

Author: Liu Jian

Time: 2020-06-06

## 机器学习10-最大熵原理与最大熵模型

### 1 最大熵原理

#### 1.1 问题的构建

#### 1.2 对偶问题

### 2 最大熵模型

#### 2.1 模型描述

#### 2.2 模型学习--极大似然估计

### 3 数值优化算法

#### 3.1 改进的迭代尺度法 (improved iterative scaling, IIS)

#### 3.2 拟牛顿法

# 机器学习10-最大熵原理与最大熵模型

## 1 最大熵原理

### 1.1 问题的构建

**最大熵原理**：在满足约束条件的模型集合中选取熵最大的模型 (可以这样理解，熵越大，编码也就越长，模型的表示能力也就越强)。下面我们从最大熵原理出发，学习某个概率分布  $P(y|\mathbf{x})$  (即  $P_{model}(y|\mathbf{x})$ )，只不过这一模型的形式我们没有指定，而是根据最大熵原理推得)， $\mathbf{x} \in \mathbb{X} \subseteq \mathbb{R}^d, y \in \mathbb{Y} \subseteq \mathbb{R}$ 。样本数据  $D = \{(\mathbf{x}_i, y_i)\}_1^N$ ，通过统计频率我们可以得到经验分布  $\tilde{P}(\mathbf{x}, y)$  (即  $P_{data}(\mathbf{x}, y)$ )，当然也可以得到  $P_{data}(\mathbf{x})$ 、 $P_{data}(y|\mathbf{x})$  等)。

首先，我们给出学习的约束条件。给定特征函数  $f(\mathbf{x}, y)$ ，基于经验分布我们可以计算特征函数的期望：

$$\mathbb{E}_{\tilde{P}}[f(\mathbf{x}, y)] = \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) f(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) \tilde{P}(y|\mathbf{x}) f(\mathbf{x}, y) d\mathbf{x} dy$$

我们将上式中的  $\tilde{P}(y|\mathbf{x})$  替换为我们所要求的分布  $P(y|\mathbf{x})$ ：

$$\int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) f(\mathbf{x}, y) d\mathbf{x} dy \triangleq \mathbb{E}_P[f(\mathbf{x}, y)]$$

自然地，我们希望替换前后二式相等，于是我们令  $\mathbb{E}_{\tilde{P}}[f(\mathbf{x}, y)] = \mathbb{E}_P[f(\mathbf{x}, y)]$ ，这也就得到了最大熵模型的约束条件。

接下来，我们给出优化的目标函数。我们可以选择最大化如下的条件熵：

$$H(y|\mathbf{x}) = - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) \ln P(y|\mathbf{x}) d\mathbf{x} dy$$

综上，我们要求解如下的约束优化问题：

$$\begin{aligned} \max_P H(y|\mathbf{x}) &\Leftrightarrow \min_P \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) \ln P(y|\mathbf{x}) d\mathbf{x} dy \\ \text{s. t. } &\mathbb{E}_{\tilde{P}}[f_i(\mathbf{x}, y)] = \mathbb{E}_P[f_i(\mathbf{x}, y)], \quad i = 1, \dots, n \\ &\int_{\mathbb{Y}} P(y|\mathbf{x}) dy = 1 \end{aligned}$$

其中,  $\tilde{P}$  由训练样本  $D$  给出,  $f_i, i = 1, \dots, n$  为给定的特征函数。

最后, 我们指出:

1. 若我们要求解的对象是  $P(\mathbf{x}, y)$ , 在构建约束时, 我们令  $\mathbb{E}_{\tilde{P}}[f(\mathbf{x}, y)] = \mathbb{E}_P[f(\mathbf{x}, y)]$ , 可得:

$$\int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) f(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbb{X} \times \mathbb{Y}} P(\mathbf{x}, y) f(\mathbf{x}, y) d\mathbf{x} dy$$

但这里我们要求解的对象是  $P(y|\mathbf{x})$ , 我们令:

$$\int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) f(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) f(\mathbf{x}, y) d\mathbf{x} dy$$

即从等号左边到右边, 我们只替换其中的  $\tilde{P}(y|\mathbf{x})$  为  $P(y|\mathbf{x})$ , 即对于等式右边, 除了我们要求解的对象--条件概率  $y|\mathbf{x}$  取  $P(y|\mathbf{x})$  (即  $P_{model}(y|\mathbf{x})$ ) 外, 其他概率均取经验分布。

2. 类似地, 这里我们选择最大化条件概率  $P(y|\mathbf{x})$  的交叉熵:

$$- \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) \ln P(y|\mathbf{x}) d\mathbf{x} dy$$

其中, 除了我们要求解的对象--条件概率  $y|\mathbf{x}$  取  $P(y|\mathbf{x})$  (即  $P_{model}(y|\mathbf{x})$ ) 外, 其他概率均取经验分布。若我们要求解的对象是  $P(\mathbf{x}, y)$ , 则我们最大化信息熵:

$$- \int_{\mathbb{X} \times \mathbb{Y}} P(\mathbf{x}, y) \ln P(\mathbf{x}, y) d\mathbf{x} dy$$

## 1.2 对偶问题

原问题:

$$\begin{aligned} \min_P \quad & \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) \ln P(y|\mathbf{x}) d\mathbf{x} dy \\ \text{s. t.} \quad & \mathbb{E}_{\tilde{P}}[f_i(\mathbf{x}, y)] = \mathbb{E}_P[f_i(\mathbf{x}, y)], \quad i = 1, \dots, n \\ & \int_{\mathbb{Y}} P(y|\mathbf{x}) dy = 1 \end{aligned}$$

上述是一个泛函问题。

构造拉格朗日函数如下:

$$\begin{aligned} L(P, \mu, \lambda) = & \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) \ln P(y|\mathbf{x}) d\mathbf{x} dy + \mu(\mathbf{x}) \left( 1 - \int_{\mathbb{Y}} P(y|\mathbf{x}) dy \right) \\ & + \sum_{i=1}^n \lambda_i \int_{\mathbb{X} \times \mathbb{Y}} \left( \tilde{P}(\mathbf{x}, y) - \tilde{P}(\mathbf{x}) P(y|\mathbf{x}) \right) f_i(\mathbf{x}, y) d\mathbf{x} dy \end{aligned}$$

其中,  $\mu(\mathbf{x}), \lambda$  为拉格朗日乘子, 约束条件  $\int_{\mathbb{Y}} P(y|\mathbf{x}) dy = 1$  是一个关于自变量  $\mathbf{x}$  的函数约束, 因此其拉格朗日乘子是一个函数  $\mu(\mathbf{x})$ 。可以这样理解, 约束条件  $\int_{\mathbb{Y}} P(y|\mathbf{x}) dy = 1$  在空间  $\mathbb{X}$  中的每个点都对应了一个约束, 也就对应了一个拉格朗日乘子, 这些拉格朗日乘子可统一地用  $\mu(\mathbf{x})$  表示。拉格朗日函数  $L(P, \lambda)$  是关于  $P$  的凸函数, 因此主问题  $\min_P \max_{\mu, \lambda} L(P, \mu, \lambda)$  与对偶问题  $\max_{\mu, \lambda} \min_P L(P, \mu, \lambda)$  等价, 我们通过求解对偶问题来得到原问题的解。

首先是  $\min_P L(P, \mu, \lambda)$ , 我们可以视  $P(\cdot|\star)$  在输入输出空间  $\mathbb{X} \times \mathbb{Y}$  每个点上的值为一个参数来进行优化, 求导并令导数为 0 (KKT 条件之 stationarity):

$$\begin{aligned}\frac{\partial L(P, \mu, \lambda)}{\partial P(y|\mathbf{x})} &= 0 \\ \Updownarrow \\ \tilde{P}(\mathbf{x})(\ln P(y|\mathbf{x}) + 1) - \mu(\mathbf{x}) - \tilde{P}(\mathbf{x}) \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) &= 0 \\ P(y|\mathbf{x}) &= \frac{\exp\left(\sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y)\right)}{\exp(1 - \mu(\mathbf{x})/\tilde{P}(\mathbf{x}))} \triangleq P_{\mu, \lambda}(y|\mathbf{x})\end{aligned}$$

其中,  $f_i(\mathbf{x}, y), \tilde{P}(\mathbf{x})$  是已知的量。因此, 对偶问题:

$$\max_{\mu, \lambda} \min_P L(P, \mu, \lambda) = \max_{\mu, \lambda} L(P_{\mu, \lambda}, \mu, \lambda)$$

事实上, 我们可以作进一步化简, 消去变量  $\mu(\mathbf{x})$ , 只保留变量  $\lambda$ 。因为由  $\int_{\mathbb{Y}} P(y|\mathbf{x}) dy = 1$  (KKT 条件之primal feasibility),  $\mu(\mathbf{x}), \lambda$  之间存在如下的约束条件:

$$\exp(1 - \mu(\mathbf{x})/\tilde{P}(\mathbf{x})) = \int_{\mathbb{Y}} \exp\left(\sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y)\right) dy \triangleq Z(\mathbf{x}; \lambda)$$

使用规范化因子  $Z(\mathbf{x}; \lambda)$ , 我们记:

$$P_{\lambda}(y|\mathbf{x}) = \frac{1}{Z(\mathbf{x}; \lambda)} \exp\left(\sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y)\right)$$

则:

$$\begin{aligned}\max_{\mu, \lambda} L(P_{\mu, \lambda}, \mu, \lambda) \\ &= \max_{\lambda} \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \ln P_{\lambda}(y|\mathbf{x}) d\mathbf{x} dy \\ &+ 0 + \sum_{i=1}^n \lambda_i \int_{\mathbb{X} \times \mathbb{Y}} \left( \tilde{P}(\mathbf{x}, y) - \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \right) f_i(\mathbf{x}, y) d\mathbf{x} dy \\ &\triangleq \max_{\lambda} \Psi(\lambda)\end{aligned}$$

通过求解  $\max_{\lambda} \Psi(\lambda)$  得到  $\lambda$  的解, 代入到  $P_{\lambda}(y|\mathbf{x})$  中即可得  $P(y|\mathbf{x})$  解。此外, 还可以根据  $\mu(\mathbf{x}), \lambda$  之间的关系得到  $\mu(\mathbf{x})$  的解。

## 2 最大熵模型

### 2.1 模型描述

根据上一章的内容, 我们抽象出如下的模型, 并称之为最大熵模型:

$$P_{\lambda}(y|\mathbf{x}) = \frac{1}{Z(\mathbf{x}; \lambda)} \exp\left(\sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y)\right)$$

其中,  $f_i(\mathbf{x}, y)$  为给定的特征函数,  $Z(\mathbf{x}; \lambda)$  为规范化因子, 待学习的量为参数  $\lambda$ 。可以看到, 和线性模型、决策树等模型一样, 最大熵模型的形式给定, 我们需要对其中的待估参数进行学习。最大熵模型属于对数线性模型。

### 2.2 模型学习--极大似然估计

给定一个含有待估参数的机器学习模型和采样数据  $D$ , 我们会采用极大似然估计对未知参数进行估计:

$$\max_{\lambda} - \sum_{i=1}^N \ln P_{\lambda}(y_i | \mathbf{x}_i)$$

而事实上，上式等价于  $\max_{\lambda} \Psi(\lambda)$ ，接下来，我们来证明这一点。

1. 我们知道，极大似然估计：

$$\begin{aligned} \max_{\lambda} - \sum_{i=1}^N \ln P_{model}(y_i | \mathbf{x}_i) \\ \Downarrow \\ \max - \int_{\mathbb{X} \times \mathbb{Y}} P_{data}(\mathbf{x}, y) \ln P_{model}(y | \mathbf{x}) d\mathbf{x} dy \end{aligned}$$

因此，

$$\begin{aligned} \max_{\lambda} - \sum_{i=1}^N \ln P_{\lambda}(y_i | \mathbf{x}_i) \\ \Downarrow \\ \min_{\lambda} \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \ln P_{\lambda}(y | \mathbf{x}) d\mathbf{x} dy \end{aligned}$$

代入模型  $P_{\lambda}(y | \mathbf{x})$  进行化简：

$$\begin{aligned} & \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \ln P_{\lambda}(y | \mathbf{x}) d\mathbf{x} dy \\ &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) - \ln Z(\mathbf{x}; \boldsymbol{\lambda}) \right) d\mathbf{x} dy \\ &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \ln Z(\mathbf{x}; \boldsymbol{\lambda}) d\mathbf{x} dy \\ &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \ln Z(\mathbf{x}; \boldsymbol{\lambda}) d\mathbf{x} \end{aligned}$$

2. 对  $\Psi(\lambda)$  变形：

$$\begin{aligned}
\Psi(\lambda) &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \ln P_{\lambda}(y|\mathbf{x}) d\mathbf{x} dy + \int_{\mathbb{X} \times \mathbb{Y}} \left( \tilde{P}(\mathbf{x}, y) - \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \right) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy \\
&= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy \\
&\quad + \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \ln P_{\lambda}(y|\mathbf{x}) d\mathbf{x} dy - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy \\
&= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy \\
&\quad + \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \left( \ln P_{\lambda}(y|\mathbf{x}) - \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy \\
&= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) (-\ln Z(\mathbf{x}; \lambda)) d\mathbf{x} dy \\
&= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \ln Z(\mathbf{x}; \lambda) d\mathbf{x}
\end{aligned}$$

3. 可见对最大熵模型 (存在具体的形式和待定参数) 进行极大似然估计就等价于依据最大熵原理 (对  $P(y|\mathbf{x})$  的形式没有作任何假设) 构建模型。

## 3 数值优化算法

由前可知，我们要解决如下的优化问题：

$$\max_{\lambda} \Psi(\lambda) = \min_{\lambda} \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \ln Z(\mathbf{x}; \lambda) d\mathbf{x}$$

其中，除待优化变量  $\lambda$  未知外，其他均已知。接下来，我们介绍两种数值优化方法--改进的迭代尺度法 (improved iterative scaling, IIS) 和拟牛顿法。

### 3.1 改进的迭代尺度法 (improved iterative scaling, IIS)

**IIS 的思路是：**假设当前参数向量是  $\lambda$ ，我们希望找到一个新的参数向量  $\lambda + \delta$ ，使得目标函数增大，重复这一过程，直到找到最大值。为此，我们构造  $\Psi(\lambda + \delta) - \Psi(\lambda)$  的下界，也就是下文中的  $B(\delta; \lambda)$ ， $\Psi(\lambda + \delta) - \Psi(\lambda) \geq B(\delta; \lambda)$ ，通过不断优化下界的值来提到目标函数的值。

1. 我们有：

$$\begin{aligned}
\Psi(\lambda + \delta) - \Psi(\lambda) &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \ln \frac{Z(\mathbf{x}; \lambda + \delta)}{Z(\mathbf{x}; \lambda)} d\mathbf{x} \\
&\geq \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \left( 1 - \frac{Z(\mathbf{x}; \lambda + \delta)}{Z(\mathbf{x}; \lambda)} \right) d\mathbf{x} \\
&\quad (\text{不等式 } -\ln x \geq 1 - x, x > 0) \\
&= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + 1 - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \frac{Z(\mathbf{x}; \lambda + \delta)}{Z(\mathbf{x}; \lambda)} d\mathbf{x}
\end{aligned}$$

又由前可知：

$$Z(\mathbf{x}; \boldsymbol{\lambda}) = \frac{1}{P_{\lambda}(y|\mathbf{x})} \exp \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) = \int_{\mathbb{Y}} \exp \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) dy$$

我们有：

$$\begin{aligned} Z(\mathbf{x}; \boldsymbol{\lambda} + \boldsymbol{\delta}) &= \int_{\mathbb{Y}} \exp \left( \sum_{i=1}^n \lambda_i f_i(\mathbf{x}, y) \right) \exp \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) dy \\ &= \int_{\mathbb{Y}} Z(\mathbf{x}; \boldsymbol{\lambda}) P_{\lambda}(y|\mathbf{x}) \exp \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) dy \\ &\quad \Updownarrow \\ \frac{Z(\mathbf{x}; \boldsymbol{\lambda} + \boldsymbol{\delta})}{Z(\mathbf{x}; \boldsymbol{\lambda})} &= \int_{\mathbb{Y}} P_{\lambda}(y|\mathbf{x}) \exp \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) dy \end{aligned}$$

则

$$\begin{aligned} &\Psi(\boldsymbol{\lambda} + \boldsymbol{\delta}) - \Psi(\boldsymbol{\lambda}) \\ &\geq \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + 1 - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \frac{Z(\mathbf{x}; \boldsymbol{\lambda} + \boldsymbol{\delta})}{Z(\mathbf{x}; \boldsymbol{\lambda})} d\mathbf{x} \\ &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + 1 - \int_{\mathbb{X}} \tilde{P}(\mathbf{x}) \int_{\mathbb{Y}} P_{\lambda}(y|\mathbf{x}) \exp \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) dy d\mathbf{x} \\ &= \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + 1 - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \exp \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy \\ &\triangleq A(\boldsymbol{\delta}; \boldsymbol{\lambda}) \end{aligned}$$

至此，我们找到了一个下界  $A(\boldsymbol{\delta}; \boldsymbol{\lambda})$ ，但对  $A(\boldsymbol{\delta}; \boldsymbol{\lambda})$  的优化依然稍显复杂，为此，下面我们作进一步缩放。

2. 记  $f^{\#}(\mathbf{x}, y) = \sum_{i=1}^n f_i(\mathbf{x}, y)$ ，由 Jensen 不等式：

$$\begin{aligned} \exp \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) &= \exp \left( \sum_{i=1}^n \frac{f_i(\mathbf{x}, y)}{f^{\#}(\mathbf{x}, y)} \delta_i f^{\#}(\mathbf{x}, y) \right) \\ &\leq \sum_{i=1}^n \frac{f_i(\mathbf{x}, y)}{f^{\#}(\mathbf{x}, y)} \exp(\delta_i f^{\#}(\mathbf{x}, y)) \end{aligned}$$

则：

$$\begin{aligned} A(\boldsymbol{\delta}; \boldsymbol{\lambda}) &\leq \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) \left( \sum_{i=1}^n \delta_i f_i(\mathbf{x}, y) \right) d\mathbf{x} dy + 1 \\ &\quad - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) \sum_{i=1}^n \frac{f_i(\mathbf{x}, y)}{f^{\#}(\mathbf{x}, y)} \exp(\delta_i f^{\#}(\mathbf{x}, y)) d\mathbf{x} dy \triangleq B(\boldsymbol{\delta}; \boldsymbol{\lambda}) \end{aligned}$$

至此，我们得到了下界  $B(\boldsymbol{\delta}; \boldsymbol{\lambda})$ ，即：

$$\Psi(\boldsymbol{\lambda} + \boldsymbol{\delta}) - \Psi(\boldsymbol{\lambda}) \geq B(\boldsymbol{\delta}; \boldsymbol{\lambda})$$

3. 下面我们说明为什么优化下界  $B(\boldsymbol{\delta}; \boldsymbol{\lambda})$  比优化下界  $A(\boldsymbol{\delta}; \boldsymbol{\lambda})$  更简单，及如何对  $B(\boldsymbol{\delta}; \boldsymbol{\lambda})$  进行优化。将  $B(\boldsymbol{\delta}; \boldsymbol{\lambda})$  对任意分量  $\delta_i$  求导，并令导数为 0，我们有：

$$\frac{\partial B(\boldsymbol{\delta}; \boldsymbol{\lambda})}{\partial \delta_i} = \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}, y) f_i(\mathbf{x}, y) d\mathbf{x} dy - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\mathbf{x}) P_{\lambda}(y|\mathbf{x}) f_i(\mathbf{x}, y) \exp(\delta_i f^{\#}(\mathbf{x}, y)) d\mathbf{x} dy = 0$$

可以看到，上式中的未知量只有分量  $\delta_i$ ，而含有其他分量  $\delta_j (j \neq i)$ ，因此对每个分量  $\delta_i$ ，单独求解对应的方程即可；但若选择优化下界  $A(\boldsymbol{\delta}; \boldsymbol{\lambda})$ ，为了求解  $\boldsymbol{\delta}$ ，我们需求解的是一个联立的方程组，难度更大。

此外，若对任意的  $\boldsymbol{x}, y$ ， $f^\#(\boldsymbol{x}, y) \equiv C$ ，则  $\delta_i$  有如下的解析解：

$$\delta_i = \frac{1}{C} \ln \frac{\int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\boldsymbol{x}, y) f_i(\boldsymbol{x}, y) d\boldsymbol{x} dy}{\int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\boldsymbol{x}) P_\lambda(y|\boldsymbol{x}) f_i(\boldsymbol{x}, y) d\boldsymbol{x} dy} = \frac{1}{C} \ln \frac{\mathbb{E}_{\tilde{P}}[f_i]}{\mathbb{E}_P[f_i]}$$

否则，我们可以采用牛顿法迭代求解方程：

$$\int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\boldsymbol{x}, y) f_i(\boldsymbol{x}, y) d\boldsymbol{x} dy - \int_{\mathbb{X} \times \mathbb{Y}} \tilde{P}(\boldsymbol{x}) P_\lambda(y|\boldsymbol{x}) f_i(\boldsymbol{x}, y) \exp(\delta_i f^\#(\boldsymbol{x}, y)) d\boldsymbol{x} dy = 0$$

上述方程有单根，牛顿法恒收敛且收敛速度很快。

## 3.2 拟牛顿法

按照拟牛顿法的格式进行求解即可，可参见笔记《7-数值优化方法》中的相关内容，《统计学习方法》中给出的是 BFGS 算法，不再赘述。