

实践指南 | 深度学习中常用的tricks

NLP从入门到放弃 今天

作者 | 我要鼓励娜扎@知乎 (已授权) 编辑 | 极市平台

大家好，我是DASOU；

图片分类是非常常见的一个工作任务，但是要想做好不是容易的事情；本文汇集了**11种**图像分类的常用技巧，希望对大家有所帮助；

1 Warmup

学习率是神经网络训练中最重要的超参数之一，针对学习率的技巧有很多。Warmup是在ResNet中提到的一种学习率预热的方法。由于刚开始训练时模型的权重(weights)是随机初始化的，此时选择一个较大的学习率，可能会带来模型的不稳定。

学习率预热就是在刚开始训练的时候先使用一个较小的学习率，训练一些epoches或iterations，等模型稳定时再修改为预先设置的学习率进行训练。ResNet论文中使用一个110层的ResNet在cifar10上训练时，先用0.01的学习率训练直到训练误差低于80%(大概训练了400个iterations)，然后使用0.1的学习率进行训练。

上述的方法是constant warmup，18年Facebook又针对上面的warmup进行了改进，因为从一个很小的学习率一下变为比较大的学习率可能会导致训练误差突然增大。论文提出了**gradual warm up**来解决这个问题，即从最开始的小学习率开始，每个iteration增大一点，直到最初设置的比较大的学习率。

2 Linear scaling learning rate

在凸优化问题中，随着批量的增加，收敛速度会降低，神经网络也有类似的实证结果。随着batch size的增大，处理相同数据量的速度会越来越快，但是达到相同精度所需要的epoch数量越来越多。也就是说，使用相同的epoch时，大batch size训练的模型与小batch size训练的模型相比，验证准确率会减小。

上面提到的gradual warmup是解决此问题的方法之一。另外，linear scaling learning rate也是一种有效的方法。在mini-batch SGD训练时，梯度下降的值是随机的，因为每一个batch的数据是随机选择的。增大batch size不会改变梯度的期望，但是会降低它的方差。也就是说，大batch size会降低梯度中的噪声，所以我们可以增大学习率来加快收敛。

具体做法很简单，比如ResNet原论文中，batch size为256时选择的学习率是0.1，当我们把batch size变为一个较大的数b时，学习率应该变为 $0.1 \times b/256$ 。

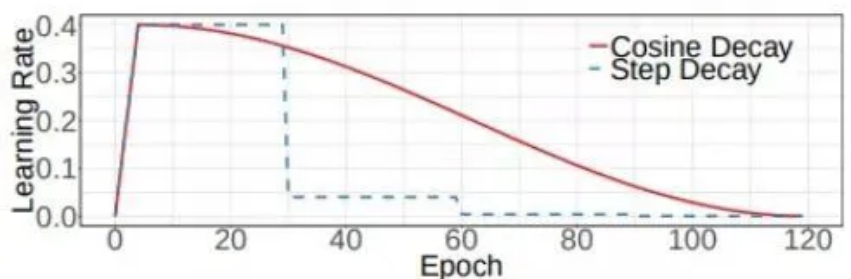
3 Cosine learning rate decay

在warmup之后的训练过程中，学习率不断衰减是一个提高精度的好方法。其中有step decay和cosine decay等，前者是随着epoch增大学习率不断减去一个小的数，后者是让学习率随着训练过程曲线下降。对于cosine decay，假设总共有T个batch（不考虑warmup阶段），在第t个batch时，学习率 η_t 为：

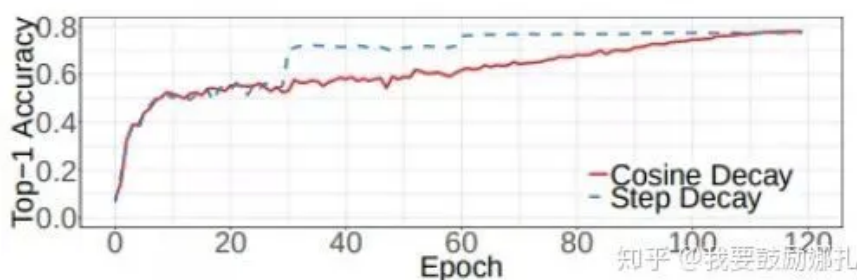
$$\eta_t = \frac{1}{2} \left(1 + \cos \left(\frac{t\pi}{T} \right) \right) \eta$$

这里， η 代表初始设置的学习率。这种学习率递减的方式称之为cosine decay。

下面是带有warmup的学习率衰减的可视化图[4]。其中，图(a)是学习率随epoch增大而下降的图，可以看出cosine decay比step decay更加平滑一点。图(b)是准确率随epoch的变化图，两者最终的准确率没有太大差别，不过cosine decay的学习过程更加平滑。



(a) Learning Rate Schedule



4 Label-smoothing

在分类问题中，我们的最后一层一般是全连接层，然后对应标签的one-hot编码，即把对应类别的值编码为1，其他为0。这种编码方式和通过降低交叉熵损失来调整参数的方式结合起来，会有一些问题。这种方式会鼓励模型对不同类别的输出分数差异非常大，或者说，模型过分相信它的判断。但是，对于一个由多人标注的数据集，不同人标注的准则可能不同，每个人的标注也可能会有一些错误。模型对标签的过分相信会导致过拟合。

标签平滑(Label-smoothing regularization,LSR)是应对该问题的有效方法之一，它的具体思想是降低我们对于标签的信任，例如我们可以将损失的目标值从1稍微降到0.9，或者将从0稍微升到0.

1. 标签平滑最早在inception-v2中被提出，它将真实的概率改造为：

$$q_i = \begin{cases} 1 - \epsilon & \text{if } i = y, \\ \epsilon / (K - 1) & \text{otherwise,} \end{cases}$$

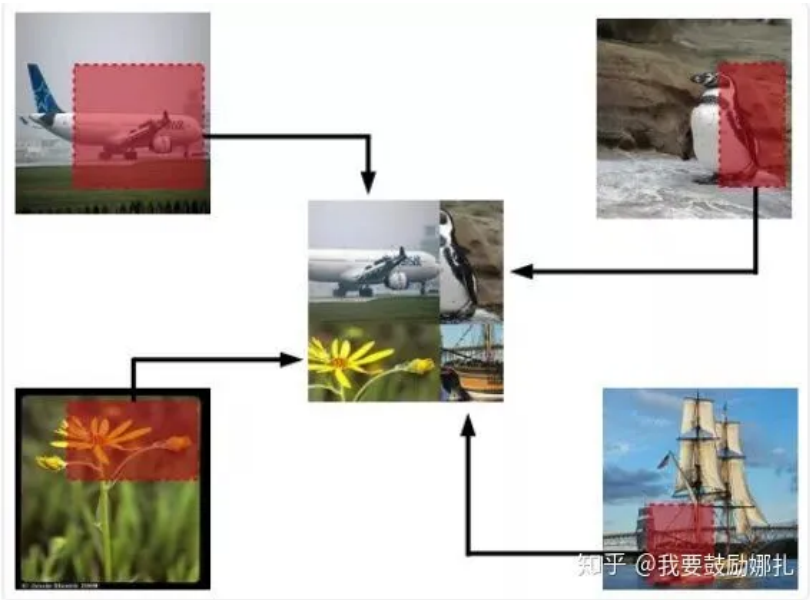
知乎 @我要鼓励娜扎

其中， ϵ 是一个小的常数， K 是类别的数目， y 是图片的真正的标签， i 代表第 i 个类别， q_i 是图片为第 i 类的概率。

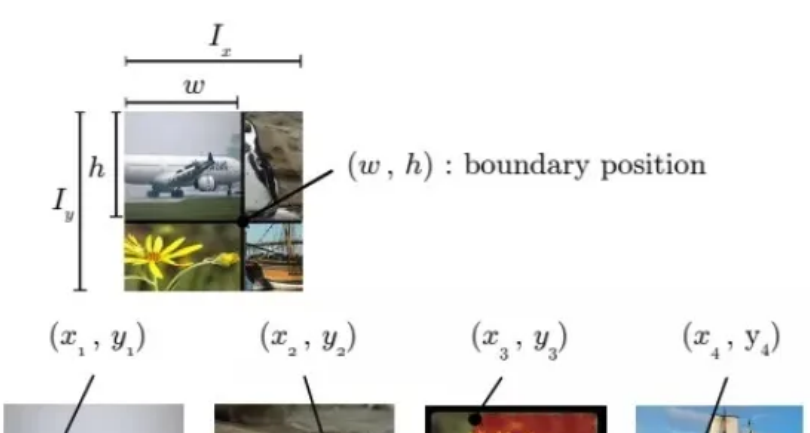
总的来说，**LSR**是一种通过在标签 y 中加入噪声，实现对模型约束，降低模型过拟合程度的一种正则化方法。

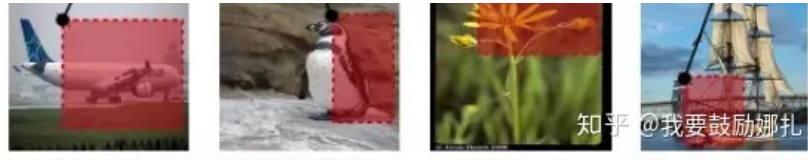
5 Random image cropping and patching

Random image cropping and patching 方法随机裁剪四个图片的中部分，然后把它们拼接为一个图片，同时混合这四个图片的标签。RICAP在caifar10上达到了2.19%的错误率。



如下图所示， I_x, I_y 是原始图片的宽和高。 w 和 h 称为boundary position，它决定了四个裁剪得到的小图片的尺寸。 w 和 h 从beta分布 $\text{Beta}(\beta, \beta)$ 中随机生成， β 也是RICAP的超参数。最终拼接的图片尺寸和原图片尺寸保持一致。





6 Knowledge Distillation

提高几乎所有机器学习算法性能的一种非常简单的方法是在相同的数据上训练许多不同的模型，然后对它们的预测进行平均。但是使用所有的模型集成进行预测是比较麻烦的，并且可能计算量太大而无法部署到大量用户。Knowledge Distillation(知识蒸馏)方法就是应对这种问题的有效方法之一。

****在知识蒸馏方法中，我们使用一个教师模型来帮助当前的模型（学生模型）训练。****教师模型是一个较高准确率的预训练模型，因此学生模型可以在保持模型复杂度不变的情况下提升准确率。比如，可以使用ResNet-152作为教师模型来帮助学生模型ResNet-50训练。在训练过程中，我们会加一个蒸馏损失来惩罚学生模型和教师模型的输出之间的差异。

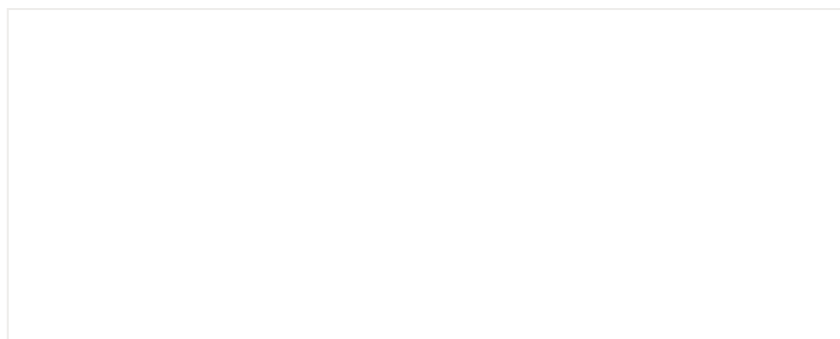
给定输入，假定 p 是真正的概率分布， z 和 r 分别是学生模型和教师模型最后一个全连接层的输出。之前我们会用交叉熵损失 $l(p, \text{softmax}(z))$ 来度量 p 和 z 之间的差异，这里的蒸馏损失同样用交叉熵。所以，使用知识蒸馏方法总的损失函数是

上式中，第一项还是原来的损失函数，第二项是添加的用来惩罚学生模型和教师模型输出差异的蒸馏损失。其中， T 是一个温度超参数，用来使softmax的输出更加平滑的。实验证明，用ResNet-152作为教师模型来训练ResNet-50，可以提高后者的准确率。

7 Cutout

Cutout是一种新的正则化方法。原理是在训练时随机把图片的一部分减掉，这样能提高模型的鲁棒性。它的来源是计算机视觉任务中经常遇到的物体遮挡问题。通过cutout生成一些类似被遮挡的物体，不仅可以让模型在遇到遮挡问题时表现更好，还能让模型在做决定时更多地考虑环境(context)。

效果如下图，每个图片的一小部分被cutout了。





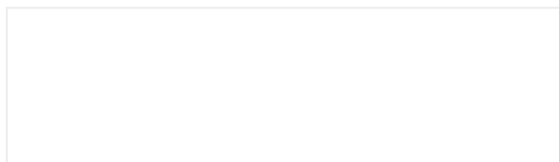
8 Random erasing

Random erasing其实和cutout非常类似，也是一种模拟物体遮挡情况的数据增强方法。区别在于，cutout是把图片中随机抽中的矩形区域的像素值置为0，相当于裁剪掉，random erasing是用随机数或者数据集中像素的平均值替换原来的像素值。而且，cutout每次裁剪掉的区域大小是固定的，Random erasing替换掉的区域大小是随机的。



9 Mixup training

Mixup是一种新的数据增强的方法。Mixup training，就是每次取出2张图片，然后将它们线性组合，得到新的图片，以此来作为新的训练样本，进行网络的训练，如下公式，其中 x 代表图像数据， y 代表标签，则得到的新的 \hat{x} , \hat{y} 。

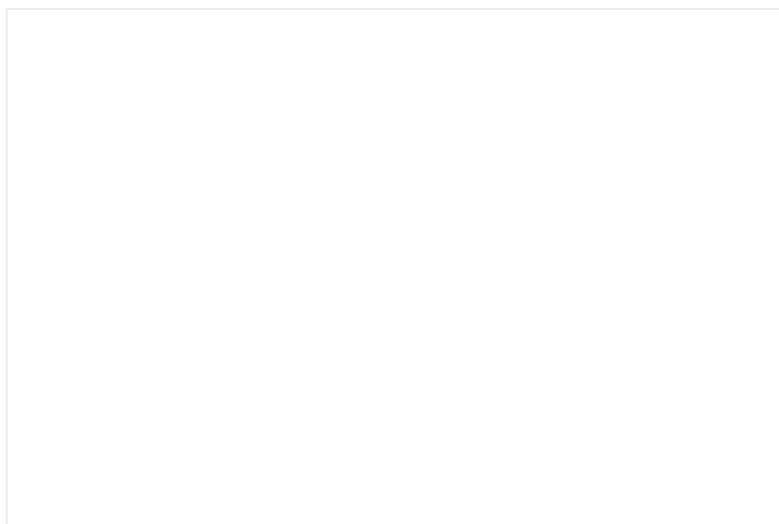


其中， λ 是从Beta(α , α)随机采样的数，在[0,1]之间。在训练过程中，仅使用(\hat{x} , \hat{y})。

Mixup方法主要增强了训练样本之间的线性表达，增强网络的泛化能力，不过mixup方法需要较长的时间才能收敛得比较好。

10 AdaBound

AdaBound，按照作者的说法，会让你的训练过程像adam一样快，并且像SGD一样好。



如下图所示，使用AdaBound会收敛速度更快，过程更平滑，结果更好。

另外，这种方法相对于SGD对超参数的变化不是那么敏感，也就是说鲁棒性更好。但是，针对不同的问题还是需要调节超参数的，只是所用的时间可能变少了。

当然，AdaBound还没有经过普遍的检验，也有可能只是对于某些问题效果好。

11 AutoAugment

数据增强在图像分类问题上有很重要的作用，但是增强的方法有很多，并非一股脑地用上所有的方法就是最好的。那么，如何选择最佳的数据增强方法呢？AutoAugment就是一种搜索适合当前问题的数据增强方法的方法。该方法创建一个数据增强策略的搜索空间，利用搜索算法选取适合特定数据集的数据增强策略。此外，从一个数据集中学到的策略能够很好地迁移到其它相似的数据集上。

- END -

[阅读原文](#)

喜欢此内容的人还喜欢

被李沐大佬虐了

NLP从入门到放弃

把家住成“小森林”，精致美好生活怎么能少了它！

小小种草囤货菌

13年前的禁忌神作，又杀疯了

网易公开课