☰ | **Navigation**

**Machine Learning Mastery**
Making Developers Awesome at Machine Learning

Click to Take the FREE XGBoost Crash-Course

Search...                                                    🔍

# A Gentle Introduction to XGBoost for Applied Machine Learning

by **Jason Brownlee** on August 17, 2016 in **XGBoost**

Tweet          Tweet          Share          **Share**

Last Updated on February 17, 2021

XGBoost is an algorithm that has recently been dominating applied machine learning and Kaggle competitions for structured or tabular data.

XGBoost is an implementation of gradient boosted decision trees designed for speed and performance.

In this post you will discover XGBoost and get a gentle introduction to what is, where it came from and how you can learn more.
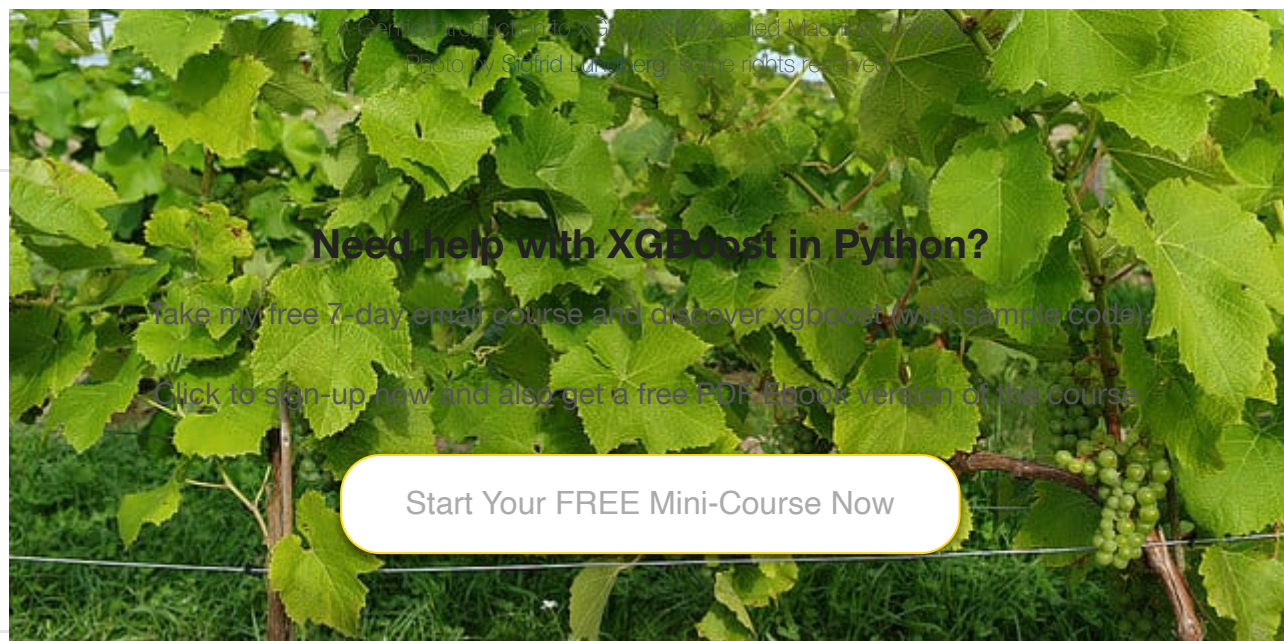
After reading this post you will know:

- What XGBoost is and the goals of the project.
- Why XGBoost must be a part of your machine learning toolkit.
- Where you can learn more to start using XGBoost on your next machine learning project.

**Kick-start your project** with my new book XGBoost With Python, including *step-by-step tutorials* and the *Python source code* files for all examples.

Let's get started.

- **Updated Feb/2021**: Fixed broken links.

**Need help with XGBoost in Python?**

Take my free 7-day email course and discover xgboost (with sample code).

Click to sign-up now and also get a free PDF Ebook version of the course.

Start Your FREE Mini-Course Now

## What is XGBoost?

XGBoost stands for e**X**treme **G**radient **B**oosting.

> *The name xgboost, though, actually refers to the engineering goal to push the limit of computations resources for boosted tree algorithms. Which is the reason why many people use xgboost.*

— Tianqi Chen, in answer to the question "What is the difference between the R gbm (gradient boosting machine) and xgboost (extreme gradient boosting)?" on Quora

It is an implementation of gradient boosting machines created by Tianqi Chen, now with contributions from many developers. It belongs to a broader collection of tools under the umbrella of the Distributed Machine Learning Community or DMLC who are also the creators of the popular mxnet deep learning library.

Tianqi Chen provides a brief and interesting back story on the creation of XGBoost in the post Story and Lessons Behind the Evolution of XGBoost.

XGBoost is a software library that you can download and install on your machine, then access from a variety of interfaces. Specifically, XGBoost supports the following main interfaces:

- Command Line Interface (CLI).
- C++ (the language in which the library is written).
- Python interface as well as a model in scikit-learn.
- R interface as well as a model in the caret package.
- Julia.
- Java and JVM languages like Scala and platforms like Hadoop.

# XGBoost Features

The library is laser focused on computational speed and model performance, as such there are few frills. Nevertheless, it does offer a number of advanced features.

## Model Features

The implementation of the model supports the features of the scikit-learn and R implementations, with new additions like regularization. Three main forms of gradient boosting are supported:

- **Gradient Boosting** algorithm also called gradient boosting machine including the learning rate.
- **Stochastic Gradient Boosting** with sub-sampling at the row, column and column per split levels.
- **Regularized Gradient Boosting** with both L1 and L2 regularization.

## System Features

The library provides a system for use in a range of computing environments, not least:

- **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 algorithm to make best use of hardware.

## Algorithm Features

The implementation of the algorithm was engineered for efficiency of compute time and memory resources. A design goal was to make the best use of available resources to train the model. Some key algorithm implementation features include:

- **Sparse Aware** implementation with automatic handling of missing data values.
- **Block Structure** to support the parallelization of tree construction.
- **Continued Training** so that you can further boost an already fitted model on new data.

XGBoost is free open source software available for use under the permissive Apache-2 license.

# Why Use XGBoost?

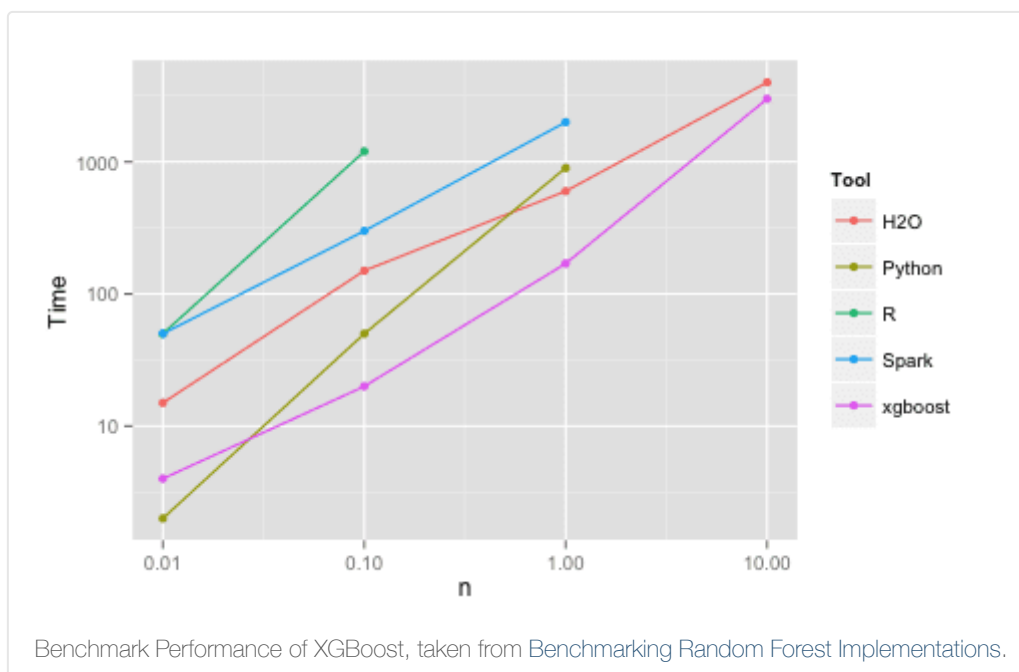The two reasons to use XGBoost are also the two goals of the project:

1. Execution Speed.
2. Model Performance.

## 1. XGBoost Execution Speed

Generally, XGBoost is fast. Really fast when compared to other implementations of gradient boosting.

Szilard Pafka performed some objective benchmarks comparing the performance of XGBoost to other implementations of gradient boosting and bagged decision trees. He wrote up his results in May 2015 in the blog post titled "Benchmarking Random Forest Implementations".

He also provides all the code on GitHub and a more extensive report of results with hard numbers.



Benchmark Performance of XGBoost, taken from Benchmarking Random Forest Implementations.

His results showed that XGBoost was almost always faster than the other benchmarked implementations from R, Python Spark and H2O.

From his experiment, he commented:

> *I also tried xgboost, a popular library for boosting which is capable to build random forests as well. It is fast, memory efficient and of high accuracy*

— Szilard Pafka, Benchmarking Random Forest Implementations.

## 2. XGBoost Model Performance

XGBoost dominates structured or tabular datasets on classification and regression predictive modeling problems.

The evidence is that it is the go-to algorithm for competition winners on the Kaggle competitive data science platform.

For example, there is an incomplete list of first, second and third place competition winners that used titled: XGBoost: Machine Learning Challenge Winning Solutions.

To make this point more tangible, below are some insightful quotes from Kaggle competition winners:

> *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' Interview: 1st place, Mad Professors

> *When in doubt, use xgboost.*

— Avito Winner's Interview: 1st place, Owen Zhang

> *I love single models that do well, and my best single model was an XGBoost that could get the 10th place by itself.*

— Caterpillar Winners' Interview: 1st place

> *I only used XGBoost.*

— Liberty Mutual Property Inspection, Winner's Interview: 1st place, Qingchen Wang

> *The only supervised learning method I used was gradient boosting, as implemented in the excellent xgboost package.*

— Recruit Coupon Purchase Winner's Interview: 2nd place, Halla Yang

## What Algorithm Does XGBoost Use?

The XGBoost library implements the gradient boosting decision tree algorithm.

This algorithm goes by lots of different names such as gradient boosting, multiple additive regression trees, stochastic gradient boosting or gradient boosting machines.

Boosting is an ensemble technique where new models are added to correct the errors made by existing models. Models are added sequentially until no further improvements can be made. A popular example is the AdaBoost algorithm that weights data points that are hard to predict.

Gradient boosting is an approach where new models are created that predict the residuals or errors of prior models and then added together to make the final prediction. It is called gradient boosting because it uses a gradient descent algorithm to minimize the loss when adding new models.

This approach supports both regression and classification predictive modeling problems.

For more on boosting and gradient boosting, see Trevor Hastie's talk on Gradient Boosting Machine Learning.

# Official XGBoost Resources

The best source of information on XGBoost is the official GitHub repository for the project.

From there you can get access to the Issue Tracker and the User Group that can be used for asking questions and reporting bugs.

A great source of links with example code and help is the Awesome XGBoost page.

There is also an official documentation page that includes a getting started guide for a range of different languages, tutorials, how-to guides and more.

There are some more formal papers on XGBoost that are worth a read for more background on the library:

- Higgs Boson Discovery with Boosted Trees, 2014.
- XGBoost: A Scalable Tree Boosting System, 2016.

# Talks on XGBoost

When getting started with a new tool like XGBoost, it can be helpful to review a few talks on the topic before diving into the code.

### XGBoost: A Scalable Tree Boosting System

Tianqi Chen, the creator of the library gave a talk to the LA Data Science group in June 2016 titled "XGBoost: A Scalable Tree Boosting System".

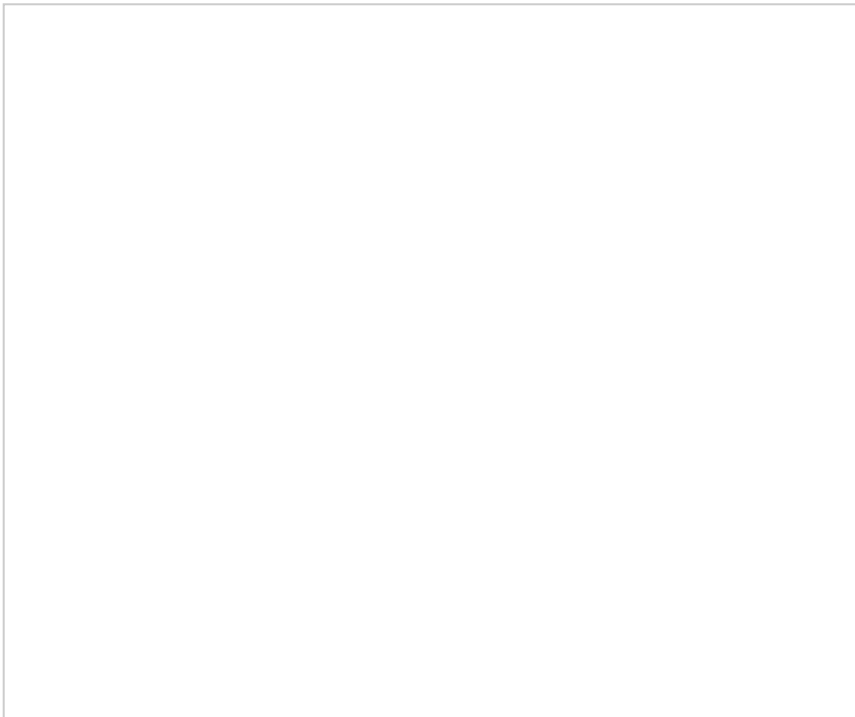You can review the slides from his talk here:

There is more information on the DataScience LA blog.

## XGBoost: eXtreme Gradient Boosting

Tong He, a contributor to XGBoost for the R interface gave a talk at the NYC Data Science Academy in December 2015 titled "XGBoost: eXtreme Gradient Boosting".

You can review the slides from his talk here:

**Xgboost** from **Vivian S. Zhang**

There is more information about this talk on the NYC Data Science Academy blog.

# Installing XGBoost

There is a comprehensive installation guide on the XGBoost documentation website.

It covers installation for Linux, Mac OS X and Windows.

It also covers installation on platforms such as R and Python.

# XGBoost in R

If you are an R user, the best place to get started is the CRAN page for the xgboost package.

From this page you can access the R vignette Package 'xgboost' [pdf].

There are also some excellent R tutorials linked from this page to get you started:

- Discover Your Data
- XGBoost Presentation
- xgboost: eXtreme Gradient Boosting [pdf]

There is also the official XGBoost R Tutorial and Understand your dataset with XGBoost tutorial.

# XGBoost in Python

Installation instructions are available on the Python section of the XGBoost installation guide.

The official Python Package Introduction is the best place to start when working with XGBoost in Python.

To get started quickly, you can type:

```
1  sudo pip install xgboost
```

There is also an excellent list of sample source code in Python on the XGBoost Python Feature Walkthrough.

# Summary

In this post you discovered the XGBoost algorithm for applied machine learning.
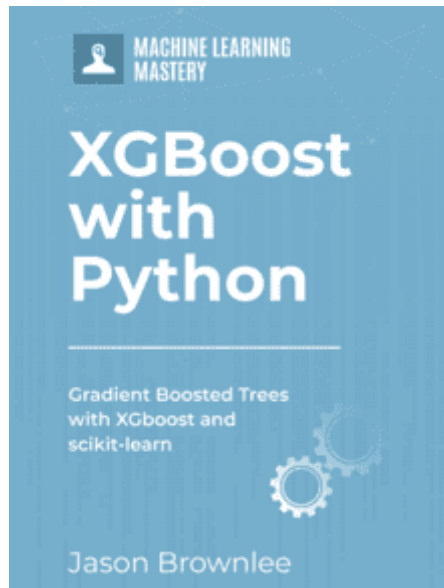
You learned:

- That XGBoost is a library for developing fast and high performance gradient boosting tree models.
- That XGBoost is achieving the best performance on a range of difficult machine learning tasks.
- That you can use this library from the command line, Python and R and how to get started.

Have you used XGBoost? Share your experiences in the comments below.

Do you have any questions about XGBoost or about this post? Ask your question in the comments below and I will do my best to answer them.

# Discover The Algorithm Winning Competitions!

### Develop Your Own XGBoost Models in Minutes

...with just a few lines of Python

Discover how in my new Ebook:
XGBoost With Python

It covers **self-study tutorials** like:
*Algorithm Fundamentals*, *Scaling*, *Hyperparameters*, and much more...

### Bring The Power of XGBoost To Your Own Projects
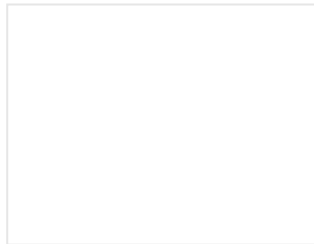
Skip the Academics. Just Results.

SEE WHAT'S INSIDE

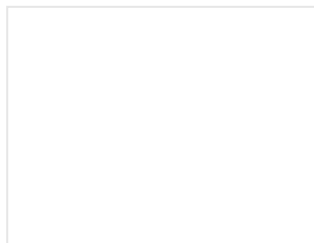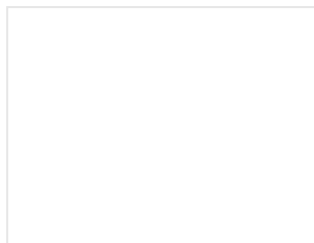Tweet          Tweet          Share          Share
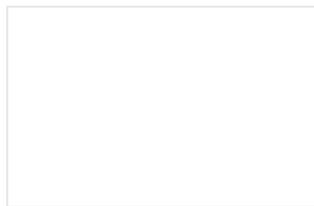
## More On This Topic



Extreme Gradient Boosting (XGBoost) Ensemble in Python
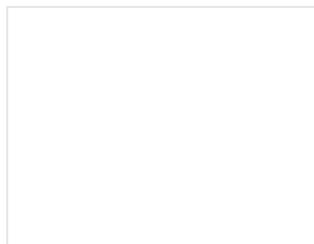


How to Develop Random Forest Ensembles With XGBoost



Tune XGBoost Performance With Learning Curves

A Gentle Introduction to XGBoost Loss Functions

XGBoost for Regression

How to Configure XGBoost for Imbalanced Classification

**About Jason Brownlee**

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on tutorials.

View all posts by Jason Brownlee →

‹ Weka Machine Learning Mini-Course          How to Develop Your First XGBoost Model in Python ›

## 67 Responses to *A Gentle Introduction to XGBoost for Applied Machine Learning*

**Seo Young Jae** July 10, 2017 at 6:25 pm #          REPLY ↩

Good information, thank you. Just one question.

Biggest difference from the gbm is normalization?

Does gbm not normalize, but does xgboost automatically normalize variables and automatically handle missing values? Did I get it right?

**Jason Brownlee** July 11, 2017 at 10:28 am #

REPLY ↩

The biggest difference is performance, not normalization.

**Seo Young Jae** July 10, 2017 at 7:46 pm #

REPLY ↩

I ran xgboost on R.

However, I found that input values can not be performed in the form of factors.

In case of gbm, it is possible to use factor type variable.

In that respect, xgboost seems to have some disadvantages.

**Jason Brownlee** July 11, 2017 at 10:29 am #

REPLY ↩

You must transform your categorical variables to be integer encoded or one hot encoded.

**Seo Young Jae** July 11, 2017 at 2:32 pm #

REPLY ↩

Is it ok to force a categorical variable to be a continuous variable?

**Jason Brownlee** July 12, 2017 at 9:39 am #

REPLY ↩

It depends on the variable. It might make sense if the variable is ordinal. If not, a one hot encoding would be the preferred approach.

**sapan** September 18, 2018 at 3:44 am #

REPLY ↩

this seems to be a limitation of the xgboost implementation you're using, not of the algorithm itself.

**dksahu** September 6, 2017 at 5:53 pm #

REPLY ↩

reference for monotonocity constraint for decision trees in xgboost?

**Jason Brownlee** September 7, 2017 at 12:50 pm #                    REPLY ↰

Sorry, I do not.

**Aman Garg** September 13, 2017 at 5:44 am #                    REPLY ↰

Could you please tell that if XGBoost can also be used for unsupervised learning – clustering of large datasets?

If yes, does XGBoost provides an edge over other unsupervised algorithms – like K means clustering, DBSCAN etc. ?

**Jason Brownlee** September 13, 2017 at 12:37 pm #                    REPLY ↰

Not as far as I know. Gradient boosting is a supervised learning algorithm.

**Petros Koulouris** September 20, 2017 at 5:14 pm #                    REPLY ↰

Jason, I would love to see how to perform repeated cross validation in order to hyper-tune model parameters. I used the caret package and it took 20-30 times longer than to train other models types ie ranger, gbm, glmnet on the same German credit dataset.

Its been touted as extremely fast which I haven't observed and most tutorials I have found employ caret.

**Jason Brownlee** September 21, 2017 at 5:37 am #                    REPLY ↰

Thanks for the suggestion Petros.

**Sasikanth** September 23, 2017 at 10:08 am #                    REPLY ↰

Hello Jason,
Have you tried to install and use LightGBM from Microsoft. It is said to be better and faster than XGboost.

REPLY ↰

**Jason Brownlee** September 24, 2017 at 5:11 am #

I have not, perhaps in the future.

REPLY ↰

**Norbert** November 10, 2017 at 6:57 pm #

Yep, XGboost rocks. One year ago I have created a quick free online course how to use it efficiently in Python – http://education.parrotprediction.teachable.com/p/practical-xgboost-in-python

REPLY ↰

**Jason Brownlee** November 11, 2017 at 9:20 am #

Cool, thanks for the ref Norbert. That is also about the time I released my book on the topic.

REPLY ↰

**Frank Ludeña** January 18, 2018 at 8:44 am #

Hello good afternoon, with respect to the fact that xgboost does not support categorical variables, I trained the following model in caret with a factor variable with xgbtree and I had no problem, (a single variable to exemplify). I am doing something wrong?

pase0.xgbTree_x=train(as.factor(PASE)~TIPO_CLIENTE,data=pase0,trControl=trainControl(method='repeat edcv',number=5,repeats=10,verboseIter = TRUE),method='xgbTree',allowParallel=TRUE,tuneGrid=xgb.tuning)

Tha categorical variable is TIPO_CLIENTE

REPLY ↰

**Jason Brownlee** January 18, 2018 at 10:15 am #

Sorry, I cannot help you with xgboost in R.

REPLY ↰

**Abhilash Menon** April 5, 2018 at 2:02 am #

Hi Dr. Brownlee, Is there a way to get all the predictions we make into the test dataset with the predictions of our model as a column in the test dataset. I am concerned about how the order of instances will be preserved in this case (as in, the prediction corresponding to an instance should be in the same row as the instance). Could you please shed some light on this issue?

Thanks!

**Jason Brownlee** April 5, 2018 at 6:14 am #

REPLY ↩

The order of the inputs will match the order of the outputs.

**Brett** April 8, 2018 at 5:00 am #

REPLY ↩

Not sure if this is the place for it, feel free to delete if not… but I just wanted to drop you a note to say thank you for the site… whenever it pops up in a search (which is often) I know I'm going to get some quality info.

**Jason Brownlee** April 8, 2018 at 6:30 am #

REPLY ↩

Thanks Brett, I really appreciate the kind words!

**IanDz** April 17, 2018 at 7:14 pm #

REPLY ↩

Jason, just wanted to thank you for all the amazing stuff you do! Your articles are some of the best online!

**Jason Brownlee** April 18, 2018 at 8:02 am #

REPLY ↩

Thanks IanDz, I really appreciate your support!

**GopalKrishna** June 3, 2018 at 12:47 pm #

REPLY ↩

Jason,

Xgb Importance output includes Split, RealCover and RealCover% in addition to Gain, Cover and Frequency when you pass add. parameters – training set ( or its subset) and label.

While Split value is understood, could you help understand/ interpret RealCover and RealCover% that appear against specific features only.

Also, in such expanded output what meaning should be derived from number of entries in the xgb importance table?

Thanks

**Jason Brownlee** June 4, 2018 at 6:22 am #                                                                                                    REPLY ↩

This documentation better explains the table:

http://xgboost.readthedocs.io/en/latest/R-package/discoverYourData.html#feature-importance

**Purvi Prajapati** December 4, 2018 at 10:09 pm #                                                                                        REPLY ↩

Could we apply XGBoost for Multi-Label Classification Problem?
Kindly reply me. I am working on Tree based approach for Multi-label classification.

**Jason Brownlee** December 5, 2018 at 6:16 am #                                                                                        REPLY ↩

Perhaps. Sorry, I don't have any examples of multi-label prediction. I hope to cover it in the
future.

**Purvi Prajapati** December 5, 2018 at 3:51 pm #                                                                                        REPLY ↩

Let me know is it applicable to Multi-Label Classification or not.

**Jason Brownlee** December 6, 2018 at 5:50 am #                                                                                        REPLY ↩

Maybe, I don't know.

**Apoorv** January 23, 2019 at 2:50 am #                                                                                                 REPLY ↩

I think you can by setting the objective function to any of the the below as per your
requirements (from xgboost documentation: https://xgboost.readthedocs.io/en/latest/parameter.html):

multi:softmax: set XGBoost to do multiclass classification using the softmax objective, you also need to
set num_class(number of classes)

multi:softprob: same as softmax, but output a vector of ndata * nclass, which can be further reshaped
to ndata * nclass matrix. The result contains predicted probability of each data point belonging to each
class.

**Mak Wai Keong** December 15, 2018 at 7:17 pm #                                                                                    REPLY ↩

Hi all

I am very keen to know how Xgb can be used in the context of Learning (such as intelligent tutoring – to pick right context of knowledge for user).

I am new in this area, but is very keen to apply AI to learning.
One way I saw was the use of Dialogs to know what is known and what is not known, and what is to be known.

Looking forward for you experts for tips and advice

---

**Jason Brownlee** December 16, 2018 at 5:22 am  #          REPLY ↩

Start with a strong definition of your problem as a supervised learning problem then apply xgboost. This framework will help:

https://machinelearningmastery.com/how-to-define-your-machine-learning-problem/

---

**Monchy** January 15, 2019 at 12:48 am  #          REPLY ↩

I see that one-hot encoding of factor variables is required. However, in my R implementation XGBoost performs without any error or warning messages when I include factors. Does the algorithm ignore these variables?

---

**Jason Brownlee** January 15, 2019 at 5:53 am  #          REPLY ↩

I think R handles the factors automatically.

No, they are not ignored.

---

**Amal** January 22, 2019 at 8:59 pm  #          REPLY ↩

Hello
I need to know what it the best to use in case of binary classification: xgboost or logistic regression with gradient discent and why
thank you so much

---

**Jason Brownlee** January 23, 2019 at 8:47 am  #          REPLY ↩

It is not knowable. You must test a suite of methods and discover what works best for a specific dataset.

**Thomas** March 22, 2019 at 5:58 am #           REPLY ↰

All the bells and whistles are there but the meat of the algorithm is extremely poorly presented. Can't believe this is listed second on Google.

**Jason Brownlee** March 22, 2019 at 8:42 am #       REPLY ↰

Sorry to hear that Thomas.

What do you think was missing exactly? What would you like to see?

**Soroosh** June 13, 2019 at 12:23 pm #       REPLY ↰

Two main points:

1) Comparing XGBoost and Spark Gradient Boosted Trees using a single node is not the right comparison. Spark GBT is designed for multi-computer processing, if you add more nodes, the processing time dramatically drops while Spark manages the cluster. XGBoost can be run on a distributed cluster, but on a Hadoop cluster.

2) XGBoost and Gradient Boosted Trees are bias-based. They reduce variance too, but not as good as variance-based models like Random Forest), so when you are dealing with Kaggle datasets XGBoost works well, but when you are dealing with the real world and data streaming problem, Random Forest is a more stable model (stability in terms of handling high variance data which happens a lot in streaming data)

**Jason Brownlee** June 13, 2019 at 2:36 pm #       REPLY ↰

Thanks.

**Abdallah Elbohy** July 22, 2019 at 12:50 pm #       REPLY ↰

Thanks for adding information. But aren't there all datasets in kaggle in a real-world? And which datasets will be more stable with random forests than in XGBoost?

**Jason Brownlee** July 22, 2019 at 2:08 pm #                    REPLY ↩

I think so.

Tabular data is often best solved with xgboost, compared to neural nets or other methods.

**Mrudhula** September 6, 2019 at 4:39 am #                    REPLY ↩

may i know the disadvantages of xgboost sir???

**Jason Brownlee** September 6, 2019 at 5:06 am #                    REPLY ↩

Good question.

It can be slow.
It can create a complex model.
…

**ralph** September 8, 2019 at 8:03 pm #                    REPLY ↩

Hi, and thanks for this very clear post!

Just to make sure I understand properly: if speed is not a concern, xgboost will bring nothing more than a classical random forest, right?

**Jason Brownlee** September 9, 2019 at 5:14 am #                    REPLY ↩

No, it is a different algorithm called stochastic gradient boosting, and it offers both performance (skill) and speed improvements over other implementations.

**Ashvin** November 14, 2019 at 1:38 am #                    REPLY ↩

Thanks for this article. Is it possible to decompose a dependent variable using XGBOOST, like coefficient times variable in a Linear Model?

**Jason Brownlee** November 14, 2019 at 8:05 am #                    REPLY ↩

Not really, no.

---

**Anthony The Koala** November 25, 2019 at 6:14 am #

Dear Dr Jason,
The "pipped" version of xgboost crashed when using the demonstration of "learning_rate on the Pima Indians Onset of Diabetes dataset" in one of your 'crash courses'.
Definition: "pipped" meaning pip install –upgrade xgboost.

Solution: while the solution worked for me, I cannot guarantee that it will work for you if your xgboost has crashed. This to get the *.whl version at https://www.lfd.uci.edu/~gohlke/pythonlibs/. Search for xgboost and obtain the suitable version of the *.whl file for the particular versoin of Python and whether you are using a 32-bit or 62-bit version of the Python interpreseter.

The direct link for example for the python v3.6 and 64-bit version
https://download.lfd.uci.edu/pythonlibs/t7epjj8p/xgboost-0.90-cp36-cp36m-win_amd64.whl

Then in your command window, you say:

```
1  pip install xgboost-0.90-cp36-cp36m-win_amd64.whl
```

Thank you,
Anthony of Sydney

---

**Jason Brownlee** November 25, 2019 at 6:33 am #

Thanks for sharing.

---

**manuela** January 14, 2020 at 5:14 pm #

several research papers use xgboost for feature engineering, is it possible to use it in resarch paper as just an algorithm for enhancing pridiction of other classifier approch?

---

**Jason Brownlee** January 15, 2020 at 8:19 am #

You can use xgboost anyway you want, e.g. for feature selection.

As for describing in the use in a research paper, I cannot comment.

---

**Apoorv Vishnoi** February 27, 2020 at 5:37 pm #

Hi Jason,

Thanks for writing this article. There is a doubt that I have not been able to clear, even after attempting to read the original paper on xgboost. Like Adaboost does XGB also weigh each sample differently for subsequent models?

**Jason Brownlee** February 28, 2020 at 6:00 am #                                           REPLY ↩

I believe so. It is key to "boosting".

**Igor Ost** April 22, 2020 at 7:00 am #                                                       REPLY ↩

IMHO " .. must be a `part` of your …" not `apart`

**Jason Brownlee** April 22, 2020 at 7:47 am #                                                REPLY ↩

Thanks! Fixed.

**Miguel** August 16, 2020 at 10:03 am #                                                      REPLY ↩

Jason, could you please explain what does "structured or tabular data" in this context? As opposed to…

Is XGBoost suitable for time series?

Thank you very much for all you excellent material.

**Jason Brownlee** August 17, 2020 at 5:44 am #                                               REPLY ↩

Data in a spreadsheet. A table of data.

Yes, xgboost can be used for time series if we change the time series data to look like tabular data. Here is an example:

https://machinelearningmastery.com/xgboost-for-time-series-forecasting/

**Hussain** November 16, 2021 at 4:06 am #                                                     REPLY ↩

is xgboost discriminative model?

**Adrian Tam** November 16, 2021 at 3:14 pm #                    REPLY ↩

Yes, as you provide data to it rather than it generate data for you.

**Jack** May 24, 2022 at 3:23 am #                    REPLY ↩

Hello, Any Ideas if we can use it in multiple instance classification, So if we have a dataset with multiple systems and each system have multiple rows.

Or would you recommend using a different approach? if so what ?

Thanks in advance.

**Thomas Hemming** July 21, 2022 at 5:12 pm #                    REPLY ↩

Thanks! Do you have, or can you guide me, to any general introduction to intuition/rationale behind XGBoost and how its different from other approaches? I'm not really looking for specific coding examples, more examples of how it can be applied in different situations. E.g. how/why is a linear regression different from a regression with XGBoost. I see many applied things, but very little on the intuition. What I'm looking for is something that will enable to apply XGBoost package to my problem (whatever it might be) based on my knowledge of how the models work etc. etc. Even though I have a background in maths and statistics, that's not really what I'm looking for. But I want to intuitively understand what is going on 🙂

**James Carmichael** July 22, 2022 at 8:25 am #                    REPLY ↩

Hi Thomas…You may find the following discussion of interest:

https://www.kaggle.com/general/196541

**Romy** September 7, 2022 at 12:00 pm #                    REPLY ↩

I tried to apply both XGBoost Classifier (XGBC) and Random Forest Classifier (RFC) on the same Pima-Indians-Diabetes data, along with data imputation to eliminate features with close to 50% missing values. After eliminating 'test' feature (close to 50% missing data), the MAE for RFC was lower than that of XGBC. However, when I test prediction based on 1 row of the original data, RFC misclassified it while

XGBC predicted it correctly. Is this part of the 'inaccuracy' portion of the model or is the MAE not good enough to make us choose the model to use? I am just wondering..appreciate your feedback. Thanks.

---

**James Carmichael** September 8, 2022 at 5:45 am  #                    REPLY ↰

Hi Romy…The following resource may prove beneficial in terms of tuning your XGBoost models:

https://machinelearningmastery.com/tune-xgboost-performance-with-learning-curves/

## Leave a Reply

Name (required)

Email (will not be published) (required)

SUBMIT COMMENT

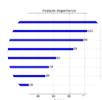**Welcome!**
I'm *Jason Brownlee* PhD
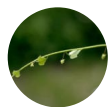and I **help developers** get results with **machine learning**.
Read more

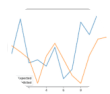## Never miss a tutorial:

## Picked for you:

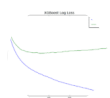Feature Importance and Feature Selection With XGBoost in Python

How to Develop Your First XGBoost Model in Python

Data Preparation for Gradient Boosting with XGBoost in Python

How to Use XGBoost for Time Series Forecasting

Avoid Overfitting By Early Stopping With XGBoost In Python

## Loving the Tutorials?

The XGBoost With Python EBook is
where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

LinkedIn | Twitter | Facebook | Newsletter | RSS

Privacy | Disclaimer | Terms | Contact | Sitemap | Search