# Cognition and computation:

- Francesco Gabriele
- student identification number: 2073224

NOTE: All the following refers to the script "Script of the assignment", except for point 7, which instead refers to the script "EMNIST-Letters restricted to ten classes".

# 1. INTRODUCTION

## 1.1 GOALS OF PAPER

The goal of the project is to analyze the efficiency of a Deep Belief Network and Feed-Forward Neural Network. Further examination on their resilience to noise in the data and to adversarial attacks will be also carried out.

## 1.2 OVERVIEW OF THE DATASET

The EMNIST-letters consists of a set of images of capital and non-capital letters, made up of 28X28 pixels each.
The dataset is splitted into two subset:
- Training set: 124800 images
- Test set: 20800 images

When displayed at the screen they look like this:



## 1.3 PREPROCESSING:

Each single pixel value falls within [0,255] range where 0 represents "black" and 255 reprepresents "white". Since the model under examination is a DBN these values will be scaled down to the range [0,1].

# 2. UNSUPERVISED LEARNING:

## 2.1 BRIEF DESCRIPTION OF THE MODEL

Deep Belief Networks also referred to with the acronym D.B.N. are generative models and as such they build representations of the data they get as input in an unsupervised learning way.

A DBN is constituted by a stack of Restricted Boltzmann Machines (R.B.M.) which are stochastic recurrent networks; more specifically they are probabilistic graphical models where each neuron activation level represents a probability. The term "restricted" refers to the absence of intralayer connections between neurons.
Each RBM is trained aiming to maximize the likelihood function via an algorithm known as "Contrastive Divergence", consequently the whole DBN is trained in a greedy way.

The representations made by the model are used to classify the original images. This classification is carried out with a linear classifier that encodes the categories trained with supervised learning.

The linear classifier that is fully connected to the deepest hidden layer investigates the task performance whereas read-outs at previous hidden layers investigate how well each of them makes representations. An increasing accuracy across the linear read-outs indicates that the representations get progressively more disentangled.
Ideally, the intention would be to structure a DBN that makes representation of the data it has been fed on in an hierarchical fashion. Therefore the representations outputted by the first RBM are expected to be higher level than the ones produced by the last one.
To facilitate this, the DBN will be structured having an increasing number of neurons at each hidden layer of the RBMs.

## 2.2 SETTING UP THE PARAMETERS

Provided that the DBN under analysis consist of a stack of three RBMs the main parameters to focus on are the:
- the number of neurons of the hidden layer of each RBM: having a network with too many neurons could require a while for the
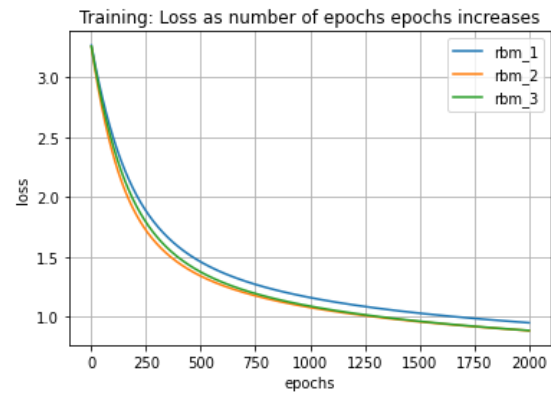
training to be completed and it could also lead to overfitting. Therefore, a reasonable number of neurons is to be chosen.
- the number of epochs the DBN is trained over
- the batch size: larger batch size may speed up the training but usually leads to worse model performance if compared to smaller ones.
- learning rate: if it is too small the training phase is faster but there is the risk of overshooting the minima, if it is too small instead the training is too slow.

The performance of the model highly depends on the structure of the network, thus some reasoning about the matter is required. It can be proved that the same DBN structured with 400 neurons in the hidden layer of the first RBM,500 in the one of the second RBM and 800 neurons in the one of the third RBM reaches a good level of accuracy over the MNIST dataset.
The MNIST dataset, though, consists of ten classes, one for each digit whereas the EMNIST- letters one has more of the double of the classes, therefore a bigger number of neurons might be required to reach a similar performance.

| number of neurons of the first hidden layer | number of neurons of the first hidden layer | number of neurons of the first hidden layer | batch size | epochs DBN | epochs linear read-outs | Accuracy of the first linear classifier | Accuracy of the second linear classifier | Accuracy of the third linear classifier |
|---|---|---|---|---|---|---|---|---|
| 400 | 600 | 900 | 50 | 50 | 2000 | 74% | 75% | 76% |


Training: Loss as number of epochs epochs increases

# 3 FURTHER INVESTIGATION OF THE MODEL
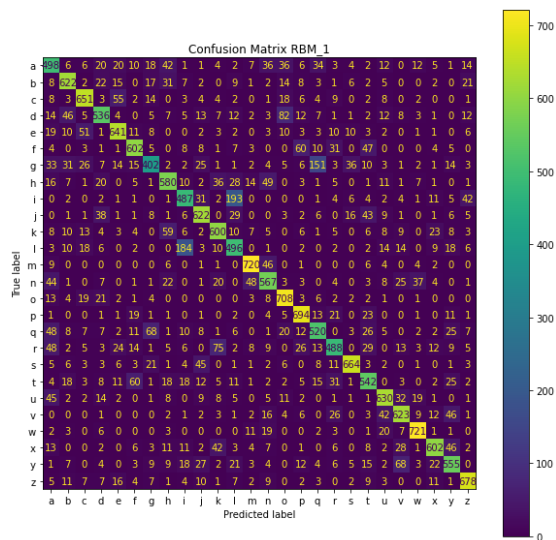
## 3.1 PREDICTION OF EACH LETTER


Prediction Accuracy of each single letter

As it can be clearly seen from the bar chart above, the letter which the model struggles the most with is the "g", while the one with it deals the best is the "w".
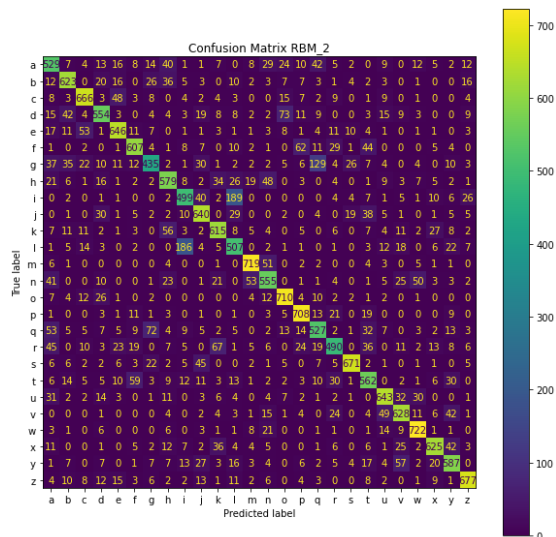
## 3.2 CONFUSION MATRICES
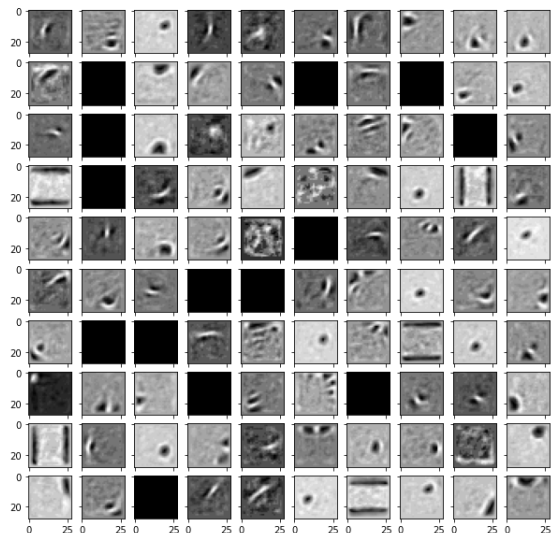An extensive description of the predictions made at test time can be represented in the form of confusion matrices.

Confusion Matrix RBM_1



Confusion Matrix RBM_2



Confusion Matrix RBM_3

DISCUSSION:

The matrices show that "i" and "l" are often mistaken for one another; the same thing happens for the letters "q" and "g".
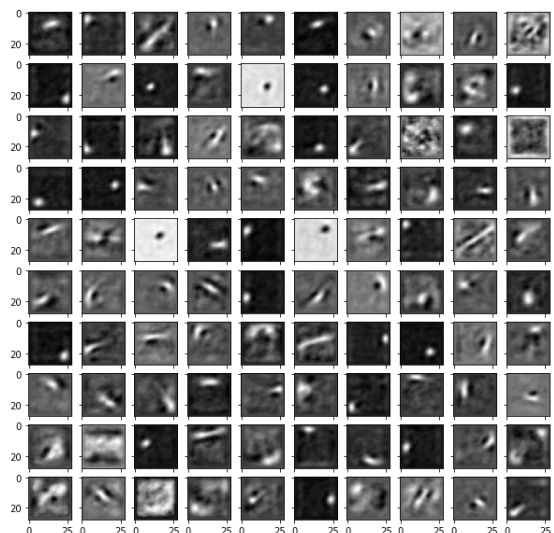
## 3.3 VISUALIZING RECEPTIVE FIELDS:

After the model has been trained it is possible to visualize the learned weights. In the images below some of the weights are plotted in cells of 28x28 pixels so that it will be possible to identify which parts of an image trigger the unit of a specific vector.



Receptive Fields RBM_1



Receptive Fields RBM_2

Receptive Fields RBM_3

In full accordance with what was mentioned before, the representations become better and better.

### 3.4 HIERARCHICAL CLUSTERING

Internal representation properties are explorable through plotting their dendrogram.



Third hidden layer

DISCUSSION:
The dendrogram above shows how similar the letters are to each other according to the DBN. The results seem to be pretty reasonable. For example the letter "o" differs from the "c" just for a small arch that would complete a "c" into an "o" that's why they are close in the dendrogram.

# 4. ROBUSTNESS TO NOISE
### 4.1 RANDOM GAUSSIAN NOISE

Let's now inject some noise in the images and then feed the models with them to examine how resilient the model is to noisy stimuli.
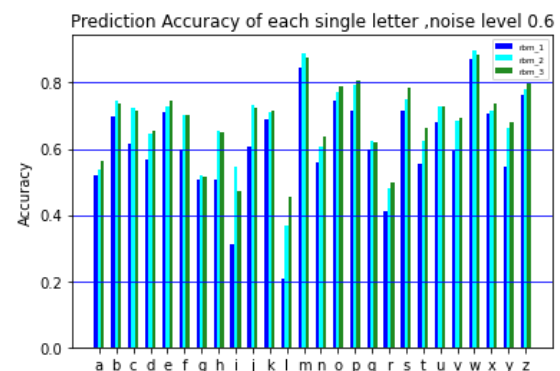
A tensor whose items come from a normal distribution of mean 0 and variance 1 is created, then it is multiplied by the noise level (ε) and finally the perturbed dataset is obtained by adding that tensor to the original data one.

### 4.2 IMPACT OF RANDOM NOISE ON THE IMAGES

In the following a string of letters affected by random gaussian noise (ε = 0.2) with the relative predictions of the DBN is displayed:



### 4.3 PREDICTION OF EACH LETTER



DISCUSSION
After the injection of the noise in the data the model really struggles with the letter "l" whereas "m" and "w" appear to be the ones it predicts more accurately, this might be related to the fact that they are closed in dendrogram.

In addition to that the bar chart highlights really well how the predictions of the first classifier are definitely more imprecise compared to the ones of the the remaining two classifiers, which instead are quite close in terms of right predictions.
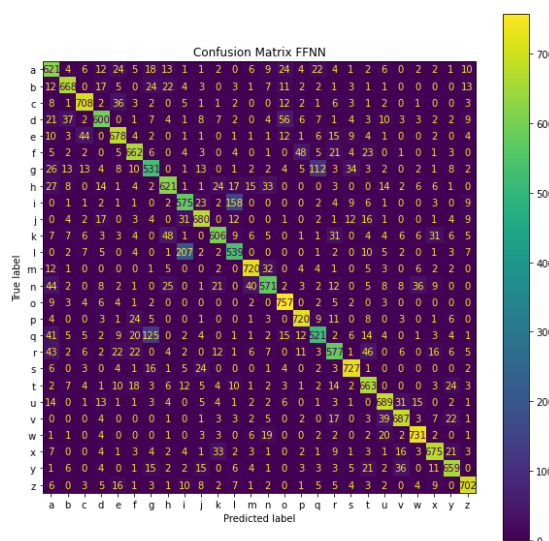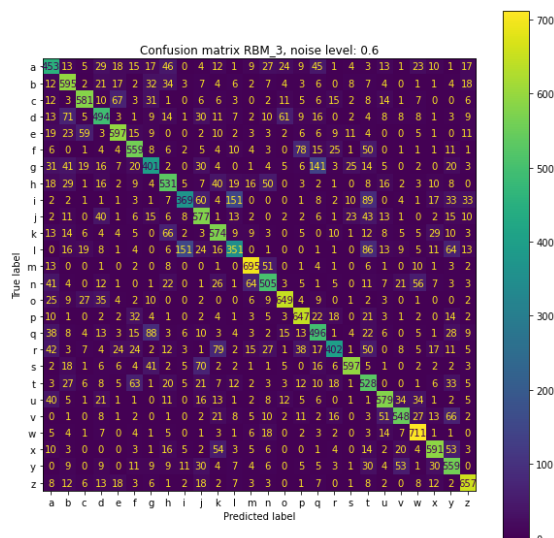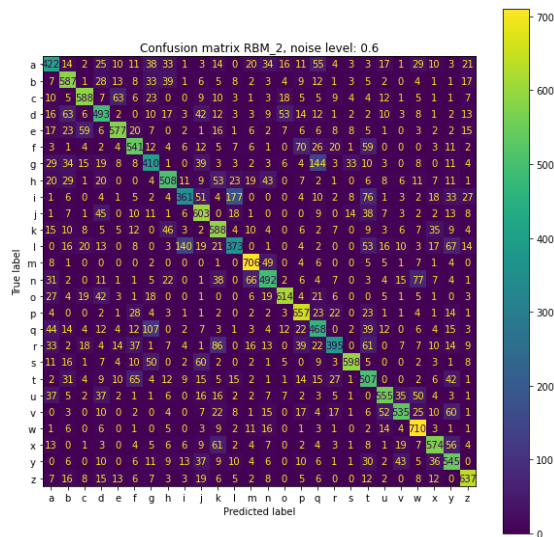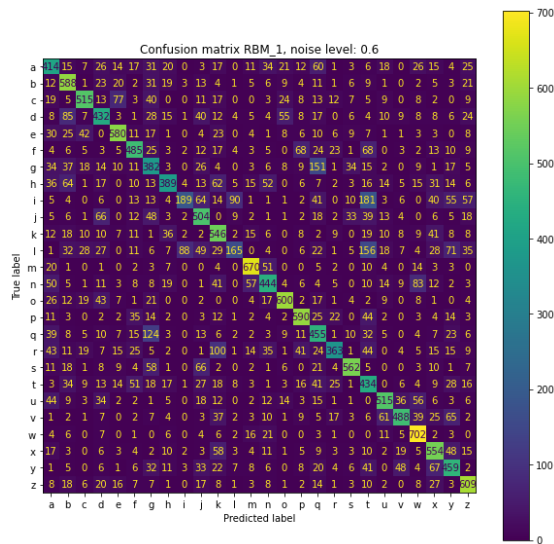
### 4.4 CONFUSION MATRICES

4

Confusion matrix RBM_1, noise level: 0.6


Confusion matrix RBM_2, noise level: 0.6


Confusion matrix RBM_3, noise level: 0.6
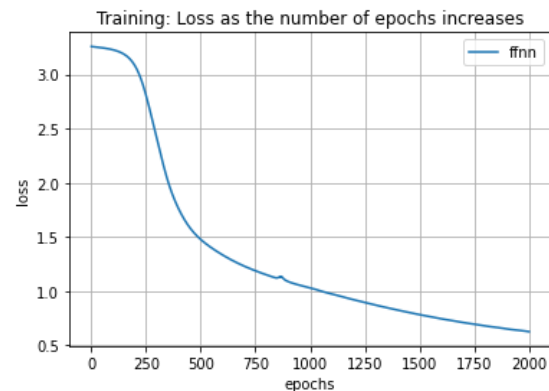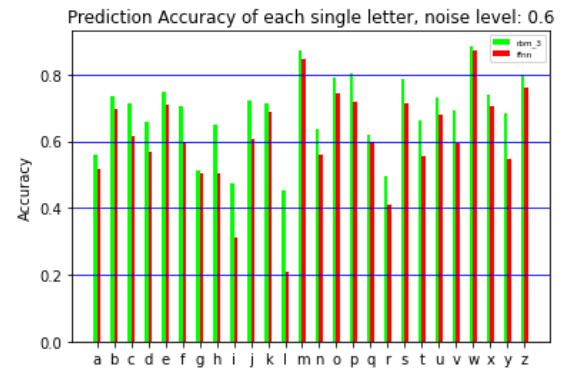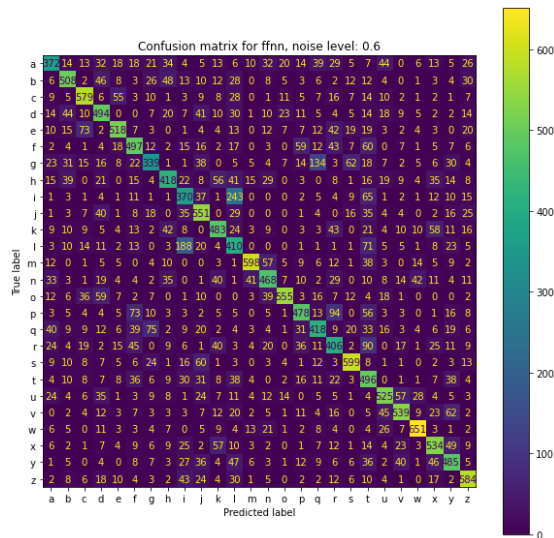
# 5. COMPARISON WITH A FEED-FORWARD NEURAL NETWORK

## 5.1 TRAINING OF FEED FORWARD NEURAL NETWORK

A feed-forward neural network is trained on the same dataset. To make a fair comparison with the DBN, the FFNN is trained with the same amount of epochs and the same size of hidden neurons. Coincely:
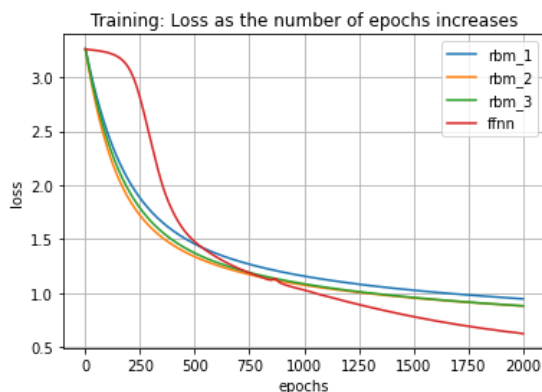
| FEED FORWARD MODEL | | | | |
|---|---|---|---|---|
| first hidden layer size | second hidden layer size | third hidden layer size | epochs | accuracy |
| 400 | 600 | 900 | 2000 | 81% |


Training: Loss as the number of epochs increases


Confusion Matrix FFNN

Confusion matrix for ffnn, noise level: 0.6

## 5.2 COMPARISON BETWEEN DEEPEST DBN READ OUT WITH THE FFNN

A five 5% gap in accuracy separates the FFNN and the DBN, but let's further investigate it.



Training: Loss as the number of epochs increases

The loss relative to the FFNN stalls for the first two hundred epochs and then decreases abruptly until the five hundredth epoch is reached, after that the loss decreases more gradually. The loss curves relative to the DBN instead, go down more smoothly. However around the thousandth epoch, these loss curves start to diverge from the red one.



Prediction Accuracy of each single letter, noise level: 0.6

With a noise random gaussian with ε = 0.6 is injected in the model the DBN perform better than the FFNN, let's make the comparison on a wide range of values of ε

.

## 5.3 COMPARISON BETWEEN EACH RBN AND THE FFNN: PSYCHOMETRIC CURVE



Robustness to noise

The accuracy of the second and third RBMs remains almost the same for values of ε smaller than 0.4, that is also the value of epsilon around which the curve of the FFNN goes below the curves of the last two RBMs. For higher levels of noise the FFNN performance is even worse than the first RBM.

## 6. ADVERSARIAL ATTACKS
### 6.1 DEFINITION OF ADVERSARIAL ATTACK

The noise added to the images in adversarial attacks is specifically crafted to confuse the model.

Random gaussian noise is random by definition so it by no means depends on the gradient of the cost function. Conversely adversarial attacks aim to inject noise in the opposite direction of the descent of the gradient.
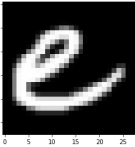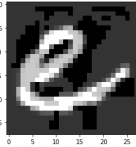The attack is implemented according to the following formula, where "x" represents the data tensor.

$$x = x + \varepsilon\, sign(\nabla x(w, x, y))$$

Adversarial attacks are meant to pass undetected so usually ε is not very large.

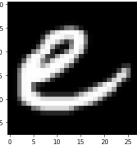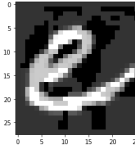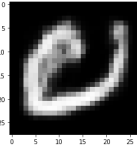## 6.2 VISUALIZE THE IMPACT OF ADVERSARIAL ATTACKS ON THE IMAGES

In the following table an image of the letter "e" is compared to its perturbed and reconstructed versions.

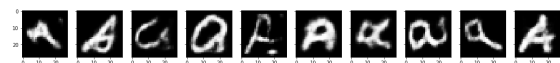| original image | perturbed image |
|---|---|
|  |  |

By taking a close look at the image in the middle it is possible to see some gray in the image due to the noise.

## 6.3 TOP DOWN RECONSTRUCTION AS A WAY TO MITIGATE THE EFFECT OF A.A.

Image reconstruction is a valid fix to adversarial attacks. In the following a top-down reconstruction will be performed.

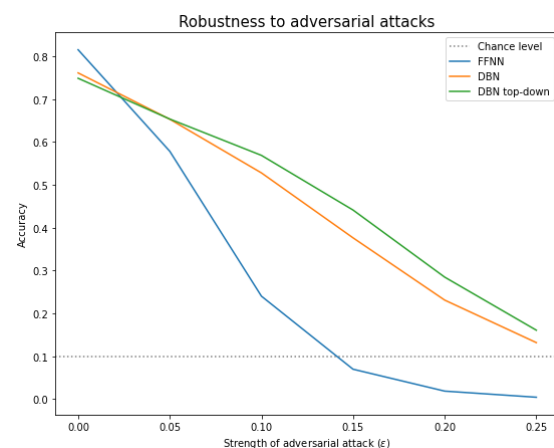| original image | perturbed image | 1 step reconstruction |
|---|---|---|
|  |  |  |

The gray in the third image is almost removed altogether even if there are still some blurry areas of the image, but the reconstruction is pretty good. Let's how good is the reconstruction on the the letter "a":



Apart from the first image the model appears to be quite good at denoising the data, in fact the reconstructions are fairly similar to the original images.

## 6.4 PSYCHOMETRICS CURVES

Let's now explore the robustness of each model to different attack strengths.



As the reconstructed images shown above suggest, the model denoise the images very well and this results in higher accuracy and more resilience to adversarial attacks.
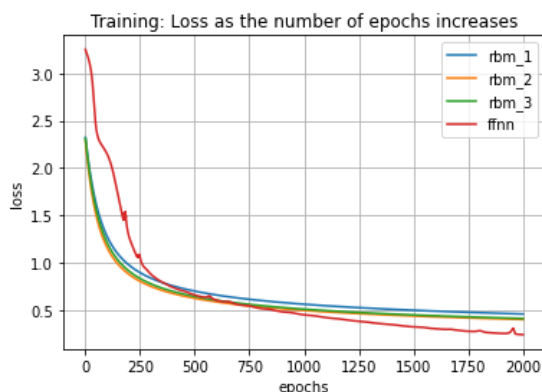Not only the DBN with top-down reconstruction outperforms the FFNN but it also outdoes the DBN without reconstructions.

# 7. ONE LAST NOTE ABOUT THE ACCURACY

It can be pointed out that the accuracy of the two models on the EMNIST-letters is not that great especially if compared to the model performance on MNIST dataset where the accuracy is found to be well above the 90%.

To make a fair comparison let's take into consideration a subset of the classes of the EMNIST- letters dataset. More specifically, let's restrict the classes to the first ten letters of the alphabet and see what happens. All the other parameters are left unaltered.

| EMNIST-letters classes restricted to the first ten letters of the alphabet | |
|---|---|
| First read-out accuracy | 86% |
| Second read-out accuracy | 87% |
| Third read-out accuracy | 88% |
| ffnn read out | 91% |


Training: Loss as the number of epochs increases

As expected the levels of accuracy of the model working on a smaller number of classes are comparable to the one recorded with MNIST dataset.

Therefore the low level of accuracy seems to be related to the number of classes.

## REFERENCE PAPERS

- [Training restricted Boltzmann machines: An introduction - ScienceDirect](#)
- [Frontiers | Deep Unsupervised Learning on a Desktop PC: A Primer for Cognitive Scientists (frontiersin.org)](#)
- [Reducing the Dimensionality of Data with Neural Networks | Science](#)
- [guideTR.pdf (toronto.edu)](#)