

A note on the instability and degeneracy of deep learning models

Andee Kaplan

Iowa State University
ajkaplan@iastate.edu

June 22, 2017

Slides available at <http://bit.ly/kaplan-private>

Joint work with D. Nordman and S. Vardeman

Introduction

- A probability model exhibits *instability* if small changes in a data outcome result in large changes in probability
- Model *degeneracy* implies placing all probability on a small portion of the sample space

Goal: Quantify instability for general probability models defined on sequences of observations, where each sequence of length N has a finite number of possible outcomes

Notation

- $\mathbf{X} = (X_1, \dots, X_N)$ a set of discrete random variables with a finite sample space, \mathcal{X}^N
- For each N , P_{θ_N} is a probability model on \mathcal{X}^N

FSFS models

Finitely Supported Finite Sequence (FSFS) model class

A series P_{θ_N} of probability models, indexed by a generic sequence of parameters θ_N , to describe data of each length $N \geq 1$ with model support of P_{θ_N} equaling the (finite) sample space \mathcal{X}^N .

- The size and structure of such parameters θ_N are without restriction
- Natural cases include $\theta_N \in \mathbb{R}^{q(N)}$ for some arbitrary integer-valued function $q(\cdot) \geq 1$

Discrete exponential family models

Exponential family model for \mathbf{X} with pmf of the form

$$p_{N,\lambda}(\mathbf{x}) = \exp \left[\boldsymbol{\eta}^T(\lambda) \mathbf{g}_N(\mathbf{x}) - \psi(\lambda) \right], \quad \mathbf{x} \in \mathcal{X}^N,$$

for fixed positive dimensions of the parameter, $\lambda \in \Lambda \subset \mathbb{R}^k$ and natural parameter $\boldsymbol{\eta} : \mathbb{R}^k \mapsto \mathbb{R}^L$ spaces, $\mathbf{g}_N : \mathcal{X}^N \mapsto \mathbb{R}^L$ a vector of sufficient statistics,

$$\psi(\lambda) = \log \sum_{\mathbf{x} \in \mathcal{X}^N} \exp \left[\boldsymbol{\eta}^T(\lambda) \mathbf{g}_N(\mathbf{x}) \right], \quad \lambda \in \Lambda,$$

the normalizing function, and $\Lambda = \{\lambda \in \mathbb{R}^k : \psi(\lambda) < \infty, k \leq q(N)\}$ is the parameter space.

Discrete exponential family models (cont'd)

- Such models arise with
 - Spatial data on a lattice (Besag 1974)
 - Network data (Wasserman and Faust 1994; Handcock 2003)
 - Binomial sampling with N iid Bernoulli random variables
- These models are special cases of the **FSFS models**
- $P_{\theta_N}(\mathbf{x}) \equiv p_{N,\lambda_N}(\mathbf{x})$ with $\theta_N = \lambda_N$ a sequence of elements of $\Lambda \subset \mathbb{R}^k$
- $P_{\theta_N}(\mathbf{x}) > 0$ for all $\mathbf{x} \in \mathcal{X}^N$
- The dimension of the parameter θ_N is the same for each N (k)
- Schweinberger (2011) considered *instability* in such exponential models

Restricted Boltzmann machines

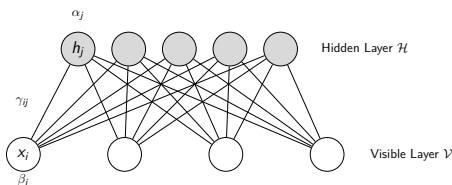


Figure 1: An example restricted Boltzmann machine (RBM). Hidden nodes are indicated by gray filled circles and the visible nodes indicated by unfilled circles.

- $\mathcal{X} = \{-1, 1\}$
- $\mathbf{X} = (X_1, \dots, X_N)$: N random variables for visibles with support \mathcal{X}^N
- $\mathbf{H} = (H_1, \dots, H_{N_H})$: N_H random variables for hiddens with support \mathcal{X}^{N_H}
- Parameters $\boldsymbol{\alpha} \in \mathbb{R}^{N_H}$, $\boldsymbol{\beta} \in \mathbb{R}^N$, Γ a matrix of size $N_H \times N$
 $(\boldsymbol{\theta}_N = (\boldsymbol{\alpha}, \boldsymbol{\beta}, \Gamma) \in \Theta_N \subset \mathbb{R}^{q(N)})$
with $q(N) = N + N_H + N * N_H$

Joint pmf:

$$P_{\boldsymbol{\theta}_N}(\tilde{\mathbf{x}}) = \exp \left[\boldsymbol{\alpha}^T \mathbf{h} + \boldsymbol{\beta}^T \mathbf{x} + \mathbf{h}^T \Gamma \mathbf{x} - \psi(\boldsymbol{\theta}_N) \right], \quad \tilde{\mathbf{x}} = (\mathbf{h}, \mathbf{x}) \in \mathcal{X}^{N+N_H}$$

Restricted Boltzmann machines (cont'd)

- The pmf for the visible variables X_1, \dots, X_N follows from marginalization:

$$P_{\theta_N}(\mathbf{x}) = \sum_{\mathbf{h} \in \mathcal{X}^{N_H}} P_{\theta_N}(\mathbf{x}, \mathbf{h}), \quad \mathbf{x} \in \mathcal{X}^N.$$

- Size of θ_N , $q(N)$, increases as a function of sample dimension N
- Can choose the number N_H of hidden variables to change with N (potentially increase)
- The RBM model specification for visibles is a **FSFS model**
- Models formed by marginalizing a base FSFS model (e.g., a type of exponential family model) is again a **FSFS model** class

Deep learning

- M the number of hidden layers included in the model
- $N_{(H,1)}, \dots, N_{((H,M))}$ the number of hidden variables within each hidden layer
- $\tilde{\mathbf{X}} = \{H_1^{(1)}, \dots, H_{N_{(H,1)}}^{(1)}, \dots, H_1^{(M)}, \dots, H_{N_{(H,M)}}^{(M)}, \mathbf{X}\}$
 - Hidden variables $\{H_i^{(j)} : i = 1, \dots, N_{(H,j)}, j = 1, \dots, M\}$
 - Visible variables $\mathbf{X} = (X_1, \dots, X_N)$
 - $H_i^{(j)}, X_k \in \mathcal{X} = \{-1, 1\}$ for all i, j, k

Two models with “deep architecture” that contain multiple hidden layers in addition to a visible layer of data

- 1 Deep Boltzmann machine (DBM)
- 2 Deep belief network (DBN)

Deep Boltzmann machine (DBM)

Stacked RBMs with conditional dependence between neighboring layers.

Joint pmf:

$$P_{\theta_N}(\tilde{\mathbf{x}}) = \exp \left[\sum_{i=1}^M \boldsymbol{\alpha}^{(i)T} \mathbf{h}^{(i)} + \boldsymbol{\beta}^T \mathbf{x} + \mathbf{h}^{(1)T} \boldsymbol{\Gamma}^{(0)} \mathbf{x} + \sum_{i=1}^{M-1} \mathbf{h}^{(i)T} \boldsymbol{\Gamma}^{(i)} \mathbf{h}^{(i+1)} - \psi(\boldsymbol{\theta}_N) \right],$$
$$\tilde{\mathbf{x}} = (\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(M)}, \mathbf{x}) \in \mathcal{X}^{N_{(H,1)} + \dots + N_{(H,M)} + N}$$

- Let $\boldsymbol{\theta}_N = (\boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(M)}, \boldsymbol{\beta}, \boldsymbol{\Gamma}^{(0)}, \dots, \boldsymbol{\Gamma}^{(M-1)}) \in \Theta_N \subset \mathbb{R}^{q(N)}$,
 $q(N) = N_{(H,1)} + \dots + N_{(H,M)} + N + N_{(H,1)} * N + N_{H,2} * H_{(H,1)} + \dots + N_{(H,M)} * H_{(H,M)-1}$
- The probability mass function for X_1, \dots, X_N follows from marginalization of the joint pmf
- The visible DBM model specification is an example of a **FSFS model**.

Deep belief network (DBN)

- **Similar** to a DBM: Multiple layers of latent random variables stacked in a deep architecture with no conditional dependence between layers
- **Difference:** all but the last stacked layer in a DBN are Bayesian networks (see Pearl 1985)
 - A “weight” parameter is placed on each interaction between visibles, X_1, \dots, X_N and the first layer of latent variables, $H_1^{(1)}, \dots, H_{N_{(H,1)}}^{(1)}$
 - The number $q(N)$ of components in the parameter vector is dependent on the dimension of the visibles
- For visibles X_1, \dots, X_N with support \mathcal{X}^N , a DBN is also a **FSFS model**

Instability results

Implications

Thank you

Questions?

- Slides – <http://bit.ly/kaplan-private>
- Contact
 - Email – ajkaplan@iastate.edu
 - Twitter – <http://twitter.com/andeekaplan>
 - GitHub – <http://github.com/andeek>

References I

Besag, Julian. 1974. "Spatial Interaction and the Statistical Analysis of Lattice Systems." *Journal of the Royal Statistical Society. Series B (Methodological)*. JSTOR, 192–236.

Handcock, Mark S. 2003. "Assessing Degeneracy in Statistical Models of Social Networks." Center for Statistics; the Social Sciences, University of Washington. <http://www.csss.washington.edu/>.

Pearl, Judea. 1985. "Bayesian Networks: A Model of Self-Activated Memory for Evidential Reasoning." UCLA Computer Science Department.

Schweinberger, Michael. 2011. "Instability, Sensitivity, and Degeneracy of Discrete Exponential Families." *Journal of the American Statistical Association* 106 (496). Taylor & Francis: 1361–70.

Wasserman, Stanley, and Katherine Faust. 1994. *Social Network Analysis: Methods and Applications*. Vol. 8. Cambridge: Cambridge University Press.