

Lab #6: Random Forests & Naïve Bayes

This lab is to deal with classification task using Random Forests and Naïve Bayes algorithms.

=====

RANDOM FORESTS

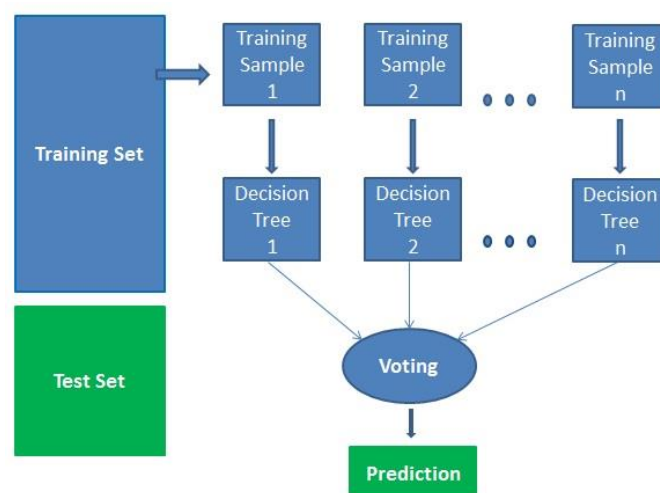
Random forests (RF) is a method by Leo Breiman (2001) for both classification and regression. Main idea: prediction is based on the combination of many decision trees, by taking the average of all individual predictions. Each tree in RF is simple but random. Each tree is grown differently, depending on the choices of the attributes and training data.

RF currently is one of the most popular and accurate methods [Fernández-Delgado et al., 2014]. It is also very general. RF can be implemented easily and efficiently. It can work with problems of very high dimensions, without overfitting.

How does the algorithm work?

It works in four steps:

1. Select random samples from a given dataset.
2. Construct a decision tree for each sample and get a prediction result from each decision tree.
3. Perform a vote for each predicted result.
4. Select the prediction result with the most votes as the final prediction.



RandomForestClassifier (<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>)

A random forest is a meta-estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Syntax:

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini',
max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0,
max_features='sqrt', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True,
oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False,
class_weight=None, ccp_alpha=0.0, max_samples=None)
```

where,

n_estimators: int, default=100

The number of trees in the forest.

Changed in version 0.22: The default value of `n_estimators` changed from 10 to 100 in 0.22.

criterion{“gini”, “entropy”, “log_loss”}, default=“gini”

The function to measure the quality of a split. Supported criteria are “gini” for the Gini impurity and “log_loss” and “entropy” both for the Shannon information gain, see [Mathematical formulation](#). Note: This parameter is tree-specific.

max_depth: int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than `min_samples_split` samples.

min_samples_split: int or float, default=2

The minimum number of samples required to split an internal node:

- If int, then consider `min_samples_split` as the minimum number.
- If float, then `min_samples_split` is a fraction and `ceil(min_samples_split * n_samples)` are the minimum number of samples for each split.

min_samples_leaf: int or float, default=1

The minimum number of samples required to be at a leaf node. A split point at any depth will only be considered if it leaves at least `min_samples_leaf` training samples

in each of the left and right branches. This may have the effect of smoothing the model, especially in regression.

- If int, then consider `min_samples_leaf` as the minimum number.
- If float, then `min_samples_leaf` is a fraction and `ceil(min_samples_leaf * n_samples)` are the minimum number of samples for each node.

min_weight_fraction_leaf: float, default=0.0

The minimum weighted fraction of the sum total of weights (of all the input samples) required to be at a leaf node. Samples have equal weight when `sample_weight` is not provided.

max_features: {"sqrt", "log2", None}, int or float, default="sqrt"

The number of features to consider when looking for the best split:

- If int, then consider `max_features` features at each split.
- If float, then `max_features` is a fraction and `round(max_features * n_features)` features are considered at each split.
- If "auto", then `max_features=sqrt(n_features)`.
- If "sqrt", then `max_features=sqrt(n_features)`.
- If "log2", then `max_features=log2(n_features)`.
- If None, then `max_features=n_features`.

Note: the search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than `max_features` features.

max_leaf_nodes: int, default=None

Grow trees with `max_leaf_nodes` in the best-first fashion. Best nodes are defined as relative reduction in impurity. If None then unlimited number of leaf nodes.

min_impurity_decrease: float, default=0.0

A node will be split if this split induces a decrease of the impurity greater than or equal to this value.

bootstrap: bool, default=True

Whether bootstrap samples are used when building trees. If False, the whole dataset is used to build each tree.

oob_score: bool, default=False

Whether to use out-of-bag samples to estimate the generalization score. Only available if `bootstrap=True`.

n_jobs: int, default=None

The number of jobs to run in parallel. [fit](#), [predict](#), [decision_path](#) and [apply](#) are all parallelized over the trees. None means 1 unless in a [joblib.parallel_backend](#) context. -1 means using all processors.

random_state: int, RandomState instance or None, default=None

Controls both the randomness of the bootstrapping of the samples used when building trees (if `bootstrap=True`) and the sampling of the features to consider when looking for the best split at each node (if `max_features < n_features`).

verbose: int, default=0

Controls the verbosity when fitting and predicting.

warm_start: bool, default=False

When set to `True`, reuse the solution of the previous call to fit and add more estimators to the ensemble, otherwise, just fit a whole new forest.

class_weight: {"balanced", "balanced_subsample"}, dict or list of dicts, default=None

Weights associated with classes in the form `{class_label: weight}`. If not given, all classes are supposed to have weight one. For multi-output problems, a list of dicts can be provided in the same order as the columns of `y`.

Note that for multioutput (including multilabel) weights should be defined for each class of every column in its own dict. For example, for four-class multilabel classification weights should be `[[0: 1, 1: 1], {0: 1, 1: 5}, {0: 1, 1: 1}, {0: 1, 1: 1}]` instead of `[[1:1], {2:5}, {3:1}, {4:1}]`.

The "balanced" mode uses the values of `y` to automatically adjust weights inversely proportional to class frequencies in the input data as `n_samples / (n_classes * np.bincount(y))`

The "balanced_subsample" mode is the same as "balanced" except that weights are computed based on the bootstrap sample for every tree grown.

For multi-output, the weights of each column of `y` will be multiplied.

Note that these weights will be multiplied with `sample_weight` (passed through the fit method) if `sample_weight` is specified.

ccp_alpha: non-negative float, default=0.0

Complexity parameter used for Minimal Cost-Complexity Pruning. The subtree with the largest cost complexity that is smaller than `ccp_alpha` will be chosen. By default, no pruning is performed.

max_samples: int or float, default=None

If bootstrap is True, the number of samples to draw from X to train each base estimator.

- If None (default), then draw `X.shape[0]` samples.
- If int, then draw `max_samples` samples.
- If float, then draw `max_samples * X.shape[0]` samples. Thus, `max_samples` should be in the interval `(0.0, 1.0]`.

Usage:

```
from sklearn.ensemble import RandomForestClassifier
#Create a Random Forest Classifier
clf=RandomForestClassifier(n_estimators=100)

#Train the model using the training sets
clf.fit(X_train,y_train)
```

NAÏVE BAYES

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of conditional independence between every pair of features given the value of the class variable. Bayes' theorem states the following relationship, given class variable y and dependent feature vector x_i through x_n :

$$P(y | x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n | y)}{P(x_1, \dots, x_n)}$$

Using the naive conditional independence assumption that

$$P(x_i | y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i | y),$$

for all i , this relationship is simplified to

$$P(y | x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i | y)}{P(x_1, \dots, x_n)}$$

- **Sklearn** provides 5 types of Naive Bayes :
 - **GaussianNB**: Naive Bayes work for **continuous features**

```
# GaussianNB
from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

- **CategoricalNB**: Naive Bayes work for **categorical features**

```
# CategoricalNB
from sklearn.naive_bayes import CategoricalNB
clf = CategoricalNB()
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
```

- **BernoulliNB**: Naive Bayes work for **binary features**. BernoulliNB implements the naive Bayes training and classification algorithms for data that is distributed according to multivariate Bernoulli distributions; i.e., there may be multiple features but each one is assumed to be a binary-valued (Bernoulli, boolean) variable.

```
# BernoulliNB
from sklearn.naive_bayes import BernoulliNB
clf = BernoulliNB()
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
```

- **MultinomialNB**: Naive Bayes work for **multinomial features**. The multinomial Naive Bayes classifier is suitable for classification with discrete features (e.g., word counts for text classification).

```
# MultinomialNB
from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB()
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
```

- **ComplementNB**: Naive Bayes work for multinomial features. It is very similar to Multinomial Naive Bayes due to the parameters but seems to be more powerful in the case of an **imbalanced dataset**.

```
# ComplementNB
from sklearn.naive_bayes import ComplementNB
model = ComplementNB()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```