



Restricted Boltzmann Machines

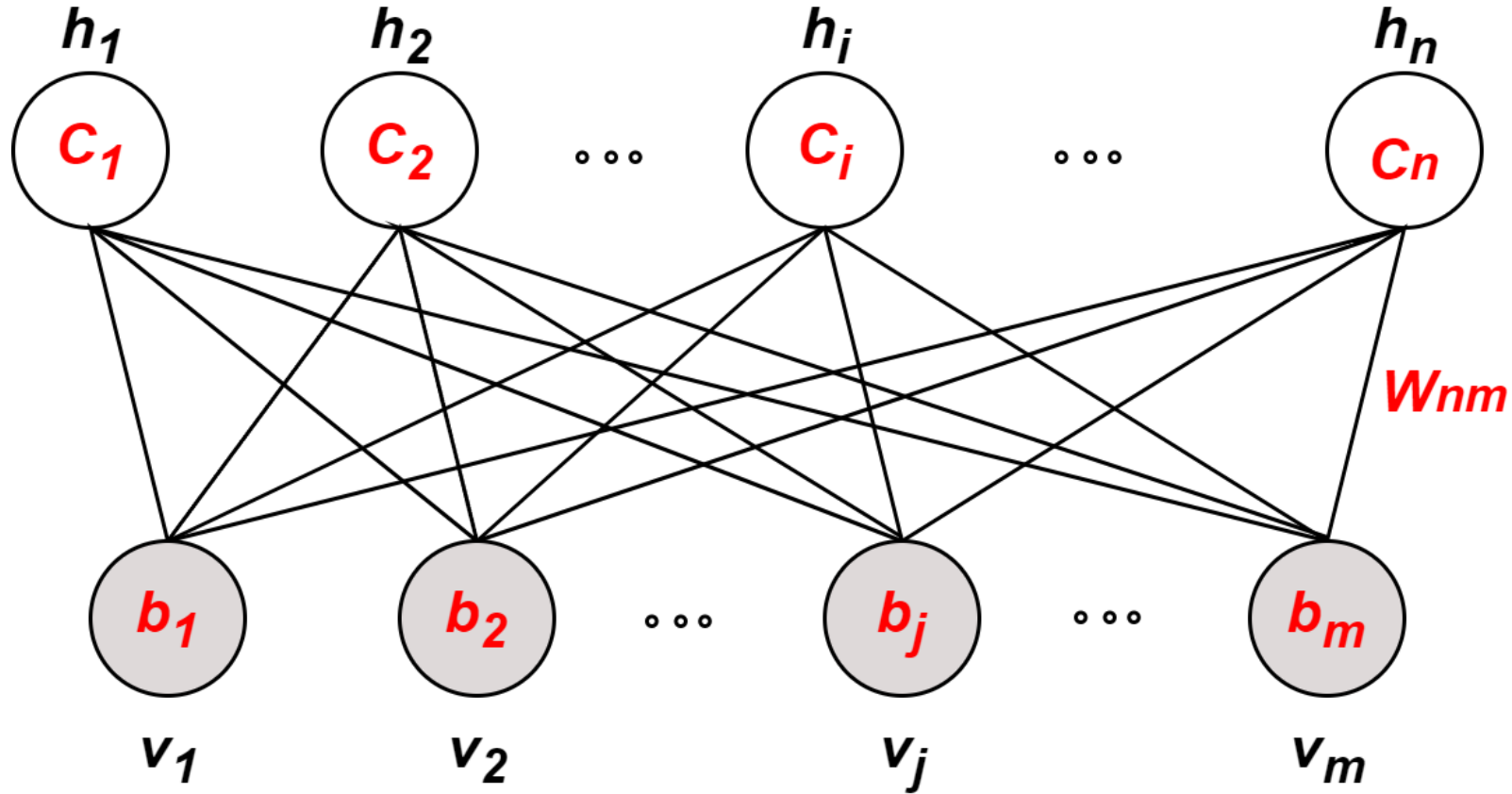
Sara Cocomello, Paolo Da Rold, Elena Rivaroli

A.Y. 2022-2023

Model Structure



Model structure



$$(\mathbf{v}, \mathbf{h}) \in \{0, 1\}^{m+n}$$

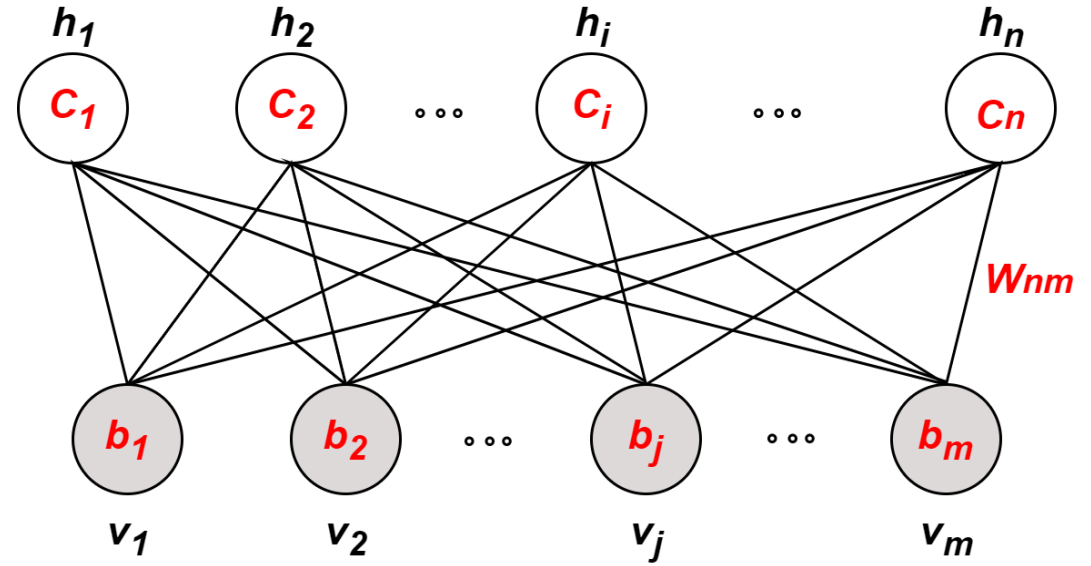
$$h_i \perp h_j | \mathbf{v}$$

$$v_i \perp v_j | \mathbf{h}$$

Fig. 1: The undirected graph of an RBM with n hidden and m visible variables

Gibbs distribution

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$$



Energy function

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i$$

Factorization

$$p(\mathbf{h} \mid \mathbf{v}) = \prod_{i=1}^n p(h_i \mid \mathbf{v}) \quad p(H_i = 1 \mid \mathbf{v}) = \sigma\left(\sum_{j=1}^m w_{ij} v_j + c_i\right)$$

$$p(\mathbf{v} \mid \mathbf{h}) = \prod_{i=1}^m p(v_i \mid \mathbf{h}) \quad p(V_j = 1 \mid \mathbf{h}) = \sigma\left(\sum_{i=1}^n w_{ij} h_i + b_j\right)$$

Gibbs sampling

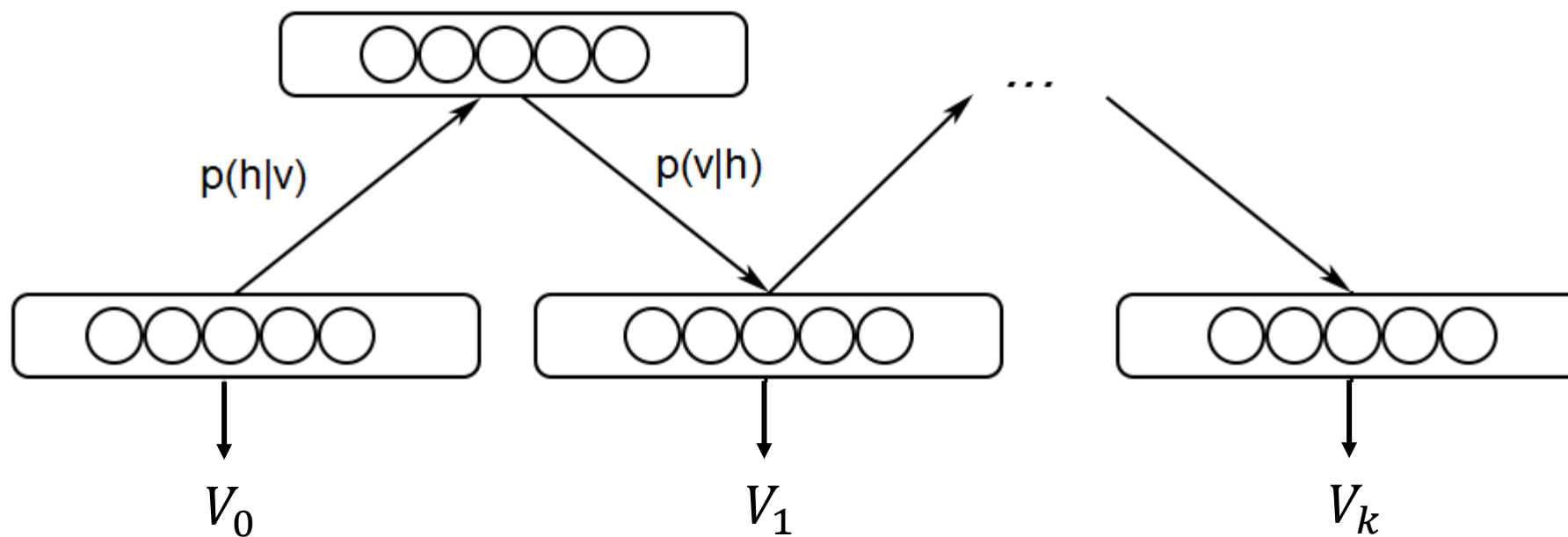


Fig. 2: Block Gibbs Sampling

RBM Training



Maximum Likelihood

The Log-Likelihood of a general MRF with latent variables is given by:

$$\ln \mathcal{L}(\boldsymbol{\theta} \mid \boldsymbol{v}) = \ln p(\boldsymbol{v} \mid \boldsymbol{\theta}) = \ln \frac{1}{Z} \sum_{\boldsymbol{h}} e^{-E(\boldsymbol{v}, \boldsymbol{h})} = \ln \sum_{\boldsymbol{h}} e^{-E(\boldsymbol{v}, \boldsymbol{h})} - \ln \sum_{\boldsymbol{v}, \boldsymbol{h}} e^{-E(\boldsymbol{v}, \boldsymbol{h})}$$

The gradient w.r.t. the parameters is:

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta} \mid \boldsymbol{v})}{\partial \boldsymbol{\theta}} = - \sum_{\boldsymbol{h}} p(\boldsymbol{h} \mid \boldsymbol{v}) \frac{\partial E(\boldsymbol{v}, \boldsymbol{h})}{\partial \boldsymbol{\theta}} + \sum_{\boldsymbol{v}, \boldsymbol{h}} p(\boldsymbol{v}, \boldsymbol{h}) \frac{\partial E(\boldsymbol{v}, \boldsymbol{h})}{\partial \boldsymbol{\theta}}$$

Maximum Likelihood for RBM

The explicit derivative w.r.t the weights w_{ij} is:

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta} | \boldsymbol{v})}{\partial w_{ij}} = - \sum_{\boldsymbol{h}} p(\boldsymbol{h} | \boldsymbol{v}) \frac{\partial E(\boldsymbol{v}, \boldsymbol{h})}{\partial w_{ij}} + \sum_{\boldsymbol{v}, \boldsymbol{h}} p(\boldsymbol{v}, \boldsymbol{h}) \frac{\partial E(\boldsymbol{v}, \boldsymbol{h})}{\partial w_{ij}}$$

Computing the derivatives and using a factorization trick we obtain:

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta} | \boldsymbol{v})}{\partial w_{ij}} = \boxed{p(H_i = 1 | \boldsymbol{v}) v_j} - \boxed{\sum_{\boldsymbol{v}} p(\boldsymbol{v}) p(H_i = 1 | \boldsymbol{v}) v_j}$$

OK NOT OK

Contrastive Divergence

To solve the problem we use the contrastive divergence technique:

$$\text{CD}_k(\boldsymbol{\theta}, \mathbf{v}^{(0)}) = - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(0)}) \frac{\partial E(\mathbf{v}^{(0)}, \mathbf{h})}{\partial \boldsymbol{\theta}} + \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(k)}) \frac{\partial E(\mathbf{v}^{(k)}, \mathbf{h})}{\partial \boldsymbol{\theta}} \quad \text{Usually } k = 1$$

Maximum Likelihood:  Minimize $[\text{KL}(q|p)]$

Contrastive divergence:  Minimize $[\text{KL}(q|p) - \text{KL}(p_k|p)]$

Pseudocode of Contrastive Divergence

Input: RBM, training batch S

Output: gradient approximations

Init $\Delta \mathbf{w} = \Delta \mathbf{b} = \Delta \mathbf{c} = \mathbf{0}$

forall v in a batch

```
┌    $v^{(0)} \leftarrow v$   
├   for  $t=0$  step  $(k-1)$   
│   ┌    $v^{(k)} \leftarrow$  sample  $h^{(t)}$  from  $p(h|v^{(t)})$  and subsequently  $v^{(t+1)}$  from  $p(h|v^{(t)})$   
│   └   update  $\Delta \mathbf{w}, \Delta \mathbf{b}, \Delta \mathbf{c}$  using CD rules
```

Code 1: k steps contrastive divergence

Application



Classification on MNIST dataset



- **PREPROCESSING** : scaling the pixel values in $[0,1]$ with grayscale
- **TRAINING** : defining the RBM architecture by setting the number of hidden units and train on MNIST
- **FEATURE EXTRACTION**: transform the data with the features extracted by RBM
- **CLASSIFICATION**: training general classifiers with extracted features

Fig. 3: Sample from MNIST dataset 28x28 images

Features extraction and classification

| MODEL | ACCURACY RBM | ACCURACY RAW |
|----------|--------------|--------------|
| Logistic | 0.94 | 0.92 |
| KNN | 0.97 | 0.98 |
| SVM | 0.97 | 0.98 |

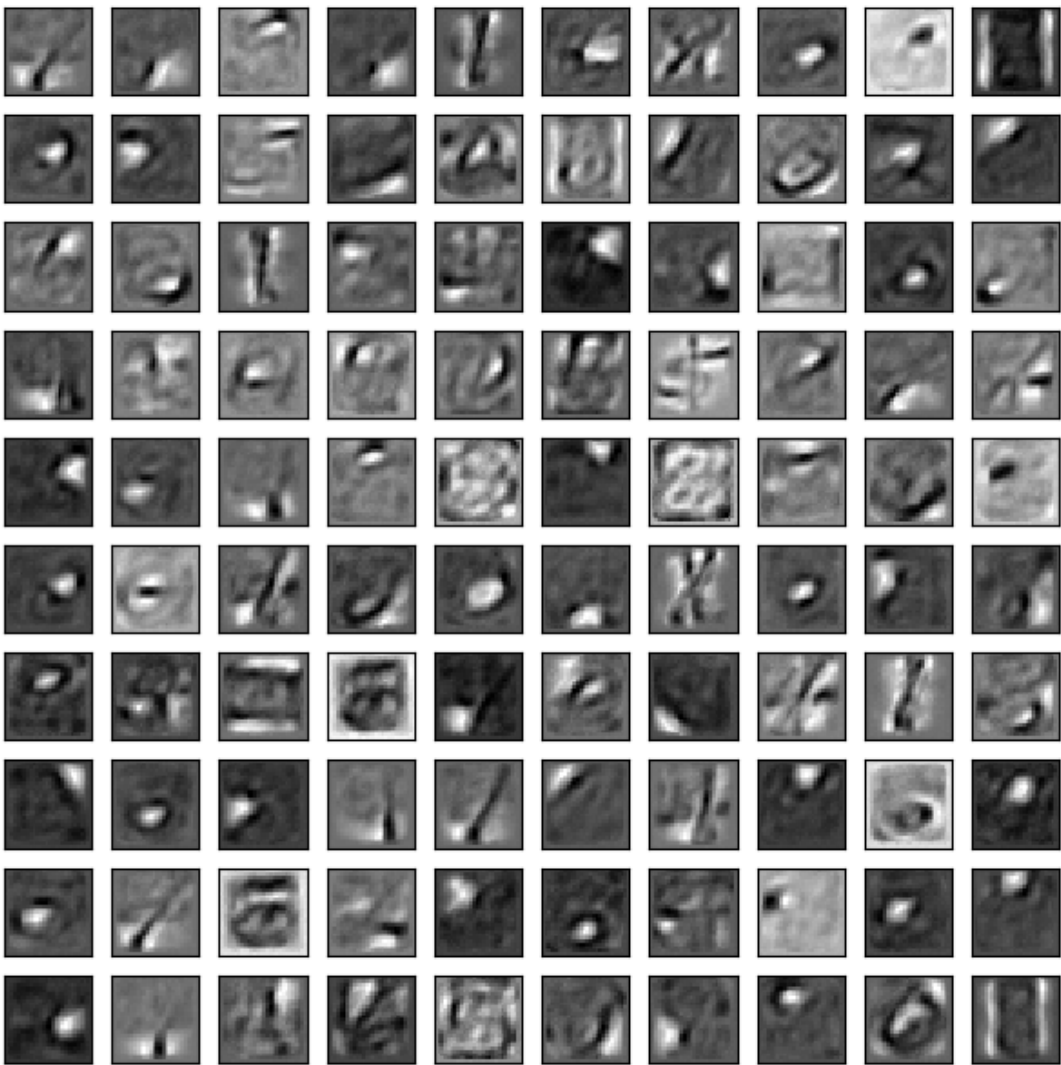


Fig. 4: 100 components

Evaluation

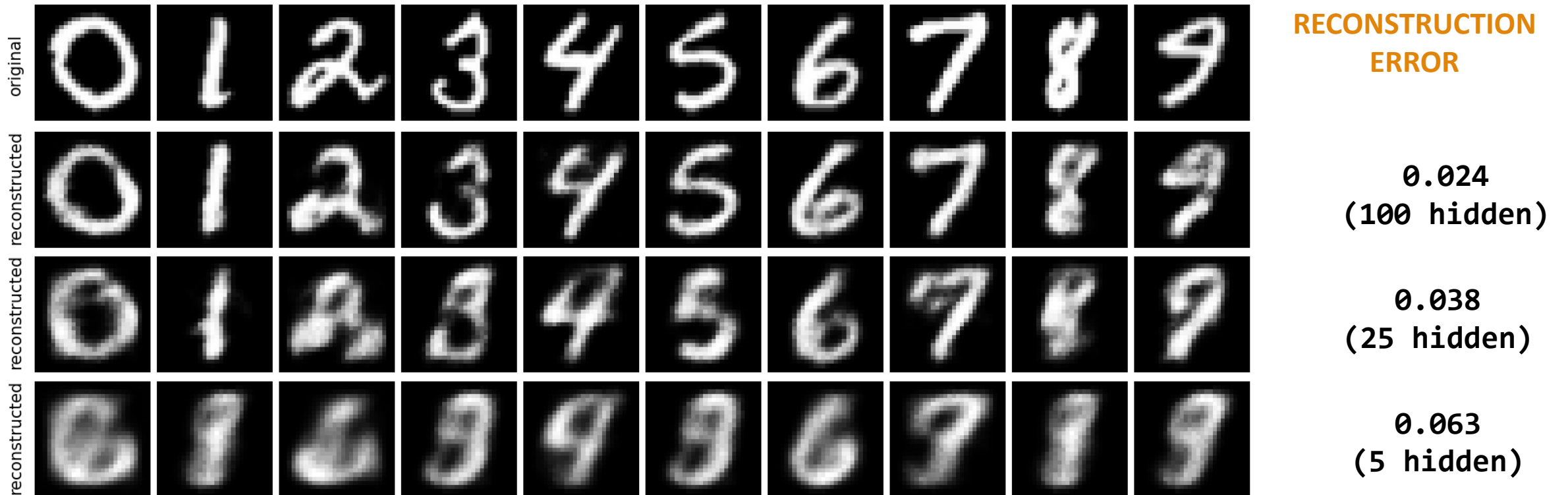


Fig. 5: Comparison of reconstruction with different numbers of hidden variables

Conclusions

RBM vs raw pixels for classification

- Feature extraction
- Dimensionality reduction
- Unsupervised learning

Current application

- Using RBM as building blocks of Deep neural networks

References

- [1] A. Fischer and C. Igel, “An introduction to restricted Boltzmannmachines”. In: Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications. Springer, Berlin Heidelberg, 2012.
- [2] Carreira-Perpiñán, M.A., Hinton, G.E.: On contrastive divergence learning. In: 10th International Workshop on Artificial Intelligence and Statistics (AISTATS 2005), pp. 59–66 (2005)
- [3] Hinton, G.E.: Training products of experts by minimizing contrastive divergence. Neural Computation 14, 1771–1800 (2002)

The image features a white background with decorative curved lines in the corners. In the top-right corner, a thick, multi-layered arc curves from the top edge towards the right, transitioning in color from a light teal at the top to a pale yellow at the bottom. In the bottom-left corner, a similar thick, multi-layered arc curves from the left edge towards the bottom, also transitioning from light teal to pale yellow.

Thank you for the attention!

Factorization trick

$$\begin{aligned}\sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} &= \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}) h_i v_j \\ &= \sum_{\mathbf{h}} \prod_{k=1}^n p(h_k | \mathbf{v}) h_i v_j = \sum_{h_i} \sum_{\mathbf{h}_{-i}} p(h_i | \mathbf{v}) p(\mathbf{h}_{-i} | \mathbf{v}) h_i v_j \\ &= \sum_{h_i} p(h_i | \mathbf{v}) h_i v_j \underbrace{\sum_{\mathbf{h}_{-i}} p(\mathbf{h}_{-i} | \mathbf{v})}_{=1} = p(H_i = 1 | \mathbf{v}) v_j = \sigma \left(\sum_{j=1}^m w_{ij} v_j + c_i \right) v_j\end{aligned}$$