

✓ 100 XP ▶

Exercise - Experiment with more powerful regression models

10 minutes

Sandbox activated! Time remaining: 1 hr 10 min

You have used 4 of 10 sandboxes for today. More sandboxes will be available tomorrow.

✓ Runtime File Edit View

💬 Comments

▶▶ Run all ∨ 🔗 azureml_py38 ∨ ↻ 🔒

🗣️ We would love to hear your feedback on the notebooks experience! Please take a few minutes to [complete our survey](#). ✕

Regression - Experimenting with additional models

In the previous notebook, we used simple regression models to look at the relationship between features of a bike rentals dataset. In this notebook, we'll experiment with more complex models to improve our regression performance.

Let's start by loading the bicycle sharing data as a **Pandas** DataFrame and viewing the first few rows. We'll also split the data into training and test datasets.

```
# Import modules we'll need for this notebook
import pandas as pd
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.model_selection import train_test_split
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

# load the training dataset
!wget https://raw.githubusercontent.com/MicrosoftDocs/mslearn-introduction-to-machine-learning/main/notebooks/daily-bike-share.csv
bike_data = pd.read_csv('daily-bike-share.csv')
bike_data['day'] = pd.DatetimeIndex(bike_data['dteday']).day
numeric_features = ['temp', 'atemp', 'hum', 'windspeed']
categorical_features = ['season', 'mnth', 'holiday', 'weekday', 'workingday', 'weathersit', 'casual', 'registered', 'rentals']
bike_data[numeric_features + categorical_features].describe()
print(bike_data.head())
```

```
# Separate features and labels
# After separating the dataset, we now have numpy arrays named **X** containing the feat
X, y = bike_data[['season', 'mnth', 'holiday', 'weekday', 'workingday', 'weathersit', 'temp',

# 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.s
```



[1] ✓ 59 sec

```
...
      instant  dteday  season  yr  mnth  holiday  weekday  workingday  \
0           1  1/1/2011      1   0     1         0         6         0
1           2  1/2/2011      1   0     1         0         0         0
2           3  1/3/2011      1   0     1         0         1         1
3           4  1/4/2011      1   0     1         0         2         1
4           5  1/5/2011      1   0     1         0         3         1

      weathersit    temp    atemp    hum  windspeed  rentals  day
0           2  0.344167  0.363625  0.805833  0.160446    331    1
1           2  0.363478  0.353739  0.696087  0.248539    131    2
2           1  0.196364  0.189405  0.437273  0.248309    120    3
3           1  0.200000  0.212122  0.590435  0.160296    108    4
4           1  0.226957  0.229270  0.436957  0.186900     82    5
Training Set: 511 rows
Test Set: 220 rows
```

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.

Experiment with Algorithms

The linear-regression algorithm we used last time to train the model has some predictive capability, but there are many kinds of regression algorithm we could try, including:

- **Linear algorithms:** Not just the Linear Regression algorithm we used above (which is technically an *Ordinary Least Squares* algorithm), but other variants such as *Lasso* and *Ridge*.
- **Tree-based algorithms:** Algorithms that build a decision tree to reach a prediction.
- **Ensemble algorithms:** Algorithms that combine the outputs of multiple base algorithms to improve generalizability.

Note: For a full list of Scikit-Learn estimators that encapsulate algorithms for supervised machine learning, see the [Scikit-Learn documentation](#). There are many algorithms from which to choose, but for most real-world scenarios, the [Scikit-Learn estimator cheat sheet](#) can help you find a suitable starting point.

```

from sklearn.linear_model import Lasso

# Fit a lasso model on the training set
model = Lasso().fit(X_train, y_train)
print (model, "\n")

# Evaluate the model using the test data
predictions = model.predict(X_test)
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)

# Plot predicted vs actual
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()

```

[2] ✓ 1 sec

```

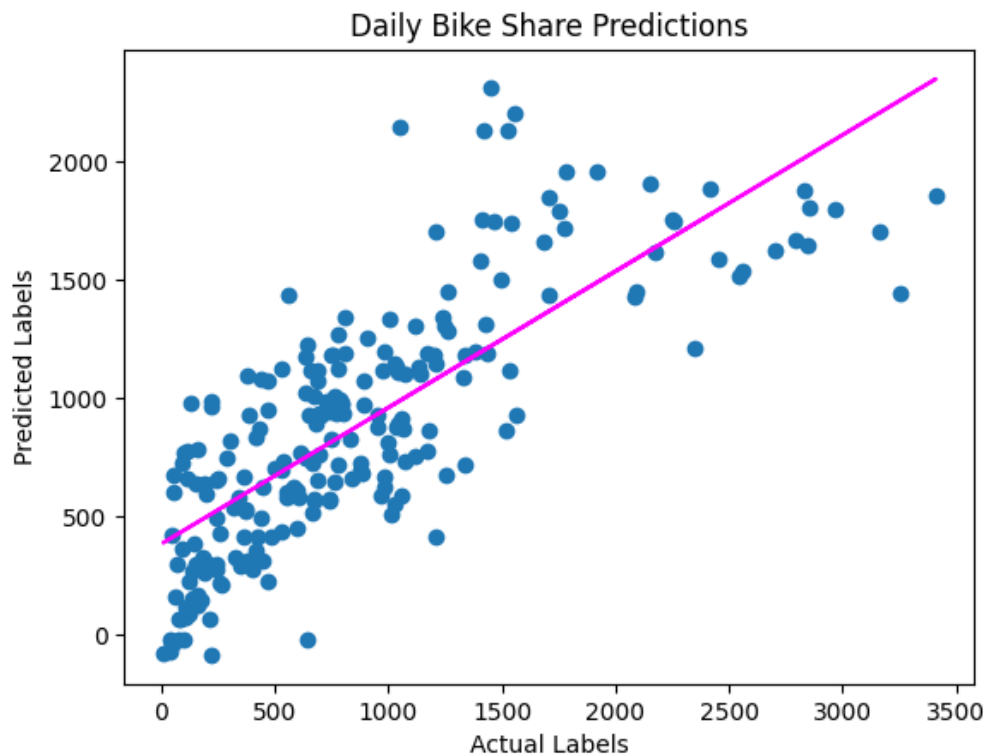
Lasso(alpha=1.0, copy_X=True, fit_intercept=True, max_iter=1000,
      normalize=False, positive=False, precompute=False, random_state=None,
      selection='cyclic', tol=0.0001, warm_start=False)

```

MSE: 201155.70593338404

RMSE: 448.5038527519959

R2: 0.6056468637824488



Try a Decision Tree Algorithm

As an alternative to a linear model, there's a category of algorithms for machine learning that uses a tree-based approach in which the features in the dataset are examined in a series of evaluations, each of which results in a *branch* in a *decision tree* based on the feature value. At the end of each series of branches are leaf-nodes with the predicted label value based on the feature values.

It's easiest to see how this works with an example. Let's train a Decision Tree regression model using the bike rental data. After training the model, the following code will print the model definition and a text representation of the tree. It also predicts label values.

```
from sklearn.tree import DecisionTreeRegressor
from sklearn.tree import export_text

# Train the model
model = DecisionTreeRegressor().fit(X_train, y_train)
print (model, "\n")

# Visualize the model tree
tree = export_text(model)
print(tree)
```

[3] ✓ <1 sec

```
DecisionTreeRegressor(ccp_alpha=0.0, criterion='mse', max_depth=None,
                      max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, presort='deprecated',
                      random_state=None, splitter='best')
```

```
|--- feature_6 <= 0.45
|   |--- feature_4 <= 0.50
|   |   |--- feature_7 <= 0.32
|   |   |   |--- feature_8 <= 0.41
|   |   |   |   |--- feature_1 <= 2.50
|   |   |   |   |   |--- feature_3 <= 3.00
|   |   |   |   |   |   |--- feature_6 <= 0.28
|   |   |   |   |   |   |   |--- value: [515.00]
|   |   |   |   |   |   |   |--- feature_6 > 0.28
```

So now we have a tree-based model, but is it any good? Let's evaluate it with the test data.

```
# Evaluate the model using the test data
predictions = model.predict(X_test)
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)

# Plot predicted vs actual
plt.scatter(y_test, predictions)
plt.xlabel('Actual Labels')
```

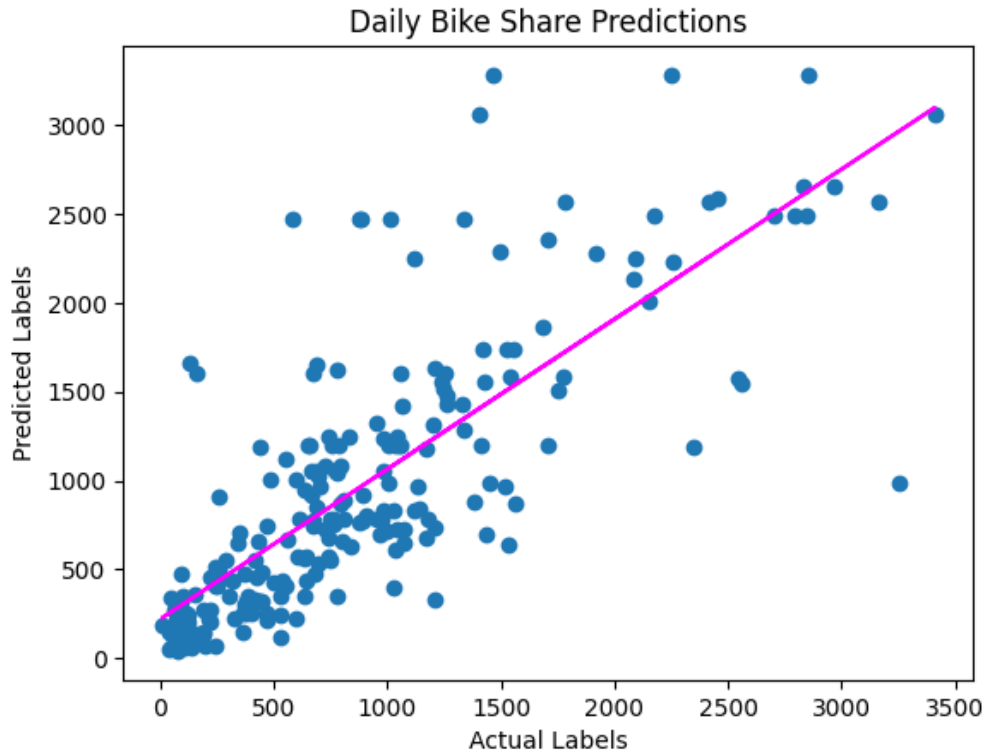
```
plt.ylabel('Predicted Labels')
plt.title('Daily Bike Share Predictions')
# overlay the regression line
z = np.polyfit(y_test, predictions, 1)
p = np.poly1d(z)
plt.plot(y_test, p(y_test), color='magenta')
plt.show()
```

[4] ✓ <1 sec

MSE: 241805.91818181818

RMSE: 491.737651783772

R2: 0.5259546740247898



The tree-based model doesn't seem to have improved over the linear model, so what else could we try?

Try an Ensemble Algorithm

Ensemble algorithms work by combining multiple base estimators to produce an optimal model, either by applying an aggregate function to a collection of base models (sometimes referred to a *bagging*) or by building a sequence of models that build on one another to improve predictive performance (referred to as *boosting*).

For example, let's try a Random Forest model. First, we'll create a random forest model:

```
from sklearn.ensemble import RandomForestRegressor

# Train the model
model = RandomForestRegressor().fit(X_train, y_train)
print (model, "\n")

# Evaluate the model using the test data
predictions = model.predict(X_test)
mse = mean_squared_error(y_test, predictions)
```

```

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)

# Plot predicted vs actual
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()

```

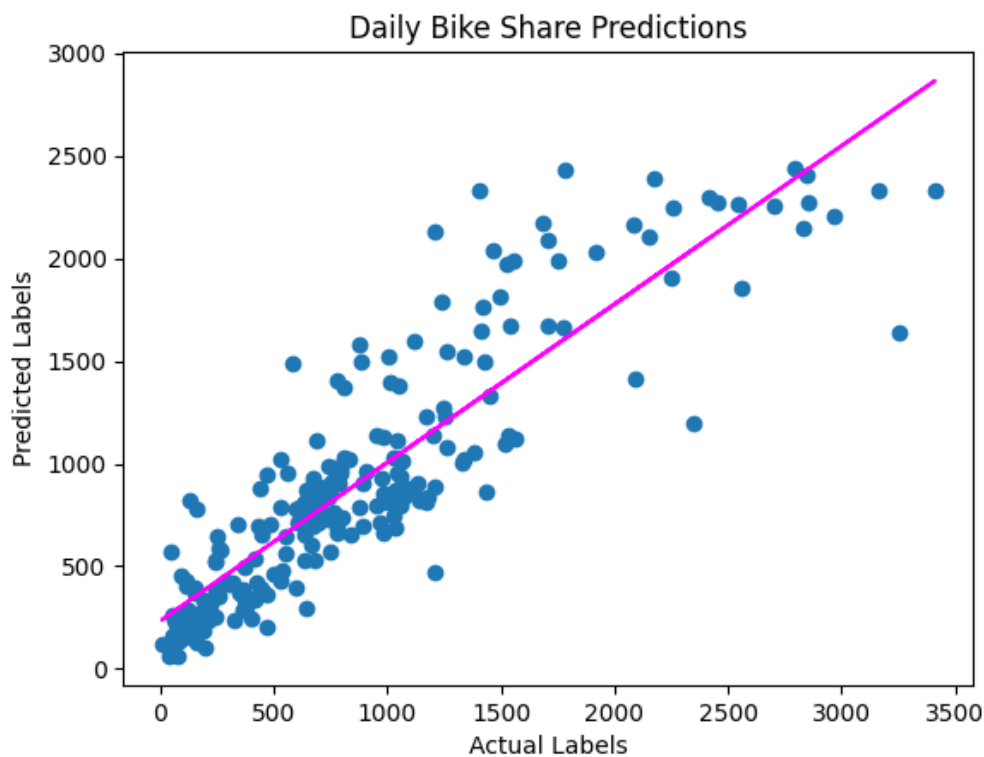
[5] ✓ 1 sec

```

RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                      max_depth=None, max_features='auto', max_leaf_nodes=None,
                      max_samples=None, min_impurity_decrease=0.0,
                      min_impurity_split=None, min_samples_leaf=1,
                      min_samples_split=2, min_weight_fraction_leaf=0.0,
                      n_estimators=100, n_jobs=None, oob_score=False,
                      random_state=None, verbose=0, warm_start=False)

```

MSE: 110297.43965181816
RMSE: 332.11058346854617
R2: 0.7837687922317003



For good measure, let's also try a *boosting* ensemble algorithm. We'll use a Gradient Boosting estimator, which like a Random Forest algorithm builds multiple trees; but instead of building

```
# Train the model
from sklearn.ensemble import GradientBoostingRegressor

# Fit a lasso model on the training set
model = GradientBoostingRegressor().fit(X_train, y_train)
print (model, "\n")

# Evaluate the model using the test data
predictions = model.predict(X_test)
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)

# Plot predicted vs actual
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()
```

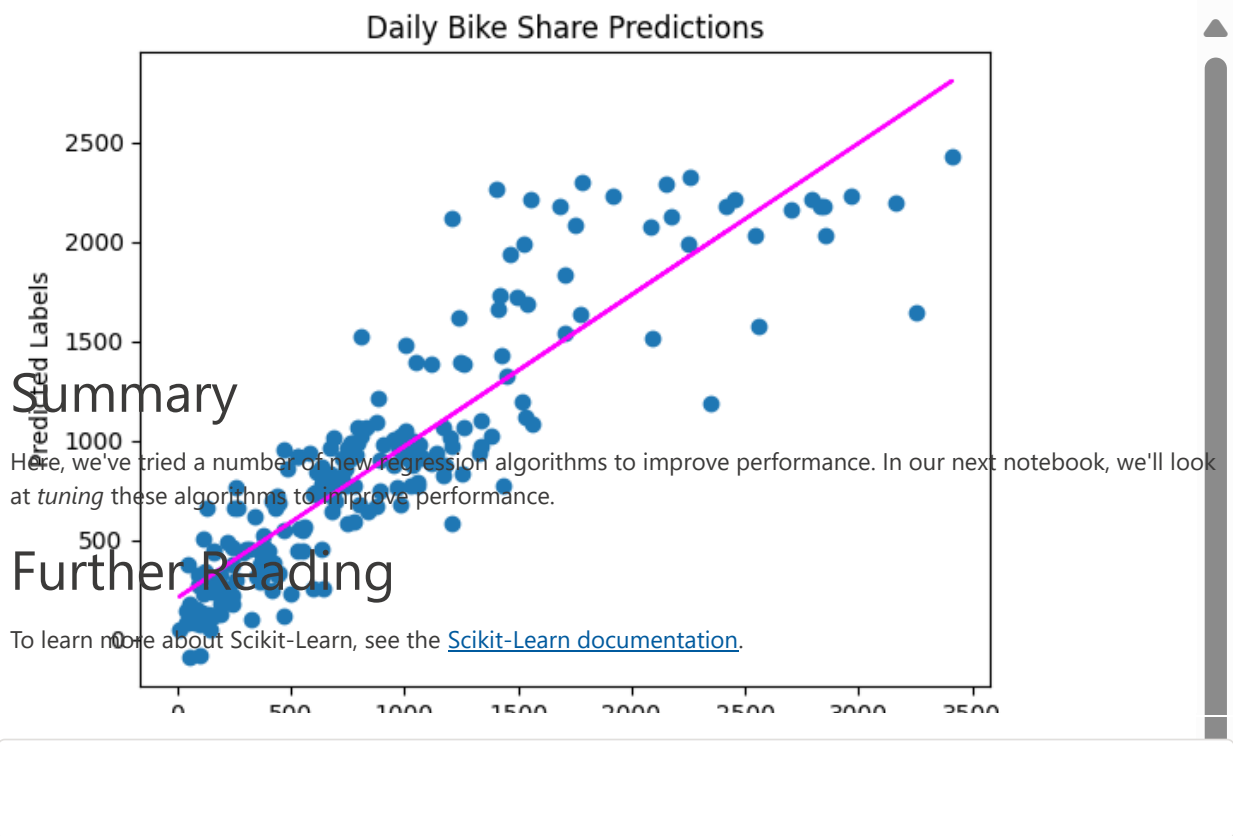
[6] ✓ <1 sec

```
... GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
                               init=None, learning_rate=0.1, loss='ls', max_depth=3,
                               max_features=None, max_leaf_nodes=None,
                               min_impurity_decrease=0.0, min_impurity_split=None,
                               min_samples_leaf=1, min_samples_split=2,
                               min_weight_fraction_leaf=0.0, n_estimators=100,
                               n_iter_no_change=None, presort='deprecated',
                               random_state=None, subsample=1.0, tol=0.0001,
                               validation_fraction=0.1, verbose=0, warm_start=False)
```

MSE: 103946.75217450749

RMSE: 322.4077421131625

R2: 0.7962189164386884



learn-notebooks-ffee7440-4204-42b1-bf43-da36245df269 Compute connected Editing Kernel idle azureml_py38

Next unit: Improve models with hyperparameters

Continue >