

Стохастическое обратное распространение
ошибки и приближенный вариационный
вывод в глубинных генеративных моделях

2014

Stochastic Backpropogation

Задача оптимизации в моделях с латентными переменными:

$$\max_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)]$$

Stochastic Backpropagation

Задача оптимизации в моделях с латентными переменными:

$$\max_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)]$$

Пример:

$$\begin{aligned} \log p(v, \theta^r, \theta^g) &= \log \sum_h q(h|\theta^r) \log p(v|h, \theta^g) \\ &= \mathbb{E}_{q(h|\theta^r)}[\log p(v|h, \theta^g)] \end{aligned}$$

Stochastic Backpropogation

Задача оптимизации в моделях с латентными переменными:

$$\max_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)]$$

Пример:

$$\begin{aligned} \log p(v, \theta^r, \theta^g) &= \log \sum_h q(h|\theta^r) \log p(v|h, \theta^g) \\ &= \mathbb{E}_{q(h|\theta^r)}[\log p(v|h, \theta^g)] \end{aligned}$$

Проблема:

$$\nabla_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)] - ?$$

Likelihood Maximization

$$\begin{aligned}\log p(\mathbf{v}) &= \sum_{\mathbf{h}} \log p(\mathbf{v})q(\mathbf{h}) \\ &= \sum_{\mathbf{h}} q(\mathbf{h}) \log p(\mathbf{v}, \mathbf{h}) - \sum_{\mathbf{h}} q(\mathbf{h}) \log p(\mathbf{h}|\mathbf{v}) \\ &= \sum_{\mathbf{h}} q(\mathbf{h}) \log \frac{p(\mathbf{v}, \mathbf{h})}{q(\mathbf{h})} - \sum_{\mathbf{h}} q(\mathbf{h}) \log \frac{p(\mathbf{h}|\mathbf{v})}{q(\mathbf{h})} \\ &= \mathcal{L}(q) + \text{KL}(q||p) \\ &= \sum_{\mathbf{h}} \log p(\mathbf{v}, \mathbf{h})q(\mathbf{h}) + \mathcal{H}(q) + \text{KL}(q||p)\end{aligned}$$

Gaussian Backpropogation

$$\nabla_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)] - ?$$

Gaussian Backpropogation

$$\nabla_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)] - ?$$

Предположим, что $q = \mathcal{N}(\xi|\mu, C)$, тогда

$$\nabla_{\mu_i} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[f(\xi)] = \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[\nabla_{\xi_i} f(\xi)]$$

$$\nabla_{C_{i,j}} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[f(\xi)] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[\nabla_{\xi_i, \xi_j}^2 f(\xi)]$$

Gaussian Backpropogation

$$\nabla_{\theta} \mathbb{E}_{q(\xi|\theta)}[f(\xi)] - ?$$

Предположим, что $q = \mathcal{N}(\xi|\mu, C)$, тогда

$$\nabla_{\mu_i} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[f(\xi)] = \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[\nabla_{\xi_i} f(\xi)]$$

$$\nabla_{C_{i,j}} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[f(\xi)] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[\nabla_{\xi_i, \xi_j}^2 f(\xi)]$$

Если μ и C зависят от параметра θ , то

$$\nabla_{\theta} \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[f(\xi)] = \mathbb{E}_{\mathcal{N}(\xi|\mu, C)} \left[\mathbf{g}^T \frac{\partial \mu}{\partial \theta} + \frac{1}{2} \text{Tr} \left(\mathbf{H} \frac{\partial C}{\partial \theta} \right) \right],$$

где \mathbf{g} — градиент $f(\xi)$, \mathbf{H} — гессиан $f(\xi)$.

Gaussian Backpropogation

$$\begin{aligned}\nabla_{\mu_i} \mathbb{E}_{\mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C})}[f(\boldsymbol{\xi})] &= \int \nabla_{\mu_i} \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C}) f(\boldsymbol{\xi}) d\boldsymbol{\xi} \\ &= - \int \nabla_{\xi_i} \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C}) f(\boldsymbol{\xi}) d\boldsymbol{\xi} \\ &= - \left[\int \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C}) f(\boldsymbol{\xi}) d\boldsymbol{\xi}_{-i} \right]_{\xi_i=-\infty}^{\xi_i=+\infty} \\ &\quad + \int \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C}) \nabla_{\xi_i} f(\boldsymbol{\xi}) d\boldsymbol{\xi} \\ &= \mathbb{E}_{\mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C})}[\nabla_{\xi_i} f(\boldsymbol{\xi})]\end{aligned}$$

$$\nabla_{\mu_i} \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C}) = -\nabla_{\xi_i} \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C})$$

$$\nabla_{C_{i,j}} \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C}) = \frac{1}{2} \nabla_{\xi_i, \xi_j} \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \mathbf{C})$$

Stochastic Backpropagation

Какие распределения q можно использовать?

Два требования:

- $\nabla_{\theta} p(x|\theta) = \nabla_x p(x|\theta) B(x, \theta)$;
- $p(x, \theta) = 0$ на границах области определения x .

Экспоненциальное семейство:

$$p(x|\theta) = h(x) \exp(\eta(\theta)^T \phi(x) - A(\theta))$$

$$B(x, \theta) = \frac{\nabla_{\theta} \eta(\theta) \phi(x) - \nabla_{\theta} A(\theta)}{\nabla_x \log h(x) + \eta(\theta)^T \nabla_x \phi(x)}$$

Преобразование координат

$$\xi \sim \mathcal{N}(\mu, C)$$

$$\xi = \mu + R\epsilon, \text{ где } \epsilon \sim \mathcal{N}(0, I) \text{ и } C = RR^T$$

$$\begin{aligned}\nabla_R \mathbb{E}_{\mathcal{N}(\xi|\mu, C)}[f(\xi)] &= \nabla_R \mathbb{E}_{\mathcal{N}(\epsilon|0, I)}[f(\mu + R\epsilon)] \\ &= \mathbb{E}_{\mathcal{N}(\epsilon|0, I)}[\epsilon g^T],\end{aligned}$$

где g — градиент f в точке $\mu + R\epsilon$.

Deep Latent Gaussian Model

$$\xi_l \sim \mathcal{N}(\xi|0, \mathbf{I}), \quad l = 1, \dots, L;$$

$$\mathbf{h}_L = \mathbf{G}_L \xi_L$$

$$\mathbf{h}_l = T_l(\mathbf{h}_{l+1}) + \mathbf{G}_l \xi_l, \quad l = 1, \dots, L;$$

$$v \sim \pi(v|T_0(\mathbf{h}_l))$$

T_l — многослойный персептрон;

\mathbf{G}_l — матрица.

Параметры генеративной части θ^g — $(\{\mathbf{G}_l\}_{l=1}^K, \{T_l\}_{l=1}^L)$

Априорное распределение на параметры:

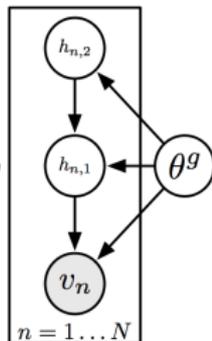
$$p(\theta^g) = \mathcal{N}(\theta|0, \kappa \mathbf{I})$$

Deep latent Gaussian Model

Совместное распределение можно записать двумя эквивалентными способами:

$$p(\mathbf{v}, \mathbf{h}) = p(\mathbf{v} | \mathbf{h}_1, \boldsymbol{\theta}^g) p(\mathbf{h}_L | \boldsymbol{\theta}^g) p(\boldsymbol{\theta}^g) \prod_{l=1}^{L-1} p(\mathbf{h}_l | \mathbf{h}_{l+1}, \boldsymbol{\theta}^g),$$

где $p(\mathbf{h}_l | \mathbf{h}_{l+1}, \boldsymbol{\theta}^g) = \mathcal{N}(\mathbf{h}_l | T_l(\mathbf{h}_{l+1}), \mathbf{G}_l \mathbf{G}_l^T)$



$$p(\mathbf{v}, \boldsymbol{\xi}) = p(\mathbf{v} | \mathbf{h}_1(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_L), \boldsymbol{\theta}^g) p(\boldsymbol{\theta}^g) \prod_{l=1}^L \mathcal{N}(\boldsymbol{\xi}_l | 0, \mathbf{I})$$

Lower bound on the marginal likelihood

\mathbf{V} — датасет.

$\mathbf{V} = \{\mathbf{v}_n\}_{n=1}^N$, где $\mathbf{v}_n = (v_{n,1}, \dots, v_{n,D})$

$$\begin{aligned} -\mathcal{L}(\mathbf{V}) &= -\log \int p(\mathbf{V}|\boldsymbol{\xi}, \boldsymbol{\theta}^g) p(\boldsymbol{\xi}, \boldsymbol{\theta}^g) \\ &= -\log \int \frac{q(\boldsymbol{\xi})}{q(\boldsymbol{\xi})} p(\mathbf{V}|\boldsymbol{\xi}, \boldsymbol{\theta}^g) p(\boldsymbol{\xi}, \boldsymbol{\theta}^g) \\ &\leq \mathcal{F}(\mathbf{V}) = KL[q(\boldsymbol{\xi})\|p(\boldsymbol{\xi})] - \mathbb{E}_{q(\boldsymbol{\xi})}[\log p(\mathbf{V}|\boldsymbol{\xi}, \boldsymbol{\theta}^g) p(\boldsymbol{\theta}^g)] \end{aligned}$$

Approximate posterior

$q(\xi|v)$ — распознающая часть модели.

$$q(\xi|V, \theta^r) = \prod_{n=1}^N \prod_{l=1}^L \mathcal{N}(\xi_{n,l} | \mu_l(v_n), C_l(v_n)),$$

где $\mu_l(\cdot)$ и $C_l(\cdot)$ — функции, представленные глубинными нейросетями. Будем обозначать параметры распределения q — θ^r .

Lower bound of marginal likelihood

$$KL[\mathcal{N}(\boldsymbol{\mu}, \mathbf{C}) \parallel \mathcal{N}(0, \mathbf{I})] = \frac{1}{2} (\text{Tr}(\mathbf{C}) + \|\boldsymbol{\mu}\|^2 - k - \log |\mathbf{C}|)$$

$$\begin{aligned} \mathcal{F}(\mathbf{V}) &= KL[q(\boldsymbol{\xi}) \parallel p(\boldsymbol{\xi})] - \mathbb{E}_{q(\boldsymbol{\xi})} [\log p(\mathbf{V} | \boldsymbol{\xi}, \boldsymbol{\theta}^g) p(\boldsymbol{\theta}^g)] \\ &= \frac{1}{2} \sum_{n,l} [\text{Tr}(\mathbf{C}_{n,l}) + \|\boldsymbol{\mu}_{n,l}\|^2 - k_l - \log |\mathbf{C}_{n,l}|] \\ &\quad - \sum_n \mathbb{E}_{q(\boldsymbol{\xi} | \mathbf{v}_n)} [\log p(\mathbf{v}_n | \mathbf{h}(\boldsymbol{\xi}))] + \frac{1}{2\kappa} \|\boldsymbol{\theta}^g\|^2 \end{aligned}$$

Gradient of lower bound

$$\begin{aligned}\mathcal{F}(\mathbf{V}) &= \frac{1}{2} \sum_{n,l} [\text{Tr}(\mathbf{C}_{n,l}) + \|\boldsymbol{\mu}_{n,l}\|^2 - k_l - \log |\mathbf{C}_{n,l}|] \\ &\quad - \sum_n \mathbb{E}_{q(\boldsymbol{\xi}|v_n)} [\log p(v_n|\mathbf{h}(\boldsymbol{\xi}))] + \frac{1}{2\kappa} \|\boldsymbol{\theta}^g\|^2\end{aligned}$$

$$\nabla_{\theta_j^g} \mathcal{F}(\mathbf{v}) = -\mathbb{E}_q \left[\nabla_{\theta_j^g} \log p(v|\mathbf{h}) \right] + \frac{1}{\kappa} \theta_j^g$$

Gradient of lower bound

$$\begin{aligned}\mathcal{F}(\mathbf{V}) &= \frac{1}{2} \sum_{n,l} [\text{Tr}(\mathbf{C}_{n,l}) + \|\boldsymbol{\mu}_{n,l}\|^2 - k_l - \log |\mathbf{C}_{n,l}|] \\ &\quad - \sum_n \mathbb{E}_{q(\boldsymbol{\xi}|v_n)} [\log p(v_n|\mathbf{h}(\boldsymbol{\xi}))] + \frac{1}{2\kappa} \|\boldsymbol{\theta}^g\|^2\end{aligned}$$

$$\nabla_{\boldsymbol{\mu}_l} \mathcal{F}(\mathbf{v}) = -\mathbb{E}_q [\nabla_{\boldsymbol{\xi}_l} \log p(v|\mathbf{h}(\boldsymbol{\xi}))] + \boldsymbol{\mu}_l$$

Gradient of lower bound

$$\begin{aligned}\mathcal{F}(\mathbf{V}) &= \frac{1}{2} \sum_{n,l} [\text{Tr}(\mathbf{C}_{n,l}) + \|\boldsymbol{\mu}_{n,l}\|^2 - k_l - \log |\mathbf{C}_{n,l}|] \\ &\quad - \sum_n \mathbb{E}_{q(\boldsymbol{\xi}|v_n)} [\log p(v_n|\mathbf{h}(\boldsymbol{\xi}))] + \frac{1}{2\kappa} \|\boldsymbol{\theta}^g\|^2\end{aligned}$$

$$\mathbf{C}_l = \mathbf{R}_l \mathbf{R}_l^T$$

$$\begin{aligned}\nabla_{R_{l,i,j}} \mathcal{F}(\mathbf{v}) &= -\mathbb{E}_q [\epsilon_{l,j} \nabla_{\xi_{l,i}} \log p(\mathbf{v}|\mathbf{h}(\boldsymbol{\xi}))] \\ &\quad + \frac{1}{2} \nabla_{R_{l,i,j}} [\text{Tr} \mathbf{C}_l - \log |\mathbf{C}_l|]\end{aligned}$$

Gradient of lower bound

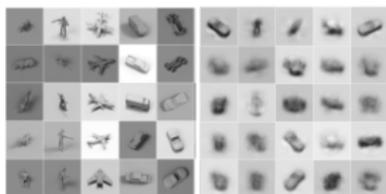
$$\begin{aligned}\mathcal{F}(\mathbf{V}) = & \frac{1}{2} \sum_{n,l} [\text{Tr}(\mathbf{C}_{n,l}) + \|\boldsymbol{\mu}_{n,l}\|^2 - k_l - \log |\mathbf{C}_{n,l}|] \\ & - \sum_n \mathbb{E}_{q(\boldsymbol{\xi}|v_n)} [\log p(v_n|\mathbf{h}(\boldsymbol{\xi}))] + \frac{1}{2\kappa} \|\boldsymbol{\theta}^g\|^2\end{aligned}$$

$$\nabla_{\boldsymbol{\theta}^r} \mathcal{F}(\mathbf{v}) = \nabla_{\boldsymbol{\mu}} \mathcal{F}(\mathbf{v})^T \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\theta}^r} + \text{Tr} \left(\nabla_{\mathbf{R}} \mathcal{F}(\mathbf{v}) \frac{\partial \mathbf{R}}{\partial \boldsymbol{\theta}^r} \right)$$

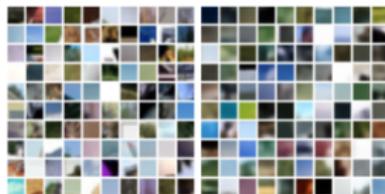
Алгоритм обучения

```
1: while hasNotConverged() do  
2:    $V \leftarrow \text{getMiniBatch}()$   
3:    $\xi_n \sim q(\xi_n | v_n)$   
4:    $h \leftarrow h(\xi)$   
5:    $\Delta \theta^{g,r} \leftarrow \text{calcGradients}()$   
6:    $\theta^{g,r} \leftarrow \theta^{g,r} + \Delta \theta^{g,r}$   
7: end while
```

Sampling from model



(a) NORB

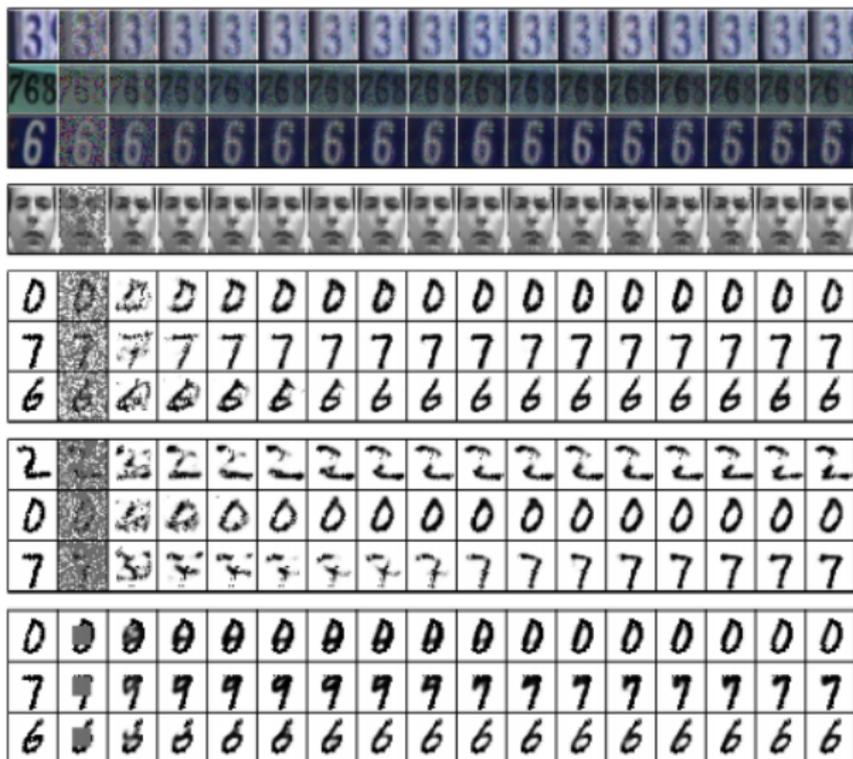


(b) CIFAR



(c) Frey

Imputation problem



Probability of MNIST test set

Table 1. Comparison of negative log-probabilities on the test set for the binarised MNIST data.

Model	$-\ln p(\mathbf{v})$
Factor Analysis	106.00
NLGBN (Frey & Hinton, 1999)	95.80
Wake-Sleep (Dayan, 2000)	91.3
DLGM diagonal covariance	87.30
DLGM rank-one covariance	86.60
<i>Results below from Uria et al. (2014)</i>	
MoBernoullis K=10	168.95
MoBernoullis K=500	137.64
RBM (500 h, 25 CD steps) approx.	86.34
DBN 2hl approx.	84.55
NADE 1hl (fixed order)	88.86
NADE 1hl (fixed order, RLU, minibatch)	88.33
EoNADE 1hl (2 orderings)	90.69
EoNADE 1hl (128 orderings)	87.71
EoNADE 2hl (2 orderings)	87.96
EoNADE 2hl (128 orderings)	85.10