

## 1 Surprise and Entropy

In this section, we will clarify the concepts of surprise and entropy. Recall that entropy is one of the standards for us to split the nodes in decision trees until we reach a certain level of homogeneity.

- (a) Suppose you have a bag of balls, all of which are black. What's the surprise of taking out a ball whose color is black?

**Solution:** 0. We aren't surprised at all when events with probability 1 occur.

- (b) With the same bag of balls, what's the surprise of taking out a white ball?

**Solution:**  $\infty$ . We are infinitely surprised when an event with probability 0 occurs.

- (c) Now we have 10 balls in the bag, each of which is black or white. Under what color distribution(s) is the entropy of the bag minimized? And under what color distribution(s) is the entropy maximized? Calculate the entropy in each case.

*Recall:* The entropy of an index set  $S$  is the expected surprise of choosing an element from  $S$ . For a set  $S$ , the entropy

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

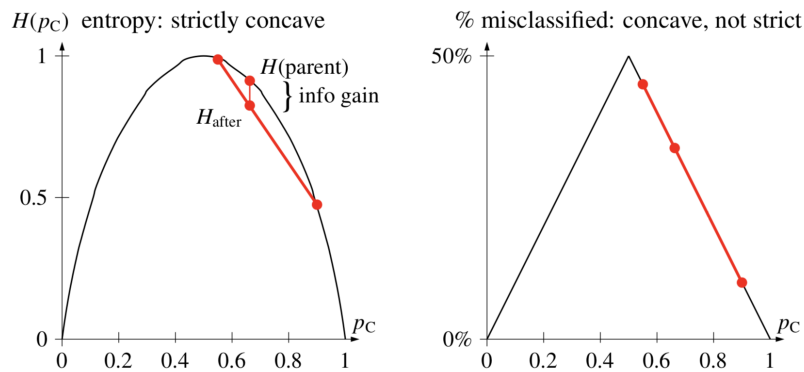
**Solution:** The entropy is minimized when, for example, all the balls are black or all the balls are white. In this case the entropy is 0. The entropy is maximized when half the balls are black and half the balls are white, in which case the entropy is  $-0.5 \log_2(0.5) - 0.5 \log_2(0.5) = 1$ .

- (d) Draw the graph of entropy  $H(p_c)$  when there are only two classes. Is the entropy function strictly concave, concave, strictly convex, or convex? Why? What is the significance?

*Hint:* For the significance, recall the information gain.

**Solution:** The function is strictly concave. Notice that the function  $-x \log x$  is strictly concave in  $[0, 1]$ , and a sum of strictly concave functions is strictly concave.

Significance: (from lecture) Suppose we pick two points on the entropy curve, then draw a line segment connecting them. Because the entropy curve is strictly concave, the interior of the line segment is strictly below the curve. Any point on that segment represents a weighted average of the two entropies for suitable weights. If you unite the two sets into one parent set, the parent set's value  $p_c$  is the weighted average of the children's  $p_c$ 's. Therefore, the point directly above that point on the curve represents the parent's entropy. The information gain is



the vertical distance between them. So the information gain is positive unless the two child sets both have exactly the same  $p_c$  and lie at the same point on the curve.

On the other hand, for the graph on the right, if we draw a line segment connecting two points on the curve, the segment might lie entirely on the curve. In that case, uniting the two child sets into one, or splitting the parent set into two, changes neither the total misclassified sample points nor the weighted average of the misclassified rate. The bigger problem, though, is that many different splits will get the same weighted average cost; this test does not distinguish the quality of different splits well.

## 2 Ensemble Learning

Ensemble learning is a general technique to combat overfitting, by combining the predictions of many varied models into a single prediction based on their average or majority vote.

- (a) **The motivation of averaging.** Consider a set of uncorrelated random variables  $\{Y_i\}_{i=1}^n$  with mean  $\mu$  and variance  $\sigma^2$ . Calculate the expectation and variance of their average. (In the context of ensemble methods, these  $Y_i$  are analogous to the prediction made by classifier  $i$ .)

**Solution:** The average of the  $Y_i$ s has the same expectation as each individual  $Y_i$ :

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n Y_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[Y_i] = \frac{1}{n} \cdot n\mu = \mu$$

but reduced variance compared to each of the individual  $Y_i$ 's:

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \left(\frac{1}{n}\right)^2 \sum_{i=1}^n \text{Var}(Y_i) = \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n}$$

- (b) **Ensemble Learning – Bagging.** In lecture, we covered bagging (Bootstrap AGGREGatING). Bagging is a randomized method for creating many different learners from the same data set.

Given a training set of size  $n$ , generate  $B$  random subsamples of size  $n'$  by sampling with replacement. Some points may be chosen multiple times, while some may not be chosen at all. If  $n' = n$ , around 63% are chosen, and the remaining 37% are called out-of-bag (OOB) samples.

(a) Why 63%?

**Solution:** Each sample has probability  $(1 - 1/n)^n$  of not being selected. For large  $n$ ,  $(1 - 1/n)^n \approx 1/e \approx .368$

(b) If we use bagging to train our model, How should we choose the hyperparameter  $B$ ? Recall,  $B$  is the number of subsamples, and typically, a few hundred to several thousand trees are used, depending on the size and nature of the training set.

**Solution:** An optimal number of trees  $B$  can be found by using cross-validation. Alternatively, we can observe the OOB error.

(c) In part (a), we see that averaging reduces variance for uncorrelated classifiers. Real world prediction will of course not be completely uncorrelated, but reducing correlation will generally reduce the final variance. Reconsider a set of correlated random variables  $\{Z_i\}_{i=1}^n$ . Suppose  $\forall i \neq j, \text{Corr}(Z_i, Z_j) = \rho$ . Calculate the variance of their average.

**Solution:**

$$\begin{aligned}\text{Var}\left(\frac{1}{n} \sum_{i=1}^n Z_i\right) &= \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n Z_i\right) = \frac{1}{n^2} \left( \sum_{i=1}^n \text{Var}(Z_i) + \sum_{i \neq j} \text{Cov}(Z_i, Z_j) \right) \\ &= \frac{n\sigma^2 + n(n-1)\sigma^2\rho}{n^2} = \rho\sigma^2 + \frac{1-\rho}{n}\sigma^2\end{aligned}$$

We can see that for large  $n$ , the first term dominates, which limits the benefit of averaging.

### 3 Decision Trees and Random Forests

Random forests are a specific ensemble method where the individual models are decision trees trained in a randomized way so as to reduce correlation among them. Because the basic decision tree building algorithm is deterministic, it will produce the same tree every time if we give it the same dataset and use the same hyperparameters (stopping conditions, etc.).

Consider constructing a decision tree on data with  $d$  features and  $n$  training points where each feature is real-valued and each label takes one of  $m$  possible values. The splits are two-way, and are chosen to maximize the information gain. We only consider splits that form a linear boundary parallel to one of the axes. In parts (a), (b) and (c) we will consider a standalone decision tree and not a random forest, so no randomization.

(a) **(From Discussion 8)** Prove or give a counter-example: In any path from the root to a leaf, the same feature will never be split on twice. If false, can you modify the conditions of the problem so that this statement is true?

**Solution:** False. Example: one dimensional feature space with training points of two classes  $x$  and  $o$  arranged as  $xxxooooxxx$ . This statement would be true if the splits were allowed to form more complex boundaries, i.e. if the splits were not binary and linear.

- (b) **(From Discussion 8)** Prove or give a counter-example: The information gain at the root is at least as much as the information gain at any other node.

*Hint:* Think about the XOR function.

**Solution:** False. Consider the XOR function, where the samples are

$$S = \{(0, 0; 0), (0, 1; 1), (1, 0; 1), (1, 1; 0)\},$$

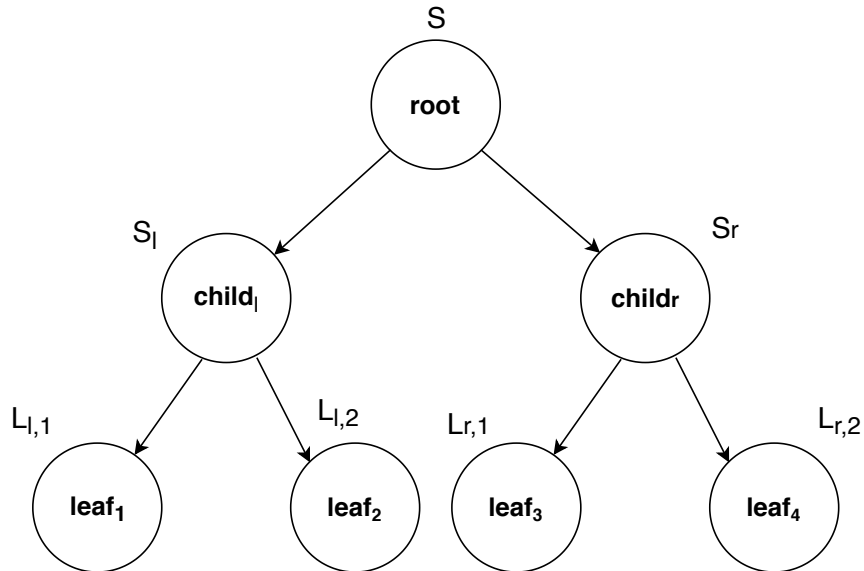
where the first two entries in every sample are features, and the last one is the label. Then,  $H(S) = 1$ . The first split is done based on the first feature, which gives  $S_l = \{(0, 0; 0), (0, 1; 1)\}$  and  $S_r = \{(1, 0; 1), (1, 1; 0)\}$ ; denote the corresponding nodes as **child<sub>l</sub>** and **child<sub>r</sub>** respectively. This gives  $H(S_l) = 1$  and  $H(S_r) = 1$ . Now we can compute the information gain of the first split:

$$IG(\mathbf{root}) = H(S) - \frac{|S_l|H(S_l) + |S_r|H(S_r)}{|S_l| + |S_r|} = 0.$$

Now we further split  $S_l$  and  $S_r$  according to the second feature, which gives 4 leaves of 1 sample each. Denote the leaf samples corresponding to  $S_r$  as  $L_{r,1}$  and  $L_{r,2}$ , and accordingly denote by  $L_{l,1}$  and  $L_{l,2}$  the leaves corresponding to  $S_l$ . Now we have

$$IG(\mathbf{child}_l) = H(S_l) - \frac{1 \cdot H(L_{l,1}) + 1 \cdot H(L_{l,2})}{1 + 1} = 1,$$

and analogously  $IG(\mathbf{child}_r) = 1$ . Therefore, the information gain at each of the child nodes is 1, while the information gain at the root is 0.



- (c) Prove or give a counter-example: For every value of  $a > 3$ , there exists some probability distribution on  $a$  objects such that its entropy is less than  $-1$ .

**Solution:** False. The entropy is always non-negative since  $-p \log p$  is non-negative when  $p \in [0, 1]$ .

- (d) One may be concerned that the randomness introduced in random forests may cause trouble. For example, some features or samples may not be considered at all. We will investigate this phenomenon in the next two parts. Consider  $n$  training points in a feature space of  $d$  dimensions. Consider building a random forest with  $T$  binary trees, each having exactly  $h$  internal nodes. Let  $m$  be the number of features randomly selected at each node. In order to simplify our calculations, we will let  $m = 1$ . For this setting, compute the probability that a certain feature (say, the first feature) is never considered for splitting.

**Solution:** The probability that it is not considered for splitting in a particular node of a particular tree is  $1 - \frac{1}{d}$ . The subsampling of  $m = 1$  features at each node is independent of all others. There are a total of  $ht$  nodes and hence the final answer is  $(1 - \frac{1}{d})^{ht}$ .

- (e) Now let us investigate the concern regarding the random selection of the samples. Suppose each tree employs  $n$  bootstrapped training samples. Compute the probability that a particular sample (say, the first sample) is never considered in any of the trees.

**Solution:** The probability that it is not considered in one of the trees is  $(1 - \frac{1}{n})^n$ . Since the choice for every tree is independent, the probability that it is not considered in any of the trees is  $(1 - \frac{1}{n})^{nT}$ , which approaches  $e^{-T}$  as  $n$  approaches infinity.

- (f) Compute the values of the probabilities you obtained in the previous two parts for the case when there are  $n = 2$  training points,  $d = 2$  dimensions,  $t = 10$  trees of depth  $h = 4$  (you may leave your answer in a fraction and exponentiated form, e.g., as  $(\frac{51}{100})^2$ ). What conclusions can you draw from your answer with regard to the concern mentioned in the beginning of the problem?

**Solution:**  $\frac{1}{2^{40}}$  and  $\frac{1}{2^{20}}$ . It is quite unlikely that a feature or a sample will be missed.