

杨威峰^{1,2} 范意兴¹ 张儒清¹ 郭嘉丰^{1,2} 程学旗^{1,2}

(1. 中国科学院计算技术研究所 网络数据科学与技术重点实验室, 北京市 100190;
2. 中国科学院大学, 北京 100190)

论文摘要

基于社交媒体数据的性格预测是指通过用户在社交媒体平台上发布的文本等信息对用户的性格进行预测。现有的研究大多采用语言特征组合加浅层机器学习模型进行性格预测, 本文采用 BERT 和 PALs 等深层预训练模型对文本进行自动特征提取和文本向量表示来预测用户的性格并取得了较大的提高。同时针对性格预测中有标注数据少的问题, 本文尝试了跨领域学习的方法, 通过同时学习不同领域间的相似分布和独有分布提高性格预测的结果, 并使得 Facebook 和 YouTube 在性格预测上的平均 F1 值分别提高了 0.4% 和 1.0%。

论文简介

性格预测是一种对用户性格特征自动分类的任务, 该任务的输入来源可以是文本或其他多媒体数据, 真实标签一般是用户填写“大五人格”(Big5)测试问卷得到。现在越来越多的领域研究需要用到性格预测, 包括社交网络分析, 推荐系统, 欺骗检测, 作者身份归属, 情感分析/观点挖掘等。

目前性格预测的一个最大难点在于数据量太少。因为隐私问题很多用户并不愿公开自己的性格测试结果和社交媒体数据, 目前研究者所采用的数据一般分为两种: 1) 之前一些项目或比赛公开的小规模数据集。2) 招募志愿者填写性格测试问卷得到用户的性格真实标签并提取用户的社交媒体数据。为了更好地和之前的研究进行对比, 本文的数据集来源于前者。

本文首先探讨了文本的输入方式, 因为每个用户都有多个社交媒体文本, 本文将用户的所有文本构成的文档分成多个长度适中的子文档进行预测, 即保证了文本信息的完整又对数据的规模进行扩充。其次, 相较于之前的研究大多通过特征提取加机器学习的方法进行性格预测, 本文采用 BERT 等深层模型对用户的社交媒体文本进行更深层次的语义提取, 并使用最终得到的语义向量进行性格预测。另外, 本文中, 我们通过使用来自三种不同社交平台的基准数据集, 对跨领域社交媒体性格预测模型进行研究。之前的跨领域研究工作中, 研究者大多从扩充训练集的角度来增加训练规模, 从而提高模型的预测准确性, 而忽略了不同领域间存在不同的分布会影响模型的学习性能。本文在增加训练集的同时, 还从模型学习的角度联合学习不同领域间的相似分布和差异分布, 从结果来看取得不错的效果。因此本文将 PALs 模型引入跨领域学习, 通过联合学习

实验分析

从实验结果来看, BERT 在 twitter、facebook 和 youtube 等三个数据集上的实验效果均优于之前的评测或竞赛中的最优结果。由此可见, 深度模型能够对文本信息进行更丰富的语义建模, 这对性格预测的结果十分重要。接下来我们使用 PALs 模型来同时学习不同领域间的关于性格方面的相似分布和领域有关的特有分布, 从结果来看 PALs 模型相比于使用 BERT 模型进行跨领域学习相比有所提升。特别是在 {YouTube + Facebook} → YouTube 和 {YouTube + Facebook} → Facebook 的实验中 F1 相比于跨领域学习前分别提高了 1.0% 和 0.4%, 证明了 PALs 模型跨领域学习的优势性。

综上所述, 本文验证了 BERT 等深度预训练模型在性格预测任务上的优势效果, 相比于原有的最优结果有了很大的提升。同时本文也探讨了跨领域学习在性格预测任务上的可行性。通过 BERT 模型和 PALs 模型在跨领域学习中的对比实验可知: 模型中完全共享多领域的特征分布并不能保证有效的特征迁移来提高预测效果, 而通过在模型中添加适配器模块, 让模型在共享与任务有关的特征学习同时分开学习与领域有关的特有分布, 可以有效提高跨领域学习的预测效果。

论文结论

基于社交媒体的性格预测是用户画像中一个非常重要的任务。本文通过引入深度学习模型更好地建模文本中的深层语义关系, 证明了其在性格预测任务上的有效性。同时为了克服性格预测任务中有标注数据样本少的问题, 本文在不同社交媒体平台的数据集间进行了跨领域学习的尝试, 通过让模型同时学习不同数据集间的任务相关的相似分布和领域相关的特有分布, 进一步提高了性格预测的效果, 取得了目前先进的表现。

不同领域间的相似分布和特有分布提高了跨领域学习在性格预测上的效果, 并为跨领域学习提出了一些新的方法和建议。

系统模型

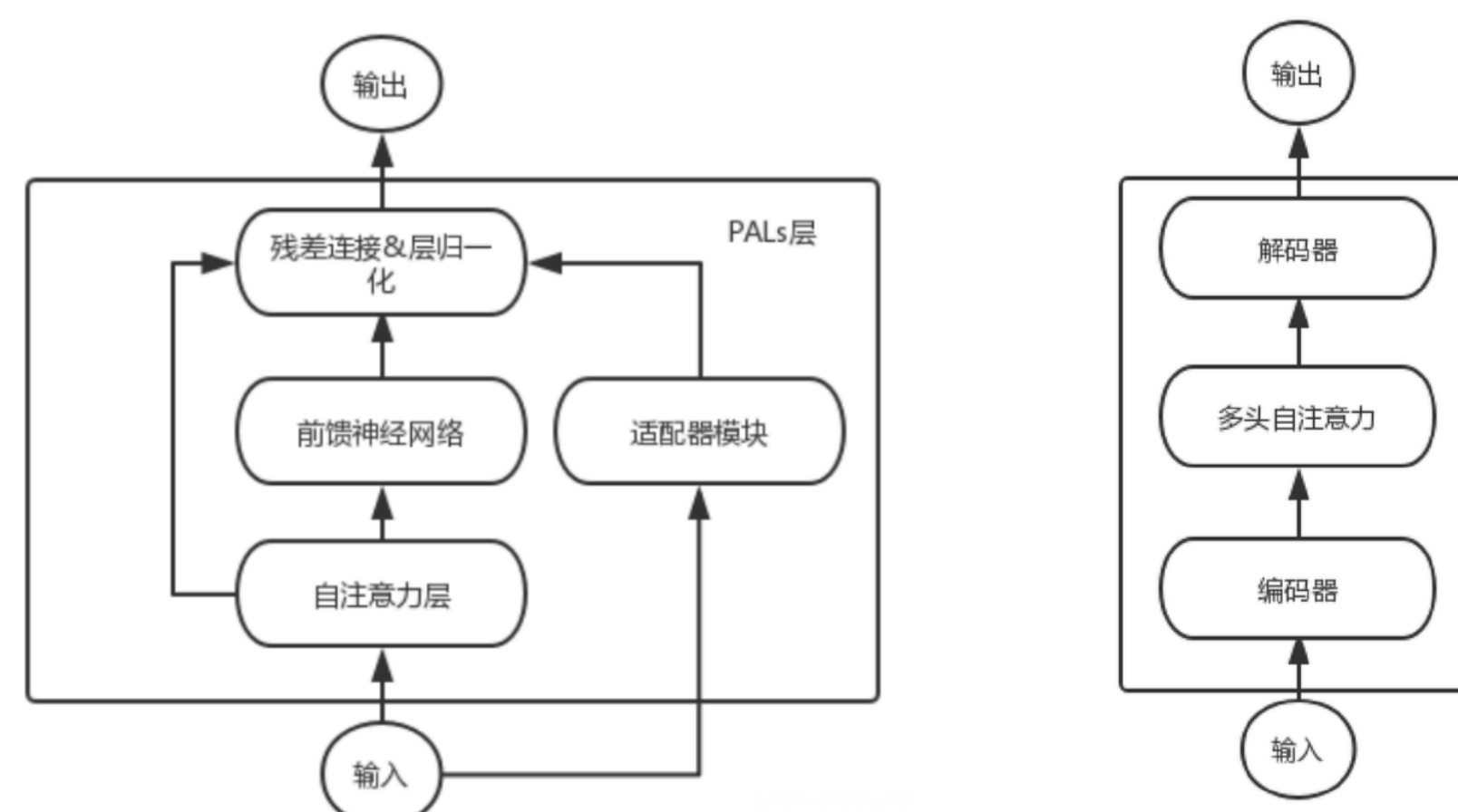


图1 PALs层的内部结构

图2 适配器的内部结构

PALs (Projected Attention Layers) 是 Stickland 等提出的用于解决多任务方法的模型。PALs 模型通过让多个任务共享 BERT 模型参数, 同时增加了少量的特定于任务的参数实现多任务学习。具体来讲, PALs 模型在 BERT 的每个层中添加了适配器模块, 每个下游任务都会有一个该任务特有的适配器来学习和特定任务相关的特征, 模型结构如图1, 2所示。

不同于原模型在多任务学习上的应用, 本文将 PALs 模型引入跨领域学习方面的研究。本文认为, 在进行跨领域学习中, 各领域间会存在和任务相关的相似分布, 同时由于不同领域数据集的噪声以及领域相关的特征, 这些领域又会存在领域特有的分布。因此, 我们希望 PALs 模型的共享参数部分学习各个数据集间关于性格分布的相似特征, 同时适配器模块中添加 n 个 Projected Attention Layers 用于学习每个领域特有的特征分布。

同时, 本文对原模型进行了一些修改。原模型中的对共享层不进行梯度传播和参数学习以提高参数学习的效率。本文中我们不再对共享层的参数进行“冻结”, 从而让模型充分学习和任务相关的特征分布, 实验证明修改后的操作能有效提高模型的预测效果。