

Notes on Anomaly Detection

Sameer Kesava PhD

Contents

1	Unsupervised Learning Algorithms	1
1.1	MeanShift Clustering	1
1.2	DBSCAN	1
2	Supervised Learning Algorithms	2
2.1	Univariate Data	2
2.2	Multivariate Data	2
2.3	Random Cut Forest	2
2.4	XGBoost	2
2.5	Isolation Forest	3
3	Improving the Accuracy	4
3.1	Hyperparameter Tuning	4

Unsupervised Learning Algorithms

For more information, please see **Deep Learning for Anomaly Detection: A Survey**; <http://arxiv.org/abs/1901.03407>.

1.1 MeanShift Clustering

- Distance-based clustering
- Has the ability to detect outliers.

1.2 DBSCAN

Density-based spatial clustering of applications with Noise

- Hyperparameters: `min_samples`, ϵ and ϵ -metric.
- Has the ability to detect outliers.
- Works well for non-linear data.
- Affected by the curse of dimensionality.

Supervised Learning Algorithms

2.1 Univariate Data

- Boxplot
- Grubbs test
- RANSAC algorithm for linear regression
- Studentized residuals and leverage points. Easy to plot for univariate data.

2.2 Multivariate Data

2.3 Random Cut Forest

From Amazon SageMaker

2.4 XGBoost

Highly popular classifier and regressor.

- Gradient boosting method
- Absolute loss and Huber loss more robust to outliers.
- Hyperparameters
 1. Max_depth
 2. Colsample_bytree
 3. Eta
 4. train-test split: 60-40/70-30/80-20.

2.5 Isolation Forest

Improving the Accuracy

3.1 Hyperparameter Tuning

- Hyperparameter optimization based on Gaussian Process Regression and Bayesian Optimization
- keras tuner in keras
- GridSearchCV or RandomSearchCV in scikit-learn
- RandomSearch can be used as the baseline against which optimization algorithms can be evaluated.

Bibliography

- [1] Pankaj Malhotra et al., Long Short Term Memory Networks for Anomaly Detection in Time Series, 2015.