



**BITS Pilani**  
Pilani Campus

# Classification: Ensemble Methods

Dr. Chetana Gavankar, Ph.D,  
IIT Bombay-Monash University Australia  
[Chetana.gavankar@pilani.bits-pilani.ac.in](mailto:Chetana.gavankar@pilani.bits-pilani.ac.in)



**BITS Pilani**  
Pilani Campus

Session 9  
Date –06/03/2022  
Time – 10 to 12.30

## **Text Book(s)**

R1	An Introduction to Data Mining – Pang-Ning Tan, Michael Steinbach, Anuj Karpatne, Vipin Kumar - 2005
----	--

These slides are prepared by the instructor, with grateful acknowledgement of and many others who made their course materials freely available online.

# Topics to be covered



## Module 9 : Ensemble Methods

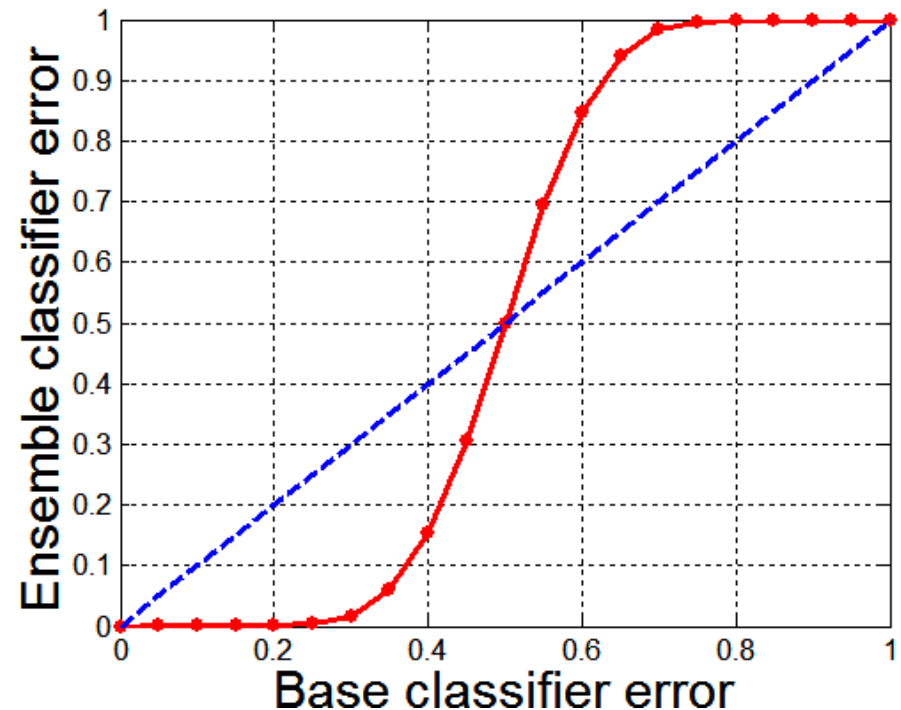
- Rational for Ensemble Method
- Methods for constructing an Ensemble Classifier
- Bagging, Boosting
- Random Forest
- AdaBoost
- eXtreme Gradient Boosting (XGBoost)
- Class Imbalance Problem & approaches to solve it
- Python Implementation of Random Forest, Adaboost and XGBoost

# Ensemble Methods

- **Ensemble methods** use multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algorithms alone
- Construct a set of classifiers from the training data
- Predict class label of test records by combining the predictions made by multiple classifiers
- Tend to reduce problems related to over-fitting of the training data.
- By combining individual models, the ensemble model tends to be more flexible (less bias) and less data-sensitive (less variance).

# Why Ensemble Methods work?

- 25 base classifiers
- Each classifier has error rate,  $\varepsilon = 0.35$
- If base classifiers are identical, then the ensemble will misclassify the same examples predicted incorrectly by the base classifiers depicted by dotted line
- Assume errors made by classifiers are uncorrelated
- ensemble makes a wrong prediction only if more than half of the base classifiers predict incorrectly
- Probability that the ensemble classifier makes a wrong prediction:



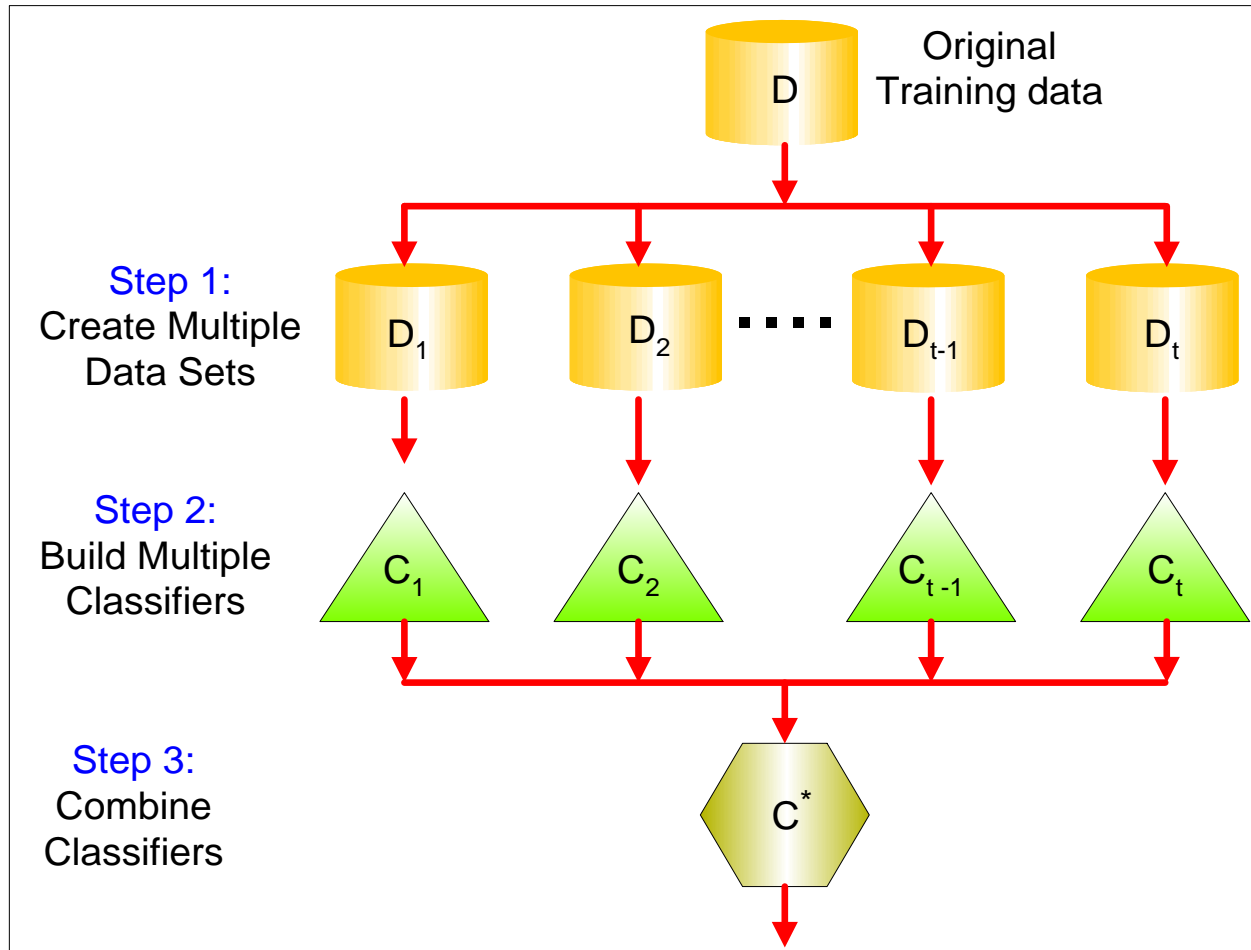
$$P(X \geq 13) = \sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06$$

# When does Ensemble work?



- Ensemble classifier performs better than the base classifiers when each classifier error is smaller than 0.5
- Necessary conditions for an ensemble classifier to perform better than a single classifier:
  - Base classifiers should be independent of each other
  - Base classifiers should do better than a classifier that performs random guessing

# General Approach



# Simple Ensemble Techniques



## Max Voting

Ex: Movie rating

The result of max voting would be something like this:

- Rating by 5 friends: 5 4 5 4 4

**Averaging-**  $(5+4+5+4+4)/5 = 4.4$  Final rating

## Weighted Average

Weight-0.23 0.23 0.18 0.18 0.18

The result is calculated as  $[(5*0.23) + (4*0.23) + (5*0.18) + (4*0.18) + (4*0.18)] = \mathbf{4.41}$ .



# Methods for constructing Ensemble Classifier



---

- Using different algorithms
- Using different parameters/hyperparameters
- Using different training sets
- By manipulating input features
- By manipulating the class labels

# Types of Ensemble Methods

---

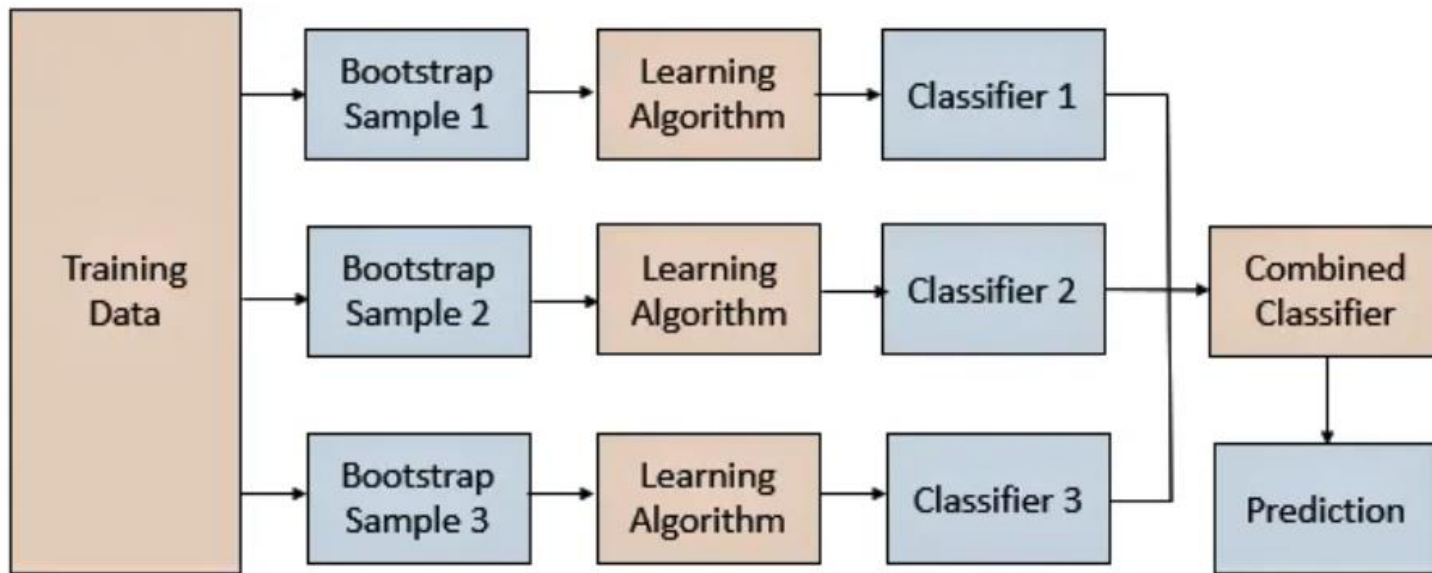
- **Simple Ensemble methods:** Max Voting, Averaging, Weighted Averaging
- **Advanced Ensemble Methods**
  - **Bagging**, often considers homogeneous weak learners, learns them independently from each other in parallel and combines them by some kind of averaging process
  - **Boosting**, often considers homogeneous weak learners, learns them sequentially in a very adaptative way (a base model depends on the previous ones) and combines them by some deterministic strategy
  - **stacking**, often considers heterogeneous weak learners, learns them in parallel and combines them by training a meta-model to output a prediction based on the different weak models predictions



# Bagging (Bootstrap Aggregating)

- Technique uses these subsets (bags) to get a fair idea of the distribution (complete set).
- The size of subsets created for bagging may be less than the original set.
- Bootstrapping is a sampling technique in which we create subsets of observations from the original dataset, **with replacement**.
- When you sample with replacement, items are independent. One item does not affect the outcome of the other. You have  $1/7$  chance of choosing the first item and a  $1/7$  chance of choosing the second item.
- If the two items are **dependent**, or linked to each other. When you choose the first item, you have a  $1/7$  probability of picking a item. Assuming you don't replace the item, you only have six items to pick from. That gives you a  $1/6$  chance of choosing a second item.

# Bagging



- Multiple subsets are created from the original dataset, selecting observations with replacement.
- A base model (weak model) is created on each of these subsets.
- The models run in parallel and are independent of each other.
- The final predictions are determined by combining the predictions from all the models.

# Bagging reduces variance (Intuition)



If each single classifier is unstable – that is, it has high variance, the aggregated classifier  $f^-$  has a smaller variance than a single original classifier

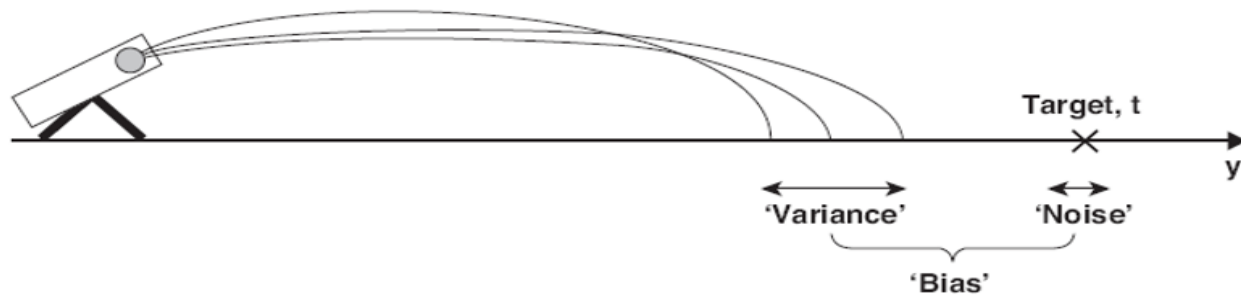


Figure 5.32. Bias-variance decomposition.

In general,

- **Bias** is contributed to by the training error; a complex model has low bias.
- **Variance** is caused by future error; a complex model has High variance.
- Bagging reduces the variance in the base classifiers.

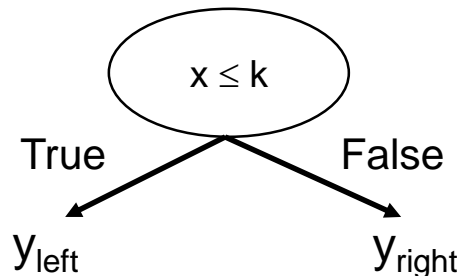
# Bagging Example

- Consider 1-dimensional data set:

Original Data:

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
y	1	1	1	-1	-1	-1	-1	1	1	1

- Classifier is a decision stump
  - Decision tree with one internal node (the root) which is immediately connected to the terminal nodes (its leaves). Decision stump makes a prediction based on the value of just a single input feature. Sometimes they are also called **1-rules**
    - Decision rule:  $x \leq k$  versus  $x > k$
    - Split point  $k$  is chosen based on entropy



# Bagging Example



Bagging Round 1:

x	0.1	0.2	0.2	0.3	0.4	0.4	0.5	0.6	0.9	0.9
y	1	1	1	1	-1	-1	-1	-1	1	1

$x \leq 0.35 \rightarrow y = 1$

$x > 0.35 \rightarrow y = -1$

Bagging Round 2:

x	0.1	0.2	0.3	0.4	0.5	0.5	0.9	1	1	1
y	1	1	1	-1	-1	-1	1	1	1	1

$x \leq 0.7 \rightarrow y = 1$

$x > 0.7 \rightarrow y = 1$

Bagging Round 3:

x	0.1	0.2	0.3	0.4	0.4	0.5	0.7	0.7	0.8	0.9
y	1	1	1	-1	-1	-1	-1	-1	1	1

$x \leq 0.35 \rightarrow y = 1$

$x > 0.35 \rightarrow y = -1$

Bagging Round 4:

x	0.1	0.1	0.2	0.4	0.4	0.5	0.5	0.7	0.8	0.9
y	1	1	1	-1	-1	-1	-1	-1	1	1

$x \leq 0.3 \rightarrow y = 1$

$x > 0.3 \rightarrow y = -1$

Bagging Round 5:

x	0.1	0.1	0.2	0.5	0.6	0.6	0.6	1	1	1
y	1	1	1	-1	-1	-1	-1	1	1	1

$x \leq 0.35 \rightarrow y = 1$

$x > 0.35 \rightarrow y = -1$

# Bagging Example



Bagging Round 6:

x	0.2	0.4	0.5	0.6	0.7	0.7	0.7	0.8	0.9	1
y	1	-1	-1	-1	-1	-1	-1	1	1	1

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

Bagging Round 7:

x	0.1	0.4	0.4	0.6	0.7	0.8	0.9	0.9	0.9	1
y	1	-1	-1	-1	-1	1	1	1	1	1

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

Bagging Round 8:

x	0.1	0.2	0.5	0.5	0.5	0.7	0.7	0.8	0.9	1
y	1	1	-1	-1	-1	-1	-1	1	1	1

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

Bagging Round 9:

x	0.1	0.3	0.4	0.4	0.6	0.7	0.7	0.8	1	1
y	1	1	-1	-1	-1	-1	-1	1	1	1

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

Bagging Round 10:

x	0.1	0.1	0.1	0.1	0.3	0.3	0.8	0.8	0.9	0.9
y	1	1	1	1	1	1	1	1	1	1

$x \leq 0.05 \rightarrow y = 1$   
 $x > 0.05 \rightarrow y = 1$



# Bagging Example

- Assume test set is the same as the original data
- Use majority vote to determine class of ensemble classifier

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	1	1	1	-1	-1	-1	-1	-1	-1	-1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	-1	-1	-1	-1	-1	-1	-1
4	1	1	1	-1	-1	-1	-1	-1	-1	-1
5	1	1	1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	1	1	1
7	-1	-1	-1	-1	-1	-1	-1	1	1	1
8	-1	-1	-1	-1	-1	-1	-1	1	1	1
9	-1	-1	-1	-1	-1	-1	-1	1	1	1
10	1	1	1	1	1	1	1	1	1	1
Sum	2	2	2	-6	-6	-6	-6	2	2	2
Sign	1	1	1	-1	-1	-1	-1	1	1	1

Predicted  
Class

# Bagging Algorithm



---

## Algorithm 5.6 Bagging Algorithm

---

- 1: Let  $k$  be the number of bootstrap samples.
  - 2: for  $i = 1$  to  $k$  do
  - 3:   Create a bootstrap sample of size  $n$ ,  $D_i$ .
  - 4:   Train a base classifier  $C_i$  on the bootstrap sample  $D_i$ .
  - 5: end for
  - 6:  $C^*(x) = \arg \max_y \sum_i \delta(C_i(x) = y)$ ,  $\{\delta(\cdot) = 1$  if its argument is true, and 0 otherwise. $\}$
-

# Boosting



- What if a data point is incorrectly predicted by the first model, and then the next (probably all models), will combining the predictions provide better results? Such situations are taken care of by boosting.
- Boosting is a sequential process, where each subsequent model attempts to correct the errors of the previous model.
- The succeeding models are dependent on the previous model.

# Boosting

- Records that are wrongly classified will have their weights increased
- Records that are classified correctly will have their weights decreased

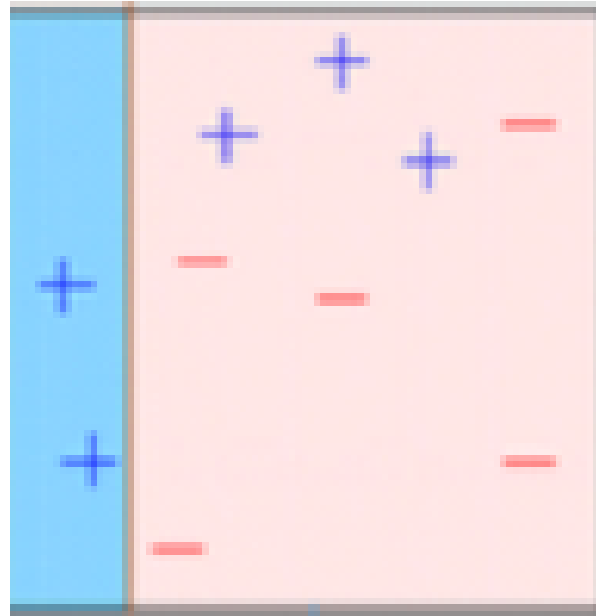
Original Data	1	2	3	4	5	6	7	8	9	10
Boosting (Round 1)	7	3	2	8	7	9	4	10	6	3
Boosting (Round 2)	5	4	9	4	2	5	1	7	4	2
Boosting (Round 3)	4	4	8	10	4	5	4	6	3	4

- Example 4 is hard to classify
- Its weight is increased, therefore it is more likely to be chosen again in subsequent rounds

# Boosting



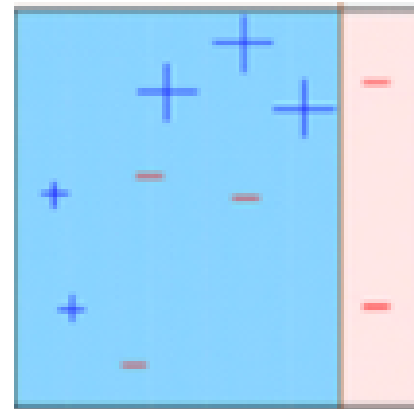
- A subset is created from the original dataset.
- Initially, all data points are given equal weights.
- A base model is created on this subset.
- This model is used to make predictions on the whole dataset.



# Boosting



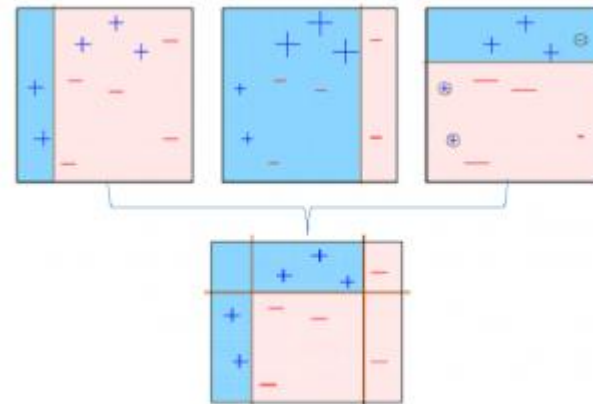
- Errors are calculated using the actual values and predicted values.
- The observations which are incorrectly predicted, are given higher weights. (Here, the three misclassified blue-plus points will be given higher weights)
- Another model is created and predictions are made on the dataset. (This model tries to correct the errors from the previous model)



# Boosting



- Similarly, multiple models are created, each correcting the errors of the previous model.
- The final model (strong learner) is the weighted mean of all the models (weak learners).



- Individual models would not perform well on the entire dataset, but they work well for some part of the dataset. Thus, each model actually boosts the performance of the ensemble.

# Algorithms based on Bagging and Boosting

---

## Bagging algorithms:

- Random forest

## Boosting algorithms:

- AdaBoost
- XGBoost



# Random Forest



- Random Forest is ensemble machine learning algorithm that follows the bagging technique.
- The base estimators in random forest are decision trees.
- Random forest randomly selects a set of features which are used to decide the best split at each node of the decision tree.

# Random Forest

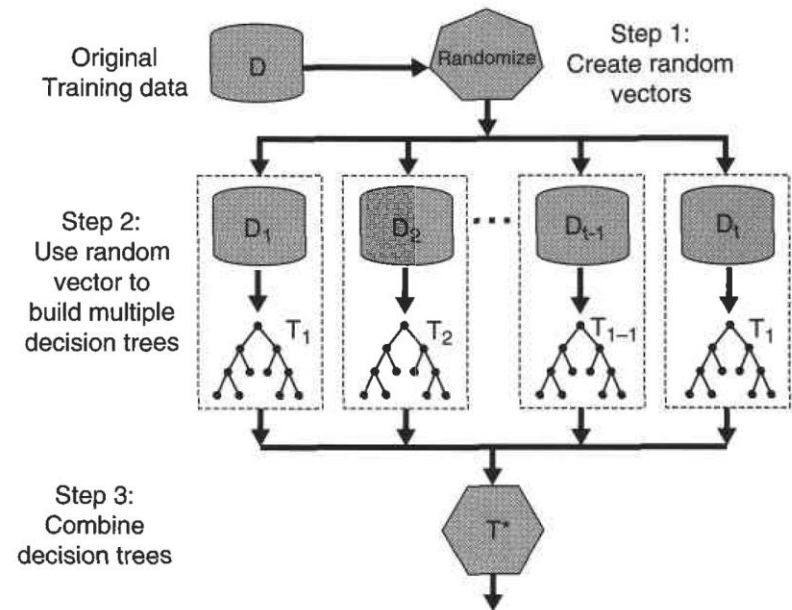
---

- Random subsets are created from the original dataset (bootstrapping).
- At each node in the decision tree, only a random set of features are considered to decide the best split.
- A decision tree model is fitted on each of the subsets.
- The final prediction is calculated by averaging the predictions from all decision trees.

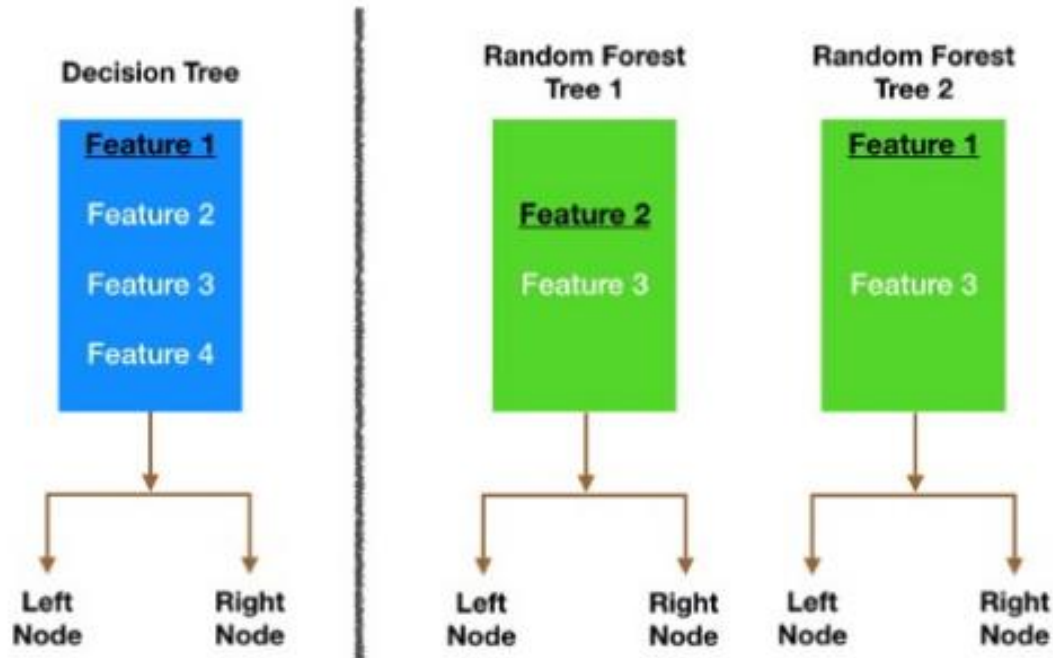
# Random Forest



- Combines the predictions made by multiple decision trees
- Each tree is generated based on the values of an independent set of random vectors
- Randomness is injected into the model-building process



# Random Forest



Node splitting in a random forest model is based on a random subset of features for each tree.

# Random Forest

Original Training Set					
Col1	Col2	Col3	Col4	Col5	Col6
1	Sdf	200	A	1	.88
3	Fg	200	A	1	.67
2	Wdv	290	A	1	.36
4	Gh	345	B	0	.85
1	J	125	AB	0	.72
3	Xcv	543	B	0	.93
2	gbn	367	A	1	.18

Training Subsets via Bootstrapping									
Col1	Col2	Col4	Col5	Col6	Col1	Col3	Col4	Col5	Col6
1	Sdf	A	1	.88	1	200	A	1	.88
3	Fg	A	1	.67	3	200	A	1	.67
Col2	Col3	Col4	Col5	Col6	Col1	Col2	Col3	Col4	Col5
Wdv	290	A	1	.36	1	Sdf	200	A	1
Gh	345	B	0	.85	2	Wdv	290	A	1
Col1	Col2	Col3	Col5	Col6	Col1	Col2	Col3	Col4	Col6
3	Fg	200	1	.67	1	Sdf	200	A	.88
2	Wdv	290	1	.36	3	Fg	200	A	.67
Col1	Col2	Col3	Col4	Col6	Col1	Col2	Col3	Col4	Col5
1	Sdf	200	A	.88	3	Fg	200	A	1
3	Fg	200	A	.67	2	Wdv	290	A	1
Col2	Col3	Col4	Col5	Col6	Col2	Col3	Col4	Col5	Col6
Sdf	200	A	1	.88	Sdf	200	A	1	.88
Wdv	290	A	1	.36	Fg	200	A	1	.67
J	125	AB	0	.72	J	125	AB	0	.72
Xcv	543	B	0	.93	Xcv	543	B	0	.93

Figure 6. The decision trees in a random forest are all slightly differently trained on a bootstrapped subset of the original dataset. The set of input features also varies for each decision tree in the random forest.

# Random Vector selection – Forest RI (Random Input selection)



- Randomly select  $F$  input features to split at each node of the decision tree and then fully grow the tree without pruning
- This helps reduce the bias present in the resulting tree
- The predictions are combined using a majority voting scheme
- To increase randomness, bagging can also be used to generate bootstrap samples
- The strength and correlation of random forests may depend on the size of  $F$  features
  - If  $F$  is sufficiently small, then the trees tend to become less correlated
- The strength of the tree classifier tends to improve with a larger number of  $F$
- Optimal number of  $F = \log_2 d + 1$  (where  $d$  is number of input features)
- Reduces the runtime of the algorithm also

# Random Vector selection – Forest RC (Random Combinations)

- Used when number of features  $d$  is small
- Increases the feature space is to create linear combinations of the input features
- Specifically, at each node, a new feature is generated by randomly selecting  $L$  of the input features
- The input features are linearly combined
- At each node,  $F$  of such randomly combined new features are generated
- The best of them is subsequently selected to split the node

# Example

## Original Dataset

Original Dataset				
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

## Test Instance to classify

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	



# Random Forest Hyperparameters

---

- `max_depth`
- `min_sample_split`
- `max_leaf_nodes`
- `min_samples_leaf`
- `n_estimators`
- `max_sample` (bootstrap sample)
- `max_features`

# Advantages of Random Forest

---

- Algorithm can solve both type of problems i.e. classification and regression
- Handles large data set with higher dimensionality.
- It can handle thousands of input variables and identify most significant variables so it is considered as one of the dimensionality reduction methods.
- Model outputs **Importance of variable**, which can be a very handy feature (on some random data set).

# Disadvantages of Random Forest

---

- May over-fit data sets that are particularly noisy.
- Random Forest can feel like a black box approach for statistical modelers – you have very little control on what the model does. You can at best – try different parameters and random seeds!

- Adaptive boosting or AdaBoost is one of the simplest boosting algorithms. Usually, decision trees are used for modelling. Multiple sequential models are created, each correcting the errors from the last model.
- AdaBoost assigns weights to the observations which are incorrectly predicted and the subsequent model works to predict these values correctly.

# Adaboost Algorithm



- Initially, all observations ( $n$ ) in the dataset are given equal weights ( $1/n$ ).
- A model is built on a subset of data.
- Using this model, predictions are made on the whole dataset.
- Errors are calculated by comparing the predictions and actual values.
- While creating the next model, higher weights are given to the data points which were predicted incorrectly.

# Adaboost Algorithm



- Weights can be determined using the error value. For instance, higher the error more is the weight assigned to the observation.
- This process is repeated until the error function does not change, or the maximum limit of the number of estimators is reached.

- Base classifiers  $C_i$ :  $C_1, C_2, \dots, C_T$
- Error rate:
  - N input samples

$$\varepsilon_i = \frac{1}{N} \sum_{j=1}^N w_j \delta(C_i(x_j) \neq y_j)$$

- Importance of a classifier:

$$\alpha_i = \frac{1}{2} \ln \left( \frac{1 - \varepsilon_i}{\varepsilon_i} \right)$$

[https://en.wikipedia.org/wiki/AdaBoost#Choosing\\_at](https://en.wikipedia.org/wiki/AdaBoost#Choosing_at)

# AdaBoost: Weight Update



Weight Update: 
$$w_i^{(j+1)} = \frac{w_i^{(j)}}{Z_j} \begin{cases} \exp^{-\alpha_j} & \text{if } C_j(x_i) = y_i \\ \exp^{\alpha_j} & \text{if } C_j(x_i) \neq y_i \end{cases} \quad \leftarrow \text{Eqn:5.88}$$

where  $Z_j$  is the normalization factor

$$C^*(x) = \arg \max_y \sum_{j=1}^T \alpha_j \delta(C_j(x) = y)$$

- Reduce weight if correctly classified else increase
- If any intermediate rounds produce error rate higher than 50%, the weights are reverted back to  $1/n$  and the resampling procedure is repeated



# Classifier weight



$$\varepsilon = 0.3$$

$$\alpha = 1/2 * \ln(1 - 0.3 / 0.3) = 0.42365$$

$$\varepsilon = 0.7$$

$$\alpha = 1/2 * \ln(1 - 0.7 / 0.7) = -0.42365$$

$$\varepsilon = 0.5$$

$$\alpha = 1/2 * \ln(1 - 0.5 / 0.5) = 0$$

Notice three interesting observations:

- 1) classifier with accuracy higher than 50% results in a positive weight for the classifier  
(in other words,  $\alpha > 0$  if  $\varepsilon \leq 0.5$ ),
- 2) classifier with exact 50% accuracy is 0, and thus, does not contribute to the final prediction, and
- 3) errors 0.3 and 0.7 lead to classifier weights with inverse signs.

# Training instance weight update

Diagram illustrating the weight update for a misclassified instance (actual = -1, predicted = 1):

$$D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

Actual: -1, Predicted: 1

Calculation:  $-0.42 \times -1 = 0.42$

Final result:  $e^{0.42} = 1.52$

Diagram illustrating the weight update for a correctly classified instance (actual = 1, predicted = 1):

$$D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

Actual: 1, Predicted: 1

Calculation:  $-0.42 \times 1 = -0.42$

Final result:  $e^{-0.42} = 0.657$

# AdaBoost Algorithm

---

## Algorithm 5.7 AdaBoost Algorithm

---

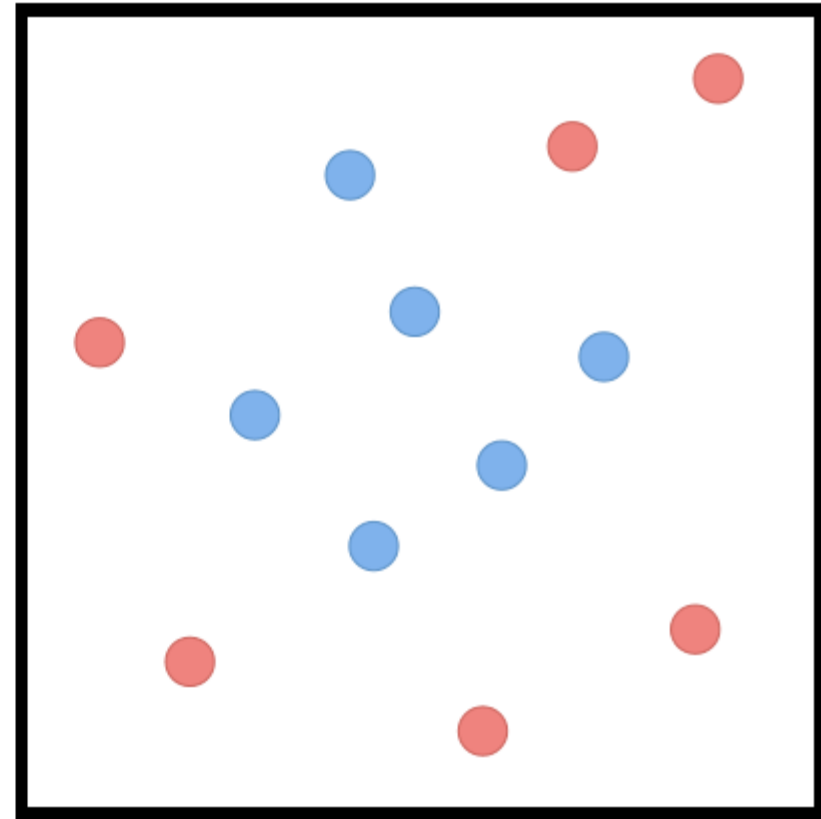
- 1:  $\mathbf{w} = \{w_j = 1/n \mid j = 1, 2, \dots, n\}$ .    {Initialize the weights for all  $n$  instances.}
  - 2: Let  $k$  be the number of boosting rounds.
  - 3: for  $i = 1$  to  $k$  do
  - 4:    Create training set  $D_i$  by sampling (with replacement) from  $D$  according to  $\mathbf{w}$ .
  - 5:    Train a base classifier  $C_i$  on  $D_i$ .
  - 6:    Apply  $C_i$  to all instances in the original training set,  $D$ .
  - 7:     $\epsilon_i = \frac{1}{n} [\sum_j w_j \delta(C_i(x_j) \neq y_j)]$     {Calculate the weighted error}
  - 8:    if  $\epsilon_i > 0.5$  then
  - 9:      $\mathbf{w} = \{w_j = 1/n \mid j = 1, 2, \dots, n\}$ .    {Reset the weights for all  $n$  instances.}
  - 10:    Go back to Step 4.
  - 11:    end if
  - 12:     $\alpha_i = \frac{1}{2} \ln \frac{1-\epsilon_i}{\epsilon_i}$ .
  - 13:    Update the weight of each instance according to equation (5.88).
  - 14: end for
  - 15:  $C^*(\mathbf{x}) = \arg \max_y \sum_{j=1}^T \alpha_j \delta(C_j(\mathbf{x}) = y)$ .
-

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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



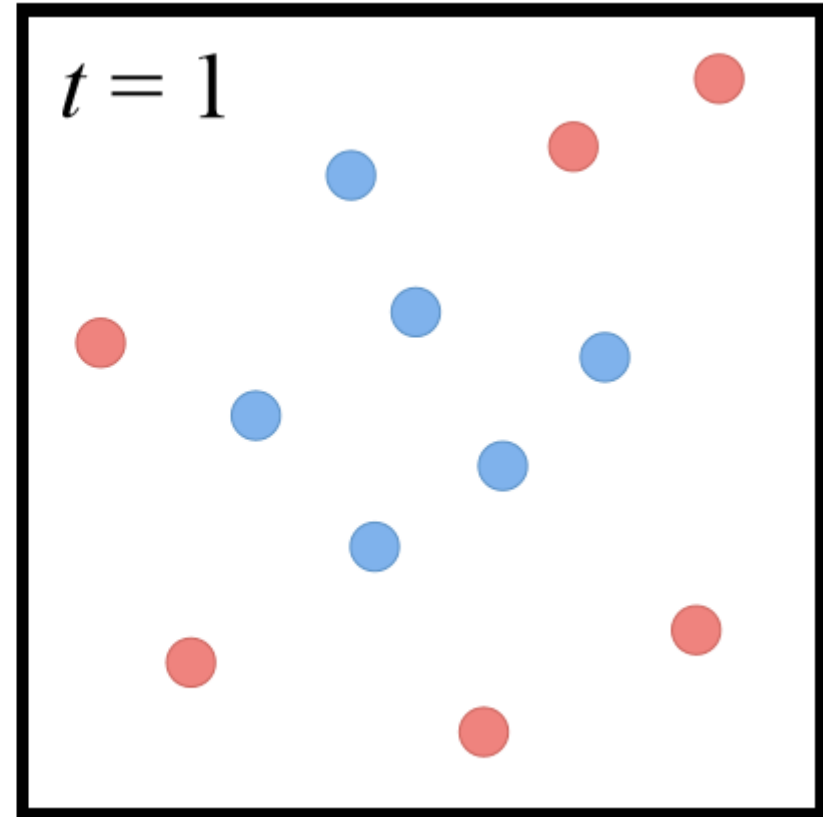
- Size of point represents the instance's weight

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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

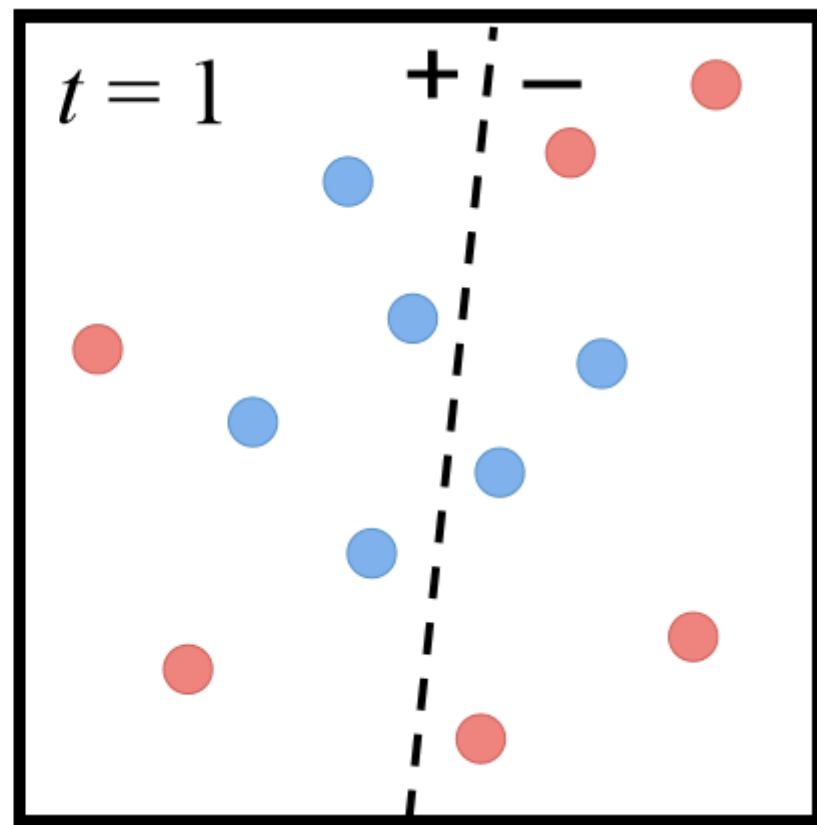


# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
     $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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

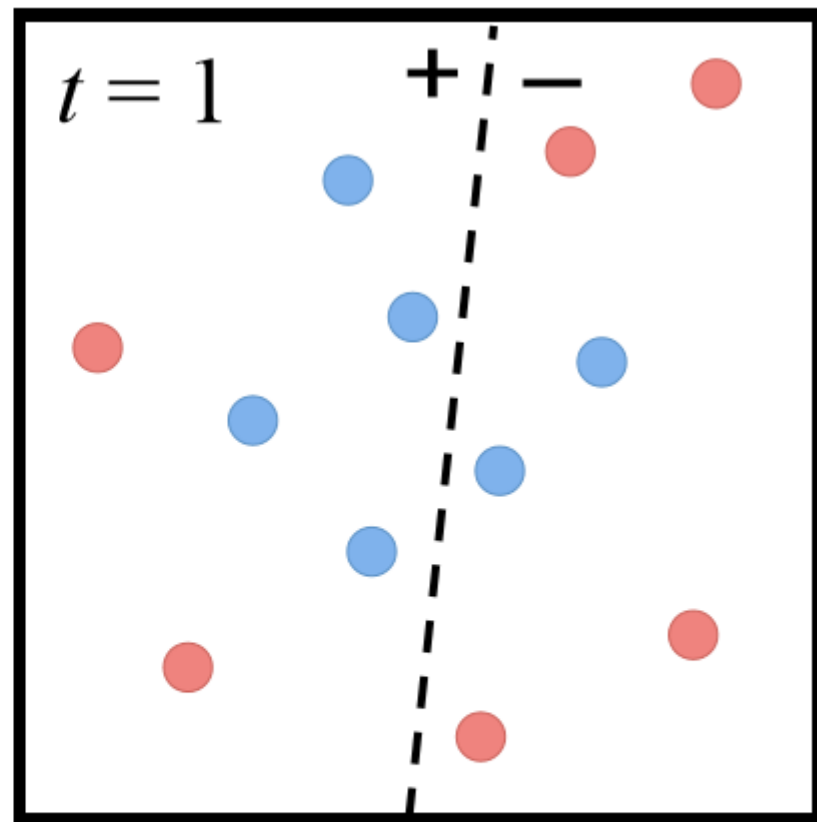


# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $w_1$
- 2: for  $t = 1, \dots, T$
- 3: Train model  $h_t$  on  $X, y$  with weights  $w_t$
- 4: Compute the weighted training error of  $h_t$
- 5: Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6: Update all instance weights:  
 $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7: Normalize  $w_{t+1}$  to be a distribution
- 8: end for
- 9: Return the hypothesis

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



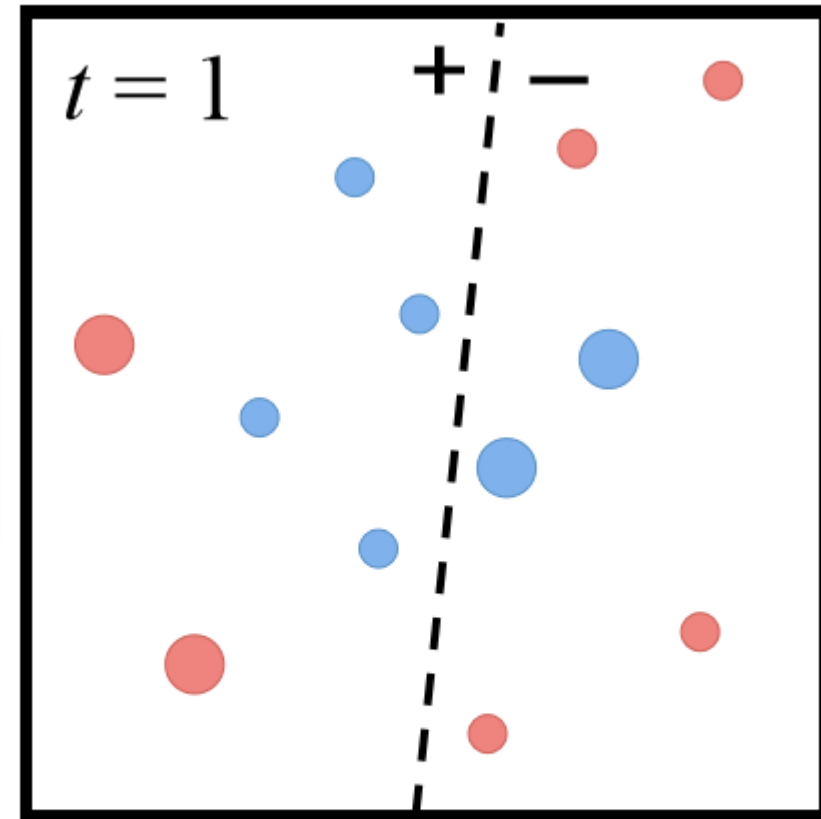
- $\beta_t$  measures the importance of  $h_t$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  ( $\beta_t$  grows as  $\epsilon_t$  gets smaller)

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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



- Weights of correct predictions are multiplied by  $e^{-\beta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{\beta_t} \geq 1$

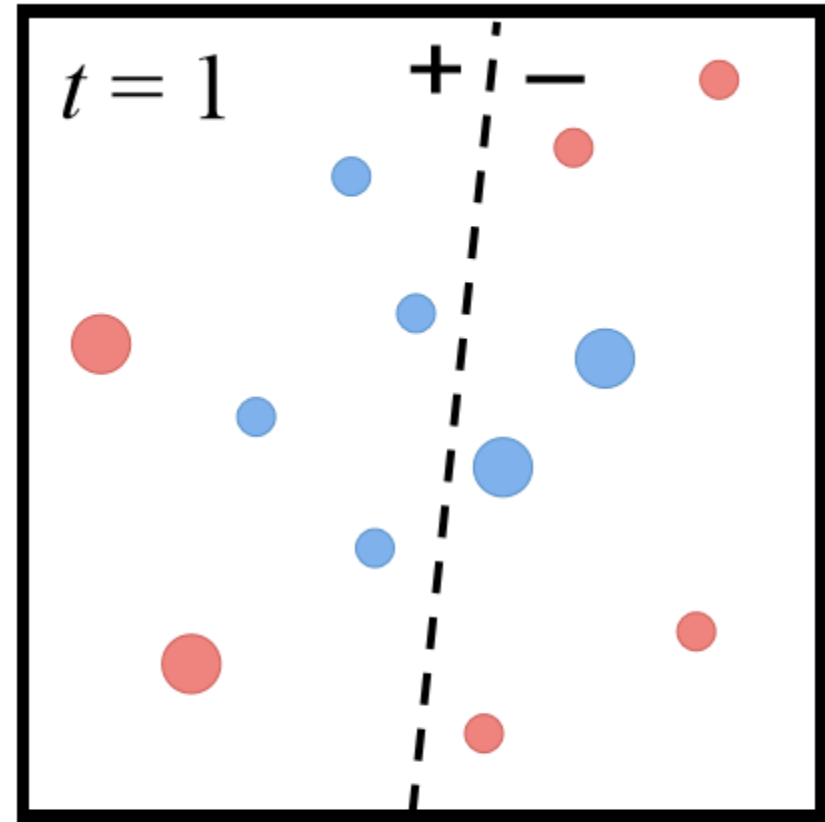


# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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



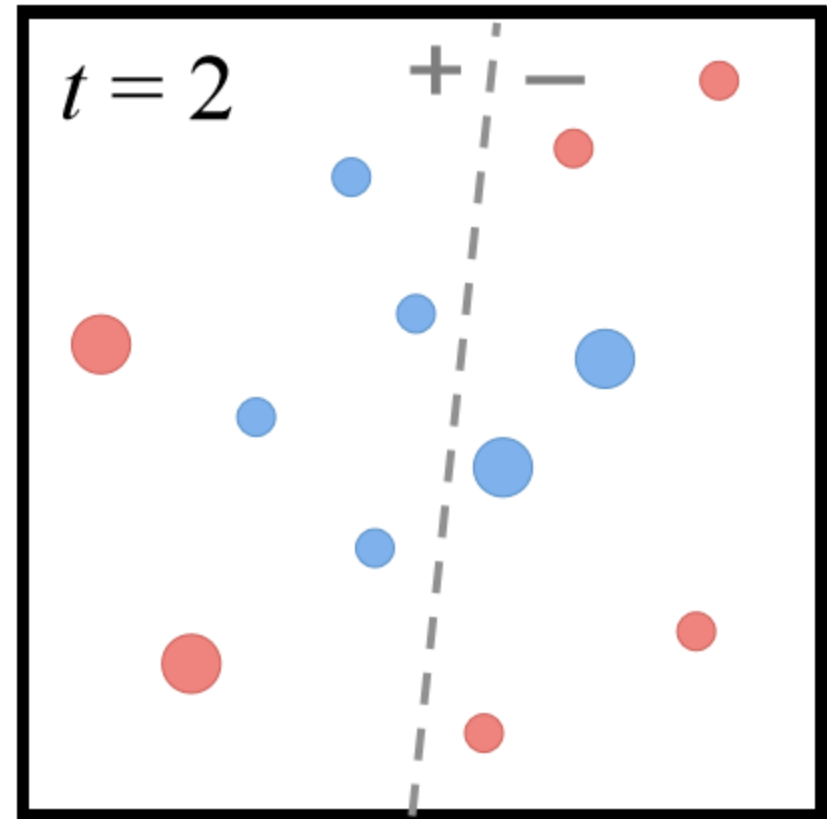
Disclaimer: Note that resized points in the illustration above are not necessarily to scale with  $\beta_t$

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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

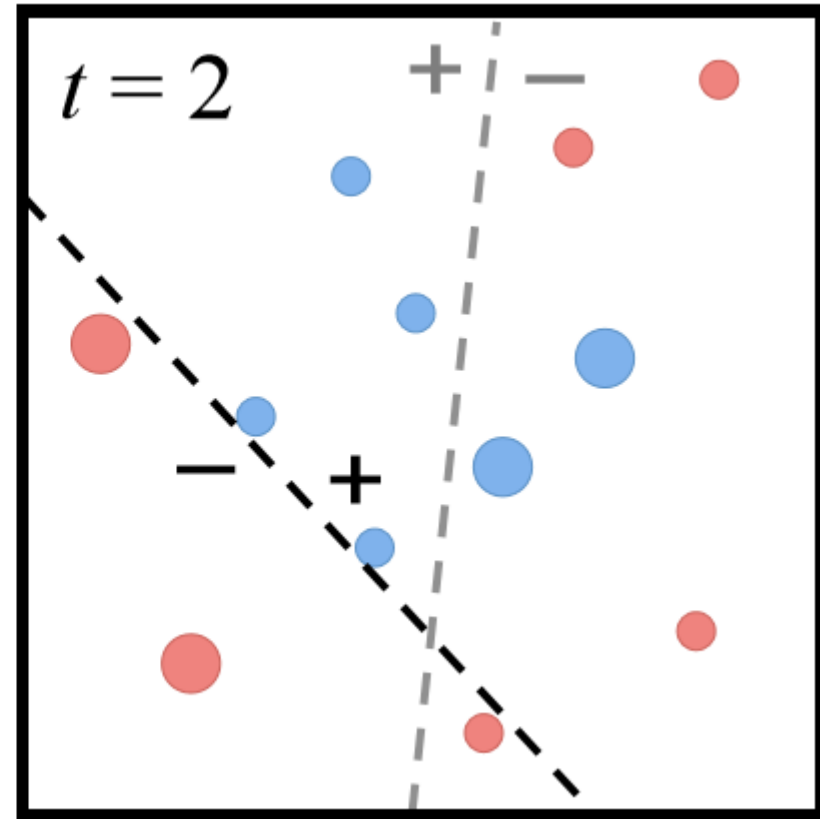


# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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

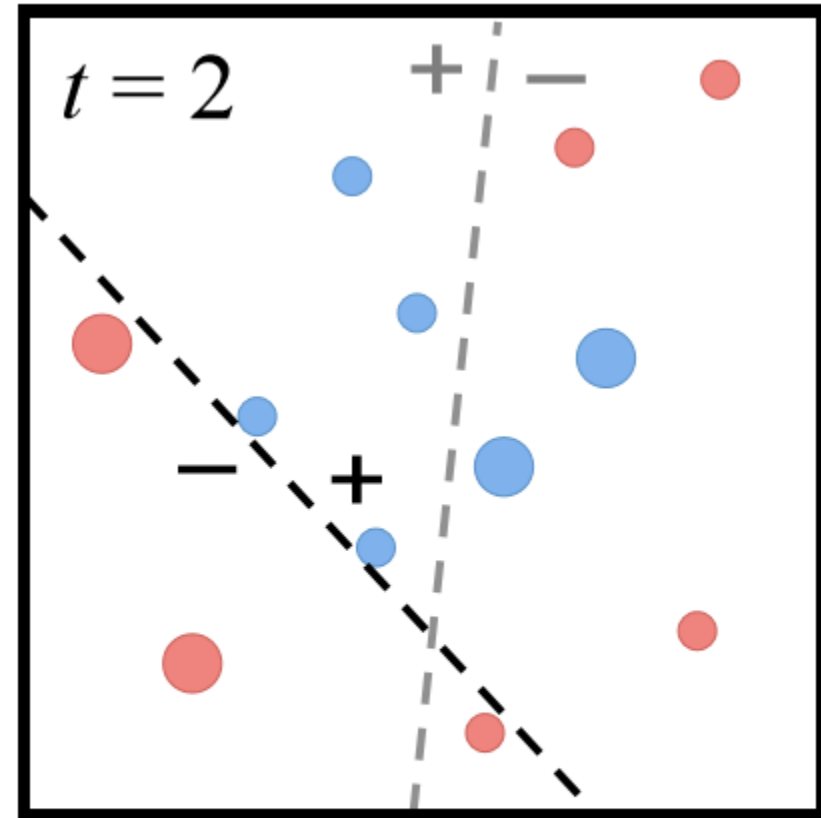


# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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



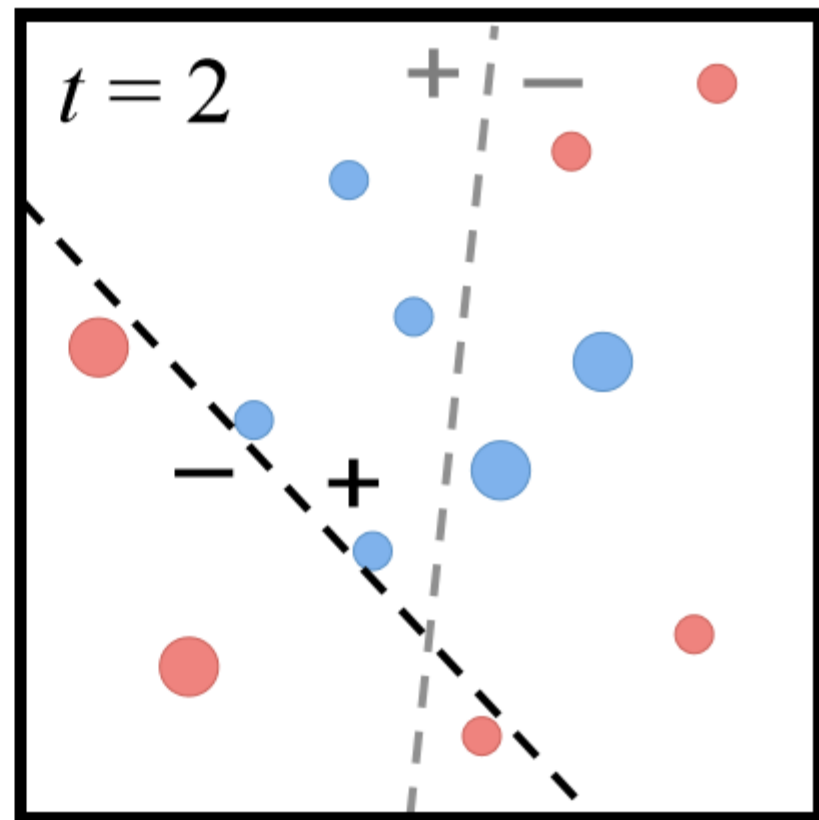
- $\beta_t$  measures the importance of  $h_t$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  ( $\beta_t$  grows as  $\epsilon_t$  gets smaller)

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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



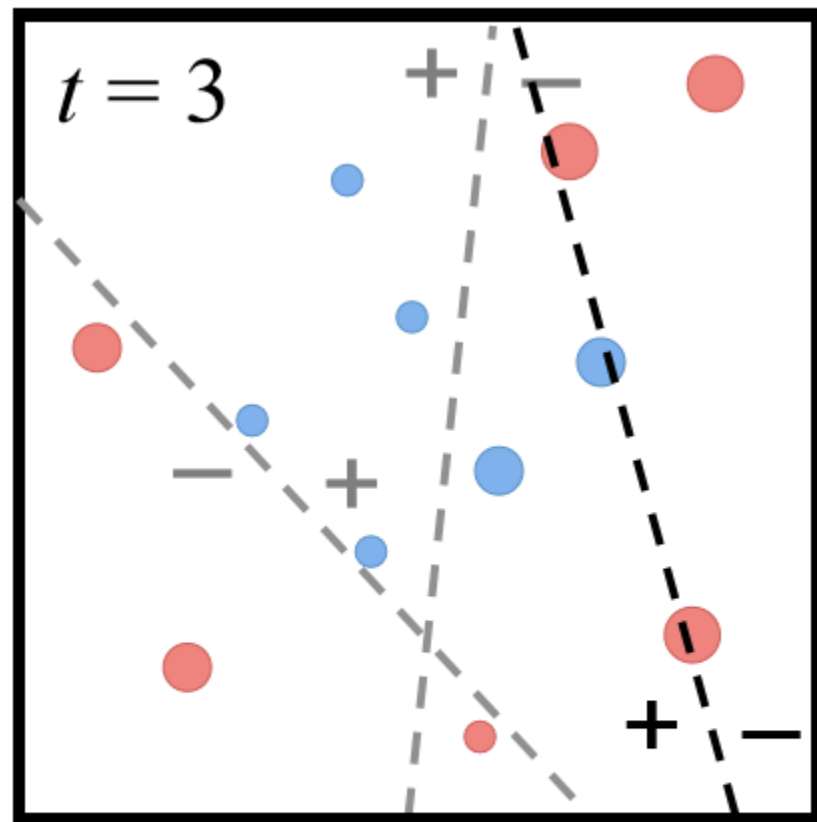
- Weights of correct predictions are multiplied by  $e^{-\beta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{\beta_t} \geq 1$

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $w_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $w_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $w_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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

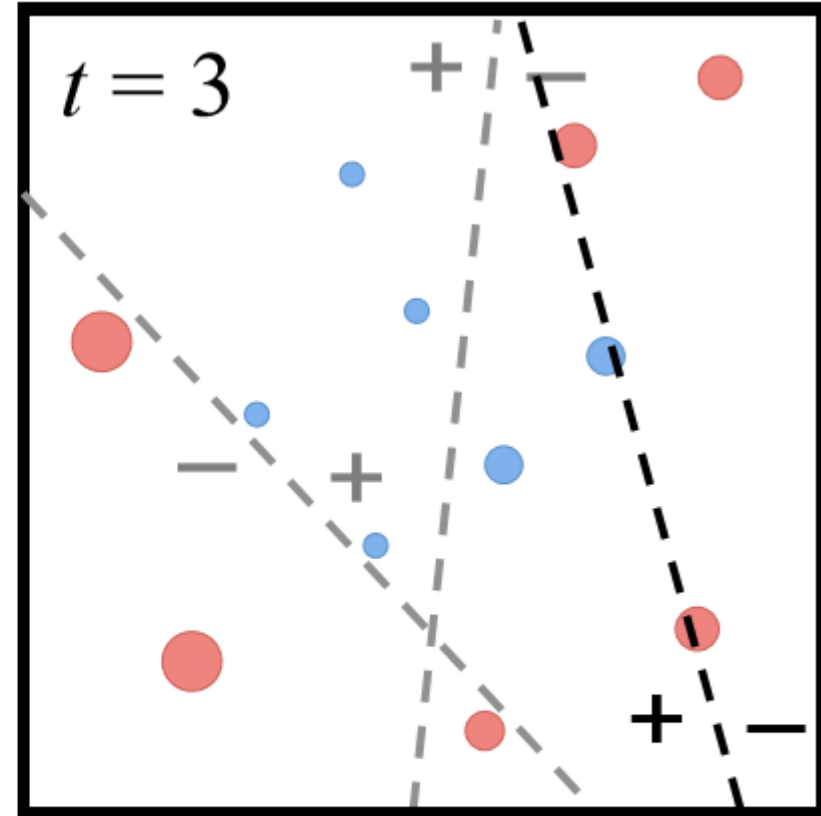


- $\beta_t$  measures the importance of  $h_t$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  ( $\beta_t$  grows as  $\epsilon_t$  gets smaller)

# AdaBoost Algorithm

- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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



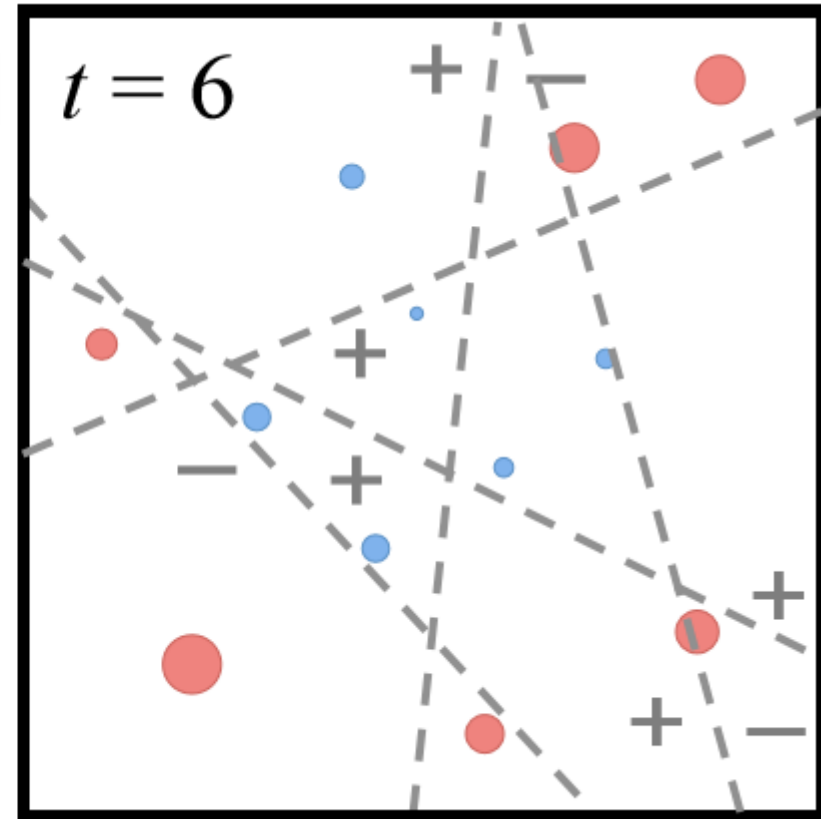
- Weights of correct predictions are multiplied by  $e^{-\beta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{\beta_t} \geq 1$

# AdaBoost Algorithm



- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**
- 9: **Return** the hypothesis

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





# AdaBoost Algorithm

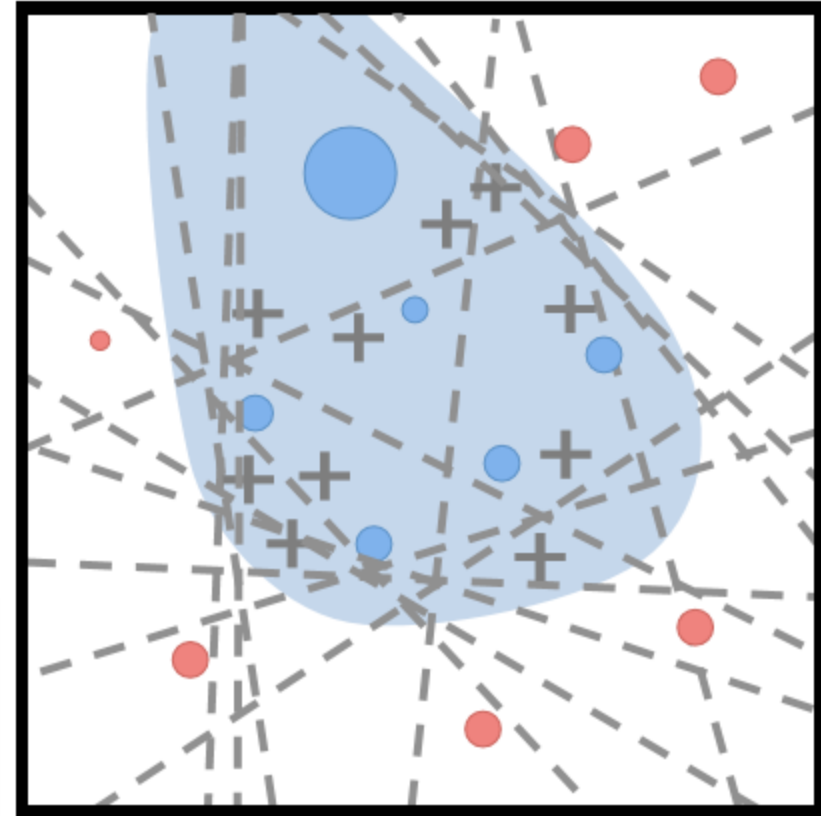


$t = T$

- 1: Initialize a vector of  $n$  uniform weights  $\mathbf{w}_1$
- 2: **for**  $t = 1, \dots, T$
- 3:   Train model  $h_t$  on  $X, y$  with weights  $\mathbf{w}_t$
- 4:   Compute the weighted training error of  $h_t$
- 5:   Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
- 6:   Update all instance weights:  
       $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7:   Normalize  $\mathbf{w}_{t+1}$  to be a distribution
- 8: **end for**

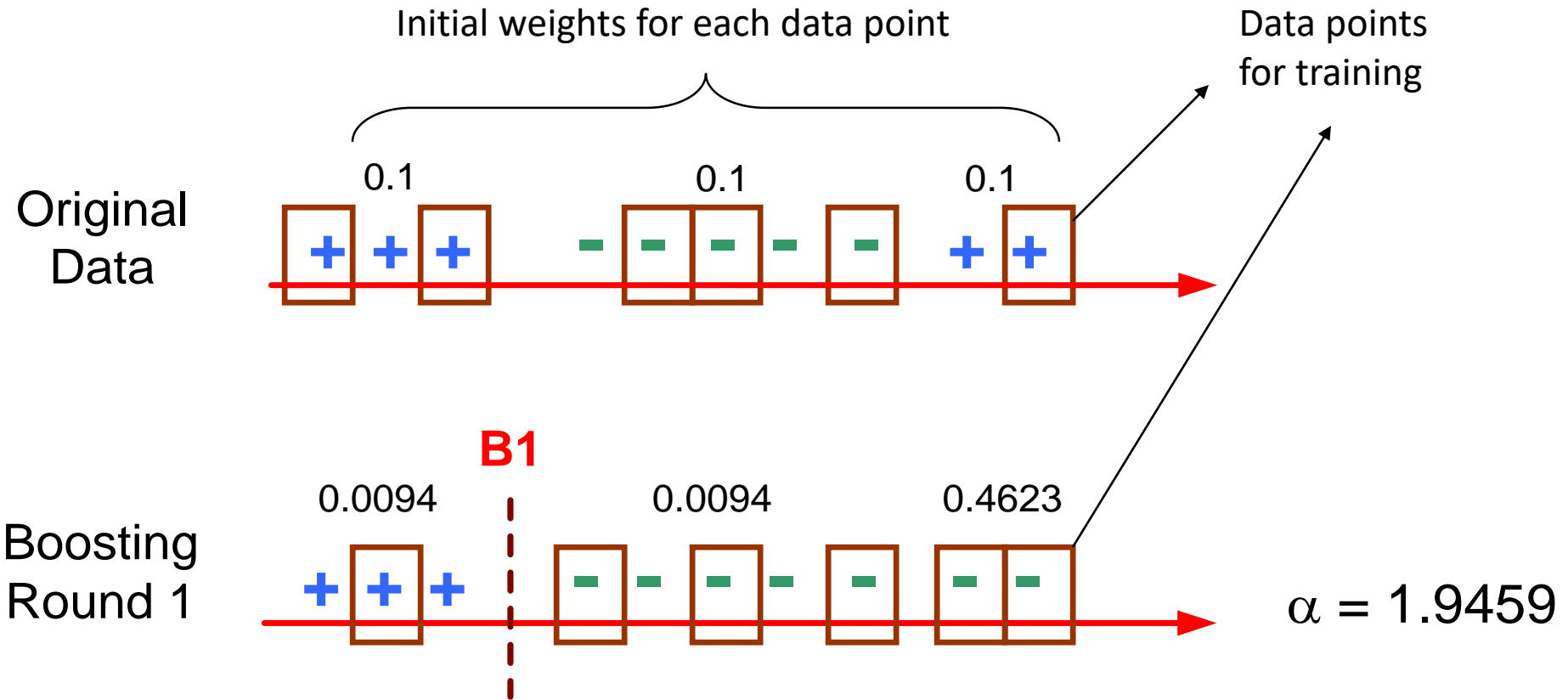
9: **Return** the hypothesis

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

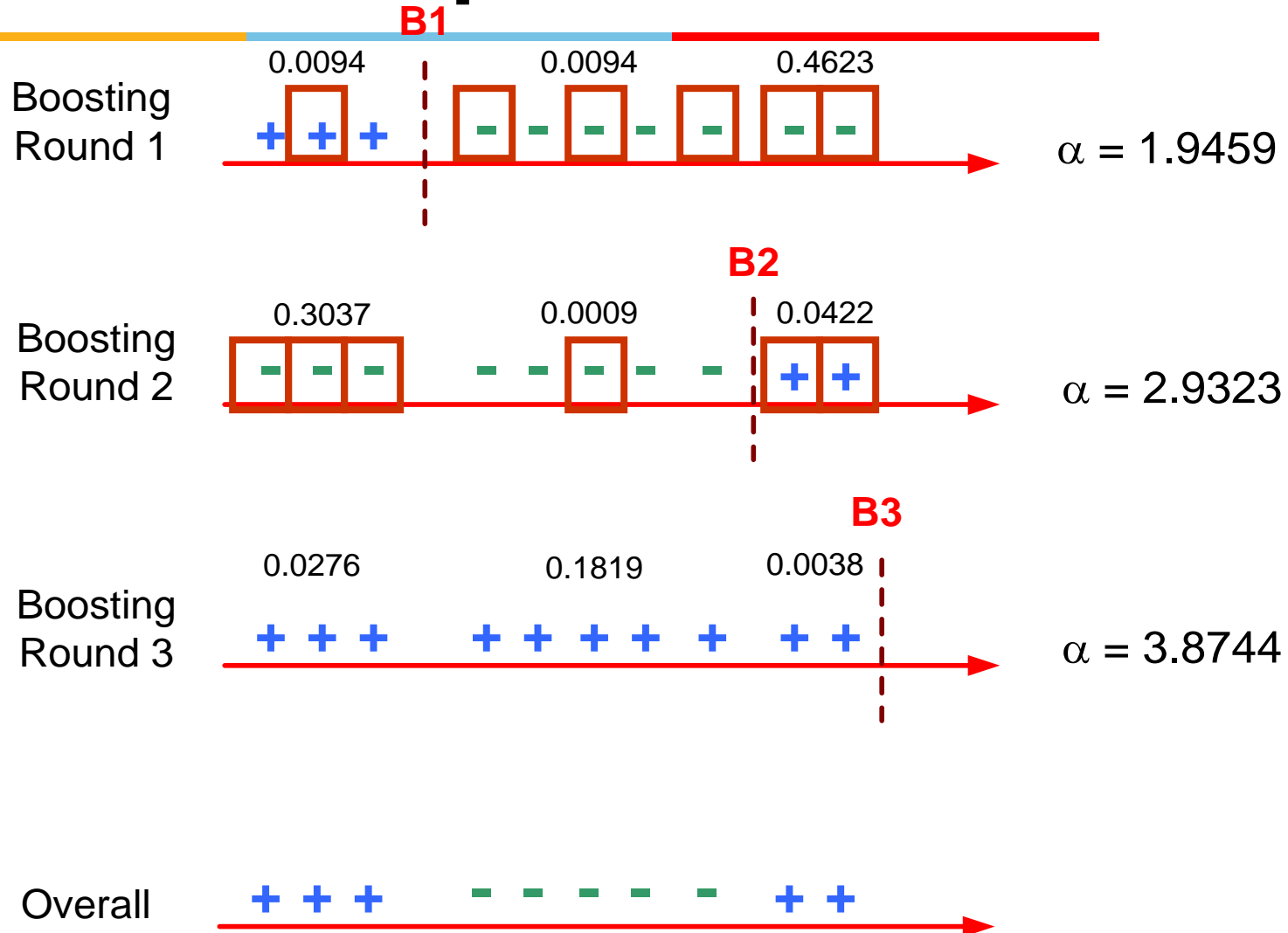


- Final model is a weighted combination of members
  - Each member weighted by its importance

# Adaboost Example



# Adaboost Example



# AdaBoost Example

- Training sets for the first 3 boosting rounds:

Boosting Round 1:

x	0.1	0.4	0.5	0.6	0.6	0.7	0.7	0.7	0.8	1
y	1	-1	-1	-1	-1	-1	-1	-1	1	1

Boosting Round 2:

x	0.1	0.1	0.2	0.2	0.2	0.2	0.3	0.3	0.3	0.3
y	1	1	1	1	1	1	1	1	1	1

Boosting Round 3:

x	0.2	0.2	0.4	0.4	0.4	0.4	0.5	0.6	0.6	0.7
y	1	1	-1	-1	-1	-1	-1	-1	-1	-1

- Summary:

Round	Split Point	Left Class	Right Class	alpha
1	0.75	-1	1	1.738
2	0.05	1	1	2.7784
3	0.3	1	-1	4.1195

# AdaBoost Example

- Weights

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
2	0.311	0.311	0.311	0.01	0.01	0.01	0.01	0.01	0.01	0.01
3	0.029	0.029	0.029	0.228	0.228	0.228	0.228	0.009	0.009	0.009

- Classification

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	-1	-1	-1	-1	-1	-1	-1	1	1	1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	-1	-1	-1	-1	-1	-1	-1
Sum	5.16	5.16	5.16	-3.08	-3.08	-3.08	-3.08	0.397	0.397	0.397
Predicted Class	1	1	1	-1	-1	-1	-1	1	1	1

AdaBoost error function takes into account the fact that only the sign of the final result is used, thus sum can be far larger than 1 without increasing error

# Gradient Boosting

- In Gradient Boosting, "shortcomings" are identified by gradients.
- Recall that, in Adaboost, "shortcomings" are identified by high-weight data points. Both high-weight data points and gradients tell us how to improve our model.

# Gradient Boosting

---

- Gradient Boosting for Different Problems  
Difficulty: regression ==> classification ==> ranking

# Gradient Boosting

- You are given  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , and the task is to fit a model  $F(x)$  to minimize square loss
- There are some mistakes:  $F(x_1) = 0.8$ , while  $y_1 = 0.9$ , and  $F(x_2) = 1.4$  while  $y_2 = 1.3$ ... How can you improve this model?
- Rules:
  - You are not allowed to remove anything from  $F$  or change any parameter in  $F$ .
  - You can add an additional model (regression tree)  $h$  to  $F$ , so the new prediction will be  $F(x) + h(x)$ .



# Gradient Boosting

- You wish to improve the model such that
  - $F(x_1) + h(x_1) = y_1$
  - $F(x_2) + h(x_2) = y_2 \dots$
  - $F(x_n) + h(x_n) = y_n$

Or, equivalently, you wish

$$h(x_1) = y_1 - F(x_1)$$

$$h(x_2) = y_2 - F(x_2) \dots$$

$$h(x_n) = y_n - F(x_n)$$

Fit a regression tree  $h$  to data

$$(x_1, y_1 - F(x_1)), (x_2, y_2 - F(x_2)), \dots, (x_n, y_n - F(x_n))$$

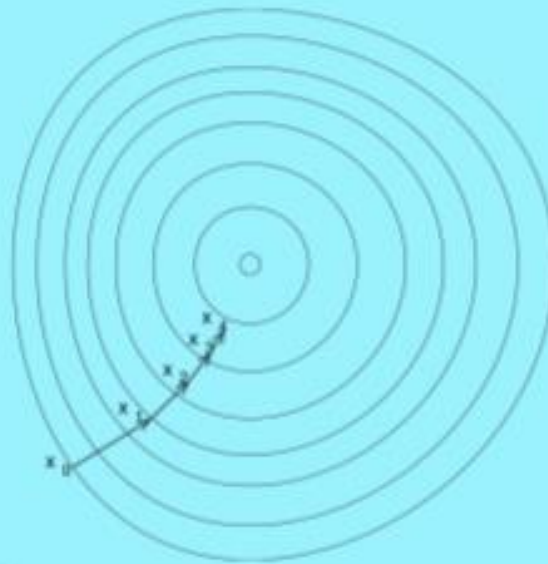
# Gradient Boosting

- Simple solution:  $y_i - F(x_i)$  are called residuals. These are the parts that existing model  $F$  cannot do well.
- The role of  $h$  is to compensate the shortcoming of existing model  $F$ .
- If the new model  $F + h$  is still not satisfactory, we can add another regression tree...
- We are improving the predictions of training data, is the procedure also useful for test data?

# Gradient Descent

Minimize a function by moving in the opposite direction of the gradient.

$$\theta_i := \theta_i - \rho \frac{\partial J}{\partial \theta_i}$$



# Gradient Boosting for regression

Loss function  $L(y, F(x)) = (y - F(x))^2/2$

We want to minimize  $J = \sum_i L(y_i, F(x_i))$  by adjusting  $F(x_1), F(x_2), \dots, F(x_n)$ .

Notice that  $F(x_1), F(x_2), \dots, F(x_n)$  are just some numbers. We can treat  $F(x_i)$  as parameters and take derivatives

$$\frac{\partial J}{\partial F(x_i)} = \frac{\partial \sum_i L(y_i, F(x_i))}{\partial F(x_i)} = \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} = F(x_i) - y_i$$

So we can interpret residuals as negative gradients.

$$y_i - F(x_i) = -\frac{\partial J}{\partial F(x_i)}$$

$$F(x_i) := F(x_i) + h(x_i)$$

$$F(x_i) := F(x_i) + y_i - F(x_i)$$

$$F(x_i) := F(x_i) - 1 \frac{\partial J}{\partial F(x_i)}$$

$$\theta_i := \theta_i - \rho \frac{\partial J}{\partial \theta_i}$$

For regression with **square loss**,

*residual  $\Leftrightarrow$  negative gradient*

*fit  $h$  to residual  $\Leftrightarrow$  fit  $h$  to negative gradient*

*update  $F$  based on residual  $\Leftrightarrow$  update  $F$  based on negative gradient*

So we are actually updating our model using **gradient descent**!

# Gradient Boosting Algorithm

- It involves three elements
  - A loss function to be optimized (minimizes expected value)
 
$$\hat{F} = \arg \min_F \mathbb{E}_{x,y} [L(y, F(x))]$$
  - Approximation of  $F(x)$  in terms of weighted sum of base(weak) learners  $h_i(x)$  to make predictions
 
$$\hat{F}(x) = \sum_{i=1}^M \gamma_i h_i(x) + \text{const}$$
  - An additive model to minimize the loss function, starting with  $F_0(x)$  and incrementally expanding in greedy fashion

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma),$$

$$F_m(x) = F_{m-1}(x) + \arg \min_{h_m \in \mathcal{H}} \left[ \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + h_m(x_i)) \right]$$

# Why XGBoost so popular?



- **Speed** : faster than other ensemble classifiers.
- **Core algorithm is parallelizable**: harness the power of multi-core computers and networks of computers enabling to train on very large datasets **Consistently outperforms other algorithm methods** : It has shown better performance on a variety of machine learning benchmark datasets.
- **Wide variety of tuning parameters** : cross-validation, regularization, missing values, tree parameters, etc
- XGBoost (Extreme Gradient Boosting) uses the gradient boosting (GBM) framework at its core.

# Class Imbalance Problem

- Find needle in haystack
- Lots of classification problems where the classes are skewed (more records from one class than another)
  - Credit card fraud
  - Intrusion detection
  - Defective products in manufacturing assembly line



# Approaches to solve Class imbalance problem

- Up-sample minority class
  - randomly duplicating observations from a minority class
- Down-sample majority class
  - removing random observations.
- Generate Synthetic Samples
  - new samples based on the distances between the point and its nearest neighbors
- Change the performance metric
  - Use Recall, Precision or ROC curves instead of accuracy
- Try different algorithms
  - Some algorithms as Support Vector Machines and Tree-Based algorithms are better to work with imbalanced classes.

# Challenges

- Evaluation measures such as accuracy is not well-suited for imbalanced class
- Detecting the rare class is like finding needle in a haystack

# Python Implementation Demo

# Good References

## Ensemble methods

<https://www.slideshare.net/hustwj/an-introduction-to-ensemble-methodsboosting-bagging-random-forests-and-more>

## Bagging and Boosting

- [https://www.youtube.com/watch?time\\_continue=2&v=m-S9Hojj1as](https://www.youtube.com/watch?time_continue=2&v=m-S9Hojj1as)

**Gradient Boosting :** <https://www.youtube.com/watch?v=jxuNLH5dXCs>

## XGBoost

- <https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/>
- <https://xgboost.readthedocs.io/en/latest/tutorials/model.html>
- [https://www.youtube.com/watch?time\\_continue=71&v=Vly8xGnNiWs](https://www.youtube.com/watch?time_continue=71&v=Vly8xGnNiWs)
- <https://www.slideshare.net/ShangxuanZhang/xgboost-55872323>
- <https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/>
- <https://towardsdatascience.com/a-beginners-guide-to-xgboost-87f5d4c30ed7>
- <https://towardsdatascience.com/machine-learning-for-diabetes-562dd7df4d42>

## Confusion Matrix

<https://towardsdatascience.com/confusion-matrix-for-your-multi-class-machine-learning-model-ff9aa3bf7826>

## ROC

<https://www.youtube.com/watch?v=OAl6eAyP-yo&feature=youtu.be>

<https://www.datacamp.com/community/tutorials/xgboost-in-python>

---

# Thank You