

带你理解朴素贝叶斯分类算法

 忆臻 · 5 个月前

贝叶斯分类是一类分类算法的总称，这类算法均以贝叶斯定理为基础，故统称为贝叶斯分类。而朴素朴素贝叶斯分类是贝叶斯分类中最简单，也是常见的一种分类方法。这篇文章我

尽可能用直白的话语总结一下我们学习会上讲到的朴素贝叶斯分类算法，希望有利于他人理解。

分类问题综述

对于分类问题，其实谁都不会陌生，日常生活中我们每天都进行着分类过程。例如，当你看到一个人，你的脑子下意识判断他是学生还是社会上的人；你可能经常会走在路上对身旁的朋友说“这个人一看就很有钱、”之类的话，其实这就是一种分类操作。

既然是贝叶斯分类算法，那么**分类的数学描述**又是什么呢？

从数学角度来说，分类问题可做如下定义：已知集合 $C = y_1, y_2, \dots, y_n$ 和 $I = x_1, x_2, x_3, \dots, x_n$ ，确定映射规则 $y = f()$ ，使得任意 $x_i \in I$ 有且仅有一个 $y_i \in C$ ，使得 $y_i \in f(x_i)$ 成立。

其中C叫做类别集合，其中每一个元素是一个类别，而叫做项集合（**特征集合**），其中每一个元素是一个待分类项，叫做分类器。**分类算法的任务就是构造分类器f。**

分类算法的内容是要求给定特征，让我们得出类别，这也是所有分类问题的关键。那么如何由指定特征，得到我们最终的类别，也是我们下面要讲的，每一个不同的分类算法，对应着不同的核心思想。

本篇文章，我会用一个具体实例，对朴素贝叶斯算法几乎所有的重要知识点进行讲解。

朴素贝叶斯分类

那么既然是朴素贝叶斯**分类算法**，它的核心算法又是什么呢？

是下面这个贝叶斯公式：

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

换个表达形式就会明朗很多，如下：

$$p(\text{类别}|\text{特征}) = \frac{p(\text{特征}|\text{类别})p(\text{类别})}{p(\text{特征})}$$

我们最终求的 $p(\text{类别}|\text{特征})$ 即可！就相当于完成了我们的任务。

例题分析

下面我先给出例子问题。

给定数据如下：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

现在给我们的问题是，如果一对男女朋友，男生想女生求婚，男生的四个特点分别是不帅，性格不好，身高矮，不上进，请你判断一下女生是嫁还是不嫁？

这是一个典型的分类问题，转为数学问题就是比较 $p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进})$ 与 $p(\text{不嫁}|\text{不帅、性格不好、身高矮、不上进})$ 的概率，谁的概率大，我就能给出嫁或者不嫁的答案！

这里我们联系到朴素贝叶斯公式：

$$p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进}) = \frac{p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁}) * p(\text{嫁})}{p(\text{不帅、性格不好、身高矮、不上进})}$$

我们需要求 $p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进})$ ，这是我们不知道的，但是通过朴素贝叶

斯公式可以转化为好求的三个量， $p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁})$ 、 $p(\text{不帅、性格不好、身高矮、不上进})$ 、 $p(\text{嫁})$ （至于为什么能求，后面会讲，那么就太好了，将待求的量转化为其它可求的值，这就相当于解决了我们的问题！）

朴素贝叶斯算法的朴素一词解释

那么这三个量是如何求得？

是根据已知训练数据统计得来，下面详细给出该例子的求解过程。

回忆一下我们要求的公式如下：

$$p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进}) = \frac{p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁}) * p(\text{嫁})}{p(\text{不帅、性格不好、身高矮、不上进})}$$

那么我只要求得 $p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁})$ 、 $p(\text{不帅、性格不好、身高矮、不上进})$ 、 $p(\text{嫁})$ 即可，好的，下面我分别求出这几个概率，最后一比，就得到最终结果。

$p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁}) = p(\text{不帅}|\text{嫁}) * p(\text{性格不好}|\text{嫁}) * p(\text{身高矮}|\text{嫁}) * p(\text{不上进}|\text{嫁})$ ，那么我就要分别统计后面几个概率，也就得到了左边的概率！

等等，为什么这个成立呢？学过概率论的同学可能有感觉了，这个等式成立的条件需要特征之间相互独立吧！

对的！这也就是为什么朴素贝叶斯分类有朴素一词的来源，朴素贝叶斯算法是假设各个特征之间相互独立，那么这个等式就成立了！

但是为什么需要假设特征之间相互独立呢？

1、我们这么想，假如没有这个假设，那么我们对右边这些概率的估计其实是不可做的，这么说，我们这个例子有4个特征，其中帅包括{帅，不帅}，性格包括{不好，好，爆好}，身高包括{高，矮，中}，上进包括{不上进，上进}，那么四个特征的联合概率分布总共是4维空间，总个数为 $2*3*3*2=36$ 个。

24个，计算机扫描统计还可以，但是现实生活中，往往有非常多的特征，每一个特征的取值也是非常之多，那么通过统计来估计后面概率的值，变得几乎不可做，这也是为什么需要假设特征之间独立的原因。

2、假如我们没有假设特征之间相互独立，那么我们统计的时候，就需要在整个特征空间中

去找，比如统计 $p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁})$ ，

我们就需要在嫁的条件下，去找四种特征全满足分别是不帅，性格不好，身高矮，不上进的人的个数，这样的话，由于数据的稀疏性，很容易统计到0的情况。这样是不合适的。

根据上面两个原因，朴素贝叶斯法对条件概率分布做了条件独立性的假设，由于这是一个较强的假设，朴素贝叶斯也由此得名！这一假设使得朴素贝叶斯法变得简单，但有时会牺牲一定的分类准确率。

好的，上面我解释了为什么可以拆成分开连乘形式。那么下面我们就开始求解！

我们将上面公式整理一下如下：

$$p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进}) = \frac{p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁}) * p(\text{嫁})}{p(\text{不帅、性格不好、身高矮、不上进})}$$

$$= \frac{p(\text{不帅}|\text{嫁}) * p(\text{性格不好}|\text{嫁}) * p(\text{身高矮}|\text{嫁}) * p(\text{不上进}|\text{嫁}) * p(\text{嫁})}{p(\text{不帅}) * p(\text{性格不好}) * p(\text{身高矮}) * p(\text{不上进})}$$

下面我将一个一个的进行统计计算（在数据量很大的时候，根据中心极限定理，频率是等于概率的，这里只是一个例子，所以我就进行统计即可）。

$p(\text{嫁}) = ?$

首先我们整理训练数据中，嫁的样本数如下：

帅？	性格好？	身高？	上进？	嫁与否
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁

则 $p(\text{嫁}) = 6/12$ （总样本数）= $1/2$

$p(\text{不帅}|\text{嫁}) = ?$ 统计满足样本数如下：

帅？	性格好？	身高？	上进？	嫁与否
不帅	好	高	上进	嫁
不帅	好	中	上进	嫁
不帅	不好	高	上进	嫁

则 $p(\text{不帅}|\text{嫁}) = 3/6 = 1/2$

$p(\text{性格不好}|\text{嫁}) = ?$ 统计满足样本数如下：

帅？	性格好？	身高？	上进？	嫁与否
不帅	不好	高	上进	嫁

则 $p(\text{性格不好}|\text{嫁}) = 1/6$

$p(\text{矮}|\text{嫁}) = ?$ 统计满足样本数如下：

帅？	性格好？	身高？	上进？	嫁与否
帅	好	矮	上进	嫁

则 $p(\text{矮}|\text{嫁}) = 1/6$

$p(\text{不上进}|\text{嫁}) = ?$ 统计满足样本数如下：

帅？	性格好？	身高？	上进？	嫁与否
帅	好	高	不上进	嫁

则 $p(\text{不上进}|\text{嫁}) = 1/6$

下面开始求分母， $p(\text{不帅})$ ， $p(\text{性格不好})$ ， $p(\text{矮})$ ， $p(\text{不上进})$

统计样本如下：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

不帅统计如上红色所示，占4个，那么 $p(\text{不帅}) = 4/12 = 1/3$

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

性格不好统计如上红色所示，占4个，那么 $p(\text{性格不好}) = 4/12 = 1/3$

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

身高矮统计如上红色所示，占7个，那么 $p(\text{身高矮}) = 7/12$

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

不上进统计如上红色所示，占4个，那么 $p(\text{不上进}) = 4/12 = 1/3$

到这里，要求 $p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁})$ 的所需项全部求出来了，下面我带进去即可，

$$\begin{aligned}
 p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进}) &= \frac{p(\text{不帅、性格不好、身高矮、不上进}|\text{嫁}) * p(\text{嫁})}{p(\text{不帅、性格不好、身高矮、不上进})} \\
 &= \frac{p(\text{不帅}|\text{嫁}) * p(\text{性格不好}|\text{嫁}) * p(\text{身高矮}|\text{嫁}) * p(\text{不上进}|\text{嫁}) * p(\text{嫁})}{p(\text{不帅}) * p(\text{性格不好}) * p(\text{身高矮}) * p(\text{不上进})}
 \end{aligned}$$

$$= (1/2 * 1/6 * 1/6 * 1/6 * 1/2) / (1/3 * 1/3 * 7/12 * 1/3)$$

下面我们根据同样的方法来求 $p(\text{不嫁}|\text{不帅, 性格不好, 身高矮, 不上进})$ ，完全一样的做法，为了方便理解，我这里也走一遍帮助理解。首先公式如下：

$$p(\text{不嫁}|\text{不帅, 性格不好, 身高矮, 不上进}) = \frac{p(\text{不帅, 性格不好, 身高矮, 不上进}|\text{不嫁}) * p(\text{不嫁})}{p(\text{不帅, 性格不好, 身高矮, 不上进})}$$

$$= \frac{p(\text{不帅}|\text{不嫁}) * p(\text{性格不好}|\text{不嫁}) * p(\text{身高矮}|\text{不嫁}) * p(\text{不上进}|\text{不嫁}) * p(\text{不嫁})}{p(\text{不帅}) * p(\text{性格不好}) * p(\text{身高矮}) * p(\text{不上进})}$$

下面我也一个一个来进行统计计算，这里与上面公式中，分母是一样的，于是我们分母不需要重新统计计算！

$p(\text{不嫁}) = ?$ 根据统计计算如下（**红色为满足条件**）：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

则 $p(\text{不嫁}) = 6/12 = 1/2$

$p(\text{不帅}|\text{不嫁}) = ?$ 统计满足条件的样本如下（**红色为满足条件**）：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

则 $p(\text{不帅}|\text{不嫁}) = 1/6$

$p(\text{性格不好}|\text{不嫁}) = ?$ 据统计计算如下（红色为满足条件）：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

则 $p(\text{性格不好}|\text{不嫁}) = 3/6 = 1/2$

$p(\text{矮}|\text{不嫁}) = ?$ 据统计计算如下（红色为满足条件）：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

则 $p(\text{矮}|\text{不嫁}) = 6/6 = 1$

$p(\text{不上进}|\text{不嫁}) = ?$ 据统计计算如下（红色为满足条件）：

帅？	性格好？	身高？	上进？	嫁与否
帅	不好	矮	不上进	不嫁
不帅	好	矮	上进	不嫁
帅	好	矮	上进	嫁
不帅	好	高	上进	嫁
帅	不好	矮	上进	不嫁
帅	不好	矮	上进	不嫁
帅	好	高	不上进	嫁
不帅	好	中	上进	嫁
帅	好	中	上进	嫁
不帅	不好	高	上进	嫁
帅	好	矮	不上进	不嫁
帅	好	矮	不上进	不嫁

则 $p(\text{不上进}|\text{不嫁}) = 3/6 = 1/2$

那么根据公式：

$$\begin{aligned}
 p(\text{不嫁}|\text{不帅、性格不好、身高矮、不上进}) &= \frac{p(\text{不帅、性格不好、身高矮、不上进}|\text{不嫁}) * p(\text{不嫁})}{p(\text{不帅、性格不好、身高矮、不上进})} \\
 &= \frac{p(\text{不帅}|\text{不嫁}) * p(\text{性格不好}|\text{不嫁}) * p(\text{身高矮}|\text{不嫁}) * p(\text{不上进}|\text{不嫁}) * p(\text{不嫁})}{p(\text{不帅}) * p(\text{性格不好}) * p(\text{身高矮}) * p(\text{不上进})}
 \end{aligned}$$

$$p(\text{不嫁}|\text{不帅、性格不好、身高矮、不上进}) = ((1/6 * 1/2 * 1 * 1/2) * 1/2) / (1/3 * 1/3 * 7/12 * 1/3)$$

很显然 $(1/6 * 1/2 * 1 * 1/2) > (1/2 * 1/6 * 1/6 * 1/6 * 1/2)$

于是有 $p(\text{不嫁}|\text{不帅、性格不好、身高矮、不上进}) > p(\text{嫁}|\text{不帅、性格不好、身高矮、不上进})$

所以我们根据朴素贝叶斯算法可以给这个女生答案，是不嫁！！！！

朴素贝叶斯分类的优缺点

优点：

(1) 算法逻辑简单,易于实现

(2) 分类过程中时空开销小

缺点：

理论上，朴素贝叶斯模型与其他分类方法相比具有最小的误差率。但是实际上并非总是如此，这是因为朴素贝叶斯模型假设属性之间相互独立，这个假设在实际应用中往往是不成立的，在属性个数比较多或者属性之间相关性较大时，分类效果不好。

而在属性相关性较小时，朴素贝叶斯性能最为良好。对于这一点，有半朴素贝叶斯之类的算法通过考虑部分关联性适度改进。

整个例子详细的讲解了朴素贝叶斯算法的分类过程，希望对大家的理解有帮助~

参考：李航博士《统计学习方法》

赞赏

算法杂货铺--分类算法之朴素贝叶斯分类(Naive Bayesian classification)

2 人赞赏

封面图来自于：算法杂货铺--分类算法之朴素贝叶斯分类(Naive Bayesian classification)



致谢：德川，皓宇，继豪，施琦

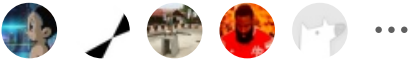
机器学习

深度学习 (Deep Learning)

自然语言处理

🔖 收藏 📄 分享 🗑 举报

👍 77



13 条评论

写下你的评论...



SJTU电院靠谱的郑建国

拉普拉斯平滑没讲到？

5 个月前



忆臻 (作者) 回复 SJTU电院靠谱的郑建国

[查看对话](#)

我打算下篇文章再用这个例子讲~

5 个月前



何晔

如果某一分类没有数据怎么办？那他的概率就是零了。

5 个月前



忆臻 (作者) 回复 何晔

[查看对话](#)

你好，下篇文章我会讲平滑~

5 个月前



张振

看懂了

5 个月前



三一

非常漂亮的讲述，物质和精神都赞一个！

5 个月前



忆臻 (作者) 回复 三一

查看对话

很开心对你有帮助，谢谢~

5 个月前

1 赞



Eric 回复 SJTU电院靠谱的郑建国

查看对话

好巧，竟然在这儿遇见你

5 个月前



如风

讲的很好，有点概率论基础的都能看的懂 学习了 谢谢

5 个月前



忆臻 (作者) 回复 如风

查看对话

很开心对你有帮助~

5 个月前

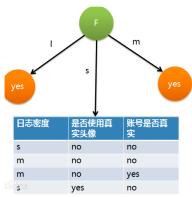
下一页

文章被以下专栏收录



进入专栏

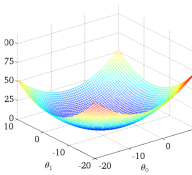
推荐阅读



理解朴素贝叶斯分类的拉普拉斯平滑

我们在上篇文章中带你理解朴素贝叶斯分类算法 - 知乎专栏已经根据朴素贝叶斯算法给出了当一... 查看全文 >

忆臻 · 5 个月前 · 发表于 机器学习算法与自然语言处理



通俗讲解平方损失函数平方形式的数学解释？

还是郭江师兄在第一次机器学习会上，给我们讲解的相关线性模型知识，这里进行总结平方损失函... 查看全文 >

忆臻 · 5 个月前 · 发表于 机器学习算法与自然语言处理



雀巢收购蓝瓶子，那么精品速溶还会远吗？

上了日报和编辑推荐，那么就原文来做个说明吧：1：专栏没有任何广告植入。2：至于国内被投资... 查看全文 >

Ray · 4 天前 · 编辑精选



你算过你一年究竟花了多少时间去排队吃饭吗？

从周五晚上的5点到9点，新一酱和朋友们花了4个多小时在上海人民广场的网红餐厅哥老官金陵东... 查看全文 >

新一酱 · 15 天前 · 编辑精选 · 发表于 新一线城市研究所