

XGBoost With Python

7-Day Mini-Course

Jason Brownlee

**MACHINE
LEARNING
MASTERY**



Jason Brownlee

XGBoost With Python

7 Day Mini-Course

XGBoost With Python

© Copyright 2016 Jason Brownlee. All Rights Reserved.

Edition: v1.0

Find the latest version of this guide online at: <http://MachineLearningMastery.com>

Contents

Before We Get Started...	1
Lesson 01: Introduction to Gradient Boosting	4
Lesson 02: Introduction to XGBoost	5
Lesson 03: Develop Your First XGBoost Model	6
Lesson 04: Monitor Performance and Early Stopping	8
Lesson 05: Feature Importance with XGBoost	10
Lesson 06: How to Configure Gradient Boosting	11
Lesson 07: XGBoost Hyperparameter Tuning	12
Final Word Before You Go...	14

Before We Get Started...

XGBoost is an implementation of gradient boosting that is being used to win machine learning competitions. It is powerful but it can be hard to get started. In this guide you will discover a 7-part crash course on XGBoost with Python. This mini course is designed for Python machine learning practitioners that are already comfortable with scikit-learn and the SciPy ecosystem. Let's get started.

Who Is This Mini-Course For?

Before we get started, let's make sure you are in the right place. The list below provides some general guidelines as to who this course was designed for. Don't panic if you don't match these points exactly, you might just need to brush up in one area or another to keep up.

- **Developers that know how to write a little code.** This means that it is not a big deal for you to get things done with Python and know how to setup the SciPy ecosystem on your workstation (a prerequisite). It does not mean you're a wizard coder, but it does mean you're not afraid to install packages and write scripts.
- **Developers that know a little machine learning.** This means you know about the basics of machine learning like cross validation, some algorithms and the bias-variance trade-off. It does not mean that you are a machine learning PhD, just that you know the landmarks or know where to look them up.

This mini-course is not a textbook on XGBoost. It will take you from a developer that knows a little machine learning in Python to a developer who can get results and bring the power of XGBoost to your own projects.

Mini-Course Overview (what to expect)

This mini-course is divided into 7 parts. Each lesson was designed to take the average developer about 30 minutes. You might finish some much sooner and others you may choose to go deeper and spend more time. You can complete each part as quickly or as slowly as you like. A comfortable schedule may be to complete one lesson per day over a one week period. Highly recommended. The topics you will cover over the next 7 lessons are as follows:

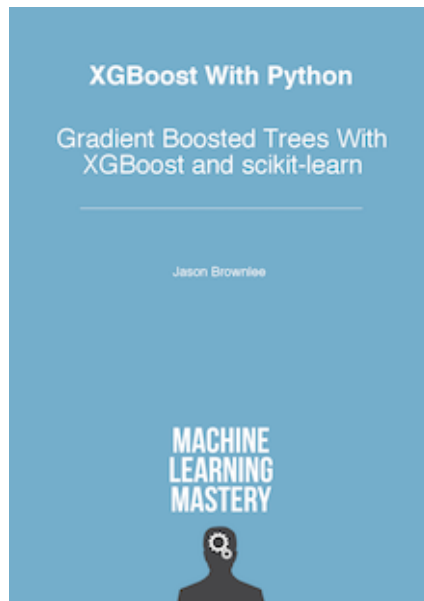
- **Lesson 01:** Introduction to Gradient Boosting.
- **Lesson 02:** Introduction to XGBoost.

- **Lesson 03:** Develop Your First XGBoost Model.
- **Lesson 04:** Monitor Performance and Early Stopping.
- **Lesson 05:** Feature Importance with XGBoost.
- **Lesson 06:** How to Configure Gradient Boosting.
- **Lesson 07:** XGBoost Hyperparameter Tuning.

This is going to be a lot of fun. You're going to have to do some work though, a little reading, a little research and a little programming. You want to learn XGBoost right? Here's a tip: All of the answers these lessons can be found on this blog <http://MachineLearningMastery.com>. Use the search feature.

Hang in there, don't give up!

If you would like me to step you through each lesson in great detail (and much more), take a look at my book: **XGBoost With Python**:



Learn more here:

<https://machinelearningmastery.com/xgboost-with-python>

Lesson 01: Introduction to Gradient Boosting

Gradient boosting is one of the most powerful techniques for building predictive models. The idea of boosting came out of the idea of whether a weak learner can be modified to become better. The first realization of boosting that saw great success in application was Adaptive Boosting or AdaBoost for short. The weak learners in AdaBoost are decision trees with a single split, called decision stumps for their shortness.

AdaBoost and related algorithms were recast in a statistical framework and became known as Gradient Boosting Machines. The statistical framework cast boosting as a numerical optimization problem where the objective is to minimize the loss of the model by adding weak learners using a gradient descent like procedure, hence the name. The Gradient Boosting algorithm involves three elements:

- **A loss function to be optimized**, such as cross entropy for classification or mean squared error for regression problems.
- **A weak learner to make predictions**, such as a greedily constructed decision tree.
- **An additive model**, used to add weak learners to minimize the loss function.

New weak learners are added to the model in an effort to correct the residual errors of all previous trees. The result is a powerful predictive modeling algorithm, perhaps more powerful than random forest. In the next lesson we will take a closer look at the XGBoost implementation of gradient boosting.

Lesson 02: Introduction to XGBoost

XGBoost is an implementation of gradient boosted decision trees designed for speed and performance. XGBoost stands for eXtreme Gradient Boosting. It was developed by Tianqi Chen and is laser focused on computational speed and model performance, as such there are few frills. In addition to supporting all key variations of the technique, the real interest is the speed provided by the careful engineering of the implementation, including:

- **Parallelization** of tree construction using all of your CPU cores during training.
- **Distributed Computing** for training very large models using a cluster of machines.
- **Out-of-Core Computing** for very large datasets that don't fit into memory.
- **Cache Optimization** of data structures and algorithms to make best use of hardware.

Traditionally, gradient boosting implementations are slow because of the sequential nature in which each tree must be constructed and added to the model. The on performance in the development of XGBoost has resulted in one of the best predictive modeling algorithms that can now harness the full capability of your hardware platform, or very large computers you might rent in the cloud. As such, XGBoost has been a cornerstone in competitive machine learning, being the technique used to win and recommended by winners. For example, here is what some recent Kaggle competition winners have said:

As the winner of an increasing amount of Kaggle competitions, XGBoost showed us again to be a great all-round algorithm worth having in your toolbox.

– Dato Winners' Interview¹

When in doubt, use xgboost.

– Avito Winner's Interview²

In the next lesson, we will develop our first XGBoost model in Python.

¹<http://goo.gl/AHkmWx>

²<http://goo.gl/sGyGtu>

Lesson 03: Develop Your First XGBoost Model

Assuming you have a working SciPy environment, XGBoost can be installed easily using pip. For example:

```
1 sudo pip install xgboost
```

Listing 1: Install XGBoost using pip.

You can learn more about installing and building XGBoost on your platform in the XGBoost Installation Instructions³. XGBoost models can be used directly in the scikit-learn framework using the wrapper classes, `XGBClassifier` for classification and `XGBRegressor` for regression problems. This is the recommended way to use XGBoost in Python. Download the Pima Indians onset of diabetes dataset from the UCI Machine Learning repository⁴. It is a good test dataset for binary classification as all input variables are numeric, meaning the problem can be modeled directly with no data preparation. We can train an XGBoost model for classification by constructing it and calling the `model.fit()` function:

```
1 model = XGBClassifier()
2 model.fit(X_train, y_train)
```

Listing 2: Train an XGBoost Model.

This model can then be used to make predictions by calling the `model.predict()` function on new data.

```
1 y_pred = model.predict(X_test)
```

Listing 3: Make Predictions with an XGBoost Model.

We can tie this all together as follows:

```
1 # First XGBoost model for Pima Indians dataset
2 from numpy import loadtxt
3 from xgboost import XGBClassifier
4 from sklearn.cross_validation import train_test_split
5 from sklearn.metrics import accuracy_score
6 # load data
7 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=",")
8 # split data into X and y
9 X = dataset[:,0:8]
10 Y = dataset[:,8]
```

³<http://xgboost.readthedocs.io/en/latest/build.html>

⁴<https://archive.ics.uci.edu/ml/datasets/Pima+Indians+Diabetes>

```
11 # split data into train and test sets
12 seed = 7
13 test_size = 0.33
14 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=test_size,
15     random_state=seed)
16 # fit model on training data
17 model = XGBClassifier()
18 model.fit(X_train, y_train)
19 # make predictions for test data
20 y_pred = model.predict(X_test)
21 predictions = [round(value) for value in y_pred]
22 # evaluate predictions
23 accuracy = accuracy_score(y_test, predictions)
24 print("Accuracy: %.2f%%" % (accuracy * 100.0))
```

Listing 4: First XGBoost Model.

In the next lesson we will look at how we can use early stopping to limit overfitting.

Lesson 04: Monitor Performance and Early Stopping

The XGBoost model can evaluate and report on the performance on a test set for the model during training. It supports this capability by specifying both an test dataset and an evaluation metric on the call to `model.fit()` when training the model and specifying verbose output (`verbose=True`). For example, we can report on the binary classification error rate (`error`) on a standalone test set (`eval_set`) while training an XGBoost model as follows:

```
1 eval_set = [(X_test, y_test)]
2 model.fit(X_train, y_train, eval_metric="error", eval_set=eval_set, verbose=True)
```

Listing 5: Monitor Performance During Training.

We can use this evaluation to stop training once no further improvements have been made to the model. We can do this by setting the `early_stopping_rounds` parameter when calling `model.fit()` to the number of iterations that no improvement is seen on the validation dataset before training is stopped. The full example using the Pima Indians Onset of Diabetes dataset is provided below.

```
1 # exmaple of early stopping
2 from numpy import loadtxt
3 from xgboost import XGBClassifier
4 from sklearn.cross_validation import train_test_split
5 from sklearn.metrics import accuracy_score
6 # load data
7 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=",")
8 # split data into X and y
9 X = dataset[:,0:8]
10 Y = dataset[:,8]
11 # split data into train and test sets
12 seed = 7
13 test_size = 0.33
14 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=test_size,
15     random_state=seed)
16 # fit model no training data
17 model = XGBClassifier()
18 eval_set = [(X_test, y_test)]
19 model.fit(X_train, y_train, early_stopping_rounds=10, eval_metric="logloss",
20     eval_set=eval_set, verbose=True)
21 # make predictions for test data
22 y_pred = model.predict(X_test)
23 predictions = [round(value) for value in y_pred]
24 # evaluate predictions
25 accuracy = accuracy_score(y_test, predictions)
```

```
24 | print("Accuracy: %.2f%%" % (accuracy * 100.0))
```

Listing 6: Example of Early Stopping.

In the next lesson we will look at how we calculate the importance of features using XGBoost.

Lesson 05: Feature Importance with XGBoost

A benefit of using ensembles of decision tree methods like gradient boosting is that they can automatically provide estimates of feature importance from a trained predictive model. A trained XGBoost model automatically calculates feature importance on your predictive modeling problem. These importance scores are available in the `feature_importances_` member variable of the trained model. For example, they can be printed directly as follows:

```
1 print(model.feature_importances_)
```

Listing 7: Print Feature Importance.

The XGBoost library provides a built-in function to plot features ordered by their importance. The function is called `plot_importance()` and can be used as follows:

```
1 plot_importance(model)
2 pyplot.show()
```

Listing 8: Plot Feature Importance.

These importance scores can help you decide what input variables to keep or discard. They can also be used as the basis for automatic feature selection techniques. The full example of plotting feature importance scores using the Pima Indians Onset of Diabetes dataset is provided below.

```
1 # plot feature importance using built-in function
2 from numpy import loadtxt
3 from xgboost import XGBClassifier
4 from xgboost import plot_importance
5 from matplotlib import pyplot
6 # load data
7 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=",")
8 # split data into X and y
9 X = dataset[:,0:8]
10 y = dataset[:,8]
11 # fit model no training data
12 model = XGBClassifier()
13 model.fit(X, y)
14 # plot feature importance
15 plot_importance(model)
16 pyplot.show()
```

Listing 9: Example of Plotting Feature Importance.

In the next lesson we will look at heuristics for best configuring the gradient boosting algorithm.

Lesson 06: How to Configure Gradient Boosting

Gradient boosting is one of the most powerful techniques for applied machine learning and as such is quickly becoming one of the most popular. But how do you configure gradient boosting on your problem? A number of configuration heuristics were published in the original gradient boosting papers. They can be summarized as:

- Learning rate or shrinkage (`learning_rate` in XGBoost) should be set to 0.1 or lower, and smaller values will require the addition of more trees.
- The depth of trees (`tree_depth` in XGBoost) should be configured in the range of 2-to-8, where not much benefit is seen with deeper trees.
- Row sampling (`subsample` in XGBoost) should be configured in the range of 30% to 80% of the training dataset, and compared to a value of 100% for no sampling.

These are a good starting points when configuring your model. A good general configuration strategy is as follows:

1. Run the default configuration and review plots of the learning curves on the training and validation datasets.
2. If the system is overlearning, decrease the learning rate and/or increase the number of trees.
3. If the system is underlearning, speed the learning up to be more aggressive by increasing the learning rate and/or decreasing the number of trees.

Owen Zhang, the former #1 ranked competitor on Kaggle and now CTO at Data Robot proposes an interesting strategy to configure XGBoost⁵. He suggests to set the number of trees to a target value such as 100 or 1000, then tune the learning rate to find the best model. This is an efficient strategy for quickly finding a good model. In the next and final lesson, we will look at an example of tuning the XGBoost hyperparameters.

⁵<http://goo.gl/0qIRIc>

Lesson 07: XGBoost Hyperparameter Tuning

The scikit-learn framework provides the capability to search combinations of parameters. This capability is provided in the GridSearchCV class and can be used to discover the best way to configure the model for top performance on your problem. For example, we can define a grid of the number of trees (`n_estimators`) and tree sizes (`max_depth`) to evaluate by defining a grid as:

```
1 n_estimators = [50, 100, 150, 200]
2 max_depth = [2, 4, 6, 8]
3 param_grid = dict(max_depth=max_depth, n_estimators=n_estimators)
```

Listing 10: Define a Grid Search.

And then evaluate each combination of parameters using 10-fold cross validation as:

```
1 kfold = StratifiedKFold(label_encoded_y, n_folds=10, shuffle=True, random_state=7)
2 grid_search = GridSearchCV(model, param_grid, scoring="log_loss", n_jobs=-1, cv=kfold,
3 verbose=1)
3 result = grid_search.fit(X, label_encoded_y)
```

Listing 11: Perform a Grid Search.

We can then review the results to determine the best combination and the general trends in varying the combinations of parameters. This is the best practice when applying XGBoost to your own problems. The parameters to consider tuning are:

- The number and size of trees (`n_estimators` and `max_depth`).
- The learning rate and number of trees (`learning_rate` and `n_estimators`).
- The row and column subsampling rates (`subsample`, `colsample_bytree` and `colsample_bylevel`).

Below is a full example of tuning just the `learning_rate` on the Pima Indians Onset of Diabetes dataset.

```
1 # Tune learning_rate
2 from numpy import loadtxt
3 from xgboost import XGBClassifier
4 from sklearn.grid_search import GridSearchCV
5 from sklearn.cross_validation import StratifiedKFold
6 # load data
7 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=",")
8 # split data into X and y
```



```

9 X = dataset[:,0:8]
10 Y = dataset[:,8]
11 # grid search
12 model = XGBClassifier()
13 learning_rate = [0.0001, 0.001, 0.01, 0.1, 0.2, 0.3]
14 param_grid = dict(learning_rate=learning_rate)
15 kfold = StratifiedKFold(Y, n_folds=10, shuffle=True, random_state=7)
16 grid_search = GridSearchCV(model, param_grid, scoring="log_loss", n_jobs=-1, cv=kfold)
17 result = grid_search.fit(X, Y)
18 # summarize results
19 print("Best: %f using %s" % (result.best_score_, result.best_params_))
20 means, stdevs = [], []
21 for params, mean_score, scores in result.grid_scores_:
22     stdev = scores.std()
23     means.append(mean_score)
24     stdevs.append(stdev)
25 print("%f (%f) with: %r" % (mean_score, stdev, params))

```

Listing 12: Example of a Grid Search.

Final Word Before You Go...

You made it. Well done! Take a moment and look back at how far you have come:

- You discovered deep learning libraries in Python including the powerful numerical libraries Theano and TensorFlow and the easy to use Keras library for applied deep learning.
- You built your first neural network using Keras and learned how to use your deep learning models with scikit-learn and how to retrieve and plot the training history for your models.
- You learned about more advanced techniques such as dropout regularization and learning rate schedules and how you can use these techniques in Keras.
- Finally, you took the next step and learned about and developed convolutional neural networks for complex computer vision tasks and learned about augmentation of image data.

Don't make light of this, you have come a long way in a short amount of time. This is just the beginning of your machine learning journey with Python. Keep practicing and developing your skills.

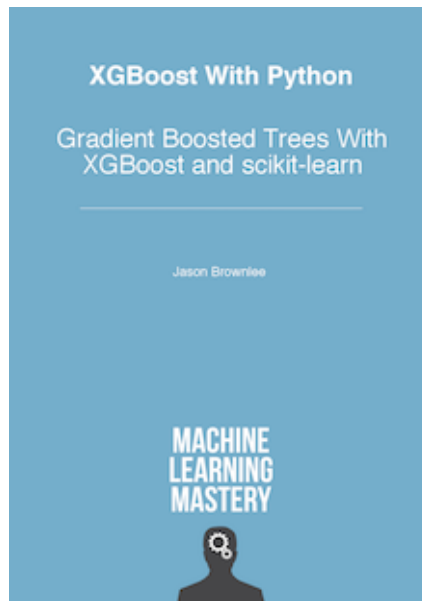
How Did You Go With The Mini-Course?

Did you enjoy this mini-course?

Do you have any questions or sticking points?

Let me know, send me an email at: **jason@MachineLearningMastery.com**

If you would like me to step you through each lesson in great detail (and much more), take a look at my book: **XGBoost With Python**:



Learn more here:

<https://machinelearningmastery.com/xgboost-with-python>