

Density estimation: kernel methods

Juneteenth, 2020

Given a sample of data X_1, \dots, X_n , we would like to estimate its underlying **probability density function** f .

$$P(a < X < b) = \int_a^b f(x) dx \quad (1)$$

Considering that for a **window width** (or a bin width) h

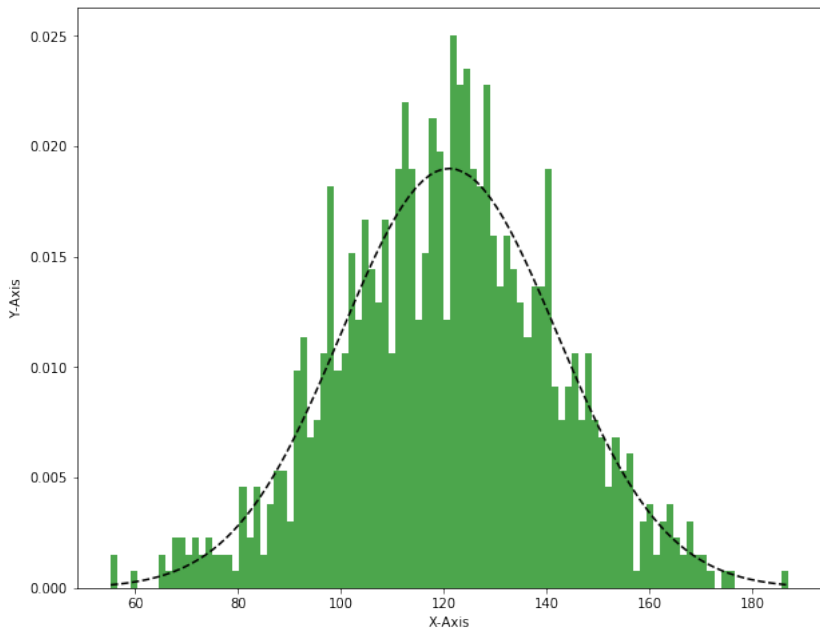
$$f(x) = \lim_{h \rightarrow 0} \frac{1}{2h} P(x - h < X < x + h) \quad (2)$$

we can define the estimate \hat{f} as

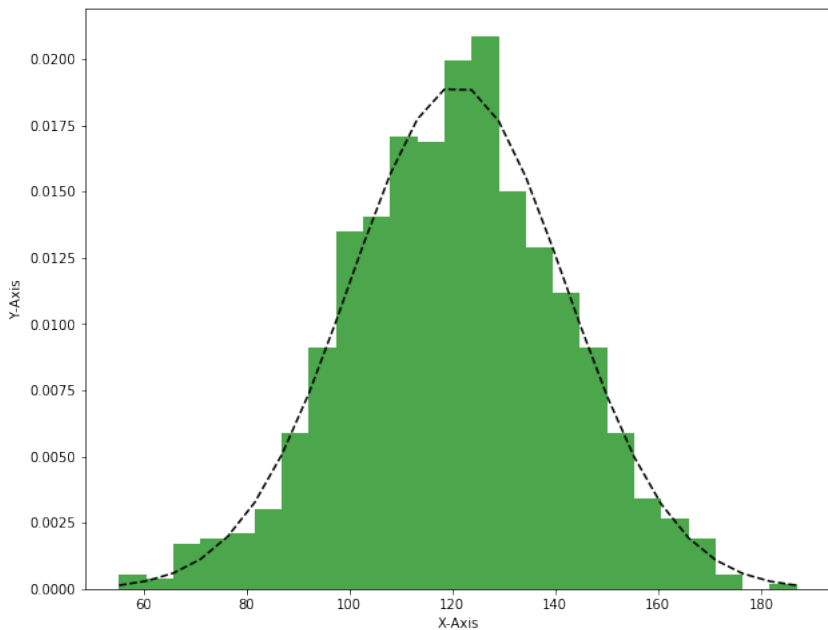
$$\hat{f}(x) = \frac{1}{2hn} m \quad (3)$$

with m being the amount of samples X_i falling in $[x - h, x + h]$

Histograms with 100 bins



Histograms with 25 bins



Define

$$d_1(t) \leq d_2(t) \leq \dots \leq d_n(t)$$

as the distance from a point t to the points of the sample.

For each $r > 0$, you expect $2rnf(t)$ samples to fall in $[t - r, t + r]$.
So in $[t - d_k(t), t + d_k(t)]$ you have k samples by definition.

Define ¹

$$\hat{f}(t) = \frac{k}{2nd_k(t)} \quad (4)$$

¹"A non-parametric estimate of a multivariate density function", by Loftsgaarden and Quesenberry, 1965

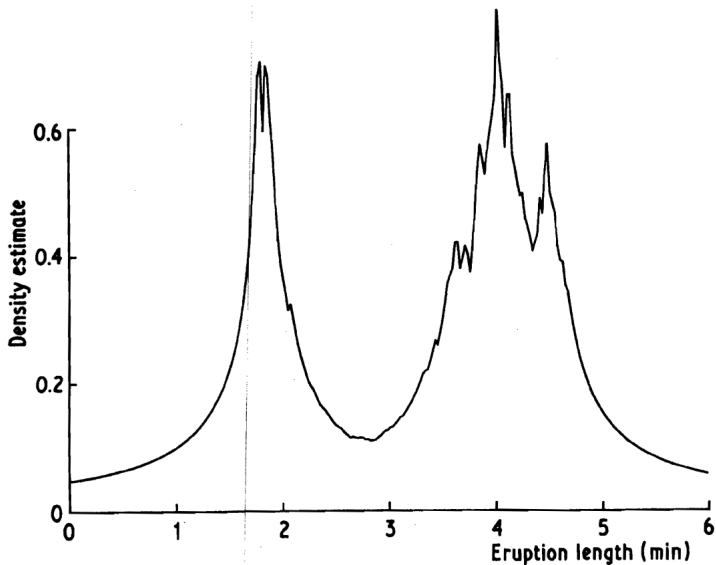


Fig. 2.10 *Nearest neighbour estimate for Old Faithful geyser data, $k = 20$.*

Assume K a function satisfying

$$\int_{-\infty}^{\infty} K(x)dx = 1 \quad (5)$$

The kernel estimator is defined ² as

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right) \quad (6)$$

²"Remarks on some non-parametric estimates of a density function", by Rosenblatt, 1956

The variable kernel estimator

$$\hat{f}(t) = \frac{1}{n} \sum_{j=1}^n \frac{1}{hd_{j,k}} K\left(\frac{t - X_j}{hd_{j,k}}\right) \quad (7)$$

or the generalized k-th nearest neighbour estimate

$$\hat{f}(t) = \frac{1}{nd_k(t)} \sum_{i=1}^n K\left(\frac{t - X_i}{d_k(t)}\right) \quad (8)$$

For ϕ_ν a set of orthogonal functions,

$$\hat{f}(t) = \sum_{\nu=0}^M \hat{f}_\nu \phi_\nu(t) \quad (9)$$

with

$$\hat{f}_\nu = E\phi_\nu(X) = \frac{1}{n} \sum_{i=1}^n \phi_\nu(X_i) \quad (10)$$

$$MSE(\hat{f}) = E[\hat{f}(x) - f(x)]^2 = [E\hat{f}(x) - f(x)]^2 + \text{var}\hat{f}(x) \quad (11)$$

$$MISE(\hat{f}) = \int MSE dx^3 \quad (12)$$

$$E\hat{f}(x) = \int \frac{1}{h} K\left(\frac{x-y}{h}\right) f(y) dy \quad (13)$$

$$n\text{var}\hat{f}(x) = \int \frac{1}{h^2} K\left(\frac{x-y}{h}\right)^2 f(y) dy - \left(\int \frac{1}{h} K\left(\frac{x-y}{h}\right) f(y) dy \right)^2 \quad (14)$$

³Also Rosenblatt, 1956

Suppose that the kernel satisfies

$$\int K(t)dt = 1 \quad (15)$$

$$\int tK(t)dt = 0 \quad (16)$$

$$\int t^2K(t)dt = k_2 \neq 0 \quad (17)$$

Doing a change of variables $y = x - ht$

$$bias = E\hat{f}(x) - f(x) = \int K(t)[f(x - ht) - f(x)]dt \quad (18)$$

Then use a Taylor's expansion on f to find

$$bias = \frac{1}{2}h^2 f''(x)k_2 + \text{higher order terms in } h \quad (19)$$

Integrating

$$\int bias(x)^2 dx \approx \frac{1}{4}h^4 k_2^2 \int f''(x)^2 dx \quad (20)$$

Using $y = x - ht$ and the bias approximation from the last slide

$$\begin{aligned} \text{var} \hat{f}(x) &= \frac{1}{n} \int \frac{1}{h^2} K \left(\frac{x-y}{h} \right)^2 f(y) dy - \frac{1}{n} [f(x) + \text{bias}(x)]^2 \\ &\approx \frac{1}{hn} \int f(x - ht) K(t)^2 dt - \frac{1}{n} [f(x) + O(h^2)] \end{aligned} \quad (21)$$

Assuming h is small, n is large and using Taylor's expansion,

$$\text{var} \hat{f}(x) \approx \frac{1}{nh} f(x) \int K(t)^2 dt \quad (22)$$

And integrating

$$\int \text{var} \hat{f}(x) dx \approx \frac{1}{nh} \int K(t)^2 dt \quad (23)$$

So far, we have

$$MISE \approx \frac{1}{4} h^4 k_2^2 \int f''(x)^2 dx + \frac{1}{nh} \int K(t)^2 dt \quad (24)$$

Somehow there is an optimal value of h that minimizes $MISE$,⁴

$$h_{opt} = k_2^{-2/5} n^{-1/5} \left(\int K(t)^2 dt \right)^{1/5} \left(\int f''(x)^2 dx \right)^{-1/5} \quad (25)$$

Substituting,

$$MISE(h_{opt}) = \frac{4}{5} C(K) \left(\int f''(x)^2 dx \right)^{1/5} n^{-4/5} \quad (26)$$

$$\text{with } C(K) = k_2^{2/5} \left(\int K(t)^2 dt \right)^{4/5} \quad (27)$$

⁴"On estimation of a probability density function and mode", by Parzen, 1962 (Lemma 4A)

So choose a Kernel that minimizes $C(K)$. Assume $k_2 = 1$ (or rescale) to see this as minimization problem with constraints. It has a solution.

$$K_e(t) = \begin{cases} \frac{3}{4\sqrt{5}} \left(1 - \frac{t^2}{5}\right) & \text{if } -\sqrt{5} \leq t \leq \sqrt{5} \\ 0 & \text{otherwise.} \end{cases}$$

K_e is called the *Epanechnikov*^{5 6} kernel and it is used to measure efficiency of different kernels by computing

$$\text{eff}(K) = \left(\frac{C(K_e)}{C(K)} \right)^{5/4} \quad (28)$$

⁵ "Nonparametric estimation of a multidimensional probability density", by Epanechnikov, 1969

⁶ "The efficiency of some nonparametric competitors of the t-test" by Hodges and Lehmann, 1956

Table 3.1 *Some kernels and their efficiencies*

<i>Kernel</i>	$K(t)$	<i>Efficiency (exact and to 4 d.p.)</i>
Epanechnikov	$\frac{3}{4}(1 - \frac{1}{5}t^2)/\sqrt{5}$ for $ t < \sqrt{5}$, 0 otherwise	1
Biweight	$\frac{15}{16}(1 - t^2)^2$ for $ t < 1$ 0 otherwise	$\left(\frac{3087}{3125}\right)^{1/2} \approx 0.9939$
Triangular	$1 - t $ for $ t < 1$, 0 otherwise	$\left(\frac{243}{250}\right)^{1/2} \approx 0.9859$
Gaussian	$\frac{1}{\sqrt{2\pi}}e^{-(1/2)t^2}$	$\left(\frac{36\pi}{125}\right)^{1/2} \approx 0.9512$
Rectangular	$\frac{1}{2}$ for $ t < 1$, 0 otherwise	$\left(\frac{108}{125}\right)^{1/2} \approx 0.9295$

Look at ISE ⁷⁸

$$\int (\hat{f} - f)^2 = \int \hat{f}^2 - 2 \int \hat{f}f + \int f^2. \quad (29)$$

Define

$$R(\hat{f}) = \int \hat{f}^2 - 2 \int \hat{f}f \quad (30)$$

$$\hat{f}_{-i}(x) = \frac{1}{h(n-1)} \sum_{j \neq i} K\left(\frac{x - X_j}{h}\right) \quad (31)$$

$$M_0(h) = \int \hat{f}^2 - \frac{2}{n} \sum_i \hat{f}_{-i}(X_i) \quad (32)$$

⁷ "Empirical choice of histograms and kernel density estimators", by Rudemo, 1982

⁸ "An alternative method of cross-validation for the smoothing of density estimates", by Bowman, 1984

Noting that

$$E \frac{1}{n} \sum_i \hat{f}_{-i}(X_i) = E \hat{f}_{-n}(X_n) = E \int \hat{f}(x) f(x) dx \quad (33)$$

implies

$$EM_0(h) = ER(\hat{f}). \quad (34)$$

So $M_0(h) + \int f^2$ is an unbiased estimator for *MISE*. Assuming the minimizers of M_0 and EM_0 are close, a good choice of h is the one that minimizes M_0 .

M_0 can be rewritten into a more computer-friendly manner.

Assuming that ⁹

- K (a bounded Borel function) satisfies $\int |K(t)|dt < \infty$ and $\int K(t)dt = 1$
- $|tK(t)| \rightarrow 0$ as $|t| \rightarrow \infty$
- h (that depends on n) satisfies $h \rightarrow 0$ and $nh \rightarrow \infty$ as $n \rightarrow \infty$

Then $\hat{f}(x) \rightarrow f(x)$ in probability as $n \rightarrow \infty$

⁹"On estimation of a probability density function and mode", by Parzen, 1962

Assuming that ¹⁰

- K is bounded, has bounded variation, the set of discontinuities has measure zero
- f is uniformly continuous on $(-\infty, \infty)$
- h (that depends on n) satisfies $h \rightarrow 0$ and $nh(\log n)^{-1} \rightarrow \infty$ as $n \rightarrow \infty$

Then $\sup_x |\hat{f}(x) - f(x)| \rightarrow 0$ with probability 1 $n \rightarrow \infty$. Conditions are necessary as well as sufficient.

¹⁰ "Convergence uniforme d'un estimateur de la densité par la méthode de noyau", by Bertrand-Retali, 1978

Assuming that ¹¹

- K is a non-negative Borel function that integrates to 1
- h (that depends on n) satisfies $h \rightarrow 0$ and $nh \rightarrow \infty$ as $n \rightarrow \infty$

Then $\int |\hat{f}(x) - f(x)| dx \rightarrow 0$ with probability 1 $n \rightarrow \infty$. Conditions are necessary as well as sufficient and no assumption is made on f .

¹¹ "Non-parametric density estimation: The L_1 view", by Devroye and Györfi, 1985

To choose the ideal window width h , you can make an assumption over f'' and draw an estimate from the h_{opt} expression, but that leads to massive oversmoothing.¹²

If the number of samples n is fixed there is no guarantee of convergence of any kind. Using Gaussian kernels, for any f in $L^2(\mathbb{R})$ and any $\epsilon > 0$ there exists $t > 0$, $N \in \mathbb{N}$ and $a_n \in \mathbb{R}$ such that¹³

$$\left\| f - \sum_{n=0}^N a_n e^{-(x-nt)^2} \right\|_2 < \epsilon \quad (35)$$

¹² "Density Estimation for Statistics and Data Analysis", by Silverman, 1986

¹³ "Approximating with Gaussians", by Calcaterra and Boldt, 2008

