

Classification with Restricted Boltzmann Machines

Projects in Machine Learning and AI

Fritjof Wolf
Katarzyna Tarnowska

Technische Universität Berlin

2015

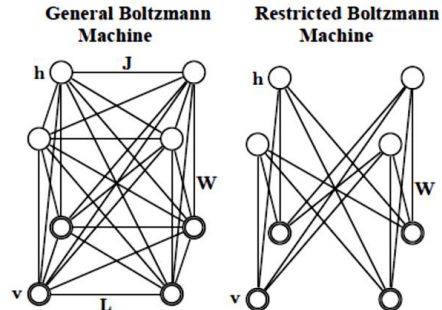
Table of Contents

- 1 Theory
 - Boltzmann Machines
 - Restricted Boltzmann Machines
 - Contrastive Divergence
 - RBM for classification
- 2 Implementation
- 3 Results
- 4 Further work



Boltzman Machine and Restricted Boltzmann Machine

- Recurrent neural network
- Hidden layer and visible layer
- Symmetric weights
- Stochastic binary neurons
- Generative Model
- In a Restricted Boltzmann Machine the joints between hidden units and also between visible units are disconnected



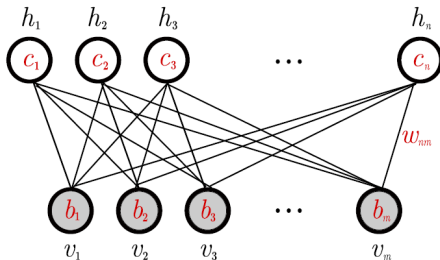
Boltzman Machine

- Energy function depends on model parameter
- Probability depends on weights and state of the other neurons
- Unsupervised learning
- Used to model probability distribution:
 - Apply random input
 - Run the model for some time to generate sample from learned distribution
- First used as an feature extractor



Restricted Boltzmann Machines

- Complete bipartite graph
- Stochastic neural network:
 - nodes - neurons
 - edges - synaptic connections



Source: A.Fischer, Ch.Igel: Training Restricted Boltzmann Machines: An Introduction

Mathematical description of the model

Energy function

$$E(\mathbf{v}, \mathbf{h}) = \sum_{i=1}^V \frac{(\mathbf{v}_i - \mathbf{b}_i^v)^2}{2\sigma_i^2} - \sum_{j=1}^H \mathbf{b}_j^h \mathbf{h}_j - \sum_{i=1}^V \sum_{j=1}^H \frac{\mathbf{v}_i}{\sigma_i} \mathbf{h}_j \mathbf{w}_{ij} \quad (1)$$

Probability of (\mathbf{v}, \mathbf{h})

$$\mathbf{p}(\mathbf{v}, \mathbf{h}) = \frac{\mathbf{e}^{-E(\mathbf{v}, \mathbf{h})}}{\sum_{\mathbf{x}} \sum_{\mathbf{k}} \mathbf{e}^{-E(\mathbf{x}, \mathbf{k})}} \quad (2)$$

Conditional distributions

$$\mathbf{p}(\mathbf{h}|\mathbf{v}) = \sum_{\mathbf{i}} \mathbf{p}(\mathbf{h}_i|\mathbf{v}) \quad (3)$$

$$\mathbf{p}(\mathbf{h}_j = 1|\mathbf{v}) = \text{sigm}(\mathbf{c}_j + \sum_{\mathbf{i}} \mathbf{W}_{ji} \mathbf{x}_i) \quad (4)$$



Contrastive Divergence

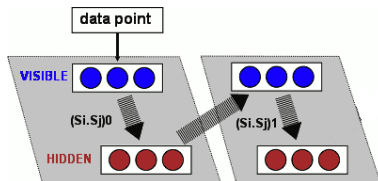
- Problem: Log likelihood gradient is hard to compute
- Run Markov chain to approximate the model distribution
- One step of Gibbs Sampling is sufficient

Remark

Training a RBM is performed by algorithm known as "Contrastive Divergence Learning"

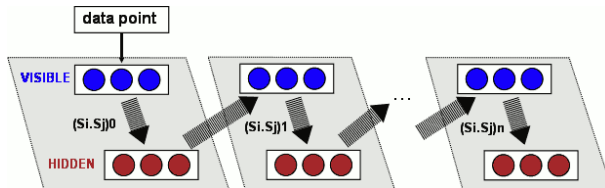


Contrastive Divergence - reconstruction step



- Get one data point from data set
- Use values of the data to set state of visible units - S_i
- Compute S_j for each hidden neuron based on S_i
- Compute $(S_i.S_j)_0$
- Reconstruction: on visible units compute S_i using the S_j
- Compute state of hidden neurons S_j again using S_i
- Use S_i and S_j to compute $(S_i.S_j)_1$

Contrastive Divergence in n steps - whole algorithm



- For each data point in data set:
 - perform reconstruction in n-steps
 - Accumulate $CD_{pos} = CD_{pos} + (S_i.S_j)_0$
 - Accumulate $CD_{neg} = CD_{neg} + (S_i.S_j)_n$
- Compute average CD_{pos} and CD_{neg} (divide by nr of points)
- Compute gradient $CD = CD_{pos} - CD_{neg}$
- Update weights and biases $W'' = W + \alpha * CD$
- Repeat for whole dataset in number of epochs (iterations)

Using RBM for classification

Three approaches (Hinton):

- Use the hidden features learned by the RBM as the inputs for some standard discriminative method
- Train a separate RBM on each class
- Train a joint density model using a single RBM (that has two sets of visible units - y for label and x for data)

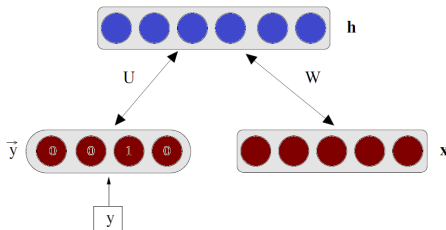


Table of Contents

- 1 Theory
 - Boltzmann Machines
 - Restricted Boltzmann Machines
 - Contrastive Divergence
 - RBM for classification
- 2 Implementation
- 3 Results
- 4 Further work



Tools



GitHub



Data loading and preprocessing module

8 2 9 4 4 6 4 9 7 0 9 2 9 5 1 5 9 1 0 3
1 3 5 9 1 7 6 2 8 2 2 5 0 7 4 9 7 8 3 2
1 1 8 3 6 1 0 3 1 0 0 1 1 2 7 3 0 4 6 5
2 6 4 7 1 8 9 9 3 0 7 1 0 2 0 3 5 4 6 5

- MNIST - handwritten digit images
- Raw data consists of greyscale normalized images (28x28 pixels, pixel is number 0-255) and their labels (0..9)
- Dataset divided into training (50000), validation (10000) and test (10000) subsets
- Loaded optionally from cPickle file or CSV
- Data-specific, help functions implemented (binarization, transformations, scaling, visualizations)



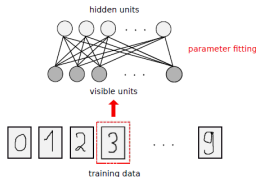
Generative and Discriminative models of RBM

RBM	
- NumOfVisibleUnits	: int
- NumOfHiddenUnits	: int
- VisibleBiases	: float
- HiddenBiases	: float
- Weights	: float
+ train ()	
+ sample ()	

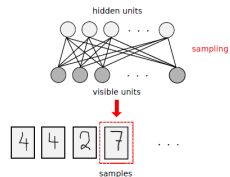


DRBM	
- NumOfTargetUnits	: int
- TargetBiases	: float
- WeightsVH	: float
- WeightsTH	: float
+ <<Override>> train ()	
+ <<Override>> sample ()	
+ predict ()	

learning

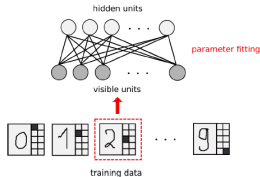


generating

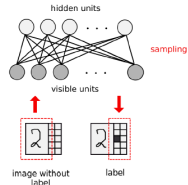


Source: A.Fischer, Ch.Igel: Training Restricted Boltzmann Machines: An Introduction

learning with labels

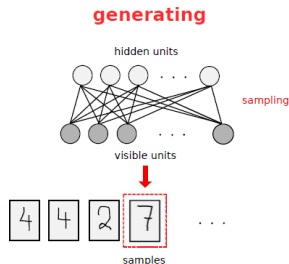
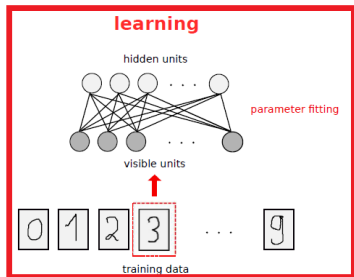


classification



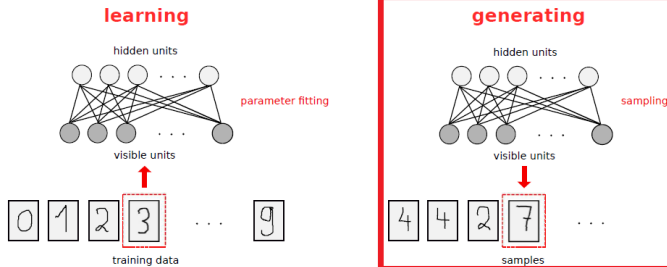
Source: A.Fischer, Ch.Igel: Training Restricted Boltzmann Machines: An Introduction

Generative model - train()



- Fitting RBM parameters so that to model distribution of the training data
- Iteratively performs one step of Contrastive Divergence (using Gibbs sampling) on data subset of one-class
- Learns until specified error threshold between data and reconstruction is reached

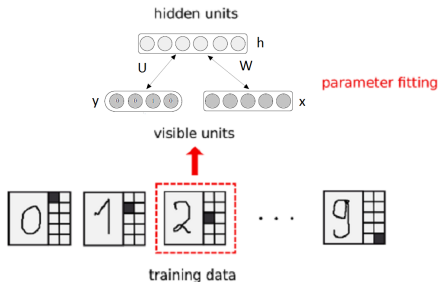
Generative model - sample()



- Trained RBM used to generate samples from learned distribution
- Shows reconstructed image for the specified digit

Discriminative model - train()

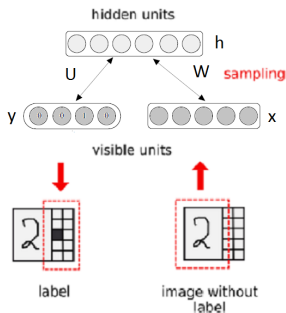
learning with labels



- DRBM models a joint distribution of inputs (x) and target classes (y)
- Two sets of visible units and two weight matrices: between x and h (W) and between y and h (U)
- Train() performs n -step Contrastive Divergence for a mini-batch

Discriminative model - predict()

classification



- Fix the visible variables corresponding to the image
- Sample target variables corresponding to the labels in chosen number of iterations
- For each datapoint in testset return probabilities of each class
- Choose the label class with highest probability
- Compare with original labels
- Count wrong predictions and compute accuracy

Table of Contents

- 1 Theory
 - Boltzmann Machines
 - Restricted Boltzmann Machines
 - Contrastive Divergence
 - RBM for classification
- 2 Implementation
- 3 Results
- 4 Further work



Testing methodology and assumptions

- Reducing to binary problem (binarization threshold = 0.5)
- Parameters possible to test:
 - size of training set,
 - size of test set,
 - learning rate,
 - initial weight distribution,
 - momentum,
 - l2 penaltization,
 - number of steps for contrastive divergence,
 - size of hidden units,
 - number of epochs for training,
 - number of iterations for sampling,
 - error threshold for traning,
 - random state

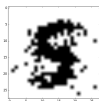


Testing generative model

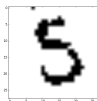
Experiments on different sizes of hidden units



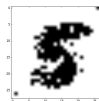
Original



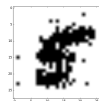
hidden=200



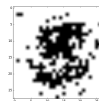
hidden=300



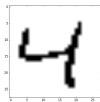
hidden=400



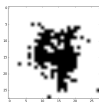
hidden=500



hidden=700



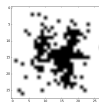
Original



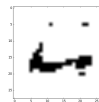
hidden=100



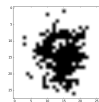
hidden=200



hidden=300



hidden=400



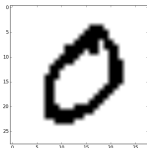
hidden=500

Remark

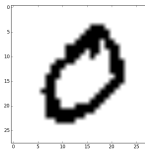
Different digit classes have different optimal hyperparameters

Testing reconstruction with DRBM

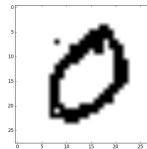
- Learning in mini-batches improved performance
- Momentum parameter for weight update other than 0.0 worsened results
- Optimized results for reconstruction after 500 epochs of training were good:



Original image



momentum=0.0

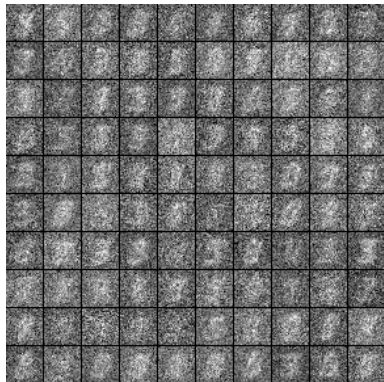
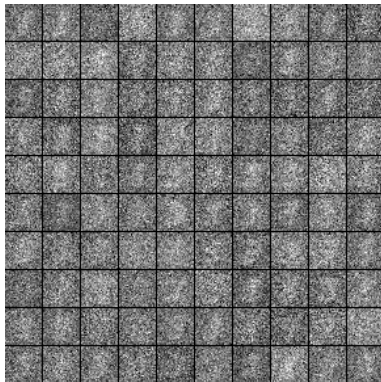


momentum=0.5

- For 500-epoch training MSE falls below 1.0 - in about 30 minutes (on 50 train size)

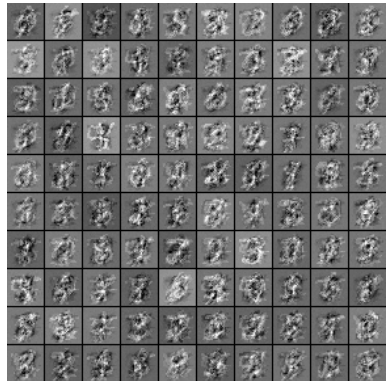
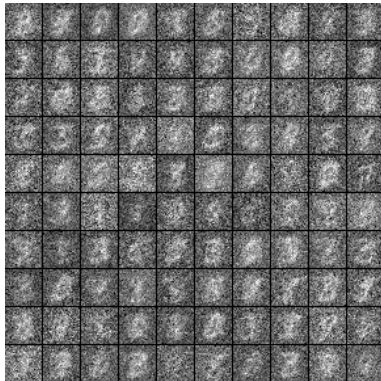


Monitoring progress of learning



Learned weights after 1 and 5 iterations

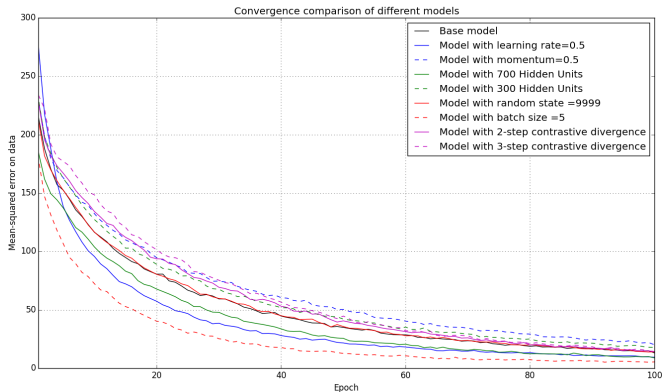
Monitoring progress of learning



Learned weights after 10 and 500 iterations

Monitoring progress of learning

Reconstruction error for 100 epochs



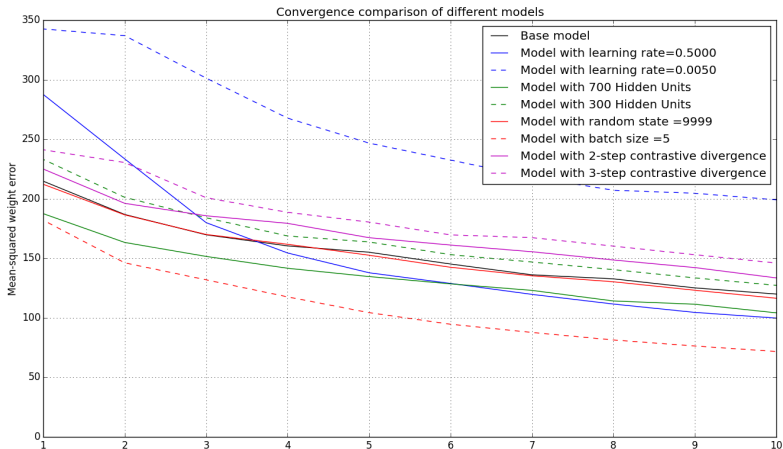
Remark

The reconstruction error on the training set falls rapidly and consistently at the start of learning and then more slowly.

Model selection

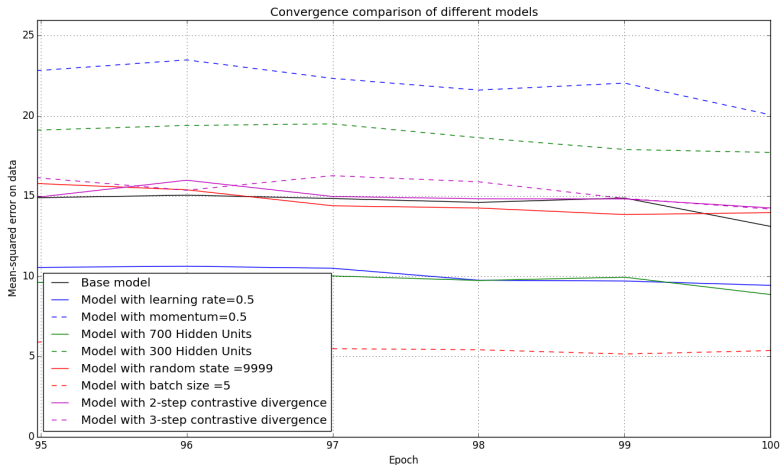
Reconstruction error for first 10 epochs

Base model: lr = 0.01, hidden units=500, random state =1234, batch size=10, 1-step contrastive divergence, no momentum



Model selection

Reconstruction error for last 5 epochs



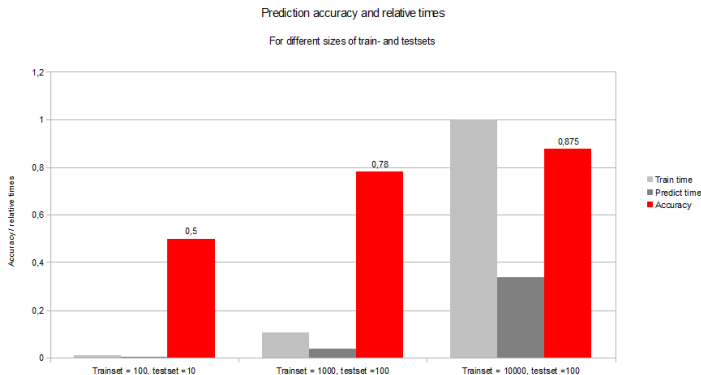
Model selection

Conclusions

- Learning in smaller mini-batches and increasing number of hidden units (from 400 to 700) improved reconstruction
- However, these changes resulted in longer train time
- Higher learning rate (0.5 instead of 0.05) caused reconstruction error to drop more sharply
- For larger train datasets higher learning rate caused instability (after some time of drop error started to increase)
- Different random states do not change reconstruction error significantly
- 1-step contrastive divergence is optimal



Testing RBM for classification (I)



- Classification accuracy depends on are train- and test sets size
- For larger data sets train and prediction times become prohibitive for personal computing

Testing RBM for classification(II)

- Optimal hyperparameters for training phase were chosen
- Performing 100 percent accurate classification on training data could be achieved even on small sets
- Classification accuracy on new data achieved so far is 95 percent(MNIST 50000 trainset and 10000-validation set)
- Better results are expected given greater computing power

Conclusion

Implementation and tests on synthetic data within this project proved that RBMs can be effectively used as standalone classifiers.



Table of Contents

- 1 Theory
 - Boltzmann Machines
 - Restricted Boltzmann Machines
 - Contrastive Divergence
 - RBM for classification
- 2 Implementation
- 3 Results
- 4 Further work



Plans for further work

- Test-against-all-labels prediction approach
- Optimizing algorithms for best performance
- Testing on gaussian values
- Another dataset, possibly CIFAR



Literature I



Hugo Larochelle, Yoshua Bengio.

Classification using Discriminative Restricted Boltzmann Machines.

Proceedings of the 25th International Conference on Machine Learning, 2008.



Geoffrey Hinton.

A Practical Guide to Training Restricted Boltzmann Machines.
UTML TR 2010-003, 2010.



Miguel A. Carreira-Perpin Geoffrey E. Hinton

On Contrastive Divergence Learning

Artificial Intelligence and Statistics, 2005.



Thank you

