☐ Feedback ☐ Edit ☐ Share ☐ Bookmark Prepare data for building a model 09/10/2019 • 6 minutes to read • 😝 📵 🦡 🦚

Data is often unclean and sparse. ML.NET machine learning algorithms expect input or features to be in a single numerical vector. Similarly, the value to predict (label), especially when it's categorical data, has to be encoded. Therefore one of the goals of data preparation is to get the data into the format expected by ML.NET algorithms. Filter data

Get Started

Is this page

✓ Yes
✓ No

In this article

Replace missing

Use normalizers

categorical data

Work with text

Filter data

Work with

values

data

Sign in

helpful?

of interest. It's important to note that because filter operations are not an IEstimator or ITransformer like those in the TransformsCatalog, they cannot be included as part of an EstimatorChain or TransformerChain data preparation pipeline. Using the following input data which is loaded into an IDataView: Copy C C#

<u>IDataView</u> containing all of the data and return an <u>IDataView</u> containing only the data points

new HomeData

```
Price=100000f
      },
      new HomeData
          NumberOfBedrooms=2f,
          Price=300000f
      new HomeData
          NumberOfBedrooms=6f,
          Price=600000f
  };
To filter data based on the value of a column, use the <u>FilterRowsByColumn</u> method.
                                                                                   Copy C
  C#
  // Apply filter
  IDataView filteredData = mlContext.Data.FilterRowsByColumn(data, "Price", lowerBound
```

```
The sample above takes rows in the dataset with a price between 200000 and 1000000. The
result of applying this filter would return only the last two rows in the data and exclude the
first row because its price is 100000 and not between the specified range.
```

Using the following input data which is loaded into an IDataView: **С**ору C#

new HomeData

meaningful value such as the mean value in the data.

NumberOfBedrooms=1f, Price=100000f }, new HomeData

}, new HomeData

```
Notice that the last element in our list has a missing value for Price. To replace the missing
values in the Price column, use the <u>ReplaceMissingValues</u> method to fill in that missing value.
         (i) Important
          ReplaceMissingValue only works with numerical data.
                                                                                                                                                                                                                                                                                                                                                                                  Сору
          C#
           // Define replacement estimator
           var replacementEstimator = mlContext.Transforms.ReplaceMissingValues("Price", replacementEstimator = mlContext.Transforms.ReplacementEstimator = mlContext.Transforms.ReplacementEstimato
          // Fit data to estimator
          // Fitting generates a transformer that applies the operations of defined by estimat
          ITransformer replacementTransformer = replacementEstimator.Fit(data);
```

detailed list and description of normalization transforms. Min-Max normalization Using the following input data which is loaded into an IDataView:

result fills in the Price property for the last element in our data with 200,000 since it's the

new HomeData

// Define min-max estimator

// Fit data to estimator

// Transform data

NumberOfBedrooms = 2f,

Price = 200000f

average of 100,000 and 300,000.

NumberOfBedrooms = 1f, Price = 100000f Normalization can be applied to columns with single numerical values as well as vectors. Normalize the data in the Price column using min-max normalization with the

normalization formula which generates output values in the range of 0-1. Binning

Using the following input data which is loaded into an IDataView:

NumberOfBedrooms=1f,

NumberOfBedrooms=6f,

Price=100000f

Price=300000f

Price=600000f

new CarData

new CarData

new CarData

},

},

Color="Red",

Color="Blue",

Color="Black",

VehicleType="SUV"

VehicleType="SUV"

VehicleType="Sedan"

}, new HomeData NumberOfBedrooms=2f,

C#

parameter enables you to specify the number of bins needed to classify your data. In this example, data will be put into two bins. Copy C

```
var binningTransformer = binningEstimator.Fit(data);
  // Transform Data
  IDataView transformedData = binningTransformer.Transform(data);
The result of binning creates bin bounds of [0,200000,Infinity]. Therefore the resulting bins
are [0,1,1] because the first observation is between 0-200000 and the others are greater
than 200000 but less than infinity.
Work with categorical data
Non-numeric categorical data needs to be converted to a number before being used to build
a machine learning model.
Using the following input data which is loaded into an <a href="IDataView">IDataView</a>:
                                                                                    Copy C
  C#
```

OneHotEncoding method. C# // Define categorical transform estimator

```
The resulting transform converts the text value of VehicleType to a number. The entries in the
VehicleType column become the following when the transform is applied:
                                                                                        Copy C
  text
       1, // SUV
       2, // Sedan
       1 // SUV
Work with text data
Text data needs to be transformed into numbers before using it to build a machine learning
model. Visit the <u>transforms page</u> for a more detailed list and description of text transforms.
Using data like the data below that has been loaded into an <a href="IDataView">IDataView</a>:</a>
```

// Define text transform estimator var textEstimator = mlContext.Transforms.Text.FeaturizeText("Description"); // Fit data to estimator // Fitting generates a transformer that applies the operations of defined by estimat

new ReviewData

new ReviewData

},

};

Rating=4.7f

Rating=2.3f

Combine complex text processing steps into an EstimatorChain to remove noise and potentially reduce the amount of required processing resources as needed. C# // Define text transform estimator var textEstimator = mlContext.Transforms.Text.NormalizeText("Description") .Append(mlContext.Transforms.Text.TokenizeIntoWords("Description"))

benefit of a more complex pipeline is control and visibility over the transformations applied to the data. Using the first entry as an example, the following is a detailed description of the results produced by the transformation steps defined by textEstimator: **Original Text: This is a good product Description Transform** Result

["good","product"] 3. Removes stopwords like is and a. RemoveDefaultStopWords 4. MapValueToKey Maps the values to keys [1,2] (categories) based on the input data 5. ProduceNGrams [1,1,1,0,0] Transforms text into sequence of consecutive words [0.577350529, 0.577350529, 6. NormalizeLpNorm Scale inputs by their lp-norm 0.577350529, 0, 0] Feedback Send feedback about This product ☑ **7** This page

∨ Tutorials

> API

> Infer.NET

→ How-to guides

Load data

the model

> Model Builder & CLI

Download PDF

Prepare data

> Train, evaluate, and explain

Use the trained model

Concepts

 \vee API

> Model Builder & CLI

Learn how to use ML.NET to prepare data for additional processing or building a model.

Sometimes, not all data in a dataset is relevant for analysis. An approach to remove irrelevant data is filtering. The <u>DataOperationsCatalog</u> contains a set of filter operations that take in an

HomeData[] homeDataList = new HomeData[] NumberOfBedrooms=1f,

Replace missing values Missing values are a common occurrence in datasets. One approach to dealing with missing values is to replace them with the default value for the given type if any or another

HomeData[] homeDataList = new HomeData[] NumberOfBedrooms=2f,

Price=300000f NumberOfBedrooms=6f, **}**;

Use normalizers Normalization is a data pre-processing technique used to standardize features that are not on the same scale which helps algorithms converge faster. For example, the ranges for values like age and income vary significantly with age generally being in the range of 0-100 and income generally being in the range of zero to thousands. Visit the transforms page for a more

new HomeData }, **}**; NormalizeMinMax method.

C#

IDataView transformedData = minMaxTransformer.Transform(data);

HomeData[] homeDataList = new HomeData[] new HomeData }, new HomeData

}; Normalize the data into bins using the NormalizeBinning method. The maximumBinCount C# // Define binning estimator var binningEstimator = mlContext.Transforms.NormalizeBinning("Price", maximumBinCoun // Fit data to estimator // Fitting generates a transformer that applies the operations of defined by estimat

CarData[] cars = new CarData[]

}; // Fit data to estimator // Fitting generates a transformer that applies the operations of defined by estimat ITransformer categoricalTransformer = categoricalEstimator.Fit(data); // Transform Data IDataView transformedData = categoricalTransformer.Transform(data);

0 Open There are no open issues

English (United States)

normalized word and character ngrams.

C# ITransformer textTransformer = textEstimator.Fit(data); // Transform data IDataView transformedData = textTransformer.Transform(data); The resulting transform would convert the text values in the Description column to a numerical vector that looks similar to the output below:

C# ReviewData[] reviews = new ReviewData[]

text [0.2041241, 0.2041241, 0.2041241, 0.4082483, 0.2041241, 0.2041241, 0.2041241, 0.204

5 Closed

Price=float.NaN

// Transform data IDataView transformedData = replacementTransformer.Transform(data); ML.NET supports various <u>replacement modes</u>. The sample above uses the Mean replacement mode which will fill in the missing value with that column's average value. The replacement 's

C# HomeData[] homeDataList = new HomeData[]

1. NormalizeText Converts all letters to lowercase this is a good product by default ["this","is","a","good","product"] 2. TokenizeWords Splits string into individual words

> View on GitHub
> □ Previous Version Docs • Blog • Contribute • Privacy & Cookies • Terms of Use • Site Feedback • Trademarks

Copy Copy var minMaxEstimator = mlContext.Transforms.NormalizeMinMax("Price"); // Fitting generates a transformer that applies the operations of defined by estimat ITransformer minMaxTransformer = minMaxEstimator.Fit(data); The original price values [200000,100000] are converted to [1, 0.5] using the MinMax Binning converts continuous values into a discrete representation of the input. For example, suppose one of your features is age. Instead of using the actual age value, binning creates ranges for that value. 0-18 could be one bin, another could be 19-35 and so on. **С**ору

The categorical VehicleType property can be converted into a number using the Copy C var categoricalEstimator = mlContext.Transforms.Categorical.OneHotEncoding("VehicleT

Description="This is a good product", Description="This is a bad product", The minimum step to convert text to a numerical vector representation is to use the FeaturizeText method. By using the FeaturizeText transform, a series of transformations is applied to the input text column resulting in a numerical vector representing the lp-

Copy C

Copy Copy

Сору

Copy Copy .Append(mlContext.Transforms.Text.RemoveDefaultStopWords("Description")) .Append(mlContext.Transforms.Conversion.MapValueToKey("Description")) .Append(mlContext.Transforms.Text.ProduceNgrams("Description")) .Append(mlContext.Transforms.NormalizeLpNorm("Description")); textEstimator contains a subset of operations performed by the FeaturizeText method. The