

# Недорогая градиентная оптимизация

Виктор Панкратов

06/12/2021г

# Мотивация

## Описание задачи

Одной из подзадач при работе с нейронными сетями является выбор или оптимизация гиперпараметров. Существуют техники для этого, когда размерность небольшая (отдельно выделяют  $< 5$  и от 5 до порядка 100), но для больших получить решение сложнее.

## Проблема

Веса модели - функции от гиперпараметров. Таким образом оптимизация гиперпараметров сводится к вычислениям якобиана и это слишком сложно для вычислений.

# Основные обозначения

$\mathcal{L}_T, \mathcal{L}_V$  – train и validation losses

$w$  – веса сети

$\lambda$  – гиперпараметры

$\lambda^* = \arg \min_{\lambda} \mathcal{L}_V^*(\lambda)$ , где

$\mathcal{L}_V^*(\lambda) = \mathcal{L}_V(\lambda, w^*(\lambda))$  и  $w^*(\lambda) = \arg \min_w \mathcal{L}_T(\lambda, w)$

# Гиперградиент

$$\begin{aligned}\frac{\partial \mathcal{L}_V^*(\lambda)}{\partial \lambda} &= \left( \frac{\partial \mathcal{L}_V}{\partial \lambda} + \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}^*} \frac{\partial \mathbf{w}^*}{\partial \lambda} \right) \Big|_{\lambda, \mathbf{w}^*(\lambda)} = \\ &= \frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \lambda} + \frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \mathbf{w}^*(\lambda)} \times \frac{\partial \mathbf{w}^*(\lambda)}{\partial \lambda}\end{aligned}$$

- $\frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \lambda}$  – прямая производная гиперпараметров
- Легко считать, но часто равна 0 - приходится обращаться к другому слагаемому
- $\frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \mathbf{w}^*(\lambda)} \times \frac{\partial \mathbf{w}^*(\lambda)}{\partial \lambda}$  – непрямая производная гиперпараметров
- $\frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \mathbf{w}^*(\lambda)}$  – прямая производная параметров
- $\frac{\partial \mathbf{w}^*(\lambda)}{\partial \lambda}$  – best response Jacobian

# IFT

В предположении существования  $\lambda', w'$ :  $\frac{\partial \mathcal{L}_T}{\partial w} \big|_{\lambda', w'} = 0$  в окрестности  $\lambda', w'$  выполнено

$$\frac{\partial w^*}{\partial \lambda} \bigg|_{\lambda'} = - \left[ \frac{\partial^2 \mathcal{L}_T}{\partial w \partial w^T} \right]^{-1} \times \frac{\partial^2 \mathcal{L}_T}{\partial w \partial \lambda} \bigg|_{\lambda', w^*(\lambda')}$$

Как считать обратный гессиан?

$$\left[ \frac{\partial^2 \mathcal{L}_T}{\partial w \partial w^T} \right]^{-1} = \lim_{i \rightarrow \infty} \sum_{j=0}^i \left[ I - \frac{\partial^2 \mathcal{L}_T}{\partial w \partial w^T} \right]^j$$

# Основной алгоритм

---

**Algorithm 1** Gradient-based HO for  $\lambda^*, \mathbf{w}^*(\lambda^*)$ 

---

```
1: Initialize hyperparameters  $\lambda'$  and weights  $\mathbf{w}'$ 
2: while not converged do
3:   for  $k = 1 \dots N$  do
4:      $\mathbf{w}' \leftarrow \alpha \cdot \frac{\partial \mathcal{L}_T}{\partial \mathbf{w}}|_{\lambda', \mathbf{w}'}$ 
5:      $\lambda' \leftarrow \text{hypergradient}(\mathcal{L}_V, \mathcal{L}_T, \lambda', \mathbf{w}')$ 
6: return  $\lambda', \mathbf{w}'$  ▷  $\lambda^*, \mathbf{w}^*(\lambda^*)$  from Eq.1
```

---

---

**Algorithm 2**  $\text{hypergradient}(\mathcal{L}_V, \mathcal{L}_T, \lambda', \mathbf{w}')$ 

---

```
1:  $\mathbf{v}_1 = \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}}|_{\lambda', \mathbf{w}'}$ 
2:  $\mathbf{v}_2 = \text{approxInverseHVP}(\mathbf{v}_1, \frac{\partial \mathcal{L}_T}{\partial \mathbf{w}})$ 
3:  $\mathbf{v}_3 = \text{grad}(\frac{\partial \mathcal{L}_T}{\partial \lambda}, \mathbf{w}, \text{grad\_outputs} = \mathbf{v}_2)$ 
4: return  $\frac{\partial \mathcal{L}_V}{\partial \lambda}|_{\lambda', \mathbf{w}'} - \mathbf{v}_3$  ▷ Return to Alg. 1
```

---

---

**Algorithm 3**  $\text{approxInverseHVP}(\mathbf{v}, \mathbf{f})$ : Neumann approximation of inverse-Hessian-vector product  $\mathbf{v}[\frac{\partial \mathbf{f}}{\partial \mathbf{w}}]^{-1}$ 

---

```
1: Initialize sum  $\mathbf{p} = \mathbf{v}$ 
2: for  $j = 1 \dots i$  do
3:    $\mathbf{v} \leftarrow \alpha \cdot \text{grad}(\mathbf{f}, \mathbf{w}, \text{grad\_outputs} = \mathbf{v})$ 
4:    $\mathbf{p} \leftarrow \mathbf{v}$ 
5: return  $\mathbf{p}$  ▷ Return to Alg. 2
```

---

# Классический подход

---

**Algorithm 1** Stochastic gradient descent with momentum

---

```
1: input: initial  $\mathbf{w}_1$ , decays  $\gamma$ , learning rates  $\alpha$ , loss function  $L(\mathbf{w}, \boldsymbol{\theta}, t)$ 
2: initialize  $\mathbf{v}_1 = \mathbf{0}$ 
3: for  $t = 1$  to  $T$  do
4:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$  ▷ evaluate gradient
5:    $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$  ▷ update velocity
6:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$  ▷ update position
7: end for
8: output trained parameters  $\mathbf{w}_T$ 
```

---

---

**Algorithm 2** Reverse-mode differentiation of SGD

---

```
1: input:  $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$ , train loss  $L(\mathbf{w}, \boldsymbol{\theta}, t)$ , loss  $f(\mathbf{w})$ 
2: initialize  $d\mathbf{v} = \mathbf{0}, d\boldsymbol{\theta} = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$ 
3: initialize  $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$ 
4: for  $t = T$  counting down to 1 do
5:    $d\alpha_t = d\mathbf{w}^T \mathbf{v}_t$ 
6:    $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$ 
7:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$ 
8:    $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$  } exactly reverse
9:    $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$  } gradient descent
10:   $d\gamma_t = d\mathbf{v}^T (\mathbf{v}_t + \mathbf{g}_t)$  } operations
11:   $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$ 
12:   $d\boldsymbol{\theta} = d\boldsymbol{\theta} - (1 - \gamma_t) d\mathbf{v} \nabla_{\boldsymbol{\theta}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \boldsymbol{\theta}, t)$ 
13:   $d\mathbf{v} = \gamma_t d\mathbf{v}$ 
14: end for
15: output gradient of  $f(\mathbf{w}_T)$  w.r.t  $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$  and  $\boldsymbol{\theta}$ 
```

---

# Классический подход

$$\mathbf{w}_{i+1}(\lambda) = \mathbf{w}_i(\lambda) - \alpha \frac{\partial \mathcal{L}_T(\lambda, \mathbf{w}_i(\lambda))}{\partial \mathbf{w}}$$

**Lemma.** *Given the recurrence from unrolling SGD optimization in Eq. 5, we have:*

$$\frac{\partial \mathbf{w}_{i+1}}{\partial \lambda} = - \sum_{j \leq i} \left[ \prod_{k < j} I - \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}^T} \Big|_{\lambda, \mathbf{w}_{i-k}(\lambda)} \right] \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \lambda^T} \Big|_{\lambda, \mathbf{w}_{i-j}(\lambda)}$$



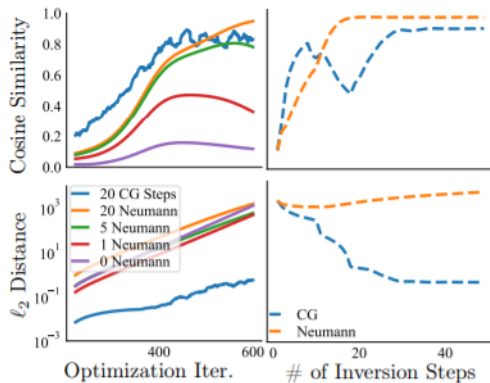


Figure 3: Comparing approximate hypergradients for inverse Hessian approximations to true hypergradients. The Neumann scheme often has greater cosine similarity than CG, but larger  $\ell_2$  distance for equal steps.

Method	Validation	Test	Time(s)
Grid Search	97.32	94.58	100k
Random Search	84.81	81.46	100k
Bayesian Opt.	72.13	69.29	100k
STN	70.30	67.68	25k
<b>No HO</b>	75.72	71.91	18.5k
<b>Ours</b>	69.22	66.40	18.5k
<b>Ours, Many</b>	<b>68.18</b>	<b>66.14</b>	18.5k

Table 4: Comparing HO methods for LSTM training on PTB. We tune millions of hyperparameters faster and with comparable memory to competitors tuning a handful. Our method competitively optimizes the same 7 hyperparameters as baselines from [26] (first four rows). We show a performance boost by tuning millions of hyperparameters, introduced with per-unit/weight dropout and DropConnect. “No HO” shows how the hyperparameter initialization affects training.