

Introduction to Machine Learning Project Phase 2 Report

Ali Sadeghian¹, Amirreza Tanevardi²

¹ Sharif University of Technology, Electrical Engineering, 400101464, ali.sadeghian@ee.sharif.edu

² Sharif University of Technology, Electrical Engineering, 400100898, amirrezatanevardi@gmail.com

Keywords:

Restricted Boltzmann Machines (RBMs) are undirected generative models that use a layer of hidden variables to model a distribution over visible variables. Typically, they are employed to model input distributions in classification tasks, but they can also capture the joint distribution of inputs and target classes. This characteristic makes them particularly useful in deep learning architectures, such as the last layer of a Deep Belief Network.

An RBM with n hidden units is a parametric model describing the joint distribution between hidden variables $h = (h_1, \dots, h_n)$ and observed variables (x, y) , where $x = (x_1, \dots, x_d)$ represents input features and y denotes the target class. The joint probability distribution takes the form:

$$p(y, x, h) \propto e^{-E(y, x, h)}$$

where the energy function is given by:

$$E(y, x, h) = -h^T W x - b^T x - c^T h - d^T \tilde{y} - h^T U \tilde{y}$$

with parameters $\Theta = (W, b, c, d, U)$ and one-hot encoded class representation $\tilde{y} = (1_{y=i})_{i=1}^C$. The model is applicable to binary, integer-valued, and continuous-valued inputs, making it flexible across various applications.

Theoretical Question 1

Explain the ElGamal encryption algorithm completely and write its mathematical relationships.

Theoretical Question 1 answer

The **ElGamal encryption algorithm** is a public-key cryptosystem based on the *Discrete Logarithm Problem (DLP)*. It consists of three main steps: **key generation**, **encryption**, and **decryption**.

Key Generation

1. Choose a large prime number p .
2. Select a primitive root g of p .
3. Choose a private key x such that $1 \leq x \leq p - 2$.

4. Compute the public key y as:

$$y = g^x \mod p$$

5. The **public key** is (p, g, y) , and the **private key** is x .

Encryption

To encrypt a message M :

1. Represent M as an integer such that $0 \leq M < p$.

2. Choose a random integer k such that $1 \leq k \leq p - 2$.

3. Compute the first part of the ciphertext:

$$C_1 = g^k \mod p$$

4. Compute the second part of the ciphertext:

$$C_2 = M \cdot y^k \mod p$$

5. The ciphertext is (C_1, C_2) .

Theoretical Question 2

Explain partial decryption in detail.

Theoretical Question 2 answer

Partial decryption in ElGamal encryption refers to computing an intermediate decryption step without fully recovering the plaintext message. This step involves calculating the shared secret S .

Computing the Shared Secret

Given the ciphertext (C_1, C_2) , the receiver computes:

$$S = C_1^x \mod p$$

Since $C_1 = g^k$, we get:

$$S = (g^k)^x \mod p = g^{kx} \mod p$$

Importance of Partial Decryption

- The shared secret S is a crucial value used in full decryption.
- Without knowing x , an attacker cannot compute S .
- Partial decryption stops at this step, leaving the final message recovery incomplete.

Theoretical Question 3

Provide an algorithm similar to the method seen above for computing the multiplication of two numbers. Assume that Party A knows the number M , and Party B knows the number N , and they want to compute MN .

Theoretical Question 3 answer

The following steps outline a secure multiplication protocol using ElGamal encryption:

Step 1: Initialization

- Party A generates a random number R .
- Party A computes:

$$m_i = M \cdot i - R$$

where i represents the possible input values of Party B .

- Each m_i is treated as plaintext and encrypted using the ElGamal encryption scheme:

$$E(m_i, r_i)$$

where r_i is a new random number for each encryption.

- Party A sends $E(m_i, r_i)$ to Party B in increasing order of i .

Step 2: Party B Picks and Rerandomizes

- Party B selects the ciphertext corresponding to its value N , i.e., $E(m_N, r_N)$.
- Party B rerandomizes the encryption by adding a new random value s :

$$E(m_N, r') = E(m_N, r_N + s)$$

- Party B sends $E(m_N, r')$ back to Party A .

Step 3: Partial Decryption by Party A

- Party A partially decrypts $E(m_N, r')$.
- The partially decrypted message is then sent to Party B .

Step 4: Final Decryption by Party B

- Party B fully decrypts the received message.
- The decrypted message is:

$$m_N = M \cdot N - R$$

- Since R is only known to Party A , and m_N is only known to Party B , the final result is computed as:

$$MN = m_N + R$$

Simulation Question 1

Using the MNIST dataset and the secure RBM as a feature extractor, perform the classification task and compare its performance against the regular RBM.

Simulation Question 1 answer

The code is presented part by part and the functionality of each section is explained. First we inspect the Regular RBM. here is the models code:

```

1 class RBM:
2     def __init__(self, input=None, n_visible=784, n_hidden=128, W=None, hbias=None, vbias=None, numpy_rng=None):
3         self.n_visible = n_visible
4         self.n_hidden = n_hidden
5
6         if numpy_rng is None:
7             numpy_rng = np.random.RandomState(1234)
8
9         if W is None:
10            a = 1. / n_visible
11            W = numpy_rng.uniform(-a, a, size=(n_visible, n_hidden))
12
13        if hbias is None:
14            hbias = np.zeros(n_hidden)
15
16        if vbias is None:
17            vbias = np.zeros(n_visible)
18
19        self.W = W
20        self.hbias = hbias
21        self.vbias = vbias
22        self.numpy_rng = numpy_rng
23
24    def sigmoid(self, x):
25        return 1 / (1 + np.exp(-x))
26
27    def sample_h_given_v(self, v):
28        h_mean = self.sigmoid(np.dot(v, self.W) + self.hbias)
29        h_sample = self.numpy_rng.binomial(size=h_mean.shape, n=1, p=h_mean)
30        return h_mean, h_sample
31
32    def sample_v_given_h(self, h):
33        v_mean = self.sigmoid(np.dot(h, self.W.T) + self.vbias)
34        v_sample = self.numpy_rng.binomial(size=v_mean.shape, n=1, p=v_mean)
35        return v_mean, v_sample
36
37    def gibbs_hvh(self, h0_sample):
38        v_mean, v_sample = self.sample_v_given_h(h0_sample)
39        h_mean, h_sample = self.sample_h_given_v(v_sample)
40        return v_mean, v_sample, h_mean, h_sample
41
42    def contrastive_divergence(self, data, k=1, lr=0.1):
43        ph_mean, ph_sample = self.sample_h_given_v(data)
44
45        chain_start = ph_sample
46        for step in range(k):
47            if step == 0:
48                nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh(chain_start)
49            else:

```

```

50         nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh(nh_sample)
51
52         self.W += lr * (np.dot(data.T, ph_sample) - np.dot(nv_sample.T, nh_mean)) / len(data)
53         self.vbias += lr * np.mean(data - nv_sample, axis=0)
54         self.hbias += lr * np.mean(ph_sample - nh_mean, axis=0)
55
56     def reconstruct(self, v):
57         h = self.sigmoid(np.dot(v, self.W) + self.hbias)
58         reconstructed_v = self.sigmoid(np.dot(h, self.W.T) + self.vbias)
59         return reconstructed_v
60
61     def save_model(self, file_path):
62         with open(file_path, 'wb') as f:
63             pickle.dump(self, f)
64
65     @staticmethod
66     def load_model(file_path):
67         with open(file_path, 'rb') as f:
68             return pickle.load(f)

```

This is the exact same model we developed in the first phase of the project. The training function is as follows:

```

1
2def train_rbm(train_loader, test_loader, n_visible, n_hidden, learning_rate, epochs, k):
3    rbm = RBM(n_visible=n_visible, n_hidden=n_hidden)
4
5    for epoch in range(epochs):
6        for batch, _ in train_loader:
7            batch = batch.view(-1, n_visible).numpy()
8            rbm.contrastive_divergence(batch, k=k, lr=learning_rate)
9
10           test_data = next(iter(test_loader))[0].view(-1, n_visible).numpy()
11           reconstruction = rbm.reconstruct(test_data)
12           loss = np.mean((test_data - reconstruction) ** 2)
13           print(f"Epoch {epoch + 1}, k={k}: Reconstruction Loss = {loss:.4f}")
14
15    return rbm

```

The explanation for these pieces of code was given before. The piece of code we wrote was this function that was used to extract features from the MNIST dataset using the RBM model. The features are considered to be the mean of the hidden layer after the first step of Gibbs sampling process.

```

1
2    from sklearn.linear_model import LogisticRegression
3from sklearn.metrics import accuracy_score
4
5# Extract features using RBM
6def extract_features_rbm(rbm, loader):
7    features = []
8    labels = []
9
10   with torch.no_grad():
11       for batch, lbls in loader:
12           batch = batch.view(-1, rbm.n_visible)
13           prob_h, _ = rbm.sample_h_given_v(batch)
14           features.append(prob_h)
15           labels.append(lbls)
16
17   return features, labels
18
19transform = transforms.Compose([transforms.ToTensor(), transforms.Lambda(lambda x: x.view(-1))])
20train_dataset = datasets.MNIST(root="./data", train=True, transform=transform, download=True)
21test_dataset = datasets.MNIST(root="./data", train=False, transform=transform, download=True)
22
23train_loader = DataLoader(train_dataset, batch_size=1, shuffle=True)

```

```

24 test_loader = DataLoader(test_dataset, batch_size=1, shuffle=False)
25 x_train_rbm, y_train_rbm = extract_features_rbm(regular_rbm, train_loader)
26 x_test_rbm, y_test_rbm = extract_features_rbm(regular_rbm, test_loader)
27
28 # Train a logistic regression model on RBM features
29 clf_rbm = LogisticRegression(max_iter=1000)
30 clf_rbm.fit(np.array(x_train_rbm).squeeze(1), np.array(y_train_rbm).squeeze(1))
31
32 # Evaluate the model
33 y_pred_rbm = clf_rbm.predict(np.array(x_test_rbm).squeeze(1))
34 accuracy_rbm = accuracy_score(np.array(y_test_rbm).squeeze(1), y_pred_rbm)
35 print(f"RBM Feature Classification Accuracy: {accuracy_rbm:.4f}")
36

```

As you can see, we used a simple Logistic Regression model as our classifier. Here is the result for this classification task:

Model	Precision	Recall	F1-score	accuracy
Regular RBM (CD=1)	0.93	0.93	0.93	0.94

Table 1: Performance of regular RBM

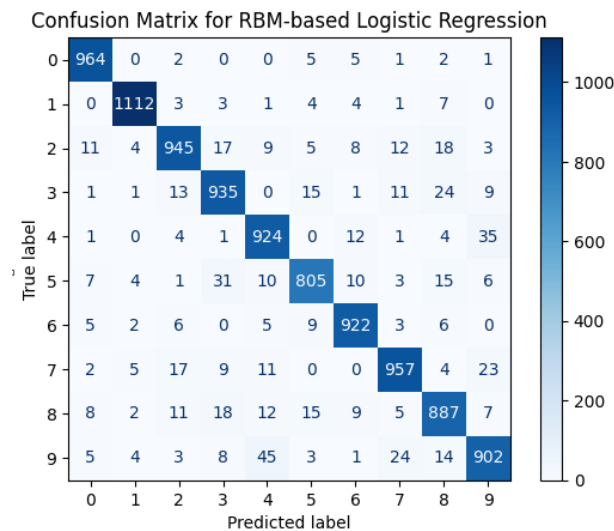


Figure 1: Confusion Matrix for Regular RBM

Now we examine the secure RBM. The model stays the same for the most part. But the data gets split vertically before passing to the sampling process. The sampling is done via the secure method described in the project sheet.

```

1 class SecureRBM:
2     def __init__(self, input=None, n_visible=784, n_hidden=128, W1=None, W2=None, hbias=None, vbias1=None, vbias2=None):
3         self.n_visible = n_visible
4         self.n_hidden = n_hidden
5
6         if numpy_rng is None:
7             numpy_rng = np.random.RandomState(1234)
8
9         if W1 is None:
10             a = 1. / n_visible
11             W1 = numpy_rng.uniform(-a, a, size=(n_visible//2, n_hidden))
12
13         if W2 is None:

```

```

14         a = 1. / n_visible
15         W2 = numpy_rng.uniform(-a, a, size=(n_visible//2, n_hidden))
16
17     if hbias is None:
18         hbias = np.zeros(n_hidden)
19
20     if vbias1 is None:
21         vbias = np.zeros(n_visible)
22
23
24     self.W1 = W1
25     self.W2 = W2
26     self.hbias = hbias
27     self.vbias = vbias
28     self.numpy_rng = numpy_rng
29
30     def sigmoid(self, x):
31         return 1 / (1 + np.exp(-x))
32
33     def sample_h_given_vA(self, v):
34         h_mean = np.dot(v, self.W1) + self.hbias
35         #h_sample = self.numpy_rng.binomial(size=h_mean.shape, n=1, p=h_mean)
36         return h_mean
37
38     def sample_h_given_vB(self, v):
39         h_mean = np.dot(v, self.W2) + self.hbias
40         #h_sample = self.numpy_rng.binomial(size=h_mean.shape, n=1, p=h_mean)
41         return h_mean
42
43     def sample_v_given_hA(self, h):
44         W = np.concatenate((self.W1, self.W2))
45         v_mean = np.dot(h, W.T) + self.vbias
46         #v_sample = self.numpy_rng.binomial(size=v_mean.shape, n=1, p=v_mean)
47         return v_mean
48
49     def sample_v_given_hB(self, h):
50         W = np.concatenate((self.W1, self.W2))
51         v_mean = np.dot(h, W.T) + self.vbias
52         #v_sample = self.numpy_rng.binomial(size=v_mean.shape, n=1, p=v_mean)
53         return v_mean
54
55
56     def gibbs_hvh_secure(self, h0_1, h0_2):
57         v1_mean1 = self.sample_v_given_hA(h0_1)
58         v1_mean2 = self.sample_v_given_hB(h0_1)
59
60         v_1 = self.secure_sig(v1_mean1, v1_mean2)
61
62         h1_mean1 = self.sample_h_given_vA(v_1[:, :self.n_visible//2])
63         h1_mean2 = self.sample_h_given_vB(v_1[:, self.n_visible//2:])
64
65         h_1 = self.secure_sig(h1_mean1, h1_mean2)
66
67         return v_1, v_1, h_1, h_1
68
69
70     def secure_sig(self, h1, h2):
71         return self.sigmoid(h1+h2)
72
73     def contrastive_divergence(self, data, k=1, lr=0.1):
74         data_A = data[:, :self.n_visible//2]
75         data_B = data[:, self.n_visible//2:]
76         ph_mean1 = self.sample_h_given_vA(data_A)
77         ph_mean2 = self.sample_h_given_vB(data_B)
78         h0 = self.secure_sigmoid(ph_mean1, ph_mean2)
79
80         nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh_secure(h0/2, h0/2)
81
82
83
84         chain_start = h0
85         for step in range(k):
86             if step == 0:

```

```

87         nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh_secure(chain_start/2 , chain_start/2)
88     else:
89         nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh_secure(nh_sample/2 , nh_sample/2)
90
91
92     self.W1 += lr * (np.dot(data_A.T, h0) - np.dot(nv_sample[: , :self.n_visible//2].T, nh_mean)) / len(da
93     self.W2 += lr * (np.dot(data_B.T, h0) - np.dot(nv_sample[: , self.n_visible//2:].T, nh_mean)) / len(da
94
95
96     self.vbias += lr * np.mean(data - nv_sample, axis=0)
97     self.hbias += lr * np.mean(h0 - nh_mean, axis=0)
98
99     def reconstruct(self, v):
100         h = self.sigmoid(np.dot(v, self.W) + self.hbias)
101         reconstructed_v = self.sigmoid(np.dot(h, self.W.T) + self.vbias)
102         return reconstructed_v
103
104     def save_model(self, file_path):
105         with open(file_path, 'wb') as f:
106             pickle.dump(self, f)
107
108     @staticmethod
109     def load_model(file_path):
110         with open(file_path, 'rb') as f:
111             return pickle.load(f)

```

The security operations are performed by these functions:

```

1 import torch
2 import random
3 from sympy import mod_inverse
4
5 # ElGamal key generation for tensors
6 def generate_keys(p, shape):
7     """Generate ElGamal public and private keys for tensors."""
8     g = torch.randint(2, p - 2, shape) # Generator tensor
9     x = torch.randint(1, p - 2, shape) # Private key tensor
10    y = torch.pow(g, x) % p # Public key tensor
11    return (p, g, y), x # Public key: (p, g, y), Private key: x
12
13 # ElGamal encryption for tensors
14 def encrypt(public_key, message):
15     p, g, y = public_key
16     k = torch.randint(1, p - 2, message.shape) # Random integer tensor k
17     c1 = torch.pow(g, k) % p
18     c2 = (message * torch.pow(y, k)) % p
19     return c1, c2
20
21 # ElGamal decryption (partial and final) for tensors
22 def decrypt_partial(private_key, public_key, ciphertext):
23     p, g, y = public_key
24     c1, c2 = ciphertext
25     s = torch.pow(c1, private_key) % p # Shared secret tensor
26     return c2, s # Partially decrypted message with the shared secret
27
28 def decrypt_final(partially_decrypted, shared_secret, p):
29     c2, s = partially_decrypted
30
31     # Handle the case where s = 0 (modular inverse of 0 does not exist)
32     s = torch.where(s == 0, torch.tensor(1, dtype=s.dtype), s)
33
34     # Compute the modular inverse of s
35     s_inv = torch.tensor([mod_inverse(int(val), p) for val in s.flatten()]).reshape(s.shape)
36
37     # Decrypt the message
38     message = (c2 * s_inv) % p
39     return message
40
41 # Sigmoid function for tensors
42 def sigmoid(z):
43     return 1 / (1 + torch.exp(-z)) # Sigmoid function for tensors
44

```



```

45 # Function to generate possible_input as a grid of tensors
46 def generate_possible_input(y):
47     """Generate the nearest tensor to y, where each element is rounded to the nearest 0.1."""
48     # Round each element of y to the nearest 0.1
49     possible_input = torch.round(y * 10) / 10
50     return [possible_input]
51
52 # Secure Sigmoid Computation for tensors
53 def secure_sigmoid(x, y):
54     """Securely compute sigmoid(x + y) using ElGamal encryption for tensors."""
55     p = 467 # A prime number
56     shape = x.shape # Shape of the input tensors
57
58     possible_input = generate_possible_input(y)
59
60     # Step 1: Party A generates a random number R and computes mi = sigmoid(x + i) - R for i in possible_input
61     # print(torch.max(y))
62     R = torch.randint(1, int(torch.max(y)) + 2, shape) # Random tensor R
63
64     # Generate keys
65     public_key, private_key = generate_keys(p, shape)
66
67     encrypted_messages = []
68
69     # Iterate over each tensor i in the set possible_input
70     for i in possible_input:
71         m_i = (sigmoid(x + i) * p - R).int() # Scale sigmoid to match finite field
72         encrypted_messages.append(encrypt(public_key, m_i))
73
74     # Step 2: Party B selects the encrypted message corresponding to its input y
75     selected_ciphertexts = encrypted_messages[0] # Reshape to match tensor shape
76
77     # Party B re-randomizes the selected ciphertexts
78     random_s = torch.randint(1, p - 1, shape)
79     randomized_c1 = (selected_ciphertexts[0] * torch.pow(public_key[1], random_s)) % p
80     randomized_c2 = (selected_ciphertexts[1] * torch.pow(public_key[2], random_s)) % p
81     randomized_encrypted = (randomized_c1, randomized_c2)
82
83     # Step 3: Party A partially decrypts and sends it back
84     partially_decrypted_c2, shared_secret_a = decrypt_partial(private_key, public_key, randomized_encrypted)
85
86     # Step 4: Party B completes decryption
87     final_message_b = decrypt_final((partially_decrypted_c2, shared_secret_a), random_s, p)
88
89     # Party B computes final sigmoid value
90     result = (final_message_b + R) / p # Reverse scaling
91
92     return torch.where(result > 1, torch.tensor(0.1), R/p), torch.where(result > 1, result - 1.1, final_message_b)
93
94 # Secure Multiplication Computation for tensors
95 def secure_multiplication(x, y):
96     """Securely compute x * y using ElGamal encryption for tensors."""
97     p = 467 # A prime number
98     shape = torch.matmul(x, y).shape # Shape of the input tensors
99
100     possible_input = generate_possible_input(y)
101
102     # Step 1: Party A generates a random number R and computes mi = (x * i) - R for i in possible_input
103     R = torch.randint(1, int(torch.max(y)) + 2, shape) # Random tensor R
104
105     # Generate keys
106     public_key, private_key = generate_keys(p, shape)
107
108     encrypted_messages = []
109
110     # Iterate over each tensor i in the set possible_input
111     for i in possible_input:
112         m_i = (torch.matmul(x, i) * 10 - R).int()
113         # print(m_i.shape) # Scale multiplication to match finite field
114         encrypted_messages.append(encrypt(public_key, m_i))
115
116     # Step 2: Party B selects the encrypted message corresponding to its input y
117     selected_ciphertexts = encrypted_messages[0]

```

```

118
119 # Party B re-randomizes the selected ciphertexts
120 random_s = torch.randint(1, p - 1, shape)
121 randomized_c1 = (selected_ciphertexts[0] * torch.pow(public_key[1], random_s)) % p
122 randomized_c2 = (selected_ciphertexts[1] * torch.pow(public_key[2], random_s)) % p
123 randomized_encrypted = (randomized_c1, randomized_c2)
124
125 # Step 3: Party A partially decrypts and sends it back
126 partially_decrypted_c2, shared_secret_a = decrypt_partial(private_key, public_key, randomized_encrypted)
127
128 # Step 4: Party B completes decryption
129 final_message_b = decrypt_final((partially_decrypted_c2, shared_secret_a), random_s, p)
130
131 # Party B computes final multiplication value
132 result = (final_message_b + R) / 10 # Reverse scaling
133
134 return torch.where(result > 1, torch.tensor(0.1), R/10), torch.where(result > 1, result - 1.1, final_messa

```

The model was only trained for 5 epochs due to the slower process of encryption and decryption. The Results are as follows:

Model	Precision	Recall	F1-score	accuracy
Secure RBM (CD=1)	0.85	0.85	0.85	0.85

Table 2: Performance of Secure RBM

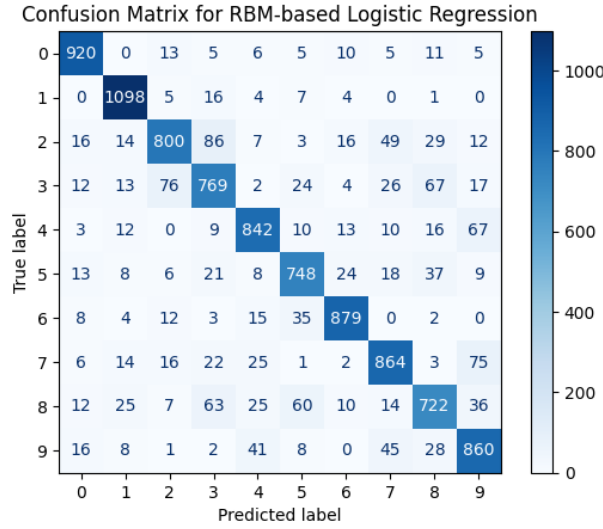


Figure 2: Confusion Matrix for Secure RBM

Theoretical Question 4 answer

1. Joint Distribution $p(y, x)$

The joint distribution $p(y, x, h)$ is defined as:

$$p(y, x, h) = \frac{e^{-E(y, x, h)}}{Z}$$

where Z is the partition function:

$$Z = \sum_{y,x,h} e^{-E(y,x,h)}$$

To obtain $p(y, x)$, we marginalize over h :

$$p(y, x) = \sum_h p(y, x, h) = \frac{1}{Z} \sum_h e^{-E(y,x,h)}$$

Using the energy function:

$$E(y, x, h) = -h^T W x - b^T x - c^T h - d^T \tilde{y} - h^T U \tilde{y}$$

we get:

$$p(y, x) = \frac{e^{b^T x + d^T \tilde{y}}}{Z} \sum_h e^{h^T (W x + c + U \tilde{y})}$$

Since $h_j \in \{0, 1\}$, the sum simplifies to:

$$\sum_h e^{h^T (W x + c + U \tilde{y})} = \prod_{j=1}^n (1 + e^{W_j x + c_j + U_j \tilde{y}})$$

Thus, the joint distribution is:

$$p(y, x) = \frac{e^{b^T x + d^T \tilde{y}}}{Z} \prod_{j=1}^n (1 + e^{W_j x + c_j + U_j \tilde{y}})$$

2. Conditional Distribution $p(y|x)$

The conditional probability follows:

$$p(y|x) = \frac{p(y, x)}{p(x)}$$

where:

$$p(x) = \sum_y p(y, x)$$

Substituting $p(y, x)$:

$$p(y|x) = \frac{e^{d^T \tilde{y}} \prod_{j=1}^n (1 + e^{W_j x + c_j + U_j \tilde{y}})}{\sum_{y^*} e^{d^T \tilde{y}^*} \prod_{j=1}^n (1 + e^{W_j x + c_j + U_j \tilde{y}^*})} \quad (1)$$

3. Log-likelihood for $p(y, x)$

The log-likelihood function is:

$$L_{\text{gen}}(D_{\text{train}}) = - \sum_{i=1}^{|D_{\text{train}}|} \log p(y_i, x_i)$$

The gradient is:

$$\frac{\partial \log p(y_i, x_i)}{\partial \theta} = -\mathbb{E}_{h|y_i, x_i} \left[\frac{\partial}{\partial \theta} E(y_i, x_i, h) \right] + \mathbb{E}_{y, x, h} \left[\frac{\partial}{\partial \theta} E(y, x, h) \right]$$

The second term is typically approximated using Contrastive Divergence (CD) or Persistent Contrastive Divergence (PCD).

4. Log-likelihood for $p(y|x)$

The conditional log-likelihood:

$$L_{\text{cond}}(D_{\text{train}}) = - \sum_{i=1}^{|D_{\text{train}}|} \log p(y_i | x_i)$$

Its gradient is:

$$\frac{\partial \log p(y_i | x_i)}{\partial \theta} = \frac{\partial}{\partial \theta} \log \left(\frac{e^{d_{y_i}} \prod_{j=1}^n (1 + e^{W_j x_i + c_j + U_j \tilde{y}_i})}{\sum_{y^*} e^{d_{y^*}} \prod_{j=1}^n (1 + e^{W_j x_i + c_j + U_j \tilde{y}^*})} \right)$$

5. Free Energy Representation $F(x, y)$

The free energy function is:

$$F(x, y) = -\log \sum_h e^{-E(y, x, h)}$$

Using the energy function:

$$F(x, y) = -b^T x - d^T \tilde{y} - \sum_{j=1}^n \log (1 + e^{W_j x + c_j + U_j \tilde{y}})$$

This representation is crucial for computing probabilities efficiently.

Theoretical Question 5

Prove the following equation:

$$\frac{\partial \log p(y_i | x_i)}{\partial \theta} = \sum_j \text{sigm}(o_{yj}(\mathbf{x}_i)) \frac{\partial o_{yj}(\mathbf{x}_i)}{\partial \theta} - \sum_{j, y^*} \text{sigm}(o_{y^*j}(\mathbf{x}_i)) p(y^* | x_i) \frac{\partial o_{y^*j}(\mathbf{x}_i)}{\partial \theta} \quad (2)$$

where

$$o_{yj}(x) = c_j + \sum_k W_{jk} x_k + U_{jy}$$

As a result, we can compute the gradient exactly.

Theoretical Question 5 answer

To derive the equality (2) we start from the conditional probability $p(y|x)$ and its log-likelihood. Let's break this down step by step.

Write the conditional probability $p(y|x)$

From the previous (1) derivation, the conditional probability $p(y|x)$ is:

$$p(y|x) = \frac{e^{d_y} \prod_{j=1}^n (1 + e^{o_{yj}(\mathbf{x})})}{\sum_{y^*} e^{d_{y^*}} \prod_{j=1}^n (1 + e^{o_{y^*j}(\mathbf{x})})},$$

where:

$$o_{yj}(\mathbf{x}) = c_j + \sum_k W_{jk} x_k + U_{jy}$$

Take the logarithm of $p(y|x)$

The log-likelihood of $p(y|x)$ is:

$$\log p(y|x) = \log \left(e^{d_y} \prod_{j=1}^n (1 + e^{o_{yj}(\mathbf{x})}) \right) - \log \left(\sum_{y^*} e^{d_{y^*}} \prod_{j=1}^n (1 + e^{o_{y^*j}(\mathbf{x})}) \right)$$

Simplify the first term:

$$\log \left(e^{d_y} \prod_{j=1}^n (1 + e^{o_{yj}(\mathbf{x})}) \right) = d_y + \sum_{j=1}^n \log (1 + e^{o_{yj}(\mathbf{x})})$$

Thus:

$$\log p(y|x) = d_y + \sum_{j=1}^n \log (1 + e^{o_{yj}(\mathbf{x})}) - \log \left(\sum_{y^*} e^{d_{y^*}} \prod_{j=1}^n (1 + e^{o_{y^*j}(\mathbf{x})}) \right)$$

Compute the gradient of $\log p(y|x)$ with respect to θ

We want to compute:

$$\frac{\partial \log p(y|x)}{\partial \theta}$$

From the expression for $\log p(y|x)$, the gradient is:

$$\frac{\partial \log p(y|x)}{\partial \theta} = \sum_{j=1}^n \frac{\partial}{\partial \theta} \log (1 + e^{o_{yj}(\mathbf{x})}) - \frac{\partial}{\partial \theta} \log \left(\sum_{y^*} e^{d_{y^*}} \prod_{j=1}^n (1 + e^{o_{y^*j}(\mathbf{x})}) \right)$$

Compute the first term

The first term is:

$$\sum_{j=1}^n \frac{\partial}{\partial \theta} \log (1 + e^{o_{yj}(\mathbf{x})})$$

Using the chain rule:

$$\frac{\partial}{\partial \theta} \log (1 + e^{o_{yj}(\mathbf{x})}) = \frac{e^{o_{yj}(\mathbf{x})}}{1 + e^{o_{yj}(\mathbf{x})}} \cdot \frac{\partial o_{yj}(\mathbf{x})}{\partial \theta}$$

Since:

$$\text{sigm}(o_{yj}(\mathbf{x})) = \frac{e^{o_{yj}(\mathbf{x})}}{1 + e^{o_{yj}(\mathbf{x})}},$$

we get:

$$\sum_{j=1}^n \frac{\partial}{\partial \theta} \log (1 + e^{o_{yj}(\mathbf{x})}) = \sum_{j=1}^n \text{sigm}(o_{yj}(\mathbf{x})) \frac{\partial o_{yj}(\mathbf{x})}{\partial \theta}$$

Compute the second term

The second term is:

$$\frac{\partial}{\partial \theta} \log \left(\sum_{y^*} e^{d_{y^*}} \prod_{j=1}^n (1 + e^{o_{y^*j}(\mathbf{x})}) \right)$$

By applying the chain rule and summing over y^* , we obtain:

$$\sum_{y^*} p(y^*|x) \sum_{j=1}^n \text{sigm}(o_{y^*j}(\mathbf{x})) \frac{\partial o_{y^*j}(\mathbf{x})}{\partial \theta}$$

Combine the terms

Thus, we arrive at:

$$\frac{\partial \log p(y_i|x_i)}{\partial \theta} = \sum_j \text{sigm}(o_{yj}(\mathbf{x}_i)) \frac{\partial o_{yj}(\mathbf{x}_i)}{\partial \theta} - \sum_{j,y^*} \text{sigm}(o_{y^*j}(\mathbf{x}_i)) p(y^*|x_i) \frac{\partial o_{y^*j}(\mathbf{x}_i)}{\partial \theta}$$

Simulation Question 2

Load the **MNIST** dataset and convert it to a binary state (for example, by applying a threshold to the pixel intensity). Then, train a **RBM** to estimate the joint distribution $p(y, x)$ using the algorithm 3. After training the model, perform final classification using $p(y|x)$ and report its accuracy.

Algorithm 1 Training update for RBM over (y, \mathbf{x}) using Contrastive Divergence

Input: training pair (y_i, \mathbf{x}_i) and learning rate λ
% Notation: $a \leftarrow b$ means a is set to value b
% $a \sim p$ means a is sampled from p

% Positive phase
 $y^0 \leftarrow y_i, \mathbf{x}^0 \leftarrow \mathbf{x}_i, \hat{\mathbf{h}}^0 \leftarrow \text{sigm}(\mathbf{c} + W\mathbf{x}^0 + U\mathbf{y}^0)$

% Negative phase
 $\mathbf{h}^0 \sim p(\mathbf{h}|y^0, \mathbf{x}^0), y^1 \sim p(y|\mathbf{h}^0), \mathbf{x}^1 \sim p(\mathbf{x}|\mathbf{h}^0)$
 $\hat{\mathbf{h}}^1 \leftarrow \text{sigm}(\mathbf{c} + W\mathbf{x}^1 + U\mathbf{y}^1)$

% Update
for $\theta \in \Theta$ **do**
 $\theta \leftarrow \theta - \lambda \left(\frac{\partial}{\partial \theta} E(y^0, \mathbf{x}^0, \hat{\mathbf{h}}^0) - \frac{\partial}{\partial \theta} E(y^1, \mathbf{x}^1, \hat{\mathbf{h}}^1) \right)$
end for

Figure 3: Training update for RBM over (y, x) using Contrastive Divergence

Simulation Question 2 answer

Dataset Preprocessing

We use the MNIST dataset and binarize it using a threshold of 0.5 after normalizing the values between 0 and 1. The processed dataset is shown in Figure 4.

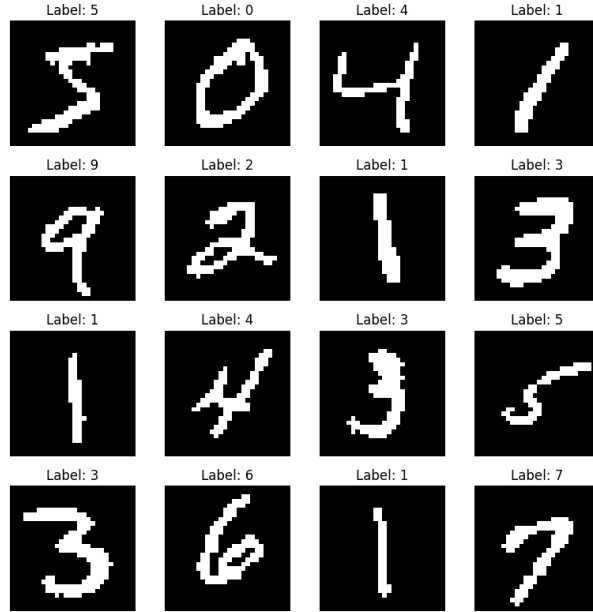


Figure 4: Binarized MNIST dataset.

RBM Model Definition

An RBM consists of visible units v and hidden units h , connected via a weight matrix W . The energy function of the RBM is defined as:

$$E(v, h) = -v^T W h - v^T b - h^T c$$

where b and c are biases for the visible and hidden layers, respectively.

Model Implementation

Below is the PyTorch implementation of our RBM model:

```

1
2 class RBM(nn.Module):
3     def __init__(self, n_visible=784+10, n_hidden=128, device='cpu'):
4         super(RBM, self).__init__()
5         self.n_visible = n_visible
6         self.n_hidden = n_hidden
7         self.device = device
8         self.W = nn.Parameter(torch.randn(n_visible, n_hidden) * 0.01)
9         self.hbias = nn.Parameter(torch.zeros(n_hidden))
10        self.vbias = nn.Parameter(torch.zeros(n_visible))
11        self.to(self.device)
12
13    def sample_h_given_v(self, v):
14        h_mean = torch.sigmoid(F.linear(v, self.W.t(), self.hbias))
15        return h_mean, torch.bernoulli(h_mean)
16
17    def sample_v_given_h(self, h):
18        v_mean = torch.sigmoid(F.linear(h, self.W, self.vbias))
19        return v_mean, torch.bernoulli(v_mean)
20
21    def gibbs_hvh(self, h0_sample):
22        v_mean, v_sample = self.sample_v_given_h(h0_sample)
23        h_mean, h_sample = self.sample_h_given_v(v_sample)
24        return v_mean, v_sample, h_mean, h_sample
25

```

```

26 def contrastive_divergence(self, data, k=1, lr=0.1):
27     data = data.to(self.device)
28     ph_mean, ph_sample = self.sample_h_given_v(data)
29     chain_start = ph_sample
30     for _ in range(k):
31         nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh(chain_start)
32         chain_start = nh_sample
33     self.W.data += lr * (torch.mm(data.t(), ph_sample) -
34                          torch.mm(nv_sample.t(), nh_mean)) / data.size(0)
35     self.vbias.data += lr * torch.mean(data - nv_sample, dim=0)
36     self.hbias.data += lr * torch.mean(ph_sample - nh_mean, dim=0)
37
38 def energy(self, v, h):
39     return -torch.matmul(v, torch.matmul(self.W, h)) -
40            torch.matmul(v, self.vbias) - torch.matmul(h, self.hbias)
41
42 def reconstruct(self, v):
43     h = torch.sigmoid(F.linear(v, self.W.t(), self.hbias))
44     return torch.sigmoid(F.linear(h, self.W, self.vbias))
45
46 def save_model(self, file_path):
47     with open(file_path, 'wb') as f:
48         pickle.dump(self, f)
49
50 @staticmethod
51 def load_model(file_path, device=None):
52     with open(file_path, 'rb') as f:
53         model = pickle.load(f)
54     model.device = device if device else torch.device("cuda" if torch.cuda.is_available() else "cpu")
55     model.to(model.device)
56     return model

```

Training the RBM

We train the RBM using Contrastive Divergence (CD- k), where $k = 1$ is commonly used for efficiency. The training function is implemented as follows:

```

1 def train_rbm(train_loader, test_loader, n_visible, n_hidden, learning_rate, epochs, k, device):
2     rbm = RBM(n_visible, n_hidden, device).to(device)
3     optimizer = torch.optim.SGD([rbm.W, rbm.hbias, rbm.vbias], lr=learning_rate)
4
5     for epoch in range(epochs):
6         for batch, label in train_loader:
7             batch = batch.view(batch.size(0), -1).to(device)
8             batch = torch.cat((batch, label.to(device)), dim=1)
9             optimizer.zero_grad()
10            rbm.contrastive_divergence(batch, k=k)
11            optimizer.step()
12    return rbm

```

Classification Using RBM Energy

The trained RBM can classify images by computing the energy of all possible labels and selecting the one with the lowest energy:

```

1 def classify_rbm(model, x, device="cpu"):
2     model.to(device)
3     x = x.flatten().to(device)
4     y_classes = torch.eye(10, device=device)
5     x_repeated = x.repeat(10, 1)
6     xy_concat = torch.cat((x_repeated, y_classes), dim=1)
7     _, h_sampled = model.sample_h_given_v(xy_concat)
8     energies = torch.tensor([model.energy(v, h) for v, h in zip(xy_concat, h_sampled)], device=device)
9     return torch.argmax(energies).item()

```

All the results have been gathered in theoretical question 6.

Simulation Question 3

Load the MNIST dataset and convert it to a binary format (e.g., by applying a threshold to pixel intensities). Then, train an RBM to estimate the distribution $p(x, y)$ using the algorithm in Figure 2. After training the model, perform classification using $p(y|x)$ and report its accuracy.

Simulation Question 3 answer

The DRBM model is similar to the RBM but differs in how it handles visible units. Instead of treating all visible units the same, it explicitly separates the feature vector x and the label vector y . During Gibbs sampling, x is kept fixed, and only y is sampled. like Algorithm 5 .

Algorithm 2 Training update for DRBM over (y, \mathbf{x}) using Contrastive Divergence

Input: training pair (y_i, \mathbf{x}_i) and learning rate λ
 % Notation: $a \leftarrow b$ means a is set to value b
 % $a \sim p$ means a is sampled from p

% Positive phase
 $y^0 \leftarrow y_i, \mathbf{x}^0 \leftarrow \mathbf{x}_i, \hat{\mathbf{h}}^0 \leftarrow \text{sigm}(\mathbf{c} + W\mathbf{x}^0 + U\mathbf{y}^0)$

% Negative phase
 $\mathbf{h}^0 \sim p(\mathbf{h}|y^0, \mathbf{x}^0), y^1 \sim p(y|\mathbf{h}^0)$
 $\hat{\mathbf{h}}^1 \leftarrow \text{sigm}(\mathbf{c} + W\mathbf{x}^0 + U\mathbf{y}^1)$

% Update
for $\theta \in \Theta$ **do**
 $\theta \leftarrow \theta - \lambda \left(\frac{\partial}{\partial \theta} E(y^0, \mathbf{x}^0, \hat{\mathbf{h}}^0) - \frac{\partial}{\partial \theta} E(y^1, \mathbf{x}^0, \hat{\mathbf{h}}^1) \right)$
end for

Figure 5: Training update for DRBM over $(y|x)$ using Contrastive Divergence

Model Implementation

```

1 class DRBM(nn.Module):
2     def __init__(self, n_visible=784, n_hidden=128, label_size=10, device='cpu'):
3         super(DRBM, self).__init__()
4         self.n_visible = n_visible
5         self.n_hidden = n_hidden
6         self.label_size = label_size
7         self.feature_size = n_visible - label_size
8         self.device = device
9
10        self.W = nn.Parameter(torch.randn(n_visible, n_hidden) * 0.01)
11        self.hbias = nn.Parameter(torch.zeros(n_hidden))
12        self.vbias = nn.Parameter(torch.zeros(n_visible))
13
14        self.to(self.device)
15
16        def sample_h_given_v(self, v):
17            h_mean = torch.sigmoid(F.linear(v, self.W.t(), self.hbias))
18            h_sample = torch.bernoulli(h_mean)
19            return h_mean, h_sample
20
21        def sample_v_given_h(self, h, x_fixed):
22            v_mean = torch.sigmoid(F.linear(h, self.W, self.vbias))
23            v_sample = v_mean.clone()

```

```

24     v_sample[:, :self.feature_size] = x_fixed
25     v_sample[:, self.feature_size:] = torch.bernoulli(v_mean[:, self.feature_size:])
26     return v_mean, v_sample
27
28     def gibbs_hvh(self, h0_sample, x_fixed):
29         v_mean, v_sample = self.sample_v_given_h(h0_sample, x_fixed)
30         h_mean, h_sample = self.sample_h_given_v(v_sample)
31         return v_mean, v_sample, h_mean, h_sample
32
33     def contrastive_divergence(self, data, k=1, lr=0.1):
34         data = data.to(self.device)
35         x_fixed = data[:, :self.feature_size]
36         ph_mean, ph_sample = self.sample_h_given_v(data)
37
38         chain_start = ph_sample
39         for _ in range(k):
40             nv_mean, nv_sample, nh_mean, nh_sample = self.gibbs_hvh(chain_start, x_fixed)
41             chain_start = nh_sample
42
43         self.W.data +=
44             lr * (torch.mm(data.t(), ph_sample) - torch.mm(nv_sample.t(), nh_mean)) / data.size(0)
45         self.vbias.data += lr * torch.mean(data - nv_sample, dim=0)
46         self.hbias.data += lr * torch.mean(ph_sample - nh_mean, dim=0)
47
48     def energy(self, v, h):
49         term1 = -torch.sum((v @ self.W) * h, dim=1)
50         term2 = -torch.matmul(v, self.vbias)
51         term3 = -torch.matmul(h, self.hbias)
52         return term1 + term2 + term3
53
54     def reconstruct(self, v):
55         v = v.to(self.device)
56         h = torch.sigmoid(F.linear(v, self.W.t(), self.hbias))
57         reconstructed_v = torch.sigmoid(F.linear(h, self.W, self.vbias))
58         reconstructed_v[:, :self.feature_size] = v[:, :self.feature_size]
59         return reconstructed_v, v[:, self.feature_size:]

```

Training and Classification Using RBM Energy

The training and classification process follows the same approach as the standard RBM, with one key difference: in DRBM, Contrastive Divergence (CD) is modified to ensure that x remains fixed while sampling y . This adjustment allows the model to learn the conditional probability $p(y|x)$ directly.

All the results have been gathered in theoretical question 6.

Simulation Question 4

Once again, perform training using the hybrid method with $\alpha \in \{1, 0.1, 0.01\}$, then classify using $p(y|x)$ and report the final accuracy.

Simulation Question 4 answer

The HybridRBM model is a modified Restricted Boltzmann Machine (RBM) that incorporates a unified learning approach by combining two types of loss functions: the standard RBM loss and the Discriminative RBM (DRBM) loss, weighted by a parameter α .

Model Implementation

Below is the implementation of the HybridRBM model:

```
1 class UnifiedRBM(nn.Module):
2     def __init__(self, n_visible=784+10, n_hidden=128, label_size=10, alpha=0.5, device='cpu'):
3         super(UnifiedRBM, self).__init__()
4
5         self.n_visible = n_visible
6         self.n_hidden = n_hidden
7         self.label_size = label_size
8         self.feature_size = n_visible - label_size
9         self.alpha = alpha
10        self.device = device
11
12        self.W = nn.Parameter(torch.randn(self.n_visible, n_hidden) * 0.01)
13        self.hbias = nn.Parameter(torch.zeros(n_hidden))
14        self.vbias = nn.Parameter(torch.zeros(self.n_visible))
15
16        self.to(self.device)
17
18    def sample_h_given_v(self, v):
19        h_mean = torch.sigmoid(F.linear(v, self.W.t(), self.hbias))
20        h_sample = torch.bernoulli(h_mean)
21        return h_mean, h_sample
22
23    def sample_v_given_h(self, h, x_fixed=None):
24        v_mean = torch.sigmoid(F.linear(h, self.W, self.vbias))
25        v_sample = torch.bernoulli(v_mean)
26
27        if x_fixed is not None:
28            v_sample[:, :self.feature_size] = x_fixed
29        return v_mean, v_sample
30
31    def gibbs_hvh(self, h0_sample, x_fixed=None):
32        v_mean, v_sample = self.sample_v_given_h(h0_sample, x_fixed)
33        h_mean, h_sample = self.sample_h_given_v(v_sample)
34        return v_mean, v_sample, h_mean, h_sample
35
36    def contrastive_divergence(self, data, k=1, lr=0.1):
37        data = data.to(self.device)
38        x_fixed = data[:, :self.feature_size]
39
40        ph_mean_dis, ph_sample_dis = self.sample_h_given_v(data)
41        chain_start_dis = ph_sample_dis
42        for _ in range(k):
43            nv_mean_dis, nv_sample_dis, nh_mean_dis, nh_sample_dis =
44                self.gibbs_hvh(chain_start_dis, x_fixed)
45            chain_start_dis = nh_sample_dis
46
47        W_update_dis = torch.mm(data.t(), ph_sample_dis) - torch.mm(nv_sample_dis.t(), nh_mean_dis)
48        vbias_update_dis = torch.mean(data - nv_sample_dis, dim=0)
49        hbias_update_dis = torch.mean(ph_sample_dis - nh_mean_dis, dim=0)
50
51        ph_mean_gen, ph_sample_gen = self.sample_h_given_v(data)
52        chain_start_gen = ph_sample_gen
53        for _ in range(k):
54            nv_mean_gen, nv_sample_gen, nh_mean_gen, nh_sample_gen = self.gibbs_hvh(chain_start_gen)
55            chain_start_gen = nh_sample_gen
56
57        W_update_gen = torch.mm(data.t(), ph_sample_gen) - torch.mm(nv_sample_gen.t(), nh_mean_gen)
58        vbias_update_gen = torch.mean(data - nv_sample_gen, dim=0)
59        hbias_update_gen = torch.mean(ph_sample_gen - nh_mean_gen, dim=0)
60
61        W_update = W_update_dis + (self.alpha) * W_update_gen
62        vbias_update = vbias_update_dis + (self.alpha) * vbias_update_gen
63        hbias_update = hbias_update_dis + (self.alpha) * hbias_update_gen
64
65        self.W.data += lr * W_update / data.size(0)
66        self.vbias.data += lr * vbias_update
67        self.hbias.data += lr * hbias_update
68
```

```

69 def energy(self, v, h):
70     term1 = -torch.matmul(v, torch.matmul(self.W, h))
71     term2 = -torch.matmul(v, self.vbias)
72     term3 = -torch.matmul(h, self.hbias)
73     return term1 + term2 + term3
74
75 def reconstruct(self, v):
76     v = v.to(self.device)
77     h = torch.sigmoid(F.linear(v, self.W.t(), self.hbias))
78     reconstructed_v = torch.sigmoid(F.linear(h, self.W, self.vbias))
79     return reconstructed_v, v[:, self.feature_size:]
80
81 def save_model(self, file_path):
82     with open(file_path, 'wb') as f:
83         pickle.dump(self, f)
84
85 @staticmethod
86 def load_model(file_path, device=None):
87     with open(file_path, 'rb') as f:
88         model = pickle.load(f)
89     model.device = device if device else torch.device("cuda" if torch.cuda.is_available() else "cpu")
90     model.to(model.device)
91     return model

```

Training and Classification

Training and classification using the HybridRBM model follows the standard RBM process, but the Contrastive Divergence (CD) method has been modified within the model's implementation.

All the results have been gathered in theoretical question 6.

Theoretical Question 6

Report the final learning outcomes and classification accuracy for all previous models and compare them. Are the results as expected?

Theoretical Question 6 answer

Model	Precision	Recall	F1-score
RBM (CD=1)	0.78	0.78	0.78
RBM (CD=5)	0.79	0.78	0.78
RBM (CD=10)	0.78	0.78	0.78
DRBM (CD=1)	0.81	0.80	0.80
DRBM (CD=5)	0.83	0.82	0.82
DRBM (CD=10)	0.83	0.83	0.83
Mix ($\alpha = 0.01$, CD=1)	0.77	0.76	0.76
Mix ($\alpha = 0.1$, CD=1)	0.80	0.80	0.80
Mix ($\alpha = 1$, CD=1)	0.82	0.82	0.82
Mix ($\alpha = 0.01$, CD=5)	0.80	0.79	0.79
Mix ($\alpha = 0.1$, CD=5)	0.81	0.81	0.81
Mix ($\alpha = 1$, CD=5)	0.83	0.82	0.82
Mix ($\alpha = 0.01$, CD=10)	0.81	0.81	0.81
Mix ($\alpha = 0.1$, CD=10)	0.82	0.82	0.82
Mix ($\alpha = 1$, CD=10)	0.83	0.83	0.83

Table 3: Performance comparison of different RBM-based models. Bold values indicate the best performance in each metric.

Analysis of Model Performance

Table 3 highlights some key insights into how different models perform:

- **RBM Performance:** The standard RBM with Contrastive Divergence (CD) shows stable performance across different CD values. Increasing CD from 1 to 5 leads to slight improvements, but beyond that (CD=10), the gains are negligible. The best performance is observed at CD=5, with a precision of 0.79.
- **DRBM Performance:** The Discriminative RBM (DRBM) performs significantly better than the standard RBM. This suggests that explicitly separating the feature vector x from the label vector y enhances classification. Performance improves as CD increases, reaching its best at CD=10, where all metrics achieve 0.83. This indicates that deeper sampling benefits the DRBM framework.
- **Mix Model Performance:** The hybrid Mix RBM model, which combines RBM and DRBM losses, is highly sensitive to the tuning parameter α . The results show that:
 - At $\alpha = 0.01$, the model underperforms compared to DRBM, likely because the RBM loss dominates.
 - At $\alpha = 0.1$, performance improves, approaching that of DRBM.
 - At $\alpha = 1$, the model matches DRBM (CD=10), suggesting that at this setting, the DRBM loss is the dominant factor.

This suggests that adjusting α allows the Mix model to balance between RBM and DRBM characteristics, making it a flexible option for different tasks.

Expected vs. Actual Results

- As anticipated, DRBM outperforms the standard RBM, reinforcing the idea that conditioning the visible layer helps classification.
- The Mix model with $\alpha = 1$ achieves performance on par with DRBM (CD=10), confirming that combining both training approaches is effective.
- Increasing CD generally enhances performance, but improvements taper off after CD=5, particularly for RBM.

Conclusion

- **Best overall model:** The DRBM (CD=10) and Mix model ($\alpha = 1, CD = 10$) both achieve the highest performance, with all metrics reaching 0.83.
- **Hybrid approach is effective:** The Mix model can achieve DRBM-level performance while offering flexibility through α .

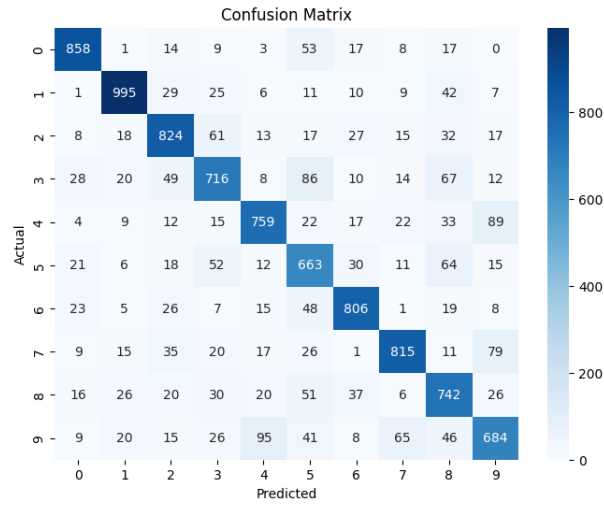


Figure 6: Confusion Matrix for RBM

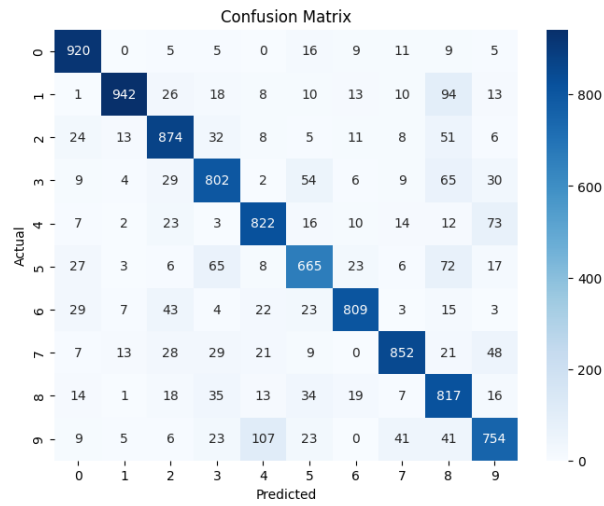


Figure 7: Confusion Matrix for DRBM

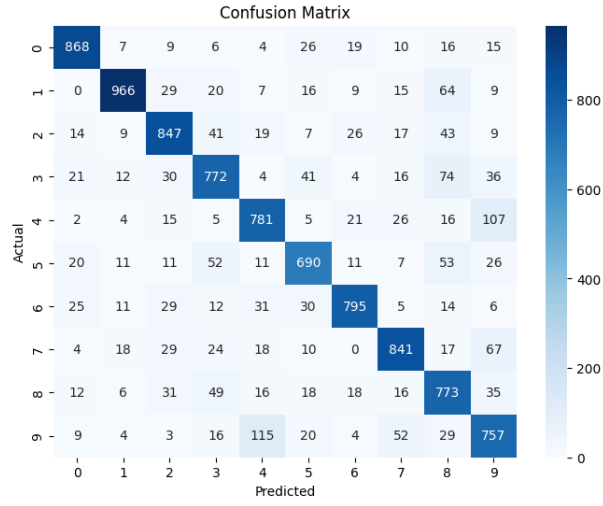


Figure 8: Confusion Matrix for Mix $\alpha = 0.01$

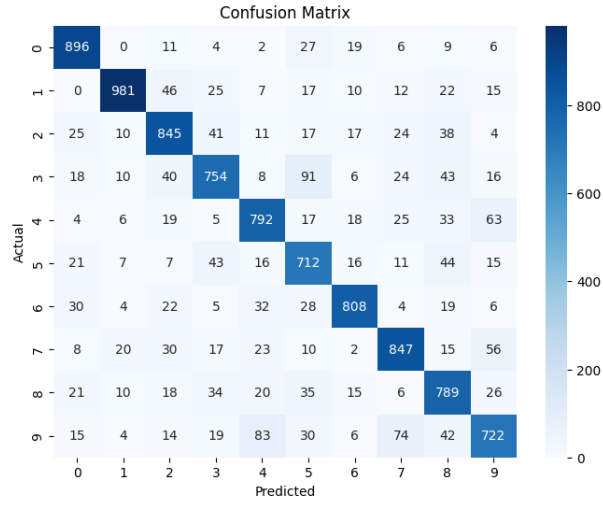


Figure 9: Confusion Matrix for Mix $\alpha = 0.1$

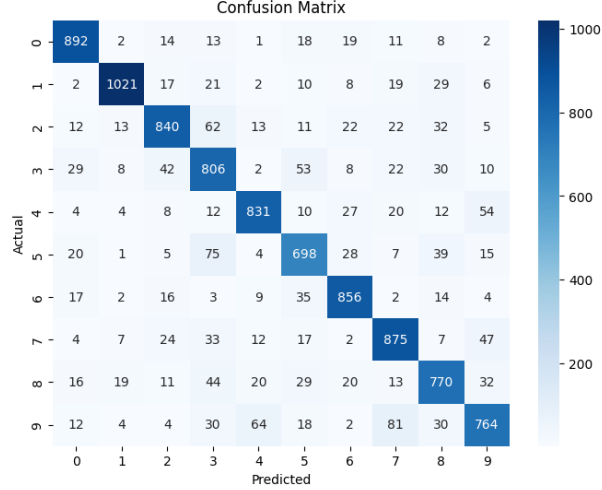


Figure 10: Confusion Matrix for Mix $\alpha = 1$

Theoretical Question 7

Can we use unlabeled data D_{unsup} alongside labeled data D_{sup} ? To be precise, we aim for a final objective function as follows:

$$\mathcal{L}_{semi-sup}(D_{sup}, D_{unsup}) = \mathcal{L}_{sup}(D_{sup}) + \beta \mathcal{L}_{unsup}(D_{unsup})$$

Propose a suitable function for \mathcal{L}_{unsup} and explain the details of its optimization.

Theoretical Question 7 answer

Yes, we can use unlabeled data D_{unsup} alongside labeled data D_{sup} in semi-supervised learning. Semi-supervised learning algorithms leverage the unlabeled data to introduce constraints on the trained model, which helps improve generalization.

Defining \mathcal{L}_{unsup}

A suitable choice for \mathcal{L}_{unsup} is the negative log-likelihood of the unlabeled data under the model's distribution:

$$\mathcal{L}_{unsup}(D_{unsup}) = - \sum_{i=1}^{|D_{unsup}|} \log p(x_i)$$

where $D_{unsup} = \{x_i\}_{i=1}^{|D_{unsup}|}$ represents the set of unlabeled inputs. This function ensures that the model learns a good generative representation of the unlabeled data, making it useful for semi-supervised learning.

Optimization Strategy

To optimize this loss, we use **contrastive divergence**, which provides an approximation of the log-likelihood gradient. The gradient is computed as:

$$\frac{\partial \log p(x_i)}{\partial \theta} = -\mathbb{E}_{y,h|x_i} \left[\frac{\partial}{\partial \theta} E(y_i, x_i, h) \right] + \mathbb{E}_{y,x,h} \left[\frac{\partial}{\partial \theta} E(y, x, h) \right]$$

where $E(y, x, h)$ is the energy function in the RBM framework. The first term can be efficiently computed using:

$$\mathbb{E}_{y|x_i} \left[\mathbb{E}_{h|y,x_i} \left[\frac{\partial}{\partial \theta} E(y_i, x_i, h) \right] \right]$$

This optimization can be performed using **sampling techniques**, where we either:

1. Compute the expected gradient over all possible class labels y , weighted by $p(y|x_i)$, or
2. Sample a class label $y \sim p(y|x_i)$ and compute the gradient for the selected y .

Final Semi-Supervised Objective

The final objective function for semi-supervised learning is:

$$\mathcal{L}_{semi-sup}(D_{sup}, D_{unsup}) = \mathcal{L}_{TYPE}(D_{sup}) + \beta \mathcal{L}_{unsup}(D_{unsup})$$

where \mathcal{L}_{TYPE} corresponds to a supervised loss, which can be generative (*gen*), discriminative (*disc*), or a hybrid of both. The parameter β controls the contribution of the unsupervised loss.