



AI Workshop @ NTNE

ML Models

Chun-Hsiang Chan

Department of Geography, National Taiwan Normal University

2024.10.24

Outline

10.24 | 機器學習演算法

- Concept of Classification
- Evaluation Metrics
- Model Optimization
- Overfitting vs Underfitting
- Bagging
- Boosting
- K-Nearest Neighbors
- Decision Tree
- Random Forest
- AdaBoost
- Gradient Boost
- XGBoost
- LightGBM
- CatBoost
- Naive Bayes
- Linear SVM
- Remarks



Open Google Colab

<https://bit.ly/1017ai>



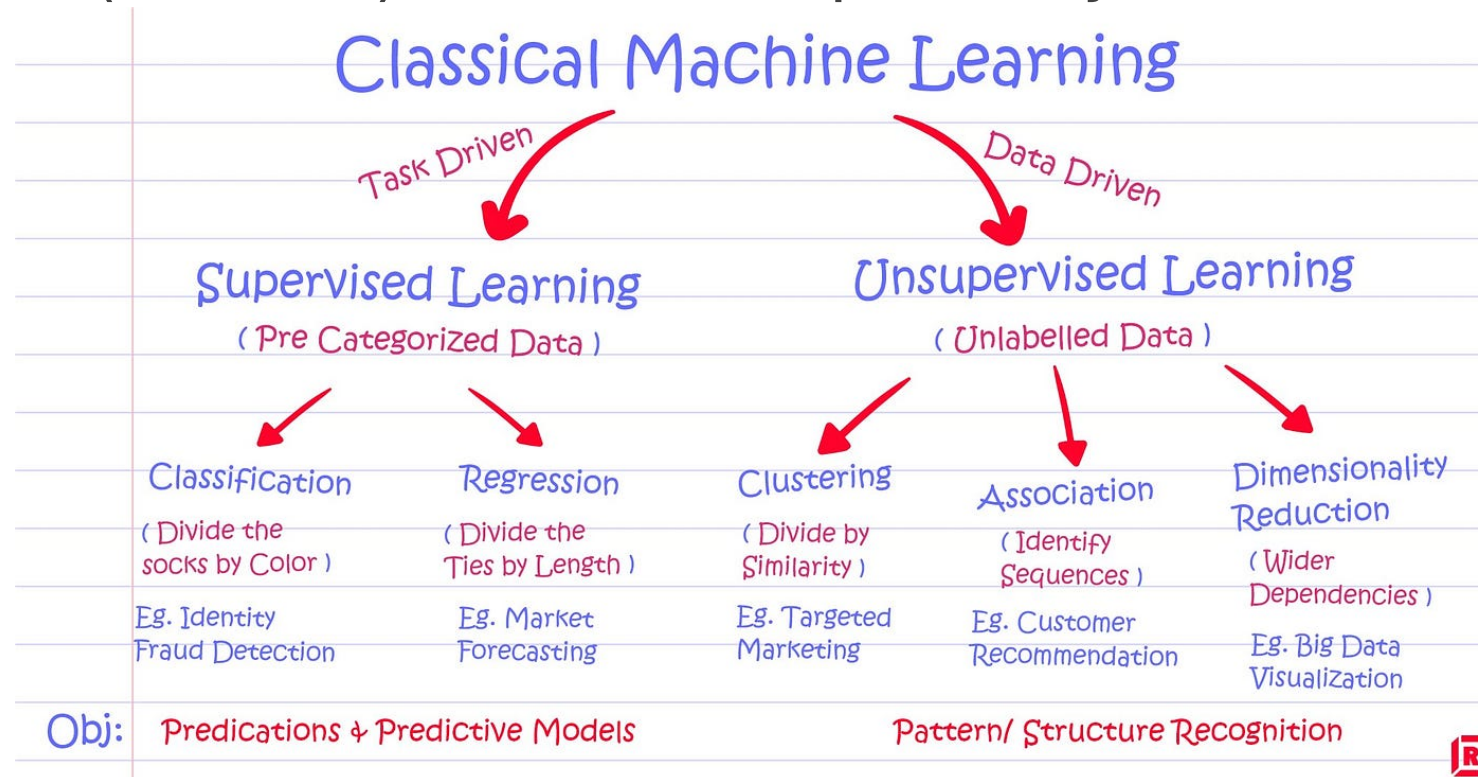
Google Colaboratory

simple**LEARN**

Concept of Classification

Supervised vs unsupervised learning

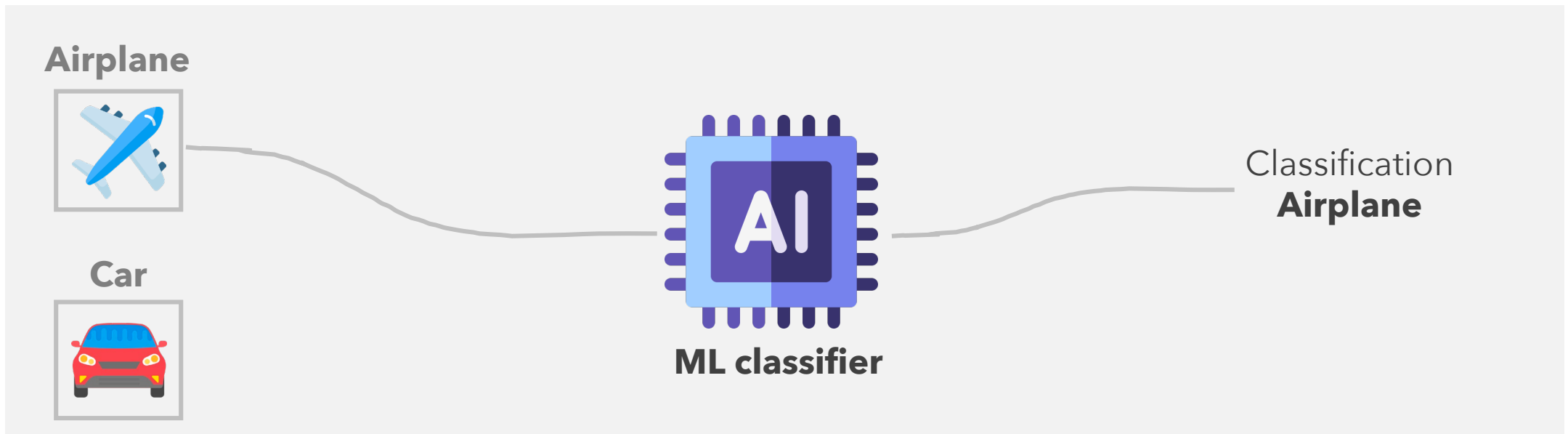
- Regression and classification are designed for continuous and categorical (discrete) datasets, respectively.



Concept of Classification

Classifier Types

- Basically, all machine learning classifiers consist of binary, multi-class, and multi-label classifications.
- **Binary Classifier:**



Concept of Classification

Classifier Types

- **Multi-class classification:** classify airplanes, cars, and ships.

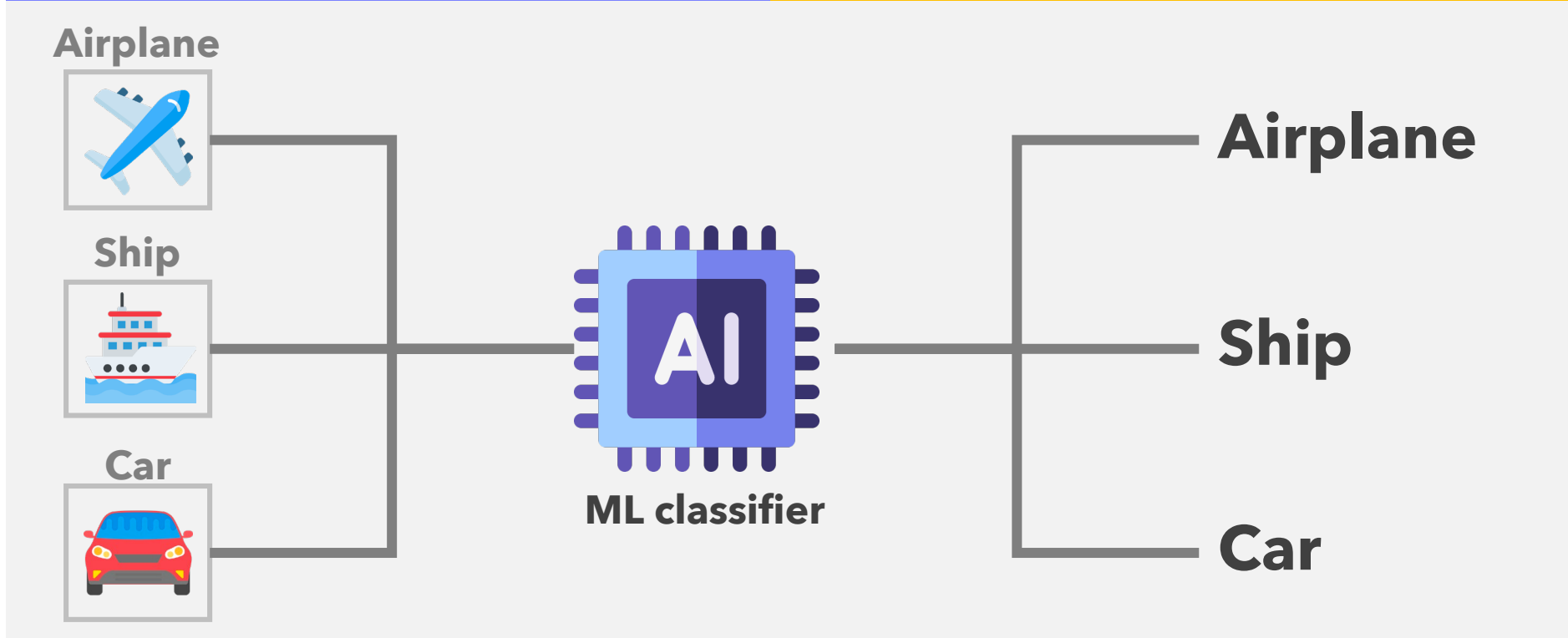


Concept of Classification

Classifier Types

- Multi-label classification:

One Model classifies all categories ...



Concept of Classification

Data Imbalance Issue

- Normally, it is too difficult to collect all categories with a huge sample size; therefore, some categories might have only a few samples, leading to a bad prediction result.
- Think about why the ML models cannot predict the small-sample-size category well?
- If so, how can we overcome this dilemma?
- Please list your solutions...

Concept of Classification

- **Possible solutions:**
 - **Cluster-based oversampling.**
 - **Random undersampling:** random elimination of examples from the majority class.
 - **SMOTE oversampling:** random replication of examples from the minority class.

SMOTE

```
class imblearn.over_sampling.SMOTE(*, sampling_strategy='auto',  
random_state=None, k_neighbors=5, n_jobs=None)
```

[\[source\]](#)

Class to perform over-sampling using SMOTE.

This object is an implementation of SMOTE - Synthetic Minority Over-sampling Technique as presented in [\[1\]](#).

Read more in the [User Guide](#).

Evaluation Metrics

- Evaluation metrics are vital elements to quantitatively assess the performance of ML model training and testing processes.
- The most common way to demonstrate the model performance is a confusion matrix.

Confusion Matrix		Actual Values	
		Positive (Yes)	Negative (No)
Predicted Values	Positive (Yes)	True Positive (TP)	False Positive (FP)
	Negative (No)	False Negative (FN)	True Negative (TN)

Type I Error

Evaluation Metrics

Items	Formula	Meanings
Accuracy	$\frac{TP + TN}{TP + TN + FN + FP}$	Out of the prediction made by the model, what percentage is correct.
Precision	$\frac{TP}{TP + FP}$	Out of all the YES predications, how many of them were correct?
Recall (Sensitivity)	$\frac{TP}{TP + FN}$	How good was the model at predicting real YES events?
Recall (Specificity)	$\frac{TN}{TN + FP}$	How good was the model at predicting real NO events?
F1 score	$\frac{2}{\frac{1}{Recall} + \frac{1}{Precision}}$	Dealing with the imbalanced dataset, we have to use F1 score, indicating the harmonic mean of the precision and recall.

Source: <https://www.datacamp.com/blog/classification-machine-learning>

Evaluation Metrics

AUC-ROC Curve

AUC: area under curve

ROC: receiver operating characteristic curve

- AUC-ROC generates probability values instead of binary 0/1 values. It should be used when your dataset is roughly balanced.
- Using ROC for imbalanced datasets lead to incorrect interpretation.
- ROC curves provide a good overview of the trade-off between the TP rate and FP rate for binary classifiers using different probability thresholds.

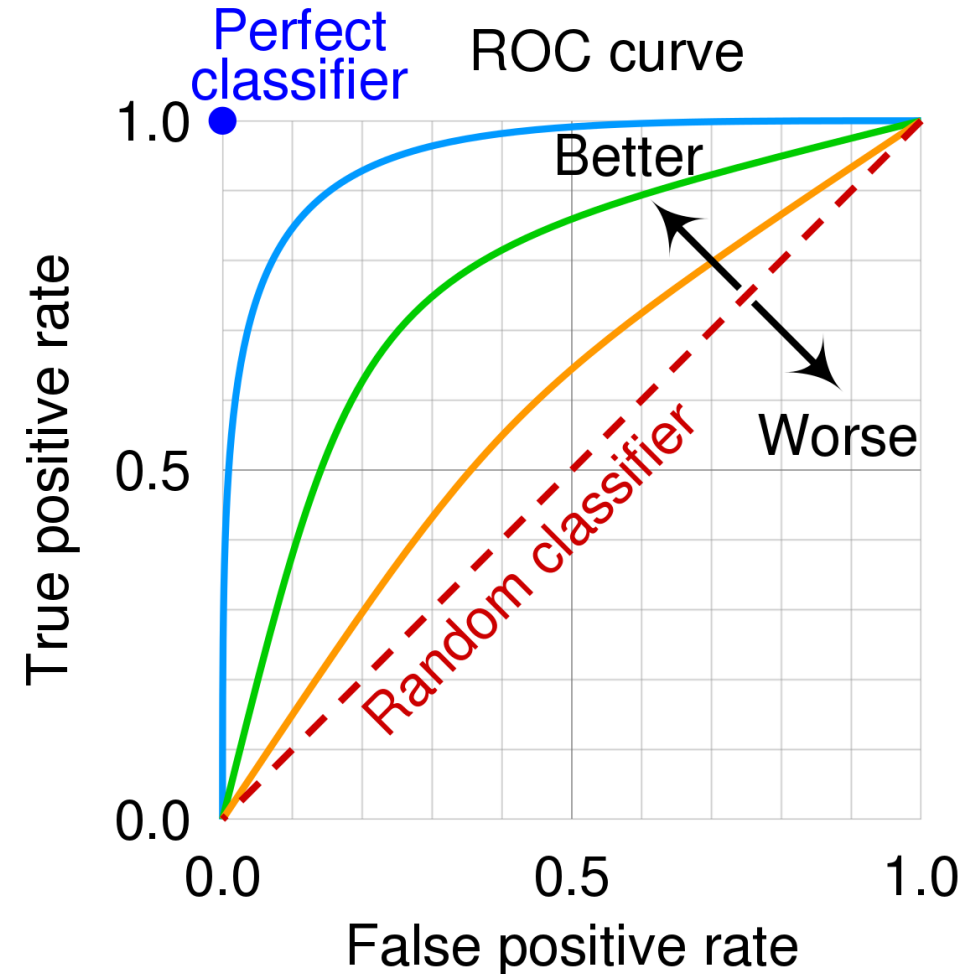
Evaluation Metrics

AUC-ROC Curve

AUC Value	Meanings
<0.5	Poor classifier
0.5~0.7	Random classifier
0.7~0.8	Good classifier
0.8~	Strong classifier
1	Perfect classifier

Source: <https://www.datacamp.com/blog/classification-machine-learning>

Oct. 24, 2024



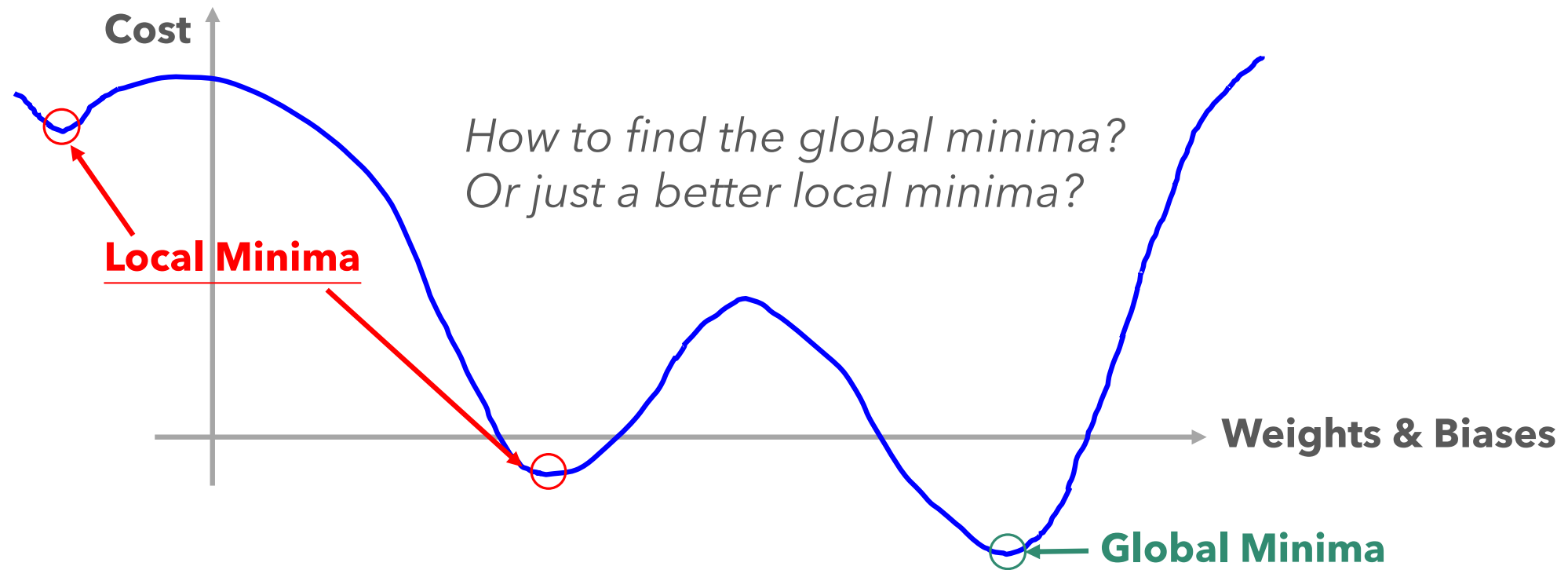
Source: <https://medium.com/@ilyurek/roc-curve-and-auc-evaluating-model-performance-c2178008b02>

Chun-Hsiang Chan (2024)

12

Model Optimization

- To reach the best model performance, it is important to optimize (tune) the hyperparameters of your model.

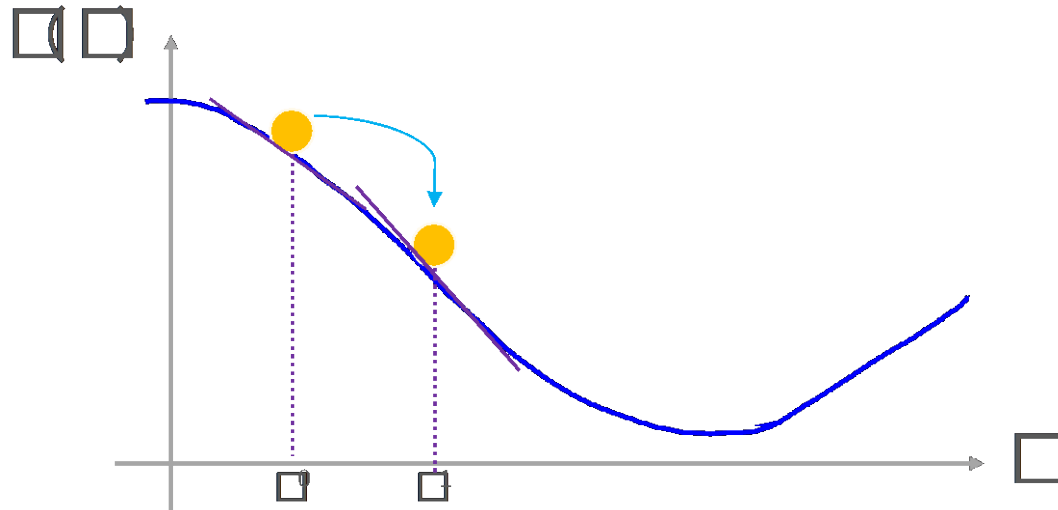


Model Optimization

- **Gradient Descent:**

Which way?

How does it work?



Assume that θ has only one variable...

- Randomly initiate at θ^0
- Compute $\frac{dC(\theta^0)}{d\theta} : \theta^1 \leftarrow \theta^0 - \eta \frac{\partial C(\theta^0)}{\partial \theta}$
- Compute $\frac{dC(\theta^1)}{d\theta} : \theta^2 \leftarrow \theta^1 - \eta \frac{\partial C(\theta^1)}{\partial \theta}$
- ...

$$\rightarrow \theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} C(\theta^i)$$

where η is the learning rate

Model Optimization

- For more information: plz see [here](#) [YouTube Link].

Deep Learning Gradient Descent

Assume that θ has only one variable...

- Randomly initiate at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$
- Compute the gradients of $C(\theta)$ at θ^0 :

$$\nabla_{\theta} C(\theta^0): \begin{bmatrix} \frac{\partial C(\theta^0)}{\partial \theta_1} \\ \frac{\partial C(\theta^0)}{\partial \theta_2} \end{bmatrix} \quad \theta^{i+1} \leftarrow \theta^i - \eta \nabla_{\theta} C(\theta^i)$$
- Update:

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial C(\theta^0)}{\partial \theta_1} \\ \frac{\partial C(\theta^0)}{\partial \theta_2} \end{bmatrix}$$
- Compute the gradients of $C(\theta)$ at θ^1 :

$$\nabla_{\theta} C(\theta^1): \begin{bmatrix} \frac{\partial C(\theta^1)}{\partial \theta_1} \\ \frac{\partial C(\theta^1)}{\partial \theta_2} \end{bmatrix}$$

Deep Learning Gradient Descent

Stochastic Gradient Descent (SGD)

Gradient Descent

$$\theta^{i+1} = \theta^i - \eta \nabla C(\theta^i), \nabla C(\theta^i) = \frac{1}{K} \sum_k \nabla C_k(\theta^i)$$

Stochastic Gradient Descent (SGD)

- Select a training sample x_k
- $\theta^{i+1} = \theta^i - \eta \nabla C_k(\theta^i)$
- If all training samples have same probability (uniform probability) to be selected

$$E[\nabla C_k(\theta^i)] = \frac{1}{K} \sum_k \nabla C_k(\theta^i)$$

The model updates after processing each selected training sample → much faster

Deep Learning Gradient Descent

Mini-Batch Stochastic Gradient Descent

Batch Gradient Descent	$\theta^{i+1} = \theta^i - \eta \frac{1}{K} \sum_k \nabla C_k(\theta^i)$	Use all K samples in each iteration
Stochastic Gradient Descent	$\theta^{i+1} = \theta^i - \eta \nabla C_k(\theta^i)$	Use 1 samples in each iteration
Mini-Batch SGD	$\theta^{i+1} = \theta^i - \eta \frac{1}{B} \sum_{x_k \in B} \nabla C_k(\theta^i)$	Use all B samples in each iteration

Chun-Hsiang Chan 13

Deep Learning Gradient Descent

Gradient Descent – Simple Case | Square Loss

- Update three parameters for $t - th$ iteration

$$w_1^{(t+1)} = w_1^{(t)} - \eta \frac{\partial C(\theta^{(t)})}{\partial w_1}$$

$$w_2^{(t+1)} = w_2^{(t)} - \eta \frac{\partial C(\theta^{(t)})}{\partial w_2}$$

$$b_1^{(t+1)} = b_1^{(t)} - \eta \frac{\partial C(\theta^{(t)})}{\partial b_1}$$

Square loss

$$C(\theta) = \sum_x \|\hat{y} - f(x; \theta)\|^2$$

$$\begin{bmatrix} w_1^{i+1} \\ w_2^{i+1} \\ b^{i+1} \end{bmatrix} \leftarrow \begin{bmatrix} w_1^i \\ w_2^i \\ b^i \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial C(\theta)}{\partial w_1} \\ \frac{\partial C(\theta)}{\partial w_2} \\ \frac{\partial C(\theta)}{\partial b} \end{bmatrix}$$

Chun-Hsiang Chan 23

Deep Learning Gradient Descent

Gradient Descent & Stochastic Gradient Descent

Gradient Descent

Stochastic Gradient Descent

Chun-Hsiang Chan 27

Deep Learning Gradient Descent

Learning Rate

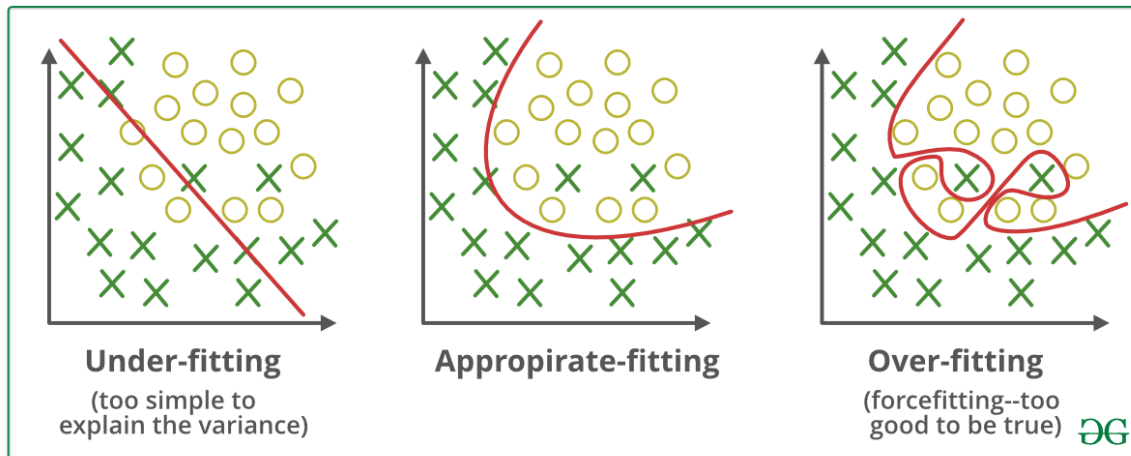
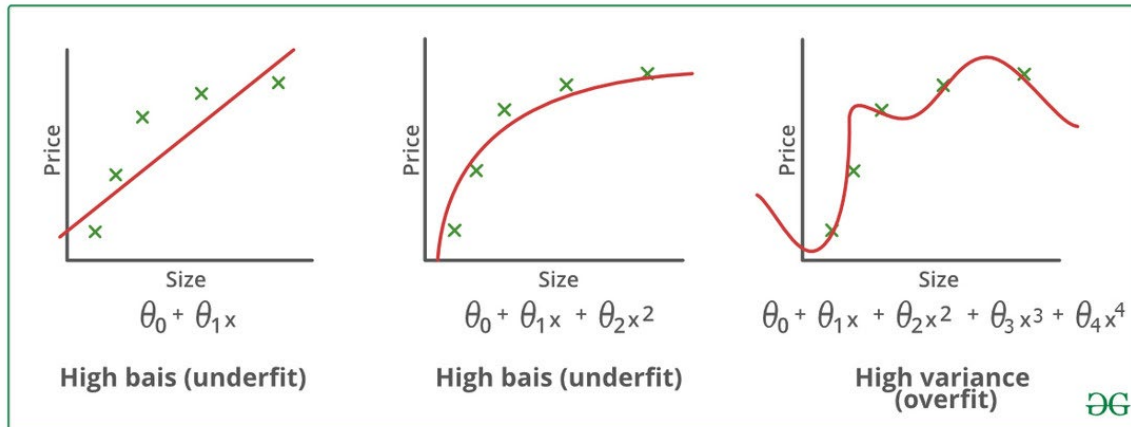
$$\theta^{i+1} = \theta^i - \eta \nabla C_k(\theta^i)$$

Single Model Standard LR Schedule

Snapshot Ensemble Cyclic LR Schedule

Left: Illustration of SGD optimization with a typical learning rate schedule. The model converges to a minimum at the end of training. **Right:** Illustration of Snapshot Ensembling. The model undergoes several learning rate annealing cycles, converging to and escaping from multiple local minima. We take a snapshot at each minimum for test-time ensembling.

Overfitting vs Underfitting



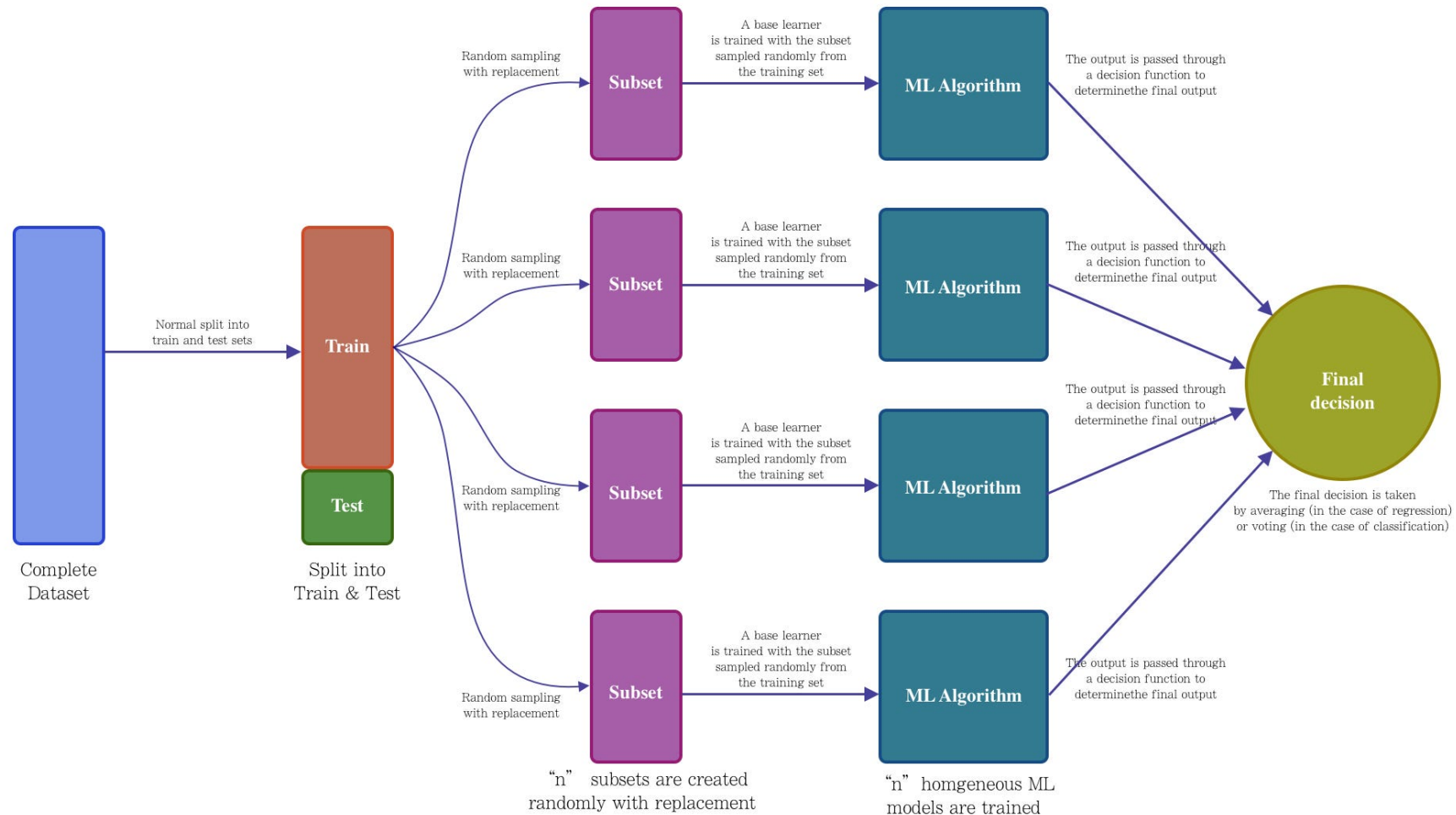
Techniques to reduce overfitting

1. Increase training data.
2. Reduce model complexity.
3. Early stopping during the training phase (keep an eye on the loss over the training period; as soon as loss increases, stop training).
4. Ridge Regularization and Lasso Regularization
5. Use dropout for neural networks to tackle overfitting.

Bagging (Bootstrap Aggregation)

- The two main components of the **bagging** technique are the *random sampling with replacement* (**bootstrapping**) and the *set of homogeneous* machine learning algorithms (**ensemble learning**).
- The **bagging** process is quite easy to understand. First, it extracts “ n ” subsets from the training set, and then these subsets are used to train “ n ” base learners of the same type.
- For making a prediction, each one of the “ n ” learners is fed with the test sample, and the output of each learner is averaged (in case of regression) or voted (majority vote).

Bagging (Bootstrap Aggregation)



Source: <https://towardsdatascience.com/ensemble-learning-bagging-boosting-3098079e5422>

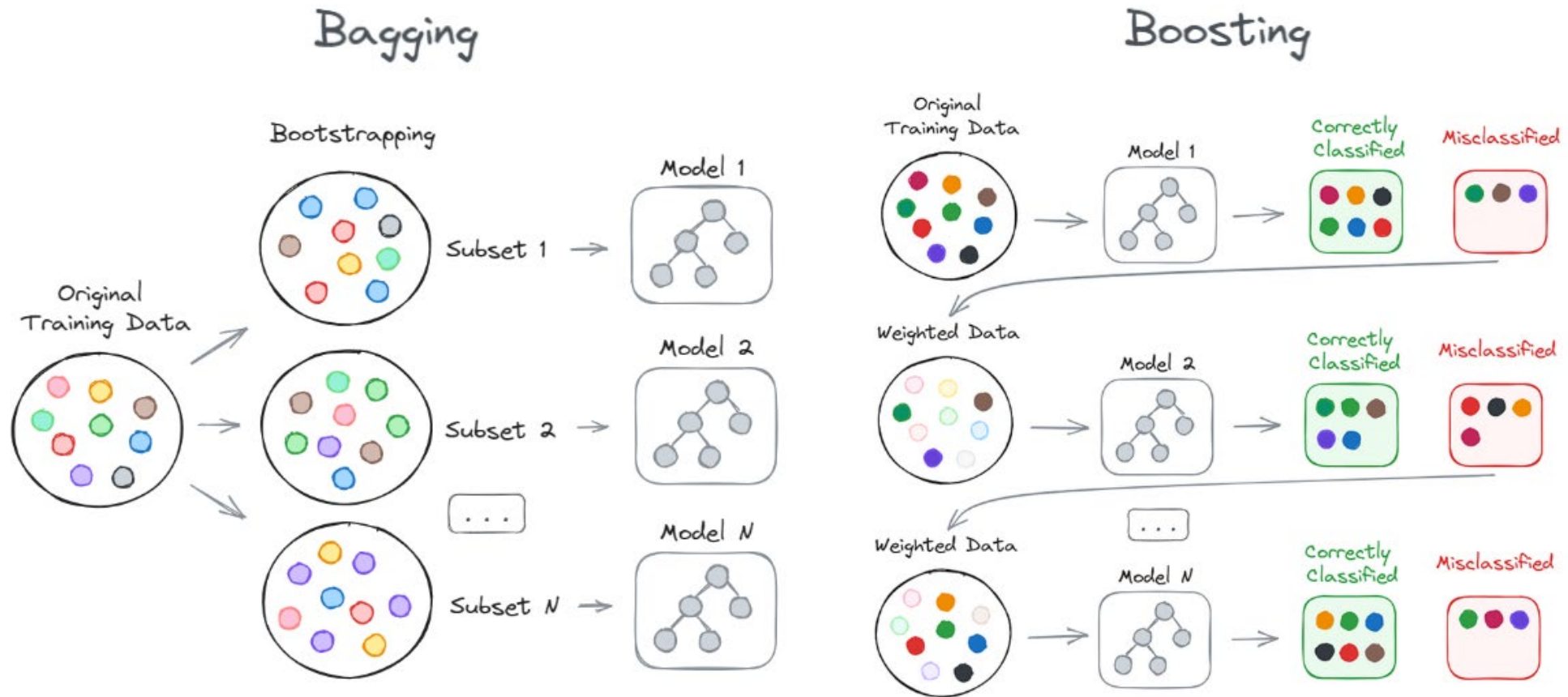
Boosting

- Boosting creates an ensemble model by combining several weak decision trees sequentially. It assigns weights to the output of individual trees. Then it gives incorrect classifications from the first decision tree a higher weight and input to the next tree. After numerous cycles, the boosting method combines these weak rules into a single powerful prediction rule.
- **Step 1:** The boosting algorithm assigns equal weight to each data sample. It feeds the data to the first machine model, called the base algorithm. The base algorithm makes predictions for each data sample.

Boosting

- **Step 2:** The boosting algorithm assesses model predictions and increases the weight of samples with a more significant error. It also assigns a weight based on model performance. A model that outputs excellent predictions will have a high amount of influence over the final decision.
- **Step 3:** The algorithm passes the weighted data to the next decision tree.
- **Step 4:** The algorithm repeats steps 2 and 3 until instances of training errors are below a certain threshold.

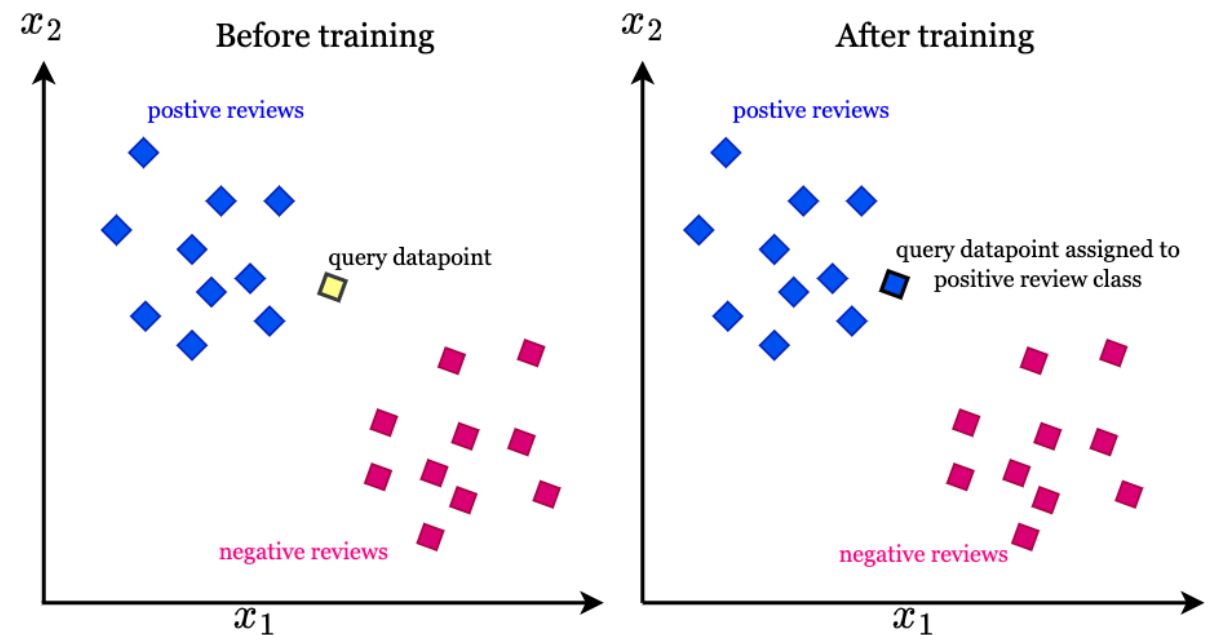
Bagging vs Boosting



Source: <https://pub.towardsai.net/bagging-vs-boosting-the-power-of-ensemble-methods-in-machine-learning-6404e33524e6>

k-Nearest Neighbors

- The k-nearest neighbors classify new data points to a particular category based on their similarity with the other data points in that category.
- Distance measurement:
 - Manhattan Distance
 - Minkowski Distance
 - Cosine Similarity
 - ...



Decision Tree

- The decision tree is the most basic machine learning classifier.
- Before introducing the decision tree concept, we first need to know "entropy" and "information gain", used to determine the tree node.
- (Information) Entropy:

$$Entropy = - \sum_{i=1}^n P(x_i) \log P(x_i)$$

- Information Gain:

$$Information\ gain = Entropy(S) - \sum_{i=1}^n w_i Entropy(i)$$

Decision Tree

Idx	Outlook	Temp	Humidity	Windy	Play
0	Sunny	Hot	High	0	0
1	Sunny	Hot	High	1	0
2	Overcast	Hot	High	0	1
3	Rainy	Mild	High	0	1
4	Rainy	Cool	Normal	0	1
5	Rainy	Cool	Normal	1	0
6	Overcast	Cool	Normal	1	1
7	Sunny	Mild	High	0	0
8	Sunny	Cool	Normal	0	1
9	Rainy	Mild	Normal	0	1
10	Sunny	Mild	Normal	1	1
11	Overcast	Mild	High	1	1
12	Overcast	Hot	Normal	0	1
13	Rainy	Mild	High	1	0

$$E(S) = 0.94$$

$$E(\text{outlook} = \text{sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$E(\text{outlook} = \text{overcast}) = -\frac{4}{4} \log_2 \frac{4}{4} = 0$$

$$E(\text{outlook} = \text{rainy}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

$$w_{\text{sunny}} = \frac{5}{14}; w_{\text{overcast}} = \frac{4}{14}; w_{\text{rainy}} = \frac{5}{14}$$

$$\text{Info Gain}(\text{outlook})$$

$$= 0.94 - \left[\frac{5}{14} * E(\text{sunny}) + \frac{4}{14} E(\text{overcast}) + \frac{5}{14} E(\text{rainy}) \right] = 0.247$$

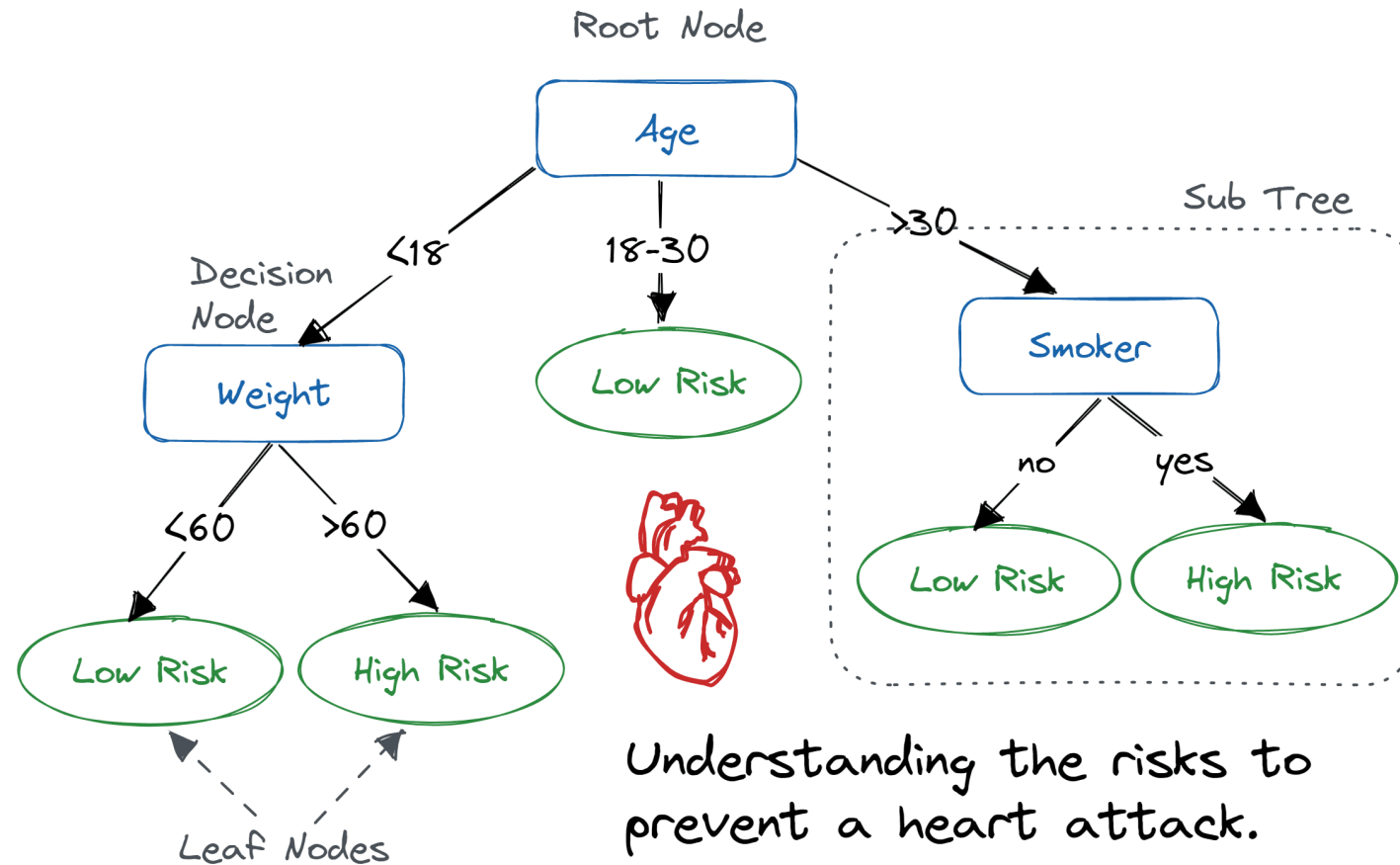
$$\text{Info Gain}(\text{temp}) = 0.029$$

$$\text{Info Gain}(\text{humidity}) = 0.152$$

$$\text{Info Gain}(\text{windy}) = 0.048$$

Select the feature with the highest information gain.

Decision Tree



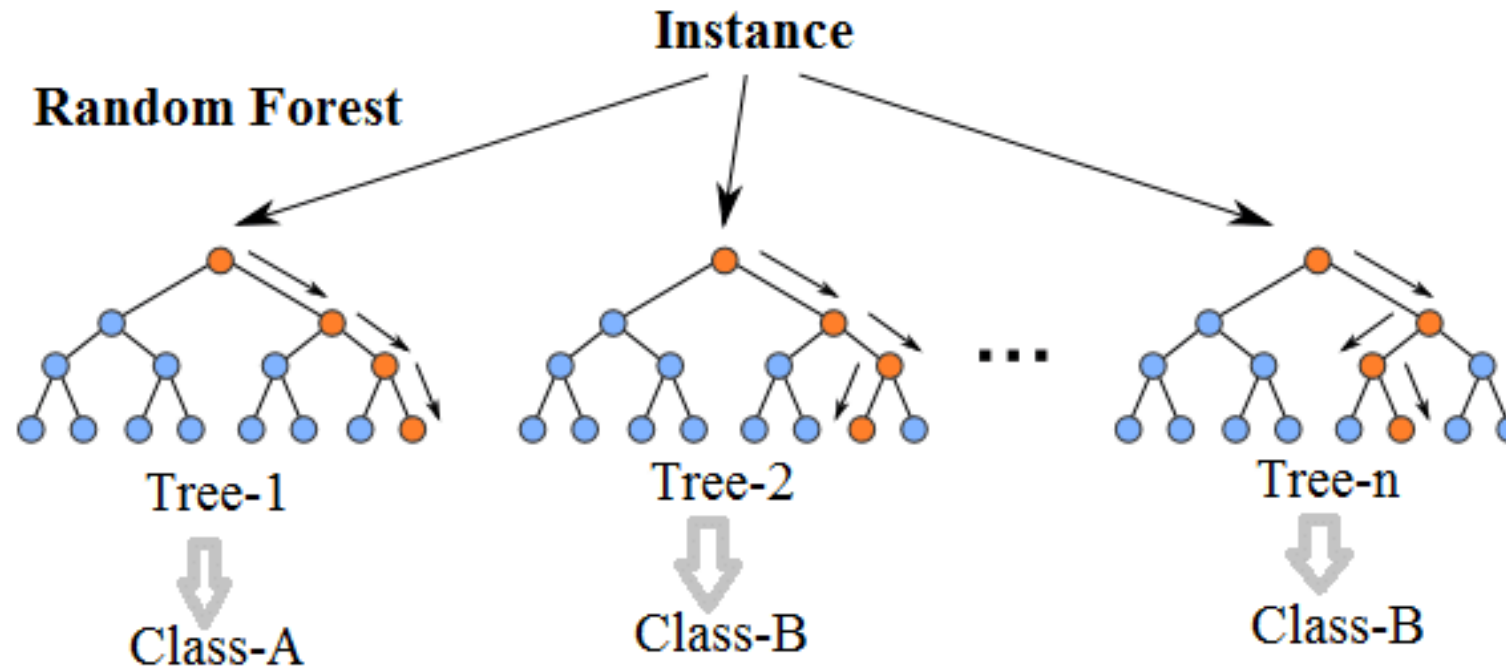
Random Forest

Random Forest = Bagging + Decision Tree

1. Randomly select n samples for a sample set (bagging).
 2. Building a decision tree depends on n samples; for each node, randomly select d features for splitting with information gain.
 3. Repeat steps 1 and 2 for k times.
- Integrate all prediction results from all trees to determine the classification result with the majority vote.

Random Forest

Random Forest Simplified



Source: <https://medium.com/chung-yi/ml%E5%85%A5%E9%96%80-%E5%8D%81%E4%B8%83-%E9%9A%A8%E6%A9%9F%E6%A3%AE%E6%9E%97-random-forest-6afc24871857>

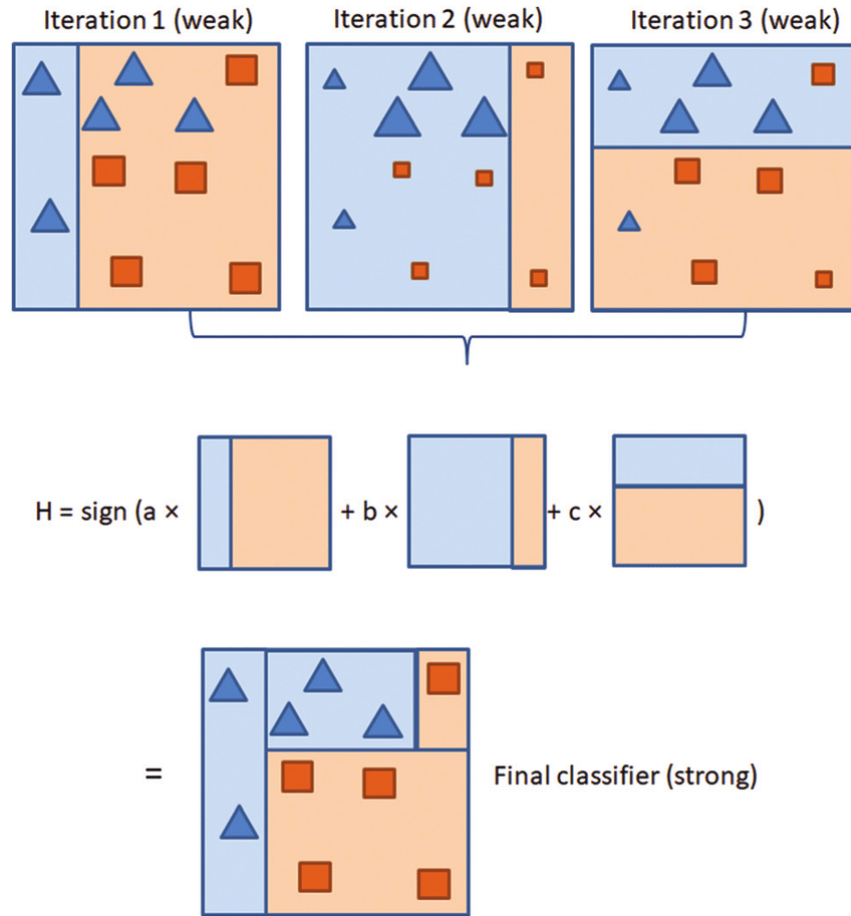
AdaBoost

- Adaptive Boosting (AdaBoost) was one of the earliest boosting models developed. It adapts and tries to self-correct in every iteration of the boosting process.
- AdaBoost initially gives the same weight to each dataset. Then, it automatically adjusts the weights of the data points after every decision tree. It gives more weight to incorrectly classified items to correct them for the next round. It repeats the process until the residual error, or the difference between actual and predicted values, falls below an acceptable threshold.

AdaBoost

- Assign every observation, x_i , an initial weight value, $w_i = \frac{1}{n}$, where n is the total number of observations.
- Train a “weak” model (most often a decision tree).
- For each observation:
 - If predicted incorrectly, w_i is increased
 - If predicted correctly, w_i is decreased
- Train a new weak model where observations with greater weights are given more priority.
- Repeat step 3 and 4 until observations perfectly predicted are present number of trees are trained.

AdaBoost



Input :

- A training set $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$.

Initialization :

- Maximum number of iterations T ;
- initialize the weight distribution $\forall i \in \{1, \dots, m\}, D^{(1)}(i) = \frac{1}{m}$.

for $t = 1, \dots, T$ **do**

- Learn a classifier $f_t : \mathbb{R}^d \rightarrow \{-1, +1\}$ using distribution $D^{(t)}$
- Set $\epsilon_t = \sum_{i: f_t(\mathbf{x}_i) \neq y_i} D^{(t)}(i)$
- Choose $a_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$
- Update the weight distribution over examples

$$\forall i \in \{1, \dots, m\}, D^{(t+1)}(i) = \frac{D^{(t)}(i)e^{-a_t y_i f_t(\mathbf{x}_i)}}{Z^{(t)}}$$

where $Z^{(t)} = \sum_{i=1}^m D^{(t)}(i)e^{-a_t y_i f_t(\mathbf{x}_i)}$ is a normalization factor such that $D^{(t+1)}$ remains a distribution.

Output : The voted classifier $\forall \mathbf{x}, F(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T a_t f_t(\mathbf{x}) \right)$

Source:

https://www.researchgate.net/publication/351719735_Exploring_AdaBoost_and_Random_Forests_machine_learning_approaches_for_infrared_pathology_on_unbalanced_data_sets/figures?lo=1&utm_source=google&utm_medium=organic

Oct. 24, 2024

Source: <https://ama.liglab.fr/~amini/AdaBoost/>

Chun-Hsiang Chan (2024)

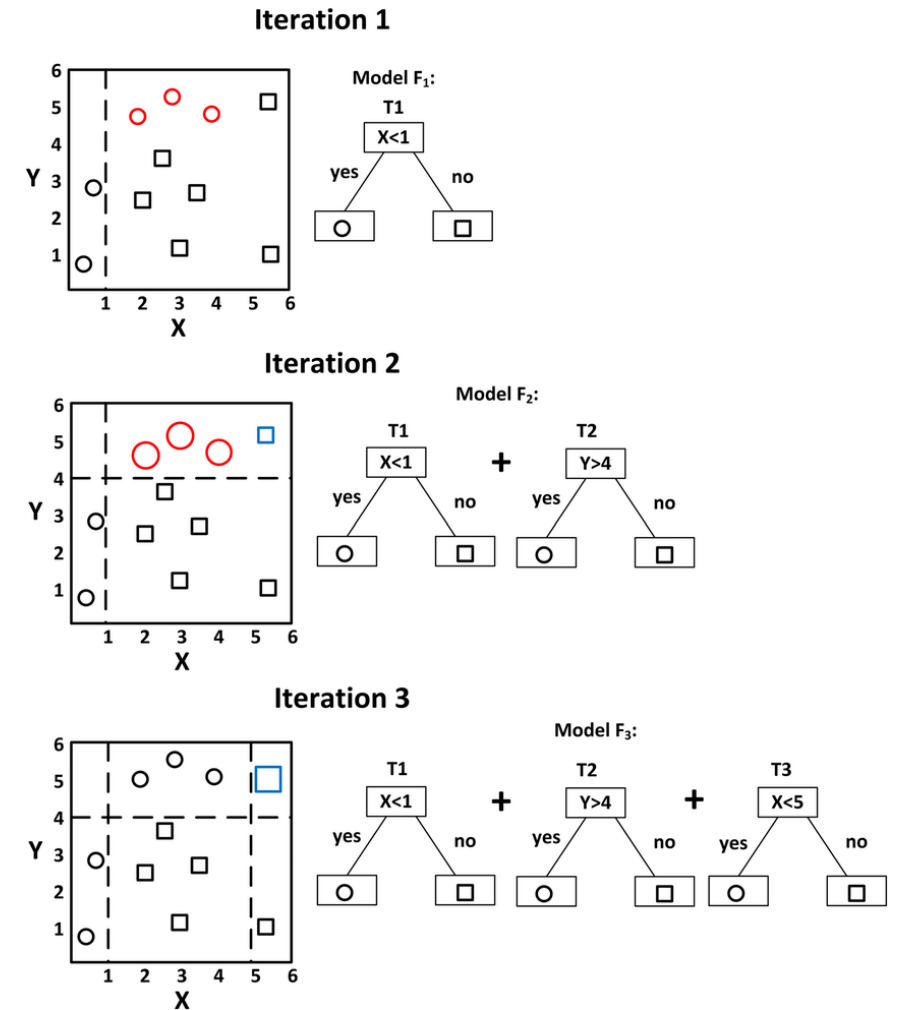
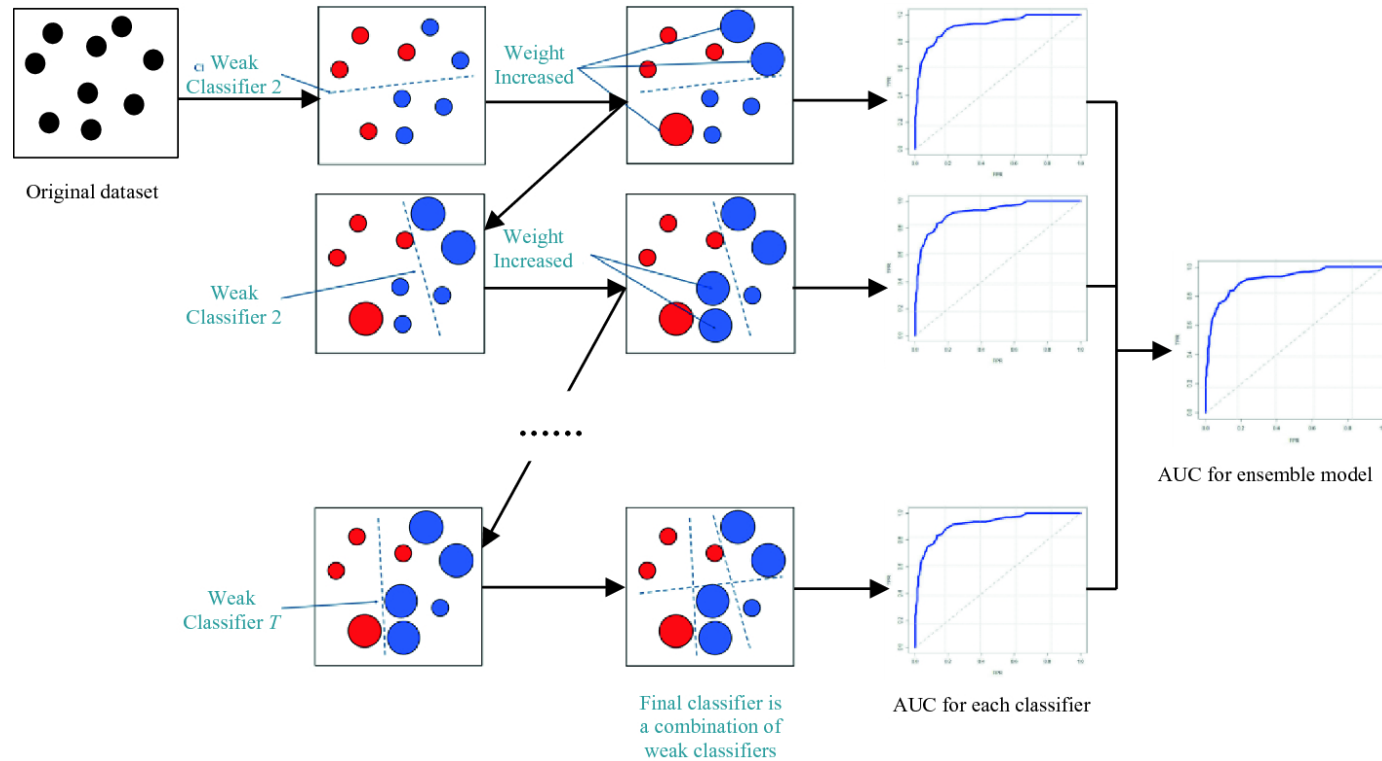
30

Gradient Boosting

- Gradient Boosting (GB) is similar to AdaBoost in that it, too, is a sequential training technique. The difference between AdaBoost and GB is that GB does not give incorrectly classified items more weight.
- Instead, GB software optimizes the loss function by generating base learners sequentially so that the present base learner is always more effective than the previous one.
- This method attempts to generate accurate results initially instead of correcting errors throughout the process, like AdaBoost.
- For this reason, GB software can lead to more accurate results. Gradient Boosting can help with both classification and regression-based problems.

Source: <https://aws.amazon.com/what-is/boosting/>

Gradient Boosting



Source: <https://www.google.com/url?sa=i&url=https%3A%2F%2Fdatascience.eu%2Fmachine-learning%2Fgradient-boosting-what-you-need-to-know%2F&psig=AOvVaw3QgdPcdywxRhcVgSQt13qn&ust=1701741598913000&source=images&cd=vfe&opi=89978449&ved=0CBQQjhXqFwoTCODL3s3X9IIDFQAAAAAdAAAAABAS>

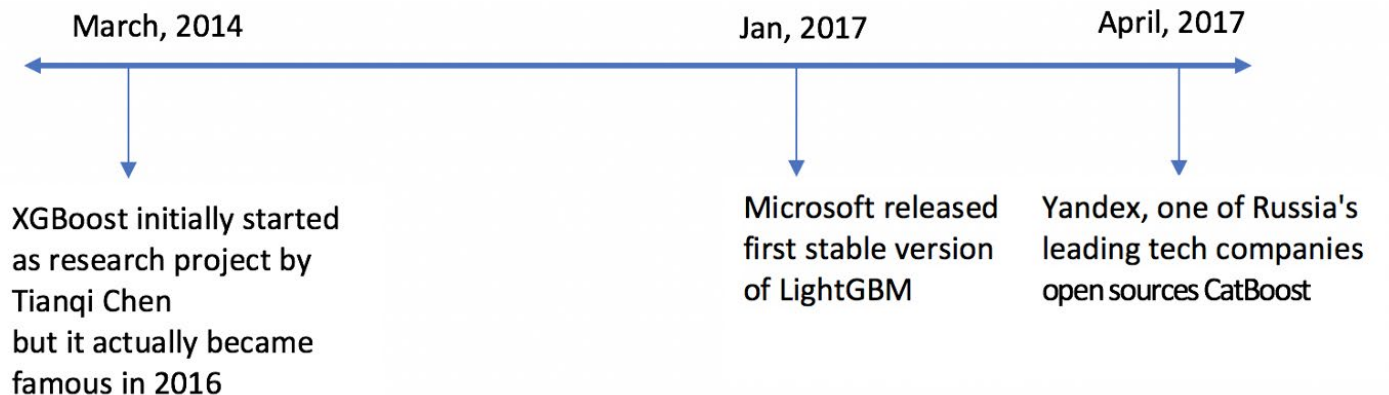
Source: https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.researchgate.net%2Ffigure%2FA-simple-example-of-visualizing-gradient-boosting_fig5_326379229&psig=AOvVaw3QgdPcdywxRhcVgSQt13qn&ust=1701741598913000&source=images&cd=vfe&opi=89978449&ved=0CBQQjhXqFwoTCODL3s3X9IIDFQAAAAAdAAAAABAJ

Extreme Gradient Boosting (XGBoost)

- Extreme Gradient Boosting (XGBoost) improves gradient boosting for computational speed and scale in several ways.
- XGBoost uses multiple cores on the CPU so that learning can occur in parallel during training. It is a boosting algorithm that can handle extensive datasets, making it attractive for big data applications.
- The key features of XGBoost are parallelization, distributed computing, cache optimization, and out-of-core processing.

LightGBM

- LightGBM - Light Gradient Boosting Machine
- LightGBM is a decision-tree based algorithm equipped with rapid, distributed, and highly efficient characteristics, developed by MS in April, 2017.



- LightGBM = Gradient Boosting Decision Tree(GBDT) + Gradient-based One-Side Sampling (GOSS) + Exclusive Feature Bundling (EFB)

LightGBM

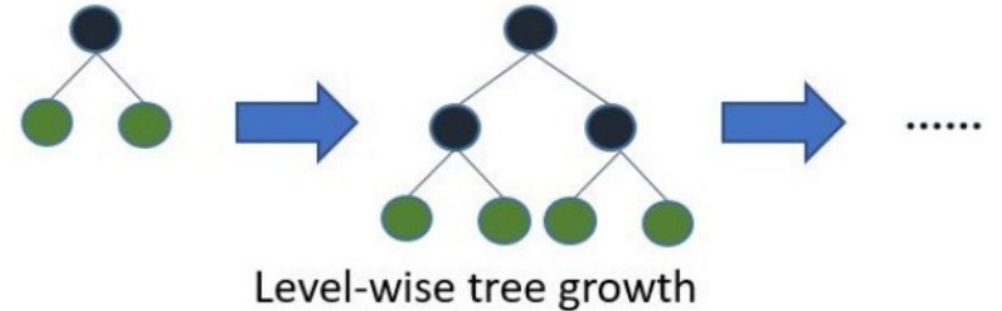
Pros:

- Fast training speed
- Low memory usage
- Support parallel computing
- Suitable for big data
- Support GPU computing
- High accuracy

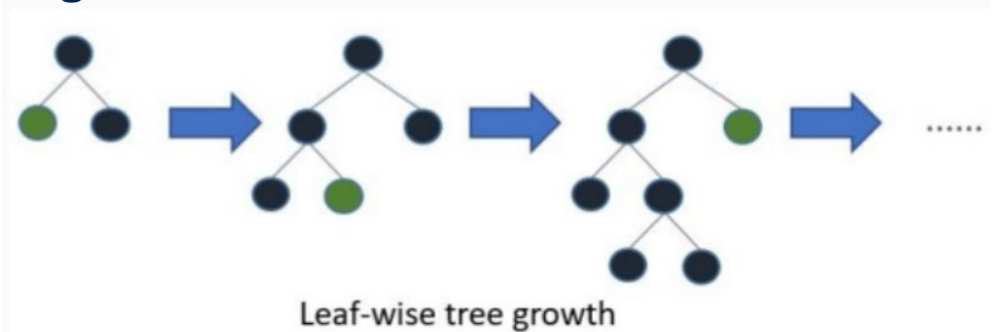
Cons:

- Not suitable for small dataset (prone to overfitting)

XGBoost



LightGBM



CatBoost

CatBoost is widely adopted in Yandex, which is the largest Russian search engine.

- It has several special functionalities
 - Automatically process categorical features
 - Automatically process missing value (minimal value)
 - Able to deal with various datasets, e.g., audio, image, and text
 - Low manual hyperparameter adjustment → low overfitting probability
 - High accuracy



Naïve Bayes

- It is a classification technique based on Bayes' Theorem with an independence assumption among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

The diagram shows the equation $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ with handwritten annotations. The term $P(A|B)$ is highlighted in a light green box and labeled "POSTERIOR" with the note "the probability of 'A' being TRUE given that 'B' is TRUE". The term $P(B|A)$ is highlighted in a light orange box and labeled "LIKELIHOOD" with the note "the probability of 'B' being TRUE given that 'A' is TRUE". The term $P(A)$ is highlighted in a light blue box and labeled "PRIOR" with the note "the probability of 'A' being TRUE". The term $P(B)$ is highlighted in a light pink box and labeled "The probability of 'B' being TRUE".

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

© luminousmen.com

Source: https://www.google.com/url?sa=i&url=https%3A%2F%2Fheartbeat.comet.ml%2Funderstanding-the-mathematics-behind-naive-bayes-ab6ee85f50d0&psig=AOvVaw1FrZcrDX7LPw6F_TsGyxns&ust=1701751731448000&source=images&cd=vfe&opi=89978449&ved=0CBiQjRxqFwoTCKDToa399IIFQAAAAAdAAAAABAS
Source: <https://www.analyticsvidhya.com/blog/2017/09/naive-bayes-explained/>

Naïve Bayes

- **Convert the data set into a frequency table**In this first step data set is converted into a frequency table
- **Create Likelihood table by finding the probabilities**Create Likelihood table by finding the probabilities like Overcast probability = 0.29 and probability of playing is 0.64.
- **Use Naive Bayesian equation to calculate the posterior probability**Now, use Naive Bayesian equation to calculate the posterior probability for each class. The class with the highest posterior probability is the outcome of the prediction.

Naïve Bayes

- Maximize Posterior Probability

$$y^* = \operatorname{argmax}_{y_i \in \mathcal{Y}} P(y_i | x_1, x_2, \dots, x_n)$$

- That is

$$y^* = \operatorname{argmax}_{y_i \in \mathcal{Y}} \frac{P(x_1, x_2, \dots, x_n | y_i) P(y_i)}{P(x_1, x_2, \dots, x_n)} = \operatorname{argmax}_{y_i \in \mathcal{Y}} P(x_1, x_2, \dots, x_n | y_i) P(y_i)$$

- Assume all features are independent, then

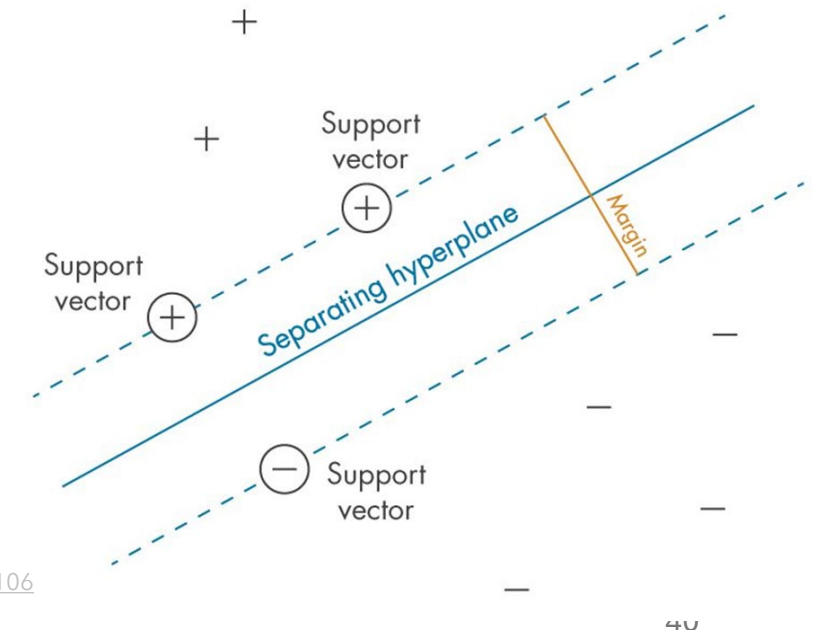
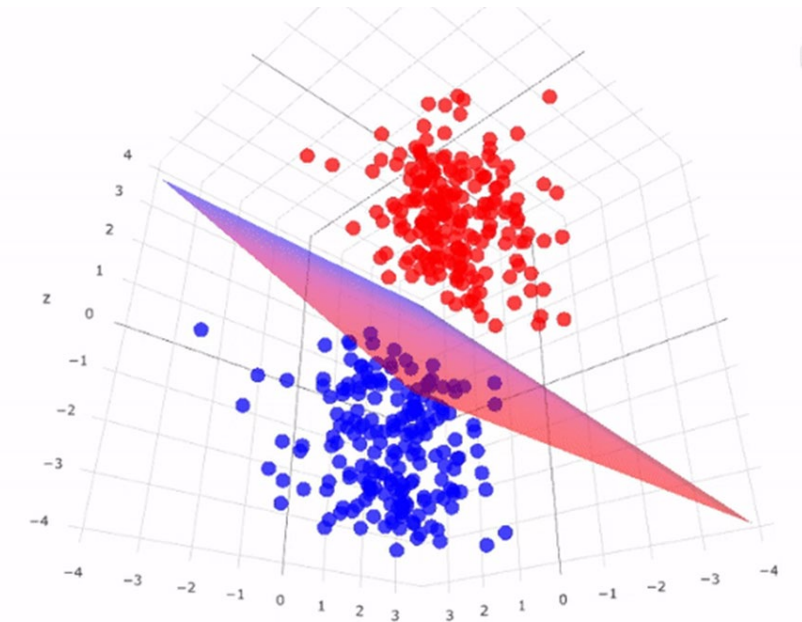
$$P(x_1, x_2, \dots, x_n | y_i) = \prod_j^n P(x_j | y_i)$$

- Therefore, posterior probability maximization could be

$$y^* = \operatorname{argmax}_{y_i \in \mathcal{Y}} P(y_i) \prod_j^n P(x_j | y_i)$$

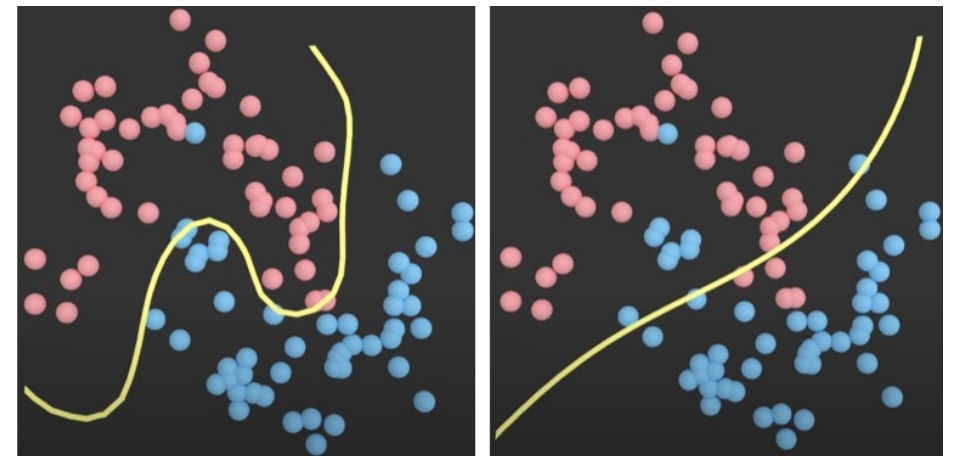
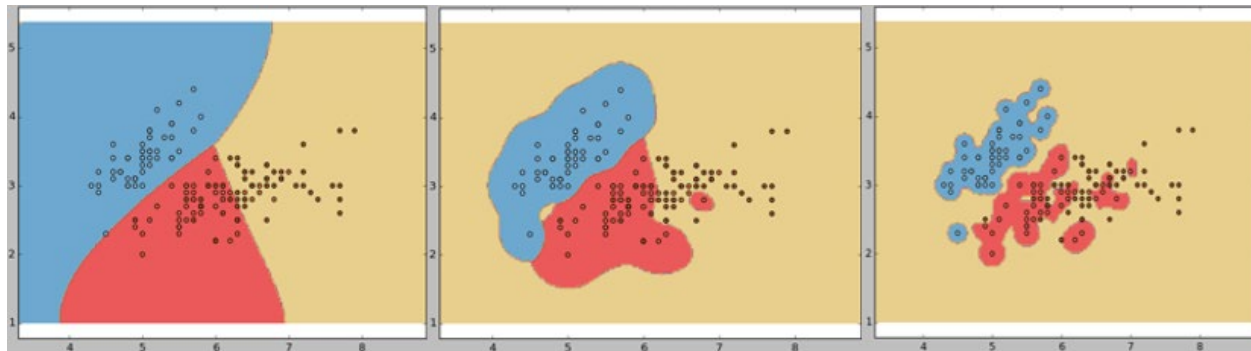
Linear SVM

- Support Vector Machines (SVMs) are a type of supervised machine learning algorithm used for classification and regression tasks. They are widely used in various fields, including pattern recognition, image analysis, and natural language processing.
- SVMs work by finding the optimal hyperplane that separates data points into different classes.



Linear SVM

- A hyperplane is a decision boundary that separates data points into different classes in a high-dimensional space. In two-dimensional space, a hyperplane is simply a line that separates the data points into two classes. In three-dimensional space, a hyperplane is a plane that separates the data points into two classes. Similarly, in N -dimensional space, a hyperplane has $(N-1)$ -dimensions.



Remarks

- Every machine learning classifier has its advantage in dealing with specific problems; therefore, it is difficult to say which model is the best one.
- In general, we will enroll all classifiers into the analysis and select the best one.
- But in fact, some of them usually perform very well in most cases: XGBoost and Random Forest.



The End

Thank you for your attention!

Email: chchan@ntnu.edu.tw

Web: toodou.github.io

