

Байесовское мультимоделирование: оптимизация гиперпараметров

Московский Физико-Технический Институт

2021

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”):

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

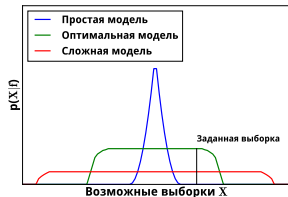
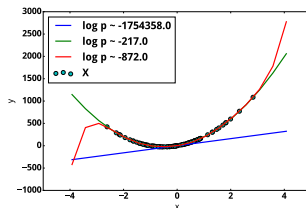


Схема выбора модели



Пример: полиномы

Что такое гиперпараметры

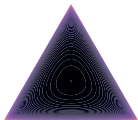
Определение

Априорным распределением $p(\mathbf{w}|\mathbf{h})$ параметров модели назовем вероятностное распределение, соответствующее предположениям о распределении параметров модели.

Определение

Гиперпараметрами $\mathbf{h} \in \mathbb{H}$ модели назовем параметры априорного распределения (параметры распределения параметров модели).

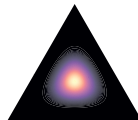
Дискретные распределения: релаксация



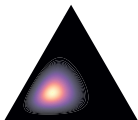
$$\bar{\alpha} = [1, 1, 1], t = 0.9$$



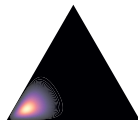
$$t = 1.0$$



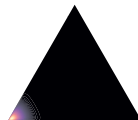
$$t = 10.0$$



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



$$[0.75, 0.125, 0.125]$$



$$[0.9, 0.05, 0.05]$$

Аппроксимация Лапласа

Нелинейный случай с m объектами и n признаками: $\mathbf{y} \sim \mathcal{N}(\mathbf{f}(\mathbf{X}, \mathbf{w}), \beta^{-1})$, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{A}^{-1})$.

Запишем интеграл:

$$p(\mathcal{D}|\mathbf{h}) = p(\mathbf{y}|\mathbf{X}, \mathbf{A}, \beta) = \frac{\sqrt{\beta \cdot |\mathbf{A}|}}{\sqrt{(2\pi)^{m+n}}} \int_{\mathbf{w}} \exp(-S(\mathbf{w})) d\mathbf{w}.$$

Разложим S в ряд Тейлора:

$$S(\mathbf{w}) \approx S(\hat{\mathbf{w}}) + \frac{1}{2} \Delta \mathbf{w}^T \mathbf{H} \Delta \mathbf{w}$$

Интеграл приводится к виду:

$$\frac{\sqrt{\beta \cdot |\mathbf{A}|}}{\sqrt{(2\pi)^{m+n}}} S(\hat{\mathbf{w}}) \int_{\mathbf{w}} \exp(-\frac{1}{2} \Delta \mathbf{w}^T \mathbf{H} \Delta \mathbf{w}) d\mathbf{w}$$

Выражение под интегралом соответствует плотности ненормированного нормального распределения.

Graves, 2011

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

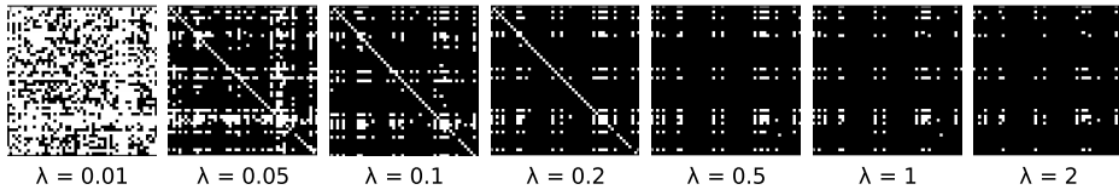
Вариационное распределение: $q(\mathbf{w}) \sim \mathcal{N}(\boldsymbol{\mu}_q, \sigma_q \mathbf{I})$.

Жадная оптимизация гиперпараметров:

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

Прунинг параметра w_i определяется относительной плотностью:

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



Постановка задачи

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

$L(\theta, \mathbf{h})$ — дифференцируемая функция потерь по которой производится оптимизация функции \mathbf{f} .

$Q(\theta, \mathbf{h})$ — дифференцируемая функция определяющая итоговое качество модели \mathbf{f} и приближающая интеграл.

Требуется найти параметры θ^* и гиперпараметры \mathbf{h}^* модели, доставляющие минимум следующему функционалу:

$$\mathbf{h}^* = \arg \max_{\mathbf{h} \in \mathbb{H}} Q(\theta^*(\mathbf{h}), \mathbf{h}),$$

$$\theta(\mathbf{h})^* = \arg \min_{\theta \in \mathbb{R}^s} L(\theta, \mathbf{h}).$$

Байесовский вывод

Пусть $\theta = [\mathbf{w}]^\top$.

Первый уровень:

$$\theta^* = \arg \max(-L(\theta, \mathbf{h})) = p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \mathbf{h}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathbf{y}|\mathbf{X}, \mathbf{h})}.$$

Второй уровень:

$$p(\mathbf{h}|\mathbf{X}, \mathbf{y}) \propto p(\mathbf{y}|\mathbf{X}, \mathbf{h})p(\mathbf{h}),$$

Полагая распределение параметров $p(\mathbf{h})$ равномерным на некоторой большой окрестности, получим задачу оптимизации гиперпараметров:

$$Q(\theta, \mathbf{h}) = p(\mathbf{y}|\mathbf{X}, \mathbf{h}) = \int_{\mathbf{w} \in \mathbb{R}^u} p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{h}) \rightarrow \max_{\mathbf{h} \in \mathbb{H}}.$$

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:

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

Пусть Q — функция качества модели:

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

Вариационная нижняя оценка

Пусть $L = -Q$:

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

где q — нормальное распределение с диагональной матрицей ковариаций:

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

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

В качестве оптимизируемых параметров $\boldsymbol{\theta}$ выступают параметры распределения q :

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

Evidence vs Кросс-валидация

Оценка Evidence:

$$\log p(\mathcal{D}|\mathbf{f}) = \log p(\mathcal{D}_1|\mathbf{f}) + \log p(\mathcal{D}_2|\mathcal{D}_1, \mathbf{f}) + \dots + \log p(\mathcal{D}_n|\mathcal{D}_1, \dots, \mathcal{D}_{n-1}, \mathbf{f}).$$

Оценка leave-one-out:

$$\text{LOU} = E \log p(\mathcal{D}_n|\mathcal{D}_1, \dots, \mathcal{D}_{n-1}, \mathbf{f}).$$

Кросс-валидация использует среднее значение последнего члена $p(\mathcal{D}_n|\mathcal{D}_1, \dots, \mathcal{D}_{n-1}, \mathbf{f})$ для оценки сложности.

Evidence учитывает **полную** сложность описания заданной выборки, определяющую предсказательную способность модели с самого начала.

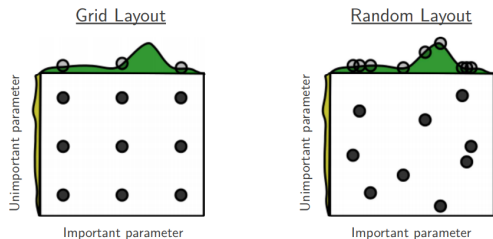
Базовые методы оптимизации гиперпараметров

Варианты:

- Поиск по решетке;
- Случайный поиск.

Оба метода страдают от проклятия размерности.

Случайный поиск может быть более эффективным, если пространство гиперпараметров вырождено.



Bergstra et al., 2012

Гауссовый процесс

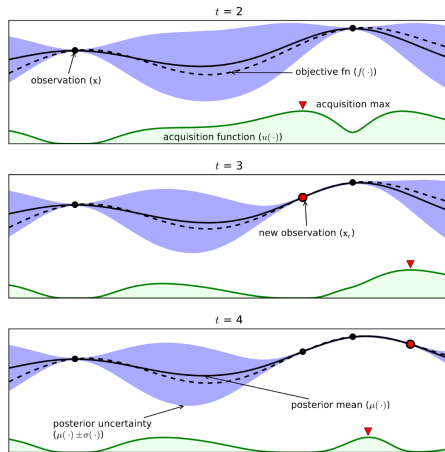
Идея:

Будем моделировать $Q(\theta(\mathbf{h})^*, \mathbf{h})$ гауссовым процессом, зависящим от \mathbf{h} .

Плюсы:

- Гибкость модели.
- Дешевле, чем обучения модели.

Минусы: кубическая сложность по количеству гиперпараметров.



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

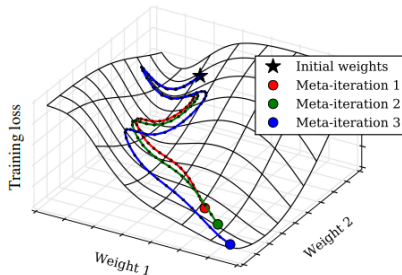
Градиентные методы

Идея: Будем производить оптимизацию вдоль всей траектории оптимизации параметров.

Плюсы:

- Оптимизация гиперпараметров будет учитывать оптимизацию параметров.
- Сложность меняется незначительно от количества гиперпараметров.

Минусы: вычислительно дорого.



Maclaurin et. al, 2015. Пример работы.

Формальная постановка задачи: градиентная оптимизация

Определение

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

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

где

$$T(\theta, \mathbf{h}) = \theta - \beta \nabla L(\theta, \mathbf{h})|_{\hat{\mathcal{D}}},$$

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

Перепишем итоговую задачу оптимизации:

$$\mathbf{h}^* = \arg \max_{\mathbf{h} \in \mathbb{H}} Q(T^\eta(\theta_0, \mathbf{h})),$$

где θ_0 — начальное значение параметров θ .

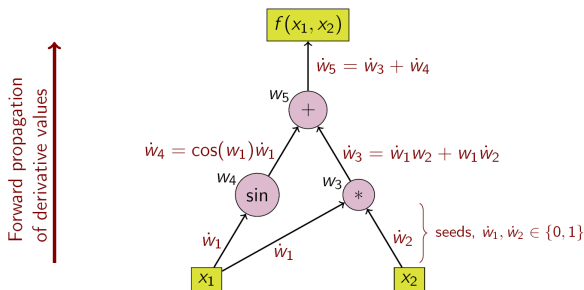
Forward-mode differentiation

Идея дифференцирования: применение формулы:

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w_{n-1}} \frac{\partial w_{n-1}}{\partial x} = \frac{\partial y}{\partial w_{n-1}} \left(\frac{\partial w_{n-1}}{\partial w_{n-2}} \frac{\partial w_{n-2}}{\partial x} \right) = \dots$$

Пример (wiki):

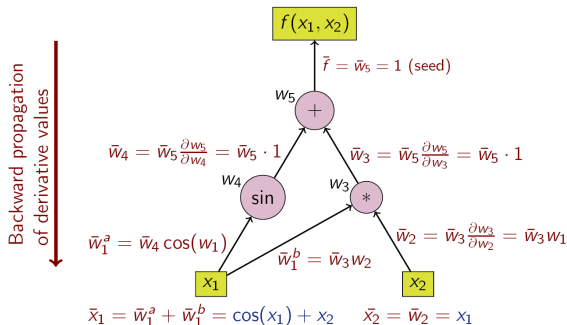
$$x_1 x_2 + \sin(x_1)$$



Reverse-mode differentiation

Идея дифференцирования: применение формулы:

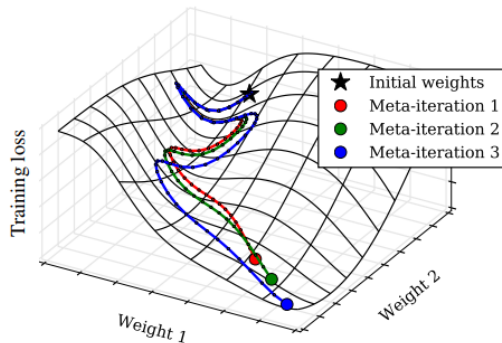
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w_1} \frac{\partial w_1}{\partial x} = \left(\frac{\partial y}{\partial w_2} \frac{\partial w_2}{\partial w_1} \right) \frac{\partial w_1}{\partial x} = \dots$$



RMAD, Maclaurin et. al, 2015

- 1 Провести η шагов оптимизации с моментом γ : $\theta = T(\theta_0, \mathbf{h})$.
- 2 Положим $\hat{\nabla} \mathbf{h} = \nabla_{\mathbf{h}} Q(\theta, \mathbf{h})$.
- 3 Положим $d\mathbf{v} = 0$.
- 4 Для $\tau = \eta \dots 1$ повторить:
 - 5 Вычислить $\theta^{\tau-1}$.
 - 6 Вычислить градиент на шаге $\tau - 1$, используя RMD.

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 с моментом γ : $\theta = T(\theta_0, \mathbf{h})$.
- 2 Положим $\hat{\nabla} \mathbf{h} = \nabla_{\mathbf{h}} Q(\theta, \mathbf{h})$.
- 3 Положим $d\mathbf{v} = 0$.
- 4 Для $\tau = \eta \dots 1$ повторить:
- 5 Вычислить $\theta^{\tau-1}$.
- 6 Вычислить градиент на шаге $\tau - 1$, используя RMD.

- 1 Провести η шагов оптимизации с моментом γ : $\theta = T(\theta_0, \mathbf{h})$.
- 2 Положим $\hat{\nabla} \mathbf{h} = \nabla_{\mathbf{h}} Q(\theta, \mathbf{h})$.
- 3 Положим $d\mathbf{v} = 0$.
- 4 Для $\tau = \eta \dots 1$ повторить:
- 5 $\theta^{\tau-1} = \theta_0 + \frac{\tau-1}{\eta} \theta^\eta$.
- 6 Вычислить градиент на шаге $\tau - 1$, используя RMD.

Аналитическая формула оптимизации параметров

Утверждение (Pedregosa, 2016)

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

Тогда

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

Жадная оптимизация гиперпараметров

На каждом шаге оптимизации параметров θ :

$$\mathbf{h}' = \mathbf{h} - \beta_{\mathbf{h}} \nabla_{\mathbf{h}} Q(T(\theta, \mathbf{h}), \mathbf{h}) = \mathbf{h} - \beta_{\mathbf{h}} \nabla_{\mathbf{h}} Q(\theta - \beta \nabla L(\theta, \mathbf{h}), \mathbf{h}),$$

где $\beta_{\mathbf{h}}$ — 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 $\mathbf{H}^{-1} \sim \mathbf{I}$.
- Сложность: $O(|\theta| \cdot |\mathbf{h}|)$
- Можно упростить, используя формулу конечных приращений, см. DARTS, $O(|\theta| + |\mathbf{h}|)$

HOAG

Численное приближение аналитической формулы:

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

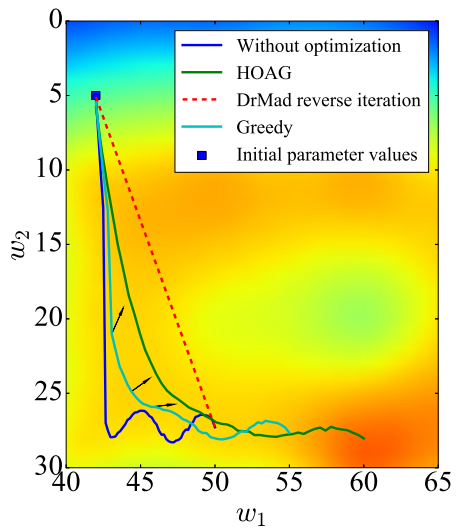
- 1 Провести η шагов оптимизации: $\boldsymbol{\theta} = T(\boldsymbol{\theta}_0, \mathbf{h})$.
- 2 Решить линейную систему для вектора $\boldsymbol{\lambda}$: $\mathbf{H}(\boldsymbol{\theta})\boldsymbol{\lambda} = \nabla_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \mathbf{h})$.
- 3 Приближенное значение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

Итоговое правило обновления:

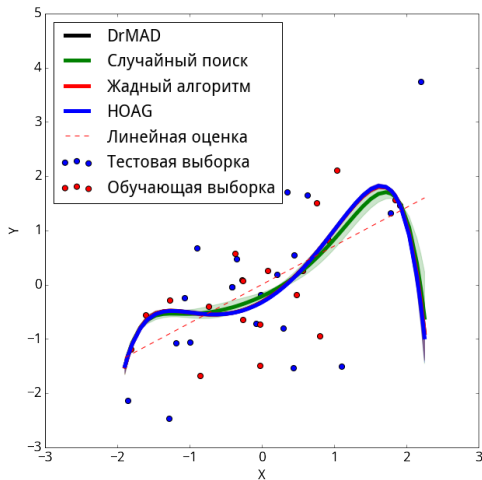
$$\mathbf{h}' = \mathbf{h} - \gamma_{\mathbf{h}} \hat{\nabla}_{\mathbf{h}} Q.$$

Сравнение алгоритмов

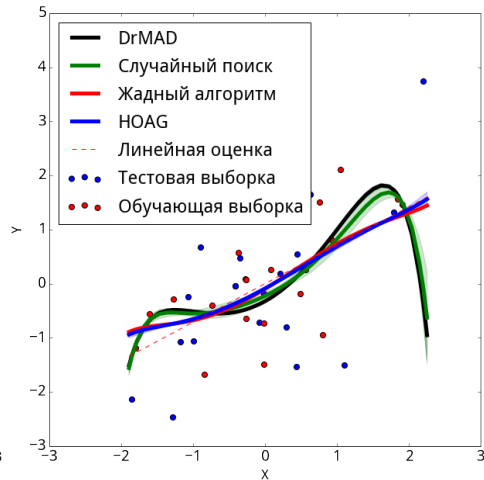
Алгоритм	+	-
Random search	Легко реализовать	Проклятие размерности
Жадная оптимизация	Оптимизация проводится внутри цикла оптимизации параметров. Легко реализовать	Жадность, неоптимальность.
HOAG	Быстрая сходимость.	Качество результатов зависит от решения линейного уравнения $\mathbf{H}(\theta)\lambda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



Эксперименты: полиномы

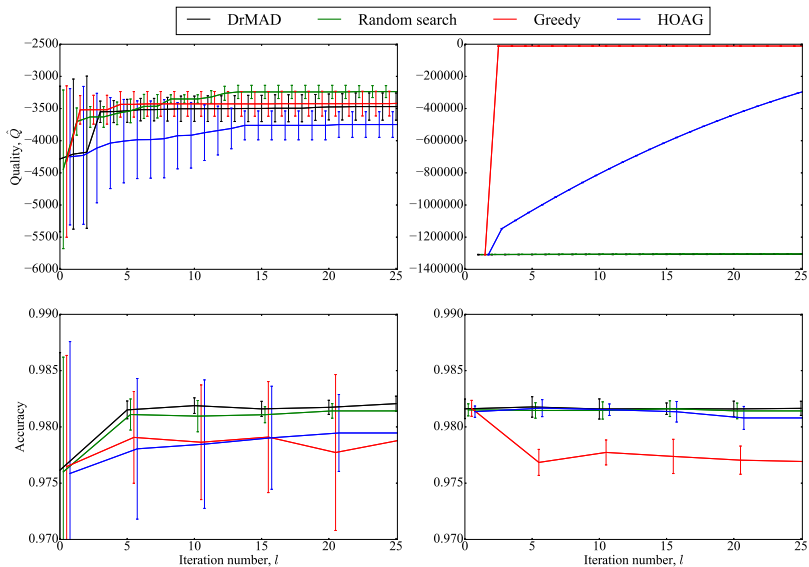


Кросс-валидация



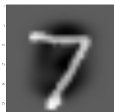
Evidence

Эксперименты: MNIST



Эксперименты: MNIST

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



Без шума



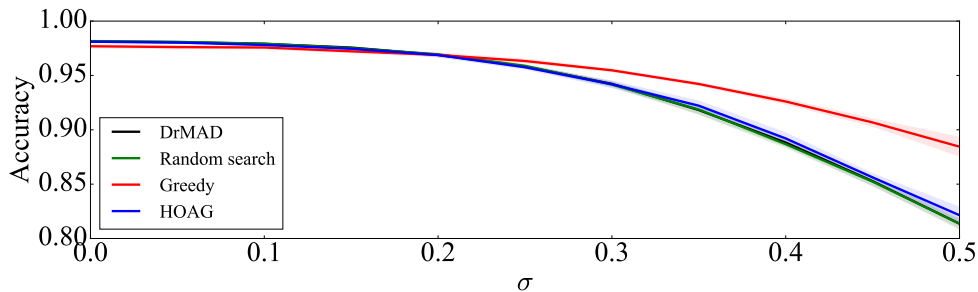
$\sigma = 0.1$



$\sigma = 0.25$



$\sigma = 0.5$



Литература и прочие ресурсы

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – Т. 128. – №. 9.
- Bakhteev O. Y., Strijov V. V. Comprehensive analysis of gradient-based hyperparameter optimization algorithms //Annals of Operations Research. – 2020. – Т. 289. – №. 1. – С. 51-65.
- Graves A. Practical variational inference for neural networks //Advances in neural information processing systems. – 2011. – Т. 24.
- Bergstra et al., Random Search for Hyper-Parameter Optimization, 2012
- Dougal Maclaurin et. al, Gradient-based Hyperparameter Optimization through Reversible Learning, 2015
- Jelena Luketina et. al, Scalable Gradient-Based Tuning of Continuous Regularization Hyperparameters, 2016
- Jie Fu et. al, DrMAD: Distilling Reverse-Mode Automatic Differentiation for Optimizing Hyperparameters of Deep Neural Networks, 2016
- Fabian Pedregosa, Hyperparameter optimization with approximate gradient, 2016
- Bobak Shahriari et. al, Taking the Human Out of the Loop: A Review of Bayesian Optimization, 2016
- Liu H., Simonyan K., Yang Y. Darts: Differentiable architecture search //arXiv preprint arXiv:1806.09055. – 2018.