

Decision Tree Variations

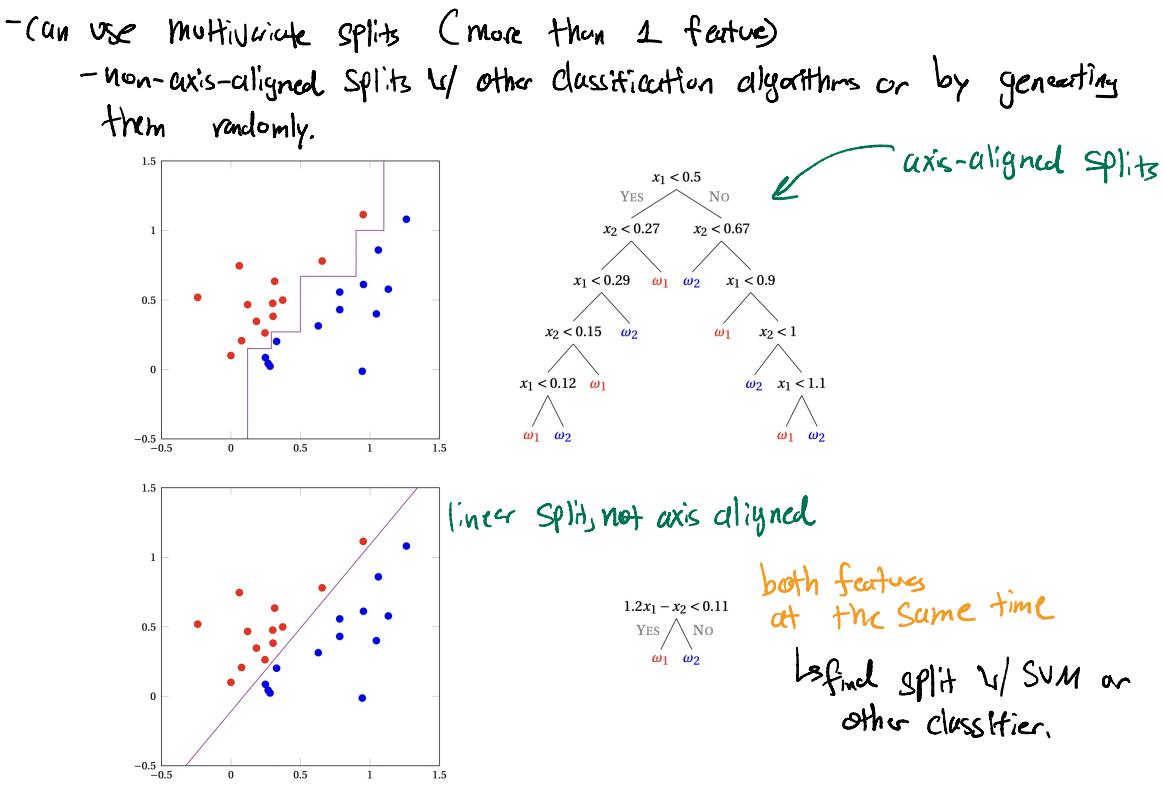


Figure 15.1: Comparison between a vanilla decision tree and a decision tree with a multivariate split

decision tree w/ linear classifier \rightarrow nonlinear decision boundaries

nonlinear classifier \rightarrow piece-wise linear boundaries

May gain better classifier at cost of:

- worse interpretability (e.g. using all features at node)
- speed

Usually, decision trees very fast

- looking at less features as you go

most decision trees in practice
only use 1 feature at a time
 \hookrightarrow assumed on final exam

Not so common

Solution: at each node, look at at most $2/3$ features

- limit # features per split: forward stepwise selection or Lasso

Decision Tree Regression:

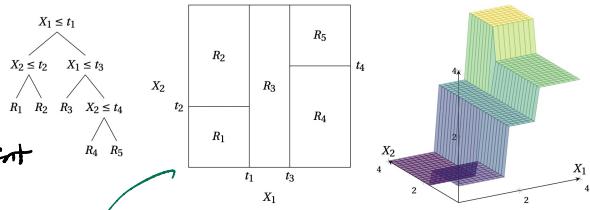
- creates a decision fn. that's piecewise constant
- each leaf stores a label $\mu_s = \frac{1}{|S|} \sum_{i \in S} Y_i$

average of all labels in index set

$$\text{Cost: } J(S) = \frac{1}{|S|} \sum_{i \in S} (Y_i - \mu_s)^2$$

↳ Variance of all labels of pts. in node
 $\text{Var}(\{Y_i : i \in S\})$

Figure 15.2: Example of regression using a decision tree



constant within each rectangle

* entropy for classification, not so good
 w/ regression (continuously-valued labels)
 choose split that minimizes this

If all training pts are identical, cost = 0, positive or/

Stopping Early

- don't always have to keep going until every leaf is pure (i.e. only contains 1 class)

Why Stop Early?

- limit tree depth (for speed)
- limit tree size (for big datasets)
- complete tree more prone to overfitting (trying to classify every single pt.)
 - data has noise or isn't separable/overlapping distribution
 - purity of leaves is counterproductive in this case (better to estimate posteriors)

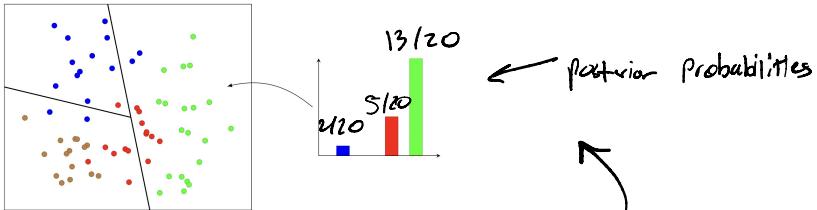


Figure 15.3: Example of a decision tree that has been stopped early, with the corresponding posterior probability distribution for one leaf

- leaves w/ multiple points, return:

- {classification}
 - majority vote
 - posterior probabilities for each class

- an average (regression)

When do we stop decision trees?

common fear split to have no info. gain, but next split makes much more progress.

- When split doesn't reduce entropy/error enough (this is dangerous, pruning better)
- When most of node's points (e.g. 95%) have the same class (prevent overfitting)
- When nodes contain too few pts. (e.g. ≤ 10)
- When cell's edges are too tiny
- When depth is too large (risky if still lots of pts. in a cell)
- use validation to compare errors
 - ↳ slowest but most effective way to know when to stop
 - ↳ to avoid overfitting - grow tree large
 - use validation to prune tree back.

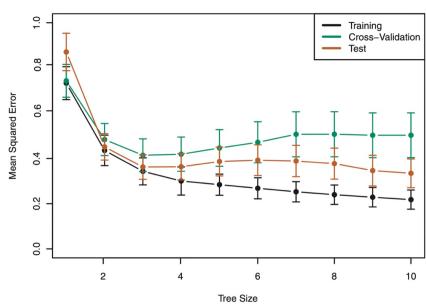
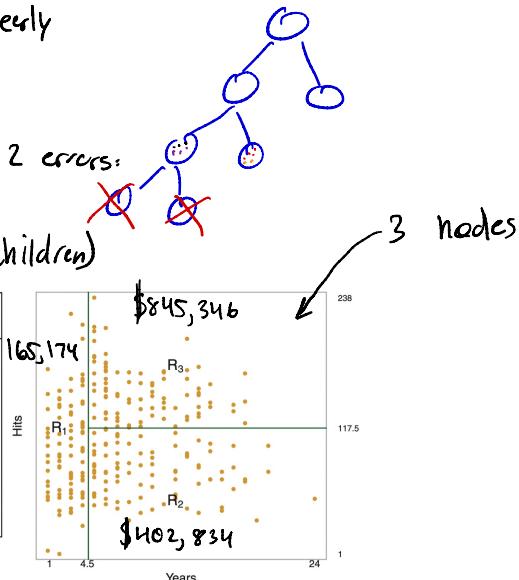
Pruning: grow tree too large then greedily remove each split whose removal improves validation

- empirically more reliable than stopping early

- don't prune node w/ no children
 ↳ work bottom-up

- when comparing validation errors, compare 2 errors:

- errors of two children
- error of parent (if we prune children)



* decision trees are fast, simple, interpretable, easy to explain, invariant under Scaling/Translation, robust to _{local features}

- but they're also not the best at prediction (compared to prev. methods)
- also high Variance (depending on split)

Ensemble Learning: combine diff. algorithms to compute collective answer
idea: reduce Variance by taking the avg. of a bunch of decision trees

Weak learning algo: does better than guessing

→ take avg. output of:

- different algorithms
- same learning algorithm on different training sets (more common)

most common {
- **bagging:** same learning algo. on random subsets of one training set (if low data)
- **random forest:** random decision trees on random subsamples of data

* averaging isn't specific to decision trees but it works particularly well for trees

for regression algos: median or mean output

for classification: majority vote or average posterior probabilities

- use learners w/ low bias (e.g. deep decision trees)

- averaging reduces variance, not necessarily bias

- high variance & some overfitting okay since averaging reduces variance

- averaging sometimes reduces bias & increases flexibility, but not reliably

e.g. averaging nonlinear decision boundaries from linear classifiers

- hyperparameter settings usually different than those for 1 learner

- # trees is a hyperparameter

↳ Variance vs. runtime tradeoff

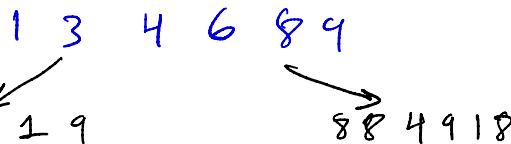
Bagging: randomized Method for creating different learners from same dataset

- works well w/ many diff. learning algorithms (except k-nearest neighbors)

given n-point training sample, generate random subsample of size n' by Sampling with replacement

↳ Some samples chosen multiple times

↳ Some won't be chosen at all



- If $n' = n$, ~63.2% are chosen

Build learner. Points chosen j times have greater weight in learning algorithm

- decision trees: pt. has j x weight in entropy

- SVMs: pt. incurs $j \times$ penalty if it violates the margin (slack var.)

- Regression: pt. incurs $j \times$ loss

repeat until T learners made

Metalearning: takes test pt., feeds it into all T classifiers, returns average/majority output of learners

Random Forests: Random Sampling w/ random classifiers

- random sampling isn't random enough

- one really strong predictor \rightarrow same feature split at top of every tree

- even more randomness \rightarrow different splits, especially at root node

Idea: - choose random sample M of d features

- choose best split from M features

- * different random sample of M features at each tree node

$M \approx \sqrt{d}$ good for classification M is a hyperparam

$M \approx d/3$ good for regression

Smaller M \rightarrow more randomness in trees, less tree correlation, more bias

Why does this work well?

- if 1 really strong predictor, only a fraction of trees can choose that predictor (M/d)

↳ decorrelates trees, taking avg \rightarrow less variance than single tree

* averaging works best when we have very strong learners that are also diverse
↳ be careful not to dumb down trees too much in order to decorrelate

disadvantage of random forest:- Slow

- loses interpretability/inference

advantage: - more accurate than a single tree