## Part 1: Required Exercises

## 1 Conditional Independence [5 pts]

Consider three binary variables $a, b, c \in \{0, 1\}$ having the joint distribution given in Table 8.2.

| $a$ | $b$ | $c$ | $p(a, b, c)$ |
|---|---|---|---|
| 0 | 0 | 0 | 0.192 |
| 0 | 0 | 1 | 0.144 |
| 0 | 1 | 0 | 0.048 |
| 0 | 1 | 1 | 0.216 |
| 1 | 0 | 0 | 0.192 |
| 1 | 0 | 1 | 0.064 |
| 1 | 1 | 0 | 0.048 |
| 1 | 1 | 1 | 0.096 |

(a) Show by direct evaluation that this distribution has the property that $a$ and $b$ are marginally dependent, so that $p(a, b) \neq p(a)p(b)$, but that they become independent when conditioned on $c$, so that $p(a, b|c) = p(a|c)p(b|c)$ for both $c = 0$ and $c = 1$.
(b) Evaluate the distribution $p(a), p(b|c)$, and $p(c|a)$ corresponding to the joint distribution given in the table. Hence show by direct evaluation that $p(a, b, c) = p(a)p(c|a)p(b|c)$. Draw the corresponding Bayesian network.
Solution goes here.

## 2 Information Gain [5 pts]

Consider the following training set with two boolean features and one continuous feature.

|  | A | B | C | Class |
|---|---|---|---|---|
| Instance 1 | F | T | 120 | Benign |
| Instance 2 | T | F | 1090 | Benign |
| Instance 3 | T | T | 245 | Malignant |
| Instance 4 | F | F | 589 | Malignant |
| Instance 5 | T | T | 877 | Malignant |

(a) How much information about the class is gained by knowing whether or not the value of feature C is less than 475?
(b) How much information about the class is gained by knowing whether or not the value of features A and B are different?
Solution goes here.

## 3 $k$-Nearest Neighbor [5 pts]

Suppose we want to learn a $k$-nearest neighbor model with the following data set and we are using Leave One Out Cross Validation (LOOCV) to select $k$. What would LOOCV pick: $k = 1$, or $k = 2$, or $k = 3$. Use Manhattan distance for calculations.

|              | Feature 1 | Feature 2 | Class    |
|--------------|-----------|-----------|----------|
| Instance 1   | 2         | 3         | Positive |
| Instance 2   | 4         | 4         | Positive |
| Instance 3   | 4         | 5         | Negative |
| Instance 4   | 6         | 3         | Positive |
| Instance 5   | 8         | 3         | Negative |
| Instance 6   | 8         | 4         | Negative |

Solution goes here.

# 4   Nearest Neighbor Regression [5 points]

Given data points $x_1 = (-1, 0), x_2 = (0, 0), x_3 = (0, 1)$ in the 2-dimensional Euclidean space and their corresponding labels $y_1 = 1, y_2 = 2, y_3 = 3$, use weighted 2-Nearest Neighbor to compute the label for $x = (1, 1)$. Here the weighted 2-Nearest Neighbor estimate is

$$f(x) = \frac{\sum_{i=1}^{2} w_i y_{(i)}}{\sum_{i=1}^{2} w_i},$$

where the weight $w_i = 1/i$ and $y_{(i)}$ is the label of the $i$-th nearest neighbor.
Solution goes here.

# 5   Evaluation [5 points]

Consider the following confusion matrix of a 2-class problem.

|                  | actual positive | actual negative |
|------------------|-----------------|-----------------|
| predict positive | 60              | 30              |
| predict negative | 50              | 60              |

Table 1: Confusion matrix of a 2-class problem. There are 200 instances in total.

Compute the following: accuracy, error, precision, recall.
Solution goes here.

# 6   Logistic Regression [10 points]

Let $f(x) = \sigma(w^\top x)$ where $w = (1, 2)$ and $\sigma$ is the sigmoid function $\sigma(z) = 1/(1 + \exp(-z))$. Compute the gradient $\nabla f$ at the point $x = (3, 4)$.

# 7   Maximum A Posterior [10 points]

Given data points $\{x_i, 1 \leq i \leq n\}$ from the Gaussian distribution $N(\mu, I)$ where the mean $\mu$ is unknown. Use the prior $p(u) = N(x_0, I)$ and compute the Maximum A Posterior estimation of $\mu$.
Solution goes here.

# 8   Bayesian Networks [10 points]
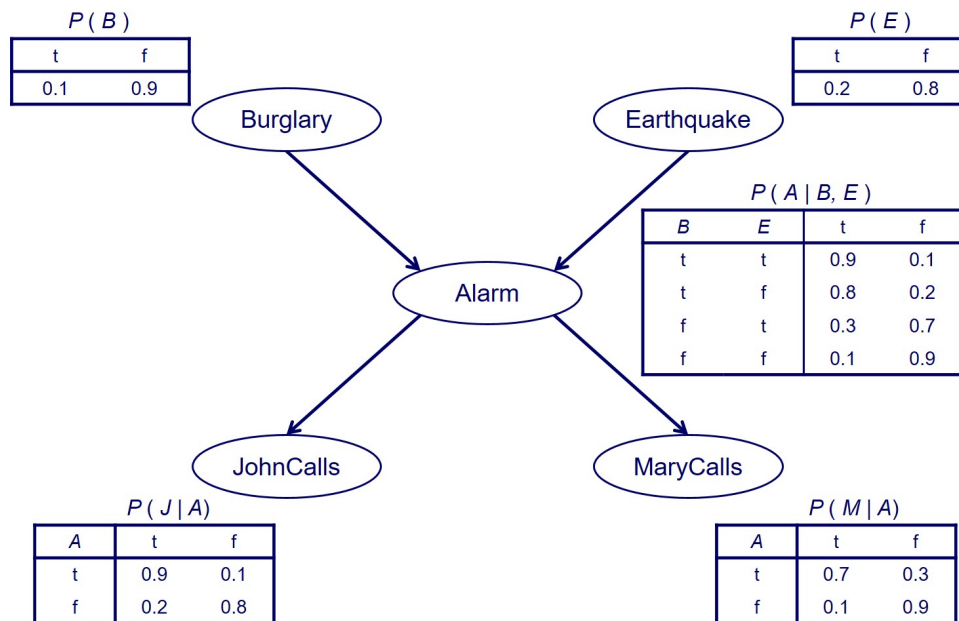
Consider the Bayesian Network in Figure 1.

| P ( B ) | |
|---|---|
| t | f |
| 0.1 | 0.9 |

| P ( E ) | |
|---|---|
| t | f |
| 0.2 | 0.8 |

| P ( A | B, E ) | | | |
|---|---|---|---|
| B | E | t | f |
| t | t | 0.9 | 0.1 |
| t | f | 0.8 | 0.2 |
| f | t | 0.3 | 0.7 |
| f | f | 0.1 | 0.9 |

| P ( J | A ) | | |
|---|---|---|
| A | t | f |
| t | 0.9 | 0.1 |
| f | 0.2 | 0.8 |

| P ( M | A ) | | |
|---|---|---|
| A | t | f |
| t | 0.7 | 0.3 |
| f | 0.1 | 0.9 |

Figure 1: A Bayesian Network example.

Compute $P(B = t, E = f, A = f, J = t, M = t)$ and $P(B = t, E = f, A = f, J = t | M = t)$.
Solution goes here.

# 9   Bayes Network: Sparse Candidate Algorithm [10 pts]

Suppose we wish to construct a Bayes Network for 3 features $X, Y$, and $Z$ using Sparse Candidate algorithm. We are given data from 100 independent experiments where each feature is binary and takes value $T$ or $F$. Below is a table summarizing the observations of the experiment:

| X | Y | Z | Count |
|---|---|---|---|
| T | T | T | 36 |
| T | T | F | 4 |
| T | F | T | 2 |
| T | F | F | 8 |
| F | T | T | 9 |
| F | T | F | 1 |
| F | F | T | 8 |
| F | F | F | 32 |

(a) Suppose we wish to compute a single candiate parent for Z. In the first round of the sparse Candidate algorithm, we compute the mutual information between $Z$ and the other random variables.

   i Compute the mutual information between $Z$ and $X$, i.e., $I(X, Z)$ based on the frequencies observed in the data.

   ii Compute the mutual information between $Z$ and $Y$, i.e., $I(Y, Z)$ based on the frequencies observed in the data.

(b) Based on your observations in part (a), which feature should be selected as candidate parent for $Z$? Why?
(c) In the first round of the algorithm, suppose that we choose $Y$ to be the parent of $Z$ in our network, $X$ to be the parent of $Y$, and that $X$ remains parent-less. Estimate the parameters of the current Bayes net, given the data.
Solution goes here.

## 10 Kernels [5 pts]

Suppose you are given the following instances in 2-D space.

| $X$ coordinate | $Y$ coordinate |
| --- | --- |
| 12 | 4 |
| 3 | 18 |
| 6 | 11 |
| 5 | 5 |

Build the Kernel Matrix for the above dataset for each of these kernels. That is, compute a matrix $K$ with entry $K_{ij}$ being the kernel value between point $i$ and point $j$.
(a) Polynomial kernel of degree 2, i.e., $k(x, z) = (x \cdot z)^2$. (b) RBF kernel with $k(x, z) = \exp(-\gamma |x_1 - x_2|^2)$ with $\gamma = 0.01$.
Solution goes here.

## 11 Kernel Methods [10 points]

Consider the following kernel

$$k(z, z') = \begin{cases} 1 & \text{if } \|z - z'\|_2 \le 1, \\ 0 & \text{otherwise.} \end{cases}$$

Given data set $x_1 = (0, 0), y_1 = 1, \quad x_2 = (0, 1), y_2 = 2, \quad x_3 = (1, 0), y_3 = 3$, define function $f(x) = \sum_{i=1}^{3} \alpha_i y_i k(x, x_i)$ where the coefficients $\alpha_i = i$. Compute $f(x)$ for $x = (1, 1)$.
Solution goes here.

## 12 Principal Component Analysis [10 points]

What is the first principal component of the following data points:

$$x_1 = (-1, 0), x_2 = (1, 0), x_3 = (0, -0.1), x_4 = (0, 0.1).$$

Solution goes here.

## 13 Reinforcement Learning [10 points]

Consider the deterministic reinforcement environment drawn below (let $\gamma = 0.1$). the number on the arcs indicate the immediate rewards. Assume we learn a Q-table. Also assume all the initial values in your Q table are 5.
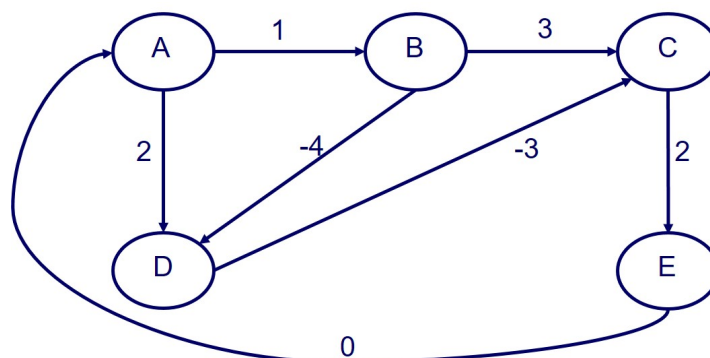


Figure 2: A deterministic reinforcement environment.

Suppose the learner follows the path $A \to D \to C \to E \to A$. Using the standard Q learning for deterministic reinforcement environment, report the final Q table on the graph above.
Solution goes here.

## Part 2: Extra Credits

## 14   Decision Tree Rank [5 points]

The rank of a decision tree is defined as follows. If the tree is a single leaf then the rank is 0. Otherwise, let $r_L$ and $r_R$ be the ranks of the left and right subtrees of the root, respectively. If $r_L = r_R$ then the rank of the tree is $r_L + 1$. Otherwise, the rank is the maximum of $r_L$ and $r_R$. Prove that a decision tree with $n$ leaves has rank at most $\log_2(n)$.

Solution goes here.

## 15   Kernel Methods [10 points]

Car-talk statistician Marge Innovera proposes the following simple kernel function:

$$k(z, z') = \begin{cases} 1 & \text{if } z = z', \\ 0 & \text{otherwise.} \end{cases}$$

Marge likes this kernel because in the $\Phi$-space, any labeling of the points in the instance space $X$ will be linearly separable. So, this should be perfect for learning any target function you want to: just run a kernelized version of SVM.

1) Why is any assignment of labels to points linearly separable?

2) Nonetheless, what is the problem with her reasoning?

Solution goes here.

## 16   Support Vector Machines [10 points]

Given data $\{(x_i, y_i), 1 \le i \le n\}$, the (hard margin) SVM objective is

$$\min_{w,b} \ \frac{1}{2}\|w\|_2^2$$
$$\text{s.t. } \ y_i(w^\top x_i + b) \ge 1(\forall i).$$

The dual is

$$\max_{\alpha} \ \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^\top x_j$$
$$\text{s.t. } \ \alpha_i \ge 0(\forall i), \ \sum_{i=1}^{n} \alpha_i y_i = 0.$$

Suppose the optimal solution for the dual is $\alpha^* = (\alpha_1^*, \alpha_2^*, \ldots, \alpha_n^*)$, and the optimal solution for the primal is $(w^*, b^*)$. Show that the margin

$$\gamma = \min_i \frac{y_i((w^*)^\top x_i + b^*)}{\|w^*\|_2}$$

satisfies

$$\frac{1}{\gamma^2} = \sum_{i=1}^{n} \alpha_i^*.$$

Hint: use the KKT conditions.
Solution goes here.