

CS189

decision trees

Midterm: everything up until this lecture

decision tree: Nonlinear method used for classification

- uses trees w/ 2 node types

- 1) Internal nodes test feature values (usually just one)
- branch after

- 2) Leaf Nodes

- predict the class or regression value $h(x)$

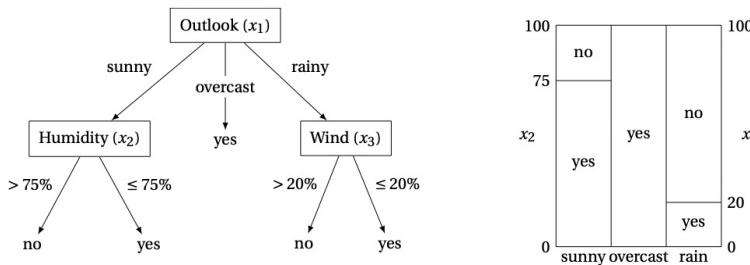
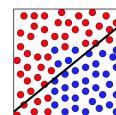
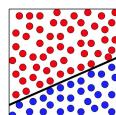


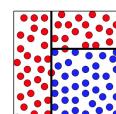
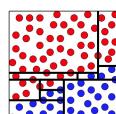
Figure 14.1: An example decision tree for when to go out for a picnic

- cuts X-space into rectangles
- works well w/ categorical and quantitative features
- interpretable results
- decision boundary can be arbitrarily complicated
- fast: only 1 feature per node

linear classifier



decision tree



need more decision rules

Classification w/ Decision Trees:

- greedy, top-down heuristic
- recursive

let $S \subseteq \{1 \dots n\}$ be set of sample point indices

top-level call has $S = \{1 \dots n\}$ (consider all nodes)

GrowTree(S):

```

if ( $\forall i \in S$  and some class  $C$ ) {
    return new leaf( $C$ )  (leaves are pure  $\rightarrow$  represent exactly 1 class)
} else {
    choose best splitting feature  $j$  and splitting value  $\beta$ 
     $S_L = \{i \in S : X_{ij} < \beta\}$     $\leq 2 >$  also work
     $S_R = \{i \in S : X_{ij} \geq \beta\}$ 
    return new node( $j, \beta, \text{GrowTree}(S_L), \text{GrowTree}(S_R)$ )
}

```

How do we choose the best split?

- try all splits:
 - for each set: $J(S) = \text{cost of } S$
 - choose split that minimizes $J(S_L) + J(S_R)$
 - or, split that minimizes weighted avg. by # training points of child

(cardinality 1:1)

$$|S_L|J(S_L) + |S_R|J(S_R)$$

$$|S_L| + |S_R| = |S|$$

How do we choose $J(S)$?

Bad idea: label S w/ the class C that labels the most points in S
 $\hookrightarrow J(S) = \# \text{ of pts. in } S \text{ not in class } C$

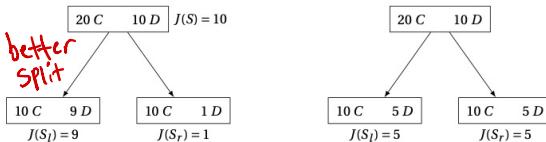


Figure 14.3: Two possible splits of a set S

Problem: $J(S_L) + J(S_R)$ same for both splits (many splits will have same total cost)

2nd idea: entropy

let Y be a random class variable, $P(Y=c) = p_c$

Surprise of Y being class C : $-\log_2 p_c$ (always nonnegative)

- ↪ event w/ prob 1 → 0 Surprise
- ↪ event w/ prob 0 → infinite Surprise

might involve fractional bits

↪ useful when compiling lots of events

Surprise: expected number of bits of information needed to transmit
 (information which events happened to a recipient who knows the probabilities of the events
 theory def.)

entropy of index set S is the average Surprise:

$$H(S) = - \sum p_c \log_2 p_c, \quad p_c = \frac{|\{i \in S : Y_i = c\}|}{|S|}$$

↪ proportion of points in S that
 are class C

example Scenarios:

- all points in S belong to same class? $H(S) = -1 \log_2 1 = 0$
- half class C , half class D $H(S) = -1/2 \log_2 1/2 - 1/2 \log_2 1/2 = 1$
- n points, all different classes: $H(S) = -\log_2 1/n = -\log_2 1 + \log_2 n = \log_2 n$

$$-\sum_{i=1}^n \frac{1}{n} \log_2 \frac{1}{n} = -\log_2 \frac{1}{n}$$

Entropy: expected # bits needed to transmit the class of sample point in S chosen
 uniformly at random

e.g. all points are class $C \rightarrow 0$ bits needed

half C & half $D \rightarrow 1$ bit

Specify one of N classes: $\log_2 N$

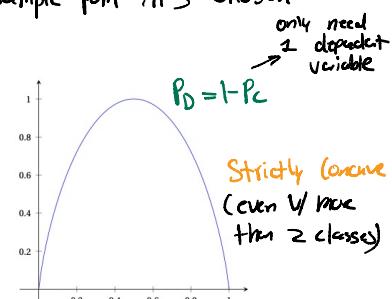


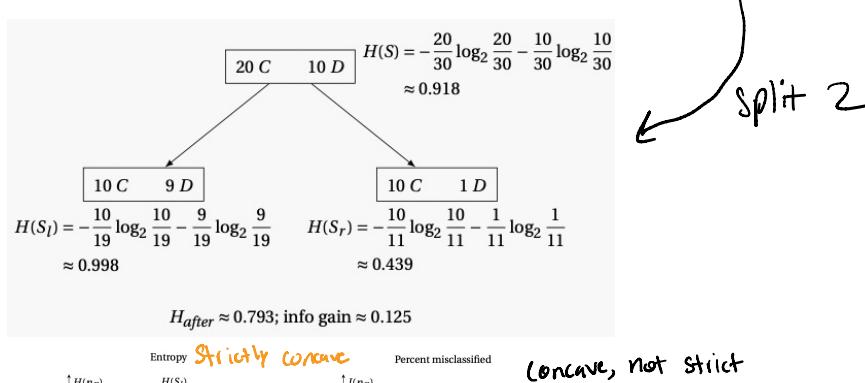
Figure 14.4: Plot of the binary entropy function $H(p_C) = -p_C \log_2 p_C - (1 - p_C) \log_2 (1 - p_C)$
 (2 classes)

Weighted avg. entropy: $H_{\text{after}} = \frac{|S_L|H(S_L) + |S_R|H(S_R)}{|S_L| + |S_R|}$
 (after a split)

→ choose split that maximizes information gain $H(S) - H_{\text{after}}$
 (Same as minimizing H_{after} since $H(S)$ is a constant w.r.t. split)

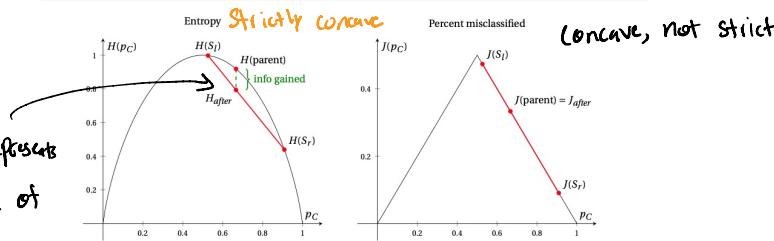
Info. gain always positive except:
 - one child is empty (\emptyset)
 or

$$P(Y_i=c | i \in S_L) = P(Y_i=c | i \in S_R)$$



split 2

any point on
line segment represents
Weighted average of
entropies for
suitable weights



info. gain: difference between weighted avg. entropy & parent entropy (directly above) } entropy
 Is always positive unless 2 child sets both have same p_C and lie at
 Same point on the curve.

W/ percent misclassified fn, curve isn't strictly concave (Segment connecting pts. may lie on curve)
 ↳ Split doesn't change # misclassified pts
 nor % pts. misclassified

biggest problem: many different splits can get same weighted avg → cost fn. doesn't distinguish quality of different splits well

* Entropy isn't the only cost fn. that works well, other concave ones like polynomial $P(1-p)$ work

Choosing a Split:

- if X_i is a binary feature, children are $X_i = 0 \ \& \ X_i = 1$
- If X_i has 3 or more discrete features, then the split depends on the application
- If X_i is quantitative:
 - 1) Sort Values of X_i in S
 - 2) Split between each pair of unequal consecutive values
 - 3) take best split at very end

Efficient: can use radix sort in linear time

↳ def use if n is huge

Note: as you scan sorted list left to right, update entropy in $O(1)$ time per point
 ↳ important for obtaining fast tree-building time

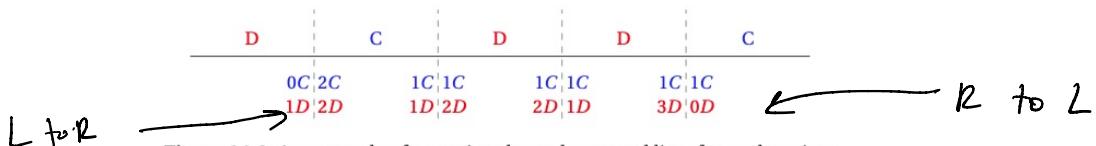


Figure 14.6: An example of scanning through a sorted list of sample points

Updating entropies: - $C := \#$ of class C sample points to the left of a potential split

- $c := \#$ sample pts. to right of split

- $D = \#$ points not in C to left of split

- $d = \#$ pts. not in C to right of split.

$$\rightarrow H(S_L) = \frac{-c}{c+d} \log_2 \frac{c}{c+d} - \frac{d}{c+d} \log_2 \frac{d}{c+d} \quad * \log_2 0 \text{ undefined, but}$$

$$H(S_R) = \frac{-c}{c+d} \log_2 \frac{c}{c+d} - \frac{d}{c+d} \log_2 \frac{d}{c+d} \quad \leftarrow \text{we suppose } 0 \log_2 0 = 0$$

choose split that minimizes weighted avg.

Weighted avg.
 $\Rightarrow -\frac{1}{n'} \left(c \log_2 \frac{c}{c+d} + d \log_2 \frac{d}{c+d} + c \log_2 \frac{c}{c+d} + d \log_2 \frac{d}{c+d} \right), n' = (c+d) + (c+d) = \# \text{pts. in tree node}$

Algorithms & Runtimes:

Classifying a test point: Walk down tree until leaf. return leaf's label

↳ Worst case: $O(\text{depth})$

↳ binary features: depth $\leq d$ (quantitative features may go deeper)

- Usually (not always): $\leq O(\log_2 n)$

Training:

Binary Features: try $O(d)$ splits at each node

- try all features $O(d)$, constant number $O(1)$ splits per feature

Quantitative features: try $O(n'd)$ splits, $n' = \# \text{pts. in node}$

- try all features $O(d)$, try all $O(n')$ splits per feature

↳ i.e. Scan through pts. linearly after radix sort

Both cases: $O(n'd)$ time to find & perform a split

↳ W binary features, Still need $O(n)$ time to split sample pts.

↳ quantitative features just as fast due to clever way
of computing entropy for each split

Each pt. participates in $O(\text{depth})$ nodes, $O(d)$ time in each node

↳ amortized analysis, charging $O(d \text{ depth})$ time to each pt.

Running Time $\leq O(nd \text{ depth})$, nd size of design matrix, logarithmic depth

