

カーネル密度推定法

氏名：西本 洋紀

2019 年 5 月 27 日

1 目的

ある分布 $p(\mathbf{x})$ に従って、 N 個の観測値が得られたとする。この時、 N 個の観測値から、分布 $p(\mathbf{x})$ を推定したい。

2 方法

\mathbf{x} を含むある小さな領域 \mathcal{R} を考える。この領域に割り当てられた確率 P は

$$P = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x} \quad (1)$$

と書ける。

領域 \mathcal{R} のデータ点の総数 K は、二項分布に従う。

$$\text{Bin}(K|N, P) = \binom{N}{K} P^K (1-P)^{N-K} \quad (2)$$

よって、 $\frac{K}{N}$ の分布の平均と分散は

$$\mathbf{E}[K/N] = P \quad (3)$$

$$\text{var}[K/N] = \frac{P(1-P)}{N} \quad (4)$$

となる事がわかる。

$N \rightarrow \infty$ のとき、 $\frac{K}{N}$ の分布は平均周りで尖った形になり、

$$K \simeq NP \quad (5)$$

となる。

また、 \mathcal{R} が小さく、確率密度 $p(\mathbf{x})$ がこの領域内で一定とみなせるとしたら、

$$P \simeq p(\mathbf{x})V \quad (6)$$

となる。ただし、 V は、 \mathcal{R} の体積である。

V を固定し、 K をデータから推定できたとすると、

$$p(\mathbf{x}) \simeq \frac{K}{NV} \quad (7)$$

により、 $p(\mathbf{x})$ を推定できる。あとは K をデータから推定をする部分を考えれば良い。

確率密度を求めたい点を \mathbf{x} 、この点を中心とする超立方体を \mathcal{R} とする。また、観測値の集合を $\{\mathbf{x}_n\}$ とする。さらに、カーネル関数

$$k(\mathbf{u}) = \begin{cases} 1 & (\|\mathbf{u}\| \leq \frac{1}{2}) \\ 0 & (o.w.) \end{cases} \quad (8)$$

を定義する。この方法では、 $k(\frac{\mathbf{x}-\mathbf{x}_n}{h})$ は \mathbf{x} を中心とする一辺が h の立方体 \mathcal{R} の内部にデータ点があれば 1 に、そうでなければ 0 になる。

立方体 \mathcal{R} 内のデータ点の総数は

$$K = \sum_{n=1}^N k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \quad (9)$$

と書け、 \mathbf{x} での推定密度は

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \quad (10)$$

となる。上式の確率密度のクラスは、カーネル密度推定法や **Parzen 推定法** と呼ばれる。

このとき、立方体の縁で、人為的な不連続が生じてしまう。より滑らかなカーネル関数を使えば、より滑らかな密度が得られる。カーネル関数 $k(\mathbf{u})$ には、以下の条件を満たす任意の関数を選ぶことができる。

$$k(\mathbf{u}) \geq 0 \quad (11)$$

$$\int k(\mathbf{u}) d\mathbf{u} = 1 \quad (12)$$

一般によく選ばれるカーネル関数はガウスカネルで、次の確率密度モデルになる。

$$p(\mathbf{x}) = \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{-\frac{(\|\mathbf{x}-\mathbf{x}_n\|)^2}{2h^2}\right\} \quad (13)$$

ただし、 h はガウス分布の標準偏差を表す。パラメータ h は平滑化パラメータの役割を果たし、 h が小さすぎると、結果はノイズの多い密度モデルになるが、大きすぎると、データの局所的な性質が捉えにくくなる。最良の密度は、中間的な h の値で得られる。