

CS 6190: Probabilistic Modeling Spring 2019

Homework 5

Handed out: 21 Nov, 2019
Due: 11:59pm, 8 Dec, 2019

- You are welcome to talk to other members of the class about the homework. I am more concerned that you understand the underlying concepts. However, you should write down your own solution. Please keep the class collaboration policy in mind.
- Feel free discuss the homework with the instructor or the TAs.
- Your written solutions should be brief and clear. You need to show your work, not just the final answer, but you do *not* need to write it in gory detail. Your assignment should be **no more than 10 pages**. Every extra page will cost a point.
- Handwritten solutions will not be accepted.
- The homework is due by **midnight of the due date**. Please submit the homework on Canvas.

Practice [100 points + 50 bonus]

1. [20 points] For warm-up, let us deal with the same scalar distribution in the last homework,

$$p(z) \propto \exp(-z^2)\sigma(10z + 3).$$

You will implement MCMC algorithms to sample from this distribution. To reach the burn-in stage, please run your chain for $100K$ iterations (i.e., generate $100K$ samples). Then continue to run $50K$ iterations, pick every 10-th sample to obtain your final samples. Set your initial sample to 0.

- (a) [9 points] Implement Metropolis-Hasting, with Gaussian proposal, $q(z_{n+1}|z_n) = \mathcal{N}(z_{n+1}|z_n, \tau)$. Vary τ from $\{0.01, 0.1, 0.2, 0.5, 1\}$. Run your chain. For each setting of τ , record the acceptance rate (i.e., how many candidate samples are accepted/the total number of candidate samples generated). Draw a figure, where the x-axis represents the setting of τ , and y-axis the acceptance rate. What do you observe? For each setting of τ , draw a figure, show a normalized hist-gram of the $5K$ samples you collected. Please set the number of bins to 50. If you use Python, please use matplotlib and look up the API at https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.hist.html. You can set the parameter “bins” to 50 and “density” to true. Also, please draw the ground-truth density curve (obtained via quadrature — you did that in the last homework). Now what do you observe?

$$\begin{aligned} p(z) &= \frac{1}{Z_p} \exp(-z^2)\sigma(10z + 3) = \frac{1}{Z_p} p'(z) \\ z_{n+1} &\sim \mathcal{N}(z|z_n, \tau) \\ A(z_{n+1}, z_n) &= \min \left(1, \frac{p'(z_{n+1})}{p'(z_n)} \right) \end{aligned} \tag{1}$$

For each candidate sample z_{n+1} , draw $u \sim U(0,1)$ and accept z_{n+1} as the next sample if $u < A(z_{n+1}, z_n)$.

We observe that as the variance increases, the rejection rate also increases because more samples are generated from the low probability region. If variance is small then, time taken for exploring the entire region will increase as small steps will be taken each time. Thus there is a trade-off between the step size and rejection rate of the generated samples.

Please see figure 1 for the experimental results.

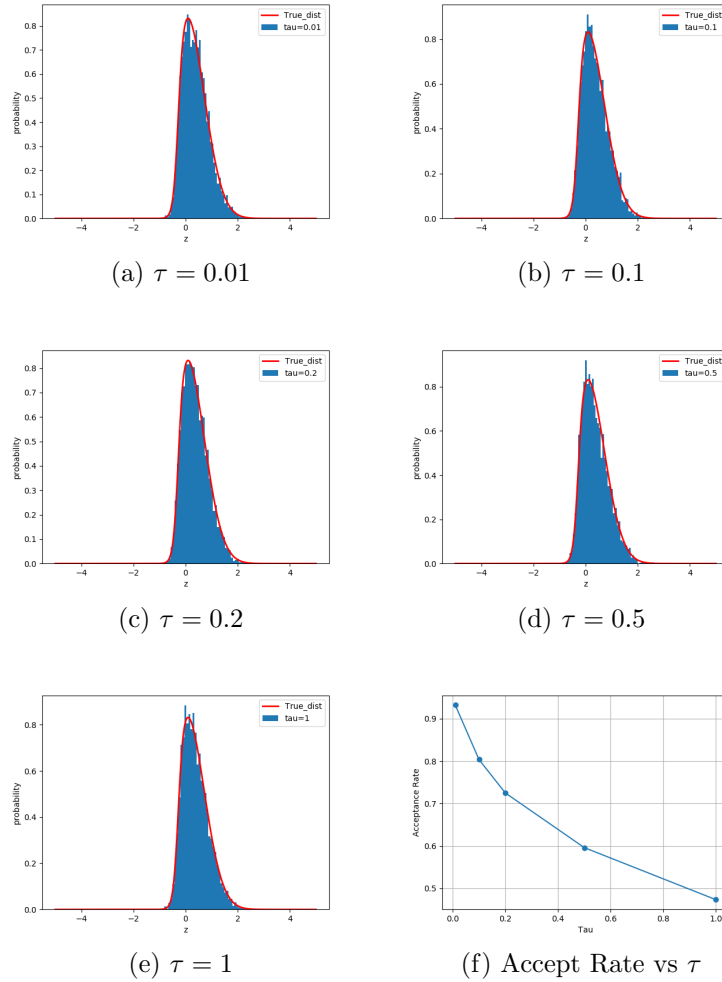


Figure 1: Metropolis-Hastings for different values of τ

- (b) [9 points] Implement Hybrid Monte-Carlo sampling with Leapfrog. Let us fix $L = 10$, and vary ϵ from $\{0.005, 0.01, 0.1, 0.2, 0.5\}$. Run your chain. Similar to the above, for each setting of ϵ , record the acceptance rate, and draw a figure showing ϵ v.s. acceptance rate. What do you observe? For each setting of ϵ , draw the normalized histogram (50 bins) of collected $5K$ samples.

what do you observe?

$$p(z) = \frac{1}{Z_p} \exp(-z^2) \sigma(10z + 3) = \frac{1}{Z_p} p'(z)$$

$$\mathcal{U}(z) = -\log(p'(z)) = z^2 - \log(\sigma(10z + 3))$$

$$\mathcal{K}(r) = \frac{1}{2} r^\top M^{-1} r$$

$$H(z, r) = \mathcal{U}(z) + \mathcal{K}(r)$$

$$p(z, r) \propto \exp(-H(z, r))$$

Observation: As we increase the value of ϵ , the acceptance rate decreases.

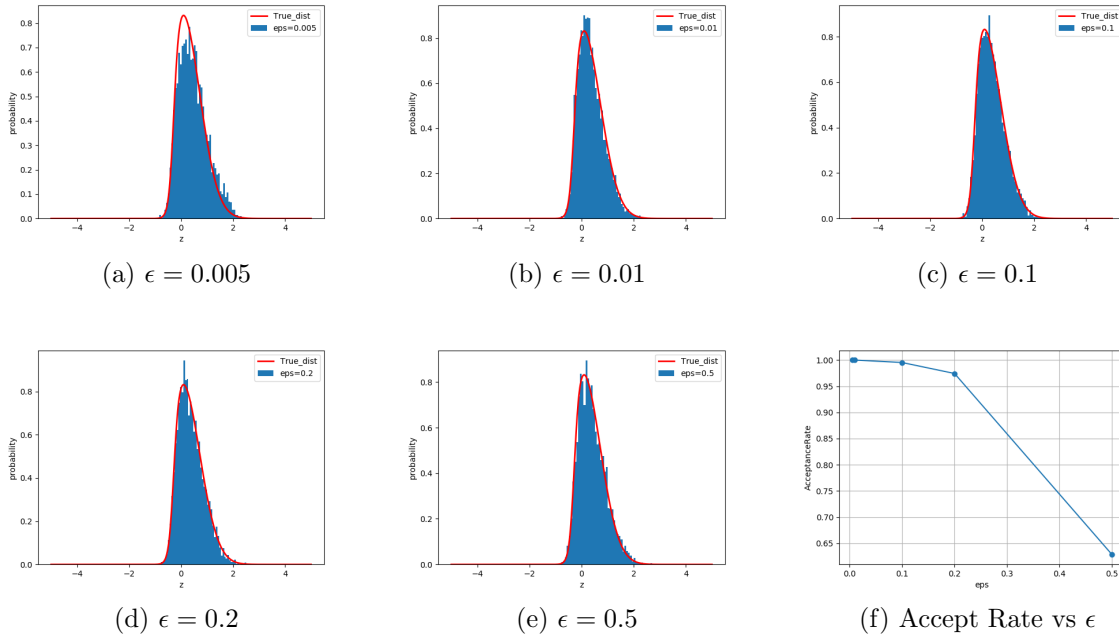


Figure 2: Hybrid Monte-Carlo sampling for different values of ϵ

Please see figure 2 for the experimental results.

- (c) [2 points] Now compare the results from the two MCMC algorithms, what do you observe and conclude?

From the above results, it can be observed that as we increase τ or ϵ , the acceptance rate decreases. However, with very small τ or ϵ , the samples generated may not approximate well the original distribution as the one generated with larger values of τ or ϵ for the same number of iterations,

2. [10 points] Let us work with a 2-dimensional Gaussian distribution,

$$p(z_1, z_2) = \mathcal{N}\left(\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \mid \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 & 2.9 \\ 2.9 & 3 \end{bmatrix}\right)$$

- (a) [1 point] Draw 500 samples from this distribution and show the scatter plot. What do you observe? As shown in figure 3, the samples drawn from the distribution $p(z_1, z_2)$ forms a

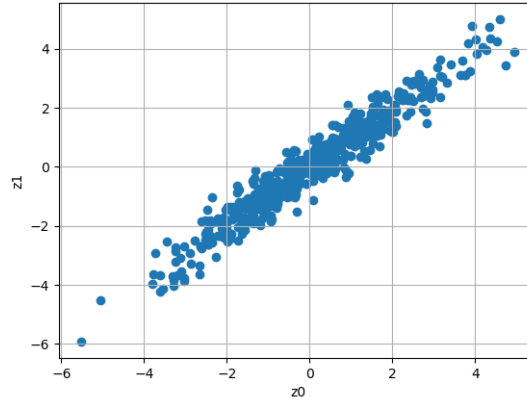


Figure 3: True Distribution

skewed ellipse showing the correlation of z_1 and z_2

- (b) [4 points] Implement Gibbs sampling to alternatively sample z_1 and z_2 . Set your initial sample to $(-4, -4)$. Run your Gibbs sampler for 100 iterations. Draw the trajectory of the samples. What do you observe?

$$p(z_1|z_2) = \mathcal{N}(z_1|\mu, \Sigma)$$

$$\mu = \mu_1 + \Sigma_{12}(\Sigma_{22})^{-1}(z_2 - \mu_2)$$

$$\Sigma = \Sigma_{11} - \Sigma_{12}(\Sigma_{22})^{-1}\Sigma_{21}$$

In Gibbs sampling, the generated samples are always accepted but it is very slow and will take

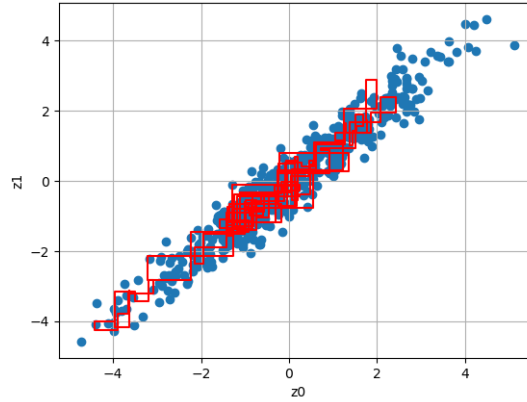


Figure 4: Gibbs Sampling

a lot of time to explore the entire distribution. The samples drawn from the Gibbs sampling are shown in figure 4

- (c) [5 points] Implement HMC with Leapfrog, set $\epsilon = 0.1$ and $L = 20$. Run your HMC for 100 iterations. Set your initial sample to $(-4, -4)$. Draw the trajectory of the samples. What do

you observe? Compare with the results of Gibbs sampling, what do you conclude?

$$\begin{aligned}
p(z) &= \mathcal{N}(z|\mu, \Sigma) \\
\mathcal{U}(z) &= -\log(p(z)) = \frac{(z - \mu)^\top \Sigma^{-1} (z - \mu)}{2} + \log(\sqrt{2\pi|\Sigma|}) \\
\mathcal{K}(r) &= \frac{1}{2} r^\top M^{-1} r \\
H(z, r) &= U(z) + K(r) \\
p(z, r) &\propto \exp(-H(z, r))
\end{aligned} \tag{2}$$

$$\begin{aligned}
p(z_1|z_2) &= \mathcal{N}(z_1|\mu, \Sigma) \\
\mu &= \mu_1 + \Sigma_{12}(\Sigma_{22})^{-1}(z_2 - \mu_2) \\
\Sigma &= \Sigma_{11} - \Sigma_{12}(\Sigma_{22})^{-1}\Sigma_{21}
\end{aligned}$$

In HMC sampling, the generated samples may not be always accepted but it is fast and may

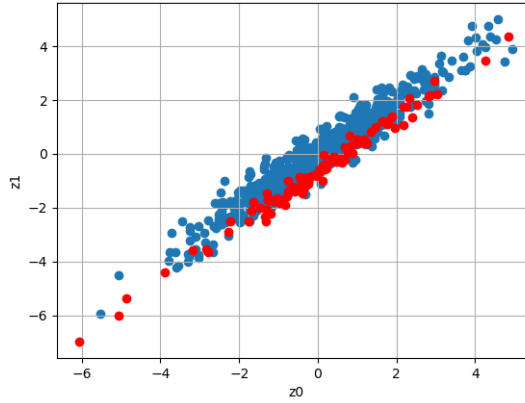


Figure 5: Hybrid Monte Carlo Sampling

explore the entire distribution in few iterations. The samples drawn from the HMC sampling are shown in figure 5

3. [70 points] Let us work on a real-world dataset we have met before. Please download the data from the folder “data/bank-note”. The features and labels are listed in the file “data-desc.txt”. The training data are stored in the file “train.csv”, consisting of 872 examples. The test data are stored in “test.csv”, and comprise of 500 examples. In both the training and testing datasets, feature values and labels are separated by commas.
 - (a) [20 points] We assign the feature weight vector \mathbf{w} a standard normal prior $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Write down the joint probability of the Bayesian logistic regression model. Now, implement HMC with Leapfrog. Set your initial sample to be $\mathbf{0}$. Run your chain for 100K iterations to reach the burn-in state; then continue to run 10K iterations, pick every 10-th sample to obtain your posterior samples. Vary ϵ from $\{0.005, 0.01, 0.02, 0.05\}$ and L from $\{10, 20, 50\}$. As we did before, we can test the predictive accuracy and predictive log-likelihood. Both can involve the posterior distribution of the weights. How? To evaluate the accuracy, for each posterior sample of \mathbf{w} in hand (1K in total), we can use it to make the predictions (1 or 0) on all the test examples and calculate the accuracy. Then we average the prediction accuracy of all the posterior samples of \mathbf{w} . In the same way, we can compute the predictive likelihood. Now you

ϵ	L	Accuracy	Likelihood
0.005	10	93.0	-0.05
0.005	20	93.0	-0.05
0.005	50	92.9	-0.06
0.01	10	92.8	-0.06
0.01	20	92.8	-0.06
0.01	50	92.8	-0.06
0.02	10	92.4	-0.06
0.02	20	92.3	-0.06
0.02	50	92.4	-0.06
0.05	10	49.1	-0.69
0.05	20	49.1	-0.69
0.05	50	49.1	-0.69

Figure 6: HMC

can sense how convenient with posterior samples — the integration over posterior distribution is turned to the average across posterior samples! List a table showing different combinations of ϵ and L and the resulting predictive accuracy, predictive log-likelihood, and acceptance rate. What do you observe and conclude?

$$\begin{aligned}
p(\mathbf{w}, \mathbf{x}, t) &= p(\mathbf{x}, t | \mathbf{w}) p(\mathbf{w}) \\
&= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}_i)^{t_i} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))^{1-t_i} p(\mathbf{w})
\end{aligned}$$

Please refer Table 6 for the results. It can be observed that the as we increase epsilon ,the error with each step increases resultilng in low accpetance rate. However, L does not affect the accuracy much.

- (b) [20 points]. We will implement Gibbs sampling for linear classification. However, Bayesian logistic regression does not allow tractable conditional posteriors. So we will use the Bayesian probit model with augmented variables. We have discussed it before in our class, if you cannot remember, please check our slides regarding generalized linear models. Again, we will assign a standard normal prior over \mathbf{w} . For each training sample n , we introduce an auxiliary variable z_n . The joint probability is given by

$$p(\mathbf{w}, \mathbf{z}, \mathcal{D}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I}) \prod_n \mathcal{N}(z_n | \mathbf{w}^T \mathbf{x}_n, 1) \mathbb{1}((2y_n - 1)z_n \geq 0) \quad (3)$$

where $\mathbf{z} = [z_1, z_2, \dots, z_N]^T$, and $\mathbb{1}(\cdot)$ is the indicator function. Note that if we marginalize out \mathbf{z} , we will recover the original Probit model,

$$\begin{aligned}
p(\mathbf{w}, \mathbf{z}, \mathcal{D}) &= \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I}) \prod_n \phi((2y_n - 1)\mathbf{w}^T \mathbf{x}_n) \\
&= \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I}) \prod_n \text{Bern}(y_n | \phi(\mathbf{w}^T \mathbf{x}_n)),
\end{aligned} \quad (4)$$

where $\phi(\cdot)$ is the CDF of the standard normal distribution, $\phi(x) = \int_{-\infty}^x \mathcal{N}(t | 0, 1) dt$. Note that Bayesian logistic regression just replaces $\phi(\cdot)$ by the Sigmoid activation function. Now implement your Gibbs sampling algorithm based on the augmented version (3). Alternatively sample \mathbf{w} and each z_n . Note that the conditional posterior of each z_n will be a truncated Gaussian. You can use `scipy.stats.truncnorm` to generate samples (or implement by yourself). Before coding,

please list your derivation of the conditional posteriors. Run your chain for $100K$ iterations to reach the burn-in stage; then continue to run $10K$ iterations, pick every 10-th sample to obtain your final posterior samples. Now compute and report the predictive accuracy and log-likelihood with your posterior samples. Note that your predictive log-likelihood should be based on the original model (4), rather than the augmented one. How does the performance compare with HMC on Bayesian logistic regression model?

$$p(W|X, Y) = \prod_i^N \phi(W^\top x_i, 1)^{y_i} (1 - \phi(W^\top x_i, 1))^{1-y_i} \mathcal{N}(W|0, I)$$

$$\text{where } \phi(W^\top x_i, 1) = \int_{-\infty}^{W^\top x_i} \mathcal{N}(t|0, I) dt$$

Augmenting the model with random variable $z_i = w^\top x_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, 1)$. Since, $w \sim \mathcal{N}(0, 1)$, $w^\top x_i$ also follows a normal distribution.

$$z_i = w^\top x_i + \epsilon_i > 0$$

$$\epsilon_i > -w^\top x_i$$

$$p(z_i > 0) = p(\epsilon_i > -w^\top x_i) = 1 - (1 - p(\epsilon_i > w^\top x_i)) = p(\epsilon_i > w^\top x_i)$$

Thus, it reduces to:

$$p(y_i = 1|x_i, w) = p(z_i > 0|x_i, w) = \phi(W^\top x_i, 1)$$

$$p(w|X, Z, Y) = p(w|X, Z) = \mathcal{N}((I + X^\top X)^{-1} X^\top Z, (X^\top X + I)^{-1})$$

$$p(z_i|w, x_i, y_i = 0) = TN(w^\top x_i, 1, -\infty, 0)$$

$$p(z_i|w, x_i, y_i = 1) = TN(w^\top x_i, 1, 0, \infty)$$

Reference : <https://people.eecs.berkeley.edu/~jordan/courses/260-spring10/lectures/lecture19.pdf>

- Predictive Accuracy: 0.986
 - Predictive Likelihood: -0.033
 - Acceptance Rate: 1
- (c) [28 points] Finally, we will implement a Bayesian neural network (NN). We will use two intermediate layers, and each layer has the same number of nodes. Vary the number of nodes from [10, 20, 50]. Vary the activation function from `{tanh, RELU}`. We will use a factorized Gaussian posterior for the NN weights (check out the slides). We will assign a standard normal prior over each weight. The output of the NN will be thrown into a Sigmoid activation function, with which we obtain a Bernoulli likelihood. Please use TensorFlow to implement the **Bayes** by BP algorithm based on reparameterization trick plus stochastic optimization. If you prefer PyTorch, that is okay. But you are NOT allowed to implement BP by yourself — it is a waste of time. Note that you do not need to sample mini-batches in this problem — because the training set is small. Please use Adam algorithm for stochastic optimization. You can tune the base learning rate from `{1e-3, 0.5e-3, 1e-4, 1e-5}` to find the best result for each layer-width and activation function setting. Initialize posterior mean and variance of each weight to be 0 and 1, respectively. Run the Adam algorithm for 1,000 iterations. For each layer width and activation function setting, calculate the predictive log-likelihood and accuracy. How? Use Monte-Carlo approximation. Check out the slides. Use the variational posterior to generate 100 samples for the weight vector; With each sample, you can calculate the prediction accuracy and log-likelihood on the test dataset; finally, you report the average results. To check the behaviour of your BNN, let us pick up one setting: the number of node is 20 and the activation function is `tanh`. For each iteration, please use the current posterior mean of the weights to compute the average log-likelihood on the training set and test set respectively. Draw two plots, one plot

showing how the training log-likelihood varies along with the number of iterations, the other showing how the test log-likelihood varies along with the number of iterations. What do you observe and conclude?

Please see figure 7

From the results, it can be observed that as we increase the number of nodes in the hidden layer, the accuracy of the model increases. Also in Bayesian Neural Network, there is a sudden improvement in the performance of the model once it has been trained for sufficient epochs. For all the experiments, the learning-rate is set to 0.001

The help is taken from the link : <https://www.nitarshan.com/bayes-by-backprop/>

- (d) [2 points] Compare the results from Bayesian linear classification and NN models. What do you observe and conclude? **Bayesian NN works better than Bayesian linear classification because of the non-linearity added at each layer. However, they may be difficult to train as they need a set of hyperparameters like learning rate to be tuned.**
4. [50 points][**Bonus**] If you feel addictive to neural networks, this is a good chance to do more state-of-the-art work. Please come to the website <http://yann.lecun.com/exdb/mnist/> and download the MNIST dataset provided by LeCun et. al. It is a famous benchmark dataset. Your goal is to classify hand-written digits. The dataset includes 60,000 training and 10,000 testing pixel images of size 28×28 . Each image is labelled with its ground-truth number (0-9 inclusive). Please preprocess the data by dividing each pixel values by 126. You will implement two versions of Bayesian neural networks. To select the hyper-parameters, please from the training data randomly select 50,000 digits for training and use the remaining 10,000 digits for validation. After the best hyper-parameters are identified, you train on the whole 60,000 digits again. Your NN output will be 10 dimensional, which are used to construct the categorical likelihood (i.e., softmax) — check out the slides if you are unclear. Please use Adam for stochastic optimization. We consider to tune the (base) learning rate from $\{1e-3, 1e-4, 1e-5\}$. We will use two hidden layers, and each layer has the same number of nodes.
- (a) [15 points] First, we consider to place an independent standard normal prior over each weight. We use factorized Gaussian posterior (check the slides). Implement your Bayes by BP algorithm. Please run your algorithm for 1,000 epochs and report the predictive accuracy. The way to compute the predictive accuracy is the same as in Problem 3. Use your variational posterior of the weights to sample 10 sets of weights, each weight set are used to predict the labels of the test samples (you choose the label with the largest probability in your model as the prediction) and calculate the accuracy; finally you average the 10 accuracy values. Vary the number of nodes in each hidden layer from $\{400, 800, 1200\}$. Try both the **tanh** and **RELU** activation functions. Report the prediction accuracy for each combination.
- (b) [20 points] Then we consider to assign a spike and slab prior over each weight,

$$p(w_i) = \pi \mathcal{N}(w_i | 0, \sigma_1^2) + (1 - \pi) \mathcal{N}(w_i | 0, \sigma_2^2)$$

where we tune $-\log(\sigma_1) \in \{0, 1, 2\}$ and $-\log(\sigma_2) \in \{6, 7, 8\}$, $\pi \in \{\frac{1}{4}, \frac{1}{2}, \frac{3}{4}\}$. We still use the factorized Gaussian posterior. Implement your Bayes by BP algorithm. Run your algorithm for 1,000 epochs. Report the predictive accuracy when the number of nodes in each hidden layer is in $\{400, 800, 1200\}$, and the activation function is **tanh** or **RELU**. Do you observe many weights have very small posterior variances and means close to 0? Why did this happen?

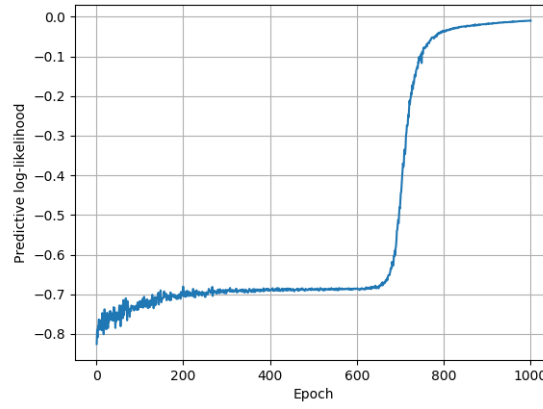
- (c) [10 points] Implement the vanilla NN training with SGD — namely, we only perform MAP estimation. Run 1,000 epochs. Report the prediction accuracy for the same settings as above. Are your results consistent with Table 1 in the paper <https://arxiv.org/pdf/1505.05424.pdf>? Compared with the two version of BNNs, what do you observe and conclude?

Nodes	10	20	50
Predictive Accuracy	96.0	98.0	99.0
Predictive Log-likelihood	-0.214	-0.107	-0.074

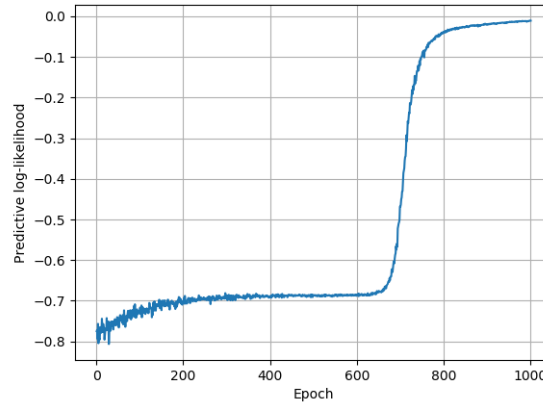
(a) Activation Function : ReLU

Number of Nodes	10	20	50
Predictive Accuracy	87.0	99.0	99.0
Predictive Log-likelihood	-0.342	-0.012	-0.004

(b) Activation Function : Tanh



(c) Train Log-Likelihood



(d) Test Log-Likelihood

Figure 7: Bayesian Neural Network

- (d) [5 points] Fix the number of nodes per layer to 1,200 and use **ReLU** activation function. For each algorithm, use the validation dataset to choose the best hyper-parameter(s) (e.g., learning rate). With the best hyper-parameter(s), re-run your algorithms on the whole training set, and draw the learning curve — how does the test accuracy vary along with the number of epochs. Is it consistent with Fig. 2 in the paper <https://arxiv.org/pdf/1505.05424.pdf>? What do

you conclude and observe?