Outline

- 1. Briefly explain random forest.
- 2. How to implement random forest in Weka.



Why use random forest? A reflection of the decision tree algorithm

- 1. Decision tree is very sensitive to small variations in the training data, which makes the model very unstable.
- 2. Very likely to overfit the training data
 - "Variance is an error resulting from sensitivity to small fluctuations in the dataset used for training. High variance will cause an algorithm to model irrelevant data, or noise, in the dataset instead of the intended outputs, called signal. This problem is called **overfitting**. An overfitted model will perform well in training, but won't be able to distinguish the noise from the signal in an actual test."



Random Forest

Random Forest is a powerful and versatile supervised machine learning algorithm that grows and combines multiple decision trees to create a "forest."

Steps of Random Forest Algorithm:

Step 1: Bootstrapping

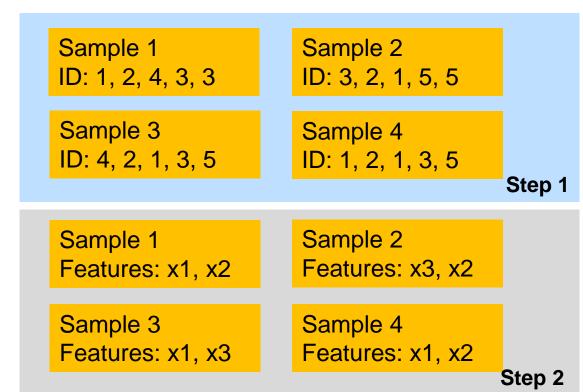
Step 2: Feature selection

Step 3: Construction of trees

Step 4: Voting and aggregation



ID	x1	x2	х3	class
1	3	3	5	1
2	4	6	3	0
3	0	6	0	1
4	1	5	6	1
5	3	2	9	0



Step 4: Voting and aggregation

