

实验四

12.26

实验内容

- 主题：上市公司财经新闻情感分析
- 实验目标：使用多种机器学习算法对文本进行情感判别，包括KNN、决策树、朴素贝叶斯、支持向量机等，学习如何进行模型训练，如何进行分类预测。要求使用至少两种分类方法。
- 要求：核心程序在MapReduce上运行，要求使用至少两种分类方法。
- 实现：
 - KNN
 - Naïve Bayes

逻辑结构

- 数据处理

- KNN

- 训练
- 预测

- Bayes

- 训练
- 预测

遇到的问题：

- Double变量精确度问题
- 导致最后统计出来的概率都是0
- 采用统计前500个高频词的方法

数据处理实现

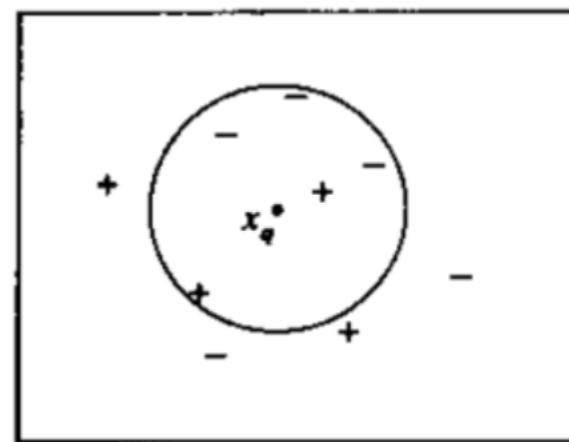
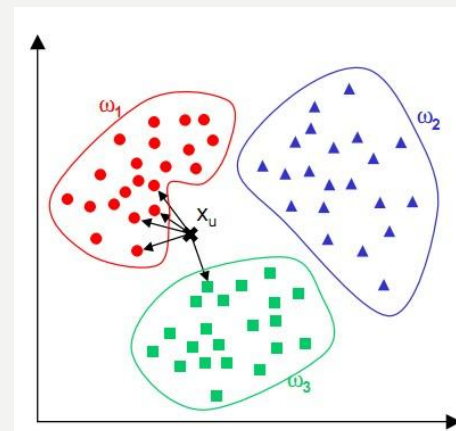
```
java
1 public class DATACREATION{
2     public static class DCMapper{
3         public void map{
4             分词, 并根据指定文件, 制作词向量;
5             保存的格式为
6             <词向量, 标签>
7         }
8     }
9     public static class DC_Another_Mapper{
10        public void setup{}
11        public void map{
12            对新闻标题数据预处理
13            分词, 制作词向量
14            保存格式为
15            <新闻标题, 词向量>
16        }
17    }
18 }
```

这里没必要用到reduce
特征选择: TF-IDF
稀疏矩阵: 统计词频,
选前500个

KNN实现

```
java
1 public class KNN{
2     public static int getdistance(String a,int[] b){}
3     //计算两个点之间的距离
4
5     public static String[] GET_TRIAN_DATA{}
6     //获得训练集数据
7
8     public static class KNNMapper(){
9         public void setup(){}
10        public void map(){
11            计算测试词向量到各个训练词向量的距离
12            取最近的三个 (k=3) 向量, 统计label
13            投票票数最多的便是预测结果
14        }
15    }
16
17
18 }
19
```

KNN介绍



```

1  java
2  public class BAYES{
3      //首先进行训练
4      public static class TrainMapper{
5          public void map(){
6              //输入训练集
7              //统计词向量中每个词的属性 (局部)
8              //比如 “下滑”: <pos,15%> <neu,15%> <neg,70%>
9          }
10     }
11     public static class TrainReducer{
12         public void reduce(){
13             //统计词向量中每个词的属性 (全局)
14         }
15     }
16     //然后进行预测
17     public static class TestReducer{
18         public void reduce(){
19             //输入测试样本数据集
20             //计算 三个词频*属性, 取最大值, 贴上标签
21             //输出分类结果
22         }
23     }
24 }

```

Bayes介绍

$$\begin{aligned}
 p(x|c_j) &= p([w_1, w_2, \dots, w_{|x|}]|c_j) = \prod_{h=1}^{|x|} p(w_h|c_j) \\
 &= \prod_{i=1}^V p(t_i|c_j)^{N(t_i,x)} = \prod_{i=1}^V \theta_{i|j}^{N(t_i,x)} \\
 p(y = c_j) &= \pi_j
 \end{aligned}$$

$$P(C | F_1) = \frac{P(CF_1)}{P(F_1)} = \frac{P(C) \cdot P(F_1 | C)}{P(F_1)}$$

港股14110才 券2400亿：银行股十年融资511亿 九八取渔	positive
上市银行屡吃罚单 中介机构：不构成重大违法违规	positive
专题报告：环杭州湾大湾区与粤港澳大湾区的对比分析	positive
估值修复遇上业绩超预期 银行股演绎“王者归来”	positive
健康元隐瞒多次环保处罚 324亿市值蒸发配股股民苦	negative
南京银行核心一级资本充足率告急 拟定增140亿元补血	positive
历史遗留问题拖累业绩 正虹科技重组预期降温	neutral
同业存单减速 中小行最受伤	negative
国家队增持南京银行 持股4.99%逼近举牌线	positive
探访茅台实体店：最严限价令下的失衡供需	negative
浦发银行同业业务违规被罚款200万 违反审慎经营规则	negative
浦发银行谢伟：银行资管业务未来发展空间巨大	positive
深交所与浦发开展战略合作 推动金融和科技深度融合	positive
茅台成25年A股第一高价股 总市值超半个贵州GDP	positive
跨界大王群兴玩具：我有故事你带钱 九年盈利不抵减持	positive
这些国家队重仓股 还被机构看好	positive
透视银行股大涨之谜：国家队增持南京银行逼近举牌线	positive
郑煤机(00564-HK)认购浦发银行2亿元人民币理财产品	positive
银行股脉冲式上演逆袭行情 工商银行股价创下9年新高	positive
麦当劳今天宣布：2700多家餐厅归中国公司了	neutral

