

DAI Assignment-3

Aditya Neeraje, Balaji Karedla, Moulik Jindal

October 10, 2024

Contents

1	Finding optimal bandwidth	1
1.1	Part 1	1
1.1.1	Part (a)	1
1.1.2	Part (b)	1
1.2	Part 2	2
1.2.1	Part (a)	2
1.2.2	Part (b)	2
1.2.3	Part (c)	2
1.2.4	Part (d)	3
1.2.5	Part (e)	3
2	Detecting Anomalous Transactions using KDE	3
2.1	Designing a custom KDE Class	3
2.2	Estimating Distribution of Transactions	3
3	Higher-Order Regression	5
3.1	Distribution of B :	5
3.2	Mean and Standard Deviation of B :	6
3.3	Part (b)	6
4	Non-parametric regression	7
4.1	Explaining the code:	7
4.2	Bandwidth corresponding to minimum estimated risk:	7
4.3	Comments on similarity and dissimilarity	8

1 Finding optimal bandwidth

1.1 Part 1

1.1.1 Part (a)

We know that the estimator for distribution function \hat{f} is given by

$$\hat{f}(x) = \sum_{j=1}^m \frac{\hat{p}_j}{h} \mathbb{1}[x \in B_j] \quad (1)$$

Also given, v_j is the number of points falling in the j^{th} bin and $\hat{p}_j = \frac{v_j}{n}$.
From this,

$$\begin{aligned} \int \hat{f}(x)^2 dx &= \int \left(\sum_{j=1}^m \frac{\hat{p}_j}{h} \mathbb{1}[x \in B_j] \right)^2 dx \\ &= \sum_{j=1}^m \int_{B_j} \left(\sum_{j=1}^m \frac{\hat{p}_j}{h} \mathbb{1}[x \in B_j] \right)^2 dx \\ &= \sum_{j=1}^m \int_{B_j} \left(\frac{\hat{p}_j}{h} \right)^2 dx \\ &= \sum_{j=1}^m \left(\frac{\hat{p}_j}{h} \right)^2 \times h \\ &= \frac{1}{n^2 h} \sum_{j=1}^m v_j^2 \end{aligned}$$

1.1.2 Part (b)

The histogram estimator after removing the i^{th} observation is given by

$$\hat{f}_{(-i)}(x) = \sum_{j=1}^m \frac{\hat{p}_{j,-i}}{h} \mathbb{1}[x \in B_j] \quad (2)$$

where $\hat{p}_{j,-i} = \frac{v_{j,-i}}{n-1}$ and $v_{j,-i}$ is the number of points falling in the j^{th} bin after removing the i^{th} observation.

Now,

$$\begin{aligned}
\sum_{i=1}^n \hat{f}_{(-1)}(X_i) &= \sum_1^n \sum_{j=1}^m \frac{\hat{p}_{j,-i}}{h} \mathbb{I}[X_i \in B_j] \\
&= \sum_{i=1}^n \sum_{j=1}^m \frac{v_{j,-i}}{(n-1)h} \mathbb{I}[X_i \in B_j] \\
&= \sum_{j=1}^m \sum_{i=1}^n \frac{v_{j,-i}}{(n-1)h} \mathbb{I}[X_i \in B_j] \\
&= \sum_{j=1}^m \frac{v_j - 1}{(n-1)h} \mathbb{I}[X_i \in B_j] \quad (\text{since whenever } \mathbb{I}[X_i \in B_j] = 1, v_{j,-i} = v_j - 1) \\
&= \sum_{j=1}^m \frac{v_j - 1}{(n-1)h} v_j \\
&= \frac{1}{(n-1)h} \sum_{j=1}^m v_j^2 - v_j
\end{aligned}$$

1.2 Part 2

1.2.1 Part (a)

The histogram of the filtered data for number of bins = 10 is shown in Figure 1.

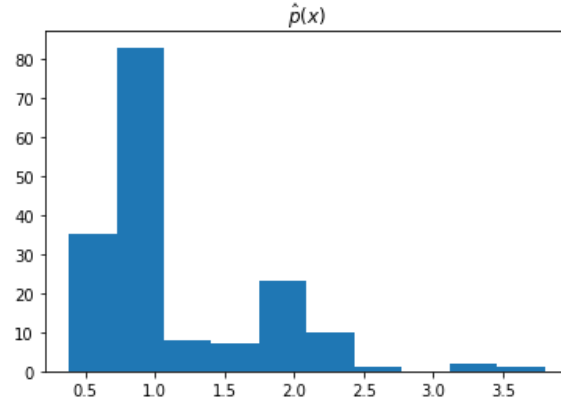


Figure 1: Histogram of filtered data with 10 bins

The values of \hat{p}_j are as follows:

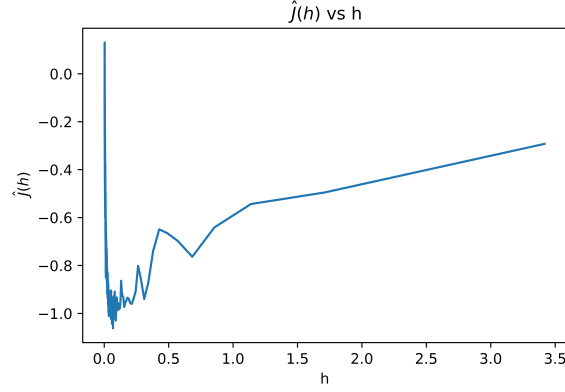
Bin	1	2	3	4	5	6	7	8	9	10
\hat{p}_j	0.602	1.428	0.138	0.12	0.396	0.172	0.017	0.0	0.034	0.017

1.2.2 Part (b)

The probability distribution is underfit as the number of bins is too low. The distribution is not smooth and the bins are too wide. This is evident from the histogram in Figure 1.

1.2.3 Part (c)

The graph of $\hat{J}(h)$ vs h for the number of bins ranging from 1 to 1000 is shown in Figure 2.

Figure 2: $\hat{J}(h)$ vs h

1.2.4 Part (d)

Based on the Cross Validation plot in Figure 2, the optimal number of bins is 50 with a value of $h^* = 0.06836$.

1.2.5 Part (e)

The histogram of the filtered data for number of bins = 50 is shown in Figure 3.

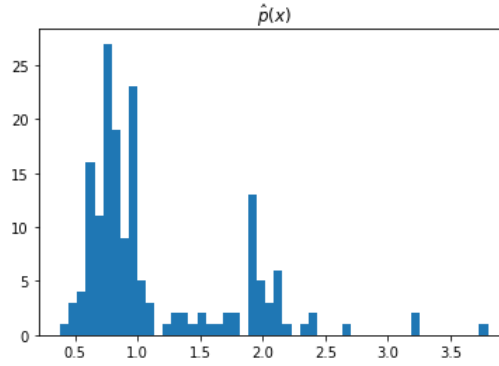


Figure 3: Histogram of filtered data with 50 bins

2 Detecting Anomalous Transactions using KDE

2.1 Designing a custom KDE Class

The KDE class with the Epachnikov kernel

$$K(x) = \frac{3}{4}(1 - \|x\|_2^2) \quad \text{if } \|x\|_2 \leq 1 \quad (3)$$

The code was implemented in the file `2.py`.

2.2 Estimating Distribution of Transactions

After estimating the distribution of transactions using the Epachnikov kernel, the estimated probability density graph is shown in Figure 4.

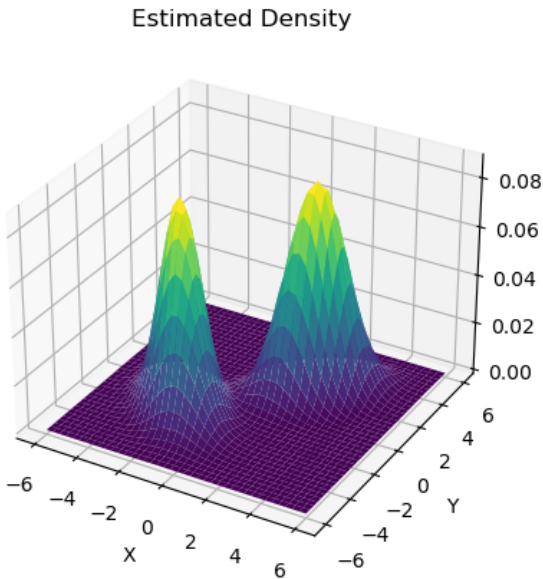


Figure 4: Estimated Probability Density of Transactions

3 Higher-Order Regression

3.1 Distribution of B :

Let us suppose we are trying to fit the dataset to a functional equation of the form:

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_m x^m + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We have to find the value of B (the estimator of $(\beta_0, \beta_1, \dots, \beta_m)$) that minimizes

$$\sum_{i=1}^n (Y_i - B_0 - B_1 x_i - B_2 x_i^2 - \cdots - B_m x_i^m)^2$$

Taking the derivative w.r.t each of the B'_i s, we get

$$\begin{aligned} \sum_{i=1}^n Y_i &= nB_0 + B_1 \sum_{i=1}^n x_i + B_2 \sum_{i=1}^n x_i^2 + \cdots + B_m \sum_{i=1}^n x_i^m \\ \sum_{i=1}^n x_i \cdot Y_i &= B_0 \sum_{i=1}^n x_i + B_1 \sum_{i=1}^n x_i^2 + B_2 \sum_{i=1}^n x_i^3 + \cdots + B_m \sum_{i=1}^n x_i^{m+1} \\ &\vdots \\ \sum_{i=1}^n x_i^m \cdot Y_i &= B_0 \sum_{i=1}^n x_i^m + B_1 \sum_{i=1}^n x_i^{m+1} + B_2 \sum_{i=1}^n x_i^{m+2} + \cdots + B_m \sum_{i=1}^n x_i^{2m} \end{aligned}$$

Let us suppose that the matrix X is defined as $X_{ij} = x_i^j$. Note that $A = X^T \cdot X$ is a matrix that satisfies $A_{ij} = \sum_{k=1}^n x_k^{i+j}$ where A_{ij} is the element in the i^{th} row and j^{th} column of A . Let Y be the vector of Y_i 's. Then, the above equations can be written as:

$$X^T \cdot Y = (X^T \cdot X) \cdot B$$

We claim that $X^T \cdot X$ is invertible as long as all values of x are distinct. To prove this, we first prove X is invertible. If that is the case, then X^T , which has the same determinant as X , is invertible, and their product also has a non-zero determinant and is invertible. For the proof of the claim that X is invertible, let us assume for the sake of contradiction that the columns of this matrix are not linearly independent.

That is, let us assume $c_0 v_0 + c_1 v_1 + \cdots + c_m v_m = 0$ where the c_i 's are the columns of X and v_i 's are real numbers, not all zero. Then we get on the k_{th} coordinate (for $k \in \{1, \dots, n\}$):

$$v_0 + v_1 x_k + \cdots + v_m x_k^m = 0$$

which means x_k is a root of the polynomial $v_0 + v_1 x + \cdots + v_m x^m$. Since this polynomial has at most m roots, we get that two of the x_k 's must be equal, which contradicts our assumption that all x 's are distinct.

Hence, $X^T \cdot X$ is invertible, and we get that B is given by:

$$B = (X^T \cdot X)^{-1} \cdot X^T \cdot Y$$

Thus, $B = C \cdot Y$ for some matrix C , with m rows and n columns.

$B_{i-1} = \sum_{j=1}^n C_{ij} Y_j$ is a linear combination of Gaussian random variables, and hence is Gaussian. B is a $(m+1)$ -tuple of Gaussian random variables.

3.2 Mean and Standard Deviation of B :

We know $Y = \beta X + \epsilon$, where $\beta = (\beta_0, \beta_1, \dots, \beta_m)$ and $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

$$\begin{aligned} E[B] &= E[(X^T \cdot X)^{-1} \cdot X^T \cdot Y] \\ &= (X^T \cdot X)^{-1} \cdot X^T \cdot E[Y] \\ &= (X^T \cdot X)^{-1} \cdot X^T \cdot \beta X \\ &= (X^T \cdot X)^{-1} \cdot X^T \cdot X \cdot \beta \\ &= \beta \end{aligned}$$

Above, notice that we have used the fact that the component-wise computation of the mean of B can be done simultaneously.

To find the variance, we again use matrix C as defined above. $B_{i-1} = \sum_{k=1}^n C_{ik} Y_k$ and $B_{j-1} = \sum_{k=1}^n C_{jk} Y_k$.

Hence, $\text{Cov}(B_{i-1}, B_{j-1}) = \text{Cov}(\sum_{k=1}^n C_{ik} Y_k, \sum_{k=1}^n C_{jk} Y_k) = \sum_{r=1}^n \sum_{l=1}^n C_{il} C_{jr} \text{Cov}(Y_l, Y_r)$.

Now, Y_l and Y_r are independent for $l \neq r$, and hence $\text{Cov}(Y_l, Y_r) = 0$ for $l \neq r$ and $\text{Var}(Y_l)$ if $l = r$.

Since $\text{Var}(Y_l) = \sigma^2$, we get that $\text{Cov}(B_{i-1}, B_{j-1}) = \sigma^2 \sum_{k=1}^n C_{ik} C_{jk}$.

The final term here is equal to the $(i, j)^{\text{th}}$ element of the matrix $C \cdot C^T$.

Thus, $\text{Cov}(B) = \sigma^2 C \cdot C^T$.

Now,

$$\begin{aligned} C^T &= ((X^T \cdot X)^{-1} \cdot X^T)^T \\ &= X \cdot ((X^T \cdot X)^{-1})^T \\ &= X \cdot (X^T \cdot X)^{-1} \end{aligned}$$

Here, the last point follows from the symmetry of $X^T \cdot X$, which implies that its inverse is also a symmetric matrix.

The above statement can be proven as follows: suppose Y is a symmetric invertible matrix, then $Y \cdot Y^{-1} = I = (Y^{-1})^T \cdot Y^T = (Y^{-1})^T \cdot Y$.

Using the fact that $Y^{-1} \cdot Y = I = (Y^{-1})^T \cdot Y^T$, we get that Y^{-1} is symmetric, since inverses are unique.

Since $\text{Cov}(B_i, B_i) = \text{Var}(B_i)$, we can get the variance of B_i as the i^{th} diagonal element of $\sigma^2 X \cdot (X^T \cdot X)^{-1}$.

The quantity σ^2 can be estimated using the sum of squares of the residuals. That is, if we let

$$SS_R = \sum_{i=1}^n (Y_i - B_0 - B_1 x_i - B_2 x_i^2 - \dots - B_m x_i^m)^2$$

then, it can be shown that $\sigma^2 = \frac{SS_R}{n-m-1}$. This is because $\frac{SS_R}{\sigma^2} \sim \chi_{n-(k+1)}^2$.

3.3 Part (b)

In computing the mean of B_i above, we concluded that $E[B_i] = \beta_i$.

Thus, $E[B] = \beta$, where β is as defined before. This is the same as the true value of the coefficients, hence B is an unbiased estimator.

4 Non-parametric regression

4.1 Explaining the code:

The code contains a `NadarayaWatsonRegressor` class and three kernel types that can be used - Gaussian, Epachnikov and Square (Tophat). Let us suppose our model is called “regressor”. We first need regressor to store our data, which is done by calling `regressor.fit(X, Y)`. `regressor.display(bandwidth)` plots the estimated KDE for the particular bandwidth.

Noticing that the data we have been given has a range of x values which are greater than 0 and less than 4, I have taken these as my limits for plotting.

Within the range of $\min(X)$ to $\max(X)$, we plot values corresponding to the y_labels given (note again that I am using the fact that the x values in the graph are already very dense, in order to save on computation). Outside this range, I plot the estimated y every 0.001 units.

`regressor.plot_cross_validator(bandwidths)` can be used to generate a plot of cross-validation scores. Assumptions made here include the assumption that bandwidths is an array of non-zero values and is in increasing order. While trying to find the optimal bandwidth, `plot_cross_validator()` was called with bandwidths being `np.linspace(0.1, 2, 200)`

Finally, `regressor.display_4_plots` can be used to generate the picture required by the assignment. The user needs to pass as arguments the bandwidths he wants to assign to the top left, top right and bottom left graphs respectively. For instance, I have called `regressor3.display_4_plots(0.01, 1.75, 0.32)` to represent that $h=0.01$ leads to overfitting and undersmoothing, $h=1.75$ leads to underfitting and oversmoothing and $h=0.32$ seems optimal based on cross validation scores.

4.2 Bandwidth corresponding to minimum estimated risk:

For a Gaussian kernel, the computed optimal bandwidth was 0.13.

For a Epachnikov kernel, the computed optimal bandwidth was 0.32.

For a Tophat/Square kernel, the computed optimal bandwidth was 0.31.

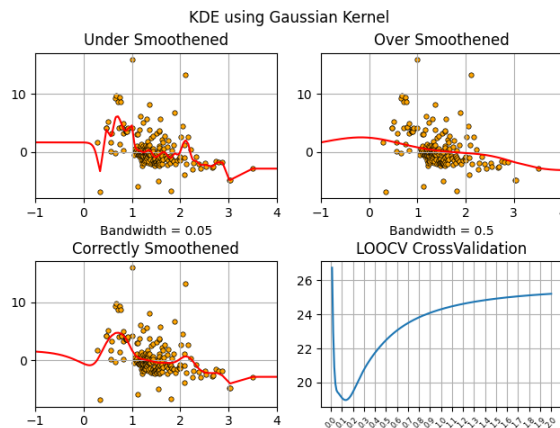


Figure 5: Gaussian Kernel Regression - undersmoothing, oversmoothing, optimal bandwidth and Cross Validation curve

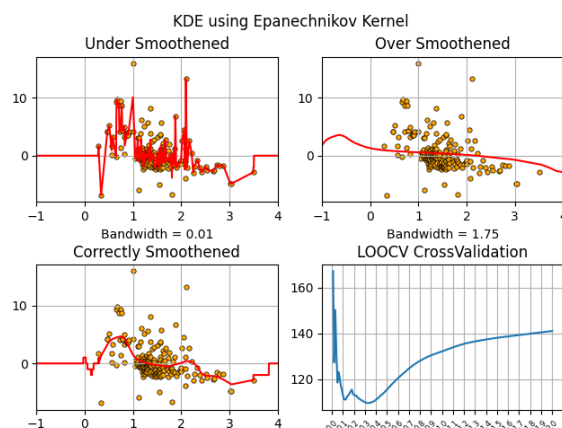


Figure 6: Epachnikov Kernel Regression - undersmoothing, oversmoothing, optimal bandwidth and Cross Validation curve

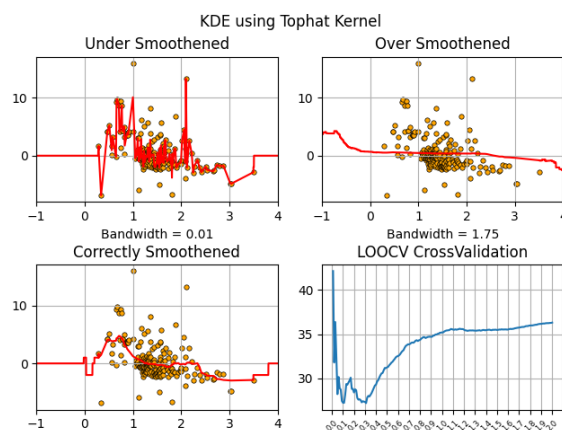


Figure 7: Tophat Kernel Regression - undersmoothing, oversmoothing, optimal bandwidth and Cross Validation curve

4.3 Comments on similarity and dissimilarity

The Gaussian estimator seems a lot smoother than the others, especially compared to Tophat, presumably because of the gaussian kernel itself being smooth and not jerky like the tophat kernel. The Gaussian estimator has a lower optimal bandwidth, presumably because it assigns a non-zero weight to all points, thus even points which are far away have some smoothing influence, whereas with Epachnikov or Tophat, both of which have finite support, we need a larger bandwidth in order to capture enough neighboring points for a smooth estimate.

The EPA and Tophat estimators do not have as smooth cross-validation curves as the Gaussian, presumably because they are abrupt jumps in the risk when a transition between bandwidths causes one object to lose all influence over the kernel at a point.