

Random Forests

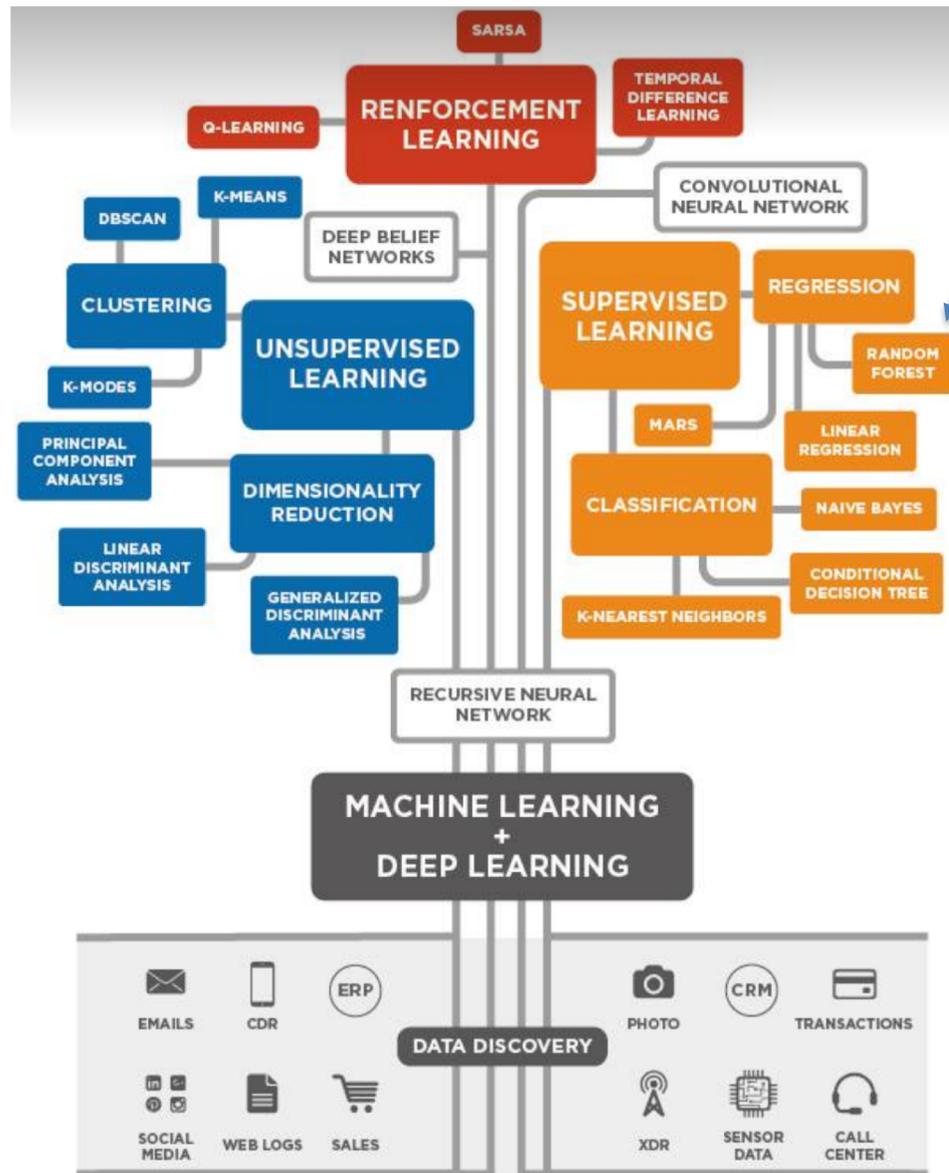
Sargur Srihari

srihari@buffalo.edu

Topics in Random Forests

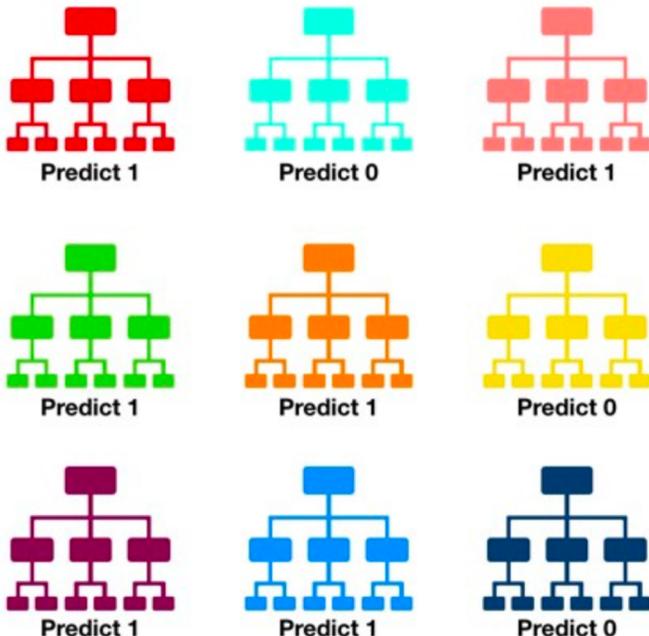
1. Random Forests in ML
2. Two types of randomness
3. Why do random forests succeed?

Random Forests in ML



Combining decision trees

- A random forest is a large number of decision trees operating as an ensemble
- Each tree outputs a class prediction
 - Class with most votes is model's prediction



With six ones and three zeroes, the prediction is 1

Regression with Random Forests

- One way to reduce the variance of an estimate is to average together many estimates
 - E.g., we can train M different trees on different random subsets of data, with replacement, and then compute the ensemble

$$f(\mathbf{x}) = \sum_{m=1}^M \frac{1}{M} f_m(\mathbf{x})$$

- where f_m is the m^{th} tree
- This technique is called bagging
 - Stands for *bootstrap aggregation*

Role of Randomness

- Unfortunately, simply re-running the same learning algorithm on different subsets of data can result in highly correlated predictors
 - Which limits the amount of variance reduction possible
- Random forests tries to decorrelate the base learners
 - By learning trees based on a randomly chosen subset of input variables, as well as a randomly chosen subset of data cases

Preconditions for Randomness

1. There needs to be some features so that models built using those features do better than random guessing
2. Predictions (and therefore the errors) made by the individual trees need to have low correlations with each other

Two types of Randomness

1. Bagging (Bootstrap Aggregation)

- Decisions trees are sensitive to training data
 - Small changes in training set can result in different trees

2. Feature Randomness

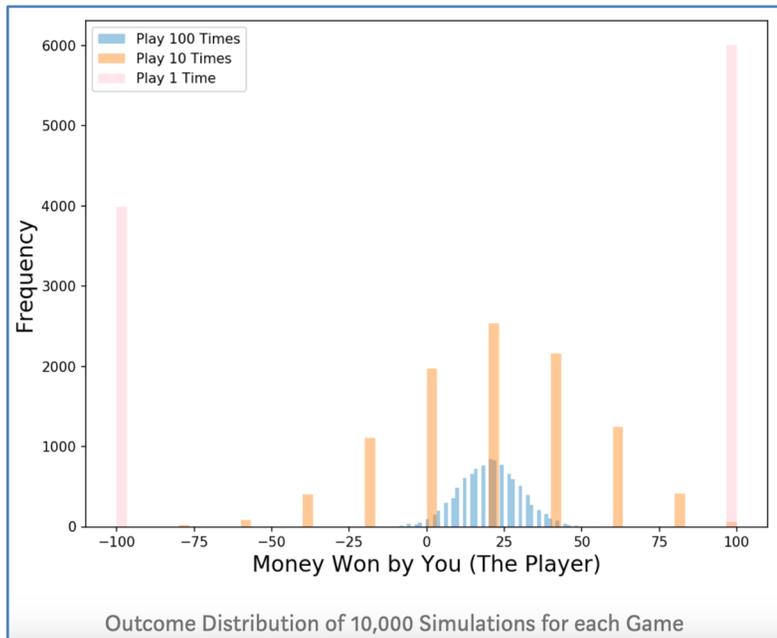
- In a normal decision tree, to split a node, we consider every feature and pick one that produces the most separation between the observations in the left node vs. those in the right node.
 - In contrast, each tree in a random forest can pick only from a random subset of features

Why does this succeed?

Game: Random no generator produces a number 1-100

You win if number greater than 40 , i.e, probability 60%. You lose otherwise

- Game 1 —play 100 times, betting \$1 each time
- Game 2— play 10 times, betting \$10 each time
- Game 3— play one time, betting \$100
- Expected Value Game 1 = $(0.60*1 + 0.40*-1)*100 = 20$
- Expected Value Game 2= $(0.60*10 + 0.40*-10)*10 = 20$
- Expected Value Game 3= $0.60*100 + 0.40*-100 = 20$



Out of 10,000 simulations
probability of making money

