

✓ 100 XP

Exercise - Train and evaluate a regression model

8 minutes

Sandbox activated! Time remaining: **1 hr 59 min**

You have used 3 of 10 sandboxes for today. More sandboxes will be available tomorrow.

 Runtime

File

Edit

View

 Save

 Export as ▾

Regression

Supervised machine learning techniques involve training a model to operate on a set of

```
import pandas as pd

# load the training dataset
!wget https://raw.githubusercontent.com/MicrosoftDocs/mslearn-introduction-to-machine-learning/main/example-data/daily-bike-share.csv
bike_data = pd.read_csv('daily-bike-share.csv')
bike_data.head()
```

In most cases, x is actually a *vector* that consists of multiple feature values, so to be a little more precise, the function could be expressed like this:

The data consists of the following columns: $y = f([x_1, x_2, x_3, \dots])$

The **instantaneous identifier** is a function that performs some kind of calculation to the data to produce the result. In this case, the data was collected daily, so there's only one possible date. The **weather** is a calculation that produces y reasonably accurately for all of the cases in the training dataset.

There are many machine learning algorithms for supervised learning, and we can broadly divide them into two types:

- **mnth**: The calendar month in which the observation was made (1:January ... 12:December)
- **Regression algorithms**: Algorithms that predict a y value that is a numeric value, such as the price of a house or the number of sales transactions
- **holiday**: A binary value indicating whether or not the observation was made on a public holiday
- **Classification algorithms**: Algorithms that predict to which category, or *class*, an observation belongs. The y value in a classification model is a vector of probability values between 0 and 1, one for each class, indicating the probability of the observation belonging to each class
- **workingday**: A binary value indicating whether or not the day is a working day (not a weekend or holiday)
- **weathersit**: A categorical value indicating the weather situation (1:clear, 2:mist/cloud, 3:light rain/snow, 4:heavy rain/hail/snow/fog)
- **temp**: The temperature in celsius (normalized)
- **atemp**: The apparent ("feels-like") temperature in celsius (normalized)
- **hum**: The humidity level (normalized)
- **windspeed**: The wind speed (normalized)
- **rentals**: The number of bicycle rentals recorded.

Explore the Data

In this dataset, **rentals** represents the label (the y value) we must train our model to

The first step in any machine-learning project is to explore the data that you'll use to train

```
bike_data['day'] = pd.DatetimeIndex(bike_data['dteday']).day
bike_data.head(32)
```

OK, let's start our analysis of the data by examining a few key descriptive statistics. We can use the dataframe's **describe** method to generate these for the numeric features as well as the rental label.

```
numeric_features = ['temp', 'atemp', 'hum', 'windspeed']
bike_data[numeric_features + ['rentals']].describe()
```

The statistics reveal some information about the distribution of the data in each of the numeric fields, including the number of observations (there are 731 records), the mean, standard deviation, minimum and maximum values, and the quartile values (the threshold values for 25%, 50%, which is also the median and 75% of the data). From this, we can see that the mean number of daily rentals is around 848; but there's a comparatively large

```
import pandas as pd
import matplotlib.pyplot as plt

# This ensures plots are displayed inline in the Jupyter notebook
%matplotlib inline

# Get the label column
label = bike_data['rentals']

# Create a figure for 2 subplots (2 rows, 1 column)
fig, ax = plt.subplots(2, 1, figsize = (9,12))

# Plot the histogram
ax[0].hist(label, bins=100)
ax[0].set_ylabel('Frequency')

# Add lines for the mean, median, and mode
ax[0].axvline(label.mean(), color='magenta', linestyle='dashed', linewidth=2)
ax[0].axvline(label.median(), color='cyan', linestyle='dashed', linewidth=2)

# Plot the boxplot
ax[1].boxplot(label, vert=False)
ax[1].set_xlabel('Rentals')

# Add a title to the Figure
fig.suptitle('Rental Distribution')

# Show the figure
fig.show()
```

The plots show that the number of daily rentals ranges from 0 to just over 3,400. However, the mean (and median) number of daily rentals is closer to the low end of that range, with most of the data between 0 and around 2,200 rentals. The few values above this are shown in the box plot as small circles, indicating that they are *outliers*; in other words, unusually high or low values beyond the typical range of most of the data.

We can do the same kind of plot for each of the numeric features. Let's create a

```
# Plot a histogram for each numeric feature
for col in numeric_features:
    fig = plt.figure(figsize=(9, 6))
    ax = fig.gca()
    feature = bike_data[col]
    feature.hist(bins=100, ax = ax)
    ax.axvline(feature.mean(), color='magenta', linestyle='dashed', linewidth=2)
    ax.axvline(feature.median(), color='cyan', linestyle='dashed', linewidth=2)
    ax.set_title(col)
plt.show()
```

The numeric features seem to be more *normally* distributed, with the mean and median nearer the middle of the range of values, coinciding with where the most commonly occurring values are.

Note: The distributions are not truly *normal* in the statistical sense, which would result in a smooth, symmetric "bell-curve" histogram with the mean and mode (the most common value) in the center; but they do generally indicate that most of the observations have a value somewhere near the middle.

```
import numpy as np

# plot a bar plot for each categorical feature count
categorical_features = ['season', 'mnth', 'holiday', 'weekday', 'workingday', 'weathersit', 'registered', 'casual']

for col in categorical_features:
    counts = bike_data[col].value_counts().sort_index()
    fig = plt.figure(figsize=(9, 6))
    ax = fig.gca()
    counts.plot.bar(ax=ax, color='steelblue')
    ax.set_title(col + ' counts')
    ax.set_xlabel(col)
    ax.set_ylabel("Frequency")
plt.show()
```

Many of the categorical features show a more or less *uniform* distribution (meaning there's roughly the same number of rows for each category). Exceptions to this include:

- **holiday:** There are many fewer days that are holidays than days that aren't.
- **workingday:** There are more working days than non-working days.
- **weathersit:** Most days are category 1 (clear), with category 2 (mist and cloud) the next most common. There are comparatively few category 3 (light rain or snow) days, and no category 4 (heavy rain, hail, or fog) days at all.

Now that we know something about the distribution of the data in our columns, we can start to look for relationships between the features and the **rentals** label we want to be able to predict.

For the numeric features, we can create scatter plots that show the relationship of features

```
for col in numeric_features:
    fig = plt.figure(figsize=(9, 6))
    ax = fig.gca()
    feature = bike_data[col]
    label = bike_data['rentals']
    correlation = feature.corr(label)
    plt.scatter(x=feature, y=label)
    plt.xlabel(col)
    plt.ylabel('Bike Rentals')
    ax.set_title('rentals vs ' + col + '- correlation: ' + str(correlation))
plt.show()
```

The results aren't conclusive, but if you look closely at the scatter plots for **temp** and **atemp**, you can see a vague diagonal trend showing that higher rental counts tend to coincide with higher temperatures, and a correlation value of just over 0.5 for both of these features supports this observation. Conversely, the plots for **hum** and **windspeed** show a slightly negative correlation, indicating that there are fewer rentals on days with high humidity or windspeed.

Next, let's compare the categorical features to the label. We'll do this by creating bar plots

```
# plot a boxplot for the label by each categorical feature
for col in categorical_features:
    fig = plt.figure(figsize=(9, 6))
    ax = fig.gca()
    bike_data.boxplot(column = 'rentals', by = col, ax = ax)
    ax.set_title('Label by ' + col)
    ax.set_ylabel("Bike Rentals")
plt.show()
```

The plots show some variance in the relationship between some category values and rentals. For example, there's a clear difference in the distribution of rentals on weekends (**weekday** 0 or 6) and those during the working week (**weekday** 1 to 5). Similarly, there are notable differences for **holiday** and **workingday** categories. There's a noticeable trend that shows different rental distributions in spring and summer months compared to winter and fall months. The **weathersit** category also seems to make a difference in rental distribution. The **temp** feature is expected for the day of the month that a little variation

Train a Regression Model

Now that we've explored the data, it's time to use it to train a regression model that uses the features we've identified as potentially predictive to predict the **rentals** label. The first thing we need to do is to separate the features we want to use to train the model from the label.

```
# Separate features and labels
X, y = bike_data[['season', 'mnth', 'holiday', 'weekday', 'workingday', 'weathersit', 'temp'],
print('Features:', X[:10], '\nLabels:', y[:10], sep='\n')
```

After separating the dataset, we now have numpy arrays named **X** containing the features and **y** containing the labels.

We *could* train a model using all of the data, but it's common practice in supervised learning to split the data into two subsets: a (typically larger) set with which to train the model, and a smaller "hold-back" set with which to validate the trained model. This allows us to evaluate how well the model performs when used with the validation dataset by comparing the predicted labels to the known labels. It's important to split the data *randomly* (rather than say, taking the first 70% of the data for training and keeping the rest for validation). This helps ensure that the two subsets of data are statistically comparable (so we validate the model with data that has a similar statistical distribution to the data on which it was trained).

```
from sklearn.model_selection import train_test_split

# Split data 70%-30% into training set and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=0)

print ('Training Set: %d rows\nTest Set: %d rows' % (X_train.shape[0], X_test.shape[0]))
```

Now we have the following four datasets:

- **X_train**: The feature values we'll use to train the model
- **y_train**: The corresponding labels we'll use to train the model
- **X_test**: The feature values we'll use to validate the model
- **y_test**: The corresponding labels we'll use to validate the model

Now we're ready to train a model by fitting a suitable regression algorithm to the training data. We'll use a *linear regression* algorithm, a common starting point for regression that works by trying to find a linear relationship between the X values and the y label. The resulting model is a function that conceptually defines a line where every possible X and y value combination intersect.

Let's fit a linear regression model to the training data and print the model:

```
# Train the model
from sklearn.linear_model import LinearRegression

# Fit a linear regression model on the training set
model = LinearRegression().fit(X_train, y_train)
print (model)
```

Evaluate the Trained Model

Now that we've trained the model, we can use it to predict rental counts for the features we held back in our validation dataset. Then we can compare these predictions to the actual label values to evaluate how well (or not!) the model is working.

Let's use the model to predict the rental counts for the test data:

```
import numpy as np

predictions = model.predict(X_test)
np.set_printoptions(suppress=True)
print('Predicted labels: ', np.round(predictions)[:10])
print('Actual labels    : ', y_test[:10])
```

Comparing each prediction with its corresponding "ground truth" actual value isn't a very efficient way to determine how well the model is predicting. Let's see if we can get a better indication by visualizing a scatter plot that compares the predictions to the actual labels. We'll also overlay a trend line to get a general sense for how well the predicted labels align with the actual labels.

Let's use matplotlib to create a scatter plot of the predictions vs. the actual labels:

```
import matplotlib.pyplot as plt

%matplotlib inline

plt.scatter(y_test, predictions)
plt.xlabel('Actual Labels')
plt.ylabel('Predicted Labels')
plt.title('Daily Bike Share Predictions')
# overlay the regression line
z = np.polyfit(y_test, predictions, 1)
p = np.poly1d(z)
plt.plot(y_test,p(y_test), color='magenta')
plt.show()
```

There's a definite diagonal trend, and the intersections of the predicted and actual values are generally following the path of the trend line, but there's a fair amount of difference between the ideal function represented by the line and the results. This variance represents the *residuals* of the model; in other words, the difference between the label predicted when the model applies the coefficients it learned during training to the validation data and the actual value of the validation label. These residuals when evaluated from the validation data indicate the expected level of *error* when the model is used with new data for which the label is unknown.

You can quantify the residuals by calculating a number of commonly used evaluation metrics. We'll focus on the following three:

- **Mean Square Error (MSE):** The mean of the squared differences between predicted and actual values. This yields a relative

```
from sklearn.metrics import mean_squared_error, r2_score

mse = mean_squared_error(y_test, predictions)
print("MSE:", mse)

rmse = np.sqrt(mse)
print("RMSE:", rmse)

r2 = r2_score(y_test, predictions)
print("R2:", r2)
```

So now, we've quantified the ability of our model to predict the number of rentals. It definitely has *some* predictive power, but we can probably do better!

Summary

Here, we've explored our data and fit a basic regression model. In the next notebook, we'll try a number of other regression algorithms to improve performance.

🔗 learn-notebooks-b4562f0e-4b02-43b6-91e7-8dbb4c8bec99 Compute not connected ✎ Viewing Kernel not connected

Next unit: Discover new regression models

Continue >

How are we doing? ☆ ☆ ☆ ☆ ☆