

# Байесовское мультимоделирование: вариационный вывод

Московский Физико-Технический Институт

2021

# Основа вариационного вывода

**Задача вариационного исчисления** — найти функцию, на которой заданный функционал достигает экстремального значения.

## Пример

Найти плотность распределения  $p$ , доставляющую максимум энтропии

$$H = - \int_w \log p(w) p(w) dw.$$

- $p$  — функция
- $H$  — функционал

Если функция задается из явно ограниченного семейства функций, то задачу вариационного исчисления можно рассматривать как задачу аппроксимации.

# Выбор модели: связанный байесовский вывод

Первый уровень: выбираем оптимальные параметры:

$$\mathbf{w} = \arg \max \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathcal{D}|\mathbf{h})},$$

Второй уровень: выбираем модель, доставляющую максимум обоснованности модели.

Обоснованность модели (“Evidence”):

$$p(\mathcal{D}|\mathbf{h}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})d\mathbf{w}.$$

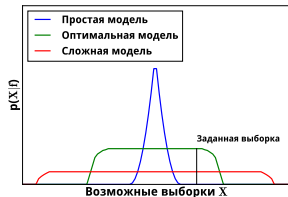
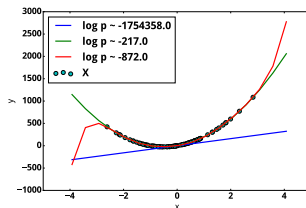


Схема выбора модели

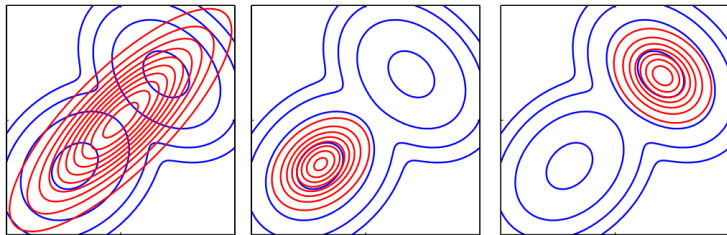


Пример: полиномы

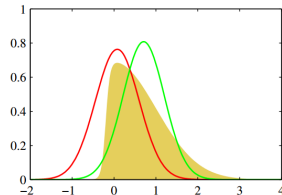
# Вариационная оценка, ELBO

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

# Получение вариационной нижней оценки

Максимизация вариационной нижней оценки

$$\int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \mathbf{h})}{q(\mathbf{w})} d\mathbf{w}$$

эквивалентна минимизации расстояния Кульбака–Лейблера между распределением  $q(\mathbf{w}) \in \mathfrak{Q}$  и апостериорным распределением параметров  $p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})$ :

$$\hat{q} = \arg \max_{q \in \mathfrak{Q}} \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \mathbf{h})}{q(\mathbf{w})} d\mathbf{w} \Leftrightarrow$$

$$\hat{q} = \arg \min_{q \in \mathfrak{Q}} D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})),$$

$$D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})) = \int_{\mathbf{w}} q(\mathbf{w}) \log \left( \frac{q(\mathbf{w})}{p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})} \right) d\mathbf{w}.$$

# Вариационная оценка и эффективный размер выборки

## Утверждение 2

Пусть  $m \gg 0$ ,  $\lambda > 0$ ,  $\frac{m}{\lambda} \in \mathbb{N}$ ,  $\frac{m}{\lambda} \gg 0$ . Тогда оптимизация функции

$$\mathbb{E}_q \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) - \lambda D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h}))$$

эквивалентна оптимизации вариационной оценки обоснованности для произвольной случайной подвыборки  $\hat{\mathbf{y}}, \hat{\mathbf{X}}$  мощности  $\frac{m}{\lambda}$  из генеральной совокупности.

См. также [ $\beta$ -VAE, Fixing Broken ELBO].

# Использование вариационной нижней оценки

Для чего используют вариационный вывод?

- получение оценок Evidence;
- получение оценок распределений моделей со скрытыми переменными (тематическое моделирование, снижение размерности).

Почему вариационный вывод?

- сводит задачу нахождения апостериорной вероятности к методам оптимизации;
- проще масштабируется, чем аппроксимация Лапласа;
- проще в использовании, чем сэмплирующие методы.

Вариационный вывод может давать сильно заниженную оценку.

## ELBO: нормальное распределение

Пусть  $q \sim \mathcal{N}(\mu_q, \mathbf{A}_q)$ .

Тогда вариационная оценка имеет вид:

$$\int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \mathbf{h}) d\mathbf{w} - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) \simeq$$
$$\sum_{i=1}^m \log p(\mathbf{y}_i|\mathbf{x}_i, \hat{\mathbf{w}}) - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) \rightarrow \max_{\mathbf{A}_q, \mu_q}, \quad \hat{\mathbf{w}} \sim q.$$

В случае, если априорное распределение параметров  $p(\mathbf{w}|\mathbf{h})$  является нормальным:

$$p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(\mu, \mathbf{A}),$$

дивергенция  $D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h}))$  вычисляется аналитически:

$$D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) = \frac{1}{2}(\text{tr}(\mathbf{A}^{-1}\mathbf{A}_q) + (\mu - \mu_q)^{\text{T}}\mathbf{A}^{-1}(\mu - \mu_q) - n + \ln |\mathbf{A}| - \ln |\mathbf{A}_q|).$$



# Graves, 2011

Априорное распределение:  $p(\mathbf{w}|\sigma) \sim \mathcal{N}(\boldsymbol{\mu}, \sigma \mathbf{I})$ .

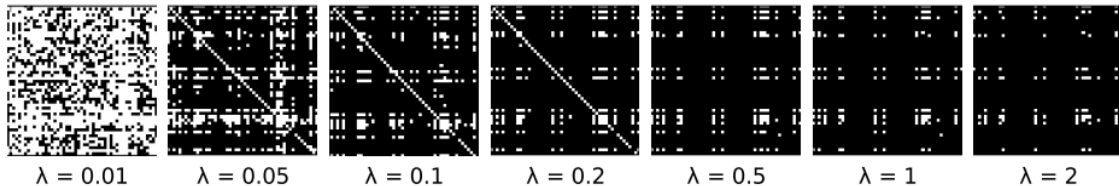
Вариационное распределение:  $q(\mathbf{w}) \sim \mathcal{N}(\boldsymbol{\mu}_q, \sigma_q \mathbf{I})$ .

Жадная оптимизация гиперпараметров:

$$\boldsymbol{\mu} = \hat{\mathbf{E}} \mathbf{w}, \quad \sigma = \hat{\mathbf{D}} \mathbf{w}.$$

Прунинг параметра  $w_i$  определяется относительной плотностью:

$$\lambda = \frac{q(0)}{q(\boldsymbol{\mu}_{i,q})} = \exp\left(-\frac{\mu_i^2}{2\sigma_i^2}\right).$$



# ELBO: нормальное распределение

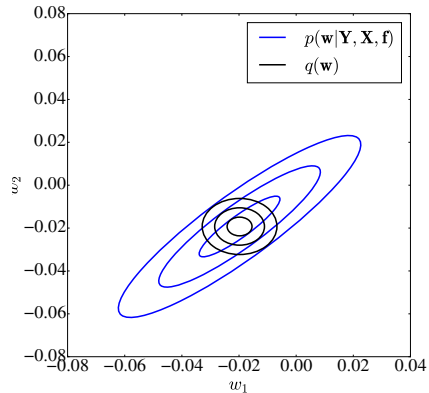
“Обычная” функция потерь:

$$L = \sum_{\mathbf{x}, \mathbf{y} \in \mathcal{D}} -\log p(\mathbf{y}|\mathbf{x}, \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2.$$

Вариационный вывод при  
( $p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(0, 1)$ ):

$$L = \sum_{\mathbf{x}, \mathbf{y}} \log p(\mathbf{y}|\mathbf{x}, \hat{\mathbf{w}}) + \\ + \frac{1}{2} (\text{tr}(\mathbf{A}_q) + \boldsymbol{\mu}_q^\top \mathbf{A}^{-1} \boldsymbol{\mu}_q - \ln |\mathbf{A}_q|).$$

Пример грубой аппроксимации нормальным  
диагональным распределением  $q$



# Локальная репараметризация

Как считать  $E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w})$ ?

- Graves, 2011: сэмплируем одну реализацию параметра на каждом шаге оптимизации. Для сэмплирования пользуемся свойством:

$$w \sim \mathcal{N}(\mu, \sigma^2) \rightarrow w \sim \varepsilon \sigma + \mu, \quad \varepsilon \sim \mathcal{N}(0, 1).$$

- ▶ Плохое приближение матожидания
- Наивный вариант: засэмплировать столько реализаций параметра, сколько у нас объектов в батче
  - ▶ BackProp будет очень медленным

# Локальная репараметризация, Kingma et al., 2015

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:

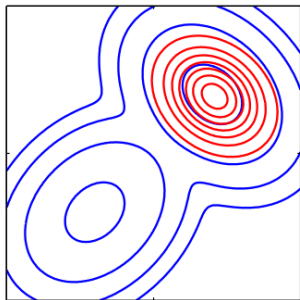
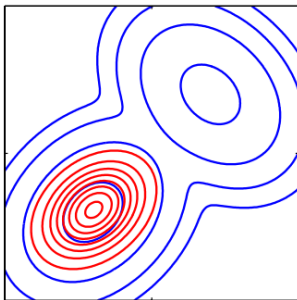
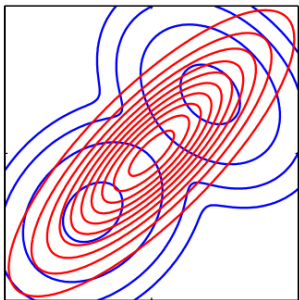
$$\mathbf{G} = \mathbf{X}\mathbf{W}, \quad G_{i,j} \sim \mathcal{N}\left(\sum_k x_{i,k} \mu_{k,j}, \sum_k x_{i,k}^2 \sigma_{k,j}^2\right).$$

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

## Пример

Пусть размер батча = 64, матрица  $\mathbf{W}$  имеет размер  $64 \times 64$ .

- Graves: сэмплируем параметры один раз,  $64 \times 64 = 4096$  элемента. Приближение матожидания одной реализацией.
- Наивный вариант: сэмплируем параметры 64 раза,  $64 \times 64 \times 64 = 262144$  элемента. Приближение матожидания 64 реализациями.
- Локальная репараметризация: сэмплируем  $\mathbf{G}$ ,  $64 \times 64 = 4096$  элемента. Приближение матожидания одной реализациями.



# Expectation propagation

Minka, 2001: представим апостериорное и аппроксимирующее распределение через произведение факторов:  $p(\mathbf{w}|\mathcal{D}) = \prod_i f_i$ ,  $q(\mathbf{w}) = \prod_i \tilde{f}_i$ .

Основная идея метода — минимизация  $KL(p(\mathbf{w}|\mathcal{D})||q(\mathbf{w}))$ .

- Выбираем фактор  $\tilde{f}_i$ , который будем приближать, «Убираем» его из рассмотрения и заменяем на истинный фактор:

$$q^i \propto f_i \prod_{j \neq i} \tilde{f}_j$$

- Приравниваем моменты  $q^i$  к моментам приближенного распределения (корректное приближение, если  $q$  из экспоненциального семейства)
- Повторять, пока не сойдемся

# Expectation propagation: плюсы и минусы

## Минусы:

- Вводится предположение на апостериорное распределение (достаточно мягкое)
- В оригинальной версии работает только для  $q$  из экспоненциального семейства
- Нет гарантий сходимости

## Плюсы:

- Приближает KL, а не его нижнюю оценку

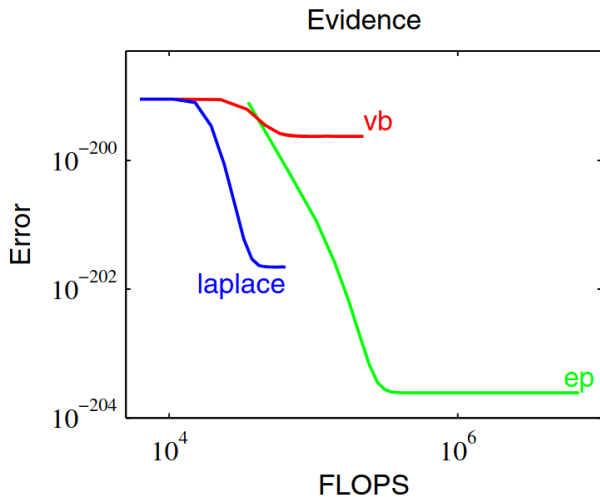


График для модели 2-х компонентной гауссовой смеси.



# Probabilistic backpropagation

Комбинация Expectation propagation и backpropagation.

**Backward pass:**

Обновляем параметры для каждого параметру по правилу Байеса:

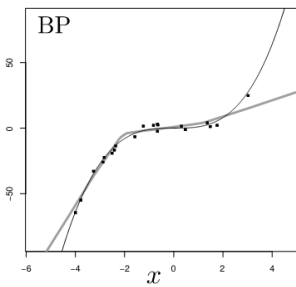
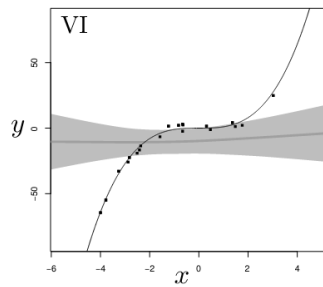
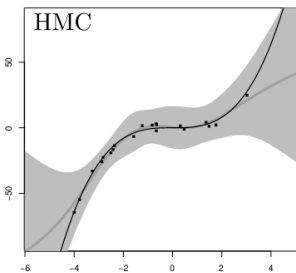
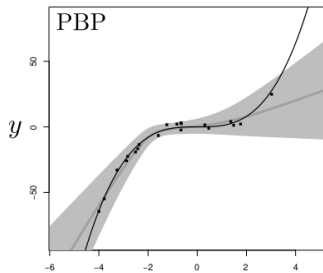
$$p(w_i|\mathcal{D}) = Z^{-1}p(\mathcal{D}|w_i, \mathbf{w}^i)p(w) \rightarrow p(w_i|\mathcal{D}) = Z^{-1}p(\mathcal{D}|w_i, \mathbf{w}^i)\mathcal{N}(w|\mu, \sigma^2).$$

**Проблема:** вычисление  $Z$ .

**Backward pass:**

$Z$  можно посчитать приближенно, в предположении  $\mathbf{f}(\mathbf{x}, \mathbf{w}) \sim \mathcal{N}(\mathbf{m}, \mathbf{v})$ .

Для подсчета  $\mathbf{m}, \mathbf{v}$  для сетей с ReLU-активацией в статье предлагается итеративный алгоритм.



# 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.
- Бахтеев О. Ю., Стрижов В. В. Выбор моделей глубокого обучения субоптимальной сложности //Автоматика и телемеханика. – 2018. – №. 8. – С. 129-147.
- Graves A. Practical variational inference for neural networks //Advances in neural information processing systems. – 2011. – Т. 24.
- Louizos C., Ullrich K., Welling M. Bayesian compression for deep learning //arXiv preprint arXiv:1705.08665. – 2017.
- Kingma D. P., Salimans T., Welling M. Variational dropout and the local reparameterization trick //Advances in neural information processing systems. – 2015. – Т. 28. – С. 2575-2583.
- Higgins I. et al. beta-vae: Learning basic visual concepts with a constrained variational framework. – 2016.
- Alemi A. et al. Fixing a broken ELBO //International Conference on Machine Learning. – PMLR, 2018. – С. 159-168.
- Minka T. P. Expectation propagation for approximate Bayesian inference //arXiv preprint arXiv:1301.2294. – 2013.
- Hernández-Lobato J. M., Adams R. Probabilistic backpropagation for scalable learning of bayesian neural networks //International conference on machine learning. – PMLR, 2015. – С. 1861-1869.