

Sign in

Get started

Follow

596K Followers

· Editors' Picks

Features

Deep Dives

Grow

Contribute

About

You have **2** free member-only stories left this month. [Sign up for Medium and get an extra one](#)

You Are Missing Out on LightGBM. It Crushes XGBoost in Every Aspect

Not anymore, XGBoost, not anymore



Bex T. · Sep 2 · 7 min read ★

Learn how to crush XGBoost in this comprehensive LightGBM tutorial.



Photo by [GR Stocks](#) on [Unsplash](#). All images are by the author unless specified otherwise.

I am confused.

So many people are drawn to XGBoost like a moth to a flame. Yes, it has seen some glorious days in prestigious competitions, and it's *still* the most widely-used ML library.

But, it has been 4 years since XGBoost lost its top spot in terms of performance. In 2017, Microsoft open-sourced LightGBM (Light Gradient Boosting Machine) that gives **equally high accuracy with 2–10 times less training speed**.

This is a game-changing advantage considering the ubiquity of massive, million-row datasets. There are other distinctions that tip the scales towards LightGBM and give it an edge over XGBoost.

By the end of this post, you will learn about these advantages, including:

- how to develop LightGBM models for classification and regression tasks
- structural differences between XGBoost and LGBM
- how to use early stopping and evaluation sets
- enabling powerful categorical feature support for up 8x times speed increase
- implementing successful cross-validation with LGBM
- hyperparameter tuning with Optuna (Part II)

XGBoost vs. LightGBM

When LGBM got released, it came with ground-breaking changes to the way it grows decision trees.

Both XGBoost and LightGBM are ensemble algorithms. They use a special type of decision trees, also called weak learners, to capture complex, non-linear patterns.

In XGBoost (and many other libraries), decision trees were built one level at a time:

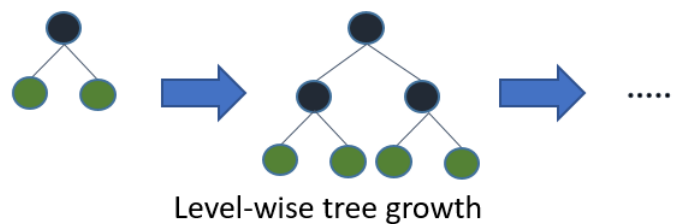


Image from LGBM documentation

This type of structure tends to result in unnecessary nodes and leaves because the trees continued to build until the `max_depth` reached. This led to higher model complexity and training cost runtime.

In contrast, LightGBM takes a leaf-wise approach:



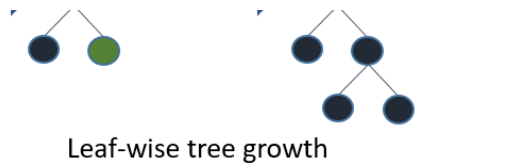


Image from LGBM documentation

The structure continues to grow with the most promising branches and leaves (nodes with the most delta loss), holding the number of the decision leaves constant. (If this doesn't make sense to you, don't sweat. This won't prevent you from effectively using LGBM).

This is one of the main reasons LGBM crushed XGBoost in terms of speed when it first came out.

The following table is the comparison of time cost:

Data	xgboost	xgboost_hist	LightGBM
Higgs	3794.34 s	165.575 s	130.094 s
Yahoo LTR	674.322 s	131.462 s	76.229 s
MS LTR	1251.27 s	98.386 s	70.417 s
Expo	1607.35 s	137.65 s	62.607 s
Allstate	2867.22 s	315.256 s	148.231 s

LightGBM ran faster than xgboost on all experiment data sets.

Image from LGBM documentation

Above is a benchmark comparison of XGBoost with traditional decision trees and LGBM with leaf-wise structure (first and last columns) on datasets with ~500k-13M samples. It shows that LGBM is orders of magnitude faster than XGB.

LGBM also uses histogram binning of continuous features, which provides even more speed-up than traditional gradient boosting. Binning numeric values significantly decrease the number of split points to consider in decision trees, and they remove the need to use sorting algorithms, which are always computation-heavy.

Inspired by LGBM, XGBoost also introduced histogram-binning, which gave massive speed-up but still not enough to match LGBM's:

Data	xgboost	xgboost_hist	LightGBM
Higgs	3794.34 s	165.575 s	130.094 s
Yahoo LTR	674.322 s	131.462 s	76.229 s
MS LTR	1251.27 s	98.386 s	70.417 s
Expo	1607.35 s	137.65 s	62.607 s
Allstate	2867.22 s	315.256 s	148.231 s

Image from LGBM documentationHistogram-binning comparison — second and third columns.

We will continue exploring the differences in the coming sections.

Model initialization and objectives

Like XGBoost, LGBM has two APIs — core learning API and Sklearn-compatible one. You know I am a big fan of Sklearn, so this tutorial will focus on that version.

Sklearn-compatible API of XGBoost and LGBM allows you to integrate their models in the Sklearn ecosystem so that you can use them inside pipelines in combination with other transformers.

Sklearn API exposes `LGBMRegressor` and `LGBMClassifier`, with the familiar `fit/predict/predict_proba` pattern:

```
1 import lightgbm as lgbm # standard alias
2
3 clf = lgbm.LGBMClassifier(objective="binary") # or 'multiclass'
4 reg = lgbm.LGBMRegressor() # default - 'regression'
```

8801.py hosted with ❤ by GitHub

[view raw](#)

objective specifies the type of learning task. Besides the common ones like `binary`, `multiclass` and `regression` tasks, there are others like `poisson`, `tweedie` regressions. See [this section](#) of the documentation for the full list of objectives.

```
>>> reg.fit(X, y)
```

```
LGBMRegressor()
```

Controlling the number of decision trees

The number of decision trees inside the ensemble significantly affects the results. You can control it using the `n_estimators` parameter in both the classifier and regressor. Below, we will fit an LGBM binary classifier on the Kaggle TPS March dataset with 1000 decision trees:

```
1 tps_march = pd.read_csv("data/train.csv")
2 tps_march.head()
```

8802.py hosted with ❤ by GitHub

[view raw](#)

	id	cat0	cat1	cat2	cat3	cat4	cat5	cat6	cat7	cat8	...	cont2	cont3	cont4	cont5	cont6	cont7	cont8	cont9	cont10	target
0	0	A	I	A	B	B	BI	A	S	Q	...	0.759439	0.795549	0.681917	0.621672	0.592164	0.791921	0.815254	0.965006	0.665915	0
1	1	A	I	A	A	E	BI	K	W	AD	...	0.386385	0.541366	0.388982	0.357778	0.600044	0.408701	0.399353	0.927406	0.493729	0
2	2	A	K	A	A	E	BI	A	E	BM	...	0.343255	0.616352	0.793687	0.552877	0.352113	0.388835	0.412303	0.292696	0.549452	0
3	3	A	K	A	C	E	BI	A	Y	AD	...	0.831147	0.807807	0.800032	0.619147	0.221789	0.897617	0.633669	0.760318	0.934242	0
4	4	A	I	G	B	E	BI	C	G	Q	...	0.338818	0.277308	0.610578	0.128291	0.578764	0.279167	0.351103	0.357084	0.328960	1

5 rows x 32 columns

```
1 >>> tps_march.shape
2 (300000, 32)
3
4 >>> tps_march.dtypes.value_counts()
5 object      19
6 float64     11
7 int64        2
8 dtype: int64
9
10 >>> X, y = tps_march.drop("target", axis=1), tps_march[["target"]].values.flatten()
```

```

11
12 # Encode categoricals
13 >>> X_enc = pd.get_dummies(X)
14
15 >>> clf = lgbm.LGBMClassifier(
16 ...     objective="binary", n_estimators=1000, random_state=1121218
17 ... )
18 >>> clf.fit(X_enc, y)
19 LGBMClassifier(n_estimators=1000, objective='binary', random_state=1121218)

```

8803.py hosted with ❤ by GitHub

[view raw](#)

Adding more trees leads to more accuracy but increases the risk of overfitting. To combat this, you can create many trees (+2000) and choose a smaller `learning_rate` (more on this later).

```

1 >>> clf = lgbm.LGBMClassifier(
2 ...     objective="binary", n_estimators=3000, learning_rate=0.1, random_state=1121218
3 ...)
4
5 >>> clf.fit(X_enc, y)
6 LGBMClassifier(n_estimators=3000, objective='binary', random_state=1121218)

```

8804.py hosted with ❤ by GitHub

[view raw](#)

Like in XGBoost, fitting a single decision tree to the data is called a **boosting round**.

Early stopping

Each tree in the ensemble builds on the predictions of the last tree — i.e., each boosting round is an improvement of the last.

If the predictions don't improve after a sequence of rounds, it is sensible to stop the training of the ensemble even if we are not at a hard stop for

`n_estimators`.

To achieve this, **LGBM provides `early_stopping_rounds` parameter inside the `fit` function**. For example, setting it to 100 means we stop the training if the predictions have not improved for the last 100 rounds.

Before looking at a code example, we should learn a couple of concepts connected to early stopping.

Eval sets and metrics

Early stopping is only enabled when you pass a set of evaluation sets to `eval_set` parameter of the `fit` method. These evaluation sets are used to keep track of the quality of the predictions from one boosting round to the next:

```

1 from sklearn.model_selection import train_test_split
2
3 X_train, X_eval, y_train, y_eval = train_test_split(X_enc, y, test_size=0.1)
4
5 clf = lgbm.LGBMClassifier(objective="binary", n_estimators=10000)

```

```

6  eval_set = [(X_eval, y_eval)]
7
8  clf.fit(
9      X_train,
10     y_train,
11     eval_set=eval_set,
12     early_stopping_rounds=100,
13     eval_metric="binary_logloss",
14 )
15
16 -----
17
18 Training until validation scores don't improve for 100 rounds
19 [1]   valid_0's binary_logloss: 0.542547
20 [2]   valid_0's binary_logloss: 0.515902
21 [3]   valid_0's binary_logloss: 0.494678
22 [4]   valid_0's binary_logloss: 0.477235
23 [5]   valid_0's binary_logloss: 0.462826
24 ...
25 [737] valid_0's binary_logloss: 0.350179
26 [738] valid_0's binary_logloss: 0.350162
27 Early stopping, best iteration is:
28 [638] valid_0's binary_logloss: 0.350073

```

8805.py hosted with ❤ by GitHub

[view raw](#)

In each round of `n_estimators`, a single decision tree is fit to `(X_train, y_train)` and predictions are made on the passed evaluation set `(X_eval, y_eval)`. The quality of predictions is measured with a passed metric in `eval_metric`.

The training stops at the 738th iteration because the validation score has not improved since the 638th one — early stopping of 100 rounds. Now, we have the luxury of creating as many trees as we want and `early_stopping_rounds` can discard the unnecessary ones.

Establish a baseline

Let's establish a baseline score with what we know so far. We will do the same for XGBoost so that we can compare the results:

```

1  import xgboost as xgb
2  from sklearn.metrics import log_loss
3
4  X, y = tps_march.drop("target", axis=1), tps_march[["target"]].values.flatten()
5
6  # Encode categoricals
7  X_enc = pd.get_dummies(X)
8  X_train, X_eval, y_train, y_eval = train_test_split(
9      X_enc, y, test_size=0.2, stratify=y
10 )

```

8806.py hosted with ❤ by GitHub

[view raw](#)

```

1  %%time
2
3  xgb_clf = xgb.XGBClassifier(
4      objective="binary:logistic",
5      random_state=1121218,
6      n_estimators=10000,
7      tree_method="hist", # enable histogram binning in XGB
8  )
9
10 xgb_clf.fit(
11     X_train,
12     y_train,
13     eval_set=[(X_eval, y_eval)],

```

```

14     eval_metric="logloss",
15     early_stopping_rounds=150,
16     verbose=False, # Disable logs
17 )
18
19 preds = xgb_clf.predict_proba(X_eval)
20 print(f"XGBoost logloss on the evaluation set: {log_loss(y_eval, preds):.5f}")
21
22 -----
23
24 XGBoost logloss on the evaluation set: 0.35482
25 Wall time: 1min 22s

```

8807.py hosted with ❤ by GitHub

[view raw](#)

```

1  %%time
2
3  lgbm_clf = lgbm.LGBMClassifier(
4      objective="binary",
5      random_state=1121218,
6      n_estimators=10000,
7      boosting="gbdt", # default histogram binning of LGBM
8      # device='gpu' # uncomment to use GPU training
9  )
10
11 lgbm_clf.fit(
12     X_train,
13     y_train,
14     eval_set=[(X_eval, y_eval)],
15     eval_metric="binary_logloss",
16     early_stopping_rounds=150,
17     verbose=False, # Disable logs
18 )
19
20 preds = lgbm_clf.predict_proba(X_eval)
21 print(f"LightGBM logloss on the evaluation set: {log_loss(y_eval, preds):.5f}")
22
23 -----
24
25 LightGBM logloss on the evaluation set: 0.35256
26 Wall time: 17.2 s

```

8808.py hosted with ❤ by GitHub

[view raw](#)

LGBM achieved a smaller loss in ~4 times less runtime. Let's see a final LGBM trick before we move on to cross-validation.

Categorical and missing values support

Histogram binning in LGBM comes with built-in support for handling missing values and categorical features. TPS March dataset contains 19 categoricals, and we have been using one-hot encoding up to this point.

This time, we will let LGBM deal with categoricals and compare the results with XGBoost once again:

```

1  X, y = tps_march.drop("target", axis=1), tps_march[["target"]].values.flatten()
2
3  # Extract categoricals and their indices
4  cat_features = X.select_dtypes(exclude=np.number).columns.to_list()
5  cat_idx = [X.columns.get_loc(col) for col in cat_features]
6
7  # Convert cat_features to pd.Categorical dtype
8  for col in cat_features:
9      X[col] = pd.Categorical(X[col])
10
11 # Unencoded train/test sets
12 X_train, X_eval, y_train, y_eval = train_test_split(

```

```

13     X, y, test_size=0.2, random_state=4, stratify=y
14 )

```

8809.py hosted with ❤ by GitHub

[view raw](#)

To specify the categorical features, pass a list of their indices to `categorical_feature` parameter in the `fit` method:

```

1  %%time
2
3  # Model initialization is the same
4  eval_set = [(X_eval, y_eval)]
5
6  _ = lgbm_clf.fit(
7      X_train,
8      y_train,
9      categorical_feature=cat_idx, # Specify the categoricals
10     eval_set=eval_set,
11     early_stopping_rounds=150,
12     eval_metric="logloss",
13     verbose=False,
14 )
15
16 preds = lgbm_clf.predict_proba(X_eval)
17 loss = log_loss(y_eval, preds)
18 print(f"LGBM logloss with default categorical feature handling: {loss:.5f}")
19
20 -----
21
22 LGBM logloss with default categorical feature handling: 0.34823
23 Wall time: 15 s

```

8810.py hosted with ❤ by GitHub

[view raw](#)

You can achieve up to 8x speed up if you use `pandas.Categorical` data type when using LGBM.

	XGB (one-hot)	LGBM (one-hot)	LGBM (cat-handling)
score	0.35482	0.35256	0.34823
runtime	1min 22s	17.2 s	15 s

The table shows the final scores and runtimes of both models. As you can see, the version with default categorical handling beat others both in accuracy and speed. Cheers!

Cross-validation with LightGBM

The most common way of doing CV with LGBM is to use Sklearn CV splitters.

I am not talking about utility functions like `cross_validate` or `cross_val_score` but splitters like `KFold` or `StratifiedKFold` with their `split` method. Doing CV in this way gives you more control over the whole process.

I have talked many times about the importance of cross-validation. You can read [this post](#) for more details.

Also, it enables you to use early stopping during cross-validation in a hassle-free manner. Here is what this looks like for the TPS March data:


```

1  import time
2
3  from sklearn.model_selection import StratifiedKFold
4
5  N_SPLITS = 7
6  strat_kf = StratifiedKFold(n_splits=N_SPLITS, shuffle=True, random_state=1121218)
7
8  scores = np.empty(N_SPLITS)
9  for idx, (train_idx, test_idx) in enumerate(strat_kf.split(X, y)):
10     print("=" * 12 + f"Training fold {idx}" + 12 * "=")
11     start = time.time()
12
13     X_train, X_val = X.iloc[train_idx], X.iloc[test_idx]
14     y_train, y_val = y[train_idx], y[test_idx]
15     eval_set = [(X_val, y_val)]
16
17     lgbm_clf = lgbm.LGBMClassifier(n_estimators=10000)
18     lgbm_clf.fit(
19         X_train,
20         y_train,
21         eval_set=eval_set,
22         categorical_feature=cat_idx,
23         early_stopping_rounds=200,
24         eval_metric="binary_logloss",
25         verbose=False,
26     )
27
28     preds = lgbm_clf.predict_proba(X_val)
29     loss = log_loss(y_val, preds)
30     scores[idx] = loss
31     runtime = time.time() - start
32     print(f"Fold {idx} finished with score: {loss:.5f} in {runtime:.2f} seconds.\n")
33
34     -----
35
36     =====Training fold 0=====
37     Fold 0 finished with score: 0.34547 in 15.46 seconds.
38
39     =====Training fold 1=====
40     Fold 1 finished with score: 0.34824 in 14.90 seconds.
41
42     =====Training fold 2=====
43     Fold 2 finished with score: 0.34665 in 12.64 seconds.
44
45     =====Training fold 3=====
46     Fold 3 finished with score: 0.35015 in 14.78 seconds.
47
48     =====Training fold 4=====
49     Fold 4 finished with score: 0.34617 in 14.61 seconds.
50
51     =====Training fold 5=====
52     Fold 5 finished with score: 0.34833 in 13.24 seconds.
53
54     =====Training fold 6=====
55     Fold 6 finished with score: 0.35155 in 13.27 seconds.
56     print(f"Average log loss: {np.mean(scores):.4f}")
57     Average log loss: 0.3481

```

8811.py hosted with  by GitHub

[view raw](#)

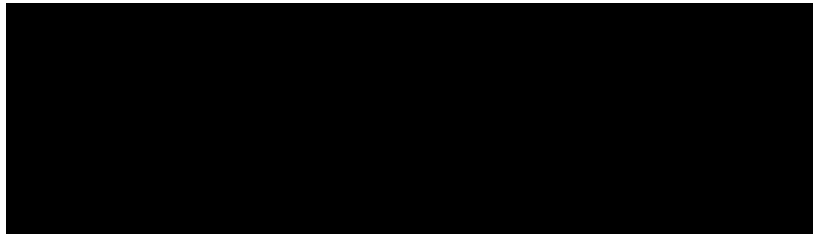
First, create a CV splitter — we are choosing `StratifiedKFold` because it is a classification problem. Then, loop through each train/test sets using `split`. In each fold, initialize and train a new LGBM model and optionally report the score and runtime. That's it! That's how most people do CV, including on Kaggle.

Conclusion

In this post, we learned pure modeling techniques with LightGBM. Next up, we will explore how to squeeze every bit of performance out of LGBM

models using Optuna.

Specifically, Part II of this article will include a detailed overview of the most important LGBM hyperparameters and introduce a well-tested hyperparameter tuning workflow. It is already out — read it [here](#).

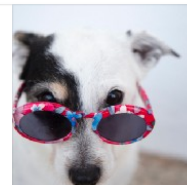


You might also be interested...

25 NumPy Functions You Never Knew Existed | P (Guarantee = 0.85)

Become a NumPy Ninja

[towardsdatascience.com](#)



Kaggle's Guide to LightGBM Hyperparameter Tuning with Optuna in 2021

[towardsdatascience.com](#)



Tired of Cliché Datasets? Here are 18 Awesome Alternatives From All Domains

[Edit description](#)

[towardsdatascience.com](#)

Love 3Blue1Brown Animations? Learn How to Create Your Own in Python in 10 Minutes

[Edit description](#)

[towardsdatascience.com](#)

7 Cool Python Packages Kagglers Are Using Without Telling You

[Edit description](#)

[towardsdatascience.com](#)

Sign up for The Variable

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

[Artificial Intelligence](#) [Data Science](#) [Machine Learning](#) [Technology](#) [Programming](#)

[About](#) [Write](#) [Help](#) [Legal](#)