

Метаоптимизация и структура

Бахтеев Олег

МФТИ

12.11.2019

В предыдущих сериях: 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).

В предыдущих сериях: гиперпараметры

Задана дифференцируемая по параметрам модель, приближающая зависимую переменную y :

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

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

Пусть также задано априорное распределение параметров $p(\mathbf{w}|\mathbf{h})$.

Пример:

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

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

В предыдущих сериях: гиперпараметры

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:

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

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

В предыдущих сериях: вариационная оценка

Вариационная оценка Evidence, Evidence lower bound — метод нахождения приближенного значения аналитически невычислимого распределения $p(\mathbf{w}|\mathcal{D}, \mathbf{h})$ распределением $q(\mathbf{w}) \in \Omega$. Получение вариационной нижней оценки обычно сводится к задаче минимизации

$$\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathcal{D})}{q(\mathbf{w})} d\mathbf{w} = \text{E}_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w}) - \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})).$$

Частным случаем вариационного распределения можно считать распределение параметров модели под действием оптимизации.

Метапараметры

Wikipedia

A parameter that controls the value of one or more others.

Определение

Метапараметрами λ модели назовем параметры оптимизации.

Чаще всего метапараметры назначаются экспертно и не подлежат оптимизации в ходе решения задачи выбора модели.

Что можно считать метапараметрами:

- параметры оператора оптимизации;
- параметры задачи оптимизации;
- структуру модели;
- функции активации слоев сети;
- вид априорного распределения и функции правдоподобия.

A neural network that embeds its own meta-levels

Предлагается разделить подмодели внутри модели сети по назначениям:

- “Normal” model: обучение и вывод.
- Evaluation model: оценка качества Q .
- Analyzing model: анализ параметров модели.
- Modifiyng model: модификация параметров.

Представлен градиентный алгоритм оптимизации нейронной сети.

Learning to learn by gradient descent by gradient descent

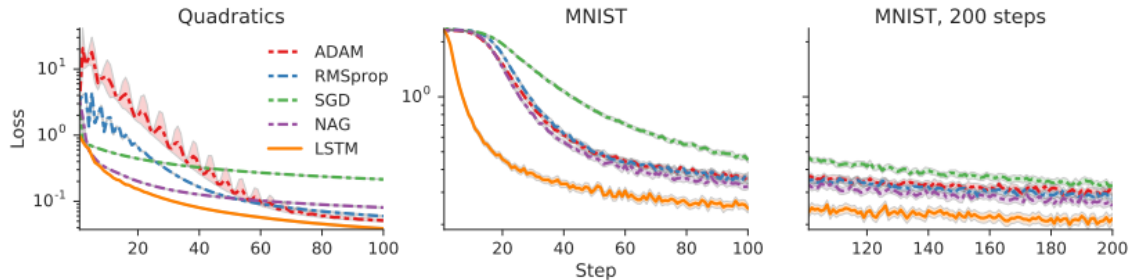
Идея: рассматривать оператор оптимизации T как дифференцируемую функцию:

$$T(\theta) = \text{LSTM}(\theta).$$

Оптимизационная задача:

$$\sum_{t=t_0}^{t_\eta} L(T^t(\theta_{t_0})) \rightarrow \max.$$

LSTM имеет небольшое число параметров и делит параметры между всеми метапараметрами оператора.



Optimal Brain Damage

Рассматривается задача удаления неинформативных параметров.

Идея метода: Разложим функцию потерь в ряд Тейлора в окрестности максимума θ^* :

$$L(\theta^* + \Delta\theta) - L(\theta^*) = -\frac{1}{2}\theta^T \mathbf{H} \theta + o(\|\Delta\theta\|^3),$$

где \mathbf{H} — гессиан функции $-L$.

Для простоты вычисления будем полагать гессиан диагональным. Задача удаления параметров сводится к рассмотрению задач условной оптимизации вида:

$$L(\theta^* + \Delta\theta) \rightarrow \max$$

при

$$\theta_i^* + \Delta\theta_i = 0.$$

Показатель информативности параметра:

$$\frac{\theta_i^2}{2[\mathbf{H}^{-1}]_{i,i}}.$$

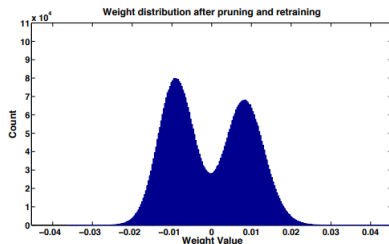
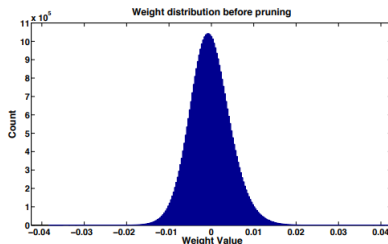
Learning both Weights and Connections for Efficient Neural Networks

Идея подхода:

- 1 Оптимизируем модель;
- 2 Удаляем наименьшие по модулю параметры;
- 3 Запускаем оптимизацию заново.

Почти очевидные факты, которые подтверждаются в статье:

- L_2 лучше для прунинга, чем L_1 в случае, если после прунинга идет оптимизация.
- Оптимизацию лучше производить из предыдущего оптимума, чем из случайной точки.
- После прунинга распределение параметров становится мультимодальным.

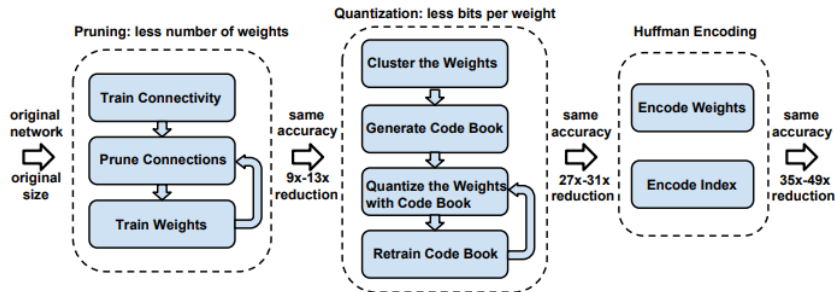


Deep Compression

Идея подхода:

- 1 Удаляем ненужные параметры модели, аналогично предыдущем подходу.
- 2 Кластеризуем параметры (K-means на каждом слое).
- 3 Производим повторную оптимизацию на центроидах.
- 4 Кодлируем индексы параметров с использованием кодов Хаффмана.

Результат: уменьшение размеров модели в 40 раз, ускорение в 3 раза.



Graves, 2011

$$\text{MDL}(\mathbf{f}, \mathcal{D}) = L(\mathbf{f}) + L(\mathcal{D}|\mathbf{f}),$$

где \mathbf{f} — модель, \mathcal{D} — выборка, L — длина описания в битах.

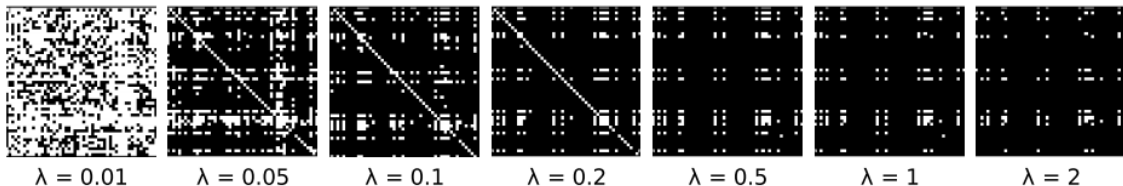
$$\text{MDL}(\mathbf{f}, \mathcal{D}) \sim L(\mathbf{f}) + L(\mathbf{w}^*|\mathbf{f}) + L(\mathcal{D}|\mathbf{w}^*, \mathbf{f}),$$

\mathbf{w}^* — оптимальные параметры модели.

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

Прунинг параметра w_i определяется относительной плотностью:

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



Bayesian Compression for Deep Learning

Модель 1:

$$\mathbf{z} \sim p(\mathbf{z}); \mathbf{w} \sim \mathcal{N}(0; \mathbf{z}^2),$$

где z поставлено в соответствие группе параметров (пример: нейронам).

Априорное распределение $p(z) \propto \frac{1}{|z|^2}$.

Вариационное распределение: $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z), \mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})$

Критерий удаления параметров:

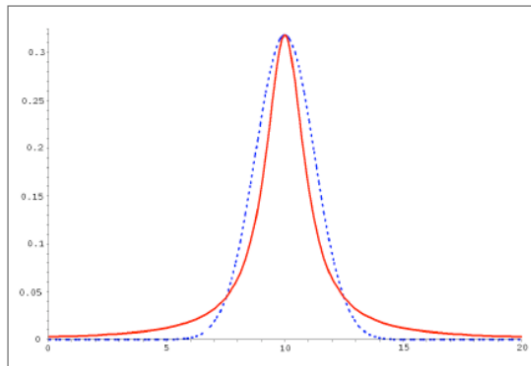
$$\log \sigma_z - \log \mu_z > \lambda.$$

Bayesian Compression for Deep Learning

Модель 2: Априорное распределение:

$$s \sim C^+(0, \lambda_0), z_i \sim C^+(0, 1), \hat{w}_{i,j} \sim \mathcal{N}(0, 1), w_{i,j} = sz_i\hat{w}_{i,j}.$$

Вариационное распределение для z_i : $\mathbf{z}_i \sim \mathcal{L}\}\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})\mathcal{L}\}\mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i)$



[Comparing the Cauchy and Gaussian (Normal) density functions, F. Masci]

Learning the structure of deep sparse graphical models

Рассматривается глубокая генеративная модель.

Структура (т.е. связи параметров) сэмплируются в соответствии с распределением Индийского буфета. Интерпретация распределения: *В ресторане находится конечное количество клиентов и бесконечное количество блюд. j -й клиент берет блюдо k с вероятностью:*

$$\frac{\eta_k}{j + \beta - 1},$$

где η_k количество выборов этого блюда предыдущими клиентами (популярность блюда), а также несколько новых блюд в соответствии с распределением Пуассона с параметром $\frac{\alpha\beta}{j+\beta-1}$.

Порождение параметров и структуры происходит с помощью MCMC.

Главный вывод: структуру можно рассматривать как случайную величину и применять вероятностные методы.

Learning the structure of deep sparse graphical models



(a) $\alpha = 1, \beta = 1$



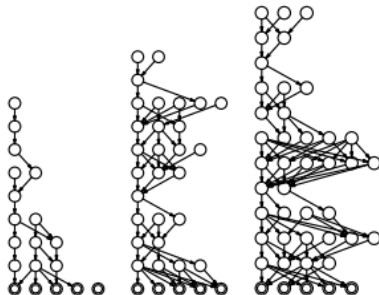
(b) $\alpha = 1, \beta = \frac{1}{2}$



(c) $\alpha = \frac{1}{2}, \beta = 1$



(d) $\alpha = 1, \beta = 2$



(e) $\alpha = \frac{3}{2}, \beta = 1$

Вид вариационной оценки

Обозначим структуру Γ . Пусть априорное распределение параметров зависит от структуры:

$$\mathbf{w} \sim p(\mathbf{w}|\Gamma, \mathbf{h}), \quad \mathbf{w} \sim^q q(\mathbf{w}|\Gamma)$$

Тогда вариационная оценка имеет вид:

$$\mathbb{E}_q \log p(\mathbf{y}|\mathbf{x}, \mathbf{w}, \Gamma) - \text{KL}(q(\Gamma) || p(\Gamma|\mathbf{h})) - \mathbb{E}_{\Gamma \sim q(\Gamma)} \text{KL}(q(\mathbf{w}|\Gamma) || p(\mathbf{w}|\Gamma, \mathbf{h})).$$

Параметрическая сложность

Относительной вариационной плотностью назовем отношение:

$$\rho(w|\mathbf{h}) = \frac{q_w(\text{mode } p(\mathbf{w}|\Gamma, \mathbf{h}, \lambda))}{q_w(\text{mode } q_w)}, \quad \rho = \prod_{w \in \mathbf{w}} \rho(w|\Gamma, \theta_w, \mathbf{h}, \lambda).$$

Определение

Параметрической сложностью модели назовем минимальную дивергенцию между априорным и вариационным распределением:

$$C_p = \min_{\mathbf{h} \in U} D_{\text{KL}}(q||p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})),$$

где U — некоторый компакт.

Утверждение

Пусть задана последовательность $\theta[1], \theta[2], \dots$, такая, что $\lim_{i \rightarrow \infty} C_p(\theta[i]) = 0$. Тогда:

$$\lim_{i \rightarrow \infty} E_{q(\Gamma)} \rho^{-1}(\mathbf{w}|\mathbf{h}[i]) = 1,$$

где $\mathbf{h} = \arg \min_{\mathbf{h} \in U} D_{\text{KL}}(q||p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}))$

Обобщающая задача оптимизации

Какие требования можно выдвинуть к “хорошей” функции оптимизации?

- 1 При некоторых значениях метапараметров функция должна приближать **метод максимального правдоподобия**.
- 2 При некоторых значениях метапараметров функция должна штрафовать **излишне сложные модели**.
- 3 При некоторых значениях метапараметров функция должна приближать обоснованность модели.
- 4 При некоторых значениях метапараметров функция должна позволять **переходить между оптимальными структурами модели**.
- 5 Функции потерь и валидации должны быть непрерывны по метапараметрам.
- 6 Область определения функции должна быть нетривиальна.

Критерии перехода между структурами

критерий 1

Существует нетривиальная константа K и метапараметры λ , такие что для любой пары локальных оптимумов $\mathbf{h}_1, \mathbf{h}_2$ и соответствующих им вариационных параметров $\theta(\mathbf{h}_1), \theta(\mathbf{h}_2)$, таких что $\text{KL}(q(\Gamma|\theta_1)||q(\Gamma|\theta_2)) > K, \text{KL}(q(\Gamma|\theta_2)||q(\Gamma|\theta_1)) > K$ и $Q(\mathbf{h}_1|\lambda) > Q(\mathbf{h}_2|\lambda)$, существует значение λ' , такое что:

- 1 соответствия $\theta(\mathbf{h}_1), \theta(\mathbf{h}_2)$ сохраняются;
- 2 $Q(\mathbf{h}_1|\lambda') < Q(\mathbf{h}_2|\lambda')$.

Критерий 2

Существует нетривиальная константа K и метапараметры λ , такие что для любой пары локальных оптимумов $\mathbf{h}_1, \mathbf{h}_2$ и соответствующих им вариационных параметров $\theta(\mathbf{h}_1), \theta(\mathbf{h}_2)$, таких что $\text{KL}(p(\Gamma|\mathbf{h}_1)||p(\Gamma|\mathbf{h}_2)) > K, \text{KL}(p(\Gamma|\mathbf{h}_2)||p(\Gamma|\mathbf{h}_1)) > K$ и $Q(\mathbf{h}_1|\lambda) > Q(\mathbf{h}_2|\lambda)$, существует значение λ' , такое что:

- 1 соответствия $\theta(\mathbf{h}_1), \theta(\mathbf{h}_2)$ сохраняются;
- 2 $Q(\mathbf{h}_1|\lambda') < Q(\mathbf{h}_2|\lambda')$.

Для $L = Q$ — вариационной нижней оценки критерий 2 не выполняется.

ДЗ: выбор задания

Дедлайн: 20 ноября, 0 часов.

```
from zlib import crc32

theory = crc32('фамилия кириллицей'.lower().encode('utf-8'))%5+1

practice = crc32('фамилия латиницей'.lower().encode('utf-8'))%3+1
```

Задания заливаются на github:

https://github.com/Intelligent-Systems-Phystech/model_selection/фамилия латиницей

ДЗ: теория

Формат: tex + pdf. Задание 1: доказать вид вариационной функции при структуре Γ (расписать дивергенцию).

ДЗ: теория

Формат: tex + pdf. Задание 2: доказать утверждение

Утверждение

Пусть

- 1 Заданы компактные множества $U_h \subset \mathbb{H}$, $U_{\theta_w} \subset \Theta_w$, $U_{\theta_\Gamma} \subset \Theta_\Gamma$.
- 2 Вариационное распределение $q_w(w|\Gamma, \theta_w)$ является абсолютно непрерывным и унимодальным на U_θ . Его мода и матожидание совпадают:
$$\text{mode } q_w(w|\Gamma, \theta_w) = E_{q_w(w|\Gamma, \theta_w)} w.$$
- 3 Априорное распределение $p(w|\Gamma, h, \lambda)$ является абсолютно непрерывным и унимодальным на U_h . Его мода и матожидание совпадают и не зависят от гиперпараметров h на U_h и структуры Γ на U_{θ_Γ} : $E_{p(w|\Gamma, h, \lambda)} w = \text{mode } p(w|\Gamma_1, h_1, \lambda) = \text{mode } p(w|\Gamma_1, h_2, \lambda) = m$ для любых $h_1, h_2 \in U_h, \Gamma_1, \Gamma_2 \in U_\Gamma$.
- 4 Параметры модели w имеют конечные вторые моменты по маргинальным распределениям: $\int_\Gamma q_\Gamma(\Gamma|\theta_\Gamma) q_w(w|\Gamma, \theta_w) d\Gamma$, $\int_\Gamma q_\Gamma(\Gamma|\theta_\Gamma) p(w|\Gamma, h, \lambda) d\Gamma$ при любых $\theta_w \in U_{\theta_w}, \theta_\Gamma \in U_{\theta_\Gamma}, h \in U_h$.
- 5 Вариационное распределение $q_w(w|\Gamma, \theta_w)$ является липшицевым по w .
- 6 Значение $q_w(w|\Gamma, \theta_w)$ не равно нулю при любых $\theta \in U_\theta, \Gamma \in \mathbb{T}$.
- 7 Точная нижняя грань $q_w(m|\Gamma, \theta_w)$ не равна нулю при $\theta_w \in U_{\theta_w}$ и $\Gamma \in \mathbb{T}$:

$$\inf_{\Gamma \in \mathbb{T}, \theta_w \in U_{\theta_w}} q_w(m|\Gamma, \theta_w) > 0.$$

Тогда

$$\left| E_{q_\Gamma(\Gamma|\theta_\Gamma)} p(w|\Gamma, \theta_w, h, \lambda)^{-1} - 1 \right| \leq \text{Const} \iint_{\Gamma, w} |w| \cdot |q_w(w|\Gamma, \theta_w) - p(w|\Gamma, h, \lambda)| q_\Gamma(\Gamma|\theta_\Gamma) dw d\Gamma.$$

ДЗ: теория

Формат: tex + pdf. Задание 3: доказать утверждение

Утверждение

Пусть

- 1 Вариационное распределение $q_w(\mathbf{w}|\Gamma, \theta_w)$ и априорное распределение $p(\mathbf{w}|\Gamma, \mathbf{h}, \lambda)$ являются абсолютно непрерывными.

- 2 Решение задачи

$$\mathbf{h}^* = \arg \min_{\mathbf{h} \in U_h} D_{\text{KL}}(q(\mathbf{w}, \Gamma|\theta) || p(\mathbf{w}, \Gamma|\mathbf{h}, \lambda)) \quad (1)$$

единственно для любого $\theta \in U_\theta$.

- 3 Задана бесконечная последовательность векторов вариационных параметров $\theta[1], \theta[2], \dots, \theta[i], \dots \in U_\theta$, такая что $\lim_{i \rightarrow \infty} C_p(\theta[i] | U_h, \lambda) = 0$.

Тогда следующее выражение стремится к нулю:

$$\iint_{\mathbf{w}, \Gamma} |p(\mathbf{w}|\Gamma, \mathbf{h}[i], \lambda) - q_w(\mathbf{w}|\Gamma, \theta_w[i])| q_\Gamma(\Gamma|\theta_\Gamma[i]) d\Gamma d\mathbf{w},$$

где $\theta[i] = [\theta_w[i], \theta_\Gamma[i]]$, $\mathbf{h}[i]$ — решение задачи (1) для $\theta[i]$.

(Воспользоваться неравенством Пинскера)

ДЗ: теория

Формат: tex + pdf. Задание 4: доказать утверждение

Утверждение

Пусть выполнены условия предыдущих двух утверждений (из д.з.). Тогда справедливо следующее выражение:

$$\lim_{i \rightarrow \infty} E_{q_{\Gamma}(\Gamma|\theta_{\Gamma}[i])} \rho(\mathbf{w}|\Gamma, \theta_{\mathbf{w}}[i], \mathbf{h}[i], \lambda)^{-1} = 1.$$

(Вариант рассуждений:

<https://math.stackexchange.com/questions/112786/convergence-in-law-and-uniformly-integrability>)

ДЗ: теория

Задание 5: доказать утверждение

Утверждение

Пусть q_{Γ} — абсолютно непрерывное распределение с дифференцируемой плотностью, такой что:

- ① Градиент плотности $\nabla_{\theta_{\Gamma}} q(\Gamma|\theta_{\Gamma})$ является ненулевым почти всюду.
- ② Выражение $\nabla_{\theta_{\Gamma}} q(\Gamma|\theta_{\Gamma}) \log p(\Gamma|\mathbf{h}, \lambda)$ ограничено на U_{θ} абсолютно непрерывной случайной величиной, не зависящей от Γ , с конечным первым моментом.
- ③ Не существует значений метапараметров λ_1, λ_2 , таких что:

$$p(\Gamma|\mathbf{h}, \lambda_1) = p(\Gamma|\mathbf{h}, \lambda_2)$$

для всех Γ .

Тогда оптимизация вариационной оценки не удовлетворяет критерию 2 перехода между структурами.

ДЗ: практика

Формат: ірυνь. Реализовать пример удаления параметров для логистической регрессии на MNIST и сравнить качество со случайным удалением параметров (ось X — процент удаленных параметров):

- ① Удаление по модулю с использованием L_1, L_2 - регуляризаций. С переобучением после прунинга и без;
- ② С использованием вариационного вывода (Graves, 2011);
- ③ Optimal Brain Damage.

При оценивании будут учитываться аккуратность кода ноутбуков и наглядность примера.

Используемые материалы

- ① Schmidhuber, Jürgen. "A neural network that embeds its own meta-levels." IEEE International Conference on Neural Networks. IEEE, 1993.
- ② Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." Advances in neural information processing systems. 2016.
- ③ LeCun, Yann, John S. Denker, and Sara A. Solla. "Optimal brain damage." Advances in neural information processing systems. 1990.
- ④ Han, Song, et al. "Learning both weights and connections for efficient neural network." Advances in neural information processing systems. 2015.
- ⑤ Han, Song, Huizi Mao, and William J. Dally. "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding." arXiv preprint arXiv:1510.00149 (2015).
- ⑥ Graves, Alex. "Practical variational inference for neural networks." Advances in neural information processing systems. 2011.
- ⑦ Louizos, Christos, Karen Ullrich, and Max Welling. "Bayesian compression for deep learning." Advances in Neural Information Processing Systems. 2017.
- ⑧ Adams, Ryan, Hanna Wallach, and Zoubin Ghahramani. "Learning the structure of deep sparse graphical models." Proceedings of the thirteenth international conference on artificial intelligence and statistics. 2010.
- ⑨ <http://web.ipac.caltech.edu/staff/fmasci/home/mystats/CauchyVsGaussian.pdf>
- ⑩ Грабовой АВ, Бахтеев ОЮ, Стрижов ВВ. Определение релевантности параметров нейросети. Информатика и её применения. 2019;13(2):62-70.