

Puzzle: Time-series forecasting

Suppose we wish to predict whether the price of a stock is going to go up or down tomorrow

- Take history over a long period of time
- Normalize the time series to zero mean, unit variance
- Form all possible input-output pairs with
 - input = previous 20 days of stock prices
 - output = price movement on the 21st day
- Randomly split data into training and testing data
- Train on training data only, test on testing data only

Based on the test data, it looks like we can consistently predict the price movement direction with accuracy ~52%

Are we going to be rich?

Decision trees

A ***decision tree*** is a method for classification/regression that aims to ask a few relatively simple questions about an input $\mathbf{x} \in \mathbb{R}^d$ and then predicts the associated output y

Decision trees are useful to a large degree because of their ***simplicity*** and ***interpretability***

The resulting output rule is something that can easily be inspected and interpreted by hand... something that is not really the case with most algorithms

We will also see later that decision trees form a useful building block for other more sophisticated algorithms

Hypothetical example

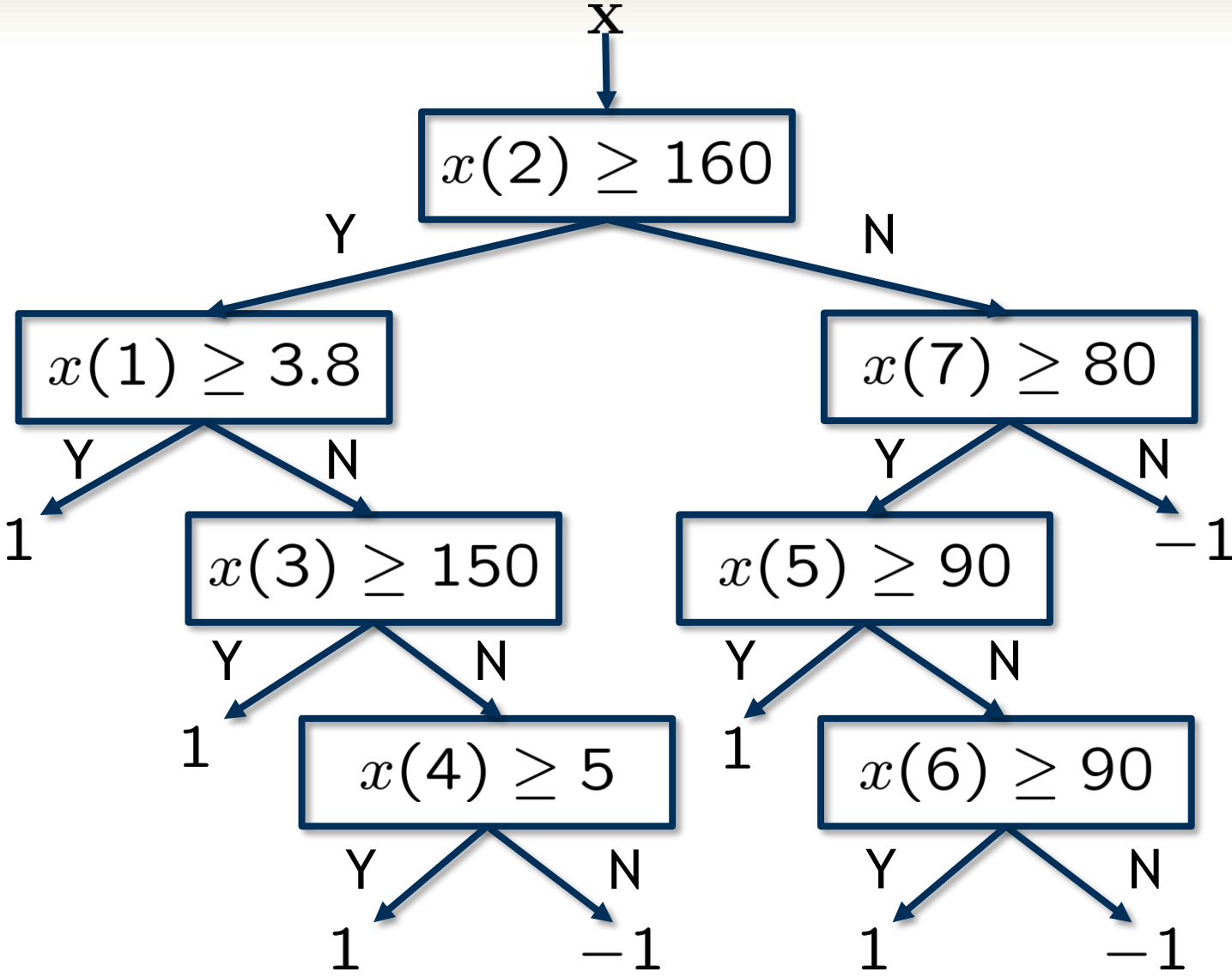
Suppose we are given a dataset where the x_i are 7-dimensional feature vectors with features

1. Undergraduate GPA (4.0 scale)
2. Quantitative GRE score (130-170)
3. Verbal GRE score (130-170)
4. Analytical Writing GRE score (0-6)
5. Undergraduate institution reputation (0-100)
6. Statement of purpose evaluation (0-100)
7. Letter of recommendation evaluation (0-100)

The labels for this dataset are:

$$y_i = \begin{cases} 1 & \text{student was accepted} \\ 0 & \text{student was denied} \end{cases}$$

Example decision tree



Generalizations

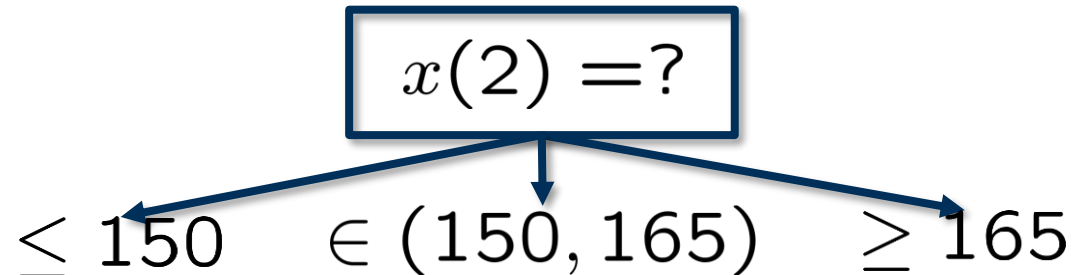
Decision trees can be applied to regression where the labels at the “leaf nodes” are real-valued (or real-valued functions)

Other generalizations include:

- splits that involve more than one feature



- Splits involving more than two outcomes

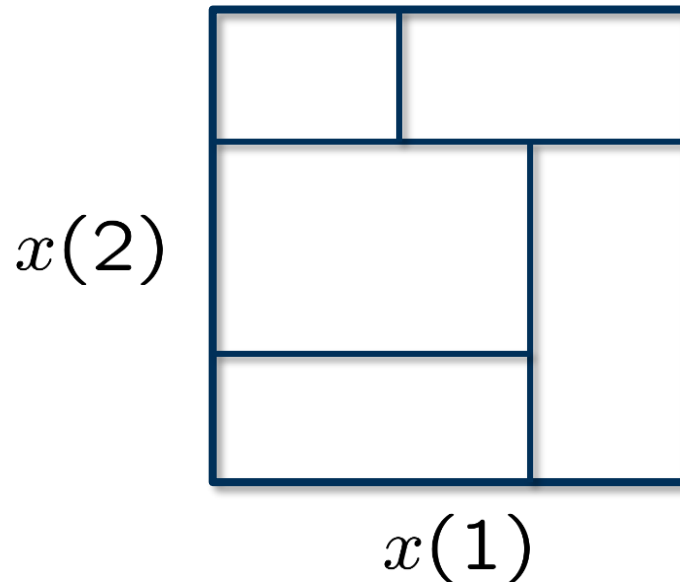


Binary decision trees

Using multiple features and allowing splits with more than two outcomes generally increases the risk of *overfitting*

We will restrict our attention to binary, single-feature splits

In this case, every (binary) decision tree is associated with a *partition* of the feature space that looks something like this example in \mathbb{R}^2



Terminology

- The elements of the partition are called *cells*
- Recall that a *graph* is a collection of *nodes* or *vertices*, some of which are joined by *edges*
- The *degree* of a vertex is the number of edges incident on that vertex
- A *tree* is a *connected* graph with *no cycles*

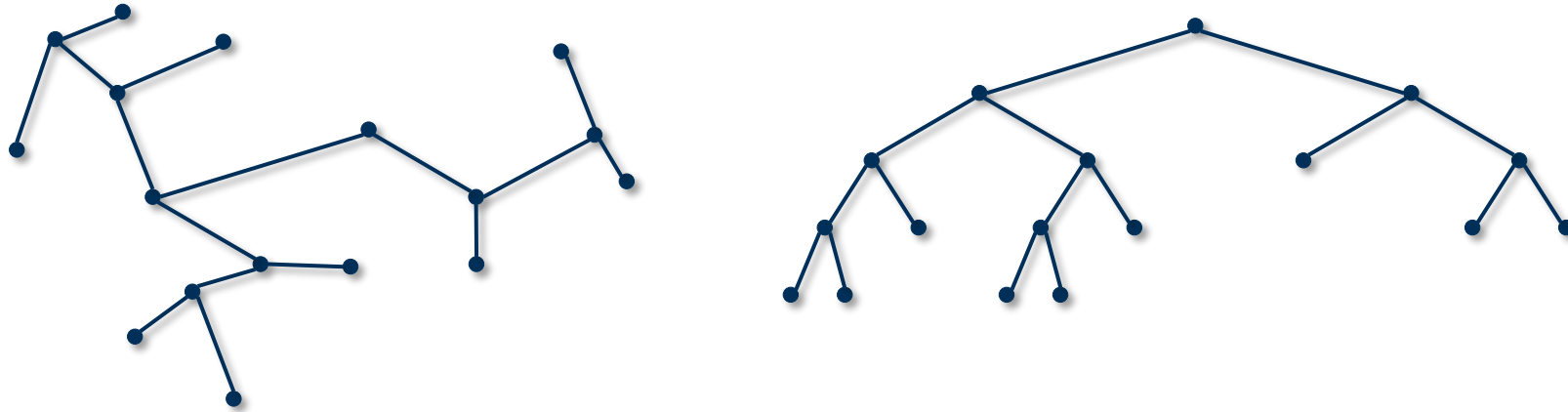
When is a tree not a tree?



Rooted binary trees

A *rooted binary tree* is a tree where

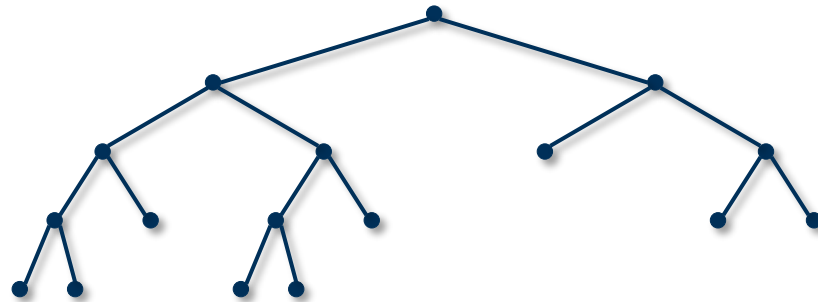
- one vertex, called the *root*, has degree 2
- and all other vertices have degree 1 or 3



- vertices with degree 1 are called *leaf* or *terminal* vertices
- vertices with degree 2 or 3 are called *internal* vertices

More terminology

- The **depth** of a vertex is the length to the root
- The **parent** of a vertex is the neighbor with depth one less
- Two vertices are **siblings** if they have the same parent
- A **subtree** is a subgraph that is also a tree
- A **rooted binary subtree** is a subtree that contains the root and is such that if any non-root vertex is in the subtree, so is its sibling



A **binary decision tree** is a rooted binary tree where each internal vertex is associated with a binary classifier, and each leaf with a label

Learning decision trees

Let \mathcal{T} denote the set of all binary decision trees

Given a dataset $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, it might be tempting to pose the learning problem as an optimization problem of the form

$$\min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, T(\mathbf{x}_i))$$

where ℓ is an appropriate loss function, e.g.,

- 0/1 loss for classification
- squared error loss for regression

Unfortunately, this will lead to massive overfitting

If we grow the tree deep enough, we can usually fit the training data perfectly!

Penalized empirical risk minimization

Since deep/complex decision trees tend to overfit, we would like to avoid such trees

A natural way to quantify the complexity of a tree is the number of leaf nodes, which we denote by $|T|$

Note that if $|T| = n$, we would very likely be overfitting

The penalized ERM approach is to consider the optimization problem

$$\min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, T(\mathbf{x}_i)) + \lambda |T|$$

Learning a tree in practice

Unfortunately, this optimization problem is intractable

A two-stage procedure is typically employed in practice

1. Grow a very large tree T_0 in a greedy fashion
2. Prune T_0 by solving

$$\min_{T \in \mathcal{T}_0} \frac{1}{n} \sum_{i=1}^n \ell(y_i, T(\mathbf{x}_i)) + \lambda |T|$$

where \mathcal{T}_0 is the set of all decision trees based on rooted binary subtrees of T_0

Growing a decision tree

The construction of T_0 follows a simple greedy (looking just one step ahead) strategy

1. Start at the root vertex (i.e., the entire feature space)
2. Decide whether to stop growing tree at current vertex
 - If yes, assign a label and stop
3. If no, consider all possible ways to split the data in the cell corresponding to the current vertex, and select the best one
4. For each branch of the split, create a new vertex and go to step 2

Implementation

To implement this strategy, we need:

- A list of possible splits

When considering splits based on a single real-valued feature, splits have the form

$$x(j) \leq t?$$

Since there are only n data points, only at most $n - 1$ values of t need to be considered for each j

For discrete or categorical features, other simple splits can be used, such as

$$x(j) = \text{“blue”}?$$

Implementation

To implement this strategy, we need:

- A labeling rule

For classification, we can just assign labels by majority vote over data in the cell corresponding to the current vertex

For regression, we can just take the average of the y_i in the cell for a piecewise constant approximation

(or least squares fit for piecewise polynomial or other approximation)

Implementation

To implement this strategy, we need:

- A rule for stopping splitting

The most common strategy is to just split until each leaf vertex contains a single data point

- A rule for selecting the best split

This is the hard part!

Split selection

Let's focus on binary classification

Suppose V is a leaf vertex at some stage in the growing process

We can think of V as also defining a cell in the feature space, allowing us to consider $\{\mathbf{x}_i : \mathbf{x}_i \in V\}$

Intuitively, a good split of V will lead to children that are more *homogenous* or *pure* than V

To define this, we will define a notion of *impurity*

Impurity measure

Assume the class labels are $\{+1, -1\}$

Let

$$q := \frac{|\{i : \mathbf{x}_i \in V, y_i = 1\}|}{|\{i : \mathbf{x}_i \in V\}|}$$

Note that $(q, (1 - q))$ defines a probability distribution on the labels (i.e., on the probability of getting each label on a randomly selected point in V)

An ***impurity measure*** is a function $i(V)$ such that

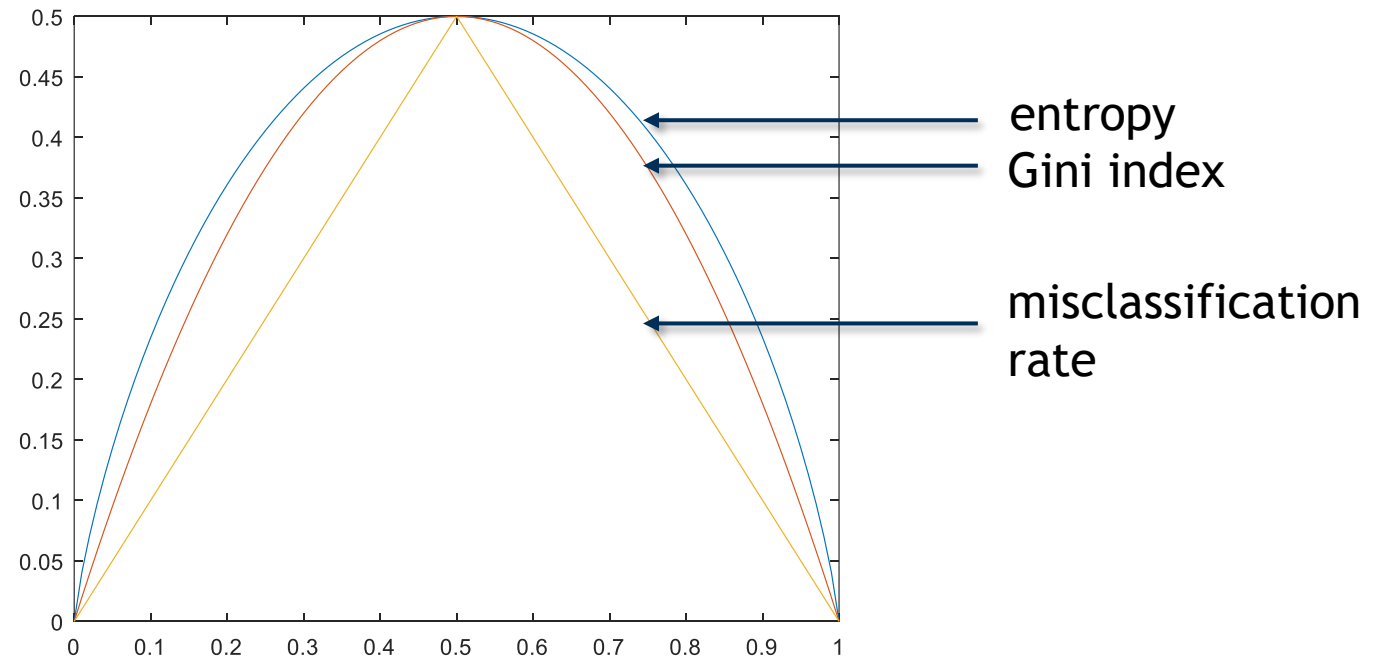
- $i(V) \geq 0$, with $i(V) = 0$ if and only if V consists of a single class
- a larger value of $i(V)$ indicates that the distribution defined by $(q, (1 - q))$ is closer to the uniform distribution

Examples

Entropy: $i(V) = -(q \log q + (1 - q) \log(1 - q))$

Gini index: $i(V) = 2q(1 - q)$

Misclassification rate: $i(V) = \min(q, 1 - q)$



Using the impurity measure

To select the best split, we aim to *maximize* the *decrease* in impurity

Let V_1 and V_2 be two children of V

Define

$$p(V_1) = \frac{|\{\mathbf{x}_i : \mathbf{x}_i \in V_1\}|}{|\{\mathbf{x}_i : \mathbf{x}_i \in V\}|} \quad p(V_2) = \frac{|\{\mathbf{x}_i : \mathbf{x}_i \in V_2\}|}{|\{\mathbf{x}_i : \mathbf{x}_i \in V\}|}$$

The decrease in impurity is

$$i(V) - \left(p(V_1)i(V_1) + p(V_2)i(V_2) \right)$$

When i is entropy, this is called the *information gain*

Pruning

Pruning the model involves solving the optimization problem:

$$\min_{T \in \mathcal{T}_0} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, T(\mathbf{x}_i)) + \lambda |T|}_{f(T)}$$

$f(T)$ is additive in the following sense:

For any tree T , let $\Pi(T) = \{T_1, T_2\}$ be the partition of the tree corresponding to the children of the root vertex

Then $f(T) = f(T_1) + f(T_2)$

Applying this idea recursively suggests a natural algorithm of ***weakest link pruning***. One can also attack this via an efficient “bottom up” dynamic programming approach.

Why grow then prune?

Why should we go to the trouble of growing the tree and then pruning?

A simpler approach would be to just grow the tree and then stop when the decrease in impurity is negligible

Answer: Ancillary splits

These are splits that have no value by themselves, but enable useful splits later on

Remarks

Advantages

- interpretable
- rapid evaluation
- easily handles
 - categorical/mixed data
 - missing data
 - multiple classes

Disadvantages

- unstable: slight perturbations of training data can drastically alter the learned tree
- jagged decision boundaries

Ensemble methods

Ensemble methods address both of these deficiencies in decision trees as well as other algorithms

The first step is to generate a number of classifiers f_1, \dots, f_T (all using the same dataset) using some method that typically involves some degree of randomness

The second step is to combine these into a single classifier

Even if *none* of the individual classifiers are particularly good, the combined result can *far* outperform any of the individual classifiers and can be surprisingly effective

“The wisdom of the crowds”

Sir Francis Galton (1822-1911)

- cousin of Charles Darwin
- statistician (introduced correlation and standard deviation)
- father of eugenics... wary of democracy and distrustful of “the mob”

How much does this ox weigh?



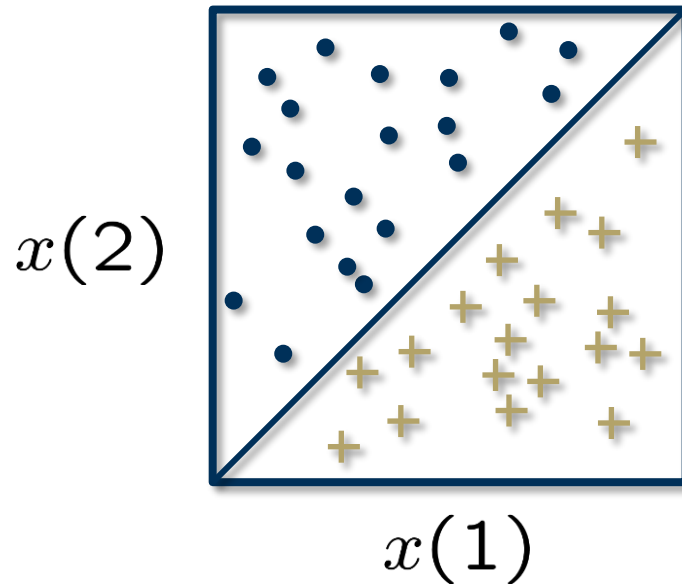
If we collect hundreds of uneducated farmers (with no particular expertise in weighing oxen), how well will they do?

Mean of the guesses: 1,197 pounds

Actual weight: 1,198 pounds

An example in classification

Suppose that the feature space is $[0, 1]^2$ and that the data looks like:



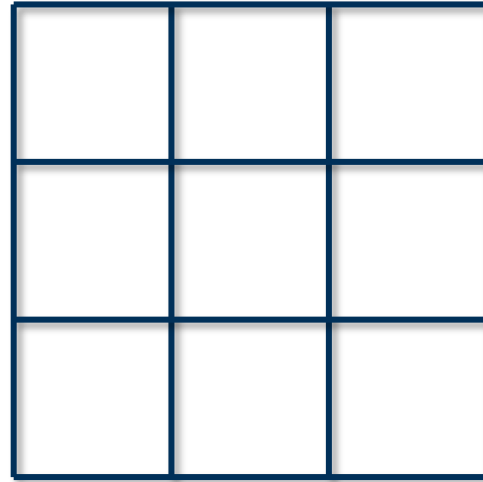
In this scenario, the Bayes risk is zero...

But the risk of certain simple classifiers can still be large

Histogram classifiers

Suppose that we are using *histogram classifiers*

In particular, we are using classifiers based on a regular partition of $[0, 1]^2$ into 9 squares



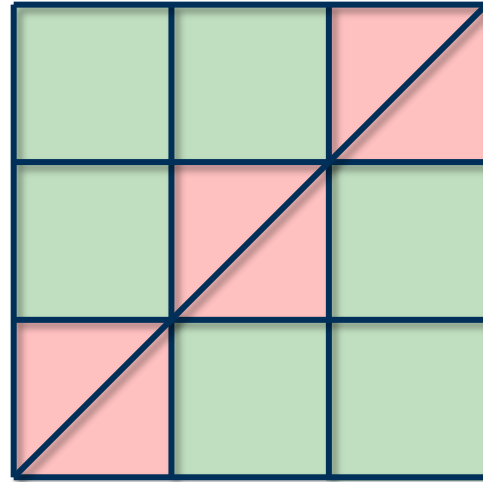
Label of each cell determined by majority vote

Histogram classifiers are pretty bad!

This classifier will not perform very well for the given distribution (or indeed, most distributions)

The risk of this classifier is:

$$R = \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{6}$$



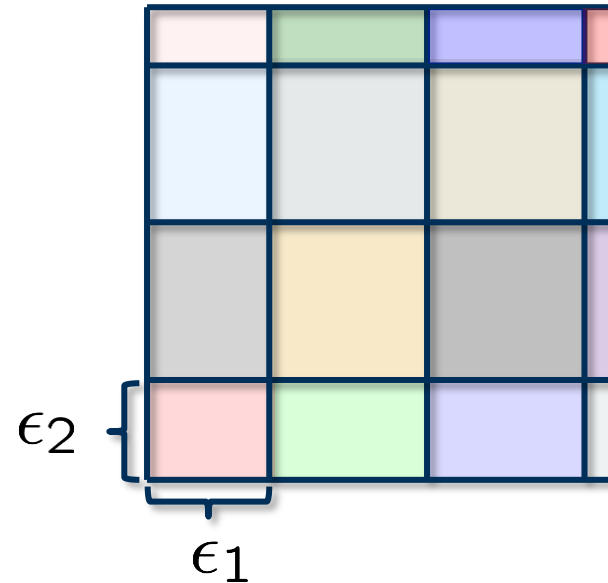
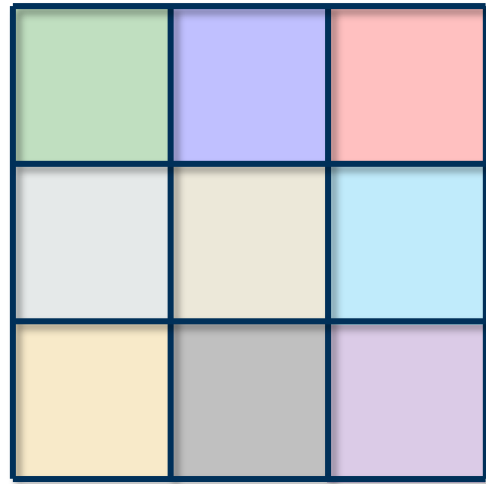
You can easily imagine that binary decision trees would have similar trouble with this example

However, we will see that with an appropriate ensemble method, we can make this classifier much more effective

Randomly shifted histogram classifiers

Suppose that we generate $\epsilon_1, \epsilon_2 \in [0, \frac{1}{3}]$ uniformly at random

Then shift the partition by $[\epsilon_1, \epsilon_2]^T$

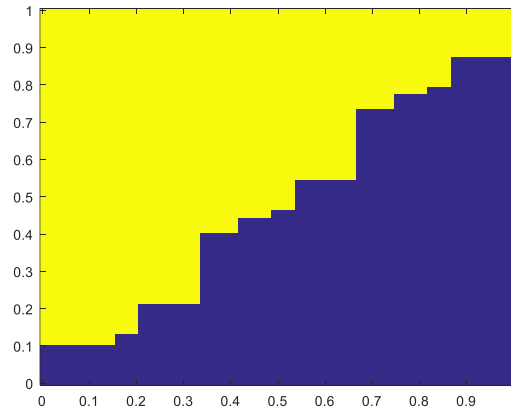


Ensemble histogram classifier

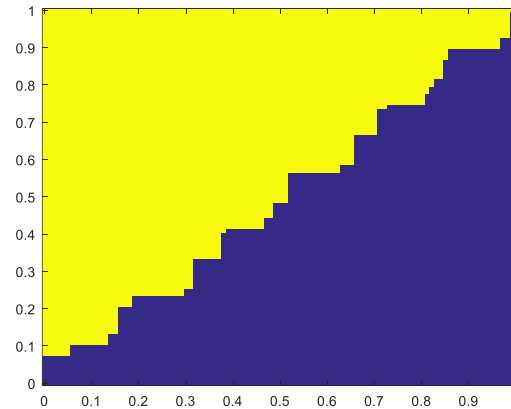
Generate f_1, \dots, f_T as independent randomly shifted histogram classifiers and take majority vote

Example: $n = 1000$

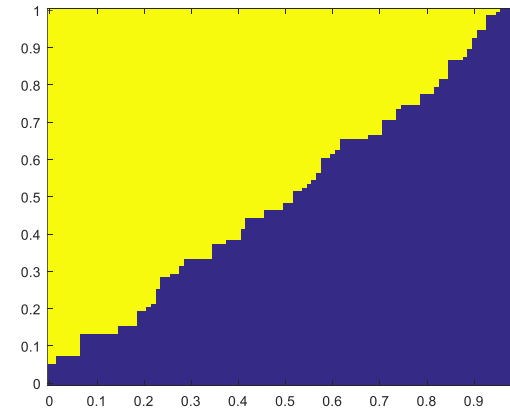
$T = 5$



$T = 11$



$T = 21$



Bagging

Another way to introduce some randomness is via **bagging**

Bagging is short for **bootstrap aggregation**

Given a training sample of size n , for $b = 1, \dots, B$ let I_b be a list of size n obtained by sampling from $\{1, \dots, n\}$ **with replacement**

Recall that I_b is called a **bootstrap sample**

Suppose we have a fixed learning algorithm

Let f_b be the classifier we obtain by applying this learning algorithm to $\{(\mathbf{x}_i, y_i)\}_{i \in I_b}$

The **bagging classifier** is just the majority vote over f_1, \dots, f_B

Random forests

A *random forest* is an ensemble of decision trees where each decision tree is (independently) randomized in some fashion

Bagging with decision trees is a simple example of a random forest

In the specific context of decision trees, bagging has one pretty big drawback

- bootstrap samples are highly correlated
- as a result, the different decision trees tend to select the same features as most informative
- this leads to partitions that tend to be highly correlated
- we would rather have partitions that are more “independent”

Random feature selection

One way to achieve this is to also incorporate *random feature selection*

- generate an classifiers by choosing random subsets of features and designing decision trees on just those features
- can be combined with bagging

Random features lead to less correlated partitions, translating to a reduced variance for the ensemble prediction

Rule of thumb: use \sqrt{d} random features

Random forests are possibly the best “off-the-shelf” method for classification

Approach also extends to regression