

人工智能导论第三次作业

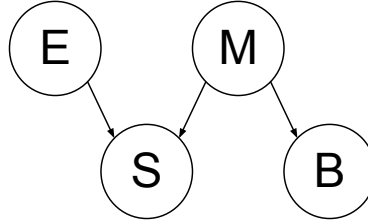
2023年5月

1 贝叶斯网络 (10分)

考虑以下含有四个变量 E, S, M, B 的贝叶斯网络。相应的条件概率表如下。

$P(E)$	
$+e$	0.4
$-e$	0.6

$P(S E, M)$			
$+e$	$+m$	$+s$	1.0
$+e$	$+m$	$-s$	0.0
$+e$	$-m$	$+s$	0.8
$+e$	$-m$	$-s$	0.2
$-e$	$+m$	$+s$	0.3
$-e$	$+m$	$-s$	0.7
$-e$	$-m$	$+s$	0.1
$-e$	$-m$	$-s$	0.9



$P(M)$	
$+m$	0.1
$-m$	0.9

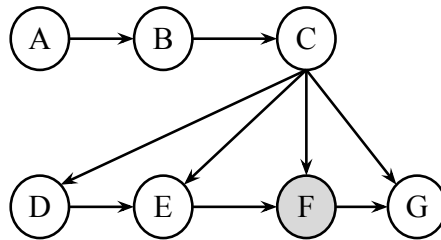
$P(B M)$		
$+m$	$+b$	1.0
$+m$	$-b$	0.0
$-m$	$+b$	0.1
$-m$	$-b$	0.9

计算以下概率:

- (1) $P(-e, -s - m, -b)$
- (2) $P(+b)$
- (3) $P(+m \mid +b)$
- (4) $P(+m \mid +s, +b, +e)$
- (5) $P(+e \mid +m)$

2 变量消除法 (10分)

考虑以下所有变量均为二元变量的贝叶斯网络。



我们想要使用变量消除法计算 $P(B, D | F = f)$ ，变量消除的顺序为 A, C, E, G 。初始的所有因子如下：

$$P(a), P(b|a), P(c|b), P(d|c), P(e|c, d), P(f|c, e), P(g|c, f)$$

当消除 A 时，我们产生一个新的因子 $\tau_1(b)$ ，如下所示

$$\tau_1(b) = \sum_a P(b|a)P(a)$$

剩余因子如下：

$$\tau_1(b), P(c|b), P(d|c), P(e|c, d), P(f|c, e), P(g|c, f)$$

请继续完成以下步骤：

- (1) 当消除 C 时，产生新的因子 τ_2 ，请写出其表达式并列出所有剩余因子。
- (2) 当消除 E 时，产生新的因子 τ_3 ，请写出其表达式并列出所有剩余因子。
- (3) 当消除 G 时，产生新的因子 τ_4 ，请写出其表达式并列出所有剩余因子。
- (4) 如何利用剩余因子计算 $P(B = b, D = d | F = f)$ ？
- (5) 因子大小是变量消除法计算复杂度的关键因素。例如，假设所有的变量都是二元变量，则因子 $P(b|a)$ 的大小是2，它有 2^2 种取值需要维护；因子 $P(e|c, d)$ 的大小是3，它有 2^3 种取值需要维护。而由于 f 的值已观测，因此 $P(g|c, f)$ 的大小只有2。你可能会发现，按照 A, C, E, G 的顺序并不是一个很好的顺序。请找出一个变量消除顺序，使得最大的因子最小。列出这个顺序，并给出在新的顺序下每次产生的因子大小。（提示：产生的最大新因子大小是2）

3 带方差的高斯线性回归（10分）

在高斯线性回归中，我们将 σ 也视作模型参数的一部分。假设 $y_n | \mathbf{w}, \mathbf{x}_n, \sigma \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)$ ，即

$$p(y_n | \mathbf{w}, \mathbf{x}_n, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (y - \mathbf{w}^T \mathbf{x}_n)^2\right),$$

其中 (\mathbf{x}_n, y_n) 为第 n 条数据， \mathbf{w} 和 σ 为模型参数。假设数据集采样自独立同分布，并使用最大似然估计来求解，即

$$\log p(\mathcal{D}_n; \mathbf{w}, \sigma) = \sum_{i=1}^N \log p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma),$$

$$\hat{\mathbf{w}}, \hat{\sigma} = \operatorname{argmax}_{\mathbf{w}, \sigma} \log p(\mathcal{D}_n; \mathbf{w}, \sigma).$$

求 $\hat{\mathbf{w}}$ 和 $\hat{\sigma}$ 。

4 采样 (10分)

考虑一个在 100×100 网格上的采样问题。对于样本 $\mathbf{x} = (x_1, x_2) \in \{1, 2, \dots, 100\}^2$, 其概率满足

$$p(x_1, x_2) \propto x_1 + \ln(x_1 x_2 + 2x_1 + 3x_2).$$

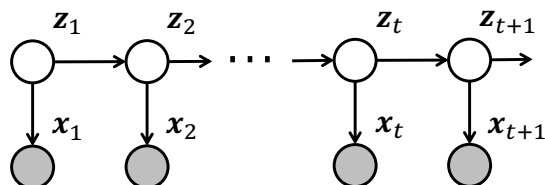
由于样本空间非常庞大, 且归一化系数难以求解, 我们使用马尔可夫链蒙特卡洛 (Markov Chain Monte Carlo) 进行采样。假设初始位于 $(2, 2)$ 。

(1) 使用Metropolis-Hastings算法进行采样, 所选的提议分布 (proposal distribution) Q 为均匀分布。若按照 Q 进行采样得到 $(1, 3)$, 则接受率 $\alpha_{(2,2),(1,3)}$ 应该是多少? 若采样得到的是 $(3, 4)$, 则接受率 $\alpha_{(2,2),(3,4)}$ 又应该是多少?

(2) 使用Gibbs采样法进行采样。第一次采样过程需要采样 x_1 , 请写出提议分布 $q(x_1)$ 的形式 (不需要计算)。在此基础上, 若采样到 $x_1 = 42$, 请写出第二次采样过程需要采样的变量及其提议分布的形式 (不需要计算)。

5 Baum-Welch算法 (20分)

Baum-Welch算法是EM算法的一种, 其解决了隐马尔科夫模型 (Hidden Markov Model, HMM) 三大主要问题中的学习问题。HMM的学习问题可以按如下方式定义: 给定观测序列 $X = \{x_1, \dots, x_T\}$, 在隐藏序列 $Z = \{z_1, \dots, z_T\}$ 未知的情况下, 如何估计模型的最佳参数 θ , 使得 $P(X | \theta)$ 最大。



参数 $\theta = \{\pi, A, B\}$, 包括初始概率分布 $\pi = [\pi_i]_N$, 转移 (Transition) 矩阵 $A = [a_{ij}]_{N \times N}$, 观测/发射 (Emission) 矩阵 $B = [b_j(k)]_{N \times M}$, 其中 N 表示隐状态总数, M 表示可观测状态总数。即:

$$P(z_1 = i | \theta) = \pi_i$$

$$P(z_{t+1} = j | z_t = i, \theta) = a_{ij}$$

$$P(x_t = k | z_t = j, \theta) = b_j(k)$$

在本题中, 我们将利用先前的知识完成该算法的推导。

(1) 首先进行E步的计算, 在这一步中我们根据当前的网络参数 θ^{old} 得到隐变量 Z 的后验概率 $q(Z) = P(Z | X, \theta^{\text{old}})$, 然后固定住 $q(Z)$, 并列出ELBO得到优化目标。请根据ELBO写出 $J(\theta)$, 并证明:

$$\arg\max_{\theta} J(\theta) = \arg\max_{\theta} \sum_Z P(X, Z | \theta^{\text{old}}) \log P(X, Z | \theta)$$

(2) 使用模型的参数 θ (包括 π, A, B) 来表示 $P(X, Z | \theta)$, 请给出相应的形式。

(3) 令 $Q(\theta, \theta^{\text{old}}) = \sum_Z P(X, Z | \theta^{\text{old}}) \log P(X, Z | \theta)$, 应用(2)中的结果, 将 $Q(\theta, \theta^{\text{old}})$ 拆分成三项之和, 其中每一项仅与 π, A, B 中的一个参数有关。

提示: 包含 π 的项为

$$\sum_Z P(X, Z | \theta^{\text{old}}) \log \pi_{z_1}.$$

(4) 上一小题实现了参数之间的解耦, 可以进入 **M** 步的计算, 此时暂时不需要考虑 θ^{old} 的处理。我们以 π 的求解为例。请首先证明

$$\sum_Z P(X, Z | \theta^{\text{old}}) \log \pi_{z_1} = \sum_{i=1}^N P(X, z_1 = i | \theta^{\text{old}}) \log \pi_i,$$

然后使用拉格朗日乘子法, 求 $\sum_Z P(X, Z | \theta^{\text{old}}) \log \pi_{z_1}$ 在

$$\sum_{i=1}^N \pi_i = 1$$

下取到极大值时 π_i 的取值。

(5) 参照上一小题的过程, 求 $Q(\theta, \theta^{\text{old}})$ 取极大值时 a_{ij} 和 $b_j(k)$ 的取值。

(6) 注意到 $\pi_i, a_{ij}, b_j(k)$ 的取值中都包含有关 θ^{old} 的项, 现在我们对它们进行求解。定义

$$\gamma_t^{\text{old}}(i) = P(z_t = i | X, \theta^{\text{old}})$$

$$\xi_t^{\text{old}}(i, j) = P(z_t = i, z_{t+1} = j | X, \theta^{\text{old}})$$

回顾课上讲过的前向-后向算法, 对于给定参数 θ^{old} 的 HMM, 我们可以求出

$$\alpha_t^{\text{old}}(i) = P(z_t = i | x_{1:t}, \theta^{\text{old}}), \beta_t^{\text{old}}(i) = P(x_{t+1:T} | z_t = i, \theta^{\text{old}}),$$

请使用 $\alpha_t^{\text{old}}(i), \beta_t^{\text{old}}(i), a_{ij}^{\text{old}}, b_j^{\text{old}}(k)$ 表示出 $\gamma_t^{\text{old}}(i)$ 和 $\xi_t^{\text{old}}(i, j)$ 。

(7) 回顾(4)(5)的结果, 请使用 $\gamma_t^{\text{old}}(i)$ 和 $\xi_t^{\text{old}}(i, j)$ 表示出 $\pi_i, a_{ij}, b_j(k)$, 并给出 Baum-Welch 算法的伪代码。

6 LDA (40分)

请使用 python 实现 Variational EM LDA。本次作业在 `./dataset` 中提供了三种不同的数据集, `dataset.txt` 是英文的小规模数据集, `dataset.cn.txt` 是中文的中等规模数据集, `dataset.cn.full.txt` 是中文的大规模数据集。建议在较小数据集上验证实现正确性之后再使用较大的数据集。以下是作业要求:

(a) 根据提供的代码框架, 写出 Variational EM LDA 的伪代码。

(b) 完成代码框架中缺失的变分推断部分。代码框架中已经实现了对于 α, β 的更新, 只需要补充 `main.py` 的两个函数, 计算 ELBO 并更新 γ, ϕ 。

(c) 设置主题个数 K 为 5, 10, 20, 使用 `dataset.cn.full.txt` 数据集, 针对不同的 K 显示每个 topic 中出现频率最高的 8 个单词。

(d) 观察结果, 找到主题分类效果最好的 K , 并分析原因。

补充说明:

1. 本次代码框架中使用了scipy, $\log(\text{gamma}(x))$ 是 gammaln 函数, $\log(\text{gamma}(x))$ 的导数是 psi 函数。
2. 本次代码框架没有引入 λ , 在变分推断更新 γ 和 ψ 时可能与课件有所出入, 同学们可以参考原论文中的这一更新过程。
3. 考虑到时间问题, 对于大规模数据集`dataset.cn_full.txt`, 最大更新轮次(`max_epochs`) 设置为10轮即可。