

# Современные методы оптимизации второго порядка

Руднев Виктор

МГУ имени М.В. Ломоносова, факультет ВМК, кафедра ММП

24 мая 2018 г.

## SGD и его модификации

- сходятся достаточно медленно,
- имеют много гиперпараметров,
- не учитывают структуры того, что оптимизируется.

Поставлена задача регрессии с выборкой  $(X, Y)$

Поставлена задача регрессии с выборкой  $(X, Y)$

$$p(y_i|x_i, \theta) = \mathcal{N}(y_i|f(x_i, \theta), I)$$

Поставлена задача регрессии с выборкой  $(X, Y)$

$$p(y_i|x_i, \theta) = \mathcal{N}(y_i|f(x_i, \theta), I)$$

$$p(Y|X, \theta) = \prod_{i=1}^N \mathcal{N}(y_i|f(x_i, \theta), I)$$

Поставлена задача регрессии с выборкой  $(X, Y)$

$$p(y_i|x_i, \theta) = \mathcal{N}(y_i|f(x_i, \theta), I)$$

$$p(Y|X, \theta) = \prod_{i=1}^N \mathcal{N}(y_i|f(x_i, \theta), I)$$

$$L(\theta) = -\log p(Y|X, \theta) = -\sum_{i=1}^N \log \mathcal{N}(y_i|f(x_i, \theta), I)$$

Поставлена задача регрессии с выборкой  $(X, Y)$

$$p(y_i|x_i, \theta) = \mathcal{N}(y_i|f(x_i, \theta), I)$$

$$p(Y|X, \theta) = \prod_{i=1}^N \mathcal{N}(y_i|f(x_i, \theta), I)$$

$$L(\theta) = -\log p(Y|X, \theta) = -\sum_{i=1}^N \log \mathcal{N}(y_i|f(x_i, \theta), I)$$

$$L(\theta) = \sum_{i=1}^N \frac{1}{2} \|y_i - f(x_i, \theta)\|^2 \rightarrow \min_{\theta}$$

Разложим  $L(\theta)$  до второй степени:



Разложим  $L(\theta)$  до второй степени:

$$L(\theta + h) = L(\theta) + (\nabla_{\theta} L(\theta))^T h + \frac{1}{2} h^T H(\theta) h + o(\|h\|^2),$$

где  $H(\theta) = \frac{\partial^2}{\partial \theta \partial \theta} L(\theta)$  – гессиан.

Разложим  $L(\theta)$  до второй степени:

$$L(\theta + h) = L(\theta) + (\nabla_{\theta} L(\theta))^T h + \frac{1}{2} h^T H(\theta) h + o(\|h\|^2),$$

где  $H(\theta) = \frac{\partial^2}{\partial \theta \partial \theta} L(\theta)$  – гессиан.

$$\arg \min_h \hat{L}(\theta + h) = -H(\theta)^{-1} \nabla_{\theta} L(\theta), \text{ если } H(\theta) \succ 0$$

Разложим  $L(\theta)$  до второй степени:

$$L(\theta + h) = L(\theta) + (\nabla_{\theta} L(\theta))^T h + \frac{1}{2} h^T H(\theta) h + o(\|h\|^2),$$

где  $H(\theta) = \frac{\partial^2}{\partial \theta \partial \theta} L(\theta)$  – гессиан.

$$\arg \min_h \hat{L}(\theta + h) = -H(\theta)^{-1} \nabla_{\theta} L(\theta), \text{ если } H(\theta) \succ 0$$

Тогда параметры лучше всего обновить так:

$$\theta_{t+1} = \theta_t - H(\theta)^{-1} \nabla_{\theta} L(\theta)$$

- сходится очень быстро по шагам, медленно по времени,
- не имеет гиперпараметров,
- не работает, если гессиан не положительно определен,
- требует очень много памяти, так как нужно хранить гессиан,
- не работает с минибатчами.

Распишем гессиан  $L(\theta)$ :

Распишем гессиан  $L(\theta)$ :

$$H(\theta)_{kl} = \frac{\partial^2}{\partial \theta_k \partial \theta_l} L(\theta) = \frac{\partial^2}{\partial \theta_k \partial \theta_l} \frac{1}{2} \sum_{i=1}^N \|y_i - f(x_i, \theta)\|^2 =$$

Распишем гессиан  $L(\theta)$ :

$$H(\theta)_{kl} = \frac{\partial^2}{\partial \theta_k \partial \theta_l} L(\theta) = \frac{\partial^2}{\partial \theta_k \partial \theta_l} \frac{1}{2} \sum_{i=1}^N \|y_i - f(x_i, \theta)\|^2 =$$
$$\sum_{i=1}^N \left\langle \frac{\partial f(x_i, \theta)}{\partial \theta_k}, \frac{\partial f(x_i, \theta)}{\partial \theta_l} \right\rangle + \left\langle f(x_i, \theta) - y_i, \frac{\partial^2 f(x_i, \theta)}{\partial \theta_k \partial \theta_l} \right\rangle$$

Распишем гессиан  $L(\theta)$ :

$$H(\theta)_{kl} = \frac{\partial^2}{\partial \theta_k \partial \theta_l} L(\theta) = \frac{\partial^2}{\partial \theta_k \partial \theta_l} \frac{1}{2} \sum_{i=1}^N \|y_i - f(x_i, \theta)\|^2 =$$
$$\sum_{i=1}^N \left\langle \frac{\partial f(x_i, \theta)}{\partial \theta_k}, \frac{\partial f(x_i, \theta)}{\partial \theta_l} \right\rangle + \left\langle f(x_i, \theta) - y_i, \frac{\partial^2 f(x_i, \theta)}{\partial \theta_k \partial \theta_l} \right\rangle$$

Без второй части — матрица Гаусса-Ньютона:

$$G(\theta) = \sum_{i=1}^N \left( \frac{\partial f(x_i, \theta)}{\partial \theta} \right) \left( \frac{\partial f(x_i, \theta)}{\partial \theta} \right)^T \succeq 0$$



Теперь  $L(\theta)$  можно расписать так:

Теперь  $L(\theta)$  можно расписать так:

$$L(\theta + h) \approx \hat{L}(\theta_h) = L(\theta) + (\nabla_{\theta} L(\theta))^T h + \frac{1}{2} h^T (G(\theta) + \alpha I) h,$$

где  $\alpha$  — коэффициент регуляризации

Теперь  $L(\theta)$  можно расписать так:

$$L(\theta + h) \approx \hat{L}(\theta_h) = L(\theta) + (\nabla_{\theta} L(\theta))^T h + \frac{1}{2} h^T (G(\theta) + \alpha I) h,$$

где  $\alpha$  — коэффициент регуляризации

$$\arg \min_h \hat{L}(\theta + h) = -(G(\theta) + \alpha I)^{-1} \nabla_{\theta} L(\theta)$$

Теперь  $L(\theta)$  можно расписать так:

$$L(\theta + h) \approx \hat{L}(\theta_h) = L(\theta) + (\nabla_{\theta} L(\theta))^T h + \frac{1}{2} h^T (G(\theta) + \alpha I) h,$$

где  $\alpha$  — коэффициент регуляризации

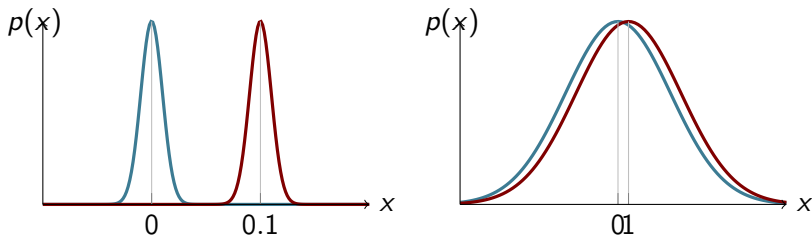
$$\arg \min_h \hat{L}(\theta + h) = -(G(\theta) + \alpha I)^{-1} \nabla_{\theta} L(\theta)$$

Тогда параметры лучше всего обновить так:

$$\theta_{t+1} = \theta_t - (G(\theta) + \alpha I)^{-1} \nabla_{\theta} L(\theta)$$

- сходится очень быстро по шагам, медленно по времени,
- почти не имеет гиперпараметров,
- всё ещё требует очень много памяти, так как нужно хранить матрицу  $G$ ,
- не работает с минибатчами.

# Проблемы с обычной оптимизацией



Две нормальных распределений – пара узких, пара широких. Если сдвинуть широкое на 1, то оно особо не изменится, но если сдвинуть узкое даже на 0.1, то оно будет совершенно другим.

Возьмем расстояние  $\rho(\theta, \theta + h) =$

$$\frac{1}{2} \left( \text{KL}[p(y|x, \theta) || p(y|x, \theta + h)] + \text{KL}[p(y|x, \theta + h) || p(y|x, \theta)] \right) =$$

Возьмем расстояние  $\rho(\theta, \theta + h) =$

$$\frac{1}{2} \left( \text{KL}[p(y|x, \theta) || p(y|x, \theta + h)] + \text{KL}[p(y|x, \theta + h) || p(y|x, \theta)] \right) =$$

$$\frac{1}{2} h^T F(\theta) h + o(\|h\|^2), \text{ где}$$

$$F(\theta) = \mathbb{E}_{p(y|x, \theta)} (\nabla_{\theta} \log p(y|x, \theta)) (\nabla_{\theta} \log p(y|x, \theta))^T -$$

матрица Фишера



# Натуральный градиент

Мы знаем, что с обычным евклидовым расстоянием направление наискорейшего спуска есть просто градиент:

$$\arg \min_{\|h\|_{\text{normal}}=1} L(\theta + h) = -\frac{\nabla_{\theta} L(\theta)}{\sqrt{(\nabla_{\theta} L(\theta))^T \nabla_{\theta} L(\theta)}}$$

# Натуральный градиент

Мы знаем, что с обычным евклидовым расстоянием направление наискорейшего спуска есть просто градиент:

$$\arg \min_{\|h\|_{\text{normal}}=1} L(\theta + h) = -\frac{\nabla_{\theta} L(\theta)}{\sqrt{(\nabla_{\theta} L(\theta))^T \nabla_{\theta} L(\theta)}}$$

С расстоянием, заданным матрицей Фишера, направление наискорейшего спуска чуть отличается:

$$\arg \min_{\|h\|_{\text{natural}}=1} L(\theta + h) = -\frac{F(\theta)^{-1} \nabla_{\theta} L(\theta)}{\sqrt{(\nabla_{\theta} L(\theta))^T F(\theta)^{-1} \nabla_{\theta} L(\theta)}}$$

# Натуральный градиент

Мы знаем, что с обычным евклидовым расстоянием направление наискорейшего спуска есть просто градиент:

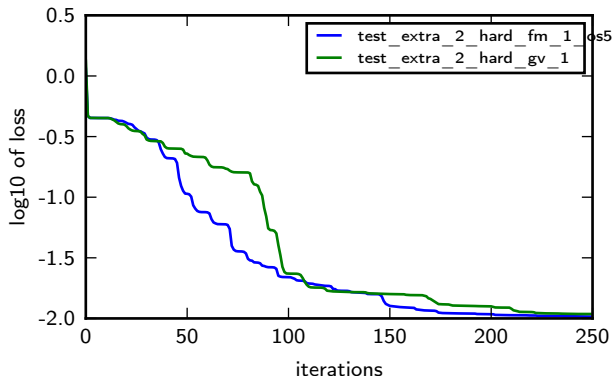
$$\arg \min_{\|h\|_{\text{normal}}=1} L(\theta + h) = -\frac{\nabla_{\theta} L(\theta)}{\sqrt{(\nabla_{\theta} L(\theta))^T \nabla_{\theta} L(\theta)}}$$

С расстоянием, заданным матрицей Фишера, направление наискорейшего спуска чуть отличается:

$$\arg \min_{\|h\|_{\text{natural}}=1} L(\theta + h) = -\frac{F(\theta)^{-1} \nabla_{\theta} L(\theta)}{\sqrt{(\nabla_{\theta} L(\theta))^T F(\theta)^{-1} \nabla_{\theta} L(\theta)}}$$

$\hat{g}(\theta) = (F(\theta_t) + \alpha I)^{-1} \nabla_{\theta_t} L(\theta_t)$  назовем натуральным градиентом, а метод  $\theta_{t+1} = \theta_t - \hat{g}(\theta_t)$  назовем методом натурального градиента.

# Натуральный градиент



Сравнение оптимизаторов на автокодировщике с полносвязной архитектурой 64-15-15-15-15-10 для DIGITS (1000 рукописных цифр размера 8x8) (синим — натуральный градиент, зелёным Гаусс-Ньютон)

# Метод натурального градиента

- сходится очень быстро по шагам, медленно по времени,
- почти не имеет гиперпараметров,
- всё ещё требует очень много памяти, так как нужно хранить матрицу  $F$ ,
- всё ещё не работает с минибатчами,
- для нейросетей есть приближения, которые не требуют памяти, не хранят полную  $F$ , работают значительно быстрее Adam по времени (и работают с минибатчами).

# Кронекерово произведение

$$A = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & \dots & b_{1q} \\ \vdots & \ddots & \vdots \\ b_{p1} & \dots & b_{pq} \end{bmatrix}$$

# Кронекерово произведение

$$A = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & \dots & b_{1q} \\ \vdots & \ddots & \vdots \\ b_{p1} & \dots & b_{pq} \end{bmatrix}$$

$$\text{vec } A = [a_{11} \quad a_{12} \quad \dots \quad a_{1m} \quad a_{21} \quad \dots \quad a_{nm}]^T$$

# Кронекерово произведение

$$A = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & \dots & b_{1q} \\ \vdots & \ddots & \vdots \\ b_{p1} & \dots & b_{pq} \end{bmatrix}$$

$$\text{vec } A = [a_{11} \ a_{12} \ \dots \ a_{1m} \ a_{21} \ \dots \ a_{nm}]^T$$

$$A \otimes B = \begin{bmatrix} a_{11} \mathbf{B} & \dots & a_{1m} \mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{n1} \mathbf{B} & \dots & a_{nm} \mathbf{B} \end{bmatrix}$$



# Кронекерово произведение

$$A \otimes B = \begin{bmatrix} a_{11}b_{11} & \dots & a_{11}b_{1q} & a_{12}b_{11} & \dots & a_{1m}b_{1q} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{11}b_{p1} & \dots & a_{11}b_{pq} & a_{12}b_{p1} & \dots & a_{1m}b_{pq} \\ a_{21}b_{11} & \dots & a_{21}b_{1q} & a_{22}b_{11} & \dots & a_{2m}b_{1q} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{21}b_{p1} & \dots & a_{21}b_{pq} & a_{22}b_{p1} & \dots & a_{2m}b_{pq} \\ a_{n1}b_{11} & \dots & a_{n1}b_{1q} & a_{n2}b_{11} & \dots & a_{nm}b_{1q} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{n1}b_{p1} & \dots & a_{n1}b_{pq} & a_{n2}b_{p1} & \dots & a_{nm}b_{pq} \end{bmatrix}$$

Зададим нейросеть  $f(x, \theta)$  с  $L$  слоями:

$$s_k = W_k a_{k-1}, \quad a_k = \sigma(s_k), \text{ где}$$

$$a_0 = x, \quad f(x, \theta) = a_L,$$

$$\theta = \left[ (\text{vec } W_1)^T; \dots; (\text{vec } W_L)^T \right]^T$$

Зададим нейросеть  $f(x, \theta)$  с  $L$  слоями:

$$s_k = W_k a_{k-1}, \quad a_k = \sigma(s_k), \text{ где}$$

$$a_0 = x, \quad f(x, \theta) = a_L,$$

$$\theta = \left[ (\text{vec } W_1)^T; \dots; (\text{vec } W_L)^T \right]^T$$

Введем обозначение  $\mathcal{D}_Z = -\frac{\partial \log p(y|x, \theta)}{\partial z} = \frac{\partial \frac{1}{2} \|y - f(x, \theta)\|^2}{\partial z}$

Зададим нейросеть  $f(x, \theta)$  с  $L$  слоями:

$$s_k = W_k a_{k-1}, \quad a_k = \sigma(s_k), \text{ где}$$

$$a_0 = x, \quad f(x, \theta) = a_L,$$

$$\theta = \left[ (\text{vec } W_1)^T; \dots; (\text{vec } W_L)^T \right]^T$$

Введем обозначение  $\mathcal{D}z = -\frac{\partial \log p(y|x, \theta)}{\partial z} = \frac{\partial \frac{1}{2} \|y - f(x, \theta)\|^2}{\partial z}$

$$\text{Тогда } F(\theta) = \mathbb{E}_{p(y|x, \theta)} \left[ \mathcal{D}\theta \mathcal{D}\theta^T \right]$$

Зададим нейросеть  $f(x, \theta)$  с  $L$  слоями:

$$s_k = W_k a_{k-1}, \quad a_k = \sigma(s_k), \text{ где}$$

$$a_0 = x, \quad f(x, \theta) = a_L,$$

$$\theta = \left[ (\text{vec } W_1)^T; \dots; (\text{vec } W_L)^T \right]^T$$

Введем обозначение  $\mathcal{D}z = -\frac{\partial \log p(y|x, \theta)}{\partial z} = \frac{\partial \frac{1}{2} \|y - f(x, \theta)\|^2}{\partial z}$

$$\text{Тогда } F(\theta) = \mathbb{E}_{p(y|x, \theta)} \left[ \mathcal{D}\theta \mathcal{D}\theta^T \right]$$

$$F(\theta)_{ij} = \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i)(\text{vec } \mathcal{D}W_j)^T \right]$$

Зададим нейросеть  $f(x, \theta)$  с  $L$  слоями:

$$s_k = W_k a_{k-1}, \quad a_k = \sigma(s_k), \text{ где}$$

$$a_0 = x, \quad f(x, \theta) = a_L,$$

$$\theta = \left[ (\text{vec } W_1)^T; \dots; (\text{vec } W_L)^T \right]^T$$

Введем обозначение  $\mathcal{D}z = -\frac{\partial \log p(y|x, \theta)}{\partial z} = \frac{\partial \frac{1}{2} \|y - f(x, \theta)\|^2}{\partial z}$

$$\text{Тогда } F(\theta) = \mathbb{E}_{p(y|x, \theta)} \left[ \mathcal{D}\theta \mathcal{D}\theta^T \right]$$

$$F(\theta)_{ij} = \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i)(\text{vec } \mathcal{D}W_j)^T \right]$$

$$\mathcal{D}W_k = (\mathcal{D}s_k) a_{k-1}^T = g_k a_{k-1}^T$$

Зададим нейросеть  $f(x, \theta)$  с  $L$  слоями:

$$s_k = W_k a_{k-1}, \quad a_k = \sigma(s_k), \text{ где}$$

$$a_0 = x, \quad f(x, \theta) = a_L,$$

$$\theta = \left[ (\text{vec } W_1)^T; \dots; (\text{vec } W_L)^T \right]^T$$

Введем обозначение  $\mathcal{D}z = -\frac{\partial \log p(y|x, \theta)}{\partial z} = \frac{\partial \frac{1}{2} \|y - f(x, \theta)\|^2}{\partial z}$

$$\text{Тогда } F(\theta) = \mathbb{E}_{p(y|x, \theta)} \left[ \mathcal{D}\theta \mathcal{D}\theta^T \right]$$

$$F(\theta)_{ij} = \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i)(\text{vec } \mathcal{D}W_j)^T \right]$$

$$\mathcal{D}W_k = (\mathcal{D}s_k) a_{k-1}^T = g_k a_{k-1}^T$$

$$\text{vec } \mathcal{D}W_k = a_{k-1} \otimes g_k$$

$$\text{vec } \mathcal{D}W_k = \text{vec } \mathbf{a}_{k-1} \otimes \text{vec } \mathbf{g}_k$$

$$F(\theta)_{ij} = \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i)(\text{vec } \mathcal{D}W_j)^T \right] =$$



$$\begin{aligned}\text{vec } \mathcal{D}W_k &= \text{vec } \mathbf{a}_{k-1} \otimes \text{vec } \mathbf{g}_k \\ F(\theta)_{ij} &= \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i)(\text{vec } \mathcal{D}W_j)^T \right] = \\ &= \mathbb{E}_{p(y|x, \theta)} \left[ (\mathbf{a}_{i-1} \otimes \mathbf{g}_i)(\mathbf{a}_{j-1} \otimes \mathbf{g}_j)^T \right] =\end{aligned}$$

$$\begin{aligned}
 \text{vec } \mathcal{D}W_k &= \text{vec } \mathbf{a}_{k-1} \otimes \text{vec } \mathbf{g}_k \\
 F(\theta)_{ij} &= \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i)(\text{vec } \mathcal{D}W_j)^T \right] = \\
 &= \mathbb{E}_{p(y|x, \theta)} \left[ (\mathbf{a}_{i-1} \otimes \mathbf{g}_i)(\mathbf{a}_{j-1} \otimes \mathbf{g}_j)^T \right] = \\
 &= \mathbb{E}_{p(y|x, \theta)} \left[ (\mathbf{a}_{i-1} \mathbf{a}_{j-1}^T) \otimes (\mathbf{g}_i \mathbf{g}_j^T) \right] \approx
 \end{aligned}$$

$$\begin{aligned}
\text{vec } \mathcal{D}W_k &= \text{vec } \mathbf{a}_{k-1} \otimes \text{vec } \mathbf{g}_k \\
F(\theta)_{ij} &= \mathbb{E}_{p(y|x, \theta)} \left[ (\text{vec } \mathcal{D}W_i) (\text{vec } \mathcal{D}W_j)^T \right] = \\
&\mathbb{E}_{p(y|x, \theta)} \left[ (\mathbf{a}_{i-1} \otimes \mathbf{g}_i) (\mathbf{a}_{j-1} \otimes \mathbf{g}_j)^T \right] = \\
&\mathbb{E}_{p(y|x, \theta)} \left[ (\mathbf{a}_{i-1} \mathbf{a}_{j-1}^T) \otimes (\mathbf{g}_i \mathbf{g}_j^T) \right] \approx \\
\mathbb{E}_{p(y|x, \theta)} \left[ \mathbf{a}_{i-1} \mathbf{a}_{j-1}^T \right] &\otimes \mathbb{E}_{p(y|x, \theta)} \left[ \mathbf{g}_i \mathbf{g}_j^T \right] = \mathbf{A}_{i-1, j-1} \otimes \mathbf{G}_{ij}
\end{aligned}$$

А если считать только блочнодиагональные элементы:

$$F_i = A_{i-1} \otimes G_i, \text{ где } A_{i-1} = \mathbb{E} a_{i-1} a_{i-1}^T, G_i = \mathbb{E} g_i g_i^T$$

А если считать только блочнодиагональные элементы:

$$F_i = A_{i-1} \otimes G_i, \text{ где } A_{i-1} = \mathbb{E} a_{i-1} a_{i-1}^T, G_i = \mathbb{E} g_i g_i^T$$

Более того, мы умеем легко обращать эту матрицу:

$$(F_i)^{-1} \text{vec } \mathcal{D}W_i = G_i^{-1} W_i A_{i-1}^{-1}$$

А если считать только блочнодиагональные элементы:

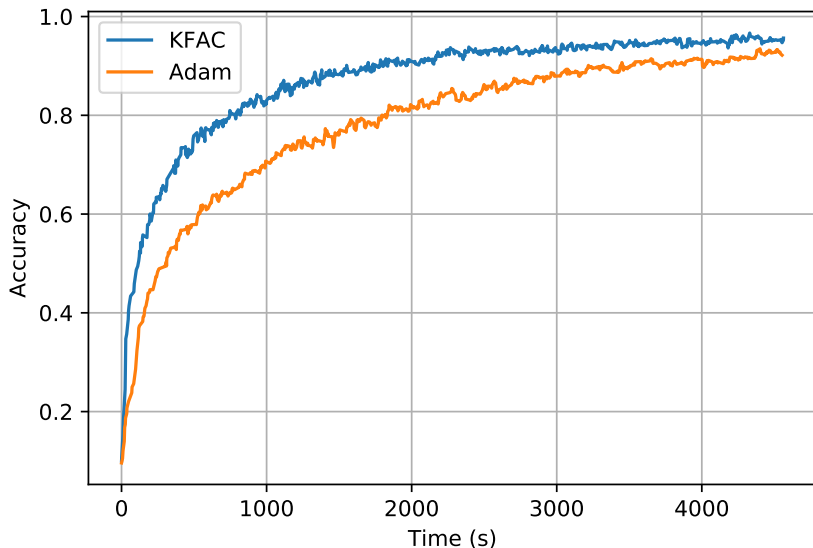
$$F_i = A_{i-1} \otimes G_i, \text{ где } A_{i-1} = \mathbb{E} a_{i-1} a_{i-1}^T, G_i = \mathbb{E} g_i g_i^T$$

Более того, мы умеем легко обращать эту матрицу:

$$(F_i)^{-1} \text{vec } \mathcal{D}W_i = G_i^{-1} W_i A_{i-1}^{-1}$$

Ура, мы умеем считать натуральный градиент по K-FAC

# K-FAC против Adam на CIFAR-10



- сходится очень быстро по шагам, быстро по времени,
- почти не имеет гиперпараметров,
- не требует много памяти,
- работает с минибатчами,
- обобщается на сверточные, рекуррентные слои.
- обобщается на задачи классификации, регрессии, автокодировщика, VAE, GAN, сэмплирования.



- Когда датасет достаточно маленький, а параметров немного, методы Ньютона, Гаусса-Ньютона, нат. градиентов сходятся значительно быстрее, и есть смысл их использовать.
- Для нейросетей есть метод K-FAC, который пока обгоняет Adam и его друзей. Заключается он просто в “отбеливании” градиентов и активаций, участвующих при подсчете градиента, для каждого слоя в отдельности.
- Более того, K-FAC доступен в TensorFlow как встроенный модуль, и его можно подключить в пару строчек.