

Never miss a tutorial:



Learning Mastery
Making Developers Awesome at Machine Learning

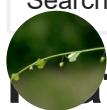
Picks just for you:



Feature Importance and Feature Selection With XGBoost in Python

[Click to Take the FREE XGBoost Crash-Course](#)

Search...



How to Develop Your First XGBoost Model in Python

Feature Importance and Feature Selection With XGBoost in Python

How to Use XGBoost for Time Series Forecasting
Browne · on August 31, 2016 in [XGBoost](#)



Tweet

Tweet

Share

Share



Data Preparation for Gradient Boosting
Updated on August 27, 2020
with XGBoost in Python

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.

Avoid Overfitting By Early Stopping With

XGBoost In Python
In this post you will discover how you can estimate the importance of features for a predictive modeling problem using the XGBoost library in Python.

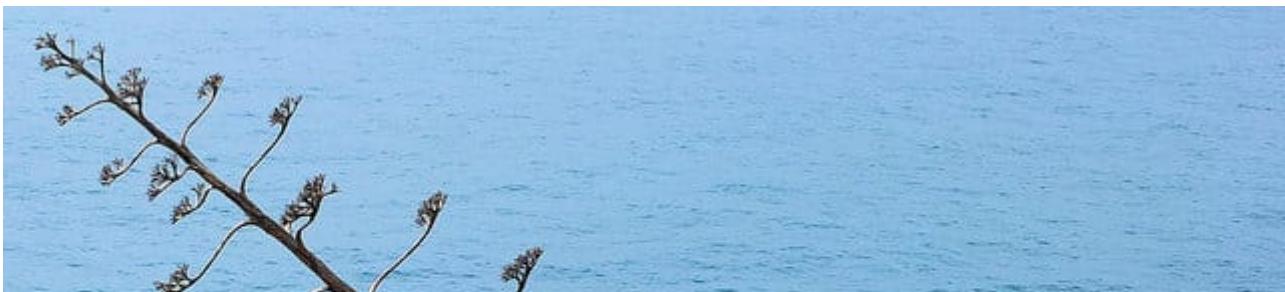
After reading this tutorial you will know:

- The XGBoost With Python EBook is where you'll find the **Really Good** stuff.
- How feature importance is calculated using the gradient boosting algorithm.
- How to plot feature importance in Python calculated by the XGBoost model.
- How [feature selection](#) is calculated by XGBoost to perform feature selection.
[>> SEE WHAT'S INSIDE](#)

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.

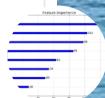
- Update Jan/2017:** Updated to reflect changes in scikit-learn API version 0.18.1.
- Update Mar/2018:** Added alternate link to download the dataset as the original appears to have been taken down.
- Update Apr/2020:** Updated example for XGBoost 1.0.2.



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



How to Use XGBoost for Time Series Forecasting

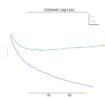
Feature Importance and Feature Selection With XGBoost in Python

Photo by Keith Roper, some rights reserved.



Data Preparation for Gradient Boosting with XGBoost in Python

Need help with XGBoost in Python?



Avoid Overfitting By Early Stopping 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.

Loving the Tutorials?

Start Your FREE Mini-Course Now

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Feature importance in Gradient Boosting

A benefit of using gradient boosting is that after the boosted trees are constructed, it is relatively straightforward to retrieve importance scores for each attribute.

Generally, importance provides a score that indicates how useful or valuable each feature was in the construction of the boosted decision trees within the model. The more an attribute is used to make key decisions with decision trees, the higher its relative importance.

This importance is calculated explicitly for each attribute in the dataset, allowing attributes to be ranked and compared to each other.

Importance is calculated for a single decision tree by the amount that each attribute split point improves the performance measure, weighted by the number of observations the node is responsible for. The performance measure may be the purity (Gini index) used to select the split points or another more specific error function.

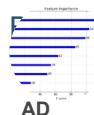
The feature importances are then averaged across all of the the decision trees within the model.

Never miss a tutorial:

For more technical information on how feature importance is calculated in boosted decision trees, see Section 13.1 “*Relative Importance of Predictor Variables*” of the book **The Elements of Statistical Learning: Data Mining, Inference, and Prediction**, page 367.

Picked for you:

Also, see Matthew Drury answer to the StackOverflow question “[Relative variable importance for Feature Selection With XGBoost in Python](#)” where he provides a very detailed and practical answer.

 Feature Importance and Feature Selection With XGBoost in Python

AD



How to Develop Your First XGBoost Model in Python

Manually Plot Feature Importance

 [How to Use XGBoost for Time Series Forecasting](#)

A trained XGBoost model automatically calculates feature importance on your predictive modeling problem.



Data Preparation for Gradient Boosting

 [with XGBoost in Python](#) 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_)
```

We can plot these scores on a bar chart directly to get a visual indication of the relative importance of each feature in the dataset. For example:

Loving the Tutorials?

```
1 # plot  The XGBoost with Python EBOOK IS
2 pyplot.bar(range(len(model.feature_importances_)), model.feature_importances_)
3 pyplot.show()
```

We can d  >> SEE WHAT'S INSIDE j an XGBoost model on the Pima Indians onset of diabetes dataset and creating a bar chart from the calculated feature importances.

Download the dataset and place it in your current working directory.

- Dataset File.
- Dataset Details.

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

```

14 print(model.feature_importances_)
15 # or miss a tutorial:
16 pyplot.bar(range(len(model.feature_importances_)), model.feature_importances_)
17 pyplot.show()

```

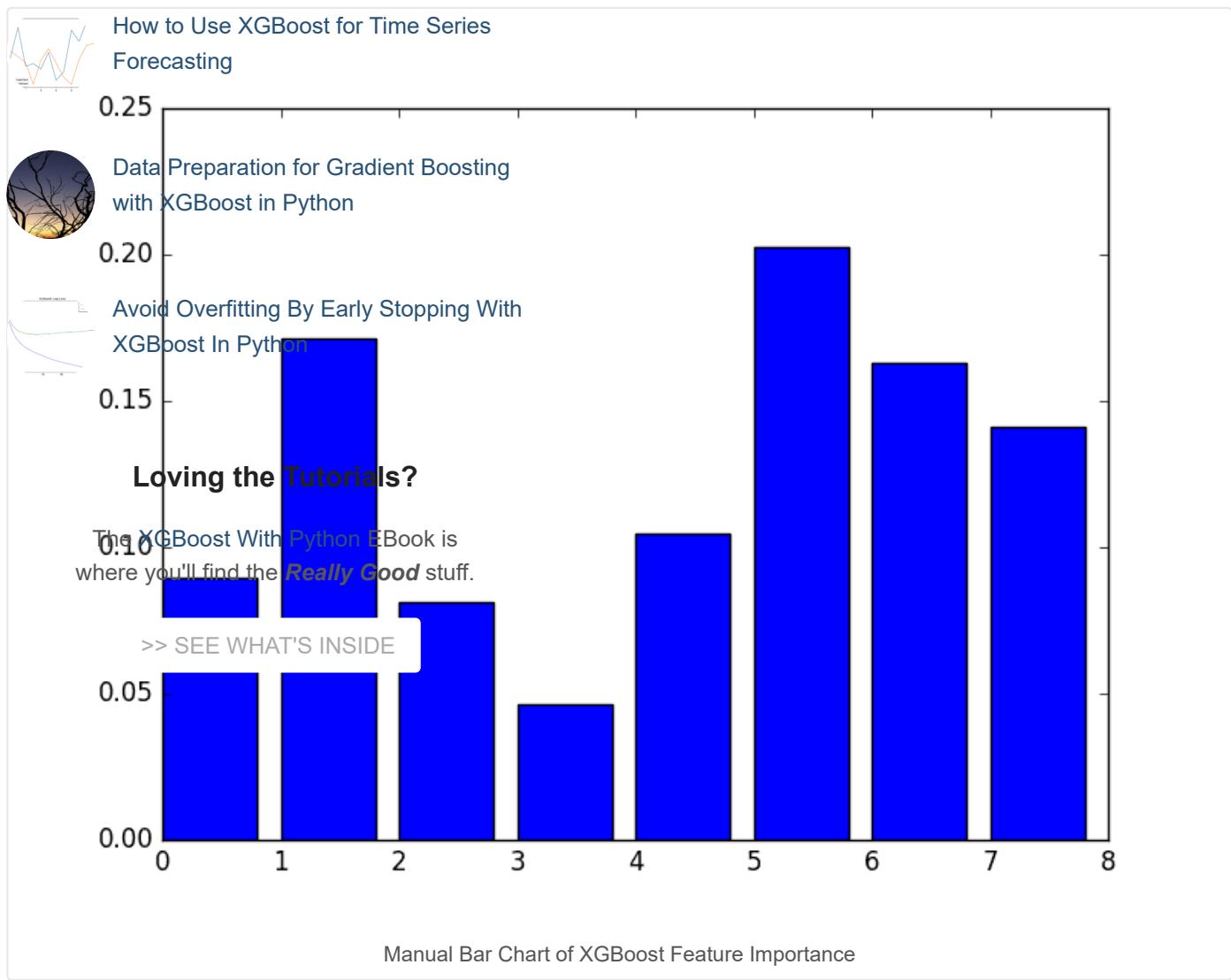


Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or picked for numerical precision. Consider running the example a few times and compare the average outcome.

Feature Importance and Feature Selection With XGBoost in Python
g this example first outputs the importance scores.



We also get a bar chart of the relative importances.



A downside of this plot is that the features are ordered by their input index rather than their importance. We could sort the features before plotting.

Thankfully, there is a built in plot function to help us.

AD Never miss a tutorial:



Using the Built-in XGBoost Feature Importance Plot

Feature Importance and Feature Selection With XGBoost in Python 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 Model in Python
2 plot_importance(model)
3 pyplot.show()
```

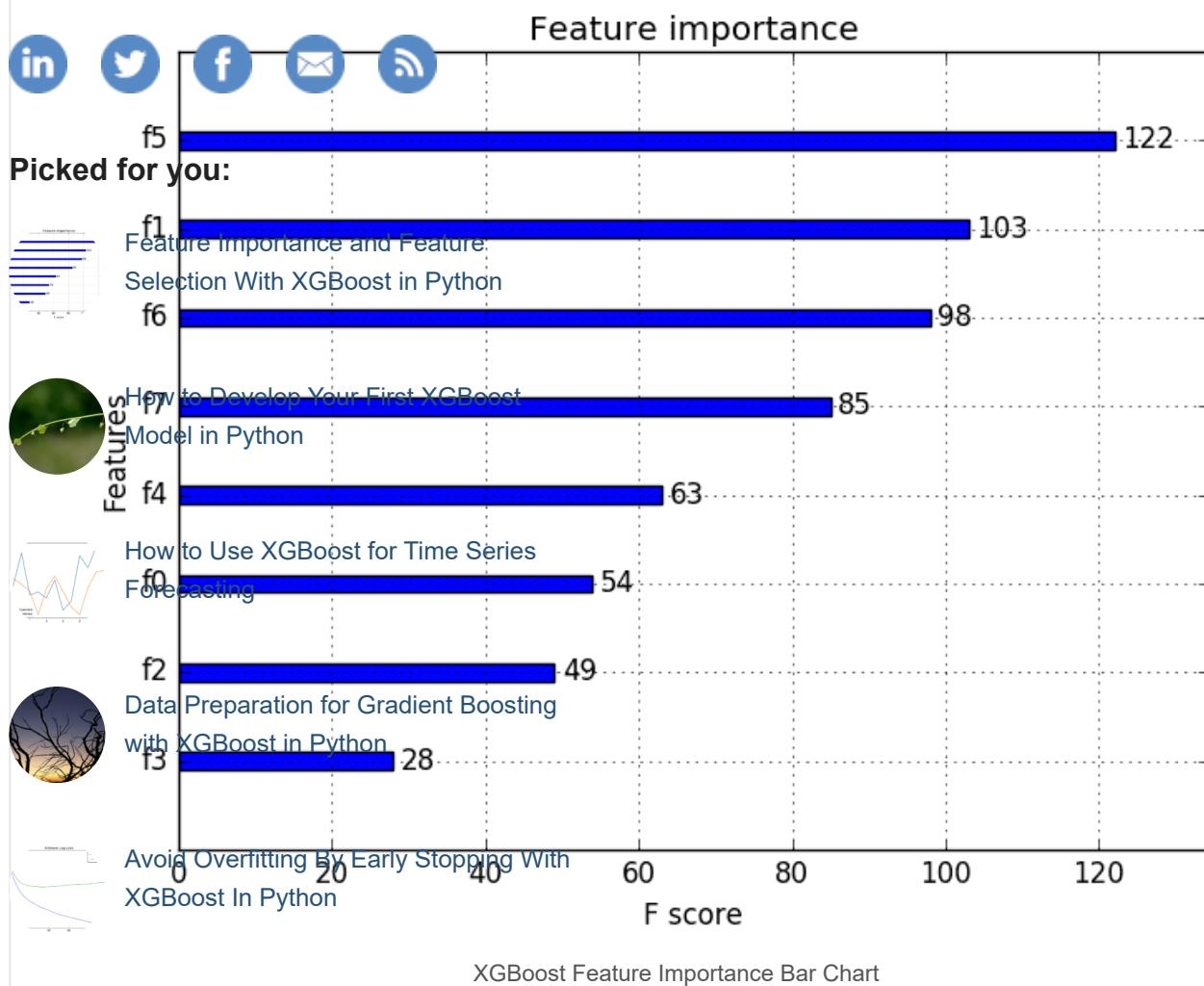
How to Use XGBoost for Time Series Forecasting Example, here is a complete code listing plotting the feature importance for the Pima Indians using the built-in **plot_importance()** function.

```
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()
```

Note: You >> SEE WHAT'S INSIDE the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

Running the example gives us a more useful bar chart.

Never miss a tutorial:



Loving the Tutorials?
You can see that features are automatically named according to their index in the input array (X) from F0 to F7.

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.
Manually mapping these indices to names in the problem description, we can see that the plot shows F5 (body fat percentage) has the highest importance and F3 (skin fold thickness) has the lowest importance.



Feature Selection with XGBoost Feature Importance Scores

Feature importance scores can be used for feature selection in scikit-learn.

This is done using the `SelectFromModel` class that takes a model and can transform a dataset into a subset with selected features.

This class can take a pre-trained model, such as one trained on the entire training dataset. It can then use a threshold to decide which features to select. This threshold is used when you call the `transform()`

method on the `SelectFromModel` instance to consistently select the same features on the training dataset and the test dataset.



In the example, we first train and then evaluate an XGBoost model on the entire training dataset and test datasets respectively.

Picked for you:

Using the feature importances calculated from the training dataset, we then wrap the model in a `SelectFromModel` instance. We use this to select features on the training dataset, train a model from the selected subset of features, then evaluate the model on the testset, subject to the same feature selection scheme.



How to Develop Your First XGBoost Model in Python

```
1 # select features using threshold
2 selection = SelectFromModel(model, threshold=thresh, prefit=True)
3 select_X_train = selection.transform(X_train)
4 # train model
5 selection_model = XGBClassifier()
6 selection_model.fit(select_X_train, y_train)
7 # eval model
8 select_X_test = selection.transform(X_test)
9 y_pred = selection_model.predict(select_X_test)
```

For interest, we can test multiple thresholds for selecting features by feature importance. Specifically, feature importance for each individual feature is supported, essentially allowing us to test each subset of features by importance, starting with all features and ending with a subset with the most important feature.

The complete code listing is provided below.

Developing the Tutorial

```
1 # use feature_importance_for_feature_selection
2 from numpy import loadtxt
3 # from numpy will find the Really Good stuff.
4 from xgboost import XGBClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from sklearn.feature_selection import SelectFromModel
8 # load data
9 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=',')
10 # split data into X and y
11 X = dataset[:,0:8]
12 Y = dataset[:,8]
13 # split data into train and test sets
14 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.33, random_state=7)
15 # fit model on all training data
16 model = XGBClassifier()
17 model.fit(X_train, y_train)
18 # make predictions for test data and evaluate
19 y_pred = model.predict(X_test)
20 predictions = [round(value) for value in y_pred]
21 accuracy = accuracy_score(y_test, predictions)
22 print("Accuracy: %.2f%%" % (accuracy * 100.0))
23 # Fit model using each importance as a threshold
24 thresholds = sort(model.feature_importances_)
25 for thresh in thresholds:
26     # select features using threshold
27     selection = SelectFromModel(model, threshold=thresh, prefit=True)
28     select_X_train = selection.transform(X_train)
29     # train model
30     selection_model = XGBClassifier()
```

```

31 selection_model.fit(select_X_train, y_train)
32 # eval model
33 select_X_test = selection.transform(X_test)
34 y_pred = selection_model.predict(select_X_test)
35 predictions = [round(value) for value in y_pred]
36 accuracy = accuracy_score(y_test, predictions)
37 print("Thresh=%f, n=%d, Accuracy: %.2f%%" % (thresh, select_X_train.shape[1], accuracy*100))

```

Picked for you:

Note, if you are using XGBoost 1.0.2 (and perhaps other versions), there is a bug in the XGBClassifier

that results in [None](#) for Feature importance and Feature Selection With XGBoost in Python

1 [KeyError: 'weight'](#)

 How to Develop Your First XGBoost Model in Python
can be fixed by using a custom `XGBClassifier` class that returns `None` for the `coef_` property.

The complete example is listed below.

```

How to Use XGBoost for Time Series Forecasting

1 # use feature importance for feature selection, with fix for xgboost 1.0.2
2 from numpy import loadtxt
3 from numpy import sort
4 from xgboost import XGBClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from sklearn.feature_selection import SelectFromModel
8
9 # define custom class to fix bug in xgboost 1.0.2
10 class MyXGBClassifier(XGBClassifier):
11     @property
12     def coef_(self):
13         return None
14
15 # load data
16 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=',')
17 # split data into X and y
18 X = dataset[:,0:8]
19 Y = dataset[:,8]
20 # split data into train and test sets
21 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.33, random_state=7)
22 # fit model on all training data
23 model = MyXGBClassifier()
24 model.fit(X_train, y_train)
25 # make predictions for test data and evaluate
26 predictions = model.predict(X_test)
27 accuracy = accuracy_score(y_test, predictions)
28 print("Accuracy: %.2f%%" % (accuracy * 100.0))
29 # Fit model using each importance as a threshold
30 thresholds = sort(model.feature_importances_)
31 for thresh in thresholds:
32     # select features using threshold
33     selection = SelectFromModel(model, threshold=thresh, prefit=True)
34     select_X_train = selection.transform(X_train)
35     # train model
36     selection_model = XGBClassifier()
37     selection_model.fit(select_X_train, y_train)
38     # eval model
39     select_X_test = selection.transform(X_test)
40     predictions = selection_model.predict(select_X_test)
41     accuracy = accuracy_score(y_test, predictions)
42     print("Thresh=%f, n=%d, Accuracy: %.2f%%" % (thresh, select_X_train.shape[1], accuracy*100.0))

```

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average

outcome.

Never miss a tutorial:

Running this example prints the following output.



```
1 Accuracy: 77.95%
2 Thresh=0.071, n=8, Accuracy: 77.95%
3 Thresh=0.073, n=7, Accuracy: 76.38%
4 Thresh=0.084, n=6, Accuracy: 77.56%
5 Thresh=0.090, n=5, Accuracy: 76.38%
6 Thresh=0.128, n=4, Accuracy: 76.38%
7 Thresh=0.160, n=3, Accuracy: 74.80%
8 Thresh=0.186, n=2, Accuracy: 71.65%
9 Thresh=0.208, n=1, Accuracy: 63.78%
```



How to Develop Your First XGBoost

I see that the performance of the model generally decreases with the number of selected features.



How to Use XGBoost for Time Series

Forecasting

problem there is a trade-off of features to test set accuracy and we could decide to take a less complex model (fewer attributes such as n=4) and accept a modest decrease in estimated accuracy from 77.95% down to 76.38%.



Data Preparation for Gradient Boosting

likely to be a wash on such a small dataset, but may be a more useful strategy on a larger

dataset and using cross validation as the model evaluation scheme.



Avoid Overfitting By Early Stopping With

XGBoost In Python



Data Mesh Architectures with Event Streams

Loving the Tutorials? Summary

The [XGBoost With Python](#) EBook is where you'll find the **Really Good Stuff**. In this post you discovered how to access features and use importance in a trained XGBoost gradient boosting model.

>> SEE WHAT'S INSIDE

Specifically, you learned:

- What feature importance is and generally how it is calculated in XGBoost.
- How to access and plot feature importance scores from an XGBoost model.
- How to use feature importance from an XGBoost model for feature selection.

Do you have any questions about feature importance in XGBoost or about this post? Ask your questions in the comments 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](#)

How to Use XGBoost for Time Series Forecasting
Tweet [Share](#) [Share](#)

 Data Preparation for Gradient Boosting with XGBoost in Python

AD

Avoid Overfitting By Early Stopping With XGBoost In Python

Children of Heroes Charity fund

Give a brighter future to Ukrainian orphans

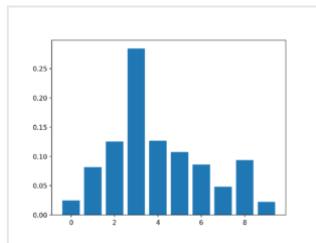
[FIND OUT HOW](#)

Loving the Tutorials?

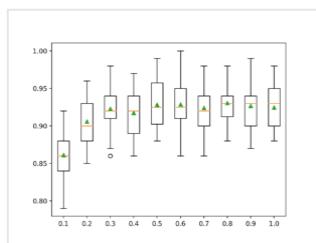
The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.

More On This Topic

[>> SEE WHAT'S INSIDE](#)



How to Calculate Feature Importance With Python



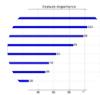
Extreme Gradient Boosting (XGBoost) Ensemble in Python

Never miss a tutorial:



Picked for you:

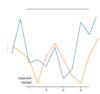
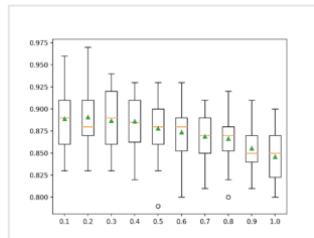
A Gentle Introduction to XGBoost for Applied Machine...



Feature Importance and Feature Selection With XGBoost in Python



How to Develop Your First XGBoost Model in Python

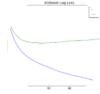
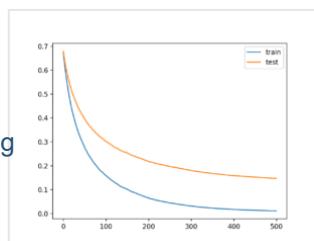


How to Develop Random Forest Ensembles With XGBoost

How to Use XGBoost for Time Series Forecasting



Data Preparation for Gradient Boosting with XGBoost in Python



Tune XGBoost Performance With Learning Curves
Avoid Overfitting By Early Stopping With XGBoost In Python



Loving the Tutorials?



>> SEE WHAT'S INSIDE **ownlee**

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 →](#)

< [How to Visualize Gradient Boosting Decision Trees With XGBoost in Python](#)

[Avoid Overfitting By Early Stopping With XGBoost In Python >](#)

208 Responses to *Feature Importance and Feature Selection With XGBoost in Python*

Trupti December 9, 2016 at 5:23 pm #

REPLY ↗

Hi. I am running “select_X_train = selection.transform(X_train)” where x_train is the data with dependent variables in few rows.

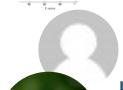


Request your help.

Picked for you:

Thanks!

 Feature Importance and Feature Selection With XGBoost in Python

 Trupti December 9, 2016 at 5:28 pm # 

 How to Develop Your First XGBoost Model in Python
“sorry the error is “TypeError: only length-1 arrays can be converted to Python scalars”.

 How to Use XGBoost for Time Series  Jason Brownlee December 10, 2016 at 8:04 am # 

Check the shape of your X_train, e.g. print(X_train.shape)

 You Data Preparation Before Gradient Boosting with XGBoost in Python

 Avoid Overfitting By Early Stopping With XGBoost In Python  July 6, 2021 at 8:04 am # 

How would you solve this? I have the same issue

Loving the Tutorials?

 Jason Brownlee July 4, 2021 at 5:58 am # 
The XGBoost With Python EBook is where you'll find the **Really Good** stuff.
This may help:

 gmastry.com/index-slice-reshape-numpy-arrays-machine-learning->> SEE WHAT'S INSIDE

 sa January 5, 2017 at 3:44 pm # 

I tried to select features for xgboost based on this post (last part which uses thresholds) but since I am using gridsearch and pipeline, this error is reported:

```
select_X_train = selection.transform(X_train)
File "C:\Users\Markazi.co\Anaconda3\lib\site-packages\sklearn\feature_selection\base.py", line 76, in
transform
mask = self.get_support()
File "C:\Users\MM.co\Anaconda3\lib\site-packages\sklearn\feature_selection\base.py", line 47, in
get_support
mask = self._get_support_mask()
File "C:\Users\Markazi.co\Anaconda3\lib\site-packages\sklearn\feature_selection\from_model.py", line
201, in _get_support_mask
scores = _get_feature_importances(estimator)
```

File "C:\Users\Markazi.co\Anaconda3\lib\site-packages\sklearn\feature_selection\from_model.py", line

Never miss a tutorial:

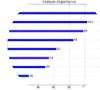
32, in _get_feature_importances

% estimator. class _name_)

value. For: The underlying estimator method has no coef_ or feature_importances_ attribute. Either pass a fitted estimator to SelectFromModel or call fit before calling transform.

Picked for you:

regards,



Feature Importance and Feature

Selection With XGBoost in Python

Jason Brownlee

January 6, 2017 at 9:05 am #

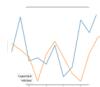
REPLY ↗



How to Develop Your First XGBoost

Model in Python

Consider trying the example without Pipelines first, get it working, then try adding in additional complexity.



How to Use XGBoost for Time Series

Forecasting



sa January 6, 2017 at 5:53 pm #

REPLY ↗

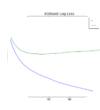


Data Preparation for Gradient Boosting

Hello Mr. Brownlee

with XGBoost in Python

Thanks



I already tried the example without Pipelines , and it works well. After adding pipeline, it could

Avoid Overfitting By Early Stopping With Extract Feature Importance but after that it fails. Thanks.

XGBoost In Python

Best regards,

Loving the Tutorials?



Johnn January 15, 2017 at 12:28 pm #

REPLY ↗

The XGBoost With Python EBook is

where you'll find the **Really Good** stuff
Thanks for the post. I don't understand what's the meaning of "F-score" in the x-axis of the feature importance plot — And what is the number next to each of the bar?

>> SEE WHAT'S INSIDE



Jason Brownlee

January 16, 2017 at 10:36 am #

REPLY ↗

Hi Johnn,

You can learn more about the F1 score here:

https://en.wikipedia.org/wiki/F1_score

The number is a scaled importance, it really only has meaning relative to other features.



Gonçalo Abreu

June 5, 2017 at 10:32 pm #

REPLY ↗

Hey Jason,

Are you sure the F score on the graph is realted to the tradicional F1-score?

I found this github page where the owner presents many ways to extract feature importance
Never miss a tutorial:
meaning from xgb. His explanation abou the F measure seems to have no relation to F1



Picked for you

 **Jason Brownlee** June 6, 2017 at 9:36 am #

REPLY ↗



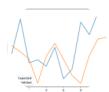
[Feature Importance and Feature Selection With XGBoost in Python](#) Importance scores are different from F scores. The above tutorial focuses on feature importance scores.



[How to Develop Your First XGBoost Model in Python](#)



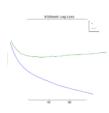
Domi April 20, 2020 at 9:29 pm #



[How to Use XGBoost for Time Series Forecasting](#)



[Data Preparation for Gradient Boosting with XGBoost in Python](#) Gonçalo has right , not the F1 score was the question. F1 score is totally different from the F score in the feature importance plot.



[Avoid Overfitting By Early Stopping With XGBoost In Python](#) F score in the feature importance context simply means the number of times a feature is used to split the data across all trees. at least, if you are using the built-in feature of Xgboost.

I hope it helped to clarify things.

Cheers,

Loving the Tutorials?

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.



Jason Brownlee April 21, 2020 at 5:54 am #

>> SEE WHAT'S INSIDE

Thanks for sharing!



tuttoaposto June 23, 2020 at 4:34 pm #

REPLY ↗

`plot_importance()` by default plots feature importance based on `importance_type = 'weight'`, which is the number of times a feature appears in a tree.

It is confusing when compared to `clf.feature_importances_`, which by default is based on normalized gain values.

You can check the correspondence between the plot and the `feature_importance_` values using this code:

```
# How to get back feature_importances_ (gain based) from plot_importance fscore
# Calculate two types of feature importance:
# Weight = number of times a feature appears in tree
# Gain = average gain of splits which use the feature = average all the gain values of the feature if it
```

appears multiple times

Never miss a tutorial:

Normalized gain = Proportion of average gain out of total average gain



k[f.groupby('Feature').agg(fscore = ('Gain', 'count'))]

group = k[k['Feature']!='Leaf'].groupby('Feature').agg(fscore = ('Gain', 'count'))

feature_importance_gain = ('Gain', 'mean'))

Picked for you:

Feature importance same as plot_importance(importance_type = 'weight'), default value

group[feature_importance_startvalue=featureending=False]

Selection With XGBoost in Python

Feature importance same as clf.feature_importance_ default = 'gain'

group['feature_importance_gain_norm'] =

group['feature_importance_gain']/group['feature_importance_gain'].sum()

How to Develop Your First XGBoost Model in Python

group['feature_importance_gain_norm'].sort_values(by='feature_importance_gain_norm', ascending=False)



Feature importance same as plot_importance(importance_type = 'gain')

How to Use XGBoost for Time Series Forecasting

Note:

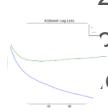


1. Features with zero feature_importance_ don't show in trees_to_dataframe(). You can check what they are with:

Data Preparation for Gradient Boosting

with XGBoost in Python

X_train.columns[[x not in k['Feature'].unique() for x in X_train.columns]]



2. The feature importance ranks for 'weight' and 'gain' types can be quite different. Be careful when choosing features based on the plot! I would choose gain over weight because gain reflects the feature's power of grouping similar instances into a more homogeneous child node at the split.

Loving the Tutorials!



June 24, 2020 at 6:23 am #

REPLY ↗

The XGBoost With Python Ebook is fantastic comment!

where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE



Soyoung Kim April 20, 2017 at 2:39 am #

REPLY ↗

Hi Jason,

Your postings are always amazing for me to learn ML techniques!

Especially this XGBoost post really helped me work on my ongoing interview project.

The task is not for the Kaggle competition but for my technical interview! 😊

I used your code to generate a feature importance ranking and some of the explanations you used to describe techniques.

You can find it here: <https://www.kaggle.com/soyoungkim/two-sigma-connect-rental-listing-inquiries/rent-interest-classifier>

I also put your link in the reference section.

Please let me know if it is not appropriate for me to use your code.

Never miss a tutorial:

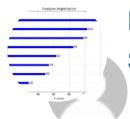
Jason Brownlee April 20, 2017 at 9:31 am #

REPLY ↗



As long as you cite the source, I am happy.

Picked for you:

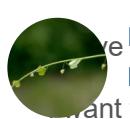


Feature Importance and Feature

Selection With XGBoost in Python

Ztara April 27, 2017 at 3:42 pm #

REPLY ↗



Hi Jason,

I have some questions about feature importance.
How to Develop Your First XGBoost Model in Python



I want to use the features that selected by XGBoost in other classification models, and I got confused on how to get the right scores of features, I mean that is it necessary to adjust parameters to get the best model and obtain the corresponding scores of features? In other words, how can I get the scores of features in the model?

Thanks a lot.



Data Preparation for Gradient Boosting
with XGBoost in Python



Jason Brownlee April 28, 2017 at 7:36 am #

REPLY ↗

Avoid Overfitting By Early Stopping With
XGBoost In Python
The scores are relative.

You can use them as a filter and select all features with a score above x, e.g. 0.5.

Loving the Tutorials?

max January 14, 2018 at 12:00 pm #

REPLY ↗

The XGBoost With Python EBook is

where you'll find the **Really Good** stuff.

Hi Jason, I know that choosing a threshold (like 0.5) is always arbitrary ...but is there a

>> SEE WHAT'S INSIDE

thanks a lot.



Jason Brownlee January 15, 2018 at 6:55 am #

REPLY ↗

Yes, start with 0.5, tune if needed.



Joe Butkovic August 26, 2019 at 10:58 pm #

Hi Jason,

Is there a score which should be discounted? For example, my highest score is 0.27, then 0.15, 0.13... Should I discount the model all together? Thanks!

Never miss a tutorial:

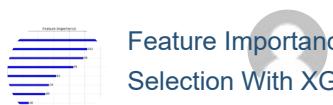


Jason Brownlee August 27, 2019 at 6:46 am #



Scores are relative. Test different cut-off values on your specific dataset.

Picked for you:



Shubham Jaiswal August 31, 2019 at 9:28 am #

Selection With XGBoost in Python

One good way to not worry about thresholds is to use something like –
CalibratedClassifierCV(clf, cv='prefit', method='sigmoid').



How to Develop Your First XGBoost Model in Python



Jason Brownlee September 1, 2019 at 5:34 am #

Nice, thanks for sharing this tip!



Data Preparation for Gradient Boosting with XGBoost in Python

<https://machinelearningmastery.com/calibrated-classification-model-in-scikit-learn/>



Avoid Overfitting By Early Stopping With XGBoost

Omogbehin May 16, 2017 at 10:23 am #

REPLY ↗

Hello sir,

For the XGBoost feature selection. How do i change the Y axis to the names of my attributes. Kind regards sir.

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.



>> SEE WHAT'S INSIDE

y 17, 2017 at 8:23 am #

REPLY ↗

Great question, I'm not sure off-hand. You may need to use the xgboost API directly.



Franco Arda October 12, 2018 at 8:27 pm #

REPLY ↗

@Omogbehin, to get the Y labels automatically, you need to switch from arrays to Pandas dataframe. By doing so, you get automatically labeled Y and X.

```
column_names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
data = pd.read_csv("diabetes.csv", names = column_names)
X = data.iloc[:,0:8]
Y = data.iloc[:,8]
model = XGBClassifier()
model.fit(X, Y)

from xgboost import plot_importance
plot_importance(model)
```

plt.show()
Never miss a tutorial:



tuttoaposto

June 23, 2020 at 3:56 pm #

REPLY ↗

Picked for you: You can plot feature_importance directly as in:

```
clf = xgb.XGBClassifier(  
    Feature Importance and Feature  
    earning_rate=0.1  
    Selection With XGBoost in Python  
    n_estimators=1000,  
    max_depth=5,  
    min_child_weight=1  
    How to Develop Your First XGBoost  
    gaModel in Python  
    subsample=0.8,  
    colsample_bytree=0.8,  
    How to Use XGBoost for Time Series  
    Forecasting  
    nthread=4,  
    # scale_pos_weight=1,  
    num_class=6,  
    Data Preparation for Gradient Boosting  
    seed=0,  
    with XGBoost in Python  
    verbosity=0).fit(X_train, y_train)  
  
%matplotlib notebook  
figAnd Overtuning By Early Stopping With  
XGBoost in Python(clf, height = 0.4, grid = False, ax=ax, importance_type='weight')  
fig.subplots_adjust(left = 0.35);
```

2. Or you can also output a list of feature importance based on normalized gain values, i.e. gain/sum of gain.

Loving the Tutorials?

pd.Series(clf.feature_importances_, index=X_train.columns,
The XGBoost With Python EBOOK is
name='Feature Importances').sort_values(ascending=False)

>> SEE WHAT'S INSIDE



Jason Brownlee

June 24, 2020 at 6:22 am #

REPLY ↗

Great tips!



Simone

June 21, 2017 at 11:14 pm #

REPLY ↗

Hi Jason,

Is it possible using "feature_importances_" in XGBRegressor() ?



Jason Brownlee

June 22, 2017 at 6:06 am #

REPLY ↗

I'm not sure off the cuff, sorry.
Never miss a tutorial:



Simone June 22, 2017 at 7:06 am #

REPLY ↗

Picked for you:

Ok, I will try another method for features selection.



Feature Importance and Feature Selection With XGBoost in Python

Thanks



How to Develop Your First XGBoost Model in Python

Richard July 22, 2017 at 8:04 pm #

REPLY ↗



Hello Jason, I use the XGBRegressor and want to do some feature selection. However, although the 'plot_importance(model)' command works, when I want to retrieve the values using [How to Use XGBoost for Time Series Forecasting](#), it says 'AttributeError: 'XGBRegressor' object has no attribute 'feature_importances_'. Any hints how to retrieve the feature importances for regression?



Data Preparation for Gradient Boosting

with XGBoost in Python

Jason Brownlee July 23, 2017 at 6:23 am #

REPLY ↗



Sorry to hear that Richard. I'm not sure of the cause.
Avoid Overfitting By Early Stopping With

XGBoost In Python

REPLY ↗

Long.Ye August 23, 2017 at 10:41 am #

Loving the Tutorials?

Hi Jason,

The [XGBoost With Python EBook](#) is
Do you know some methods to qualify variable importance in RNN or LSTM? Could the XGBoost
where you'll find the **Really Good** stuff.
method be used in regression problems of RNN or LSTM? Thanks a lot.

>> SEE WHAT'S INSIDE



Jason Brownlee August 23, 2017 at 4:23 pm #

REPLY ↗

Perhaps, I have not tried.



Edward August 25, 2017 at 3:57 pm #

REPLY ↗

Can you explain how the decision trees feature importance also works?



Biswajit September 9, 2017 at 10:36 pm #

REPLY ↗

Hi Jason while trying to fit my model in Xgboost object it is showing the below error

OSError: [WinError -529697949] Windows Error 0xe06d7363

i am using 32 bit anaconda
Never miss a tutorial:

import platform



Please suggest how to get over this issue

Feature Importance and Feature Selection With XGBoost in Python

 **Jason Brownlee** September 11, 2017 at 12:01 pm #

REPLY ↗

Sorry, I have not seen this error.

How to Develop Your First XGBoost



Perhaps you can post to stackoverflow?

How to Use XGBoost for Time Series Forecasting

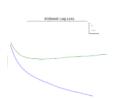
 **kim tae in** September 23, 2017 at 4:51 pm #

REPLY ↗

Hi Jason.



Data Preparation for Gradient Boosting
SelectFromModel(model, threshold=thresh, prefit=True)
with XGBoost in Python



I wonder what prefit = true means in this section. I checked on the sklearn site, but I do not understand.

Avoid Overfitting By Early Stopping With
XGBoost In Python



Jason Brownlee September 24, 2017 at 5:15 am #

REPLY ↗

Loving the Tutorials?

It specifies not to fit the model again, that we have already fit it prior.

The [XGBoost With Python](#) EBook is
where you'll find the **Really Good** stuff.



>> SEE WHAT'S INSIDE at 2:14 am #

REPLY ↗

Hi, Jason

Can you get feature importance of artificial neural network?

If you can, how?

Thanks very much.



Jason Brownlee January 19, 2018 at 6:35 am #

REPLY ↗

Perhaps, I have not seen this done.



Zhang January 25, 2018 at 11:41 pm #

REPLY ↗

Hi, Jason. I am doing a project with Stochastic gradient boosting. My database is clinical data and I think the ranking of feature importance can feed clinicians back with clinical knowledge, i.e., machine can tell us which clinical features are most important in distinguishing phenotypes of the diseases. What I did is to predict the phenotypes of diseases with all the variables of the database using SGB in the training set, and then test the performance of the model in testing set. If the testing is good (e.g., high accuracy and kappa), then I would like to say the ranking of the feature importance is reasonable as machine can make good prediction using this ranking information (i.e., the feature importance is the knowledge). Feature Importance and Feature Selection With XGBoost in Python
versa, if the prediction is poor I would like to say the ranking of feature importance is bad or even wrong. In this case we cannot trust the 'knowledge' feed back by the machine. In other words, it wastes time to do feature selection in this case because the feature importance is not correct (either because of the poor data quality or the machine learning algorithm is not suitable). May I ask whether my thinking above is reasonable?

My second question is that I did not do feature selection to identify a subset of features as you did in your post. I just treat the few features on the top of the ranking list as the most important clinical features. How to Use XGBoost for Time Series Forecasting
then did classical analysis like t test to confirm these features are statistically different in different phenotypes. Can I still name it as feature selection or feature extraction? I am little bit confused about these terms. Thanks and I am waiting for your reply.

Data Preparation for Gradient Boosting with XGBoost in Python

Jason Brownlee January 26, 2018 at 5:42 am #

REPLY ↗

Avoid Overfitting By Early Stopping With XGBoost In Python

Sorry, I'm not sure I follow. Perhaps you can distil your question into one or two lines?

Yes, you could still call this feature selection.

Loving the Tutorials?

Zhang January 26, 2018 at 8:09 pm #

REPLY ↗

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

As you >> SEE WHAT'S INSIDE << gradient boosting (SGB) is a model with built-in feature selection, which is thought to be more efficient in feature selection than wrapper methods and filter methods. But I doubt whether we can always trust the feature selected by SGB because the importance (relative influence) of the features are still provided by the model when the model has bad performance (e.g., very poor accuracy in testing). In this case, the model may be even wrong, so the selected features may be also wrong. So I would like to hear some comment from you regarding to this issue.

Thanks.

Jason Brownlee January 27, 2018 at 5:57 am #

REPLY ↗

Perhaps compare models fit with different subsets of features to see if it is lifting skill.

Try using an ensemble of models fit on different subsets of features to see if you can lift skill further.

Never miss a tutorial:

Sa January 29, 2018 at 3:25 pm #

REPLY ↗



Hi, n. Could you please let me know if the feature selection method that you used here, is classified as filter, Picked for you embedded feature selection method?



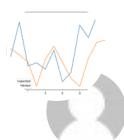
Feature Importance and Feature Selection With XGBoost in Python



Jason Brownlee [@Brownlee](#) January 30, 2018 at 9:47 am #

REPLY ↗

Here we are doing feature importance or feature scoring. It would be a filter.



How to Use XGBoost for Time Series

Forecasting

Youcai Wang February 2, 2018 at 9:37 am #

REPLY ↗



Hi Brownlee, if I have a dataset with 118 variables, but the target variable is in 116, and I want to use 6,115 and 117-118 variables as dependent variables, how can I modify the code `X = dataset[:,0:8]` `dataset[:,8:]` to get X and Y?



Avoid Overfitting By Early Stopping With XGBoost In Python
not figure out this simple question. Please help

Thanks,

Loving the Tutorials?



Nick March 18, 2018 at 9:47 am #

REPLY ↗

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.

Hi Jason,

Thanks >> SEE WHAT'S INSIDE

Did you notice that the values of the importances were very different when you used `model.get_importances_` versus `xgb.plot_importance(model)`?

I used these two methods on a model I just trained and it looks like they are completely different. Moreover, the numpy array `feature_importances` do not directly correspond to the indexes that are returned from the `plot_importance` function.

In other words, these two methods give me qualitatively different results. Any idea why?



Jason Brownlee March 19, 2018 at 6:03 am #

REPLY ↗

I have not noticed that. Perhaps post a ticket on the xgboost user group or on the project? Sounds like a fault?

Never miss a tutorial:

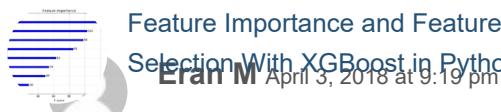
Nick March 31, 2018 at 4:01 am #

REPLY ↗



It should be `model.feature_importances`, not `model.get_importances_`.

Picked for you:



REPLY ↗

Better importance estimation:

[How to Develop Your First XGBoost Model in Python](#)
Model feature_importances_ uses the
Booster.get_fscore() which uses

`Booster.get_score(importance_type='weight')`



which is an estimation to "gain" (as of how many times all trees represented a certain feature).
I think it would be better to use `Booster.get_score(importance_type='gain')` to get a more precise evaluation of how important a feature is.



REPLY ↗



Thanks for sharing.

Loving the Tutorials?



John Markson November 22, 2018 at 12:42 am #

REPLY ↗

The [XGBoost With Python](#) EBook is
where you'll find [the Really Good](#) stuff.

>> SEE WHAT'S INSIDE e posts. Regarding the feature importance in Xgboost (or more generally gradient boosting trees), how do you feel about the SHAP? I am not sure if you already had any post discussing SHAP, but it is definitely interesting to people who need gradient boosting tree models for feature selections.

Thanks!



Jason Brownlee November 22, 2018 at 6:25 am #

REPLY ↗

What is SHAP?



dasgupso May 24, 2018 at 10:01 pm #

REPLY ↗

Hi Jason

I need to know the feature importance calculations by different methods like "weight", "gain", or "cover" etc. in Xgboost.

Never miss a tutorial:

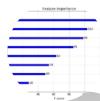
Jason Brownlee September 9, 2018 at 6:01 am #

REPLY ↗



Specifying how much time and resources you have and the goals of your project.

Perhaps a comparison of the same configuration of model with different input features would be a good step (w.g. without the grid search).



Feature Importance and Feature Selection With XGBoost in Python

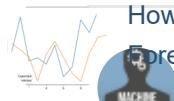
James September 26, 2018 at 12:33 pm #

REPLY ↗



How to Develop Your First XGBoost

Model in Python How can we use let's say top 10 features to train the model? I can not find a parameter to do so while initiating.



How to Use XGBoost for Time Series Forecasting

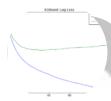
Jason Brownlee September 26, 2018 at 2:24 pm #

REPLY ↗



You must use feature selection methods to select the features you want to use. There is no best feature selection method, just different perspectives on what might be useful.

Data Preparation for Gradient Boosting with XGBoost in Python



Avoid Overfitting By Early Stopping With

Sinan Ozdemir September 28, 2018 at 6:56 am #

REPLY ↗

Hi Jason,

After reading your book, I was able to implement a model successfully. However, I have a few questions and I will appreciate if you can give me feedback:

Q1 – In the XGBoost with Python Tutorials we apply PCA (Principal Component Analysis), LDA (Linear Discriminant Analysis) or t-SNE when we use XGBOOST to determine the most important features?

>> SEE WHAT'S INSIDE

Q2 – Do we need to apply standard scaling after one hot encoding the categorical values? Again, some people say that this is not necessary in decision tree like models, but I would like to get your opinion.

Q3 – Do we need to be concerned with the dummy variable trap when we use XGBOOST? I couldn't find a good source about how XGBOOST handles the dummy variable trap meaning if it is necessary to drop a column.

As always I really appreciate your feedback.

Thank you.



Jason Brownlee September 28, 2018 at 2:58 pm #

REPLY ↗

You can try dimensionality reduction methods, it really depends on the dataset and the configuration of the model as to whether they will be beneficial.

No real need to rescale data for xgboost. Standardizing might be useful for Gaussian variables. Test and see.



especially if they expose a grouping of levels not obvious from the data (e.g. the addition of flag variables)

Picked for you:



 [Sitan Ondrejka](#) October 5, 2018 at 6:57 am #

REPLY ↗

[Selection With XGBoost in Python](#)

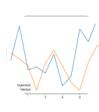
As always, thank you so much Jason.



For people who are interested in my experiment:

[How to Develop Your First XGBoost](#)

 [Michael Shor](#) Dimensionality reduction method didn't really help much. XGBOOST feature selection method was way better in my case.



Standardizing didn't really change neither the accuracy score or the predicting results.

[How to Use XGBoost for Time Series](#)

 [Forecasting](#)

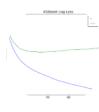
Keeping dummy variable increased the accuracy by about 2%, I used KFold to measure the accuracy.



 [Thanks](#)

[Data Preparation for Gradient Boosting with XGBoost in Python](#)

Thanks

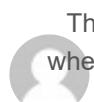


 [Jason Brownlee](#) Avoid Overfitting By Early Stopping With XGBoost in Python

REPLY ↗

Nice work!

Loving the Tutorials?



The [XGBoost With Python](#) EBook is

where you find the [Really Good Stuff](#) #

REPLY ↗

>> SEE WHAT'S INSIDE

Thanks for your post. It is really helpful.

But I am still confused about "Importance is calculated for a single decision tree by the amount that each attribute split point improves the performance measure".

How to calculate the "amount that each attribute split point improves the performance measure"?



 [Jason Brownlee](#) December 13, 2018 at 7:41 am #

REPLY ↗

This is calculated as part of constructing each individual tree. The final importance scores are an average of these scores.



 [Yang Song](#) December 27, 2018 at 7:09 pm #

REPLY ↗

Hi Jason, I have encountered a problem when I try to reimplement the python trained xgboost model by c++. I built the same decision trees as the python trained(use the 'model.dump_model' function) but I got the different scores. I didn't know why and can't figure that, can you give me several tips? thanks!



Picked for you:



Jason Brownlee December 28, 2018 at 5:54 am #

REPLY ↗



Feature Importance and Feature

Selection With XGBoost in Python
Perhaps there was a difference in your implementation? It could be one of a million things – impossible for me to diagnose sorry.



How to Develop Your First XGBoost Model in Python

Kamil February 12, 2019 at 7:14 am #

REPLY ↗



Hi How to Use XGBoost for Time Series Forecasting
n I'm running this code:



plot_importance(model)

plot_importance? preparation for Gradient Boosting

with XGBoost in Python

getting an error:



ValueError: tree must be Booster, XGBModel or dict instance

Avoid Overfitting By Early Stopping With
can I deal with that?

XGBoost In Python



Loving the Tutorials?

February 12, 2019 at 8:09 am #

REPLY ↗

The XGBoost With Python Error

I have not seen this error, I have some suggestions here:

where you'll find the Really Good stuff
<https://machinelearningmastery.com/faq/single-faq/why-does-the-code-in-the-tutorial-not-work-for-me>

>> SEE WHAT'S INSIDE



AAV March 16, 2019 at 9:16 am #

REPLY ↗

Is there any way to get sign of the features to understand if the impact is positive or negative.



Jason Brownlee March 17, 2019 at 6:14 am #

REPLY ↗

Not as far as I know, sorry.



Charles Brauer March 22, 2019 at 4:17 am #

REPLY ↗

When I click on the link: "names in the problem description" I get a 404 error.
Never miss a tutorial:
The "f1, f2..." names are not useful. I want the real column names.



 **Jason Brownlee** March 22, 2019 at 8:37 am #

REPLY ↗

Thanks, I have updated the link to:

<https://github.com/brownlee/Datasets/blob/master/pima-indians-diabetes.names>

 **Abhinav** May 7, 2019 at 7:42 pm #

REPLY ↗

Hi Jason,


Thanks again for an awesome post. Just like there are some tips which we keep in mind while feature selection using Random Forest.

Like – The categorical variable with high cardinality/ continuous variable are given preference over others

correlation is not visible in case of RF feature importance.

Do XGBoost have similar cons similar to Random Forest??

Avoid Overfitting By Early Stopping With
XGBoost In Python



Jason Brownlee May 8, 2019 at 6:43 am #

REPLY ↗

Loving the Tutorials?

Yes, perhaps this post will help:
<https://machinelearningmastery.com/configure-gradient-boosting-algorithm/>
The XGBoost With Python EBook is
where you'll find the **Really Good** stuff.



>> SEE WHAT'S INSIDE

 **Starting** May 14, 2019 at 2:10 am #

REPLY ↗

Hello!

Given feature importance is a very interesting property, I wanted to ask if this is a feature that can be found in other models, like Linear regression (along with its regularized partners), in Support Vector Regressors or Neural Networks, or if it is a concept solely defined solely for tree-based models. I ask because I am not sure whether I can consider eg Linear Regression's coefficients as the analog for feature importance.

Thanks!



Jason Brownlee May 14, 2019 at 7:49 am #

REPLY ↗

Yes, coefficient size in linear regression can be a sign of importance.

SVM, less so.
Never miss a tutorial:



Constantine

May 16, 2019 at 4:39 am #

REPLY ↗

Picked for you: thanks!



Feature Importance and Feature

Selection With XGBoost in Python



Jason Brownlee May 16, 2019 at 6:36 am #

REPLY ↗



How to Develop Your First XGBoost
Model in Python

You're welcome!



How to Use XGBoost for Time Series
Forecasting

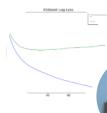
Hiro May 19, 2019 at 1:38 am #

REPLY ↗

Thanks for all of your posts. I use your blog to study a lot.



I have a question. If you had a large number of features, do you want to use all of them? I have a set with over 1,000 features but not all of them are meaningful for this classification problem I am working on. Should I reduce the number of features before applying XGBoost? If so, how can I do so?



Avoid Overfitting By Early Stopping With

XGBoost In Python

Jason Brownlee May 19, 2019 at 8:05 am #

REPLY ↗

Try modeling with all features and compare results to models fit on subsets of selected features to see if it improves performance.

Loving the Tutorials?

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.



Jonathan May 21, 2019 at 4:08 am #

REPLY ↗

>> SEE WHAT'S INSIDE

Hi Jason,

Does multicollinearity affect feature importance for boosted regression trees? If so, how would you suggest to treat this problem?

Thanks!



Jason Brownlee May 21, 2019 at 6:40 am #

REPLY ↗

Probably not.

Try modeling with an without the colinear features and compare results.



Grzegorz Kępisty July 16, 2019 at 5:03 pm #

REPLY ↗

Hello Jason,

Never miss a tutorial:

Concerning default feature importance in similar method from sklearn (Random Forest) I recommend



<https://explained.ai/rf-importance/>

The authors show that the default feature importance implementation using Gini is biased.

Picked for you:

I observed this kind of bias several times, that is overestimation of importance of artificial random
variables and importance of Feature Selection With XGBoost in Python

However, there are other methods like “drop-col importance” (described in same source). Interestingly,
while working with production data, I observed that some variables occur in head of sorted distribution or
tail – depending which method of 2 above I applied.



[How to Develop Your First XGBoost Model In Python](#)

This is somehow confusing and now I am cautious in using RF for feature selection.

Do you have some experience in this field or some best practices to share?

[How to Use XGBoost for Time Series Forecasting](#)
regards!

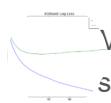


Data Preparation for Gradient Boosting

[With XGBoost](#) Jason Brownlee July 17, 2019 at 8:20 am #

REPLY ↗

Thanks for sharing.



[Avoid Overfitting By Early Stopping With](#)

XGBoost In Python My best advice is to use importance as a suggestion but remain skeptical. Test many methods, many subsets, make features earn the use in the model with hard evidence.



Loving the Tutorials?

[new_to_modelling](#) July 17, 2019 at 1:08 am #

REPLY ↗

The XGBoost With Python EBook is

where i want to predict one of those columns so remaining 5. Out of which 2 are categorical variable and 3 are numerical variable. So, i used https://scikit-learn.org/stable/modules/plot_column_transformer_mixed_types.html to workout a mixed data type issues. But when i the feature_importance size does not match with the original number of columns? The size of feature_importances_ array is 918

I mean its generating extra feature or is it creating a feature value for one_hot_encoding of the categorical variable.

I checked my data has 1665 unique brand values. So, its not the same as feature_importances_ array size



[Jason Brownlee](#) July 17, 2019 at 8:28 am #

REPLY ↗

Performing feature selection on the categorical data might be confusing as it is probably one hot encoded.

Perhaps create a subset of the data with just the numerical features and perform feature selection
Never miss a tutorial:
on that?



 **new_to_modelling** July 17, 2019 at 6:45 pm #
Picked for you:

REPLY ↗



Thanks, but i found it was working once i tried dummies in place of the above
Feature Importance and Feature Selection With XGBoost in Python
mentioned column transformer approach seems like during transformation there is some loss of information when the xgboost booster picks up the feature names

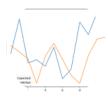


How to Develop Your First XGBoost Model in Python



Jason Brownlee July 18, 2019 at 8:24 am #

REPLY ↗



Glad to hear it.
How to Use XGBoost for Time Series Forecasting



Mike Sishi July 28, 2019 at 4:22 am #
Data Preparation for Gradient Boosting with XGBoost in Python

Interesting article, thanks a lot!!

REPLY ↗

Ho can I reverse-engineer a Decision Tree? That is, change the target variable and consequently have feature A variable adjust themselves

 **Avoid Overfitting By Early Stopping With XGBoost In Python**

Actually, I want to set a target variable value and get all possible values of feature variables that can yield the target variable value.

Loving the Tutorials?

 **Jason Brownlee** July 28, 2019 at 6:49 am #
where you'll find the **Really Good** stuff.

REPLY ↗

Reverse ML/predictive modeling is very hard if not entirely intractable.

>> SEE WHAT'S INSIDE
You can do this and do this and do this and give many "results".

It would might not make sense for an ensemble of trees.



Mike Sishi July 28, 2019 at 6:27 pm #

REPLY ↗

Hi Jason,

Thanks for your prompt response. I will try to work on the solution and let you know how it goes.

Kind Regards



Jason Brownlee July 29, 2019 at 6:11 am #

REPLY ↗

Good luck!
Never miss a tutorial:



Abdoul August 17, 2019 at 3:00 am #

REPLY ↗

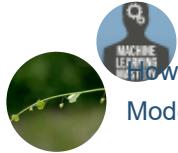
Picked for you:

How to extract the n best attributes at the end?



Feature Importance and Feature

Selection With XGBoost in Python



Jason Brownlee August 17, 2019 at 5:58 am #

REPLY ↗

How to Develop Your First XGBoost

Model in Python I give an example in the above tutorial.



How to Use XGBoost for Time Series

Forecasting **Robert Feyerharm** August 28, 2019 at 11:49 pm #

REPLY ↗

Thanks Jason, very helpful!



Is there a way to determine if a feature has a net positive or negative correlation with the outcome variable? Data Preparation for Gradient Boosting with XGBoost in Python



Avoid Overfitting By Early Stopping With

XGBoost In Python **Jason Brownlee** August 29, 2019 at 6:12 am #

REPLY ↗

Yes, you can calculate the correlation between them.

Loving the Tutorials?

The XGBoost With Python EBook is
where you can find the really good stuff! **Roger** August 29, 2019 at 10:10 am #

REPLY ↗

>> SEE WHAT'S INSIDE

Thresh=0.030, n=10, precision: 40.81%
Thresh=0.031, n=9, precision: 50.00%
Thresh=0.032, n=8, precision: 47.83%
Thresh=0.033, n=7, precision: 51.11%
Thresh=0.035, n=6, precision: 48.78%
Thresh=0.041, n=5, precision: 41.86%
Thresh=0.042, n=4, precision: 58.62%
Thresh=0.043, n=3, precision: 68.97%
Thresh=0.045, n=2, precision: 62.96%
Thresh=0.059, n=1, precision: 0.00%

Hi Jason, Thank you for your post, and I am so happy to read this kind of useful ML articles. I have a question: the above output is from my example. As you can see, when thresh = 0.043 and n = 3, the precision dramatically goes up. So, I want to take a closer look at that thresh and wants to find out the names and corresponding feature importances of those 3 features. How can I achieve this goal?

Never miss a tutorial:

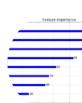
Jason Brownlee September 6, 2019 at 1:57 pm #

REPLY ↗



In feature have unique index of the column in the dataset from 0 to n. If you know the names of the columns, you can map the column index to names.

Picked for you: [How to do this in Python to automate it.](#)



I hope that helps
Feature Importance and Feature Selection With XGBoost in Python

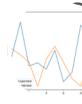


Ralph September 6, 2019 at 1:57 pm #

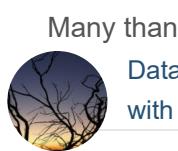
REPLY ↗

Model in Python

Hi! I am using instead the xgb.train command instead of XGBClassifier because this is much faster. By the way you have any idea why, and if it possible to obtain the same performance with XGBClassifier (might be related to the number of threads)?

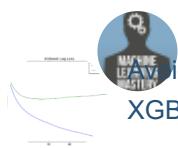


[How to Use XGBoost for Time Series Forecasting](#)
way, you have any idea of how to get importance feature with xgb.train?



Many thanks

Data Preparation for Gradient Boosting with XGBoost in Python



Jason Brownlee September 22, 2019 at 9:27 am #

REPLY ↗

Avoid Overfitting By Early Stopping With XGBoost in Python

Are they any Python? Any Python? It is faster? It should be identical in speed.



Loving the Tutorials?

abstract September 23, 2019 at 1:56 pm #

REPLY ↗

The [XGBoost With Python](#) EBook is Any reason why the Accuracy has increased from 76.38 at n=7 to 77.56 at n=6 ? where you'll find the **Really Good** stuff.

—> SEE WHAT'S INSIDE



Jason Brownlee September 26, 2019 at 6:27 am #

REPLY ↗

Perhaps the change in inputs or perhaps the stochastic nature of the learning algorithm.

A fair comparison would use repeated k-fold cross validation and perhaps a significance test.



Maria September 27, 2019 at 11:50 pm #

REPLY ↗

Hello Jason,

I work on an imbalanced dataset for anomaly detection in machines. I have 590 features and 1567 observations. I tried this approach for reducing the number of features since I noticed there was multicollinearity, however, there is no important shift in the results for my precision and recall and sometimes the results get really weird. I was wondering what could that be an indication of? Here are the results of the features selection

Thresh=0.000, n=211, f1_score: 5.71%
Never miss a tutorial:
precision_score: 50.00%



Thresh=0.000, n=210, f1_score: 5.71%
precision_score: 50.00%

Picked for you:
recall_score: 3.03%

 Feature importance and Feature Selection With XGBoost in Python
precision_score: 50.00%

 How to Develop Your First XGBoost Model in Python
accuracy_score: 91.22%
Thresh=0.000, n=208, f1_score: 5.71%
precision_score: 50.00%

 How to Use XGBoost for Time Series Forecasting
Thresh=0.000, n=207, f1_score: 5.71%
precision_score: 50.00%

 Data Preparation for Gradient Boosting With XGBoost in Python
accuracy_score: 91.22%
Thresh=0.000, n=206, f1_score: 5.71%
precision_score: 50.00%

 Avoid Overfitting By Early Stopping With XGBoost in Python
Thresh=0.006, n=55, f1_score: 11.11%
precision_score: 66.67%

recall_score: 6.06%
accuracy_score: 91.49%
Thresh=0.000, n=34, f1_score: 5.88%
precision_score: 100.00%
The XGBoost With Python EBook is recall_score: 3.03%
where you'll find the **Really Good** stuff.
accuracy_score: 91.49%

Thresh >> SEE WHAT'S INSIDE 8%
precision_score: 100.00%
recall_score: 3.03%

accuracy_score: 91.49%
Thresh=0.007, n=52, f1_score: 5.88%
precision_score: 100.00%
recall_score: 3.03%
accuracy_score: 91.49%

Thresh=0.007, n=47, f1_score: 0.00%
precision_score: 0.00%
recall_score: 0.00%
accuracy_score: 91.22%

UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.
'precision', 'predicted', average, warn_for)

Precision is ill-defined and being set to 0.0 due to no predicted samples.
Never miss a tutorial:
(precision', 'predicted', average, warn_for)



 **Jason Brownlee** September 28, 2019 at 6:20 am # REPLY ↗
Picked for you:

Interesting, I'm not sure.
 **Feature Importance and Feature Selection With XGBoost In Python**
You may need to dig into the specifics of the data to what is going on. If you're using CV, then perhaps some folds don't have examples of the target class – use stratified CV.

 **How to Develop Your First XGBoost Model in Python**
 **Maria** September 28, 2019 at 8:41 pm # REPLY ↗

 **How to Use XGBoost for Time Series Forecasting**
You're right, I have 104 examples of the minority class and 1463 of the other one. Would you suggest oversampling in that case?

 **Data Preparation for Gradient Boosting with XGBoost in Python**
 **Jason Brownlee** September 29, 2019 at 6:11 am # REPLY ↗

Perhaps test it and see.
 **Avoid Overfitting By Early Stopping With More ideas here: XGBoost In Python**
<https://machinelearningmastery.com/tactics-to-combat-imbalanced-classes-in-your-machine-learning-dataset/>

Loving the Tutorials?

 The **XGBoost With Python EBook** is November 1, 2019 at 6:51 pm # REPLY ↗
where you'll find the **Really Good** stuff.

Hi
>> SEE WHAT'S INSIDE
I have a question: how can we know the names of the features that are selected in: model using each importance as a threshold.

 **Jason Brownlee** November 2, 2019 at 6:41 am # REPLY ↗

Each column in the array of loaded data will map to the column in your raw data.
If you know column names in the raw data, you can figure out the names of columns in your loaded data, model, or visualization.

 **krs reddy** November 22, 2019 at 12:46 am # REPLY ↗

Can you try plotting model interpretation using shap library for tree based algorithms??

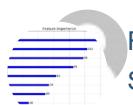
Never miss a tutorial:

Jason Brownlee November 22, 2019 at 6:05 am #

REPLY ↗



Picked for you:

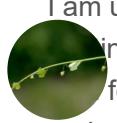


Kamal December 16, 2019 at 6:39 pm #

REPLY ↗

Selection With XGBoost in Python

Variable of Importance in Xgboost for multilinear features –



I am using 60 obseravation*90features data (all continuous variables) and the response variable is also

How to Develop Your First XGBoost

feature importance in python(xgb.feature_importances_), that sums up 1. I run xgboost 100 times

and select features based on the rank of mean variable importance in 100 runs. Let's say I choose 10

factors and then, again run xgboost with the same hyperparameters on these 10 features, surprisingly

How to Use XGBoost for Time Series

most important feature becomes least important in these 10 variables. Any feasible explanation for
Forecasting ?



Data Preparation for Gradient Boosting

with XGBoost in Python

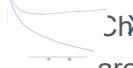
Jason Brownlee December 17, 2019 at 6:31 am #

REPLY ↗



Not off hand, does it matter though?

Avoid Overfitting By Early Stopping With



XGBoost In Python features that gives the best results/most skillful model – any importance scores
are a “suggestion” at best.

Loving the Tutorials?



Sathya Bhat January 25, 2020 at 3:48 am #

REPLY ↗

The XGBoost With Python EBook is

where you'll find the **Really Good** stuff.

In a XGBoost model, the top features we derive shows which feature is more influential than the
rest.

>> SEE WHAT'S INSIDE

For example if the top feature is tenure days, how do i determine if “more tenure days” or “less tenure
days” increase the rating in the output..

How do I determine if it is a positive influence or negative influence?



Jason Brownlee January 25, 2020 at 8:40 am #

REPLY ↗

Not sure I follow. Importance is not positive or negative.

If you're in doubt: build a model with and without it and compare the performance of the model.



Sahil Basera February 11, 2020 at 7:28 pm #

REPLY ↗

I don't understand the F -score in the feature importance plot, who can the value be 100+. Also,
if this is not the traditional F-score, could you point to the definition/explanation of it? (can't find it in the

xgb documentation)
Never miss a tutorial:

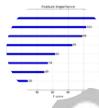


Jason Brownlee

February 12, 2020 at 5:44 am #

REPLY ↗

Picked for you The scores are relative.



Feature Importance and Feature Selection With XGBoost in Python

Shreya February 26, 2020 at 7:34 am #

REPLY ↗



How to Develop Your First XGBoost Model in Python
Is it possible to plot important features on model ensembled using Voting Classifier ?
What then what could be the alternative to plot important features in an ensembled technique ?

Thank You



How to Use XGBoost for Time Series Forecasting

Jason Brownlee February 26, 2020 at 8:30 am #

REPLY ↗



Data Preparation for Gradient Boosting
Not really with XGBoost in Python

Most ensembles of decision trees can give you feature importance.



Avoid Overfitting By Early Stopping With XGBoost In Python



Shreya February 27, 2020 at 12:57 am #

REPLY ↗

Thank You !
Loving the Tutorials?

But what about ensemble using Voting Classifier consisting of Random Forest, Decision Tree,
The [XGBoost With Python EBook](#) is
XGBoost and Logistic Regression ?
where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE



Jason Brownlee February 27, 2020 at 5:55 am #

REPLY ↗

Voting ensemble does not offer a way to get importance scores (as far as I know),
regardless of what is being combined.



Shreya March 1, 2020 at 2:29 am #

ok.

Thanks a lot !!



Jason Brownlee March 1, 2020 at 5:25 am #

You're welcome.
Never miss a tutorial:



Daniel Madsen March 23, 2020 at 9:45 pm #

REPLY ↗

Picked for you:

Hi and thanks for the codes first of all.



In I run the "select_X_train = selection.transform(X_train)" I receive the following error: "ValueError: It contains NaN, infinity or a value too large for dtype('float64')."

My features do contain some NaNs, dummy variables and categorial variables. Do you know any way around this without having to change my data?

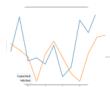


How to Develop Your First XGBoost Model in Python

Model in Python

works in advance,

Daniel



How to Use XGBoost for Time Series

Forecasting



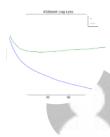
Jason Brownlee March 24, 2020 at 6:02 am #

REPLY ↗

Data Preparation for Gradient Boosting

You will need to impute the nan values first, or remove rows with nan values:
with XGBoost in Python

<https://machinelearningmastery.com/handle-missing-data-python/>



Avoid Overfitting By Early Stopping With

XGBoost In Python

Joseph Hall April 5, 2020 at 5:10 pm #

REPLY ↗

Hi,

Loving the Tutorials?

I got an error at this line of the code: "select_X_train = selection.transform(X_train)". The error is simply "KeyError: weight"

The [XGBoost With Python](#) EBook is

where you'll find the [Really Good stuff](#).
I did some research and found out that SelectFromModel expects an estimator having coef_ or feature importances. Obviously XGBoostClassifier does have this attribute. Why is it not working for me but >> SEE WHAT'S INSIDE

Please help!



Jason Brownlee April 6, 2020 at 6:02 am #

REPLY ↗

Perhaps confirm that your version of xgboost is up to date?



Gustavo April 7, 2020 at 4:04 am #

REPLY ↗

I am having this same error. I am with xgboost 1.0.2 installed through pip.

d8veone April 7, 2020 at 12:50 am #

REPLY ↗

it works in xgboost 0.90, but not 1.0.2
Never miss a tutorial:



Jason Brownlee April 7, 2020 at 5:51 am #

REPLY ↗

Picked for you:

Thanks, I will investigate!



Feature Importance and Feature Selection With XGBoost in Python

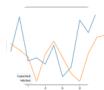


Jason Brownlee April 8, 2020 at 10:12 am #

REPLY ↗



How to Develop Your First XGBoost Model in Python
I have added a work around.



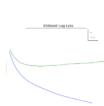
How to Use XGBoost for Time Series Forecasting
d8veone April 9, 2020 at 12:28 am #

REPLY ↗

Awesome! you're a true master. Thank you.



Data Preparation for Gradient Boosting with XGBoost in Python



Jason Brownlee April 9, 2020 at 8:04 am #

Avoid Overfitting By Early Stopping With XGBoost In Python
Thanks for your kind words.

REPLY ↗



Loving the Tutorials? 19 am #

The XGBoost With Python Ebook is Please, remove my last post... xgboost 0.90 worked where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE



Jason Brownlee April 7, 2020 at 1:30 pm #

REPLY ↗

Done.



John April 14, 2020 at 8:43 am #

REPLY ↗

Hi,

I would like to use the “Feature Selection with XGBoost Feature Importance Scores” approach with model selection in my research. How can I cite it in paper/thesis?
Thank you



Jason Brownlee April 14, 2020 at 10:39 am #

REPLY ↗

This will help:
Never miss a tutorial:

<https://machinelearningmastery.com/faq/single-faq/how-do-i-reference-or-cite-a-book-or-blog-post>



 **Mohie** April 26, 2020 at 2:22 pm #

REPLY ↗

my xgb model is taking too long for one fit and i want to try many thresholds so can i use
Feature Importance and Feature
Selection With XGBoost In Python
another simple model to know the best threshold and is yes what do you recommend ?

 **How to Develop Your First XGBoost Model**  **Jason Brownlee** April 27, 2020 at 5:28 am #

REPLY ↗

You can try, but the threshold should be calculated for the specific model.

 **How to Use XGBoost for Time Series Forecasting** 

REPLY ↗

 **kim** May 6, 2020 at 12:35 pm #

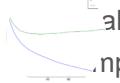
Data Preparation for Gradient Boosting

Hello Sir

with XGBoost in Python

I tried to run print(model.feature_importances_)

but it give an array with all 'nan' like [nan nan nan ... nan nan nan]

 **Avoid Overfitting By Early Stopping With XGBoost In Python** 
also, when i tried to plot the model with plot_importance(model), it return Booster.get_score() results npty

do you have any advice? thank you very much

Loving the Tutorials?

The XGBoost With Python EBook is
 **Jason Brownlee** May 6, 2020 at 1:38 pm #
where you'll find the *Really Good* stuff.

REPLY ↗

>> SEE WHAT'S INSIDE check that you fit the model?

 **kim** May 6, 2020 at 2:11 pm #

REPLY ↗

yes it return like this

```
XGBClassifier(base_score=0.5, booster=None, colsample_bylevel=1,  
colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,  
importance_type='gain', interaction_constraints=None,  
learning_rate=0.300000012, max_delta_step=0, max_depth=6,  
min_child_weight=1, missing=nan, monotone_constraints=None,  
n_estimators=100, n_jobs=0, num_parallel_tree=1,  
objective='binary:logistic', random_state=0, reg_alpha=0,  
reg_lambda=1, scale_pos_weight=1, subsample=1, tree_method=None,  
validate_parameters=False, verbosity=None)
```

Never miss a tutorial:



Jason Brownlee May 7, 2020 at 6:37 am #

REPLY ↗



I'm not sure off the cuff, you might have to try varying the training data and review the

Picked for you:



Feature Importance and Feature Selection With XGBoost in Python

Apo May 12, 2020 at 3:56 am #

REPLY ↗



How to Develop Your First XGBoost Model in Python

I'm testing your idea with feature importance of XGB and thresholds in a problem that I survey these days. I'm dealing with some weird results and I wonder if you could help.



Hi, run a part of code similar to yours to see different metrics results on each threshold (beginning all features to end up with 1). After that I check these metrics and note the best outcomes and the number of features resulting in these (best) metrics. Finally, I'm taking these features and use XGB algorithm with only these features but this time the results are different with results I got in the previous Data Preparation for Gradient Boosting with XGBoost in Python

Any good explanation of this side effect?

Thanks for your time



Avoid Overfitting By Early Stopping With XGBoost In Python



Jason Brownlee May 12, 2020 at 6:51 am #

REPLY ↗

Perhaps the difference in results is due to the stochastic nature of the learning algorithm or test harness.

Loving the Tutorials?

The [XGBoost With Python](#) Ebook is perhaps design a robust test harness and perform feature selection within the modeling pipeline. where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE



babak June 6, 2020 at 5:33 pm #

REPLY ↗

dear Jason

thank you for your program

I have 2 questions

1)if my target data are not categorical or binary for example so as Boston housing price has many price target so I encoding the price first before feature selection?

2) does the feature selection and correlation must have the same results?

thank first for your time



Jason Brownlee June 7, 2020 at 6:20 am #

REPLY ↗

No, that is a regression problem:

<https://machinelearningmastery.com/faq/single-faq/what-is-the-difference-between-classification-and-regression>

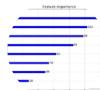
No, each technique will give you a different idea of what features may be important. You can fit a model from each suggestion and discover what actually results in a skillful model.



 **babak** June 7, 2020 at 8:57 pm #

REPLY ↗

Picked for you:



thank you for your answer

[Feature Importance and Feature](#)

[Selection With XGBoost in Python](#)



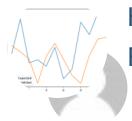
Jason Brownlee June 8, 2020 at 6:10 am #

REPLY ↗

[How to Develop Your First XGBoost](#)

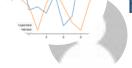
[Model in Python](#)

You're welcome.



[How to Use XGBoost for Time Series](#)

[Forecasting](#)

 **alireza** July 1, 2020 at 4:30 am #

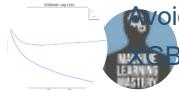
REPLY ↗

hi.



[Data Preparation for Gradient Boosting](#)

[Get Effect \(percentage\) of the Input Variables on the Output Variable with XGBoost in Python](#)



[Avoid Overfitting By Early Stopping With](#)

Jason Brownlee July 1, 2020 at 5:56 am #

REPLY ↗

One approach would be to convert each score to a ratio of the sum of the scores.

Loving the Tutorials?



James Hutton July 12, 2020 at 8:43 am #

REPLY ↗

Where you'll find the **Really Good** stuff.

Hi Jason,

Thank >> SEE WHAT'S INSIDE trial on this – I learn a lot.

I have one question, in the Feature Selection with XGBoost Feature Importance Scores section, you used

`thresholds = sort(model.feature_importances_)`

and do the for loop along these threshold values to evaluate the possible models.

Is there anyway how to do similar by using the values from `plot_importance()` results as the thresholds?

I looked at the data type from `plot_importance()` return, it is a matplotlib object instead of an array like the ones from `model.feature_importances_`

Can you please guide me on how to implement this?

Thank you



Jason Brownlee July 12, 2020 at 11:28 am #

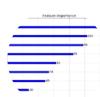
REPLY ↗

Good question James, yes there must be, but I'm not sure off hand. Perhaps check the xgboost library API for the appropriate function?



Picked for you: James Hutton July 12, 2020 at 8:27 pm #

REPLY ↗



Alright Feature Importance and Feature

Selection With XGBoost in Python

One more thing, In the results of different thresholds and respective different n number of features, how to pull in which features are in each scenario of threshold or in this n number of features? Means, which features are they? Can you show perhaps?



How to Develop Your First XGBoost

Model in Python



How to Use XGBoost for Time Series

Forecasting

Sorry, I don't follow your questions. Can you please restate or elaborate?

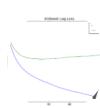


Data Preparation for Gradient Boosting

with XGBoost in Python

Julie July 30, 2020 at 1:27 am #

REPLY ↗



Jason,

Avoid Overfitting By Early Stopping With XGBoost In Python

way of explaining is very simple and straitprint(classification_report(y_test, predicted_xgb))ght ward. Please keep doing this!!!

As for this subject, I've done both manual feature importance and xgboost built-in one but got different rankings. I'm not sure why ??

Loving the Tutorials?

The XGBoost With Python EBook is

where you'll find the **Really Good stuff**

Jason Brownlee July 30, 2020 at 6:24 am #

REPLY ↗

>> SEE WHAT'S INSIDE

I believe they use a different evaluation function for the plot vs automatic.



Antonio September 29, 2020 at 1:12 am #

REPLY ↗

Amazing job Jason, Very helpful...!

If I may ask about the difference between the two ways of calculating feature importance, as I'm having contradictory results and non-matching numbers.

Exp: first way is giving output in [0,1], and the second way is giving results >1, can you explain the difference please 😊

In addition to that, if we take feature importance as ranking and setting apart the different scale issue between the two approaches, I encountered contradictory results where the number 1 important feature in the first method isn't the number 1 in the second method.

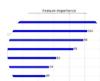
Never miss a tutorial:

Jason Brownlee September 29, 2020 at 5:42 am #

REPLY ↗



I believe the built-in method uses a different scoring system, you can change it to be consistent with
Picked for you: [any argument to the function.](#)



[Feature Importance and Feature Selection With XGBoost in Python](#)

Jacob November 17, 2020 at 5:45 am #

REPLY ↗



[How to Develop Your First XGBoost Model in Python](#)



looks like the feature importance results from the `model.feature_importances_` and the built in `xgboost.plot_importance` are different if your sort the importance weight for `model.feature_importances_`.
link [How to Use XGBoost for Time Series Forecasting](#) to determine the importance as xgboost use fs score to generate feature importance plots.



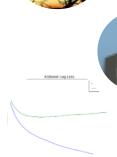
-Jacob



[Data Preparation for Gradient Boosting with XGBoost in Python](#)

Jason Brownlee November 17, 2020 at 6:34 am #

REPLY ↗



[Avoid Overfitting By Early Stopping With XGBoost In Python](#)

Thanks for sharing.

XGBoost In Python



Sarah November 18, 2020 at 2:20 am #

REPLY ↗

The [XGBoost In Python](#) difference between feature importance and feature selection methods?
where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Jason Brownlee November 18, 2020 at 6:45 am #

REPLY ↗

Good question, I answer it here:

<https://machinelearningmastery.com/faq/single-faq/what-is-the-difference-between-feature-selection-and-feature-importance>



Sarah November 19, 2020 at 2:49 am #

REPLY ↗

Great explanation, thanks. So, when we run feature selection should we expect the most important variables to be selected?

Jason Brownlee November 19, 2020 at 7:47 am #

REPLY ↗

Thanks!



XGBoost performs feature selection automatically as part of fitting the model.
Never miss a tutorial:



Saran November 19, 2020 at 10:30 am #

REPLY ↗

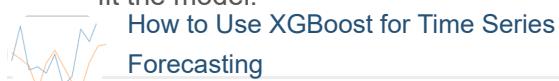
Picked for you: thank you so much for the clarification about the XG-Boost. In the XGBoost, I used `xgb.plot_importance` which plots all the features by their F score. How I can plot the selected features which are used as part of fitting the model.?
[Feature Importance and Feature Selection With XGBoost in Python](#)



How to Develop Your First XGBoost Model in Python Jason Brownlee November 19, 2020 at 1:38 pm #

REPLY ↗

The importance score itself is a reflection of the degree to which the features were used to fit the model.



[How to Use XGBoost for Time Series Forecasting](#)



Diana November 27, 2020 at 3:38 am #
[Data Preparation for Gradient Boosting with XGBoost in Python](#)

REPLY ↗

Hi,
Thank you for your wisdom.

[Using Python and the Recursive Feature Elimination \(RFE\) to select features](#). I'm trying different types of models such as XGBClassifier, Decision Trees, or KNN.

However, the RFE gives me the following error when the model is XGBClassifier or KNN. The KNN does not provide logic to do feature selection, but the XGBClassifier does.

I'm doing something wrong? Is there an explanation for this error with XGBClassifier?

Error: The classifier does not expose "coef_" or "feature_importances_" attributes
[The XGBoost With Python EBook](#) is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE



Jason Brownlee November 27, 2020 at 6:44 am #

REPLY ↗

You're welcome.

Interesting. Perhaps check of your xgboost library is up to date?

Otherwise, perhaps xgboost cannot be used in this way – which is a shame.

Jean-Marc December 7, 2020 at 10:28 pm #

REPLY ↗

Hello,

I am getting an empty `select_X_train` when using the smallest threshold (So normally I will get the same for all other thresholds). Can someone please help me find out why?

Regards!

Never miss a tutorial:

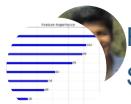
Jason Brownlee December 8, 2020 at 7:43 am #

REPLY ↗



, if the threshold is too low, you will not select any features. Increase it.

Picked for you:



mhr007 December 11, 2020 at 7:26 pm #

REPLY ↗

Selection With XGBoost in Python

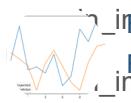
Why are you using 'SelectFromModel' here ?

can't we just do something like this ?



How to Develop Your First XGBoost Model in Python

```
model2 = xgb.XGBRegressor(**tuned_params)
model2.fit(X_imp_train,y_train,eval_set = [(X_imp_train,y_train),
(X_imp_test,y_test)],verbose=False)
```



Gain_Importance_XGBoost For Time Series

```
model2.get_booster().get_score(importance_type='gain')
gain_importance_dict2temp = sorted(gain_importance_dict2temp.items(), key=lambda x: x[1],
reverse=True)
```



Data Preparation for Gradient Boosting

Gain_Importance_in_Python(gain_importance_dict2temp)

```
temmae = 10000.0
```



Avoid Overfitting By Early Stopping With

XGBoost In Python

```
list_of_feature = [x for x,y in gain_importance_dict2temp[:feature_importance_len-i]]
print(list_of_feature)
```

X_imp_train3=X_imp_train[list_of_feature]

X_imp_test3=X_imp_test[list_of_feature]

The XGBoost With Python EBook is

where you'll find the XGBoostessor(**tuned_params)

```
regression_model.fit(X_imp_train3,y_train,eval_set = [(X_imp_train3,y_train),
(X_imp_test3,y_test)],early_stopping_rounds=se)
```

```
yPred= regression_model.predict(X_imp_test3)
```



Jason Brownlee December 12, 2020 at 6:24 am #

REPLY ↗

We are using select from model because the xgboost model has feature importance scores.

You have implemented essentially what the select from model does automatically.



mhr007 December 14, 2020 at 8:41 pm #

REPLY ↗

Thanks, you are so great, I didn't expect an answer from you for small things like this.
Thanks.

Never miss a tutorial:

Jason Brownlee December 15, 2020 at 6:19 am #

REPLY ↗



Picked for you:



Eric January 22, 2021 at 4:11 am #

REPLY ↗

Feature Importance and Feature Selection With XGBoost in Python

Followed exact same code but got "ValueError: X has a different shape than during fitting." in

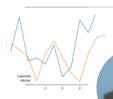
line "select_x_train = selection.transform(x_train)" after projecting the first few lines of results of the features selection.



How to Develop Your First XGBoost

Model in Python

Please help, many thanks



How to Use XGBoost for Time Series Forecasting

Jason Brownlee January 22, 2021 at 7:25 am #

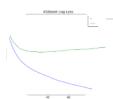
REPLY ↗

Sorry to hear that, perhaps these tips will help:



Data Preparation for Gradient Boosting

<https://machinelearningmastery.com/faq/single-faq/why-does-the-code-in-the-tutorial-not-work-for-me-with-xgboost-in-python>



Avoid Overfitting By Early Stopping With XGBoost In Python

Eric January 22, 2021 at 8:06 pm #

REPLY ↗

thanks for the reply. I add the np.sort of the threshold and problem solved

Loving the Tutorials?
Learn more about feature_importances_)

The [XGBoost With Python](#) EBook is where you'll find the **Really Good** stuff.



>> SEE WHAT'S INSIDE

Jason Brownlee January 23, 2021 at 7:02 am #

REPLY ↗

Happy to hear that Eric!



manuela April 3, 2021 at 8:54 pm #

REPLY ↗

Hi jason, I have used a standard version of Algorithm A which has features x, y, and z then I had used feature engineering to add for algorithm A new features (10 new features)

I would like to use the feature importance method to select the most important features between only the 10 features without removing any of the (x, y, z features)

can I identify first the list of features on which I would like to apply the feature importance method??



Jason Brownlee April 4, 2021 at 6:51 am #

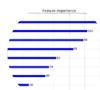
REPLY ↗

Sure. You can use any features you like, e.g. a combination of those selected by an algorithm and those you select.



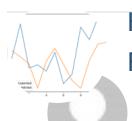
Picked for you: manuela April 4, 2021 at 3:43 pm #

REPLY ↗

 thank you for your replay
Feature Importance and Feature Selection With XGBoost in Python

 Jason Brownlee April 15, 2021 at 6:09 am #
How to Develop an XGBoost Model in Python
You're welcome.

REPLY ↗

 How to Use XGBoost for Time Series Forecasting
Amr J April 29, 2021 at 5:21 pm #

REPLY ↗

 Thank you for the article
Data Preparation for Gradient Boosting with XGBoost in Python
I decided to read in the pima Indian data using DF and put inthe feature names so that I can see those when plotting the feature importance.

 We used the following to add the feature names to the scores of model.feature_importances_
Avoid Overfitting By Early Stopping With XGBoost in Python plot:

```
from pandas import DataFrame
cols=X.columns
new_df = DataFrame (cols)
importance = model.feature_importances_*100
The XGBoost With Python EBook is
importance = importance.round(2)
Where you'll find the Really Good stuff.
new_df2 = DataFrame (importance)
fi=pd.c   >> SEE WHAT'S INSIDE   s=1)
fi.columns=[feature , score ]
fi.set_index('Feature',inplace=True)
fi.sort_values(by='score', ascending=False, inplace=True)
```

When comparing this plot to the one produced by plot_importance(model), you will notice the two do not rank the features in the same order. Any idea why?

I have tried the same thing with the famous wine data and again the two plots gave different orders to the feature importance.

One interesting thing to note is that when using catboost (as compared to xgboost) and then using SHAP to understand the impact of the features on the model, the graph is very similar to the (model.feature_importances_) method.



Jason Brownlee April 30, 2021 at 6:02 am #

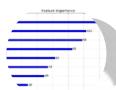
REPLY ↗

I believe that the plot_importance() uses a different metric for importance scores than feature_importances_.
Never miss a tutorial:



I believe I can config the plot function to use the same score to make the scores equivalent. I recommend checking the API.

Picked for you:



Feature Importance and Feature Selection With XGBoost in Python
Vasiliki Voukelatou September 4, 2021 at 12:36 am #

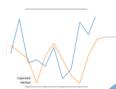
REPLY ↗

Hi Jason,



Thanks for the tutorial! Which is the default type for the feature_importances_ , i.e. weight, gain, etc? It is mentioned in the documentation.

Thank you in advance.

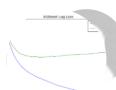


How to Use XGBoost for Time Series Forecasting
Jason Brownlee September 4, 2021 at 5:23 am #

REPLY ↗



Data Preparation for Gradient Boosting are not clear, I recommend dipping into the code.
with XGBoost in Python



Avoid Overfitting By Early Stopping With XGBoost In Python
Shengyin September 28, 2021 at 9:35 pm #

REPLY ↗

Hi Jason,

Thank you for the tutorial, it's really useful! However, I have been encountering this error (ValueError: Shape of passed values is (59372, 40), indices imply (59372, 41)) with the transform part, by any chance do you know how can I solve it? I also posted my question on Stack Overflow, but no luck 😞

The XGBoost With Python EBook is <https://stackoverflow.com/questions/69362344/valueerror-shape-of-passed-values-is-59372-40-indices-imply-59372-41> where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE



Adrian Tam September 29, 2021 at 11:58 pm #

REPLY ↗

Seems an off-by-one error. Check how you preprocess your data.



Hadi.94 December 7, 2021 at 11:35 pm #

REPLY ↗

Hey Mr. Jason .. thank you so much for your amazing article.

I'm using Feature Selection with XGBoost Feature Importance Scores with KNN based module and until now it has shown me great results.

I have one question, when I run the loop responsible of Feature Selection, I want to see the features that are involved in each iteration. Is there a simple way to do so ?

Never miss a tutorial:

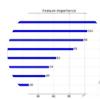
Adrian Tam December 8, 2021 at 8:07 am #

REPLY ↗



sim, way. It the best would be to drill into the XGBoost code to add a line or two to print that out.

Picked for you:



Feature Importance and Feature
SwapPy Selection With XGBoost In Python #

REPLY ↗



UserWarning: X has feature names, but SelectFromModel was fitted without feature names
arnings warn(

[How to Develop Your First XGBoost](#)



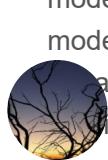
I am getting above mentioned error while I am trying to find the feature importance scores. DF has features with names in it. Below is the code I have used. Could you please mention a solution. Thank you.



[How to Use XGBoost for Time Series](#)

[Forecasting](#)


```
model = XGBClassifier()
```



```
model.fit(X_train, y_train)
```


[Take Data Preparation for Gradient Boosting](#)

[Actions - Model.predict\(X_test\)](#)


```
accuracy = accuracy_score(y_test, predictions)
```



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


[Avoid Overfitting By Early Stopping With](#)

[XGBoost In Python](#)


```
thresholds = sort(model.feature_importances_)
```


for thresh in thresholds:


```
# select features using threshold
```


[Loving the Tutorials?](#)


```
selection = SelectFromModel(model, threshold=thresh, prefit=True)
```



```
select_X_train = selection.transform(X_train)
```


[The XGBoost With Python EBook is](#)

where you'll find the [Really Good](#) stuff.


```
selection_model = selection.fit(select_X_train, y_train)
```



```
# eval | >> SEE WHAT'S INSIDE
```



```
select_X_test = selection.transform(X_test)
```



```
predictions = selection_model.predict(select_X_test)
```



```
accuracy = accuracy_score(y_test, predictions)
```



```
print("Thresh=%f, n=%d, Accuracy: %.2f%%" % (thresh, select_X_train.shape[1], accuracy*100.0))
```



James Carmichael June 15, 2022 at 7:19 am #

REPLY ↗

Hi Swappy...It looks like you are just using a code sample and not a full program listing.



Romy Opeña July 31, 2022 at 7:00 pm #

REPLY ↗

I am currently applying the XGBoost Classifier on the Kaggle mushroom classification data, replicating your codes in this article. When using XGBClassifier, the number of important features could

be reduced from original 22 variables down to 6-8 with still a high accuracy scores.
Never miss a tutorial:

I then tried to use the XGBRFClassifier on the same data and this further cut down another variable from 100 to lower values for the next 2 reductions and then it went back up to 100 from which it resumed the downward trend.

Picked for you:

It seems unusual but Is this normal or something is wrong with the module?



 How to Develop Your First XGBoost

Model in Python

James Carmichael August 1, 2022 at 8:44 am #

REPLY ↗

Hi Romy .. The following may be of interest to you:

[How to Use XGBoost for Time Series](#)

[Feature Selection](#)

<https://medium.com/towardsdatascience/techwarrior.com/why-does-the-loss-accuracy-fluctuate-during-the-training/>

 Data Preparation for Gradient Boosting

with XGBoost in Python

Romy Opena August 1, 2022 at 10:38 am #

REPLY ↗

Thanks, I will check on it. Meanwhile, I have decided to stick with XGBClassifier because I am

Avoid Overfitting By Early Stopping With

XGBoost In Python

 **Joe Salter** October 4, 2022 at 10:40 pm #

REPLY ↗

Hi Jason,

The [XGBoost With Python](#) EBook is

where you'll find the [Really Good](#) stuff.

Thanks for your great content.

Imagin >> SEE WHAT'S INSIDE and one target (y).

After building the model, I want to see what happens if I only change one of these predictors and keep the rest constant. In other words, I want to see only the effect of that specific predictor on the target.

Is there a specific way to do that?

I was thinking about making a mock dataset with all other predictors kept the same and just changing the one that I am interested in. Then predict y and plot changes in that specific predictor and changes in y.

Does that make sense? Because when I do it, then the predicted values of the mock data are the same...

 **James Carmichael** October 5, 2022 at 7:32 am #

REPLY ↗

Hi Joe...You are very welcome! In general, your suggestion is a valid one for small feature sets. The following may be of interest:



Leave a Reply

Feature Importance and Feature Selection With XGBoost in Python

How to Develop Your First XGBoost Model in Python

How to Use XGBoost for Time Series Forecasting

Name (required)

Data Preparation for Gradient Boosting with XGBoost in Python

Avoid Overfitting By Early Stopping With XGBoost In Python

Welcome!
I'm Jason Brownlee, PhD
and I help developers get results with machine learning.
The [XGBoost Best With Python](#) EBook is
where you'll find the **Really Good** stuff.

[>> SEE WHAT'S INSIDE](#)

AD

CONFLUENT

Data Mesh Architected with Event Streams

Download Now

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



How to Use XGBoost for Time Series Forecasting



Data Preparation for Gradient Boosting with XGBoost in Python



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.

© 2023 Guidina Tech Media. All Rights Reserved.

[LinkedIn](#) | [r](#) >> SEE WHAT'S INSIDE | [RSS](#)

[Privacy](#) | [Disclaimer](#) | [Terms](#) | [Contact](#) | [Sitemap](#) | [Search](#)