# GradInit: Learning to Initialize Neural Networks for Stable and Efficient Training

Докладчик: Еленик Константин

Рецензент: Ким Михаил

Практик-исследователь: Михненко Наталья

Хакер: Сафонов Иван

#### **Abstract**

**GradInit** – метод автоматической инициализации нейронных сетей, работающий для любых архитектур.

Метод основан на эвристике: норма каждого слоя сети корректируется так, чтобы один шаг оптимизатора давал минимальное значение лосса. Корректировка достигается за счет умножения каждой группы параметров сети на обучаемый скаляр.

Авторы показали, что метод ускоряет сходимость и улучшает качество на тесте для различных конволюционных сетей, а также позволяет обучить оригинальный трансформер без warmup.

#### Общая схема

- Все матрицы весов\* изначально инициализируются из нормального распределения с нулевым средним.
- 2. Полученные матрицы фиксируются. Для каждой из них обучается scale factor.
- 3. Каждая из матриц масштабируется
- 4. Модель обучается как обычно (без scale factors)

$$\{W_1,\ldots,W_m\}$$

$$\{\alpha_1 W_1, \ldots, \alpha_m W_m\}$$

$$W_i' \leftarrow \alpha_i W_i$$

$$\{W_1',\ldots,W_m'\}$$

#### Задача оптимизации

Хотим чтобы один шаг оптимизатора давал минимальное значение лосса

$$egin{cases} L(\hat{S}; heta_m - \eta \mathcal{A}[g_{S, heta_m}]) 
ightarrow \min_m \ \|g_{S, heta_m}\|_{p_{\mathcal{A}}} \leq \gamma \end{cases}$$

не даем градиенту стать слишком большим

$$m=\{lpha_1,\dots,lpha_M\}$$
  $S,\hat{S}$  — два разных мини-батча  $heta_m=\{lpha_1W_1,\dots,lpha_MW_M\}$   $\eta,\gamma$  — константные гиперпараметры  $g_{S, heta}=
abla_{ heta}L(S, heta)$   $\mathcal{A}$  — алгоритм оптимизации (Adam/SGD)

### Алгоритм оптимизации

(По сути просто SGD)

#### **Algorithm 1** *GradInit* for learning the initialization of neural networks.

1: **Input:** Target optimization algorithm  $\mathcal{A}$  and learning rate  $\eta$  for model training, initial model parameters  $\theta_0$ , learning rate  $\tau$  of the GradInit scales m, total iterations T, upper bound of the gradient  $\gamma$ , lower bound for the initialization scalars  $\underline{\alpha} = 0.01$ .

```
2: m_1 \leftarrow 1
 3: for t=1 to T do
            Sample S_t from training set.
                                                                                                        оптимизируем либо
           L_t \leftarrow \frac{1}{|S_t|} \sum_{x_k \in S_t} \ell(x_k; \boldsymbol{\theta}_{m_t}), \ \boldsymbol{g}_t \leftarrow \nabla_{\boldsymbol{\theta}} L_t
                                                                                                          норму градиента
 6:
           if \|g_t\|_{p_A} > \gamma then
                 m_{t+1} \leftarrow m_t - \tau \nabla_{m_t} \|\boldsymbol{q}_t\|_{p_A}
                                                                                                       L(\hat{S}; \theta_m - \eta \mathcal{A}[g_{S,\theta_m}]) \to \min_m 
 ||g_{S,\theta_m}||_{p,q} \le \gamma 
            else
                  Sample S_t from training set.
 9:
                  \tilde{L}_{t+1} \leftarrow \frac{1}{|\tilde{S}_t|} \sum_{x_k \in \tilde{S}_t} \ell(x_k; \boldsymbol{\theta}_{m_t} - \eta \mathcal{A}[\boldsymbol{g}_t])
10:
                                                                                                             либо значение лосса
                  m_{t+1} \leftarrow m_t - \tau \nabla_{m_t} \tilde{L}_{t+1}
11:
                                                                                                               после одного шага
            Clamp m_{t+1} using \alpha — — — не даем альфам стать слишком маленькими
12:
```

### Выбор батчей

$$\begin{cases} L(\hat{S}; \underline{\theta_m - \eta \mathcal{A}[g_{S,\theta_m}]}) \to \min_m \\ \|g_{S,\theta_m}\|_{p_{\mathcal{A}}} \le \gamma \end{cases}$$

- Значение градиента может сильно зависеть от батча в начале обучение, так что сэмплировать S и Ŝ независимо плохо.
- Если S = Ŝ, то выгодно просто увеличить норму градиента.
- Поэтому авторы формируют новый батч (Ŝ) половину сэпмплируя из старого батча (S), а половину из остальных данных.

## Выбор констант

$$\begin{cases} L(\hat{S}; \theta_m - \eta \mathcal{A}[g_{S,\theta_m}]) \to \min_m \\ \|g_{S,\theta_m}\|_{p_{\mathcal{A}}} \le \gamma \end{cases}$$

Ограничение нужно, чтобы уменьшение лосса достигалось за счет выбора правильного направления, а не за счет большого шага в менее оптимальном направление.

$$L(S; \theta_{m} - \eta \mathcal{A}[\boldsymbol{g}_{S,\theta_{m}}]) - L(S; \theta_{m}) \approx -\eta \mathcal{A}[\boldsymbol{g}_{S,\theta_{m}}]^{T} \boldsymbol{g}_{S,\theta_{m}} = \begin{cases} -\eta \|\boldsymbol{g}_{S,\theta_{m}}\|_{2}^{2}, & \text{if } \mathcal{A} \text{ is SGD,} \\ -\eta \|\boldsymbol{g}_{S,\theta_{m}}\|_{1}, & \text{if } \mathcal{A} \text{ is Adam.} \end{cases}$$

Рекомендуется выбирать гамму из такого соотношения:

$$\eta \gamma = 0.1$$
 для Adam  $\eta \gamma^2 = 0.1$  для SGD

#### Почему не ...?

$$L(\hat{S}; \theta_m - \eta \mathcal{A}[g_{S,\theta_m}]) + \lambda ||g_{S,\theta_m}|| \to \min_m$$

- градиенты второго порядка
- сложно выбрать универсальную лямбду

## Эксперименты

Table 3: First epoch  $(Acc_1)$  and best test accuracy over all epochs  $(Acc_{best})$  for models on CIFAR-10. We report the mean and standard error of the test accuracies in 4 experiments with different random seeds. Best results in each group are in bold.

Model (# Params)		VGG-19 w/o BN (20.03M)	VGG-19 w/ BN (20.04M)	ResNet-110 w/o BN (1.72M)	ResNet-110 w/ BN (1.73M)	ResNet-1202 w/ BN (19.42M)
Kaiming	$\begin{array}{ c c c } Acc_1 \\ Acc_{best} \end{array}$	$29.1 \pm 1.5$ $94.5 \pm 0.1$	$12.6 \pm 0.6$ $94.4 \pm 0.1$	$16.1 \pm 2.1$ $94.2 \pm 0.1$	$23.2 \pm 0.9$ $95.0 \pm 0.2$	$12.9 \pm 2.8$ $94.4 \pm 0.6$
+1 epoch (Const. LR)	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$37.2 \pm 1.1$ $94.4 \pm 0.1$	$19.6 \pm 4.0$ $94.5 \pm 0.1$	$21.0 \pm 3.8$ $93.9 \pm 0.4$	$32.5 \pm 3.8$ $94.7 \pm 0.3$	$12.6 \pm 2.8$ $94.0 \pm 0.4$
+1 epoch (Warmup)	$\begin{array}{ c c c } Acc_1 \\ Acc_{best} \end{array}$	$37.4 \pm 1.2$ $94.4 \pm 0.1$	$53.5 \pm 2.9$ $94.7 \pm 0.1$	$19.8 \pm 0.5$ $94.1 \pm 0.1$	$48.7 \pm 1.1$ $95.1 \pm 0.1$	$28.1 \pm 1.3$ $95.4 \pm 0.2$
MetaInit	$\begin{array}{ c c c } Acc_1 \\ Acc_{best} \end{array}$	$30.5 \pm 0.9$ $94.6 \pm 0.1$	$35.1 \pm 0.6$ $94.6 \pm 0.1$	$14.6 \pm 2.2$ $94.2 \pm 0.1$	$29.0 \pm 1.5$ $94.8 \pm 0.1$	$11.7 \pm 1.6$ $95.0 \pm 0.5$
GradInit	$\begin{array}{ c c c } Acc_1 \\ Acc_{best} \end{array}$	$29.3 \pm 0.6$ <b>94.7</b> $\pm 0.1$	$47.8 \pm 1.8$ <b>95.1</b> $\pm 0.1$	$36.2 \pm 0.8$ <b>94.6</b> $\pm 0.1$	$38.2 \pm 0.9$ <b>95.4</b> $\pm 0.1$	$29.0 \pm 1.1$ <b>96.2</b> $\pm 0.1$

## Стабилизация

Table 3: First epoch  $(Acc_1)$  and best test accuracy over all epochs  $(Acc_{best})$  for models on CIFAR-10. We report the mean and standard error of the test accuracies in 4 experiments with different random seeds. Best results in each group are in bold.

Model (# Params)		VGG-19 w/o BN (20.03M)	VGG-19 w/ BN (20.04M)	ResNet-110 w/o BN (1.72M)	ResNet-110 w/ BN (1.73M)	ResNet-1202 w/ BN (19.42M)
Kaiming	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$29.1 \pm 1.5$ $94.5 \pm 0.1$	$12.6 \pm 0.6 \\ 94.4 \pm 0.1$	$16.1 \pm 2.1$ $94.2 \pm 0.1$	$23.2 \pm 0.9$ $95.0 \pm 0.2$	$12.9 \pm 2.8$ $94.4 \pm 0.6$
+1 epoch (Const. LR)	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$37.2 \pm 1.1$ $94.4 \pm 0.1$	$19.6 \pm 4.0$ $94.5 \pm 0.1$	$21.0 \pm 3.8$ $93.9 \pm 0.4$	$32.5 \pm 3.8$ $94.7 \pm 0.3$	$12.6 \pm 2.8$ $94.0 \pm 0.4$
+1 epoch (Warmup)	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$37.4 \pm 1.2$ $94.4 \pm 0.1$	$53.5 \pm 2.9$ $94.7 \pm 0.1$	$19.8 \pm 0.5$ $94.1 \pm 0.1$	$48.7 \pm 1.1$ $95.1 \pm 0.1$	$28.1 \pm 1.3$ $95.4 \pm 0.2$
MetaInit	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$30.5 \pm 0.9$ $94.6 \pm 0.1$	$35.1 \pm 0.6$ $94.6 \pm 0.1$	$14.6 \pm 2.2$ $94.2 \pm 0.1$	$29.0 \pm 1.5$ $94.8 \pm 0.1$	$11.7 \pm 1.6$ $95.0 \pm 0.5$
GradInit	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$29.3 \pm 0.6$ <b>94.7</b> $\pm 0.1$	$47.8 \pm 1.8$ <b>95.1</b> $\pm 0.1$	$36.2 \pm 0.8$ <b>94.6</b> $\pm 0.1$	$38.2 \pm 0.9$ <b>95.4</b> $\pm 0.1$	$29.0 \pm 1.1$ <b>96.2</b> $\pm 0.1$

### Deep residual networks

Table 3: First epoch  $(Acc_1)$  and best test accuracy over all epochs  $(Acc_{best})$  for models on CIFAR-10. We report the mean and standard error of the test accuracies in 4 experiments with different random seeds. Best results in each group are in bold.

Model (# Params)		VGG-19 w/o BN (20.03M)	VGG-19 w/ BN (20.04M)	ResNet-110 w/o BN (1.72M)	ResNet-110 w/ BN (1.73M)	ResNet-1202 w/ BN (19.42M)
Kaiming	$\begin{vmatrix} Acc_1 \\ Acc_{best} \end{vmatrix}$	$29.1 \pm 1.5$ $94.5 \pm 0.1$	$12.6 \pm 0.6$ $94.4 \pm 0.1$	$16.1 \pm 2.1$ $94.2 \pm 0.1$	$23.2 \pm 0.9$ $95.0 \pm 0.2$	$12.9 \pm 2.8$ $94.4 \pm 0.6$
+1 epoch (Const. LR)	$\begin{array}{c c} Acc_1 \\ Acc_{best} \end{array}$	$37.2 \pm 1.1$ $94.4 \pm 0.1$	$19.6 \pm 4.0$ $94.5 \pm 0.1$	$21.0 \pm 3.8$ $93.9 \pm 0.4$	$32.5 \pm 3.8$ $94.7 \pm 0.3$	$12.6 \pm 2.8$ $94.0 \pm 0.4$
+1 epoch (Warmup)	$\begin{array}{c c} Acc_1 \\ Acc_{best} \end{array}$	$37.4 \pm 1.2$ $94.4 \pm 0.1$	$53.5 \pm 2.9$ $94.7 \pm 0.1$	$19.8 \pm 0.5$ $94.1 \pm 0.1$	$48.7 \pm 1.1$ $95.1 \pm 0.1$	$28.1 \pm 1.3$ $95.4 \pm 0.2$
MetaInit	$\begin{array}{ c c } Acc_1 \\ Acc_{best} \end{array}$	$30.5 \pm 0.9$ $94.6 \pm 0.1$	$35.1 \pm 0.6$ $94.6 \pm 0.1$	$14.6 \pm 2.2$ $94.2 \pm 0.1$	$29.0 \pm 1.5$ $94.8 \pm 0.1$	$11.7 \pm 1.6$ $95.0 \pm 0.5$
GradInit	$\begin{array}{c c} Acc_1 \\ Acc_{best} \end{array}$	$29.3 \pm 0.6$ <b>94.7</b> $\pm 0.1$	$47.8 \pm 1.8$ <b>95.1</b> $\pm 0.1$	$36.2 \pm 0.8$ <b>94.6</b> $\pm 0.1$	$38.2 \pm 0.9$ <b>95.4</b> $\pm 0.1$	$29.0 \pm 1.1$ $96.2 \pm 0.1$

#### Что-то еще

- рескалирование весов для слоев перед батчнормом влияет на эффективность обучения
- просто добавить scale factors и обучать сеть не работает
- авторы показывают применимость своего метода к трансформерам

#### Итоги

**GradInit** – метод автоматической инициализации нейронных сетей, работающий для любых архитектур.

С помощью модифицированного SGD, метод находит scale factors для каждого блока параметров сети так, чтобы один шаг оптимизатора давал минимальное значение лосса.

Авторы показали, что метод позволяет уменьшить дисперсию градиента, ускоряет сходимость и улучшает качество на тесте для различных конволюционных сетей, а также дает возможность обучить оригинальный трансформер без warmup и даже с использование SGD.

В статье предложен автоматический метод инициализации весов, метод учитывает специфику оптимизатора, величину шага и данные, при этом является model-agnostic.

#### Сильные стороны:

- 1. Предложен эффективный алгоритм решения оптимизационной задачи
- 2. Предложен автоматический метод подбора гиперпараметров
- 3. С помощью большого количества экспериментов было показано, что:
  - Метод уменьшает дисперсию градиентов
  - Ускоряет сходимость
  - Позволяет избавиться от warmup schedule, normalization layers

#### Слабые стороны:

- 1. Слабо исследованы случаи, когда метод неэффективен
- 2. Эксперименты покрывают небольшой набор задач: классификация и машинный перевод
- 3. Непонятна целесообразность в случае трансформеров
- 4. Не исследована применимость для GAN, Inverse RL и других моделей, в которых стабильность обучения критична

#### Актуальность:

- 1. Модели усложняются, процедуры их обучения тоже
- 2. Большинство предшествующих методов не учитывают специфику задачи и(или) специфику оптимизатора
- 3. Вычислительно эффективен

- 1. Идеи в статье изложены доходчиво, однако некоторые графики довольно сложно разглядеть
- 2. Код метода доступен
- 3. В статье перечислены гиперпараметры, необходимые для воспроизведения экспериментов.

**Оценка**: 8

**Уверенность**: 5

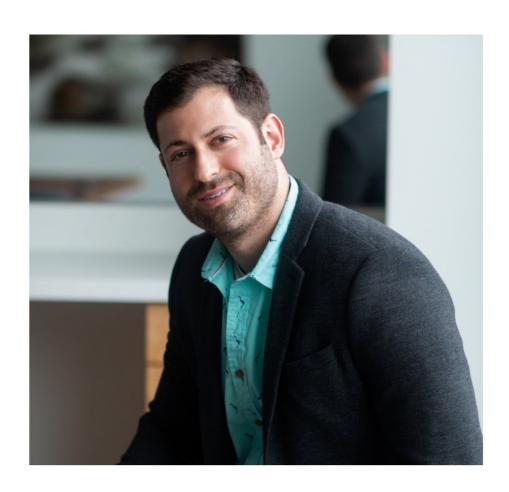
# GradInit

Практик-исследователь: Михненко Наталья

## Информация о статье

- Написана в феврале 2021
- Представлена на конференции NeurLPS (Neural Information Processing Systems)
- Авторы: Chen Zhu, Renkun Ni, Zheng Xu, Kezhi Kong, W. Ronny Huang, Tom Goldstein

## Известные авторы



#### Томас Голдштейн

- Исследования лежат на стыке машинного обучения и оптимизации и нацелены на приложения в области компьютерного зрения и обработки сигналов.
- В числе достижений в 2012 году получил премию Ричарда ДиПрима, Young Faculty Award (от управление перспективных исследовательских проектов Министерства обороны США) и стипендию Слоана

# Конкуренты

Model (# Params)		VGG-19 w/o BN (20.03M)	VGG-19 w/ BN (20.04M)	ResNet-110 w/o BN (1.72M)	ResNet-110 w/ BN (1.73M)	ResNet-1202 w/ BN (19.42M)
Kaiming	$Acc_1 Acc_{best}$	$\begin{vmatrix} 29.1 \pm 1.5 \\ 94.5 \pm 0.1 \end{vmatrix}$	$12.6 \pm 0.6 \\ 94.4 \pm 0.1$	$16.1 \pm 2.1 \\ 94.2 \pm 0.1$	$23.2 \pm 0.9$ $95.0 \pm 0.2$	$12.9 \pm 2.8$ $94.4 \pm 0.6$
+1 epoch (Const. LR)	$Acc_1 Acc_{best}$	$\begin{vmatrix} 37.2 \pm 1.1 \\ 94.4 \pm 0.1 \end{vmatrix}$	$19.6 \pm 4.0$ $94.5 \pm 0.1$	$21.0 \pm 3.8$ $93.9 \pm 0.4$	$32.5 \pm 3.8 \\ 94.7 \pm 0.3$	$12.6 \pm 2.8$ $94.0 \pm 0.4$
+1 epoch (Warmup)	$Acc_1 Acc_{best}$	$\begin{vmatrix} 37.4 \pm 1.2 \\ 94.4 \pm 0.1 \end{vmatrix}$	$53.5 \pm 2.9$ $94.7 \pm 0.1$	$19.8 \pm 0.5$ $94.1 \pm 0.1$	$48.7 \pm 1.1$ $95.1 \pm 0.1$	$28.1 \pm 1.3$ $95.4 \pm 0.2$
MetaInit	$Acc_1 Acc_{best}$	$\begin{vmatrix} 30.5 \pm 0.9 \\ 94.6 \pm 0.1 \end{vmatrix}$	$35.1 \pm 0.6$ $94.6 \pm 0.1$	$14.6 \pm 2.2 \\ 94.2 \pm 0.1$	$29.0 \pm 1.5$ $94.8 \pm 0.1$	$11.7 \pm 1.6$ $95.0 \pm 0.5$
GradInit	$Acc_1 Acc_{best}$	$29.3 \pm 0.6$ <b>94.7</b> $\pm 0.1$	$47.8 \pm 1.8$ <b>95.1</b> $\pm 0.1$	$36.2 \pm 0.8$ $94.6 \pm 0.1$	$38.2 \pm 0.9$ <b>95.4</b> $\pm 0.1$	$29.0 \pm 1.1$ $96.2 \pm 0.1$

## Вклад

- Цитат работы очень мало
- Возможно дальнейшее исследование, может ли алгоритм не зависеть от оптимизатора
- Возможность использования в архитектурах, которые раньше не сходились