

CENTRAL UNIVERSITY OF FINANCE AND ECONOMICS



中央财经大学

大数据分析统计基础课程

---

## 第七次课作业

---

吴宇翀

2021210793

EMAIL@WUYUCHONG.COM

指导老师：盖玉洁

2021 年 11 月 8 日

电影评论情感分析

## 目录

1	摘要	2
2	数据集介绍	2
3	深度学习	2
3.1	数据处理 . . . . .	2

## 1 摘要

我们使用 IMDB 数据集进行文本分类。在文本预处理阶段，我们尝试使用词编码和词向量的方式，在训练阶段，我们构建了 DNN、LSTM、BERT 等多个深度学习模型进行训练，并进行了模型比较，最高达到了 99% 的准确率。最后，为了进一步实现在超大文本集上进行训练，我们使用基于 Spark 的分布式算法在集群服务器上进行训练测试。

模型	计算配置	用时	准确率	可拓展性
tokenize + DNN	阿里云服务器 Xeon 8 核 CPU 32G 内存	10 分钟	60%	低-单机
Word2Vec + LSTM	阿里云服务器 Xeon 8 核 CPU 32G 内存	2 小时	80%	低-单机
bert - 小型	阿里云服务器 Xeon 8 核 CPU 32G 内存	1 小时	90%	低-单机
bert - AL	阿里云服务器 Xeon 8 核 CPU 32G 内存	1.5 小时	90%	低-单机
bert - 标准	阿里云服务器 Xeon 8 核 CPU 32G 内存	3 小时	92%	低-单机
spark	中央财经大学大数据高性能分布式集群	- 分钟	-%	高-集群

分布式模型在该小型数据集上没有优势，进行此项的意义在于对大型文本数据集可拓展性的技术储备，仅有在文本量级超过单机可承载上限时，分布式计算才具备意义

## 2 数据集介绍

我们选择了 IMDB 的电影评论文本数据进行大数据建模研究。

IMDB 是一个隶属于亚马逊公司旗下的世界著名互联网电影资料库 (Internet Movie Database)。它有着关于电影演员、电影、电视节目、电视明星和电影制作的在线数据，包括了影片的众多信息、演员、片长、内容介绍、分级、评论等，在电影评论评分时被广泛使用。IMDB 的论坛也十分活跃，除每个数据库条目都有留言板之外，还有关于多种多样的主题的各种综合讨论版。

我们将 IMDB 的电影评论文本用于自然语言处理的二元情感分类。我们使用 5 万条标有积极和消极标签的真实用户电影评论文本构建情感分类模型。即使用深度学习算法预测评论为正面或是负面。

我们使用的文本为多语言文本，其中英文文本数量占绝大多数比例。

## 3 深度学习

### 3.1 数据处理

1. 标签处理：分类标签由类别名称转为数字。
2. 数据集划分：在总体 5 万条文本中随机划分 20% 的测试集，再从训练集中划分 20% 的验证集。
  - 使用训练集的文本进行模型训练

- 使用验证集的文本进行模型超参数的调整
  - 使用测试集的文本进行模型效果评价
3. 使用自动的缓冲区大小，使用 32 的 `batch size`