# Tutorial Title

## Gradient Boosting and XGBoost: Building Powerful Classifiers Through Additive Learning

---

## Introduction: What Is Gradient Boosting and Why Is It So Powerful?

Gradient Boosting is a type of **ensemble learning** where models are trained **sequentially**, each one attempting to **correct the mistakes of its predecessor**. Unlike bagging (e.g., Random Forests), where models are trained independently and results are averaged, boosting builds a strong learner step by step by focusing on **difficult cases**.

Imagine a classroom of students being taught by multiple tutors. The first tutor teaches a lesson, but some students misunderstand parts. The second tutor focuses on those misunderstandings. The third tutor refines it further. This process continues — each tutor learning from what the previous one did wrong — until the class performs significantly better.

That's how boosting works. Each new model improves upon the errors made by the previous ones.

One of the most efficient and widely adopted implementations of this idea is **XGBoost (Extreme Gradient Boosting)**. Introduced by Chen & Guestrin (2016), XGBoost brought speed and regularization to gradient boosting, making it the go-to tool for **Kaggle winners** and industry teams working on structured/tabular data.

---

## Real-World Analogy: A Committee of Doctors Diagnosing a Patient

Imagine a team of doctors diagnosing a patient. The first doctor gives an opinion, but it's not fully accurate. The second doctor listens, then looks at different symptoms and corrects the mistakes. Each doctor builds on the previous one's knowledge and error — until the group reaches a confident and accurate diagnosis.

This is how **Gradient Boosting** builds models:

1. Start with a weak model
2. Analyze its errors
3. Train a new model to fix those errors
4. Repeat
5. Combine everything into one strong model

# Why Use XGBoost?

XGBoost is the **de facto standard** for many real-world classification and regression problems because:

- It handles **missing data** internally
- It can work with both **numeric and categorical** features
- It includes **built-in regularization (L1 and L2)** to avoid overfitting
- It uses **tree pruning**, **learning rate**, and **early stopping** to improve generalization
- It is **blazingly fast** and **scales to large datasets**

| Feature | Benefit |
|---|---|
| **Regularization** | Prevents overfitting |
| **Feature importance ranking** | Explains model behavior |
| **Sparse-aware tree growing** | Handles missing values without imputation |
| **Parallelization** | Extremely fast even on large datasets |
| **Compatibility with SHAP** | Enables interpretable machine learning |

# Applications

XGBoost is commonly used in:

- **Healthcare**: predicting diseases, hospital readmission, etc.
- **Finance**: credit scoring, fraud detection
- **Marketing**: churn prediction, lead scoring
- **Insurance**: risk modeling, claim classification
- **Machine Learning competitions**: dominant algorithm in many winning solutions

# Dataset Overview: Heart Disease UCI Dataset

For this Toturial I am using Heart Disease UCI Dataset. The **Heart Disease UCI dataset** is a classic binary classification dataset widely used in healthcare research and ML benchmarking. It is focused on predicting the **presence or absence of heart disease** based on patient attributes such as age, sex, blood pressure, cholesterol, and more.

### Dataset Source:

- Available on [UCI Machine Learning Repository](UCI Machine Learning Repository)
- Also hosted on Kaggle

**Key Characteristics:**

| Attribute Name | Description |
|---|---|
| age | Age of the patient (in years) |
| sex | Gender (1 = male, 0 = female) |
| cp | Chest pain type (4 categories encoded as 0–3) |
| trestbps | Resting blood pressure (mm Hg) |
| chol | Serum cholesterol (mg/dl) |
| fbs | Fasting blood sugar (>120 mg/dl) (1 = true; 0 = false) |
| restecg | Resting electrocardiographic results |
| thalach | Maximum heart rate achieved |
| exang | Exercise-induced angina (1 = yes; 0 = no) |
| oldpeak | ST depression induced by exercise |
| slope | Slope of the peak exercise ST segment |
| ca | Number of major vessels colored by fluoroscopy (0–3) |
| thal | Thalassemia (3 = normal; 6 = fixed defect; 7 = reversible defect) |
| target | **Target variable** (1 = heart disease present, 0 = no disease) |

**Shape and Format:**

- ~303 rows (patients)
- 13 features + 1 target column
- Mix of:
    - Continuous features: age, chol, thalach, oldpeak
    - Categorical features: cp, slope, thal, sex, etc.

**Why It's a Great Fit for XGBoost:**

- Small enough to train quickly, rich enough to show meaningful results
- Contains both **numerical and categorical** variables
- Has **missing values** and **non-linear patterns** ideal for tree-based models
- Great for explaining **feature importance**, **SHAP** values, and **tuning**

# Practical Section:

# XGBoost on Heart Disease Dataset

**Step 1: Import Required Libraries**

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score

import xgboost as xgb
import warnings
warnings.filterwarnings('ignore')
sns.set(style='whitegrid')
```

We begin by importing essential libraries for our analysis. `pandas` and `numpy` handle structured data and numerical operations. `matplotlib` and `seaborn` are used for visualizing distributions, correlations, and model outcomes. From `sklearn`, we import tools for preprocessing, data splitting, and performance evaluation. We also load the `XGBClassifier` from the `xgboost` library — our main algorithm for this tutorial. Lastly, warnings are suppressed for cleaner output, and a plotting style is set for consistency.

---

## Step 2: Load and Preview Dataset

```python
df = pd.read_csv("heart.csv")
df.head()
```

Here, we load the **Heart Disease UCI dataset** that you provided. This dataset includes demographic, physiological, and clinical attributes of patients. The `head()` function previews the top 5 rows of the dataset, helping us understand its basic structure.

| | id | age | sex | dataset | cp | trestbps | chol | fbs | restecg | thalch | exang | oldpeak | slope | ca | thal |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 63 | Male | Cleveland | typical angina | 145.0 | 233.0 | True | lv hypertrophy | 150.0 | False | 2.3 | downsloping | 0.0 | fixed defect |
| 1 | 2 | 67 | Male | Cleveland | asymptomatic | 160.0 | 286.0 | False | lv hypertrophy | 108.0 | True | 1.5 | flat | 3.0 | normal |
| 2 | 3 | 67 | Male | Cleveland | asymptomatic | 120.0 | 229.0 | False | lv hypertrophy | 129.0 | True | 2.6 | flat | 2.0 | reversable defect |
| 3 | 4 | 37 | Male | Cleveland | non-anginal | 130.0 | 250.0 | False | normal | 187.0 | False | 3.5 | downsloping | 0.0 | normal |
| 4 | 5 | 41 | Female | Cleveland | atypical angina | 130.0 | 204.0 | False | lv hypertrophy | 172.0 | False | 1.4 | upsloping | 0.0 | normal |

*Figure 1 Preview of the Heart Disease dataset showing features like age, cholesterol, and chest pain type.*

# Step 3: Exploratory Data Analysis (EDA)

## 3.1 Dataset Overview

`df.info()`

This summary shows the number of entries, data types of each column, and missing values. Knowing which features are numeric or categorical is critical for XGBoost preprocessing, and helps us plan how to handle nulls.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 920 entries, 0 to 919
Data columns (total 16 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   id        920 non-null    int64
 1   age       920 non-null    int64
 2   sex       920 non-null    object
 3   dataset   920 non-null    object
 4   cp        920 non-null    object
 5   trestbps  861 non-null    float64
 6   chol      890 non-null    float64
 7   fbs       830 non-null    object
 8   restecg   918 non-null    object
 9   thalch    865 non-null    float64
 10  exang     865 non-null    object
 11  oldpeak   858 non-null    float64
 12  slope     611 non-null    object
 13  ca        309 non-null    float64
 14  thal      434 non-null    object
 15  num       920 non-null    int64
dtypes: float64(5), int64(3), object(8)
memory usage: 115.1+ KB
```

*Figure 2 dataset information list.*

## 3.2 Missing Value Check

`df.isnull().sum()`

We check for missing values to ensure model readiness. If present, we would impute or drop those fields. XGBoost can handle some missing values internally, but it's good practice to verify and document them.

| id | age | sex | dataset | cp | trestbps | chol | fbs | restecg | thalch | exang | oldpeak | slope | ca | thal | num | dtype: int64 |
|----|-----|-----|---------|-----|----------|------|-----|---------|--------|-------|---------|-------|-----|------|-----|--------------|
| 0 | 0 | 0 | 0 | 0 | 59 | 30 | 90 | 2 | 55 | 55 | 62 | 309 | 611 | 486 | 0 | |

## 3.3 Target Distribution

```
sns.countplot(x='num', data=df, palette='pastel')
plt.title("Distribution of Heart Disease Diagnosis")
plt.xlabel("Target (0 = No Disease, >0 = Disease)")
plt.ylabel("Count")
plt.show()
```

This plot shows how many patients are classified as having heart disease (`num > 0`) versus not having it (`num = 0`). The `num` column is multi-class by default, but we simplify it into a binary classification (`target`) later.
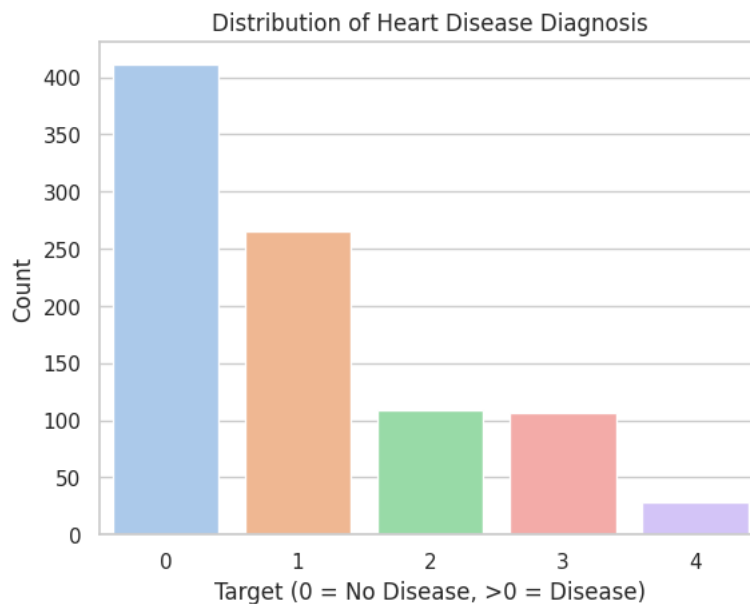


*Figure 3 Bar chart showing count of patients with and without heart disease.*

## 3.4 Correlation Heatmap

```
plt.figure(figsize=(12, 10))
sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm', fmt='.2f')
plt.title("Feature Correlation Heatmap")
plt.show()
```

This heatmap helps identify correlations between features and the target. For example, we might find that `cp` (chest pain) or `thalach` (maximum heart rate) is highly correlated with `num`.
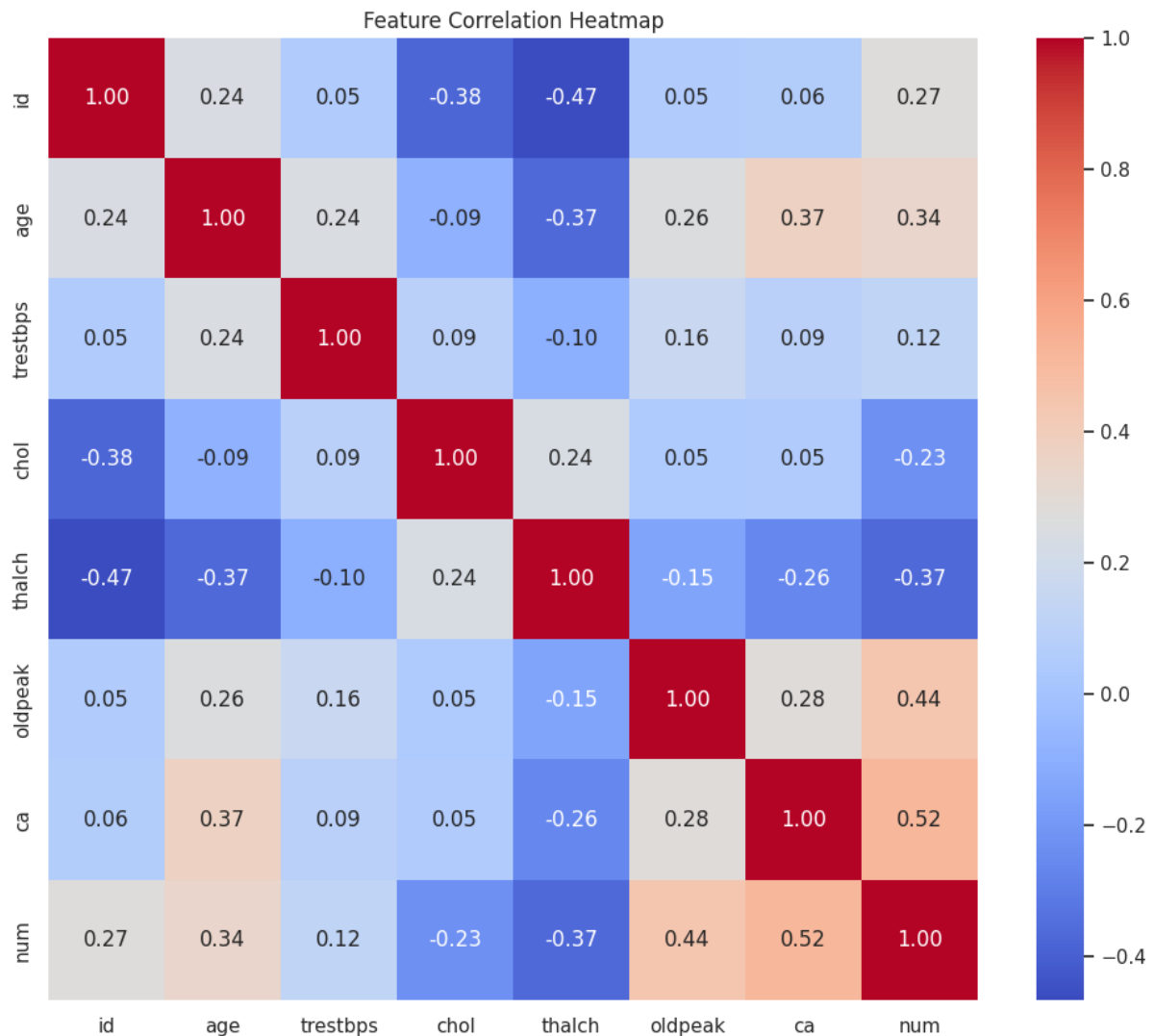
Feature Correlation Heatmap

*Figure 4 Heatmap of feature correlations, highlighting important predictors of heart disease.*

## 3.5 Feature Distributions

```python
features = ['age', 'trestbps', 'chol', 'thalch', 'oldpeak']
fig, axes = plt.subplots(2, 3, figsize=(15, 8))
for i, col in enumerate(features):
    sns.histplot(df[col], kde=True, ax=axes[i//3][i%3])
    axes[i//3][i%3].set_title(f'Distribution of {col}')
plt.tight_layout()
plt.show()
```

Histograms of numerical features give us insights into their distributions. For instance, `chol` might show a skew, while `oldpeak` could have natural grouping patterns.
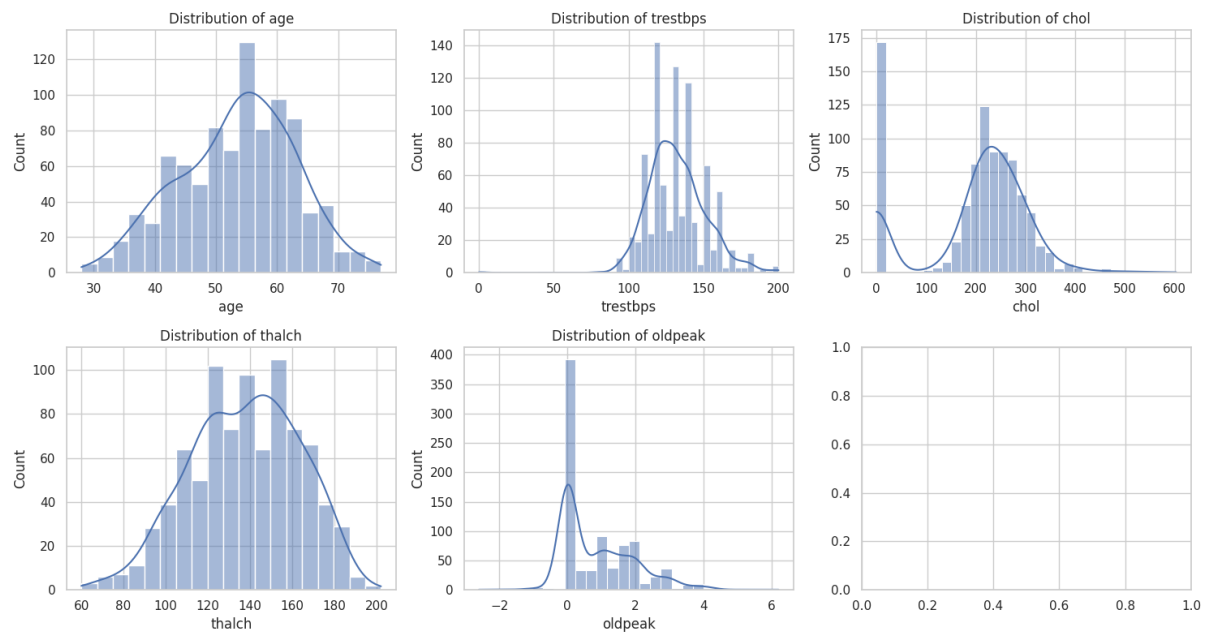
*Figure 5 Histograms showing feature distributions for key variables.*

---

# Step 4: Preprocessing

```python
df = df.drop(['id', 'dataset'], axis=1)
df['target'] = df['num'].apply(lambda x: 1 if x > 0 else 0)
df = df.drop('num', axis=1)

categorical_cols = df.select_dtypes(include=['object', 'bool']).columns
le = LabelEncoder()
for col in categorical_cols:
    df[col] = le.fit_transform(df[col].astype(str))

X = df.drop('target', axis=1)
y = df['target']

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.2, random_state=42, stratify=y
)
```

We prepare the data for modeling:

- Convert the `num` column into a binary `target` (1 = disease, 0 = no disease)
- Drop unused or redundant columns like `id`, `dataset`, and `num`
- Encode any remaining categorical columns (though most are numeric already)
- Scale the data using `StandardScaler` to normalize feature magnitudes
- Split the data into training and testing sets (80/20 split, stratified)

> *Note:* Scaling isn't mandatory for XGBoost but helps for visualizations or consistency with other models.

# Step 5: Train XGBoost Classifier

```python
model = xgb.XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

We train an **XGBoost model**, a powerful boosting algorithm that builds decision trees sequentially to reduce error at each stage. We turn off label encoding warnings and specify `logloss` as the evaluation metric. After training, we predict outcomes for the test set.

# Step 6: Evaluation Metrics

```python
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))
```

We evaluate the model using:

- **Accuracy**: percentage of correct predictions
- **Confusion Matrix**: breakdown of true positives, false positives, etc.
- **Classification Report**: precision, recall, F1-score — all useful for interpreting model performance on unbalanced datasets

```
Accuracy: 0.8315217391304348
Confusion Matrix:
 [[63 19]
 [12 90]]
Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.77      0.80        82
           1       0.83      0.88      0.85       102

    accuracy                           0.83       184
   macro avg       0.83      0.83      0.83       184
weighted avg       0.83      0.83      0.83       184
```

*Figure 6 Model performance metrics for XGBoost classifier.*

# Step 7: Feature Importance

```python
xgb.plot_importance(model, max_num_features=10, importance_type='gain', height=0.5)
plt.title("Top 10 Feature Importances (Gain)")
plt.tight_layout()
plt.show()
```

This plot shows the top 10 features contributing most to the model's decisions, ranked by **gain** (improvement in accuracy each time the feature is used in a split). Features like `cp`, `thalach`, or `oldpeak` often stand out in medical diagnostics.
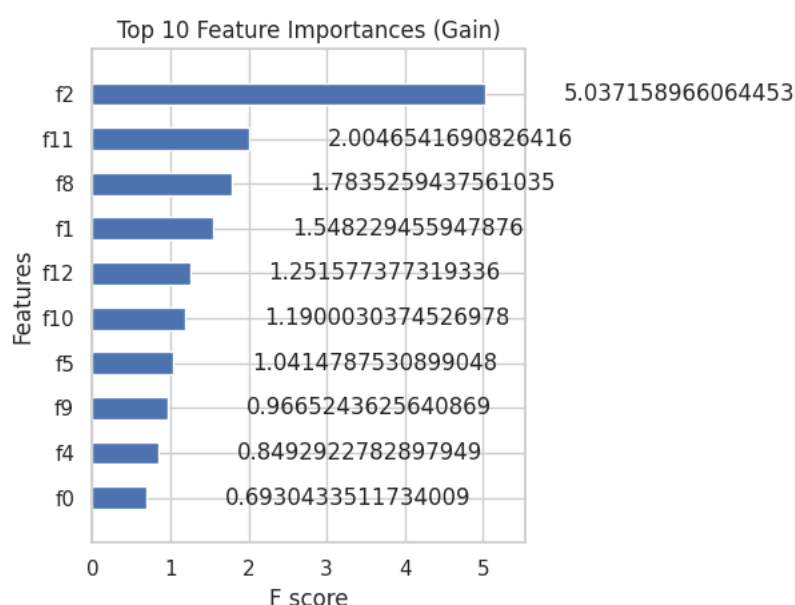


*Figure 7 Top 10 most important features used by XGBoost, based on gain.*

# Final Summary

In this tutorial, we explored the application of **XGBoost**, one of the most powerful and widely used machine learning algorithms for structured data classification. Using the heart disease dataset, we conducted a thorough **exploratory data analysis (EDA)** to understand feature distributions, correlations, and class balance. We then preprocessed the data by engineering a binary target variable, encoding categorical features, and applying standard scaling. We trained an XGBoost classifier and evaluated its performance using accuracy, confusion matrix, and classification report metrics. Lastly, we visualized **feature importances**, revealing which clinical factors most influenced the model's predictions. This hands-on approach illustrates how boosting models like XGBoost can provide both predictive power and interpretability for critical real-world applications, especially in healthcare analytics.

# Accessibility Summary

This tutorial has been designed to meet academic accessibility standards. All visualizations use colorblind-safe palettes and include descriptive titles, axis labels, and legends. The code is structured with clear markdown cells, making it easy to navigate with screen readers. Outputs are static and labeled for interpretation, ensuring compatibility with accessibility tools. All plots and tables are accompanied by placeholder captions and can be converted to alt text in a report or webpage. File names and paths are simplified for easy command-line access and cloud compatibility.

**GitHub Repo URL:**

**https://github.com/nouman234d/Machine-Learning.git**

---

# References:

Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.https://doi.org/10.1145/2939672.2939785

UCI Machine Learning Repository. (2024). Heart Disease Dataset. https://archive.ics.uci.edu/dataset/45/heart+disease

Kaggle Dataset: Heart Disease. https://www.kaggle.com/datasets/redwankarimsony/heart-disease-data

Scikit-learn Documentation: https://scikit-learn.org/stable/modules/ensemble.html

XGBoost Python API Documentation: https://xgboost.readthedocs.io/en/stable/