

Оптимизация гиперпараметров градиентными методами

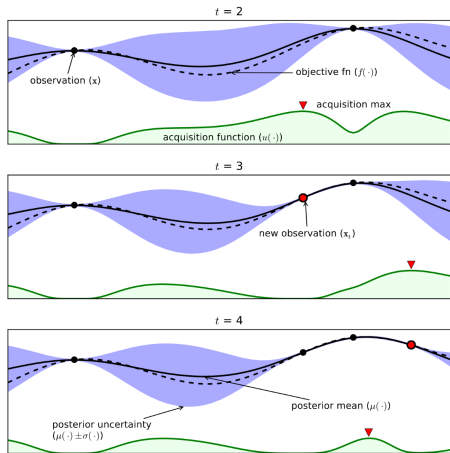
Бахтеев Олег

МФТИ

04.04.2018

Градиентные методы: зачем?

- Гиперпараметры — параметры распределения параметров модели.
- Основные методы оптимизации не позволяют проводить оптимизацию большого количества гиперпараметров (>10).
- Решение проблемы — использование градиентного спуска для гиперпараметров.



Shahriari et. al, 2016. Пример работы гауссового процесса.

Постановка задачи

Задана дифференцируемая по параметрам модель, приближающая зависимую переменную y :

$$f : \mathbb{R}^n \rightarrow \mathbb{Y}, \quad \mathbf{w} \in \mathbb{R}^u.$$

Функция f задает правдоподобие выборки $\log p(\mathbf{y}|\mathbf{X}, f)$.

Пусть также задано априорное распределение параметров:

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}),$$

где $\mathbf{A}^{-1} = \text{diag}[\alpha_1, \dots, \alpha_u]^{-1}$ — матрица ковариаций диагонального вида, определяемая гиперпараметрами $[\alpha_1, \dots, \alpha_u]$.

Кросс-валидация

Разобьем выборку \mathfrak{D} на k равных частей:

$$\mathfrak{D} = \mathfrak{D}_1 \sqcup \dots \sqcup \mathfrak{D}_k.$$

Запустим k оптимизаций модели, каждую на своей части выборки. Положим $\theta = [\mathbf{w}_1, \dots, \mathbf{w}_k]$, где $\mathbf{w}_1, \dots, \mathbf{w}_k$ — параметры модели при оптимизации k .

Пусть L — функция потерь:

$$L(\theta, \mathbf{A}^{-1}) = -\frac{1}{k} \sum_{q=1}^k \left(\frac{k}{k-1} \log p(\mathbf{y} \setminus \mathbf{y}_q | \mathbf{X} \setminus \mathbf{X}_q, \mathbf{w}_q) + \log p(\mathbf{w}_q | \mathbf{A}) \right). \quad (1)$$

Пусть Q — функция качества модели:

$$Q(\theta, \mathbf{A}^{-1}) = \frac{1}{k} \sum_{q=1}^k k \log p(\mathbf{y}_q | \mathbf{X}_q, \mathbf{w}_q).$$

Формальная постановка задачи

Задана дифференцируемая по параметрам модель, приближающая зависимую переменную y :

$$f : \mathbb{R}^n \rightarrow \mathbb{Y}, \quad \mathbf{w} \in \mathbb{R}^u.$$

Пусть $\boldsymbol{\theta} \in \mathbb{R}^s$ — множество всех оптимизируемых параметров.

$L(\boldsymbol{\theta}, \mathbf{A}^{-1})$ — дифференцируемая функция потерь по которой производится оптимизация функции f .

$Q(\boldsymbol{\theta}, \mathbf{A}^{-1})$ — дифференцируемая функция определяющая итоговое качество модели f и приближающая интеграл.

Требуется найти параметры $\hat{\boldsymbol{\theta}}$ и гиперпараметры $\hat{\mathbf{A}}^{-1}$ модели, доставляющие минимум следующему функционалу:

$$\hat{\mathbf{A}}^{-1} = \arg \max_{\mathbf{A}^{-1} \in \mathbb{R}^h} Q(\hat{\boldsymbol{\theta}}(\mathbf{A}^{-1}), \mathbf{A}^{-1}),$$

$$\hat{\boldsymbol{\theta}}(\mathbf{A}^{-1}) = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^s} L(\boldsymbol{\theta}, \mathbf{A}^{-1}).$$

Байесовский подход к сложности

Правдоподобие модели (“Evidence”):

$$p(\mathbf{y}|\mathbf{f}) = \int_{\mathbf{w}} p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{A})d\mathbf{w}.$$

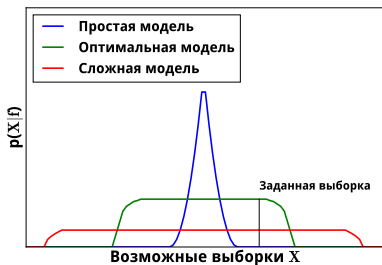
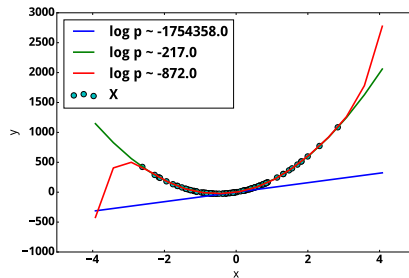


Схема выбора модели по правдоподобию



Пример: полиномы

Вариационная нижняя оценка

Пусть задано непрерывное распределение q . Тогда

$$\begin{aligned}\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A})}{q(\mathbf{w})} d\mathbf{w} + D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \mathbf{A})) \geq \\ &\geq \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A})}{q(\mathbf{w})} d\mathbf{w} = \\ &= -D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{A})) + \int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{A}) d\mathbf{w},\end{aligned}$$

где

$$D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{A})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathbf{A})}{q(\mathbf{w})} d\mathbf{w}.$$

Evidence: 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 $p(\mathbf{w}|\mathbf{f}) \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$:

$$L = \sum_{\mathbf{x}, \mathbf{y} \in \mathcal{D}} \log p(\mathbf{y}|\mathbf{x}, \hat{\mathbf{w}}) + \frac{1}{2} (\text{tr}(\mathbf{A}_q^{-1}) + \mu_q^T \mathbf{A}_q \mu_q - \ln |\mathbf{A}_q^{-1}|),$$

$$\hat{\mathbf{w}} \sim q = \mathcal{N}(\mu_q, \mathbf{A}_q^{-1}).$$

Вариационная оценка: оптимизация гиперпараметров

Пусть $L = -Q$:

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{A}) \geq \sum_{\mathbf{x}, y} \log p(y|\mathbf{x}, \hat{\mathbf{w}}) - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{A})) = -L(\theta, \mathbf{A}^{-1}) = Q(\theta, \mathbf{A}^{-1}),$$

где q — нормальное распределение с диагональной матрицей ковариаций:

$$q \sim \mathcal{N}(\mu_q, \mathbf{A}_q^{-1}),$$

$$D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{f})) = \frac{1}{2}(\text{Tr}[\mathbf{A}\mathbf{A}_q^{-1}] + (\mu - \mu_q)^T \mathbf{A}(\mu - \mu_q) - u + \ln |\mathbf{A}^{-1}| - \ln |\mathbf{A}_q^{-1}|).$$

В качестве оптимизируемых параметров θ выступают параметры распределения q :

$$\theta = [\alpha_1, \dots, \alpha_u, \mu_1, \dots, \mu_u].$$

Формальная постановка задачи: градиентная оптимизация

Определение

Оператором T назовем оператор стохастического градиентного спуска, производящий η шагов оптимизации:

$$\hat{\theta} = T \circ T \circ \dots \circ T(\theta_0, \mathbf{A}^{-1}) = T^\eta(\theta_0, \mathbf{A}^{-1}), \quad (2)$$

где

$$T(\theta, \mathbf{A}^{-1}) = \theta - \gamma \nabla L(\theta, \mathbf{A}^{-1})|_{\hat{\mathcal{D}}},$$

γ — длина шага градиентного спуска, θ_0 — начальное значение параметров θ , $\hat{\mathcal{D}}$ — случайная подвыборка исходной выборки \mathcal{D} .

Перепишем итоговую задачу оптимизации:

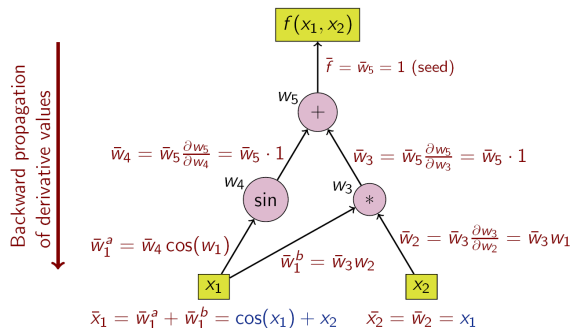
$$\hat{\mathbf{A}}^{-1} = \arg \max_{\mathbf{A}^{-1} \in \mathbb{R}^h} Q(T^\eta(\theta_0, \mathbf{A}^{-1})),$$

где θ_0 — начальное значение параметров θ .

RMAD, Maclaurin et. al, 2015

- 1 Провести η шагов оптимизации:
 $\theta = T(\theta_0, \mathbf{A}^{-1})$.
- 2 Положим $\hat{\nabla} \mathbf{A}^{-1} = \nabla_{\mathbf{A}}^{-1} Q(\theta, \mathbf{A}^{-1})$.
- 3 Положим $d\mathbf{v} = \mathbf{0}$.
- 4 Для $\tau = \eta \dots 1$ повторить:
- 5 $\theta^{\tau-1} = \theta^{\tau} - \gamma \mathbf{v}^{\tau}$.
- 6 $\mathbf{v}^{\tau-1} = \mathbf{v}^{\tau} + \gamma \hat{\nabla} \theta$.
- 7 $d\mathbf{v} = \gamma \hat{\nabla} \theta$.
- 8 $\hat{\nabla} \mathbf{A}^{-1} = \hat{\nabla} \mathbf{A}^{-1} - d\mathbf{v} \nabla_{\mathbf{A}^{-1}} \nabla_{\theta} Q$.
- 9 $\hat{\nabla} \theta = \hat{\nabla} \theta - d\mathbf{v} \nabla_{\theta} \nabla_{\theta} Q$.

Алгоритм RMAD основывается на Reverse-mode differentiation.



DrMAD

Алгоритм DrMad — упрощенный RMAD. Вводится предположение о линейности траектории обновления параметров θ .

① Провести η шагов оптимизации:
 $\theta = T(\theta_0, \mathbf{A}^{-1})$.

② Положим $\hat{\nabla} \mathbf{A}^{-1} = \nabla_{\mathbf{A}}^{-1} Q(\theta, \mathbf{A}^{-1})$.

③ Положим $d\mathbf{v} = \mathbf{0}$.

④ Для $\tau = \eta \dots 1$ повторить:

⑤ $\theta^{\tau-1} = \theta^{\tau} - \gamma \mathbf{v}^{\tau}$.

⑥ $\mathbf{v}^{\tau-1} = \mathbf{v}^{\tau} + \gamma \hat{\nabla}_{\theta}$.

⑦ $d\mathbf{v} = \gamma \hat{\nabla}_{\theta}$.

⑧ $\hat{\nabla} \mathbf{A}^{-1} = \hat{\nabla} \mathbf{A}^{-1} - d\mathbf{v} \nabla_{\mathbf{A}^{-1}} \nabla_{\theta} Q$.

⑨ $\hat{\nabla} \theta = \hat{\nabla} \theta - d\mathbf{v} \nabla_{\theta} \nabla_{\theta} Q$.

① Провести η шагов оптимизации:
 $\theta = T(\theta_0, \mathbf{A}^{-1})$.

② Положим $\hat{\nabla} \mathbf{A}^{-1} = \nabla_{\mathbf{A}}^{-1} Q(\theta, \mathbf{A}^{-1})$.

③ Положим $d\mathbf{v} = \mathbf{0}$.

④ Для $\tau = \eta \dots 1$ повторить:

⑤ $\theta^{\tau-1} = \theta_0 + \frac{\tau-1}{\eta} \theta^{\eta}$.

⑥

⑦ $d\mathbf{v} = \gamma \hat{\nabla}_{\theta}$.

⑧ $\hat{\nabla} \mathbf{A}^{-1} = \hat{\nabla} \mathbf{A}^{-1} - d\mathbf{v} \nabla_{\mathbf{A}^{-1}} \nabla_{\theta} Q$.

⑨ $\hat{\nabla} \theta = \hat{\nabla} \theta - d\mathbf{v} \nabla_{\theta} \nabla_{\theta} Q$.

Аналитическая формула оптимизации параметров

Утверждение (Pedregosa, 2016).

Пусть L — дифференцируемая функция, такая что все стационарные точки L являются локальными минимумами. Пусть также гессиан \mathbf{H}^{-1} функции потерь L является обратимым в каждой стационарной точке.

Тогда

$$\nabla_{\mathbf{A}^{-1}} Q(T(\theta_0), \mathbf{A}^{-1}) = \nabla_{\mathbf{A}^{-1}} Q(\theta^\eta, \mathbf{A}^{-1}) - \nabla_{\mathbf{A}^{-1}} \nabla_{\theta} L(\theta^\eta, \mathbf{A}^{-1})^T \mathbf{H}^{-1} \nabla_{\theta} Q(\theta^\eta, \mathbf{A}^{-1}).$$

Схема доказательства

- 1 Т.к. точка θ^η стационарна, то $\nabla_{\theta} L(\theta^\eta, \mathbf{A}^{-1}) = 0$.
- 2 Продифференцируем выражение по \mathbf{A}^{-1} :

$$\nabla_{\mathbf{A}^{-1}} \nabla_{\theta} L(\theta^\eta, \mathbf{A}^{-1}) + \mathbf{A}^{-1} \nabla_{\mathbf{A}^{-1}} T(\theta_0).$$

- 3 По правилу дифференцирования сложной функции:

$$\nabla_{\mathbf{A}^{-1}} Q(T(\theta_0), \mathbf{A}^{-1}) = \nabla_{\mathbf{A}^{-1}} Q(\theta^\eta, \mathbf{A}^{-1}) + \nabla_{\mathbf{A}^{-1}} T(\theta_0)^T \nabla_{\theta} Q(\theta^\eta, \mathbf{A}^{-1}).$$

- 4 Подставим в выражение 3 выражение 2 и получим искомое.

Жадная оптимизация гиперпараметров

На каждом шаге оптимизации параметров θ :

$$\mathbf{A}'^{-1} = \mathbf{A}^{-1} - \gamma_{\mathbf{A}^{-1}} \nabla_{\mathbf{A}^{-1}} Q(T(\theta, \mathbf{A}^{-1}), \mathbf{A}^{-1}) = \mathbf{A}^{-1} - \gamma_{\mathbf{A}^{-1}} \nabla_{\mathbf{A}^{-1}} Q(\theta - \gamma \nabla L(\theta, \mathbf{A}^{-1}), \mathbf{A}^{-1}),$$

где $\gamma_{\mathbf{A}^{-1}}$ — длина шага оптимизации гиперпараметров.

- Можно рассматривать как упрощение алгоритма RMAD, использующее только один элемент истории обновления параметров.
- Является приближением к решению аналитической формуле в случае $\mathbf{H}^{-1} \sim \mathbf{I}$.

Численное приближение аналитической формулы:

$$\nabla_{\mathbf{A}^{-1}} Q(\boldsymbol{\theta}^\eta, \mathbf{A}^{-1}) - \nabla_{\mathbf{A}^{-1}} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^\eta, \mathbf{A}^{-1})^\top \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}^\eta, \mathbf{A}^{-1}).$$

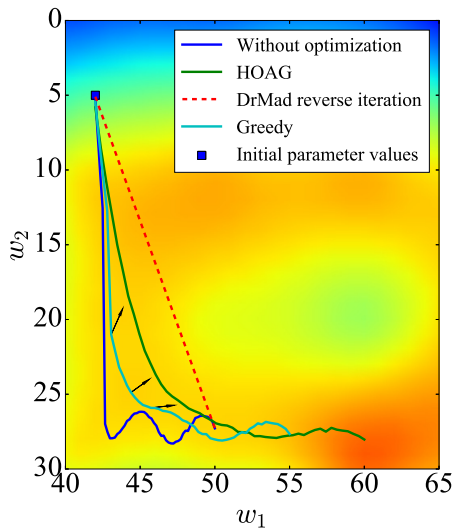
- ① Провести η шагов оптимизации: $\boldsymbol{\theta} = T(\boldsymbol{\theta}_0, \mathbf{A}^{-1})$.
- ② Решить линейную систему для вектора $\boldsymbol{\lambda}$: $\mathbf{H}^{-1}(\boldsymbol{\theta})\boldsymbol{\lambda} = \nabla_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \mathbf{A}^{-1})$.
- ③ Приближенное значение градиентов гиперпараметра вычисляется как:
 $\hat{\nabla}_{\mathbf{A}^{-1}} Q = \nabla_{\mathbf{A}^{-1}} Q(\boldsymbol{\theta}, \mathbf{A}^{-1}) - \nabla_{\boldsymbol{\theta}, \mathbf{A}^{-1}} L(\boldsymbol{\theta}, \mathbf{A}^{-1})^\top \boldsymbol{\lambda}$.

Итоговое правило обновления:

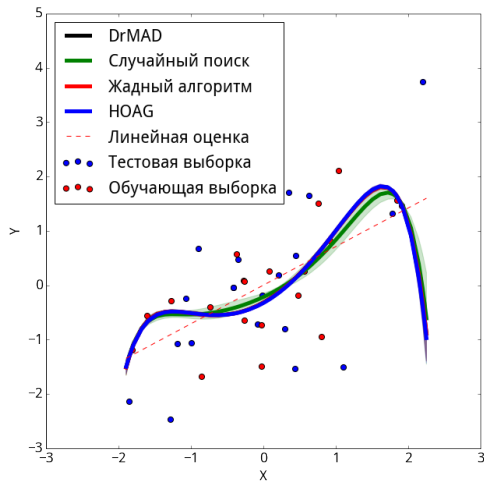
$$\mathbf{A}'^{-1} = \mathbf{A}^{-1} - \gamma_{\mathbf{A}^{-1}} \hat{\nabla}_{\mathbf{A}^{-1}} Q.$$

Сравнение алгоритмов

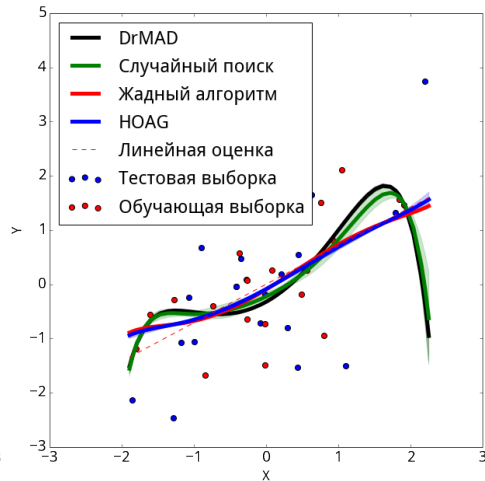
Алгоритм	+	-
Random search	Легко реализовать	Проклятие размерности
Жадная оптимизация	Оптимизация проводится внутри цикла оптимизации параметров. Легко реализовать	Жадность, неоптимальность.
HOAG	Быстрая сходимость.	Качество результатов зависит от решения линейного уравнения $\mathbf{H}^{-1}(\theta)\lambda = \nabla_{\theta}Q(\theta, \mathbf{A}^{-1})$.
DrMAD	Учитывает особенности оператора оптимизации. Можно использовать для оптимизации мета-параметров.	Неустойчив при больших значениях длины градиентного шага $\gamma_{\mathbf{A}}^{-1}$. Качество оптимизации зависит от кривизны траектории обновления параметров.



Эксперименты: полиномы

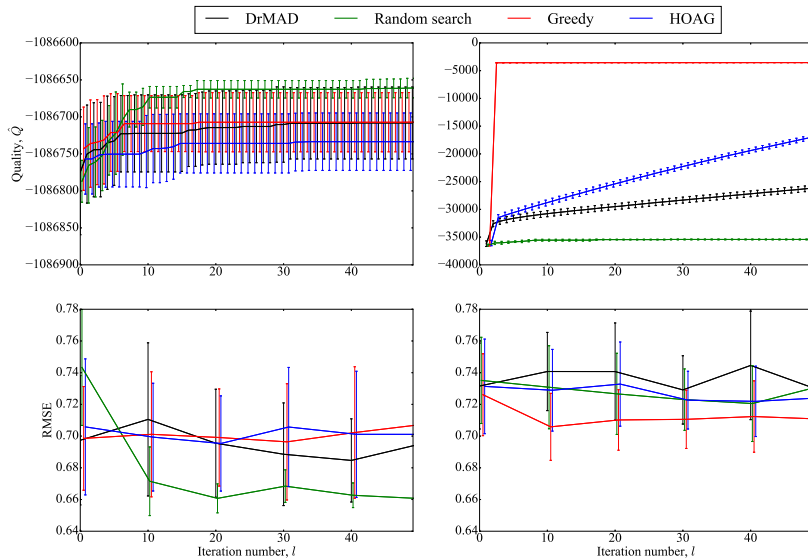


Кросс-валидация

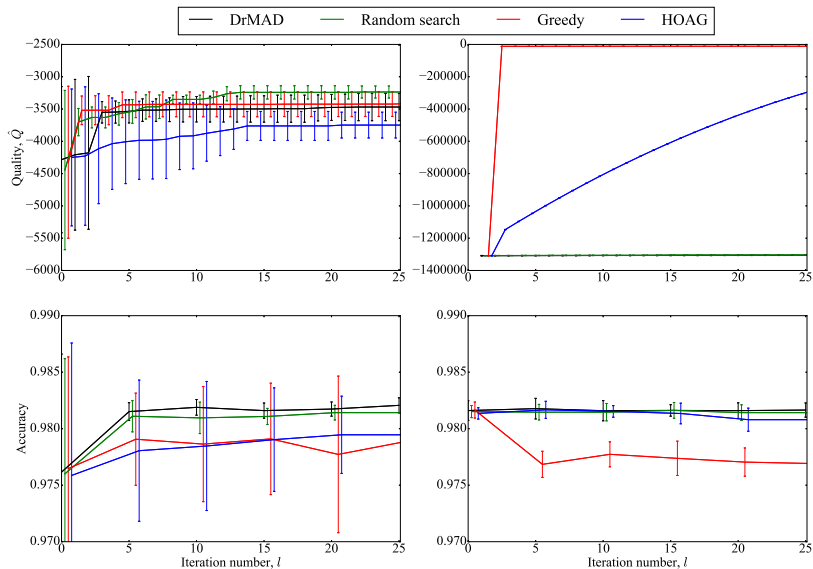


Evidence

Эксперименты: WISDM



Эксперименты: MNIST



Эксперименты: MNIST

Добавление гауссового шума $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$:



Без шума



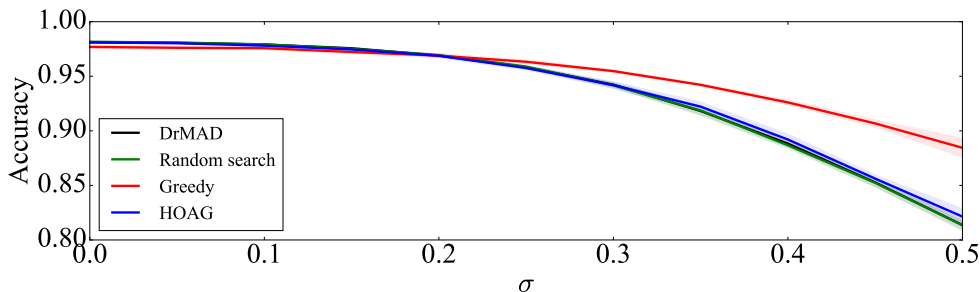
$\sigma = 0.1$



$\sigma = 0.25$



$\sigma = 0.5$



Используемые материалы

- ① David J. C. MacKay, Information Theory, Inference & Learning Algorithms, 2003
- ② Christopher Bishop, Pattern Recognition and Machine Learning, 2006
- ③ Dougal Maclaurin et. al, Gradient-based Hyperparameter Optimization through Reversible Learning, 2015
- ④ Jelena Luketina et. al, Scalable Gradient-Based Tuning of Continuous Regularization Hyperparameters, 2016
- ⑤ Jie Fu et. al, DrMAD: Distilling Reverse-Mode Automatic Differentiation for Optimizing Hyperparameters of Deep Neural Networks, 2016
- ⑥ Fabian Pedregosa, Hyperparameter optimization with approximate gradient, 2016
- ⑦ Bobak Shahriari et. al, Taking the Human Out of the Loop: A Review of Bayesian Optimization, 2016