

An exposition on the propriety of restricted Boltzmann machines

Andee Kaplan, Daniel Nordman, Stephen Vardeman
Department of Statistics, Iowa State University

Deep learning

Three layer deep Boltzmann machine, with visible-to-hidden and hidden-to-hidden connections but no within-layer connections. This can be considered as multiple single layer restricted Boltzmann machines with the lower stack hidden layer acting as the visible layer for the higher stacked model. Claimed ability to learn "internal representations that become increasingly complex" [4], used in classification problems.

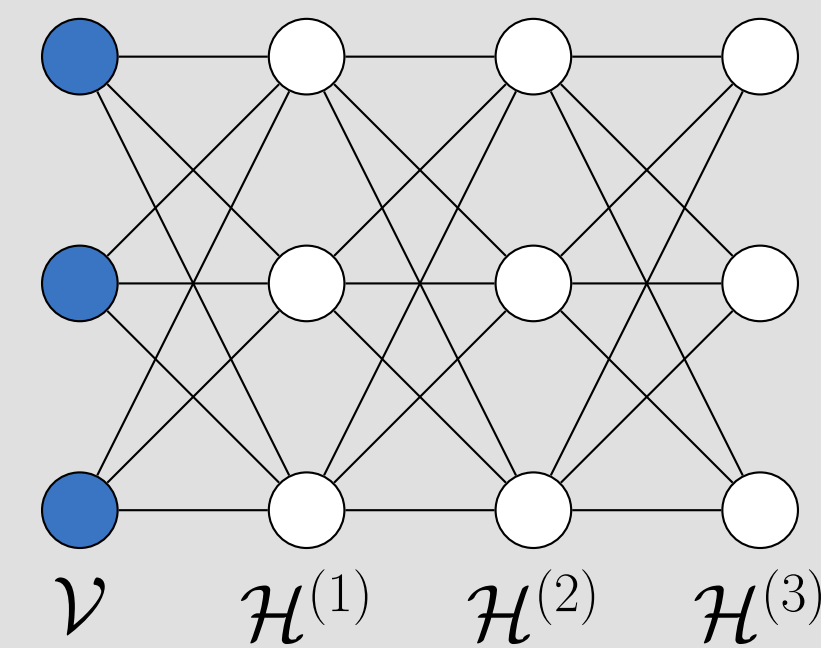


Figure 3: Deep RBM example.

Degeneracy, instability, and uninterpretability... Oh my!

The highly flexible nature of the RBM ($H+V+HV$ parameters) makes the following characteristics of model impropriety of particular concern.

| Characteristic | Detection |
|---|---|
| Near-degeneracy. Occurs when there is a disproportionate amount of probability placed on only a few elements of the sample space by the model [2]. | If the mean parametrization on the model parameters, $\mu(\theta)$, is close to the boundary of the convex hull of the set of statistics in the neg-potential function $Q(\mathbf{x})$. |
| Instability. Small changes in natural parameters result in large changes of the pmf, excessive sensitivity [5]. | If for any $C > 0$ there exists $N_C > 0$ such that $\max_{\mathbf{x}_N \in \mathcal{X}_N} [Q(\mathbf{x}_N)] > CN$ for all $N > N_C$. Where $Q(\cdot)$ is the neg-potential function of the model. |
| Uninterpretability. Due to the existence of dependence, marginal mean-structure no longer maintained [3]. | If the magnitude of the difference between model expectations and expectations under independence, $ \mathbb{E}(\mathbf{X} \theta) - \mathbb{E}(\mathbf{X} \emptyset) $, is large. |

Table 1: Table of "improper model" characteristics.

Avoiding degeneracy

For the $\{-1, 1\}$ encoding of \mathcal{V} and \mathcal{H} , the origin is the center of the parameter space. In particular, at $\theta = \mathbf{0}$, the RBM is equivalent to elements of \mathbf{X} being distributed as iid Bernoulli($\frac{1}{2}$) r.v.s. \Rightarrow No *near-degeneracy*, *instability*, or *uninterpretability*!

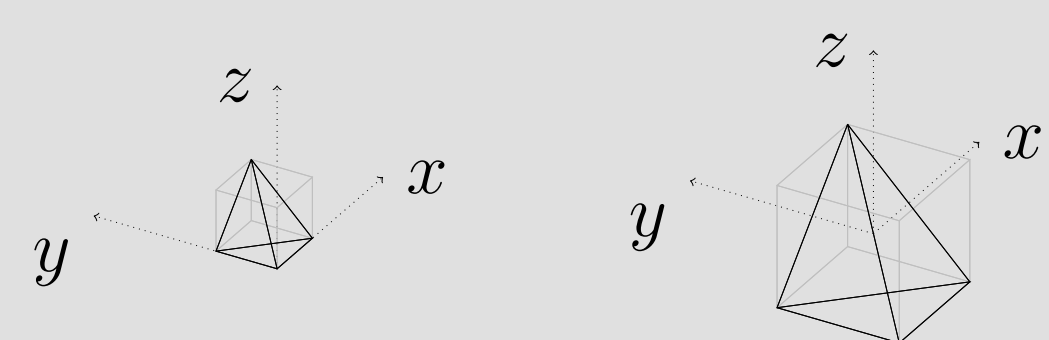


Figure 4: The convex hulls of the statistic space in three dimensions for a toy RBM with $|\mathcal{V}| = |\mathcal{H}| = 1$ for $\{0, 1\}$ -encoding (left) and $\{-1, 1\}$ -encoding (right) enclosed by an unrestricted hull of 3-space.

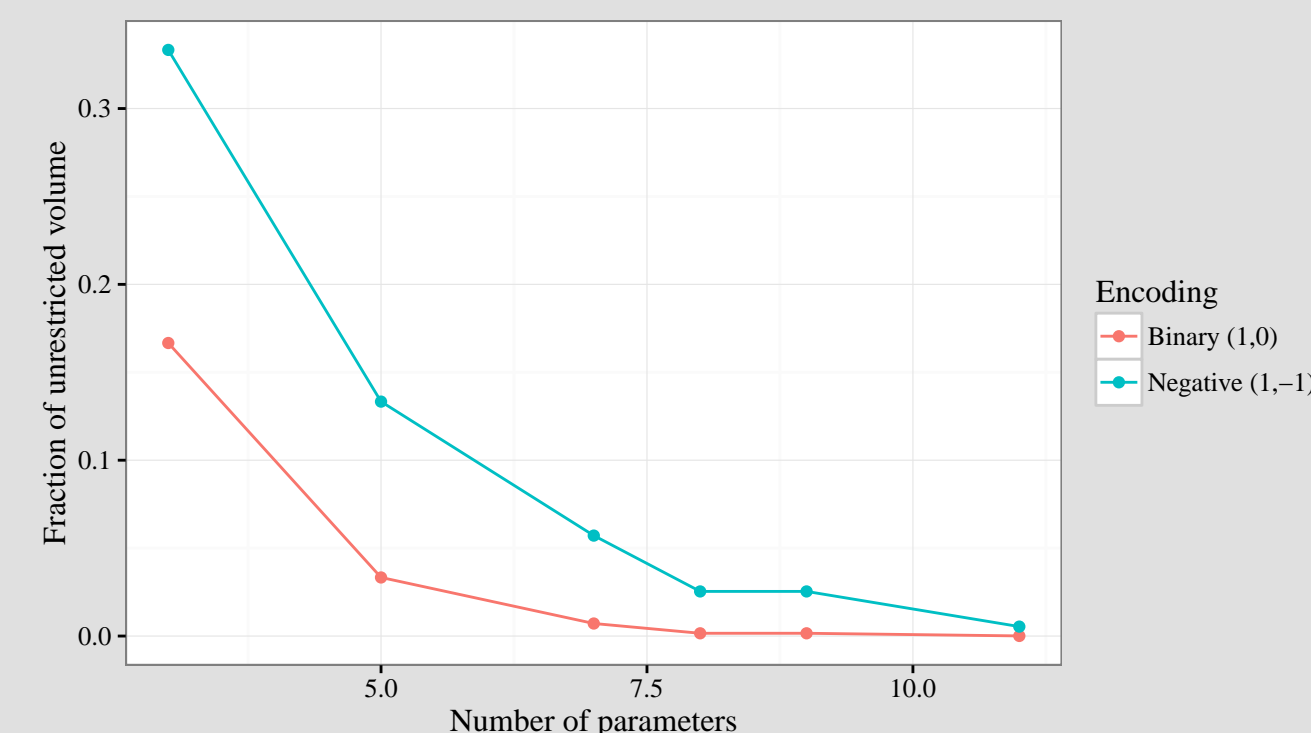


Figure 5: Volume relationship for the convex hulls of statistics in $Q(\cdot)$ vs. unrestricted space.

Restricted Boltzmann machine (RBM)

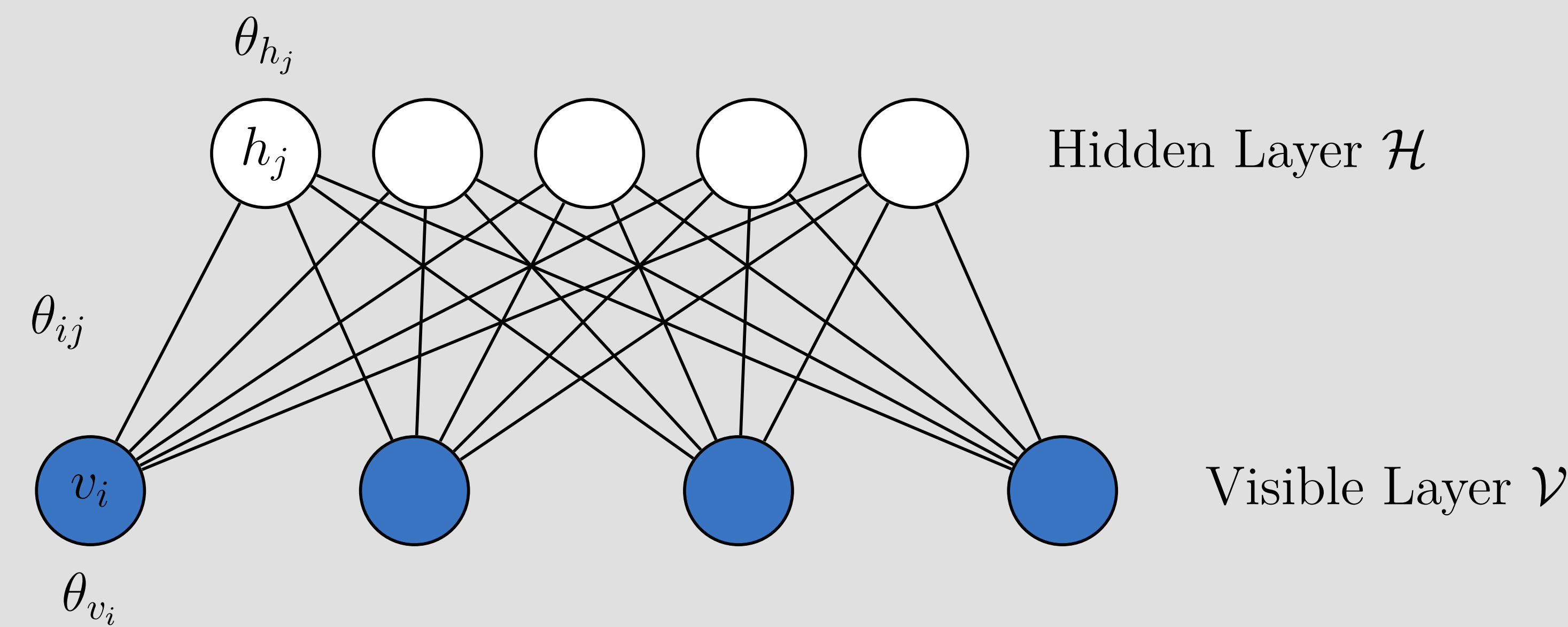


Figure 1: An example restricted Boltzmann machine (RBM), which consists of two layers, a hidden (\mathcal{H}) and a visible layer (\mathcal{V}), with no connections within a layer. Hidden nodes indicated by white circles and the visible nodes indicated by blue circles [1].

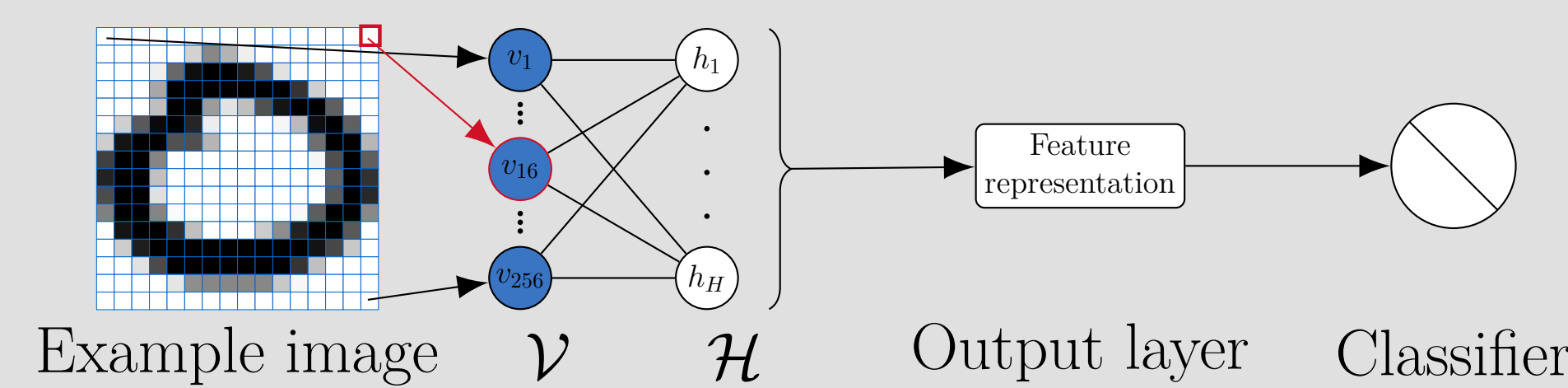


Figure 2: Image classification using a RBM. On the left, each image pixel comprises a node in the visible layer, \mathcal{V} . On the right, the output of the RBM is used to create features which are then passed to a supervised learning algorithm.

Joint distribution

Let $\mathbf{x} = \{h_1, \dots, h_H, v_1, \dots, v_V\}$ represent the states of the visible and hidden nodes in an RBM. Then the probability each node taking the the value corresponding to \mathbf{x} is:

$$f_{\theta}(\mathbf{x}) = \frac{\exp\left(\sum_{i=1}^V \sum_{j=1}^H \theta_{ij} v_i h_j + \sum_{i=1}^V \theta_{v_i} v_i + \sum_{j=1}^H \theta_{h_j} h_j\right)}{\sum_{\mathbf{x} \in \mathcal{X}} \exp\left(\sum_{i=1}^V \sum_{j=1}^H \theta_{ij} v_i h_j + \sum_{i=1}^V \theta_{v_i} v_i + \sum_{j=1}^H \theta_{h_j} h_j\right)} \quad (1)$$

References

- [1] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Vol. 1. Springer series in statistics Springer, Berlin, 2001.
- [2] Mark S Handcock et al. *Assessing degeneracy in statistical models of social networks*. Tech. rep. Working paper, 2003.
- [3] Mark S Kaiser. "Statistical Dependence in Markov Random Field Models". In: *Statistics Preprints Paper 57* (2007). URL: http://lib.dr.iastate.edu/stat_las_preprints/57.
- [4] Ruslan Salakhutdinov and Geoffrey E Hinton. "Deep boltzmann machines". In: *International Conference on Artificial Intelligence and Statistics*. 2009, pp. 448–455.
- [5] Michael Schweinberger. "Instability, sensitivity, and degeneracy of discrete exponential families". In: *Journal of the American Statistical Association* 106.496 (2011), pp. 1361–1370.

Manageable (a.k.a. small) examples

As the magnitude of θ grows, so does the occurrence of near-degeneracy, instability, and uninterpretability for RBMs of varying sizes.

Model fitting

Fitting is really hard to get right (Bayesian or ML). If you are able to fit this model properly, the result is the empirical distribution of images. You spent a lot of time just to get the empirical distribution back. If you do not have an instance of every possible image in your training set, certain parameter must $\rightarrow \infty$, making fitting more complicated (shrink) In reality it's practically impossible to have at least one of every possible image in your training set (large images, color scales, etc.)

Discussion

These models very easily are degenerate, unstable, and uninterpretable. Meaning, the space of possible fitted parameter values that leads to a proper model is highly restricted. To further complicate things, the proper fitting these models is very intricate and easily leads to parameter values running off to ∞ . If you're thinking of using RBMs, or stacked RBMs in a deep architecture, don't. Instead, fit a less flexible and more sensible model.