

This is the title

Author One

Author Two

29 February 2004

## 1 Energy-Based Model (EBM)

EBM mengaitkan sebuah energi skalar pada setiap konfigurasi variable yang diinginkan. Proses learning bertujuan untuk memodifikasi fungsi energi sehingga bentuknya memiliki sifat yang diinginkan. Sebagai contoh, misalnya diinginkan sebuah bentuk konfigurasi yang memiliki energi yang rendah, maka model probabilistik dari EBM didefinisikan sebagai distribusi probabilitas melalui fungsi energi sebagai berikut:

$$p(x) = \frac{e^{-E(x)}}{Z}. \quad (1)$$

$Z$  adalah faktor normalisasi yang disebut sebagai fungsi partisi untuk menganalogikan dengan sistem fisika.

$$Z = \sum_x e^{-E(x)}$$

EBM bisa dilatih dengan cara melakukan (stochastic) gradient descent pada negative log-likelihood (NLL)-nya secara empiris pada data training. Adapun untuk logistic regression akan didefinisikan terlebih dahulu log-likelihood  $\mathcal{L}(\theta, \mathcal{D})$  dan fungsi loss-nya sebagai NLL  $\ell(\theta, \mathcal{D})$  sebagai berikut:

$$\begin{aligned} \mathcal{L}(\theta, \mathcal{D}) &= \frac{1}{N} \sum_{x^{(i)} \in \mathcal{D}} \log p(x^{(i)}) \\ \ell(\theta, \mathcal{D}) &= -\mathcal{L}(\theta, \mathcal{D}) \end{aligned}$$

Menggunakan stochastic gradient  $-\frac{\partial \log p(x^{(i)})}{\partial \theta}$ , dimana  $\theta$  adalah parameter dari modelnya.

## 1.1 EBM dengan Hidden Units

Pada banyak kasus, sampel  $x$  biasanya tidak terobservasi secara penuh, atau akan ditambahkan variabel yang tidak teobservasi secara langsung yang disebut dengan hidden unit, dimana hal ini berguna untuk meningkatkan ekspresivitas dari model. Sehingga dikenalkan bagian yang terobservasi disini dilambangkan dengan  $x$ , dan sebuah bagian yang tersembunyi  $h$ . Sehingga bisa ditulis sebagai:

$$P(x) = \sum_h P(x, h) = \sum_h \frac{e^{-E(x, h)}}{Z}. \quad (2)$$

Pada kasus ini, untuk melakukan pemetaan rumus yang mirip dengan rumus (1), akan dikenalkan notasi (yang merupakan inspirasi dari fisika) yaitu free energy  $\mathcal{F}(x)$ , yang didefinisikan sebagai berikut:

$$\mathcal{F}(x) = -\log \sum_h e^{-E(x, h)} \quad (3)$$

Sehingga bisa diturunkan sebagai :

$$P(x) = \frac{e^{-\mathcal{F}(x)}}{Z} \text{ dengan } Z = \sum_x e^{-\mathcal{F}(x)}.$$

Data dari gradien NLL kemudian memiliki bentuk yang menarik yaitu:

$$-\frac{\partial \log p(x)}{\partial \theta} = \frac{\partial \mathcal{F}(x)}{\partial \theta} - \sum_{\tilde{x}} p(\tilde{x}) \frac{\partial \mathcal{F}(\tilde{x})}{\partial \theta}. \quad (4)$$

Notice that the above gradient contains two terms, which are referred to as the positive and negative phase. The terms positive and negative do not refer to the sign of each term in the equation, but rather reflect their effect on the probability density defined by the model. The first term increases the probability of training data (by reducing the corresponding free energy), while the second term decreases the probability of samples generated by the model.

It is usually difficult to determine this gradient analytically, as it involves the computation of  $E_P[\frac{\partial \mathcal{F}(x)}{\partial \theta}]$ . This is nothing less than an expectation over all possible configurations of the input  $x$  (under the distribution  $P$  formed by the model) !

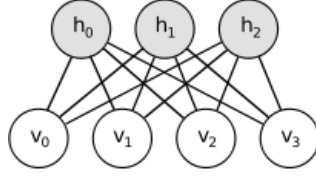


Figure 1: Grafik yang Menggambarkan RBM

The first step in making this computation tractable is to estimate the expectation using a fixed number of model samples. Samples used to estimate the negative phase gradient are referred to as negative particles, which are denoted as  $\mathcal{N}$ . The gradient can then be written as:

$$-\frac{\partial \log p(x)}{\partial \theta} \approx \frac{\partial \mathcal{F}(x)}{\partial \theta} - \frac{1}{|\mathcal{N}|} \sum_{\tilde{x} \in \mathcal{N}} \frac{\partial \mathcal{F}(\tilde{x})}{\partial \theta}. \quad (5)$$

where we would ideally like elements  $\tilde{x}$  of  $\mathcal{N}$  to be sampled according to  $P$  (i.e. we are doing Monte-Carlo). With the above formula, we almost have a practical, stochastic algorithm for learning an EBM. The only missing ingredient is how to extract these negative particles  $\mathcal{N}$ . While the statistical literature abounds with sampling methods, Markov Chain Monte Carlo methods are especially well suited for models such as the Restricted Boltzmann Machines (RBM), a specific type of EBM.

## 2 Restricted Boltzmann Machines (RBM)

Boltzmann Machines (BM) are a particular form of log-linear Markov Random Field (MRF), i.e., for which the energy function is linear in its free parameters. To make them powerful enough to represent complicated distributions (i.e., go from the limited parametric setting to a non-parametric one), we consider that some of the variables are never observed (they are called hidden). By having more hidden variables (also called hidden units), we can increase the modeling capacity of the Boltzmann Machine (BM). Restricted Boltzmann Machines further restrict BMs to those without visible-visible and hidden-hidden connections. A graphical depiction of an RBM is shown below.

The energy function  $E(v,h)$  of an RBM is defined as:

$$E(v, h) = -b'v - c'h - h'Wv \quad (6)$$

where  $W$  represents the weights connecting hidden and visible units and  $b, c$  are the offsets of the visible and hidden layers respectively.

This translates directly to the following free energy formula:

$$\mathcal{F}(v) = -b'v - \sum_i \log \sum_{h_i} e^{h_i(c_i + W_i v)}.$$

Because of the specific structure of RBMs, visible and hidden units are conditionally independent given one-another. Using this property, we can write:

$$p(h|v) = \prod_i p(h_i|v)$$

$$p(v|h) = \prod_j p(v_j|h).$$

## 2.1 RBMs with binary units

In the commonly studied case of using binary units (where  $v_j$  and  $h_i \in \{0, 1\}$ ), we obtain from Eq. (6) and (2), a probabilistic version of the usual neuron activation function:

$$P(h_i = 1|v) = \text{sigm}(c_i + W_i v) \quad (7)$$

$$P(v_j = 1|h) = \text{sigm}(b_j + W'_j h) \quad (8)$$

The free energy of an RBM with binary units further simplifies to:

$$\mathcal{F}(v) = -b'v - \sum_i \log(1 + e^{(c_i + W_i v)}). \quad (9)$$

## 2.2 Update Equations with Binary Units

Combining Eqs. (5) with (9), we obtain the following log-likelihood gradients for an RBM with binary units:

$$-\frac{\partial \log p(v)}{\partial W_{ij}} = E_v[p(h_i|v) \cdot v_j] - v_j^{(i)} \cdot \text{sigm}(W_i \cdot v^{(i)} + c_i) \quad (10)$$

$$-\frac{\partial \log p(v)}{\partial c_i} = E_v[p(h_i|v)] - \text{sigm}(W_i \cdot v^{(i)})$$

$$-\frac{\partial \log p(v)}{\partial b_j} = E_v[p(v_j|h)] - v_j^{(i)}$$

For a more detailed derivation of these equations, we refer the reader to the following page, or to section 5 of Learning Deep Architectures for AI. We will however not use these formulas, but rather get the gradient using Theano T.grad from equation (4).

### 3 Sampling in an RBM

Samples of  $p(x)$  can be obtained by running a Markov chain to convergence, using Gibbs sampling as the transition operator.

Gibbs sampling of the joint of  $N$  random variables  $S = (S_1, \dots, S_N)$  is done through a sequence of  $N$  sampling sub-steps of the form  $S_i \sim p(S_i|S_{-i})$  where  $S_{-i}$  contains the  $N - 1$  other random variables in  $S$  excluding  $S_i$ .

For RBMs,  $S$  consists of the set of visible and hidden units. However, since they are conditionally independent, one can perform block Gibbs sampling. In this setting, visible units are sampled simultaneously given fixed values of the hidden units. Similarly, hidden units are sampled simultaneously given the visibles. A step in the Markov chain is thus taken as follows:

$$h^{(n+1)} \sim \text{sigm}(W'v^{(n)} + c)$$

$$v^{(n+1)} \sim \text{sigm}(Wh^{(n+1)} + b),$$

where  $h^{(n)}$  refers to the set of all hidden units at the  $n - th$  step of the Markov chain. What it means is that, for example,  $h_i^{(n+1)}$  is randomly chosen to be 1 (versus 0) with probability  $\text{sigm}(W'_i v^{(n)} + c_i)$ , and similarly,  $v_j^{(n+1)}$  is randomly chosen to be 1 (versus 0) with probability  $\text{sigm}(W_j h^{(n+1)} + b_j)$ .

This can be illustrated graphically :

As  $t \rightarrow \infty$ , samples  $(v^{(t)}, h^{(t)})$  are guaranteed to be accurate samples of  $p(v, h)$ .

In theory, each parameter update in the learning process would require running one such chain to convergence. It is needless to say that doing so would be prohibitively expensive. As such, several algorithms have been

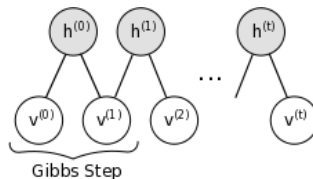


Figure 2: Gibbs Sampling

devised for RBMs, in order to efficiently sample from  $p(v, h)$  during the learning process.

## 4 Contrastive Divergence (CD-k)

Contrastive Divergence uses two tricks to speed up the sampling process:

since we eventually want  $p(v) \approx p_{train}(v)$  (the true, underlying distribution of the data), we initialize the Markov chain with a training example (i.e., from a distribution that is expected to be close to  $p$ , so that the chain will be already close to having converged to its final distribution  $p$ ). CD does not wait for the chain to converge. Samples are obtained after only  $k$ -steps of Gibbs sampling. In practice,  $k = 1$  has been shown to work surprisingly well.

## 5 Persistent CD

Persistent CD [Tieleman08] uses another approximation for sampling from  $p(v, h)$ . It relies on a single Markov chain, which has a persistent state (i.e., not restarting a chain for each observed example). For each parameter update, we extract new samples by simply running the chain for  $k$ -steps. The state of the chain is then preserved for subsequent updates.

The general intuition is that if parameter updates are small enough compared to the mixing rate of the chain, the Markov chain should be able to catch up to changes in the model.