

# Deep Learning with Stacked AEs & RBMs

## DD2437 - Artificial Neural Networks & Deep Architectures - Lab 4

Niels Agerskov	Lukas Bjarre	Gabriel Carrizo
agerskov@kth.se	lbjarre@kth.se	gabcar@kth.se

---

## 1 Introduction

This lab will examine two different artificial neural network structures, Auto Encoders (AE) and Restricted Boltzmann Machines (RBM). Their effectiveness in a learning task and the effect of the layer depth of the models will be tested and evaluated.

### 1.1 Used libraries

All of the models were implemented in Python. For the AE models the Keras framework [1] was used, which in turn is built upon Tensorflow. Keras did not have any RBMs though, so for the RBMs and DBNs Scikit-learn [2] was used.

## 2 Feature learning

In this first task shallow versions of both models are trained as benchmarks for the later deeper versions. The dataset used is a subset of the MNIST dataset containing  $28 \times 28$  images of handwritten digits from 0 to 9 together with correct labels of the written digit. All the pixel values has for simplicity's sake been converted to binary values via simple thresholding. A total of 10000 images are used from the dataset, which has been further subdivided into a training set of size 8000 and a validation set of size 2000.

### 2.1 Hidden unit size

The input size hyperparameter for both of the models is decided by the shape of the data. In our case we require  $28 \times 28 = 784$  input nodes, one for each image pixel. We do however have a choice in the number of hidden units,  $n_h$ .

Both models were trained with  $n_h = 50, 75, 100, 150$  hidden units. The error curves on the validation set during the training are displayed in fig. 1 for the RBM, and in fig. 2 for the AE using Stochastic Gradient Descent (SGD). However, the AE clearly converges to the same values no matter the number of hidden units. Complementary to using SGD to train the AE ADADELTA [3] was also used, which error curve can be seen in fig. 3. ADADELTA shows a similar improvement given more hidden units as the error curves for the RBM. The errors are also on one order of magnitude smaller compared to the errors using SGD, which is why the ADADELTA trained AE is only used furthermore.

The quality of the models can be seen in fig. 4, where one image of each class have been used to get both models recalled versions.

### 2.2 Learned features

The learned weights of each model can be examined to get an idea what and how the models are learning. By reshaping the weight vectors back into  $28 \times 28$  grids each hidden units weights can be represented as images where each pixel value corresponds to the strength of that weight from the given pixel to the hidden unit. Plots for these are found in fig. 5. The AE models seems to in general have better trained hidden nodes as more of the components seems to have learned a specific shape instead of something that resembles white noise.

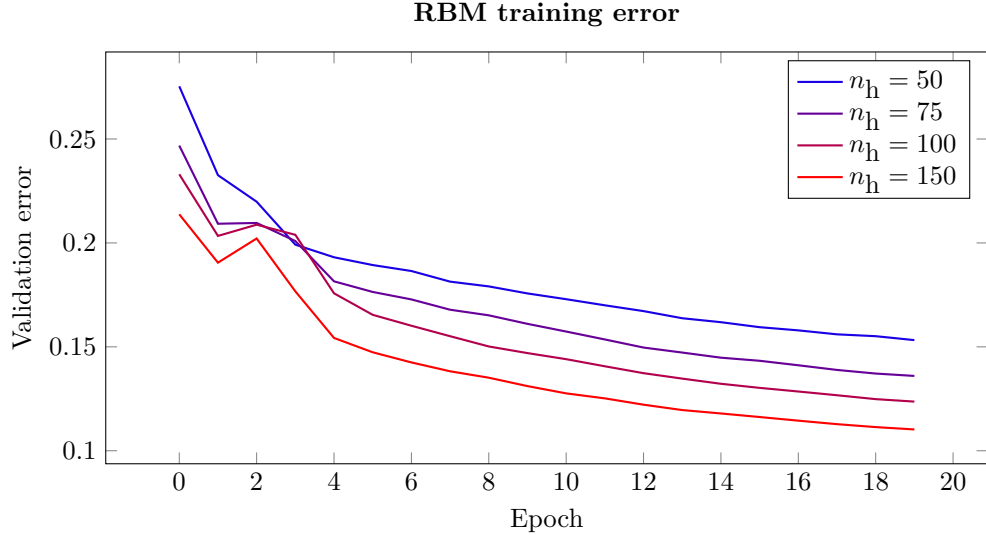


Figure 1: Error curves on the validation set for the RBM.

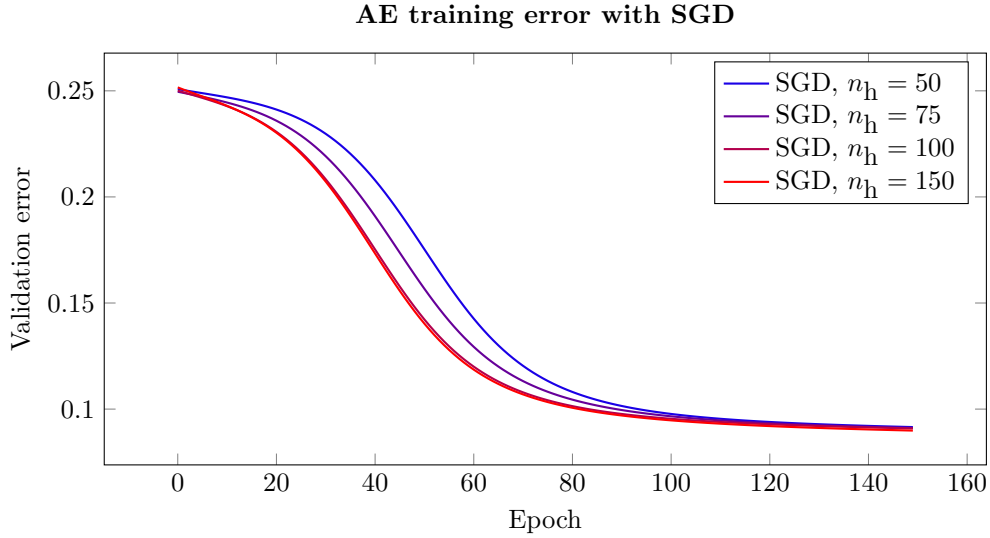


Figure 2: Error curves on the validation set for the AE using SGD.

### 3 Deep architectures

The next step is to make the architectures deep by introducing more hidden layers. For the Auto Encoders this corresponds to adding more encoding layers after the first encoder layer, resulting in a structure called Stacked Auto Encoders (SAE). The deeper version of the RBM is also constructed by adding more hidden layers after the first, resulting in a Deep Belief Network (DBN).

The task for the deep architectures was to classify a given image to the corresponding digit. To achieve this one additional perceptron layer was added to the end of all the deep models. The idea of the structure then becomes that the deep layers are supposed to find and model hierarchical representations of the images in lower and lower dimensions, which the final perceptron layer can use to classify the digit.

The method used to train the networks was to pre-train the hidden layers one-by-one using the unsupervised methods used previously. After one layer was trained the dataset was transformed using the trained weights and then the next layer used this transformed set to train its weights. One final supervised training run on the entire architecture was also used at the end to fine-tune the weights.

Models with different layer depths ranging from one to three were trained and evaluated. The size

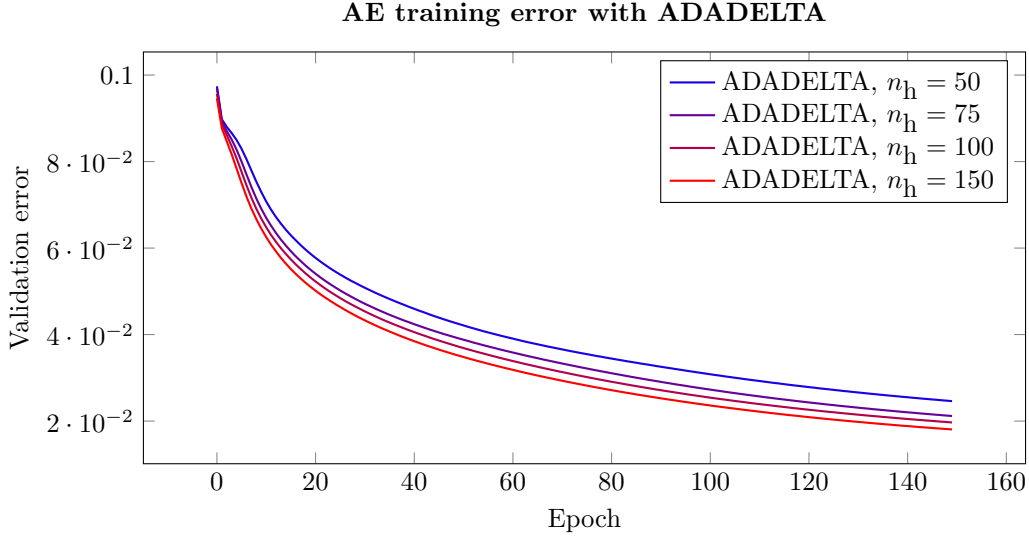


Figure 3: Error curves on the validation set for the AE using ADADELTA.

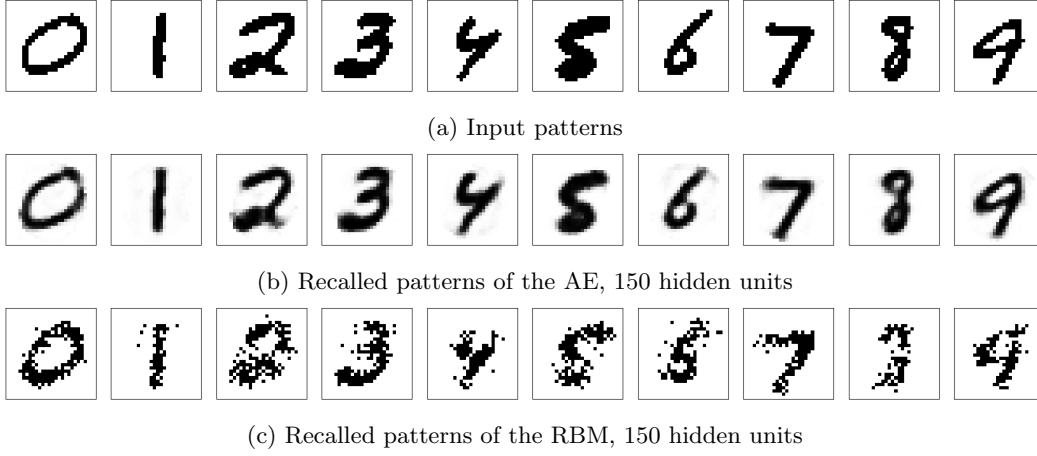


Figure 4: Recalled images from the models.

of each hidden layer was determined first from the best performing versions of the shallow models then decreasing this size for the consecutive layers. The best hidden unit size can be seen from fig. 1 and fig. 3 which in both models is  $n_h = 150$ . This was set to be the first hidden layer size. To decrease the consecutive layer sizes the two other sizes was set to 100 and 50.

### 3.1 Performance

Table 1: Accuracy of different deep architectures

Layers	DBN	SAE
0	0.902	
1	0.937	0.937
2	0.907	0.943
3	0.859	0.866

In table 1 the accuracy of the trained models on the validation set is presented. The 0 layer size represents a simple perceptron layer trained directly on the dataset. We can see that the performance in general is similar for the different setups, with the difference of seeing a slight decrease in performance

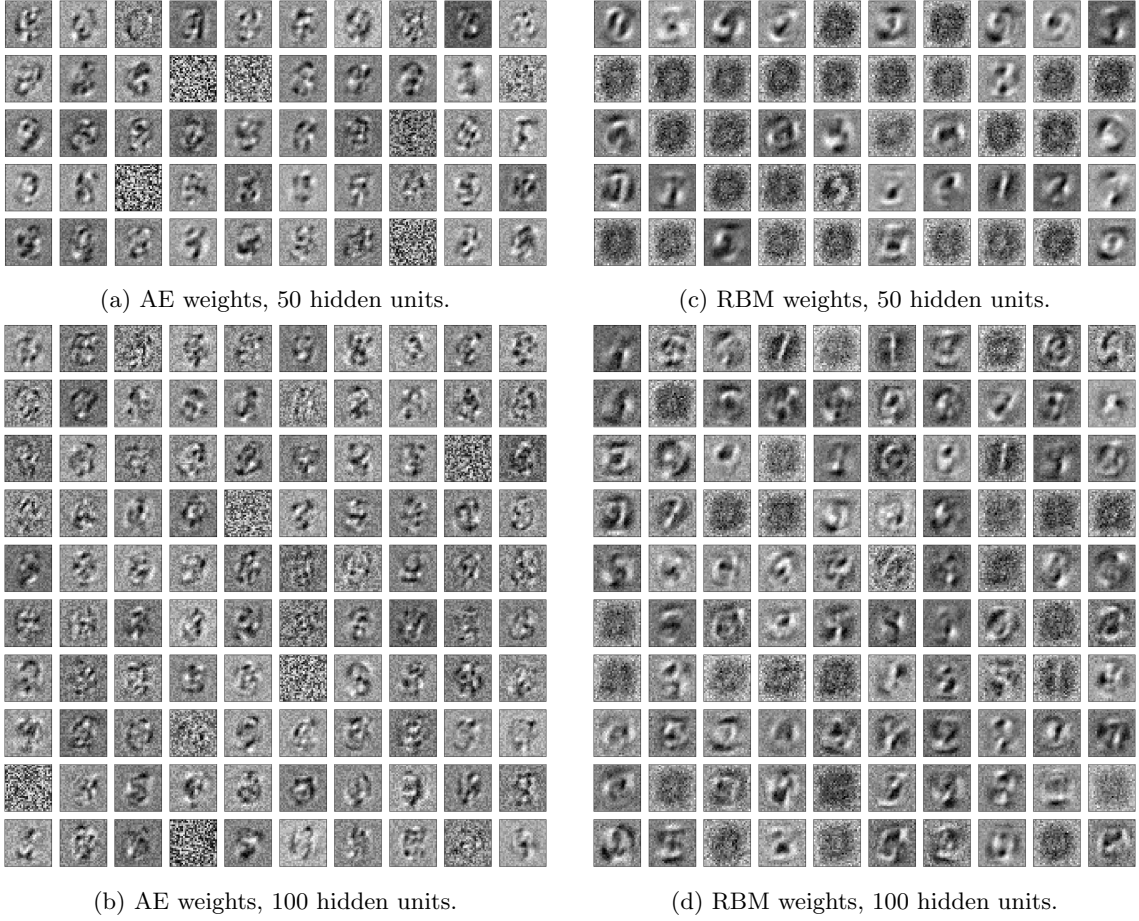


Figure 5: Image representations of the hidden weights in the AE and RBM models for 50 and 100 hidden units.

for the deepest setups. Overall the best performance was seen with the Stacked Auto Encoder with two hidden layers.

### 3.2 Learned feature

Similarly to how the learned features was studied earlier by looking at the learned weights the weights between the different hidden layers in the deep architectures can also be viewed. Figure 6 plots the weights between deep hidden layers into images. However, unlike the weights of the first layers the hidden units represent a transform from one abstract representation to another which is not easily represented in the original  $28 \times 28$  pixel space. The visualisation is still done by pixels, where each pixel represents the weight from one previous layer unit to the represented hidden unit, but since the previous data representation does not necessarily makes sence in a pixel space no special patterns can be seen. However, to the neural network this might still represent something useful even though we humans cannot observe anything useful from it.

Another way to visualise the learned features is to observe the activations of the hidden units when entering input patterns from different classes. Figure 7 plots these activations for both models when images of the classes 0 and 8 are entered into the classifier. The visualisation is once again done by images, each pixel represents one hidden unit and the different images represents different layers. The problem of the abstract feature space is once again present since the pixel arrangement does not reflect any meaningful pattern in the representation. The more important part to observe in these pictures however is that the activation for different image classes produce different patterns in the representation space, which is what the final perceptron layer is then using to classify the digits.

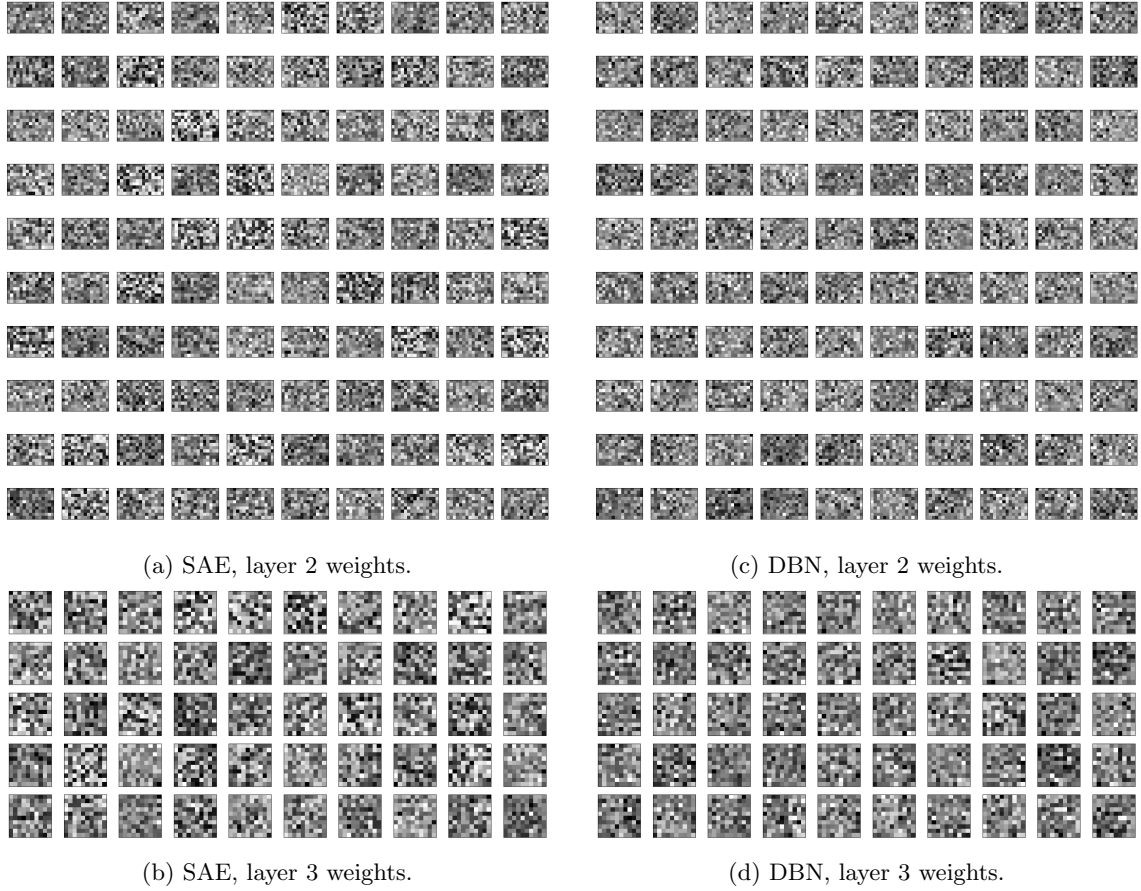


Figure 6: Image representations of the hidden weights in the SAE and DBN networks.

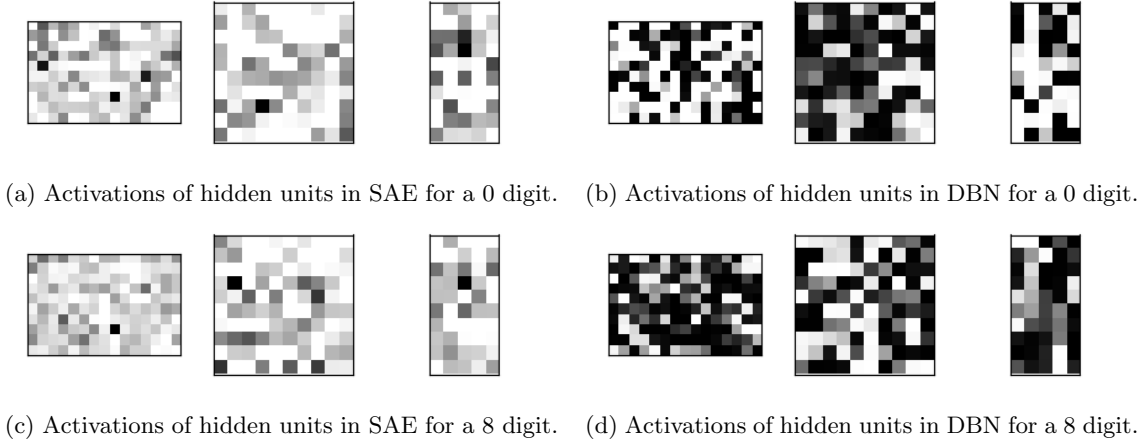


Figure 7: Image representations of the activations in each hidden layers in the deep models when fed with an image of a 0 and an 8.

### 3.3 Comparison to deep MLP

In addition to the pre-trained deep models a comparison to a equivalent deep MLP with the same hidden layer setup but trained in a regular supervised way. Since the best performing pre-trained network setup was the SAE with two layers, first layer with 150 hidden units and second layer with 100 hidden units, this setup was used for the supervised MLP. Since two different optimizer was used for the AE in the start, SGD and ADADELTA, both of these optimizers were tested for the supervised

training. Table 2 shows the validation set accuracy for the final trained networks. The supervised trained network with ADADELTA performs extremely similar to the pre-trained network, but using regular SGD shows a much worse performing network.

Table 2: Accuracy of deep pre-trained and regular supervised networks

Network	SAE	MLP, SGD	MLP, ADADELTA
Accuracy	0.943	0.599	0.942

## References

- [1] François Chollet et al. Keras. <https://keras.io>, 2015.
- [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [3] Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.