

# CS207 Algorithm Analysis: Kernel Density Estimation

Gioia Dominedò    Nicolas Drizard    Kendrick Lo    Malcolm Mason Rodriguez

May 4, 2016

## 1 Introduction

Kernel Density Estimation (KDE) is a non-parametrical statistical technique that is used to estimate the probability density function (PDF) of a random variable based on a finite data sample<sup>1</sup>. Mathematically, we can denote the KDE of the PDF of independent and identically distributed (*iid*) random variables  $X = \{x_1, x_2, \dots, x_n\}$  as:

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

where  $K(\cdot)$  represents the *kernel* and  $h$  represents the *bandwidth* (i.e. a smoothing parameter). We note that  $\hat{f}_h(x)$  denotes the estimated PDF at a particular point  $x$ , and that the entire data set  $X$  is used to calculate it.

Various kernels are used in practice (e.g. Gaussian, Epanechnikov, exponential, etc.), all of which meet the criteria of being non-negative functions with mean zero that integrate to one. Similarly, different (strictly positive) values of  $h$  can be chosen depending on the data set in question. Figure 1 illustrates how the choice of these two parameters affects the shape of the estimated PDF.

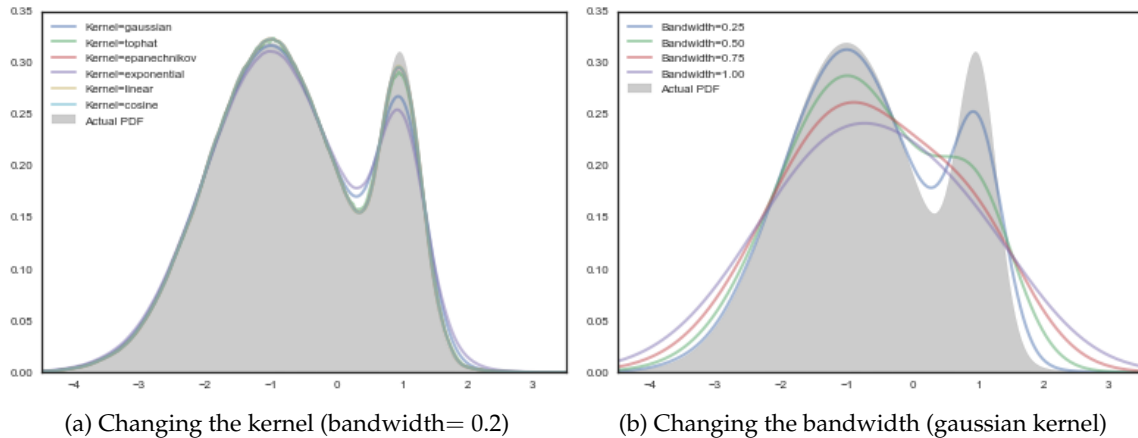


Figure 1: Impact of KDE parameters

<sup>1</sup>[https://en.wikipedia.org/wiki/Kernel\\_density\\_estimation](https://en.wikipedia.org/wiki/Kernel_density_estimation)

## 2 Algorithm

A discussion on the “best” parameters for KDE is beyond the scope of this paper. For ease of comparability, all our experiments are based on a Gaussian kernel and a bandwidth of 0.2. Figure 2 shows the result of running the Scikit Learn KDE algorithm with these parameters on 100K one-dimensional random variables drawn from a bimodal distribution, where the PDF is plotted at 1K evenly-spaced values of  $x$ .

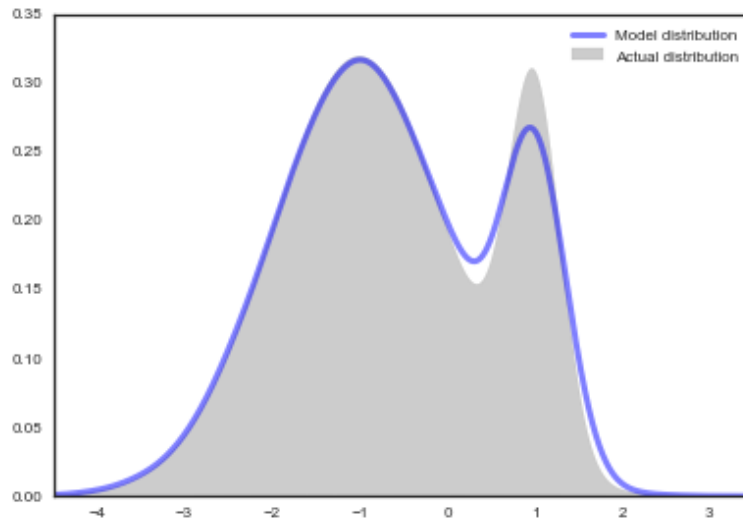


Figure 2: Sample KDE output

We started by implementing a “pure” Python version of KDE; in keeping with that objective, we only used lists. We can see from the code below that the KDE algorithm is a direct application of Equation 1. The double loop, which is necessary due to the use of lists, immediately stands out as a significant source of inefficiencies.

```
def gaussian(x):  
    """  
    Returns the Gaussian function for a given input variable.  
  
    Parameters  
    -----  
    x : float  
        the x-value at which to evaluate the Gaussian function  
  
    Returns  
    -----  
    gaussian : float  
        the result of evaluating the Gaussian function at x  
    """  
    return np.exp(-1.0 * (x**2) / 2) * (1.0 / np.sqrt(2.0 * np.pi))
```

```

def naive_kde(x, x_grid, h):
    """
    Calculates the KDE PDF estimate, based on a Gaussian kernel.

    Parameters
    -----
    x : list
        iid random variables that are used to calculate the KDE estimate
    x_grid : list
        values along which to calculate the KDE estimate
    h : float
        bandwidth parameter for the KDE estimate

    Returns
    -----
    estimates: list
        KDE estimates for all values of x_grid
    """
    global N # number of iid random variables
    estimates = [0] * len(x_grid)

    for j, gridpt in enumerate(x_grid):
        for i in range(N):
            val_in_sum = gaussian((x[i] - gridpt) / h)
            estimates[j] += val_in_sum

        estimates[j] = estimates[j] / (N * h)

    return estimates

```

Before moving on to the Scikit Learn implementation, we tested the effect of switching from lists to Numpy arrays and vectorizing the inner loop. This simple change, which is shown in the code below, reduced the runtime from nearly 300 seconds to just over 1 second!

```

def vectorized_kde(x, x_grid, h):
    """
    Calculates the KDE PDF estimate, based on a Gaussian kernel.
    Inputs
    -----
    x      : An array of iid random variables that are used to calculate the KDE estimate
    x_grid : An array of values along which to calculate the KDE estimate
    h      : Float, the bandwidth parameter for the KDE estimate
    Returns
    -----
    An array of KDE estimates for all values of x_grid.
    """
    global N # number of iid random variables

    estimates = np.zeros(len(x_grid))

    for j, gridpt in enumerate(x_grid):
        val_in_sum = gaussian((x - gridpt)/h) # returns vector of dim N
        estimates[j] = np.sum(val_in_sum) / (N * h)

    return estimates

```

### 3 Scikit Learn Optimizations

After completing our Python implementation, we next moved on to examine the Scikit Learn implementation<sup>2</sup> and analyze its optimizations.

Implementation	Runtime (seconds)
Naive Python	298.28
Vectorized Python	1.34
Scikit Learn	5.73

*Table 1: Comparing one-dimensional KDE implementations*

---

<sup>2</sup>The full source code is available at <https://github.com/scikit-learn/scikit-learn/tree/master/sklearn/neighbors>.