清华大学软件学院机器学习实验室,专注于迁移学习、深度学习、知识学习等基础理...

THUML





# 迁移学习的预训练-微调范式来有效降低训练成本。迁移学习使得深度神经网络以预训练模型

来越大。幸运的是,在计算机视觉、自然语言处理等人工智能应用的主要领域,人们能够采用

的形式走进千家万户,不用上千块TPU,我们也能够使用BERT、EfficientNet等大型模型。 如今,对于深度学习框架来说,丰富的预训练模型库已经是标配了(例如TensorFlow Hub, Torchvision Models)。在一些研究领域(比如2020年非常热门的自监督学习),研究成果最终 也是以预训练模型的方式呈现给社区。在深度学习社区里,一些热门领域已经积累了成百上千

个预训练模型。 面对众多预训练模型,我们在进行迁移时,该用哪一个好呢?这个重要问题很少有人研究,因 此人们目前只好使用一些简单粗暴的办法: ● 使用常见的预训练模型(例如ResNet50)

如果想要准确地选择最好的预训练模型,我们需要把每一个候选模型都做一遍微调。因为微调 涉及到模型训练,时间至少几个小时起步。有些预训练模型的微调还需要进行超参数搜索,想

- 要决定一个预训练模型的迁移效果就需要将近50个小时!
- 针对这一问题,我们进行了深入探究,提出了一种名为LogME的方法。它能极大地加速预训练 模型选择的过程,将衡量单个预训练模型的时间从50个小时减少到一分钟,疯狂提速三千倍!

论文标题: LogME: Practical Assessment of Pre-trained Models for Transfer Learning 论文链接:

https://arxiv.org/abs/2102.11005 GitHub链接:

的预训练模型。

Target Data

Modality

好。

Pre-train

classification

**预训练模型选择**问题,就是针对用户给定的数据集,从预训练模型库中选择一个最适合的预训 练模型用于迁移学习。其流程可以简单概括为下图,核心就是要对每一个预训练模型进行迁移

Transferability

Assessment

**LEEP** 

NCE

LogME

Potentially Best Model

性评估(Transferability Assessment),简单来说就是为每个模型打分,然后选择出打分最高

无疑是最好的,总是能选出最合适的预训练模型。然而,它的时间开销太大(每个预训练模型 需要50小时), 因此无法实用。 一个好的打分标准,需要在保持与ground-truth打分的高度相关性的同时,尽可能降低时间开 销,才能满足实际使用的要求。除了ground-truth方法之外,目前还有两种打分方法(LEEP和

▲预训练模型选择问题

classification regression vision classification contrastive contrastive regression classification LM language ▲应用场景比较, LogME能胜任几乎所有常见场景

**Target** 

classification

能够多大程度上用于预测这些标注。 征与标注的关系。 说到这里,很多人会想到,一种直观的方法是通过Logistic Regression或者Linear

Regression得到最优权重 $w^*$ ,然后使用似然函数 $p(y \mid F, w^*)$ 作为打分标准。但是这样容易导

致过拟合问题,而且这些方法也有很多超参数需要选择,这使得它们的时间开销很大且效果不

我们选用的是统计学中的证据(evidence,也叫marginalized likelihood)来衡量特征与标注的关

系。它不使用某个特定的 $w^*$ 的值,而是使用的分布来得到边缘化似然的值

 $p(y|F) = \int p(w)p(y|F,w)dw$ 。它相当于取遍了所有可能的w值,能够更加准确地反映特征与

标注的关系,不会有过拟合的问题。其中,p(w)与 $p(y \mid F)$ 分别由超参数 $\alpha$ 和 $\beta$ 决定,但是它们

不需要grid search,可以通过最大化evidence来直接求解。于是,我们就得到了对数最大证 据(Log Maximum Evidence, 缩写LogME)标准来作为预训练模型选择的依据。具体数学细节 不在这里赘述, 感兴趣的读者可以阅读底部的论文。算法的具体细节在下图中给出了。注意, 虽然LogME计算过程中将预训练模型 $\phi$ 视作特征提取器,但是LogME可以用于衡量 $\phi$ 被用于迁 移学习(微调)的性能。

2: Output: logarithm of maximum evidence (LogME) 3: Extract features using pre-trained model  $\phi$ :  $F \in \mathbb{R}^{n \times D}, f_i = \phi(x_i), Y \in \mathbb{R}^{n \times K}$ 4: Compute SVD  $F^TF = V \operatorname{diag} \{\sigma\} V^T$ 5: **for** k = 1 to K **do** Let  $y = Y^{(k)} \in \mathbb{R}^n$ , initialize  $\alpha = 1, \beta = 1$ while  $\alpha, \beta$  not converge do 7: Compute  $\gamma = \sum_{i=1}^{D} \frac{\beta \sigma_i}{\alpha + \beta \sigma_i}$ ,  $\Lambda = \text{diag}\{(\alpha + \beta \sigma)\}$ 8: Naïve:  $A = \alpha I + \beta F^T F, m = \beta A^{-1} F^T y$ 9:

▲LogME算法具体流程 值得一提的是,LogME算法涉及到很多矩阵分解、求逆、相乘操作,因此一不小心就容易使得 算法的复杂度很高(例如上图第9行,粗糙的实现方式)。我们在深入研究该算法后发现,很多 矩阵运算的开销可以通过巧妙的计算优化手段大大降低,因此将计算流程优化为上图第10行,

15: Return LogME  $\frac{1}{K} \sum_{k=1}^{K} \mathcal{L}_k$ 

14: **end for** 

合成数据

-0.5

💋 实验 💋 在实验部分,我们用合成数据、真实数据等多种方式方式,测试了LogME在17个数据集、14 **个预训练模型**上的效果,LogME在这么多数据集、预训练模型上都表现得很好,展现了它优异 的性能。

首先让我们看看,LogME给出的打分标准与人的主观感觉是否一致。我们为分类问题和回归问

题分别设计了一个toy实验,使用生成数据来测量LogME的值。从下图中可以看出,不管是分

类任务还是回归任务,**当特征质量越来越差时,LogME的值也越来越低**,说明LogME可以很

好地衡量特征与标注的关系,从而作为预训练模型选择的标准。

## 15 20 25 5 10 standard deviation of noise 1.0 0.5 0.0 TogM -0.5

接下来,我们用LogME来进行预训练模型选择。我们使用若干个常用预训练模型,通过耗时的

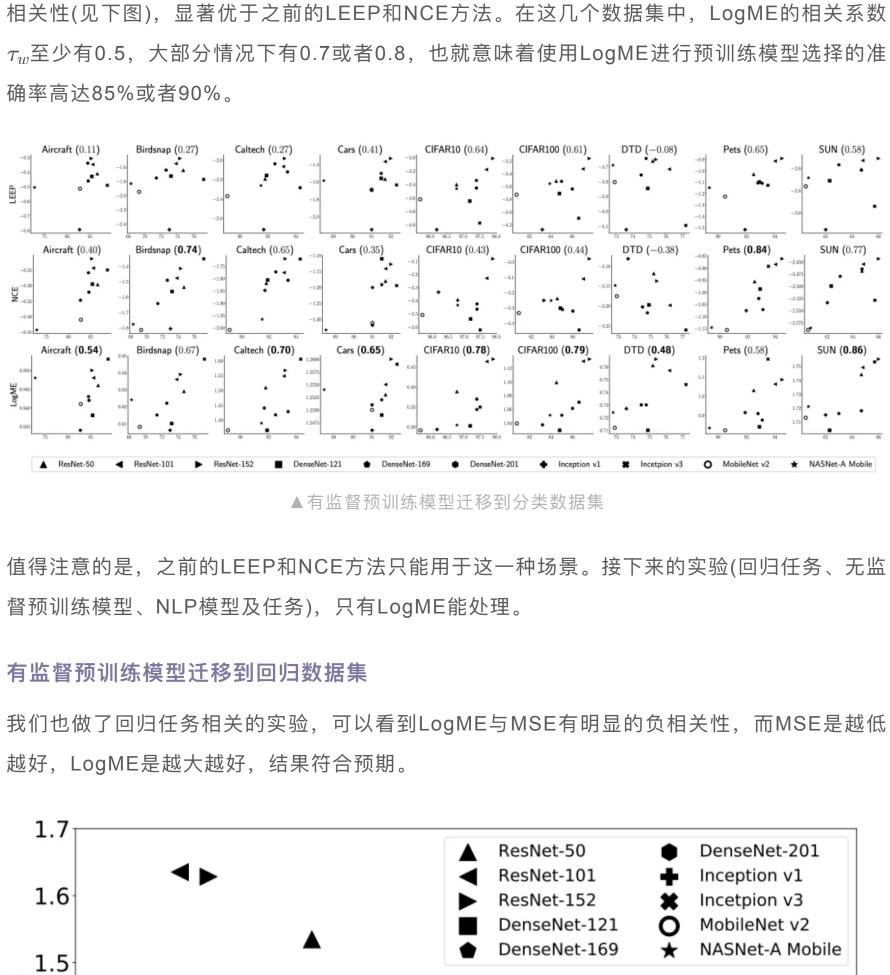
微调过程得到它们的迁移性指标,然后衡量LogME与迁移性指标的相关性。相关性指标 $au_w$ 为

加权肯达尔系数,它的取值范围是[-1,1]。相关系数为 $\tau_w$ 意味着如果LogME认为预训练模型

我们将10个常用预训练模型迁移到9个常见分类数据集中,发现LogME与微调准确率有很高的

 $\phi_1$ 比 $\phi_2$ 好,那么确实 $\phi_1$ 比 $\phi_2$ 好的概率是 $\frac{\tau_w+1}{2}$ 。也就是说, $\tau_w$ 越大越好。

有监督预训练模型迁移到分类数据集



 $\tau_w$ : 1.0  $\tau_w$ : 1.0 ▲使用LogME来衡量无监督预训练模型 自然语言处理任务

LogME并不局限于视觉模型与任务,我们还测试了它对NLP预训练模型的评价能力。可以看

到,在五个任务上,LogME完美地预测了四个预训练模型的表现的相对大小,在另外两个任务

时间加速 🥏

LogME方法不仅效果好,更难得的是它所需要的时间非常短,可以快速评价预训练模型。如果

Accuracy (%)

81.68

84.16

86.99

在2020年,视觉领域的重要进展之一就是无监督预训练模型。因此我们也尝试了使用LogME

来判断无监督预训练模型的质量。从下图的结果来看,不论是分类任务(Aircraft)还是回归任务

Aircraft

LogME

0.93

0.94

0.95

dSprites

LogME

1.52

1.64

1.58

**MSE** 

0.069

0.047

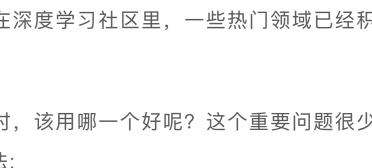
0.050

(dSprites), LogME都能准确衡量无监督预训练模型的质量。

将直接微调的时间作为基准, LogME只需要0.31‰的时间(注意不是百分号, 是千分号), 也就 是说加速了3000倍!而之前的方法如LEEP和NCE,虽然耗时更少,但是效果很差,适用范围 也很有限,完全不如我们的LogME方法。 Wall-clock time (second) Proportion fine-tune (upper bound)  $(1.61 \pm 0.06) \times 10^5$ 1000% $37.3 \pm 0.6$ extract feature (lower bound) 0.23%LEEP (Nguyen et al., 2020)  $37.3 \pm 0.6$ NCE (Tran et al., 2019)  $37.5 \pm 0.6$ LogME (naïve implementation)  $839.8 \pm 5.6$ 

RoBERTa-D

其它领域有所作为。例如,在无监督预训练中,评估一次预训练模型就需要在整个ImageNet 数据集上进行linear protocol evaluation,整个过程需要几个小时。若采用LogME,则只需 要一分钟不到,因此可以在训练过程中将LogME作为early stopping的准则。 后台回复关键词【入群】



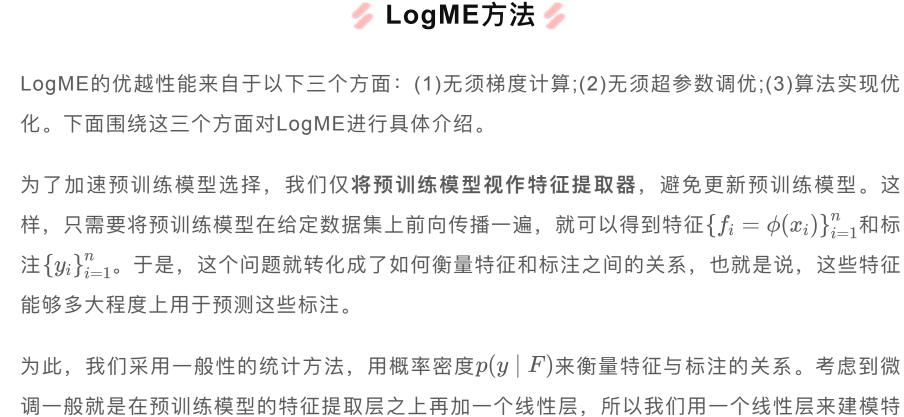
● 使用预训练指标(例如ImageNet准确率)高的模型

目前该论文已被ICML2021接受。

https://github.com/thuml/LogME 🥠 问题描述 🥠

Model Zoo

最直接的打分方法,就是将预训练模型在给定数据集上进行调参、微调,将最终的准确率或者 其它衡量指标作为预训练模型的分数。我们将这种方法称为ground-truth方法,它的选择效果 NCE), 但是它们的使用范围非常有限,只能用于有监督预训练模型迁移到分类任务的场景, 如下表所示,而我们提出的**LogME则能够胜任几乎所有常见的场景**,覆盖了视觉、NLP、分 类、回归、有监督预训练模型、无监督预训练模型等方向。

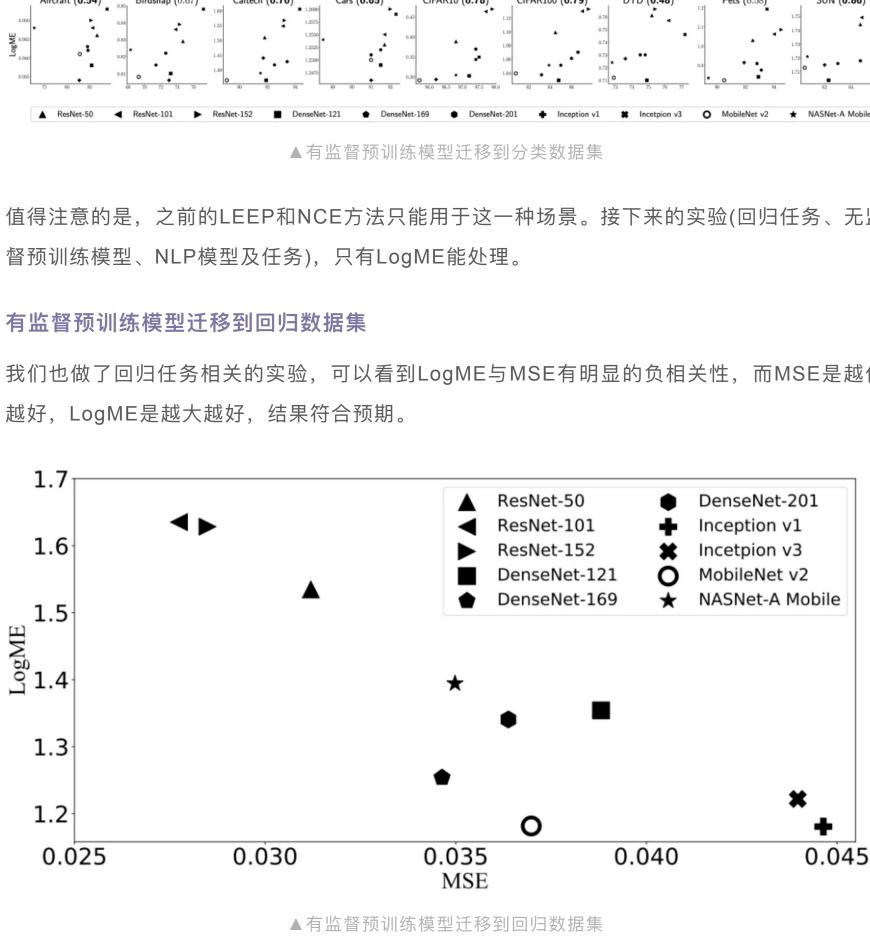


# **Algorithm 1** LogME 1: **Input:** Pre-trained model $\phi$ Target dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Optimized:  $m = \beta(V(\Lambda^{-1}(V^T(F^Ty))))$ 10: Update  $\alpha \leftarrow \frac{\gamma}{m^T m}, \beta \leftarrow \frac{n-\gamma}{||Fm-y||_2^2}$ 11: end while 12: Compute  $\mathcal{L}_k = \frac{1}{n}\mathcal{L}(\alpha,\beta)$  using Eq. 2 13:

整体的计算复杂度降低了一个阶,从四次方降低为三次方(见下表),使得该算法在数秒内就能 处理常见情况。 Complexity per for-loop Overall complexity  $\mathcal{O}(KD^3 + nKD^2)$  $\mathcal{O}(D^3 + nD^2)$   $\qquad \mathcal{O}(KD^3 + nKD^2)$   $\qquad \mathcal{O}(D^2 + nD)$   $\qquad \mathcal{O}(KD^2 + nKD + D^3 + nD^2)$ naïve optimized ▲优化前后复杂度对比, n是数据量, D是特征维度, K是类别数目

-0.6-0.8-1.0-1.5standard deviation of noise ▲特征质量越来越差时, LogME也越来越低。



喜欢此内容的人还喜欢

夕小瑶的卖萌屋

若被制裁,中国AI会雪崩吗?

上的表现也不错。

![LogME衡量NLP预训练模型]

无监督预训练模型

Pre-trained Network

MoCo V1

MoCo V2

**MoCo 800** 

0.23%0.23%5.22%LogME (optimized)  $50.4 \pm 0.7$ 0.31%▲各种方法耗时比较,LogME加速3000倍 值得注意的是,像LogME这种根据概率公式计算的方法,一般效果更好,但是耗时也更高。事

实上,如果我们采用简单粗暴的实现,评估一个模型就需要八百多秒。正是有了精心优化的版

🥠 展望 🥠

因为它的准确、快速、实用性,我们相信LogME除了能够用作预训练模型选择之外,还能够在

本,我们才能够**既有概率方法的优越效果,又有简单高效的实现**。



