tds  Published in Towards Data Science

Bex T.  Follow

Apr 7, 2021 · 8 min read ★ · ▶ Listen

🔖 Save    🐦    🅵    in    🔗

# Beginner's Guide to XGBoost for Classification Problems

Utilize the hottest ML library for state-of-the-art performance



🏠            🔍            👤

Let me introduce you to the hottest Machine Learning library in the ML community — XGBoost. In recent years, it has been the main driving force behind the algorithms that win massive ML competitions. Its speed and performance are unparalleled and it consistently outperforms any other algorithms aimed at supervised learning tasks.

The library is parallelizable which means the core algorithm can run on clusters of GPUs or even across a network of computers. This makes it feasible to solve ML tasks by training on hundreds of millions of training examples with high performance.

Originally, it was written in C++ as a command-line application. After winning a huge competition in the field of physics, it started being widely adopted by the ML community. As a result, now the library has its APIs in several other languages including Python, R, and Julia.

In this post, you will learn the fundamentals of XGBoost to solve classification tasks, an overview of the massive list of XGBoost's hyperparameters and how to tune them.

**Refresher on Terminology**

Before we move on to code examples of XGBoost, let's refresh on some of the terms we will be using throughout the post.

**Classification task**: a supervised machine learning task in which one should predict if an instance is in some category by studying the instance's features. For example, by looking at the body measurements, patient history, and glucose levels of a person, you can predict whether a person belongs to the 'Has diabetes' or 'Does not have diabetes' group.

**Binary classification**: One type of classification where the target instance can only belong to either one of two classes. For example, predicting whether an email is a spam or not, whether a customer purchases some product or not, etc.

**Multi-class classification**: Another type of classification problem where the target can

If you find yourself confused by other terminology, I have written a small ML dictionary for beginners:

---

**Codeless Machine Learning Dictionary For Dummies**

Edit description

towardsdatascience.com

---

**How to preprocess your datasets for XGBoost**

Apart from basic data cleaning operations, there are some requirements for XGBoost to achieve top performance. Mainly:

- Numeric features should be scaled

- Categorical features should be encoded

To show how these steps are done, we will be using the Rain in Australia dataset from Kaggle where we will predict whether it will rain today or not based on some weather measurements. In this section, we will focus on preprocessing by utilizing Scikit-Learn Pipelines.

```python
1   import pandas as pd
2
3   rain = pd.read_csv("data/weatherAUS.csv")
4
5   >>> rain.head()
```
5001.py hosted with ❤ by GitHub                                                    view raw

```
1    >>> rain.info()
2
3    <class 'pandas.core.frame.DataFrame'>
4    RangeIndex: 145460 entries, 0 to 145459
5    Data columns (total 23 columns):
6     #   Column         Non-Null Count   Dtype
7    ---  ------         --------------   -----
8     0   Date           145460 non-null  object
9     1   Location       145460 non-null  object
10    2   MinTemp        143975 non-null  float64
11    3   MaxTemp        144199 non-null  float64
12    4   Rainfall       142199 non-null  float64
13    5   Evaporation    82670 non-null   float64
14    6   Sunshine       75625 non-null   float64
15    7   WindGustDir    135134 non-null  object
16    8   WindGustSpeed  135197 non-null  float64
17    9   WindDir9am     134894 non-null  object
18    10  WindDir3pm     141232 non-null  object
19    11  WindSpeed9am   143693 non-null  float64
20    12  WindSpeed3pm   142398 non-null  float64
21    13  Humidity9am    142806 non-null  float64
22    14  Humidity3pm    140953 non-null  float64
23    15  Pressure9am    130395 non-null  float64
24    16  Pressure3pm    130432 non-null  float64
25    17  Cloud9am       89572 non-null   float64
26    18  Cloud3pm       86102 non-null   float64
27    19  Temp9am        143693 non-null  float64
28    20  Temp3pm        141851 non-null  float64
29    21  RainToday      142199 non-null  object
30    22  RainTomorrow   142193 non-null  object
31   dtypes: float64(16), object(7)
32   memory usage: 25.5+ MB
```

**5002.py** hosted with ❤ by **GitHub**                                                        view raw

The dataset contains weather measures of 10 years from multiple weather stations in Australia. You can either predict whether it will rain tomorrow or today, so there are two targets in the dataset named `RainToday`, `RainTomorrow`.

```
cols_to_drop = ["Date", "Location", "RainTomorrow", "Rainfall"]

rain.drop(cols_to_drop, axis=1, inplace=True)
```

> *Dropping the* `Rainfall` *column is a must because it records the amount of rain in millimeters.*

Next, let's deal with missing values starting by looking at their proportions in each column:

```
1   missing_props = rain.isna().mean(axis=0)
2
3   >>> missing_props
4
5   MinTemp           0.010209
6   MaxTemp           0.008669
7   Evaporation       0.431665
8   Sunshine          0.480098
9   WindGustDir       0.070989
10  WindGustSpeed     0.070555
11  WindDir9am        0.072639
12  WindDir3pm        0.029066
13  WindSpeed9am      0.012148
14  WindSpeed3pm      0.021050
15  Humidity9am       0.018246
16  Humidity3pm       0.030984
17  Pressure9am       0.103568
18  Pressure3pm       0.103314
19  Cloud9am          0.384216
20  Cloud3pm          0.408071
21  Temp9am           0.012148
22  Temp3pm           0.024811
23  RainToday         0.022419
24  dtype: float64
```

**5003.py** hosted with 🧡 by **GitHub**                                    view raw

If the proportion is higher than 40% we will drop the column:

```
4
5    Evaporation      0.431665
6    Sunshine         0.480098
7    Cloud3pm         0.408071
8    dtype: float64
```

Three columns contain more than 40% missing values. We will drop them:

```
1    rain.drop(over_threshold.index,
2             axis=1,
3             inplace=True)
```

Now, before we move on to pipelines, let's divide the data into feature and target arrays beforehand:

```
1    X = rain.drop("RainToday", axis=1)
2    y = rain.RainToday
```

Next, there are both categorical and numeric features. We will build two separate pipelines and combine them later.

> The next code examples will heavily use Sklearn-Pipelines. If you are not familiar with them, check out my separate article for the _complete guide_ on them.

For the categorical features, we will impute the missing values with the mode of the column and encode them with One-Hot encoding:

```
1    from sklearn.impute import SimpleImputer
2    from sklearn.pipeline import Pipeline
3    from sklearn.preprocessing import OneHotEncoder
```

```
 9        ]
10    )
```

For the numeric features, I will choose the mean as an imputer and `StandardScaler` so that the features have 0 mean and a variance of 1:

```
1    from sklearn.preprocessing import StandardScaler
2
3    numeric_pipeline = Pipeline(
4        steps=[("impute", SimpleImputer(strategy="mean")),
5                ("scale", StandardScaler())]
6    )
```

Finally, we will combine the two pipelines with a column transformer. To specify which columns the pipelines are designed for, we should first isolate the categorical and numeric feature names:

```
1    cat_cols = X.select_dtypes(exclude="number").columns
2    num_cols = X.select_dtypes(include="number").columns
```

Next, we will input these along with their corresponding pipelines into a `ColumnTransFormer` instance:

```
1    from sklearn.compose import ColumnTransformer
2
3    full_processor = ColumnTransformer(
4        transformers=[
5            ("numeric", numeric_pipeline, num_cols),
6            ("categorical", categorical_pipeline, cat_cols),
7        ]
```

The full pipeline is finally ready. The only thing missing is the XGBoost classifier, which we will add in the next section.

**An Example of XGBoost For a Classification Problem**

To get started with `xgboost`, just install it either with `pip` or `conda`:

```
# pip
pip install xgboost

# conda
conda install -c conda-forge xgboost
```

After installation, you can import it under its standard alias — `xgb`. For classification problems, the library provides `XGBClassifier` class:

```
1    import xgboost as xgb
2
3    xgb_cl = xgb.XGBClassifier()
4
5    >>> print(type(xgb_cl))
6
7    <class 'xgboost.sklearn.XGBClassifier'>
```

5011.py hosted with ❤ by **GitHub**                                    view raw

Fortunately, the classifier follows the familiar fit-predict pattern of `sklearn` meaning we can freely use it as any `sklearn` model.

Before we train the classifier, let's preprocess the data and divide it into train and test sets:

```
1    # Apply preprocessing
```

```
7    from sklearn.model_selection import train_test_split
8
9    X_train, X_test, y_train, y_test = train_test_split(
10       X_processed, y_processed, stratify=y_processed, random_state=1121218
11   )
```

**5012.py** hosted with ❤ by **GitHub**                                    view raw

Since the target contains `NaN` , I imputed it by hand. Also, it is important to pass `y_processed` to `stratify` so that the split contains the same proportion of categories in both sets.

Now, we fit the classifier with default parameters and evaluate its performance:

```
1    from sklearn.metrics import accuracy_score
2
3    # Init classifier
4    xgb_cl = xgb.XGBClassifier()
5
6    # Fit
7    xgb_cl.fit(X_train, y_train)
8
9    # Predict
10   preds = xgb_cl.predict(X_test)
11
12   # Score
13   >>> accuracy_score(y_test, preds)
14   0.8507080984463082
```

**5013.py** hosted with ❤ by **GitHub**                                    view raw

Even with default parameters, we got an 85% accuracy which is reasonably good. In the next sections, we will try to improve the model even further by using `GridSearchCV` offered by Scikit-learn.

Unlike many other algorithms, XGBoost is an ensemble learning algorithm meaning that it combines the results of many models, called **base learners** to make a prediction.

Just like in Random Forests, XGBoost uses Decision Trees as base learners:

# Visualizing a Decision Tree of Rain

Yes    ( Sunny? )    No

No rain

Yes    ( Cloudy? )    No

No rain

Yes    ( Temp below 10 °C )    No
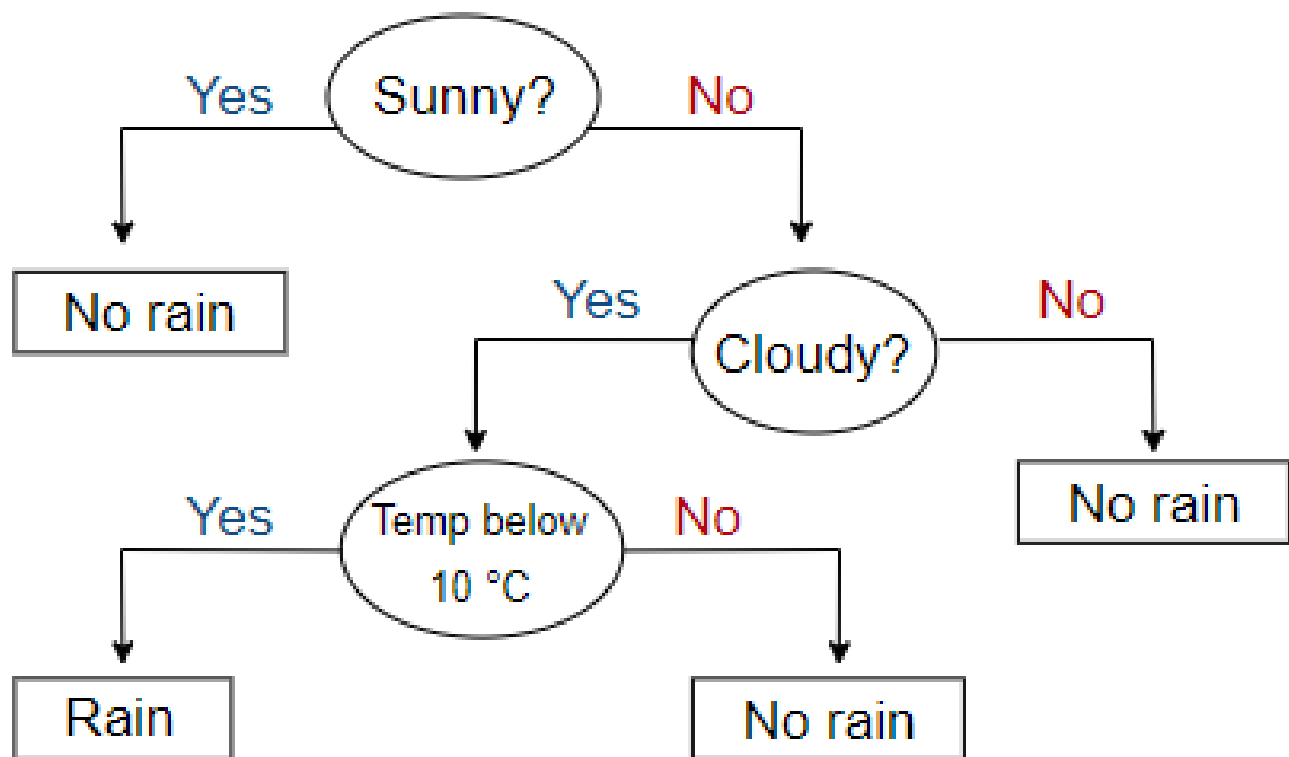
Rain                No rain

Image by the author. Decision tree to predict rain

An example of a decision tree can be seen above. In each decision node (circles), there is a single question that is being asked with only two possible answers. At the bottom of each tree, there is a single decision (rectangles). In the above tree, the first question is whether it is sunny or not. If yes, you immediately decide that it is not going to rain. If otherwise, you continue to ask more binary (yes/no) questions that ultimately will lead to some decision at the last "leaf" (rectangle).

unseen data.

However, the trees used by XGBoost are a bit different than traditional decision trees. They are called CART trees (Classification and Regression trees) and instead of containing a single decision in each "leaf" node, they contain real-value scores of whether an instance belongs to a group. After the tree reaches max depth, the decision can be made by converting the scores into categories using a certain threshold.

I am in no way an expert when it comes to the internals of XGBoost. That's why I recommend you to check out this awesome YouTube playlist entirely on XGBoost and another one solely aimed at Gradient Boosting which I did not mention at all.

**Overview of XGBoost Classifier Hyperparameters**

So far, we have been using only the default hyperparameters of the XGBoost Classifier:

```
1    >>> xgb_cl
2
3    XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
4                  colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
5                  importance_type='gain', interaction_constraints='',
6                  learning_rate=0.300000012, max_delta_step=0, max_depth=6,
7                  min_child_weight=1, missing=nan, monotone_constraints='()',
8                  n_estimators=100, n_jobs=4, num_parallel_tree=1, random_state=0,
9                  reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
10                 tree_method='exact', validate_parameters=1, verbosity=None)
```

**5014.py** hosted with ♥ by **GitHub**                                 view raw

> *Terminology refresher:* hyperparameters *of a model are the settings of that model which should be provided by the user. The model itself cannot learn these from the given training data.*

1. `learning_rate` : also called *eta*, it specifies how quickly the model fits the residual errors by using additional base learners.

- typical values: 0.01–0.2

2. `gamma, reg_alpha`, `reg_lambda` : these 3 parameters specify the values for 3 types of regularization done by XGBoost - minimum loss reduction to create a new split, L1 reg on leaf weights, L2 reg leaf weights respectively

- typical values for `gamma` : 0 - 0.5 but highly dependent on the data

- typical values for `reg_alpha` and `reg_lambda` : 0 - 1 is a good starting point but again, depends on the data

3. `max_depth` - how deep the tree's decision nodes can go. Must be a positive integer

- typical values: 1–10

4. `subsample` - fraction of the training set that can be used to train each tree. If this value is low, it may lead to underfitting or if it is too high, it may lead to overfitting

- typical values: 0.5–0.9

5. `colsample_bytree` - fraction of the features that can be used to train each tree. A large value means almost all features can be used to build the decision tree

- typical values: 0.5–0.9

The above are the main hyperparameters people often tune. It is perfectly OK if you don't understand them all completely (like me) but you can refer to this post which gives a thorough overview of how each of the above parameters works and how to tune them.

separate models on the given data for each combination of hyperparameters. I won't go into detail about how GridSearch works but you can check out my separate comprehensive article on the topic:

**Automatic Hyperparameter Tuning with Sklearn GridSearchCV and RandomizedSearchCV**

Edit description

towardsdatascience.com

We will be tuning only a few of the parameters in two rounds because of how tuning is both computationally and time-expensive. Let's create the parameter grid for the first round:

```
1   param_grid = {
2       "max_depth": [3, 4, 5, 7],
3       "learning_rate": [0.1, 0.01, 0.05],
4       "gamma": [0, 0.25, 1],
5       "reg_lambda": [0, 1, 10],
6       "scale_pos_weight": [1, 3, 5],
7       "subsample": [0.8],
8       "colsample_bytree": [0.5],
9   }
```

**5015.py** hosted with ❤ by **GitHub**                              view raw

In the grid, I fixed `subsample` and `colsample_bytree` to recommended values to speed things up and prevent overfitting.

We will import `GridSearchCV` from `sklearn.model_selection`, instantiate and fit it to our preprocessed data:

```
1   from sklearn.model_selection import GridSearchCV
```

```
7    grid_cv = GridSearchCV(xgb_cl, param_grid, n_jobs=-1, cv=3, scoring="roc_auc")

8

9    # Fit
10   _ = grid_cv.fit(X_processed, y_processed)
```

**5016.py** hosted with 🖤 by **GitHub**                    view raw

After an excruciatingly long time, we finally got the best params and best score:

```
1    >>> grid_cv.best_score_
2    0.853810061069393
```

**5017.py** hosted with 🖤 by **GitHub**                    view raw

This time, I chose `roc_auc` metric which calculates the area under the ROC (receiver operating characteristic) curve. It is one of the most popular and robust evaluation metrics for unbalanced classification problems. You can learn more about it here. Let's see the best params:

```
1    >>> grid_cv.best_params_
2
3    {'gamma': 1,
4     'learning_rate': 0.1,
5     'max_depth': 7,
6     'reg_lambda': 10,
7     'scale_pos_weight': 3}
```

**5018.py** hosted with 🖤 by **GitHub**                    view raw

As you can see, only `scale_pos_weight` is in the middle of its provided range. The other parameters are at the end of their ranges meaning that we have to keep exploring:

```
1    # Insert the new fixed values to the grid
2    param_grid["scale_pos_weight"] = [3]
3    param_grid["subsample"] = [0.8]
4    param_grid["colsample_bytree"] = [0.5]
```

```
10    param_grid["learning_rate"] = [0.3, 0.5, 0.7, 1]
```

We will fit a new GridSearch object to the data with the updated param grid and see if we got an improvement on the best score:

```
1    grid_cv_2 = GridSearchCV(xgb_cl, param_grid,
2                             cv=3, scoring="roc_auc", n_jobs=-1)
3
4    _ = grid_cv_2.fit(X_processed, y_processed)
5
6    >>> grid_cv_2.best_score_
7    0.8501002182976786
```

Looks like the second round of tuning resulted in a slight decrease in performance. We have got no choice but to stick with the first set of parameters which were:

```
1    >>> grid_cv.best_params_
2
3    {'gamma': 1,
4     'learning_rate': 0.1,
5     'max_depth': 7,
6     'reg_lambda': 10,
7     'scale_pos_weight': 3}
```

Let's create a final classifier with the above parameters:

```
1    final_cl = xgb.XGBClassifier(
2        **grid_cv.best_params_,
3        objective="binary:logistic",
4        colsample_bytree=0.5,
5        subsample=0.8
```

Finally, make predictions on the test set:

```python
from sklearn.metrics import roc_auc_score

_ = final_cl.fit(X_train, y_train)

preds = final_cl.predict(X_test)
```

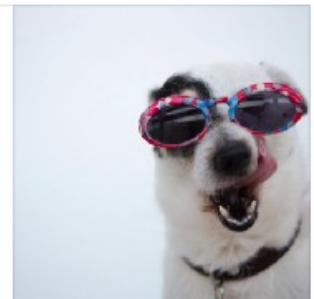**5023.py** hosted with ❤️ by **GitHub**                    view raw

## Conclusion

We have made it to the end of this introductory guide on XGBoost for classification problems. Even though we covered a lot, there are still many topics to explore in terms of XGBoost itself and on the topic of classification.

I strongly recommend you to check out the links I provided as additional sources to learn XGBoost and suggest reading more on how to tackle classification problems.

How To Touch Into the Heart of Matplotlib And Create Amazing Plots

Use Matplotlib like an absolute boss!

towardsdatascience.com

👏 159  |  💬 2

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. Take a look.

Get this newsletter