

朴素贝叶斯

最适合简单的文本分析的算法

朴素贝叶斯 — 核心思想

请在此刀隔线以上回复内容

您提交的 # 3152号工单：来自于李先生的留言 有更新。

请点击以下链接查看工单处理进度：

<https://tingyun.kf5.com/hc/request/view/3152/>

要添加另外的工单评论，请回复此邮件。

在垃圾邮件里经常出现“**链接**”，“**点击**”这种单词。



假如一个邮件里包含了这些单词，这个邮件很可能是垃圾邮件。

一个简单的例子：垃圾邮件分类

Part 1: 朴素贝叶斯核心

垃圾邮件分类



正常邮件

垃圾邮件



Q: 新的邮件，属于哪一个分类？

对于“购买”单词

假设：每个邮件包含10个单词

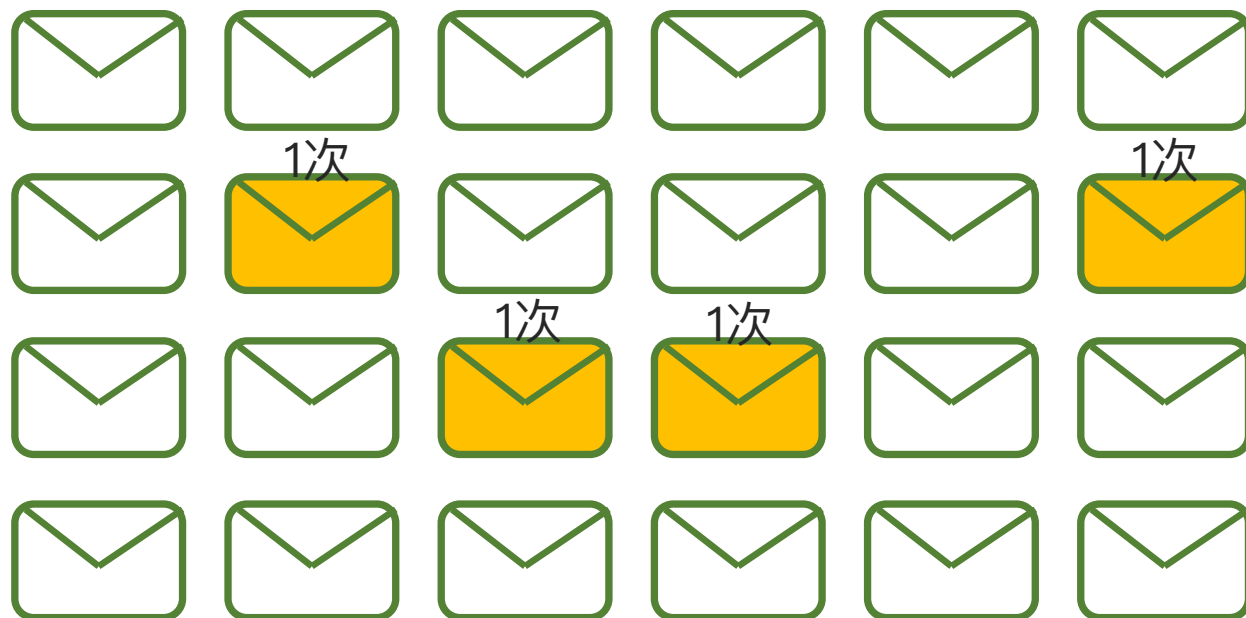


正常邮件含有“购买”词的概率多少？ $p(\text{“购买”}|\text{正常}) = 3/(24*10) = 1/80$

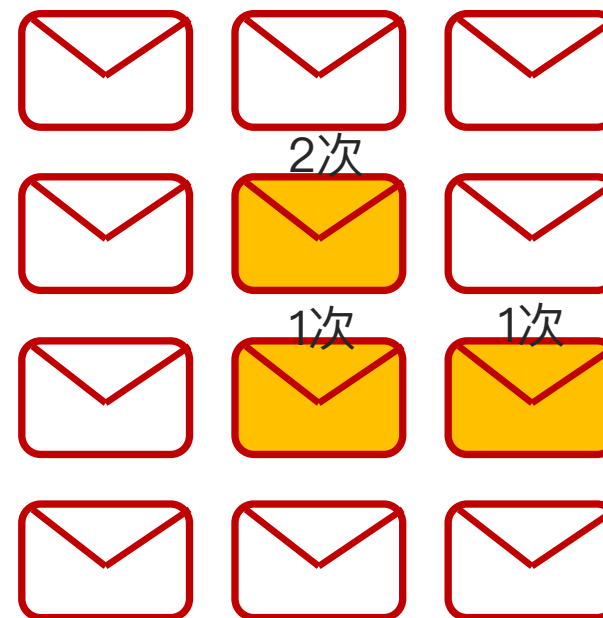
垃圾邮件含有“购买”词的概率多少？ $p(\text{“购买”}|\text{垃圾}) = 7/(12*10) = 7/120$

对于“物品”单词

正常邮件



垃圾邮件



正常邮件含有“物品”词的概率多少？

垃圾邮件含有“物品”词的概率多少？

对于“不是”单词



正常邮件含有“不是”词的概率多少？

垃圾邮件含有“不是”词的概率多少？

对于“广告”单词



正常邮件含有“广告”词的概率多少？

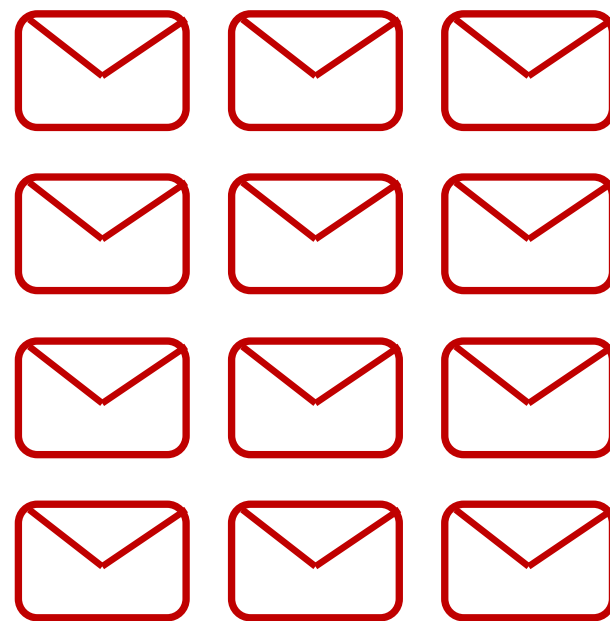
垃圾邮件含有“广告”词的概率多少？

还有一个概率需要统计

正常邮件



垃圾邮件



有多少邮件是正常邮件（百分比）？

有多少邮件是垃圾邮件（百分比）？

利用朴素贝叶斯识别垃圾邮件

前面计算出了很多零零散散的概率，怎么整合这些信息来完成识别任务？

从概率统计的角度

$P(\text{垃圾}|\text{邮件内容})$: 一个邮件内容为垃圾邮件的概率

$P(\text{正常}|\text{邮件内容})$: 一个邮件内容为正常邮件的概率

如何做判断？

如果 $P(\text{垃圾}|\text{邮件内容}) > P(\text{正常}|\text{邮件内容})$, 则可以认为是垃圾邮件

如果 $P(\text{垃圾}|\text{邮件内容}) \leq P(\text{正常}|\text{邮件内容})$, 则可以认为是正常邮件

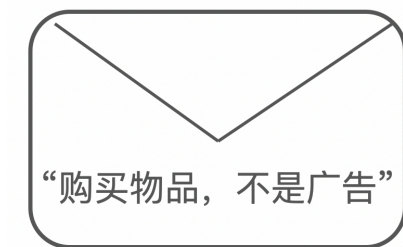
新的问题: $P(\text{垃圾}|\text{邮件内容})$, $P(\text{正常}|\text{邮件内容})$ 怎么计算???

回顾：贝叶斯定理



$P(\text{"购买"} \text{正常}) =$	$P(\text{"购买"} \text{垃圾}) =$
$P(\text{"物品"} \text{正常}) =$	$P(\text{"物品"} \text{垃圾}) =$
$P(\text{"不是"} \text{正常}) =$	$P(\text{"不是"} \text{垃圾}) =$
$P(\text{"广告"} \text{正常}) =$	$P(\text{"广告"} \text{垃圾}) =$

$P(\text{正常}) =$ $P(\text{垃圾}) =$



新邮件

是垃圾邮件还是正常邮件?

$$\begin{aligned}
 P(\text{正常} \mid \text{邮件内容}) &= \frac{P(\text{邮件内容}|\text{正常}) * p(\text{正常})}{p(\text{邮件内容})} \\
 &= \frac{P(\text{"购买", "物品", "不是", "广告"}|\text{正常}) * p(\text{正常})}{p(\text{邮件内容})} \\
 &= \frac{P(\text{"购买"}|\text{正常}) * P(\text{"物品"}|\text{正常}) * P(\text{"不是"}|\text{正常}) * P(\text{"广告"}|\text{正常}) * p(\text{正常})}{p(\text{邮件内容})}
 \end{aligned}$$

$$P(\text{垃圾} \mid \text{邮件内容}) = \frac{P(\text{邮件内容}|\text{垃圾}) * p(\text{垃圾})}{p(\text{邮件内容})}$$

怎么处理概率为0的情况？

$$P(\text{正常} \mid \text{邮件内容}) = \frac{P(\text{"购买"} \mid \text{正常}) * P(\text{"物品"} \mid \text{正常}) * P(\text{"不是"} \mid \text{正常}) * P(\text{"广告"} \mid \text{正常}) p(\text{正常})}{p(\text{邮件内容})}$$

手推一个完整的例子

垃圾邮件

1. 点击 更多 信息
2. 最新 产品
3. 信息 点击 链接

正常邮件

1. 开会
2. 信息 详见 邮件
3. 最新 信息

新邮件

最新 产品 实惠 点击 链接



属于正常邮件还是垃圾邮件?

手推一个完整的例子

利用严格的数学来表示朴素贝叶斯过程

给定一个文本向量 $x = (x_1, \dots, x_N)$, 以及两种类型标签 $y = \{0, 1\}$,

其中 N 为词典库的大小, θ 为模型的参数, 其中 θ_{i0} 和 θ_{i1}

分别表示第 i 个参数出现在分类 0 和分类 1 的概率。 x_i 表示词典中第 i 个单词出现的次数。

Coding Time-1

Part 2: 文本表示

单词的表示

词典：[我们，去，爬山，今天，你们，昨天，跑步]

每个单词的表示：

我们：

爬山：

跑步：

昨天：

句子的表示 (boolean)

词典：[我们，又，去，爬山，今天，你们，昨天，跑步]

每个句子的表示

我们 今天 去 爬山：

你们 昨天 跑步：

你们 又 去 爬山 又 去 跑步：

句子的表示 (count)

词典：[我们，又，去，爬山，今天，你们，昨天，跑步]

每个句子的表示

我们 今天 去 爬山：

你们 昨天 跑步：

你们 又 去 爬山 又 去 跑步：

句子的表示

句子1: He is going from Beijing to Shanghai

句子2: He denied my request, but he actually lied.

句子3: Mike lost the phone, and phone was in the car

句子1: (0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0)

句子2: (1, 0, 0, 1, 0, 1, 0, 0, 2, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0)

句子3: (0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 2, 0, 0, 2, 0, 1)

denied

he

这种表示有什么缺点？

句子的表示

句子1: He is going from Beijing to Shanghai

句子2: He denied my request, but he actually lied.

句子3: Mike lost the phone, and phone was in the car

句子1: (0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0)

句子2: (1, 0, 0, 1, 0, **1**, 0, 0, **2**, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0)


句子3: (0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 2, 0, 0, 2, 0, 1)

denied

he

并不是出现的越多就越重要!
并不是出现的越少就越不重要!

Tf-idf 表示

$$tfidf(w) = tf(d, w) * idf(w)$$


文档 d 中 w 的词频

$\log \frac{N}{N(w)}$

N : 语料库中的文档总数

$N(w)$: 词语 w 出现在多少个文档?

$$tfidf(w) = tf(d, w) * idf(w)$$

今天 上 机器学习 课程

今天 的 课程 有 意思

数据 课程 也 有 意思

Coding Time-1

Part 3: Extensions (可选)

1. 当特征为实数型的时候

2. 为什么叫“朴素”? — 条件独立

$$\begin{aligned}P(\text{正常} \mid \text{邮件内容}) &= \frac{P(\text{邮件内容} \mid \text{正常}) * p(\text{正常})}{p(\text{邮件内容})} \\&= \frac{P(\text{“购买”, “物品”, “不是”, “广告”} \mid \text{正常}) * p(\text{正常})}{p(\text{邮件内容})} \\&= \frac{P(\text{“购买”} \mid \text{正常}) * P(\text{“物品”} \mid \text{正常}) * P(\text{“不是”} \mid \text{正常}) * P(\text{“广告”} \mid \text{正常}) * p(\text{正常})}{p(\text{邮件内容})}\end{aligned}$$

3. 朴素贝叶斯的极大似然

$$\begin{aligned} p(D) &= \prod_{i=1}^N p(x^i, y^i) = \prod_{i=1}^N p(x^i | y^i) p(y^i) \\ &= \prod_{i=1}^N p(x_1^i, x_2^i, \dots, x_{m_i}^i | y^i) p(y^i) \\ &= \prod_{i=1}^N \prod_{j=1}^{m_i} p(x_j^i | y^i) p(y^i) \end{aligned}$$

<https://zhuanlan.zhihu.com/p/71960086>

4. 生成模型与判别模型

逻辑回归

朴素贝叶斯

目标函数

$$\prod_{i=1}^m p(y_i | x_i)$$

$$\prod_{i=1}^m p(x_i, y_i)$$

