# AdaBoost

# DISCLAIMER

In a **Random Forest**, each time you make a tree, you make a full sized tree.

In contrast, in a **Forest of Trees** made with **AdaBoost**, the trees are usually just a **node** and two **leaves**.

Some trees might be bigger than others, but there is no predetermined maximum depth.

A tree with just one node and two leaves is called a *stump*.

...so this is really a **Forest of Stumps** rather than trees.

| Chest Pain | Good Blood Circ. | Blocked Arteries | Weight | Heart Disease |
|------------|------------------|------------------|--------|---------------|
| No | No | No | 125 | No |
| Yes | Yes | Yes | 180 | Yes |
| Yes | Yes | No | 210 | No |
| Yes | No | Yes | 167 | Yes |

…then a full sized **Decision Tree** would take advantage of all **4** variables that we measured (**Chest Pain**, **Blood Circulation**, **Blocked Arteries** and **Weight**) to make a decision…

| Chest Pain | Good Blood Circ. | Blocked Arteries | Weight | Heart Disease |
|---|---|---|---|---|
| No | No | No | 125 | No |
| Yes | Yes | Yes | 180 | Yes |
| Yes | Yes | No | 210 | No |
| Yes | No | Yes | 167 | Yes |

…but a **Stump** can only use one variable to make a decision.



Thus, **Stumps** are technically "weak learners".

However, that's the way **AdaBoost** likes it, and it's one of the reasons why they are so commonly combined.

In a **Random Forest**, each tree has an equal vote on the final classification.

In contrast, in a **Forest of Stumps** made with **AdaBoost**, some stumps get more say in the final classification than others.

Lastly, in a **Random Forest**, each decision tree is made independently of the others.

In contrast, in a **Forest of Stumps** made with **AdaBoost**, order is important.

The errors that the first stump makes...

...influence how the second stump is made...

To review, the three ideas behind **AdaBoost** are…

1) **AdaBoost** combines a lot of "weak learners" to make classifications. The weak learners are almost aways **stumps**.

2) Some **stumps** get more say in the classification than others.

3) Each **stump** is made by taking the previous **stump's** mistakes into account.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease |
|---|---|---|---|
| Yes | Yes | 205 | Yes |
| No | Yes | 180 | Yes |
| Yes | No | 210 | Yes |
| Yes | Yes | 167 | Yes |
| No | Yes | 156 | No |
| No | Yes | 125 | No |
| Yes | No | 168 | No |
| Yes | Yes | 172 | No |

We create a **Forest of Stumps** with **AdaBoost** to predict if a patient has heart disease.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease |
|------------|------------------|----------------|---------------|
| Yes | Yes | 205 | Yes |
| No | Yes | 180 | Yes |
| Yes | No | 210 | Yes |
| Yes | Yes | 167 | Yes |
| No | Yes | 156 | No |
| No | Yes | 125 | No |
| Yes | No | 168 | No |
| Yes | Yes | 172 | No |

We will make these predictions based on a patient's **Chest Pain** and **Blocked Artery** status and their **Weight**.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|:---:|:---:|:---:|:---:|:---:|
| Yes | Yes | 205 | Yes | |
| No | Yes | 180 | Yes | |
| Yes | No | 210 | Yes | |
| Yes | Yes | 167 | Yes | |
| No | Yes | 156 | No | |
| No | Yes | 125 | No | |
| Yes | No | 168 | No | |
| Yes | Yes | 172 | No | |

The first thing we do is give each sample a weight that indicates how important it is to be correctly classified.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

At the start, all samples get the same weight…

$$\frac{1}{\text{total number of samples}} = \frac{1}{8}$$

However, after we make the first stump, these weights will change in order to guide how the next stump is created.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|------------|------------------|----------------|---------------|---------------|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

This is done finding the variable, **Chest Pain**, **Blocked Arteries** or **Patient Weight**, that does the best job classifying the samples.

Chest Pain

| Yes Heart Disease | | No Heart Disease | |
|---|---|---|---|
| Correct | Incorrect | Correct | Incorrect |
| 3 | 2 | 2 | 1 |

Gini Index 0.47

Blocked Arteries

| Yes Heart Disease | | No Heart Disease | |
|---|---|---|---|
| Correct | Incorrect | Correct | Incorrect |
| 3 | 3 | 1 | 1 |

Gini Index 0.5

The **Gini Index** for **Patient Weight** is the lowest...

Weight > 176

| Yes Heart Disease | | No Heart Disease | |
|---|---|---|---|
| Correct | Incorrect | Correct | Incorrect |
| 3 | 0 | 4 | 1 |

Gini Index 0.2

...so this will be the first stump in the forest.

Now we need to determine how much say this stump will have in the final classification.

**Weight > 176**

**Yes Heart Disease**

Correct | Incorrect
--- | ---
3 | 0

**No Heart Disease**

Correct | Incorrect
--- | ---
4 | 1

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

This patient, who weighs less than 176, *has* heart disease, but the stump says they do not.

**Weight > 176**

| Yes Heart Disease | |
|---|---|
| Correct | Incorrect |
| 3 | 0 |

| No Heart Disease | |
|---|---|
| Correct | Incorrect |
| 4 | 1 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

We use the **Total Error** to determine **Amount of Say** this stump has in the final classification with the following formula:

$$\text{Amount of Say} = \frac{1}{2} \log\left(\frac{1 - \text{Total Error}}{\text{Total Error}}\right)$$

**Weight > 176**

**Yes Heart Disease**

| Correct | Incorrect |
|---|---|
| 3 | 0 |

**No Heart Disease**

| Correct | Incorrect |
|---|---|
| 4 | 1 |

The **Blue Line** tells us the **Amount of Say** for **Total Error** values between **0** and **1**.

$$\text{Amount of Say} = \frac{1}{2} \log\left(\frac{1 - \text{Total Error}}{\text{Total Error}}\right)$$

**Weight > 176**

**Yes Heart Disease**

| Correct | Incorrect |
| --- | --- |
| 3 | 0 |

**No Heart Disease**

| Correct | Incorrect |
| --- | --- |
| 4 | 1 |

…and the **Amount of Say** that this stump has on the final classification is **0.97**.

$$\text{Amount of Say} = \frac{1}{2} \log( 7 ) = 0.97$$

**Weight > 176**

**Yes Heart Disease**

| Correct | Incorrect |
|---------|-----------|
| 3 | 0 |

**No Heart Disease**

| Correct | Incorrect |
|---------|-----------|
| 4 | 1 |

Total Error

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

…we will emphasize the need for the next stump to correctly classify it by increasing its **Sample Weight**…

**Weight > 176**

**Yes Heart Disease**

| Correct | Incorrect |
|---|---|
| 3 | 0 |

**No Heart Disease**

| Correct | Incorrect |
|---|---|
| 4 | 1 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|:---:|:---:|:---:|:---:|:---:|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

…and decreasing all of the other **Sample Weights**.

**Weight > 176**

**Yes Heart Disease**

| Correct | Incorrect |
|:---:|:---:|
| 3 | 0 |

**No Heart Disease**

| Correct | Incorrect |
|:---:|:---:|
| 4 | 1 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|------------|------------------|----------------|---------------|---------------|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

$$\text{New Sample Weight} = \text{sample weight} \times e^{\text{amount of say}}$$

This is the formula we will use to *increase* the **Sample Weight** for the sample that was *incorrectly* classified.

New Sample Weight $= \text{sample weight} \times e^{\text{amount of say}}$

$$= \frac{1}{8} e^{\text{amount of say}}$$

$e^{\text{amount of say}}$

…then the previous **Sample Weight** is scaled by a relatively small number.

Amount of Say

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

New Sample Weight $=$ sample weight $\times\ e^{\text{-amount of say}}$

The big difference is the *negative* sign in front of **Amount of Say**

$$\text{New Sample Weight} = \text{sample weight} \times e^{-\text{amount of say}}$$

$$= \frac{1}{8} \boxed{e^{-\text{amount of say}}}$$

$e^{-\text{amount of say}}$

...then we will scale the **Sample Weight** by a value close to **1**.

This means that the **New Sample Weight** will be just a little smaller than the old one.

Amount of Say

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight | New Weight |
|---|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 | 0.05 |
| No | Yes | 180 | Yes | 1/8 | 0.05 |
| Yes | No | 210 | Yes | 1/8 | 0.05 |
| Yes | Yes | 167 | Yes | 1/8 | 0.33 |
| No | Yes | 156 | No | 1/8 | 0.05 |
| No | Yes | 125 | No | 1/8 | 0.05 |
| Yes | No | 168 | No | 1/8 | 0.05 |
| Yes | Yes | 172 | No | 1/8 | 0.05 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight | New Weight |
|---|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 | 0.05 |
| No | Yes | 180 | Yes | 1/8 | 0.05 |
| Yes | No | 210 | Yes | 1/8 | 0.05 |
| Yes | Yes | 167 | Yes | 1/8 | 0.33 |
| No | Yes | 156 | No | 1/8 | 0.05 |
| No | Yes | 125 | No | 1/8 | 0.05 |
| Yes | No | 168 | No | 1/8 | 0.05 |
| Yes | Yes | 172 | No | 1/8 | 0.05 |

Now we need to normalize the **New Sample Weights** so that they will add up to 1.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight | New Weight | Norm. Weight |
|---|---|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 | 0.05 | 0.07 |
| No | Yes | 180 | Yes | 1/8 | 0.05 | 0.07 |
| Yes | No | 210 | Yes | 1/8 | 0.05 | 0.07 |
| Yes | Yes | 167 | Yes | 1/8 | 0.33 | 0.49 |
| No | Yes | 156 | No | 1/8 | 0.05 | 0.07 |
| No | Yes | 125 | No | 1/8 | 0.05 | 0.07 |
| Yes | No | 168 | No | 1/8 | 0.05 | 0.07 |
| Yes | Yes | 172 | No | 1/8 | 0.05 | 0.07 |

So we divide each **New Sample Weight** by **0.68** to get the normalized values.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | No | 0.07 |
| Yes | No | 168 | No | 0.07 |
| Yes | Yes | 172 | No | 0.07 |

Now we just transfer the **Normalized Sample Weights** to the **Sample Weights** column, since those are what we will use for the next stump.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | No | 0.07 |
| Yes | No | 168 | No | 0.07 |
| Yes | Yes | 172 | No | 0.07 |

In theory, we could use the **Sample Weights** to calculate **Weighted Gini Indexes** to determine which variable should split the next stump.

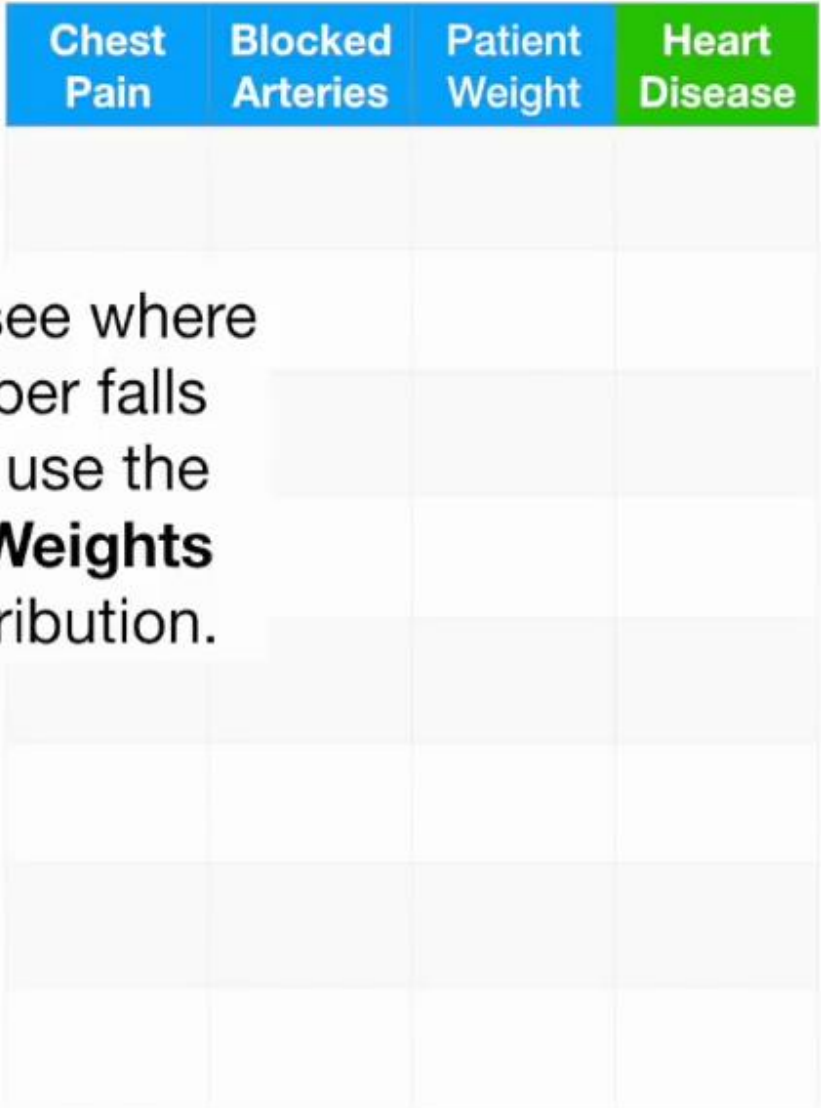| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | No | 0.07 |
| Yes | No | 168 | No | 0.07 |
| Yes | Yes | 172 | No | 0.07 |

Alternatively, instead of using a **Weighted Gini Index**, we can make a new collection of samples that contains duplicate copies of the samples with the largest **Sample Weights**.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | No | 0.07 |
| Yes | No | 168 | No | 0.07 |
| Yes | Yes | 172 | No | 0.07 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease |
|---|---|---|---|

…and we see where that number falls when we use the **Sample Weights** like a distribution.

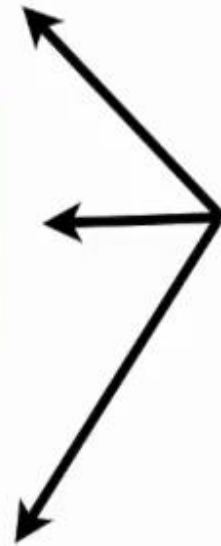| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | | |
| Yes | No | 168 | | |
| Yes | Yes | 172 | | |

Ultimately, this sample was added to the new collection of samples **4** times, reflecting its larger **Sample Weight**.
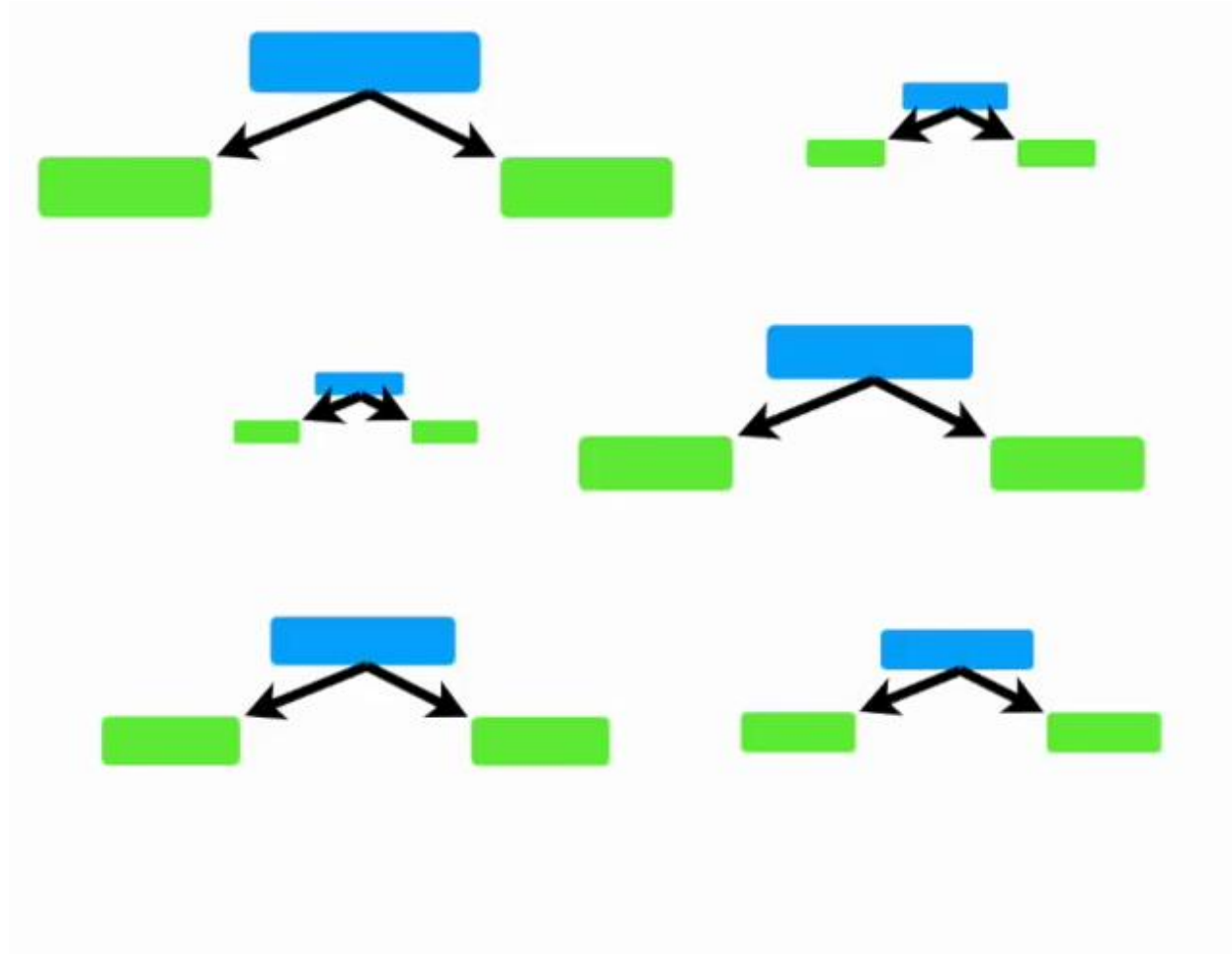
| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease |
|---|---|---|---|
| No | Yes | 156 | No |
| Yes | Yes | 167 | Yes |
| No | Yes | 125 | No |
| Yes | Yes | 167 | Yes |
| Yes | Yes | 167 | Yes |
| Yes | Yes | 172 | No |
| Yes | Yes | 205 | Yes |
| Yes | Yes | 167 | Yes |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| No | Yes | 156 | No | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| Yes | Yes | 172 | No | 1/8 |
| Yes | Yes | 205 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |

Lastly, we give all the samples equal **Sample Weights**, just like before.

Because these samples are all the same, they will be treated as a block, creating a large penalty for being misclassified.
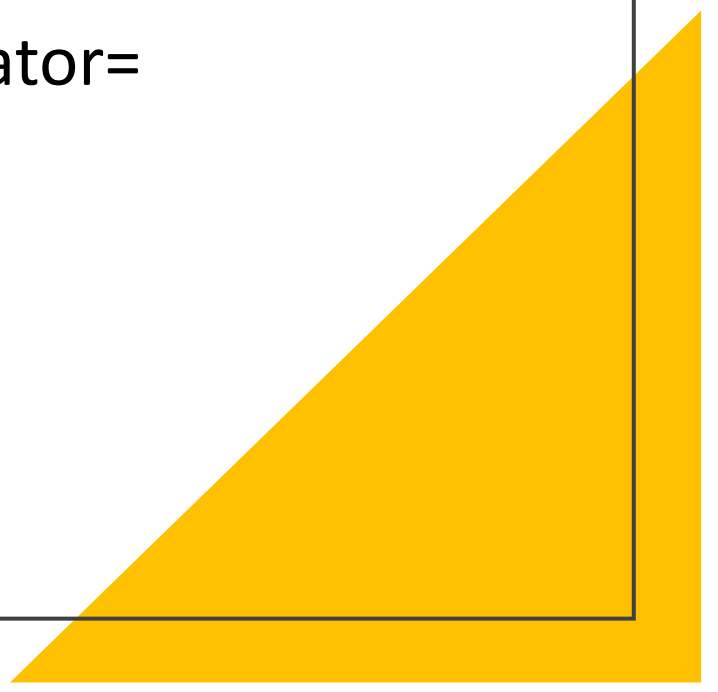
This goes on and on....
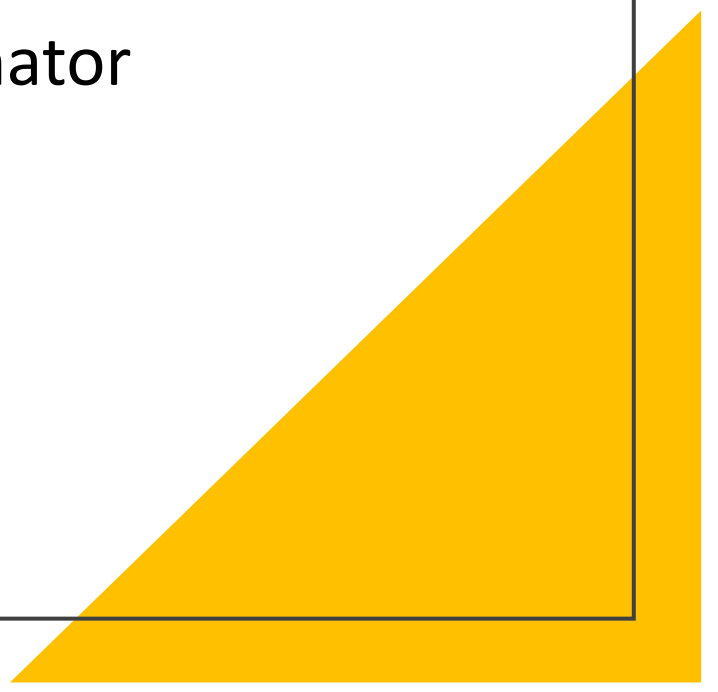Until you train sufficient number of stumps

# AdaBoostClassifier in Sklearn

class
sklearn.ensemble.AdaBoostClassifier(base_estimator=
None, *, n_estimators=50, learning_rate=1.0,
algorithm='SAMME.R', random_state=None)

# AdaBoostRegressor in Sklearn

class
sklearn.ensemble.AdaBoostRegressor(base_estimator
=None, *, n_estimators=50, learning_rate=1.0,
loss='linear', random_state=None)[source]

# Why Adaboost works?

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | No | 0.07 |
| Yes | No | 168 | No | 0.07 |
| Yes | Yes | 172 | No | 0.07 |