

Иерархические модели

Московский Физико-Технический Институт

2022

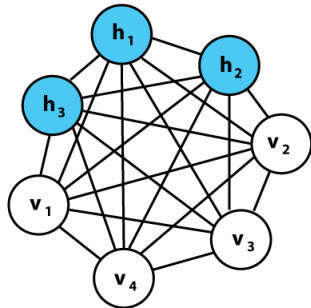
Машина Больцмана

Energy-based модель: $p(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}$,

$$E(\mathbf{x}) = -\mathbf{x}^T \mathbf{W} \mathbf{x} - \mathbf{w}_b^T \mathbf{x}.$$

Вариант со скрытыми переменными:

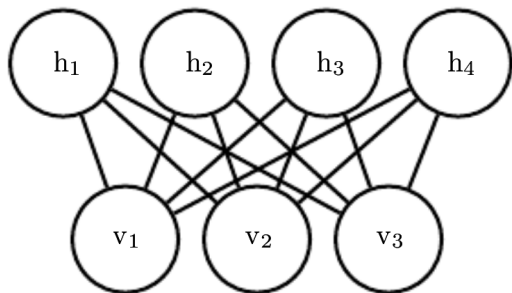
$$E(\mathbf{x}) = -\mathbf{x}^T \mathbf{W}_v \mathbf{x} - \mathbf{x}^T \mathbf{W}_{vh} \mathbf{h} - \mathbf{h}^T \mathbf{W}_h \mathbf{h} - \mathbf{w}_{bh}^T \mathbf{h} - \mathbf{w}_{bv}^T \mathbf{x}.$$



Ограниченная машина Больцмана

Частный случай машины Больцмана — модель представима в виде двудольного графа:

$$E(\mathbf{x}) = -\mathbf{h}^T \mathbf{W} \mathbf{h} - \mathbf{w}_{bh}^T \mathbf{h} - \mathbf{w}_{bb}^T \mathbf{x}.$$

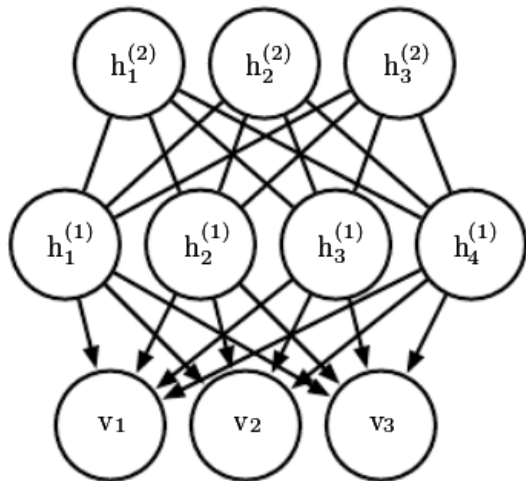


Свойства RBM

- Модель ненаправленная
- $p(\mathbf{h}|\mathbf{x}) = \prod_{j=1}^{n_h} \sigma((2\mathbf{h} - 1) \cdot (\mathbf{w}_{bh} + \mathbf{W}^T \mathbf{x}))_j$
- $p(\mathbf{x}|\mathbf{h}) = \prod_{j=1}^n \sigma((2\mathbf{x} - 1) \cdot (\mathbf{w}_{bv} + \mathbf{W}\mathbf{h}))_j$

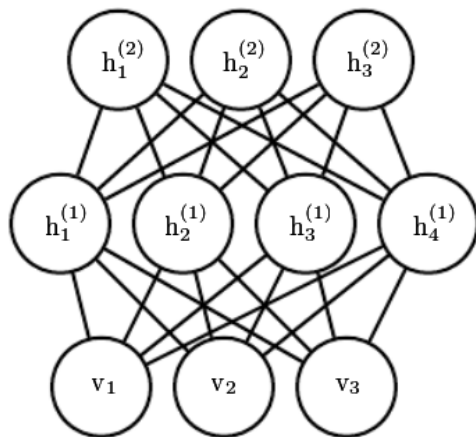
Глубокие сети доверия

- Стэк из нескольких RBM
- Оптимизируется послойно
- Модель направленная

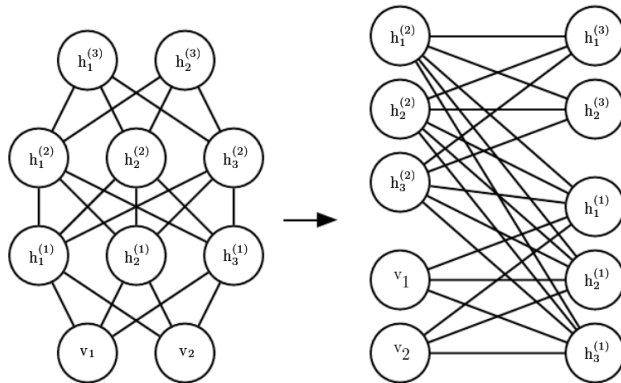


Глубокие машины Больцмана

- RBM с одним видимым слоем и несколькими скрытыми
- Модель ненаправленная

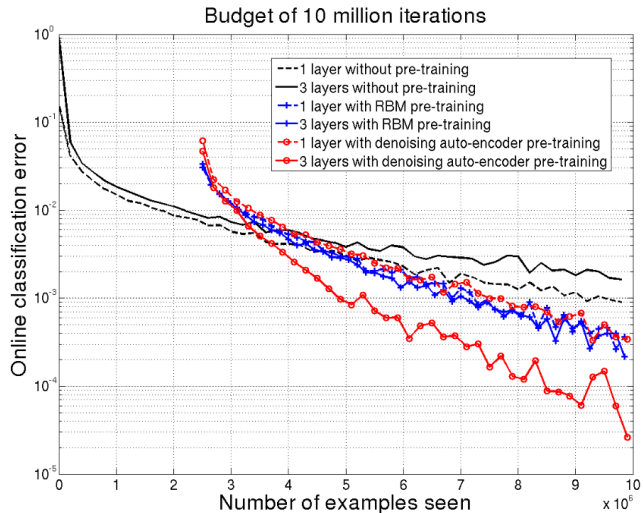


Глубокие машины Больцмана

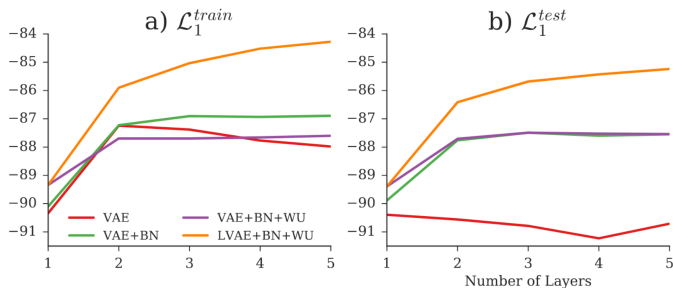


Обуславливаемся на одну долю графа → получаем независимые переменные во второй доле.

Жадное послойное обучение моделей



Жадное послойное обучение моделей: не всегда работает



Contrastive Divergence: идея

Energy-based model:

$$p(x|w) = \frac{\exp(-E_w(x))}{Z(w)}, \quad Z = \int_x \exp(-E_w(x)),$$

$$\frac{\partial \log p(x|w)}{\partial w} = E_{x'} \frac{\partial E(x')}{\partial w} - \frac{\partial E(x)}{\partial w}$$

Алгоритм для RBM:

- Берем x из выборки
- $h_0 \sim p(h_0|x)$
- $x_1 \sim p(x|h_0)$
- ...
- Получаем x_k
- $\frac{\partial \log p(x|w)}{\partial w} = \frac{\partial E(x_k)}{\partial w} - \frac{\partial E(x)}{\partial w}$

Дискриминативная модель как EBM

Our key observation in this work is that one can slightly re-interpret the logits obtained from f_θ to define $p(\mathbf{x}, y)$ and $p(\mathbf{x})$ as well. Without changing f_θ , one can re-use the logits to define an energy based model of the joint distribution of data point \mathbf{x} and labels y via:

$$p_\theta(\mathbf{x}, y) = \frac{\exp(f_\theta(\mathbf{x})[y])}{Z(\theta)}, \quad (5)$$

where $Z(\theta)$ is the unknown normalizing constant and $E_\theta(\mathbf{x}, y) = -f_\theta(\mathbf{x})[y]$.

By marginalizing out y , we obtain an unnormalized density model for \mathbf{x} as well,

$$p_\theta(\mathbf{x}) = \sum_y p_\theta(\mathbf{x}, y) = \frac{\sum_y \exp(f_\theta(\mathbf{x})[y])}{Z(\theta)}. \quad (6)$$

Notice now that the $\text{LogSumExp}(\cdot)$ of the logits of *any* classifier can be re-used to define the energy function at a data point \mathbf{x} as

$$E_\theta(\mathbf{x}) = -\text{LogSumExp}_y(f_\theta(\mathbf{x})[y]) = -\log \sum_y \exp(f_\theta(\mathbf{x})[y]). \quad (7)$$

Оптимизация:

$$\log p_\theta(\mathbf{x}, y) = \log p_\theta(\mathbf{x}) + \log p_\theta(y|\mathbf{x}).$$

второе слагаемое — обычная кросс-энтропия.

Как оптимизировать первое слагаемое?

Стохастическая динамика Ланжевена

Модификация стохастического градиентного спуска:

$$T = \mathbf{x} - \lambda \nabla_{\mathbf{x}} L(\mathbf{x}, \mathbf{w}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \frac{\lambda}{2})$$

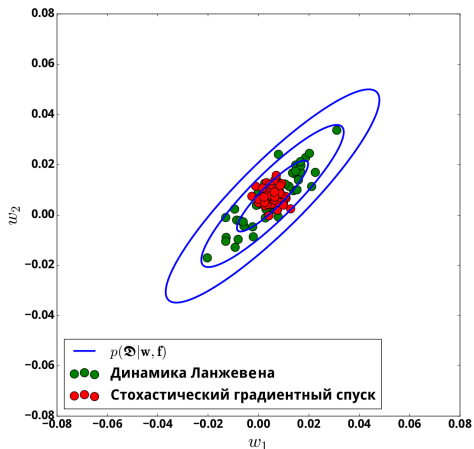
где шаг оптимизации λ изменяется с количеством итераций:

$$\sum_{\tau=1}^{\infty} \lambda_{\tau} = \infty, \quad \sum_{\tau=1}^{\infty} \lambda_{\tau}^2 < \infty.$$

Утверждение [Welling, 2011]. Распределение $T \circ T \circ \dots T$ сходится к распределению $p(\mathbf{x})$.

Стохастическая динамика Ланжевена

Распределение параметров после 2000 итераций:



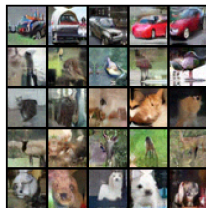
Алгоритм

Algorithm 1 JEM training: Given network f_θ , SGLD step-size α , SGLD noise σ , replay buffer B , SGLD steps η , reinitialization frequency ρ

```
1: while not converged do
2:   Sample  $\mathbf{x}$  and  $y$  from dataset
3:    $L_{\text{clf}}(\theta) = \text{xent}(f_\theta(\mathbf{x}), y)$ 
4:   Sample  $\hat{\mathbf{x}}_0 \sim B$  with probability  $1 - \rho$ , else  $\hat{\mathbf{x}}_0 \sim \mathcal{U}(-1, 1)$  ▷ Initialize SGLD
5:   for  $t \in [1, 2, \dots, \eta]$  do ▷ SGLD
6:      $\hat{\mathbf{x}}_t = \hat{\mathbf{x}}_{t-1} + \alpha \cdot \frac{\partial \text{LogSumExp}_{y'}(f_\theta(\hat{\mathbf{x}}_{t-1})[y'])}{\partial \hat{\mathbf{x}}_{t-1}} + \sigma \cdot \mathcal{N}(0, I)$ 
7:   end for
8:    $L_{\text{gen}}(\theta) = \text{LogSumExp}_{y'}(f(\mathbf{x})[y']) - \text{LogSumExp}_{y'}(f(\hat{\mathbf{x}}_t)[y'])$  ▷ Surrogate for Eq 2
9:    $L(\theta) = L_{\text{clf}}(\theta) + L_{\text{gen}}(\theta)$ 
10:  Obtain gradients  $\frac{\partial L(\theta)}{\partial \theta}$  for training
11:  Add  $\hat{\mathbf{x}}_t$  to  $B$ 
12: end while
```

Качество моделей

Class	Model	Accuracy% \uparrow	IS \uparrow	FID \downarrow
Hybrid	Residual Flow	70.3	3.6	46.4
	Glow	67.6	3.92	48.9
	IGEBM	49.1	8.3	37.9
	JEM $p(\mathbf{x} y)$ factored	30.1	6.36	61.8
	JEM (Ours)	92.9	8.76	38.4
Disc.	Wide-Resnet	95.8	N/A	N/A
Gen.	SNGAN	N/A	8.59	25.5
	NCSN	N/A	8.91	25.32



Литература

- Goodfellow I., Bengio Y., Courville A. Deep learning. – MIT press, 2016.
- Carreira-Perpinan M. A., Hinton G. On contrastive divergence learning //International workshop on artificial intelligence and statistics. – PMLR, 2005. – C. 33-40.
- Restricted Boltzmann Machine, a complete analysis:
<https://medium.com/datatype/restricted-boltzmann-machine-a-complete-analysis-part-3-contrastive-divergence-algorithm-3d06bbebb10c>
- Sønderby C. K. et al. Ladder variational autoencoders //Advances in neural information processing systems. – 2016. – T. 29.
- Grathwohl W. et al. Your Classifier is Secretly an Energy Based Model and You Should Treat it Like One. – 2020.
- Welling M., Teh Y. W. Bayesian learning via stochastic gradient Langevin dynamics //Proceedings of the 28th international conference on machine learning (ICML-11). – 2011. – C. 681-688.