

Online Density Estimation of Heterogeneous Data Streams in Higher Dimensions

Michael Geilke, Andreas Karwath, and Stefan Kramer

Johannes Gutenberg University Mainz, Germany

September 22, 2016

















 \Rightarrow about 2 GBs of data

















Query:

Return the probability distribution for sensors in the living room during the week days.



$$f(X_1, \ldots, X_n)$$





$$f(X_1) \cdot \prod_{i=2}^n f(X_i | X_1, \dots, X_{i-1})$$





$$f(X_1) \cdot \prod_{i=2}^n f(X_i | X_1, \dots, X_{i-1})$$









only for discrete random variables



Goals

A density estimator that

- estimates joint densities from data streams
- is able to deal with heterogeneous data, and
- and works for higher dimensional data.

For density estimation, 100 variables is high dimensional.



Main Idea



Main Idea





 $(temperature = 20, humidity = 50, ...) = \blacksquare \in \mathbb{R}^{n}$



$$\mathbb{R}^n \ni \vec{x} = (x_1, x_2, \dots, x_n) \longrightarrow \vec{v} = (v_1, v_2, v_3, v_4) \in \mathbb{R}^m$$



$$I = \{ \vec{x} \in \mathbb{R}^n \mid h_L(\vec{x}) = \vec{v} = (v_1, v_2, v_3, v_4) \}$$

 $\hat{g}(\vec{v}) = \sum_{\vec{x} \in I} \hat{f}(\vec{x})$







- representative
- instance

$$(temperature = 10, humidity = 80, ...) = \bullet \in \mathbb{R}^n$$







- representative
- instance





Mahalanobis distance:
$$\sqrt{(\vec{x} - \vec{v})^T \Sigma^{-1} (\vec{x} - \vec{v})}$$



$$\vec{v} = (v_1, v_2, v_3, v_4)$$





$$\vec{v} = (v_1, v_2, v_3, v_4)$$











Choice of Landmarks

Main idea:

- theoretical foundation
- landmarks are orthogonal to each other
- if |*L*|= d + 1, then consistent estimator
- back translation by system of linear equations



Evaluation: Parameter Setting

Parameters:

- $\theta_{C \to R} = 100$
- Euclidean norm
- $|L| \in \{2, 3, 5, 10, 20\}$
- $M \in \{0.1, 0.5, 1.0, 2.0, 5.0, 10.0\}$

Datasets
Synthetic
Gaussian mixtures
Real-World
Covertype Electricity Letter Shuttle



Online Density Estimation of Heterogeneous Data Streams in Higher Dimensions







Evaluation: Parameter Setting



- |L| depends on dimensionality of data
- small M partition the space better
- but at some point too few instances per region



Evaluation: Performance

oKDE:

- online Kernel Density Estimator
- for multi-variate densities
- for continuous variables
- by Kristan et al. (2011)

Datasets
Synthetic
Gaussian mixtures
Real-World
Covertype Electricity Letter Shuttle



Online Density Estimation of Heterogeneous Data Streams in Higher Dimensions



Conclusions

- online density estimation in higher dimensions
- heterogeneous data stream
- theoretical foundation
- comparable to the state of the art

Future Work:

- new strategies for landmarks selection
- outlier detection
- detection of emerging trends





Thank you for your attention















instance

$$I = \{ \vec{x} \in \mathbb{R}^n \mid h_L(\vec{x}) = \vec{v} = (v_1, v_2, v_3, v_4) \}$$

