

Restricted Boltzmann Machines

GRAHAM TAYLOR

VECTOR INSTITUTE

SCHOOL OF ENGINEERING
UNIVERSITY OF GUELPH

CANADIAN INSTITUTE
FOR ADVANCED RESEARCH

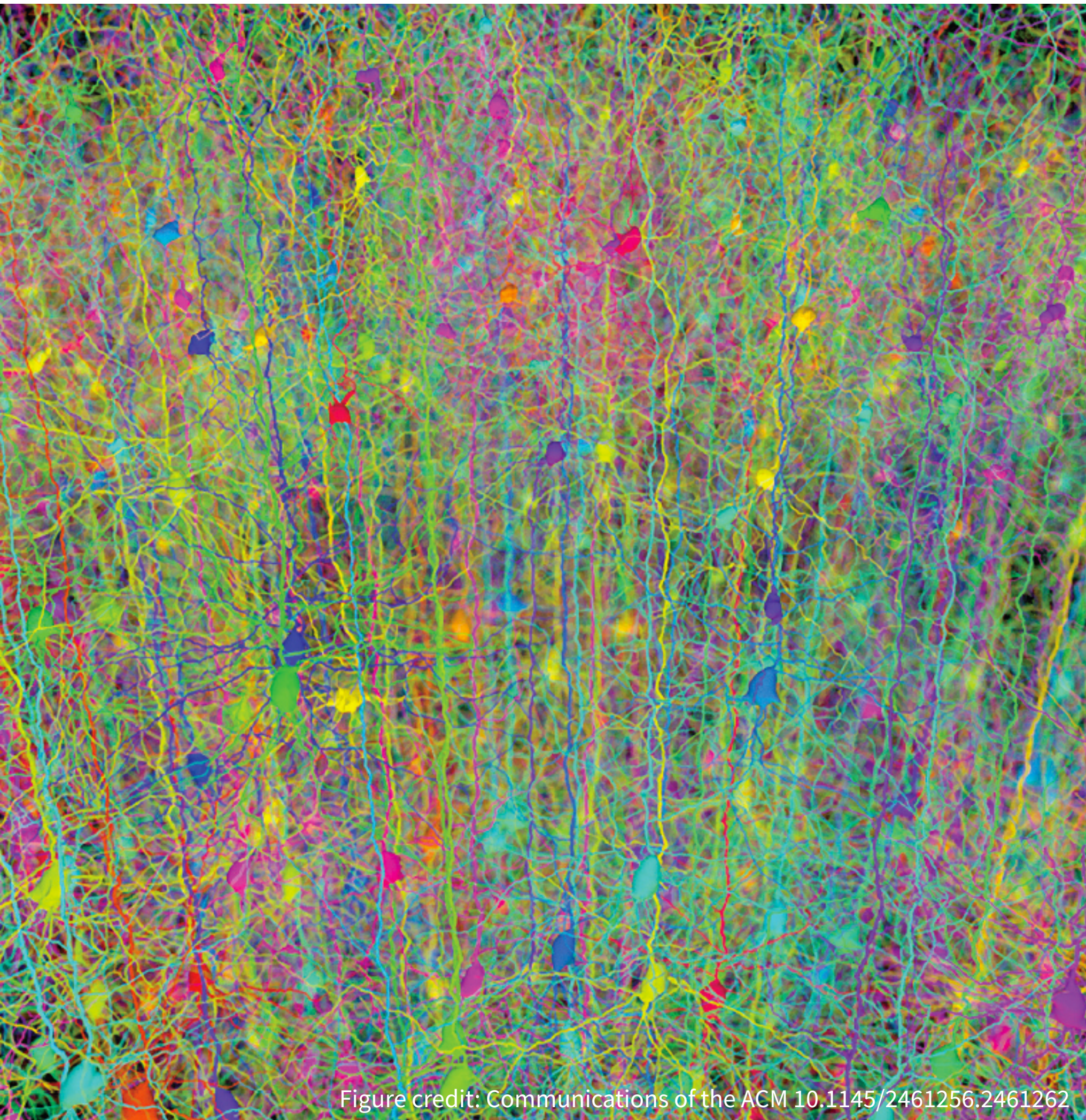


Figure credit: Communications of the ACM 10.1145/2461256.2461262

UNIVERSITY
of GUELPH

CHANGING LIVES
IMPROVING LIFE

CIFAR
CANADIAN
INSTITUTE
FOR
ADVANCED
RESEARCH

Scaling ML to the Challenges of AI

Classification algorithms take an input from a high-dimensional distribution and **summarize** it with a category label

- During this process, the classifier **discards most of the information** in the input and produces a single output

It is possible to ask our ML models to do many other tasks

- Some of which require them to produce **multiple outputs**
- Most require a complete understanding of the **entire structure of the input**, with no option to ignore sections of it

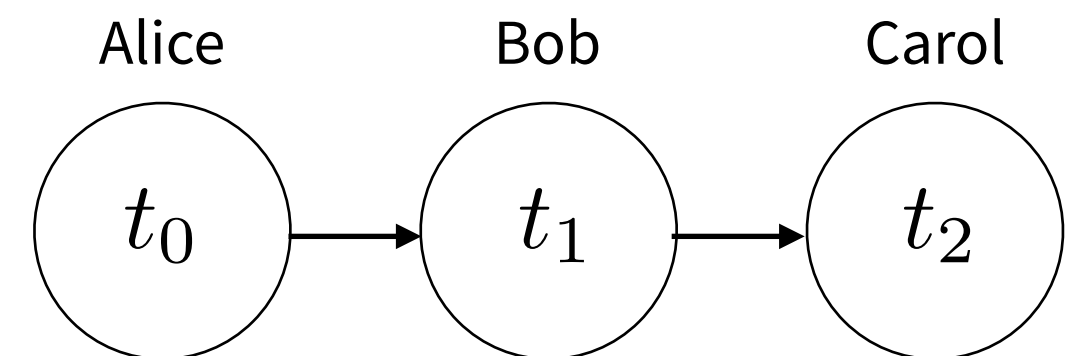
Structured Probabilistic Models for Machine Learning

Modeling a rich distribution over random variables is a challenging task, both **computationally** and **statistically**

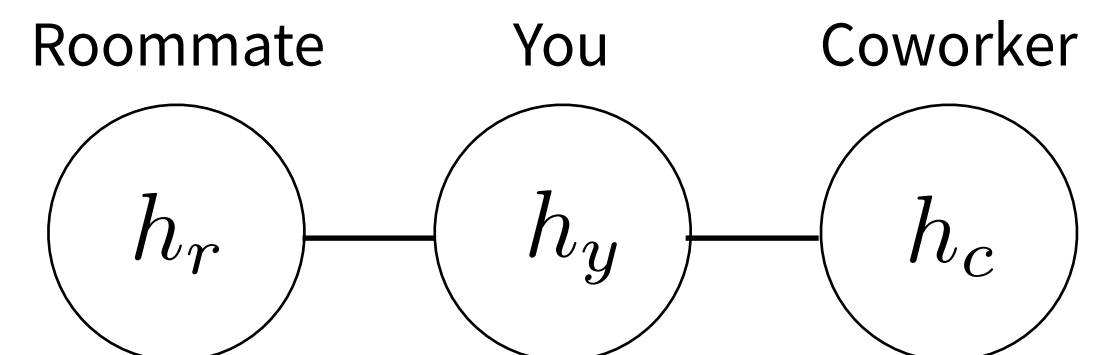
Structured probabilistic models (graphical models) provide a formal framework for modeling **only direct interactions** between random variables

- significantly fewer parameters
- estimated reliably from less data

directed graphical model

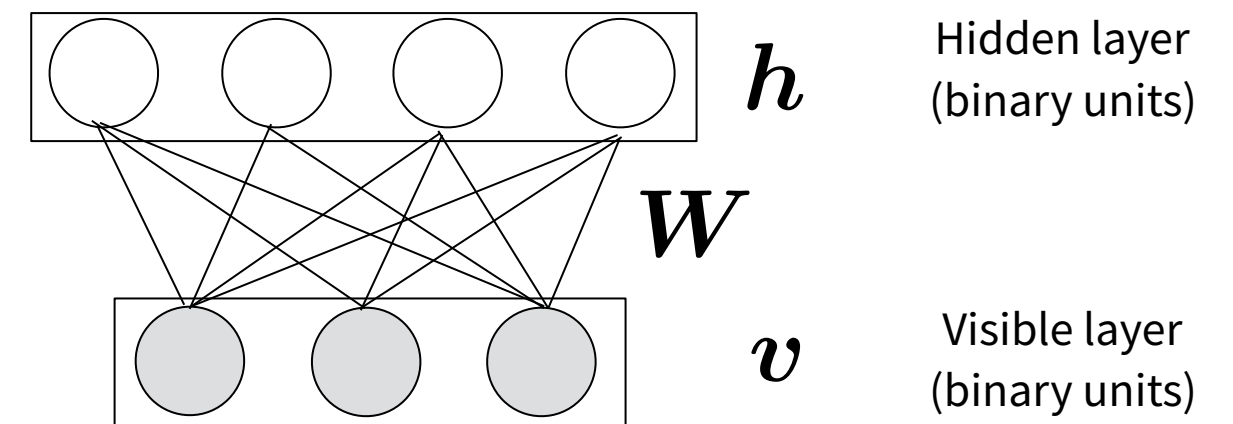


undirected graphical model



Restricted Boltzmann Machine

- Undirected graphical model with **bipartite structure** and **restricted connectivity** to make inference and learning easier



- Energy function:

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}$$

- Distribution:

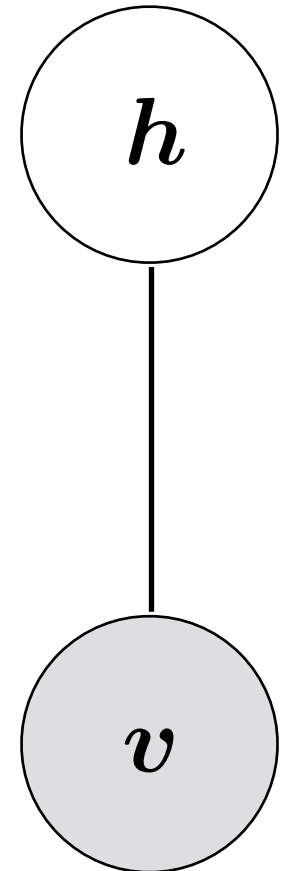
$$P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp \{-E(\mathbf{v}, \mathbf{h})\}$$

Markov Network View

Markov network with **vector nodes**:

$$\begin{aligned} P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) &= \exp(-E(\mathbf{v}, \mathbf{h}))/Z \\ &= \exp(\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h})/Z \\ &= \underbrace{\exp(\mathbf{b}^\top \mathbf{v}) \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}_{\text{factors}}/Z \end{aligned}$$

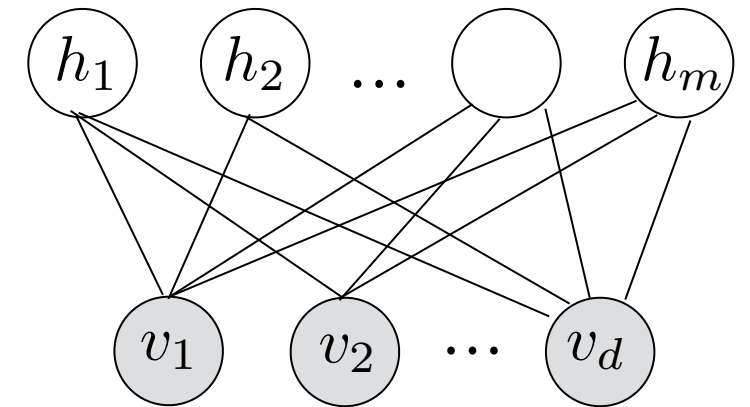


The notation based on an energy function is simply an **alternative** to the representation as the product of factors

Markov Network View (2)

Markov network with **scalar nodes**:

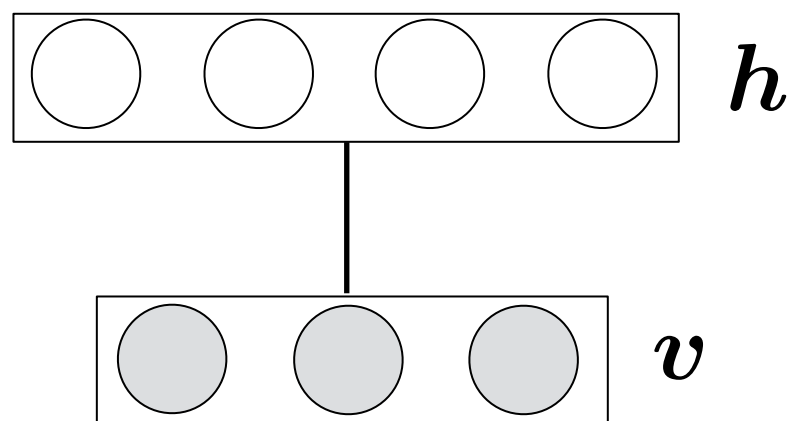
$$P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) = \frac{1}{Z} \underbrace{\prod_j \prod_k \exp(W_{j,k} h_j v_k)}_{\text{pair-wise factors}} \underbrace{\prod_k \exp(b_k v_k) \prod_j \exp(c_j h_j)}_{\text{unary factors}}$$



This scalar visualization is more informative of the **structure within the vectors**

Inference in RBM

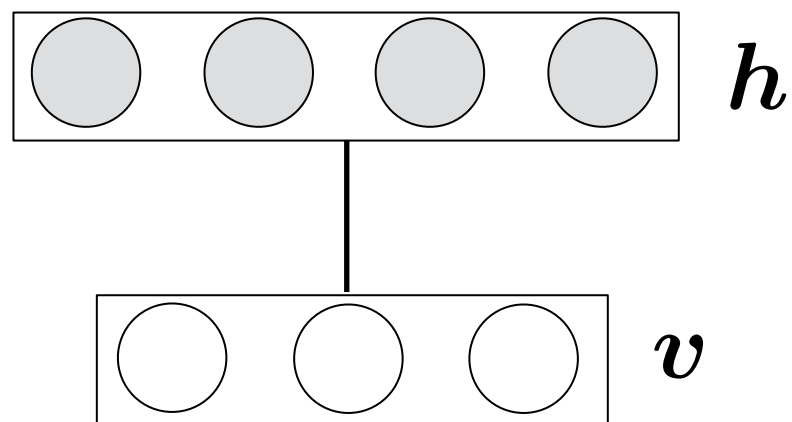
Though $P(\mathbf{v})$ is intractable, the bipartite graph structure of the RBM has the special property of its conditional distributions $P(\mathbf{h}|\mathbf{v})$ and $P(\mathbf{v}|\mathbf{h})$ being factorial and relatively **simple to compute and sample**:



$$P(\mathbf{h}|\mathbf{v}) = \prod_j P(h_j = 1|\mathbf{v})$$

$$P(h_j = 1|\mathbf{v}) = \frac{1}{1 + \exp(-(c_j + \mathbf{v}^\top \mathbf{W}_{:,j}))} = \sigma(c_j + \mathbf{v}^\top \mathbf{W}_{:,j})$$

↖ jth column of \mathbf{W}



$$P(\mathbf{v}|\mathbf{h}) = \prod_k P(v_k = 1|\mathbf{h})$$

$$P(v_k = 1|\mathbf{h}) = \frac{1}{1 + \exp(-(b_k + \mathbf{W}_{k,:} \mathbf{h}))} = \sigma(b_k + \mathbf{W}_{k,:} \mathbf{h})$$

↖ kth row of \mathbf{W}

Free Energy

- Many algorithms that operate on probabilistic models need to compute not $p_{\text{model}}(\mathbf{x})$ but only $\log \tilde{p}_{\text{model}}(\mathbf{x})$
- For energy-based models with latent variables \mathbf{h} , these algorithms are sometimes phrased in terms of the negative of this quantity, called the **free energy**:

$$\mathcal{F}(\mathbf{x}) = -\log \sum_{\mathbf{h}} \exp(-E(\mathbf{x}, \mathbf{h}))$$

- For the RBM:

$$\begin{aligned} P(\mathbf{v}) &= \sum_{\mathbf{h} \in \{0,1\}^m} P(\mathbf{v}, \mathbf{h}) = \sum_{\mathbf{h} \in \{0,1\}^m} \exp(-E(\mathbf{v}, \mathbf{h})) / Z \\ &= \exp \left(\mathbf{b}^\top \mathbf{v} + \sum_{j=1}^m \log (1 + \exp (c_j + \mathbf{v}^\top \mathbf{W}_{:,j})) \right) / Z \\ &= \exp (-\mathcal{F}(\mathbf{v})) / Z \end{aligned}$$