

Boosting

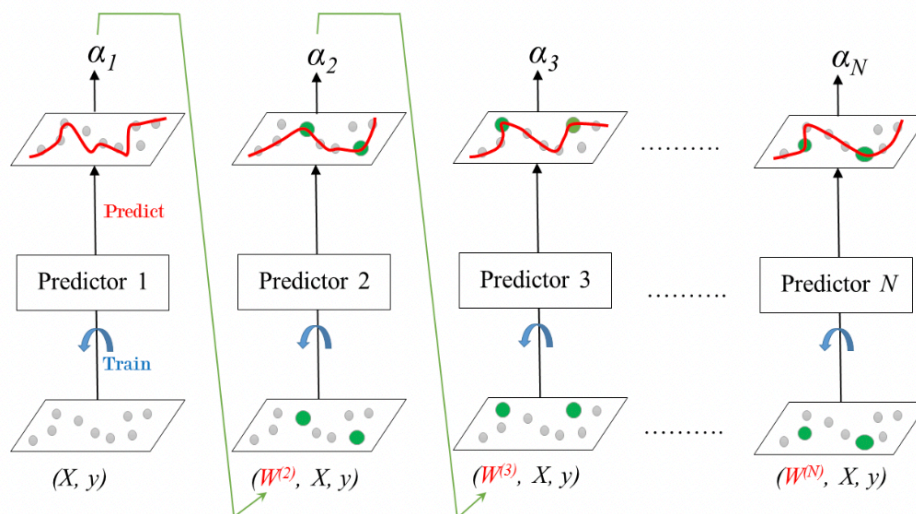
Boosting

- **Boosting:** Ensemble method combining several weak learners to form a strong learner.
- **Weak learner:** Model doing slightly better than random guessing.
- Example of weak learner: Decision stump (CART whose maximum depth is 1).

Boosting

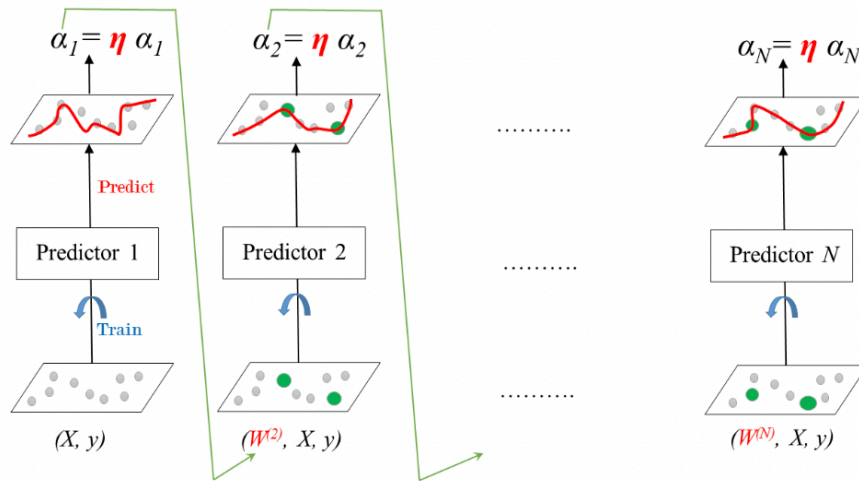
- Train an ensemble of predictors sequentially.
- Each predictor tries to correct its predecessor.
- Most popular boosting methods:
 - AdaBoost,
 - Gradient Boosting.

AdaBoost: Training



Learning Rate

Learning rate: $0 < \eta \leq 1$



- Smaller value for learning rate should be compensated by an increase in the number of estimators

AdaBoost: Prediction

- Classification:
 - Weighted majority voting.
 - In sklearn: `AdaBoostClassifier` .
- Regression:
 - Weighted average.
 - In sklearn: `AdaBoostRegressor` .

AdaBoost Classification in sklearn (Breast Cancer dataset)

```
# Import models and utility functions
from sklearn.ensemble import AdaBoostClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import train_test_split

# Set seed for reproducibility
SEED = 1

# Split data into 70% train and 30% test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
                                                    stratify=y,
                                                    random_state=SEED)
```

```
# Instantiate a classification-tree 'dt'
dt = DecisionTreeClassifier(max_depth=1, random_state=SEED)

# Instantiate an AdaBoost classifier 'adb_clf'
adb_clf = AdaBoostClassifier(base_estimator=dt, n_estimators=100)

# Fit 'adb_clf' to the training set
adb_clf.fit(X_train, y_train)

# Predict the test set probabilities of positive class
y_pred_proba = adb_clf.predict_proba(X_test)[:,-1]

# Evaluate test-set roc_auc_score
adb_clf_roc_auc_score = roc_auc_score(y_test, y_pred_proba)
```

AdaBoost Classification in sklearn (Breast Cancer dataset)

```
# Print adb_clf_roc_auc_score
print('ROC AUC score: {:.2f}'.format(adb_clf_roc_auc_score))
```

```
ROC AUC score: 0.99
```

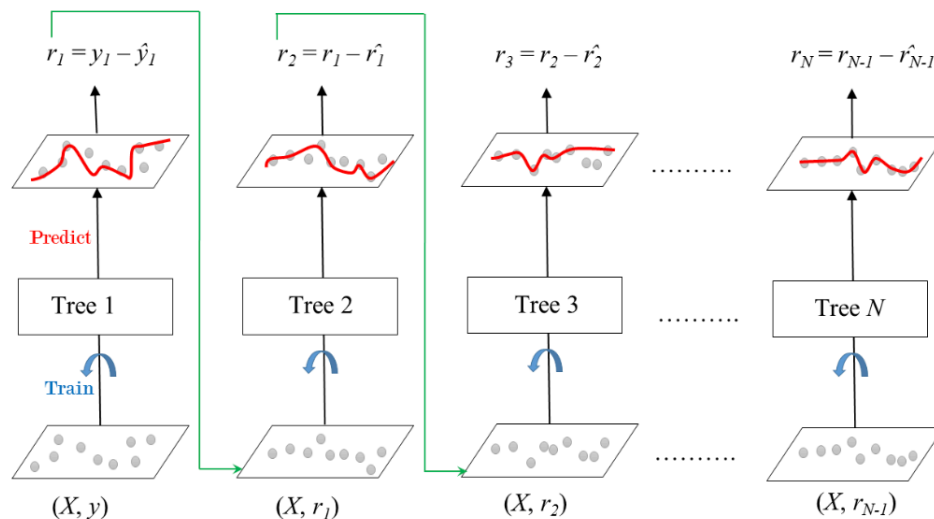
AUC - Area Under the Curve

Gradient Boosting (GB)

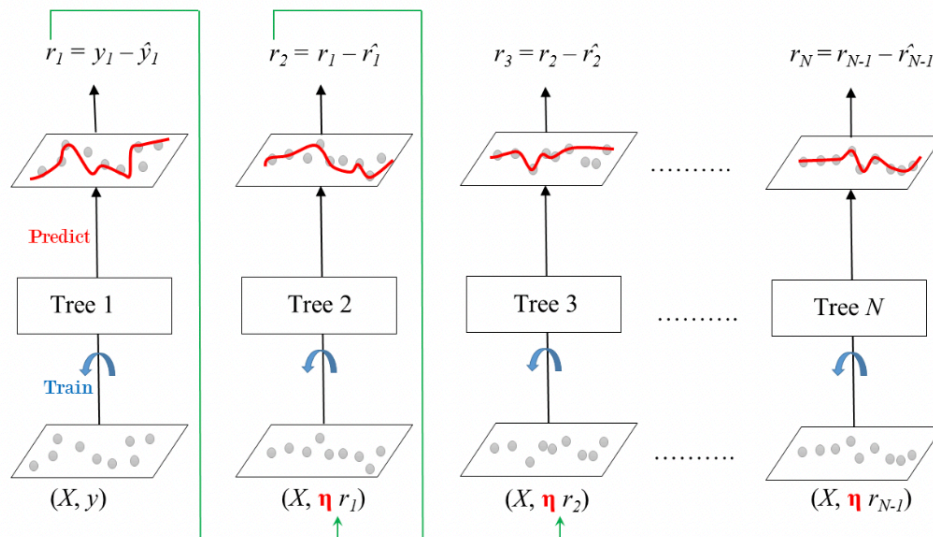
Gradient Boosted Trees

- Sequential correction of predecessor's errors.
- Does not tweak the weights of training instances.
- Fit each predictor is trained using its predecessor's residual errors as labels.
- Gradient Boosted Trees: a CART is used as a base learner.

Gradient Boosted Trees for Regression: Training



Shrinkage



Gradient Boosted Trees: Prediction

- Regression:
 - $y_{pred} = y_1 + \eta r_1 + \dots + \eta r_N$
 - In sklearn: `GradientBoostingRegressor`.
- Classification:
 - In sklearn: `GradientBoostingClassifier`.

Gradient Boosting in sklearn (auto dataset)

```
# Import models and utility functions
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error as MSE

# Set seed for reproducibility
SEED = 1

# Split dataset into 70% train and 30% test
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.3,
                                                    random_state=SEED)
```

```
# Instantiate a GradientBoostingRegressor 'gbr'
gbr = GradientBoostingRegressor(n_estimators=300, max_depth=1, random_state=SEED)

# Fit 'gbr' to the training set
gbr.fit(X_train, y_train)

# Predict the test set labels
y_pred = gbr.predict(X_test)

# Evaluate the test set RMSE
rmse_test = MSE(y_test, y_pred)**(1/2)

# Print the test set RMSE
print('Test set RMSE: {:.2f}'.format(rmse_test))
```

```
Test set RMSE: 4.01
```

Stochastic Gradient Boosting (SGB)

Gradient Boosting: Cons

- GB involves an exhaustive search procedure.
- Each CART is trained to find the best split points and features.
- May lead to CARTs using the same split points and maybe the same features.

- Each tree is trained on a random subset of rows of the training data.
- The sampled instances (40%-80% of the training set) are sampled without replacement.
- Features are sampled (without replacement) when choosing split points.
- Result: further ensemble diversity.
- Effect: adding further variance to the ensemble of trees.

The diagram illustrates the training process for a single tree in a boosting framework. It shows a sequence of operations:

- Input Data:** A set of sample instances (X, y) is shown at the bottom.
- Sampling:** A subset of these instances is sampled to form $(X_{sampled}, y_{sampled})$.
- Training:** The sampled data is used to **Train** **Tree 1**, indicated by a blue circular arrow.
- Feature Extraction:** From the trained tree, **Sample features at each split** are extracted, represented by orange circles labeled f_1, f_2, \dots, f_N .
- Prediction:** The tree is used to **Predict** the output \hat{y}_l for the original data X .
- Residual Calculation:** The residual $r_l = y_l - \hat{y}_l$ is calculated for each instance.

 A green line on the right side of the diagram indicates the progression to the next step in the boosting process, where the next tree is trained on the residuals r_l and the previous model's output \hat{y}_l is combined to form the next model's output \hat{y}_{l+1} .

[illegible]

Stochastic Gradient Boosting in sklearn (auto dataset)

```
# Instantiate a stochastic GradientBoostingRegressor 'sgbt'
sgbt = GradientBoostingRegressor(max_depth=1,
                                subsample=0.8,
                                max_features=0.2,
                                n_estimators=300,
                                random_state=SEED)

# Fit 'sgbt' to the training set
sgbt.fit(X_train, y_train)

# Predict the test set labels
y_pred = sgbt.predict(X_test)
```

max_features = 0.2 - Each tree uses a maximum of 20% of the available features to perform the best split

Stochastic Gradient Boosting in sklearn (auto dataset)

```
# Evaluate test set RMSE 'rmse_test'
rmse_test = MSE(y_test, y_pred)**(1/2)

# Print 'rmse_test'
print('Test set RMSE: {:.2f}'.format(rmse_test))
```

```
Test set RMSE: 3.95
```

Tuning a CART's Hyperparameters

Hyperparameters

Machine learning model:

- **parameters:** learned from data
 - CART example: split-point of a node, split-feature of a node, ...
- **hyperparameters:** not learned from data, set prior to training
 - CART example: max_depth, min_samples_leaf, splitting criterion ...

What is hyperparameter tuning?

- **Problem:** search for a set of optimal hyperparameters for a learning algorithm.
- **Solution:** find a set of optimal hyperparameters that results in an optimal model.
- **Optimal model:** yields an optimal **score**.
- **Score:** in sklearn defaults to accuracy (classification) and R^2 (regression).
- Cross validation is used to estimate the generalization performance.

Why tune hyperparameters?

- In `sklearn`, a model's default hyperparameters are not optimal for all problems.
- Hyperparameters should be tuned to obtain the best model performance.

Approaches to hyperparameter tuning

- Grid Search
- Random Search
- Bayesian Optimization
- Genetic Algorithms
-

Grid search cross validation

- Manually set a grid of discrete hyperparameter values.
- Set a metric for scoring model performance.
- Search exhaustively through the grid.
- For each set of hyperparameters, evaluate each model's CV score.
- The optimal hyperparameters are those of the model achieving the best CV score.

Grid search cross validation: example

- Hyperparameters grids:
 - `max_depth` = {2,3,4},
 - `min_samples_leaf` = {0.05, 0.1}
- hyperparameter space = { (2,0.05) , (2,0.1) , (3,0.05), ... }
- CV scores = { $score_{(2,0.05)}$, ... }
- optimal hyperparameters = set of hyperparameters corresponding to the best CV score.

Inspecting the hyperparameters of a CART in sklearn

```
# Import DecisionTreeClassifier
from sklearn.tree import DecisionTreeClassifier

# Set seed to 1 for reproducibility
SEED = 1

# Instantiate a DecisionTreeClassifier 'dt'
dt = DecisionTreeClassifier(random_state=SEED)
```

Inspecting the hyperparameters of a CART in sklearn

```
# Print out 'dt's hyperparameters
print(dt.get_params())
```

```
{'class_weight': None,
 'criterion': 'gini',
 'max_depth': None,
 'max_features': None,
 'max_leaf_nodes': None,
 'min_impurity_decrease': 0.0,
 'min_impurity_split': None,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'min_weight_fraction_leaf': 0.0,
 'presort': False,
 'random_state': 1,
 'splitter': 'best'}
```

```

# Import GridSearchCV
from sklearn.model_selection import GridSearchCV
# Define the grid of hyperparameters 'params_dt'
params_dt = {
    'max_depth': [3, 4, 5, 6],
    'min_samples_leaf': [0.04, 0.06, 0.08],
    'max_features': [0.2, 0.4, 0.6, 0.8]
}
# Instantiate a 10-fold CV grid search object 'grid_dt'
grid_dt = GridSearchCV(estimator=dt,
                        param_grid=params_dt,
                        scoring='accuracy',
                        cv=10,
                        n_jobs=-1)
# Fit 'grid_dt' to the training data
grid_dt.fit(X_train, y_train)

```

Extracting the best hyperparameters

```

# Extract best hyperparameters from 'grid_dt'
best_hyperparams = grid_dt.best_params_
print('Best hyperparameters:\n', best_hyperparams)

```

```

Best hyperparameters:
{'max_depth': 3, 'max_features': 0.4, 'min_samples_leaf': 0.06}

```

```

# Extract best CV score from 'grid_dt'
best_cv_score = grid_dt.best_score_
print('Best CV accuracy'.format(best_cv_score))

```

```

Best CV accuracy: 0.938

```

Extracting the best estimator

```
# Extract best model from 'grid_dt'
best_model = grid_dt.best_estimator_

# Evaluate test set accuracy
test_acc = best_model.score(X_test, y_test)

# Print test set accuracy
print("Test set accuracy of best model: {:.3f}".format(test_acc))
```

```
Test set accuracy of best model: 0.947
```

Tuning a RF's Hyperparameters

Random Forests Hyperparameters

- CART hyperparameters
- number of estimators
- bootstrap
-

Tuning is expensive

Hyperparameter tuning:

- computationally expensive,
- sometimes leads to very slight improvement,

Weight the impact of tuning on the whole project.

Inspecting RF Hyperparameters in sklearn

```
# Import RandomForestRegressor
from sklearn.ensemble import RandomForestRegressor

# Set seed for reproducibility
SEED = 1

# Instantiate a random forests regressor 'rf'
rf = RandomForestRegressor(random_state= SEED)
```

```
# Inspect rf' s hyperparameters
rf.get_params()
```

```
{'bootstrap': True,
 'criterion': 'mse',
 'max_depth': None,
 'max_features': 'auto',
 'max_leaf_nodes': None,
 'min_impurity_decrease': 0.0,
 'min_impurity_split': None,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'min_weight_fraction_leaf': 0.0,
 'n_estimators': 10,
 'n_jobs': -1,
 'oob_score': False,
 'random_state': 1,
 'verbose': 0,
 'warm_start': False}
```

```

# Basic imports
from sklearn.metrics import mean_squared_error as MSE
from sklearn.model_selection import GridSearchCV
# Define a grid of hyperparameter 'params_rf'
params_rf = {
    'n_estimators': [300, 400, 500],
    'max_depth': [4, 6, 8],
    'min_samples_leaf': [0.1, 0.2],
    'max_features': ['log2', 'sqrt']
}
# Instantiate 'grid_rf'
grid_rf = GridSearchCV(estimator=rf,
                       param_grid=params_rf,
                       cv=3,
                       scoring='neg_mean_squared_error',
                       verbose=1,
                       n_jobs=-1)

```

Searching for the best hyperparameters

```

# Fit 'grid_rf' to the training set
grid_rf.fit(X_train, y_train)

```

```

Fitting 3 folds for each of 36 candidates, totalling 108 fits
[Parallel(n_jobs=-1)]: Done 42 tasks      | elapsed: 10.0s
[Parallel(n_jobs=-1)]: Done 108 out of 108 | elapsed: 24.3s finished
RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=4,
                       max_features='log2', max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=0.1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, n_estimators=400, n_jobs=1,
                       oob_score=False, random_state=1, verbose=0, warm_start=False)

```

Extracting the best hyperparameters

```
# Extract the best hyperparameters from 'grid_rf'
best_hyperparams = grid_rf.best_params_

print('Best hyperparameters:\n', best_hyperparams)
```

```
Best hyperparameters:
{'max_depth': 4,
 'max_features': 'log2',
 'min_samples_leaf': 0.1,
 'n_estimators': 400}
```

Evaluating the best model performance

```
# Extract the best model from 'grid_rf'
best_model = grid_rf.best_estimator_
# Predict the test set labels
y_pred = best_model.predict(X_test)
# Evaluate the test set RMSE
rmse_test = MSE(y_test, y_pred)**(1/2)
# Print the test set RMSE
print('Test set RMSE of rf: {:.2f}'.format(rmse_test))
```

```
Test set RMSE of rf: 3.89
```