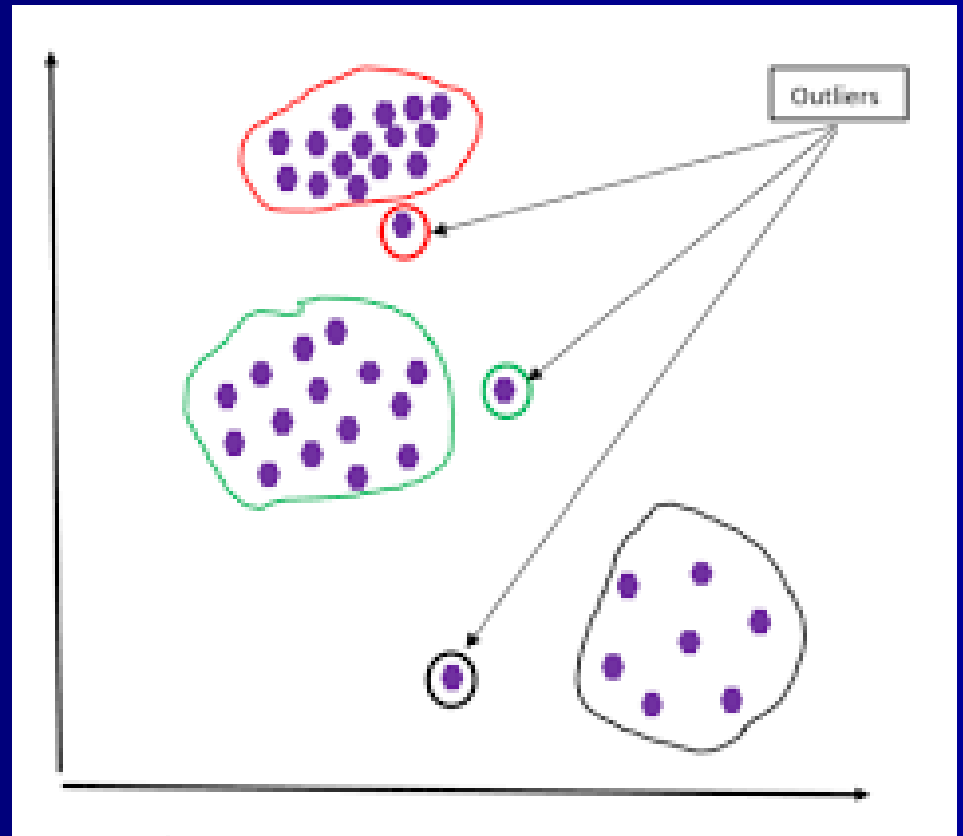
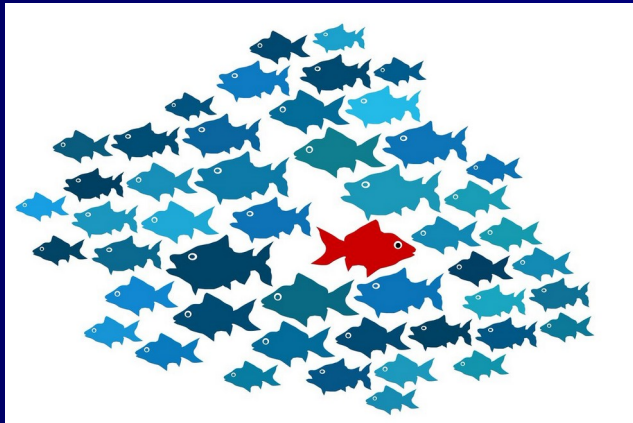


ONE-CLASS CLASSIFICATION

– OUTLIER DETECTION –



OUTLINE

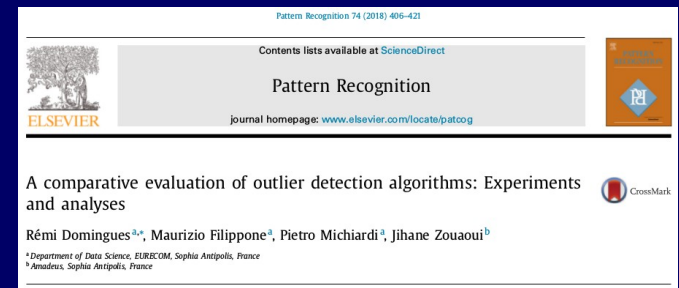
- The terms
- One-Class Classification
- Outlier detection
- Other scenarios: anomaly detection, novelty detection
- Main “one-class classification” algorithms
- References and software

One-class classification: Concept learning in the absence of counter-examples.

[DMJ Tax - 2002 - elibrary.ru](#)

Degree: Dr. DegreeYear: 2001 Institute: Technische Universiteit Delft (The Netherlands)
Publisher: Print partners Ipskamp, Capitool 25, Postbus 333, 7500 AH Enschede, The Netherlands. This thesis treats the problem of one-class classification. It starts with an ...

☆ 97 Citado por 1338 [Artículos relacionados](#) [Las 5 versiones](#)

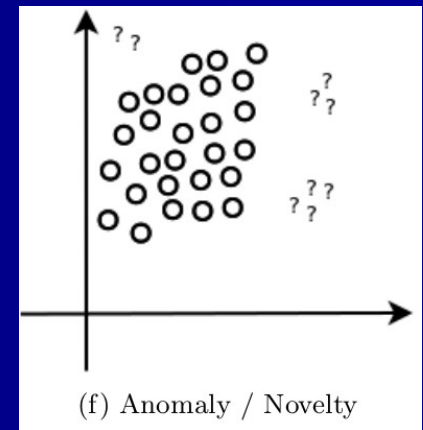
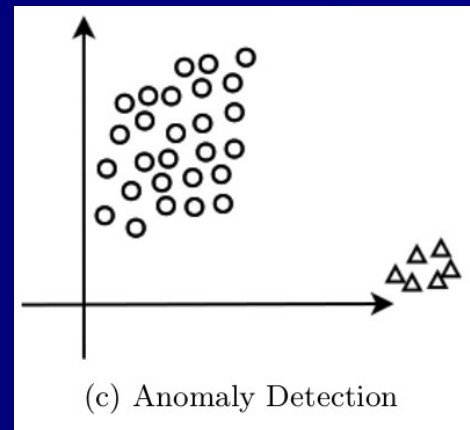
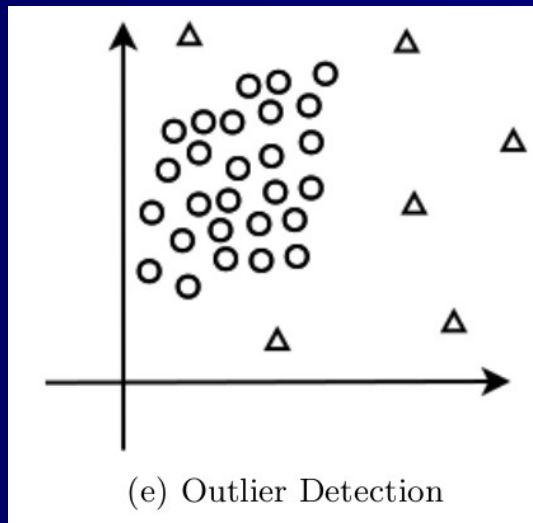


On the evaluation of outlier detection and one-class classification: a comparative study of algorithms, model selection, and ensembles

Henrique O. Marques, Lorne Swersky, Jörg Sander, Ricardo J. G. B. Campello & Arthur Zimek

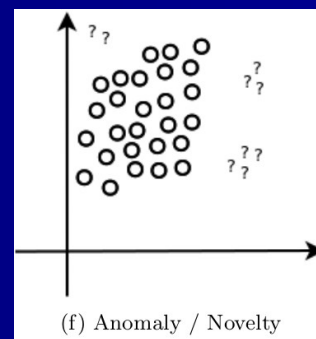
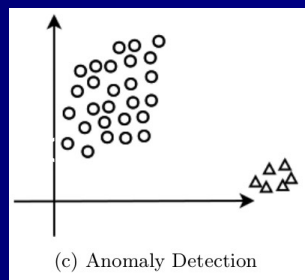
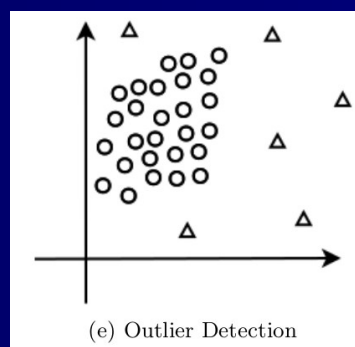
[Data Mining and Knowledge Discovery](#) (2023) | [Cite this article](#)

OUTLIER – ANOMALY – NOVELTY



- "Normal" category → common, majority
- Outlier samples → sparse, minority, category
- Anomaly samples → minoritaria, categoría ~ imbalanced supervised learning
- Novelty samples → anomaly, unknown in training time

OUTLIER – ANOMALY – NOVELTY

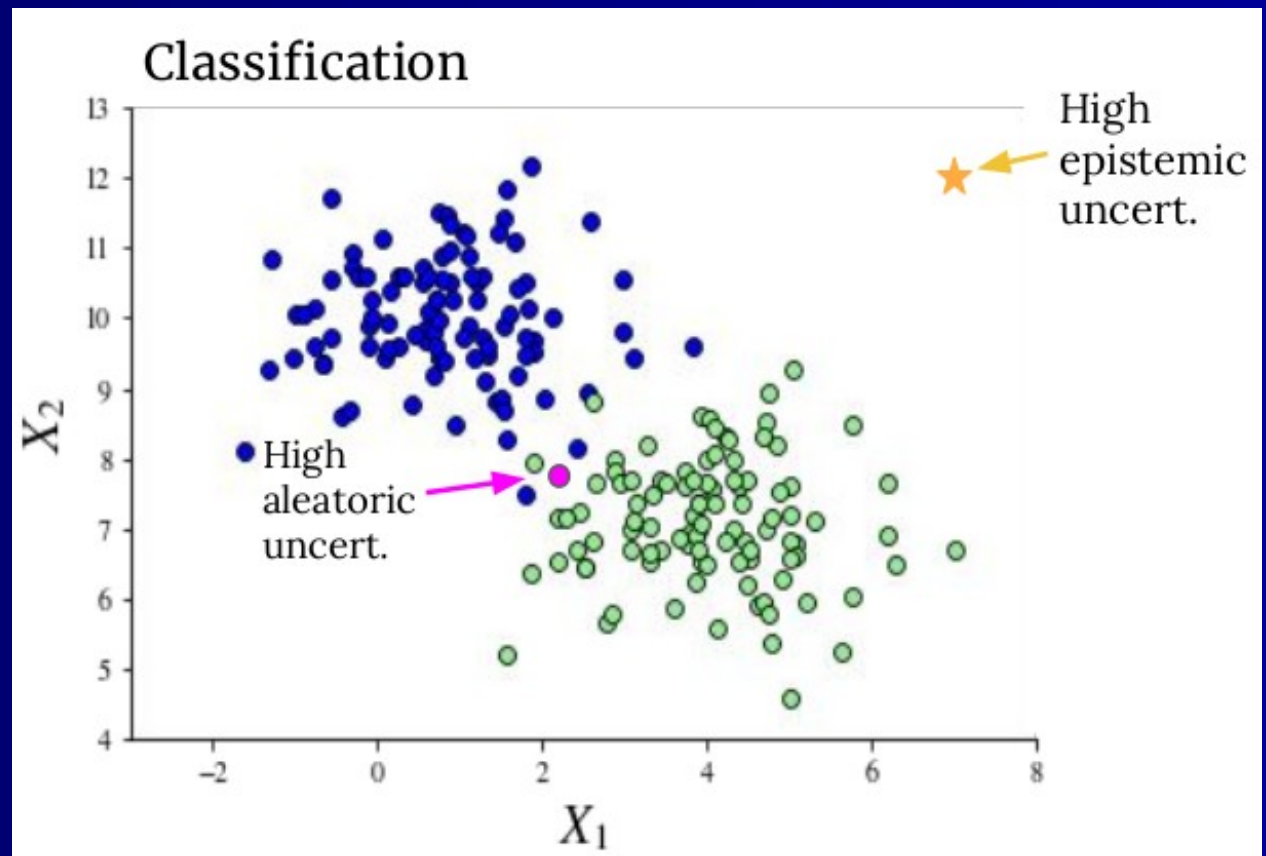
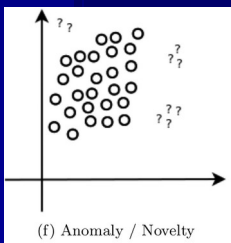
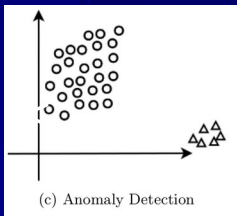
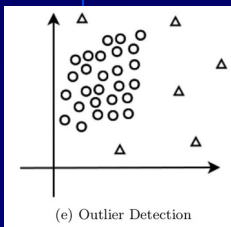


Artificial Intelligence Review (2020) 53:3575–3594
<https://doi.org/10.1007/s10462-019-09771-y>

Analyzing rare event, anomaly, novelty and outlier detection terms under the supervised classification framework

Ander Carreño¹  · Iñaki Inza¹ · Jose A. Lozano^{1,2}

ALEATORIC versus EPISTEMIC UNCERTAINTY

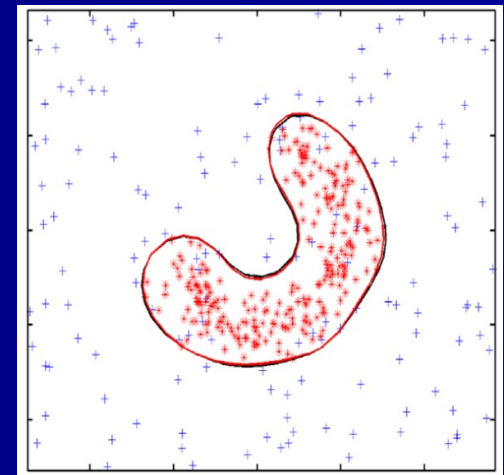
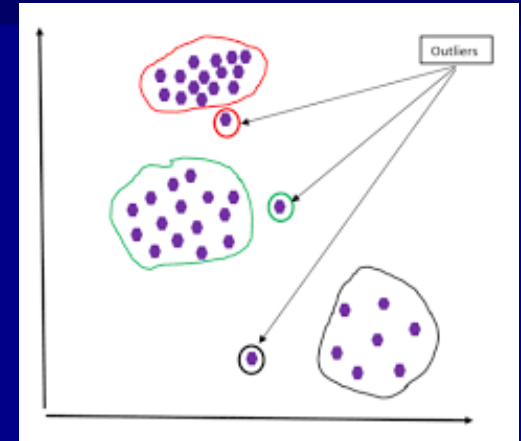


OUTLIER ~~ ONE-CLASS CLASSIFICATION

- Outlier samples → sparse, minority, category
- Outlier scenario → multi-class scenario + more than one class

most cited ones is Hawkins' definition (Hawkins 1980), which refers to an outlier as “an observation which deviates so much from other observations as to arouse suspicions that it was generated by a different mechanism”. Detecting such patterns is important because they might represent extraordinary behaviors that deserve special

- One-class classification → single class
- Single class → model its boundary + isolate the “rest”
- Non-single-class samples → non-modeled



OUTLIER DETECTION

– THE UNIVARIATE APPROACH –

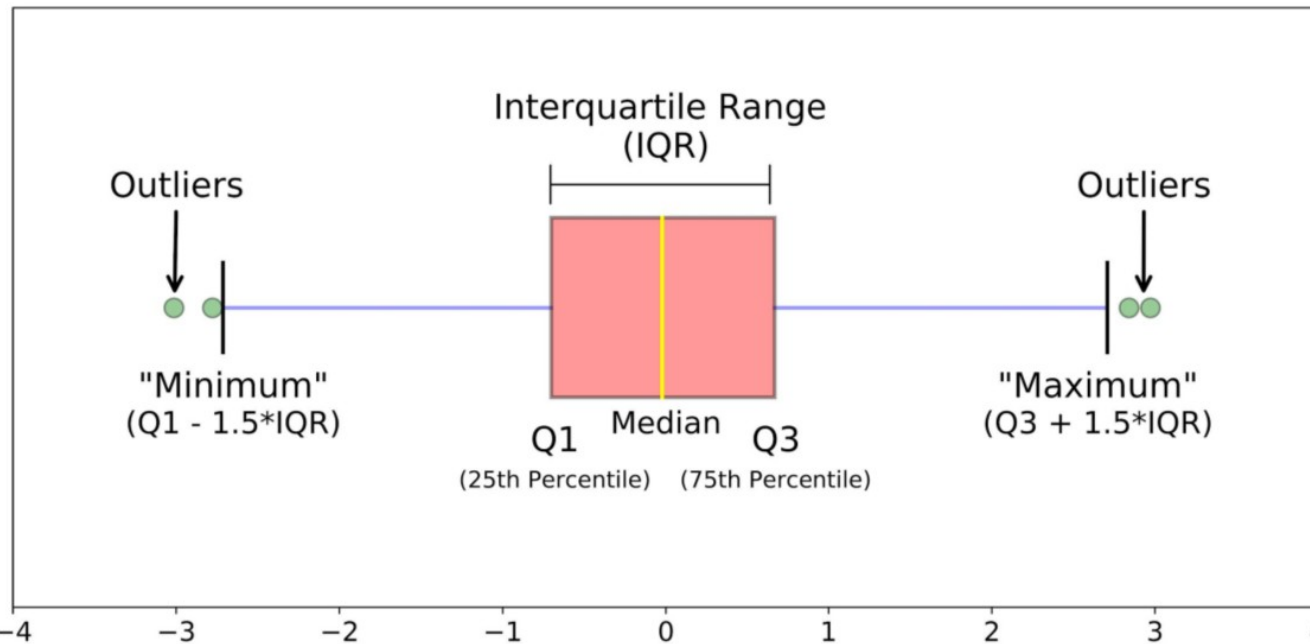
Whenever the goal is to identify univariate outliers, the statistical methods are among the simplest ones. Assuming a Gaussian distribution and learning the parameters from the data, parametric methods identify the points with low probability as outliers. One of the methods used to spot such outliers is the boxplot method, introduced by [Tukey \(1977\)](#). Based on the first quartile ($Q1$), the third quartile ($Q3$) and the interquartile range ($IQR = Q3 - Q1$) of the data, it determines that the interval $[Q1 - 1.5 * IQR, Q3 + 1.5 * IQR]$ contains 99.3% of data. Therefore, points outside that interval are considered as mild outliers, and points outside the interval $[Q1 - 3 * IQR, Q3 + 3 * IQR]$ are considered extreme outliers.

Regarding the unsupervised learning approach, we have used a boxplot-based method (boxplotEns) as follows: we generate five boxplots, one for each attribute, on the training data. For each new observation, we check if the value of each attribute is considered to be an outlier by the respective boxplot. The observation is considered to be an **Abnormal cycle** if at least one of the five attributes values is signalled as an outlier. This approach corresponds to an ensemble of boxplots to deal with multivariate data.

OUTLIER DETECTION – THE UNIVARIATE APPROACH –

Whenever the goal is to identify univariate outliers, the statistical methods are among the simplest ones. Assuming a Gaussian distribution and learning the parameters from the data, methods

used to :
first qua
the data,
of data.
outside



the data,
methods
d on the
- Q1) of
is 99.3%
d points
ers.

F
(box
data
outl

at least one of the five attributes values is signified as an outlier. This approach corresponds to an ensemble of boxplots to deal with multivariate data.

method
training
to be an
cycle if

OUTLIER DETECTION DATASETS – BENCHMARKS –

- *What do you mean by “outlier”?*
- **Multi-dimensional point datasets** [* ← ours]
- Time series graph datasets for event detection
- Time series point datasets (uni/multi -variate time series)
- Cyber-attack scenarios – security datasets
- Crowded scene video for anomaly detection

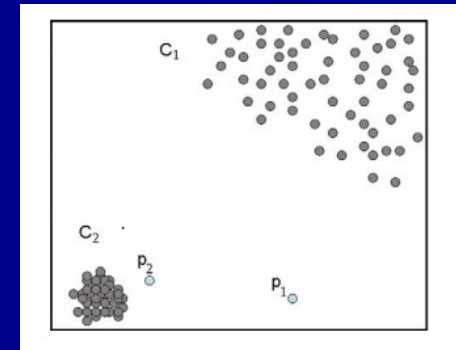
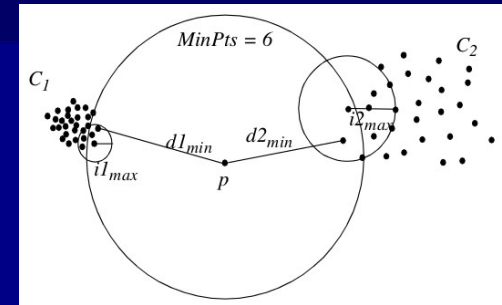


Outlier Detection DataSets (ODDS)

In ODDS, we openly provide access to a large collection of outlier detection datasets with ground truth (if available). Our focus is to provide datasets from different domains and present them under a single umbrella for the research community. As such, we arrange the datasets based on their types into different tables in the order as listed below. [\[read more about ODDS\]](#)

LOCAL OUTLIER FACTOR - LOF

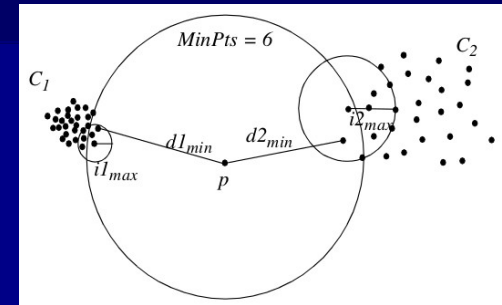
- Distance-based algorithm
 - To decide "outlier"
 - → by local neighborhood
 - → by local density
 - Parameter → k , number of neighbours
 - Calculate the neighborhood
-
- Outlier → defined "locally"
 - Outlierness → compute density of its local k -neighborhood



```
library(DDOutlier)
# 1860 Daily Closing Prices of Major European Stock Indices
# https://stat.ethz.ch/R-manual/R-devel/library/
# datasets/html/EuStockMarkets.html
data("EuStockMarkets")
colnames(EuStockMarkets)
# calculate "outlierness" score, by LOF
outlierness = LOF(dataset=EuStockMarkets, k=5)
# assign an index to outlierness values
names(outlierness) <- 1:nrow(EuStockMarkets)
sort(outlierness, decreasing=TRUE)
hist(outlierness)
which(outlierness > 2.0)
```

LOCAL OUTLIER FACTOR - LOF

- Distance-based algorithm
- To decide “outlier”
 - → by local neighborhood
 - → by local density
- Parameter → k , number of neighbours
- Calculate the neighborhood
- Outlier → defined “locally”
- Outlierness → compute density of its local k -neighborhood



Proc. ACM SIGMOD 2000 Int. Conf. On Management of Data, Dallas, TX, 2000

LOF: Identifying Density-Based Local Outliers

Markus M. Breunig[†], Hans-Peter Kriegel[†], Raymond T. Ng[‡], Jörg Sander[†]

[†] Institute for Computer Science
University of Munich
Oettingenstr. 67, D-80538 Munich, Germany

Department of Computer Science
University of British Columbia
Vancouver, BC V6T 1Z4 Canada

{ breunig | kriegel | sander }
@dbs.informatik.uni-muenchen.de

rng@cs.ubc.ca

```
library(DDOutlier)
# 1860 Daily Closing Prices of Major European Stock Indices
# https://stat.ethz.ch/R-manual/R-devel/library/
# datasets/html/EuStockMarkets.html
data("EuStockMarkets")
colnames(EuStockMarkets)
# calculate "outlierness" score, by LOF
outlierness = LOF(dataset=EuStockMarkets, k=5)
# assign an index to outlierness values
names(outlierness) <- 1:nrow(EuStockMarkets)
sort(outlierness, decreasing=TRUE)
hist(outlierness)
which(outlierness > 2.0)
```

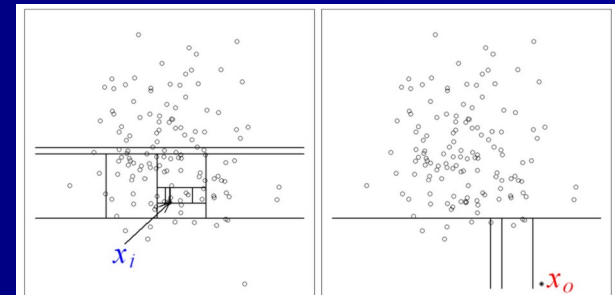
ISOLATION FOREST

- Compute “isolation score” *per sample*
- Construct a tree per sample by:
 - Random selection of
 - {attribute-split and attribute-split-value}
- Isolate the sample from the rest
- “Outliers easier to isolate...”
 - with fewer splits
- Path length from root to node
- ~ “isolation score” = “outlierness”
- “low path length” ~ “high outlierness”
- → easy to isolate point
- → graph “outlierness” values → threshold

Isolation Forest

Fei Tony Liu, Kai Ming Ting
Gippsland School of Information Technology
Monash University, Victoria, Australia
{tony.liu},{kaiming.ting}@infotech.monash.edu.au

Zhi-Hua Zhou
National Key Laboratory
for Novel Software Technology
Nanjing University, Nanjing 210093, China
zhouzh@lamda.nju.edu.cn



(a) Isolating x_i

(b) Isolating x_o

```
# package with benchmark datasets
library(mlbench)
# Census data for 506 Boston houses
data("BostonHousing", package = "mlbench")
# Package with IsolationForest implementation
library(solitude)
# Empty tree structure
iso <- isolationForest$new()
# Learn the IsolationTree for our data
iso$fit(BostonHousing)
p <- iso$predict(BostonHousing)
print(p)
sort(p$anomaly_score)
plot(density(p$anomaly_score))
# based on the plot, decide the cut-off point:
# indexes of samples with Outlierness > 0.63
which(p$anomaly_score > 0.63)
```

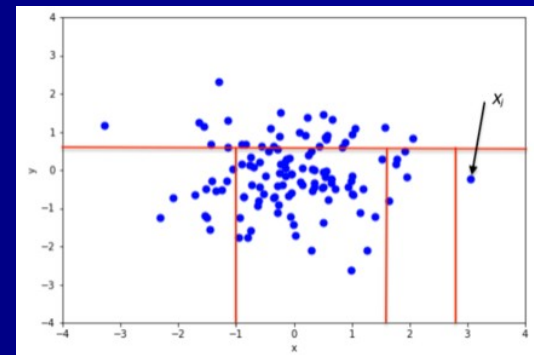
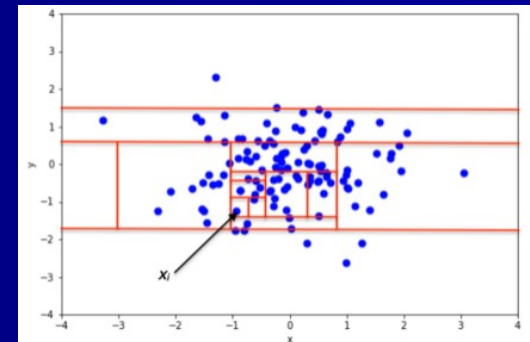
ISOLATION FOREST

- Compute “isolation score” *per sample*
- Construct a tree per sample by:
 - Random selection of
 - {attribute-split and attribute-split-value}
- Isolate the sample from the rest
- “Outliers easier to isolate...”
 - with fewer splits
- Path length from root to node
- ~ “isolation score” = “outlierness”
- “low path length” ~ “high outlierness”
- → easy to isolate point
- → graph “outlierness” values → threshold

Isolation Forest

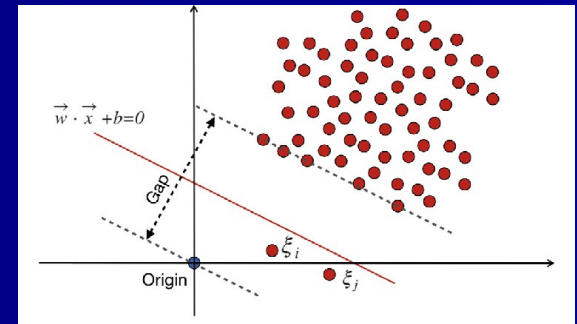
Fei Tony Liu, Kai Ming Ting
Gippsland School of Information Technology
Monash University, Victoria, Australia
{tony.liu},{kaiming.ting}@infotech.monash.edu.au

Zhi-Hua Zhou
National Key Laboratory
for Novel Software Technology
Nanjing University, Nanjing 210093, China
zhouzh@lamda.nju.edu.cn



OneClass SVMs – OCVSM

- Learn a SVM with single-class samples
- Map to higher dimension space
- Separating hyperplane
- Maximize margin between origin and data
- Outliers → points outside boundary



```
library(e1071)
# Daily air quality measurements in New York,
# May to September 1973.
# https://stat.ethz.ch/R-manual/R-devel/library/
# datasets/html/airquality.html
data(airquality)
df <- airquality

# all variables to be numerical

#train a SVM one-classification model
model <- svm(df, y=NULL, type='one-classification')

print(model)
summary(model) #print summary

# test on the whole set
# TRUE values mean suspect outliers
pred <- predict(model, df)
which(pred==TRUE)
```

Support Vector Method for Novelty Detection

Bernhard Schölkopf*, Robert Williamson[§],
Alex Smola[§], John Shawe-Taylor[†], John Platt*

* Microsoft Research Ltd., 1 Guildhall Street, Cambridge, UK

[§] Department of Engineering, Australian National University, Canberra 0200

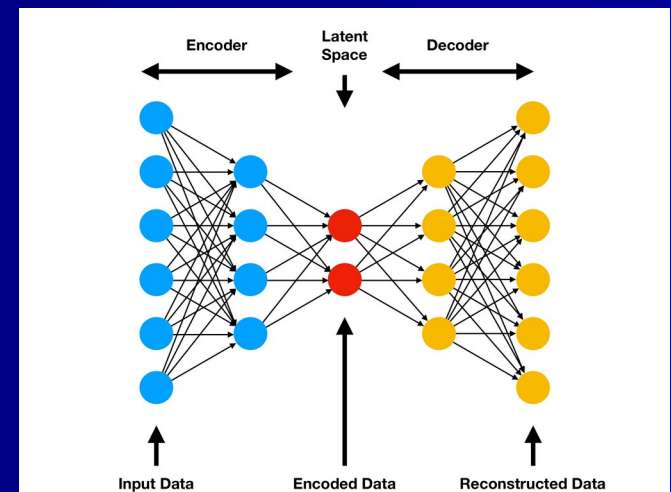
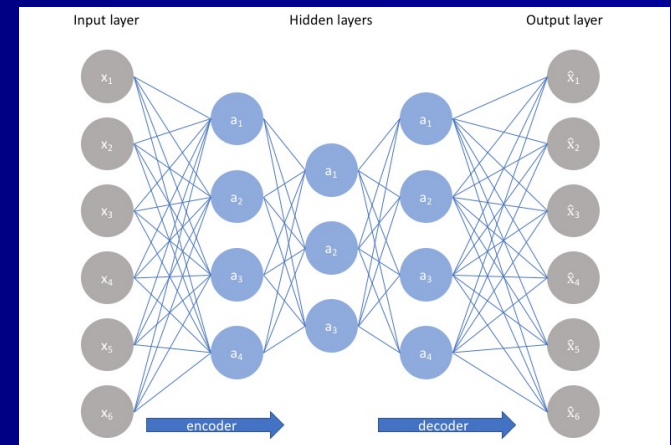
[†] Royal Holloway, University of London, Egham, UK

* Microsoft, 1 Microsoft Way, Redmond, WA, USA

bsc/jplatt@microsoft.com, Bob.Williamson/Alex.Smola@anu.edu.au, john@dc.srbnrc.ac.uk

AUTOENCODERS – DEEP LEARNING –

- Learn “encoded” data representation
- Reducing to non-linear dimensions in hidden layers
- {Encode + Decode} 1-class data
- Check for anomalies
- Does the autoencoder “reconstruct” the input data in the output?
- → “reconstruction error”
- → high value indicative of outlierness
- Hidden layers' features
- Non-linear, compact representation
- → learn with them a supervised model?

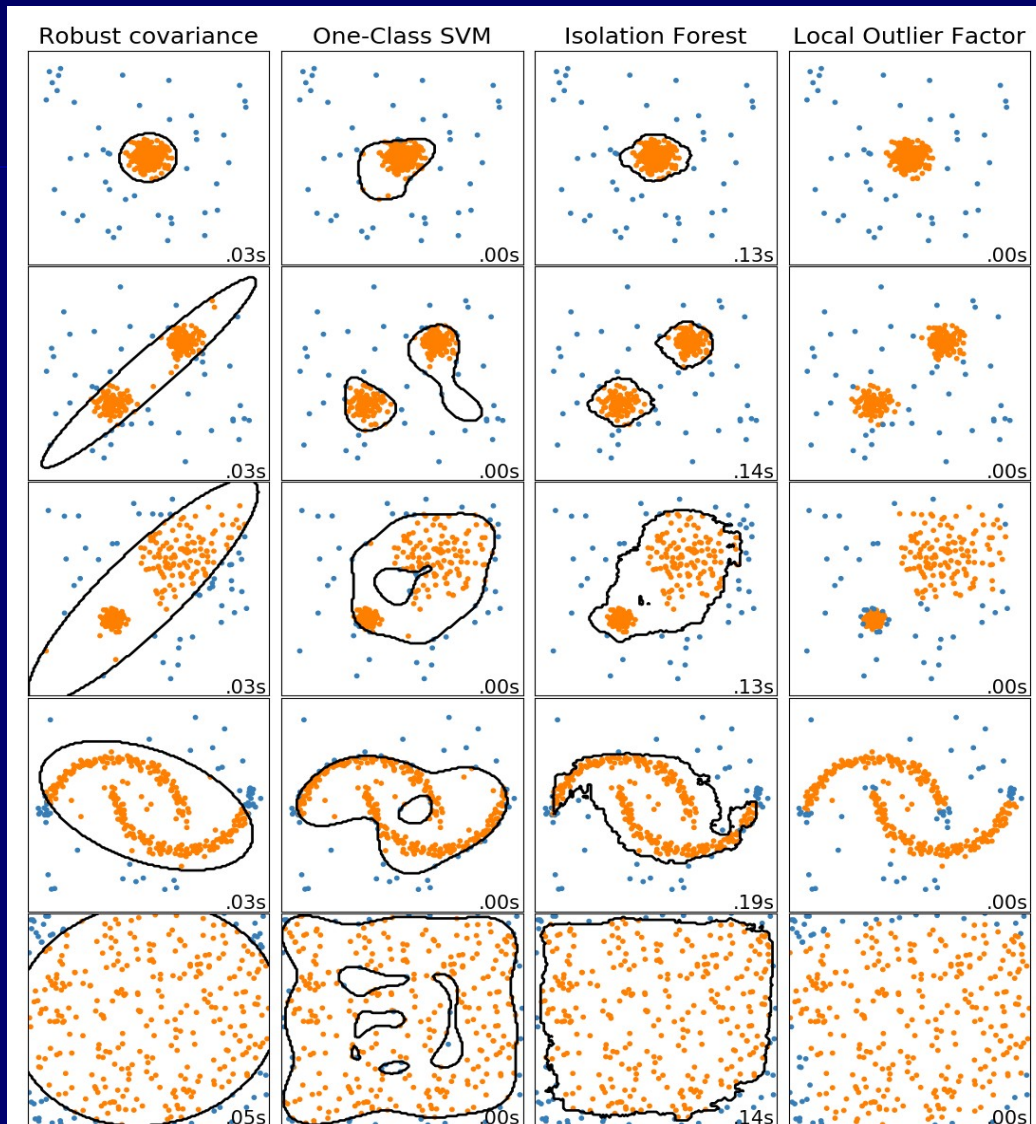


AUTOENCODERS – DEEP LEARNING –

- Learn representation of data
- Reducing to non-linear dimensions in hidden layers
- {Encode + Decode} 1-class data
- Check for anomalies
- Does the autoencoder “reconstruct” the input data in the output?
- → “reconstruction error”
- → high value indicative of outlierness
- Hidden layers' features
- Compact, non-linear representation
- → learn with them a supervised model?

```
library(h2o)
h2o.init()
prostate_path = system.file("extdata",
                             "prostate.csv", package = "h2o")
prostate = h2o.importFile(path = prostate_path)
colnames(prostate)
dim(prostate)
# learn autoencoder with 2 hidden layers of 10 units each
autoencoder_model = h2o.deeplearning(x = 3:9,
                                     training_frame = prostate, autoencoder = TRUE,
                                     hidden = c(10, 10), epochs = 5)
# features in the autoencoder's first hidden layer
deep_features_layer1 = h2o.deepfeatures(autoencoder_model,
                                         prostate, layer=1)
# further supervised models can be trained with these features
head(deep_features_layer1)
# reconstruction error per sample ~ outlierness indicative
reconstruction_error = h2o.anomaly(autoencoder_model, prostate)
head(reconstruction_error)
reconstruction_error = as.data.frame(reconstruction_error)
plot(sort(reconstruction_error$Reconstruction.MSE),
     main='Reconstruction Error')
which(reconstruction_error > 0.15)
```


ONE-CLASS CLASSIFICATION





Contents lists available at ScienceDirect

Pattern Recognition

journal homepage: www.elsevier.com/locate/patcog



A comparative evaluation of outlier detection algorithms: Experiments and analyses



Rémi Domingues^{a,*}, Maurizio Filippone^a, Pietro Michiardi^a, Jihane Zouaoui^b

^a Department of Data Science, EURECOM, Sophia Antipolis, France

^b Amadeus, Sophia Antipolis, France

On the evaluation of outlier detection and one-class classification: a comparative study of algorithms, model selection, and ensembles

[Henrique O. Marques](#) , [Lorne Swersky](#), [Jörg Sander](#), [Ricardo J. G. B. Campello](#) & [Arthur Zimek](#)

Data Mining and Knowledge Discovery (2023) | [Cite this article](#)

“BASQUE” APPLICATION INDUSTRY 4.0

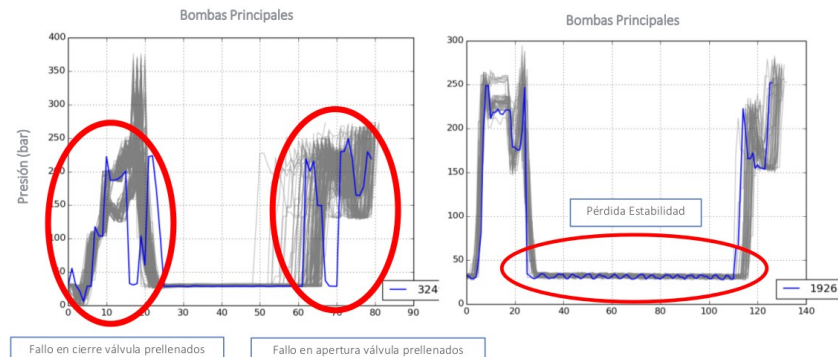


MACHINE LEARNING: OUTLIER DETECTION Y CLUSTERING



Técnicas de Machine Learning

- Definición de comportamientos normales → Desarrollo de Patrones.
- Análisis de ciclos en Tiempo real.
- Detección de desviaciones y análisis de causas → búsqueda del origen.
- Comportamientos que revelan síntomas de fallo en otros elementos.



- “Machine-tool” manufacturers
- Non-availability of “failure-class data”
- Predictive maintenance – “early prediction”
- “Do not arrive to failure and avoid machine stop”

EXERCISE

- Choose a publication → describing a previous method
- Read abstract + show the paper
- Find a software package which develops it
- Study + resume the parameters of the method
- Choose a supervised dataset (e.g. iris)
- Choose one of its classes (e.g. "iris setosa")
- Apply the one-class method over it
- Be careful !! → methods may only work with numerical features
- → remove the class column!!
- Graphs showing "outlierness" distribution
- Cut-off value in "outlierness" to decide outliers?