

2_问题的定义

1. 数据回顾

浏览了一下数据，数据中自变量 X 包含了一个小短句，如：“陈丽婷在踢足球”。因变量分别为：接受度，喜爱度，交友意愿。由于这几个变量高度相关，我们可以暂时看成一个变量：受欢迎程度 $Y = (\text{接受度} + \text{喜爱度} + \text{交友意愿})/3$ 。还有一些条件变量，不过我们可以暂时忽略。

序号	1.陈立婷在踢足球。			2.胡玉蝶在搬重物。			3.杨华峰在看娱乐节目。			4.万丽萍在看英雄、武侠、动作类影		
	接受度	喜爱度	交友意愿	接受度	喜爱度	交友意愿	接受度	喜爱度	交友意愿	接受度	喜爱度	交友意愿
1	6	4	5	7	5	5	5	3	3	4	4	4
2	4	2	1	4	3	3	4	4	4	6	4	4
3	6	4	5	6	4	5	6	4	5	6	3	5
4	7	7	7	7	7	7	7	4	7	7	7	7
5	4	4	3	4	4	3	6	4	4	6	6	5
6	7	5	4	7	6	4	7	4	4	7	4	4
7	5	3	4	6	4	4	6	4	4	2	2	2
8	5	3	3	5	4	4	6	5	4	6	5	6
9	7	4	4	4	4	4	4	4	4	6	4	4
10	5	5	5	4	4	5	4	1	2	3	3	2
11	4	4	4	4	4	4	3	3	3	4	4	4
12	2	2	2	2	2	2	2	2	2	2	2	2
13	1	1	1	4	4	4	1	1	1	1	1	1
14	7	4	4	4	4	4	6	4	4	7	4	5
15	7	5	5	7	6	6	7	5	5	7	5	4
16	6	6	7	7	7	7	7	7	7	6	6	7
17	7	6	6	7	7	7	7	6	7	7	7	7
18	3	3	2	2	2	2	4	3	3	3	3	4
19	4	3	3	5	4	4	1	5	2	1	1	1
20	4	2	2	7	5	4	7	5	4	5	4	4
21	5	5	4	6	6	6	4	5	5	7	7	7

图1: 部分数据

基于这个数据，我们需要对自变量和因变量进行建模： $p(Y|X)$ 。但是直接做相关就有点没有技术含量了，我们需要做一些比较复杂的贡献。比如我们可以猜测，从人们看到输入的句子，到他们写下受欢迎程度的评分的过程中发生了什么。



图2: 在这个过程中人们想了什么呢?

如果从概率的角度来看，我们就是把问题看成一个贝叶斯推测问题。而解贝叶斯推断的思路为

1. 建立生成模型
2. 用一些算法做推断

2.数学模型

2.1 变量的定义

- 定义句中主语为随机变量 N 取值空间为 $\mathcal{N} = (\text{"陈立婷"}, \text{"胡玉蝶"}, \text{"杨华峰"}, \dots)$ 。
- 定义句子谓语为随机变量 V 取值空间为 $\mathcal{V} = (\text{"踢足球"}, \text{"搬重物"}, \text{"看娱乐节目"}, \dots)$ 。
- 定义性别为随机变量 G ，取值空间为 $\mathcal{G} = (\text{"男"}, \text{"女"})$ 。之后为了简便，我们定义男为0，女为1、
- 定义某人是男（是女）的概率为随机变量 Q ，取值空间为 $\mathcal{Q} = [0 - 1]$ 。
- 定义*喜欢程度为随机变量 Y ，取值空间为 $\mathcal{Y} = (\text{"1"}, \text{"2"}, \text{"3"}, \text{"4"}, \text{"5"}, \text{"6"}, \text{"7"}, \text{"8"}, \text{"9"})$ 。
- 定义我们描述的思维模型为 M ，输入为 (N, V) ，输出为 Y 。

2.2 模型结构

图2为各种概念的生成模型。其中蓝色为观测变量，即我们从数据中收集到的。红色概念为隐变量，即我们脑中的概念。

这个模型的意义是，好比游戏中的捏人。我们脑中本来就有一个人的模板。这个模板说了一个人大概是男的(女的)的可能性是 Q 。根据这个分布,我们生成对我们要捏的个体的性别 G 。根据这个性别我们对这个角色进行命名 N 以及希望他干什么 V .另一方面，我们对模板 Q 的分布进行偏好的推测。（需要组织一下语言，不知道怎么从数学语言翻译过来）。

如果用这个模型描述一个人从词语到喜欢程度，那其实我们在反转这个生成模型来进行猜测。对这个推断过程，我们可以有以理解。当给予对一个人的描述，“名字”+“行为”，我们推测这个描述背后人的性别是什么，而通过推测的难易程度，我们产生对这个描述的喜好程度

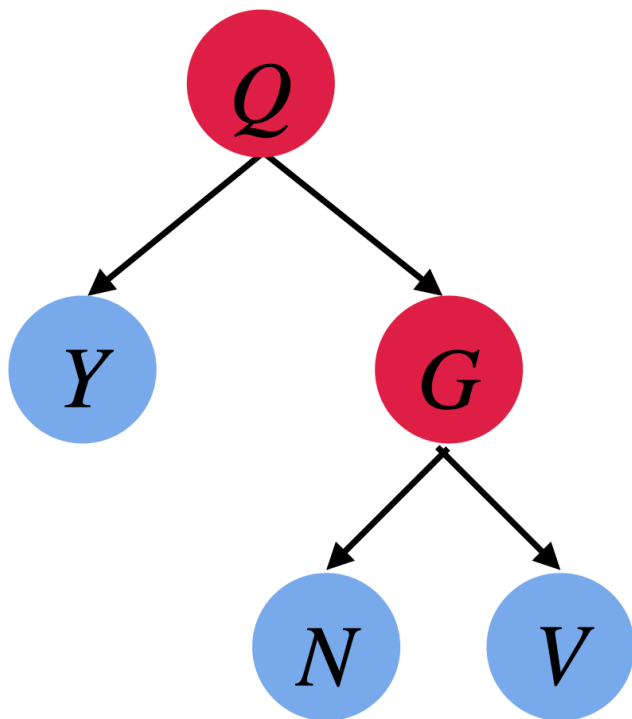


图3: 受欢迎程度的生成模型

2.3 模型参数

1. $p(G|Q)$:
 - $p(G = 0|Q) = \theta$ (可学习)
 - $p(G = 1|Q) = 1 - \theta$
2. $p(N|G)$ 需要实验收集
3. $p(V|G)$ 需要模型收集
4. $p(Y|Q)$ 这个是我们的假设。我们的假设其服从Dirichlet分布。其实就是线性相关，只是从连续变量映射到离散变量的方法，不需要学习参数。

2.4 具体模型的使用

这个模型描述的过程是：一个被试给定某人的描述（名字+动作），基于这个描述推断这个人的受欢迎程度。

抽象的描述是: $p(Y|N = n, V = v)$

具体例子来说是，对图1中第一行的数据。我们需要对我们的模型输入（"陈立婷"，"足球"）输出（7）作为受喜欢程度的平均值。

3.其他事项

3.1 可能需要补充实验

”假如是一个男生，他可能取以下哪些名字，有多少把握。“

3.2 做实验之前还需要做什么？

- 对数据可能的结果进行模拟，以及对模型进行调试。需要成功率比较高才进行数据收集。（10/30/2020 前）
- 实验控制过程。
- 讨论这个模型的意义在于什么？
- 怎么加入词语的效价？