

Метод Монте-Карло

Московский Физико-Технический Институт

2022

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”):

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

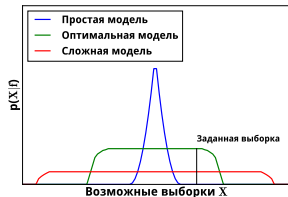
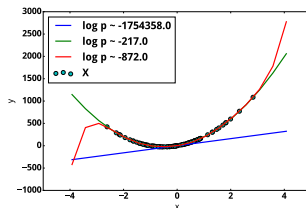


Схема выбора модели



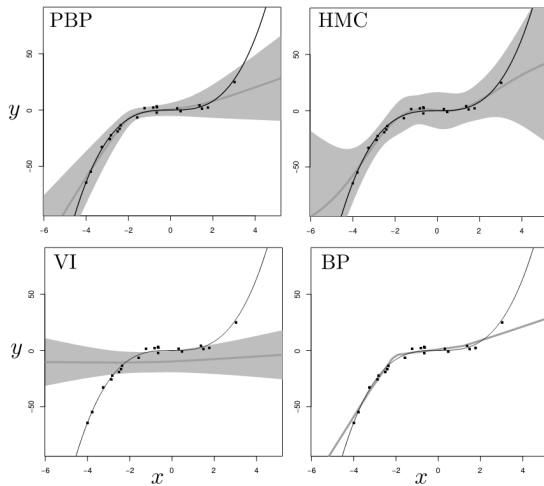
Пример: полиномы

Оценки интеграла

$$Ef = \int_{\mathbf{w}} f(\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$

- Аппроксимация Лапласа
 - ▶ Негибкий, аппроксимирует нормальным распределением
 - ▶ Есть проблемы с большими размерностями пространства
- Вариационный вывод
 - ▶ Гибкий, позволяет аппроксимировать неизвестное распределением широким классом распределений
 - ▶ Нижняя оценка интеграла, при некорректном выборе вариационного распределения может давать сильно заниженную оценку
- МС
 - ▶ Гибкий
 - ▶ Точный
 - ▶ Медленный

VI vs MC



Наивный метод

$$I = \mathbb{E}f = \int_{\mathbf{w}} f(\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$

Аппроксимируем:

$$\hat{I} = \frac{1}{N} \sum_{\mathbf{w} \sim p(\mathbf{w})} f(\mathbf{w}).$$

Свойства

Оценка интеграла:

- Сильно состоятельна: $\hat{I} \xrightarrow{\text{п.н.}} I$
- Несмещена: $E\hat{I} = I$
- Асимптотически нормальна;
- $D\hat{I} = O(\frac{1}{N})$.
- **Проблема:** нужно уметь сэмплировать из p .

Наивный метод

Пусть существует обратимая функция T из $\mathcal{U}(0, 1)$ в некоторое распределение w . Тогда

$$F_w(t) = p(w \leq t) = p(T(w') \leq t) = p(w' \leq T^{-1}(t)) = T^{-1}(w).$$

Отсюда $F_w^{-1} = T$.

Пример

$$w = \lambda \exp(-\lambda t).$$

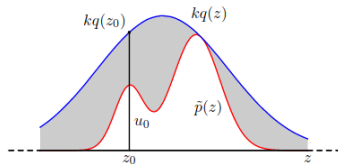
$$F_w(t) = 1 - \exp(-\lambda t).$$

$$F_w^{-1}(t') = -1 \frac{1}{\lambda} \log(1 - t').$$

Сэмплирование с отклонением

- Задана плотность $p(w)$ (может быть задана с точностью до нормировочной константы)
- Введем распределение q
- Подберем множитель k таким образом, чтобы $kq(w) \geq p(z)$ для всех z
- В цикле
 - ▶ Просэмплируем $w_0 \sim q$
 - ▶ Просэмплируем $u \sim \mathcal{U}(0, kq(w_0))$
 - ▶ Если $u \leq p(w_0)$ — считать его сэмплом из $p(w)$

Идея метода: сэмплы u равномерно распределены в регионе, ограниченном кривой $p(w)$.



Bishop, 2006

Сэмплирование по значимости

Пусть мы не можем сэмплировать из $p(w)$, но можем оценивать правдоподобие в каждой точке, и хотим получить интерал

$$Ef = \int f(w)p(w)dw.$$

Тогда введем распределение q :

$$Ef = \int f(w)p(w)dw = \int f(w)\frac{p(w)}{q(w)}dz \approx \frac{1}{L} \sum_{l=1}^L \frac{p(w^l)}{q(w^l)} f(w^l).$$

Nested sampling

Приближение интеграла вида:

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

- В базовом виде — приближение через априорное распределение $p(\mathbf{w})$
- Можно ли сэмплировать из апостериорное распределение?

Идея:

- Введем $I_t = \int_{\mathbf{w}:p(\mathcal{D}|\mathbf{w})>t} p(\mathbf{w})d\mathbf{w}$
- Вычисление интеграла сводится к интегрированию по контуру: $I = \int_0^\infty I_t dt$.

МСМС

Основная идея: Сэмплируем аналогично сэмплированию с отклонениями, но q — марковское распределение, обусловленное на предыдущий успешный шаг
Хотим, чтобы предельное (стационарное) распределение соответствовало нашему распределению $p(w)$.
Достаточное условие

$$p(w) T(w|w') = p(w') T(w'|w).$$

Алгоритм Метрополиса — Гастингса

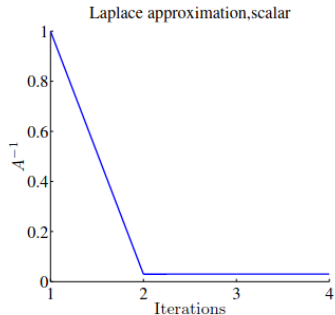
- Сэмплируем новое значение $w' \sim q(w|w^t)$.
- Принимаем его с вероятностью $A(w'|w^t) = \min \left(1, \frac{p(w')q(w^t|w')}{p(w^t)q(w'|w^t)} \right)$.
- Если приняли: $w^{t+1} = w'$,
- иначе: $w^{t+1} = w^t$.

Условие предельного распределения выполняется:

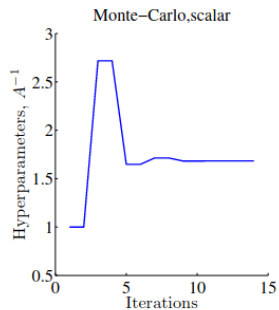
$$\begin{aligned} p(w)T(w|w') &= p(w)T(w'|w) = p(w')T(w'|w^t) = p(w')q(w'|w^t)A(w'|w^t) = \\ &= p(w^t)q(w^t|w')A(w^t|w'). \end{aligned}$$

- Сэмплы скоррелированы. Если требуется декоррелировать сэмплы, можно брать каждый k -й сэмпл.
- Работает в пространствах высокой размерности значительно лучше, чем сэмплирование с отклонением.
- Выбор адекватного распределения q — главная оптимизационная задача алгоритма.

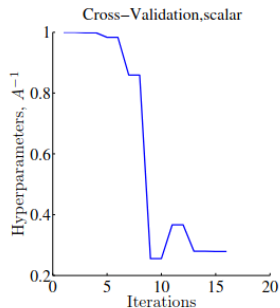
Пример: линейная модель



(a) Laplace approximation

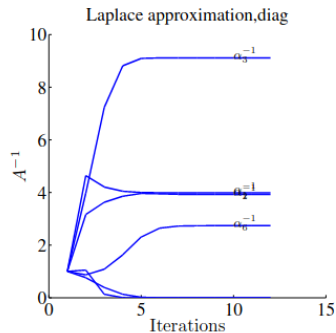


(b) Monte-Carlo

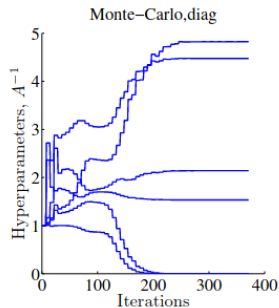


(c) Cross validation

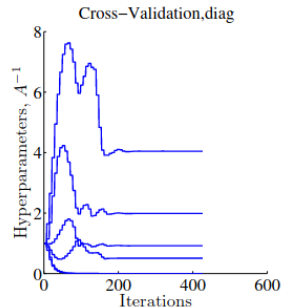
Пример: линейная модель



(a) Laplace approximation

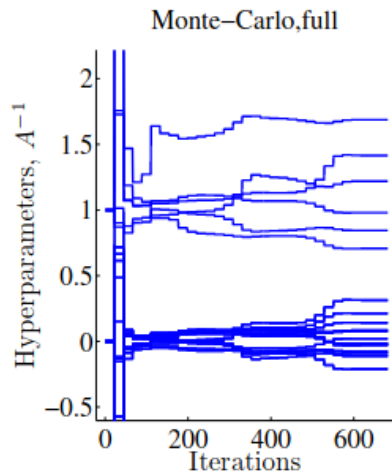


(b) Monte-Carlo

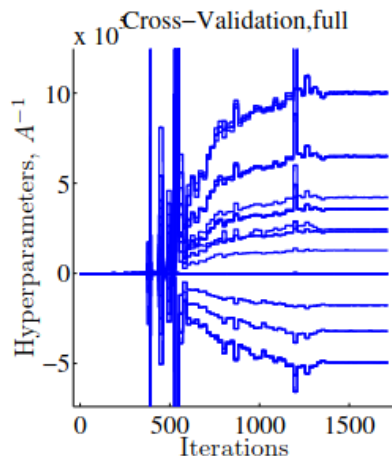


(c) Cross validation

Пример: линейная модель



(a) Monte-Carlo



(b) Cross validation

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:

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

Precondition-матрица для динамики Ланжевена

SGLD:

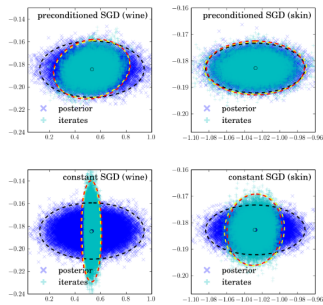
$$T = \mathbf{w} - \beta \nabla L + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \frac{\alpha}{2})$$

pSGLD:

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

матрица \mathbf{M} выбирается так, чтобы сделать шаг градиентного спуска равномерным по всем направлениям (с учетом дисперсии градиента).

Пример для SGD, для SGLD — аналогично.



Сэмплирование по Гиббсу

Пусть задана графическая модель для набора переменных w_1, \dots, w_n .
Тогда в цикле по всем переменным:

$$\hat{w}_i \sim p(w_i | w_1, w_{i-1}, w_{i+1}, w_n), w_i := \hat{w}_i$$

Сэмплирование Гиббса — частный случай МН-алгоритма.

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

Автокодировщик: порождающая модель?

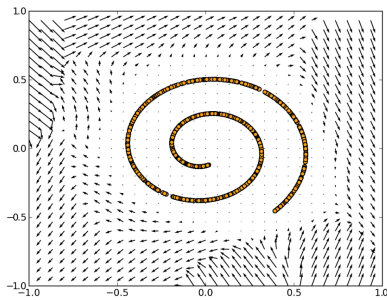
(Alain, Bengio 2012): рассмотрим модель автокодировщика с регуляризацией:

$$||\mathbf{f}(\mathbf{x}, \sigma) - \mathbf{x}||^2,$$

где σ — уровень шума, подаваемого на вход модели кодирования. Тогда

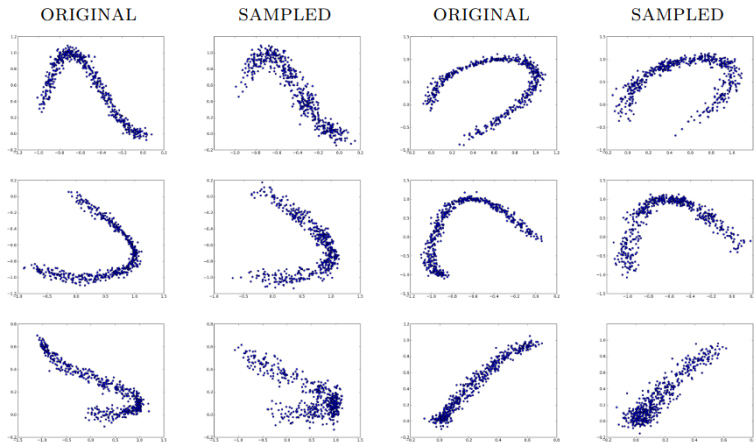
$$\frac{\partial \log p(\mathbf{x})}{\partial \mathbf{x}} = \frac{||\mathbf{f}(\mathbf{x}, \sigma) - \mathbf{x}||^2}{\sigma^2} + o(1) \text{ при } \sigma \rightarrow 0.$$

Векторное поле, индуцированное ошибкой реконструкции автокодировщика



Автокодировщик для оценки сэмплирования

$$A = \frac{p(x^*)}{p(x)} = \exp(E(x) - E(x^*)) \approx \frac{\partial E(x)^T}{\partial x} (x^* - x) + o(\|x - x^*\|).$$



Оптимизация распределения q

Распределение q можно задавать не как статическое распределение, а как сложное нейросетевое распределение.

- **Основное требование:** определенность распределений $p(x|x')$, $p(x'|x) \rightarrow$ распределение должно быть инвертируемо.
- Нейросеть вида $\mathbf{f}(\mathbf{x}, \mathbf{w}) = \mathbf{x} + \mathbf{g}(\mathbf{x}, \mathbf{w})$ является потоком и инвертируемо.

Варианты оптимизации:

- Энтропия * Acceptance rate (Li et al., 2020)
- Состязательная оптимизация между эмпирическим распределением и распределением q (Song et al., 2017).

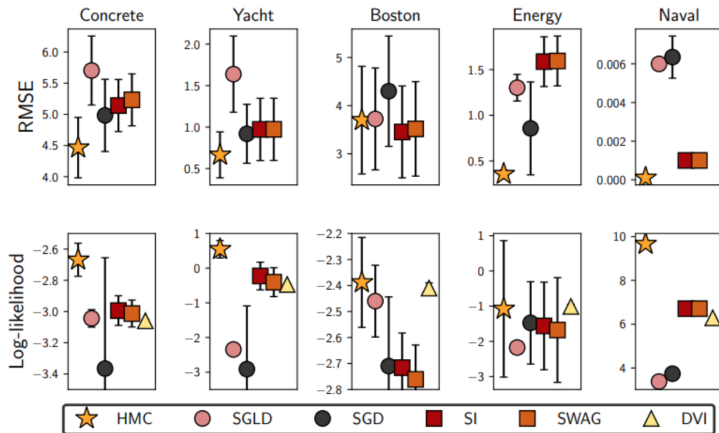
What Are Bayesian Neural Network Posteriors Really Like?

Izmailov et al., 2021:

- Методами НМС решается задача нахождения апостериорного распределения для глубоких моделей на стандартных датасетах.
- Вычисления: 512 TPU

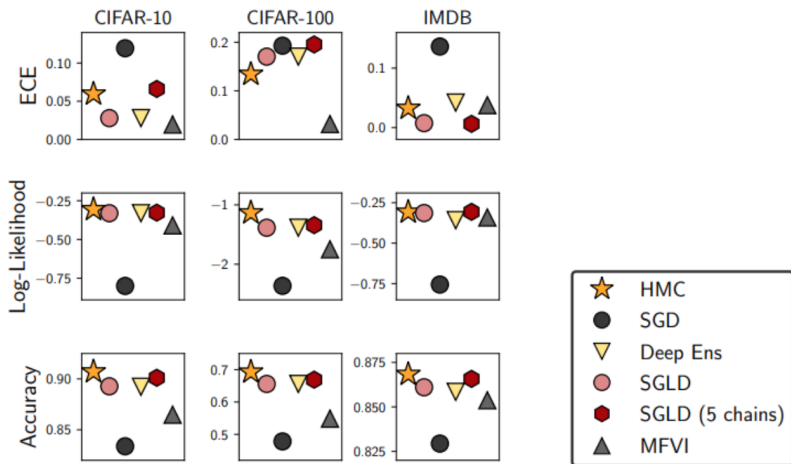
What Are Bayesian Neural Network Posteriors Really Like?

BNN evaluation: UCI



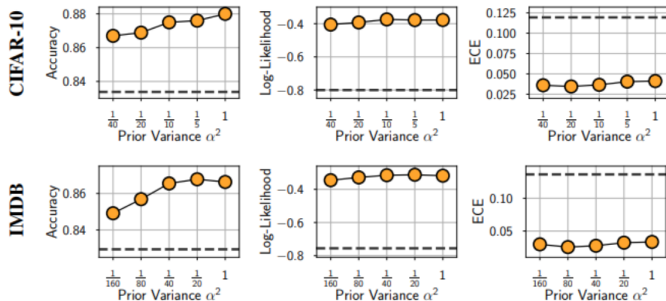
What Are Bayesian Neural Network Posteriors Really Like?

BNN evaluation: CIFAR and IMDB



What Are Bayesian Neural Network Posteriors Really Like?

Effect of priors



HMC BNNs are fairly robust to Gaussian prior variance.

Литература

- Bishop C. M., Nasrabadi N. M. Pattern recognition and machine learning. – New York : springer, 2006. – Т. 4. – №. 4. – С. 738.
- 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.
- Kuznetsov M., Tokmakova A., Strijov V. Analytic and stochastic methods of structure parameter estimation //Informatica. – 2016. – Т. 27. – №. 3. – С. 607-624.
- Welling M., Teh Y. W. Bayesian learning via stochastic gradient Langevin dynamics //Proceedings of the 28th international conference on machine learning (ICML-11). – 2011. – С. 681-688.
- Mandt S., Hoffman M., Blei D. A variational analysis of stochastic gradient algorithms //International conference on machine learning. – PMLR, 2016. – С. 354-363.
- Alain G., Bengio Y. What regularized auto-encoders learn from the data-generating distribution //The Journal of Machine Learning Research. – 2014. – Т. 15. – №. 1. – С. 3563-3593.
- Li Z., Chen Y., Sommer F. T. A neural network mcmc sampler that maximizes proposal entropy //arXiv preprint arXiv:2010.03587. – 2020.
- Song J., Zhao S., Ermon S. A-nice-mc: Adversarial training for mcmc //Advances in Neural Information Processing Systems. – 2017. – Т. 30.
- Izmailov P. et al. What are Bayesian neural network posteriors really like? //International Conference on Machine Learning. – PMLR, 2021. – С. 4629-4640.