

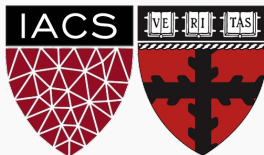
Lecture #14: Decision Trees

Data Science 1

CS 109A, STAT 121A, AC 209A, E-109A

Pavlos Protopapas Kevin Rader

Margo Levine Rahul Dave



Lecture Outline

Announcements

Motivation

Decision Trees

Splitting Criteria

Stopping Conditions & Pruning

Announcements

Monday 10/30: **Intro to algorithmic accountability and ethics in data science**

Rachel Schutt, IACS

Wednesday 11/1: **Storytelling with Data: Finding the Narrative in the Numbers**

Maggie Mulvihill, BU

The fusion of data and storytelling has created an entire new job of data storyteller today. Boston University Computational Journalism Professor Maggie Mulvihill explains how investigative and data journalists work with data, how they uncover hard-to-find, high-value datasets and best practices for presenting stories that move the needle for an audience?

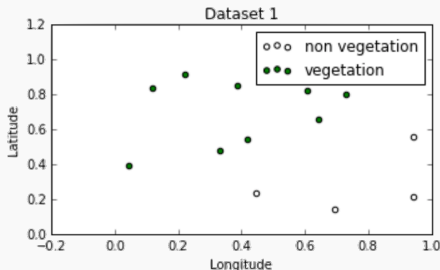
Motivation

Geometry of Data

Recall that **logistic regression, LDA/QDA** for classification works best when the classes are well-separated in the feature space by a decision boundary defined by some equation

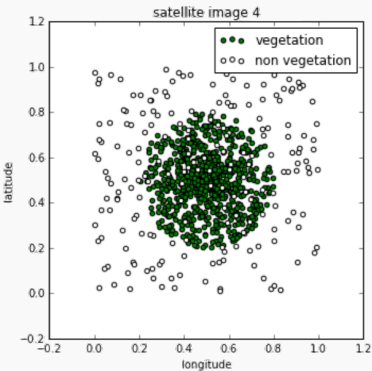
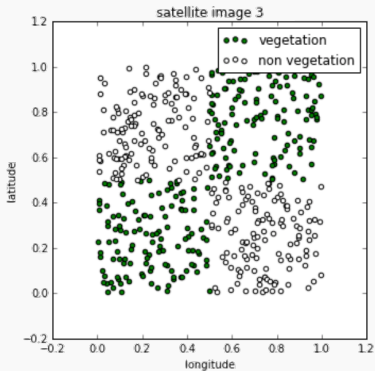
$$f(x_1, \dots, x_J) = 0$$

The following is a typical dataset for logistic regression with a linear boundary:



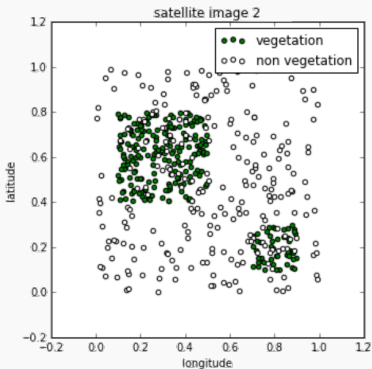
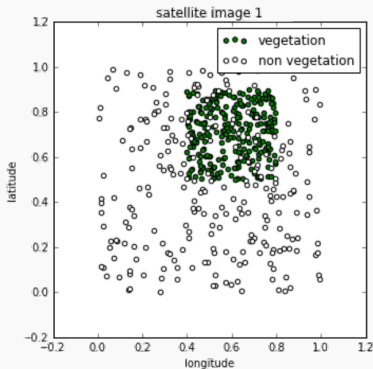
Geometry of Data

Discuss the suitability of the following datasets for logistic regression and LDA/QDA:



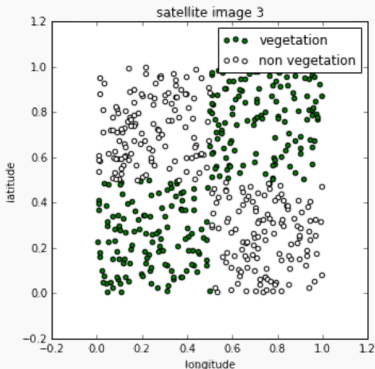
Geometry of Data

Discuss the suitability of the following datasets for logistic regression and LDA/QDA:



Geometry of Data

Notice that in all of the datasets the classes are still well-separated in the feature space, but ***the decision boundaries cannot be described by single equations:***



Interpretable Models

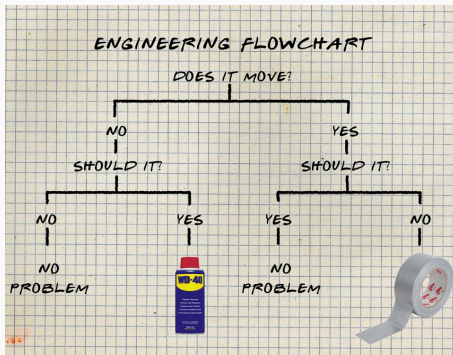
While logistic regression models with linear boundaries are intuitive to interpret by examining the impact of each predictor on the log-odds of a positive classification, it is less straightforward to interpret nonlinear decision boundaries in context:

$$(x_3 + 2x_2)^2 - x_1 + 10 = 0$$

It would be desirable to build models with complex decision boundaries that are also easy to interpret.

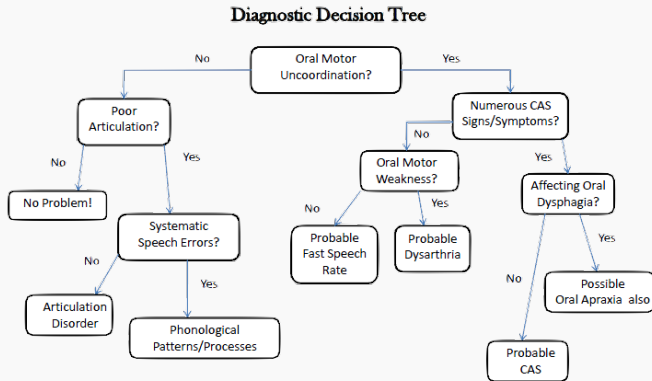
Interpretable Models

But people in every walk of life have long been using interpretable models for differentiating between classes of objects and phenomena:



Interpretable Models

But people in every walk of life have long been using interpretable models for differentiating between classes of objects and phenomena:



It turns out that the simple flow charts in our examples can be formulated as mathematical models for classification and these models have the properties we desire; they are:

1. interpretable by humans
2. have sufficiently complex decision boundaries
3. the decision boundaries are locally linear, each component of the decision boundary is simple to describe mathematically.

Decision Trees

The Geometry of Flow Charts

Flow charts whose graph is a tree (connected and no cycles) represents a model called a **decision tree**.

Formally, a **decision tree model** is one in which the final outcome of the model is based on a series of comparisons of the values of predictors against threshold values.

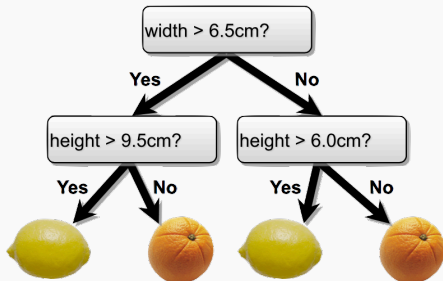
In a graphical representation (flow chart),

- ▶ the internal nodes of the tree represent attribute testing
- ▶ branching in the next level is determined by attribute value
- ▶ leaf nodes represent class assignments

The Geometry of Flow Charts

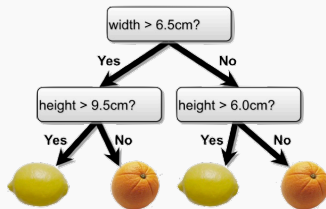
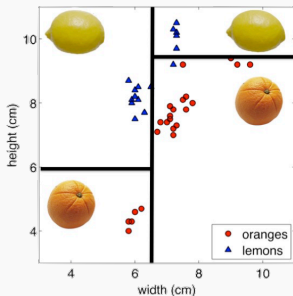
Flow charts whose graph is a tree (connected and no cycles) represents a model called a **decision tree**.

Formally, a **decision tree model** is one in which the final outcome of the model is based on a series of comparisons of the values of predictors against threshold values.



The Geometry of Flow Charts

Every flow chart tree corresponds to a partition of the feature space by lines aligned to the axis or (hyper) planes. Conversely, every such partition can be written as a flow chart tree.



Each comparison and branching represents splitting a region in the feature space. Typically, at each iteration, we split once along one dimension (one predictor).

Learning the Model

Given a training set, learning a decision tree model for binary classification means to produce an 'optimal' partition of the feature space with axis aligned linear boundaries, wherein each region is given a class label based on the largest class of the training points in that region.

Learning the Model

Learning the smallest ‘optimal’ decision tree for any given set of data is NP complete for numerous simple definitions of ‘optimal’. Instead, we will seek a reasonably model using a greedy algorithm.

1. Start with an empty decision tree (undivided feature space)
2. Choose the ‘optimal’ predictor on which to split and choose the ‘optimal’ threshold value for splitting.
3. Recurse on on each new node until **stopping condition** is met

Now, we need only define our splitting criterion and stopping condition.

Numerical vs Categorical Attributes

Note that the compare and branch method by which we defined regression tree works well for numerical features.

However, if a feature is categorical (with more than two possible values), comparisons like *feature < threshold* does not make sense.

A simple solution is to encode the values of a categorical feature using numbers and treat this feature like a numerical variable.

This is indeed what some computational libraries (e.g. `sklearn`) do, however, this method has drawbacks.

Numerical vs Categorical Attributes

Example

Suppose the feature we want to split on is **color**, and the values are: Red, Blue and Yellow. If we encode the categories numerically as:

$$\text{Red} = 0, \text{Blue} = 1, \text{Yellow} = 2$$

Then the possible non-trivial splits on **color** are

$$\{\{\text{Red}\}, \{\text{Blue}, \text{Yellow}\}\}, \quad \{\{\text{Red}, \text{Blue}\}, \{\text{Yellow}\}\}$$

But if we encode the categories numerically as:

$$\text{Red} = 2, \text{Blue} = 0, \text{Yellow} = 1$$

The possible splits are

$$\{\{\text{Blue}\}, \{\text{Yellow}, \text{Red}\}\}, \quad \{\{\text{Blue}, \text{Yellow}\}, \{\text{Red}\}\}$$

Depending on the encoding, the splits we can optimize over can be different!

Numerical vs Categorical Attributes

In practice, the effect of our choice of naive encoding of categorical variables are often negligible - models resulting from different choices of encoding will perform comparably.

In cases where you might worry about encoding, there is a more sophisticated way to numerically encode the values of categorical variables so that one can optimize over all possible partitions of the values of the variable.

This more principled encoding scheme is computationally more expensive but is implemented in a number of computational libraries (e.g. R's `randomForest`, H2O packages, etc).

Splitting Criteria

Optimality of Splitting

While there is no ‘correct’ way to define an optimal split, there are some common sensical guidelines for every splitting criterion:

- ▶ the regions in the feature space should grow progressively more pure with the number of splits. That is, we should see each region ‘specialize’ towards a single class.
- ▶ the fitness metric of a split should take a differentiable form (making optimization possible)
- ▶ we shouldn’t end up with empty regions - regions containing no training points.

Classification Error

Suppose we have J number of predictors and K classes.

Suppose we select the j -th predictor and split a region containing N number of training points along the threshold $t_j \in \mathbb{R}$.

We can assess the quality of this split by measuring the classification error made by each newly created region, R_1, R_2 :

$$\text{Error}(i|j, t_j) = 1 - \max_k p(k|R_i)$$

where $p(k|R_i)$ is the proportion of training points in R_i that are labeled class k .

Example

	Class 1	Class 2	Error($i j, t_j$)
R_1	0	6	$1 - \max\{6/6, 0/6\} = 0$
R_2	5	8	$1 - \max\{5/13, 8/13\} = 5/13$

We can now try to find the predictor j and the threshold t_j that minimizes the average classification error over the two regions, weighted by the population of the regions:

$$\min_{j, t_j} \left\{ \frac{N_1}{N} \text{Error}(1|j, t_j) + \frac{N_2}{N} \text{Error}(2|j, t_j) \right\}$$

where N_i is the number of training points inside region R_i .

Suppose we have J number of predictors, N number of training points and K classes.

Suppose we select the j -th predictor and split a region containing N number of training points along the threshold $t_j \in \mathbb{R}$.

We can assess the quality of this split by measuring the purity of each newly created region, R_1, R_2 . This metric is called the **Gini Index**:

$$\text{Gini}(i|j, t_j) = 1 - \sum_k p(k|R_i)^2$$

Question: What is the effect of squaring the proportions of each class? What is the effect of summing the squared proportions of classes within each region?

Example

	Class 1	Class 2	Gini($i j, t_j$)
R_1	0	6	$1 - (6/6^2 + 0/6^2) = 0$
R_2	5	8	$1 - [(5/13)^2 + (8/13)^2] = 80/169$

We can now try to find the predictor j and the threshold t_j that minimizes the average Gini Index over the two regions, weighted by the population of the regions:

$$\min_{j, t_j} \left\{ \frac{N_1}{N} \text{Gini}(1|j, t_j) + \frac{N_2}{N} \text{Gini}(2|j, t_j) \right\}$$

where N_i is the number of training points inside region R_i .

Information Theory

The last metric for evaluating the quality of a split is motivated by metrics of uncertainty in information theory.

Ideally, our decision tree should split the feature space into regions such that each region represents a single class. In practice, the training points in each region is distributed over multiple classes, e.g.:

	Class 1	Class 2
R_1	1	6
R_2	5	6

However, though both imperfect, R_1 is clearly sending a stronger 'signal' for a single class (Class 2) than R_2 .

Information Theory

One way to quantify the strength of a signal in a particular region is to analyze the distribution of classes within the region. We compute the **entropy** of this distribution.

For a random variable with a discrete distribution, the entropy is computed by

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

Higher entropy means the distribution is uniform-like (flat histogram) and thus values sampled it are ‘less predictable’ (all possible values are equally probable).

Lower entropy means the distribution has more defined peaks and valleys and thus values sampled from it are ‘more predictable’ (values around the peaks are more probable).

Entropy

Suppose we have J number of predictors, N number of training points and K classes.

Suppose we select the j -th predictor and split a region containing N number of training points along the threshold $t_j \in \mathbb{R}$.

We can assess the quality of this split by measuring the entropy of the class distribution in each newly created region, R_1, R_2 :

$$\text{Entropy}(i|j, t_j) = - \sum_k p(k|R_i) \log_2[p(k|R_i)]$$

Note: we are actually computing the conditional entropy of the distribution of training points amongst the K classes given that the point is in region i .

Example

	Class 1	Class 2	Entropy($i j, t_j$)
R_1	0	6	$-\left(\frac{6}{6} \log_2 \frac{6}{6} + \frac{0}{6} \log_2 \frac{0}{6}\right) = 0$
R_2	5	8	$-\left(\frac{5}{13} \log_2 \frac{5}{13} + \frac{8}{13} \log_2 \frac{8}{13}\right) \approx 1.38$

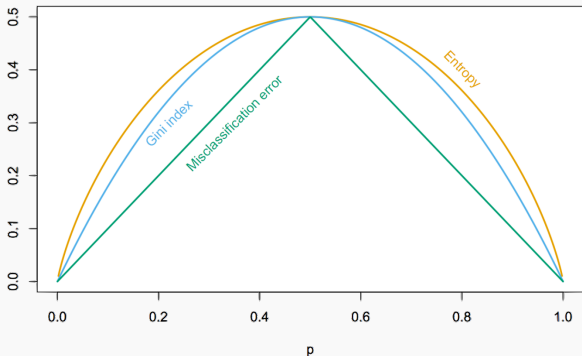
We can now try to find the predictor j and the threshold t_j that minimizes the average entropy over the two regions, weighted by the population of the regions:

$$\min \left\{ \frac{N_1}{N} \text{Entropy}(1|j, t_j) + \frac{N_2}{N} \text{Entropy}(2|j, t_j) \right\}$$

Comparison of Criteria

Recall our intuitive guidelines for splitting criteria, which of the three criteria fits our guideline the best?

We have the following comparison of the value of the three criteria at different levels of purity (from 0 to 1) in a single region.



Comparison of Criteria

Recall our intuitive guidelines for splitting criteria, which of the three criteria fits our guideline the best?

To note that entropy penalizes impurity the most is not to say that it is the best splitting criteria. For one, a model with purer leaf nodes on a training set may not perform better on the testing test.

Another factor to consider is the size of the tree (i.e. model complexity) each criteria tends to promote.

To compare different decision tree models, we need to first discuss **stopping conditions**.

Stopping Conditions & Pruning

Variance vs Bias

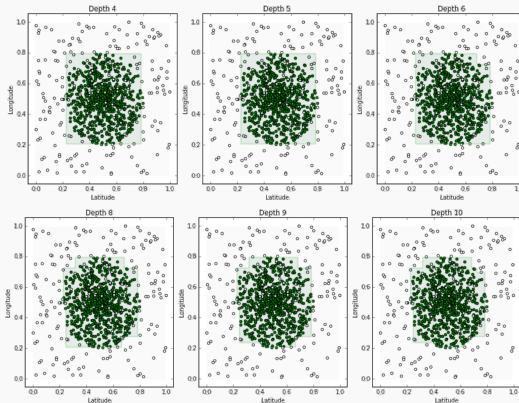
If we don't terminate the decision tree learning algorithm manually, the tree will continue to grow until each region defined by the model contains exact one training point (and the model attains 100% training accuracy).

To prevent this from happening, we can simply stop the algorithm at a particular depth.

But how do we determine the appropriate depth?

Variance vs Bias

Consider the result of training a decision tree of various depths on a previous example dataset:



Variance vs Bias

We make some observations about our models:

- ▶ **(Bias)** A tree of depth 4 is not a good fit for the training data - it's unable to capture the nonlinear boundary separating the two classes.
- ▶ **(Bias)** With an extremely high depth, we can obtain a model that correctly classifies all points on the boundary (by zig-zagging around each point).
- ▶ **(Variance)** The tree of depth 4 is robust to slight perturbations in the training data - the square carved out by the model is stable if you move the boundary points a bit.
- ▶ **(Variance)** Trees of high depth are sensitive to perturbations in the training data, especially to changes in the boundary points.

Not surprisingly, complex ones have low bias (able to capture more complex geometry in the data) but high variance (can over fit). Complex trees are also harder to interpret and more computationally expensive to train.

Stopping Conditions

Common simple stopping conditions:

- ▶ Don't split a region if all instances in the region belong to the same class
- ▶ Don't split a region if the number of instances in the sub-region will fall below pre-defined threshold
- ▶ Don't split a region if the total number of leaves in the tree will exceed pre-defined threshold

The appropriate thresholds can be determined by evaluating the model on a held-out data set or, better yet, via cross-validation.

Stopping Conditions

More restrictive stopping conditions:

- ▶ Don't split a region if the class distribution of the training points inside the region are independent of the predictors
- ▶ Compute the gain in purity, information or reduction in entropy of splitting a region R

$$\text{Gain}(R) = \Delta(R) = m(R) - \frac{N_1}{N}m(R_1) - \frac{N_2}{N}m(R_2)$$

where m is a metric like the Gini Index or entropy.
Don't split if the gain is less than some pre-defined threshold.

Pruning

Rather than preventing a complex tree from growing, we can obtain a simpler tree by ‘pruning’ a complex one.

There are many method of pruning, a common one is **cost complexity pruning**, where by we select from a array of smaller subtrees of the full model that optimizes a balance of performance and efficiency.

That is, we measure

$$C(T) = \text{Error}(T) + \alpha|T|$$

where T is a decision (sub) tree, $|T|$ is the number of leaves in the tree and α is the parameter for penalizing model complexity.

Pruning

The pruning algorithm:

1. Start with a full tree T_0 (each leaf node contains exactly one training point)
2. Replace a subtree in T_0 with a leaf node to obtain a pruned tree T_1 . This subtree should be selected to minimize

$$\frac{\text{Error}(T_0) - \text{Error}(T_1)}{|T_0| - |T_1|}$$

3. Iterate this pruning process to obtain T_0, T_1, \dots, T_L , where T_L is the tree containing just the root of T_0 .
4. Select the optimal tree T_i by cross validation.

Note: you might wonder where we are computing the cost-complexity $C(T_i)$. One can prove that this process is equivalent to explicitly optimizing C .

An Example

[demonstrate difference between different splitting criteria] [demonstrate difference between different stopping conditions] [demonstrate overfitting and variance]