

Problem Set 3

Instructor: Kamalika Chaudhuri

Due on: never

Problem 1: 8 points

A group of biologists would like to determine which genes are associated with a certain form of liver cancer. After much research, they have narrowed the possibilities down to two genes, let us call them A and B. After analyzing a lot of data, they have also calculated the following joint probabilities.

	Cancer	No Cancer		Cancer	No Cancer
Gene A	$\frac{1}{2}$	$\frac{1}{10}$	Gene B	$\frac{2}{5}$	$\frac{3}{20}$
No Gene A	$\frac{1}{5}$	$\frac{1}{5}$	No Gene B	$\frac{3}{10}$	$\frac{3}{20}$

- Let X denote the 0/1 random variable which is 1 when a patient has cancer and 0 otherwise. Let Y denote the 0/1 random variable which is 1 when gene A is present, 0 otherwise, and let Z denote the 0/1 random variable which is 1 when gene B is present and 0 otherwise. Write down the conditional distributions of $X|Y = y$ for $y = 0, 1$ and $X|Z = z$, for $z = 0, 1$.
- Calculate the conditional entropies $H(X|Y)$ and $H(X|Z)$.
- Based on these calculations, which of these genes is more informative about cancer?

Problem 2: 8 points

Since a decision tree is a classifier, it can be thought of as a function that maps a feature vector x in some set \mathcal{X} to a label y in some set \mathcal{Y} . We say two decision trees T and T' are *equal* if for all $x \in \mathcal{X}$, $T(x) = T'(x)$.

The following are some statements about decision trees. For these statements, assume that $\mathcal{X} = \mathbb{R}^d$, that is, the set of all d -dimensional feature vectors. Also assume that $\mathcal{Y} = \{1, 2, \dots, k\}$. Write down if each of these statements are correct or not. If they are correct, provide a brief justification or proof; if they are incorrect, provide a counterexample to illustrate a case when they are incorrect.

- If the decision trees T and T' do not have exactly the same structure, then they can never be equal.
- If T and T' are any two decision trees that produce zero error on the same training set, then they are equal.

Problem 3: 8 points

- A fair coin (that is, a coin with equal probability of coming up heads and tails) is flipped until the first head occurs. Let X denote the number of flips required. What is the entropy $H(X)$ of X ? You may find the following expressions useful:

$$\sum_{j=0}^{\infty} r^j = \frac{1}{1-r}, \quad \sum_{j=0}^{\infty} jr^j = \frac{r}{(1-r)^2}$$

- Let X be a discrete random variable which takes values x_1, \dots, x_m and let Y be a discrete random variable which takes values x_{m+1}, \dots, x_{m+n} . (That is, the values taken by X and the values taken by Y are disjoint.) Let:

$$\begin{aligned} Z &= X \text{ with probability } \alpha \\ &= Y \text{ with probability } 1 - \alpha \end{aligned}$$

Find $H(Z)$ as a function of $H(X)$, $H(Y)$ and α .