



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

关键词提取

汇报人：宋策、冀温瑾、何妙、叶姿逸、邹媛婷、冷晓晗

时间：2021/11/01



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

目录

CONTENTS

1

关键词提取介绍及应用

2

关键词提取基本原理

3

经典算法介绍

4

开源工具与实例展示

5

关键词提取前沿研究介绍



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

1

关键词提取介绍及应用

汇报人：冀温瑾

关键词提取

关键词：是指能包含一个词素(语言中最小的有意义的单位)的词或语言里最小的可以自由运用的单位，切能够是表达文本主题内容的词，包括单词，术语和短语，在含义上是独立非复合的。包含一定的信息量，对文本内容的理解有作用。

联合
建设 现代 建筑 电子
财经 管理 理工
贸易 首都
工业 科技 工商
职业 工程
信息 经济 技术
国际 经贸 金融 师范
商贸 商务



文献检索

由论文作者给出论文的关键词，用户可通过一个或多个关键词匹配查找到相关文献，简化搜索结果。

自动摘要

通过查找关键词最多的句子可以自动形成摘要。

推荐系统

通过用户历史浏览记录等用户标签，基于关键词匹配为用户推荐相关的产品信息。

文本分类

文本分类的核心问题是从文本中提取出关键词，然后基于一定的规则对文本分类。

搜索引擎

当我们搜索时，算法会从输入的语句中提取关键词并对文本内容进行相关性匹配。

机器翻译

机器翻译离不开关键词的抽取。



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

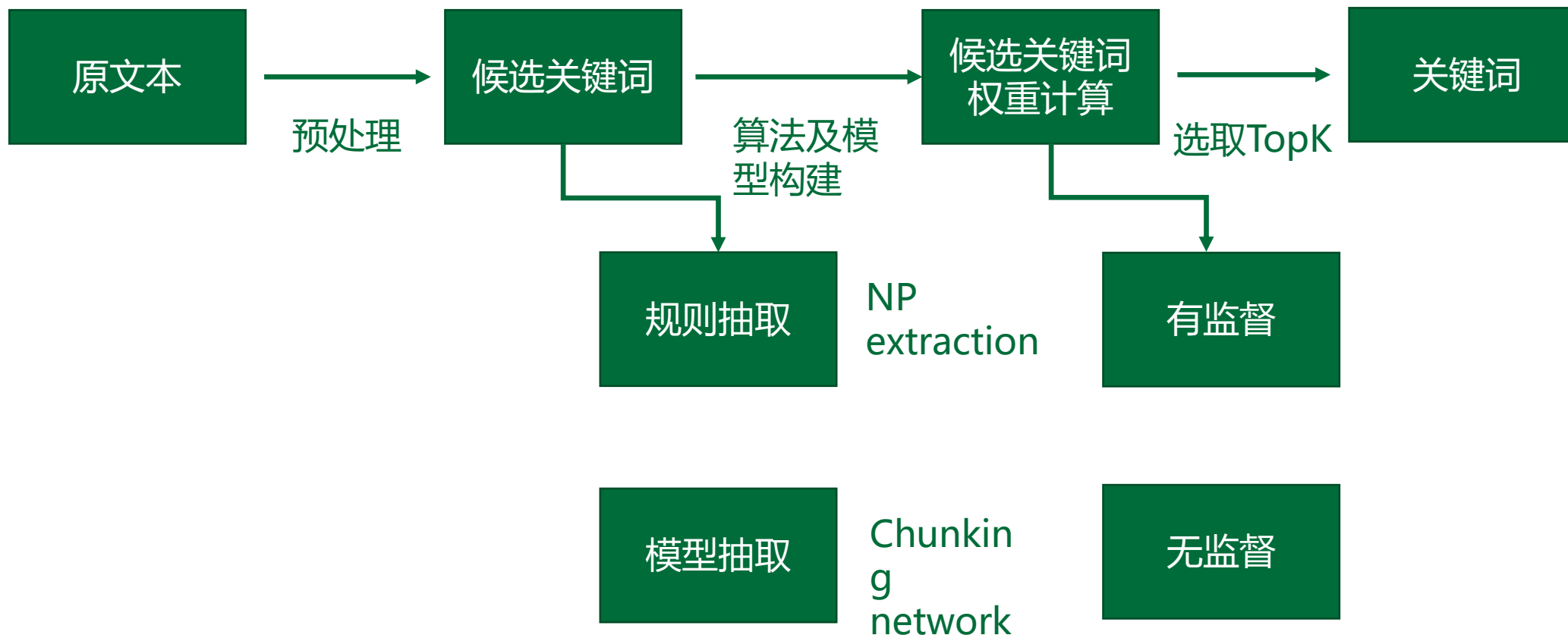
2

关键词提取基本原理

汇报人：冀温瑾

关键词抽取存在最基本的问题，即什么是“关键词”，而且如何提取“关键词”。





关于文本的关键词提取方法分为**有监督**和**无监督**两种：

1

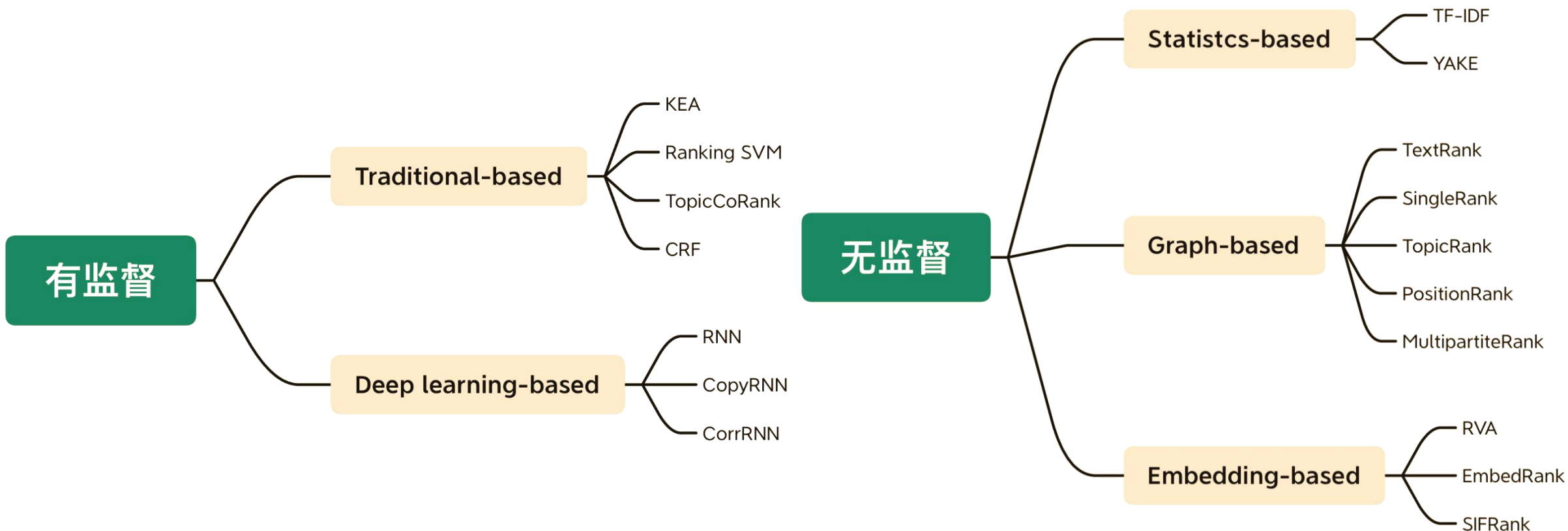
有监督

它是建关键词抽取算法看作是二分类问题，判断文档中的词或者短语是或者不是关键词。既然是分类问题，就需要提供已经标注好的训练语料，利用训练语料训练关键词提取模型，根据模型对需要抽取关键词的文档进行关键词抽取。

2

无监督

不需要人工标注的语料，利用某些方法发现文本中比较重要的词作为关键词，进行关键词抽取。





北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

3

经典算法介绍

汇报人：叶姿逸、何妙



目的：用以评估一字词对于一个文件集或一个语料库中的其中一份文件的重要程度。字词的重要性随着它在文件中出现的**次数**成正比增加，但同时会随着它在语料库中出现的**频率**成反比下降。

·TF(Term Frequency)词频：词条（关键字）在文本中出现的频率

$$TF = \frac{\text{在某一类中词条}w\text{出现的次数}}{\text{该类中所有的词条数目}}$$

·IDF(Inverse Document Frequency)逆向文件频率：由总文件数目除以包含该词语的文件的数目，再将得到的商取对数得到。

·如果包含词条 w 的文档越少，IDF越大，则说明词条具有很好的类别区分能力。

$$IDF = \log \frac{\text{语料库的文档总数}}{(\text{包含词条}w\text{的文档数} + 1)} \quad (\text{分母之所以要} + 1 \text{ 是避免分母为} 0)$$



·某一特定文件内的高词语频率，以及该词语在整个文件集中的低文件频率，可以产生出高权重的TF-IDF。因此，TF-IDF倾向于过滤掉常见的词语，保留重要的词语。

$$TF - IDF = TF * IDF$$

某个词对文章的重要性越高，它的TF-IDF值就越大。所以，排在最前面的几个词，就是这篇文章的关键词。

TF-IDF应用

- (1) 搜索引擎； (2) 关键词提取； (3) 文本相似性； (4) 文本摘要

TextRank算法：正规的TextRank公式在PageRank的公式的基础上，引入了边的权值的概念，代表两个句子的相似度。

TextRank 一般模型可以表示为一个有向有权图 $G = (V, E)$ ，由点集合 V 和边集合 E 组成， E 是 $V \times V$ 的子集。图中任两点 V_i, V_j 之间边的权重为 w_{ji} ，对于一个给定的点 V_i ， $In(V_i)$ 为指向该点的点集合， $Out(V_i)$ 为点 V_i 指向的点集合。点 V_i 的得分定义如下：

$$WS(V_i) = (1 - d) + d * \sum_{V_j \in In(V_i)} \frac{w_{ji}}{\sum_{V_k \in Out(V_j)} w_{jk}} WS(V_j)$$

其中， d 为阻尼系数，取值范围为 0 到 1，代表从图中某一特定点指向其他任意点的概率，一般取值为 0.85。



样例:

程序员(英文Programmer)是从事程序开发、维护的专业人员。一般将程序员分为程序设计人员和程序编码人员,但两者的界限并不非常清楚,特别是在中国。软件从业人员分为初级程序员、高级程序员、系统分析员和项目经理四大类。

首先对这句话分词,得出分词结果

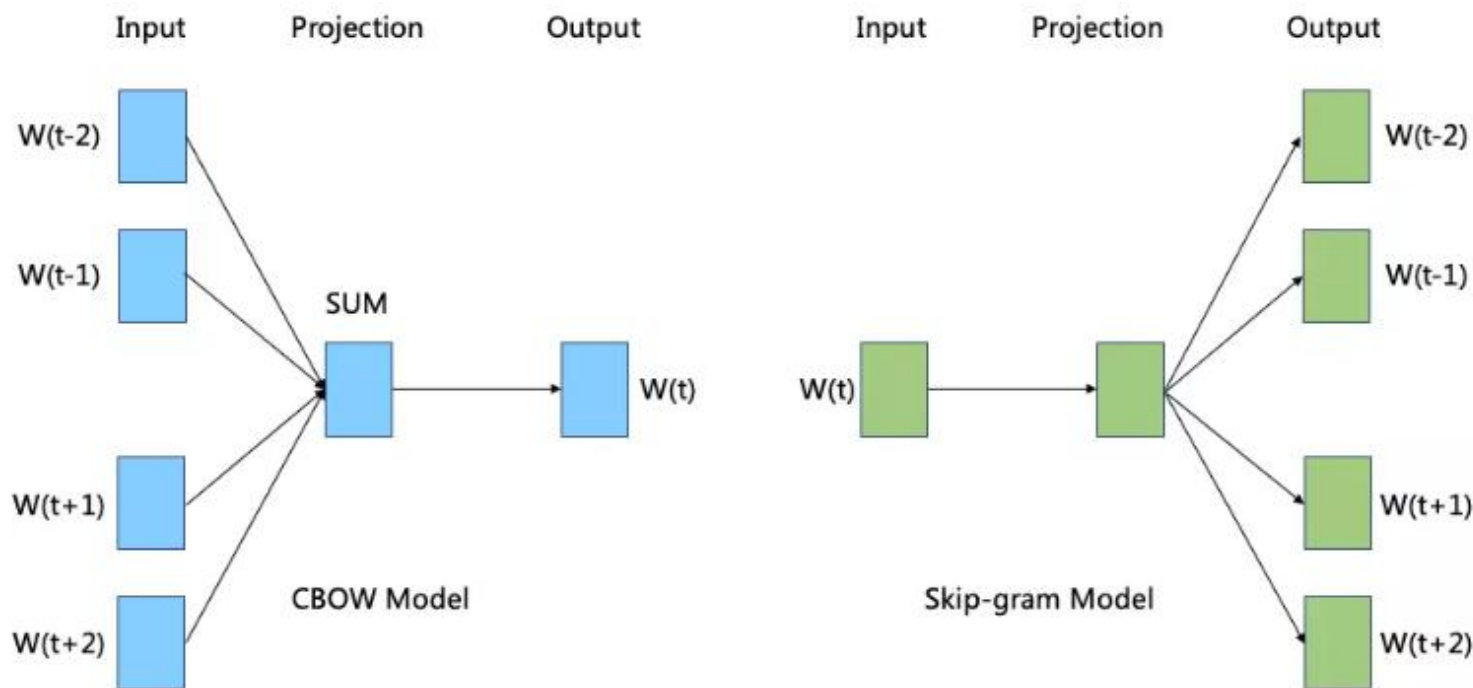
[程序员/n, (, 英文/nz, programmer/en,), 是/v, 从事/v, 程序/n, 开发/v, /w, 维护/v, 的/uj, 专业/n, 人员/n, 。 /w, 一般/a, 将/d, 程序员/n, 分为/v, 程序/n, 设计/vn, 人员/n, 和/c, 程序/n, 编码/n, 人员/n, , /w, 但/c, 两者/r, 的/uj, 界限/n, 并/c, 不/d, 非常/d, 清楚/a, , /w, 特别/d, 是/v, 在/p, 中国/ns, 。 /w, 软件/n, 从业/b, 人员/n, 分为/v, 初级/b, 程序员/n, 、 /w, 高级/a, 程序员/n, 、 /w, 系统/n, 分析员/n, 和/c, 项目/n, 经理/n, 四/m, 大/a, 类/q, 。 /w]

然后去掉里面的停用词,比如标点符号、常用词、以及“名词、动词、形容词、副词之外的词”。得出实际有用的词语:

[程序员, 英文, 程序, 开发, 维护, 专业, 人员, 程序员, 分为, 程序, 设计, 人员, 程序, 编码, 人员, 界限, 特别, 中国, 软件, 人员, 分为, 程序员, 高级, 程序员, 系统, 分析员, 项目, 经理]

然后我们首先计算距离每个词距离不大于窗口大小的词的集合,然后根据公式进行迭代,直到所有词的重要性收敛到某一个值的时候,就可以停止迭代并输出结果。

word2vec工具主要包含两个模型：连续词袋模型（CBOW, continuous bag of words）和跳字模型（skip-gram）



CBOW是根据上下文去预测目标词来训练得到词向量

https://blog.csdn.net/qq_30189255
Skip-gram是根据目标词去预测周围词来训练得到词向量



Latent Dirichlet Allocation, 简称LDA, 由David M. Blei、Andrew Y. Ng和Michael I. Jordan在2003年提出。

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI



The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

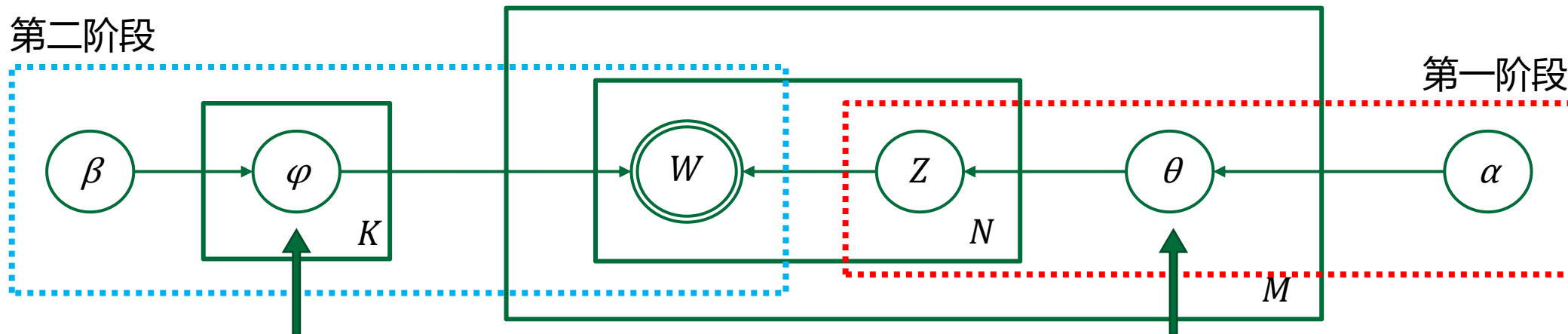


LDA



核心思想：①文档是若干主题的混合分布 ②每个主题又是词的概率分布

$$P(\text{词}|\text{文档}) = \sum P(\text{词}|\text{主题}) \times P(\text{主题}|\text{文档})$$

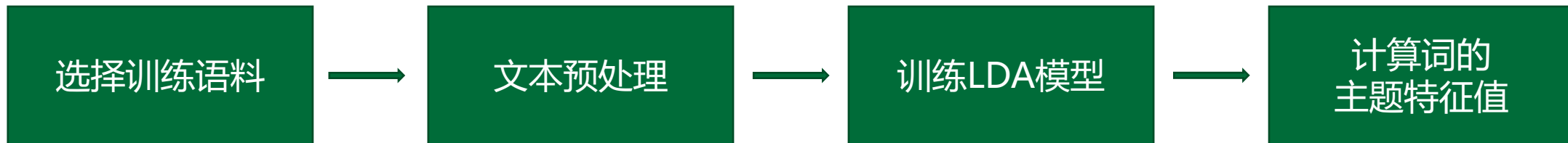


文档集总词汇数为 V ，一篇由 n 个词组成的文档，每个词的生成都服从多项分布。

不同文档之间的骰子不是同一个，每次为文档选一个主题骰子，这个过程也服从多项分布。

- **阶段一**：选取一个参数为 θ_m 的文档—主题分布，然后对第 m 篇文档的第 n 个词的主题，生成 $Z_{m,n}$ 的编号。从一个参数为 α 的Dirichlet分布中采样出一个多项分布 θ_m ，作为该文档在 k 个主题上的分布。
- **阶段二**：生成第 m 篇文档的第 n 个词。对该文档中的每个词，根据上步中的 θ_m 分布，采样出一个主题编号，然后根据主题—词分布对应参数为 β 的Dirichlet分布中采样出一个多项分布，作为该主题下词的分布。

算法流程:



- **方法一**: 假定每个词只能代表一个主题, 取模型中各主题下权重高的 TOP K 个词作为该主题的词。
- **方法二**: 假定主题区分度大的词应该是那些在某个主题下的权重高、而在其他主题中出现频率少的词语, 每个词都只是代表其最能代表的那个主题。计算出词的主题特征值之后, 即可以作为关键词, 也可以将词的主题特征值作为其他自动关键词抽取方法中词的一个特征使用。

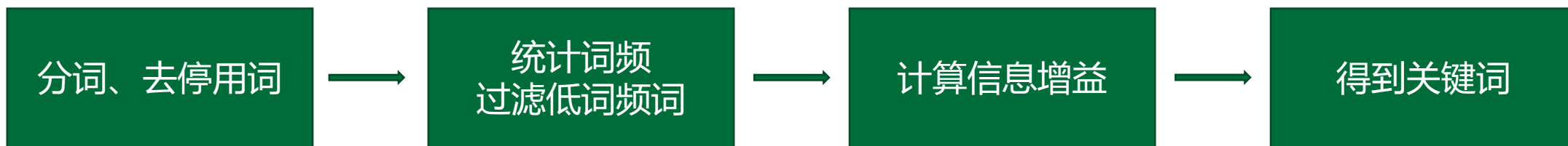
信息增益：是一种基于熵的度量方式。在文本特征提取中，某个特征的信息增益是指有该特征和没有该特征时，为整个分类所能提供的信息量的差别。

$$\begin{aligned}
 \text{Gain}(t_i) &= \overset{\text{信息熵}}{\text{Entropy}(S)} - \overset{\text{条件熵}}{\text{ExpectedEntropy}(S)} \\
 &= \left\{ -\sum_{j=1}^M P(C_j) \log P(C_j) \right\} \\
 &\quad - P(t_i) \left\{ -\sum_{j=1}^M P(C_j | t_i) \log P(C_j | t_i) \right\} \\
 &\quad - P(\bar{t}_i) \left\{ -\sum_{j=1}^M P(C_j | \bar{t}_i) \log P(C_j | \bar{t}_i) \right\}
 \end{aligned}$$



核心思想：把关键词的提取过程看成是文本分类中的特征提取的过程。在已知分类的情况下，通过将专业语料库和通用语料库进行对比，计算词条对类别信息量的贡献值，贡献值越大，其成为关键词的可能性就越大。

算法流程：





卡方检验：卡方是数理统计中用于检验两个变量独立性的方法，是一种确定两个分类变量之间是否存在相关性的统计方法，经典的卡方检验是检验定性自变量对定性因变量的相关性。

$$\chi^2(w) = \sum_{g \in G} \frac{(\text{freq}(w, g) - n_w p_g)^2}{n_w p_g}$$

词w与高频词集G中的词g共现的分布偏差

表示 w 与 g 的共现频数

高频词g在文档中出现的期望概率

共现词的期望频率

词w出现的次数

阈值 χ^2_α

核心思想：呈特殊分布特征的词语可能是关键词，根据该规则，通过卡方检验计算词语与高频词共现的分布偏差，当偏差大于阈值时，认定该词语为关键词。

算法流程：





北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

4

开源工具与实例展示

汇报人：邹媛婷、冷晓晗、宋策



HanLP

<https://github.com/hankcs/HanLP>

- 面向生产环境的多语种自然语言处理工具包，基于PyTorch和TensorFlow 2.x双引擎。
- 轻量级RESTful API，仅数KB，适合敏捷开发、移动APP等场景。服务器算力有限，匿名用户配额较少，建议申请公益API密钥auth
- 海量级native API，依赖PyTorch、TensorFlow等深度学习技术，适合专业NLP工程师、研究者以及本地海量数据场景。要求Python 3.6以上。可以在CPU上运行，推荐GPU/TPU。

主要功能:

1. 分词 (粗分、细分2个标准, 强制、合并、校正3种词典模式)
2. 词性标注 (PKU、863、CTB、UD四套词性规范)
3. 命名实体识别 (PKU、MSRA、OntoNotes三套规范)
4. 依存句法分析 (SD、UD规范)
5. 成分句法分析
6. 语义依存分析 (SemEval16、DM、PAS、PSD四套规范)
7. 语义角色标注
8. 词干提取
9. 词法语法特征提取
10. 抽象意义表示 (AMR)
11. 指代消解
12. 语义文本相似度
13. 文本风格转换



HanLP ——关键词提取

- 内部采用**TextRank算法**实现，用户可以直接调用：
`TextRankKeyword.getKeywordList(document, size)`

```
String content = "程序员(英文Programmer)是从事程序开发、维护的专业人员。一般将程序员分为程序设计人员和程序编码人员，但两者的界限并不非常清楚，特别是在中国。软件从业人员分为初级程序员、高级程序员、系统分析员和项目经理四大类。";  
List<String> keywordList = HanLP.extractKeyword(content, 5);  
System.out.println(keywordList);
```

```
[程序员, 程序, 分为, 人员, 软件]
```

步骤:

1. 分词

把文本通过一个的分词算法进行分词，这里采用的是HMM算法。

2. 构造窗口

为分个词构造窗口，这个词前后各四个词就是这个词的窗口

3. 迭代投票

每个词最后的投票得分由这个词的窗口进行多次迭代投票决定，迭代的结束条件就是大于最大迭代次数这里是 200次，或者两轮之前某个词的权重小于某一值，这里是 0.001f



小明NLP

<https://github.com/seanlee97/xmnlp/>

- 轻量级中文自然语言处理工具
- 当前版本：2019年9月发布
- **主要功能：**
 - 中文分词 & 词性标注
 - 支持繁體
 - 支持自定义词典
 - 中文拼写检查
 - 文本摘要 & **关键词提取(TextRank)**
 - 情感分析
 - 文本转拼音
 - 获取汉字偏旁部首

从文本中提取关键词:

```
xmnlp.keyword(text: str, k: int = 10, stopwords: bool = True, allowPOS: Optional[List[str]] = None) -> List[Tuple[str, float]]
```

参数:

- text: 文本输入
- k: 返回关键词的个数
- stopwords: 是否去除停用词
- allowPOS: 配置允许的词性

结果返回:

由关键词和权重组成的列表

并行处理关键词提取:

```
xmnlp.keyphrase_parallel(texts: List[str], k: int = 10, stopwords: bool = False, n_jobs: int = 2) -> Generator[List[str], None, None]
```

参数:

- texts: 文本列表
- n_jobs: 线程 worker 数



SnowNLP

<https://github.com/isnowfy/snownlp>

- 一个处理中文文本的 Python 类库

主要功能

1. 中文分词 (Character-Based Generative Model)
2. 词性标注 (TnT 3-gram)
3. 情感分析 (现在训练数据主要是买卖东西时的评价, 所以对其他的一些可能效果不是很好, 待解决)
4. 文本分类 (Naive Bayes)
5. 转换成拼音 (Trie树实现的最大匹配)
6. 繁体转简体 (Trie树实现的最大匹配)
7. 提取文本关键词 (TextRank算法)
8. 提取文本摘要 (TextRank算法)
9. tf-idf
10. Tokenization (分割成句子)
11. 文本相似 (BM25)

关键词提取:

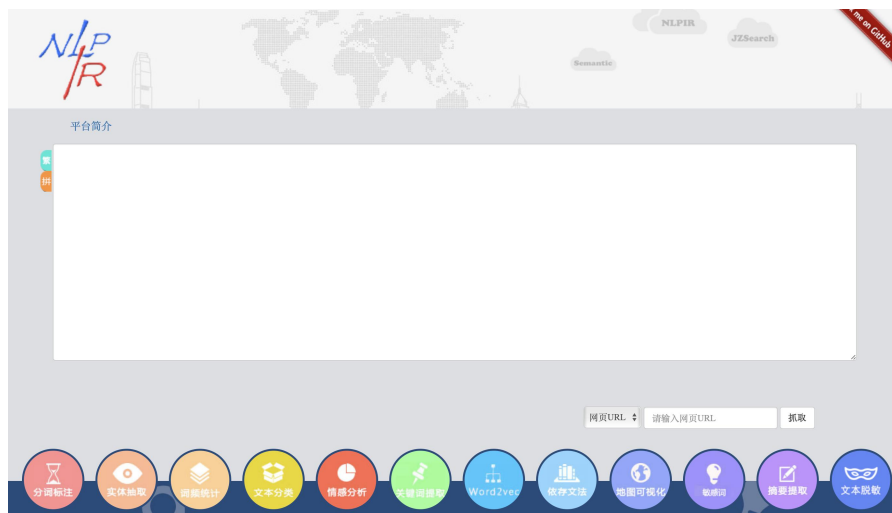
```
keywords(self, limit=5,  
merge=False)
```

1. 输入文段
2. 按换行符和中文符号 (, . ? ! ;) 划分成句子: list[句子]
3. 对每个句子分词: list[[词]]
4. 统计该词的关联词, 去重并排序 (关联词, 该词所在句的其他词都为关联词)
5. 计算每个词的关键度 (与该词的关联词的数量相关): $m[i] += 0.85 / \text{len}(\text{list}[\text{关联词}]) * m[i](\text{old})$ 。默认计算200次, 期间前后两次 $m[i]$ 相减绝对值小于等于0.001退出计算
6. 排序, 截取

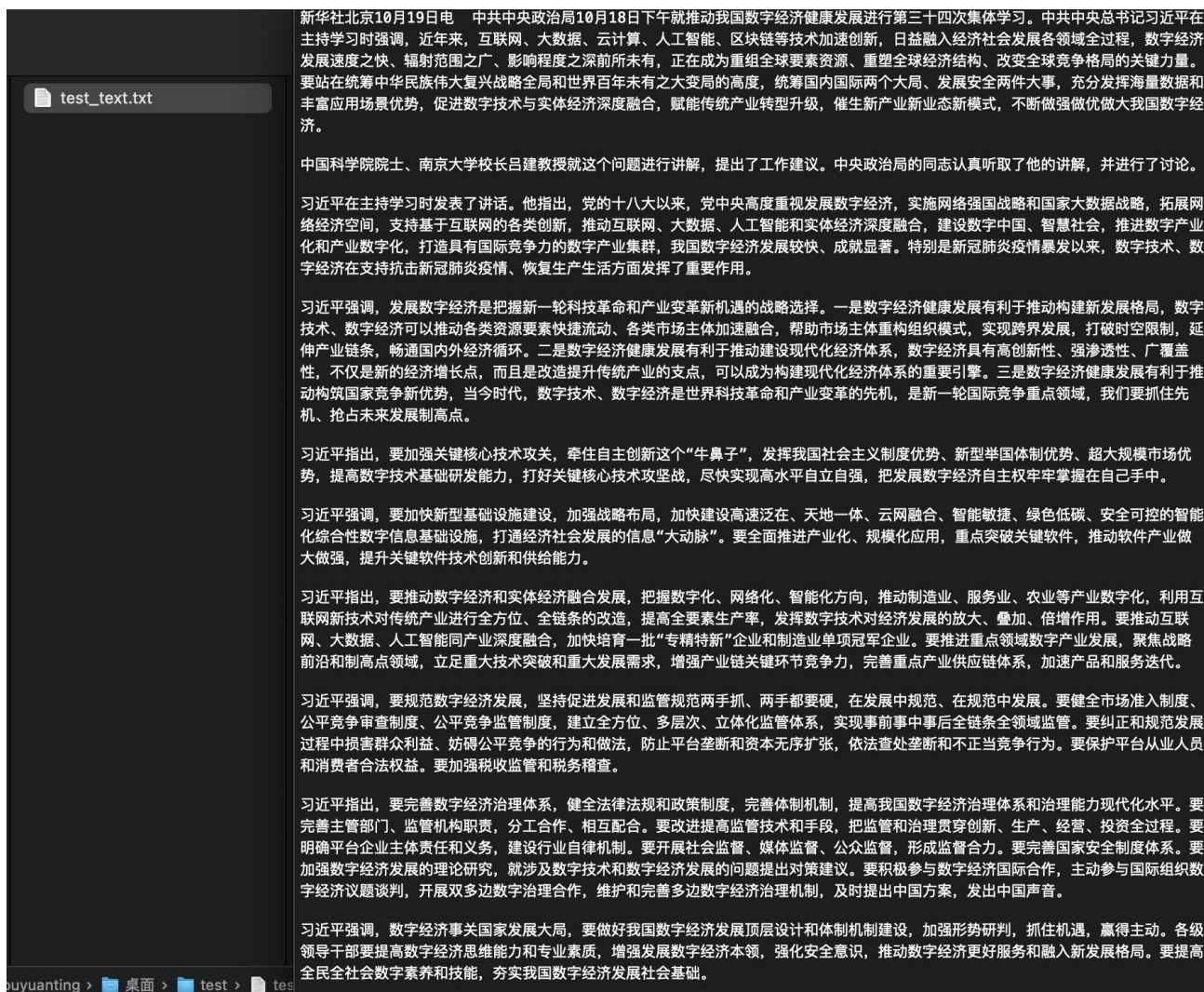
NLPIR-ICTCLAS汉语分词系统

<https://github.com/NLPIR-team/NLPIR>

- **主要功能**包括中文分词；英文分词；词性标注；命名实体识别；新词识别；**关键词提取（交叉信息熵算法）**；支持用户专业词典与微博分析。
- 支持多种操作系统（Windows、Linux等），多种开发语言与平台（C/C++/C#、Java、Python、Hadoop等）



- `nlpir.key_extract.import_dict(word_list: list) → list`
- `nlpir.key_extract.delete_user_word(word_list: list)`
- `nlpir.key_extract.import_blacklist(filename: str, pos_blacklist=typing.List[str]) → bool`
- ...
- `nlpir.key_extract.get_key_words(text: str, max_key: int = 50) → List[dict]`
 - 获取文本对应的关键词,以及对应的权值,词性,词频等信息



测试文本

习近平 系列重要讲话数据库



人民网 >> 中国共产党新闻网 >> 习近平系列重要讲话数据库

习近平在中共中央政治局第三十四次集体学习时强调

把握数字经济发展趋势和规律 推动我国数字经济健康发展

来源：人民网-人民日报 发布时间：2021-10-20



```
In [1]: import time
from snownlp import SnowNLP
import xmnlp
from pyhanlp import *
from nlpir import key_extract
```

```
In [2]: with open("g:/Desktop/test/test_text.txt", "r") as f:
text = f.read()
```

```
s = SnowNLP(text)
t1SnowNLP = time.time()
snowNLPtags = s.keywords(5)
t2SnowNLP = time.time()
print("SnowNLP: ", end="")
print(",".join(snowNLPtags), end="")
print("")
```

```
xmnlp.set_model('g:/Downloads/xmnlp-onnx-models')
t1xmNLP = time.time()
xmNLPtags = xmnlp.keyword(text, 5)
t2xmNLP = time.time()
print("xmNLP: ", end="")
for var in xmNLPtags:
    print(var)
print("")
```

```
t1HanLP = time.time()
HanLPtags = HanLP.extractKeyword(text, 5)
t2HanLP = time.time()
print("HanLP: ", end="")
print(",".join(HanLPtags), end="")
print("\n")
```

```
t1nlp = time.time()
nlpTags = key_extract.get_key_words(text, 5)
t2nlp = time.time()
print("NLPIR: ", end="")
for var in nlpTags:
    print(var)
```

关键词对比:

SnowNLP: 经济,数字,发展,产业,新

xmNLP: ('发展', 8.211491499369146)
('数字经济', 7.3708887710746)
('数字', 4.2225105950089326)
('习近平', 3.6402442137621684)
('推动', 3.329270354844655)

HanLP: 数字,经济,发展,产业,技术

NLPIR: {'freq': 29, 'pos': 'n_new', 'weight': 30.199034255113137, 'word': '数字经济'}
{'freq': 6, 'pos': 'n_new', 'weight': 21.691608909786563, 'word': '我国数字经济'}
{'freq': 31, 'pos': 'v', 'weight': 20.31373449776709, 'word': '发展'}
{'freq': 3, 'pos': 'n_new', 'weight': 14.551162767576342, 'word': '数字经济治理'}
{'freq': 7, 'pos': 'n_new', 'weight': 14.332681159996882, 'word': '数字技术'}

耗时对比:

SnowNLP: 0.4926900863647461

xmNLP: 1.3275861740112305

HanLP: 1.4350521564483643

NLPIR: 0.14505410194396973



YAKE

<https://github.com/LIAAD/yake>

- 论文: YAKE! Collection-Independent Automatic Keyword Extractor。
- 轻量级无监督自动关键字提取方法, 它基于从单个文档中提取的文本统计特征来选择文本中最重要的关键字
- 不能直接作用于中文的关键词提取

关键词提取使用的特征:

- **大写term (Casing)**: 大写字母的term (除了每句话的开头单词) 的重要程度比那些小写字母的term重要程度要大
 - **词的位置 (Word Position)**: 文本越开头的部分句子的重要程度比后面的句子重要程度要大
 - **词频 (Term Frequency)**: 一个词在文本中出现的频率越大, 相对来说越重要
 - **上下文关系 (Term Related to Context)**: 一个词与越多不相同的词共现, 该词的重要程度越低
 - **词在句子中出现的频率 (Term Different Sentence)**: 一个词在越多句子中出现, 相对更重要
- $S(t)$ 表示的是单词t的分值情况, 其中 $s(t)$ 分值越小, 表示的单词t越重要

$$S(t) = \frac{T_{Rel} * T_{Position}}{T_{case} + \frac{TF_{norm}}{T_{Rel}} + \frac{T_{Sentence}}{T_{Rel}}}$$

assuming default parameters

```
kw_extractor = yake.KeywordExtractor()
keywords = kw_extractor.extract_keywords(text)

for kw in keywords:
    print(kw)
```

specifying parameters

```
language = "en"
max_ngram_size = 3
deduplication_threshold = 0.9
deduplication_algo = 'seqm'
window_size = 1
num_of_keywords = 20

custom_kw_extractor = yake.KeywordExtractor(lan=language, n=max_ngram_size, dedupLim=deduplication_threshold)
keywords = custom_kw_extractor.extract_keywords(text)

for kw in keywords:
    print(kw)
```

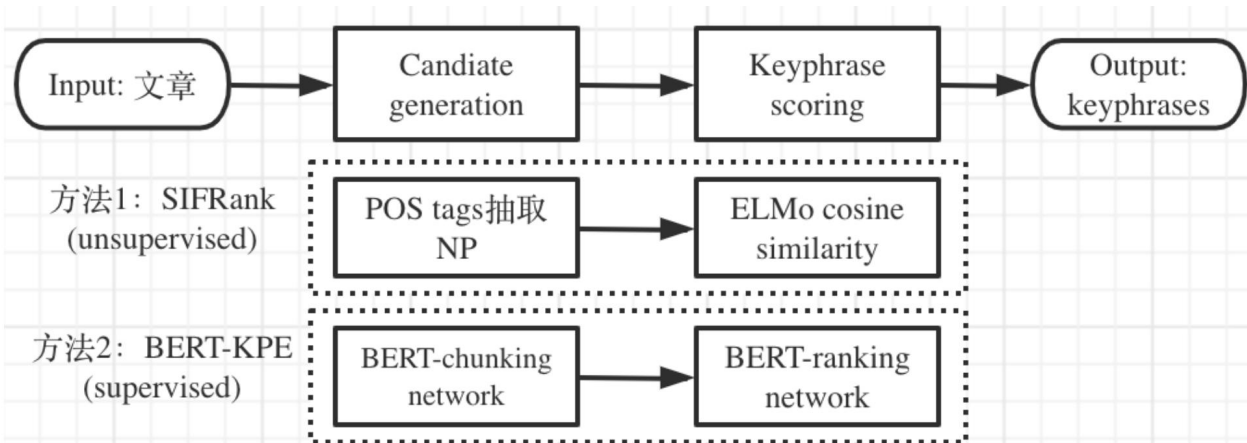

BERT-KPE

<https://github.com/thunlp/BERT-KPE>

- 论文: Capturing Global Informativeness in Open Domain Keyphrase Extraction <https://arxiv.org/abs/2004.13639>.
- 使用预训练模型BERT, 采用多任务学习范式, 同时训练分块网络chunking network和排名网络ranking network, 分别完成candidate chunking和ranking两个任务。此外, 该方法也非常灵活, 可以根据需要, 删除任意一个子网络, 仅训练另外一个子网络。

方法步骤:

1. Token embedding: 采用BERT将文本编码成向量。
2. N-gram representation: 使用CNN融合窗口大小为k的字对应的token embeddings, 得到相应的n-gram representation
3. Multi-task:
 1. Chunking network: 预测该n-gram是否是关键词候选
 2. Ranking network: 对n-gram进行排序





- TFIDF类的初始化

```
def __init__(self, idf_path=None):  
    # 加载  
    self.tokenizer = jieba.dt  
    self.postokenizer = jieba.posseg.dt  
    self.stop_words = self.STOP_WORDS.copy()  
    self.idf_loader = IDFLoader(idf_path or DEFAULT_IDF)  
    self.idf_freq, self.median_idf = self.idf_loader.get_idf()
```



- 主调函数TFIDF.extract_tags

```
def extract_tags(self, sentence, topK=20, withWeight=False, allowPOS=(), withFlag=False):  
    # 传入了词性限制集合  
    if allowPOS:  
        allowPOS = frozenset(allowPOS)  
        # 调用词性标注接口  
        words = self.postokenizer.cut(sentence)  
    # 没有传入词性限制集合  
    else:  
        # 调用分词接口  
        words = self.tokenizer.cut(sentence)  
    freq = {}  
    for w in words:  
        if allowPOS:  
            if w.flag not in allowPOS:  
                continue  
            elif not withFlag:  
                w = w.word  
        wc = w.word if allowPOS and withFlag else w  
        # 判断词的长度是否小于2, 或者词是否为停用词  
        if len(wc.strip()) < 2 or wc.lower() in self.stop_words:  
            continue  
        # 将其添加到词频词典中, 次数加1  
        freq[w] = freq.get(w, 0.0) + 1.0
```



- 主调函数TFIDF.extract_tags

```
# 统计词频词典中的总次数
total = sum(freq.values())
for k in freq:
    kw = k.word if allowPOS and withFlag else k
    # 计算每个词的tf-idf值
    freq[k] *= self.idf_freq.get(kw, self.median_idf) / total

# 根据tf-idf值进行排序
if withWeight:
    tags = sorted(freq.items(), key=itemgetter(1), reverse=True)
else:
    tags = sorted(freq, key=freq.__getitem__, reverse=True)
# 输出topK个词作为关键词
if topK:
    return tags[:topK]
else:
    return tags
```

- TextRank类的初始化

```
def __init__(self):  
    self.tokenizer = self.posttokenizer = jieba.posseg.dt  
    self.stop_words = self.STOP_WORDS.copy()  
    self.pos_filt = frozenset(('ns', 'n', 'vn', 'v'))  
    self.span = 5
```



- 主调函数TextRank.texrank

```
def texrank(self, sentence, topK=20, withWeight=False, allowPOS=('ns', 'n', 'vn', 'v'), withFlag=False):  
  
    self.pos_filt = frozenset(allowPOS)  
    #定义无向有权图  
    g = UndirectWeightedGraph()  
    #定义共现词典  
    cm = defaultdict(int)  
    #分词  
    words = tuple(self.tokenizer.cut(sentence))  
    #一次遍历每个词  
    for i, wp in enumerate(words):  
        #词i满足过滤条件  
        if self.pairfilter(wp):  
            # 依次遍历词i 之后窗口范围内的词  
            for j in xrange(i + 1, i + self.span):  
                # 词j 不能超出整个句子  
                if j >= len(words):  
                    break  
                # 词j不满足过滤条件, 则跳过  
                if not self.pairfilter(words[j]):  
                    continue  
                # 将词i和词j作为key, 出现的次数作为value, 添加到共现词典中  
            if allowPOS and withFlag:  
                cm[(wp, words[j])] += 1  
            else:  
                cm[(wp.word, words[j].word)] += 1
```



- 主调函数TextRank.textrank

```
# 依次遍历共现词典的每个元素，将词i，词j作为一条边起始点和终止点，共现的次数作为边的权重
for terms, w in cm.items():
    g.addEdge(terms[0], terms[1], w)
    # 运行textrank算法
nodes_rank = g.rank()
# 根据指标值进行排序
if withWeight:
    tags = sorted(nodes_rank.items(), key=itemgetter(1), reverse=True)
else:
    tags = sorted(nodes_rank, key=nodes_rank.__getitem__, reverse=True)
# 输出topK个词作为关键词
if topK:
    return tags[:topK]
else:
    return tags
```



- UndirectWeightedGraph类的初始化及添加边函数

```
def __init__(self):  
    self.graph = defaultdict(list)#这是进行分词后的一个词典  
  
def addEdge(self, start, end, weight):  
    # use a tuple (start, end, weight) instead of a Edge object  
    self.graph[start].append((start, end, weight))  
    self.graph[end].append((end, start, weight))
```



- 权重迭代收敛函数rank

```
def rank(self):
    ws = defaultdict(float)#权值list表
    outSum = defaultdict(float)
    # 初始化各个结点的权值
    # 统计各个结点的出度的次数之和
    wsdef = 1.0 / (len(self.graph) or 1.0)
    for n, out in self.graph.items():
        ws[n] = wsdef
        outSum[n] = sum((e[2] for e in out), 0.0)#e[2]是什么?

    # this line for build stable iteration
    sorted_keys = sorted(self.graph.keys())
    # 遍历若干次
    for x in xrange(10): # 10 iters
        #遍历各个节点
        for n in sorted_keys:
            s = 0
            # 遍历结点的入度结点
            for e in self.graph[n]:
                # 将这些入度结点贡献后的权值相加
                # 贡献率 = 入度结点与结点n的共现次数 / 入度结点的所有出度的次数
                s += e[2] / outSum[e[1]] * ws[e[1]]
            # 更新结点n的权值
            ws[n] = (1 - self.d) + self.d * s
```



- 权重迭代收敛函数rank

```
(min_rank, max_rank) = (sys.float_info[0], sys.float_info[3])
# 获取权值的最大值和最小值
for w in itervalues(ws):
    if w < min_rank:
        min_rank = w
    if w > max_rank:
        max_rank = w
# 对权值进行归一化
for n, w in ws.items():
    # to unify the weights, don't *100.
    ws[n] = (w - min_rank / 10.0) / (max_rank - min_rank / 10.0)

return ws
```

爬虫与关键词提取器结合对中国新闻进行分析

无监督

分为两步：

- 使用爬虫爬取相应网站的新闻信息
- 通过关键词提取器提取关键词

国内城市新闻分析：基于tfidf算法，利用jieba所提供的语料库进行计算
国际新闻分析：基于textrank



爬虫与关键词提取器结合对中国新闻关注点近些年发展的分析

对于不同地区新闻的分析：
样例分别为重庆和甘肃

chongqing
gansu

2021/10/23 20:58

2021/10/23 21:06

文本文档

文本文档



爬虫所爬到的数据文档

Analyzechongqing
Analyzegansu

2021/10/23 21:01

2021/10/23 21:07

文本文档

文本文档



关键词提取器的分析结果



爬虫与关

对于不同地
样例分别为

gansu - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

21年甘肃省新冠肺炎疫情防控工作(第五场)新闻发布会23日在兰州召开。会议披露, 22日0时
对当前兰州地区疫情防控要求和团体就餐需求加大、居民外出就餐不便等实际, 充分发挥中
配送企业和餐饮企业, 重点面向隔离小区、防疫点位、团体单位及社会大众, 每日提供主、
送食品种类、配送范围、配送方式、联系人、联系电话。杨丽平说, 各团体单位和市民, 根
厅启动12类72种商品价格及供应监测, 建立与全省重点保供企业“一对一”联系机制, 保障
小组办公室副主任、省政府副秘书长梁朝阳通报了该省新冠肺炎疫情最新情况以及疫情防控
急引入第三方核酸检测机构, 调配7台大型方舱实验室, 加快检测进度。下一步, 甘肃将继续
市工业和信息化局(以下简称“兰州市工信局”)22日披露, 根据电煤供需紧缺程度, 兰州严密
四场)新闻发布会22日在兰州召开, 会议披露, 自2021年10月21日16时至24时, 甘肃省新增
物资应急保障、保障通信供电稳定等方面为抗击疫情保驾护航。在保障电煤供应稳定的过程
0万吨煤炭采购意向, 作为兰州煤炭短缺时的备用煤源。督促华能兰州范坪热电公司、华能
兰州市工信局督促三大电厂和靖煤、华煤近期补签电煤供应合同, 并派专人到新炭供应企
季电煤保供运输协调组, 制定了《兰州市2021年-2022年供暖期电煤保供预案》, 会同相关
煤底线, 按程序中断部分高能耗工业企业煤炭供应, 优先保障兰州地区发电供热用煤之外, 以
形势, 兰州市工信局将积极协调、联动调度, 多措并举加强供需保障, 落实各项应对措施, 在
当息, 10月22日0-24时, 甘肃省新增确诊病例17例。其中, 兰州13例(12例为云南旅行团成
现住云南省昆明市盘龙区穿金路金林碧水小区, 与确诊病例105、106、107同为昆明来甘
昆明市盘龙区穿金路金林碧水小区。10月10日-16日与确诊病例111同程在内蒙古额济纳旗,
省昆明市西山区白马小区。10月10日-16日与确诊病例111同程在内蒙古额济纳旗, 甘肃省金
划纲要》发
蒙, 融入大循环
含5个方面内容。一是坚持发展共享。齐心协力把联盟打造成为集资源集聚、协同创新、技术

第 19 行, 第 1 列 100% Windows (CRLF) UTF-8

第 14 行, 第 21 列 100% Windows (CRLF) UTF-8

爬虫与关键词提取器结合对中国新闻关注点近些年发展的分析

Analyzechongqing - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

重庆/生态/创新/发展/成渝/双城/经济圈/国际/绿色/产业/

→ 经济发展、生态保护

Analyzegansu - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

病例/核酸/确诊/兰州/疫情/兰州市/检测/新冠/隔离/甘肃省/

→ 疫情反复



爬虫与关键词提取器结合对中国新闻关注点近些年发展的分析

对中国新闻网国际模块的新闻分析：

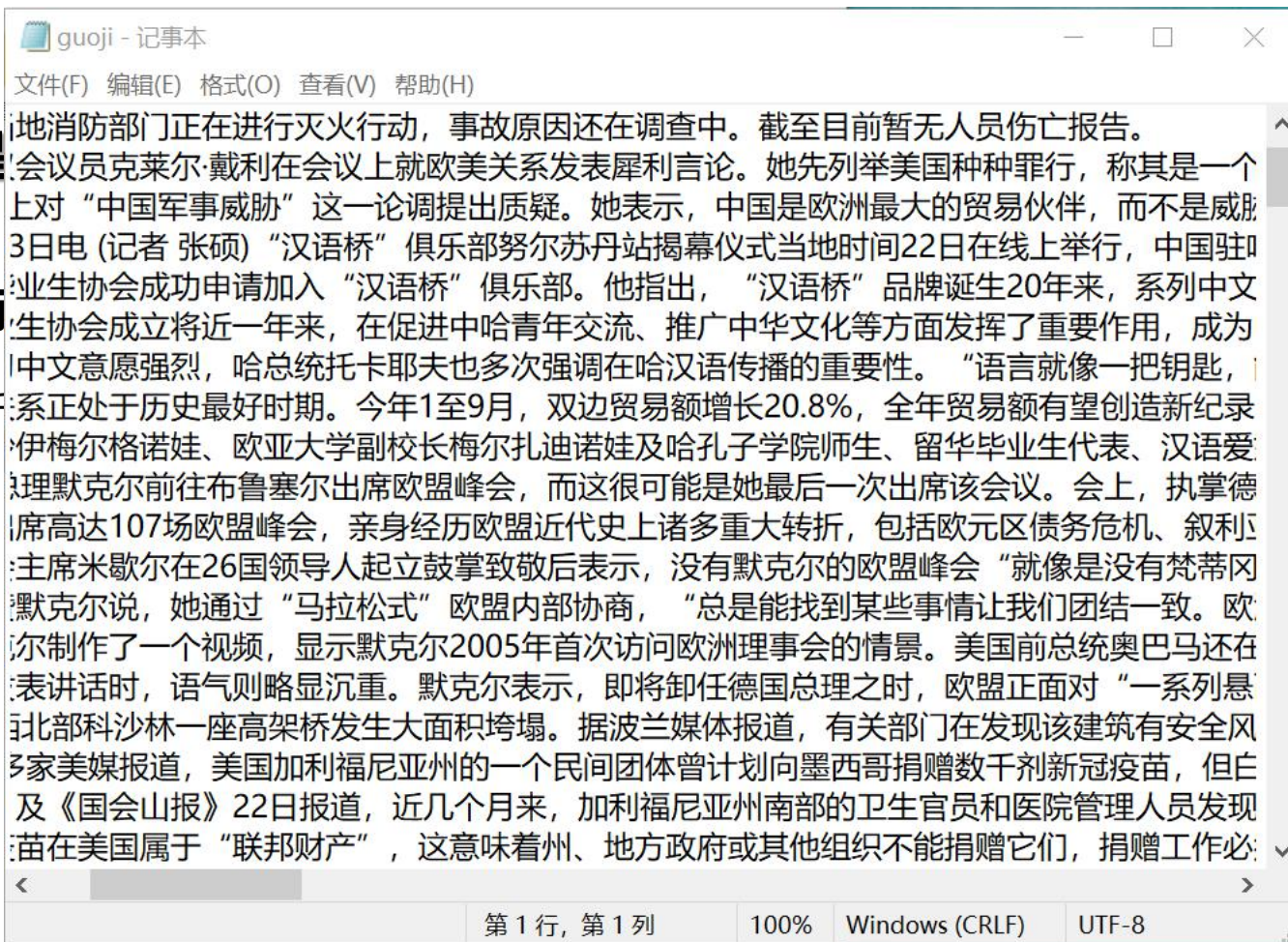
由于对该模块的分析主要是想研究近期与中国冲突较多的国家，与国际近期发生频率较高的事件。采用textrank算法。



爬虫与关键词提

对中国新闻网国际

由于对该模块的分析
textrank算法。



事件。采用

爬虫与关键词提取器结合对中国新闻关注点近些年发展的分析

对中国新闻网国际模块的新闻分析：

关键词分析所得到的结果：

```
keywords by textRank:
```

```
美国/新冠/疫苗/病例/问题/新疆/疫情/供应链/表示/中国/
```

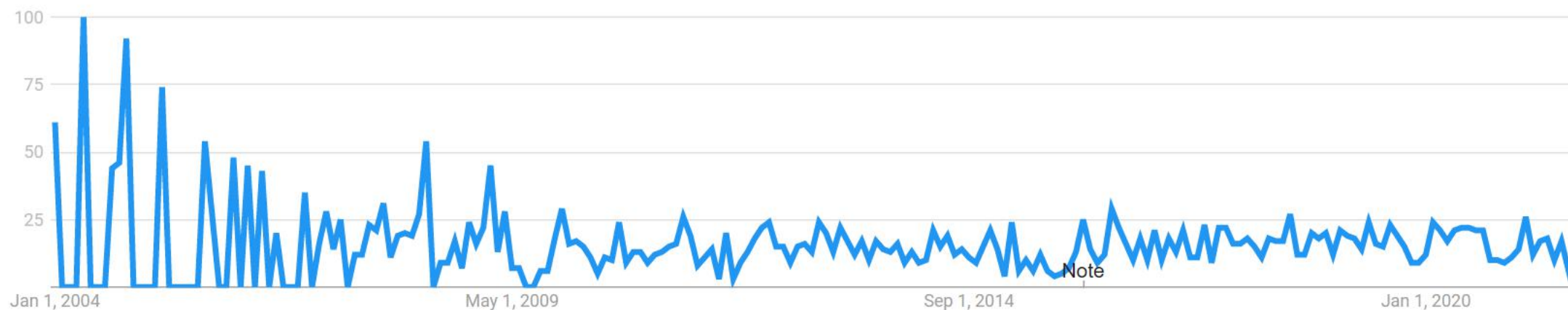


北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

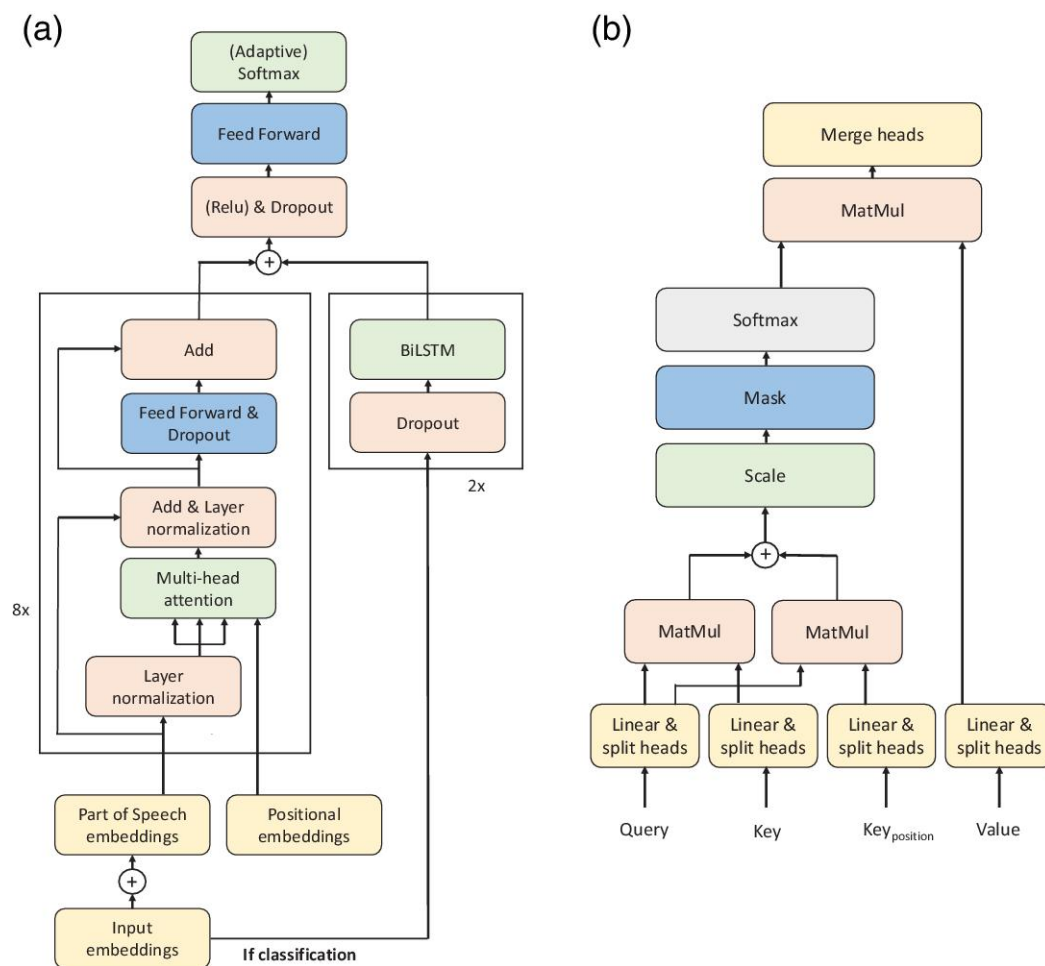
5

关键词提取前沿研究介绍

汇报人：冀温瑾

Interest over time 

- Data were collected from Google Scholar by searching 'keyword extraction' and 'keyphrase extraction'.



[1] Martinc M , Krlj B , Pollak S . TNT-KID: Transformer-based neural tagger for keyword identification[J]. Natural Language Engineering, 2021:1-40.

	KP20k	Inspec	Krapivin	NUS	SemEval	KPTimes	JPTimes	DUC	Average
Unsupervised algorithms									
Tfidf									
F1@5	0.072	0.160	0.067	0.112	0.088	0.179*	0.266*	0.098*	0.130
F1@10	0.094	0.244	0.093	0.140	0.147	0.151*	0.229*	0.120*	0.152
TextRank									
F1@5	0.181	0.286	0.185	0.230	0.217	0.022*	0.012*	0.120*	0.157
F1@10	0.151	0.339	0.160	0.216	0.226	0.030*	0.026*	0.181*	0.166
YAKE									
F1@5	0.141*	0.204*	0.215*	0.159*	0.151*	0.105*	0.109*	0.106*	0.149
F1@10	0.146*	0.223*	0.196*	0.196*	0.212*	0.118*	0.135*	0.132*	0.170
RaKUn									
F1@5	0.177*	0.101*	0.127*	0.224*	0.167*	0.168*	0.225*	0.189*	0.172
F1@10	0.160*	0.108*	0.106*	0.193*	0.159*	0.139*	0.185*	0.172*	0.153
Key2Vec									
F1@5	0.080*	0.121*	0.068*	0.109*	0.081*	0.126*	0.158*	0.062*	0.101
F1@10	0.090*	0.181*	0.082*	0.121*	0.126*	0.116*	0.145*	0.078*	0.117
EmbedRank									
F1@5	0.135*	0.345*	0.149*	0.173*	0.189*	0.063*	0.081*	0.219*	0.169
F1@10	0.134*	0.394*	0.158*	0.190*	0.217*	0.057*	0.074*	0.246*	0.184

Supervised algorithms

KEA									
F1@5	0.046	0.022	0.018	0.073	0.068	/	/	/	/
F1@10	0.044	0.022	0.017	0.071	0.065	/	/	/	/
Maui									
F1@5	0.005	0.035	0.005	0.004	0.011	/	/	/	/
F1@10	0.005	0.046	0.007	0.006	0.014	/	/	/	/
Semi-supervised CopyRNN									
F1@5	0.308	0.326	0.296	0.356	0.322	/	/	/	/
F1@10	0.245	0.334	0.240	0.320	0.294	/	/	/	/
CopyRNN									
F1@5	0.317	0.244	0.305	0.376	0.318	0.406*	0.256*	0.083	0.288
F1@10	0.273	0.289	0.266	0.352	0.318	0.393	0.246	0.105	0.280
CatSeqD									
F1@5	0.348	0.276	0.325	0.374	0.327	0.424*	0.238*	0.063*	0.297
F1@10	0.298	0.333	0.285	0.366	0.352	0.424*	0.238*	0.063*	0.295
CorrRNN									
F1@5	/	/	0.318	0.361	0.320	/	/	/	/
F1@10	/	/	0.278	0.335	0.320	/	/	/	/
GPT-2									
F1@5	0.275*	0.413*	0.253*	0.318*	0.257*	0.421*	0.331*	0.298*	0.321
F1@10	0.278*	0.469*	0.253*	0.323*	0.278*	0.423*	0.336*	0.312*	0.334
GPT-2 + BiLSTM-CRF									
F1@5	0.355*	0.462*	0.287*	0.329*	0.246*	0.478*	0.386*	0.333*	0.360
F1@10	0.360*	0.524*	0.288*	0.336*	0.274*	0.479*	0.389*	0.371*	0.378
TNT-KID									
F1@5	0.336*	0.460*	0.310*	0.350*	0.283*	0.485*	0.359*	0.318*	0.363
F1@10	0.338*	0.536*	0.320*	0.358*	0.337*	0.485*	0.361*	0.373*	0.389



近年来，中译英、英译中等“翻译抄袭”作为一种较为隐蔽的学术不端行为，引起了高校、期刊编辑部等越来越多学术科研机构的关注。针对这一现象，知网学术不端文献检测系统在国内首个实现了跨语言检测功能，为抄袭检测提供了更加严谨科学的检测手段，自发布以来得到论文诚信管理部门的一致认可。

跨语言关键字检测需要模型对于少数据样本训练，有相关论文想法是在多语言语料库上预训练模型，在一种语言上对其进行fine-tune，然后在第二种语言上对模型进行零样本跨语言测试。



- [1] Mihalcea R, Tarau P. TextRank: Bringing order into texts[C]. Association for Computational Linguistics, 2004.
- [2] https://pdf.hanspub.org/CSA20130100000_81882762.pdf
- [3] <https://github.com/bitcarmanlee/easy-algorithm-interview-and-practice>
- [4] 时永宾,余青松.基于共现词卡方值的关键词提取算法[J].计算机工程,2016,42(06):191-195.
- [5] 徐涛,蓝传铤.基于卡方统计量的藏文新闻网页关键词提取方法[J].电脑知识与技术,2017,13(26):171-173.
- [6] 叶秋永,吴华琼,宋继华. 基于信息增益的中文术语抽取[A]. 中文教学现代化学会.数字化对外汉语教学实践与反思[C].中文教学现代化学会:中文教学现代化学会,2010:7.
- [7] 赵京胜,朱巧明,周国栋,张丽.自动关键词抽取研究综述[J].软件学报,2017,28(09):2431-2449.
- [8] <https://github.com/fxsjy/jieba/blob/master/jieba/analyse/tfidf.py>
- [9] <https://github.com/fxsjy/jieba/blob/master/jieba/analyse/textrank.py>
- [10] Martinc M , Krlj B , Pollak S . TNT-KID: Transformer-based neural tagger for keyword identification[J]. Natural Language Engineering, 2021:1-40.



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

谢谢观看
敬请老师批评指正