
isotree Documentation

David Cortes

Jun 14, 2022

CONTENTS

1	Installation	3
2	Introduction to the library and methods	5
3	Quick example notebooks	7
4	Methods	9
5	IsolationForest	11
6	Indices and tables	43
	Index	45

This is the documentation page for the *isotree* Python package, which provides Isolation Forest models for outlier/anomaly detection and other purposes. See project's GitHub for more details:

<https://www.github.com/david-cortes/isotree>

For the R version, see the CRAN webpage:

<https://cran.r-project.org/web/packages/isotree/index.html>

INSTALLATION

The Python version of this package can be easily installed from PyPI

```
pip install isotree
```

(See the GitHub page for more details, esp. section “Reducing library size and compilation times”)

Note that it is only available in source form (not in binary wheel form), which means you will need a toolchain for compiling C++ source code (e.g. GCC in linux, msys2 that comes with anaconda on windows, clang in mac).

INTRODUCTION TO THE LIBRARY AND METHODS

- An introduction to Isolation Forests.

QUICK EXAMPLE NOTEBOOKS

- An introduction to Isolation Forests.
- General library usage.
- As missing value imputer.
- As kernel for SVMs.
- Converting to treelite for faster predictions.

METHODS

- *IsolationForest*
- *append_trees*
- *build_indexer*
- *copy*
- *decision_function*
- *drop_imputer*
- *drop_indexer*
- *drop_reference_points*
- *export_model*
- *fit*
- *fit_predict*
- *fit_transform*
- *generate_sql*
- *get_num_nodes*
- *get_params*
- *import_model*
- *partial_fit*
- *predict*
- *predict_distance*
- *predict_kernel*
- *set_params*
- *set_reference_points*
- *subset_trees*
- *to_treelite*
- *transform*

ISOLATIONFOREST

```
class isotree.IsolationForest(sample_size='auto', ntrees=500, ndim=3, ntry=1, categ_cols=None,  
max_depth='auto', ncols_per_tree=None, prob_pick_pooled_gain=0.0,  
prob_pick_avg_gain=0.0, prob_pick_full_gain=0.0, prob_pick_dens=0.0,  
prob_pick_col_by_range=0.0, prob_pick_col_by_var=0.0,  
prob_pick_col_by_kurt=0.0, min_gain=0.0, missing_action='auto',  
new_categ_action='auto', categ_split_type='auto', all_perm=False,  
coef_by_prop=False, recode_categ=False, weights_as_sample_prob=True,  
sample_with_replacement=False, penalize_range=False,  
standardize_data=True, scoring_metric='depth', fast_bratio=True,  
weigh_by_kurtosis=False, coefs='uniform', assume_full_distr=True,  
build_imputer=False, min_imp_obs=3, depth_imp='higher',  
weigh_imp_rows='inverse', random_seed=1, use_long_double=False,  
nthreads=- 1, n_estimators=None, max_samples=None, n_jobs=None,  
random_state=None, bootstrap=None)
```

Bases: object

Isolation Forest model

Isolation Forest is an algorithm originally developed for outlier detection that consists in splitting sub-samples of the data according to some attribute/feature/column at random. The idea is that, the rarer the observation, the more likely it is that a random uniform split on some feature would put outliers alone in one branch, and the fewer splits it will take to isolate an outlier observation like this. The concept is extended to splitting hyperplanes in the extended model (i.e. splitting by more than one column at a time), and to guided (not entirely random) splits in the SCiForest and FCF models that aim at isolating outliers faster and/or finding clustered outliers.

This version adds heuristics to handle missing data and categorical variables. Can be used to approximate pairwise distances by checking the depth after which two observations become separated, and to approximate densities by fitting trees beyond balanced-tree limit. Offers options to vary between randomized and deterministic splits too.

Note: The default parameters in this software do not correspond to the suggested parameters in any of the references. In particular, the following default values are likely to cause huge differences when compared to the defaults in other software: `ndim`, `sample_size`, `ntrees`. The defaults here are nevertheless more likely to result in better models. In order to mimic scikit-learn for example, one would need to pass `ndim=1`, `sample_size=256`, `ntrees=100`, `missing_action="fail"`, `nthreads=1`.

Note: Shorthands for parameter combinations that match some of the references:

'iForest' (reference¹):

`ndim=1, sample_size=256, max_depth=8, ntrees=100, missing_action="fail"`.

¹ Liu, Fei Tony, Kai Ming Ting, and Zhi-Hua Zhou. "Isolation forest." 2008 Eighth IEEE International Conference on Data Mining. IEEE, 2008.

'EIF' (reference^{Page 12, 3}):

`ndim=2, sample_size=256, max_depth=8, ntrees=100, missing_action="fail",`
`coefs="uniform", standardize_data=False` (plus standardizing the data **before** passing it).

'SCiForest' (reference⁴):

`ndim=2, sample_size=256, max_depth=8, ntrees=100, missing_action="fail",`
`coefs="normal", ntry=10, prob_pick_avg_gain=1, penalize_range=True`. Might provide much better results with `max_depth=None` despite the reference's recommendation.

'FCF' (reference¹¹):

`ndim=2, sample_size=256, max_depth=None, ntrees=200, missing_action="fail",`
`coefs="normal", ntry=1, prob_pick_pooled_gain=1`. Might provide similar or better results with `ndim=1` and/or sample size as low as 32. For the FCF model aimed at imputing missing values, might give better results with `ntry=10` or higher and much larger sample sizes.

'RRCF' (reference¹²):

`ndim=1, prob_pick_col_by_range=1, sample_size=256 or more, max_depth=None, ntrees=100`
or more, `missing_action="fail"`. Note however that reference¹² proposed a different method for calculation of anomaly scores, while this library uses isolation depth just like for 'iForest', so results might differ significantly from those of other libraries. Nevertheless, experiments in reference¹¹ suggest that isolation depth might be a better scoring metric for this model.

Note: The model offers many tunable parameters (see reference¹¹ for a comparison). The most likely candidate to tune is `prob_pick_pooled_gain`, for which higher values tend to result in a better ability to flag outliers in multimodal datasets, at the expense of poorer generalizability to inputs with values outside the variables' ranges to which the model was fit (see plots generated from the examples in GitHub notebook for a better idea of the difference). The next candidate to tune is `sample_size` - the default is to use all rows, but in some datasets introducing sub-sampling can help, especially for the single-variable model. In smaller datasets, one might also want to experiment with `weigh_by_kurtosis` and perhaps lower `ndim`. If using `prob_pick_pooled_gain`, models are likely to benefit from deeper trees (controlled by `max_depth`), but using large samples and/or deeper trees can result in significantly slower model fitting and predictions - in such cases, using `min_gain` (with a value like 0.25) with `max_depth=None` can offer a better speed/performance trade-off than changing `max_depth`.

If the data has categorical variables and these are more important for determining outlieriness compared to numerical columns, one might want to experiment with `ndim=1, categ_split_type="single_categ", and scoring_metric="density"`.

For small datasets, one might also want to experiment with `ndim=1, scoring_metric="adj_depth" and penalize_range=True`.

Note: The default parameters will not scale to large datasets. In particular, if the amount of data is large, it's suggested to set a smaller sample size for each tree (parameter `sample_size`) and to fit fewer of them (parameter `ntrees`). As well, the default option for 'missing_action' might slow things down significantly. See the documentation of the parameters for more details. These defaults can also result in very big model sizes in memory and as serialized files (e.g. models that weight over 10GB) when the number of rows in the data is large. Using fewer trees, smaller sample sizes, and shallower trees can help to reduce model sizes if that becomes a problem.

³ Hariri, Sahand, Matias Carrasco Kind, and Robert J. Brunner. "Extended Isolation Forest." arXiv preprint arXiv:1811.02141 (2018).

⁴ Liu, Fei Tony, Kai Ming Ting, and Zhi-Hua Zhou. "On detecting clustered anomalies using SCiForest." Joint European Conference on Machine Learning and Knowledge Discovery in Databases. Springer, Berlin, Heidelberg, 2010.

¹¹ Cortes, David. "Revisiting randomized choices in isolation forests." arXiv preprint arXiv:2110.13402 (2021).

¹² Guha, Sudipto, et al. "Robust random cut forest based anomaly detection on streams." International conference on machine learning. PMLR, 2016.

Note: See the documentation of `predict` for some considerations when serving models generated through this library.

Parameters

- **sample_size** (*str* “auto”, *int*, *float*(0,1), or *None*) – Sample size of the data sub-samples with which each binary tree will be built. If passing ‘None’, each tree will be built using the full data. Recommended value in^{Page 11, 1},²,^{Page 12, 3} is 256, while the default value in the author’s code in⁵ is ‘None’ here.

If passing “auto”, will use the full number of rows in the data, up to 10,000 (i.e. will take ‘sample_size=min(nrows(X), 10000)’) **when calling fit**, and the full amount of rows in the data **when calling the variants fit_predict** or **fit_transform**.

If passing None, will take the full number of rows in the data (no sub-sampling).

If passing a number between zero and one, will assume it means taking a sample size that represents that proportion of the rows in the data.

Hint: seeing a distribution of scores which is on average too far below 0.5 could mean that the model needs more trees and/or bigger samples to reach convergence (unless using non-random splits, in which case the distribution is likely to be centered around a much lower number), or that the distributions in the data are too skewed for random uniform splits.

- **ntrees** (*int*) – Number of binary trees to build for the model. Recommended value in^{Page 11, 1} is 100, while the default value in the author’s code in⁵ is 10. In general, the number of trees required for good results is higher when (a) there are many columns, (b) there are categorical variables, (c) categorical variables have many categories, (d) *ndim* is high, (e) `prob_pick_pooled_gain` is used, (f) `scoring_metric="density"` or `scoring_metric="boxed_density"` are used.

Hint: seeing a distribution of scores which is on average too far below 0.5 could mean that the model needs more trees and/or bigger samples to reach convergence (unless using non-random splits, in which case the distribution is likely to be centered around a much lower number), or that the distributions in the data are too skewed for random uniform splits.

- **ndim** (*int*) – Number of columns to combine to produce a split. If passing 1, will produce the single-variable model described in^{Page 11, 1} and², while if passing values greater than 1, will produce the extended model described in^{Page 12, 3} and^{Page 12, 4}. Recommended value in^{Page 12, 4} is 2, while^{Page 12, 3} recommends a low value such as 2 or 3. Models with values higher than 1 are referred hereafter as the extended model (as in^{Page 12, 3}).

Note that, when using `ndim>1` plus `standardize_data=True`, the variables are standardized at each step as suggested in^{Page 12, 4}, which makes the models slightly different than in^{Page 12, 3}.

In general, when the data has categorical variables, models with `ndim=1` plus `categ_split_type="single_categ"` tend to produce better results, while models `ndim>1` tend to produce better results for numerical-only data, especially in the presence of missing values.

- **ntry** (*int*) – When using any of `prob_pick_pooled_gain`, `prob_pick_avg_gain`, `prob_pick_full_gain`, `prob_pick_dens`, how many variables (with `ndim=1`) or linear combinations (with `ndim>1`) to try for determining the best one according to gain.

² Liu, Fei Tony, Kai Ming Ting, and Zhi-Hua Zhou. “Isolation-based anomaly detection.” *ACM Transactions on Knowledge Discovery from Data* (TKDD) 6.1 (2012): 3.

⁵ <https://sourceforge.net/projects/iforest/>

Recommended value in reference^{Page 12, 4} is 10 (with `prob_pick_avg_gain`, for outlier detection), while the recommended value in reference^{Page 12, 11} is 1 (with `prob_pick_pooled_gain`, for outlier detection), and the recommended value in reference⁹ is 10 to 20 (with `prob_pick_pooled_gain`, for missing value imputations).

- **`categ_cols`** (*None or array-like*) – Columns that hold categorical features, when the data is passed as an array or matrix. Categorical columns should contain only integer values with a continuous numeration starting at zero, with negative values and NaN taken as missing, and the array or list passed here should correspond to the column numbers, with numeration starting at zero. The maximum categorical value should not exceed 'INT_MAX' (typically $2^{31} - 1$). This might be passed either at construction time or when calling `fit` or variations of `fit`.

This is ignored when the input is passed as a `DataFrame` as then it will consider columns as categorical depending on their dtype (see the documentation for `fit` for details).

- **`max_depth`** (*int, None, or str "auto"*) – Maximum depth of the binary trees to grow. If passing `None`, will build trees until each observation ends alone in a terminal node or until no further split is possible. If using "auto", will limit it to the corresponding depth of a balanced binary tree with number of terminal nodes corresponding to the sub-sample size (the reason being that, if trying to detect outliers, an outlier will only be so if it turns out to be isolated with shorter average depth than usual, which corresponds to a balanced tree depth). When a terminal node has more than 1 observation, the remaining isolation depth for them is estimated assuming the data and splits are both uniformly random (separation depth follows a similar process with expected value calculated as in⁶). Default setting for^{Page 11, 1}, ^{Page 13, 2}, ^{Page 12, 3}, ^{Page 12, 4} is "auto", but it's recommended to pass higher values if using the model for purposes other than outlier detection.

Note that models that use `prob_pick_pooled_gain` or `prob_pick_avg_gain` are likely to benefit from deeper trees (larger `max_depth`), but deeper trees can result in much slower model fitting and predictions.

If using pooled gain, one might want to substitute `max_depth` with `min_gain`.

- **`ncols_per_tree`** (*None, int, or float(0,1]*) – Number of columns to use (have as potential candidates for splitting at each iteration) in each tree, somewhat similar to the 'mtry' parameter of random forests. In general, this is only relevant when using non-random splits and/or weighted column choices.

If passing a number between zero and one, will assume it means taking a sample size that represents that proportion of the columns in the data. If passing exactly 1, will assume it means taking 100% of the columns rather than taking 1 column.

If passing `None` (the default) or zero, will use the full number of available columns.

- **`prob_pick_pooled_gain`** (*float[0, 1]*) – This parameter indicates the probability of choosing the threshold on which to split a variable (with `ndim=1`) or a linear combination of variables (when using `ndim>1`) as the threshold that maximizes a pooled standard deviation gain criterion (see references⁹ and^{Page 12, 11}) on the same variable or linear combination, similarly to regression trees such as CART.

If using `ntry>1`, will try several variables or linear combinations thereof and choose the one in which the largest standardized gain can be achieved.

For categorical variables with `ndim=1`, will use shannon entropy instead (like in⁷).

⁹ Cortes, David. "Imputing missing values with unsupervised random trees." arXiv preprint arXiv:1911.06646 (2019).

⁶ <https://math.stackexchange.com/questions/3388518/expected-number-of-paths-required-to-separate-elements-in-a-binary-tree>

⁷ Quinlan, J. Ross. C4. 5: programs for machine learning. Elsevier, 2014.

Compared to a simple averaged gain, this tends to result in more evenly-divided splits and more clustered groups when they are smaller. Recommended to pass higher values when used for imputation of missing values. When used for outlier detection, datasets with multimodal distributions usually see better performance under this type of splits.

Note that, since this makes the trees more even and thus it takes more steps to produce isolated nodes, the resulting object will be heavier. When splits are not made according to any of `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`, both the column and the split point are decided at random. Note that, if passing value 1 (100%) with no sub-sampling and using the single-variable model, every single tree will have the exact same splits.

Be aware that `penalize_range` can also have a large impact when using `prob_pick_pooled_gain`.

Under this option, models are likely to produce better results when increasing `max_depth`. Alternatively, one can also control the depth through `min_gain` (for which one might want to set `max_depth=None`).

Important detail: if using any of `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`, the distribution of outlier scores is unlikely to be centered around 0.5.

- **`prob_pick_avg_gain`** (*float*[0, 1]) – This parameter indicates the probability of choosing the threshold on which to split a variable (with `ndim=1`) or a linear combination of variables (when using `ndim>1`) as the threshold that maximizes an averaged standard deviation gain criterion (see references [Page 12, 4](#) and [Page 12, 11](#)) on the same variable or linear combination.

If using `ntry>1`, will try several variables or linear combinations thereof and choose the one in which the largest standardized gain can be achieved.

For categorical variables with `ndim=1`, will take the expected standard deviation that would be gotten if the column were converted to numerical by assigning to each category a random number $\sim \text{Unif}(0, 1)$ and calculate gain with those assumed standard deviations.

Compared to a pooled gain, this tends to result in more cases in which a single observation or very few of them are put into one branch. Typically, datasets with outliers defined by extreme values in some column more or less independently of the rest, usually see better performance under this type of split. Recommended to use sub-samples (parameter `sample_size`) when passing this parameter. Note that, since this will create isolated nodes faster, the resulting object will be lighter (use less memory).

When splits are not made according to any of `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`, both the column and the split point are decided at random. Default setting for [Page 11, 1](#), [Page 13, 2](#), [Page 12, 3](#) is zero, and default for [Page 12, 4](#) is 1. This is the randomization parameter that can be passed to the author's original code in [Page 13, 5](#), but note that the code in [Page 13, 5](#) suffers from a mathematical error in the calculation of running standard deviations, so the results from it might not match with this library's.

Be aware that, if passing a value of 1 (100%) with no sub-sampling and using the single-variable model, every single tree will have the exact same splits.

Under this option, models are likely to produce better results when increasing `max_depth`.

Important detail: if using any of `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`, the distribution of outlier scores is unlikely to be centered around 0.5.

- **prob_pick_full_gain** (*float[0,1]*) – This parameter indicates the probability of choosing the threshold on which to split a variable (with `ndim=1`) or a linear combination of variables (when using `ndim>1`) as the threshold that minimizes the pooled sums of variances of all columns (or a subset of them if using `ncols_per_tree`).

In general, this is much slower to evaluate than the other gain types, and does not tend to lead to better results. When using this option, one might want to use a different scoring metric (particulary "density", "boxed_density2" or "boxed_ratio"). Note that the calculations are all done through the (exact) sorted-indices approach, while is much slower than the (approximate) histogram approach used by other decision tree software.

Be aware that the data is not standardized in any way for the variance calculations, thus the scales of features will make a large difference under this option, which might not make it suitable for all types of data.

This option is not compatible with categorical data, and `min_gain` does not apply to it.

When splits are not made according to any of `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`, both the column and the split point are decided at random. Default setting for [Page 11, 1](#), [Page 13, 2](#), [Page 12, 3](#), [Page 12, 4](#) is zero.

- **prob_pick_dens** (*float[0,1]*) – This parameter indicates the probability of choosing the threshold on which to split a variable (with `ndim=1`) or a linear combination of variables (when using `ndim>1`) as the threshold that maximizes the pooled densities of the branch distributions.

The `min_gain` option does not apply to this type of splits.

When splits are not made according to any of `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`, both the column and the split point are decided at random. Default setting for [Page 11, 1](#), [Page 13, 2](#), [Page 12, 3](#), [Page 12, 4](#) is zero.

- **prob_pick_col_by_range** (*float[0, 1]*) – When using `ndim=1`, this denotes the probability of choosing the column to split with a probability proportional to the range spanned by each column within a node as proposed in reference [Page 12, 12](#).

When using `ndim>1`, this denotes the probability of choosing columns to create a hyperplane with a probability proportional to the range spanned by each column within a node.

This option is not compatible with categorical data. If passing column weights, the effect will be multiplicative.

Be aware that the data is not standardized in any way for the range calculations, thus the scales of features will make a large difference under this option, which might not make it suitable for all types of data.

If there are infinite values, all columns having infinite values will be treated as having the same weight, and will be chosen before every other column with non-infinite values.

Note that the proposed RRCF model from [Page 12, 12](#) uses a different scoring metric for producing anomaly scores, while this library uses isolation depth regardless of how columns are chosen, thus results are likely to be different from those of other software implementations. Nevertheless, as explored in [Page 12, 11](#), isolation depth as a scoring metric typically provides better results than the “co-displacement” metric from [Page 12, 12](#) under these split types.

- **prob_pick_col_by_var** (*float[0, 1]*) – When using `ndim=1`, this denotes the probability of choosing the column to split with a probability proportional to the variance of each column within a node.

When using `ndim>1`, this denotes the probability of choosing columns to create a hyperplane with a probability proportional to the variance of each column within a node.

For categorical data, it will calculate the expected variance if the column were converted to numerical by assigning to each category a random number $\sim \text{Unif}(0, 1)$, which depending on the number of categories and their distribution, produces numbers typically a bit smaller than standardized numerical variables.

Note that when using sparse matrices, the calculation of variance will rely on a procedure that uses sums of squares, which has less numerical precision than the calculation used for dense inputs, and as such, the results might differ slightly.

Be aware that this calculated variance is not standardized in any way, so the scales of features will make a large difference under this option.

If passing column weights, the effect will be multiplicative.

If passing a `missing_action` different than “fail”, infinite values will be ignored for the variance calculation. Otherwise, all columns with infinite values will have the same probability and will be chosen before columns with non-infinite values.

- **prob_pick_col_by_kurt** (*float*[0, 1]) – When using `ndim=1`, this denotes the probability of choosing the column to split with a probability proportional to the kurtosis of each column **within a node** (unlike the option `weigh_by_kurtosis` which calculates this metric only at the root).

When using `ndim>1`, this denotes the probability of choosing columns to create a hyperplane with a probability proportional to the kurtosis of each column within a node.

For categorical data, it will calculate the expected kurtosis if the column were converted to numerical by assigning to each category a random number $\sim \text{Unif}(0, 1)$.

Note that when using sparse matrices, the calculation of kurtosis will rely on a procedure that uses sums of squares and higher-power numbers, which has less numerical precision than the calculation used for dense inputs, and as such, the results might differ slightly.

If passing column weights, the effect will be multiplicative. This option is not compatible with `weigh_by_kurtosis`.

If passing a `missing_action` different than “fail”, infinite values will be ignored for the kurtosis calculation. Otherwise, all columns with infinite values will have the same probability and will be chosen before columns with non-infinite values.

If using `missing_action="impute"`, the calculation of kurtosis will not use imputed values in order not to favor columns with missing values (which would increase kurtosis by all having the same central value).

Be aware that kurtosis can be a rather slow metric to calculate.

- **min_gain** (*float* > 0) – Minimum gain that a split threshold needs to produce in order to proceed with a split. Only used when the splits are decided by a variance gain criterion (`prob_pick_pooled_gain` or `prob_pick_avg_gain`, but not `prob_pick_full_gain` nor `prob_pick_dens`). If the highest possible gain in the evaluated splits at a node is below this threshold, that node becomes a terminal node.

This can be used as a more sophisticated depth control when using pooled gain (note that `max_depth` still applies on top of this heuristic).

- **missing_action** (*str*, one of “divide” (single-variable only), “impute”, “fail”, “auto”) – How to handle missing data at both fitting and prediction time. Options are:

"divide":

(For the single-variable model only, recommended) Will follow both branches and combine the result with the weight given by the fraction of the data that went to each branch when fitting the model.

"impute":

Will assign observations to the branch with the most observations in the single-variable model, or fill in missing values with the median of each column of the sample from which the split was made in the extended model (recommended for the extended model) (but note that the calculation of medians does not take into account sample weights when using `weights_as_sample_prob=False`). When using `ndim=1`, gain calculations will use median-imputed values for missing data under this option.

"fail":

Will assume there are no missing values and will trigger undefined behavior if it encounters any.

"auto":

Will use “divide” for the single-variable model and “impute” for the extended model.

In the extended model, infinite values will be treated as missing. Passing “fail” will produce faster fitting and prediction times along with decreased model object sizes.

Models from [Page 11, 1](#), [Page 13, 2](#), [Page 12, 3](#), [Page 12, 4](#) correspond to “fail” here.

Typically, models with `'ndim>1'` are less affected by missing data than models with `'ndim=1'`.

- **new_categ_action** (*str*, one of “weighted” (single-variable only), “impute” (extended only), “smallest”, “random”) – What to do after splitting a categorical feature when new data that reaches that split has categories that the sub-sample from which the split was done did not have. Options are:

"weighted":

(For the single-variable model only, recommended) Will follow both branches and combine the result with weight given by the fraction of the data that went to each branch when fitting the model.

"impute":

(For the extended model only, recommended) Will assign them the median value for that column that was added to the linear combination of features (but note that this median calculation does not use sample weights when using `weights_as_sample_prob=False`).

"smallest":

In the single-variable case will assign all observations with unseen categories in the split to the branch that had fewer observations when fitting the model, and in the extended case will assign them the coefficient of the least common category.

"random":

Will assign a branch (coefficient in the extended model) at random for each category beforehand, even if no observations had that category when fitting the model. Note that this can produce biased results when deciding splits by a gain criterion.

Important: under this option, if the model is fitted to a `DataFrame`, when calling `predict` on new data which contains new categories (unseen in the data to which the model was fitted), they will be added to the model’s state on-the-fly. This means that, if calling `predict` on data which has new categories, there might be inconsistencies in the results if predictions are done in parallel or if passing the same data in batches or with different row orders. It also means that the `predict` function will not be thread-safe (e.g. cannot be used alongside `joblib` with a backend that uses shared memory).

"auto":

Will select “weighted” for the single-variable model and “impute” for the extended model.

Ignored when passing ‘categ_split_type’ = ‘single_categ’.

- **categ_split_type** (*str*; one of “auto”, “subset”, or “single_categ”) – Whether to split categorical features by assigning sub-sets of them to each branch, or by assigning a single category to a branch and the rest to the other branch. For the extended model, whether to give each category a coefficient, or only one while the rest get zero.

If passing "auto", will select "subset" for the extended model and "single_categ" for the single-variable model.

- **all_perm** (*bool*) – When doing categorical variable splits by pooled gain with `ndim=1` (single-variable model), whether to consider all possible permutations of variables to assign to each branch or not. If `False`, will sort the categories by their frequency and make a grouping in this sorted order. Note that the number of combinations evaluated (if `True`) is the factorial of the number of present categories in a given column (minus 2). For averaged gain, the best split is always to put the second most-frequent category in a separate branch, so not evaluating all permutations (passing `False`) will make it possible to select other splits that respect the sorted frequency order. Ignored when not using categorical variables or not doing splits by pooled gain or using `ndim > 1`.
- **coef_by_prop** (*bool*) – In the extended model, whether to sort the randomly-generated coefficients for categories according to their relative frequency in the tree node. This might provide better results when using categorical variables with too many categories, but is not recommended, and not reflective of real “categorical-ness”. Ignored for the single-variable model (`ndim=1`) and/or when not using categorical variables.
- **recode_categ** (*bool*) – Whether to re-encode categorical variables even in case they are already passed as `pd.Categorical`. This is recommended as it will eliminate potentially redundant categorical levels if they have no observations, but if the categorical variables are already of type `pd.Categorical` with only the levels that are present, it can be skipped for slightly faster fitting times. You’ll likely want to pass `False` here if merging several models into one through `append_trees`.
- **weights_as_sample_prob** (*bool*) – If passing sample (row) weights when fitting the model, whether to consider those weights as row sampling weights (i.e. the higher the weights, the more likely the observation will end up included in each tree sub-sample), or as distribution density weights (i.e. putting a weight of two is the same as if the row appeared twice, thus higher weight makes it less of an outlier, but does not give it a higher chance of being sampled if the data uses sub-sampling).
- **sample_with_replacement** (*bool*) – Whether to sample rows with replacement or not (not recommended). Note that distance calculations, if desired, don’t work well with duplicate rows.

Note that it is not possible to call `fit_predict` or `fit_transform` when using this option.

- **penalize_range** (*bool*) – Whether to penalize (add -1 to the terminal depth) observations at prediction time that have a value of the chosen split variable (linear combination in extended model) that falls outside of a pre-determined reasonable range in the data being split (given by `2 * range` in data and centered around the split point), as proposed in [Page 12, 4](#) and implemented in the authors’ original code in [Page 13, 5](#). Not used in single-variable model when splitting by categorical variables.

This option is not supported when using density-based outlier scoring metrics.

It’s recommended to turn this off for faster predictions on sparse CSC matrices.

Note that this can make a very large difference in the results when using `prob_pick_pooled_gain`.

Be aware that this option can make the distribution of outlier scores a bit different (i.e. not centered around 0.5).

- **scoring_metric** (*str*) – Metric to use for determining outlier scores (see reference¹³). Options are:

"depth"

Will use isolation depth as proposed in reference^{Page 11, 1}. This is typically the safest choice and plays well with all model types offered by this library.

"density"

Will set scores for each terminal node as the ratio between the fraction of points in the sub-sample that end up in that node and the fraction of the volume in the feature space which defines the node according to the splits that lead to it. If using `ndim=1`, for categorical variables, this is defined in terms of number of categories that go towards each side of the split divided by number of categories in the observations that reached that node.

The standardized outlier score from density for a given observation is calculated as the negative of the logarithm of the geometric mean from the per-tree densities, which unlike the standardized score produced from depth, is unbounded, but just like the standardized score from depth, has a natural threshold for defining outlierness, which in this case is zero instead of 0.5. The non-standardized outlier score is calculated as the geometric mean, while the per-tree scores are calculated as the density values.

This might lead to better predictions when using `ndim=1`, particularly in the presence of categorical variables. Note however that using density requires more trees for convergence of scores (i.e. good results) compared to isolation-based metrics.

This option is incompatible with `penalize_range`.

"adj_depth"

Will use an adjusted isolation depth that takes into account the number of points that go to each side of a given split vs. the fraction of the range of that feature that each side of the split occupies, by a metric as follows:

$$d = \frac{2}{1 + \frac{1}{2p}}$$

Where p is defined as:

$$p = \frac{n_s / r_s}{n_t / r_t}$$

With n_t being the number of points that reach a given node, n_s the number of points that are sent to a given side of the split/branch at that node, r_t being the range (maximum minus minimum) of the splitting feature or linear combination among the points that reached the node, and r_s being the range of the same feature or linear combination among the points that are sent to this same side of the split/branch. This makes each split add a number between zero and two to the isolation depth, with this number's probabilistic distribution being centered around 1 and thus the expected isolation depth remaining the same as in the original "depth" metric, but having more variability around the extremes.

Scores (standardized, non-standardized, per-tree) are aggregated in the same way as for "depth".

This might lead to better predictions when using `ndim=1`, particularly in the presence of categorical variables and for smaller datasets, and for smaller datasets, might make sense to combine it with `penalize_range=True`.

¹³ Cortes, David. "Isolation forests: looking beyond tree depth." arXiv preprint arXiv:2111.11639 (2021).

"adj_density"

Will use the same metric from "adj_depth", but applied multiplicatively instead of additively. The expected value for this adjusted density is not strictly the same as for isolation, but using the expected isolation depth as standardizing criterion tends to produce similar standardized score distributions (centered around 0.5).

Scores (standardized, non-standardized, per-tree) are aggregated in the same way as for "depth".

This option is incompatible with `penalize_range`.

"boxed_ratio"

Will set the scores for each terminal node as the ratio between the volume of the boxed feature space for the node as defined by the smallest and largest values from the split conditions for each column (bounded by the variable ranges in the sample) and the variable ranges in the tree sample. If using `ndim=1`, for categorical variables this is defined in terms of number of categories. If using `ndim>1`, this is defined in terms of the maximum achievable value for the splitting linear combination determined from the minimum and maximum values for each variable among the points in the sample, and as such, it has a rather different meaning compared to the score obtained with `ndim=1` - boxed ratio scores with `ndim>1` typically provide very poor quality results and this metric is thus not recommended to use in the extended model. With 'ndim>1', it also has a tendency of producing too small values which round to zero.

The standardized outlier score from boxed ratio for a given observation is calculated simply as the the average from the per-tree boxed ratios. This metric has a lower bound of zero and a theoretical upper bound of one, but in practice the scores tend to be very small numbers close to zero, and its distribution across different datasets is rather unpredictable. In order to keep rankings comparable with the rest of the metrics, the non-standardized outlier scores are calculated as the negative of the average instead. The per-tree scores are calculated as the ratios.

This metric can be calculated in a fast-but-not-so-precise way, and in a low-but-precise way, which is controlled by parameter `fast_bratio`. Usually, both should give the same results, but in some fatasets, the fast way can lead to numerical inaccuracies due to round-offs very close to zero.

This metric might lead to better predictions in datasets with many rows when using `ndim=1` and a relatively small `sample_size`. Note that more trees are required for convergence of scores when using this metric. In some datasets, this metric might result in very bad predictions, to the point that taking its inverse produces a much better ranking of outliers.

This option is incompatible with `penalize_range`.

"boxed_density2"

Will set the score as the ratio between the fraction of points within the sample that end up in a given terminal node and the boxed ratio metric.

Aggregation of scores (standardized, non-standardized, per-tree) is done in the same way as for density, and it also has a natural threshold at zero for determining outliers and inliers.

This metric is typically usable with 'ndim>1', but tends to produce much bigger values compared to 'ndim=1'.

Albeit unintuitively, in many datasets, one can usually get better results with metric "boxed_density" instead.

The calculation of this metric is also controlled by `fast_bratio`.

This option is incompatible with `penalize_range`.

"boxed_density"

Will set the score as the ratio between the fraction of points within the sample that end up in a given terminal node and the ratio between the boxed volume of the feature space in the sample and the boxed volume of a node given by the split conditions (inverse as in "boxed_density2"). This metric does not have any theoretical or intuitive justification behind its existence, and it is perhaps illogical to use it as a scoring metric, but tends to produce good results in some datasets.

The standardized outlier scores are defined as the negative of the geometric mean of this metric, while the non-standardized scores are the geometric mean, and the per-tree scores are simply the 'density' values.

The calculation of this metric is also controlled by `fast_bratio`.

This option is incompatible with `penalize_range`.

- **fast_bratio** (*bool*) – When using “boxed” metrics for scoring, whether to calculate them in a fast way through cumulative sum of logarithms of ratios after each split, or in a slower way as sum of logarithms of a single ratio per column for each terminal node.

Usually, both methods should give the same results, but in some datasets, particularly when variables have too small or too large ranges, the first method can be prone to numerical inaccuracies due to roundoff close to zero.

Note that this does not affect calculations for models with 'ndim>1', since given the split types, the calculation for them is different.

- **standardize_data** (*bool*) – Whether to standardize the features at each node before creating a linear combination of them as suggested in [Page 12, 4](#). This is ignored when using `ndim=1`.
- **weigh_by_kurtosis** (*bool*) – Whether to weigh each column according to the kurtosis obtained in the sub-sample that is selected for each tree as briefly proposed in [Page 11, 1](#). Note that this is only done at the beginning of each tree sample. For categorical columns, will calculate expected kurtosis if the column were converted to numerical by assigning to each category a random number $\sim \text{Unif}(0, 1)$.

Note that when using sparse matrices, the calculation of kurtosis will rely on a procedure that uses sums of squares and higher-power numbers, which has less numerical precision than the calculation used for dense inputs, and as such, the results might differ slightly.

Using this option makes the model more likely to pick the columns that have anomalous values when viewed as a 1-d distribution, and can bring a large improvement in some datasets.

This is intended as a cheap feature selector, while the parameter `prob_pick_col_by_kurt` provides the option to do this at each node in the tree for a different overall type of model.

If passing column weights or using weighted column choices proportional to some other metric (`prob_pick_col_by_range`, `prob_pick_col_by_var`), the effect will be multiplicative.

If passing `missing_action="fail"` and the data has infinite values, columns with rows having infinite values will get a weight of zero. If passing a different value for missing action, infinite values will be ignored in the kurtosis calculation.

If using `missing_action="impute"`, the calculation of kurtosis will not use imputed values in order not to favor columns with missing values (which would increase kurtosis by all having the same central value).

- **coefs** (*str, one of “normal” or “uniform”*) – For the extended model, whether to sample random coefficients according to a normal distribution $\sim \text{Normal}(0, 1)$ (as proposed in [Page 12, 4](#)) or according to a uniform distribution $\sim \text{Unif}(-1, +1)$ as proposed in [Page 12, 3](#). Ignored for

the single-variable model. Note that, for categorical variables, the coefficients will be sampled $\sim N(0,1)$ regardless - in order for both types of variables to have transformations in similar ranges (which will tend to boost the importance of categorical variables), pass "uniform" here.

- **assume_full_distr** (*bool*) – When calculating pairwise distances (see⁸), whether to assume that the fitted model represents a full population distribution (will use a standardizing criterion assuming infinite sample, and the results of the similarity between two points at prediction time will not depend on the presence of any third point that is similar to them, but will differ more compared to the pairwise distances between points from which the model was fit). If passing 'False', will calculate pairwise distances as if the new observations at prediction time were added to the sample to which each tree was fit, which will make the distances between two points potentially vary according to other newly introduced points. This will not be assumed when the distances are calculated as the model is being fit (see documentation for method 'fit_transform').

This was added for experimentation purposes only and it's not recommended to pass False. Note that when calculating distances using a tree indexer (after calling `build_index`), there might be slight discrepancies between the numbers produced with or without the indexer due to what are considered "additional" observations in this calculation.

- **build_imputer** (*bool*) – Whether to construct missing-value imputers so that later this same model could be used to impute missing values of new (or the same) observations. Be aware that this will significantly increase the memory requirements and serialized object sizes. Note that this is not related to 'missing_action' as missing values inside the model are treated differently and follow their own imputation or division strategy.
- **min_imp_obs** (*int*) – Minimum number of observations with which an imputation value can be produced. Ignored if passing 'build_imputer' = 'False'.
- **depth_imp** (*str, one of "higher", "lower", "same"*) – How to weight observations according to their depth when used for imputing missing values. Passing "higher" will weigh observations higher the further down the tree (away from the root node) the terminal node is, while "lower" will do the opposite, and "same" will not modify the weights according to node depth in the tree. Implemented for testing purposes and not recommended to change from the default. Ignored when passing 'build_imputer' = 'False'.
- **weigh_imp_rows** (*str, one of "inverse", "prop", "flat"*) – How to weight node sizes when used for imputing missing values. Passing "inverse" will weigh a node inversely proportional to the number of observations that end up there, while "proportional" will weigh them heavier the more observations there are, and "flat" will weigh all nodes the same in this regard regardless of how many observations end up there. Implemented for testing purposes and not recommended to change from the default. Ignored when passing 'build_imputer' = 'False'.
- **random_seed** (*int*) – Seed that will be used for random number generation.
- **use_long_double** (*bool*) – Whether to use 'long double' (extended precision) type for more precise calculations about standard deviations, means, ratios, weights, gain, and other potential aggregates. This makes such calculations accurate to a larger number of decimals (provided that the compiler used has wider long doubles than doubles) and it is highly recommended to use when the input data has a number of rows or columns exceeding 2^{53} (an unlikely scenario), and also highly recommended to use when the input data has problematic scales (e.g. numbers that differ from each other by something like 10^{-100} or columns that include values like 10^{100} , 10^{-10} , and 10^{-100} and still need to be sensitive to a difference of 10^{-10}), but will make the calculations slower, the more so in platforms in which 'long double' is a software-emulated type (e.g. Power8 platforms). Note that some platforms (most

⁸ Cortes, David. "Distance approximation using Isolation Forests." arXiv preprint arXiv:1910.12362 (2019).

notably windows with the msvc compiler) do not make any difference between ‘double’ and ‘long double’.

If ‘long double’ is not going to be used, the library can be compiled without support for it (making the library size smaller) by defining an environment variable `NO_LONG_DOUBLE` before installing this package (e.g. through `export NO_LONG_DOUBLE=1` before running the `pip` command).

This option is not available on Windows, due to lack of support in some compilers (e.g. msvc) and lack of thread-safety in the calculations in others (e.g. mingw).

- **nthreads** (*int*) – Number of parallel threads to use. If passing a negative number, will use the same formula as joblib does for calculating number of threads (which is `n_cpus + 1 + n_jobs` - i.e. pass -1 to use all available threads). Note that, the more threads, the more memory will be allocated, even if the thread does not end up being used. Be aware that most of the operations are bound by memory bandwidth, which means that adding more threads will not result in a linear speed-up. For some types of data (e.g. large sparse matrices with small sample sizes), adding more threads might result in only a very modest speed up (e.g. 1.5x faster with 4x more threads), even if all threads look fully utilized.
- **n_estimators** (*None or int*) – Synonym for `ntrees`, kept for better compatibility with scikit-learn.
- **max_samples** (*None or int*) – Synonym for `sample_size`, kept for better compatibility with scikit-learn.
- **n_jobs** (*None or int*) – Synonym for `nthreads`, kept for better compatibility with scikit-learn.
- **random_state** (*None, int, or RandomState*) – Synonym for `random_seed`, kept for better compatibility with scikit-learn.
- **bootstrap** (*None or bool*) – Synonym for `sample_with_replacement`, kept for better compatibility with scikit-learn.

Variables

- **cols_numeric** (*array(n_num_features,)*) – Array with the names of the columns that were taken as numerical (Only when fitting the model to a DataFrame object).
- **cols_categ** (*array(n_categ_features,)*) – Array with the names of the columns that were taken as categorical (Only when fitting the model to a DataFrame object).
- **is_fitted** (*bool*) – Indicator telling whether the model has been fit to data or not.

References

`append_trees` (*other*)

Appends isolation trees from another Isolation Forest model into this one

This function is intended for merging models **that use the same hyperparameters** but were fitted to different subsets of data.

In order for this to work, both models must have been fit to data in the same format - that is, same number of columns, same order of the columns, and same column types, although not necessarily same object classes (e.g. can mix `np.array` and `scipy.sparse.csc_matrix`).

If the data has categorical variables, the models should have been built with parameter `recode_categ=False` in the class constructor, and the categorical columns passed as type `pd.Categorical` with the same encoding - otherwise different models might be using different encodings

for each categorical column, which will not be preserved as only the trees will be appended without any associated metadata.

Note: This function will not perform any checks on the inputs, and passing two incompatible models (e.g. fit to different numbers of columns) will result in wrong results and potentially crashing the Python process when using it.

Note: This function is not thread-safe - that is, it will produce problems if one tries to call this function on the same model object in parallel through e.g. `joblib` with a shared-memory backend (which is not the default for `joblib`).

Parameters

other (*IsolationForest*) – Another Isolation Forest model from which trees will be appended to this model. It will not be modified during the call to this function.

Returns

self – This object.

Return type

obj

build_indexer(*with_distances=False*)

Build indexer for faster terminal node predictions and/or distance calculations

Builds an index of terminal nodes for faster prediction of terminal node numbers (calling `predict` with `output="tree_num"`).

Optionally, can also pre-calculate terminal node distances in order to speed up distance calculations (calling `predict_distance`).

Note: This feature is not available for models that use `missing_action="divide"` or `new_categ_action="weighted"` (which are the defaults when passing `ndim=1`).

Parameters

with_distances (*bool*) – Whether to also pre-calculate node distances in order to speed up `predict_distance`. Note that this will consume a lot more memory and make the resulting object significantly heavier.

Returns

self – This object

Return type

obj

copy()

Get a deep copy of this object

Returns

copied – A deep copy of this object

Return type

obj

decision_function(X)

Wrapper for 'predict' with 'output=score'

This function is kept for compatibility with Scikit-Learn.

Parameters

X (*array or array-like (n_samples, n_features)*) – Observations for which to predict outlier-ness or average isolation depth. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC or CSR matrix.

Returns

score – Outlier scores for the rows in 'X' (the higher, the most anomalous).

Return type

array(n_samples,)

drop_imputer()

Drops the imputer sub-object from this model object

Drops the imputer sub-object from this model object, if it was fitted with data imputation capabilities. The imputer, if constructed, is likely to be a very heavy object which might not be needed for all purposes.

Returns

self – This object

Return type

obj

drop_indexer()

Drops the indexer sub-object from this model object

Drops the indexer sub-object from this model object, if it was constructed. The indexer, if constructed, is likely to be a very heavy object which might not be needed for all purposes.

Note that reference points as added through `set_reference_points` are associated with the indexer object and will also be dropped if any were added.

Returns

self – This object

Return type

obj

drop_reference_points()

Drops reference points from this model

Drops any reference points used for distance and/or kernel calculations from the model object, if any were set through `set_reference_points`.

Returns

self – This object

Return type

obj

export_model(file, add_metada_file=False)

Export Isolation Forest model

Save Isolation Forest model to a serialized file along with its metadata, in order to be re-used in Python or in the R or the C++ versions of this package.

This function is not suggested to be used for passing models to and from Python - in such case, one can use `pickle` instead, although the function still works correctly for serializing objects between Python processes.

Note that, if the model was fitted to a `DataFrame`, the column names must be something exportable as JSON, and must be something that R could use as column names (for example, using integers as column names is valid in pandas but not in R).

Can optionally generate a JSON file with metadata such as the column names and the levels of categorical variables, which can be inspected visually in order to detect potential issues (e.g. character encoding) or to make sure that the columns are of the right types.

The metadata file, if produced, will contain, among other things, the encoding that was used for categorical columns - this is under `data_info.cat_levels`, as an array of arrays by column, with the first entry for each column corresponding to category 0, second to category 1, and so on (the C++ version takes them as integers). When passing `categ_cols`, there will be no encoding but it will save the maximum category integer and the column numbers instead of names.

The serialized file can be used in the C++ version by reading it as a binary file and de-serializing its contents using the C++ function 'deserialize_combined' (recommended to use 'inspect_serialized_object' beforehand).

Be aware that this function will write raw bytes from memory as-is without compression, so the file sizes can end up being much larger than when using `pickle`.

The metadata is not used in the C++ version, but is necessary for the R and Python versions.

Note: While in earlier versions of this library this functionality used to be faster than `pickle`, starting with version 0.3.0, this function and `pickle` should have similar timings and it's recommended to use `pickle` for serializing objects across Python processes.

Note: Important: The model treats boolean variables as categorical. Thus, if the model was fit to a `DataFrame` with boolean columns, when importing this model into C++, they need to be encoded in the same order - e.g. the model might encode `True` as zero and `False` as one - you need to look at the metadata for this. Also, if using some of Pandas' own Boolean types, these might end up as non-boolean categorical, and if importing the model into R, you might need to pass values as e.g. "True" instead of `TRUE` (look at the `.metadata` file to determine this).

Note: The files produced by this function will be compatible between:

- Different operating systems.
- Different compilers.
- Different Python/R versions.
- Systems with different 'size_t' width (e.g. 32-bit and 64-bit), as long as the file was produced on a system that was either 32-bit or 64-bit, and as long as each saved value fits within the range of the machine's 'size_t' type.
- Systems with different 'int' width, as long as the file was produced on a system that was 16-bit, 32-bit, or 64-bit, and as long as each saved value fits within the range of the machine's int type.
- Systems with different bit endianness (e.g. x86 and PPC64 in non-le mode).
- Versions of this package from 0.3.0 onwards, **but only forwards compatible** (e.g. a model saved with versions 0.3.0 to 0.3.5 can be loaded under version 0.3.6, but not the other way around, and attempting to do so will cause crashes and memory corruptions without an informative error message). **This last point applies also to models saved through pickle.** Note that loading a model produced by an earlier version of the library might be slightly slower.

But will not be compatible between:

- Systems with different floating point numeric representations (e.g. standard IEEE754 vs. a base-10 system).
- Versions of this package earlier than 0.3.0.

This pretty much guarantees that a given file can be serialized and de-serialized in the same machine in which it was built, regardless of how the library was compiled.

Reading a serialized model that was produced in a platform with different characteristics (e.g. 32-bit vs. 64-bit) will be much slower.

Note: On Windows, if compiling this library with a compiler other than MSVC or MINGW, there might be issues exporting models larger than 2GB.

Parameters

- **file** (*str*) – The output file path into which to export the model. Must be a file name, not a file handle.
- **add_metada_file** (*bool*) – Whether to generate a JSON file with metadata, which will have the same name as the model but will end in ‘.metadata’. This file is not used by the de-serialization function, it’s only meant to be inspected manually, since such contents will already be written in the produced model file.

Returns

self – This object.

Return type

obj

fit(*X, y=None, sample_weights=None, column_weights=None, categ_cols=None*)

Fit isolation forest model to data

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Data to which to fit the model. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix. If passing a DataFrame, will assume that columns are:
 - Numeric, if their dtype is a subtype of NumPy’s ‘number’ or ‘datetime64’.
 - Categorical, if their dtype is ‘object’, ‘Categorical’, or ‘bool’. Note that, if *Categorical* dtypes are ordered, the order will be ignored here.

Other dtypes are not supported.

Note that, if passing NumPy arrays, they are used in column-major order (a.k.a. “Fortran arrays”), and if they are not already in column-major format, will need to create a copy of the data.

If passing a DataFrame with categorical columns, then column names must be unique.

- **y** (*None*) – Not used. Kept as argument for compatibility with Scikit-Learn pipelining.
- **sample_weights** (*None or array(n_samples,)*) – Sample observation weights for each row of ‘X’, with higher weights indicating either higher sampling probability (i.e. the observation has a larger effect on the fitted model, if using sub-samples), or distribution density

(i.e. if the weight is two, it has the same effect of including the same data point twice), according to parameter 'weights_as_sample_prob' in the model constructor method.

- **column_weights** (*None or array(n_features,)*) – Sampling weights for each column in 'X'. Ignored when picking columns by deterministic criterion. If passing None, each column will have a uniform weight. If used along with kurtosis weights, the effect is multiplicative.
- **categ_cols** (*None or array-like*) – Columns that hold categorical features, when the data is passed as an array or matrix. Categorical columns should contain only integer values with a continuous numeration starting at zero, with negative values and NaN taken as missing, and the array or list passed here should correspond to the column numbers, with numeration starting at zero. The maximum categorical value should not exceed 'INT_MAX' (typically $2^{31} - 1$). This might be passed either at construction time or when calling `fit` or variations of `fit`.

This is ignored when the input is passed as a `DataFrame` as then it will consider columns as categorical depending on their dtype.

Returns

self – This object.

Return type

obj

fit_predict(*X, column_weights=None, output_outlierness='score', output_distance=None, square_mat=False, output_imputed=False, categ_cols=None*)

Fit the model in-place and produce isolation or separation depths along the way

See the documentation of other methods ('init', 'fit', 'predict', 'predict_distance') for details.

Note: The data must NOT contain any duplicate rows.

Note: This function will be faster at predicting average depths than calling 'fit' + 'predict' separately when using full row samples.

Note: If using 'penalize_range' = 'True', the resulting scores/depths from this function might differ a bit from those of 'fit' + 'predict' ran separately.

Note: Sample weights are not supported for this method.

Note: When using multiple threads, there can be small differences in the predicted scores or average depth or separation/distance between runs due to roundoff error.

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Data to which to fit the model. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix. If passing a DataFrame, will assume that columns are:
 - Numeric, if their dtype is a subtype of NumPy's 'number' or 'datetime64'.

- Categorical, if their dtype is ‘object’, ‘Categorical’, or ‘bool’. Note that, if *Categorical* dtypes are ordered, the order will be ignored here.

Other dtypes are not supported.

If passing a DataFrame with categorical columns, then column names must be unique.

- **column_weights** (*None or array(n_features,)*) – Sampling weights for each column in ‘X’. Ignored when picking columns by deterministic criterion. If passing None, each column will have a uniform weight. If used along with kurtosis weights, the effect is multiplicative. Note that, if passing a DataFrame with both numeric and categorical columns, the column names must not be repeated, otherwise the column weights passed here will not end up matching.
 - **output_outlierness** (*None or str in [“score”, “avg_depth”]*) – Desired type of outlierness output. If passing “score”, will output standardized outlier score. If passing “avg_depth” will output average isolation depth without standardizing. If passing ‘None’, will skip outlierness calculations.
 - **output_distance** (*None or str in [“dist”, “avg_sep”]*) – Type of distance output to produce. If passing “dist”, will standardize the average separation depths. If passing “avg_sep”, will output the average separation depth without standardizing it (note that lower separation depth means furthest distance). If passing ‘None’, will skip distance calculations.
- Note that it might be much faster to calculate distances through a fitted object with `build_indexer` instead of calling this method.
- **square_mat** (*bool*) – Whether to produce a full square matrix with the distances. If passing ‘False’, will output only the upper triangular part as a 1-d array in which entry (i,j) with $0 \leq i < j < n$ is located at position $p(i,j) = (i * (n - (i+1)/2) + j - i - 1)$. Ignored when passing ‘output_distance’ = ‘None’.
 - **output_imputed** (*bool*) – Whether to output the data with imputed missing values. Model object must have been initialized with ‘build_imputer’ = ‘True’.
 - **categ_cols** (*None or array-like*) – Columns that hold categorical features, when the data is passed as an array or matrix. Categorical columns should contain only integer values with a continuous numeration starting at zero, with negative values and NaN taken as missing, and the array or list passed here should correspond to the column numbers, with numeration starting at zero. The maximum categorical value should not exceed ‘INT_MAX’ (typically $2^{31} - 1$). This might be passed either at construction time or when calling `fit` or variations of `fit`.

This is ignored when the input is passed as a DataFrame as then it will consider columns as categorical depending on their dtype.

Returns

output – Requested outputs about isolation depth (outlierness), pairwise separation depth (distance), and/or imputed missing values. If passing either ‘output_distance’ or ‘output_imputed’, will return a dictionary with keys “pred” (array(n_samples,)), “dist” (array(n_samples * (n_samples - 1) / 2,)) or array(n_samples, n_samples)), “imputed” (array-like(n_samples, n_columns)), according to whether each output type is present.

Return type

array(n_samples,), or dict

fit_transform(*X, y=None, column_weights=None, categ_cols=None*)

Scikit-Learn pipeline-compatible version of ‘fit_predict’

Will fit the model and output imputed missing values. Intended to be used as part of Scikit-learn pipelining. Note that this is just a wrapper over ‘fit_predict’ with parameter ‘output_imputed’ = ‘True’. See the documentation of ‘fit_predict’ for details.

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Data to which to fit the model and whose missing values need to be imputed. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix (see the documentation of `fit` for more details).
- **y** (*None*) – Not used. Kept for compatibility with Scikit-Learn.
- **column_weights** (*None or array(n_features,)*) – Sampling weights for each column in ‘X’. Ignored when picking columns by deterministic criterion. If passing `None`, each column will have a uniform weight. If used along with kurtosis weights, the effect is multiplicative. Note that, if passing a DataFrame with both numeric and categorical columns, the column names must not be repeated, otherwise the column weights passed here will not end up matching.
- **categ_cols** (*None or array-like*) – Columns that hold categorical features, when the data is passed as an array or matrix. Categorical columns should contain only integer values with a continuous numeration starting at zero, with negative values and `NaN` taken as missing, and the array or list passed here should correspond to the column numbers, with numeration starting at zero. The maximum categorical value should not exceed ‘INT_MAX’ (typically $2^{31} - 1$). This might be passed either at construction time or when calling `fit` or variations of `fit`.

This is ignored when the input is passed as a `DataFrame` as then it will consider columns as categorical depending on their dtype.

Returns

imputed – Input data ‘X’ with missing values imputed according to the model.

Return type

array-like(n_samples, n_columns)

generate_sql (*enclose='doublequotes', output_tree_num=False, tree=None, table_from=None, select_as='outlier_score', column_names=None, column_names_categ=None*)

Generate SQL statements representing the model prediction function

Generate SQL statements - either separately per tree (the default), for a single tree if needed (if passing `tree`), or for all trees concatenated together (if passing `table_from`). Can also be made to output terminal node numbers (numeration starting at zero).

Note: Making predictions through SQL is much less efficient than from the model itself, as each terminal node will have to check all of the conditions that lead to it instead of passing observations down a tree.

Note: If constructed with the default arguments, the model will not perform any sub-sampling, which can lead to very big trees. If it was fit to a large dataset, the generated SQL might consist of gigabytes of text, and might lay well beyond the character limit of commands accepted by SQL vendors.

Note: The generated SQL statements will not include range penalizations, thus predictions might differ from calls to `predict` when using `penalize_range=True`.

Note: The generated SQL statements will only include handling of missing values when using `missing_action="impute"`. When using the single-variable model with categorical variables + subset splits, the rule buckets might be incomplete due to not including categories that were not present in a given node - this last point can be avoided by using `new_categ_action="smallest"`, `new_categ_action="random"`, or `missing_action="impute"` (in the latter case will treat them as missing, but the `predict` function might treat them differently).

Note: The resulting statements will include all the tree conditions as-is, with no simplification. Thus, there might be lots of redundant conditions in a given terminal node (e.g. “`X > 2`” and “`X > 1`”, the second of which is redundant).

Note: If using `scoring_metric="density"` or `scoring_metric="boxed_ratio"` plus `output_tree_num=False`, the outputs will correspond to the logarithm of the density rather than the density.

Parameters

- **enclose** (*str*) – With which symbols to enclose the column names in the select statement so as to make them SQL compatible in case they include characters like dots. Options are:
 - "doublequotes":**
Will enclose them as "column_name" - this will work for e.g. PostgreSQL.
 - "squarebraces":**
Will enclose them as [column_name] - this will work for e.g. SQL Server.
 - "none":**
Will output the column names as-is (e.g. column_name)
- **output_tree_num** (*bool*) – Whether to make the statements return the terminal node number instead of the isolation depth. The numeration will start at zero.
- **tree** (*int or None*) – Tree for which to generate SQL statements. If passed, will generate the statements only for that single tree. If passing 'None', will generate statements for all trees in the model.
- **table_from** (*str or None*) – If passing this, will generate a single select statement for the outlier score from all trees, selecting the data from the table name passed here. In this case, will always output the outlier score, regardless of what is passed under `output_tree_num`.
- **select_as** (*str*) – Alias to give to the generated outlier score in the select statement. Ignored when not passing `table_from`.
- **column_names** (*None or list[str]*) – Column names to use for the **numeric** columns. If not passed and the model was fit to a `DataFrame`, will use the column names from that `DataFrame`, which can be found under `self.cols_numeric_`. If not passing it and the model was fit to data in a format other than `DataFrame`, the columns will be named "column_N" in the resulting SQL statement. Note that the names will be taken verbatim - this function will not do any checks for whether they constitute valid SQL or not, and will not escape characters such as double quotation marks.
- **column_names_categ** (*None or list[str]*) – Column names to use for the **categorical** columns. If not passed, will use the column names from the `DataFrame` to which the model was fit. These can be found under `self.cols_categ_`.

Returns

sql – A list of SQL statements for each tree as strings, or the SQL statement for a single tree if passing ‘tree’, or a single select-from SQL statement with all the trees concatenated if passing `table_from`.

Return type

list[str] or str

get_num_nodes()

Get number of nodes per tree

Gets the number of nodes per tree, along with the number of terminal nodes.

Returns

nodes – A tuple in which the first element denotes the total number of nodes in each tree, and the second element denotes the number of terminal nodes. Both are returned as arrays having one entry per tree.

Return type

tuple(array(n_trees,), array(n_trees,))

get_params(*deep=True*)

Get parameters for this estimator.

Kept for compatibility with scikit-learn.

Parameters

deep (*bool*) – Ignored.

Returns

params – Parameter names mapped to their values.

Return type

dict

property has_indexer_**property has_reference_points_****static import_model(*file*)**

Load an Isolation Forest model exported from R or Python

Loads a serialized Isolation Forest model as produced and exported by the function `export_model` or by the R version of this package. Note that the metadata must be something importable in Python - e.g. column names must be valid for Pandas.

It’s recommended to generate a ‘.metadata’ file (passing `add_metadata_file=True`) and to visually inspect said file in any case.

See the documentation for `export_model` for details about compatibility of the generated files across different machines and versions.

Note:

This is a static class method - that is, it should be called like this:

```
iso = IsolationForest.import_model(...)
```

(i.e. no parentheses after *IsolationForest*)

Note: While in earlier versions of this library this functionality used to be faster than `pickle`, starting with version 0.3.0, this function and `pickle` should have similar timings and it's recommended to use `pickle` for serializing objects across Python processes.

Parameters

file (*str*) – The input file path containing an exported model along with its metadata file. Must be a file name, not a file handle.

Returns

iso – An Isolation Forest model object reconstructed from the serialized file and ready to use.

Return type

IsolationForest

partial_fit(*X*, *sample_weights=None*, *column_weights=None*, *X_ref=None*)

Add additional (single) tree to isolation forest model

Adds a single tree fit to the full (non-subsampled) data passed here. Must have the same columns as previously-fitted data.

Note: If constructing trees with different sample sizes, the outlier scores with depth-based metrics will not be centered around 0.5 and might have a very skewed distribution. The standardizing constant for the scores will be taken according to the sample size passed in the construction argument (if that is `None` or `"auto"`, will then set it as the sample size of the first tree).

If trees are going to be fit to samples of different sizes, it's strongly recommended to use density-based scoring metrics instead.

Note: This function is not thread-safe - that is, it will produce problems if one tries to call this function on the same model object in parallel through e.g. `joblib` with a shared-memory backend (which is not the default for `joblib`).

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Data to which to fit the new tree. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix. If passing a DataFrame, will assume that columns are:

- Numeric, if their dtype is a subtype of NumPy's 'number' or 'datetime64'.

- Categorical, if their dtype is 'object', 'Categorical', or 'bool'. Note that, if *Categorical* dtypes are ordered, the order will be ignored here. Categorical columns, if any, may have new categories.

Other dtypes are not supported.

If passing an array and the array is not in column-major format, will be forcibly converted to column-major, which implies an extra data copy.

If passing a DataFrame with categorical columns, then column names must be unique.

- **sample_weights** (*None or array(n_samples,)*) – Sample observation weights for each row of 'X', with higher weights indicating distribution density (i.e. if the weight is two, it has

the same effect of including the same data point twice). If not 'None', model must have been built with 'weights_as_sample_prob' = 'False'.

- **column_weights** (*None or array(n_features,)*) – Sampling weights for each column in 'X'. Ignored when picking columns by deterministic criterion. If passing None, each column will have a uniform weight. If used along with kurtosis weights, the effect is multiplicative.
- **X_ref** (*array or array-like (n_references, n_features)*) – Reference points for distance and/or kernel calculations, if these were previously added to the model object through `set_reference_points`. Must correspond to the same points that were passed to the call to `set_reference_points`.

Might be passed in either row-major (preferred) or column-major order. If sparse, only CSC format is supported.

This is ignored if the model has no stored reference points.

Returns

self – This object.

Return type

obj

predict (*X, output='score'*)

Predict outlieriness based on average isolation depth or density

Calculates the approximate depth that it takes to isolate an observation according to the fitted model splits, or the average density of the branches in which observations fall. Can output either the average depth/density, or a standardized outlier score based on whether it takes more or fewer splits than average to isolate observations. In the standardized outlier score for density-based metrics, values closer to 1 indicate more outlieriness, while values closer to 0.5 indicate average outlieriness, and close to 0 more averageness (harder to isolate). When using `scoring_metric="density"`, the standardized outlier scores are instead unbounded, with larger values indicating more outlieriness and a natural threshold of zero for determining inliers and outliers.

Note: For multi-threaded predictions on many rows, it is recommended to set the number of threads to the number of physical cores of the CPU rather than the number of logical cores, as it will typically have better performance that way. Assuming a typical x86-64 desktop CPU, this typically involves dividing the number of threads by 2 - for example:

```
import multiprocessing;model.set_params(nthreads=multiprocessing.
cpu_count()/2)
```

Note: Depending on the model parameters, it might be possible to convert the models to 'treelite' format for faster predictions or for easier model serving. See method `to_treelite` for details.

Note: If the model was built with 'nthreads>1', this prediction function will use OpenMP for parallelization. In a linux setup, one usually has GNU's "gomp" as OpenMP as backend, which will hang when used in a forked process - for example, if one tries to call this prediction function from 'flask'+ 'unicorn', which uses process forking for parallelization, it will cause the whole application to freeze; and if using kubernetes on top of a different backend such as 'falcon', might cause it to run slower than needed or to hang too. A potential fix in these cases is to set the number of threads to 1 in the object (e.g. 'model.nthreads = 1'), or to use a different version of this library compiled without OpenMP (requires manually altering the 'setup.py')

file), or to use a non-GNU OpenMP backend. This should not be an issue when using this library normally in e.g. a jupyter notebook.

Note: For model serving purposes, in order to have a smaller and leaner library, it is recommended to compile this library without support for 'long double' type, which can be done by setting up an environment variable "NO_LONG_DOUBLE" before installation of this package (see the GitHub page of this library for more details).

Note: The more threads that are set for the model, the higher the memory requirements will be as each thread will allocate an array with one entry per row.

Note: In order to save memory when fitting and serializing models, the functionality for outputting terminal node number will generate index mappings on the fly for all tree nodes, even if passing only 1 row, so it's only recommended for batch predictions. If this type of prediction is desired, it can be sped up by building an index of terminal nodes through `build_indexer`.

Note: The outlier scores/depth predict functionality is optimized for making predictions on one or a few rows at a time - for making large batches of predictions, it might be faster to use the 'fit_predict' functionality.

Note: If using non-random splits (parameters `prob_pick_avg_gain`, `prob_pick_pooled_gain`, `prob_pick_full_gain`, `prob_pick_dens`) and/or range penalizations (which are off by default), the distribution of scores might not be centered around 0.5.

Note: When making predictions on CSC matrices with many rows using multiple threads, there can be small differences between runs due to roundoff error.

Parameters

- **X** (array or array-like (*n_samples*, *n_features*)) – Observations for which to predict outlieriness or average isolation depth. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC or CSR matrix.

If 'X' is sparse and one wants to obtain the outlier score or average depth or tree numbers, it's highly recommended to pass it in CSC format as it will be much faster when the number of trees or rows is large.

While the 'X' used by `fit` always needs to be in column-major order, predictions can be done on data that is in either row-major or column-major orders, with row-major being faster for dense data.

- **output** (*str*, one of "score", "avg_depth", "tree_num", "tree_depths") – Desired type of output. Options are:

"score":

Will output standardized outlier scores. For all scoring metrics, higher values indicate

more outlierness.

"avg_depth":

Will output unstandardized average isolation depths. For `scoring_metric="density"`, will output the geometric mean instead. See the documentation for `scoring_metric`, for more details about the calculation for other metrics. For all scoring metrics, higher values indicate less outlierness.

"tree_num":

Will output the index of the terminal node under each tree in the model. If this calculation is going to be performed frequently, it's recommended to build node indices through `build_indexer`.

"tree_depths":

Will output non-standardized per-tree isolation depths or densities or log-densities (note that they will not include range penalties from `penalize_range=True`). See the documentation for `scoring_metric` for details about the calculation for each metrics.

Returns

score – Requested output type for each row according to parameter 'output' (outlier scores, average isolation depth, terminal node indices, or per-tree isolation depths).

Return type

array(n_samples,) or array(n_samples, n_trees)

predict_distance (*X*, *output='dist'*, *square_mat=True*, *X_ref=None*, *use_reference_points=True*)

Predict approximate distances or isolation kernels/proximities between points

Predict approximate pairwise distances between points, or individual distances between two sets of points based on how many splits it takes to separate them, or isolation kernels (a.k.a. proximity matrix, which for example can be used for a generalized least-squares regressions as a rough estimate of residual correlations) from the model based on the number of trees in which two observations end up in the same terminal node. Can output either the average number of paths/steps it takes to separate two observations, or a standardized metric (in the same way as the outlier score) in which values closer to zero indicate nearer points, closer to one further away points, and closer to 0.5 average distance, or a kernel/proximity metric, either standardized (values between zero and one) or raw (values ranging from zero to the number of trees in the model).

Note: The more threads that are set for the model, the higher the memory requirement will be as each thread will allocate an array with one entry per combination (with an exception being calculation of distances to reference points, which do not do this).

Note: Separation depths are very slow to calculate. By default, it will do it through a procedure that counts steps as observations are passed down the trees, which is especially slow and not recommended for more than a few thousand observations. If this function is going to be called repeatedly and/or it is going to be called for a large number of rows, it's highly recommended to build node distance indexes beforehand through `build_indexer` with option `with_distances=True`, as then the computation will be done based on terminal node indices instead, which is a much faster procedure. If the calculations are all going to be performed with respect to a fixed set of points, it's highly recommended to set those points as references through `set_reference_points`.

Note: If using `assume_full_distr=False` (not recommended to use such option), predictions with and without an indexer will differ slightly due to differences in what they count towards "additional" observa-

tions in the calculation.

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Observations for which to calculate approximate pairwise distances or kernels, or first group for distances/kernels between sets of points. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix.
- **output** (*str, one of “dist”, “avg_sep”, “kernel”, “kernel_raw”*) – Type of output to produce. If passing “dist”, will standardize the average separation depths. If passing “avg_sep”, will output the average separation depth without standardizing it (note that lower separation depth means furthest distance). If passing “kernel”, will output the fraction of the trees in which two observations end up in the same terminal node. If passing “kernel_raw”, will output the number (not fraction) of trees in which two observations end up in the same terminal node.

Note that for “kernel” and “kernel_raw”, having an indexer without reference points will not speed up calculations, and if such calculations are going to be done frequently, it is highly recommended to set reference points in the model object.

- **square_mat** (*bool*) – Whether to produce a full square matrix with the pairwise distances or kernels. If passing ‘False’, will output only the upper triangular part as a 1-d array in which entry (i,j) with $0 \leq i < j < n$ is located at position $p(i,j) = (i * (n - (i+1)/2) + j - i - 1)$.

Ignored when passing X_ref or use_reference_points=True plus having reference points.

- **X_ref** (*array or array-like (n_ref, n_features)*) – Second group of observations. If passing it, will calculate distances/kernels between each point in X and each point in X_ref. If passing None (the default), will calculate pairwise distances/kernels between the points in X. Must be of the same type as X (e.g. array, DataFrame, CSC).

Note that, if X_ref is passed and the model object has an indexer with reference points added (through set_reference_points), those reference points will be ignored for the calculation.

- **use_reference_points** (*bool*) – When the model object has an indexer with reference points (which can be added through set_reference_points), whether to calculate the distances/kernels from X to those reference points instead of the pairwise distances/kernels between points in X.

This is ignored when passing X_ref or when the model object does not contain an indexer or the indexer does not contain reference points.

Returns

dist – Approximate distances or average separation depth or kernels/proximities between points, according to parameter ‘output’. Shape and size depends on parameters square_mat, use_reference_points, and whether X_ref is passed.

Return type

array(n_samples * (n_samples - 1) / 2,) or array(n_samples, n_samples) or array(n_samples, n_ref)

predict_kernel(X, square_mat=True, X_ref=None, use_reference_points=True)

Predict isolation kernel between points

This is a shorthand for predict_distance with output="kernel".

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Observations for which to calculate approximate pairwise kernels/proximities, or first group for kernels between sets of points. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix.
- **square_mat** (*bool*) – Whether to produce a full square matrix with the pairwise kernels. If passing ‘False’, will output only the upper triangular part as a 1-d array in which entry (i,j) with $0 \leq i < j < n$ is located at position $p(i,j) = (i * (n - (i+1)/2) + j - i - 1)$. Ignored when passing X_ref.
- **X_ref** (*array or array-like (n_ref, n_features)*) – Second group of observations. If passing it, will calculate kernels between each point in X and each point in X_ref. If passing None (the default), will calculate pairwise kernels between the points in X. Must be of the same type as X (e.g. array, DataFrame, CSC).

Note that, if X_ref is passed and the model object has an indexer with reference points added (through `set_reference_points`), those reference points will be ignored for the calculation.

- **use_reference_points** (*bool*) – When the model object has an indexer with reference points (which can be added through `set_reference_points`), whether to calculate the kernels from X to those reference points instead of the pairwise kernels between points in X.

This is ignored when passing X_ref or when the model object does not contain an indexer or the indexer does not contain reference points.

Returns

dist – Approximate kernels between points, according to parameter ‘output’. Shape and size depends on parameter `square_mat`, and whether X_ref is passed.

Return type

`array(n_samples * (n_samples - 1) / 2,)` or `array(n_samples, n_samples)` or `array(n_samples, n_ref)`

`set_params(**params)`

Set the parameters of this estimator.

Kept for compatibility with scikit-learn.

Note: Setting any parameter other than the number of threads, will reset the model object to a blank state - that is, if it was fitted to some data, the fitted model will be lost, and it will need to be refitted before being able to make predictions.

Parameters

****params** (*dict*) – Estimator parameters.

Returns

self – Estimator instance.

Return type

estimator instance

`set_reference_points(X, with_distances=False)`

Set reference points to calculate distances or kernels with

Sets some points as pre-defined landmarks with respect to which distances and/or isolation kernel values will be calculated for arbitrary new points in calls to `predict_distance` and/or `predict_kernel`. If any points have already been set as references in the model object, they will be overwritten with the new points passed here.

Note that points are added in terms of their terminal node indices, but the raw data about them is not kept - thus, calling `partial_fit` later on a model with reference points requires passing those reference points again to add their node indices to the new tree.

Be aware that adding reference points requires building a tree indexer.

Parameters

- **X** (*array or array-like (n_samples, n_features)*) – Observations to set as references for future distance and/or isolation kernel calculations. Can pass a NumPy array, Pandas DataFrame, or SciPy sparse CSC matrix.
- **with_distances** (*bool*) – Whether to pre-calculate node distances (this is required to calculate distance from arbitrary points to the reference points).

Note that reference points for distances can only be set when using `assume_full_distr=False` (which is the default).

Returns

self – This object

Return type

obj

`subset_trees(trees_take)`

Subset trees of a given model

Creates a new model containing only selected trees of this model object.

Parameters

trees_take (*array_like(n,)*) – Indices of the trees of this model to copy over to the new model. Must be integers with numeration starting at zero.

Returns

new_model – A new IsolationForest model object, containing only the subset of trees from this object that was specified under ‘trees_take’.

Return type

obj

`to_treelite(use_float32=False)`

Convert model to ‘treelite’ format

Converts an IsolationForest model to a ‘treelite’ object, which can be compiled into a small standalone runtime library for smaller models and usually faster predictions:

<https://treelite.readthedocs.io/en/latest/index.html>

A couple notes about this conversion:

- It is only possible to convert to ‘treelite’ when using `ndim=1` (which is not the default).
- The ‘treelite’ and ‘treelite_runtime’ libraries must be installed for this to work.
- The options for handling missing values in ‘treelite’ are more limited. This function will always produce models that force `missing_action="impute"`, regardless of how the IsolationForest model itself handles them.
- The options for handling unseen categories in categorical variables are also more limited in ‘treelite’. It’s not possible to convert models that use `new_categ_action="weighted"`, and categories that were not present within the training data (which are not meant to be passed to ‘treelite’) will always be sent to the right side of the split, which might produce different results from `predict`.
- Some features such as range penalizations will not be kept in the ‘treelite’ model.

- While this library always uses C ‘double’ precision (typically ‘float64’) for model objects and prediction outputs, ‘treelite’ (a) can use ‘float32’ precision, (b) converts floating point numbers to a decimal representation and back to floating point; which combined can result in some precision loss which leads to producing slightly different predictions from the `predict` function in this package.
- If the model was fit to a `DataFrame` having a mixture of numerical and categorical columns, the resulting ‘treelite’ object will be built assuming all the numerical columns come before the categorical columns, regardless of which order they originally had in the data that was passed to ‘fit’. In such cases, it is possible to check the order of the columns under attributes `self.cols_numeric_` and `self.cols_categ_`.
- Categorical columns in ‘treelite’ are passed as integer values. if the model was fit to a `DataFrame` with categorical columns, the encoding that is used can be found under `self._cat_mapping`.
- The ‘treelite’ object returned by this function will not yet have been compiled. It’s necessary to call `compile` and `export_lib` afterwards in order to be able to use it.

Parameters

use_float32 (*bool*) – Whether to use ‘float32’ type for the model. This is typically faster but has less precision than the typical ‘float64’ (outside of this conversion, models from this library always use ‘float64’).

Returns

model – A ‘treelite’ model object.

Return type

obj

transform(X)

Impute missing values in the data using isolation forest model

Note: In order to use this functionality, the model must have been built with imputation capabilities (`‘build_imputer’ = ‘True’`).

Note: Categorical columns, if imputed with a model fit to a `DataFrame`, will always come out with pandas categorical dtype.

Note: The input may contain new columns (i.e. not present when the model was fitted), which will be output as-is.

Parameters

X (*array or array-like (n_samples, n_features)*) – Data for which missing values should be imputed. Can pass a NumPy array, Pandas `DataFrame`, or SciPy sparse CSR matrix.

If the model was fit to a `DataFrame` with categorical columns, must also be a `DataFrame`.

Returns

X_imputed – Object of the same type and dimensions as ‘X’, but with missing values already imputed. Categorical columns will be output as pandas’s ‘Categorical’ regardless of their dtype in ‘X’.

Return type

array or array-like (n_samples, n_features)

INDICES AND TABLES

- genindex
- modindex
- search

A

append_trees() (*isotree.IsolationForest* method), 24

B

build_indexer() (*isotree.IsolationForest* method), 25

C

copy() (*isotree.IsolationForest* method), 25

D

decision_function() (*isotree.IsolationForest* method), 25

drop_imputer() (*isotree.IsolationForest* method), 26

drop_indexer() (*isotree.IsolationForest* method), 26

drop_reference_points() (*isotree.IsolationForest* method), 26

E

export_model() (*isotree.IsolationForest* method), 26

F

fit() (*isotree.IsolationForest* method), 28

fit_predict() (*isotree.IsolationForest* method), 29

fit_transform() (*isotree.IsolationForest* method), 30

G

generate_sql() (*isotree.IsolationForest* method), 31

get_num_nodes() (*isotree.IsolationForest* method), 33

get_params() (*isotree.IsolationForest* method), 33

H

has_indexer_ (*isotree.IsolationForest* property), 33

has_reference_points_ (*isotree.IsolationForest* property), 33

I

import_model() (*isotree.IsolationForest* static method), 33

IsolationForest (class in *isotree*), 11

isotree

module, 11

M

module

isotree, 11

P

partial_fit() (*isotree.IsolationForest* method), 34

predict() (*isotree.IsolationForest* method), 35

predict_distance() (*isotree.IsolationForest* method), 37

predict_kernel() (*isotree.IsolationForest* method), 38

S

set_params() (*isotree.IsolationForest* method), 39

set_reference_points() (*isotree.IsolationForest* method), 39

subset_trees() (*isotree.IsolationForest* method), 40

T

to_treelite() (*isotree.IsolationForest* method), 40

transform() (*isotree.IsolationForest* method), 41