

Градиентные методы оптимизации гиперпараметров моделей глубокого обучения

О. Ю. Бахтеев

Научный руководитель: д.ф.-м.н. В.В. Стрижов
Московский Физико-Технический Институт (Государственный Университет)

ММРО-2017

09.10.2017

Цель работы

Исследуются

Методы автоматического порождения и прореживания моделей глубокого обучения.

Требуется

Предложить алгоритм нахождения оптимальных значений гиперпараметров (параметров распределения параметров) модели.

Проблемы нахождения оптимальных значений гиперпараметров

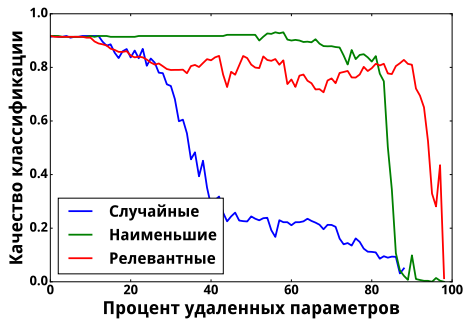
- Многоэкстремальность задачи оптимизации параметров модели,
- Вычислительная сложность оптимизации,
- Большое количество гиперпараметров.

Решение

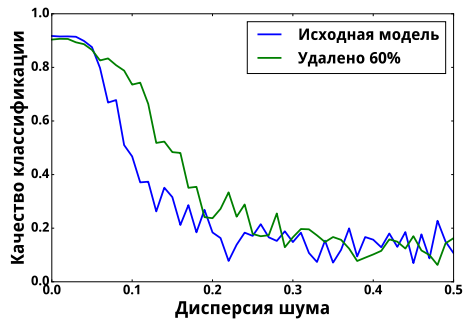
Предлагается оптимизировать параметры и гиперпараметры модели в единой процедуре с использованием градиентных методов. В качестве критерия оптимальности модели рассматривается нижняя оценка правдоподобия модели.

Проблемы обучения сетей

Правдоподобие моделей с избыточным количеством параметров не меняется при удалении параметров.

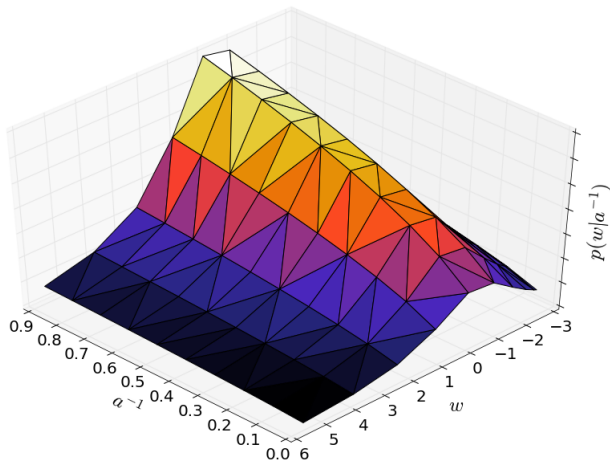


Избыточность параметров модели



Неустойчивость модели

Зависимость правдоподобия от гиперпараметров



Формальная постановка задачи

Задана дифференцируемая по параметрам модель $\mathbf{f}(\mathbf{w}, \mathbf{x})$, задающая правдоподобие выборки $p(\mathcal{D}|\mathbf{w})$. Задано априорное распределение параметров модели:

$$p(\mathbf{w}|\mathbf{A}^{-1}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}) = \mathcal{N}(\mathbf{0}, \text{diag}[\alpha_1, \dots, \alpha_n]^T).$$

Пусть $\boldsymbol{\theta} \in \mathbb{R}^s$ — множество параметров, подлежащих оптимизации (соответствует параметрам модели \mathbf{w}),

L — оптимизируемая функция потерь,

Q — критерий качества модели.

Итоговая задача оптимизации:

$$\hat{\mathbf{A}}^{-1} = \arg \max_{[\alpha_1, \dots, \alpha_n]^T \in \mathbb{R}^n} Q(\hat{\boldsymbol{\theta}}(\mathbf{A}^{-1}), \mathbf{A}^{-1}, \mathcal{D}),$$

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

L и Q: Кросс-валидация

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

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

Запустим k оптимизаций модели, r -я модель обучается на выборках $\mathcal{D}^r = \mathcal{D}_1, \dots, \mathcal{D}_{r-1}, \mathcal{D}_{r+1}, \dots, \mathcal{D}_k$.

Положим $\theta = [\mathbf{w}_1, \dots, \mathbf{w}_k]$ — параметры всех запусков модели.

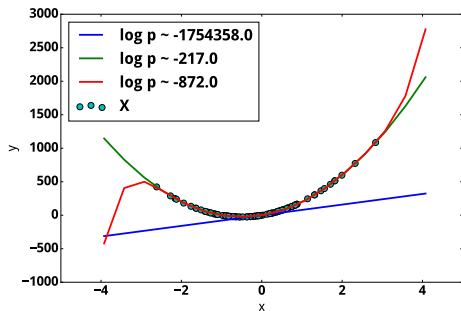
$$L(\theta, \mathbf{A}^{-1}, \mathcal{D}) = -\frac{1}{k} \sum_{r=1}^k \left(\frac{k}{k-1} \log p(\mathcal{D}^r | \mathbf{w}_r) - \log p(\mathbf{w}_r | \mathbf{A}^{-1}) \right).$$

$$Q(\theta, \mathbf{A}^{-1}, \mathcal{D}) = \frac{1}{k} \sum_{r=1}^k k \log p(\mathcal{D}_r | \mathbf{w}_r).$$

Правдоподобие модели

Модель $\mathbf{f} \in \mathfrak{F}$ оптимальна, если достигается максимум правдоподобия модели:

$$p(\mathcal{D}|\mathbf{A}^{-1}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{A}^{-1})d\mathbf{w}.$$



Пусть q — непрерывное распределение.

$$\log p(\mathcal{D}|\mathbf{A}^{-1}) \geq \int q(\mathbf{w})\log \frac{p(\mathcal{D}, \mathbf{w}|\mathbf{A}^{-1})}{q(\mathbf{w})}d\mathbf{w} =$$

$$= \int q(\mathbf{w})\log p(\mathcal{D}|\mathbf{w}, \mathbf{A}^{-1})d\mathbf{w} - D_{\text{KL}},$$

где

$$D_{\text{KL}} = - \int q(\mathbf{w})\log \frac{p(\mathbf{w}|\mathbf{A}^{-1})}{q(\mathbf{w})}d\mathbf{w}.$$

Вариационная оценка на основе мултистарта

$$\log p(\mathcal{D}|\mathbf{A}) \geq \mathbb{E}_{q(\mathbf{w})}[\log p(\mathcal{D}, \mathbf{w}|\mathbf{A}^{-1})] - S(q(\mathbf{w})),$$

S — энтропия.

Теорема [Бахтеев, 2016]. Пусть L — функция потерь, градиент которой — непрерывно-дифференцируемая функция с константой Липшица C . Пусть $\boldsymbol{\theta} = [\mathbf{w}^1, \dots, \mathbf{w}^k]$ — начальные приближения оптимизации модели. Пусть γ — шаг градиентного спуска, такой что:

- $\gamma < \frac{1}{C}$,
- $\gamma^{(-1)} > \max_{r \in \{1, \dots, k\}} \lambda_{\max}(\mathbf{H}(\mathbf{w}^r))$.

Тогда

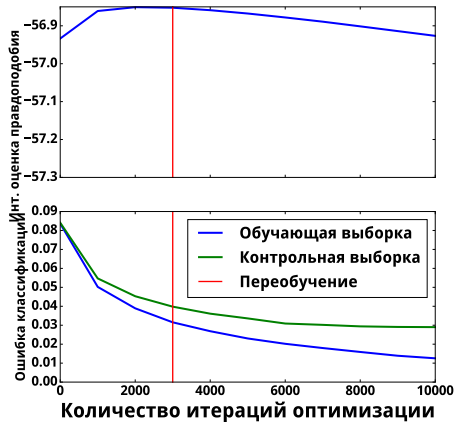
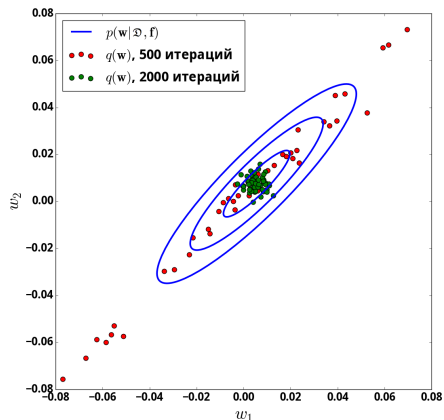
$$S(q^\tau(\mathbf{w})) - S(q^{\tau-1}(\mathbf{w})) \sim \frac{1}{k} \sum_{r=1}^k (\gamma \text{Tr}[\mathbf{H}(\mathbf{w}^r)] - \gamma^2 \text{Tr}[\mathbf{H}(\mathbf{w}^r)\mathbf{H}(\mathbf{w}^r)]) + o_{\gamma \rightarrow 0}(1),$$

где \mathbf{H} — гессиан функции потерь L , q^τ — распределение $q(\mathbf{w})$ в момент оптимизации τ .

Вариационная оценка с использованием градиентного спуска

Максимизация вариационной оценки эквивалентна минимизации

$D_{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}, \mathbf{A}^{-1}))$. Градиентный спуск не минимизирует $D_{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}, \mathbf{A}^{-1}))$.



L и Q: Вариационная оценка

Пусть $q = \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q^{-1})$, $\boldsymbol{\theta} = [\boldsymbol{\mu}_q, \mathbf{A}_q^{-1}]$.

Тогда вариационная оценка имеет вид:

$$\int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathcal{D}, \mathbf{w}, \mathbf{A}^{-1}) d\mathbf{w} - D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{A}^{-1})) \simeq$$

$$\sum_{i=1}^m \log p(\mathbf{x}_i | \mathbf{w}_i) - D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{A}^{-1})) = -L(\boldsymbol{\theta}, \mathbf{A}^{-1}, \mathcal{D}) = Q(\boldsymbol{\theta}, \mathbf{A}^{-1}, \mathcal{D}),$$

где $\mathbf{w}_i \sim q$.

Дивергенция $D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{A}^{-1}))$ вычисляется аналитически:

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

Общая схема алгоритма оптимизации

Вход: количество итераций оптимизации гиперпараметров ℓ , количество итераций оптимизации параметров τ , длина шага градиентного спуска γ .

① Повторять в цикле от $1, \dots, \ell$:

① Инициализировать параметры θ_0 .

② Провести оптимизацию параметров с использованием стохастического градиентного спуска:

$$\hat{\theta} = T \circ T \circ \dots \circ T(\theta_0) = T^\tau(\theta_0),$$

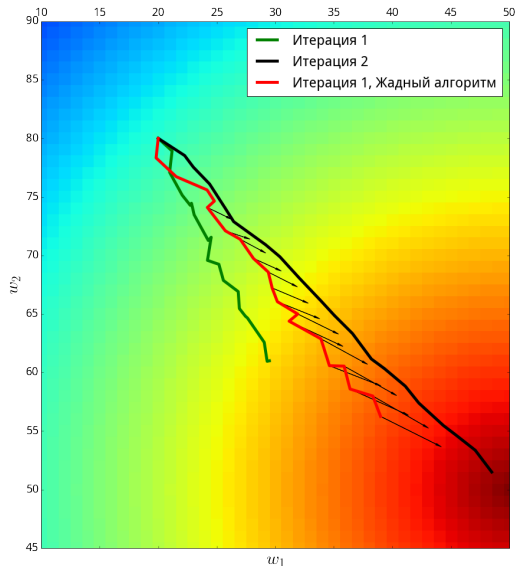
где

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

$\hat{\mathcal{D}}$ — случайная подвыборка \mathcal{D} .

③ Провести оптимизацию гиперпараметров: $Q(\hat{\theta}(\mathbf{A}^{-1}), \mathbf{A}^{-1}, \mathcal{D}) \rightarrow \max$.

НОАГ и Жадный алгоритм



Жадный алгоритм

$$\mathbf{A}'^{-1} = \mathbf{A}^{-1} - \gamma_{\mathbf{A}}(\nabla_{\mathbf{A}^{-1}} Q(T(\theta), \mathbf{A}^{-1}, \mathcal{D})),$$

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

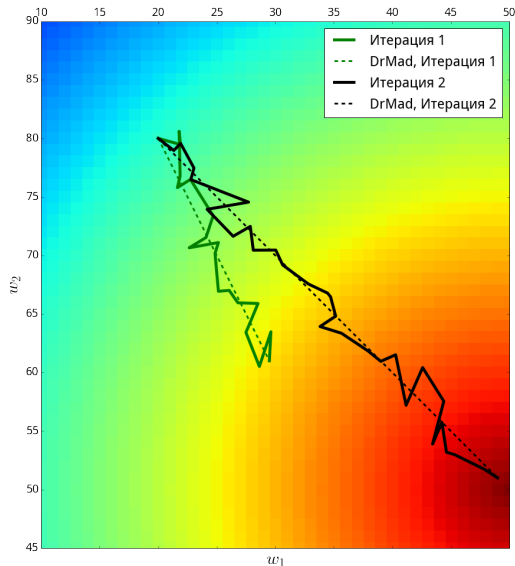
НОАГ

$$\mathbf{A}'^{-1} = \mathbf{A}^{-1} - \gamma_{\mathbf{A}} \hat{\nabla}_{\mathbf{A}^{-1}} Q(T^{\tau}(\theta_0), \mathbf{A}^{-1}, \mathcal{D})),$$

где $\hat{\nabla}_{\mathbf{A}^{-1}}$ — численное приближение градиента:

$$\nabla_{\mathbf{A}^{-1}} Q(\hat{\theta}) = \left(\frac{\partial^2 L}{\partial \theta \partial \mathbf{A}^{-1}} \right)^{\top} \mathbf{H}(L(\theta))^{-1} \nabla_{\theta} Q.$$

Алгоритм DrMad



Рассматривается оптимизация функции $Q(T^T(\theta_0), A^{-1}, \mathcal{D})$ по всей истории оптимизации параметров.

Вводятся предположения о линейности траектории оптимизации параметров. Градиент $\nabla_{A^{-1}} Q$ аккумулируется по правилу:

$$\nabla_{A^{-1}} Q = \nabla_{A^{-1}} Q - \gamma \nabla_{\theta} Q(\theta^T) \nabla_{A^{-1}} \nabla_{\theta} L(\theta^T).$$

Вычислительный эксперимент

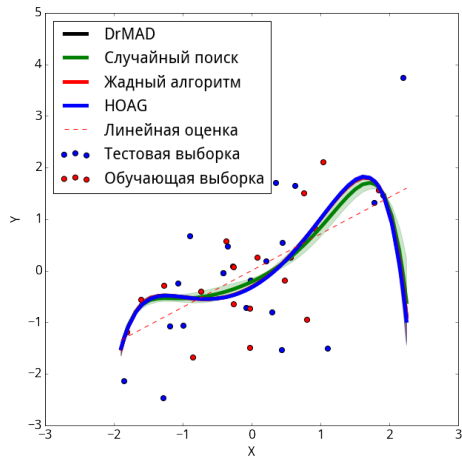
Цель эксперимента: анализ рассматриваемых алгоритмов и итоговых моделей.

Данные:

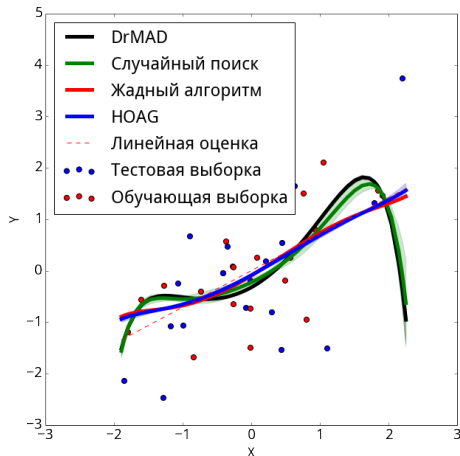
- Синтетические данные: 40 точек на плоскости. Модель: полином 12 степени.
- Набор записей акселерометра WISDM. Рассматривается задача регрессии с нейронной сетью с одним скрытым слоем (10 нейронов).
- Набор рукописных цифр MNIST. Рассматривается задача регрессии с нейронной сетью с одним скрытым слоем (300 нейронов).

В качестве Q и L рассматривается кросс-валидация ($k = 4$) и вариационная оценка.

Синтетические данные: результат

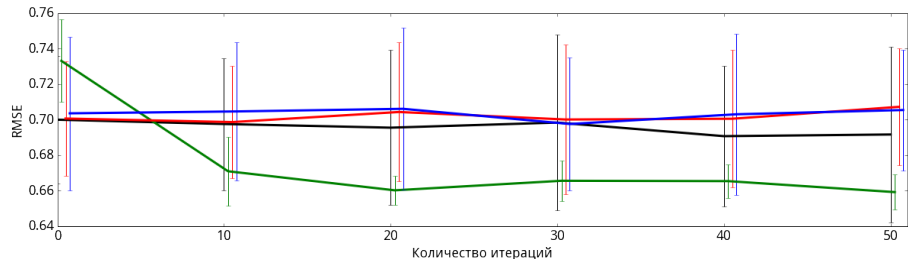
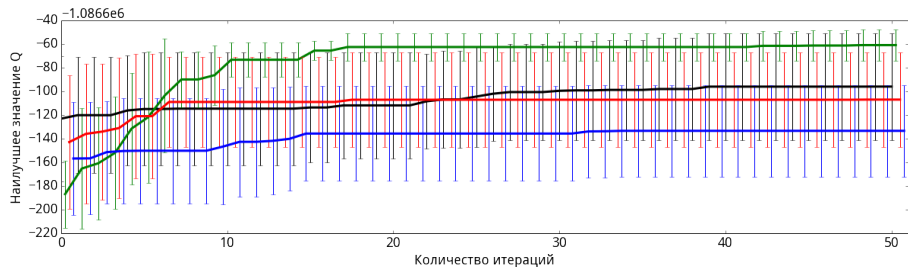


Кросс-Валидация

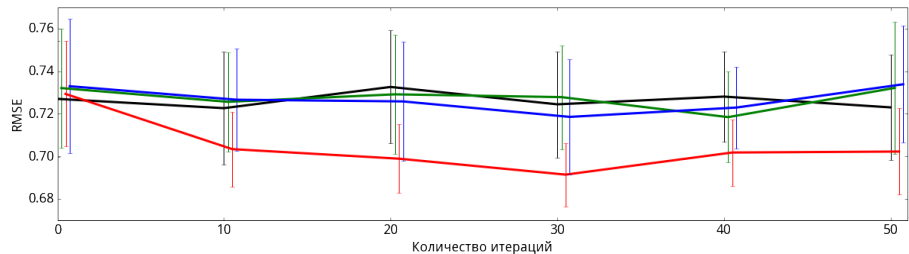
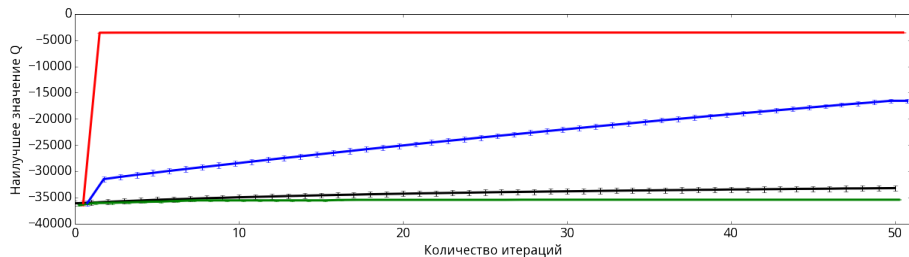


Вариационная оценка

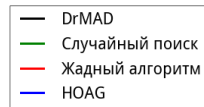
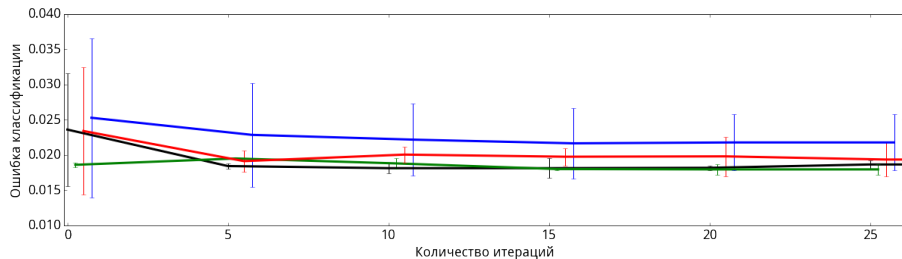
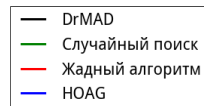
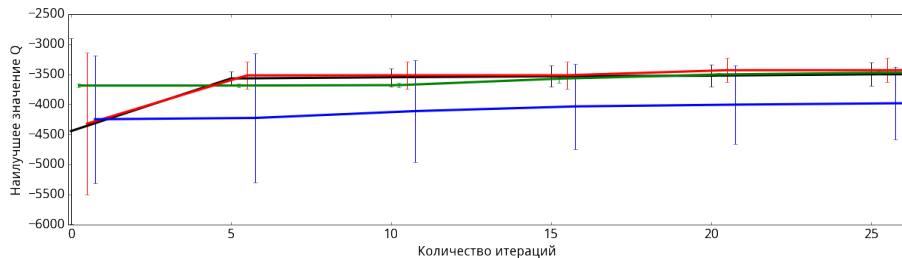
WISDM: кросс-валидация



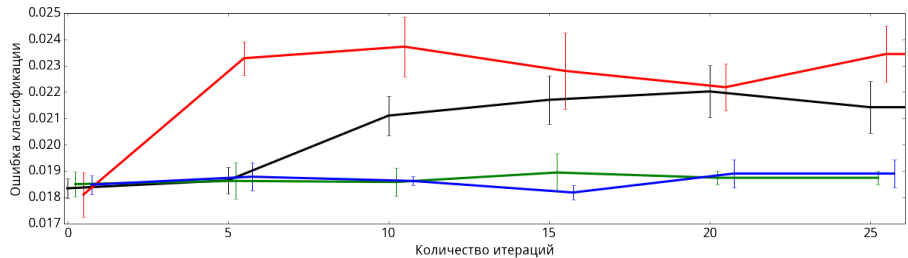
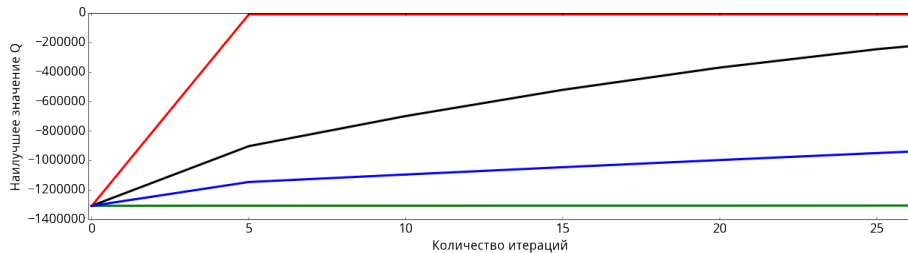
WISDM: вариационная оценка



MNIST: кросс-валидация



MNIST: вариационная оценка



MNIST: добавление шума

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



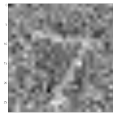
Без шума



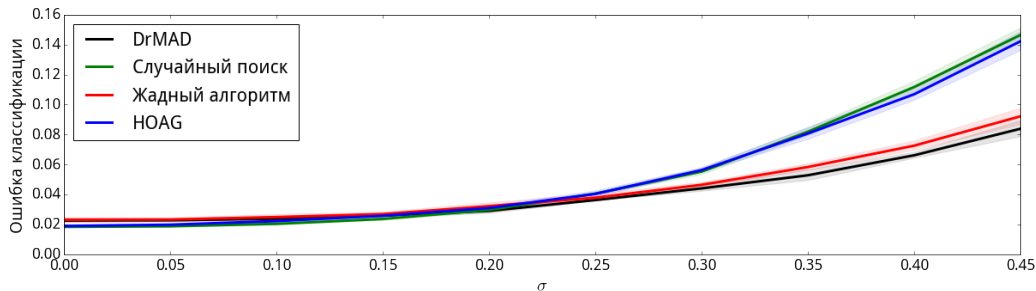
$\sigma = 0.1$



$\sigma = 0.25$



$\sigma = 0.5$



Заключение

- Предложен критерий оптимальной сложности модели глубокого обучения.
- Исследована зависимость интегральной оценки правдоподобия от устойчивости модели и возможности переобучения.
- Предложен алгоритм выбора субоптимальной модели классификации без использования кросс-валидации.
- Рассмотрены градиентные алгоритмы оптимизации гиперпараметров.
- Проведены эксперименты на ряде выборок в задачах классификации и регрессии.
- Наилучшие результаты показали алгоритмы жадной оптимизации.

Исследование основывается на следующих работах

- Graves A. Practical variational inference for neural networks //Advances in Neural Information Processing Systems. – 2011
- Maclaurin D., Duvenaud D., Adams R. Gradient-based hyperparameter optimization through reversible learning //International Conference on Machine Learning. – 2015
- Luketina J. et al. Scalable gradient-based tuning of continuous regularization hyperparameters //International Conference on Machine Learning. - 2016
- J. Fu et al., DrMAD: Distilling Reverse-Mode Automatic Differentiation for Optimizing Hyperparameters of Deep Neural Networks // IJCAI - 2016
- Pedregosa F. Hyperparameter optimization with approximate gradient //International Conference on Machine Learning. – 2016. –