

# Оптимизация гиперпараметров градиентными методами

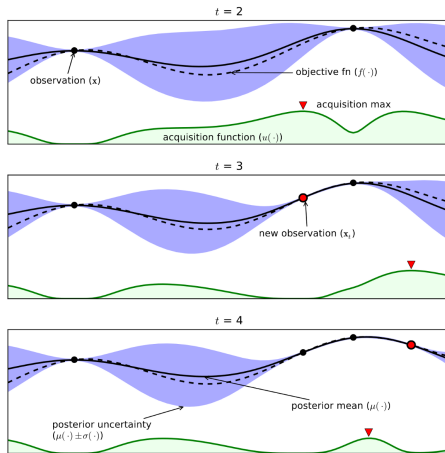
Бахтеев Олег

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:

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

Пусть  $Q$  — функция качества модели:

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

# 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.

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

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

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

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

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

# Байесовый подход к сложности

Правдоподобие модели (“Evidence”):

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

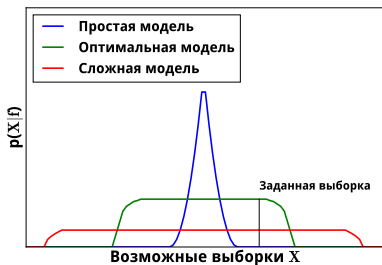
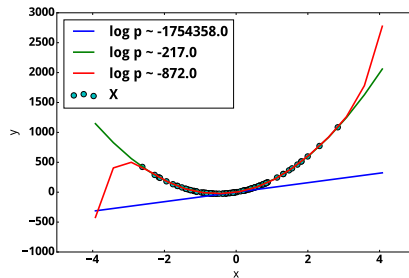


Схема выбора модели по правдоподобию



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

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

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



# Вариационная оценка: оптимизация гиперпараметров

Пусть  $L = -Q$ :

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{A}) \geq \sum_{\mathbf{x}, y} \log p(y|\mathbf{x}, \hat{\mathbf{w}}) - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{A})) = -L(\boldsymbol{\theta}, \mathbf{h}) = Q((\boldsymbol{\theta}, \mathbf{h}), \quad (2)$$

где  $q$  — нормальное распределение с диагональной матрицей ковариаций:

$$q \sim \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q^{-1}), \quad (3)$$

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

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

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

# Формальная постановка задачи: градиентная оптимизация

## Определение

Оператором  $T$  назовем оператор стохастического градиентного спуска, производящий  $\eta$  шагов оптимизации:

$$\hat{\theta} = T \circ T \circ \dots \circ T(\theta_0, \mathbf{h}) = T^\eta(\theta_0, \mathbf{h}), \quad (4)$$

где

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

$\gamma$  — длина шага градиентного спуска,  $\theta_0$  — начальное значение параметров  $\theta$ ,  $\hat{\mathcal{D}}$  — случайная подвыборка исходной выборки  $\mathcal{D}$ .

Перепишем итоговую задачу оптимизации:

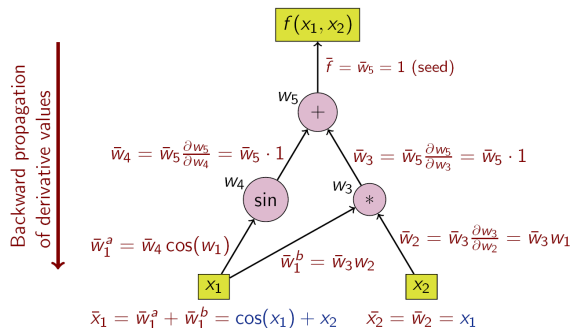
$$\hat{\mathbf{h}} = \arg \max_{\mathbf{h} \in \mathbb{R}^h} Q(T^\eta(\theta_0, \mathbf{h})),$$

где  $\theta_0$  — начальное значение параметров  $\theta$ .

# RMAD, Maclaurin et. al, 2015

- 1 Провести  $\eta$  шагов оптимизации:  
 $\theta = T(\theta_0, \mathbf{h})$ .
- 2 Положим  $\hat{\nabla} \mathbf{h} = \nabla_{\mathbf{h}} Q(\theta, \mathbf{h})$ .
- 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  $\hat{\nabla} \mathbf{h} = \hat{\nabla} \mathbf{h} - d\mathbf{v} \nabla_{\mathbf{h}} \nabla_{\theta} Q$ .
- 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  $\theta$ .

① Провести  $\eta$  шагов оптимизации:  
 $\theta = T(\theta_0, \mathbf{h})$ .

② Положим  $\hat{\nabla} \mathbf{h} = \nabla_{\mathbf{h}} Q(\theta, \mathbf{h})$ .

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

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

① Провести  $\eta$  шагов оптимизации:  
 $\theta = T(\theta_0, \mathbf{h})$ .

② Положим  $\hat{\nabla} \mathbf{h} = \nabla_{\mathbf{h}} Q(\theta, \mathbf{h})$ .

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

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

# Закрытая форма оптимизации параметров

Утверждение (Pedregosa, 2016).

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

Тогда

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

## Доказательство

① Т.к. точка  $\theta^\eta$  стационарна, то  $\nabla_{\theta} L(\theta^\eta, \mathbf{A}^{-1}) = 0$ .

② Продифференцируем выражение по  $\mathbf{A}^{-1}$  :

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

③ По цепному правилу:

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

④ Подставим в выражение 3 выражение 2 и получим искомое.

# Жадная оптимизация гиперпараметров

На каждом шаге оптимизации параметров  $\theta$ :

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

где  $\gamma_{\mathbf{h}}$  — длина шага оптимизации гиперпараметров.

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

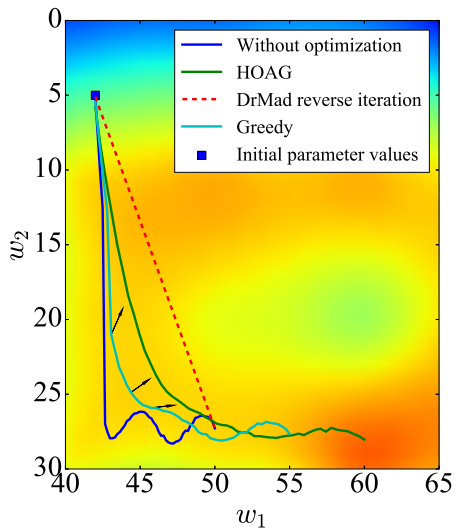
Численное приближение закрытой формы:

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

- ① Провести  $\eta$  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:  
 $\hat{\nabla}_{\mathbf{h}} Q = \nabla_{\mathbf{h}} Q(\boldsymbol{\theta}, \mathbf{h}) - \nabla_{\boldsymbol{\theta}, \mathbf{h}} L(\boldsymbol{\theta}, \mathbf{h})^\top \boldsymbol{\lambda}$ .

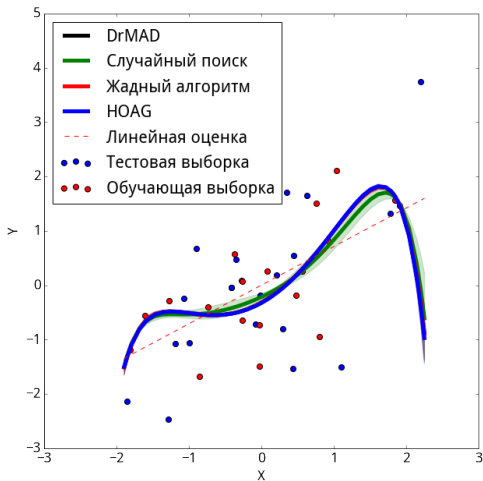
Итоговое правило обновления:

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

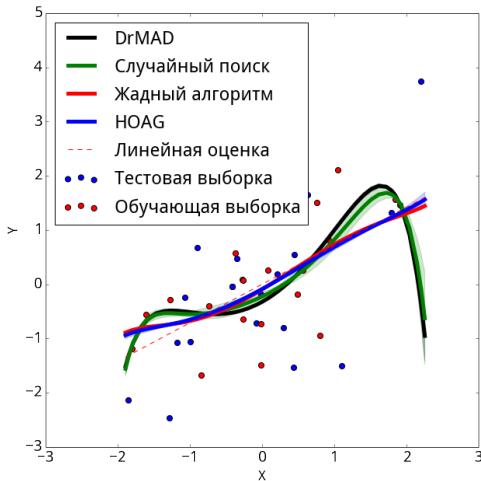




# Эксперименты: полиномы

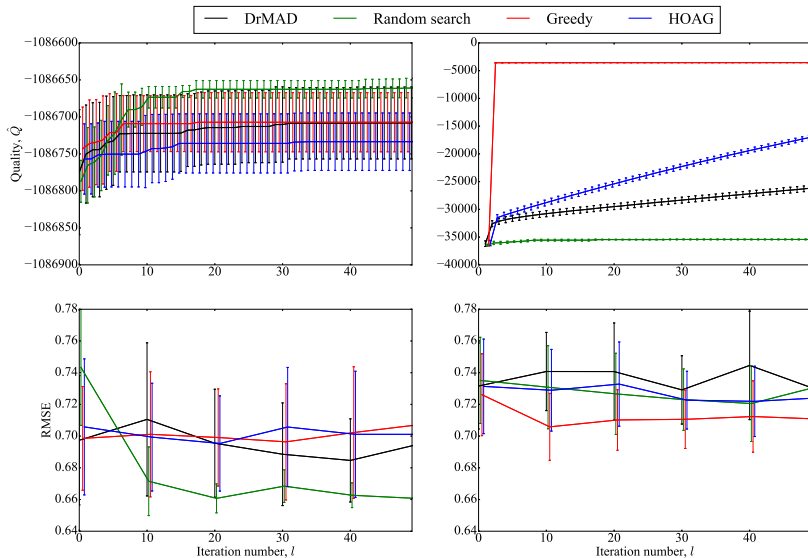


Кросс-валидация

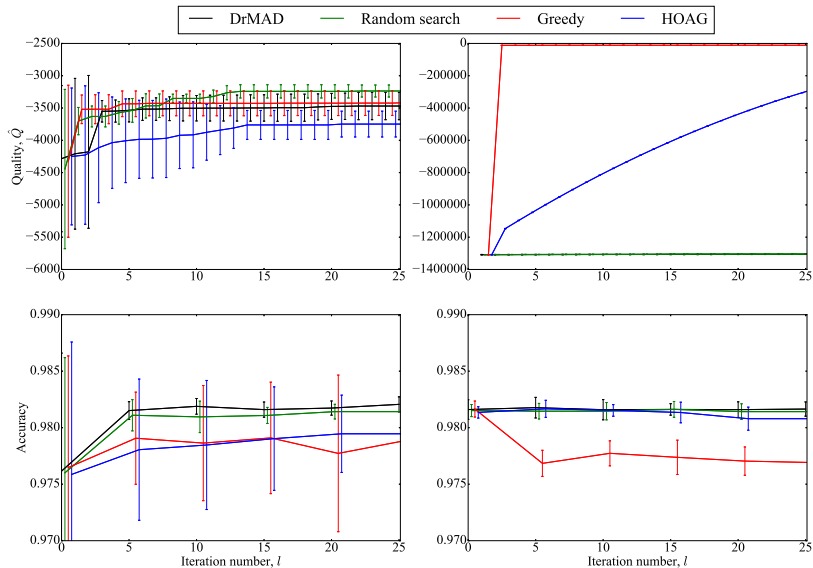


Evidence

# Эксперименты: WISDM



# Эксперименты: MNIST



# Эксперименты: MNIST

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



Без шума



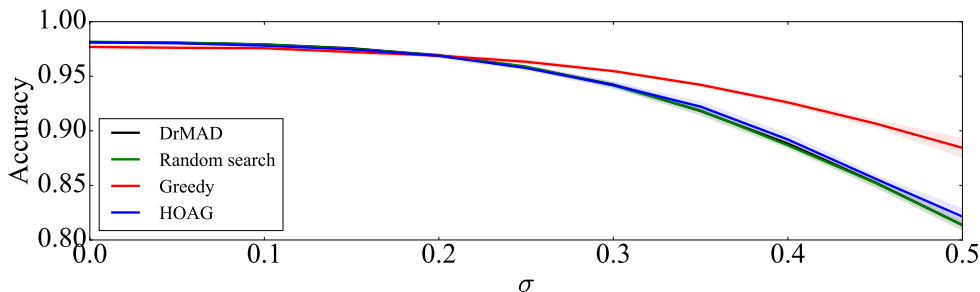
$\sigma = 0.1$



$\sigma = 0.25$



$\sigma = 0.5$



# 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