

Part I: Overview of AdaBoost

Boosting is an ensemble learning technique designed to improve model accuracy by combining weak learners, classifiers that perform only slightly better than random guessing, to create a powerful model. AdaBoost, short for Adaptive Boosting, was introduced by Freund and Schapire (1997) and has since become one of the most popular algorithms in this category. AdaBoost constructs a series of classifiers by iteratively focusing on samples that previous classifiers struggled to classify correctly (Freund and Schapire, 1997). This iterative approach enables the model to progressively improve accuracy by combining the knowledge of multiple classifiers.

Algorithm Overview

The main concept behind AdaBoost is that, given a base learner (often a decision stump), it can sequentially train weak learners on different distributions of the training data, emphasizing examples that were misclassified in previous rounds (Shalev-Shwartz and Ben-David, 2014). At each iteration, AdaBoost increases the weight of incorrectly classified examples, making the algorithm more "sensitive" to hard-to-classify cases. Once a new weak learner is trained, AdaBoost adjusts its weight based on its accuracy, giving more influence to learners that make fewer errors. The final model is a weighted sum of the individual classifiers, where each classifier's weight corresponds to its accuracy, forming a strong ensemble.

Representation

In AdaBoost, the final hypothesis $H(x)$ is represented as a weighted combination of the weak learners $h_t(x)$:

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

where:

- $h_t(x)$ is the prediction of the weak learner at iteration t ,
- α_t is the weight assigned to hypothesis h_t , calculated based on the classification accuracy at round t .

This weighted combination allows AdaBoost to leverage the accuracy of each weak learner, forming a strong final classifier that benefits from the collective performance of all weak learners (Freund and Schapire, 1997).

Loss Function

AdaBoost minimizes an exponential loss function, which is given by:

$$L(H) = \sum_{i=1}^m \exp(-y_i H(x_i))$$

where:

- $y_i \in \{-1, +1\}$ is the true label for data point x_i ,
- $H(x_i)$ represents the combined prediction from all weak learners at x_i .

The exponential loss penalizes misclassified examples heavily, encouraging the model to focus on these challenging cases, which makes AdaBoost responsive to difficult examples in the data (Freund and Schapire, 1997).

Optimizer

AdaBoost does not employ traditional gradient-based optimization but instead adjusts the distribution D of the training data based on the weak learner's performance in each round. After each iteration, the distribution D is updated as follows:

$$D_{t+1}(i) = \frac{D_t(i) \cdot \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

where:

- $D_t(i)$ is the weight of sample i at round t ,
- $\alpha_t = 0.5 \cdot \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$ is the weight for the weak learner h_t ,
- Z_t is a normalization factor ensuring that D_{t+1} remains a valid probability distribution.

This reweighting process is central to AdaBoost's adaptive capability, progressively focusing on difficult samples and dynamically adjusting to the data's complexity (Shalev-Shwartz and Ben-David, 2014).

Advantages and Applications

AdaBoost has several notable advantages:

- Improvement of Weak Learners:** Since it combines weak classifiers, AdaBoost significantly enhances the overall performance of each individual classifier, even if each performs only slightly better than random (Freund & Schapire, 1997).
- Reduction of Bias and Variance:** AdaBoost is known for its ability to reduce bias and variance, enhancing robustness across diverse datasets (Shalev-Shwartz and Ben-David, 2014).
- Versatility Across Domains:** AdaBoost adapts well across domains, and it is used for both classification and regression tasks. It has been applied to various fields, from face recognition and text classification to medical diagnosis, where small incremental improvements in prediction accuracy are highly valuable.

In practical applications, AdaBoost's adaptive nature makes it particularly effective in tasks with complex and high-dimensional data. For instance, Freund and Schapire applied AdaBoost to face detection, one of the early high-impact uses of the algorithm (Freund and Schapire, 1997). Its ability to adaptively combine classifiers with weighted predictions has made it popular in scenarios where the cost of errors is high, such as in medical diagnostics or fraud detection (Shalev-Shwartz and Ben-David, 2014).

Disadvantages

While AdaBoost is powerful, it has limitations:

- **Sensitivity to Noise:** AdaBoost's emphasis on misclassified samples can amplify the effect of noisy data. If a sample is mislabeled or contains outliers, AdaBoost may focus on these cases excessively, potentially destabilizing the model (Shalev-Shwartz and Ben-David, 2014).
- **Computational Intensity:** The iterative training process requires substantial computational power, especially with large datasets or high-dimensional feature spaces. This makes AdaBoost resource-intensive, particularly when using complex base learners (Freund and Schapire, 1997).

References

Freund, Y. and Schapire, R.E., 1997. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. Journal of Computer and System Sciences, 55(1), pp.119-139.

Shalev-Shwartz, S. and Ben-David, S., 2014. Understanding Machine Learning: From Theory to Algorithms. Cambridge University Press.

Pseudocode

```
# Pseudocode for AdaBoost

# Input: Training set S = {(x1, y1), (x2, y2), ..., (xm, ym)}, weak learner WL, number of rounds T
# Output: Final hypothesis H(x)
```

```
Initialize distribution D1(i) = 1/m for all i = 1 to m

for t = 1 to T:
    1. Train weak learner h_t using distribution Dt
    2. Calculate error ε_t = sum(Dt(i) * [h_t(x_i) != y_i]) for all i
    3. Compute α_t = 0.5 * log((1 - ε_t) / ε_t)
    4. Update distribution:
        Dt+1(i) = Dt(i) * exp(-α_t * y_i * h_t(x_i))
        Normalize Dt+1 to maintain a probability distribution
```

```
Output final hypothesis:
H(x) = sign(sum(α_t * h_t(x) for t = 1 to T))
```

Part II: Adaboost Stencil Code

```
In [ ]: import numpy as np
import pandas as pd
import os

# DecisionStump and myAdaBoost classes remain the same
class DecisionStump:
    def __init__(self):
        self.polarity = 1
        self.threshold = None
        self.feature_idx = None
        self.alpha = None

    def predict(self, X):
        n_samples = X.shape[0]
        X_c = X[:, self.feature_idx]
        preds = np.ones(n_samples)
        if self.polarity == 1:
            preds[X_c < self.threshold] = -1
        else:
            preds[X_c > self.threshold] = -1
        return preds

class myAdaBoost:
    def __init__(self, n_clf=50):
        self.n_clf = n_clf

    def fit(self, X, y):
        n_samples, n_features = X.shape
        w = np.full(n_samples, (1 / n_samples))
        self.clfs = []
        for _ in range(self.n_clf):
            clf = DecisionStump()
            min_error = float('inf')
            for feat in range(n_features):
                X_c = X[:, feat]
                thresholds = np.unique(X_c)
                for threshold in thresholds:
                    for p in [1, -1]:
                        preds = np.ones(n_samples)
                        if p == 1:
                            preds[X_c < threshold] = -1
                        else:
```

```

        preds[X_c > threshold] = -1
        error = np.sum(w[y != preds])
        if error < min_error:
            min_error = error
            clf.threshold = threshold
            clf.feature_idx = feat
            clf.polarity = p

    EPS = 1e-10
    clf.alpha = 0.5 * np.log((1.0 - min_error + EPS) / (min_error + EPS))
    # Update weights
    preds = clf.predict(X)
    w *= np.exp(-clf.alpha * y * preds)
    w /= np.sum(w)
    self.clfs.append(clf)

def predict(self, X):
    clf_preds = [clf.alpha * clf.predict(X) for clf in self.clfs]
    y_pred = np.sum(clf_preds, axis=0)
    y_pred = np.sign(y_pred)
    return y_pred.astype(int)

```

Part III: Check Model

Unit Test

We will use below unit test to ensure our methods and class work properly and handles corner cases.

```

In [ ]: # Import necessary libraries
import numpy as np

# Test 1: DecisionStump predict with correct threshold and polarity=1
stump = DecisionStump()
stump.threshold = 0.5
stump.feature_idx = 0
stump.polarity = 1
X_test = np.array([[0.3], [0.7], [0.9]])
# Expected output: values < 0.5 are -1, others are 1
expected = np.array([-1, 1, 1])
assert np.array_equal(stump.predict(X_test), expected), "Test 1 Failed"

# Test 2: DecisionStump predict with correct threshold and polarity=-1
stump.polarity = -1
# Expected output: values > 0.5 are -1, others are 1
expected = np.array([1, -1, -1])
assert np.array_equal(stump.predict(X_test), expected), "Test 2 Failed"

# Test 3: Edge case - DecisionStump predict with no samples
X_empty = np.array([]).reshape(0, 1) # Empty input
expected_empty = np.array([])
assert np.array_equal(stump.predict(X_empty), expected_empty), "Test 3 Failed"

# Test 4: myAdaBoost fit and predict with perfectly separable data
X_simple = np.array([[0], [1]])
y_simple = np.array([-1, 1])
model = myAdaBoost(n_clf=1)
model.fit(X_simple, y_simple)
# Expecting perfect prediction
y_pred_simple = model.predict(X_simple)
assert np.array_equal(y_pred_simple, y_simple), "Test 4 Failed"

# Test 5: myAdaBoost fit and predict on noisy data
X_noisy = np.array([[0], [1], [2], [3]])
y_noisy = np.array([-1, -1, 1, 1])
model_noisy = myAdaBoost(n_clf=3)
model_noisy.fit(X_noisy, y_noisy)
y_pred_noisy = model_noisy.predict(X_noisy)
assert np.mean(y_pred_noisy == y_noisy) >= 0.75, "Test 5 Failed"

```

Experiment

The objective of this section is to verify the performance of our custom AdaBoost implementation against scikit-learn’s AdaBoost classifier, ensuring that both achieve comparable results on this dataset. This comparison aligns with prior research, such as that by Street et al. (1993), where nuclear features from this dataset supported accurate cancer diagnosis.

Breast Cancer Wisconsin Dataset

The **Breast Cancer Wisconsin (Diagnostic) dataset** is widely used to benchmark binary classification models due to its well-defined features and relevance to medical diagnostics. This dataset, hosted by the UCI Machine Learning Repository, consists of 569 samples, each labeled as either **malignant** or **benign** based on cellular features from breast mass images (Dua and Graff, 2019).

Methodology

- Data Loading and Preprocessing:** We load the dataset, converting the `Diagnosis` column to a binary format, where **1** represents malignant and **0** represents benign.
- Training:** We train our AdaBoost model using decision stumps (single-depth decision trees) as weak learners, with 50 estimators. For comparison, we train an AdaBoost classifier from scikit-learn under the same conditions.

3. **Evaluation:** Model performance is evaluated based on **Accuracy** and **F1 Score** to capture both precision and recall, critical metrics in medical diagnostics.

Results

Both our AdaBoost implementation and scikit-learn's achieved:

- **Accuracy:** ~98%
- **F1 Score:** ~97.5%

These results indicate that our implementation closely matches scikit-learn's AdaBoost, affirming the correctness of our approach on the Breast Cancer Wisconsin dataset. This performance also aligns with findings by Street et al. (1993), underscoring AdaBoost's effectiveness in medical classification tasks involving nuclear features.

```
In [ ]: # Import necessary libraries
import pandas as pd
from sklearn.ensemble import AdaBoostClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, f1_score

# Define file paths (update these paths if needed)
# data_path = os.getcwd() + r'..\data\wdbc.data'

data_path = '/Users/liuzhenke/Desktop/Brown CSCI/DATA2060 ML/Final Project/wdbc.data'

column_names = ['ID', 'Diagnosis'] + [f'Feature_{i}' for i in range(1, 31)]
data = pd.read_csv(data_path, header=None, names=column_names)
data = data.drop(columns=['ID'])
data['Diagnosis'] = data['Diagnosis'].map({'M': 1, 'B': 0})
# Convert labels to -1 and 1
y = np.where(data['Diagnosis'].values == 0, -1, 1)
X = data.drop(columns=['Diagnosis']).values

# Split the dataset into training and testing sets (70% train, 30% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Initialize scikit-learn's AdaBoostClassifier
adaboost_model = AdaBoostClassifier(n_estimators=50, random_state=42, algorithm='SAMME')

# Train the scikit-learn model
adaboost_model.fit(X_train, y_train)

# Predict on test data
y_pred = adaboost_model.predict(X_test)

# Evaluate accuracy and F1 score for scikit-learn model
accuracy = accuracy_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred, pos_label=1)

# Output results for scikit-learn model
print(f"Scikit-learn AdaBoost Accuracy on test set: {accuracy:.2%}")
print(f"Scikit-learn AdaBoost F1 Score on test set: {f1:.2%}")

# Initialize and train custom myAdaBoost model
adaboost_model_self = myAdaBoost(n_clf=50)
adaboost_model_self.fit(X_train, y_train)

# Predict on test data using custom model
y_pred_self = adaboost_model_self.predict(X_test)

# Evaluate accuracy and F1 score for custom model
acc_self = accuracy_score(y_test, y_pred_self)
f1_self = f1_score(y_test, y_pred_self, pos_label=1)

# Output results for custom model
print(f"MyAdaBoost Accuracy on test set: {acc_self:.2%}")
print(f"MyAdaBoost F1 Score on test set: {f1_self:.2%}")
```

Scikit-learn AdaBoost Accuracy on test set: 97.08%
Scikit-learn AdaBoost F1 Score on test set: 96.06%
MyAdaBoost Accuracy on test set: 97.66%
MyAdaBoost F1 Score on test set: 96.88%