

- Additional supervised learning techniques
- Tune tree-based models
- 📺

Video: Bootstrap aggregation

4 min

📖

Reading: Bagging: How it works and why to use it

20 min

📺

Video: Explains a random forest

2 min

📖

Reading: More about random forests

20 min

📺

Video: Tuning a random forest

3 min

📖

Reading: Following instructions: Build and cross-validate a random forest model with Python

20 min

📄

Lab: Annotated following guide: Build and cross-validate a random forest model

20 min

📖

Video: Build and cross-validate a random forest model with Python

4 min

📺

Video: Build and validate a random forest model using a validation data set

7 min

📖

Reading: Reference guide: Random forest tuning

20 min

📄

Lab: Activity: Build a random forest model

1h

📄

Lab: Exemplar: Build a random forest model

20 min

📖

Reading: Case Study: Machine Learning model unearths recurring insights for Booz Allen Hamilton

20 min

📖

Practice Quiz: Test your knowledge: Bagging

4 questions

Boosting

Review: Tree-based modeling
- ## Reference guide: Random forest tuning
- ### Reference guide: Random forest tuning
- Previously, you learned about random forest models and studied how to build and tune them. This reading is a quick-reference guide to help you when you're building models of your own. It includes:
- Import statements
  - Hyperparameters
- ### Save this course item
- You may want to save a copy of this guide for future reference. You can use it as a resource for additional practice or in your future professional projects. To access a downloadable version of this course item, click the link below and select "Use Template."
- [Reference guide: Random forest tuning](#)
- OR
- If you don't have a Google account, you can download the item directly from the attachment below.
- [Reference guide\\_Random forest tuning](#)  
DOCX File
- ### Import statements
- #### Models
- For classification tasks:
- ```
from sklearn.ensemble import RandomForestClassifier
```
- For regression tasks:
- ```
from sklearn.ensemble import RandomForestRegressor
```
- #### Evaluation metrics
- For classification tasks:
- ```
from sklearn.metrics import
```
- |                                                                        |                                                                                                  |
|------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------|
| <b>accuracy_score</b> $\mathcal{O}(y\_true, y\_pred, "1", ...)$        | Accuracy classification score.                                                                   |
| <b>average_precision_score</b> $\mathcal{O}(y\_true, ...)$             | Compute average precision (AP) from prediction scores.                                           |
| <b>confusion_matrix</b> $\mathcal{O}(y\_true, y\_pred, *)$             | Compute confusion matrix to evaluate the performance of the training of a model.                 |
| <b>f1_score</b> $\mathcal{O}(y\_true, y\_pred, "1", ...)$              | Compute the F1 score, also known as balanced F-score or F-measure.                               |
| <b>fbeta_score</b> $\mathcal{O}(y\_true, y\_pred, *, beta)$            | Compute the F-beta score.                                                                        |
| <b>metrics.log_loss</b> $\mathcal{O}(y\_true, y\_pred, "1", ops, ...)$ | Log loss, aka logistic loss or cross-entropy loss.                                               |
| <b>multilabel_confusion_matrix</b> $\mathcal{O}(y\_true, ...)$         | Compute a confusion matrix for each class or sample.                                             |
| <b>precision_recall_curve</b> $\mathcal{O}(y\_true, ...)$              | Compute precision-recall pairs for different probability thresholds.                             |
| <b>precision_score</b> $\mathcal{O}(y\_true, y\_pred, "1", ...)$       | Compute the precision.                                                                           |
| <b>recall_score</b> $\mathcal{O}(y\_true, y\_pred, "1", ...)$          | Compute the recall.                                                                              |
| <b>roc_auc_score</b> $\mathcal{O}(y\_true, y\_score, "1", ...)$        | Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores. |
- For regression tasks:
- ```
from sklearn.metrics import
```
- |  |   |
|--|---|
| <b>mean_absolute_error</b> $\mathcal{O}(y\_true, y\_pred, *)$    | Mean absolute error regression loss.                            |
| <b>mean_squared_error</b> $\mathcal{O}(y\_true, y\_pred, *)$     | Mean squared error regression loss.                             |
| <b>mean_squared_log_error</b> $\mathcal{O}(y\_true, y\_pred, *)$ | Mean squared logarithmic error regression loss.                 |
| <b>median_absolute_error</b> $\mathcal{O}(y\_true, y\_pred, *)$  | Median absolute error regression loss.                          |
| <b>mean_absolute_percentage_error</b> $\mathcal{O}(...)$         | Mean absolute percentage error (MAPE) regression loss.          |
| <b>r2_score</b> $\mathcal{O}(y\_true, y\_pred, "1", ...)$        | $R^2$ (coefficient of determination) regression score function. |
- ### Hyperparameters
- The following are some of the most important hyperparameters for random forest classification models in scikit-learn.
- | Hyperparameter          | What it does   | Input type                          | Default Value | Considerations  |
|-------------------------|--|-------------------------------------|---------------|---|
| <b>n_estimators</b>     | Specifies the number of trees your model will build in its ensemble  | int                                 | 100           | A typical range is 50-500. Consider how much data you have, how deep the trees are allowed to grow and how many samples are bootstrapped to grow each tree (you generally need more trees if they're shallow, and more trees if your bootstrap sample size is smaller). Also consider if your use case has latency requirements.  |
| <b>max_depth</b>        | Specifies how many levels your tree can have.<br><br>If None, trees grow until leaves are pure or until all leaves contain less than <b>min_sample_split</b> samples.  | int                                 | None          | Random forest models often use base learners that are fully grown, but restricting tree depth can reduce train/latency times and prevent overfitting. If not None, consider values 3-8.   |
| <b>min_sample_split</b> | Controls threshold below which nodes become leaves<br><br>If float, then it represents a percentage (0-1) of <b>max_samples</b> .  | int or float                        | 2             | Consider (a) how many samples are in your dataset, and (b) how much of that data you're allowing each base learner to use (i.e., the value of the <b>max_samples</b> hyperparameter). The fewer samples available, the lesser the number of samples may need to be allowed in each leaf node (otherwise the tree would be very shallow).  |
| <b>min_samples_leaf</b> | A split can only occur if it guarantees a minimum of this number of observations in each resulting node.<br><br>If float, then it represents a percentage (0-1) of <b>max_samples</b> .  | int or float                        | 1             | Consider (a) how many samples are in your dataset, and (b) how much of that data you're allowing each base learner to use (i.e., the value of the <b>max_samples</b> hyperparameter). The fewer samples available, the lesser the number of samples may need to be allowed in each leaf node (otherwise the tree would be very shallow).  |
| <b>max_features</b>     | Specifies the number of features that each tree randomly selects during training.<br><br>If int, then consider <b>max_features</b> features at each split.<br><br>If float, then <b>max_features</b> is a fraction and $\text{round}(\text{max\_features} * \text{n\_features})$ features are considered at each split.<br><br>If "sqrt", then <b>max_features</b> = $\text{sqrt}(\text{n\_features})$ .<br><br>If "log2", then <b>max_features</b> = $\text{log2}(\text{n\_features})$ .<br><br>If None, then <b>max_features</b> = <b>n_features</b> . | "sqrt", "log2", None, int or float, | "sqrt"        | Consider how many features the dataset has and how many trees will be grown. Fewer features sampled during each bootstrap means more base learners would be needed. Small <b>max_features</b> values on datasets with many features mean more unrepresentative trees in the ensemble.   |
| <b>max_samples</b>      | Specifies the number of samples bootstrapped from the dataset to train each base model<br><br>If float, then it represents a percentage (0-1) of the dataset.<br><br>If None, then draw <b>X.shape[0]</b> samples.   | int or float                        | None          | Consider the size of your dataset. When working with large datasets, it can be beneficial to limit the number of samples in each tree, because doing so can greatly reduce training time and yet still result in a robust model. For example, 20% of 1 billion may be enough to capture patterns in the data, but if you only have 1,000 samples in your dataset, you'll probably need to use them all. |
- \* Note that **max\_samples** was not used in the "Build and validate/cross-validate a random forest model" videos and notebook, but is included here so you can use this hyperparameter in your own work. Remember that using fractions of the data to train each base learner can possibly improve model predictions and certainly speed up execution times.
- ### Key takeaways
- When building machine learning models, it's essential to have the right tools and to understand how to use them. While there are numerous other hyperparameters to explore, the ones in this reference guide are among the most important. Be inquisitive and try different approaches. Discovering ways to improve your model is a lot of fun!
- ### Resources for more information
- More detailed information about random forest tuning can be found here:
- scikit-learn documentation:
    - [Model metrics](#)
    - [Random forest classifier](#): documentation for model used for classification tasks
    - [Random forest regressor](#): documentation for model used for regression tasks
- Mark as completed
-