



2.5.0

概率密度直方图方法：① 为估计在某个特定位置的概率密度。
应考虑位于那个点的某个领域内的数据点。
② 为获得好的结果，平滑参数的值既不能太大也不能太小。

2.5.1. 假设有一未知的概率密度分布 $p(x)$ ，收集了 N 次观测，每个数据点落入 R 中的概率为

$$\text{Bin}(k, N, p) = \frac{N!}{k!(N-k)!} p^k (1-p)^{N-k}$$

假定 N 较大，则 $p \approx \frac{k}{N}$ 即 $k \approx NP$ ①

假定 R 足够小，使在这个区域的概率密度 $p(x)$ 大致为常数

则 $p \approx p(x)V$ 其中 V 为区域 R 的体积。②

由 ①. ② 得 $p(x) = \frac{k}{NV}$

此式依赖于两个相互矛盾的假设即区域 R 要足够小，使这个区域内的概率密度近似为常数，但也要足够大，使 k 足够让二项分布达到尖峰。

核方法：固定 V 然后从数据中确定 k 。

R 取为以 x 为中心的小立方体

定义 $k(u) = \begin{cases} 1 & |u_i| \leq \frac{1}{2} \quad i=1, \dots, D \\ 0 & \text{其他} \end{cases}$

当边长为 h 时 $k(\frac{x-x_0}{h}) = \begin{cases} 1 \\ 0 \end{cases}$

位于立方体内数据点总数为

$$k = \sum_{i=1}^N k(\frac{x-x_i}{h})$$

即 $p(x) = \frac{1}{N} \sum_{i=1}^N \frac{1}{h^D} k(\frac{x-x_i}{h})$

D 维边长为 h 的立方体 $V = h^D$

核密度估计或 Parzen 估计

问题：人为带来的非连续性。且训练阶段只需存储数据。



2.5.2

一个以 x 为中心的小球体 且半径可以自由增长. 直到他精确的包含 k 个数据点 利用公式 $p(x) = \frac{k}{NkV}$ 计算概率分布

近邻方法: 固定 k 的值然后确定 V



推广到分类.

一数据集 N_k 个数据属于类别 C_k , 数据点总数为 N . $\therefore \sum N_k = N$.

现对一新的数据点进行分类. 画一个以 x 为中心的球 且精确地包含 k 个数据点 (无论属于哪个类别) 设球体体积为 V 且包含 C_k 类的 k_k 个数据点 则 $p(x|C_k) = \frac{k_k}{NkV}$ ~ 每个类别关联的概率密度估计.

\therefore 无条件概率密度为 $p(x) = \frac{k}{NkV}$

先验为 $p(C_k) = \frac{N_k}{N}$

由贝叶斯定理得 $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)} = \frac{k_k}{k}$

因为要最小化错误分类的概率. \therefore 应把测试点 x 分配给有最大后验概率的类别. 对应最大的 $\frac{k_k}{k}$.

最近邻 ($k=1$) 分类器 在 $N \rightarrow \infty$ 时 错误率 不会超过最优分类器

可达到的最小错误率的 2 倍



即用真实概率分布的分类器.

