



学生姓名: \_\_\_\_\_ 学号: \_\_\_\_\_ 培养单位: \_\_\_\_\_ 分数: \_\_\_\_\_

一、选择题（单选，每题 2 分，共 20 分）

1. 关于布尔检索，以下说法正确的是（ ）
  - A) 在布尔检索系统中，进行词干还原从不降低正确率
  - B) 在布尔检索系统中，进行词干还原从不降低召回率
  - C) 词干还原会增加词项词典的大小
  - D) 词干还原应该在构建索引时调用，而不应在查询处理时调用
2. 关于倒排索引，以下说法错误的是（ ）
  - A) 词典的开销通常小于倒排记录表的开销 ✓
  - B) 倒排索引中往往会存储词项的频率 ✓
  - C) 词典中存储的都是词
  - D) 倒排索引可以加快搜索的速度
3. 关于向量空间模型的特点，下面说法不正确的是（ ）
  - A) 支持部分匹配和近似匹配，结果可以排序 ✓
  - B) 理论上不够严谨，往往基于直觉的经验性公式 ✓
  - C) 词项之间的独立性假设与实际不符 ✓
  - D) 应用于检索问题时，效果不如概率检索模型
4. 关于tf-idf，以下说法不正确的是（ ）
  - A) 一个罕见词的idf往往很高
  - B) idf的大小是无限的
  - C) 词项的tf-idf权重可以超过1 ✓
  - D) 词项t在所有文档中出现，则权重取值很小 ✓
5. 关于隐式相关反馈，下面说法错误的是（ ）
  - A) 不需要用户显式参与，减轻用户负担 ✓
  - B) 对行为分析有较高要求 ✓
  - C) 一定能提升检索准确率
  - D) 某些情况下需要增加额外设备 ✓



9. 伪相关反馈一定能提高每个查询的检索效果。( ) **✗**

10. 倒排索引中的VB编码压缩方法是一种无损压缩方法。( ) **✓**

三、计算题（每题 8 分，共 32 分）

1. 给定矩阵  $C$ ， $C = \begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}$ ，利用SVD分解，得到分解结果为  $U = \begin{pmatrix} -0.816 & 0.000 \\ -0.408 & -0.707 \\ -0.408 & 0.707 \end{pmatrix}$ ， $\Sigma = \begin{pmatrix} 1.732 & 0.000 \\ 0.000 & 1.000 \end{pmatrix}$ ， $V^T = \begin{pmatrix} -0.707 & -0.707 \\ 0.707 & -0.707 \end{pmatrix}$  计算矩阵  $C$  的1-秩逼近  $C_1$ ，并给出该逼近下的F范数误差值。

2. 已知采用欧式归一化方法对表1中三个文档的tf值进行归一化，得到Doc1、Doc2和Doc3的  $\sqrt{\sum_{i=1}^M V_i^2(d)}$  值分别是30.56、46.84和41.3，

	Doc1	Doc2	Doc3
北京	27	4	24
中国	3	33	0
天安门	0	33	29
科技	14	0	17

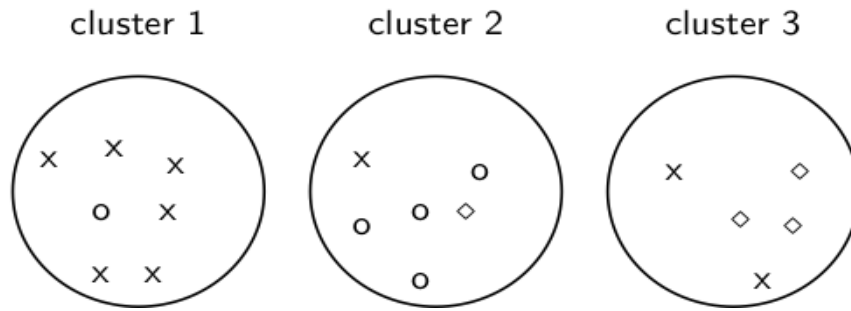
表1 tf值

- a) 请计算tf归一化的结果，填入下表：

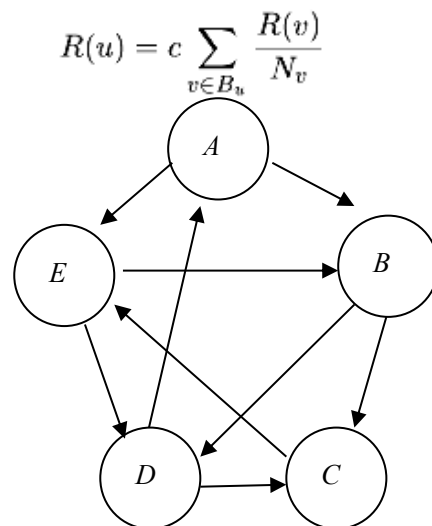
	Doc1	Doc2	Doc3
北京	0.884	0.085	0.581
中国	0.098	0.705	0
天安门	0	0.705	0.702
科技	0.458	0	0.412

- b) 已知Doc1、Doc2和Doc3的静态得分分别是0.25，0.5和0.8，画出当使用静态得分与欧式归一化tf值求和的结果进行排序的倒排记录表。

3. 有18篇已经标注好的文档，所属类别分别用×、○、◇表示。通过聚类算法得到了三个簇，如下图所示，请计算兰迪指数（Rand Index），并附计算过程。



4. 假设五个网页A、B、C、D、E构成如下的链接关系，试采用如下公式( $u$ 指当前页面， $B_u$ 是所有链接到 $u$ 的页面集合， $N_v$ 是页面 $v$ 的出链(outlink)总数， $R(u)$ 和 $R(v)$ 分别是 $u$ 和 $v$ 的PageRank值，为计算方便，假定 $c=1$ )计算每个网页的归一化PageRank值，即五个网页的PageRank总和为1。



#### 四、简答题（共 28 分）

1. 通过对信息检索的结果进行聚类，可以为用户提供更有效的展示。若采用 K-均值聚类算法，为了达到较好的聚类效果，请描述下可能采用的初始质心向量（种子）的选择方法。（8 分）
2. 为了提高查询召回率，请设计一种不带相关反馈的查询扩展方法。请详细描述你的思路和做法。（10分）
3. 请设计一种语义检索模型，即查询和文档之间用词可能完全不同，但是语义近似。请详细描述你的思路和做法。（10分）