

Layer-wise Training in Deep Neural Network

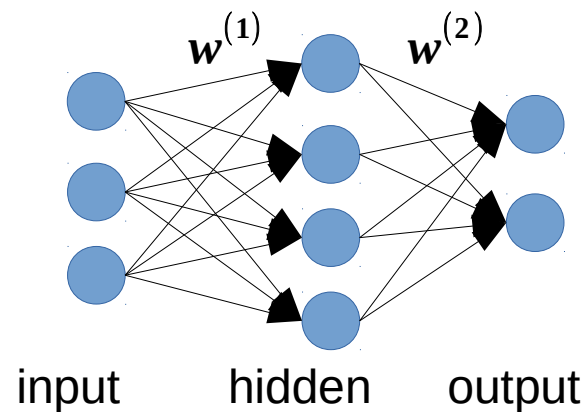
Do Quoc Truong
Nara Institute of Science and Technology (NAIST)
1/6/15

Neural Networks

(Quick overview)

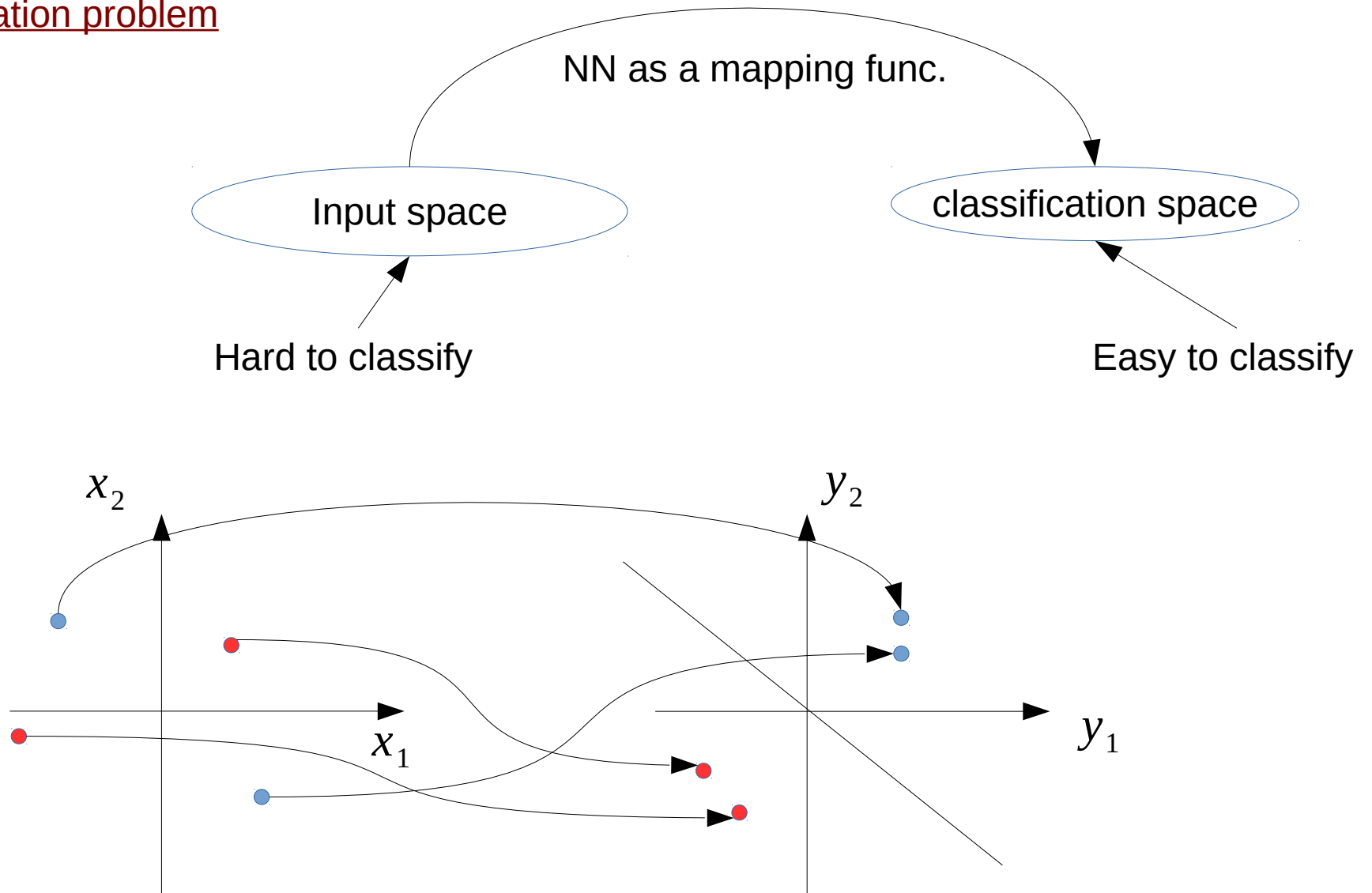
Neural Networks

- A network consists of perceptions and their connection, divided into layers.
- Solve either classification or regression problem



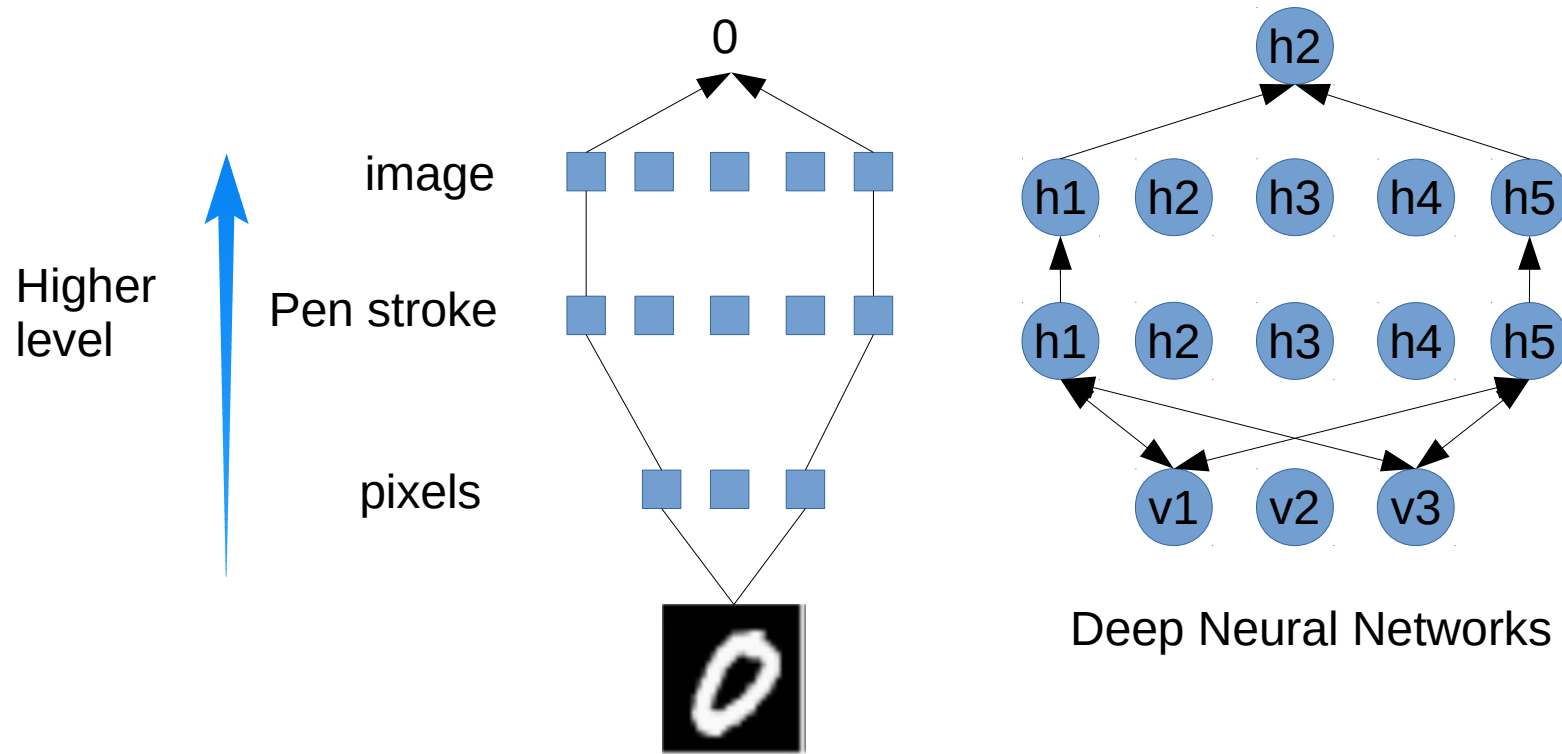
Neural Networks

Classification problem



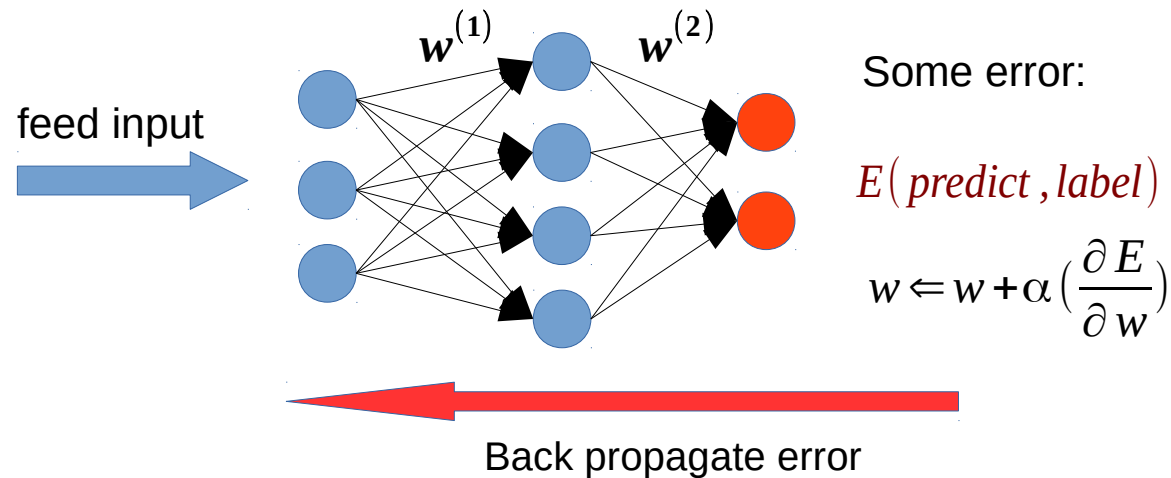
Deep Neural Network

- Deep neural network involve more hidden layers
 - Model non-linear function
 - Learn higher level of representation of data



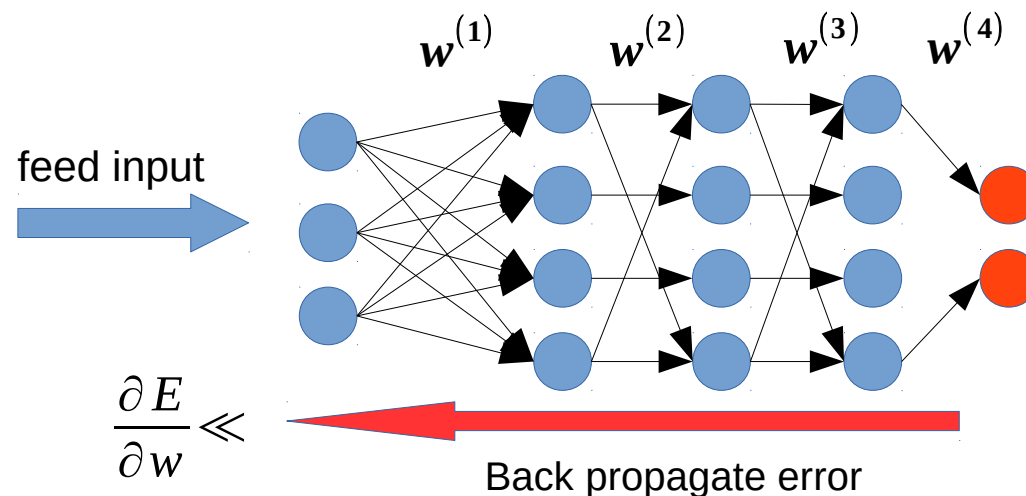
Neural Networks

- Training:
 - stochastic gradient descent + back propagation [Yann LeCun et al. 1989]



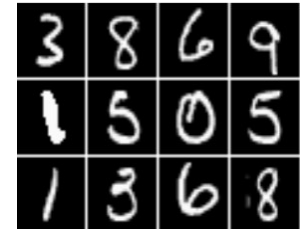
Deep NN

- Trap in local optimal
- Gradient vanishing

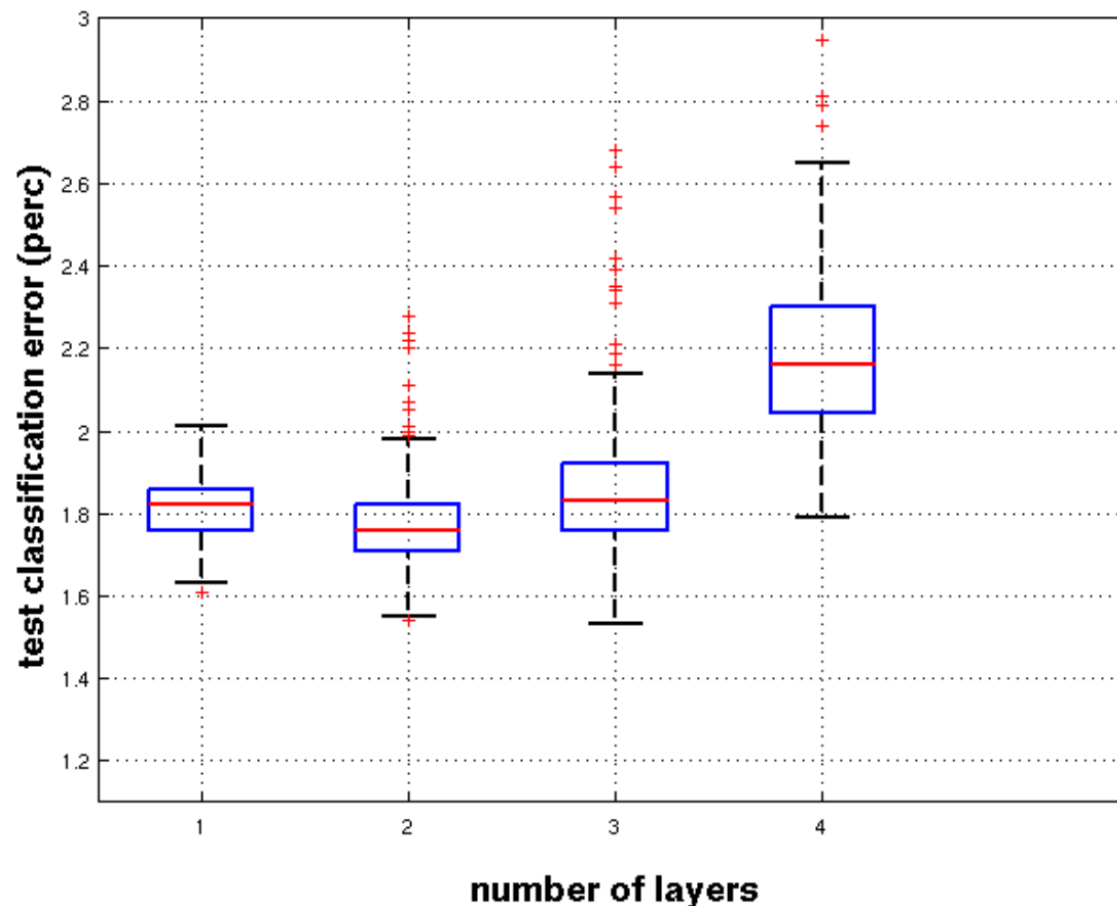


Deep Neural Networks

- Difficulty in training DNN
 - MNIST digit classification task

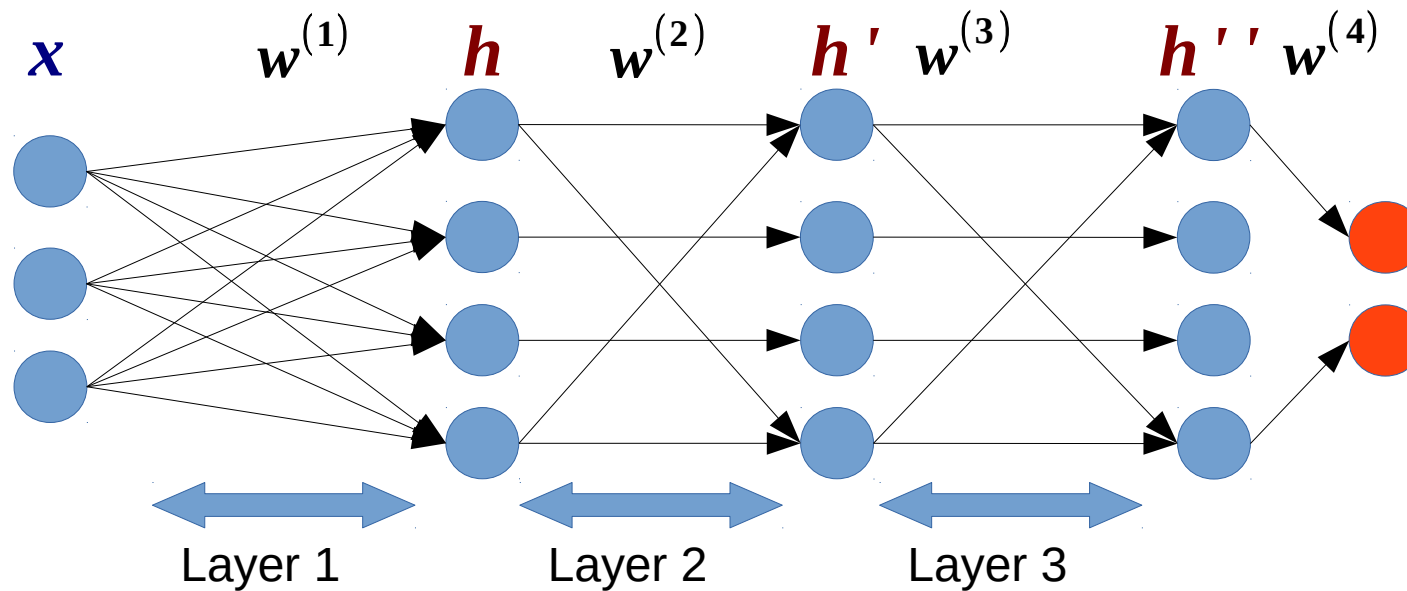


[Erhan et al., 2009]



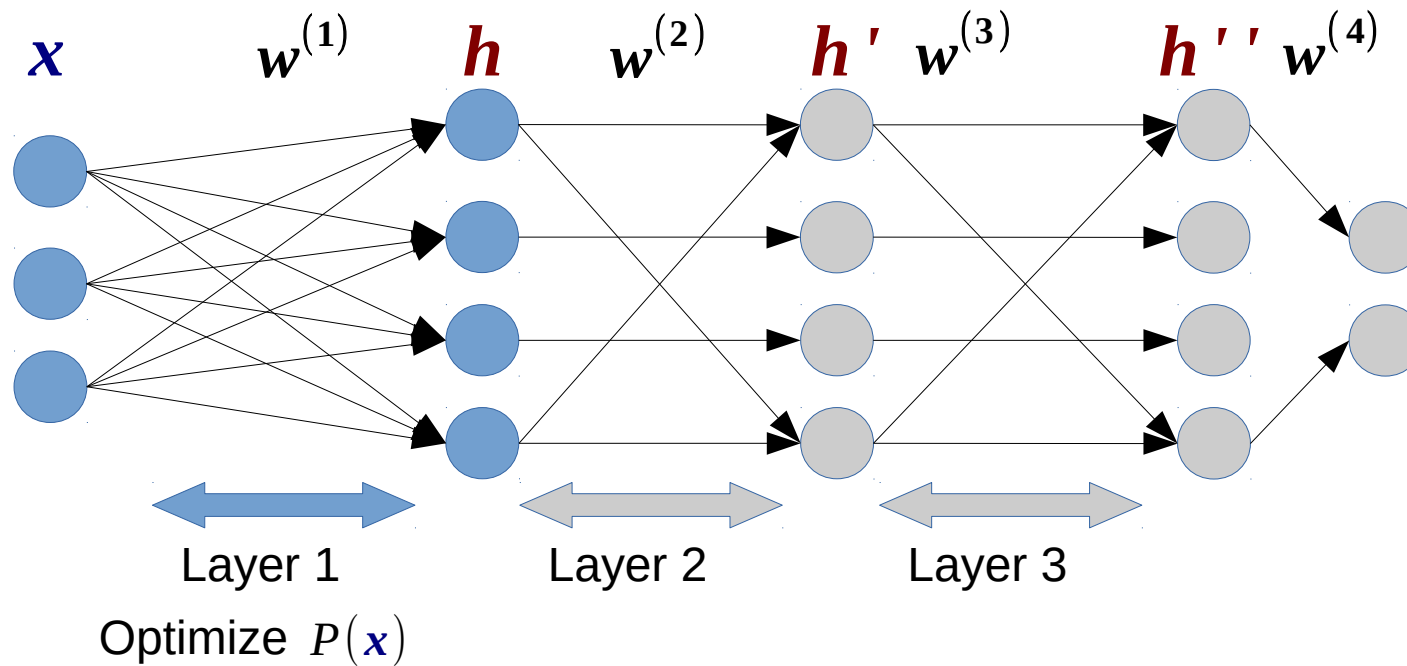
Layer-wise training

Train **one layer** at a time



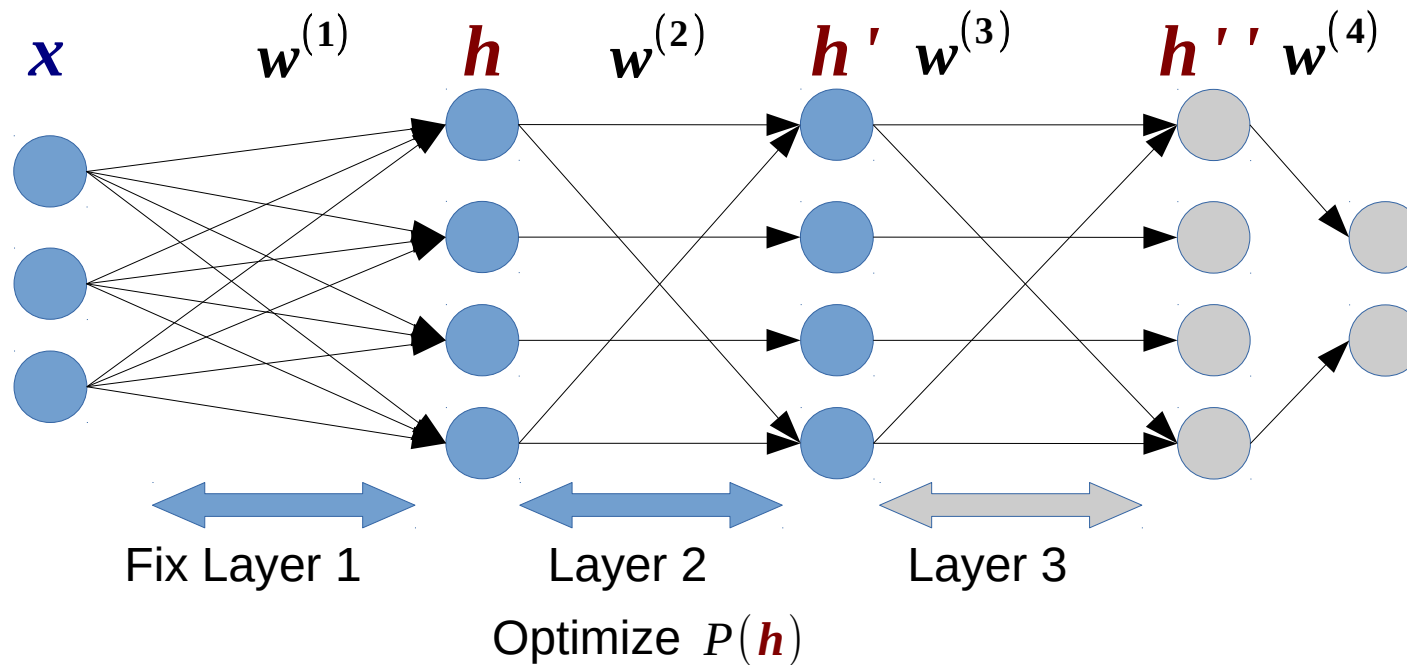
Layer-wise training

Train **one layer** at a time



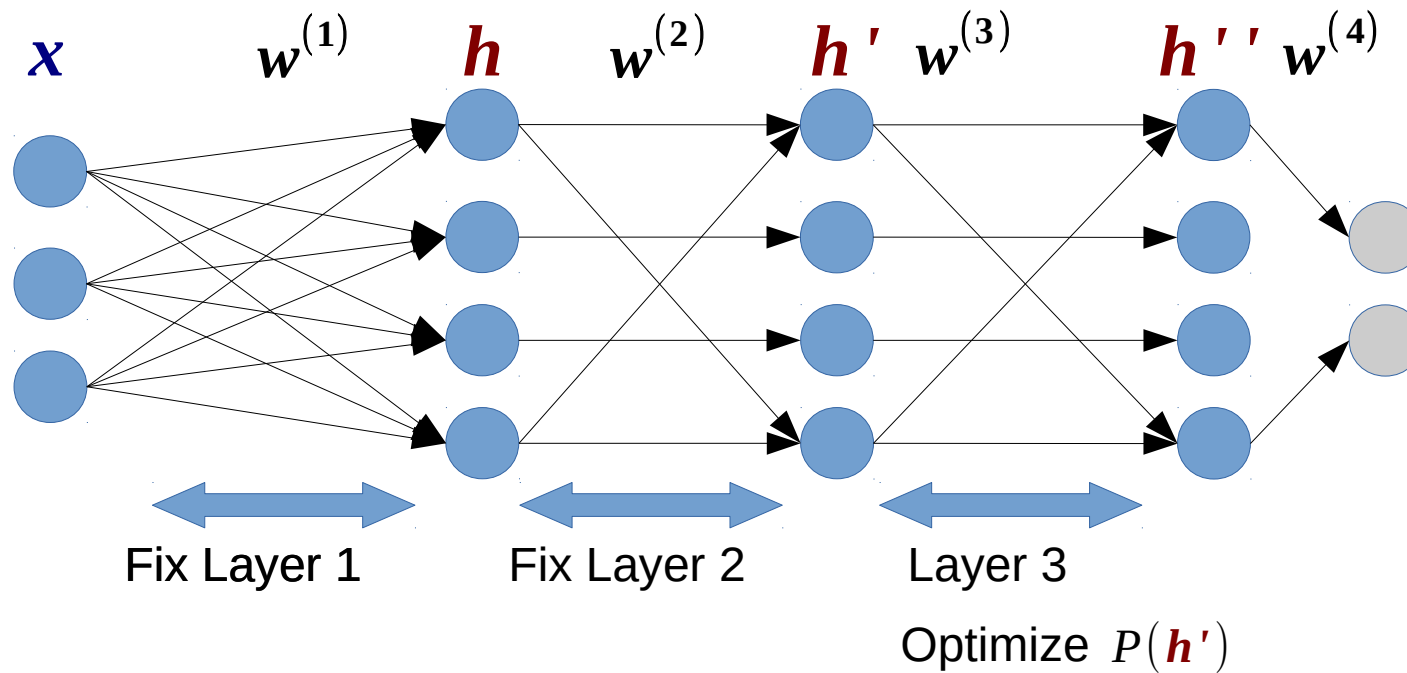
Layer-wise training

Train **one layer** at a time



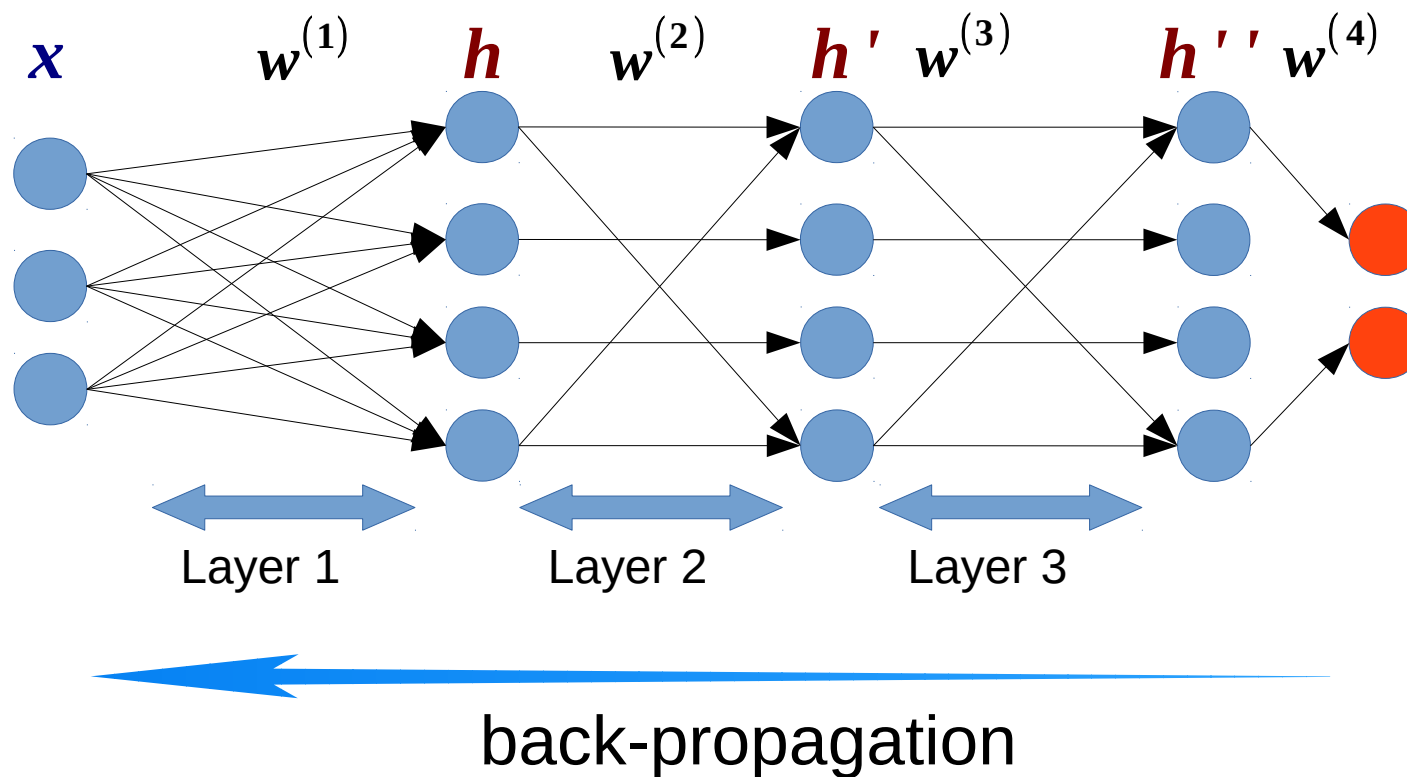
Layer-wise training

Train **one layer** at a time



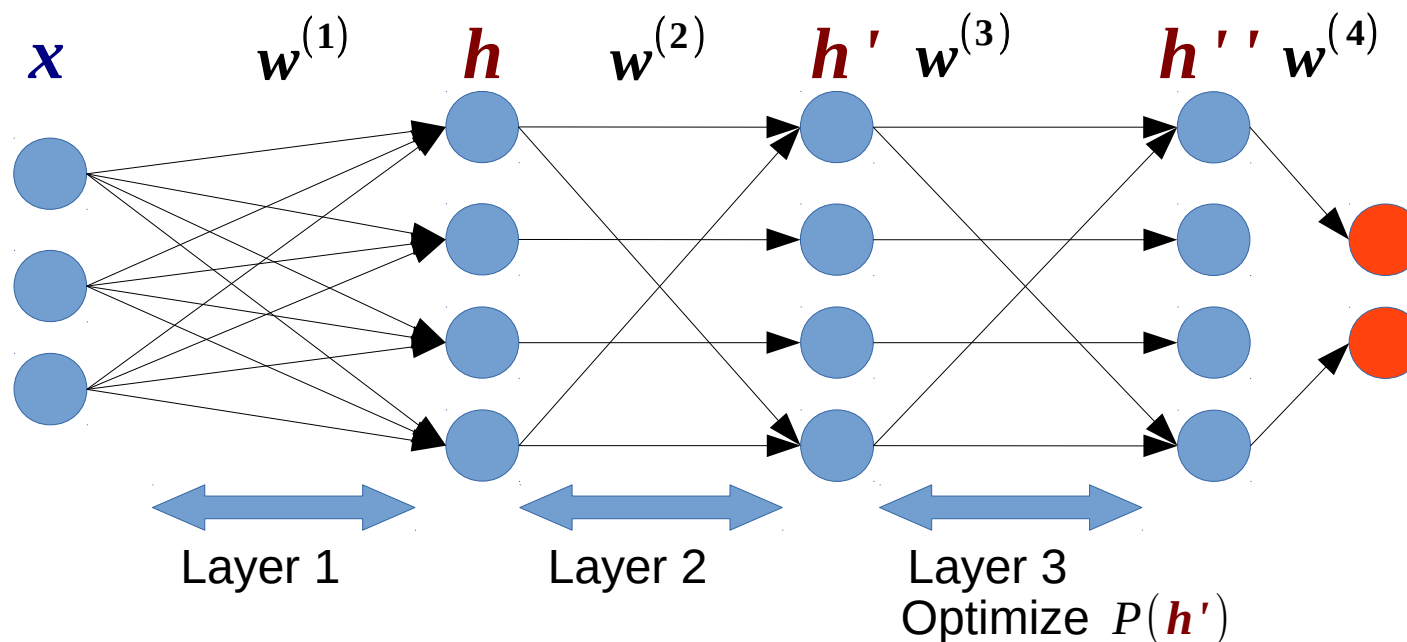
Layer-wise training

Finally, optimize $P(\text{label} \mid \text{input})$ with back-propagation



Layer-wise training

Finally, optimize $P(\text{label} \mid \text{input})$ with back-propagation



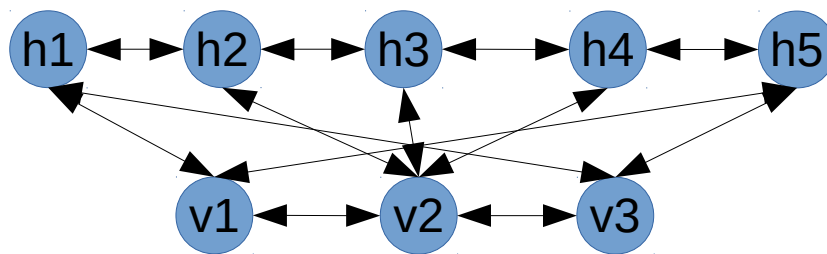
We need to understand our data first.
In other words, model $P(\text{input})$ before model $P(\text{label} \mid \text{input})$

Layer-wise training

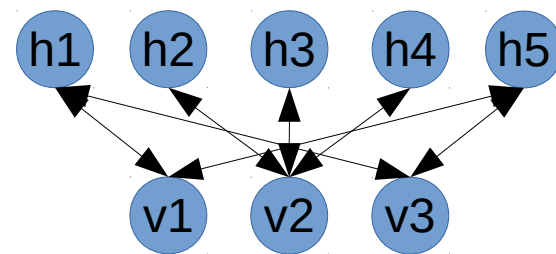
- Restricted Boltzmann Machine (RBM)
 - How to model $P(\text{input})$?
 - How it can help DNN training ?

Restricted Boltzmann Machine (RBM)

- Boltzmann Machine ?
 - a network that connect **binary neurons** using **symmetric connection**
- Restricted Boltzmann Machine ?
 - 2 layers: one hidden, one input
 - **No connection between hidden nodes** (Restricted)



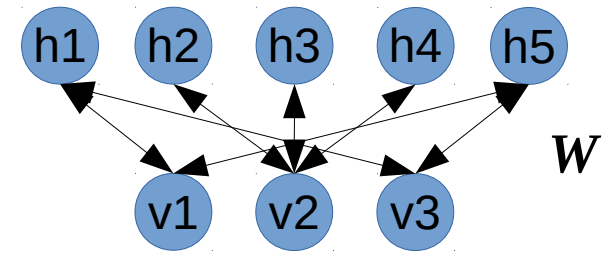
Boltzmann Machine



Restricted Boltzmann Machine

Restricted Boltzmann Machine (RBM)

Weights \rightarrow **Energy** \rightarrow **Probabilities**
 W $E(\mathbf{x}, \mathbf{h})$ $P(\mathbf{x}, \mathbf{h})$



Restricted Boltzmann Machine

Each possible joint configuration has an energy:

$$-E(\mathbf{x}, \mathbf{h}) = \mathbf{x}^T W \mathbf{h} + b^T \mathbf{x} + d^T \mathbf{h} = \sum_{i,j} x_i W_{ij} h_j + \sum_i b_i x_i + \sum_j d_j h_j$$

Probability of a joint configuration:

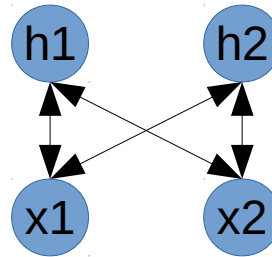
$$P(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{x}, \mathbf{h})} \quad \longrightarrow \quad P(\mathbf{x}) = \sum_{\mathbf{h}} P(\mathbf{x}, \mathbf{h})$$

Where:

- Z is partition function (normalizer) $Z = - \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$
- b, and d are bias terms

Restricted Boltzmann Machine (RBM)

Examples



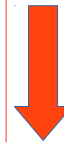
RBM with 2 visible, 2 hidden units

Observe: $x_1 = 1, x_2 = 1$

never observe: $x_1 = 0, x_2 = 0$

$$\begin{aligned} &P(v_1=1, v_2=1, h_1=1, h_2=1) \\ &P(v_1=1, v_2=1, h_1=1, h_2=0) \\ &P(v_1=1, v_2=1, h_1=0, h_2=1) \\ &P(v_1=1, v_2=1, h_1=0, h_2=0) \end{aligned}$$

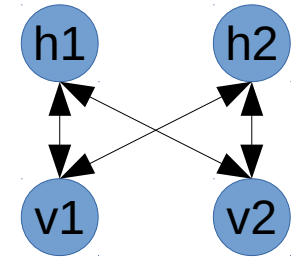
$$\begin{aligned} &P(v_1=0, v_2=0, h_1=1, h_2=0) \\ &P(v_1=0, v_2=0, h_1=0, h_2=1) \\ &\dots \end{aligned}$$



Restricted Boltzmann Machine

• Inference

- Hidden nodes are conditional independent given observation x



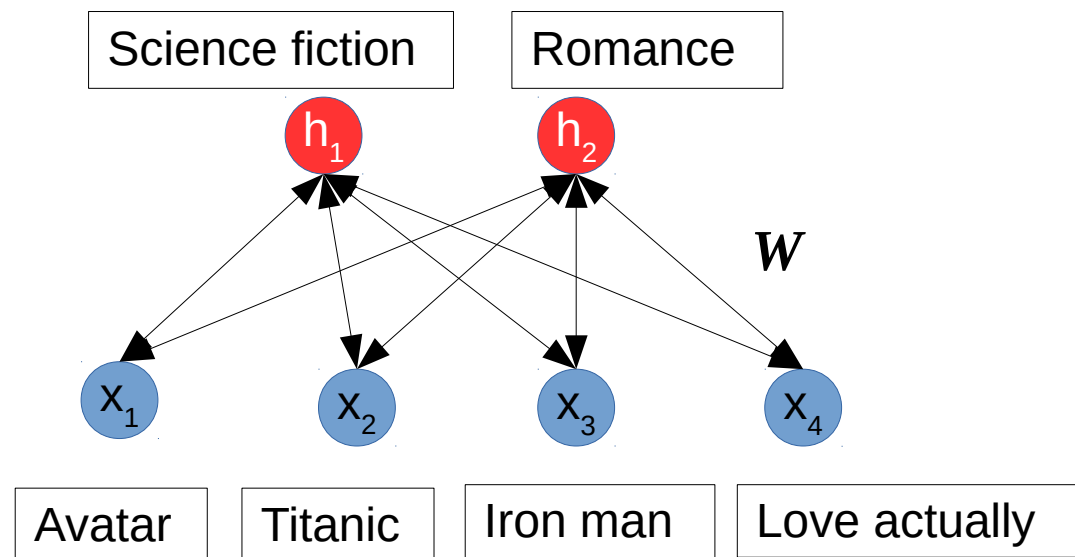
$$\begin{aligned}
 \Rightarrow P(\mathbf{h}|\mathbf{x}_m) &= \frac{P(\mathbf{x}_m, \mathbf{h})}{p(\mathbf{x}_m)} = \frac{e^{-E(\mathbf{x}_m, \mathbf{h})} / Z}{\sum_{\mathbf{h}'} e^{-E(\mathbf{x}_m, \mathbf{h}')} / Z} = \frac{e^{(\mathbf{x}_m^T \mathbf{W} \mathbf{h} + \mathbf{x}_m^T \mathbf{b} + \mathbf{h}^T \mathbf{d})}}{\sum_{\mathbf{h}'} e^{(\mathbf{x}_m^T \mathbf{W} \mathbf{h}' + \mathbf{x}_m^T \mathbf{b} + \mathbf{h}'^T \mathbf{d})}} \\
 &= \frac{e^{(\mathbf{x}_m^T \mathbf{W} \mathbf{h} + \mathbf{h}^T \mathbf{d})}}{\sum_{\mathbf{h}'} e^{(\mathbf{x}_m^T \mathbf{W} \mathbf{h}' + \mathbf{h}'^T \mathbf{d})}} = \frac{e^{(\sum_j \mathbf{x}_m^T \mathbf{W}_j \mathbf{h}_j + \mathbf{h}_j^T \mathbf{d}_j)}}{\sum_{\mathbf{h}'} e^{(\sum_j \mathbf{x}_m^T \mathbf{W}_j \mathbf{h}'_j + \mathbf{h}'_j^T \mathbf{d}_j)}} \\
 &= \frac{\prod_j e^{\mathbf{x}_m^T \mathbf{W}_j \mathbf{h}_j + \mathbf{h}_j^T \mathbf{d}_j}}{\prod_j \sum_{h_j'} e^{\mathbf{x}_m^T \mathbf{W}_j \mathbf{h}'_j + \mathbf{h}'_j^T \mathbf{d}_j}} = \prod_j \frac{e^{\mathbf{x}_m^T \mathbf{W}_j \mathbf{h}_j + \mathbf{h}_j^T \mathbf{d}_j}}{\sum_{h_j'} e^{\mathbf{x}_m^T \mathbf{W}_j \mathbf{h}'_j + \mathbf{h}'_j^T \mathbf{d}_j}} = \prod_j P(h_j | \mathbf{x}_m)
 \end{aligned}$$

$$\Rightarrow P(\mathbf{x}_m | \mathbf{h}) = \prod_j P(x_m^{(j)} | \mathbf{h})$$

$$\Rightarrow P(h_j = 1 | \mathbf{x}_m) = \frac{e^{\mathbf{x}_m^T \mathbf{W}_j + d_j}}{1 + e^{\mathbf{x}_m^T \mathbf{W}_j + d_j}} = \frac{1}{1 + e^{-(\mathbf{x}_m^T \mathbf{W}_j + d_j)}} = \text{sigmoid}(\mathbf{x}_m^T \mathbf{W}_j + d_j)$$

Restricted Boltzmann Machine

- Examples



➡ Observations

1 0 1 0

$P(h_1=1|\mathbf{x}_m) = \text{sigmoid}(\mathbf{x}_m^T \mathbf{W}_j + d_j) = 0.8$ ➡ This person likes SF with $p = 0.8$

$P(h_1=0, \mathbf{x}_m) = 1 - 0.8 = 0.2$

In practice: $h_1 = 1$ if $P(h_1|\mathbf{x}_m) > U(0, 1)$

➡ Generate data

Given a person likes SF: $h_1 = 1, h_2 = 0$

Generate which movies they should watch with $P(\mathbf{x}|h_1=1, h_2=0)$

Restricted Boltzmann Machine

Training process:

Given training sample \mathbf{x}_m , we want to maximize the likelihood

$$P(X = \mathbf{x}_m) = \sum_{\mathbf{h}} P(\mathbf{x}_m, \mathbf{y})$$

$$\hat{W} = \underset{W}{\operatorname{argmax}} \log(P(X = \mathbf{x}_m)) = \underset{W}{\operatorname{argmax}} \log\left(\sum_{\mathbf{h}} P(\mathbf{x}_m, \mathbf{y})\right)$$

Take derivative of Log-Likelihood w.r.t w_{ij}

$$\frac{\partial \log(P(X = \mathbf{x}_m))}{\partial w_{ij}} = \sum_{\mathbf{x}, \mathbf{h}} P(\mathbf{x}, \mathbf{h}) \frac{\partial(E(\mathbf{x}, \mathbf{h}))}{\partial w_{ij}} - \sum_{\mathbf{h}} P(\mathbf{h} | \mathbf{x}_m) \frac{\partial E(\mathbf{x}_m, \mathbf{h})}{\partial w_{ij}}$$

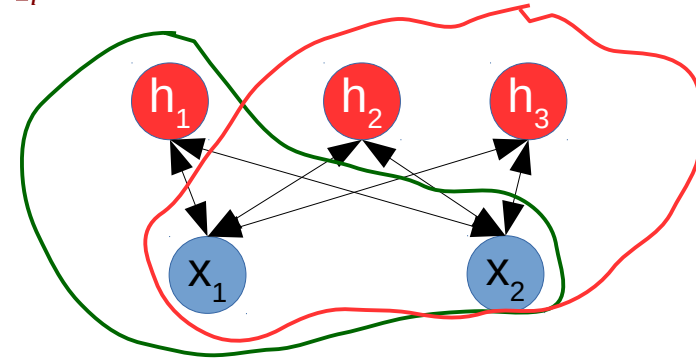
$$W \leftarrow W - \alpha \left(\frac{\partial \log(P(X = \mathbf{x}_m))}{\partial w_{ij}} \right) = W + \underbrace{\alpha \left(\sum_{\mathbf{h}} P(\mathbf{h} | \mathbf{x}_m) \frac{\partial E(\mathbf{x}_m, \mathbf{h})}{\partial w_{ij}} \right)}_{\text{Positive terms}} - \underbrace{\alpha \left(\sum_{\mathbf{x}, \mathbf{h}} P(\mathbf{x}, \mathbf{h}) \frac{\partial(E(\mathbf{x}, \mathbf{h}))}{\partial w_{ij}} \right)}_{\text{Negative terms}}$$

Restricted Boltzmann Machine

Training process:

The positive **term** is easy to compute:

$$\begin{aligned}\sum_{\mathbf{h}} P(\mathbf{h}|\mathbf{x}_m) \frac{\partial E(\mathbf{x}_m, \mathbf{h})}{\partial w_{ij}} &= \sum_{\mathbf{h}} P(\mathbf{h}|\mathbf{x}_m) h_i x_j = \sum_{h_i} \sum_{\mathbf{h}_{-i}} P(h_i|\mathbf{x}_m) P(\mathbf{h}_{-i}|\mathbf{x}_m) h_i x_j \\ &= \sum_{h_i} P(h_i|\mathbf{x}_m) h_i x_j \underbrace{\sum_{\mathbf{h}_{-i}} P(\mathbf{h}_{-i}|\mathbf{x}_m)}_{=1} \\ &= P(h_i=1, \mathbf{x}_m) x_j = \sigma\left(\sum_{j=1} w_{ji} x_j + d_i\right) x_j\end{aligned}$$



The **negative term** is hard when big hidden units

$$\sum_{\mathbf{x}, \mathbf{h}} P(\mathbf{x}, \mathbf{h}) \frac{\partial (E(\mathbf{x}, \mathbf{h}))}{\partial w_{ij}}$$

↓
of configurations \mathbf{x}, \mathbf{h} is exponential $2^{\text{num of units}}$

2 visible, 3 hidden $\rightarrow 2^2 \cdot 2^3 = 32$

➡ Can not compute $\frac{\partial \log(P(X=\mathbf{x}_m))}{\partial w_{ij}}$ directly

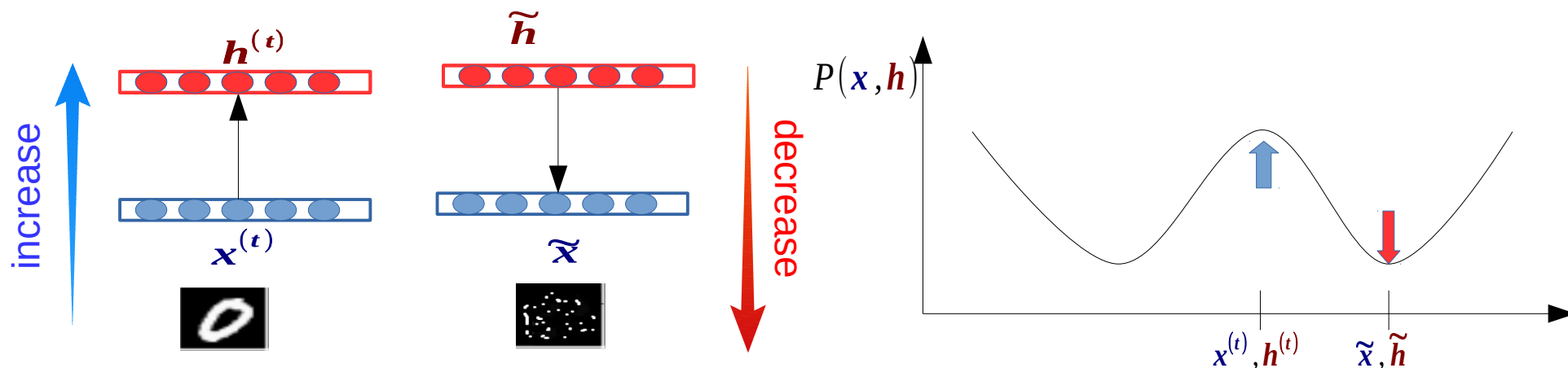
CONTRASTIVE DIVERGENCE [HINTON et al. 2002]

CONTRASTIVE DIVERGENCE [HINTON et al. 2002]

$$\text{ML} \quad \underbrace{\sum_{\mathbf{h}} P(\mathbf{h}|\mathbf{x}^{(t)}) \frac{\partial E(\mathbf{x}^{(t)}, \mathbf{h})}{\partial w_{ij}}}_{\mathcal{L}} - \underbrace{\sum_{\mathbf{x}, \mathbf{h}} P(\mathbf{x}, \mathbf{h}) \frac{\partial (E(\mathbf{x}, \mathbf{h}))}{\partial w_{ij}}}_{\mathcal{L}}$$

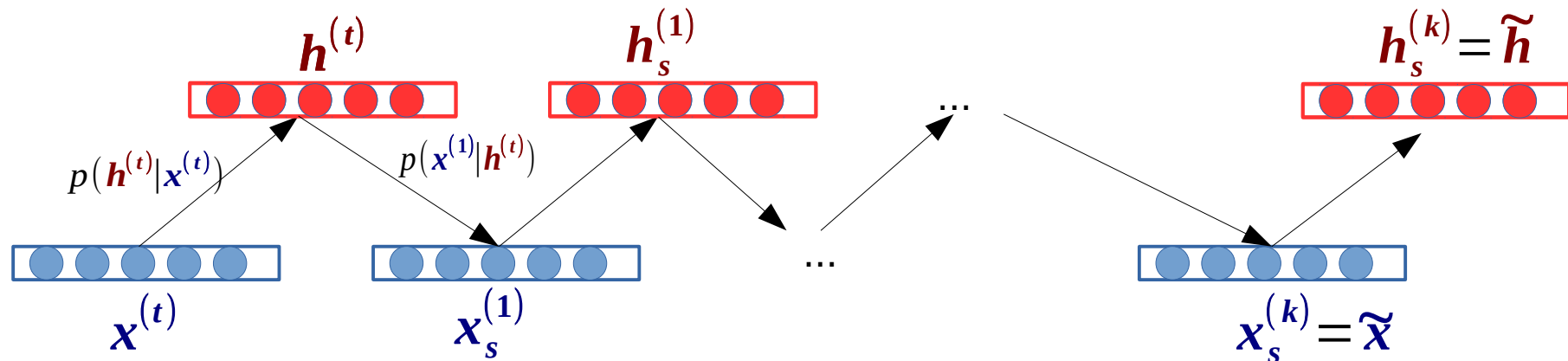
$$\text{CD} \quad \frac{\partial E(\mathbf{x}^{(t)}, \mathbf{h}^{(t)})}{\partial (w_{ij})} - \frac{\partial (E(\tilde{\mathbf{x}}, \tilde{\mathbf{h}}))}{\partial w_{ij}} = \mathbf{x}^{(t)T} \mathbf{h}^{(t)} - \tilde{\mathbf{x}}^T \tilde{\mathbf{h}}$$

$$\Rightarrow W \leftarrow W + \alpha (\mathbf{x}^{(t)T} \mathbf{h}^{(t)} - \tilde{\mathbf{x}}^T \tilde{\mathbf{h}})$$



CONTRASTIVE DIVERGENCE [HINTON et al. 2002]

Gibbs sampling



For each training sample $x^{(t)}$

1. Generate \tilde{x} using k-step Gibbs sampling from $x^{(t)}$
2. Update parameters:

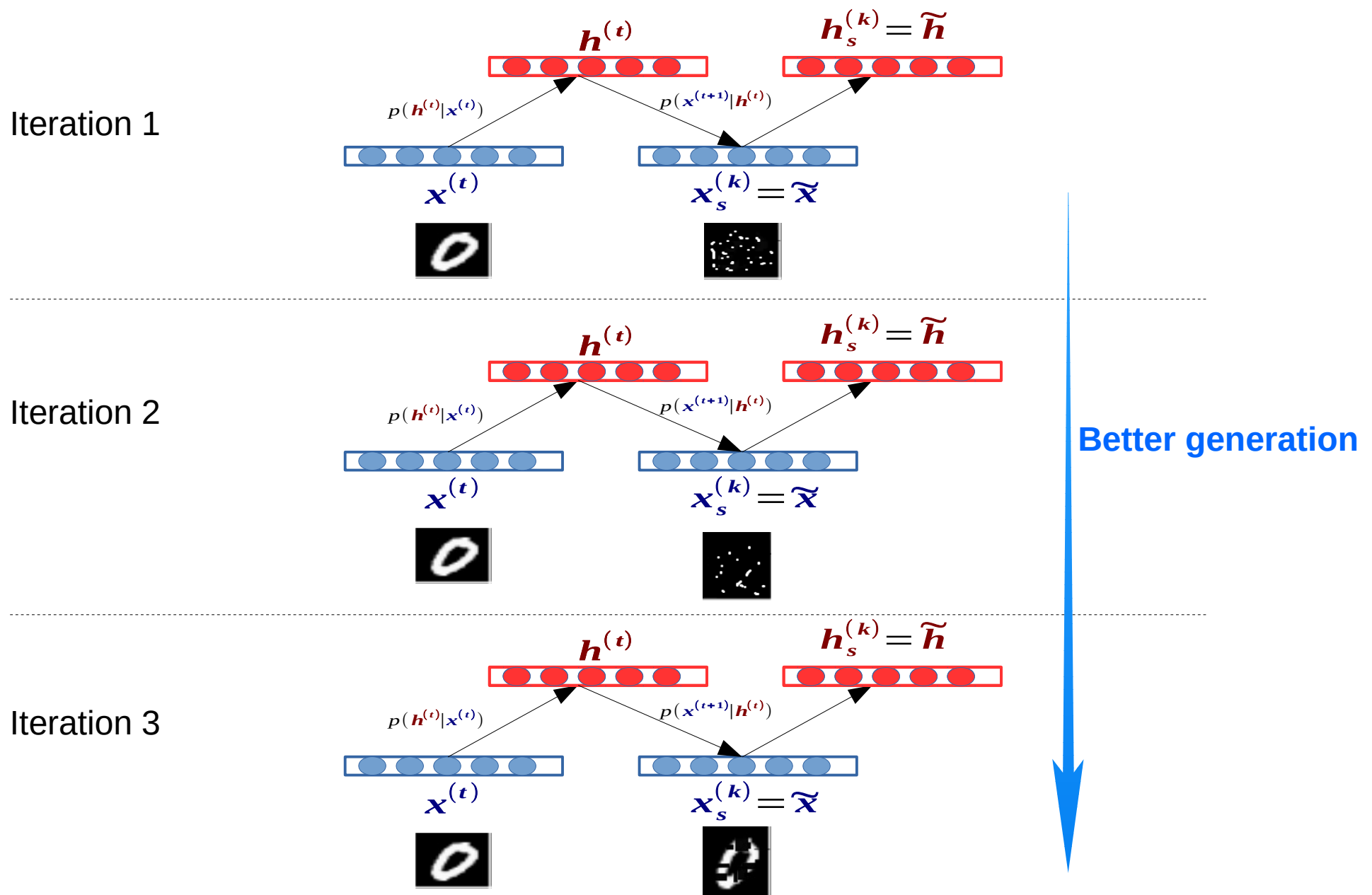
$$W \leftarrow W + \alpha (x^{(t)T} h^{(t)} - \tilde{x}^T \tilde{h})$$

$$b \leftarrow b + \alpha (x^{(t)} - \tilde{x})$$

$$d \leftarrow d + \alpha (h^{(t)} - \tilde{h})$$

3. Repeat until stopping criteria

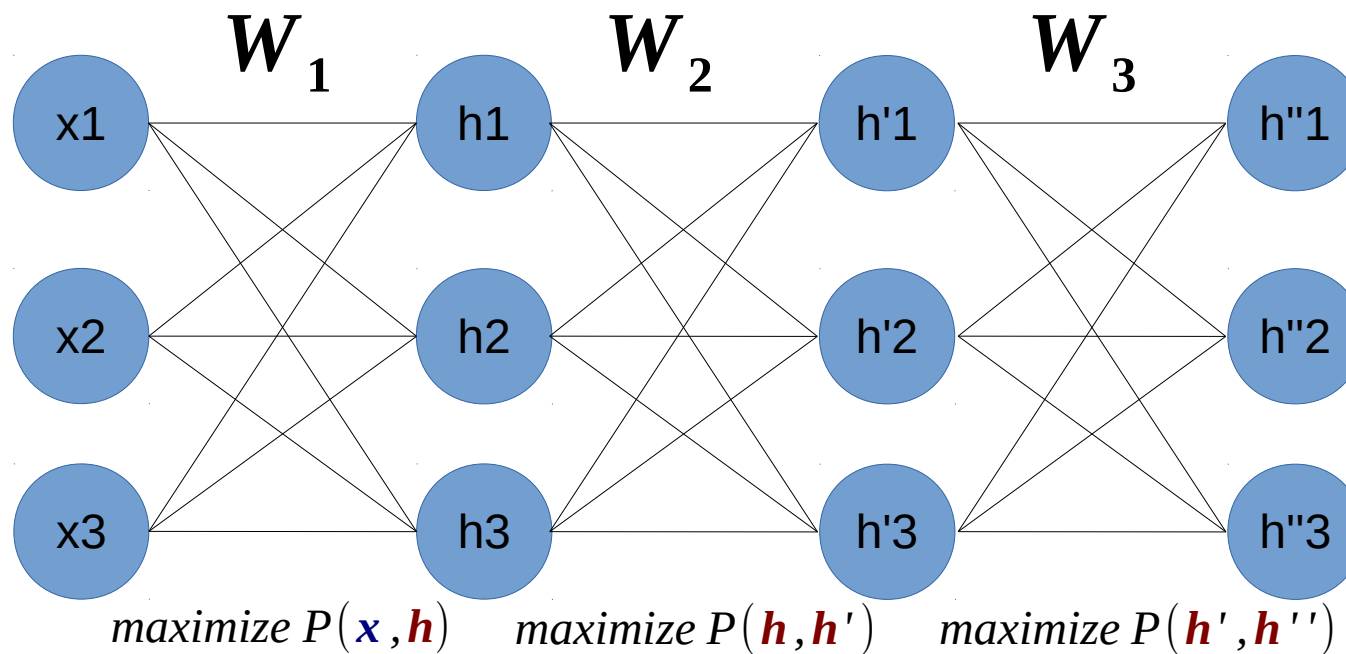
CONTRASTIVE DIVERGENCE [HINTON et al. 2002]



CONTRASTIVE DIVERGENCE [HINTON et al. 2002]

- CD-k: k iterations of Gibbs sampling
- The bigger k is, the better the estimate of gradient will be (Law of large number)
- In practice, k=1 is good for pre-training:
 - Optimize $P(x)$ is not the goal
 - Later fine-tuning with BP: $P(\text{label} | x)$

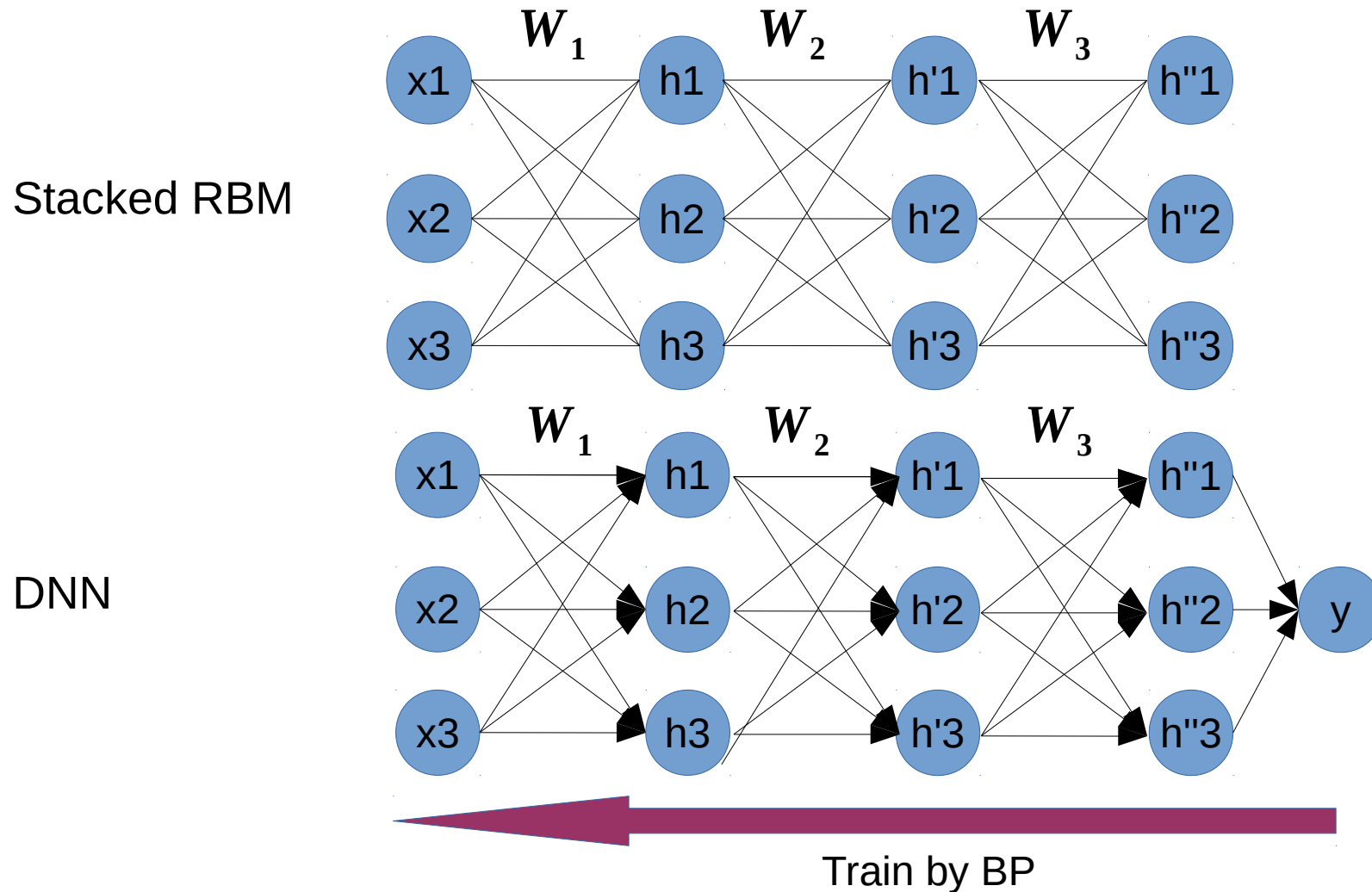
Layer-wise Pre-training RBM



Layer-wise Pre-training RBM

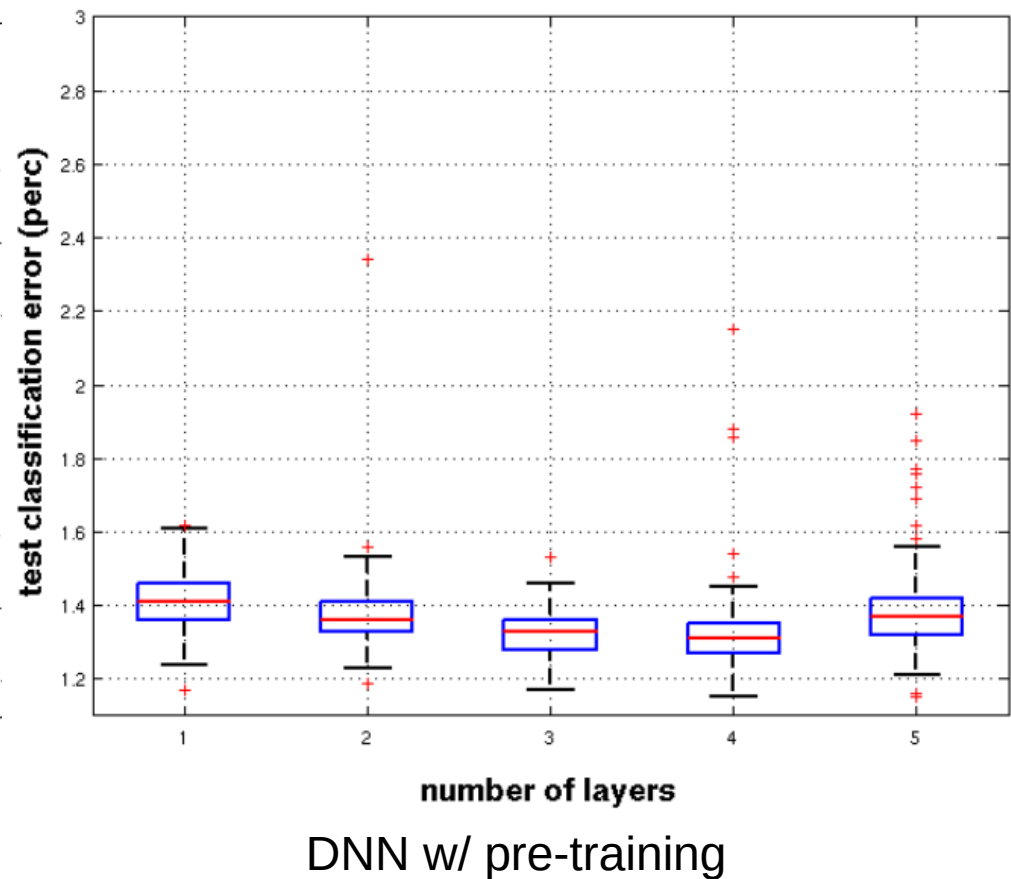
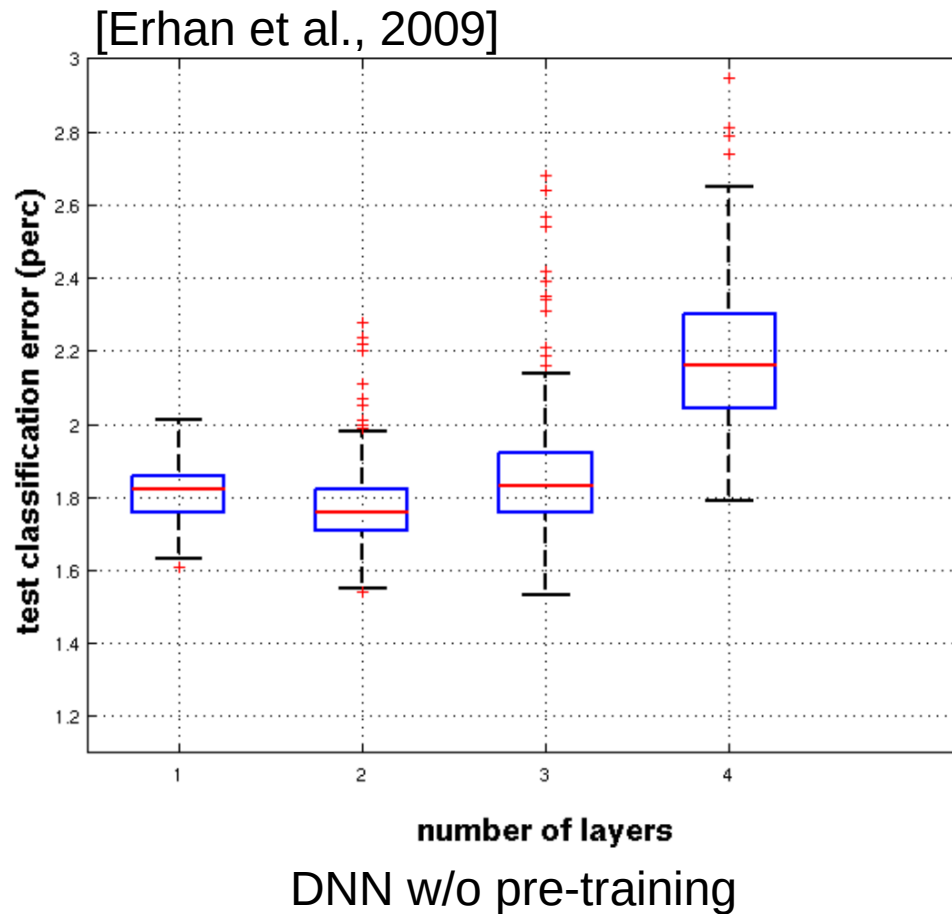
Application:

- Stacked RBM – Deep Belief Network: Generative model
- Use as initial parameter for DNN

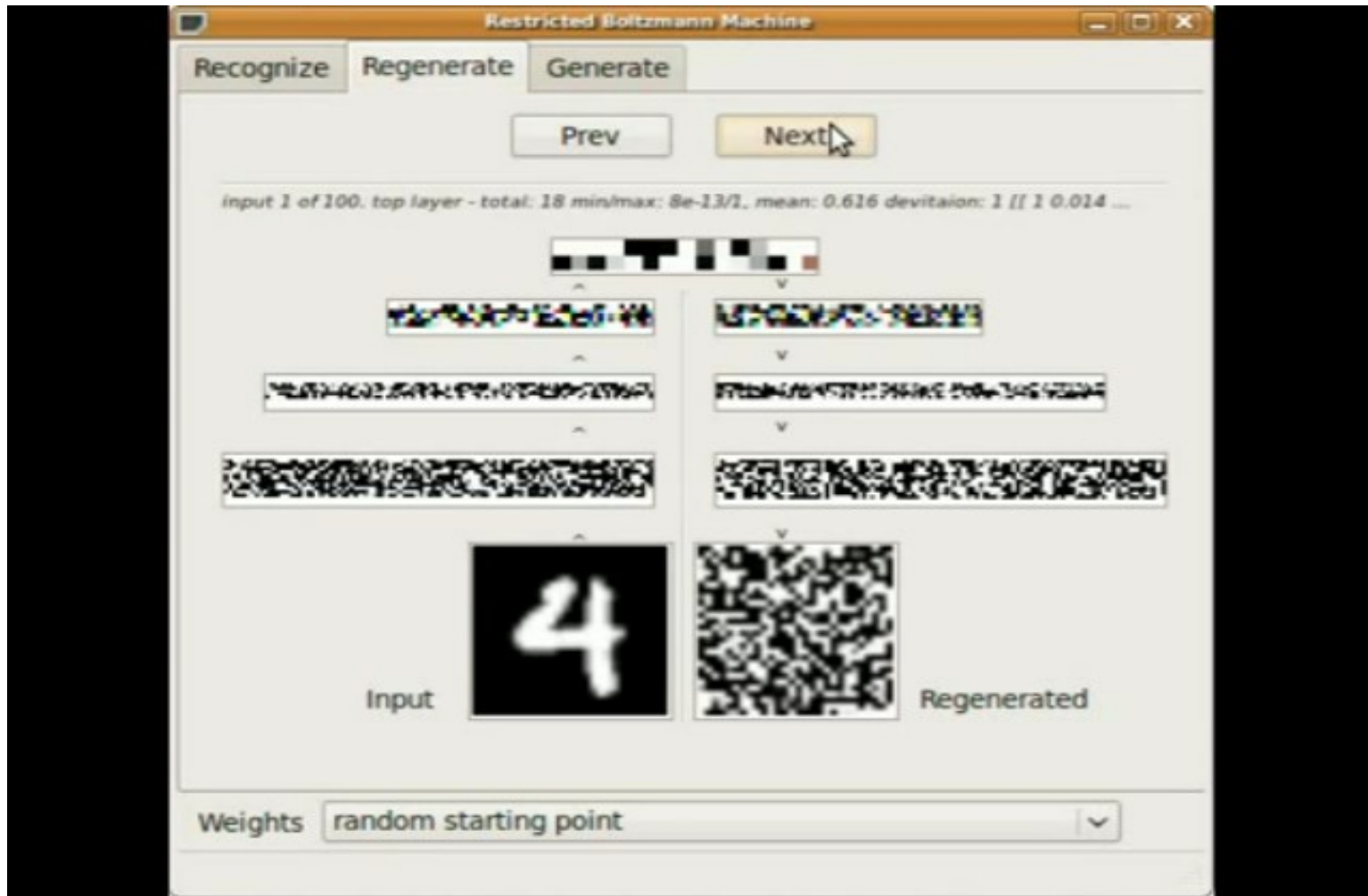


Deep Neural Networks

- Pre-train help DNN works better



RBM examples



<https://www.youtube.com/watch?v=0LTG64s6Xuc>