

Report of Deep Learning for Natural Language Processing

冷元飞
ffxdds@163.com

Abstract

本实验旨在通过实际的中文金庸小说语料库数据，验证著名的 Zipf's Law，并进一步计算以词和字为基本单位的中文文本的平均信息熵，以揭示语言的内在统计规律及其信息复杂度。

Introduction

语言作为一种复杂而有序的社会现象，蕴含着丰富的统计规律和信息结构。理解和量化这些规律对于语言学研究、信息论应用及自然语言处理技术的发展具有重要意义。本实验聚焦于两个关键语言统计特性：Zipf's Law 与信息熵，并通过实际的中文语料库进行深入探究。

Zipf's Law 由美国语言学家 George Kingsley Zipf 于 1935 年提出，是描述自然语言中词汇分布非均匀性的一种经验定律。该定律指出，在一个大规模文本中，单词（或在汉字书写体系中的字符）的出现频率大致与其在频率表中的排名成反比的幂律关系，即第 r 位最常见的单词（或汉字）出现的频率约为第 1 位单词（或汉字）频率的 $\frac{1}{r}$ ，其中通常接近于 1。此存在揭示了语言使用的经济原则和人类认知的局限性，它在不同语言、不同文本类型中展现出惊人的普适性，是语言统计学的重要基石。在语言学中，信息熵被用来衡量语言表达的平均信息复杂度，即一个语言系统中词语或字符序列的平均不确定性。对于给定的词汇表或字符集，其信息熵越大，表明该语言系统在表达信息时的平均选择自由度越高，语言的复杂度和多样性越显著。针对特定语言（如中文）的验证依然具有学术价值，因为不同的语言系统可能存在特定的词汇使用模式和频率分布特征。对中文语料库的验证有助于深入理解汉语词汇使用的独特性，同时也能为语言模型构建、信息检索算法优化等应用提供更为精准的语言学基础。同时，分别计算词级和字级的信息熵有助于从不同粒度层面揭示汉语的信息组织特性。词级信息熵反映词汇组合的复杂程度和词汇表征信息的能力，而字级信息熵则直接体现了单个汉字作为语言基本单元的信息承载能力。对比两者可以洞察汉语在词汇层面与字符层面的信息分布差异，这对于理解汉字的表意功能、词法构造以及汉语信息处理算法的设计均有重要指导意义。

综上所述，本实验报告旨在通过严谨的数据分析和理论探讨，揭示 Zipf's Law 在中文语料库中的表现以及汉语在词、字两个层次上的信息熵特征，从而增进对汉语内在规律和信息复杂性的认识，为语言学研究、信息论应用及自然语言处理技术提供理论依据和实践指导。

Methodology

一、Zipf's Law 验证

Zipf's Law 是语言学和信息科学中的一项重要统计规律，它指出在一个自然语言文本中，单词（或汉字）的出现频率与其在频率表中的排名大致呈幂律关系，即第 r 位最常见的单词（或汉字）出现的频率约为第 1 位单词（或汉字）频率的 $\frac{1}{r}$ ，通常接近于 1。实验首先对选定的大型中文语料库进行了预处理，包括分词（以词为单位验证时）、去除停用词和标点符号等，以获得纯净的语言单元频数数据。然后，对处理后的数据进行排序并计算每个位置的单词（或汉字）频率，绘制频率与排名的双对数图，观察是否呈现出典型的直线趋势，即验证 Zipf's Law 的适用性。同时，通过拟合幂律模型，量化中文语料库中词汇分布的幂律特征。

二、信息熵计算

信息熵作为衡量信息不确定性和复杂度的关键指标，被应用于评估语言的内在结构和多样性。本实验分别计算了以信息熵作为衡量信息不确定性和复杂度的关键指标，被应用于评估语言的内在结构和多样性。实验分别计算了以词和字为单位的中文文本的信息熵

Experimental Studies

本实验旨在通过实际的中文语料库数据，一方面验证 Zipf's Law 在汉语环境中的适用性，通过绘制频率-排名双对数图和估计 Zipf 指数，直观展现并量化汉语词汇分布的幂律特征；另一方面，计算词级和字级的信息熵，以量化汉语文本的平均信息复杂度。实验采用标准的文本预处理步骤，包括分词、停用词过滤等，确保数据的准确性和代表性。所得结果将为深入理解汉语的统计规律和信息结构提供实证支持，同时也为相关领域的研究和应用提供有价值的参考数据、实验报告将总结 Zipf's Law 在中文语料库中的验证结果，包括直观的图形展示和 Zipf 指数估计值，分析其对中文语言特性的解释力。同时，报告将呈现词级和字级信息熵的具体数值，对比两者差异。

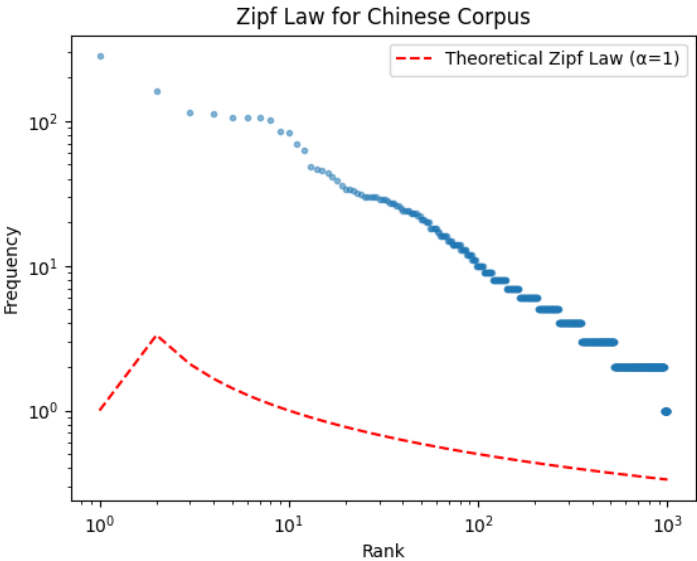


图 1. 频率与排名的对数图

通过计算 16 个文本库下中文平均信息熵如下所示：

语料库	词单位平均熵 (bits)	字单位平均熵 (bits)
三十三剑客图	0.0014	0.0042
书剑恩仇录	0.0003	0.0024
侠客行	0.0004	0.0027
倚天屠龙记	0.0006	0.0031
天龙八部	0.0002	0.0022
射雕英雄传	0.0010	0.0035
白马啸西风	0.0002	0.0023
碧血剑	0.0002	0.0024
神雕侠侣	0.0003	0.0025
笑傲江湖	0.0002	0.0022
越女剑	0.0005	0.0028
连城诀	0.0003	0.0024
雪山飞狐	0.0008	0.0034

飞狐外传	0.0003	0.0024
鸳鸯刀	0.0020	0.0047
鹿鼎记	0.0034	0.0062
总语料库	0.0034	0.0062

Conclusions

通过本次实验，我们期望不仅验证 Zipf's Law 在中文语境下的普适性，还能够通过信息熵的定量计算，深入理解中文语言系统的内在规律与复杂特性