

---

# **pygbm Documentation**

*Release 0.1.0.dev0*

**Olivier Grisel, Nicolas Hug**

**Dec 19, 2018**



<b>1</b>	<b>Gradient Boosting Estimators</b>	<b>3</b>
<b>2</b>	<b>Losses</b>	<b>9</b>
<b>3</b>	<b>Binning</b>	<b>13</b>
<b>4</b>	<b>Grower</b>	<b>15</b>
<b>5</b>	<b>Splitting</b>	<b>19</b>
<b>6</b>	<b>Indices and tables</b>	<b>21</b>
	<b>Python Module Index</b>	<b>23</b>



**Warning:** Pygbm's API and default values are likely to be changed in future version, without any deprecation cycle.

**Warning:** Pygbm's API and default values are likely to be changed in future version, without any deprecation cycle.



---

## Gradient Boosting Estimators

---

Gradient Boosting decision trees for classification and regression.

```
class pygbm.gradient_boosting.GradientBoostingClassifier (loss='auto', learning_rate=0.1,
                                                         max_iter=100,
                                                         max_leaf_nodes=31,
                                                         max_depth=None,
                                                         min_samples_leaf=20,
                                                         l2_regularization=0.0,
                                                         max_bins=256, scoring=None, validation_split=0.1,
                                                         n_iter_no_change=5,
                                                         tol=1e-07, verbose=0,
                                                         random_state=None)
```

Scikit-learn compatible Gradient Boosting Tree for classification.

### Parameters

- **loss** (*{'auto', 'binary\_crossentropy', 'categorical\_crossentropy'}, optional (default='auto')*) – The loss function to use in the boosting process. ‘binary\_crossentropy’ (also known as logistic loss) is used for binary classification and generalizes to ‘categorical\_crossentropy’ for multiclass classification. ‘auto’ will automatically choose either loss depending on the nature of the problem.
- **learning\_rate** (*float, optional (default=1)*) – The learning rate, also known as *shrinkage*. This is used as a multiplicative factor for the leaves values. Use 1 for no shrinkage.
- **max\_iter** (*int, optional (default=100)*) – The maximum number of iterations of the boosting process, i.e. the maximum number of trees for binary classification. For multiclass classification, *n\_classes* trees per iteration are built.
- **max\_leaf\_nodes** (*int or None, optional (default=None)*) – The maximum number of leaves for each tree. If None, there is no maximum limit.

- **max\_depth** (*int or None, optional (default=None)*) – The maximum depth of each tree. The depth of a tree is the number of nodes to go from the root to the deepest leaf.
- **min\_samples\_leaf** (*int, optional (default=20)*) – The minimum number of samples per leaf.
- **l2\_regularization** (*float, optional (default=0)*) – The L2 regularization parameter. Use 0 for no regularization.
- **max\_bins** (*int, optional (default=256)*) – The maximum number of bins to use. Before training, each feature of the input array *X* is binned into at most `max_bins` bins, which allows for a much faster training stage. Features with a small number of unique values may use less than `max_bins` bins. Must be no larger than 256.
- **scoring** (*str or callable or None, optional (default=None)*) – Scoring parameter to use for early stopping (see `sklearn.metrics` for available options). If `None`, early stopping is check w.r.t the loss value.
- **validation\_split** (*int or float or None, optional (default=0.1)*) – Proportion (or absolute size) of training data to set aside as validation data for early stopping. If `None`, early stopping is done on the training data.
- **n\_iter\_no\_change** (*int or None, optional (default=5)*) – Used to determine when to “early stop”. The fitting process is stopped when none of the last `n_iter_no_change` scores are better than the “`n_iter_no_change - 1`”th-to-last one, up to some tolerance. If `None` or 0, no early-stopping is done.
- **tol** (*float or None optional (default=1e-7)*) – The absolute tolerance to use when comparing scores. The higher the tolerance, the more likely we are to early stop: higher tolerance means that it will be harder for subsequent iterations to be considered an improvement upon the reference score.
- **verbose** (*int, optional (default=0)*) – The verbosity level. If not zero, print some information about the fitting process.
- **random\_state** (*int, np.random.RandomStateInstance or None, optional (default=None)*) – Pseudo-random number generator to control the subsampling in the binning process, and the train/validation data split if early stopping is enabled. See [scikit-learn glossary](#).

## Examples

```
>>> from sklearn.datasets import load_iris
>>> from pygbm import GradientBoostingClassifier
>>> X, y = load_iris(return_X_y=True)
>>> clf = GradientBoostingClassifier().fit(X, y)
>>> clf.score(X, y)
0.97...
```

**fit** (*X, y*)

Fit the gradient boosting model.

### Parameters

- **X** (*array-like, shape=(n\_samples, n\_features)*) – The input samples. If `X.dtype == np.uint8`, the data is assumed to be pre-binned and the prediction methods (`predict`, `predict_proba`) will only accept pre-binned data as well.



- **y** (*array-like*, *shape=(n\_samples,)*) – Target values.

**Returns** **self**

**Return type** object

**get\_params** (*deep=True*)

Get parameters for this estimator.

**Parameters** **deep** (*boolean, optional*) – If True, will return the parameters for this estimator and contained subobjects that are estimators.

**Returns** **params** – Parameter names mapped to their values.

**Return type** mapping of string to any

**predict** (*X*)

Predict classes for X.

**Parameters** **X** (*array-like, shape=(n\_samples, n\_features)*) – The input samples. If `X.dtype == np.uint8`, the data is assumed to be pre-binned and the estimator must have been fitted with pre-binned data.

**Returns** **y** – The predicted classes.

**Return type** array, shape (n\_samples,)

**predict\_proba** (*X*)

Predict class probabilities for X.

**Parameters** **X** (*array-like, shape=(n\_samples, n\_features)*) – The input samples. If `X.dtype == np.uint8`, the data is assumed to be pre-binned and the estimator must have been fitted with pre-binned data.

**Returns** **p** – The class probabilities of the input samples.

**Return type** array, shape (n\_samples, n\_classes)

**score** (*X, y, sample\_weight=None*)

Returns the mean accuracy on the given test data and labels.

In multi-label classification, this is the subset accuracy which is a harsh metric since you require for each sample that each label set be correctly predicted.

**Parameters**

- **X** (*array-like, shape = (n\_samples, n\_features)*) – Test samples.
- **y** (*array-like, shape = (n\_samples) or (n\_samples, n\_outputs)*) – True labels for X.
- **sample\_weight** (*array-like, shape = [n\_samples], optional*) – Sample weights.

**Returns** **score** – Mean accuracy of `self.predict(X)` wrt. `y`.

**Return type** float

**set\_params** (*\*\*params*)

Set the parameters of this estimator.

The method works on simple estimators as well as on nested objects (such as pipelines). The latter have parameters of the form `<component>__<parameter>` so that it's possible to update each component of a nested object.

**Returns**

Return type self

```
class pygbm.gradient_boosting.GradientBoostingRegressor (loss='least_squares',  
learning_rate=0.1,  
max_iter=100,  
max_leaf_nodes=31,  
max_depth=None,  
min_samples_leaf=20,  
l2_regularization=0.0,  
max_bins=256, scoring=None, validation_split=0.1,  
n_iter_no_change=5,  
tol=1e-07, verbose=0,  
random_state=None)
```

Scikit-learn compatible Gradient Boosting Tree for regression.

### Parameters

- **loss** (*{'least\_squares'}*, *optional (default='least\_squares')*) – The loss function to use in the boosting process.
- **learning\_rate** (*float, optional (default=0.1)*) – The learning rate, also known as *shrinkage*. This is used as a multiplicative factor for the leaves values. Use 1 for no shrinkage.
- **max\_iter** (*int, optional (default=100)*) – The maximum number of iterations of the boosting process, i.e. the maximum number of trees.
- **max\_leaf\_nodes** (*int or None, optional (default=None)*) – The maximum number of leaves for each tree. If None, there is no maximum limit.
- **max\_depth** (*int or None, optional (default=None)*) – The maximum depth of each tree. The depth of a tree is the number of nodes to go from the root to the deepest leaf.
- **min\_samples\_leaf** (*int, optional (default=20)*) – The minimum number of samples per leaf.
- **l2\_regularization** (*float, optional (default=0)*) – The L2 regularization parameter. Use 0 for no regularization.
- **max\_bins** (*int, optional (default=256)*) – The maximum number of bins to use. Before training, each feature of the input array X is binned into at most `max_bins` bins, which allows for a much faster training stage. Features with a small number of unique values may use less than `max_bins` bins. Must be no larger than 256.
- **scoring** (*str or callable or None, optional (default=None)*) – Scoring parameter to use for early stopping (see `sklearn.metrics` for available options). If None, early stopping is check w.r.t the loss value.
- **validation\_split** (*int or float or None, optional (default=0.1)*) – Proportion (or absolute size) of training data to set aside as validation data for early stopping. If None, early stopping is done on the training data.
- **n\_iter\_no\_change** (*int or None, optional (default=5)*) – Used to determine when to “early stop”. The fitting process is stopped when none of the last `n_iter_no_change` scores are better than the “`n_iter_no_change - 1`”-th-to-last one, up to some tolerance. If None or 0, no early-stopping is done.

- **tol** (*float or None optional (default=1e-7)*) – The absolute tolerance to use when comparing scores. The higher the tolerance, the more likely we are to early stop: higher tolerance means that it will be harder for subsequent iterations to be considered an improvement upon the reference score.
- **verbose** (*int, optional (default=0)*) – The verbosity level. If not zero, print some information about the fitting process.
- **random\_state** (*int, np.random.RandomStateInstance or None, optional (default=None)*) – Pseudo-random number generator to control the subsampling in the binning process, and the train/validation data split if early stopping is enabled. See [scikit-learn glossary](#).

## Examples

```
>>> from sklearn.datasets import load_boston
>>> from pygbm import GradientBoostingRegressor
>>> X, y = load_boston(return_X_y=True)
>>> est = GradientBoostingRegressor().fit(X, y)
>>> est.score(X, y)
0.92...
```

**fit** (*X, y*)

Fit the gradient boosting model.

### Parameters

- **X** (*array-like, shape=(n\_samples, n\_features)*) – The input samples. If `X.dtype == np.uint8`, the data is assumed to be pre-binned and the prediction methods (`predict`, `predict_proba`) will only accept pre-binned data as well.
- **y** (*array-like, shape=(n\_samples,)*) – Target values.

**Returns** `self`

**Return type** object

**get\_params** (*deep=True*)

Get parameters for this estimator.

**Parameters** **deep** (*boolean, optional*) – If `True`, will return the parameters for this estimator and contained subobjects that are estimators.

**Returns** **params** – Parameter names mapped to their values.

**Return type** mapping of string to any

**predict** (*X*)

Predict values for *X*.

**Parameters** **X** (*array-like, shape=(n\_samples, n\_features)*) – The input samples. If `X.dtype == np.uint8`, the data is assumed to be pre-binned and the estimator must have been fitted with pre-binned data.

**Returns** **y** – The predicted values.

**Return type** array, shape (n\_samples,)

**score** (*X, y, sample\_weight=None*)

Returns the coefficient of determination  $R^2$  of the prediction.

The coefficient  $R^2$  is defined as  $(1 - u/v)$ , where  $u$  is the residual sum of squares  $((y\_true - y\_pred) ** 2).sum()$  and  $v$  is the total sum of squares  $((y\_true - y\_true.mean()) ** 2).sum()$ . The best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse). A constant model that always predicts the expected value of  $y$ , disregarding the input features, would get a  $R^2$  score of 0.0.

#### Parameters

- **X** (*array-like, shape = (n\_samples, n\_features)*) – Test samples. For some estimators this may be a precomputed kernel matrix instead, *shape = (n\_samples, n\_samples\_fitted]*, where *n\_samples\_fitted* is the number of samples used in the fitting for the estimator.
- **y** (*array-like, shape = (n\_samples) or (n\_samples, n\_outputs)*) – True values for X.
- **sample\_weight** (*array-like, shape = [n\_samples], optional*) – Sample weights.

**Returns score** –  $R^2$  of `self.predict(X)` wrt. `y`.

**Return type** float

**set\_params** (*\*\*params*)

Set the parameters of this estimator.

The method works on simple estimators as well as on nested objects (such as pipelines). The latter have parameters of the form `<component>__<parameter>` so that it's possible to update each component of a nested object.

#### Returns

**Return type** self

**Warning:** Pygbm's API and default values are likely to be changed in future version, without any deprecation cycle.

This module contains the loss classes.

Specific losses are used for regression, binary classification or multiclass classification.

**class** `pygbm.loss.BinaryCrossEntropy`  
Binary cross-entropy loss, for binary classification.

For a given sample  $x_i$ , the binary cross-entropy loss is defined as the negative log-likelihood of the model which can be expressed as:

$$\text{loss}(x_i) = \log(1 + \exp(\text{raw\_pred}_i)) - y_{\text{true}_i} * \text{raw\_pred}_i$$

See The Elements of Statistical Learning, by Hastie, Tibshirani, Friedman.

**get\_baseline\_prediction** (*y\_train*, *prediction\_dim*)  
Return initial predictions (before the first iteration).

**Parameters**

- **y\_train** (*array-like*, *shape=(n\_samples,)*) – The target training values.
- **prediction\_dim** (*int*) – The dimension of one prediction: 1 for binary classification and regression, *n\_classes* for multiclass classification.

**Returns** **baseline\_prediction** – The baseline prediction.

**Return type** float or array of shape (1, *prediction\_dim*)

**init\_gradients\_and\_hessians** (*n\_samples*, *prediction\_dim*)  
Return initial gradients and hessians.

Unless hessians are constant, arrays are initialized with undefined values.

**Parameters**

- **n\_samples** (*int*) – The number of samples passed to *fit()*

- **prediction\_dim** (*int*) – The dimension of a raw prediction, i.e. the number of trees built at each iteration. Equals 1 for regression and binary classification, or K where K is the number of classes for multiclass classification.

#### Returns

- **gradients** (*array-like, shape=(n\_samples \* prediction\_dim)*)
- **hessians** (*array-like, shape=(n\_samples \* prediction\_dim).*) – If hessians are constant (e.g. for LeastSquares loss, shape is (1,)) and the array is initialized to 1.

**inverse\_link\_function** = <ufunc 'expit'>

**update\_gradients\_and\_hessians** (*gradients, hessians, y\_true, raw\_predictions*)

Update gradients and hessians arrays, inplace.

The gradients (resp. hessians) are the first (resp. second) order derivatives of the loss for each sample with respect to the predictions of model, evaluated at iteration  $i - 1$ .

#### Parameters

- **gradients** (*array-like, shape=(n\_samples \* prediction\_dim)*) – The gradients (treated as OUT array).
- **hessians** (*array-like, shape=(n\_samples \* prediction\_dim) or (1,)*) – The hessians (treated as OUT array).
- **y\_true** (*array-like, shape=(n\_samples,)*) – The true target values or each training sample.
- **raw\_predictions** (*array-like, shape=(n\_samples, prediction\_dim)*) – The raw\_predictions (i.e. values from the trees) of the tree ensemble at iteration  $i - 1$ .

**class** pygbm.loss.CategoricalCrossEntropy

Categorical cross-entropy loss, for multiclass classification.

For a given sample  $x_i$ , the categorical cross-entropy loss is defined as the negative log-likelihood of the model and generalizes the binary cross-entropy to more than 2 classes.

**get\_baseline\_prediction** (*y\_train, prediction\_dim*)

Return initial predictions (before the first iteration).

#### Parameters

- **y\_train** (*array-like, shape=(n\_samples,)*) – The target training values.
- **prediction\_dim** (*int*) – The dimension of one prediction: 1 for binary classification and regression, n\_classes for multiclass classification.

**Returns** **baseline\_prediction** – The baseline prediction.

**Return type** float or array of shape (1, prediction\_dim)

**init\_gradients\_and\_hessians** (*n\_samples, prediction\_dim*)

Return initial gradients and hessians.

Unless hessians are constant, arrays are initialized with undefined values.

#### Parameters

- **n\_samples** (*int*) – The number of samples passed to *fit()*
- **prediction\_dim** (*int*) – The dimension of a raw prediction, i.e. the number of trees built at each iteration. Equals 1 for regression and binary classification, or K where K is the number of classes for multiclass classification.

**Returns**

- **gradients** (*array-like, shape=(n\_samples \* prediction\_dim)*)
- **hessians** (*array-like, shape=(n\_samples \* prediction\_dim).*) – If hessians are constant (e.g. for LeastSquares loss, shape is (1,)) and the array is initialized to 1.

**update\_gradients\_and\_hessians** (*gradients, hessians, y\_true, raw\_predictions*)

Update gradients and hessians arrays, inplace.

The gradients (resp. hessians) are the first (resp. second) order derivatives of the loss for each sample with respect to the predictions of model, evaluated at iteration  $i - 1$ .

**Parameters**

- **gradients** (*array-like, shape=(n\_samples \* prediction\_dim)*) – The gradients (treated as OUT array).
- **hessians** (*array-like, shape=(n\_samples \* prediction\_dim) or (1,)*) – The hessians (treated as OUT array).
- **y\_true** (*array-like, shape=(n\_samples,)*) – The true target values or each training sample.
- **raw\_predictions** (*array-like, shape=(n\_samples, prediction\_dim)*) – The raw\_predictions (i.e. values from the trees) of the tree ensemble at iteration  $i - 1$ .

**class** pygbm.loss.LeastSquares

Least squares loss, for regression.

For a given sample  $x_i$ , least squares loss is defined as:

$$\text{loss}(x_i) = (y_{\text{true}_i} - \text{raw\_pred}_i)**2$$

**get\_baseline\_prediction** (*y\_train, prediction\_dim*)

Return initial predictions (before the first iteration).

**Parameters**

- **y\_train** (*array-like, shape=(n\_samples,)*) – The target training values.
- **prediction\_dim** (*int*) – The dimension of one prediction: 1 for binary classification and regression,  $n_{\text{classes}}$  for multiclass classification.

**Returns** **baseline\_prediction** – The baseline prediction.

**Return type** float or array of shape (1, prediction\_dim)

**init\_gradients\_and\_hessians** (*n\_samples, prediction\_dim*)

Return initial gradients and hessians.

Unless hessians are constant, arrays are initialized with undefined values.

**Parameters**

- **n\_samples** (*int*) – The number of samples passed to *fit()*
- **prediction\_dim** (*int*) – The dimension of a raw prediction, i.e. the number of trees built at each iteration. Equals 1 for regression and binary classification, or  $K$  where  $K$  is the number of classes for multiclass classification.

**Returns**

- **gradients** (*array-like, shape=(n\_samples \* prediction\_dim)*)

- **hessians** (*array-like, shape=(n\_samples \* prediction\_dim).*) – If hessians are constant (e.g. for LeastSquares loss, shape is (1,)) and the array is initialized to 1.

**update\_gradients\_and\_hessians** (*gradients, hessians, y\_true, raw\_predictions*)

Update gradients and hessians arrays, inplace.

The gradients (resp. hessians) are the first (resp. second) order derivatives of the loss for each sample with respect to the predictions of model, evaluated at iteration  $i - 1$ .

#### Parameters

- **gradients** (*array-like, shape=(n\_samples \* prediction\_dim)*) – The gradients (treated as OUT array).
- **hessians** (*array-like, shape=(n\_samples \* prediction\_dim) or (1,)*) – The hessians (treated as OUT array).
- **y\_true** (*array-like, shape=(n\_samples,)*) – The true target values or each training sample.
- **raw\_predictions** (*array-like, shape=(n\_samples, prediction\_dim)*) – The raw\_predictions (i.e. values from the trees) of the tree ensemble at iteration  $i - 1$ .



This module contains the `BinMapper` class.

`BinMapper` is used for mapping a real-valued dataset into integer-valued bins with equally-spaced thresholds.

**class** `pygbm.binning.BinMapper` (*max\_bins=256, subsample=100000, random\_state=None*)  
Transformer that maps a dataset into integer-valued bins.

The bins are created in a feature-wise fashion, with equally-spaced quantiles.

Large datasets are subsampled, but the feature-wise quantiles should remain stable.

If the number of unique values for a given feature is less than `max_bins`, then the unique values of this feature are used instead of the quantiles.

#### Parameters

- **max\_bins** (*int, optional (default=256)*) – The maximum number of bins to use. If for a given feature the number of unique values is less than `max_bins`, then those unique values will be used to compute the bin thresholds, instead of the quantiles.
- **subsample** (*int or None, optional (default=1e5)*) – If `n_samples > subsample`, then `sub_samples` samples will be randomly chosen to compute the quantiles. If `None`, the whole data is used.
- **random\_state** (*int or numpy.random.RandomState or None, optional (default=None)*) – Pseudo-random number generator to control the random sub-sampling. See [scikit-learn glossary](#).

**fit** (*X, y=None*)

Fit data `X` by computing the binning thresholds.

**Parameters** `x` (*array-like*) – The data to bin

**Returns** `self`

**Return type** `object`

**transform** (*X*)

Bin data `X`.

**Parameters** **x** (*array-like*) – The data to bin

**Returns** **X\_binned** – The binned data

**Return type** array-like

This module contains the `TreeGrower` class.

`TreeGrower` builds a regression tree fitting a Newton-Raphson step, based on the gradients and Hessians of the training data.

```
class pygbm.grower.TreeGrower(X_binned, gradients, hessians, max_leaf_nodes=None,
                               max_depth=None, min_samples_leaf=20, min_gain_to_split=0.0,
                               max_bins=256, n_bins_per_feature=None,
                               l2_regularization=0.0, min_hessian_to_split=0.001, shrinkage=1.0)
```

Tree grower class used to build a tree.

The tree is fitted to predict the values of a Newton-Raphson step. The splits are considered in a best-first fashion, and the quality of a split is defined in `splitting._split_gain`.

#### Parameters

- **`X_binned`** (*array-like of int, shape=(n\_samples, n\_features)*) – The binned input samples. Must be Fortran-aligned.
- **`gradients`** (*array-like, shape=(n\_samples,)*) – The gradients of each training sample. Those are the gradients of the loss w.r.t the predictions, evaluated at iteration `i - 1`.
- **`hessians`** (*array-like, shape=(n\_samples,)*) – The Hessians of each training sample. Those are the Hessians of the loss w.r.t the predictions, evaluated at iteration `i - 1`.
- **`max_leaf_nodes`** (*int or None, optional (default=None)*) – The maximum number of leaves for each tree. If `None`, there is no maximum limit.
- **`max_depth`** (*int or None, optional (default=None)*) – The maximum depth of each tree. The depth of a tree is the number of nodes to go from the root to the deepest leaf.
- **`min_samples_leaf`** (*int, optional (default=20)*) – The minimum number of samples per leaf.

- **min\_gain\_to\_split** (*float, optional (default=0.)*) – The minimum gain needed to split a node. Splits with lower gain will be ignored.
- **max\_bins** (*int, optional (default=256)*) – The maximum number of bins. Used to define the shape of the histograms.
- **n\_bins\_per\_feature** (*array-like of int or int, optional (default=None)*) – The actual number of bins needed for each feature, which is lower or equal to `max_bins`. If it's an int, all features are considered to have the same number of bins. If None, all features are considered to have `max_bins` bins.
- **l2\_regularization** (*float, optional (default=0)*) – The L2 regularization parameter.
- **min\_hessian\_to\_split** (*float, optional (default=1e-3)*) – The minimum sum of Hessians needed in each node. Splits that result in at least one child having a sum of Hessians less than `min_hessian_to_split` are discarded.
- **shrinkage** (*float, optional (default=1)*) – The shrinkage parameter to apply to the leaves values, also known as learning rate.

**can\_split\_further** ()

Return True if there are still nodes to split.

**grow** ()

Grow the tree, from root to leaves.

**make\_predictor** (*numerical\_thresholds=None*)

Make a `TreePredictor` object out of the current tree.

**Parameters** **numerical\_thresholds** (*array-like of floats, optional (default=None)*) – The actual thresholds values of each bin, expected to be in sorted increasing order. None if the training data was pre-binned.

**Returns**

**Return type** A `TreePredictor` object.

**split\_next** ()

Split the node with highest potential gain.

**Returns**

- **left** (*TreeNode*) – The resulting left child.
- **right** (*TreeNode*) – The resulting right child.

**class** `pygbm.grower.TreeNode` (*depth, sample\_indices, sum\_gradients, sum\_hessians, parent=None*)  
Tree Node class used in `TreeGrower`.

This isn't used for prediction purposes, only for training (see `TreePredictor`).

**Parameters**

- **depth** (*int*) – The depth of the node, i.e. its distance from the root
- **samples\_indices** (*array of int*) – The indices of the samples at the node
- **sum\_gradients** (*float*) – The sum of the gradients of the samples at the node
- **sum\_hessians** (*float*) – The sum of the Hessians of the samples at the node
- **parent** (*TreeNode or None, optional (default=None)*) – The parent of the node. None for root.

**depth**

*int* – The depth of the node, i.e. its distance from the root

**samples\_indices**

*array of int* – The indices of the samples at the node

**sum\_gradients**

*float* – The sum of the gradients of the samples at the node

**sum\_hessians**

*float* – The sum of the hessians of the samples at the node

**parent**

*TreeNode or None, optional(default=None)* – The parent of the node. None for root.

**split\_info**

*SplitInfo or None* – The result of the split evaluation

**left\_child**

*TreeNode or None* – The left child of the node. None for leaves.

**right\_child**

*TreeNode or None* – The right child of the node. None for leaves.

**value**

*float or None* – The value of the leaf, as computed in `finalize_leaf()`. None for non-leaf nodes

**find\_split\_time**

*float* – The total time spent computing the histogram and finding the best split at the node.

**construction\_speed**

*float* – The Number of samples at the node divided `find_split_time`.

**apply\_split\_time**

*float* – The total time spent actually splitting the node, e.g. splitting `samples_indices` into left and right child.

**hist\_subtraction**

*bool* – Whether the subtraction method was used for computing the histograms.



This module contains njitted routines and data structures to:

- Find the best possible split of a node. For a given node, a split is characterized by a feature and a bin.
- Apply a split to a node, i.e. split the indices of the samples at the node into the newly created left and right childs.

`pygbm.splitting.find_node_split`

For each feature, find the best bin to split on at a given node.

Returns the best split info among all features, and the histograms of all the features. The histograms are computed by scanning the whole data.

#### Parameters

- **context** (*SplittingContext*) – The splitting context
- **sample\_indices** (*array of int*) – The indices of the samples at the node to split.

#### Returns

- **best\_split\_info** (*SplitInfo*) – The info about the best possible split among all features.
- **histograms** (*array of HISTOGRAM\_DTYPE, shape=(n\_features, max\_bins)*) – The histograms of each feature. A histogram is an array of HISTOGRAM\_DTYPE of size `max_bins` (only `n_bins_per_features[feature]` entries are relevant).

`pygbm.splitting.find_node_split_subtraction`

For each feature, find the best bin to split on at a given node.

Returns the best split info among all features, and the histograms of all the features.

This does the same job as `find_node_split()` but uses the histograms of the parent and sibling of the node to split. This allows to use the identity: `histogram(parent) = histogram(node) - histogram(sibling)`, which is significantly faster than computing the histograms from data.

Returns the best `SplitInfo` among all features, along with all the feature histograms that can be latter used to compute the sibling or children histograms by subtraction.

#### Parameters

- **context** (*SplittingContext*) – The splitting context
- **sample\_indices** (*array of int*) – The indices of the samples at the node to split.
- **parent\_histograms** (*array of HISTOGRAM\_DTYPE of shape (n\_features, max\_bins)*) – The histograms of the parent
- **sibling\_histograms** (*array of HISTOGRAM\_DTYPE of shape (n\_features, max\_bins)*) – The histograms of the sibling

#### Returns

- **best\_split\_info** (*SplitInfo*) – The info about the best possible split among all features.
- **histograms** (*array of HISTOGRAM\_DTYPE, shape=(n\_features, max\_bins)*) – The histograms of each feature. A histogram is an array of HISTOGRAM\_DTYPE of size max\_bins (only n\_bins\_per\_features[feature] entries are relevant).

`pygbm.splitting.split_indices`

Split samples into left and right arrays.

#### Parameters

- **context** (*SplittingContext*) – The splitting context
- **split\_ingo** (*SplitInfo*) – The SplitInfo of the node to split
- **sample\_indices** (*array of int*) – The indices of the samples at the node to split. This is a view on context.partition, and it is modified inplace by placing the indices of the left child at the beginning, and the indices of the right child at the end.

#### Returns

- **left\_indices** (*array of int*) – The indices of the samples in the left child. This is a view on context.partition.
- **right\_indices** (*array of int*) – The indices of the samples in the right child. This is a view on context.partition.



## CHAPTER 6

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



**p**

`pygbm.binning`, 13  
`pygbm.gradient_boosting`, 3  
`pygbm.grower`, 15  
`pygbm.loss`, 9  
`pygbm.splitting`, 19



**A**

apply\_split\_time (pygbm.grower.TreeNode attribute), 17

**B**

BinaryCrossEntropy (class in pygbm.loss), 9

BinMapper (class in pygbm.binning), 13

**C**

can\_split\_further() (pygbm.grower.TreeGrower method), 16

CategoricalCrossEntropy (class in pygbm.loss), 10

construction\_speed (pygbm.grower.TreeNode attribute), 17

**D**

depth (pygbm.grower.TreeNode attribute), 16

**F**

find\_node\_split (in module pygbm.splitting), 19

find\_node\_split\_subtraction (in module pygbm.splitting), 19

find\_split\_time (pygbm.grower.TreeNode attribute), 17

fit() (pygbm.binning.BinMapper method), 13

fit() (pygbm.gradient\_boosting.GradientBoostingClassifier method), 4

fit() (pygbm.gradient\_boosting.GradientBoostingRegressor method), 7

**G**

get\_baseline\_prediction() (pygbm.loss.BinaryCrossEntropy method), 9

get\_baseline\_prediction() (pygbm.loss.CategoricalCrossEntropy method), 10

get\_baseline\_prediction() (pygbm.loss.LeastSquares method), 11

get\_params() (pygbm.gradient\_boosting.GradientBoostingClassifier method), 5

get\_params() (pygbm.gradient\_boosting.GradientBoostingRegressor method), 7

GradientBoostingClassifier (class in pygbm.gradient\_boosting), 3

GradientBoostingRegressor (class in pygbm.gradient\_boosting), 6

grow() (pygbm.grower.TreeGrower method), 16

**H**

hist\_subtraction (pygbm.grower.TreeNode attribute), 17

**I**

init\_gradients\_and\_hessians() (pygbm.loss.BinaryCrossEntropy method), 9

init\_gradients\_and\_hessians() (pygbm.loss.CategoricalCrossEntropy method), 10

init\_gradients\_and\_hessians() (pygbm.loss.LeastSquares method), 11

inverse\_link\_function (pygbm.loss.BinaryCrossEntropy attribute), 10

**L**

LeastSquares (class in pygbm.loss), 11

left\_child (pygbm.grower.TreeNode attribute), 17

**M**

make\_predictor() (pygbm.grower.TreeGrower method), 16

**P**

parent (pygbm.grower.TreeNode attribute), 17

predict() (pygbm.gradient\_boosting.GradientBoostingClassifier method), 5

predict() (pygbm.gradient\_boosting.GradientBoostingRegressor method), 7

predict\_proba() (pygbm.gradient\_boosting.GradientBoostingClassifier method), 5

pygbm.binning (module), 13  
pygbm.gradient\_boosting (module), 3  
pygbm.grower (module), 15  
pygbm.loss (module), 9  
pygbm.splitting (module), 19

## R

right\_child (pygbm.grower.TreeNode attribute), 17

## S

samples\_indices (pygbm.grower.TreeNode attribute), 17  
score() (pygbm.gradient\_boosting.GradientBoostingClassifier method), 5  
score() (pygbm.gradient\_boosting.GradientBoostingRegressor method), 7  
set\_params() (pygbm.gradient\_boosting.GradientBoostingClassifier method), 5  
set\_params() (pygbm.gradient\_boosting.GradientBoostingRegressor method), 8  
split\_indices (in module pygbm.splitting), 20  
split\_info (pygbm.grower.TreeNode attribute), 17  
split\_next() (pygbm.grower.TreeGrower method), 16  
sum\_gradients (pygbm.grower.TreeNode attribute), 17  
sum\_hessians (pygbm.grower.TreeNode attribute), 17

## T

transform() (pygbm.binning.BinMapper method), 13  
TreeGrower (class in pygbm.grower), 15  
TreeNode (class in pygbm.grower), 16

## U

update\_gradients\_and\_hessians()  
(pygbm.loss.BinaryCrossEntropy method),  
10  
update\_gradients\_and\_hessians()  
(pygbm.loss.CategoricalCrossEntropy  
method), 11  
update\_gradients\_and\_hessians()  
(pygbm.loss.LeastSquares method), 12

## V

value (pygbm.grower.TreeNode attribute), 17