1. **Markov Decision Processes**
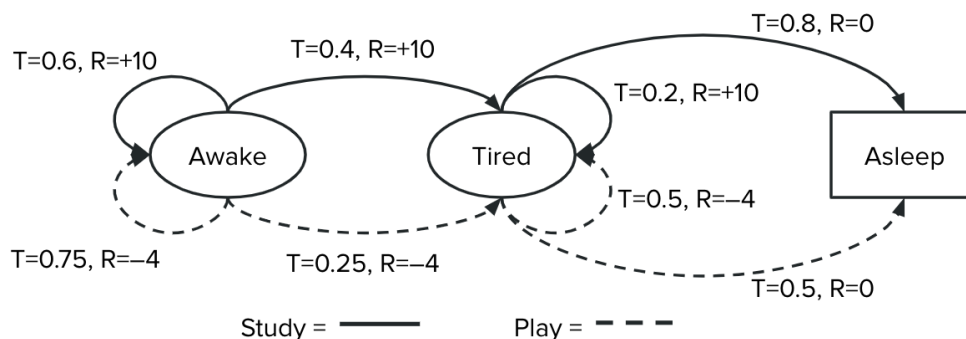
   The Markov Decision Process below models three states − Awake, Tired, and Asleep − and two actions − Study (the solid arrows) and Play (the dashed arrows). The transition value (i.e. probability) and reward for each $(s, a, s')$ triple is provided on the arc of each arrow. The Asleep state is a terminal state with no actions. (Alternatively, you could think of the Asleep state as having a single non-action arc looping back to Asleep with probability 1.0 and reward 0.)

   

   Perform value iteration on this MDP up to $V_2(s)$ for states Awake and Tired (Sleep will always have value 0). Remember that $V_0(s)$ is initialized to be 0 and the update function is given by:

   $$V_{k+1}(s) = \max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma V_k(s)]$$

   To keep things simple, assume $\gamma = 1$.

2. **Reinforcement Learning**

   Consider the following gridworld with states $A, B, G1,$ and $G2$.

   

   From state $A$, possible actions are Right ($\rightarrow$) and Down ($\downarrow$). From state $B$, possible actions are Left ($\leftarrow$) and Down ($\downarrow$). For states $G1$ and $G2$, the only possible action is Exit. Upon exiting, we receive a reward ($+10$ and $+1$ respectively) and end at the end-of-game absorbing state $X$ (which always has value $V(X) = 0$). The discount $\gamma = 1$ and there is no noise (i.e. all actions are deterministic).

   Consider the following episodes. The following questions will use various sequences of these episodes as examples.

   Episode 1 ($E1$)

   | $s$ | $a$ | $s'$ | $r$ |
   |-----|-----|------|-----|
   | $A$ | $\downarrow$ | $G1$ | 0 |
   | $G1$ | exit | $X$ | 10 |

   Episode 2 ($E2$)

   | $s$ | $a$ | $s'$ | $r$ |
   |-----|-----|------|-----|
   | $B$ | $\downarrow$ | $G2$ | 0 |
   | $G2$ | exit | $X$ | 1 |

   Episode 3 ($E3$)

   | $s$ | $a$ | $s'$ | $r$ |
   |-----|-----|------|-----|
   | $A$ | $\rightarrow$ | $B$ | 0 |
   | $B$ | $\downarrow$ | $G2$ | 0 |
   | $G2$ | exit | $X$ | 1 |

   Episode 4 ($E4$)

   | $s$ | $a$ | $s'$ | $r$ |
   |-----|-----|------|-----|
   | $B$ | $\leftarrow$ | $A$ | 0 |
   | $A$ | $\downarrow$ | $G1$ | 0 |
   | $G1$ | exit | $X$ | 10 |

(a) Suppose the observed sequence is $\langle E1, E2, E3, E4 \rangle$. Determine the value $V(s)$ of each state using **direct evaluation** under policy $\pi$ where $\pi(A) = $ Down and $\pi(B) = $ Down.

(b) Suppose the observed sequence is $\langle E3, E4, E1, E2 \rangle$. Determine the value $V(s)$ of each state using **temporal difference learning**. All values are initialized to zero and the learning rate is $\alpha = 0.5$. For reference, here is the update for each sample:

$$\text{sample} = R(s, \pi(s), s') + \gamma V^\pi(s')$$
$$V^\pi(s) = (1 - \alpha)V^\pi(s) + (\alpha)\text{sample}$$

(c) Consider using **Q-Learning** to learn the q-values of this gridworld. For which of the following sequences, if repeated an infinite number of times, would the q-values for *all* state-action pairs $(s, a)$ converge to their optimal value $Q^*(s, a)$?

  i. $\langle E1, E2, E1, E2, ... \rangle$

  ii. $\langle E3, E4, E3, E4, ... \rangle$

  iii. $\langle E1, E2, E3, E4, ... \rangle$

3. **Naive Bayes**

Consider the following labeled corpus of text messages, with punctuation and capitalization removed for simplicity:

(Spam) `you have a chance to win $100`
(Spam) `send your love with an exclusive offer`
(Spam) `you have an offer for free tickets`
(Ham) `you up`
(Ham) `i am here`
(Ham) `have you seen crazy stupid love`
(Ham) `yo terron canceled the quiz you wanna get drinks`

For the following problems, you need only represent probabilities with fractions, not computed decimal values.

(a) Compute the following probabilities:

    Prior probability of spam: $P(Y = spam)$
    Prior probability of ham: $P(Y = ham)$
    Probability of the word `you` given spam: $P(W = you|Y = spam)$
    Probability of the word `you` given ham: $P(W = you|Y = ham)$

(b) What is the probability of spam and the probability of ham given the following text: `love you`. That is, what are $P(Y = spam|X = \text{love you})$ and $P(Y = ham|X = \text{love you})$? You answer should be written as a product of probabilities (i.e. you do not have to compute the exact value of the product).

(c) Consider the text: `are you crazy`. What are the values of $P(Y = spam|X = \text{are you crazy})$ and $P(Y = ham|X = \text{are you crazy})$? (Hint: You shouldn't need to write out the entire product to determine this.) Why is this the case and how can we account for the issue observed here?

4. **Multiclass Perceptron**

Suppose for some multiclass classification problem (e.g. genres of music, categories for news articles, etc.) we have three labels, simplified here to the labels 1, 2, and 3. As such, we have three weight vectors, one per label. There are three features, so each vector contains three values. At some iteration during the update process, the weight vectors have the following values:

$$w_1 = [1, 2, -2], w_2 = [3, -2, -1], w_3 = [-1, 2, 4]$$

(a) On the next iteration, we sample a random data point $x_i$ with label $y = 1$ and extract the feature vector $f(x_i) = [2, 2.5, -1]$. Perform the multiclass perceptron update for $x_i$ and compute the values of any updated weight vectors.

(b) After performing the update in part (a), we sample another point $x_j$ with label $y = 2$ and feature vector $f(x_j) = [1, -0.5, 3]$. Perform the multiclass perceptron update for $x_j$ and compute the values of any updated weight vectors.

5. **Bias**

It is easy to think that machine learning algorithms are objectively correct because they are based on data. However, while the algorithm may optimally separate, cluster, or otherwise detect patterns in data, a key thing to keep in mind is how biased your data can be. As an extreme example, if a classifier is trained on a data set that only contains pictures of dogs, it will not be able to recognize pictures of cats.

For each of the following data sets, state how the data set may be biased for the task at hand:

(a) Data: A hand-selected set of student essays for an English class
Task: Autograding English essays

(b) Data: Collection of hand-drawn sketches by people who were told to draw a "shoe"
Task: Classifying a sketch as being a shoe or not

(c) Data: Former presidents of the United States
Task: Predicting who will be president next