# 1  Entropy and Information

In this problem, we try to build intuition as to why entropy of a random variable corresponds to the amount of information that variable transmits. In particular, it determines the number of 0's and 1's needed to "efficiently" encode a random variable.

A coin with bias $b \in (0, 1)$ is flipped until the first head occurs, meaning that each flip gives heads with probability $b$. Let $X$ denote the number of flips required. Recall that the entropy of a random variable $Y$ is defined as:

$$H(Y) = - \sum_y \mathbb{P}(Y = y) \log(\mathbb{P}(Y = y)).$$

(a) Find the entropy $H(X)$. Assuming the logarithm in the definition of entropy has base 2, then the entropy is measured in *bits*.
   *Hint*: The following expressions might be useful:

$$\sum_{n=0}^{\infty} b^n = \frac{1}{1 - b}, \quad \sum_{n=1}^{\infty} nb^n = \frac{b}{(1 - b)^2}.$$

(b) Let $b = \frac{1}{2}$. Find an "efficient" sequence of yes-no questions of the form, "Is $X$ contained in the set $S$?", such that $X$ is determined as fast as possible. Compare $H(X)$ to the expected number of asked questions.

# 2 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. How surprised are you if you take out a black ball?

(b) With the same bag of balls, how surprised are you if you take out a white ball?

(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 a measure of expected surprise from choosing an element from $S$; that is,

$$H(S) = -\sum_C p_C \log_2(p_C), \text{ where } p_C = \frac{|i \in S : y_i = C|}{|S|}.$$

(d) Draw the graph of entropy $H(p_c)$ when there are only two classes C and D, with $p_D = 1 - p_C$. Is the entropy function strictly concave, concave, strictly convex, or convex? Why? What is the significance?

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

# 3  Decision Trees

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. We will only consider a standalone decision tree and not a random forest (hence no randomization). Recall the definition of information gain:

$$IG(\mathbf{node}) = H(S) - \frac{|S_l|H(S_l) + |S_r|H(S_r)}{|S_l| + |S_r|},$$

where $S$ is set of samples considered at **node**, $S_l$ is the set of samples remaining in the left subtree after **node**, and $S_r$ is the set of samples remaining in the right subtree after **node**. Intuitively, this is the difference between the entropy of a node's elements if no split is performed and the entropy of the node's elements if a split is performed. Thus, information gain essentially measures the reduction in entropy achieved by the split.

(a) 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?

(b) 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.

(c) Intuitively, how is the depth of a decision tree related to overfitting and underfitting?

(d) Suppose that a learning algorithm is trying to find a consistent hypothesis when the labels are actually being generated randomly. There are $d$ Boolean features and 1 Boolean label, and examples are drawn uniformly from the set of $2^{d+1}$ possible examples. Calculate the number of samples required before the probability of finding a contradiction in the data reaches $\frac{1}{2}$.
   (A contradiction is reached if two samples with identical features but different labels are drawn.)