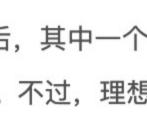
卖萌屋@深度学习炼丹技巧

来自专辑





多大的batch size。 这时候通常要往优化过程着手,比如使用混合精度训练(tensorflow下还可以使用一种叫做bfloat16的 新型浮点格式),即省显存又加速训练;又或者使用更省显存的优化器,比如RMSProp就比Adam更 省显存。本文则介绍AdaFactor,一个由Google提出来的新型优化器,首发论文为《Adafactor: Adaptive Learning Rates with Sublinear Memory Cost .

些缺陷。 Adam

首先我们来回顾一下常用的Adam优化器的更新过程。设t为迭代步数, α_t 为当前学习率, $L(\theta)$ 是损失 函数, θ 是待优化参数, ϵ 则是防止溢出的小正数,那么Adam的更新过程为

$$eta_2 v_{t-1} + (1- i) \ = m_t ig/ ig(1- eta_2^t) \ v_t ig/ ig(1- eta_2^t)$$

要省显存,就首先得知道显存花在哪里的。首先,计算量和显存的大头肯定都是 $\nabla_{\theta}L(\theta_t)$,也就是 说,计算梯度是很费资源的,这也是为啥"ALBERT相比BERT参数量虽然少了那么多,但训练速度也 没见快多少"的原因了;除此之外,显存的消耗主要是m,v了,我们要维护两组缓存变量,来滑动计算

梯度的前两阶矩(也就是
$$m$$
和 v),用以计算参数的更新量。这两组变量每一组都跟训练参数本身一样大,因此对于参数比较多的模型,两组缓存变量所消耗的显存也不少。

AdaFactor

我们知道,CV模型很多时候要靠"SGD+动量"来炼出最优效果来,自适应学习率优化器通常训练不出 最好的效果。但对于NLP模型来说,情况有点相反,自适应学习率显得更重要一些,很少听到由纯靠 SGD调NLP模型的案例。因此,作为省显存的第一步,我们可以抛弃Adam里边的动量,这样就少一 组缓存参数了,自然也就省了显存:

广义KL散度

特别重要。

在SGD中,所有参数都是共用一个标量学习率;在Adam中,则是每一个参数都有自己的学习率 $lpha_t ig/ \sqrt{\hat{v}_t + \epsilon}$ 。我们知道通过精调学习率,SGD其实也能有不错的效果,这表明"每一个参数都有自 己的学习率"这件事情都不是特别重要,或者换一种说法,就是"精调每一个参数自己的学习率"并不是

 $l = \sum_{i,j} c_{i,j} \log rac{c_{i,j}}{a_i b_j} - c_{i,j} + a_i b_j$ 这个度量源自不等式 $x\log x \geq x-1 (\forall x>0)$, 当且仅当x=1时等号成立。所以代入 x = p/q (p, q > 0),然后两端乘以q,我们有

 $a_i b_i \approx c_{i,j}$

 $\sum_i (a_i b_j - c_{i,j})^2$

但在这个距离之下, a_i, b_i 并没有解析解;此外,在优化过程中 $c_{i,j}$ (即 \hat{v}_t)是非负的,而通过上述目

标优化出来的 a_ib_j 无法保证非负,因此很可能扰乱优化过程。原论文的作者们很机智地换了一个度量

标准,使得 a_i, b_j 有解析解。具体来说,它使用了"广义KL散度",又称"I散度",其形式为:

$$egin{cases} rac{\partial l}{\partial a_i} = \sum_j -rac{c_{i,j}}{a_i} + b_j = 0 \ rac{\partial l}{\partial b_j} = \sum_i -rac{c_{i,j}}{b_j} + a_i = 0 \end{cases}$$

直观理解

布,那么 \hat{a}_i , \hat{b}_j 就相当于它们的边缘分布,即

找 $\{a_i\}_{i=1}^m$ 和 $\{b_j\}_{j=1}^n$,使得

 $p\log\frac{p}{q}-p+q\geq 0$ 当且仅当p=q成立,如果p,q有多个分量,那么对多个分量的结果求和即可,这就得到了度量。显 然,广义KL散度是概率的KL散度的自然推广,但它不要求 $c_{i,j}$ 和 a_ib_j 满足归一化,只要求它们非负,

$$egin{cases} a_i \sum_j b_j = \sum_j c_{i,j} \ b_j \sum_i a_i = \sum_i c_{i,j} \end{cases}$$

注意到如果 (a_i,b_j) 是一组最优解,那么 $(\lambda a_i,b_j/\lambda)$ 也是,说白了,所有的 a_i 乘以一个常数,所有的 b_j

也除以这个常数, a_ib_j 是不变的。那么我们就可以随意指定 $\sum_i a_i$ 或 $\sum_j b_j$,因为它们就只是一个缩放

我们也可以从另一个角度理解结果。由于 $c_{i,j}$ 是非负的,我们可以将它归一化,变成具有概率分布的特 性,即 $\hat{c}_{i,j}=rac{c_{i,j}}{\sum c_{i,j}}$,然后我们试图完成分解 $\hat{c}_{i,j}pprox \hat{a}_i\hat{b}_j$,由于 $\hat{c}_{i,j}$ 现在相当于一个二元联合概率分

$$\hat{a}_i = \sum_j \hat{c}_{i,j} = \frac{\sum_j c_{i,j}}{\sum_{i,j} c_{i,j}}, \quad \hat{b}_j = \sum_i \hat{c}_{i,j} = \frac{\sum_i c_{i,j}}{\sum_{i,j} c_{i,j}}$$
 现在 $\hat{c}_{i,j}$ 到 $c_{i,j}$ 还需要乘上一个 $\sum_{i,j} c_{i,j}$,我们可以把它乘到 \hat{a}_i 或 \hat{b}_j 中,不失一般性,我们假设乘到 \hat{a}_i 上,那么就得到(5)。

有了结果(5)后,我们就可以用它来构建更省内存的优化器了,这就是AdaFactor的雏形。简单来说,

当参数 θ 是普通一维向量时,优化过程保持不变;但 θ 是 $m \times n$ 的矩阵时,算出来的梯度 g_t 也是矩阵,

 $\left\{ v_{i,j;t} = v_{i;t}^{(r)} v_{j;t}^{(c)} \middle/ \sum_{j} v_{j;t}^{(c)}
ight.$ (把 ϵ 加到 g_t^2 上去而不是 \hat{v}_t 上去,这是AdaFactor整出来的形式,不是笔者的锅~).

在Adam以及上述AdaFactor雏形中,滑动权重 β_2 都是恒为常数,AdaFactor指出这是不科学的,并提

练不稳定,因为训练后期梯度变小,训练本身趋于稳定,校正学习率的意义就不大了,因此学习率的 校正力度应该变小,并且 $t \to \infty$,学习率最好恒定为常数(这时候相当于退化为SGD),这就要求 为了达到这个目的,AdaFactor采用如下的衰减策略 $\hat{\beta}_{2,t} = 1 - \frac{1}{t^c}$ 它满足 $\hat{eta}_{2,0}=0,\lim_{t\to\infty}\hat{eta}_{2,t}=1$ 。但即便如此,也不是任何c都适合,必须有0< c<1。c>0好理 解,那为什么要c < 1呢?原论文包含了对它的分析,大家可以去读读,但笔者觉得原论文的推导过 于晦涩, 所以这里给出自己的理解。

最后,我们还可以进一步根据参数的模长来校正更新量,这个思路来自LAMB优化器,在之前的文章

《6个派生优化器的简单介绍及其实现》中也介绍过。简单来说,它就是将最后的更新量标准化,然

后乘以参数的模长,说白了,就是不管你怎么折腾,最后的更新量我只要你的方向,而大小由参数本

身的模长和预先设置学习率共同决定,使得所有层所有参数的相对变化程度保持一致。

至此,我们终于可以写出完整版AdaFactor的更新过程了:

 $RMS(u_t) > d$ 时才执行归一化。原论文中的默认参数为

本文介绍了Google提出来的AdaFactor优化器,一个旨在减少显存占用的优化器,并且针对性地分析

关注&星标小夕,带你解锁AI秘籍 订阅号主页下方「撩一下」有惊喜哦

···· 点击查看精选留言

batch_size过小,那么梯度估算本身也带来较大的误差,两者叠加优化过程可能还不收敛。对于预训 练模型来说,batch_size通常还是很大的,所以现在不少预训练模型开始用AdaFactor优化器了;对于 普通的下游任务来说,AdaFactor也可以尝试,但可能需要多炼炼丹,才能搞出由于无脑Adam的效 果。

• Google|突破瓶颈,打造更强大的Transformer • 推荐系统的发展与简单回顾 • ACL2020|FastBERT: 放飞BERT的推理速度 • LayerNorm是Transformer的最优解吗?

阅读原文

一只小狐狸带你解锁炼丹术&NLP秘籍 作者: 苏剑林(来自追一科技,人称"苏神") 前言 自从GPT、BERT等预训练模型流行起来后,其中一个明显的趋势是模型越做越大,因为更大的模型

配合更充分的预训练通常能更有效地刷榜。不过,理想可以无限远,现实通常很局促,有时候模型太 大了,大到哪怕你拥有了大显存的GPU甚至TPU,依然会感到很绝望。比如GPT2最大的版本有15亿 参数,最大版本的T5模型参数量甚至去到了110亿,这等规模的模型,哪怕在TPU集群上也没法跑到

AdaFactor具有自适应学习率的特性,但比RMSProp还要省显存,并且还针对性地解决了Adam的一

 $igg| heta_t = heta_{t-1} - lpha_t \hat{m}_t \Big/ \sqrt{\hat{v}_t + \epsilon}$

样大,因此对于参数比较多的模型,两组缓存变量所消耗的显存也不少。 在这一节中,我们会相对详细地介绍一些AdaFactor优化器,介绍中会设计比较多的公式和推导。如 果只求一个大致了解的读者,可以自行跳过部分数学内容~ 抛弃动量

这其实就是RMSProp的变种,比RMSProp多了 $\hat{v}_t = v_t / (1 - eta_2^t)$ 这一步。 低秩分解 去掉m之后,缓存变量直接减少了一半,但AdaFactor还不满意,它希望保留自适应学习率功能,但

既然要近似,就要有一个度量的标准。很容易想到的标准是欧氏距离,即

推导过程 直接对求偏导数并让偏导数等于0,得 整理得

从而 g_t^2 也是矩阵,这时候我们对 g_t^2 做低秩分解,然后维护两组缓存变量 $v_t^{(r)} \in \mathbb{R}^m, v_t^{(c)} \in \mathbb{R}^n$,分别

滑动权重

出新的策略。

等价形式

为了认识到这一点,我们重写一下Adam的 \hat{v}_t 的更新过程:

 $\hat{v}_t = v_t/(1-eta_2^t)$

上,那么就得到(5)。

AdaFactor雏形

首先,对于 \hat{v}_t 来说,一个很容易想到的方案是所有梯度平方的平均,即: 时, $1-\frac{1}{t^c}<1-\frac{1}{t}$,因此一个简洁的方案是在式中取c<1,AdaFactor默认的c是0.8。

层自适应

 $t o\infty$ 时, $\hat{eta}_{2,t} o 1$ 。

新的衰减策略

其中 $RMS(x) = \sqrt{\frac{1}{n}\sum_{i=1}^n x_i^2}$ 是模长的变种, $\max(1,RMS(u_t)/d)$ 这一步相当于做了个截断,即

置在bert4keras中,方便大家调用。 需要提醒的是,用AdaFactor的时候,batch_size最好大一些,因为本身低秩分解会带来误差,而如果

并解决了Adam的一些缺陷。笔者认为,AdaFactor针对Adam所做的分析相当经典,值得我们认真琢

夕小瑶的卖萌屋

• 卖萌屋上线Arxiv论文速刷神器,直达学术最前沿! • 13个offer, 8家SSP, 谈谈我的秋招经验

 $egin{cases} g_t =
abla_{ heta} L(heta_t) \ m_t = eta_1 m_{t-1} + (1-eta_1) g_t \ v_t = eta_2 v_{t-1} + (1-eta_2) g_t^2 \ \hat{m}_t = m_t / ig(1-eta_1^tig) \ \hat{v}_t = v_t / ig(1-eta_2^tig) \ \end{cases}$

 $egin{cases} g_t =
abla_ heta L(heta_t) \ v_t = eta_2 v_{t-1} + (1-eta_2) g_t^2 \ \hat{v}_t = v_t ig/ ig(1-eta_2^tig) \ heta_t = heta_{t-1} - lpha_t g_t ig/ \sqrt{\hat{v}_t + \epsilon} \end{cases}$

这启发我们,将 \hat{v}_t 换一种参数更少的近似可能也就足够了。而"参数更少的近似",我们就不难想到低 秩分解了。对于 $m \times n$ 的矩阵C,我们希望找到 $m \times k$ 的矩阵A和 $k \times n$ 的矩阵B,使得 $AB \approx C$ 当k足够小时,A、B的参数总量就小于C的参数量。为了"省"到极致,AdaFactor直接让k=1,即寻

把缓存变量v的参数量再压一压。这一次,它用到了矩阵的低秩分解。

这正好对应了AdaFactor的场景。而且巧妙的是,这种情形配上这个目标,刚好有解析解: $a_i = \sum_j c_{i,j}, \quad b_j = rac{\sum\limits_i c_{i,j}}{\sum\limits_i c_{i,j}} \cdots (5)$

其实这个解析解也很形象,就是行、列分别求和,然后相乘,再除以全体的和。

标量而已。不失一般性,我们指定
$$\sum_{j}b_{j}=1$$
,那么就解得 (5) 。 **直观理解**

滑动平均低秩分解后的结果,最后用 $v_t^{(r)},v_t^{(c)}$ 共同调整学习率: $egin{cases} g_{i,j;t} =
abla_{ heta} L(heta_{i,j;t}) \ v_{i;t}^{(r)} = eta_2 v_{t-1;i}^{(r)} + (1-eta_2) \sum_j ig(g_{i,j;t}^2 + \epsilonig) \ v_{j;t}^{(c)} = eta_2 v_{t-1;j}^{(c)} + (1-eta_2) \sum_i ig(g_{i,j;t}^2 + \epsilonig) \end{cases}$

$$\begin{split} &=\frac{\beta_2 v_{t-1} + (1-\beta_2)g_t^2}{1-\beta_2^t} \\ &=\frac{\beta_2 \hat{v}_{t-1} \left(1-\beta_2^{t-1}\right) + (1-\beta_2)g_t^2}{1-\beta_2^t} \\ &=\frac{\beta_2 \frac{1-\beta_2^{t-1}}{1-\beta_2^t}}{1-\beta_2^t} \hat{v}_{t-1} + \left(1-\beta_2\frac{1-\beta_2^{t-1}}{1-\beta_2^t}\right)g_t^2 \\ &\text{所以如果设} \hat{\beta}_{2,t} = \beta_2 \frac{1-\beta_2^{t-1}}{1-\beta_2^t}, \;\; \text{那么更新公式就是} \\ &\qquad \qquad \hat{v}_t = \hat{\beta}_{2,t} \hat{v}_{t-1} + \left(1-\hat{\beta}_{2,t}\right)g_t^2 \\ &\qquad \qquad \hat{\rho}_{2,t} \text{够不够合理呢?} \;\; \text{答案是可能不大够。} \\ &= 0 \text{ 时} \hat{\beta}_{2,t} = 0, \;\; \text{这时候} \hat{v}_t \text{就是} g_t^2, \;\; \text{也就是用} \end{split}$$

实时梯度来校正学习率,这时候校正力度最大;当 $t o \infty$ 时, $\hat{\beta}_{2,t} o \beta_2$,这时候 v_t 是累积梯度平方

与当前梯度平方的加权平均,由于 $\beta_2 < 1$,所以意味着当前梯度的权重 $1 - \beta_2$ 不为0,这可能导致训

 $\hat{v}_t = rac{1}{t} \sum_{i=1}^t g_i^2 = rac{t-1}{t} \hat{v}_{t-1} + rac{1}{t} g_t^2$ 所以这等价于让 $\hat{\beta}_{2,t}=1-\frac{1}{t}$ 。这个方案美中不足的一点是,每一步梯度都是平权的,这不符合直 觉,因为正常来说<mark>越久远的梯度应该越不重要</mark>才对,所以应该适当降低历史部分权重,而当c < 1

$$egin{aligned} v_{j;t}^{(c)} = \hat{eta}_{2,t} v_{t-1;j}^{(c)} + \left(1 - \hat{eta}_{2,t}
ight) \sum_{i} \left(g_{i,j;t}^2 + \epsilon_1
ight) \ \hat{v}_{i,j;t} = v_{i;t}^{(r)} v_{j;t}^{(c)} igg/\sum_{j} v_{j;t}^{(c)} \ u_t = g_t igg/\sqrt{\hat{v}_t} \ \hat{u}_t = u_t / ext{max}(1, RMS(u_t)/d) imes ext{max}(\epsilon_2, RMS(heta_{t-1})) \ heta_t = heta_{t-1} - lpha_t \hat{u}_t \end{aligned}$$

 $egin{aligned} v_{i;t}^{(r)} &= \hat{eta}_{2,t} v_{t-1;i}^{(r)} + \left(1 - \hat{eta}_{2,t}
ight) \sum_{j} \left(g_{i,j;t}^2 + \epsilon_1
ight) \end{aligned}$

磨体味,对有兴趣研究优化问题的读者来说,更是一个不可多得的分析案例。 当然,没有什么绝对能有效的方法,有的只是方法虽好,要想实际有效,依然要用心炼丹。