

# **Artificial Intelligence in Industry**

Last update: 24 September 2024

Academic Year 2024 – 2025

Alma Mater Studiorum · University of Bologna

# Contents

<b>1</b>	<b>Preliminaries</b>	<b>1</b>
<b>2</b>	<b>Low dimensional anomaly detection: Taxi calls</b>	<b>2</b>
2.1	Data . . . . .	2
2.2	Approaches . . . . .	2
2.2.1	Gaussian assumption . . . . .	2
2.2.2	Characterize data distribution . . . . .	2
2.2.2.1	Univariate kernel density estimation . . . . .	3
2.2.2.2	Multivariate kernel density estimation . . . . .	4
2.2.2.3	Time-dependent estimator . . . . .	6
2.2.2.4	Time-indexed model . . . . .	7
<b>3</b>	<b>High dimensional anomaly detection: HPC centers</b>	<b>8</b>
3.1	Data . . . . .	8
3.1.1	High-dimensional data visualization . . . . .	8
3.2	Approaches . . . . .	9
3.2.1	Multivariate KDE . . . . .	9

# 1 Preliminaries

**Problem formalization** Defines the ideal goal.

**Solution formalization** Defines the actual possible approaches to solve a problem.

**Occam's razor** Principle for which, between two hypotheses, the simpler one is usually correct.

|**Remark.** This approach has less variance and more bias, making it more robust.

Problem  
formalization  
Solution  
formalization  
Occam's razor

## 2 Low dimensional anomaly detection: Taxi calls

**Anomaly** Event that deviates from the usual pattern.

Anomaly

**Time series** Data with an ordering (e.g., chronological).

Time series

### 2.1 Data

The dataset is a time series and it is a **DataFrame** with the following fields:

**timestamp** with a 30 minutes granularity.

**value** number of calls.

The label is a **Series** containing the timestamps of the anomalies.

An additional **DataFrame** contains information about the time window in which the anomalies happen:

**begin** acceptable moment from which an anomaly can be detected.

**end** acceptable moment from which there are no anomalies anymore.

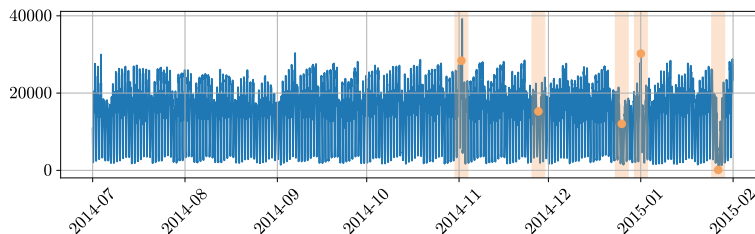


Figure 2.1: Plot of the time series, anomalies, and windows

### 2.2 Approaches

#### 2.2.1 Gaussian assumption

Assuming that the data follows a Gaussian distribution, mean and variance can be used to determine anomalies through a threshold.  $z$ -score can also be used.

#### 2.2.2 Characterize data distribution

Classify a data point as an anomaly if it is too unlikely.

**Problem formalization** Given a random variable  $X$  with values  $x$  to represent the number of taxi calls, we want to find its probability density function (PDF)  $f(x)$ .

An anomaly is determined whether:

$$f(x) \leq \varepsilon$$

where  $\varepsilon$  is a threshold.

**Remark.** A PDF can be reasonably used even though the dataset is discrete if its data points are sufficiently fine-grained.

**Remark.** It is handy to use negated log probabilities as:

- The logarithm adds numerical stability.
- The negation makes the probability an alarm signal, which is a more common measure.

Therefore, the detection condition becomes:

$$-\log f(x) \geq \varepsilon$$

**Solution formalization** The problem can be tackled using a density estimation technique.

### 2.2.2.1 Univariate kernel density estimation

**Kernel density estimation (KDE)** Based on the assumption that whether there is a data point, there are more around it. Therefore, each data point is the center of a density kernel.

Kernel density estimation (KDE)

**Density kernel** A kernel  $K(x, h)$  is defined by:

- The input variable  $x$ .
- The bandwidth  $h$ .

**Gaussian kernel** Kernel defined as:

$$K(x, h) \propto e^{-\frac{x^2}{2h^2}}$$

where:

- The mean is 0.
- $h$  is the standard deviation.

As the mean is 0, an affine transformation can be used to center the kernel on a data point  $\mu$  as  $K(x - \mu, h)$ .

Given  $m$  training data points  $\bar{x}_i$ , the density of any point  $x$  can be computed as the kernel average:

$$f(x, \bar{x}, h) = \frac{1}{m} \sum_{i=0}^m K(x - \bar{x}_i, h)$$

Therefore, the train data themselves are used as the parameters of the model while the bandwidth  $h$  has to be estimated.

**Data split** Time series are usually split chronologically:

**Train** Should ideally contain only data representing the normal pattern. A small amount of anomalies might be tolerated as they have low probabilities.

**Validation** Used to find the threshold  $\varepsilon$ .

**Test** Used to evaluate the model.

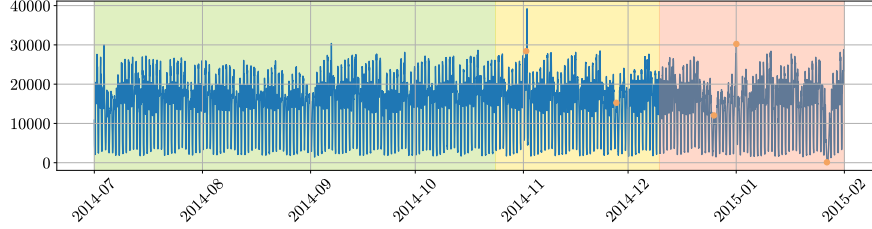


Figure 2.2: Train, validation, and test splits

**Metrics** It is not straightforward to define a metric for anomaly detection. A cost model to measure the benefits of a prediction is more suited. A simple cost model can be based on:

**True positives (TP)** Windows for which at least an anomaly is detected;

**False positives (FP)** Detections that are not actually anomalies;

**False negatives (FN)** Undetected anomalies;

**Advance (adv)** Time between an anomaly and when it is first detected;

and is computed as:

$$(c_{\text{false}} \cdot \text{FP}) + (c_{\text{miss}} \cdot \text{FN}) + (c_{\text{late}} \cdot \text{adv}_{\leq 0})$$

where  $c_{\text{false}}$ ,  $c_{\text{miss}}$ , and  $c_{\text{late}}$  are hyperparameters.

**Bandwidth estimation** According to some statistical arguments, a rule-of-thumb to estimate  $h$  in the univariate case is the following:

$$h = 0.9 \cdot \min \left\{ \hat{\sigma}, \frac{\text{IQR}}{1.34} \right\} \cdot m^{-\frac{1}{5}}$$

where:

- IQR is the inter-quartile range.
- $\hat{\sigma}$  is the standard deviation computed over the whole dataset.

**Threshold optimization** Using the train and validation set, it is possible to find the best threshold  $\varepsilon$  that minimizes the cost model through linear search.

**Remark.** The train set can be used alongside the validation set to estimate  $\varepsilon$  as this operation is not used to prevent overfitting.

**Remark.** The evaluation data should be representative of the real world distribution. Therefore, in this case, to evaluate the model the whole dataset can be used.

**Remark.** KDE assumes that the Markov property holds. Therefore, each data point is considered independent to the others.

#### 2.2.2.2 Multivariate kernel density estimation

**Remark.** In this dataset, nearby points tend to have similar values.

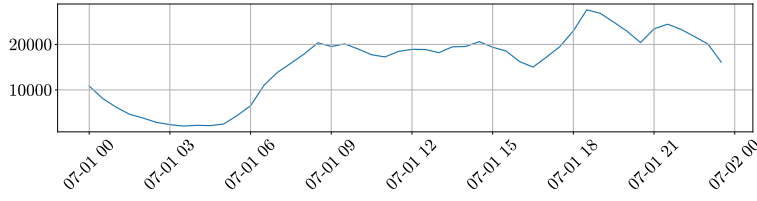


Figure 2.3: Subset of the dataset

**Autocorrelation plot** Plot to visualize the correlation between nearby points of a series. Given the original series, it is duplicated, shifted by a lag  $l$ , and the Pearson correlation coefficient is then computed between the two series. This operation is repeated over different values of  $l$ .

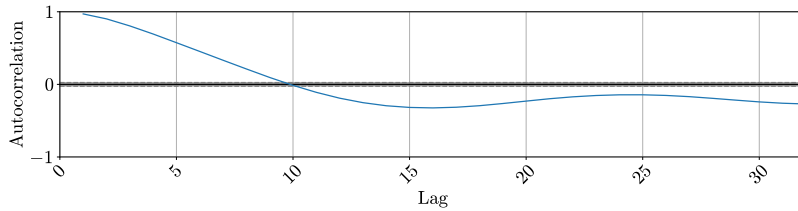


Figure 2.4: Autocorrelation plot of the subset of the dataset.

There is strong correlation up to 4-5 lags.

**Sliding window** Given a window size  $w$  and a stride  $s$ , the dataset is split into sequences of  $w$  continuous elements.

**Remark.** Incomplete sequences at the start and end of the dataset are ignored.

**Remark.** In `pandas`, the `rolling` method of `Dataframe` allows to create a slicing window iterator. This approach creates the windows row-wise and also considers incomplete windows. However, a usually more efficient approach is to construct the sequences column-wise by hand.

**Multivariate KDE** Extension of KDE to vector variables.

**Window size estimation** By analyzing the autocorrelation plot, an ideal window size can be picked as the lag with a low correlation (e.g., 10 according to Figure 2.4).

**Bandwidth estimation** Differently from the univariate case, the bandwidth has to be estimated by maximizing the log-likelihood on the validation set. Given:

- The validation set  $\tilde{x}$ ,
- The observations  $x$ ,
- The bandwidth  $h$ ,
- The density estimator  $\hat{f}$

the likelihood is computed, by assuming independent observations, as:

$$L(h, x, \tilde{x}) = \prod_{i=1}^m \hat{f}(x_i, \tilde{x}_i, h)$$

Maximum likelihood estimation is defined as:

$$\arg \max \mathbb{E}_{x \sim f(x), \tilde{x} \sim f(x)} [L(h, x, \tilde{x})]$$

where  $f(x)$  is the true distribution.

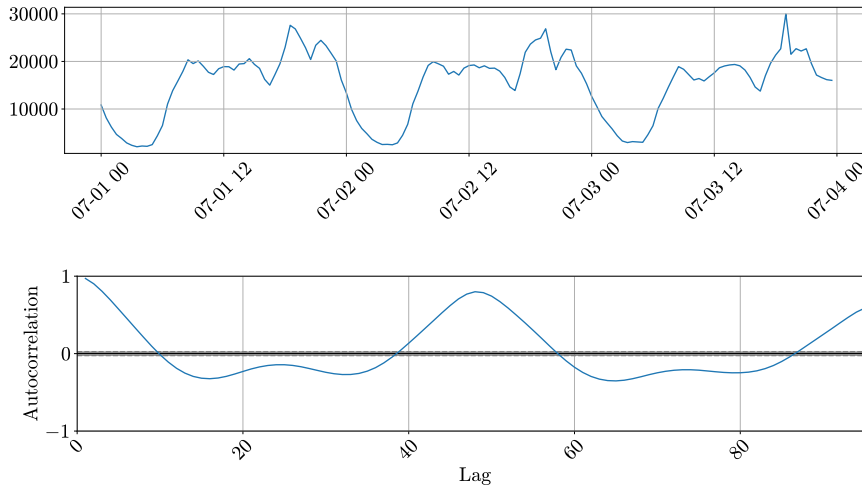
In practice, the expected value is sampled from multiple validation and training sets through cross-validation and grid-search.

**Remark.** As different folds for cross-validation must be tried, grid-search is an expensive operation.

**Threshold optimization** The threshold can be determined as in Section 2.2.2.1.

### 2.2.2.3 Time-dependent estimator

The approach taken in Sections 2.2.2.1 and 2.2.2.2 is based on the complete timestamp of the dataset. Therefore, the same times on different days are treated differently. However, it can be seen that this time series is approximately periodic, making the approach used so far more complicated than necessary.



**Time input** It is possible to take time into consideration by adding it as a parameter of the density function:

$$f(t, x)$$

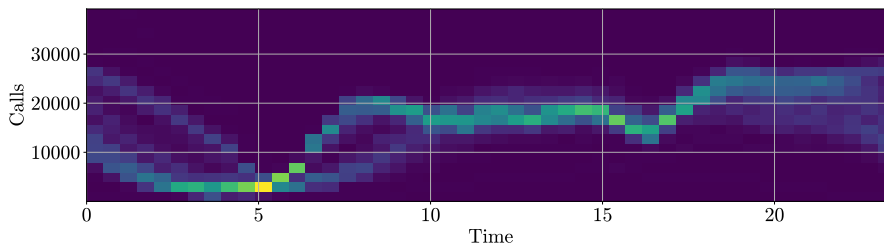


Figure 2.5: 2D histogram of the distribution. Lighter colors indicate a higher frequency of occurrence.



However,  $t$  is a controlled variable and is completely predictable (i.e., if times are samples with different frequencies, they should not affect the overall likelihood). Therefore, a conditional density should be used:

$$f(x | t) = \frac{f(t, x)}{f(t)}$$

where  $f(t)$  can be easily computed using KDE on the time data.

**Remark.** For this dataset, the time distribution is uniform. Therefore,  $f(t)$  is a constant and it is correct to use  $f(t, x)$  as the estimator.

**Bandwidth estimation** The bandwidth can be estimated as in Section 2.2.2.2.

**Remark.** When using the `scikit-learn`, the dataset should be normalized as the implementation of KDE use the same bandwidth for all features.

**Threshold optimization** The threshold can be determined as in Section 2.2.2.1.

#### 2.2.2.4 Time-indexed model

**Ensemble model** Model defined as:

Ensemble model

$$f_{g(t)}(x)$$

where:

- $f_i$  are estimators, each working on subsets of the dataset and solving a smaller problem.
- $g(t)$  determines which particular  $f_i$  should solve the input.

**Time-indexed model** Consider both time and sequence inputs by using an ensemble model. Each estimator is specialized on a single time value (i.e., an estimator for 00:00, one for 00:30, ...).

**Bandwidth estimation** The bandwidth can be estimated as in Section 2.2.2.2.

**Threshold optimization** The threshold can be determined as in Section 2.2.2.1.

## 3 High dimensional anomaly detection: HPC centers

### 3.1 Data

The dataset is a time series with the following fields:

`timestamp` with a 5 minutes granularity.

**HPC data** technical data related to the cluster.

`anomaly` indicates if there is an anomaly.

In practice, the cluster has three operational modes:

**Normal** the frequency is proportional to the workload.

**Power-saving** the frequency is always at the minimum.

**Performance** the frequency is always at the maximum.

For this dataset, both power-saving and performance are considered anomalies.

#### 3.1.1 High-dimensional data visualization

**Individual plots** Plot individual columns.

**Statistics** Show overall statistics (e.g., `pandas describe` method)

**Heatmap** Heatmap with standardized data.

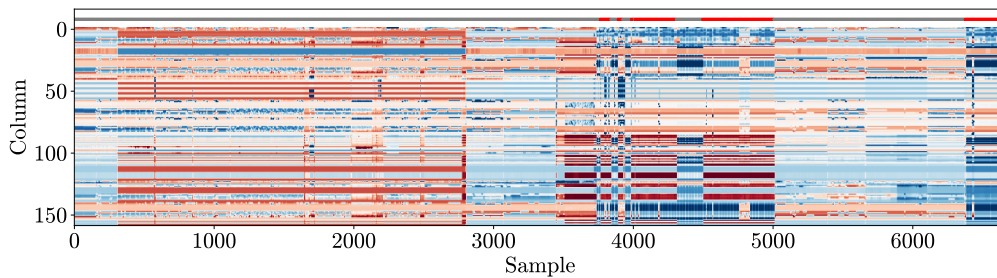


Figure 3.1: Heatmap of the dataset. On the top line, gray points represent normal behavior while red points are anomalies. In the heatmap, white tiles represent the mean, red tiles represent values below average, and blue tiles represent values above average.

## 3.2 Approaches

### 3.2.1 Multivariate KDE

Given the train, validation, and test splits, the dataset is standardized using the training set alone. This avoids leaks from the test set.

The KDE model, bandwidth, and threshold are fitted as in Chapter 2.

**Cost model** A simple cost model can be based on:

- False positives,
- False negatives,
- A time tolerance.

**Problems** KDE on this dataset has the following problems:

- It is highly subject to the curse of dimensionality and requires more data to be reliable.
- As the dataset is used during inference, a large dataset is computationally expensive.
- It only provides an alarm signal and not an explanation of the anomaly. In low-dimensionality, this is still acceptable, but with high-dimensional data it is harder to find an explanation.