

# Report of Deep Learning for Natural Language Processing

冯佳铂

775193638@qq.com

## Abstract

利用给定语料库金庸小说语料, 利用 1 ~ 2 种神经语言模型 (如: 基于 Word2Vec, LSTM, GloVe 等模型) 来训练词向量, 通过计算词向量之间的语意距离、某一类词语的聚类、某些段落直接的语意关联、或者其他方法来验证词向量的有效性。

## Introduction

使用 Word2Vec 模型训练词向量。word2vec 是 google 推出的一个 NLP 工具, 它可以将所有的词向量化, 可以定量的去度量词与词之间的关系, 挖掘词之间的联系。

最早期的词向量编码方式为 one hot representation, 这种编码简单, 但却维度过高, 同时词与词之间相互正交, 无法较好的衡量词与词之间的关系。通过神经语言模型训练, 将每个词映射到一个较短的词向量上, 以此构成新的向量空间, 在新的向量空间上, 便可以通过普通的统计学方法研究词与词之间的关系。

word2vec 之前, 一般使用 DNN 模型训练, 其中又包含 CBOW(Continuous Bag-of-Words) 与 Skip-Gram 两种模型。CBOW 模型的训练输入是某一个特征词的上下文相关的词对应的词向量, 输出是该特定的词的词向量; Skip-Gram 模型与 CBOW 的思路相反, 即输入是一个特定词的词向量, 输出的是特定词对应的上下文词向量。

当词汇达百万级别时 DNN 输出层需要进行 softmax 计算各个词的输出概率计算量过大。使用霍夫曼树来代替隐藏层和输出层的神经元, 霍夫曼树的叶子节点起到输出层神经元的作用, 叶子节点的个数为词汇表的大小, 内部节点起隐藏层神经元的作用。

在此基础上, 进一步使用 Negative Sampling (负采样) 方法, 进一步减少计算量。

gensim 封装了 google 的 C 语言版的 word2vec, 通过它可实现词向量训练。

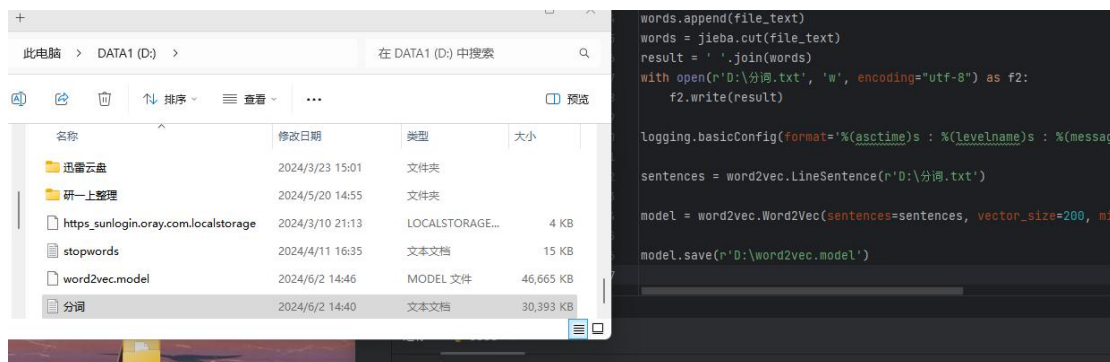
## Methodology

流程如下：

1. 语料库建立：将 16 部金庸小说合并建立语料库。
2. 文本预处理：将语料库中多余符号进行删除，如“本书来自 www.cr173.com 免费 txt 小说下载站”、“更多更新免费电子书请关注 www.cr173.com”等，一般的 NLP 处理中，需要去除停用词，但 word2vec 算法依赖于上下文，上下文可能就是停用词，可不去除停用词。使用 jieba 库函数进行分词。
3. 训练模型：使用 gensim 库函数对分词后的语料库进行训练，产生对应模型。
4. 分别计算词向量之间的语意距离、某一类词语的聚类、某些段落直接的语意关联，验证词向量的有效性。

## Experimental Studies

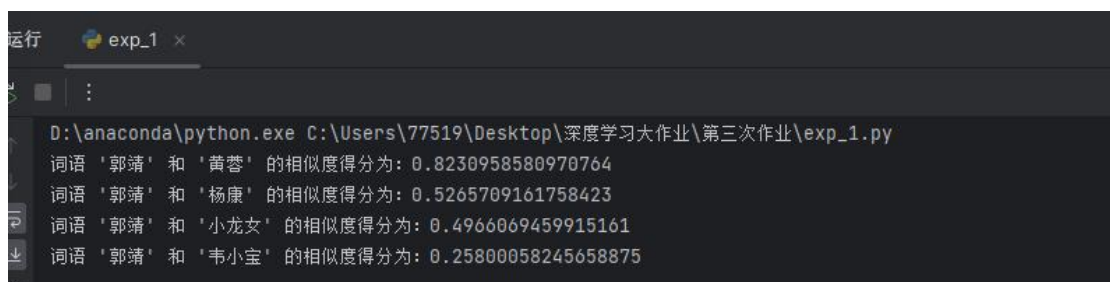
将训练所得模型“word2vec.model”保存至 D 盘。



### T1: 计算词向量之间的语意距离

所选词：

- ①“郭靖” “黄蓉”：词意距离：0.8231
- ②“郭靖” “杨康”：词意距离：0.5266
- ③“郭靖” “小龙女”：词意距离：0.4966
- ④“郭靖” “韦小宝”：词意距离：0.2580



所计算词意距离满足关系，郭靖与黄蓉为夫妻，关系最为密切，词意距离为 0.8231；郭靖与杨康虽为异姓兄弟，但道不同，词意距离 0.5266；小龙女与郭靖之间由于杨过存在一定联系，但关系应略低于郭靖与杨康，词意距离为 0.4966，满足关系。郭靖与韦小宝为 2 部不同小说主角，关联程度很低，词意距离为 0.258。

## T2：某一类词语的聚类

所选词：使用 K-means 聚类算法进行聚类分析

所得结果中大部分聚类比较成功，但同样存在效果较差部分。

较好的：

①与身体相关：

Cluster 0: 身上, 内力, 全身, 穴道, 气, 受伤, 鲜血, 穴, 解药, 解开, 伤口, 呼吸, 体内, 晕, 毒, 运气, 疼痛, 手足, 伤势, 剧毒, 毒药, 发作, 动弹不得, 但觉, 剧痛, 血, 练功, 吐, 散, 过得, 中毒, 吸, 难当, 服, 经脉, 真气, 九阳, 痛, 痛得, 丹田, 气息, 毒蛇, 毙命, 点中, 解, 内息, 掌心, 神智, 动弹, 毒性, 肌肤, 知觉, 自身, 呻吟, 察觉, 内伤, 酸软, 头脑, 毒针, 病, 喷出, 周身, 腹中, 流, 热气, 医治, 支持不住, 四肢, 运, 浑厚, 治, 药物, 流血, 克制, 痛楚, 毒物, 透, 无影无踪, 人身, 解救, 疗伤, 解毒, 晕倒, 鼻中, 寒气, 气绝, 两处, 生出, 瞳中, 支撑, 相抗, 符, 沾, 有毒, 精花, 传到, 热, 之毒, 运功, 七日, 丹药, 丹, 筋, 寒毒, 治伤, 不知不觉, 身体, 鼻息, 伤处, 昏迷, 毒气, 指力, 肋骨, 昏迷, 痛苦, 充沛, 关节, 丸, 冰冷, 便觉, 痊愈, 微弱, 创口, 无碍, 猛烈, 化为, 脉搏, 凉出, 推拿, 吸星, 跳动, 蛇, 筋骨, 大穴, 腋下, 着手, 全失, 伤痕, 胸膛, 包扎, 之法, 气血, 十成, 甚重, 致命, 救治, 送入, 传入, 渗出, 血痕, 清醒, 毒质, 剧烈, 解了, 冰凉, 灵药, 废人, 上升, 皮肉, 三处, 之伤, 灵台, 天竺僧, 舒畅, 抵受, 敷, 寒冷, 服食, 调匀, 闭住, 敷上, 蛤蟆功, 股, 冰蚕, 蜜蜂, 细微, 相触, 阴, 依照, 没伤, 依法, 断臂, 断筋, 身中, 麻痒, 走火入魔, 金创药, 逐步, 翻涌, 昏昏沉沉, 痒, 输入, 两穴, 份量, 华辉, 迷糊, 驱除, 丸药, 骨髓, 精力, 刺伤, 脑子, 五毒, 真力, 逐渐, 一倍, 伤药, 针刺, 呃, 郁光标, 放松, 断腿, 大碍, 毒发, 断绝, 枕, 发烧, 凝聚, 冰窖, 麻木, 冰, 汗珠, 清凉,

②与武功相关：

Cluster 10: 内功, 精微, 敬惜, 凶性, 外功, 断臂, 朝思暮想, 举重, 下功, 全到, 打心  
Cluster 10: 武功, 派, 弟子, 功夫, 剑法, 天下, 少林, 当年, 高手, 最, 高, 武林中, 内功, 颇, 人物, 武林, 掌门人, 传, 练, 华山派, 佩服, 功力, 指点, 第一, 生平, 神功, 门下, 五岳, 寻常, 家, 好手, 胜, 学, 武当派, 不同, 武艺, 道理, 本门, 武当, 武学, 聪明, 大有, 乃, 高僧, 一派, 了不起, 高强, 高明, 辟邪, 数十年, 传授, 本领, 一手, 剑术, 全真教, 真正, 绝技, 本派, 昆仑, 来历, 所说, 以及, 见识, 深厚, 泰山, 之一, 外号, 上乘, 全真, 门派, 名头, 奇, 九阴真经, 识, 西域, 一套, 法门, 诸般, 自称, 四大, 各派, 是从, 修习, 高人, 祖师, 不识, 这套, 胜过, 门人, 中土, 一门, 七, 共有, 天下第一, 这门, 过招, 修为, 较量, 所有, 平平, 学会, 打败, 九, 崆峒, 太师父, 用心, 所授, 点穴, 所学, 功, 秘笈, 真有, 惊人, 学武, 大法, 成名, 任何, 之士, 各有, 号称, 阴毒,

③与吃喝相关：

Cluster 7: 吃, 喝, 银子, 一口, 酒, 咬, 喝酒, 买, 钱, 水, 换, 挑, 碗, 一碗, 肉, 酒杯, 闻到, 吃饭, 骰子, 茶, 店小二, 一大, 干净, 药, 装, 干干净净, 饿, 饭, 饮, 端, 舒服, 饮酒, 葫芦, 一杯, 滋味, 骨, 金银, 熟, 拿到, 香, 之类, 斟, 筷子, 大碗, 药丸, 杯酒, 铜牌, 蜈蚣, 窝, 香气, 牛肉, 侍候, 金子, 菜, 汤, 不吃, 店伴, 拣, 去取, 几口, 煮, 送来, 厨房, 酒保, 缸, 酒壶, 这口, 老板, 肚里, 美酒, 点心, 洗, 饭菜, 菜肴, 珠宝, 酒菜, 清水, 有的是, 饮食, 鸡, 店伙, 泡, 酒来, 毒虫, 杯, 掌柜的, 明珠, 馒头, 入口, 玉瓶, 剥, 醉, 草, 供, 煮, 银两, 宣伊璜, 粽子, 酒饭, 剩, 人参, 老鼠, 鱼, 酒店, 翡翠, 馍, 羊, 蒙汗药, 喝, 一连, 财主, 食物, 粮食, 毛, 两口, 盛, 鲨鱼, 送饭, 一锭, 猪, 擀骰子, 伙计, 端起, 一包, 米, 找些, 叫化子, 酒里, 西瓜, 牛, 肚中, 种, 精致, 嚼, 一百, 几只, 茶碗, 蝎子, 好酒, 坛, 珍贵, 食, 烤, 嘴边, 金银珠宝, 一叠, 肚, 尿, 口里, 味道, 新鲜, 斟酒, 送上, 贵重, 举杯, 分给, 酒水, 拨开, 一古脑儿, 一大口, 元宝, 毒酒, 喝一杯, 好吃, 一盆, 柴, 三杯, 牛羊, 玉龙杯, 晚饭, 等物, 亲随, 熬, 冷水, 厨子, 吞, 盘中, 扑鼻, 饥饿, 王如意, 切, 酒碗, 酒香, 一桌, 一小, 一两, 整治, 仆役, 酒席, 犯人, 喝茶, 药材, 席上, 干粮, 豆腐, 药粉, 油, 甜, 几杯, 采, 吃饱, 玉杯, 下肚, 肚皮, 迷, 一箇, 灌, 漫, 喜酒, 粥, 馒头, 舌, 盐, 葵, 酒肉, 羹, 药箱, 一饮而尽, 吃些, 战战兢兢, 二千两, 热水, 炸, 有钱, 大里, 船家, 咀嚼, 店家, 膳八粥, 嘟嘟哈哈, 咬了一口, 三夜, 药酒, 煮饭, 吃喝, 羊肉, 几块, 尝尝, 清香, 药方, 啧啧, 老鸹, 津津有味, 做饭, 三口, 磨, 木盘, 钵, 鸡蛋, 五十两, 茶壶, 老仆, 放到, 馄饨, 碗碟, 水缸,

④与武功门派相关：

Cluster 10: 武功, 派, 弟子, 功夫, 剑法, 天下, 少林, 当年, 高手, 最, 高, 武林中, 内功, 颇, 人物, 武林, 掌门人, 传, 练, 华山派, 佩服, 功力, 指点, 第一, 生平, 神功, 门下, 五岳, 寻常, 家, 好手, 胜, 学, 武当派, 不同, 武艺, 道理, 本门, 武当, 武学, 聪明, 大有, 乃, 高僧, 一派, 了不起, 高强, 高明, 辟邪, 数十年, 传授, 本领, 一手, 剑术, 全真教, 真正, 绝技, 本派, 昆仑, 来历, 所说, 以及, 见识, 深厚, 泰山, 之一, 外号, 上乘, 全真, 门派, 名头, 奇, 九阴真经, 识, 西域, 一套, 法门, 诸般, 自称, 四大, 各派, 是从, 修习, 高人, 祖师, 不识, 这套, 胜过, 门人, 中土, 一门, 七, 共有, 天下第一, 这门, 过招, 修为, 较量, 所有, 平平, 学会, 打败, 九, 崆峒, 太师父, 用心, 所授, 点穴, 所学, 功, 秘笈, 真有, 惊人, 学武, 大法, 成名, 任何, 之士, 各有, 号称, 阴毒, 世间, 心法, 自来, 绝顶, 先师, 大不相同, 名, 真经, 岳先生, 神剑, 一阳指, 可比, 之术, 世, 钦佩, 甚高, 见过, 讲究, 地步, 十八, 太师, 学到, 练武, 相同, 兼, 平常, 真有, 当今, 王重阳, 段氏, 极少, 天山, 西藏, 两派, 本寺, 轻身, 学得, 之功, 人才, 相比, 口诀, 威震, 金刚, 八卦, 所知, 传给, 常人, 数百年, 以此, 苦练, 之道, 佛门, 记住, 最高, 当世, 名门, 单打独斗, 七十二, 罕见, 见之, 名, 远胜, 深湛, 六脉, 属, 造诣, 自负, 胜得, 天竺, 天龙, 钻研, 根本, 教中,

⑤人物相关：

Cluster 12: 韦小宝, 令狐冲, 张无忌, 杨过, 众人, 郭靖, 袁承志, 黄蓉, 陈家洛, 小龙女, 段誉, 石破天, 三人, 胡斐, 虚竹, 汉子, 公主, 丐帮, 欧阳锋, 萧峰, 周伯通, 洪七公, 点, 岳不群, 侍卫, 张翠山, 黄药师, 太后, 老畜, 盈盈, 群雄, 李莫愁, 谢逊, 郭, 三, 林平之, 岳灵珊, 青青, 王语嫣, 乔峰, 马, 四人, 慕容复, 周芷若, 丘处机, 师弟, 恒山, 僧, 狄云, 杨, 程灵素, 张召重, 赵敏, 阿, 法王, 段正淳, 木婉清, 二, 双儿, 田伯光, 婆婆, 徐天宏, 仪琳, 段, 那人, 霍青桐, 大汉, 周, 乾隆, 余鱼同, 文泰来, 柯镇恶, 陆无双, 高山, 南海, 鸠摩智, 笑, 殷素素, 苗人凤, 白衣, 女, 天地会, 女郎, 海老公, 张, 郭襄, 李沅芷, 老人, 向问天, 任我行, 游坦之, 红花, 陆菲青, 酒, 欧阳克, 喇嘛, 梅超风, 赵志敬, 赵半山, 契丹, 誉, 左冷禅, 陈近南, 阿朱, 雪山, 鳌拜, 张三丰, 道人, 星宿, 四, 裘千仞, 丁当, 白万剑, 僧人, 李文秀, 灭绝师太, 丁春秋, 袁紫衣, 师娘, 郑克爽, 总舵主, 石青, 周仲英, 包不同, 郡主, 洪, 风, 阿珂, 茅十八, 長老, 尹志平, 无忌, 碧冰, 香香公主, 水笙, 宫女, 俞莲舟, 余沧海, 岳夫人, 全轮法王, 碧神, 阿紫, 保定, 沐, 绮, 锤, 杨康, 丁不四, 指, 金花婆婆, 那姓, 谢烟客, 福康安, 穆念慈, 万震山, 马春花, 丁典, 童姥, 青城派, 袁千仞, 蓉, 完颜洪烈, 东方不败, 五人, 桃谷六仙, 峨嵋派, 何铁手, 道士, 完颜康, 乌老大, 段延庆, 殷梨亭, 戚芳, 无尘, 朱, 黄, 王夫人, 施琅, 俞岱岩, 老妇, 史婆婆, 齐, 万圭, 空, 吴应彪, 钟灵, 胡一刀, 胡青牛, 胖子, 冲虚, 一灯, 沐剑屏, 令狐大哥, 耶律齐, 血刀, 宋青书, 石, 彭连虎, 陆冠英, 萧湘子, 方怡, 琪, 霍都, 刘正风, 韦一笑, 周颠, 绿萼, 徐天川, 风天南, 王处一, 田归农, 武宣, 杨铁心, 定逸, 胖头陀, 花铁干, 头陀, 贝海石, 段公子, 闵柔, 铁掌, 黑白子, 黑衣, 拖雷, 林震南, 六怪, 何太冲, 帝, 金蛇, 李自成, 武三通, 洪夫人, 神雕, 梁子翁, 心砚, 尼摩星, 公孙止, 说不

较差的：

Cluster 8: 拿, 用, 地下, 一条, 之上, 一只, 取出, 一件, 一张, 一块, 放在, 一把, 柄, 怀中, 桌上, 一根, 露出, 匕首, 头上, 摸, 衣衫, 接过, 衣服, 经书, 一对, 交给, 掉, 头发, 插, 一看, 一枚, 衣袖, 随手, 一遍, 石凳, 两只, 藏, 捧, 怀里, 抬起, 烧, 下面, 一颗, 兵器, 两, 一头, 脚, 黄金, 物事, 一团, 钉, 绳索, 开来, 挂, 捏, 缚, 尸体, 蛇, 衣襟, 摸, 抹, 割, 树枝, 画, 一层, 铁, 包袱, 尸首, 之物, 上面, 底下, 脸颊, 一枝, 绣, 两根, 递, 长袍, 包裹, 挖, 靴, 宝剑, 包, 看时, 短, 缠, 牙齿, 鼻子, 盖, 银针, 托, 金针, 眼珠, 两枚, 墙上, 布袋, 提了, 细看, 肌肉, 大小, 取过, 放入, 壁上, 顶, 石头, 地图, 火焰, 呈, 高高, 手帕, 渔网, 握, 抚摸, 图形, 递给, 撕下, 皮, 铺, 两块, 半边, 一摸, 塞, 投入, 棺材, 长长的, 裹, 三枚, 割断, 死尸, 一粒, 颜色, 半, 罩, 石上, 粗, 面颊, 掌中, 铁盒, 首级, 一堆, 银票, 烂, 溅, 拉开, 凑, 敲, 深入, 跌落, 两边, 物, 赫然, 数, 小小的, 拿出, 底, 堆, 埋, 帽子, 套, 字迹, 裤子, 腰带, 黑色, 圆圈, 瓷瓶, 鞋子, 掏出, 两件, 利器, 白色, 一支, 锋利, 捡, 一幅, 贴, 宝物, 口子, 缚住, 绳子, 铜钱, 入怀, 甲, 一寸, 黑黝黝, 骸骨, 一具, 鼻, 厚, 件, 龙眼, 两行, 放着, 图画, 柔软, 人头, 碎, 袍子, 这块, 红色, 几条, 琵琶, 辫子, 火烧, 小指, 明晃晃, 铁板, 中心, 裂开, 满地, 摸, 箱子, 两张, 镯, 宝石, 棋子, 背负, 拾, 块, 断, 血迹, 手铐, 那条, 缝, 泥土, 玉, 剪刀, 扣, 珍珠, 三只, 白布, 几块, 利刃, 火石, 石块, 钥匙, 抄, 小孔, 追住, 珍室, 钢针, 铁, 几根, 僧袍, 银, 两把, 尖刀, 尾巴, 窟窿, 骷髅, 棋盘, 扇, 盒子, 袍, 一大块, 揉, 揩, 铸, 利剑, 碾, 一朵, 树皮, 根, 一顶, 剪, 灰尘, 石灰, 撕成, 几件, 模花, 人形, 系, 一具, 光芒, 膝下, 交在, 一条, 一根, 小帕, 道袍, 烧饼, 上刺, 眉毛, 结, 化成, 石板, 另一, 金丝, 写定, 贴身, 塞入, 暗入, 第

屠戮, 书来, 天聪, 实录, 例如, 圣教, 习惯, 中堂, 越王, 万物, 九月, 原在, 昆仑奴, 难以捉摸, 珍重, 加赋, 郝, 明报, 丰邵, 改用, 方式, 吾人, 志, 未尝, 狼, 一在, 历, 所奏, 形似, 一关, S, 八门, 时说, 主角, 初步, 古来, 释迦牟尼, 评论, 磨勒, 八句, 古诗, 四卷, 金玉, 嘘, 真言, 于外, 参透, 艰深, ●, 匈奴, 歌颂, 盲, 背熟, 揣摩, 启, G, 玉洞, 莲花, 翼, ◎, ◎, 拳诀, 呜呼, 鲜卑, 类, 二十二, 不符, 之曲, 汉文, 圆圆, 修改, 发表, 四张机, 双飞, 朱笔, 四行, 喝过, 及其, 欣赏, 候, 雁门, 不实, 谈得, 谓之, 零, 屑, 道德, 丧乱, 熟读, 昨夜, 神剑剑, 摩, 研读, 见解, 有言, 梁山泊, 探索, 昼夜, 中均, 真迹, 并称, 孔子, 第五个, 希, 三十九, @, 神剑经, 鄱阳, 崇焕, 情是, 九宫, 精品, 五斗, 无忧, 瞠目, 袁枚, 事宜, 无妄, 期, 兹, 观念, 孙子兵法, 最早, 甚少, 银骨簪, 其下, 改写, 下联, 碧莹莹, 三十七, 插图, 笔力, 之水, 太难, 平沙, 初五, 阿修罗, 交通, 笔墨, 三卷, 尽属, 戏剧, 落花, 中写, 第二道, 宋史, 善用, 又称, 语录, 调子, 岱, 甘, 后记, 中所, 连载, 文绉绉, 贴切, 诗经, 琴曲, 圆满, 言词, 少阴, 俄, 八荒, 卿, 书架上, 龙虎, 孟子, 明书辑, 乌拉, 书铺, ①, 宋藩, 一案, 昏迷, 表现, 算题, 一章, 易筋, 近代, 批, 夜雨, 杜甫, 菩菩, 数千里, 愁康, ", 默写, 闻, 4, 无所, 两首, 细字, 一册, 续, 由此而来, 重视, 玛米儿, 语句, 厚薄, 虚幻, 互通, 升天, 顿首, 种类, 辟, 御制, 讳, 一知半解, 不古, 当头棒喝, 周瑜, 上联, 卢生, 李白, 这出, 唐时, 黎干, 钱, 笔录, 武经, 资治通鉴

T3: 计算某些段落直接的语意关联

基于 word2vec 的句向量,选取最简单的相加求均值：

计算公式：

$$\text{sen\_vec} = \frac{\sum_{i=1}^m \text{vec}_i}{m}$$

sen\_vec 为句向量，m 为句中词数量，vec<sub>i</sub>为词的词向量。

①选取两端语义相近的段落进行计算：（选取倚天屠龙记中围攻光明顶中有关七伤拳的两段讨论）

宗维侠强道：“七伤拳是我崆峒绝技，怎能说有害无益？当年我掌门师祖木灵子以七伤拳威震天下，名扬四海，寿至九十一岁，怎么说会损害自身？你这不是胡说八道麼？”张无忌道：“木灵子前辈想必内功深湛，自然能练，不但无害，反而强壮肝腑。依晚辈之见，宗前辈的内功如不到那个境界，若要强练，只怕终归无用。”

宗维侠是崆峒名宿，虽知他所说的不无有理，但在各派高手之前，被这少年指摘本派的镇山绝技无用，如何不恼？大声喝道：“凭你也配说我崆峒绝技有用无用。既说无用，那就来试试。”张无忌淡淡一笑，说道：“七伤拳自是神妙精奥的绝技，拳力刚中有柔，柔中有刚，七般拳劲各不相同，吞吐闪烁，变幻百端，敌手委实难防难挡……”宗维侠听他赞誉七伤拳的神妙，说来语语中肯，不禁脸露微笑，不住点头，却听他继续说道：“……晚辈只是说内功修为倘若不到，那便练之有害无益。”

计算结果：段落间语义相似度为：0.9089

Prefix dict has been built successfully.  
段落间的语义相似度：0.908865749835968



②选取两篇小说中不相关的两端进行计算：

原来黄药师对妻子情深意重，兼之爱妻为他而死，当时一意便要以死相殉。他自知武功深湛，上吊服毒，一时都不得便死，死了之后，尸身又不免受岛上哑仆糟蹋，于是去大陆捕拿造船巧匠，打造了这艘花船。这船的龙骨和寻常船只无异，但船底木材却并非用铁钉钉结，而是以生胶绳索胶缠在一起，泊在港中之时固是一艘极为华丽的花船，但如驶入大海，给浪涛一打，必致沉没。他本拟将妻子遗体放入船中，驾船出海，当波涌舟碎之际，按玉箫吹起《碧海潮生曲》，与妻子一齐葬身万丈洪涛之中，如此潇洒倜傥以终此一生，方不辱没了当世武学大宗匠的身分，但每次临到出海，总是既不忍携女同行，又不忍将她抛下不顾，终于造了墓室，先将妻子的棺木厝下。这艘船却是每年油漆，历时常新。要待女儿长大，有了妥善归宿，再行此事。

宗维侠强道：“七伤拳是我崆峒绝技，怎能说有害无益？当年我掌门师祖木灵子以七伤拳威震天下，名扬四海，寿至九十一岁，怎么说会损害自身？你这不是胡说八道麽？”张无忌道：“木灵子前辈想必内功深湛，自然能练，不但无害，反而强壮肝腑。依晚辈之见，宗前辈的内功如不到那个境界，若要强练，只怕终归无用。”

计算结果：段落间语义相似度为：0.7259

段落间的语义相似度：0.7258730530738831

有所下降，但感觉下降的不够多，排除是因为同一个人写的小说缘故

选用一篇其他小说段落

这时，李达康的专车也驶临了收费站。侯亮平下车，当道一站，举起手掌做出停车的手势。李达康的轿车缓缓停下。后面跟踪而来的张华华的警车也在李达康专车的左侧停了下来。李达康的司机下了车，走到侯亮平和陆亦可面前：你们想干啥？知道这是谁的车吗？侯亮平说：我只知道被传唤人欧阳菁在这台车上。司机满脸的不屑：我说同志，欧阳菁是谁的夫人，你不会不知道吧？侯亮平说：欧阳菁是谁的夫人与我们检察院办案没关系！陆亦可解释道：有人举报了欧阳菁副行长，我们要请她去谈一谈！司机不无傲慢地说：知道吗？欧阳副行长是市委李达康书记的夫人，这台车也是李达康书记的专车！

计算结果为：0.7963

段落间的语义相似度：0.7963464260101318

可见计算段落之间语义关联，具备一定参考价值，具体关系值得商榷。

## Conclusions

通过 word2vec 模型训练了一个具备一定意义的词向量模型，可在一定程度上验证词向量的有效性。

## References

- [1] 用 gensim 学习 word2vec - 刘建平 Pinard - 博客园 (cnblogs.com)
- [2] 基于 word2vec 的中文词向量训练\_word2vec 能用汉字吗-CSDN 博客
- [3] 基于 Word2vec 词聚类的关键词实现\_word2vec 词聚类的关键词提取算法-CSDN 博客
- [4] 基于 Word2vec 文本聚类\_词向量模型 word2vec 与 聚类-CSDN 博客
- [5] NLP 词向量和句向量方法总结及实现-CSDN 博客

Part 2 : 计算中文平均信息熵

## Abstract

对金庸先生 16 篇小说建立中文语料库，并计算其中文平均信息熵。

## Introduction

信息是个很抽象的概念。人们常常说信息很多，或者信息较少，但却很难说清楚信息到底有多少。直至香农提出了“信息熵”的概念，才解决了对信息的量化度量问题。分别利用基于词的一元模型、二元模型和三元模型计算所给中文语料库的信息熵。

## Methodology

### M1: 一元模型

$H(x) = - \sum_{x \in X} P(x) \log P(x)$ ，其中 $P(x)$ 近似等于每个词在语料库中出现的频率。

### M2: 二元模型

$H(X|Y) = - \sum_{x \in X, y \in Y} P(x, y) \log P(x|y)$ ，其中 $P(x, y)$ 近似等于每个二元词在语料库中出现的频率，条件概率 $P(x|y)$ 近似等于每个二元词组在语料库中出现的频数与该二元词组的第一个词为词首的二元词组的频数的比值。

### M3: 三元模型

$H(X|Y, Z) = - \sum_{x \in X, y \in Y, z \in Z} P(x, y, z) \log P(x|y, z)$ ，其中 $P(x, y, z)$ 近似等于每个三元词组在语料库中出现的频率，条件概率 $P(x|y, z)$ 近似等于每个三元词组在语料库中出现的频数与以该三元词组的前两个词为词首的三元词组的频数的比值。

## Experimental Studies

## T1: 验证齐夫定律

模型	词信息熵
一元模型	8.8132
二元模型	5.4837
三元模型	3.0216

```
一元信息熵为: 8.81318171981199
二元信息熵为: 5.483662467190433
三元信息熵为: 3.021636809115951
进程已结束, 退出代码为 0
```

## Conclusions

随着模型复杂度的增加, 信息熵下降, 信息熵的值越大, 文本中不同字符或词组的出现频率相差越大, 文本的特征值越分散; 模型变复杂后, 考虑词与词之间的搭配, 上下文之间的联系, 能够更准确的扑捉到某些语言现象, 从而使得这些现象的出现变得更加可预测, 所以信息熵降低。

## References

- [1] 刘鸿博.《Python 自然语言处理》.[Z].undefined.电子工业出版社.2022.
- [2][中文信息熵的计算 汉语信息熵-CSDN 博客](#)
- [3]Brown, P. F., Pietra, S. D. A., Pietra, V. D. J., Lai, J. C., & Mercer, R. L. . (1992). An estimate of an upper bound for the. Computational Lingus, 18(1), 31-40.