

Байесовское мультимоделирование: вариационный вывод-2

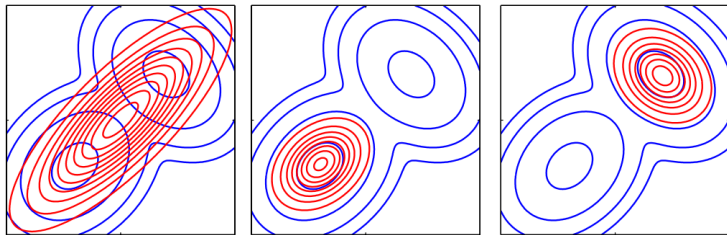
Московский Физико-Технический Институт

2021

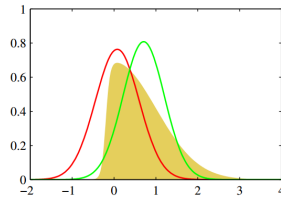
Вариационная оценка, ELBO

Вариационная оценка Evidence, Evidence lower bound — метод нахождения приближенного значения аналитически невычислимого распределения $p(\mathbf{w}|\mathcal{D}, \mathbf{h})$ распределением $q(\mathbf{w}) \in \mathcal{Q}$. Получение вариационной нижней оценки обычно сводится к задаче минимизации

$$\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathcal{D})}{q(\mathbf{w})} d\mathbf{w} = \mathbb{E}_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w}) - \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})).$$



Вариационный вывод и expectation propagation (Bishop)



Аппроксимация Лапласа и
вариационная оценка

Вариационный вывод: наивная интерпретация

Вариационная оценка обоснованности получается путем приближения заранее определенным **параметрическим непрерывным** распределением (**нормальным**) апостериорного распределения. Задача получения вариационной оценки **сводится к минимизации**

$$\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathcal{D})}{q(\mathbf{w})} d\mathbf{w} = \mathbb{E}_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w}) - \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})).$$

МСМС и вариационный вывод

Идея МСМС: Порождаем сэмплы из простого распределения и принимаем их, если заданное отношение больше порога:

$$\min \left(1, \frac{p(\mathbf{w}^\tau | \mathbf{y}, \mathbf{X}, \mathbf{h})}{p(\mathbf{w}^{\tau-1} | \mathbf{y}, \mathbf{X}, \mathbf{h})} \right),$$

где \mathbf{w}^τ выбирается на основе предыдущего сэмпла:

$$\mathbf{w}^\tau = T(\mathbf{w}^{\tau-1}).$$

Salimans et al., 2014: будем интерпретировать последовательность применения оператора T как оптимизацию вариационной оценки:

$$T^1 \circ \dots \circ T^n(\mathbf{w}) \rightarrow p(\mathbf{w}^\tau | \mathbf{y}, \mathbf{X}, \mathbf{h}).$$

Maclaurin et. al, 2015: в качестве оператора T будем рассматривать оператор оптимизации. Откажемся от отклонения сэмплов по порогу.

Оператор оптимизации, Maclaurin et. al, 2015

Определение

Назовем оператором оптимизации алгоритм T выбора вектора параметров \mathbf{w}' по параметрам предыдущего шага \mathbf{w} :

$$\mathbf{w}' = T(\mathbf{w}).$$

Определение

Пусть L — дифференцируемая функция потерь.

Оператором градиентного спуска назовем следующий оператор:

$$T(\mathbf{w}) = \mathbf{w} - \beta \nabla L(\mathbf{w}, \mathbf{y}, \mathcal{D}).$$

Градиентный спуск для оценки правдоподобия

Рассмотрим максимизацию совместного распределения параметров:

$$L = -\log p(\mathfrak{D}, \mathbf{w}|\mathbf{h}) = - \sum_{\mathfrak{D} \in \mathfrak{D}} \log p(\mathfrak{D}|\mathbf{w}, \mathbf{h})p(\mathbf{w}|\mathbf{h})$$

Проведем оптимизацию нейросети из r различных начальных приближений $\mathbf{w}_1, \dots, \mathbf{w}_r$ с использованием градиентного спуска:

$$\mathbf{w}' = T(\mathbf{w}).$$

Векторы параметров $\mathbf{w}_1, \dots, \mathbf{w}_r$ соответствуют некоторому скрытому распределению $q(\mathbf{w})$.

Энтропия

Формулу вариационной оценки можно переписать с использованием энтропии:

$$\log p(\mathcal{D}|\mathbf{f}) \geq \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathcal{D}, \mathbf{w}|\mathbf{h})}{q(\mathbf{w})} d\mathbf{w} = \\ \mathbb{E}_{q(\mathbf{w})}[\log p(\mathcal{D}, \mathbf{w}|\mathbf{h})] + S(q(\mathbf{w})),$$

где $S(q(\mathbf{w}))$ — энтропия:

$$S(q(\mathbf{w})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log q(\mathbf{w}) d\mathbf{w}.$$

Градиентный спуск для оценки правдоподобия

Утверждение 3

Пусть L — липшицева функция, оператор оптимизации — биекция. Тогда разность энтропии на различных шагах оптимизации вычисляется как:

$$S(q'(\mathbf{w})) - S(q(\mathbf{w})) \simeq \frac{1}{r} \sum_{g=1}^r (-\beta \text{Tr}[\mathbf{H}(\mathbf{w}'^g)] - \beta^2 \text{Tr}[\mathbf{H}(\mathbf{w}'^g)\mathbf{H}(\mathbf{w}'^g)]).$$

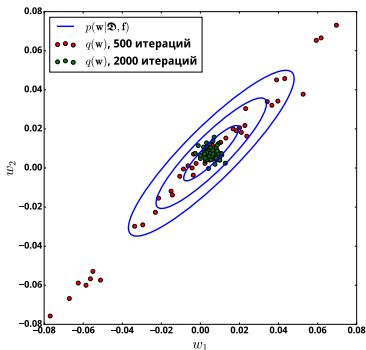
Итоговая оценка на шаге оптимизации τ :

$$\begin{aligned} \log \hat{p}(\mathbf{Y}|\mathcal{D}, \mathbf{h}) &\sim \frac{1}{r} \sum_{g=1}^r L(\mathbf{w}_{\tau}^g, \mathcal{D}, \mathbf{Y}) + S(q^0(\mathbf{w})) + \\ &+ \frac{1}{r} \sum_{b=1}^{\tau} \sum_{g=1}^r (-\beta \text{Tr}[\mathbf{H}(\mathbf{w}_b^g)] - \beta^2 \text{Tr}[\mathbf{H}(\mathbf{w}_b^g)\mathbf{H}(\mathbf{w}_b^g)]), \end{aligned}$$

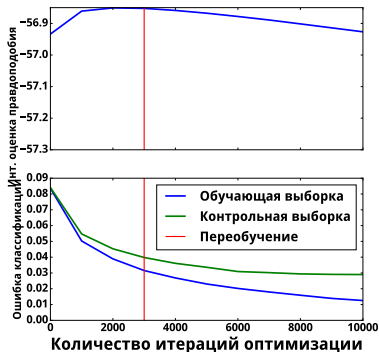
\mathbf{w}_b^g — вектор параметров старта g на шаге b , $S(q^0(\mathbf{w}))$ — начальная энтропия.

Переобучение, Maclaurin et. al, 2015

Градиентный спуск не минимизирует дивергенцию $KL(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}, \mathbf{h}))$. При приближении к моде распределения снижается оценка Evidence, что интерпретируется как переобучение модели.



Схождение распределения к моде



Оценка начала переобучения

Стохастическая динамика Ланжевена

Модификация стохастического градиентного спуска:

$$T = \mathbf{w} - \beta \nabla L + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \frac{\alpha}{2})$$

где шаг оптимизации α изменяется с количеством итераций:

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

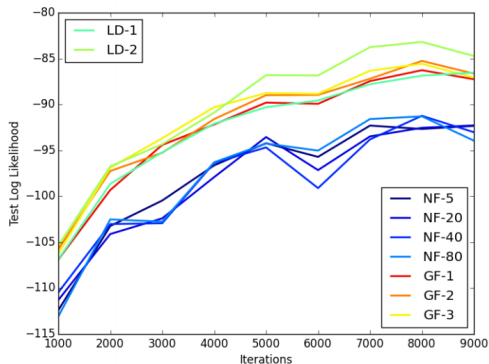
Утверждение [Welling, 2011]. Распределение $q^{\tau}(\mathbf{w})$ сходится к апостериорному распределению $p(\mathbf{w}|\mathbf{X}, \mathbf{f})$.

Изменение энтропии с учетом добавленного шума:

$$\hat{S}(q^{\tau}(\mathbf{w})) \geq \frac{1}{2} |\mathbf{w}| \log \left(\exp\left(\frac{2S(q^{\tau}(\mathbf{w}))}{|\mathbf{w}|}\right) + \exp\left(\frac{2S(\epsilon)}{|\mathbf{w}|}\right) \right).$$

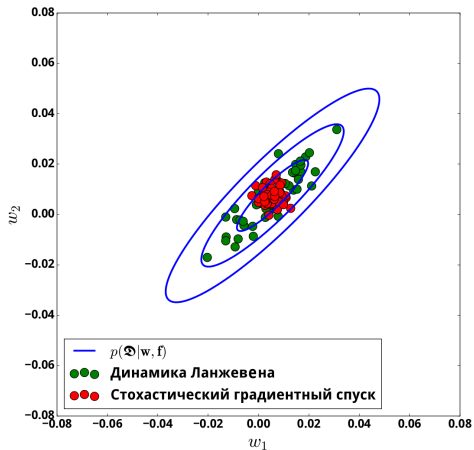
Стохастическая динамика Ланжевена в генеративных моделях

Altieri et al., 2015: будем сэмплировать скрытую переменную z и приближать его распределение к максимуму вариационной оценки с использованием динамики Ланжевена.



Стохастическая динамика Ланжевена

Распределение параметров после 2000 итераций:



Reparametrization trick: проблемы

Идея репараметризации:

$$\varepsilon = S_{\theta}(\mathbf{w}), \quad \mathbf{w} = S_{\theta}^{-1}(\varepsilon).$$

Тогда:

$$\nabla_{\theta} E_q f(\mathbf{w}) = E_q \nabla_{\theta} f(S_{\theta}^{-1}(\varepsilon)) = E_q \nabla_{\mathbf{w}} f(S_{\theta}^{-1}(\varepsilon)) \nabla_{\theta} S^{-1}(\varepsilon).$$

Пример:

$$w \sim \mathcal{N}(\mu, \sigma^2) \rightarrow S(w) = \frac{w - \mu}{\sigma} \sim \mathcal{N}(0, 1).$$

Проблема: часто получение S^{-1} вычислительно дорого.

Implicit reparametrization trick

$$\nabla_{\theta} E_q f(\mathbf{w}) = E_q \nabla_{\mathbf{w}} f(\mathbf{w}) \nabla_{\theta} \mathbf{w}.$$

Применим формулу полной производной к $\varepsilon = S_{\theta}(\mathbf{w})$:

$$\nabla_{\mathbf{w}} S_{\theta}(\mathbf{w}) \nabla_{\theta} \mathbf{w} + \nabla_{\theta} S_{\theta}(\mathbf{w}) = 0 \rightarrow$$

$$\rightarrow \nabla_{\theta} \mathbf{w} = -(\nabla_{\mathbf{w}} S_{\theta}(\mathbf{w}))^{-1} \nabla_{\theta} S_{\theta}.$$

Получили формулу без обратной функции к S .

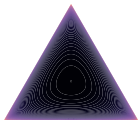
Для одномерных пространств в качестве S можно взять, например:

$$S(\mathbf{w}) = F(\mathbf{w}|\theta) \sim \mathcal{U}(0, 1).$$

Table 4: Test negative log-likelihood (lower is better) for VAE on MNIST. Mean \pm standard deviation over 5 runs. The von Mises-Fisher results are from [9].

Prior	Variational posterior	$D = 2$	$D = 5$	$D = 10$	$D = 20$	$D = 40$
$\mathcal{N}(0, 1)$	$\mathcal{N}(\mu, \sigma^2)$	131.1 ± 0.6	107.9 ± 0.4	92.5 ± 0.2	88.1 ± 0.2	88.1 ± 0.0
Gamma(0.3, 0.3)	Gamma(α, β)	132.4 ± 0.3	108.0 ± 0.3	94.0 ± 0.3	90.3 ± 0.2	90.6 ± 0.2
Gamma(10, 10)	Gamma(α, β)	135.0 ± 0.2	107.0 ± 0.2	92.3 ± 0.2	88.3 ± 0.2	88.3 ± 0.1
Uniform(0, 1)	Beta(α, β)	128.3 ± 0.2	107.4 ± 0.2	94.1 ± 0.1	88.9 ± 0.1	88.6 ± 0.1
Beta(10, 10)	Beta(α, β)	131.1 ± 0.4	106.7 ± 0.1	92.1 ± 0.2	87.8 ± 0.1	87.7 ± 0.1
Uniform($-\pi, \pi$)	vonMises(μ, κ)	127.6 ± 0.4	107.5 ± 0.4	94.4 ± 0.5	90.9 ± 0.1	91.5 ± 0.4
vonMises(0, 10)	vonMises(μ, κ)	130.7 ± 0.8	107.5 ± 0.5	92.3 ± 0.2	87.8 ± 0.2	87.9 ± 0.3
Uniform(S^D)	vonMisesFisher($\boldsymbol{\mu}, \kappa$)	132.5 ± 0.7	108.4 ± 0.1	93.2 ± 0.1	89.0 ± 0.3	90.9 ± 0.3

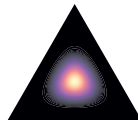
Дискретные распределения: релаксация



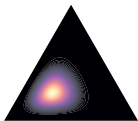
$$\bar{\alpha} = [1, 1, 1], t = 0.9$$



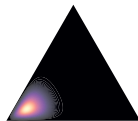
$$t = 1.0$$



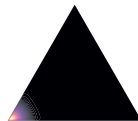
$$t = 10.0$$



$$\bar{\alpha} = [0.5, 0.25, 0.25], t = 30.0$$



$$[0.75, 0.125, 0.125]$$



$$[0.9, 0.05, 0.05]$$

Дискретные распределения в вариационном выводе

Релаксация:

- Распределение Дирихле (если использовать Implicit reparametrization trick)
- Gumbel-softmax:

$$p(\mathbf{w}) = \Gamma(k) \tau^{k-1} \left(\sum_{i=1}^k \alpha_i / w_i \right)^{-k} \prod_{i=1}^k (\alpha_i / w_i \tau + 1)$$

- ▶ Возможна репараметризация
- ▶ нет аналитической формы для KL
- Invertible Gaussian reparametrization:

$$p(\mathbf{w}) = \text{softmax}(\alpha), \quad \alpha \sim \mathcal{N},$$

(в знаменатель softmax добавляется константа для гарантии обратимости функции)

- ▶ Возможна репараметризация
- ▶ $KL(\mathbf{w}_1 | \mathbf{w}_2) = KL(\alpha_1 | \alpha_2)$
- ▶ Интерпретация параметров сильно сложнее, чем у GS или Дирихле

Локальная репараметризация

Пусть $y = \text{ReLU}(\mathbf{X}\mathbf{W})$ и матрица параметров \mathbf{W} распределена нормально:
 $w_{i,j} \sim \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$.

Тогда результатом линейной операции $\mathbf{X}\mathbf{W}$ будет гауссовая матрица:

$$\mathbf{G} = \mathbf{X}\mathbf{W}, \quad G_{i,j} \sim \mathcal{N}\left(\sum_k x_{i,k} \mu_{k,j}, \sum_k x_{i,k}^2 \sigma_{k,j}^2\right).$$

Вместо сэмлирования полноценного вектора параметров для каждого элемента батча на каждом шаге оптимизации, сэмплируем элементы из \mathbf{G} (то, что идет перед ReLU).

Если \mathbf{w} принимает значения из дискретного набора значений, то пользуясь ЦПТ (в формулировке Ляпунова):

$$\sum_k x_{i,k} w_{k,j} \sim \mathcal{N}(\cdot, \cdot).$$

Значит, к дискретным значениям также можно применить локальную репараметризацию.

Дивергенция Реньи

$$D_{\alpha}(p(\mathbf{w})||q(\mathbf{w})) = \frac{1}{\alpha - 1} \log \int p(\mathbf{w})^{\alpha} q(\mathbf{w})^{1-\alpha} d\mathbf{w}.$$

Table 1: Special cases in the Rényi divergence family.

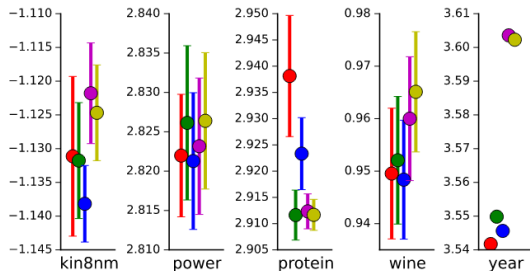
α	Definition	Notes
$\alpha \rightarrow 1$	$\int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}$	<i>Kullback-Leibler (KL) divergence</i> , used in VI ($\text{KL}[q p]$) and EP ($\text{KL}[p q]$)
$\alpha = 0.5$	$-2 \log(1 - \text{Hel}^2[p q])$	function of the square <i>Hellinger distance</i>
$\alpha \rightarrow 0$	$-\log \int_{p(\boldsymbol{\theta}) > 0} q(\boldsymbol{\theta}) d\boldsymbol{\theta}$	zero when $\text{supp}(q) \subseteq \text{supp}(p)$ (not a divergence)
$\alpha = 2$	$-\log(1 - \chi^2[p q])$	proportional to the χ^2 -divergence
$\alpha \rightarrow +\infty$	$\log \max_{\boldsymbol{\theta} \in \Theta} \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})}$	<i>worst-case regret</i> in <i>minimum description length principle</i>

mass-covering ←

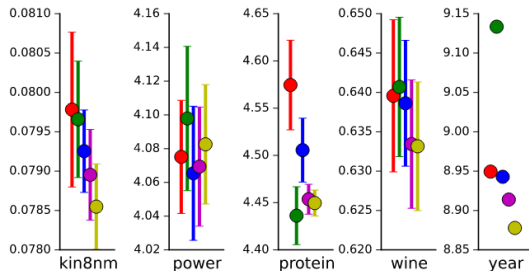
→ zero-forcing

● $\alpha \rightarrow -\infty$ (max) ● $\alpha=0.0$ ● $\alpha=0.5$ ● $\alpha=1.0$ (VI) ● $\alpha \rightarrow +\infty$

average negative test LL/nats



average test RMSE



Литература и прочие ресурсы

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – Т. 128. – №. 9.
- MacKay D. J. C., Mac Kay D. J. C. Information theory, inference and learning algorithms. – Cambridge university press, 2003.
- Salimans, Tim, Diederik Kingma, and Max Welling, 2015. Markov chain monte carlo and variational inference: Bridging the gap
- Altieri: <http://approximateinference.org/accepted/AltieriDuvenaud2015.pdf>
- Stephan Mandt, Matthew D. Hoffman, David M. Blei, 2017. Stochastic Gradient Descent as Approximate Bayesian Inference
- Бахтеев О. Ю., Стрижов В. В. Выбор моделей глубокого обучения субоптимальной сложности //Автоматика и телемеханика. – 2018. – №. 8. – С. 129-147.
- Figurnov M., Mohamed S., Mnih A. Implicit reparameterization gradients //arXiv preprint arXiv:1805.08498. – 2018.
- Jang E., Gu S., Poole B. Categorical reparameterization with gumbel-softmax //arXiv preprint arXiv:1611.01144. – 2016.
- Potapczynski A., Loaiza-Ganem G., Cunningham J. P. Invertible gaussian reparameterization: Revisiting the gumbel-softmax //arXiv preprint arXiv:1912.09588. – 2019.
- Maddison C. J., Mnih A., Teh Y. W. The concrete distribution: A continuous relaxation of discrete random variables //arXiv preprint arXiv:1611.00712. – 2016.
- Shayer O., Levi D., Fetaya E. Learning discrete weights using the local reparameterization trick //arXiv preprint arXiv:1710.07739. – 2017.
- Li Y., Turner R. E. Rényi Divergence Variational Inference //arXiv preprint arXiv:1602.02311. – 2016.