

Chapter 8: Sampling

采样的作用

- 从特定的概率分布中抽取对应的样本点
 - 采样得到的样本集可以看作一个非参数模型，而用少量的样本点（经验分布）来近似总体分布
 - 例子：随机森林的弱分类器训练，每个弱分类器使用训练集的一个子集来训练，然后组合起来得到整体分布的一个样本集，作为训练集来训练，称之为bagging。
 - 训练模型的时候是使最小化模型在训练集上的损失函数(过拟合风险)，在评估模型时，使用另一个总体分布的样本集作为测试集
 - 用于信息可视化，帮助人们决策，直观地了解总体分布中数据的结构和特性
- 信息浓缩
 - 对当前数据集重采样，充分利用已有数据集，挖掘更多信息
- 随机模拟
 - 含有独立变量的联合模型，可以模拟含有独立变量的分布进行采样
 - 指数模型和近似求解方法，不适用于一般确定性的近似求解方法，如蒙特卡罗方法、期望传播等

不均采样本集的采样

- 基于数据的方法
 - 对数据集进行重采样，使原本不均匀的样本变得均衡
 - 过采样Over Sampling
 - 随机从少数类样本集中重复抽取样本(有放回)以得到更多的样本
 - 对少数类样本进行了多次重采，扩大了数据规模，缓解了模型的分层度，缓解了过拟合
 - SMOTE
 - 为每个少数类的样本合成相同数量的新样本，这可能会增加类间重叠度，并且会生成一些不能提供有益信息的样本
 - Borderline-SMOTE
 - 只给那些分类边界上的少数类样本合成新样本
 - ADASYN
 - ADASYN根据不同的少数类样本合成不同个数的新样本
 - 欠采样Under Sampling
 - 随机从多数类样本集中选取较少的样本，有放回或无放回
 - 会丢弃一些样本，可能会丢失部分有用信息，造成模型只学到了数据模式的一部分，导致欠拟合
- 基于算法的方法
 - 将问题转化为单类学习One-class learning，异常检测Anomaly detection
 - 改变模型训练时的目标函数，如代价敏感学习中不同类别有不同的权重来纠正这种不平衡性

贝叶斯网络的采样

- 概率图模型描述多个随机变量的联合概率分布，又称联合网络或有向无环图模型
- 贝叶斯网络是利用有向无环图来刻画一组随机变量之间的条件概率分布关系
- 祖先采样
- 逻辑采样
- 粒度和投票采样
- MCMC采样法

高斯分布的采样

- 逆变换法
 - Box-Muller
 - 需要计算三角函数
 - Marsaglia Polar Method
- 拒绝采样法

均匀分布随机数

- 计算机是确定性的，存储和计算单元只能处理离散状态值
- 计算机不能产生完全的均匀分布随机数，只能产生离散分布的随机数来逼近连续分布，使用很大的离散空间来提供足够的精度
- 生成方法
 - 利用当前生成的随机数 x_i 来产生下一个随机数 x_{i+1} ，初始值 x_0 称为种子
 - 得到区间 $[0, m-1]$ 上的随机数，除以 m 得到区间 $[0, 1]$
 - 随机数并不严格相互独立，最多产生 m 个不同的随机数
 - 线性同余法 Linear Congruential Generator
 - $x_{i+1} \equiv a \cdot x_i + c \pmod{m}$
 - 主要产生方式
 - 生成方式
 - $m = 2^{31} - 1$
 - $a = 1103515245$
 - $c = 12345$
 - gcc取用版本

常见采样方法

- 几乎所有的采样方法都通过均匀分布随机数作为基本操作
- 简单分布无离散分布可以称做连续法来采样
- 函数变换法
 - 如果变换关系 $g(\cdot)$ 是 x 的累积分布函数的话，则得到逆函数 g^{-1} 来采样(Inverse Transform Sampling)
 - 假设待采样的目标分布的概率密度函数为 $p(x)$ ，它的累积分布函数为
 - 1. 从均匀分布 $U(0, 1)$ 产生一个随机数 u_i
 - 2. 计算 $x_i = g^{-1}(u_i)$ ，其中 $g^{-1}(\cdot)$ 是累积分布函数的逆函数
 - 如果带采样的目标分布的概率密度函数无法求解或者不容易计算，则不适用于函数变换法
- 拒绝采样法
 - 拒绝采样(Accept-Reject Sampling)
 - 自适应拒绝采样
 - 在拒绝采样中，如果在第一步中采样被拒绝，则逐步不会产生新样本，需要重新进行采样
- 重要性采样(Importance Sampling)
 - 可以构造一个容易采样的参考分布，先对参考分布进行采样，然后对得到的样本进行一定的后处理操作，使得最终的样本服从目标分布
 - 对于类别估计的随机性误差，重要性采样会降低方差分布，采样效率低下，样本的误差随样本量减小而显著性降低
- 基于采样的数值近似求解方法
 - 蒙特卡诺法
 - 用于进行采样
 - 马尔科夫链
 - 针对待采样的目标分布，构造一个马尔科夫链，使得该马尔科夫链的平稳分布就是目标分布
 - 然后，从任何一个初始状态出发，沿着马尔科夫链进行状态转移，最终得到的状态转移序列会收敛到目标分布，由此可以得到目标分布的一系列样本
 - Metropolis-Hastings采样法
 - 马尔科夫链对应的不同的MCMC采样法
 - 吉布斯采样法
 - MCMC采样法每一步都会产生一个样本，只有当这个样本与之前的样本一样时才会停止
 - MCMC采样法满足不能法代过程中收敛到平稳分布的，因此实际应用中一般会得到收敛样本序列 $\{x_1, \dots, x_n\}$ 来逼近，即根据收敛到平稳分布的样本，以收敛后的样本
 - MCMC产生的样本不是独立的，因为后一个样本是由前一个样本根据特定的转移概率得到的，或者有一道概率度量前一个样本，如果仅靠采样，并不需要样本之间相互独立
 - 产生独立样本
 - 如果确实需要产生独立同分布的样本，可以同时运行多条马尔科夫链，这样不同链上的样本是相互独立的
 - 或者用一条马尔科夫链上每隔若干个样本抽取一个，这样抽取出来的样本也是相互独立的。

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

2018.12.20 Thursday. By KuKuXia@githu.com