

Never miss a tutorial:



Machine Learning Mastery
Making Developers Awesome at Machine Learning

Picked for you:



How to Choose a Feature Selection
Method For Machine Learning

[Click to Take the FREE Data Preparation Crash-Course](#)

Search...



Data Preparation for Machine Learning (7-
Day Mini-Course)

Automatic Outlier Detection Algorithms in Python

by **Jason Brownlee** on July 8, 2020 in **Data Preparation**

How to Calculate Feature Importance



With Python

Share

Last Updated on August 17, 2020



Recursive Feature Elimination (RFE) for

Feature Selection in Python

Identifying and **removing outliers** is challenging with simple statistical methods for most machine learning datasets given the large number of input variables. Instead, **model-based outlier detection** methods can be used in the modeling pipeline and compared, just like other data preparation transforms may be applied to the dataset.

In this tutorial, you will discover how to use automatic outlier detection and removal to improve machine learning predictive modeling performance.

Loving the Tutorials?

After completing this tutorial, you will know:

- The [Data Preparation](#) EBook is where you'll find the **Really Good** stuff.
- Automatic outlier detection models provide an alternative to statistical techniques with a larger number of input variables with complex and unknown inter-relationships.
- How to use automatic outlier detection and removal to the training dataset only to avoid data leakage.
- How to evaluate and compare predictive modeling pipelines with outliers removed from the training dataset.

[Start Machine Learning](#)

Kick-start your project with my new book [Data Preparation for Machine Learning](#), including *step-by-step tutorials* and the *Python source code* files for all examples.



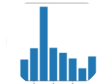
Picked for you:



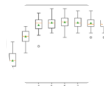
How to Choose a Feature Selection Method For Machine Learning



Data Preparation for Machine Learning (7 Day Mini-Course)



How to Calculate Feature Importance With Python



Recursive Feature Elimination (RFE) for Feature Selection in Python



How to Remove Outliers for Machine Learning

Loving the Tutorials?

The [Data Preparation EBook](#) is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Model-Based Outlier Detection and Removal in Python
Photo by Zoltán Vörös, some rights reserved.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Tutorial Overview

This tutorial is divided into the parts; they are:

1. Outlier Detection and Removal

Picked for you:

1. House Price Regression Dataset
2. How to Choose a Feature Selection Method For Machine Learning
3. Automatic Outlier Detection
4. Isolation Forest
5. Minimum Covariance Determinant
6. Data Preparation for Machine Learning (7- Day Mini Course)
7. Local Outlier Factor
8. One-Class SVM

Outlier Detection and Removal

Outliers are observations in a dataset that don't fit in some way.

The most common or familiar type of outlier is the observations that are far from the rest of the observations.

This is easy to understand when we have one or two variables and we can visualize the data as a histogram. When we have many input variables defining a high-dimensional input feature space.

In this case, simple statistical methods for identifying outliers can break down, such as methods that use standard deviations or the interquartile range.

It can be important to identify and remove outliers from data when training machine learning algorithms for predictive modeling.

Loving the Tutorials?

Outliers can skew statistical measures and data distributions, providing a misleading representation of the underlying data and relationships. Removing outliers from training data prior to modeling can result in a better fit of the data and, in turn, more skillful predictions.

Thankfully, there are several automatic model-based methods for identifying outliers in input data. Importantly, each method approaches the definition of an outlier is slightly different ways, providing alternate approaches to preparing a training dataset.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

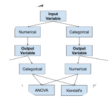
like any other data preparation step in a modeling pipeline.

Never miss a tutorial:

Before we dive into automatic outlier detection methods, let's first select a standard machine learning dataset that we can use as the basis for our investigation.



Picked for you:



How to Choose a Feature Selection Method For Machine Learning

Want to Get Started With Data Preparation?

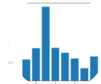
Take my free 7-day email crash course now (with sample code).



Data Preparation for Machine Learning (7 Day Mini-Course)

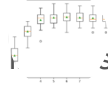
Click to sign-up and also get a free PDF Ebook version of the course.

Download Your FREE Mini-Course



How to Calculate Feature Importance With Python

Dataset and Performance Baseline



Recursive Feature Elimination (RFE) for Feature Selection in Python

In this section, we will first select a standard machine learning dataset and establish a baseline in performance.

This will provide the context for exploring the outlier identification and removal method of data preparation.



How to Remove Outliers for Machine Learning

House Price Regression Dataset

We will use the house price regression dataset.

Loving the Tutorials?

This dataset has 13 input variables that describe the properties of the house and suburb and requires the prediction of the median value of houses in the suburb in thousands of dollars.

The Data Preparation Book is where you'll find the **Really Good** stuff.

You can learn more about the dataset here:

>> SEE WHAT'S INSIDE

- [House Price Dataset \(housing.csv\)](#)
- [House Price Dataset Description \(housing.names\)](#)

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:

Open the dataset and review the raw data. The first few rows of data are listed below.



We can see that it is a regression predictive modeling problem with numerical input variables, each of which has different scales.

Picked for you:

[illegible]

The dataset is a **Set of 100,000 machine learning (7 Day Mini-Course)** examples that have unknown and complex relationships. We don't know that outliers exist in this dataset, but we may guess that some outliers may be present.

The example below loads the dataset and splits it into the input and output columns, splits it into training and testing sets, and calculates the feature importance.



How to Calculate Feature Importance With Python

```

1 # load and summarize the dataset
2 from pandas import read_csv
3 from sklearn.model_selection import train_test_split
4 # load the dataset
5 url = https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv
6 df = read_csv(url, header=None)
7 # retrieve the array
8 data = df.values
9 # split into input and output elements
10 X, y = data[:, :-1], data[:, -1]
11 # summarize the shape of the dataset
12 print(X.shape, y.shape)
13 # split into train and test sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
15 # summarize the shape of the train and test sets
16 print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)

```

Running the example, we can see that the dataset was loaded correctly and that there are 506 rows of data with 13 input variables and a single target variable.

>> SEE WHAT'S INSIDE

The dataset is split into two sets with 339 rows used for model training and 167 for model evaluation.

1	(506, 13)	(506,)
---	-----------	--------

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Start Machine Learning

```
2 (339, 13) (167, 13) (339,) (167,)
```

Never miss a tutorial:

Next, let's evaluate a model on this dataset and establish a baseline in performance.




Baseline Model Performance

Picked for you: predictive modeling problem, meaning that we will be predicting a numeric value. All input variables are also numeric.

In this case, we will fit a linear regression algorithm and evaluate model performance by training the model on the test dataset and making a prediction on test data and evaluate the predictions using the mean absolute error (MAE).

The complete example of evaluating a linear regression model on the dataset is listed below.

 **Data Preparation for Machine Learning (7-Day Mini-Course)**

```
1 # evaluate model on the raw dataset
2 from pandas import read_csv
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_absolute_error
6 # load the dataset
7 url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv'
8 df = read_csv(url, header=None)
9 # retrieve the array
10 data = df.values
11 # split into input and output elements
12 X, y = data[:, :-1], data[:, -1]
13 # split into train and test sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
15 # fit the model
16 model = LinearRegression()
17 model.fit(X_train, y_train)
18 # evaluate the model
19 yhat = model.predict(X_test)
20 # evaluate predictions
21 mae = mean_absolute_error(y_test, yhat)
22 print('MAE: %.3f' % mae)
```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Running the Data Preparation Book evaluates the model, then reports the MAE. where you'll find the **Really Good** stuff.

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the examp >> SEE WHAT'S INSIDE are the average outcome.

Start Machine Learning

In this case, we can see that the model achieved a MAE of about 3.417. This provides a baseline in performance to which we can compare different outlier identification and removal procedures.

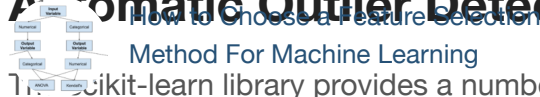


1 MAE: 3.417

Next, we can try removing outliers from the training dataset.

Picked for you:

Automatic Outlier Detection



How to Choose a Feature Selection

Method For Machine Learning

The scikit-learn library provides a number of built-in automatic methods for identifying outliers in data.

In this section, we will review four methods and compare their performance on the house price dataset.



Data Preparation for Machine Learning (7

Day Mini-Course)

Each method will be defined, then fit on the training dataset. The fit model will then predict which examples are not (so-called inliers). The outliers will then be removed from the training dataset, then the model



How to Calculate Feature Importance

With Python

It would be invalid to fit the outlier detection method on the entire training dataset as this would result in access to data (or information about the data) in the test set not used to train the model. This may result



Recursive Feature Elimination (RFE) for

Feature Selection in Python

to completely detect outliers on “new data” such as the test set prior to making a prediction,

One approach might be to return a “None” indicating that the model is unable to make a prediction or to explore that may be appropriate for your project.



How to Remove Outliers for Machine

Learning

Isolation Forest

Isolation Forest, or iForest for short, is a tree-based anomaly detection algorithm.

Loving the Tutorials?

It is based on modeling the normal data in such a way as to isolate anomalies that are both few in number and different in the feature space.

The Data Preparation Ebook is where you'll find the **Really Good** stuff.



... our proposed method takes advantage of two anomalies' quantitative properties: i) they are the minority consisting of fewer instances and ii) the >> SEE WHAT'S INSIDE that are very different from those of normal instances.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

— Isolation Forest, 2008.
Never miss a tutorial:

The scikit-learn library provides an implementation of Isolation Forest in the `IsolationForest` class.

Perhaps the most important hyperparameter in the model is the “*contamination*” argument, which is used to help estimate the number of outliers in the dataset. This is a value between 0.0 and 0.5 and by default is set to 0.1.

Picked for you:

```
1 # ... How to Choose a Feature Selection
2 # identify outliers in the training dataset
3 iso = IsolationForest(contamination=0.1)
4 yhat = iso.fit_predict(X_train)
```

Once identified, we can remove the outliers from the training dataset.

Data Preparation for Machine Learning (7-Day Mini-Course)

```
1
2 # select all rows that are not outliers
3 mask = yhat != -1
4 X_train, y_train = X_train[mask, :], y_train[mask]
```

is together, the complete example of evaluating the linear model on the housing dataset with is listed below.

```
1 # Recursive Feature Elimination (RFE) for
2 from pandas import read_csv
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.ensemble import IsolationForest
6 from sklearn.metrics import mean_absolute_error
7 # load the dataset
8 url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv'
9 df = read_csv(url, header=None)
10 # retrieve the array
11 data = df.values
12 # split into input and output elements
13 X, y = data[:, :-1], data[:, -1]
14 # split into train and test sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
16 # summarize the shape of the training dataset
17 print(X_train.shape, y_train.shape)
18 # identify outliers in the training dataset
19 iso = IsolationForest(contamination=0.1)
20 yhat = iso.fit_predict(X_train)
21 # select all rows that are not outliers
```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning


```

22 mask = yhat != -1
23 # remove outliers: X_train[mask, :], y_train[mask]
24 # summarize the shape of the updated training dataset
25 print(X_train.shape, y_train.shape)
26 # fit the model
27 model = LinearRegression()
28 model.fit(X_train, y_train)
29 # evaluate the model
30 yhat = model.predict(X_test)
31 # evaluate predictions
32 mae = mean_absolute_error(y_test, yhat)
33 print('MAE: %f' % mae)

```

Running the example fits and evaluates the model, then reports the MAE.



Data Preparation for Machine Learning (7-Day Mini-Course)

Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

In this case, we can see that that model identified and removed 34 outliers and achieved a MAE of about 3.417.



How to Calculate Feature Importance with Python

```

1 (339, 13) (339,)
2 (305, 13) (305,)
3 MAE: 3.417

```

Recursive Feature Elimination (RFE) for Feature Selection in Python

Minimum Covariance Determinant



How to Remove Outliers for Machine Learning

If input variables have a Gaussian distribution, then simple statistical methods can be used to detect outliers.

For example, if the dataset has two input variables and both are Gaussian, then the feature space forms a multi-dimensional Gaussian and knowledge of this distribution can be used to identify values far from the distribution.

Loving the Tutorials?

This approach can be generalized by defining a hypersphere (ellipsoid) that covers the normal data, and data that falls outside this shape is considered an outlier. An efficient implementation of this technique for multivariate data is known as the Minimum Covariance Determinant, or MCD for short.

The Data Preparation Ebook is where you'll find the **Really Good** stuff.



The Minimum Covariance Determinant (MCD) method is a highly robust estimator of multivariate location and scatter, for which a fast algorithm also serves as a convenient and efficient tool for outlier detection.

>> SEE WHAT'S INSIDE

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

— Minimum Covariance Determinant and Extensions, 2017.
Never miss a tutorial:

The scikit-learn library provides access to this method via the `EllipticEnvelope` class.

It provides the “contamination” argument that defines the expected ratio of outliers to be observed in practice. In this case, we will set it to a value of 0.01, found with a little trial and error.
Picked for you:

```
1 # ... How to Choose a Feature Selection
2 # identify outliers in the training dataset
3 ee = EllipticEnvelope(contamination=0.01)
4 yhat = ee.fit_predict(X_train)
```

Once identified, the outliers can be removed from the training dataset as we did in the prior example.



Data Preparation for Machine Learning (7-

Day Mini-Course)

By putting this together, the complete example of identifying and removing outliers from the housing data (using the minimum covariance determinant) method is listed below.

```
1 # evaluate model performance with outliers removed using elliptical envelope
2 from pandas import read_csv
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.covariance import EllipticEnvelope
6 from sklearn.metrics import mean_absolute_error
7 # load the dataset
8 url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv'
9 df = read_csv(url, header=None)
10 # retrieve the array
11 data = df.values
12 # split into input and output elements
13 X, y = data[:, :-1], data[:, -1]
14 # split into train and test sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
16 # summarize the shape of the training dataset
17 print(X_train.shape, y_train.shape)
18 # identify outliers in the training dataset
19 ee = EllipticEnvelope(contamination=0.01)
20 yhat = ee.fit_predict(X_train)
21 # select all rows that are not outliers
22 mask = yhat != -1
23 X_train, y_train = X_train[mask, :], y_train[mask]
24 # summarize the shape of the updated training dataset
25 print(X_train.shape, y_train.shape)
26 # fit the model
```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.
 Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

```

27 model = LinearRegression()
28 model.fit(X_train, y_train)
29 # evaluate the model
30 yhat = model.predict(X_test)
31 # evaluate predictions
32 mae = mean_absolute_error(y_test, yhat)
33 print('MAE: %.3f' % mae)

```

Picked for you:

Running the example fits and evaluates the model, then reports the MAE.



How to Choose a Feature Selection

Method For Machine Learning

Our results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.



Data Preparation for Machine Learning (

Day 3)

```

1 (339, 13) (339,)
2 (335, 13) (335,)
3 MAE: 3.388

```



How to Calculate Feature Importance

With Python

Local Outlier Factor



Recursive Feature Elimination (RFE) for

Feature Selection in Python

The approach to identifying outliers is to locate those examples that are far from the other examples.

This can work well for feature spaces with low dimensionality (few features), although it can become referred to as the curse of dimensionality.



How to Remove Outliers for Machine

Learning

Local outlier factor, or LOF for short, is a technique that attempts to harness the idea of nearest neighbors for outlier detection. Each example is assigned a scoring of how isolated or how likely it is to be outliers based on the size of its local neighborhood. Those examples with the largest score are more likely to be outliers.

Loving the Tutorials?



We introduce a local outlier (LOF) for each object in the dataset, indicating its degree of outlier-ness.

The Data Preparation Ebook is

where you'll find the **Really Good** stuff.

— LOF: Identifying Density-based Local Outliers, 2000.

>> SEE WHAT'S INSIDE

The scikit-learn library provides an implementation of this approach in the `LocalOutlierFactor` class.

Start Machine Learning

You can master applied Machine Learning

without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

The model provides the “contamination” argument, that is the expected percentage of outliers in the dataset, be indicated and defaults to 0.1.

Never miss a tutorial:

```
1 ...
2 # identify outliers in the training dataset
3 lof = LocalOutlierFactor()
4 yhat = lof.fit_predict(X_train)
```

Picked for you:

Tying this together, the complete example of identifying and removing outliers from the housing dataset using the local outlier factor method is listed below.

How to Choose a Feature Selection
Method For Machine Learning

```
1 # evaluate model performance with outliers removed using local outlier factor
2 from pandas import read_csv
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.neighbors import LocalOutlierFactor
6 from sklearn.metrics import mean_absolute_error
7 # load the dataset
8 url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv'
9 df = read_csv(url, header=None)
10 # retrieve the array
11 data = df.values
12 # split into input and output elements
13 X, y = data[:, :-1], data[:, -1]
14 # split into train and test sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
16 # summarize the shape of the training dataset
17 print(X_train.shape, y_train.shape)
18 # identify outliers in the training dataset
19 lof = LocalOutlierFactor()
20 yhat = lof.fit_predict(X_train)
21 # select all rows that are not outliers
22 mask = yhat != -1
23 X_train, y_train = X_train[mask, :], y_train[mask]
24 # summarize the shape of the updated training dataset
25 print(X_train.shape, y_train.shape)
26 # fit the model
27 model = LinearRegression()
28 model.fit(X_train, y_train)
29 # evaluate the model
30 yhat = model.predict(X_test)
31 # evaluate predictions
32 mae = mean_absolute_error(y_test, yhat)
33 print('MAE: %.3f' % mae)
```

Loving the Tutorials?

The Data Preparation EBook is

SEARCH THIS BLOG

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Running the example fits and evaluates the model, then reports the MAE.

Start Machine Learning

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.



In this case, we can see that the local outlier factor method identified and removed 34 outliers, the same number as isolation forest, resulting in a drop in MAE from 3.417 with the baseline to 3.356. Better, but not as good as isolation forest, suggesting a different set of outliers were identified and removed.

Picked for you:

- 1 (339, 13) (339,)
 - 2 (305, 13) (305,)
 - 3 MAE: 3.356
- How to Choose a Feature Selection Method For Machine Learning

One-Class SVM



Data Preparation for Machine Learning (7 Day Mini-Course)

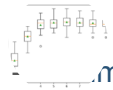
Support vector machine, or SVM, algorithm developed initially for binary classification can be used for one-class classification.

When modeling one class, the algorithm captures the density of the majority class and classifies ex outliers. This modification of SVM is referred to as One-Class SVM.



How to Calculate Feature Importance With Python

... an algorithm that computes a binary function that is supposed to capture regions in input support), that is, a function such that most of the data will live in the region where the function



Recursive Feature Elimination (RFE) for Feature Selection in Python

Estimating the Support of a High-Dimensional Distribution, 2001.

Although SVM is a classification algorithm and One-Class SVM is also a classification algorithm, it is not a classification algorithm and classification datasets.



How to Remove Outliers for Machine Learning

The scikit-learn library provides an implementation of one-class SVM in the `OneClassSVM` class.

The class provides the “nu” argument that specifies the approximate ratio of outliers in the dataset, which defaults to 0.1. In this case, we will set it to 0.01, found with a little trial and error.

The Data Preparation EBook is

```
1 .. where you'll find the Really Good stuff.
2 # identify outliers in the training dataset
3 ee = OneClassSVM(nu=0.01)
4 yhat = ee.fit_predict(X_train)
```

Tying this together, the complete example of identifying and removing outliers from the housing data

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE


```

1 # evaluate model performance with outliers removed using one class SVM
2 from pandas import read_csv
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.svm import OneClassSVM
6 from sklearn.metrics import mean_absolute_error
7 # load the dataset
8 url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv'
9 df = read_csv(url, header=None)
10 # retrieve the array
11 data = df.values
12 # split into input and output elements
13 X, y = data[:, :-1], data[:, -1]
14 # split into train and test sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
16 # summarize the shape of the training dataset
17 print(X_train.shape, y_train.shape)
18 # identify outliers in the training dataset
19 ee = OneClassSVM(nu=0.01)
20 yhat = ee.fit_predict(X_train)
21 # select all rows that are not outliers
22 mask = yhat != -1
23 X_train, y_train = X_train[mask, :], y_train[mask]
24 # summarize the shape of the updated training dataset
25 print(X_train.shape, y_train.shape)
26 # fit the model
27 model = LinearRegression()
28 model.fit(X_train, y_train)
29 # evaluate the model
30 yhat = model.predict(X_test)
31 # evaluate predictions
32 mae = mean_absolute_error(y_test, yhat)
33 print('MAE: %.3f' % mae)

```

Running the example fits and evaluates the model, then reports the MAE.

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

Loving the Tutorials?

The Data Preparation EBook is

In this case, we did not see *Really Good* outliers were identified and removed and the model achieved a MAE of about 3.431, which is not better than the baseline model that achieved 3.417. Perhaps better performance can be achieved with more tuning.

>> SEE WHAT'S INSIDE

```

1 (339, 13) (339,)
2 (336, 13) (336,)
3 MAE: 3.431

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial: Further Reading



This section provides more resources on the topic if you are looking to go deeper.

Picked for you: Related Tutorials



Papers



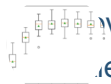
Data Preparation for Machine Learning (7-Day Mini Course), 2008.

- Minimum Covariance Determinant and Extensions, 2017.
- LOF: Identifying Density-based Local Outliers, 2000.



How to Calculate Feature Importance
Estimating the Support of a High-Dimensional Distribution, 2001.
With Python

APIs



Recursive Feature Elimination (RFE) for
Feature Selection in Python
Scikit-learn user guide.
scikit-learn.covariance.EllipticEnvelope API.

- sklearn.svm.OneClassSVM API.
- sklearn.neighbors.LocalOutlierFactor API.



How to Remove Outliers for Machine
Learning
scikit-learn.ensemble.IsolationForest API.

Summary

In this tutorial, you learned how to use automatic outlier detection and removal to improve machine learning predictive modeling performance.

Specifically, you learned:
The Data Preparation EBook is where you'll find the **Really Good** stuff.

- AutoML models provide an alternative to statistical techniques with a larger number of input variables with complex and unknown relationships.
- How to correctly apply automatic outlier detection and removal to the training dataset only to avoid overfitting.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

- How to evaluate and compare predictive modeling pipelines with outliers removed from the training dataset.
- Never miss a tutorial:**

Do you have any questions?

Ask your questions in the comments below and I will do my best to answer.

Picked for you:



How to Choose a Feature Selection Method For Machine Learning

Data Preparation for Machine Learning (7-Day Mini-Course)
Data Cleaning, Feature Selection, and Data Transforms in Python

How to Calculate Feature Importance With Python

Recursive Feature Elimination (RFE) for Feature Selection in Python

How to Remove Outliers for Machine Learning

Get a Handle on Modern Data Preparation!

Prepare Your Machine Learning Data in Minutes

...with just a few lines of python code

Discover how in my new
Data Preparation for Machine Learning

It provides **self-study tutorials** with
Feature Selection, RFE, Data Cleaning, Data Transforms, Scaling

Bring Modern Data Preparation to Your Machine Learning

SEE WHAT'S INSIDE

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.**
Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Tweet

Share
Loving the Tutorials?



The Data Preparation EBook is
where you'll find the *Really Good* stuff.

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on

>> SEE WHAT'S INSIDE

[view all posts by Jason Brownlee →](#)

Start Machine Learning

< [How to Use Feature Extraction on Tabular Data for Machine Learning](#)
Never miss a tutorial:

[6 Dimensionality Reduction Algorithms With Python](#) >



38 Responses to 4 Automatic Outlier Detection Algorithms in Python



[How to Choose a Feature Selection](#)

[Method For Machine Learning](#)

Joseph July 8, 2020 at 7:00 pm #

REPLY ↩

Hi Jason, thanks for one more great article!



[Data Preparation for Machine Learning \(7 Day Mini Course\)](#)

Hi Jason, thanks for one more great article!



[How to Calculate Feature Importance](#)

[With Python](#)

Jason Brownlee July 9, 2020 at 6:39 am #

I think trees are pretty robust to outliers. Test for your dataset.



[Recursive Feature Elimination \(RFE\) for](#)

[Feature Selection in Python](#)

JParzival July 9, 2020 at 10:42 am #

[How to Remove Outliers for Machine](#)

[Learning](#)

Thank you for sharing your experience!

Loving the Tutorials?



[The Data Preparation EBook is](#)

where you'll find the **Really Good** stuff.

You're welcome.

>> SEE WHAT'S INSIDE

REPLY ↩

Start Machine Learning

You can master applied Machine Learning

without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:

Nagdev. A

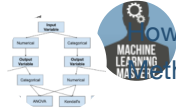
July 10, 2020 at 9:39 am #

REPLY ↩



Two more tutorials on autoencoders and PCA

Picked for you:



Jason Brownlee

How to Choose a Feature Selection Method For Machine Learning

For outlier detection? How so?

REPLY ↩



Data Preparation for Machine Learning (7-Day Mini-Course)



Ali November 19, 2020 at 4:20 pm #



Actually, autoencoders can provide best performance for anomaly detection problems

How to Calculate Feature Importance With Python



Recursive Feature Elimination (RFE) for Feature Selection in Python

Jason Brownlee

November 20, 2020 at 6:43 am #

Depends on the specific dataset.



How to Remove Outliers for Machine Learning



Ali November 19, 2020 at 4:27 pm #

Both Autoencoder and PCA are dimensionality reduction techniques. Interestingly, during the process of dimensionality reduction outliers are identified

Loving the Tutorials?

The Data Preparation EBook is

where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

pm #

REPLY ↩

Start Machine Learning



You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

REPLY ↩

Start Machine Learning

Hello sir,
Never miss a tutorial:

It was a great article. Just one doubt:



MCD technique doesn't perform well when the data has very large dimensions like >1000. In that case, it is a good option to feed the model with principal components of the data. The paper that you mentioned in the link says:

"For large p we can still make a rough estimate of the scatter as follows. First compute the first $q < p$ robust principal components of the data. For this we can use the MCD-based ROBPCA method⁵³, which requires that the number of components q be set rather low."



the ROBPCA is not a Feature Selection Method For Machine Learning

Thank you



Data Preparation for Machine Learning (7-Day Mini-Course)

Jason Brownlee July 11, 2020 at 6:13 am #



Great tip, thanks.
 How to Calculate Feature Importance

With Python

Perhaps find a different platform that implements the method?

Perhaps implement it yourself?

Perhaps use a different method entirely?

Recursive Feature Elimination (RFE) for
 Feature Selection in Python



fabou July 10, 2020 at 11:08 pm #

How to Remove Outliers for Machine Learning

Hi Jason,

as usual great educational article.

How could automatic outlier detection be integrated into a cross validation loop? Does it have to be part of a pipeline which steps would be : outlier detection > outlier removal (transformer) > modeling?

In this case Data Preparation is a transformer "outlier remover" be created?

where you'll find the **Really Good** stuff.
 Thanks

>> SEE WHAT'S INSIDE

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.
 Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:

Jason Brownlee July 11, 2020 at 6:16 am #

REPLY ↩

you would have to run the CV loop manually and apply the method to the data prior to fitting/evaluating a model or pipeline.

It's disappointing that sklearn does not support methods in pipelines that add/remove rows. imbalanced learn can do this kind of thing...

Picked for you:

How to Choose a Feature Selection

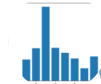
Method For Machine Learning

Chayma July 15, 2020 at 1:16 am #

REPLY ↩



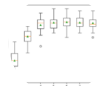
Thank you for the great article.
Data Preparation for Machine Learning (7-Day Mini Course)
Which algorithm is the most suitable for outlier detection in time series data?



How to Calculate Feature Importance

with Python **Jason Brownlee** July 15, 2020 at 8:27 am #

I don't know off hand, I hope to write about that topic in the future.



**Recursive Feature Elimination (RFE) for
Feature Selection in Python**

Nagesh August 27, 2020 at 12:33 pm #



How to Remove Outliers for Machine Learning
Thoughts on this one ? <https://github.com/arundo/adtk>



Loving the Tutorials? **Jason Brownlee** August 27, 2020 at 1:36 pm #

REPLY ↩

The **Data Preparation EBook** is
I'm not familiar with it.
where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Allen21 September 1, 2020 at 5:54 am #

Start Machine Learning



You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

If anyone is getting a `TypeError` with `X_train[mask, :]`, just change it to `X_train[mask]`. Another great article BTW

Never miss a tutorial:



Jason Brownlee September 1, 2020 at 6:38 am #

REPLY ↩

Picked for you:

Sorry to hear that.



How to Choose a Feature Selection Method For Machine Learning

Perhaps these tips will help.

<https://machinelearningmastery.com/faq/single-faq/why-does-the-code-in-the-tutorial-not-work-for-me>



Data Preparation for Machine Learning (7-Day Mini-Course)

Vidya Manu Shankar September 10, 2020 at 1:51 pm #

Hi Jason .

How to Calculate Feature Importance



Thanks for this post.

Couple of questions though:

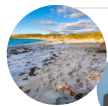
1. How do we validate the output of the outlier detection algorithms mentioned in this post , whether the



Recursive Feature Elimination (RFE) for

Feature Selection in Python

so, why don't we include the target variable as input to the outlier algorithms ? I missed this point.



How to Remove Outliers for Machine Learning

Jason Brownlee September 11, 2020 at 5:48 am #

Thanks.

Good question, you can validate the model by either evaluating predictions on dataset with known outliers or inspecting identified outliers and using a subject matter expert to determine if they are true outliers or not.

The [Data Preparation](#) EBook is

The algorithms are one class algorithms, no target variable is required. where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Vidya September 12, 2020 at 2:38 am #

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Awesome, thank you !
Never miss a tutorial:



Jason Brownlee September 12, 2020 at 6:18 am #

REPLY ↩

Picked for you:

You're welcome.



[How to Choose a Feature Selection Method For Machine Learning](#)

arnan maipradit September 13, 2020 at 10:29 am #

REPLY ↩



[Data Preparation for Machine Learning \(7-Day Mini-Course\)](#)

do you have any example of outlier detection using Q-learning, I found that Q-learning almost using in case of many actions (robot move up down left right so it has 4 actions) but in the case of outlier detection it has only 2 actions (normal behavior and abnormal behavior) can be used on outlier detection (anomaly detection) or not ? . If you could make an example or suggest



[How to Calculate Feature Importance With Python](#)



Jason Brownlee September 14, 2020 at 6:44 am #



[Recursive Feature Elimination \(RFE\) for Feature Selection in Python](#)

Sorry, I do not have any examples or RL at this stage.

Thanks for the suggestion.



[How to Remove Outliers for Machine Learning](#)

Arushi Mahajan October 14, 2020 at 1:45 pm #

REPLY ↩

Hi, amazing tutorial.

Loving the Tutorials?

Just one question. How can you see all the rows that were dropped?

The [Data Preparation](#) EBook is where you'll find the **Really Good** stuff.



>> SEE WHAT'S INSIDE October 14, 2020 at 1:51 pm #

REPLY ↩

What do you mean by "dropped rows"?

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:



Ma October 17, 2020 at 3:57 am #

REPLY ↩

Hey Jason,

Picked for you:

I've read about hyperparameter tuning of Isolation Forests etc. When all models/removing the detected outliers doesn't really add value or doesn't improve baseline model's score, I think it makes sense to invest time into hyperparameter tuning of these anomaly detection models?



Method For Machine Learning

: From my point of view those outliers seem to be legit to me...

Cheers again,



Data Preparation for Machine Learning (7-Day Mini-Course)



How to Calculate Feature Importance

With Python

Jason Brownlee October 17, 2020 at 6:13 am #

If it not improving performance, no.



Recursive Feature Elimination (RFE) for Feature Selection in Python

Mohammad Ali Shahlaei

November 30, 2020 at 6:38 am #



How to Remove Outliers for Machine

Learning
I think he meant that the rows were identified as outliers (dropped rows)!

Loving the Tutorials?

September 27, 2020 at 4:06 am #

REPLY ↩

The Data Preparation eBook is a great article. I have a question that is why we don't apply the outlier detection algorithm to the whole dataset rather than only the training dataset? It will not bother the accuracy of the model if there are outlier data in the test dataset?

>> SEE WHAT'S INSIDE

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:



Jason Brownlee

December 7, 2020 at 6:22 am #

REPLY ↩

don't use example only applies the automatic methods to the training dataset.

Picked for you:



Mitra December 20, 2020 at 1:42 am #

REPLY ↩

How to Choose a Feature Selection Method For Machine Learning

Hi sir! Thank you for the amazing content, Just wanted to point out one thing. In the Isolation Forests, documentation of Scikit learn I read that the default value for contamination is no longer 0.1 and it's turned to auto. You can correct that part 😊



Data Preparation for Machine Learning (7-Day Mini-Course)

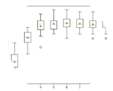


Jason Brownlee December 20, 2020 at 5:58 am #

How to Calculate Feature Importance With Python

You're welcome.

Thanks for the suggestion.



Recursive Feature Elimination (RFE) for Feature Selection in Python

Mitra December 20, 2020 at 4:49 am #



How to Remove Outliers for Machine Learning

One quick note! In the Minimum Covariance Determination method, you said we can use this method well in the dataset you're using the features don't have such shape. Most of them are skewed. I think we should first apply a transformation(log, box-cox, etc.) and then use this method on features with little or no skewness. I'm actually writing a Kaggle kernel on this and would love to hear what you think about it when it's done!

Loving the Tutorials?

The Data Preparation EBook is where you'll find the **Really Good** stuff.

REPLY ↩

>> SEE WHAT'S INSIDE need evaluating the approach with and without the data prep and use the approach that results in the best performance.

Start Machine Learning



You can master applied Machine Learning **without math or fancy degrees.** Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:

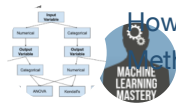
Divya Sami December 29, 2020 at 2:18 am #

REPLY ↩



Question- Should we always drop the rows containing outliers? Will outlier imputation work better in some cases?

Picked for you:



How to Choose a Feature Selection

Jason Brownlee

Method For Machine Learning December 29, 2020 at 5:15 am #

REPLY ↩

Thanks.



Data Preparation for Machine Learning (7

Day Mini-Course)



How to Calculate Feature Importance

With Python



Recursive Feature Elimination (RFE) for
Feature Selection in Python



How to Remove Outliers for Machine
Learning

Start Machine Learning



You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Loving the Tutorials?

Name (required)

The [Data Preparation](#) EBook is
where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

website

Start Machine Learning

Never miss a tutorial:

SUBMIT COMMENT



Welcome!

Picked for you!

I'm Jason Brownlee PhD

and I **help developers** get results with **machine learning**.

[How to Calculate Feature Importance With Python](#)
[Recursive Feature Elimination \(RFE\) for Feature Selection in Python](#)
[How to Remove Outliers for Machine Learning](#)



Data Preparation for Machine Learning (7-Day Mini-Course)



How to Calculate Feature Importance With Python



Recursive Feature Elimination (RFE) for Feature Selection in Python



How to Remove Outliers for Machine Learning

Loving the Tutorials?

The [Data Preparation](#) EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Start Machine Learning



You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning

Never miss a tutorial:



Picked for you:



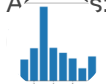
How to Choose a Feature Selection
Method For Machine Learning



Data Preparation for Machine Learning (7-
Day Mini-Course)

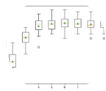
© 2020 Machine Learning Mastery Pty. Ltd. All Rights Reserved.

Address: PO Box 206, Vermont Victoria 3133, Australia. | ACN: 626 223 336.



How to Calculate Feature Importance
With Python

[Privacy](#) | [Disclaimer](#) | [Terms](#) | [Contact](#) | [Sitemap](#) | [Search](#)



Recursive Feature Elimination (RFE) for
Feature Selection in Python



How to Remove Outliers for Machine
Learning

Loving the Tutorials?

The [Data Preparation](#) EBook is
where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Start Machine Learning



You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Start Machine Learning