

Chapter 6: Probabilistic Graphical Model

2018.12.25 Tuesday. By KuKuXia@github.com

ML/DL/RL深度交流讨论群: 622545311, 加群时请备注: 姓名-方向-公司OR高校

用贝叶斯节点表示观测到的数据, 用隐含节点表示潜在的知识, 用边来描述知识与数据的相关关系, 最好基于这样的关系而获得一个概念分布

有向图 ○ 贝叶斯网络(Bayesian Network)

无向图 ○ 马尔科夫网络(Markov Network)

主要类别

PGM

假设可以观测到的变量集合为X, 需要预测的变量集合为Y, 其他变量集合为Z

$$P(Y|X) = \frac{P(X, Y)}{P(X)} = \sum_Z \frac{P(X, Y, Z)}{\sum_Z P(X, Y, Z)}$$

对联合概率分布P(X, Y, Z)进行建模, 在给定观测集合X的条件下, 通过计算边缘分布率得到对变量集合Y的推断

朴素贝叶斯

贝叶斯网络

pLSA

LDA

生成式模型

隐马尔科夫模型

直接对条件概率分布和P(Y, Z|X)进行建模, 然后向未知变量Z就可以得到对变量集合Y的预测

$$P(Y|X) = \sum_Z P(Y, Z|X)$$

判别式模型

最大熵模型

条件随机场

对序列数据进行建模

满足无后效性的随机过程

$$P(x_n|x_1, x_2, \dots, x_{n-1}) = P(x_n|x_{n-1})$$

假设一个随机过程中, $t+1$ 时刻的状态 x_t 的条件分布, 仅仅与其前一个状态 x_{t-1} 有关, 和距离为马尔科夫过程

时间和状态取值都是离散的马尔科夫过程也称为马尔科夫链

马尔科夫过程

在鲁棒的马尔科夫模型中, 所有状态对于观测者都是可见的, 因此在马尔科夫模型中只包含状态间的转移概率

隐马尔科夫模型是对含有未知参数状态状态的马尔科夫链进行建模的生成模型

隐状态 x_t 对于观测者而言是不可见的, 能观测到的只有每个隐状态 x_t 对应的输出 y_t , 而观测状态 y_t 的概率分布仅仅取决于对应的隐状态 x_t

假定观测状态两两相互独立

隐状态间的转移概率

隐状态到观测状态的输出概率

隐状态 x_t 的取值空间

观测状态 y_t 的取值空间

初始状态的概率分布

参数

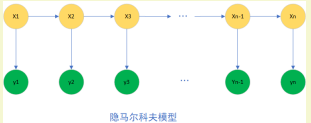
概率计算问题

预测问题

学习问题

基本问题

隐马尔科夫模型



去掉了隐马尔科夫模型中观测状态相互独立的假设, 考虑了整个观测序列, 更准确的表达能力

直接对标注后的后验概率P(y|x)进行建模的判别式模型

$$P(x_{1:n}|y_{1:n}) = \prod_{i=1}^n P(x_i|x_{1:n-1}, y_{1:n})$$

建模公式

最大熵隐马尔科夫模型MEMM: Maximum Entropy Markov Model

由于局部性一致的影响, 隐状态会倾向于转移到相邻的隐状态引致更少的状态上, 以提高整体的后验概率, 这就是标注偏置问题

在MEMM的基础上, 进行了全局归一化, 解决了局部归一化带来的标注偏置问题

$$P(x_{1:n}|y_{1:n}) = \frac{1}{Z(y_{1:n})} \prod_{i=1}^n \exp(P(x_i, x_{i-1}, y_{1:n}))$$

Z(y_{1:n})是在全局范围进行归一化

条件随机场CRF: Conditional Random Field

$y = \max_m P(y|x)$

通过预测指定样本属于特定类别的概率 $P(y_i|x)$ 来预测该样本的所属类别

$$P(y_i|x) = \frac{P(x|y_i)P(y_i)}{P(x)}$$

贝叶斯公式

$$P(y_i|x) \propto P(x|y_i)P(y_i) = P(x_1|y_i)P(x_2|y_i) \dots P(x_n|y_i)P(y_i)$$

后验概率 $P(x_i|y_i)$ 的值决定了分类的结果, 并且任意特征 x_i 都受 y_i 的取值影响

朴素贝叶斯模型的生成模型

上图是图式记法, 表示变量 y 同时对 x_1, x_2, \dots, x_n 这 N 个变量产生影响

信息是指人们对于事物的理解不确定性的降低或消除, 而不确定性是指不确定性, 熵越大, 不确定性就越大

最大熵原理: 概率分布学习的一个准则, 指导模型是在满足约束条件的模型集合中选取熵最大的模型, 即不确定性最大的模型, 就是最接近联合分布 $P(y|x)$, 受约束条件 $\sum_i P(y_i|x)$ 的约束

$$H(P) = - \sum_{ij} P(x_i) \log P(x_i)$$

假设离散随机变量 x 的分布 $P(x)$, 则关于 x 的熵的定义如图, 当 x 服从均匀分布时, 对应的熵最大, 也就是不确定性最高

$$H(P) = - \sum_{ij} P(x_i) P(y_j|x_i) \log P(y_j|x_i)$$

给定离散随机变量 x 和 y 上的条件概率分布 $P(y|x)$, 定义在条件概率分布上的条件熵

$$H(y|x) = - \sum_{ij} P(x_i) P(y_j|x_i) \log P(y_j|x_i)$$

最大熵模型的学习等价于约束最优化问题

$$P_w(y|x) = \frac{1}{Z} \exp \left(\sum_{ij} w_{ij} f_{ij}(x, y) \right)$$

最大熵模型的表达式形式, 即最终归结为学习最佳的参数 w , 使得 $w(y|x)$ 最大化

最大熵模型的概率图模型

概率图模型: $P_w(y|x)$ 的表达式形式非常类似于势函数为指数函数的马尔科夫网络, 其中变量 x 和 y 构成了最大团

基于词频模型或N-gram模型的文本表示模型有一个明显缺点, 就是无法识别出两个不同的短语但具有相同的话题

主题模型是一种能够识别具有相同主题的词频词矩阵到一个潜在主题

一个生成模型来建模文档的生成过程, 概率图模型, 将每个文档对应的主题分布 $p(\theta_i | d, m)$ 和每个主题对应的词分布 $p(w_i | \theta_i, k)$ 看成确定的未知参数, 并可以求解出来

利用似然函数表示整个语料库中的文本生成概率, 但是由于参数中包含有隐变量, 因此无法用最大似然法直接求解, 可以利用最大期望算法来求解

贝叶斯推理思想, 认为待估计的参数(主题分布和词分布)不是一个确定的参数, 而是服从一定分布的随机变量, 这个分布符合一定的先验概率分布, 即狄利克雷分布(Dirichlet), 并且在观察到样本信息之后, 可以对先验分布进行修正, 从而得到后验分布

确定K的第一种方法: 可以利用训练集、验证集进行训练, 然后用测试集对超参数进行选择

评估指标 ○ 困惑度Perplexity

主题个数K是一个预先指定的超参数

是一种非参数主题模型HDP-LDA

不需要预先指定主题的个数, 模型可以随着文档数目的变化而自动对主题个数进行调整

缺点是使得模型更加复杂, 训练速度更加缓慢, 实际中经常采用另一种方法确定主题个数K

分层狄利克雷过程HDP: Hierarchical Dirichlet Process

是而在没有大量用户数据的情况下如何给用户进行个性化推荐, 目的是优化点击率、转化率或者用户停留时间(用户停留时间、留存率等)

用户冷启动 ○ 指对一个之前没有行为或行为很少的新用户进行推荐

物品冷启动 ○ 指为一个新上市的商品或电影(这时没有与之相关的评分或用户行为数据), 寻找具有潜在兴趣的用户

系统冷启动 ○ 指如何为一个新开发的网站设计个性化推荐系统

解决方法 ○ 一般解决方法是基于内容的推荐