

# Parameter estimation for text analysis - 笔记

冯柏淋

2018 年 10 月 23 日

## 目录

<b>1 引入</b>	<b>2</b>
<b>2 基本问题与贝叶斯公式</b>	<b>2</b>
2.1 两类基本问题 . . . . .	2
2.2 贝叶斯公式 . . . . .	2
2.3 记号 . . . . .	2
<b>3 极大似然估计 MLE</b>	<b>2</b>
3.1 参数估计 . . . . .	2
3.2 预测 . . . . .	3
3.3 一个例子 . . . . .	3
<b>4 最大后验概率 MAP</b>	<b>3</b>
4.1 参数估计 . . . . .	3
4.2 预测 . . . . .	4
4.3 一个例子 . . . . .	4
<b>5 贝叶斯估计</b>	<b>4</b>
5.1 参数估计 . . . . .	4
5.2 预测 . . . . .	5
5.3 一个例子 . . . . .	5

# 1 引入

此笔记是本人在阅读论文 Gregor Heinrich: *Parameter estimation for text analysis* 的过程的一点笔记，大多是对于原文的翻译和摘录。另有一部分自己的理解，因而可能存在错误。

## 2 基本问题与贝叶斯公式

### 2.1 两类基本问题

- **估计** 估计分布的参数以解释观测值
- **预测/回归** 基于以往观测值，预测新的观测值的概率

### 2.2 贝叶斯公式

$$p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) \cdot p(\theta)}{p(\mathcal{X})} \quad (1)$$

对以上公式中的每一部分，我们有如下名称：

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}} \quad (2)$$

### 2.3 记号

数据集：

$$\mathcal{X} \triangleq \{x\}_i^{|\mathcal{X}|}$$

## 3 极大似然估计 MLE

### 3.1 参数估计

极大似然函数：

$$L(\theta; \mathcal{X}) \triangleq p(\mathcal{X} | \theta) = \prod_{x \in \mathcal{X}} p(x | \theta) \quad (3)$$

极大似然估计 (MLE)：

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \mathcal{L}(\theta; \mathcal{X}) = \arg \max_{\theta} \sum_{x \in \mathcal{X}} \log p(x | \theta) \quad (4)$$

### 3.2 预测

对于新的观测值的概率估计过程如下：<sup>1</sup>

$$p(\tilde{x} | \mathcal{X}) = \int_{\theta \in \Theta} p(\tilde{x} | \theta) \cdot p(\theta | \mathcal{X}) d\theta \quad (5)$$

$$\approx \int_{\theta \in \Theta} p(\tilde{x} | \hat{\theta}_{\text{MLE}}) \cdot p(\theta | \mathcal{X}) d\theta \quad (6)$$

$$= p(\tilde{x} | \hat{\theta}_{\text{MLE}}). \quad (7)$$

### 3.3 一个例子

暂时略。

## 4 最大后验概率 MAP

### 4.1 参数估计

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta | \mathcal{X}) \quad (8)$$

$$= \arg \max_{\theta} \frac{p(\mathcal{X} | \theta) \cdot p(\theta)}{p(\mathcal{X})} \quad (9)$$

$$= \arg \max_{\theta} p(\mathcal{X} | \theta) \cdot p(\theta) \quad (10)$$

$$= \arg \max_{\theta} \{\mathcal{L}(\theta | \mathcal{X}) + \log p(\theta)\} \quad (11)$$

$$= \arg \max_{\theta} \left\{ \sum_{x \in \mathcal{X}} \log p(x | \theta) + \log p(\theta) \right\} \quad (12)$$

MAP 相比于 MLE 而言，添加了关于先验分布的额外信息  $p(\theta)$ 。实践中，可以通过假定一些简单的先验分布  $p(\theta)$  来防止过拟合。<sup>2</sup>

<sup>1</sup> 此处， $p(\tilde{x} | \hat{\theta}_{\text{MLE}})$  视为常数，并且  $\int_{\theta \in \Theta} p(\theta | \mathcal{X}) d\theta = 1$

<sup>2</sup> Occam's razor 原理

## 4.2 预测

$$p(\tilde{x} | \mathcal{X}) \approx \int_{\theta \in \Theta} p(\tilde{x} | \hat{\theta}_{\text{MAP}}) \cdot p(\theta | \mathcal{X}) d\theta \quad (13)$$

$$= p(\tilde{x} | \hat{\theta}_{\text{MAP}}). \quad (14)$$

## 4.3 一个例子

暂时略。

# 5 贝叶斯估计

## 5.1 参数估计

贝叶斯估计是对 MAP 的一个拓展。MAP 对于参数只产生一个确定的估计值  $\hat{\theta}_{\text{MAP}}$ ，而贝叶斯估计则产生对于参数的一个分布。从而，贝叶斯分布则提供了关于参数的额外信息，比如期望和方差。

主要步骤为计算后验分布：

$$p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) \cdot p(\theta)}{p(\mathcal{X})} \quad (15)$$

在前述的 MAP 方法中，只需要最大化此式，则无需求出分母  $p(\mathcal{X})$ ，而在此处则需要计算，且：

$$p(\mathcal{X}) = \int_{\theta \in \Theta} p(\mathcal{X} | \theta) \cdot p(\theta) d\theta \quad (16)$$

因为没最大化参数的过程，参数估计的过程也就不是一个具体的值，而是求取参数的分布：

$$p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) \cdot p(\theta)}{p(\mathcal{X})} \quad (17)$$

$$= \frac{p(\mathcal{X} | \theta) \cdot p(\theta)}{\int_{\theta \in \Theta} p(\mathcal{X} | \theta) \cdot p(\theta) d\theta} \quad (18)$$

## 5.2 预测

$$p(\tilde{x} \mid \mathcal{X}) = \int_{\theta \in \Theta} p(\tilde{x} \mid \theta) \cdot p(\theta \mid \mathcal{X}) \, \mathrm{d}\theta \quad (19)$$

$$= \int_{\theta \in \Theta} p(\tilde{x} \mid \theta) \cdot \frac{p(\mathcal{X} \mid \theta) \cdot p(\theta)}{p(\mathcal{X})} \, \mathrm{d}\theta. \quad (20)$$

## 5.3 一个例子

暂时略。