



Санкт-Петербургский государственный университет

Оптимизационные алгоритмы для нейросетевых гидроаэромеханических расчётов

Егоров Павел Алексеевич

Научный руководитель:
к.ф.-м.н. Гориховский В. И.

Санкт-Петербург
2024

Понятие оптимизатора нейронной сети

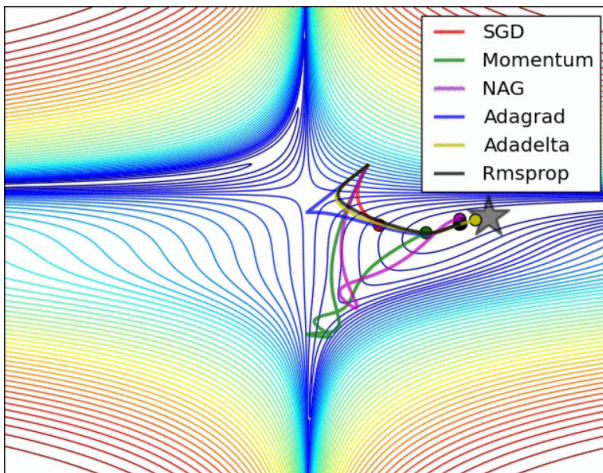


Рис. 1: Принцип работы оптимизационных алгоритмов

- Цель работы: применение различных оптимизационных схем: первого и второго порядков, а также быстрообучаемых оптимизаторов для нахождения оптимальных вычислительных конфигураций с заданным уровнем точности проводимых расчетов
- Используемые технологии: python 3.10, TensorFlow 2.12

Тестовая модель

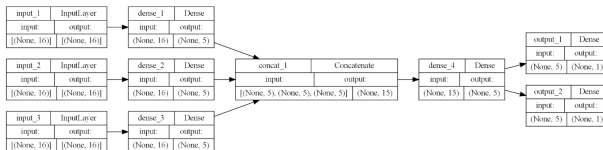


Рис. 2: Тестовая модель

- В качестве функций потерь моделей использовалась средняя абсолютная ошибка
- Выборка делилась на тестовую и обучающую в соотношении один к четырем
- Обучение проводилось на тридцати эпохах

Реализованные алгоритмы

Algorithm 2 AdaMod

Input: initial parameter θ_0 , step sizes $\{\alpha_t\}_{t=1}^T$, moment decay $\{\beta_1, \beta_2, \beta_3\}$, regularization constant ϵ , stochastic objective function $f(\theta_0)$

```
1: Initialize  $m_0 = 0, v_0 = 0, s_0 = 0$ 
2: for  $t = 1$  to  $T$  do
3:    $g_t = \nabla f_t(\theta_{t-1})$ 
4:    $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$ 
5:    $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ 
6:    $\hat{m}_t = m_t / (1 - \beta_1^t)$ 
7:    $\hat{v}_t = v_t / (1 - \beta_2^t)$ 
8:    $\eta_t = \alpha_t / (\sqrt{\hat{v}_t} + \epsilon)$ 
9:    $s_t = \beta_3 s_{t-1} + (1 - \beta_3) \eta_t$ 
10:   $\hat{\eta}_t = \min(\eta_t, s_t)$ 
11:   $\theta_t = \theta_{t-1} - \hat{\eta}_t \hat{m}_t$ 
12: end for
```

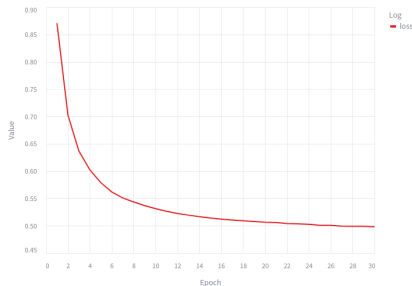


Рис. 3: Adamod

Реализованные алгоритмы

Algorithm 1 MADGRAD

Require: γ_k stepsize sequence, c_k momentum sequence, initial point x_0 , epsilon ϵ

```
1:  $s_0 : d = 0, \nu_0 : d = 0$   
2: for  $k = 0, \dots, T$  do  
3:   Sample  $\xi_k$  and set  $g_k = \nabla f(x_k, \xi_k)$   
4:    $\lambda_k = \gamma_k \sqrt{k+1}$   
5:    $s_{k+1} = s_k + \lambda_k g_k$   
6:    $\nu_{k+1} = \nu_k + \lambda_k (g_k \circ g_k)$   
7:   
$$z_{k+1} = x_0 - \frac{1}{\sqrt[3]{\nu_{k+1}} + \epsilon} \odot s_{k+1}$$
  
8:    $x_{k+1} = (1 - c_{k+1}) x_k + c_{k+1} z_{k+1}$   
9: end for  
10: return  $x_T$ 
```

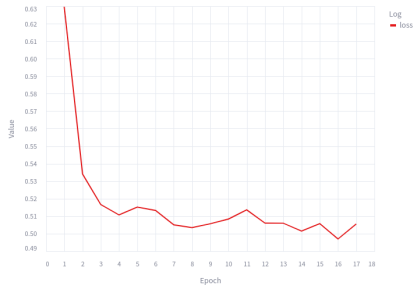


Рис. 4: MADGRAD

Реализованные алгоритмы

Algorithm 1 SGD with LARS. Example with weight decay, momentum and polynomial LR decay.

Parameters: base LR γ_0 , momentum m , weight decay β , LARS coefficient η , number of steps T

Init: $t = 0, v = 0$. Init weight w_0^l for each layer l

while $t < T$ for each layer l **do**

$g_t^l \leftarrow \nabla L(w_t^l)$ (obtain a stochastic gradient for the current mini-batch)

$\gamma_t \leftarrow \gamma_0 * (1 - \frac{t}{T})^2$ (compute the global learning rate)

$\lambda^l \leftarrow \frac{\|w_t^l\|}{\|g_t^l\| + \beta \|w_t^l\|}$ (compute the local LR λ^l)

$v_{t+1}^l \leftarrow mv_t^l + \gamma_{t+1} * \lambda^l * (g_t^l + \beta w_t^l)$ (update the momentum)

$w_{t+1}^l \leftarrow w_t^l - v_{t+1}^l$ (update the weights)

end while

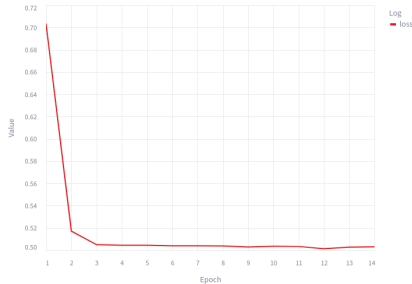


Рис. 5: LARS

Реализованные алгоритмы

Algorithm 1: APOLLO, our proposed algorithm for nonconvex stochastic optimization. All operations on vectors are element-wise. Good default settings are $\beta = 0.9$ and $\epsilon = 1e^{-4}$.

```

Initial:  $m_0, d_0, B_0 \leftarrow 0, 0, 0$  // Initialize  $m_0, d_0, B_0$  to zero
while  $t \in \{0, \dots, T\}$  do
  for  $\theta \in \{\theta^1, \dots, \theta^L\}$  do
     $g_{t+1} \leftarrow \nabla f_t(\theta_t)$  // Calculate gradient at step  $t$ 
     $m_{t+1} \leftarrow \frac{\beta(1-\beta^t)}{1-\beta^{t+1}}m_t + \frac{1-\beta}{1-\beta^{t+1}}g_{t+1}$  // Update bias-corrected moving average
     $\alpha \leftarrow \frac{d_t^T(m_{t+1}-m_t)+d_t^T B_t d_t}{(\|d_t\|+t)^4}$  // Calculate coefficient of  $B$  update
     $B_{t+1} \leftarrow B_t - \alpha \cdot \text{Diag}(d_t^2)$  // Update diagonal Hessian
     $D_{t+1} \leftarrow \text{rectify}(B_{t+1}, 0.01)$  // Handle nonconvexity
     $d_{t+1} \leftarrow D_{t+1}^{-1}m_{t+1}$  // Calculate update direction
     $\theta_{t+1} \leftarrow \theta_t - \eta_{t+1}d_{t+1}$  // Update parameters
  end
end
return  $\theta_T$ 

```

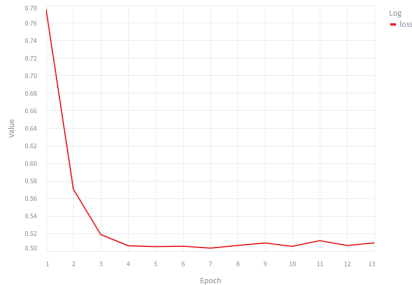


Рис. 6: Apollo

Реализованные алгоритмы

Algorithm 1: ADAHESSIAN

Require: Initial Parameter: θ_0

Require: Learning rate: η

Require: Exponential decay rates: β_1, β_2

Require: Block size: b

Require: Hessian Power: k

Set: $m_0 = 0, v_0 = 0$

for $t = 1, 2, \dots$ **do** // Training Iterations

$\mathbf{g}_t \leftarrow$ current step gradient

$\mathbf{D}_t \leftarrow$ current step estimated diagonal Hessian

 Compute $\mathbf{D}_t^{(s)}$ based on Eq. 10

 Update \mathbf{D}_t based on Eq. 11

 Update m_t, v_t based on Eq. 12

$\theta_t = \theta_{t-1} - \eta m_t / v_t$

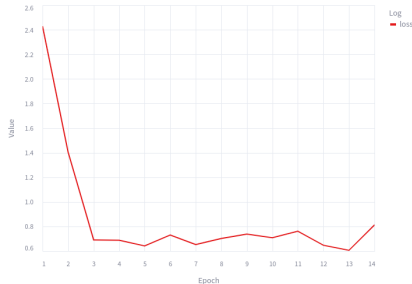


Рис. 7: AdaHessian

Встроенные алгоритмы

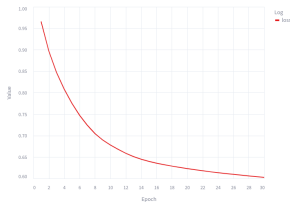


Рис. 8: SGD

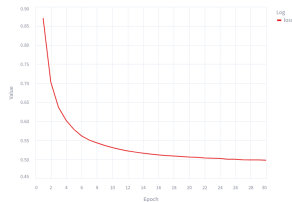


Рис. 9: Adam

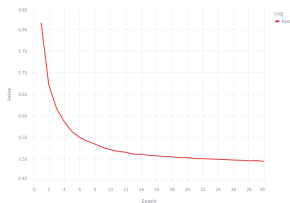


Рис. 10: RMSProp

В ходе учебной практики были реализованы различные оптимизационные алгоритмы:

- 1 AdaMod
- 2 Apollo
- 3 LARS
- 4 MADGRAD
- 5 AdaHessian

Все алгоритмы были интегрированы в приложение и протестированы.