

Assignment 3

CS215: Data Structures and Algorithms

Shaik Awez Mehtab, 23B1080

Satyam Sinoliya, 23B0958

Vaibhav Singh, 23B1068

Solutions

SOLUTION 1

Detecting Anomalous Transactions using KDE

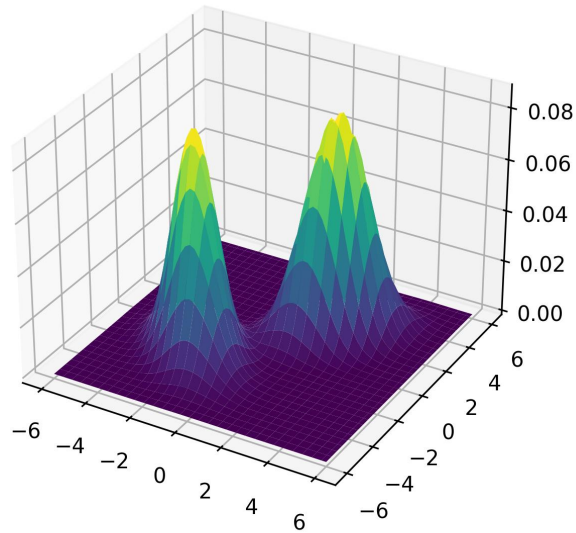


Figure 1.1: Distribution of transactions

As can be seen in the given figure, the resulting estimated distribution contains two nodes

SOLUTION 2

Higher-Order Regression

Part 1

Suppose our estimates for α and β are A and B respectively, then these values of A and B minimize

$$\sum_{i=1}^n (y_i - A - Bx_i)^2 \quad (1.1)$$

$$\Rightarrow \frac{\partial}{\partial A} \sum_{i=1}^n (y_i - A - Bx_i)^2 = 0 \quad (1.2)$$

$$\sum_{i=1}^n -2(y_i - A - Bx_i) = 0 \quad (1.3)$$

$$n\bar{y} - nA - nB\bar{x} = 0 \quad (1.4)$$

$$\bar{y} = A + B\bar{x} \quad (1.5)$$

Least square regression line is given by $y = A + Bx$. Thus by (1.5), (\bar{x}, \bar{y}) lies on the regression line.

Part 2

Suppose our estimates for β_0^* and β_1^* are A^* and B^* respectively, then A^* and B^* minimize $\sum_{i=1}^n (y_i - A^* - B^* z_i)^2$

$$\implies \frac{\partial}{\partial A^*} \sum_{i=1}^n (y_i - A^* - B^* z_i)^2 = 0 \quad \frac{\partial}{\partial B^*} \sum_{i=1}^n (y_i - A^* - B^* z_i)^2 = 0 \quad (1.6)$$

$$\sum_{i=1}^n -2(y_i - A^* - B^* z_i) = 0 \quad \sum_{i=1}^n -2z_i(y_i - A^* - B^* z_i) = 0 \quad (1.7)$$

$$n\bar{y} - nA^* - nB^*\bar{z} = 0 \quad \sum z_i y_i - A^* n\bar{z} - B^* \sum z_i^2 = 0 \quad (1.8)$$

$$\sum y_i z_i - n(\bar{y} - B^* \bar{z})\bar{z} - B^* \sum z_i^2 = 0 \quad (1.9)$$

$$B^* = \frac{\sum y_i z_i - n\bar{y}\bar{z}}{n\bar{z}^2 - \sum z_i^2} \quad A^* = \bar{y} - B^* \bar{z} \quad (1.10)$$

SOLUTION 3

Non-parametric regression

1. Report Bandwidth Corresponding to Minimum Estimated Risk

After running the Nadaraya-Watson kernel regression using the Epanechnikov and Gaussian kernel and performing cross-validation for bandwidth selection, the optimal bandwidth corresponding to the minimum estimated risk is:

Optimal Bandwidth of Gaussian kernel: 0.180

Optimal Bandwidth of Gaussian kernel: 0.164

2. Comment on Similarities and Differences Due to Choice of Different Kernel Functions

Similarities

- **General Functionality:** Both kernels assign weights to data points based on their distance from the query point, resulting in similar predictions in regions with high data density.
- **Smoothing:** As the bandwidth increases, all kernel functions produce smoother estimates. At very large bandwidths, all kernels oversmooth the data, giving too much influence to distant points.
- **Cross-validation Behavior:** Both kernels display a similar behavior during cross-validation, and the corresponding risk curves follow the same trend with bandwidth changes.

Differences

- **Shape of the Weights:**
 - **Epanechnikov Kernel:** This kernel assigns zero weight to points farther than the bandwidth due to its quadratic form, creating a more localized effect.
 - **Gaussian Kernel:** This kernel assigns non-zero weight to every point, regardless of distance, due to its exponential decay. It results in smoother estimates, but it is more sensitive to distant points.
- **Sensitivity to Outliers:**
 - **Epanechnikov Kernel:** This kernel is more resilient to outliers because they assign zero or reduced weight to distant points, decreasing the influence of outliers on the prediction.

- **Gaussian Kernel:** The Gaussian kernel is more prone to incorporating outliers, as it assigns non-zero weights even to far-away points, making it less resilient in the presence of outliers.
- **Plots**
 - **Epanechnikov Kernel:** This kernel produces more precise and localized estimates, with a good balance between bias and variance when using the optimal bandwidth.
 - **Gaussian Kernel:** The Gaussian kernel leads to smoother curves but gives undue influence to distant points, which can result in overfitting or oversmoothing depending on the bandwidth.