Tikrit University Computer Science Dept.

> Master Degree Lecture -5-



Assistant Professor Dr. Eng. Zaidoon.T.AL-Qaysi

College of Computer Science and Mathematics (2023-2024)

- Many machine learning algorithms perform better when numerical input variables are scaled to a standard range. This includes algorithms that use a weighted sum of the input, like linear regression, and algorithms that use distance measures, like k-nearest neighbors. The two most popular techniques for scaling numerical data prior to modeling are normalization and standardization.
- Normalization scales each input variable separately to the range 0-1, which is the range for floating-point values where we have the most precision.
- Standardization scales each input variable separately by subtracting the mean (called centering) and dividing by the standard deviation to shift the distribution to have a mean of zero and a standard deviation of one.
- Data scaling is a recommended pre-processing step when working with many machine learning algorithms. It can be achieved by normalizing or standardizing real-valued input and output variables.
- Differences in the scales across input variables may increase the difficulty of the problem being modeled. An example of this is that large input values (e.g. a spread of hundreds or thousands of units) can result in a model that learns large weight values. A model with large weight values is often unstable, meaning that it may suffer from poor performance during learning and sensitivity to input values resulting in higher generalization error.
- The difference in scale for input variables does not affect all machine learning algorithms. For example, algorithms that fit a model that use a weighted sum of input variables are affected, such as linear regression, logistic regression, and artificial neural networks (deep learning).
- When the distance or dot products between predictors are used (such as K-nearest neighbors or support vector machines) or when the variables are required to be a common scale in order to apply a penalty, a standardization procedure is essential.
- Algorithms that use distance measures between examples are affected, such as k-nearest neighbors and support vector machines. There are also algorithms that are unaffected by the scale of numerical input variables, most notably decision trees and ensembles of trees, like random forest.
- It can also be a good idea to scale the target variable for regression predictive modeling problems to make the problem easier to learn, most notably in the case of neural network models. A target variable with a large spread of values, in turn, may result in large error gradient values causing weight values to change dramatically, making the learning process unstable. Scaling input and output variables is a critical step in using neural network models.
- Both normalization and standardization can be achieved using the scikit-learn library.

• Normalization is a rescaling of the data from the original range so that all values are within the new range of 0 and 1. Normalization requires that you know or are able to accurately estimate the minimum and maximum observable values. You may be able to estimate these values from your available data. Attributes are often normalized to lie in a fixed range - usually from zero to oneby dividing all values by the maximum value encountered or by subtracting the minimum value and dividing by the range between the maximum and minimum values.

A value is normalized as follows:

$$y = \frac{x - \min}{\max - \min}$$

Where the minimum and maximum values pertain to the value x being normalized. For example, for a dataset, we could guesstimate the min and max observable values as 30 and -10. We can then normalize any value, like 18.8, as follows:

$$y = \frac{x - \min}{\max - \min} \\ = \frac{18.8 - -10}{30 - -10} \\ = \frac{28.8}{40} \\ = 0.72$$

```
# example of a normalization
from numpy import asarray
from sklearn.preprocessing import MinMaxScaler
# define data
data = asarray([[100, 0.001],
       [8, 0.05],
       [50, 0.005],
       [88, 0.07],
       [4, 0.1]])
print(data)
# define min max scaler
scaler = MinMaxScaler()
# transform data
scaled = scaler.fit_transform(data)
print(scaled)
```

```
[[1.0e+02 1.0e-03]
[8.0e+00 5.0e-02]
[5.0e+01 5.0e-03]
[8.8e+01 7.0e-02]
[4.0e+00 1.0e-01]]
[[1. 0. ]
[0.04166667 0.49494949]
[0.47916667 0.04040404]
[0.875 0.6969697 ]
[0. 1. ]]
```

Example 1: Normalizing values in a dataset

• Standardizing a dataset involves rescaling the distribution of values so that the mean of observed values is 0 and the standard deviation is 1. This can be thought of as subtracting the mean value or centering the data. Like normalization, standardization can be useful and even required in some machine learning algorithms when your data has input values with differing scales. Standardization assumes that your observations fit a Gaussian distribution (bell curve) with a well-behaved mean and standard deviation.

A value is standardized as follows:

$$y = \frac{x - \text{mean}}{\text{standard_deviation}}$$

Where the mean is calculated as:

$$\mathrm{mean} = \frac{1}{N} \times \sum_{i=1}^{N} x_i$$

And the standard_deviation is calculated as:

standard_deviation =
$$\sqrt{\frac{\sum_{i=1}^{N} (x_i - \text{mean})^2}{N-1}}$$

We can guesstimate a mean of 10.0 and a standard deviation of about 5.0. Using these values, we can standardize the first value of 20.7 as follows:

$$y = \frac{x - \text{mean}}{\frac{\text{standard_deviation}}{5}}$$
$$= \frac{\frac{10.7}{5}}{\frac{10.7}{5}}$$
$$= 2.14$$

```
# example of a standardization
from numpy import asarray
from sklearn.preprocessing import StandardScaler
# define data
data = asarray([[100, 0.001],
       [8, 0.05],
       [50, 0.005],
       [88, 0.07],
       [4, 0.1]])
print(data)
# define standard scaler
scaler = StandardScaler()
# transform data
scaled = scaler.fit_transform(data)
print(scaled)
```

```
[[1.0e+02 1.0e-03]
[8.0e+00 5.0e-02]
[5.0e+01 5.0e-03]
[8.8e+01 7.0e-02]
[4.0e+00 1.0e-01]]
[[ 1.26398112 -1.16389967]
[-1.06174414 0.12639634]
[ 0. -1.05856939]
```

Example 2: Standardizing values in a dataset

•

```
# load and summarize the diabetes dataset
from pandas import read_csv
from matplotlib import pyplot
# load the dataset
dataset = read_csv('pima-indians-diabetes.csv', header=None)
# summarize the shape of the dataset
print(dataset.shape)
# summarize each variable
print(dataset.describe())
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

(768,	9)					
	0	1	2	6	7	8
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	0.331329	11.760232	0.476951
min	0.000000	0.00000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	2.420000	81.000000	1.000000

Example 3: Summarizing and Histogram Plots for the Diabetes Dataset.



Figure 1: Histogram Plots for the Diabetes Dataset.



Example 4: Evaluating model performance on the diabetes dataset

<pre>visualize a minmax scaler transform of the diabetes dataset rom pandas import read_csv rom pandas import DataFrame rom sklearn.preprocessing import MinMaxScaler rom matplotlib import pyplot load the dataset ataset = read_csv('pima-indians-diabetes.csv', header=None) retrieve just the numeric input values ata = dataset.values[:, :-1] perform a robust scaler transform of the dataset rans = MinMaxScaler() ata = trans.fit_transform(data) convert the array back to a dataframe ataset = DataFrame(data) summarize rint(dataset.describe()) histograms of the variables ig = dataset.hist(xlabelsize=4, ylabelsize=4) x.title.set_size(4) for x in fig.ravel()] show the plot yplot.show()</pre>
0 1 2 5 6 7 ount 768.000000 768.000000 768.000000 768.000000 768.000000
ean 0.226180 0.607510 0.566438 0.476790 0.168179 0.204015
td 0.198210 0.160666 0.158654 0.117499 0.141473 0.196004
in 0.000000 0.000000 0.000000 0.000000
5% 0.058824 0.497487 0.508197 0.406855 0.070773 0.050000
0% 0.176471 0.587940 0.590164 0.476900 0.125747 0.133333
5% 0.352941 0.704774 0.655738 0.545455 0.234095 0.333333
ax 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000

Example 5: Summarizing the variables from the diabetes dataset after a MinMaxScaler transform



Figure 2: Histogram Plots of MinMaxScaler Transformed Input Variables for the Diabetes Dataset

```
# evaluate knn on the diabetes dataset with minmax scaler transform
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import MinMaxScaler
from sklearn.pipeline import Pipeline
# load the dataset
dataset = read_csv('pima-indians-diabetes.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
trans = MinMaxScaler()
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('t', trans), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

```
Accuracy: 0.739 (0.053)
```

Example 6: Evaluating model performance after a MinMaxScaler transform



Example 7: Reviewing the data after a StandardScaler transform.

Histogram plots of the variables are created, although the distributions don't look much different from their original distributions seen in the previous section other than their scale on the x-axis. We can see that the center of mass for each distribution is centered on zero, which is more obvious for some variables than others.



Figure 3: Histogram Plots of StandardScaler Transformed Input Variables for the Diabetes Dataset

```
# evaluate knn on the diabetes dataset with standard scaler transform
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
# load the dataset
dataset = read_csv('pima-indians-diabetes.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
trans = StandardScaler()
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('t', trans), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.741 (0.050)
```

Example 8: Evaluating model performance after a StandardScaler transform

```
# load and summarize the diabetes dataset
from pandas import read_csv
from matplotlib import pyplot
# load dataset
dataset = read_csv('pima-indians-diabetes.csv', header=None)
# summarize the shape of the dataset
print(dataset.shape)
# summarize each variable
print(dataset.describe())
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

(768, 9)

0 1 2 ... 6 7 8 count 768.000000 768.000000 768.000000 ... 768.000000 768.000000 768.000000 3.845052 120.894531 69.105469 ... 0.471876 33.240885 0.348958 mean 3.369578 31.972618 19.355807 ... 0.331329 11.760232 0.476951 std0.000000 ... 0.078000 21.000000 0.000000 0.00000 0.00000 min 1.000000 99.000000 62.000000 ... 0.243750 24.000000 25% 0.000000 50% 3.000000 117.000000 72.000000 ... 0.372500 29.000000 0.000000 6.000000 140.250000 80.000000 ... 0.626250 41.000000 1.000000 75% 17.000000 199.000000 122.000000 ... 2.420000 81.000000 1.000000 max

Example 9: Example of loading and summarizing the diabetes dataset



Figure 4: Histogram Plots of Input Variables for the Diabetes Dataset.

Sometimes an input variable may have outlier values. These are values on the edge of the distribution that may have a low probability of occurrence, yet are overrepresented for some reason. Outliers can skew a probability distribution and make data scaling using standardization difficult as the calculated mean and standard deviation will be skewed by the presence of the outliers. One approach to standardizing input variables in the presence of outliers is to ignore the outliers from the calculation of the mean and standard deviation, then use the calculated values to scale the variable. This is called robust standardization or robust data scaling. This can be achieved by calculating the median (50th percentile) and the 25th and 75th percentiles. The values of each variable then have their median subtracted and are divided by the interquartile range (IQR) which is the difference between the 75th and 25th percentiles

value =
$$\frac{\text{value} - \text{median}}{p_{75} - p_{25}}$$

The resulting variable has a zero mean and median and a standard deviation of 1, although not skewed by outliers and the outliers are still present with the same relative relationships to other values.

Interestingly, the definition of the scaling range can be specified via the quantile range argument. It takes a tuple of two integers between 0 and 100 and defaults to the percentile values of the IQR, specifically (25, 75).



Example 10: Evaluating model performance after a RobustScaler transform.

• Histogram plots of the variables are created, although the distributions don't look much different from their original distributions seen in the previous section. We can see that the center of mass for each distribution is now close to zero



Figure 5: Histogram Plots of Robust Scaler Transformed Input Variables for the Diabetes Dataset

```
# explore the scaling range of the robust scaler transform
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import RobustScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.pipeline import Pipeline
from matplotlib import pyplot
# get the dataset
def get_dataset():
 # load dataset
  dataset = read_csv('pima-indians-diabetes.csv', header=None)
  data = dataset.values
  # ensure inputs are floats and output is an integer label
  X = X.astype('float32')
  y = LabelEncoder().fit_transform(y.astype('str'))
  return X, y
# get a list of models to evaluate
def get_models():
 models = dict()
  for value in [1, 5, 10, 15, 20, 25, 30]:
   # define the pipeline
   trans = RobustScaler(quantile_range=(value, 100-value))
   model = KNeighborsClassifier()
   models[str(value)] = Pipeline(steps=[('t', trans), ('m', model)])
  return models
# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
 cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
  scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
 return scores
# define dataset
X, y = get_dataset()
# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
 scores = evaluate_model(model, X, y)
 results.append(scores)
 names.append(name)
 print('>%s %.3f (%.3f)' % (name, mean(scores), std(scores)))
# plot model performance for comparison
pyplot.boxplot(results, labels=names, showmeans=True)
pyplot.show()
>1 0.734 (0.054)
>5 0.736 (0.051)
>10 0.739 (0.047)
>15 0.740 (0.045)
>20 0.734 (0.050)
>25 0.734 (0.044)
>30 0.735 (0.042)
```

Example 11: Comparing model Performance with different ranges for the RobustScaler transform



Figure 6: Box Plots of Robust Scaler IQR Range vs Classification Accuracy of KNN

- Numerical data, as its name suggests, involves features that are only composed of numbers, such as integers or floating-point values. Categorical data are variables that contain label values rather than numeric values. The number of possible values is often limited to a fixed set. Categorical variables are often called nominal. Some examples include: ^
 - 1. A pet variable with the values: dog and cat.
 - 2. A color variable with the values: red, green, and blue.
 - 3. A place variable with the values: first, second, and third.

```
# example of a ordinal encoding
from numpy import asarray
from sklearn.preprocessing import OrdinalEncoder
# define data
data = asarray([['red'], ['green'], ['blue']])
print(data)
# define ordinal encoding
encoder = OrdinalEncoder()
# transform data
result = encoder.fit_transform(data)
print(result)
[['red']
['green']
['blue']]
```

['blue']] [[2.] [1.] [0.]]

Example 12: Demonstrating an ordinal encoding of color categories

```
# example of a one hot encoding
from numpy import asarray
from sklearn.preprocessing import OneHotEncoder
# define data
data = asarray([['red'], ['green'], ['blue']])
print(data)
# define one hot encoding
encoder = OneHotEncoder(sparse=False)
# transform data
onehot = encoder.fit_transform(data)
print(onehot)
```

[['red']
['green']
['blue']]
[[0. 0. 1.]
[0. 1. 0.]
[1. 0. 0.]]

Example 13: Demonstrating a one hot encoding of color categories

```
# ordinal encode the breast cancer dataset
from pandas import read_csv
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OrdinalEncoder
# load the dataset
dataset = read_csv('breast-cancer.csv', header=None)
# retrieve the array of data
data = dataset.values
# separate into input and output columns
X = data[:, :-1].astype(str)
y = data[:, -1].astype(str)
# ordinal encode input variables
ordinal_encoder = OrdinalEncoder()
X = ordinal_encoder.fit_transform(X)
# ordinal encode target variable
label_encoder = LabelEncoder()
y = label_encoder.fit_transform(y)
# summarize the transformed data
print('Input', X.shape)
print(X[:5, :])
print('Output', y.shape)
print(y[:5])
Input (286, 9)
```

[[2. 2. 2. 0. 1. 2. 1. 2. 0.] [3. 0. 2. 0. 0. 0. 1. 0. 0.] [3. 0. 6. 0. 0. 1. 0. 1. 0.] [2. 2. 6. 0. 1. 2. 1. 1. 1.] [2. 2. 5. 4. 1. 1. 0. 4. 0.]] Output (286,) [1 0 1 0 1]

Example 14: Ordinal encoding of the breast cancer dataset.



Example 15: Evaluating a model on the breast cancer dataset with an ordinal encoding



Example 16: Evaluating a model on the breast cancer dataset with an one hot encoding

- Machine learning algorithms like Linear Regression and Gaussian Naive Bayes assume the numerical variables have a Gaussian probability distribution. Your data may not have a Gaussian distribution and instead may have a Gaussian-like distribution (e.g. nearly Gaussian but with outliers or a skew) or a totally different distribution (e.g. exponential). As such, you may be able to achieve better performance on a wide range of machine learning algorithms by transforming input and/or output variables to have a Gaussian or more Gaussian distribution. Power transforms like the Box-Cox transform and the Yeo-Johnson transform provide an automatic way of performing these transforms on your data and are provided in the scikit-learn Python machine learning library.
- A power transform will make the probability distribution of a variable more Gaussian. This is often described as removing a skew in the distribution, although more generally is described as stabilizing the variance of the distribution.

```
# demonstration of the power transform on data with a skew
from numpy import exp
from numpy.random import randn
from sklearn.preprocessing import PowerTransformer
from matplotlib import pyplot
# generate gaussian data sample
data = randn(1000)
# add a skew to the data distribution
data = exp(data)
# histogram of the raw data with a skew
pyplot.hist(data, bins=25)
pyplot.show()
# reshape data to have rows and columns
data = data.reshape((len(data),1))
# power transform the raw data
power = PowerTransformer(method='yeo-johnson', standardize=True)
data_trans = power.fit_transform(data)
# histogram of the transformed data
pyplot.hist(data_trans, bins=25)
pyplot.show()
```

Example 17: Demonstration of the effect of the power transform on a skewed data distribution



Figure 7: Histogram of Skewed Gaussian Distribution.



Figure 8: Histogram of Skewed Gaussian Data After Power Transform.

```
# load and summarize the sonar dataset
from pandas import read_csv
from matplotlib import pyplot
# load dataset
dataset = read_csv('sonar.csv', header=None)
# summarize the shape of the dataset
print(dataset.shape)
# summarize each variable
print(dataset.describe())
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

Running the example first summarizes the shape of the loaded dataset. This confirms the 60 input variables, one output variable, and 208 rows of data. A statistical summary of the input variables is provided showing that values are numeric and range approximately from 0 to 1.

(208,	61)					
	0	1	2	57	58	59
count	208.000000	208.000000	208.000000	208.000000	208.000000	208.000000
mean	0.029164	0.038437	0.043832	0.007949	0.007941	0.006507
std	0.022991	0.032960	0.038428	0.006470	0.006181	0.005031
min	0.001500	0.000600	0.001500	0.000300	0.000100	0.000600
25%	0.013350	0.016450	0.018950	0.003600	0.003675	0.003100
50%	0.022800	0.030800	0.034300	0.005800	0.006400	0.005300
75%	0.035550	0.047950	0.057950	0.010350	0.010325	0.008525
max	0.137100	0.233900	0.305900	0.044000	0.036400	0.043900

Example 18: Summarizing the variables from the sonar dataset

In **Figure 9**, a histogram is created for each input variable. If we ignore the clutter of the plots and focus on the histograms themselves, we can see that many variables have a skewed distribution. The dataset provides a good candidate for using a power transform to make the variables more Gaussian.



Figure 9: Histogram Plots of Input Variables for the Sonar Binary Classification Dataset

```
# evaluate knn on the raw sonar dataset
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define and configure the model
model = KNeighborsClassifier()
# evaluate the model
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report model performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.797 (0.073)
```

Example 19: Example of evaluating model performance on the sonar dataset

- The Box-Cox transform is named for the two authors of the method. It is a power transform that assumes the values of the input variable to which it is applied are strictly positive. That means 0 and negative values are not supported.
- The Sonar dataset does not have negative values but may have zero values. This may cause a problem. Let's try anyway. The complete example of creating a Box-Cox transform of the sonar dataset and plotting histograms of the result is listed below.

```
# visualize a box-cox transform of the sonar dataset
from pandas import read_csv
from pandas import DataFrame
from sklearn.preprocessing import PowerTransformer
from matplotlib import pyplot
# Load dataset
dataset = read_csv('sonar.csv', header=None)
# retrieve just the numeric input values
data = dataset.values[:, :-1]
# perform a box-cox transform of the dataset
pt = PowerTransformer(method='box-cox')
# NOTE: we expect this to cause an error !!!
data = pt.fit_transform(data)
# convert the array back to a dataframe
dataset = DataFrame(data)
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

ValueError: The Box-Cox transformation can only be applied to strictly positive data

Example 20: Applying the Box-Cox transform in a way that results in an error.

```
# visualize a box-cox transform of the scaled sonar dataset
from pandas import read_csv
from pandas import DataFrame
from sklearn.preprocessing import PowerTransformer
from sklearn.preprocessing import MinMaxScaler
from sklearn.pipeline import Pipeline
from matplotlib import pyplot
# Load dataset
dataset = read_csv('sonar.csv', header=None)
# retrieve just the numeric input values
data = dataset.values[:, :-1]
# perform a box-cox transform of the dataset
scaler = MinMaxScaler(feature_range=(1, 2))
power = PowerTransformer(method='box-cox')
pipeline = Pipeline(steps=[('s', scaler),('p', power)])
data = pipeline.fit_transform(data)
# convert the array back to a dataframe
dataset = DataFrame(data)
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

Example 21: Summarizing the sonar dataset after applying a Box-Cox transform



Figure 10: Histogram Plots of Box-Cox Transformed Input Variables for the Sonar Dataset.

```
# evaluate knn on the box-cox sonar dataset
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import PowerTransformer
from sklearn.preprocessing import MinMaxScaler
from sklearn.pipeline import Pipeline
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
scaler = MinMaxScaler(feature_range=(1, 2))
power = PowerTransformer(method='box-cox')
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('s', scaler),('p', power), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

Accuracy: 0.811 (0.085)

Example 22: Evaluating a model on the sonar dataset after applying a Box-Cox transform

• The Yeo-Johnson transform is also named for the authors. Unlike the Box-Cox transform, it does not require the values for each input variable to be strictly positive. It supports zero values and negative values. This means we can apply it to our dataset without scaling it first. We can apply the transform by defining a PowerTransformer object and setting the method argument to 'yeo-johnson' (the default).

<pre># visualize a yeo-johnson transform of the sonar dataset</pre>						
from pandas import read_csv						
from pandas import DataFrame						
from sklearn.preprocessing import PowerTransformer						
from matplotlib import pyplot						
# Load dataset						
<pre>dataset = read_csv('sonar.csv', header=None)</pre>						
<pre># retrieve just the numeric input values</pre>						
<pre>data = dataset.values[:, :-1]</pre>						
<pre># perform a yeo-johnson transform of the dataset</pre>						
<pre>pt = PowerTransformer(method='yeo-johnson')</pre>						
<pre>data = pt.fit_transform(data)</pre>						
# convert the array back to a dataframe						
dataset = DataFrame(data)						
# histograms of the variables						
<pre>fig = dataset.hist(xlabelsize=4, ylabelsize=4)</pre>						
<pre>[x.title.set_size(4) for x in fig.ravel()]</pre>						
# show the plot						
<pre>pyplot.show()</pre>						





Figure 11: Histogram Plots of Yeo-Johnson Transformed Input Variables for the Sonar Dataset

```
# evaluate knn on the yeo-johnson sonar dataset
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import PowerTransformer
from sklearn.pipeline import Pipeline
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
power = PowerTransformer(method='yeo-johnson')
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('p', power), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

Accuracy: 0.808 (0.082)

Example 24: Evaluating a model on the sonar dataset after applying a Yeo-Johnson transform

```
# evaluate knn on the yeo-johnson standardized sonar dataset
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import PowerTransformer
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
scaler = StandardScaler()
power = PowerTransformer(method='yeo-johnson')
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('s', scaler), ('p', power), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

```
Accuracy: 0.816 (0.077)
```

Example 25: Evaluating a model on after applying a StandardScaler and Yeo-Johnson transforms

- Running the example, we can see that the Yeo-Johnson transform results in a lift in performance from 79.7 percent accuracy without the transform to about 80.8 percent with the transform, less than the Box-Cox transform that achieved about 81.1 percent. Sometimes a lift in performance can be achieved by first standardizing the raw dataset prior to performing a Yeo-Johnson transform. We can explore this by adding a StandardScaler as a first step in the pipeline. The complete example is listed below.
- Running the example, we can see that standardizing the data prior to the Yeo-Johnson transform resulted in a small lift in performance from about 80.8 percent to about 81.6 percent, a small lift over the results for the Box-Cox transform.
- A quantile transform will map a variable's probability distribution to another probability distribution. Recall that a quantile function, also called a percent-point function (PPF), is the inverse of the cumulative probability distribution (CDF). A CDF is a function that returns the probability of a value at or below a given value. The PPF is the inverse of this function and returns the value at or below a given probability.
- The quantile function ranks or smooths out the relationship between observations and can be mapped onto other distributions, such as the uniform or normal distribution. The transformation can be applied to each numeric input variable in the training dataset and then provided as input to a machine learning model to learn a predictive modeling task. This quantile transform is available in the scikit-learn Python machine learning library via the QuantileTransformer class.
- The class has an output distribution argument that can be set to 'uniform' or 'normal' and defaults to 'uniform'. It also provides a n quantiles that determines the resolution of the mapping or ranking of the observations in the dataset. This must be set to a value less than the number of observations in the dataset and defaults to 1,000.

```
# demonstration of the quantile transform
from numpy import exp
from numpy.random import randn
from sklearn.preprocessing import QuantileTransformer
from matplotlib import pyplot
# histogram of the raw data with a skew
pyplot.hist(data, bins=25)
pyplot.show()
# reshape data to have rows and columns
data = data.reshape((len(data),1))
# quantile transform the raw data
quantile = QuantileTransformer(output_distribution='normal')
data_trans = quantile.fit_transform(data)
# histogram of the transformed data
pyplot.hist(data_trans, bins=25)
pyplot.show()
```

Example 26: Demonstration of the effect of the quantile transform on a skewed data distribution



Figure 12: Histogram of Skewed Gaussian Distribution



Figure 13: Histogram of Skewed Gaussian Data After Quantile Transform

```
# evaluate knn on the raw sonar dataset
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define and configure the model
model = KNeighborsClassifier()
# evaluate the model
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report model performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

Accuracy: 0.797 (0.073)

Example 27: Evaluating model performance on the sonar dataset

```
# visualize a normal quantile transform of the sonar dataset
from pandas import read_csv
from pandas import DataFrame
from sklearn.preprocessing import QuantileTransformer
from matplotlib import pyplot
# load dataset
dataset = read_csv('sonar.csv', header=None)
# retrieve just the numeric input values
data = dataset.values[:, :-1]
# perform a normal quantile transform of the dataset
trans = QuantileTransformer(n_quantiles=100, output_distribution='normal')
data = trans.fit_transform(data)
# convert the array back to a dataframe
dataset = DataFrame(data)
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

Example 28: Summarizing the sonar dataset after applying a normal quantile transform



Figure 14: Histogram Plots of Normal Quantile Transformed Input Variables for the Sonar Dataset

```
# evaluate knn on the sonar dataset with normal quantile transform
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import QuantileTransformer
from sklearn.pipeline import Pipeline
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
trans = QuantileTransformer(n_quantiles=100, output_distribution='normal')
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('t', trans), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

```
Accuracy: 0.817 (0.087)
```

Example 29: Evaluating a model on the sonar dataset after applying a normal quantile transform

```
# visualize a uniform quantile transform of the sonar dataset
from pandas import read_csv
from pandas import DataFrame
from sklearn.preprocessing import QuantileTransformer
from matplotlib import pyplot
# load dataset
dataset = read_csv('sonar.csv', header=None)
# retrieve just the numeric input values
data = dataset.values[:, :-1]
# perform a uniform quantile transform of the dataset
trans = QuantileTransformer(n_quantiles=100, output_distribution='uniform')
data = trans.fit_transform(data)
# convert the array back to a dataframe
dataset = DataFrame(data)
# histograms of the variables
fig = dataset.hist(xlabelsize=4, ylabelsize=4)
[x.title.set_size(4) for x in fig.ravel()]
# show the plot
pyplot.show()
```

Example 30: Summarizing the sonar dataset after applying a uniform quantile transform.



Figure 15: Histogram Plots of Uniform Quantile Transformed Input Variables for the Sonar Dataset

26

```
# explore number of quantiles on classification accuracy
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import QuantileTransformer
from sklearn.preprocessing import LabelEncoder
from sklearn.pipeline import Pipeline
from matplotlib import pyplot
# get the dataset
def get_dataset(filename):
 # load dataset
 dataset = read_csv(filename, header=None)
 data = dataset.values
  # separate into input and output columns
 X, y = data[:, :-1], data[:, -1]
  # ensure inputs are floats and output is an integer label
 X = X.astype('float32')
 y = LabelEncoder().fit_transform(y.astype('str'))
 return X, y
# get a list of models to evaluate
def get_models():
 models = dict()
 for i in range(1,100):
   # define the pipeline
   trans = QuantileTransformer(n_quantiles=i, output_distribution='uniform')
   model = KNeighborsClassifier()
   models[str(i)] = Pipeline(steps=[('t', trans), ('m', model)])
 return models
# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
 cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
 scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
 return scores
# define dataset
X, y = get_dataset('sonar.csv')
# get the models to evaluate
models = get_models()
# evaluate the models and store results
results = list()
for name, model in models.items():
 scores = evaluate_model(model, X, y)
 results.append(mean(scores))
 print('>%s %.3f (%.3f)' % (name, mean(scores), std(scores)))
# plot model performance for comparison
pyplot.plot(results)
pyplot.show()
>1 0.466 (0.016)
>2 0.813 (0.085)
>3 0.840 (0.080)
>4 0.854 (0.075)
>5 0.848 (0.072)
>6 0.851 (0.071)
>7 0.845 (0.071)
```

>10 0.843 (0.074)

>8 0.848 (0.066) >9 0.848 (0.071)

Example 31: Comparing the number of partitions for the dataset in a uniform quantile transform



Figure 16: Line Plot of Number of Quantiles vs. Classification Accuracy of KNN

```
# evaluate knn on the sonar dataset with uniform quantile transform
from numpy import mean
from numpy import std
from pandas import read_csv
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import QuantileTransformer
from sklearn.pipeline import Pipeline
# load dataset
dataset = read_csv('sonar.csv', header=None)
data = dataset.values
# separate into input and output columns
X, y = data[:, :-1], data[:, -1]
# ensure inputs are floats and output is an integer label
X = X.astype('float32')
y = LabelEncoder().fit_transform(y.astype('str'))
# define the pipeline
trans = QuantileTransformer(n_quantiles=100, output_distribution='uniform')
model = KNeighborsClassifier()
pipeline = Pipeline(steps=[('t', trans), ('m', model)])
# evaluate the pipeline
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report pipeline performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
Accuracy: 0.845 (0.074)
```

Example 32: Evaluating a model on the sonar dataset after applying a uniform quantile transform.