

random forests & gradient boosting

Overview

- A decision tree is learnt by splitting a subset of data (that maps to a current leaf) by the most useful features to split across.
- A random forest is an ensemble of decision trees → ensemble to reduce variance
- each decision tree in the forest considers a random subset of features when forming questions and only has access to a random set of the training data points.
- in case of regression, we average the outputs of each tree in the ensemble
- in case of classification, we let the trees vote and take a majority
- Bagging vs boosting
 - **Bagging** is a simple ensembling technique in which we build many *independent* predictors/models/learners and combine them using some model averaging techniques. (e.g. weighted average, majority vote or normal average)
 - Example of bagging ensemble is **Random Forest models**.
 - **Boosting** is an ensemble technique in which the predictors are not made independently, but sequentially.
 - This technique employs the logic in which the subsequent predictors learn from the mistakes of the previous predictors.
 - Gradient boosting is an example of a boosting ensemble vs random forests which is a bagging ensemble.
 - Because we keep training on the mistakes in the boosting model, we need to be wary of overfitting and have a good stopping criterion
- **Boosting**
 - Focus new learners on examples that others get wrong
 - Train learners sequentially
 - Errors of early predictions indicate the “hard” examples
 - Focus later predictions on getting these examples right
 - Combine the whole set in the end
 - Convert many “weak” learners into a complex predictor
- math for gradient boosting
 - pick the next leaf to split based on which feature gives a good split → sequentially picking the worst performing leaf next

- we use an additive strategy: fix what we have learned, and add one new tree at a time.
- for an overview of gradient boosting, see [xgboost intro](#)

Random forests

- the tree is built by choosing splitting on which feature gives us the best information gain. if the feature is a set of continuous valued functions, sample at random.



- One measure for information gain is entropy minimization: $-\sum p_i \log(1/p_i)$ where p_i is the probability (by counting) of a class i . Note that we get better information gain (lower entropy) if the split is something like $[3/4, 1/8, 1/8]$ rather than $[2/3, 2/3, 2/3]$.
- Another could be $\sum_i p_i^2$ over all classes in that split.
- using something like the above measures, we choose a split
- how to build a random tree? lets say we have a dataset of 5 inputs with 3 features in each input. pick a random set of features \rightarrow lets say two features, use that to generate some candidate split points (here we use the projection onto axes as candidates). pick the one that gives the most information gain

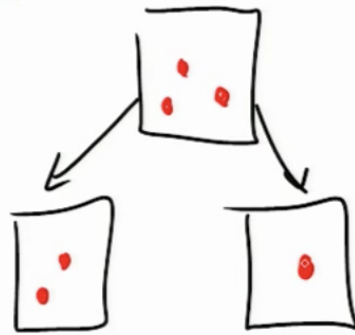
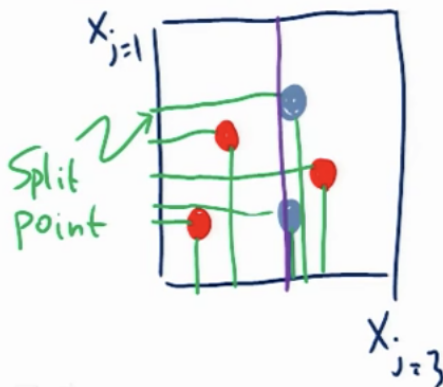
$d=3$ features
 $n=5$ data

Building a random tree

$$X = \begin{matrix} & \begin{matrix} i=1 & i=2 & i=3 & i=4 & i=5 \end{matrix} \\ \begin{matrix} j=1 \\ j=2 \\ j=3 \end{matrix} & \begin{bmatrix} \textcircled{1} & 3 & 0 & 8 & 5 \\ 0 & 6 & 2 & 9 & 5 \\ \textcircled{2} & 1 & 4 & 0 & 1 \end{bmatrix} \end{matrix}$$

$$Y = \begin{matrix} & \begin{matrix} i=1 & i=2 & i=3 & i=4 & i=5 \end{matrix} \\ \begin{matrix} j=1 \\ j=2 \\ j=3 \end{matrix} & \begin{bmatrix} 0 & 1 & 0 & 1 & 0 \end{bmatrix} \end{matrix}$$

Pick 2 features at random.



- why pick at random? we have tons of features and tons of split points. lets pick at random and build many trees (like a 1000 trees). over time we hope to pick up good split points.

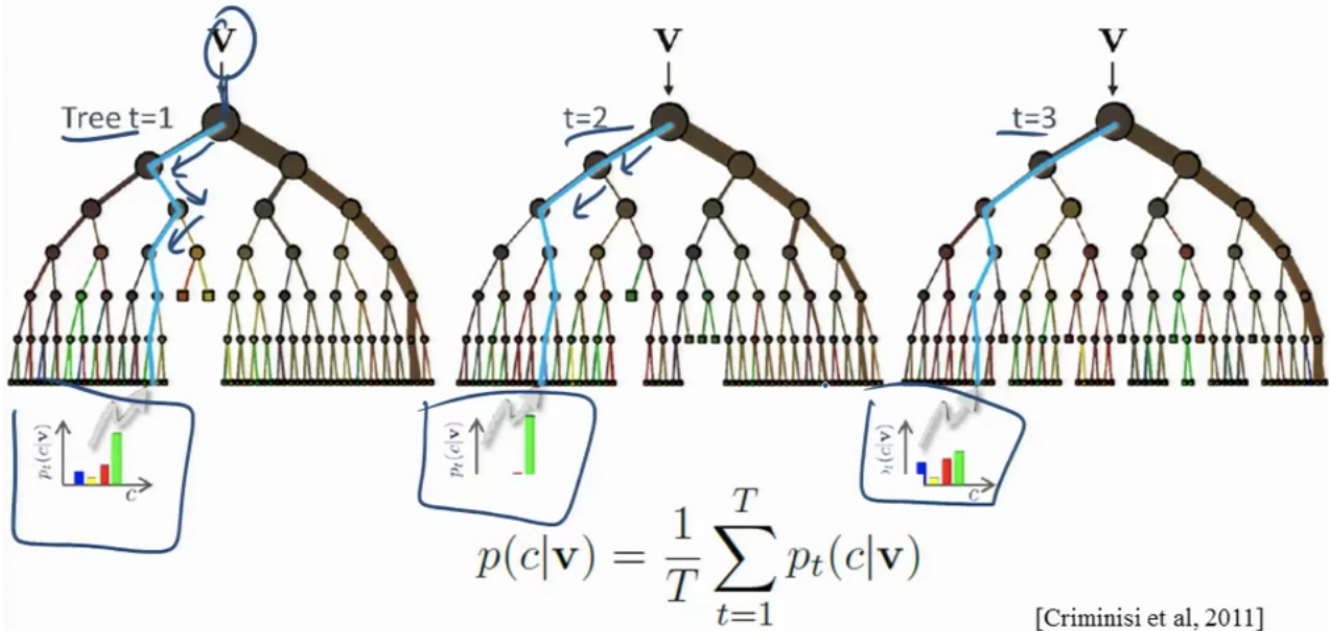
Random Forests algorithm

- For $b = 1$ to B :
 - Draw a bootstrap sample \mathbf{Z}^* of size N from the training data.
 - Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.
 - Select m variables at random from the p variables.
 - Pick the best variable/split-point among the m .
 - Split the node into two daughter nodes.
- Output the ensemble of trees $\{T_b\}_1^B$.

- note that the algorithm has randomness both in the data chosen to train each tree but also randomness in the feature picked as split point.
- because we build many independent trees, this is trivially parallelizable and if we have a cluster of machines we can take advantage of it.
- **Bootstrapping** is the idea of sampling data at random and training a set of classifiers with the thus sampled subsets of data. we use bootstrapping here (see algorithm above)
- each of these trees has very high variance. like really high
- for a ML algorithm to benefit from bagging, they need to be uncorrelated errors.

Building a forest (ensemble)

In a forest with T trees we have $t \in \{1, \dots, T\}$. All trees are trained independently (and possibly in parallel). During testing, each test point \mathbf{v} is simultaneously pushed through all trees (starting at the root) until it reaches the corresponding leaves.



Text classification example

In news categorization, a possible term is *Bill Clinton*. A corresponding **weak learner (node)** is: If the term *Bill Clinton* appears in the document predict that the document belongs to *News*.

