



北京航空航天大学
BEIHANG UNIVERSITY

深度学习与自然语言处理 第二次大作业

基于 EM 算法的参数估计

院（系）名称	自动化科学与电气工程学院
--------	--------------

学 生 学 号	ZY2103810
---------	-----------

学 生 姓 名	游虎杰
---------	-----

2022 年 4 月

一、问题描述

现有一堆硬币，包括三种硬币，每种硬币的比例分别为 s_1 , s_2 , s_3 ，对应的抛硬币正面朝上概率分别为 p , q , r 。现在进行 N 次独立实验，每次实验随机从硬币堆中抽取某个硬币并抛掷，记录该硬币是正面朝上（结果记为 1）还是反面朝上（结果记为 0），然后放回硬币， N 次实验后，根据所记录观测结果，利用 EM 算法估计参数 s_1 , s_2 , p , q , r 。

二、问题表达

算法用到的各项符号与对应参数如表 1 所示：

表 1 相关符号及对应参数

符号	对应参数
$x^{(j)}$	模拟投掷硬币时生成的第 j 个随机数
$y^{(j)}$	第 j 次实验结果
$s1_i^{(j)}$	在第 i 轮迭代的 E 步中，根据第 j 次实验结果，对参数 s_1 的期望值
$s2_i^{(j)}$	在第 i 轮迭代中 E 步中，根据第 j 次实验结果，对参数 s_2 的期望值
$s1_i$	在第 i 轮迭代的 M 步中，对参数 s_1 的估计值
$s2_i$	在第 i 轮迭代的 M 步中，对参数 s_2 的估计值
p_i	在第 i 轮迭代的 M 步中，对参数 p 的估计值
q_i	在第 i 轮迭代的 M 步中，对参数 q 的估计值
r_i	在第 i 轮迭代的 M 步中，对参数 r 的估计值
$\widehat{s1}_i$	EM 算法收敛时，对 s_1 的最终估计值
$\widehat{s2}_i$	EM 算法收敛时，对 s_2 的最终估计值
\hat{p}	EM 算法收敛时，对 p 的最终估计值
\hat{q}	EM 算法收敛时，对 q 的最终估计值

\hat{r}	EM 算法收敛时，对 r 的最终估计值
m	模拟抛掷硬币实验中硬币正面出现的频率
\hat{m}	按照参数估计值计算出的硬币证明出现概率

三、具体算法实现

1. 模拟投掷结果生成

经计算可知，每次投掷出证明的概率应当为 $s1 * p + s2 * q + s3 * r$ ，因此可通过随机生成 N 个随机数来生成实验结果。若生成的第 j 个随机数 $x^{(j)} < s1 * p + s2 * q + s3 * r$ ，则判断第 j 次实验硬币投掷出证明，否则判断投掷出反面。利用列表记录 N 次模拟投掷的结果。

2. 基于 EM 算法的参数估计

2.1 对各参数进行初始化设置，并令 $i=0$ 。

2.2 期望值求解（E 步）

根据每次实验结果，计算 $s1$ ， $s2$ 的对应期望值，具体公式如下：

$$s1_i^{(j)} = \frac{s1_i(p_i)^{y^{(j)}}(1-p_i)^{(1-y^{(j)})}}{s1_i(p_i)^{y^{(j)}}(1-p_i)^{(1-y^{(j)})} + s2_i(q_i)^{y^{(j)}}(1-q_i)^{(1-y^{(j)})} + (1-s1_i-s2_i)(r_i)^{y^{(j)}}(1-r_i)^{(1-y^{(j)})}}$$

$$s2_i^{(j)} = \frac{s2_i(q_i)^{y^{(j)}}(1-q_i)^{(1-y^{(j)})}}{s1_i(p_i)^{y^{(j)}}(1-p_i)^{(1-y^{(j)})} + s2_i(q_i)^{y^{(j)}}(1-q_i)^{(1-y^{(j)})} + (1-s1_i-s2_i)(r_i)^{y^{(j)}}(1-r_i)^{(1-y^{(j)})}}$$

2.3 似然函数最大化

更新各参数的估计值，使得似然函数趋于局部最大，具体公式如下：

$$s1_{i+1} = \frac{1}{N} \sum_{j=1}^N s1_i^{(j)}$$

$$s2_{i+1} = \frac{1}{N} \sum_{j=1}^N s2_i^{(j)}$$

$$p_{i+1} = \frac{\sum_{j=1}^N s1_i^{(j)} y^{(j)}}{\sum_{j=1}^N s1_i^{(j)}}$$

$$q_{i+1} = \frac{\sum_{j=1}^N s2_i^{(j)} y^{(j)}}{\sum_{j=1}^N s2_i^{(j)}}$$

$$r_{i+1} = \frac{\sum_{j=1}^N (1 - s1_i^{(j)} - s2_i^{(j)}) y^{(j)}}{\sum_{j=1}^N (1 - s1_i^{(j)} - s2_i^{(j)})}$$

$$i = i + 1$$

重复步骤 2.2~2.3，直至各参数的绝对误差和小于 0.001，此时 $\hat{p} = p_{i+1}$ ， $\hat{q} = q_{i+1}$ ， $\hat{r} = r_{i+1}$

四、运行结果

进行 10 组基于 EM 算法的参数估计实验，每组实验抛掷硬币的次数 N 均为 1000，结果如表 2 所示：

表 2 实验结果

序号	s1	$\widehat{s1}_i$	s2	$\widehat{s2}_i$	p	\hat{p}	q	\hat{q}	r	\hat{r}	m	\hat{m}
1	0.30	0.30	0.50	0.49	0.40	0.34	0.30	0.34	0.70	0.70	0.416	0.416
2	0.20	0.28	0.20	0.47	0.40	0.50	0.30	0.50	0.70	0.80	0.558	0.558
3	0.20	0.31	0.50	0.5	0.40	0.23	0.20	0.23	0.30	0.56	0.287	0.287
4	0.20	0.32	0.30	0.54	0.10	0.09	0.20	0.09	0.10	0.30	0.121	0.121
5	0.50	0.31	0.20	0.52	0.30	0.17	0.30	0.17	0.10	0.47	0.216	0.216
6	0.10	0.22	0.20	0.55	0.30	0.64	0.80	0.86	0.90	0.88	0.815	0.815
7	0.40	0.36	0.50	0.36	0.30	0.14	0.10	0.36	0.90	0.41	0.288	0.288

8	0.80	0.24	0.10	0.54	0.90	0.53	0.10	0.80	0.20	0.83	0.743	0.743
9	0.10	0.22	0.90	0.55	0.10	0.67	0.90	0.88	0.10	0.90	0.873	0.873
10	0.20	0.30	0.60	0.50	0.40	0.33	0.60	0.63	0.60	0.68	0.553	0.553

由实验结果可知，每组实验完成后，根据 EM 算法估计的参数计算硬币正面出现概率 \hat{m} ，都和实验过程中硬币抛掷出正面的频率 m 相等，但是估计出的参数与真实值往往相差很大，这是因为 EM 算法在迭代过程中会收敛到初始值附近的局部最优值。事实上，最终收敛结果一般在初始值附近。并且，在实验 1~5 中， p 与 q 的估计值完全一致，这是因为在实验 1~5 中， p 与 q 的初始值完全一致。以上实验结果说明，EM 算法对参数的估计值，能够准确计算出实验结果，但是未必与真实值相同，并且很大程度上受初始值的影响，这些在选用 EM 算法进行参数估计时都要充分考虑。

五、个人总结和体会

通过这次作业，我对 EM 算法的理解有了进一步的加深，对 EM 算法的计算过程有了更为熟练的掌握。同时，在编写代码的过程中，我对 python 的应用水平有了进一步提升。

六、作业代码

https://github.com/youlll/DP_NLP2.git