

# HOMEWORK 2

Nikhil Malik  
nmalik1@andrew.cmu.edu

## Problem 1 (10 pts)

Loss Function

$$J = -\log(y_m) \quad (1)$$

$$\frac{\partial J}{\partial \mathbf{y}} = [0, 0, \dots, -1/y_m, \dots, 0] \quad (2)$$

$$(3)$$

Softmax Normalization

$$y_m = \exp(x_m) / (\sum_k \exp(x_k)) \quad (4)$$

$$\frac{\partial y_i}{\partial x_j} = y_i * [(i == j) - y_j] \quad (5)$$

$$(6)$$

Softmax Weights: Where  $d_2$  is the dimensionality of two post max pool feature maps

$$x_i = \beta_i + \sum_{l=1}^2 \sum_{e=1}^{d_2} \sum_{f=1}^{d_2} w_{efi}^l * v_{ef}^l \quad (7)$$

$$\frac{\partial x_i}{\partial \beta_i} = 1 \quad (8)$$

$$\frac{\partial x_i}{\partial w_{efi}^l} = v_{ef}^l \quad (9)$$

$$\frac{\partial x_i}{\partial v_{ef}^l} = w_{efi}^l \quad (10)$$

$$(11)$$

Maxpool Operation

$$v_{ef}^l = \max((a_{e*p+c, f*p+d}^l)_{d=0}^{p-1})_{c=0}^{p-1} \quad (12)$$

$$(13)$$

$$\frac{\partial v_{ef}^l}{\partial a_{gh}^l} = \begin{cases} 1, & a_{gh}^l \text{ was the max in max pool window.} \\ 0, & \text{otherwise.} \end{cases} \quad (14)$$

$$(15)$$

## Convolution Operation

$$a_{gh}^l = \theta_l + \sum_{u=0}^{k-1} \sum_{v=0}^{k-1} \theta_{uv}^l * input_{g*s+u, h*s+v} \quad (16)$$

$$\frac{\partial a_{gh}^l}{\partial \theta_{ab}^l} = \sum_{g,h} input_{g*s+a, h*s+b} \quad (17)$$

$$(18)$$

Chain rule expressions, where all the individuals components have been calculated above

$$\frac{\partial J}{\partial \beta_i} = \frac{\partial J}{\partial \mathbf{y}} * \sum_j \left( \frac{\partial y_j}{\partial x_i} * \frac{\partial x_i}{\partial \beta_i} \right) \quad (19)$$

$$\frac{\partial J}{\partial w_{efi}^l} = \frac{\partial J}{\partial \mathbf{y}} * \sum_j \left( \frac{\partial y_j}{\partial x_i} * \frac{\partial x_i}{\partial w_{efi}^l} \right) \quad (20)$$

$$\frac{\partial J}{\partial \theta_{ab}^l} = \frac{\partial J}{\partial \mathbf{y}} * \sum_j \left( \frac{\partial y_j}{\partial x_i} * \frac{\partial x_i}{\partial v_{ef}^l} * \sum_{e,f} \left( \frac{\partial v_{ef}^l}{\partial a_{gh}^l} * \sum_{g,h} \frac{\partial a_{gh}^l}{\partial \theta_{ab}^l} \right) \right) \quad (21)$$

**Problem 2 (10 pts)**

$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k | \mathbf{pa}_k) \quad (22)$$

We need to prove

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = 1 \quad (23)$$

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = \sum_{x_1, x_2, \dots, x_K} \prod_{k=1}^K p(x_k | \mathbf{pa}_k) \quad (24)$$

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = \sum_{x_1, x_2, \dots, x_K} p(x_K | \mathbf{pa}_K) * \prod_{k=1}^{K-1} p(x_k | \mathbf{pa}_k) \quad (25)$$

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = \sum_{x_1, x_2, \dots, x_{K-1}} \left( \sum_{x_K} p(x_K | \mathbf{pa}_K) \right) * \prod_{k=1}^{K-1} p(x_k | \mathbf{pa}_k) \quad (26)$$

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = \sum_{x_1, x_2, \dots, x_{K-1}} (1) * \prod_{k=1}^{K-1} p(x_k | \mathbf{pa}_k) \quad (27)$$

Repeating the same steps as above K-1 times.

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = (1) * (1) * \dots * (1) \quad (28)$$

$$\sum_{x_1, x_2, \dots, x_K} p(\mathbf{x}) = 1 \quad (29)$$

This proves that the original representation was correctly normalized.

**Problem 3 (10pts)**

$$E(\mathbf{v}, \mathbf{h}; \theta) = - \sum_{i=1}^D \sum_{j=1}^P W_{ij} v_i h_j - \sum_{i=1}^D v_i b_i - \sum_{j=1}^P h_j a_j. \quad (30)$$

$$P_\theta(\mathbf{v}, \mathbf{h}) = \frac{1}{\mathcal{Z}} \exp(-E(\mathbf{v}, \mathbf{h}; \theta)) \quad (31)$$

$$\mathcal{Z} = \sum_{\mathbf{v}, \mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}; \theta)) \quad (32)$$

Since  $h_2, h_3, \dots, h_P$  are not part of markov blanket of  $h_1$ , we can write

$$P_\theta(h_1, h_2, \dots, h_P | \mathbf{v}) = P_\theta(h_1 | \mathbf{v}) * P_\theta(h_2, \dots, h_P | \mathbf{v}) \quad (33)$$

Repeating the same step P times

$$P_\theta(h_1, h_2, \dots, h_P | \mathbf{v}) = P_\theta(h_1 | \mathbf{v}) * P_\theta(h_2 | \mathbf{v}) * \dots * P_\theta(h_P | \mathbf{v}) \quad (34)$$

$$P_\theta(h_1, h_2, \dots, h_P | \mathbf{v}) = \prod_j P_\theta(h_j | \mathbf{v}) \quad (35)$$

$$P_\theta(h_j | \mathbf{v}) = \frac{1}{\mathcal{Z}_j} \exp(-E(\mathbf{v}, h_j; \theta)) \quad (36)$$

$$P_\theta(h_j | \mathbf{v}) = \frac{1}{\mathcal{Z}_j} \exp\left(\sum_{i=1}^D W_{ij} v_i h_j + \sum_{i=1}^D v_i b_i + h_j a_j\right) \quad (37)$$

$$P_\theta(h_j = 1 | \mathbf{v}) = \frac{1}{\mathcal{Z}_j} \exp\left(\sum_{i=1}^D W_{ij} v_i + \sum_{i=1}^D v_i b_i + a_j\right) \quad (38)$$

$$P_\theta(h_j = 0 | \mathbf{v}) = \frac{1}{\mathcal{Z}_j} \exp\left(\sum_{i=1}^D v_i b_i\right) \quad (39)$$

$$P_\theta(h_j = 1 | \mathbf{v}) = P_\theta(h_j = 1 | \mathbf{v}) / (P_\theta(h_j = 1 | \mathbf{v}) + P_\theta(h_j = 0 | \mathbf{v})) \quad (40)$$

$$P_\theta(h_j = 1 | \mathbf{v}) = \exp\left(\sum_{i=1}^D W_{ij} v_i + \sum_{i=1}^D v_i b_i + a_j\right) / \left(\exp\left(\sum_{i=1}^D W_{ij} v_i + \sum_{i=1}^D v_i b_i + a_j\right) + \exp\left(\sum_{i=1}^D v_i b_i\right)\right) \quad (41)$$

$$P_\theta(h_j = 1 | \mathbf{v}) = \exp\left(\sum_{i=1}^D W_{ij} v_i + a_j\right) / \left(\exp\left(\sum_{i=1}^D W_{ij} v_i + a_j\right) + 1\right) \quad (42)$$

$$P_\theta(h_j = 1 | \mathbf{v}) = 1 / (1 + \exp(-\sum_{i=1}^D W_{ij} v_i - a_j)) \quad (43)$$

$$P_\theta(h_j = 1 | \mathbf{v}) = \text{Sigm}\left(\sum_{i=1}^D W_{ij} v_i + a_j\right) \quad (44)$$

$$(45)$$

This proves the original expression

**Problem 4 (10 pts)**

$$E(\mathbf{x}, \mathbf{y}) = h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i, \quad (46)$$

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{y})) \quad (47)$$

- (5 pts)

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) = h \sum_{i \neq k} x_i - \beta \sum_{i \neq k, j} x_i x_j - \eta \sum_{i \neq k} x_i y_i + h - \beta \sum_{j \in ne(x_k)} x_j - \eta y_k \quad (48)$$

$$E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) = h \sum_{i \neq k} x_i - \beta \sum_{i \neq k, j} x_i x_j - \eta \sum_{i \neq k} x_i y_i - h + \beta \sum_{j \in ne(x_k)} x_j + \eta y_k \quad (49)$$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) - E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) = 2h - 2\beta \sum_{j \in ne(x_k)} x_j - 2\eta y_k \quad (50)$$

As per this expression the difference in energy depends only on neighbourhood of particular element of  $\mathbf{x}$ .

- (5 pts)

Substituting  $\beta$  and  $h$  in expression above

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) - E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) = -2\eta y_k \quad (51)$$

Given  $\eta \geq 0$

If  $y_k = 1$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) - E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) = -2\eta \quad (52)$$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) - E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) < 0 \quad (53)$$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) < E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) \quad (54)$$

$$(55)$$

Lower value of energy corresponds to higher likelihood. Therefore  $x_k = 1$

Similarly, if  $y_k = -1$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) - E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) = 2\eta \quad (56)$$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) - E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) > 0 \quad (57)$$

$$E(x_1, \dots, x_k = 1, \dots, x_K, \mathbf{y}) > E(x_1, \dots, x_k = -1, \dots, x_K, \mathbf{y}) \quad (58)$$

$$(59)$$

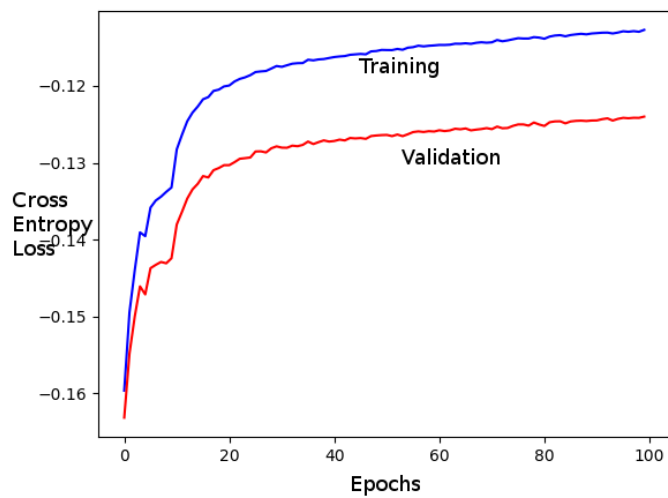
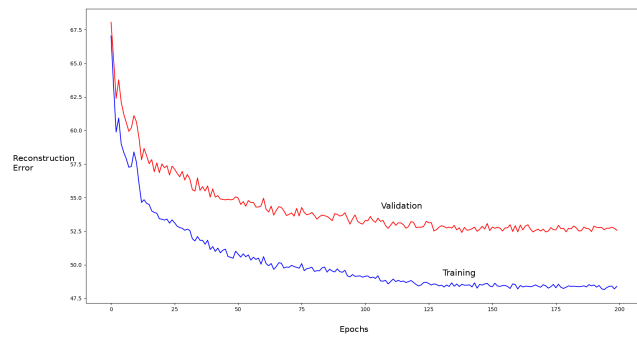
Higher value of energy corresponds to lower likelihood. Therefore  $x_k = -1$

From the two results above the most probable configuration of latent variables corresponds to  $y_k = x_k$

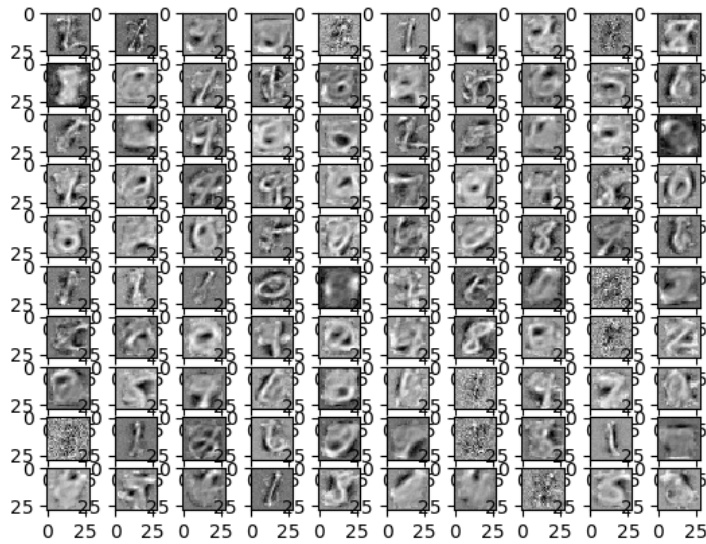
## Problem 5 (60 pts)

### Contrastive Divergence (CD), Autoencoders

#### a) Basic generalization [20 points]

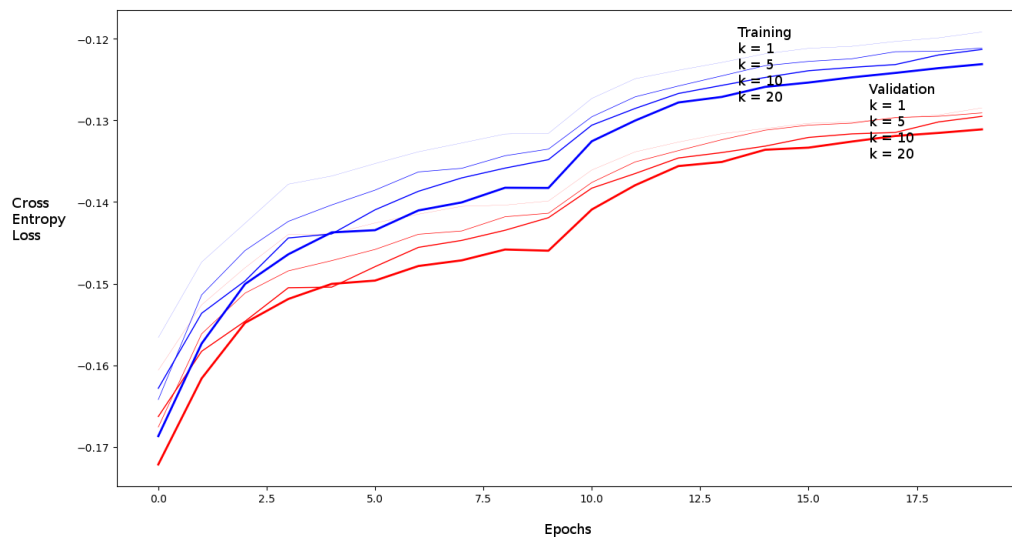


Training set error goes down faster compared to validation set as expected. After 125 epochs the validation error stops reducing representing the generalization error region.



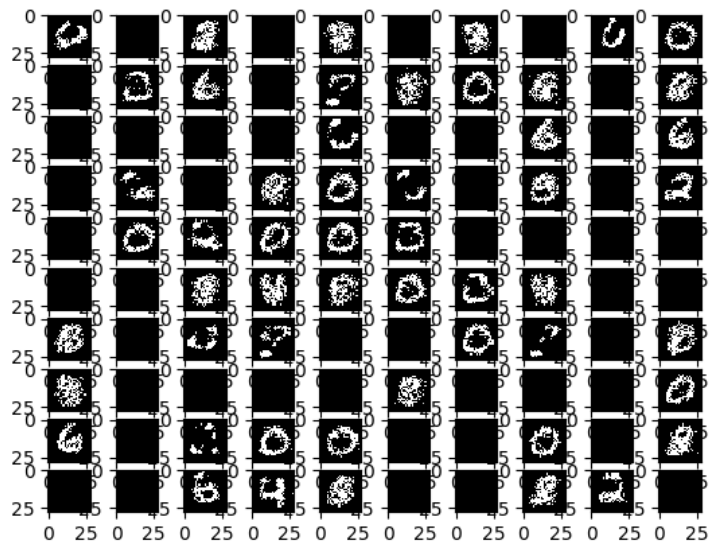
The visualized Weights  $W$  do represent structure similar to original digits or components such as curves that look similar to components of the original digits.

**b) Number of CD steps [5 points]**



In terms of convergence of training accuracy, higher value of  $k$  (indicated by thicker lines in the figure) converges slowly compared to smaller values (indicated by thinner lines in the figure) i.e.  $k=1$  converges the fastest.

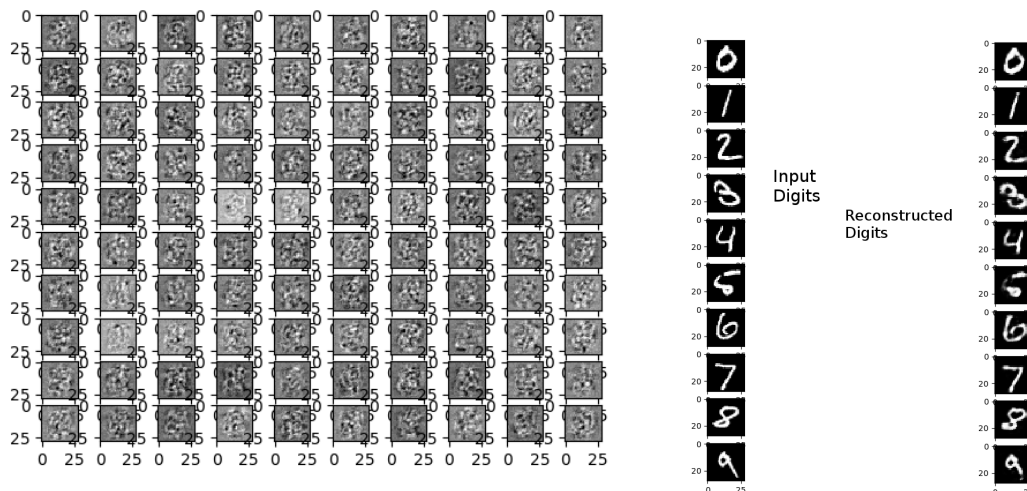
However in terms of generalization error,  $k=5$  does the best.  $k=10$  is comparable to  $k=5$ . While  $k=1$  has the highest generalization error.  $k=20$  underfits, at least until 50 epochs compared to smaller values of  $k$ . Discribe Observation

**c) Sampling from the RBM model [5 points]**

They do show a very rough structure of handwritten digits.

**d) Unsupervised Learning as pretraining [5 points]**

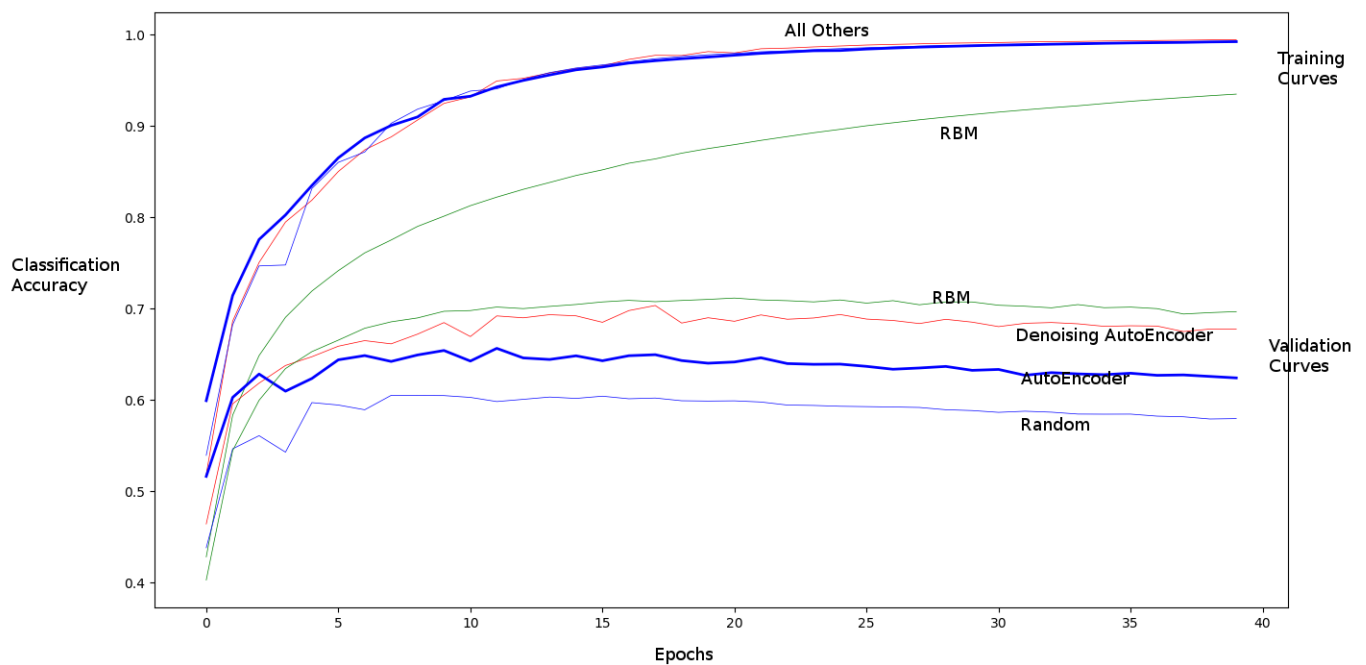
Yes, rbm initialization improves the model generalizability. It achieves the best validation accuracy, and tends to not overfit unlike all other initialization methods. (see Figure A)

**e) Autoencoder [5 points]**

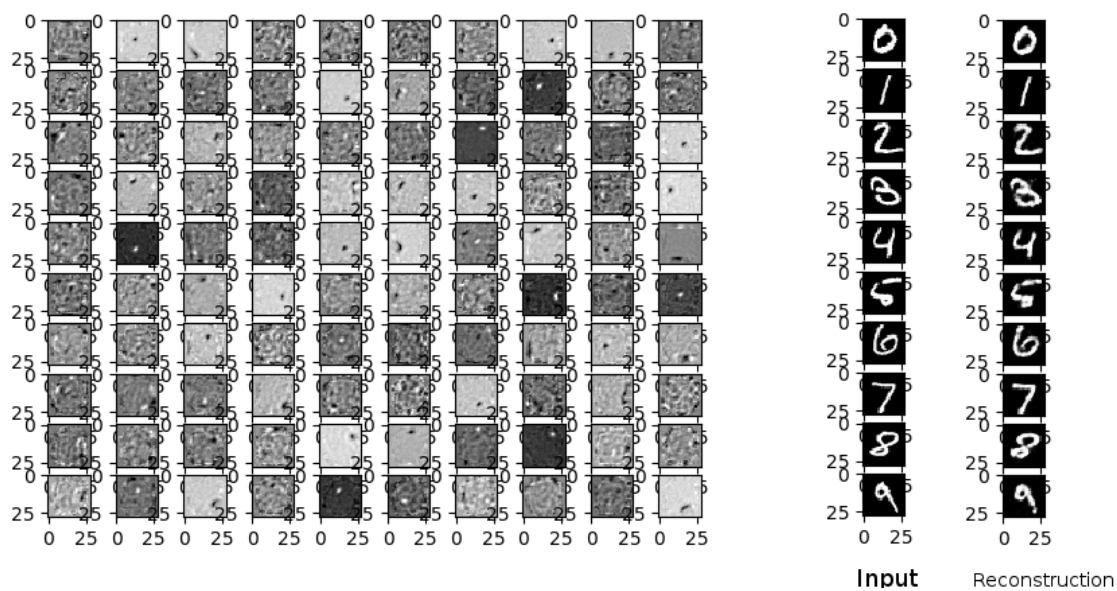
Unlike rbm, there is very limited structure to discern in the visualization of Auto encoder weights. Even so, the auto encoder is able to accomplish a fairly good reconstruction of the input digit.

The standard autoencoder does better than random initialization in terms of validation accuracy. However it still starts to overfit after 10 epochs or so. It doesn't offer the same generalizability as the rbm initialization. (see Figure A)

Figure 1: Figure A: Comparing Generalization for different initialization strategies



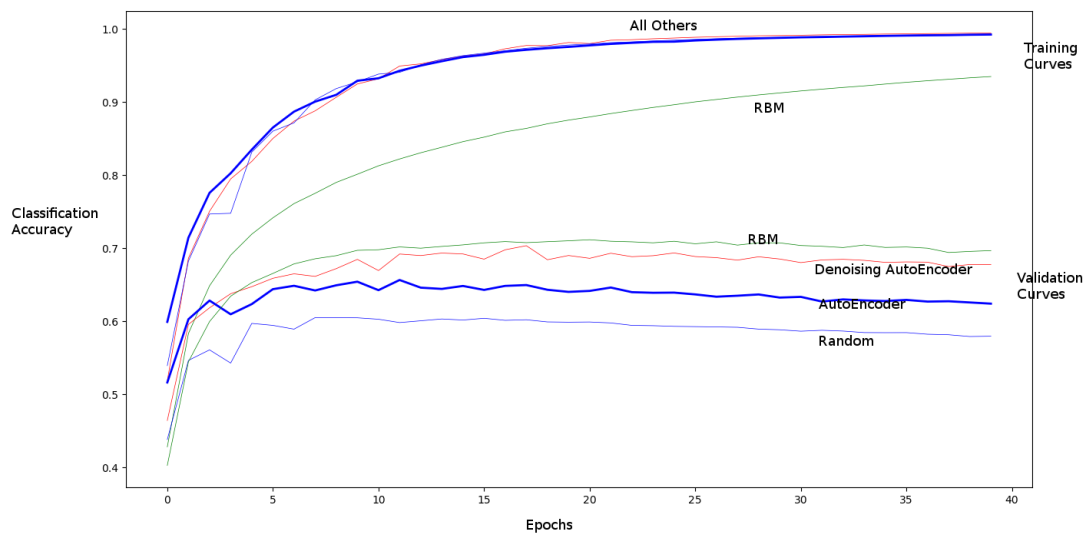
## f) Denoising Autoencoder [10 points]



Compared to no drop out auto encoder, there is greater structure visible in the weight visualization. The reconstructed digits are also sharper compared to the standard autoencoder. The filters are still not as crisp as the rbm case. The denoising autoencoder pre-training is an improvement over both random initialization and standard auto encoder.

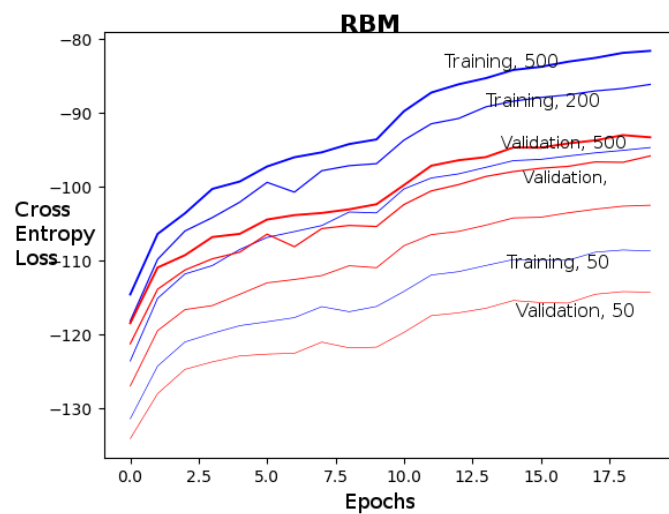


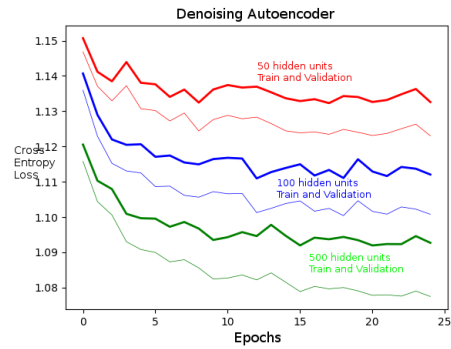
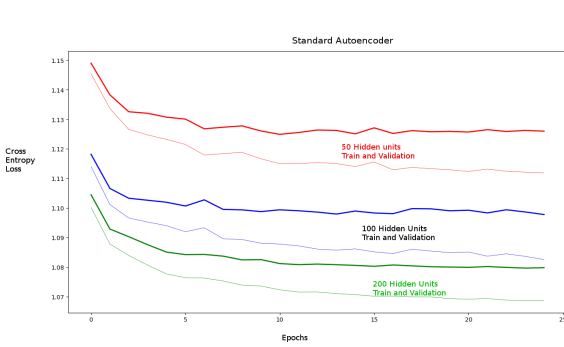
Figure 2: Figure A (repeated): Comparing Generalization for different initialization strategies



Compared to rbm initializations it performs as good or slightly worse on multiple attempts in terms of validation accuracy and generalizability. (see Figure A)

**g) Number of hidden units [10 points]**





(Note: Some of the loss plots are negative log likelihood.)

The larger the hidden units, the better training and validation accuracy across all 3 settings. It follows from intuition that larger hidden units allow preserving more information about the input making the reconstruction easier. Obviously, with larger number of hidden units the gradient descent takes longer to converge given the larger number of parameters to be learnt.