

算法 1：基于关键字的正样本资讯筛选算法

输入：关键字文件列表： pos_files^1 , neg_files^2 , $pos_mixed_files^3$, $neg_mixed_files^4$
待标记的资讯文件： $unlabel_file^5$

输出：正样例集合 $S_1 = \{news_1, news_2, \dots, news_m\}$
负样例集合 $S_{-1} = \{news_1, news_2, \dots, news_n\}$

算法步骤：

(1) 读取各个文件关键字：

```
POS_KEYS = read_key_files6(pos_files)
NEG_KEYS = read_key_files(neg_files)
POS_MIXED_KEYS = read_key_files(pos_mixed_files)
NEG_MIXED_KEYS = read_key_files(neg_mixed_files)
UNLABEL_NEWS = read_news_file(unlabel_file)
 $S_1 = \emptyset$ ,  $S_{-1} = \emptyset$ 
```

(2) **FOR** 每一个 $NEWS \in UNLABEL_NEWS$:

```
IF contain_key7(NEWS, NEG_KEYS):
     $S_{-1} = S_{-1} \cup NEWS$ 
ELSE IF contain_mixed_key8(NEWS, POS_MIXED_KEYS):
     $S_1 = S_1 \cup NEWS$ 
ELSE IF contain_mixed_key(NEWS, NEG_MIXED_KEYS):
     $S_{-1} = S_{-1} \cup NEWS$ 
ELSE IF contain_key(NEWS, POS_KEYS):
     $S_{-1} = S_{-1} \cup NEWS$ 
ELSE IF:
     $S_{-1} = S_{-1} \cup NEWS$ 
END
```

¹ $pos_files = [\text{'keywords_EVT.txt'}, \text{'keywords_GOV.txt'}, \text{'keywords_IDX.txt'}, \text{'keywords_VIP.txt'}, \text{'keywords_POS.txt'}]$

² $neg_files = [\text{'keywords_BLK.txt'}, \text{'keywords_ORG.txt'}, \text{'keywords_PER.txt'}, \text{'keywords_GPE.txt'}, \text{'keywords_NEG.txt'}]$

³ $pos_mixed_files = [\text{'mixed_keywords_POS.txt'}]$

⁴ $neg_mixed_files = [\text{'mixed_keywords_NEG.txt'}]$

⁵ $unlabel_file$: 需要标注资讯，文件格式：每行一条资讯

⁶ 读取文件集函数，每个文件存储一类关键字。文件格式：一个关键字占一行，详见文件 $keywords_EVT.txt$

⁷ 判断是否含有关键字。匹配规则：NEWS（单个）匹配到 KEYS 中一个关键字即为成功

⁸ 逻辑同 7，主要区别在关键字组成：每个关键字由 '+' 连接，如 '黄金' + '美元'，必须同时匹配到才算成功。组合关键字组成详见文件 $mixed_keywords_POS.txt$