

MLE与MAP

1. MLE for Discrete Random Variable

假设我们有训练集 X ，设其为离散随机变量，且其可能的取值为 $\{x_1, \dots, x_n\}$ 。我们另有模型 $f(x; \theta)$ ，且 $P(x | \theta) = f(x; \theta)$ ，同时我们假设样本是*i. i. d.*的。MLE是一个用来估计参数 θ 取值的方法，它旨在寻找一个最优的 θ_{MLE} ，使得 $P(X | \theta_{MLE})$ 最大化。也就是说，在 θ_{MLE} 参数化下，训练集 X 出现的概率会是最大的。

具体来说，当 X 为离散随机变量时，我们有似然函数：

$$\begin{aligned} L(\theta) &= P(X | \theta) \\ &= P(x_1 | \theta) \cdot P(x_2 | \theta) \cdots P(x_n | \theta) \\ &= \prod_{i=1}^n P(x_i | \theta). \end{aligned} \tag{1}$$

显然，

$$\begin{aligned} \theta_{MLE} &= \operatorname{argmax}_{\theta} L(\theta) \\ &= \operatorname{argmax}_{\theta} \log L(\theta) \\ &= \operatorname{argmin}_{\theta} -\log L(\theta) \\ &= \operatorname{argmin}_{\theta} -\sum_{i=1}^n \log P(x_i | \theta). \end{aligned} \tag{2}$$

当 X 为连续随机变量时，这篇[blog](#)用一个很好的角度将离散/连续两种情况联系在一起了。

这里提一句，为什么我们将 $L(\theta) = P(X | \theta)$ 叫作Likelihood，而不是Probability呢？其实两个不同的叫法只是为了区分出哪一个是 $P(X | \theta)$ 中的变量。

1. 给定一个 $P(X | \theta)$ ，如果我们称之为Probability，那么我们可以推出，其中 θ 已知，而 X 是变量， $P(X | \theta)$ 是一个关于 X 的函数。在这个称呼下，我们感兴趣的问题是：给定 θ 后，对任意的 X ，得到其出现的概率。
2. 给定一个 $P(X | \theta)$ ，如果我们称之为Likelihood，那么我们可以推出，其中 X 已知，而 θ 是变量， $P(X | \theta)$ 是一个关于 θ 的函数。在这个称呼下，我们感兴趣的问题是：给定 X ，寻找一个 θ ，使得在该 θ 下， X 出现的概率最大。

2. MAP for Discrete Random Variable

MLE通过优化参数 θ ，使得训练集 X 出现的概率最大化。但有些情况下我们会更加希望加入人工先验知识。

举个例子，如果我们掷硬币10次，结果全都是正面。我们若认为该实验是一个以 θ 参数化的伯努利分布，那么MLE就会得出 $\theta = 1$ 的结论。这种估计有可能是过激的。有可能这个硬币本身是均匀的，只是恰好这10次实验全是正面而已。如果能在估计中加入先验知识，那么我们估计的结果可能会更鲁棒。

MAP则是一种在估计中嵌入先验知识的方法。具体来说，若我们假设 θ 也是一个随机变量，同时它服从某个先验分布，我们可以令 $\theta_{MAP} = \operatorname{argmax}_{\theta} P(\theta | X)$ 。也就是说，给定 θ 的先验分布和 X 后，我们希望找到在这个情况下概率最大的取值 θ_{MAP} 。

显然，在 $P(\theta | X)$ 中， X 是已观测到的随机变量， θ 是未观测到的随机变量，因此我们无法直接对其进行计算。所以我们这里应用贝叶斯公式来规避 $P(\theta | X)$ 的计算，具体来说，

$$\begin{aligned}
 \theta_{MAP} &= \operatorname{argmax}_{\theta} P(\theta | X) \\
 &= \operatorname{argmax}_{\theta} \frac{P(X | \theta)P(\theta)}{P(X)} && \text{贝叶斯公} \\
 &= \operatorname{argmax}_{\theta} P(X | \theta)P(\theta) && P(X) \text{与} \theta \text{无} \\
 &= \operatorname{argmax}_{\theta} P(\theta) \prod_{i=1}^n P(x_i | \theta) \\
 &= \operatorname{argmax}_{\theta} \log \left(P(\theta) \prod_{i=1}^n P(x_i | \theta) \right) \\
 &= \operatorname{argmax}_{\theta} \left(\log(P(\theta)) + \sum_{i=1}^n \log(P(x_i | \theta)) \right) \\
 &= \operatorname{argmin}_{\theta} \left(-\log(P(\theta)) - \sum_{i=1}^n \log(P(x_i | \theta)) \right).
 \end{aligned}$$

因此，我们可以通过预设 θ 的先验分布 $P(\theta)$ ，来嵌入先验知识。

比如，对于上面掷硬币的例子，我们可以假设 θ 服从*Bernoulli*分布的共轭先验分布*Beta*分布，即 $\theta \sim \text{Beta}(\alpha, \beta)$ ，从而使得估计更加鲁棒。当然这里也可以选择其它的先验分布，只是如果我们采用共轭先验分布的话，可以使得后验分布 $P(\theta | X)$ 与先验分布 $P(\theta)$ 有着相同的数学形式，在计算和理解上带来方便。

MAP中两个值得注意的地方：

1. 若假设先验分布 θ 服从均匀分布，则 $P(\theta)$ 为常数。此时 $\theta_{MAP} = \theta_{MLE}$ ，即MAP等价于MLE。

-
2. 从机器学习的角度来看，MAP可以看作是MLE增加了一个关于参数的先验分布的正则项 $\log(P(\theta))$ 。

References

[1] [Link1](#)

[2] [Link2](#)

[3] [Link3](#)

[4] [Link4](#)