sigma_i)$ . Further, we assume that the network parameters are independent, that is

$$q(\mathbf{w}) = \prod_i q(w_i) = \prod_i \mathcal{N}(w_i|\mu_i, \sigma_i) \quad (2)$$

The parameter vector of the variational distribution  $\lambda = (\mu_i, \sigma_i), i = 1, 2, \dots, N$ , with  $N$  being the number of parameters. As usual, we approximate the posterior network parameter distribution  $p(\mathbf{w}|\mathcal{D})$  with the variational distribution  $q(\mathbf{w})$  and aim to find the parameter set  $\lambda$  which minimizes the KL distance between the posterior and variational distributions.

As described in the lecture, KL distance is minimized indirectly by maximizing the Evidence Lower Bound (ELBO) which is given by

$$\hat{\mathcal{L}}(\lambda) = \frac{1}{S} \sum_{s=1}^S [\ln p(\mathbf{w}(\lambda, \boldsymbol{\epsilon}^s), \mathcal{D}) - \ln q(\mathbf{w}(\lambda, \boldsymbol{\epsilon}^s), \lambda)]$$

where each element of  $\boldsymbol{\epsilon}^s$  is drawn from a standard normal distribution and  $\mathbf{w}^s = \mathbf{w}(\lambda, \boldsymbol{\epsilon}^s)$  represents the re-parameterization trick. The above expression can also be written as

$$-\hat{\mathcal{L}}(\lambda) = \frac{1}{S} \sum_{s=1}^S [\ln q(\mathbf{w}^s, \lambda) - \ln p(\mathcal{D}|\mathbf{w}^s) - \ln p(\mathbf{w}^s)] \quad (3)$$

Then we use back-propagation to minimize the loss  $-\hat{\mathcal{L}}(\lambda)$ .

The first term and the last term of equation 3 is evaluated using equation 2 and equation 1 respectively. These terms are data independent and therefore they are added to the model losses. The second term in equation 3 is data dependent and it is the negative log probability of data for the given parameter vector  $\mathbf{w}^s$ . This quantity is typically available after a forward pass through the network.

## 4 Task

1. Run `mc_train_test.py` and comment on the variance of the output as depicted in the generated figure. Hint: Run the script several times to get a feeling of how it looks like in general. Note the difference in test and training sets, i.e. training set has 32 samples in the range of  $(-0.5, 0.5)$  while the test set has 1500 samples in the range of  $(-1.5, 1.5)$ .
2. Run `vi_train_test.py` and comment on the variance of the output. Compare it with that of MC-dropout.
3. Find the line numbers of the code in `densevariational.py` which implements reparameterization-trick  $\boldsymbol{w}^s = \boldsymbol{w}(\lambda, \boldsymbol{\epsilon}^s)$ .
4. Find the line numbers of the code in `densevariational.py` which adds the first and the third loss components in equation 3 to the computational graph.
5. Find the line numbers of the code where the second loss component in equation 3 is added to the computation graph.
6. Markov Chain Monte Carlo (MCMC) is an alternative to variational inference. Give an advantage and a disadvantage of MCMC.