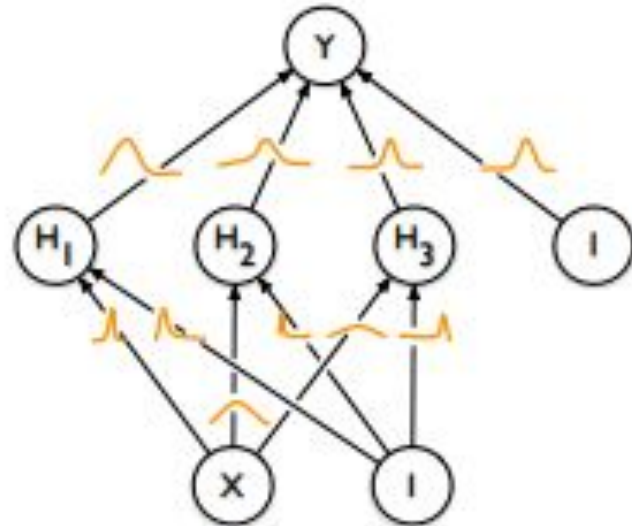
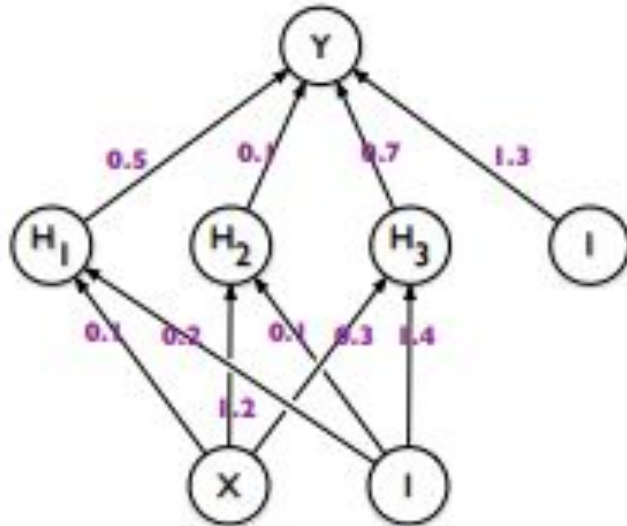


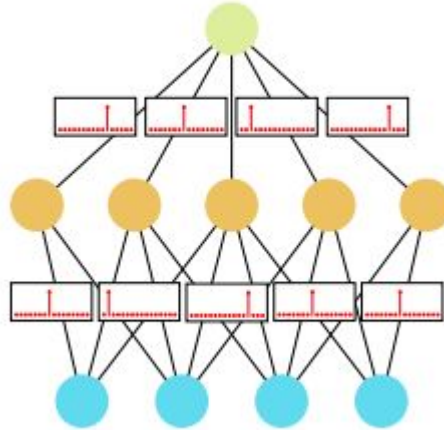
Bayesian Neural Network

Present by Linjun Huang

Neural Networks

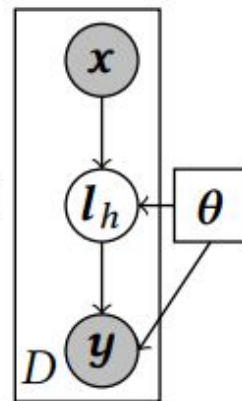
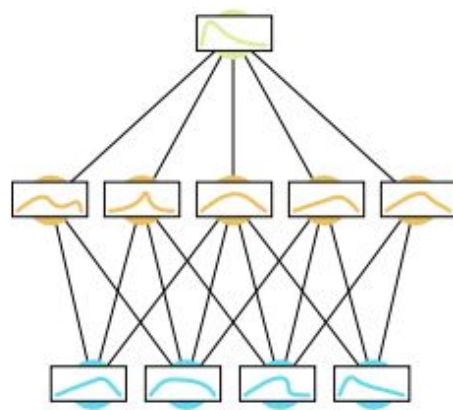
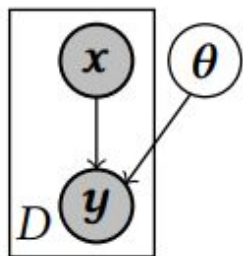
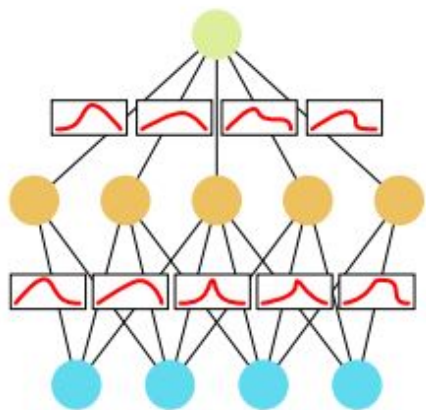


Regular Neural Network



$$\begin{aligned}l_0 &= \mathbf{x}, \\l_i &= \text{nl}_i(\mathbf{W}_i l_{i-1} + \mathbf{b}_i) \quad \forall i \in [1, n], \\ \mathbf{y} &= l_n.\end{aligned}$$

Bayesian Neural Network

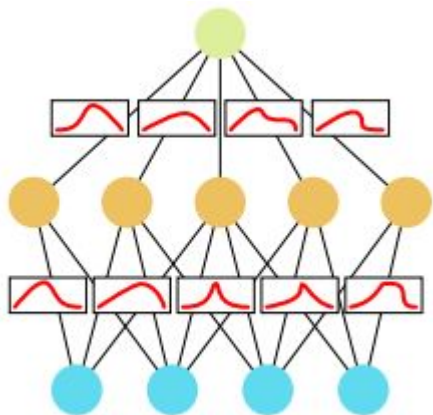


$$\theta \sim p(\theta),$$

$$\mathbf{y} = \text{NN}_{\theta}(\mathbf{x}) + \epsilon$$

$$p(\theta|D) = \frac{p(D_{\mathbf{y}}|D_{\mathbf{x}}, \theta)p(\theta)}{\int_{\theta} p(D_{\mathbf{y}}|D_{\mathbf{x}}, \theta')p(\theta')d\theta'} \propto p(D_{\mathbf{y}}|D_{\mathbf{x}}, \theta)p(\theta).$$

Define BNN from a PGM



$$\theta \sim p(\theta) = \mathcal{N}(\mu, \Sigma),$$

$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\text{NN}_\theta(\mathbf{x}), \Sigma).$$

$$p(D_{\mathbf{y}}|D_{\mathbf{x}}, \theta) = \prod_{(\mathbf{x}, \mathbf{y}) \in D} p(\mathbf{y}|\mathbf{x}, \theta).$$

$$l_0 = \mathbf{x},$$

$$l_i \sim p(l_i|l_{i-1}) = \text{nl}_i(\mathcal{N}(\mathbf{W}_i l_{i-1} + \mathbf{b}_i, \Sigma)) \quad \forall i \in [1, n],$$

$$\mathbf{y} = l_n.$$

$$p(D_{\mathbf{y}}, l_{[1, n-1]}|D_{\mathbf{x}}) = \prod_{(l_0, l_n) \in D} \left(\prod_{i=1}^n p(l_i|l_{i-1}) \right).$$

$$\mathbf{W} \sim \mathcal{N}(\mu_{\mathbf{W}}, \Sigma_{\mathbf{W}}),$$

$$\mathbf{b} \sim \mathcal{N}(\mu_{\mathbf{b}}, \Sigma_{\mathbf{b}}),$$

$$l = \text{nl}(\mathbf{W}l_{-1} + \mathbf{b})$$

Optimization on Neural Networks

Does not place a prior upon weights \mathbf{w}

$$\begin{aligned}\mathbf{w}^{\text{MLE}} &= \arg \max_{\mathbf{w}} \log P(\mathcal{D}|\mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \sum_i \log P(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}).\end{aligned}$$

Place a prior upon weights \mathbf{w}

$$\begin{aligned}\mathbf{w}^{\text{MAP}} &= \arg \max_{\mathbf{w}} \log P(\mathbf{w}|\mathcal{D}) \\ &= \arg \max_{\mathbf{w}} \log P(\mathcal{D}|\mathbf{w}) + \log P(\mathbf{w}).\end{aligned}$$

Bayesian Estimation

Posterior Distribution $P(\mathbf{w}|\mathcal{D}) = \frac{P(\mathbf{w}, \mathcal{D})}{P(\mathcal{D})} = \frac{P(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})}$

Predictive Distribution $P(\hat{\mathbf{y}}|\hat{\mathbf{x}}) = \mathbb{E}_{P(\mathbf{w}|\mathcal{D})}[P(\hat{\mathbf{y}}|\hat{\mathbf{x}}, \mathbf{w})]$

Optimization via Kullback-Leibler divergence

Use a distribution $q(\mathbf{w}|\theta)$ to approximate the posterior distribution $P(\mathbf{w}|\mathcal{D})$

$$\begin{aligned}\theta^* &= \arg \min_{\theta} D_{\text{KL}}[q(\mathbf{w}|\theta) || P(\mathbf{w}|\mathcal{D})] \\ &= \arg \min_{\theta} \int q(\mathbf{w}|\theta) \log \frac{q(\mathbf{w}|\theta)}{P(\mathbf{w})P(\mathcal{D}|\mathbf{w})} d\mathbf{w} \\ &= \arg \min_{\theta} D_{\text{KL}}[q(\mathbf{w}|\theta) || P(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}|\theta)}[\log P(\mathcal{D}|\mathbf{w})]\end{aligned}$$

Denote objective function: $\mathcal{F}(\mathcal{D}, \theta) = D_{\text{KL}}[q(\mathbf{w}|\theta) || P(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}|\theta)}[\log P(\mathcal{D}|\mathbf{w})]$

Prior Distribution

Scale mixture gaussian prior:

$$P(\mathbf{w}) = \prod_j \pi \mathcal{N}(\mathbf{w}_j | 0, \sigma_1^2) + (1 - \pi) \mathcal{N}(\mathbf{w}_j | 0, \sigma_2^2)$$

gaussian prior:

$$P(\mathbf{w}) = \prod_j \mathcal{N}(\mathbf{w}_j | 0, \sigma^2)$$

Reparametrize

Reparameterize process: $z \sim \mathcal{N}(\mu, \sigma^2)$

$$z = \sigma\epsilon + \mu$$

$$\epsilon \sim \mathcal{N}(0, 1)$$

The distribution of ϵ satisfies $q(\epsilon)d\epsilon = q(\mathbf{w}|\theta)d\mathbf{w}$

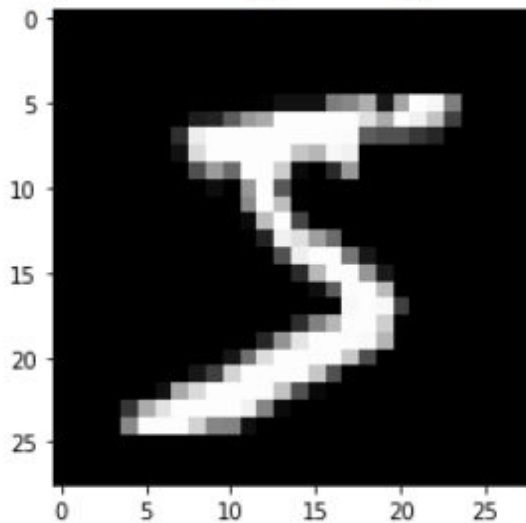
Unbiased Monte Carlo gradients

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbb{E}_{q(\mathbf{w}|\theta)}[f(\mathbf{w}, \theta)] &= \frac{\partial}{\partial \theta} \int f(\mathbf{w}, \theta) q(\mathbf{w}|\theta) d\mathbf{w} \\ &= \frac{\partial}{\partial \theta} \int f(\mathbf{w}, \theta) q(\epsilon) d\epsilon \\ &= \mathbb{E}_{q(\epsilon)} \left[\frac{\partial f(\mathbf{w}, \theta)}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial \theta} + \frac{\partial f(\mathbf{w}, \theta)}{\partial \theta} \right]\end{aligned}$$

$$\mathcal{F}(\mathcal{D}, \theta) \approx \sum_{i=1}^n \log q(\mathbf{w}^{(i)}|\theta) - \log P(\mathbf{w}^{(i)}) - \log P(\mathcal{D}|\mathbf{w}^{(i)})$$

Numerical Results

```
Loading MNIST dataset  
X_train.shape = (60000, 28, 28, 1)  
y_train.shape = (60000, 10)  
X_test.shape = (10000, 28, 28, 1)  
y_test.shape = (10000, 10)  
<matplotlib.image.AxesImage at 0x7f9d0d465cc0>
```

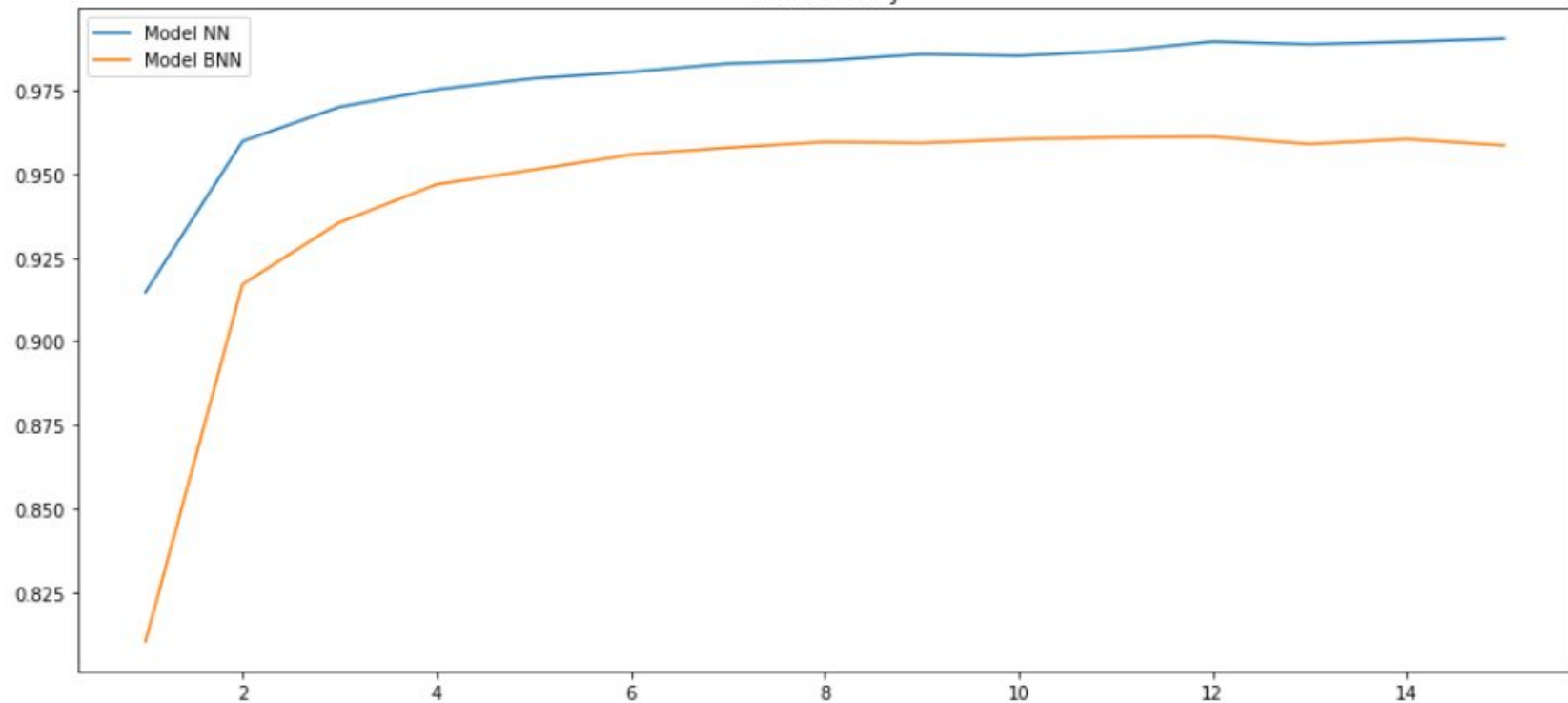


Numerical Results

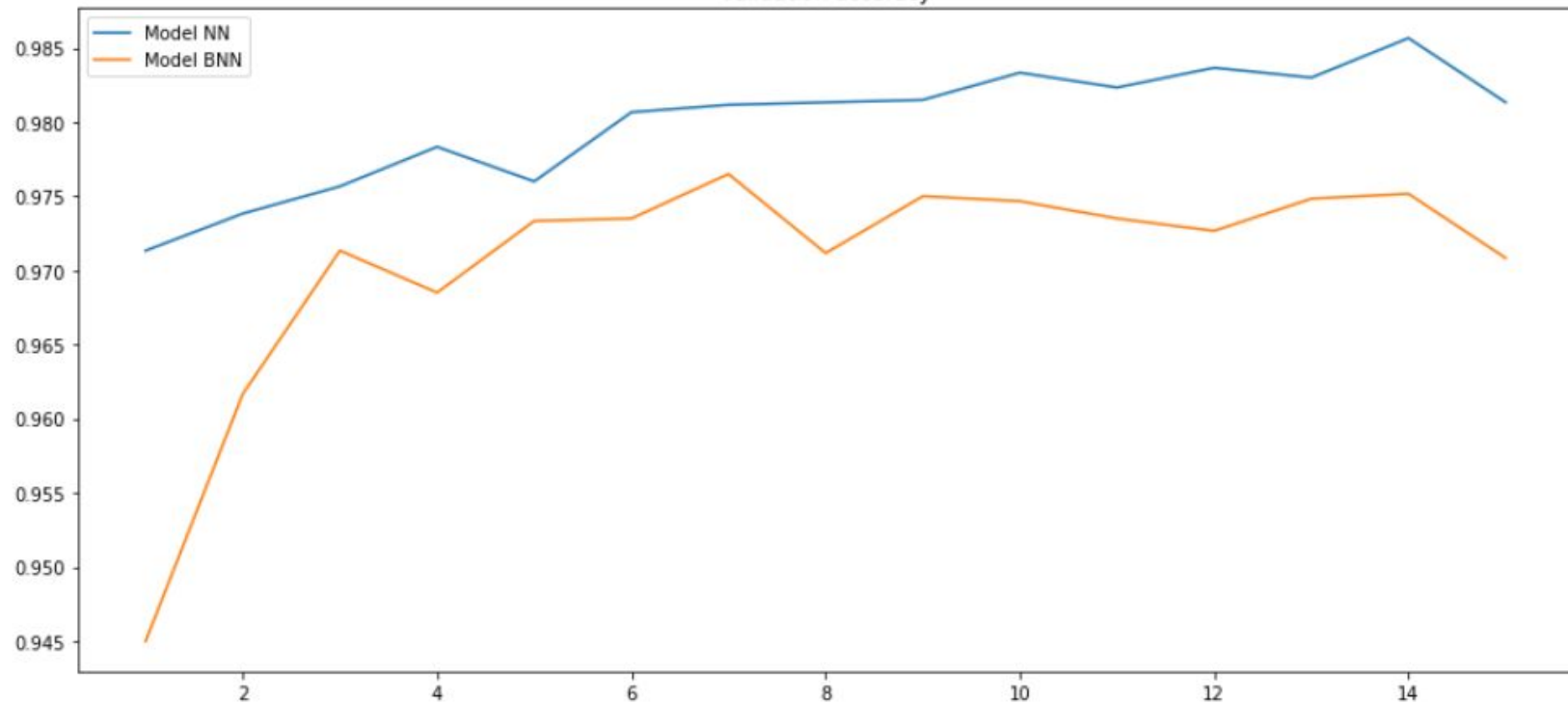
Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 28, 28, 1)]	0
flatten (Flatten)	(None, 784)	0
dense_flipout (DenseFlipout)	(None, 512)	803329
batch_normalization (Batch Normalization)	(None, 512)	2048
dropout (Dropout)	(None, 512)	0
dense_flipout_1 (DenseFlipout)	(None, 256)	262401
batch_normalization_1 (Batch Normalization)	(None, 256)	1024
dropout_1 (Dropout)	(None, 256)	0
dense_flipout_2 (DenseFlipout)	(None, 128)	65665
dropout_2 (Dropout)	(None, 128)	0
dense_flipout_3 (DenseFlipout)	(None, 10)	2571
Total params: 1,137,038		
Trainable params: 1,135,498		
Non-trainable params: 1,540		

Layer (type)	Output Shape	Param #
input_2 (InputLayer)	[(None, 28, 28, 1)]	0
flatten_1 (Flatten)	(None, 784)	0
dense (Dense)	(None, 512)	401920
batch_normalization_2 (Batch Normalization)	(None, 512)	2048
dropout_3 (Dropout)	(None, 512)	0
dense_1 (Dense)	(None, 256)	131328
batch_normalization_3 (Batch Normalization)	(None, 256)	1024
dropout_4 (Dropout)	(None, 256)	0
dense_2 (Dense)	(None, 128)	32896
dropout_5 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 10)	1290
Total params: 570,506		
Trainable params: 568,970		
Non-trainable params: 1,536		

train accuracy



validation accuracy



Test accuracy

```
2 preds_nn = nn_model.evaluate(X_test, y_test)
3 print ("Loss NN = {}".format(preds_nn[0]))
4 print ("Test Accuracy NN = {}".format(preds_nn[1]))
```

```
313/313 [=====] - 1s 2ms/step - loss:
Loss NN = {}0.07821856439113617
Test Accuracy NN = 0.9800000190734863
```

```
1 preds_bnn = bnn_model.evaluate(X_test, y_test)
2 print ("Loss BNN = {}".format(preds_bnn[0]))
3 print ("Test Accuracy BNN = {}".format(preds_bnn[1]))
```

```
313/313 [=====] - 1s 4ms/step - loss:
Loss BNN = {}8.128716468811035
Test Accuracy BNN = 0.9689000248908997
```


Quantify the uncertainty in predictions for bnn

```
1 n_mc_run = 100
2 med_prob_thres = 0.2
3
4
5 y_pred_logits_list = [bnn_model.predict(X_test) for _ in range(n_mc_run)] # a list of predicted logits
6 y_pred_prob_all = np.concatenate([softmax(y, axis=-1)[:, :, np.newaxis] for y in y_pred_logits_list], axis=-1)
7 y_pred = [[int(np.median(y) >= med_prob_thres) for y in y_pred_prob] for y_pred_prob in y_pred_prob_all]
8 y_pred = np.array(y_pred)
9
10 idx_valid = [any(y) for y in y_pred]
11 print('Number of recognizable samples:', sum(idx_valid))
12
13 idx_invalid = [not any(y) for y in y_pred]
14 print('Unrecognizable samples:', np.where(idx_invalid)[0])
15
16 print('Test accuracy on MNIST (recognizable samples):',
17       sum(np.equal(np.argmax(y_test[idx_valid], axis=-1), np.argmax(y_pred[idx_valid], axis=-1))) / len(y_test[idx_valid]))
18
19 print('Test accuracy on MNIST (unrecognizable samples):',
20       sum(np.equal(np.argmax(y_test[idx_invalid], axis=-1), np.argmax(y_pred[idx_invalid], axis=-1))) / len(y_test[idx_invalid]))
```

Number of recognizable samples: 9996

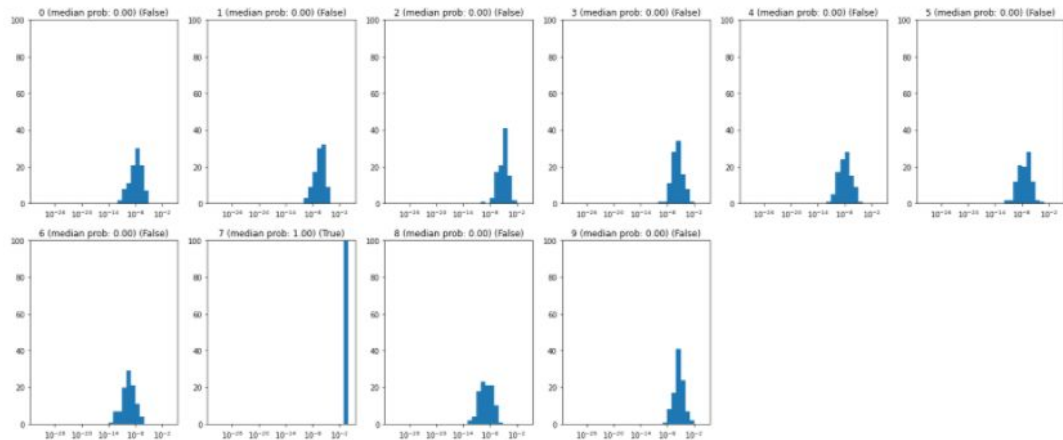
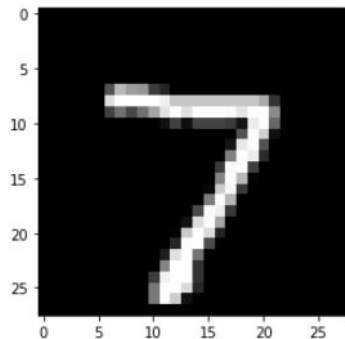
Unrecognizable samples: [1941 2266 6625 9634]

Test accuracy on MNIST (recognizable samples): 0.9746898759503801

Test accuracy on MNIST (unrecognizable samples): 0.25

A recognizable example

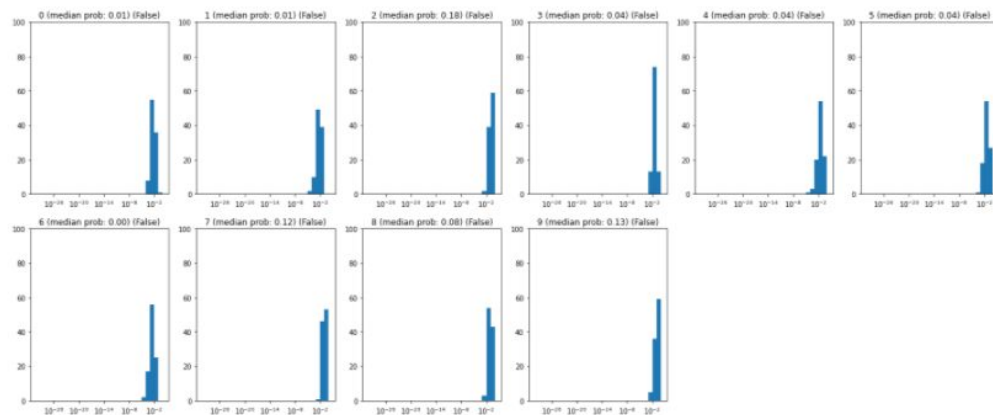
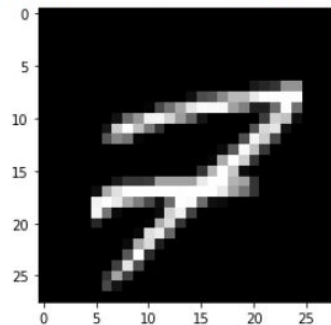
True label of the test sample 0: 7



Predicted label of the test sample 0: 7

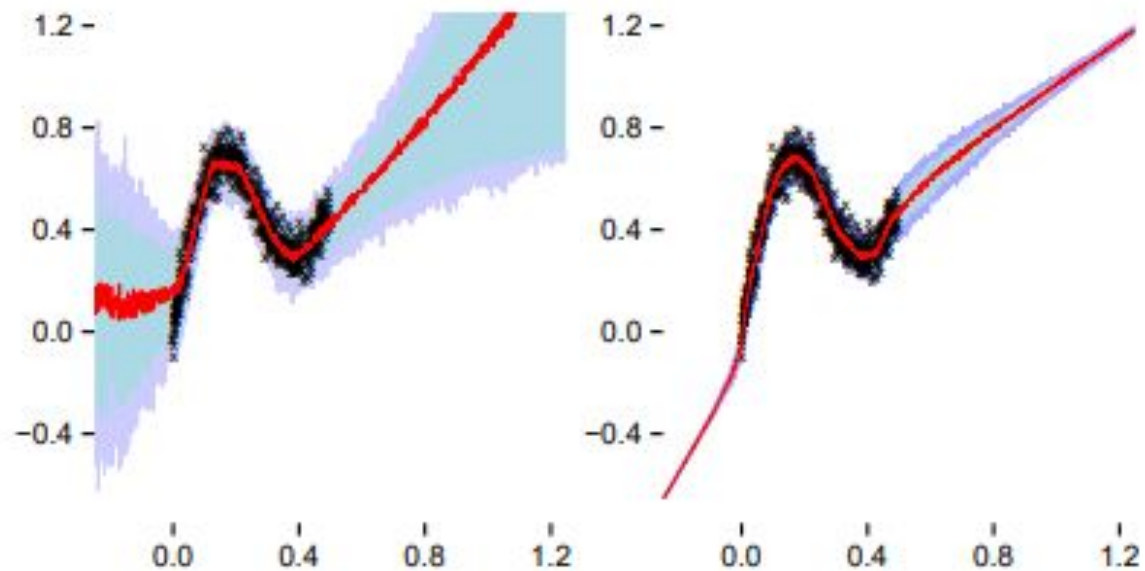
Unrecognizable example

True label of the test sample 1941: 7



I don't know!

Example in paper



Reference

[1] Laurent Valentin Jospin, etc, Hands-on Bayesian Neural Network - a Tutorial for Deep Learning Users.

[2] Charles Blundell, etc, Weight Uncertainty in Neural Networks.