

DS310 HW2

ZhiXuan Yang

October 15, 2019

1. Let  $\mathbf{W}_t$  be the weight vector after  $t$  weight updates.

Let  $\mathbf{w}^*$  be such that

$$\forall \mathbf{x}_n \in S^+, \mathbf{w}^* \cdot \mathbf{x}_n \geq \gamma \text{ and } \forall \mathbf{x}_n \in S^-, \mathbf{w}^* \cdot \mathbf{x}_n \leq -\gamma.$$

We also assume that  $\mathbf{w}^* \cdot \mathbf{x} = 0$  passes through the origin.

$\therefore$  training set  $E$  is linearly separable iff

$$\exists \mathbf{w}^*, \text{such that } \forall \mathbf{x}_p \in S^+, \mathbf{w}^* \cdot \mathbf{x}_p \geq \gamma$$

$$\text{and } \forall \mathbf{x}_p \in S^-, \mathbf{w}^* \cdot \mathbf{x}_p \leq -\gamma$$

$\therefore$  whenever  $\mathbf{w}^* \in \mathbb{R}^{n+1}$  and  $\gamma > 0$  exist, (Training set  $E$  is linearly separable) (A)

$$\text{Let } \forall \mathbf{x}_n \in S^+, z_n = \mathbf{x}_n, \forall \mathbf{x}_n \in S^-, z_n$$

$$\forall \mathbf{x}_n \in S^-, z_n = -\mathbf{x}_n,$$

$$Z = \{z_n\}$$

$$(\forall \mathbf{x}_n \in S^+, \mathbf{w}^* \cdot \mathbf{x}_n \geq \gamma \text{ and } \forall \mathbf{x}_n \in S^-, \mathbf{w}^* \cdot \mathbf{x}_n \leq -\gamma)$$

$$\Leftrightarrow (\forall z_n \in Z, \mathbf{w}^* \cdot z_n \geq \gamma)$$

Let  $E' = \{(z_n, d)\}$ ,  $\mathbf{W}_{t+1} = \mathbf{W}_t + \eta(d_n - y_n)z_n$ , where  $\eta$  is the learning rate and  $d$  is the desired output,  $y$  is the real output., learning rate  $\eta > 0$ , and  $\mathbf{W}_0 = [0 \ 0 \ \dots \ 0]^T$

since [Weight update based on example  $(z_n, d) \Leftrightarrow [(d_n=1) \vee (y_n=-1)]$ ]

$$\therefore \mathbf{w}^* \cdot \mathbf{W}_{t+1} = \mathbf{w}^* \cdot (\mathbf{W}_t + \eta z_n) = (\mathbf{w}^* \cdot \mathbf{W}_t) + \eta (\mathbf{w}^* \cdot z_n)$$

$$\because \forall z_n \in Z, (\mathbf{w}^* \cdot z_n \geq \gamma), \mathbf{w}^* \cdot \mathbf{W}_{t+1} \geq \mathbf{w}^* \cdot \mathbf{W}_t + \eta \gamma$$

$$\therefore \forall t, \mathbf{w}^* \cdot \mathbf{W}_t \geq 2 + \eta \gamma \quad (B)$$

$$\|W_{t+1}\|^2 = W_{t+1} \cdot W_{t+1}$$

$$= (W_t + 2nZ_H) \cdot (W_t + Z_H Z_H)$$

$$= (W_t \cdot W_t) + 4n(W_t \cdot Z_H) + 4n^2(Z_H \cdot Z_H)$$

Note that given any finite training set, the length of the training patterns is bounded. That is,  $\forall z_H, \|z_H\| \leq L$ . where  $L = \max_{x_{HS}} \|x_H\|$  and  $S = S^+ \cup S^-$

$$\therefore \|W_{t+1}\|^2 \leq \|W_t\|^2 + 4n^2 \|Z_H\|^2 \leq \|W_t\|^2 + 4n^2 L^2$$

$$\text{Hence } \|W_t\|^2 \leq 4n^2 L^2$$

$$\therefore \forall t \quad \|W_t\| \leq 2nL\sqrt{t} \quad (C)$$

From (b) we have:  $\forall t \quad \hat{W}_t (W^* \cdot W_t) \geq 2tn\delta \Rightarrow \{\forall t \quad 2tn\delta \leq W^* \cdot W_t\} \Rightarrow$

$$\{\forall t \quad 2tn\delta \leq \|W^*\| \|W_t\| \cos\theta\} \Rightarrow \{\forall t \quad 2tn\delta \leq \|W^*\| \|W_t\|\}$$

$$\because \forall \theta \quad \cos\theta \leq 1, \text{ substituting from for an upper bound on } \|W_t\| \text{ from (C)}$$

$$\forall t \quad \{2tn\delta \leq \|W^*\| 2nL\sqrt{t}\} \Rightarrow \{\forall t \quad (\delta\sqrt{t} \leq \|W^*\| L)\} \Rightarrow t \leq \left(\frac{\|W^*\| L}{\delta}\right)^2 \quad (D)$$

From formula (D), we have prove that the bound on the number of weight updates does not depend on the learning rate, and the training set is linear separable from (A). This convergence theorem holds for any bounded learning rate  $n > 0$ , and there is a condition that " $0 < A \leq n(t) \leq B$ ", thus this theorem is always suitable for this case.

Above, we have shown that learning rate does not directly influence the number of weight updates. Now, we have to prove that variant learning rates is between  $0 < A \leq n(t) \leq B$ , where A and B are fixed lower and upper bounds.

The Perception algorithm changes  $W$  as follows:

$$W_{t+1} \leftarrow W_t + n X_t$$

Suppose we choose  $n$  such that at each weight update,

$$W_{t+1} X_t = (W_t + n_t X_t) \cdot X_t$$

$W_{t+1} \cdot X_t > 0$  if  $W_t \cdot X_t < 0$  and  $C_t = 1$ .

$$\therefore n_t = \begin{cases} 1 & W_t \cdot X_t \leq 0 \\ \frac{W_t \cdot X_t}{X_t \cdot X_t} & X_t \cdot X_t \neq 0 \end{cases} \quad (E)$$

Since the training set is finite, therefore, this equation proves that  $0 < A \leq n(t) \leq B$ , where  $A$  and  $B$  are fixed lower and upper bounds respectively.

All in All, We have prove that the training data are linearly separable from (A), and the number of weight updates does not depend on learning rate from (D), and  $0 < A \leq n(t) \leq B$  from (E), Thus, the resulting variant of the perception learning algorithm is guaranteed to terminate with a weight vector.

Additionally, the perception algorithm is quite robust as a learning model. It can be shown that the algorithm converges even when  $n$  is allowed to fluctuate arbitrarily over time as long as  $0 < n_t \leq B$  where  $B$  is the upper bound on the learning rate, since We have prove that  $0 < A \leq n(t) \leq B$ , therefore, this proves that  $n$  is allowed to fluctuate arbitrarily.

2. Ordered 4-tuples of attributes values corresponding to

outlook (sunny, overcast, rain)

Temperature (hot, mild, cool)

Humidity (high, normal)

Wind (strong, weak)

and we have a class label PlayTennis, and it corresponding to

PlayTennis (Yes(+), No(-))

and we have 14 instances from Training Data.

Instance

class label

$$\hat{H}(X) = -\frac{5}{14} \log_2 \frac{5}{14} - \frac{9}{14} \log_2 \frac{9}{14} \approx 0.940$$

I<sub>1</sub>(S, h, i, w)

-

$$\hat{H}(X | \text{outlook} = S) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

I<sub>2</sub>(S, h, i, t)

-

$$\hat{H}(X | \text{outlook} = O) = -\frac{0}{4} \log_2 \frac{0}{4} - \frac{4}{4} \log_2 \frac{4}{4} = 0$$

I<sub>3</sub>(O, h, i, w)

+

$$\hat{H}(X | \text{outlook} = R) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

I<sub>4</sub>(R, m, i, w)

+

$$\hat{H}(X | \text{outlook}) = \frac{5}{14}(0.971) + \frac{4}{14}(0) + \frac{5}{14}(0.971) = 0.694$$

I<sub>5</sub>(R, c, n, w)

-

$$(I(\text{outlook}) = 0.940 - 0.694 = 0.246)$$

I<sub>6</sub>(R, c, n, t)

-

$$\hat{H}(X | \text{Temperature} = h) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$$

I<sub>7</sub>(O, c, h, t)

+

$$\hat{H}(X | \text{Temperature} = m) = -\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6} = 0.918$$

I<sub>8</sub>(S, m, i, w)

-

$$\hat{H}(X | \text{Temperature} = C) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.811$$

I<sub>9</sub>(S, c, h, w)

+

$$\hat{H}(X | \text{Temperature}) = \frac{4}{14}(1) + \frac{6}{14}(0.918) + \frac{4}{14}(0.811) = 0.911$$

I<sub>10</sub>(R, m, h, w)

+

$$(I(\text{temperature}) = 0.940 - 0.911 = 0.029)$$

I<sub>11</sub>(S, m, n, t)

+

$$\hat{H}(X | \text{Humidity} = i) = -\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7} = 0.985$$

I<sub>12</sub>(O, m, i, t)

-

$$\hat{H}(X | \text{Humidity} = n) = -\frac{1}{7} \log_2 \frac{1}{7} - \frac{6}{7} \log_2 \frac{6}{7} = 0.592$$

$$\hat{H}(\text{Humidity}) = \frac{7}{14}(0.985) + \frac{7}{14}(0.592) = 0.789$$

$$(I(\text{Humidity}) = 0.940 - 0.789 = 0.151)$$

$$\hat{H}(X|\text{Wind} = t) = -\frac{3}{6}\log_2 \frac{3}{6} - \frac{3}{6}\log_2 \frac{3}{6} = 1$$

$$\hat{H}(X|\text{Wind} = w) = -\frac{6}{8}\log_2 \frac{6}{8} - \frac{2}{8}\log_2 \frac{2}{8} = 0.811$$

$$\hat{H}(X|\text{Wind}) = \frac{1}{4} * (1) + \frac{1}{4} (0.811) = 0.892$$

$$I(\text{Wind}) = 0.940 - 0.892 = 0.048$$

Based on the results, Outlook has the largest information gain, therefore Outlook is the root of this tree.

since overcast of Outlook is pure, which means it has 1 for Yes, and 0 for No, and then we don't need to think about this.

Then, we have find the child under sunny,

$$\hat{H}(\text{temperature} = h|S) = -\frac{2}{2}\log_2 \frac{2}{2} - \frac{0}{2}\log_2 \frac{0}{2} = 0.$$

$$\hat{H}(\text{temperature} = m|S) = -\frac{1}{2}\log_2 \frac{1}{2} - \frac{1}{2}\log_2 \frac{1}{2} = 1.$$

$$\hat{H}(\text{temperature} = c|S) = -\frac{0}{1}\log_2 \frac{0}{1} - \frac{1}{1}\log_2 \frac{0}{1} = 0$$

$$I(\text{temperature}|S) = 0.971 - \frac{2}{5} * 0 - \frac{2}{5} * 1 - \frac{1}{5} * 0 = 0.571$$

$$\hat{H}(\text{Humidity} = i|S) = -\frac{3}{3}\log_2 \frac{3}{3} - \frac{0}{3}\log_2 \frac{0}{3} = 0.$$

$$\hat{H}(\text{Humidity} = n|S) = -\frac{0}{2}\log_2 \frac{0}{2} - \frac{2}{2}\log_2 \frac{2}{2} = 0$$

$$I(\text{Humidity}|S) = 0.971 - 0 - 0 = 0.971$$

$$\hat{H}(\text{Wind} = t|S) = -\frac{1}{2}\log_2 \frac{1}{2} - \frac{1}{2}\log_2 \frac{1}{2} = 1$$

$$\hat{H}(\text{Wind} = w|S) = -\frac{2}{3}\log_2 \frac{2}{3} - \frac{1}{3}\log_2 \frac{1}{3} = 0.918$$

$$I(\text{Wind}|S) = 0.971 - (\frac{2}{5}) * 1 - (\frac{3}{5})(0.918) = 0.0202$$

Therefore, Humidity is the child of sunny because it has largest information gain.

Next, we have to find child of rain.

$$H(\text{temperature} = h|r) = 0$$

$$H(\text{temperature} = m|r) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.918$$

$$H(\text{temperature} = c|r) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$I(\text{temperature}|r) = 0.971 - \frac{3}{5}(0.918) - \left(\frac{2}{5}\right)(1) = 0.0202$$

$$H(\text{Humidity} = i|r) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$H(\text{Humidity} = n|r) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.918$$

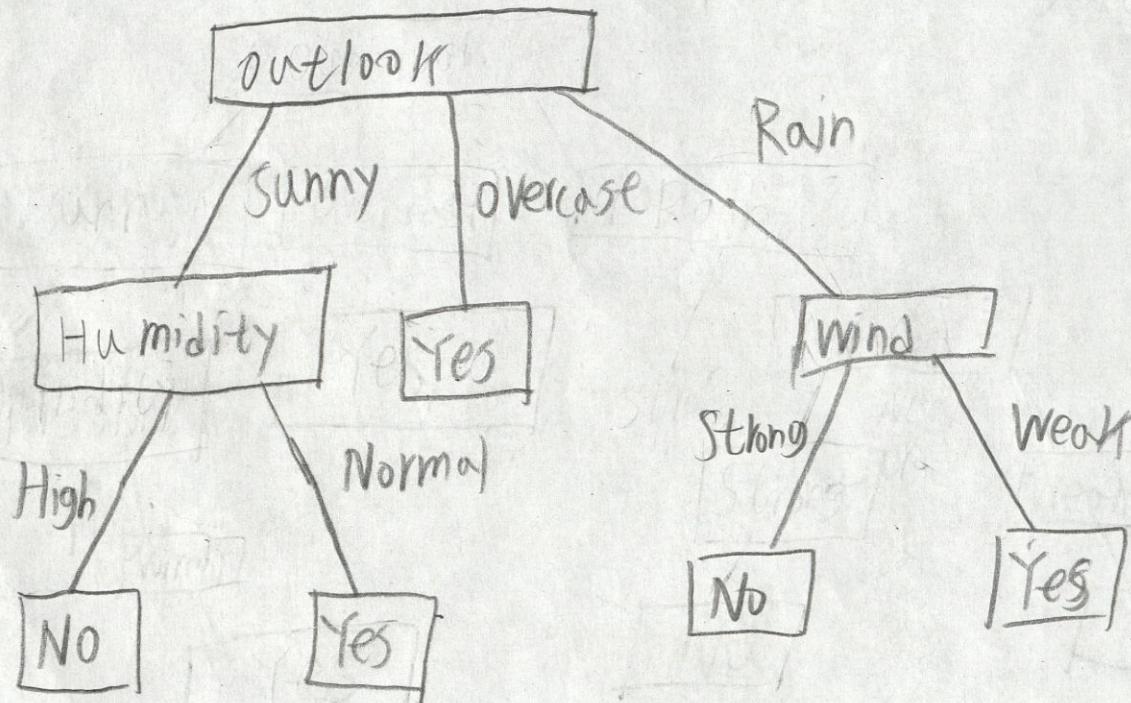
$$I(\text{Humidity}|r) = 0.971 - \frac{3}{5}(0.918) - \frac{2}{5}(1) = 0.0202$$

$$H(\text{wind} = t|r) = -\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$$

$$H(\text{wind} = w|r) = -\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$I(\text{wind}|r) = 0.971 - 0 - 0 = 0.971$$

Based on the results, wind is the child of rain because it has biggest information gain. Finally, our decision tree should be the following



3. To construct a decision tree to minimize the cost of correct classification of data samples, my assumption is that we use marginal information gain instead of information gain. The function of marginal information gain is  $\frac{I(F)^2}{C(F)}$ , where  $C$  is the cost of measuring feature  $F$ , and

$I$  is the information gain of feature  $F$ . The run-time cost of a feature at a position is the sum of cost of seeing the feature and the cost of moving to its position. And below is the original information gain of each feature we got from Problem 2.

$$I(\text{outlook}) = 0.246 \quad I(\text{temperature}) = 0.029 \quad I(\text{Humidity}) = 0.15$$

$$I(\text{wind}) = 0.048$$

And cost of each feature is given.

$$\text{Cost}(\text{outlook}) = 1 \quad \text{cost}(\text{temperature}) = 2 \quad \text{cost}(\text{Humidity}) = 2$$

$$\text{cost}(\text{wind}) = 3$$

The marginal information gains of each feature are below

$$MI(\text{outlook}) = \frac{(0.246)^2}{1} = 0.0605 \quad MI(\text{temperature}) = \frac{0.029^2}{2} = 0.00042$$

$$MI(\text{Humidity}) = \frac{(0.15)^2}{2} = 0.0114 \quad MI(\text{wind}) = \frac{0.048^2}{3} = 0.00077$$

Based on the above results, outlook is the root of cost-effective decision tree because it has biggest marginal information gain.

since overcast of outlook has 4 Yes and 0 No, so, the "Yes" would be the branch of the tree and stop here.

Let's find the first child of Sunny. and we have already got the information gain under sunny, and they are,

$$I(\text{temperature}/S) = 0.571 \quad I(\text{Humidity}/S) = 0.971 \quad I(\text{wind}/S) = 0.0202$$

And because of the formula of marginal information gain  $MI(S/A)$ , the marginal information gain of each attribute under sunny is,

$$MI(\text{temperature}/S) = \frac{(0.571)^2}{1} = 0.326 \quad MI(\text{Humidity}/S) = \frac{(0.971)^2}{1} = 0.943$$

$$MI(\text{Wind} | S) = \frac{(0.0202)^2}{1} = 0.00041$$

Based on the above results, Humidity should be the child of sunny.  
Let's calculate the child of rain. the original information gains of each features under rain from problem 2. are.

$$I(\text{temperature} | r) = 0.0202, I(\text{Humidity} | r) = 0.0202, I(\text{wind} | r) = 0.971$$

And, the marginal information gains are

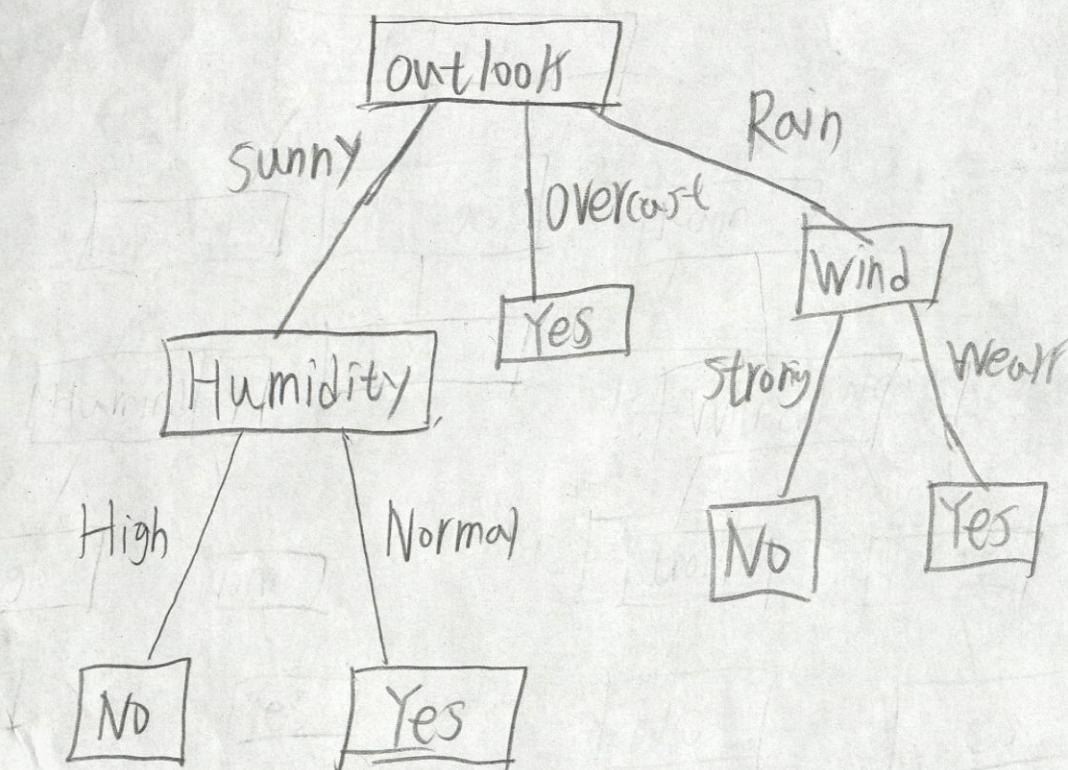
$$MI(\text{temperature} | r) = \frac{(0.0202)^2}{1} = 0.00041$$

$$MI(\text{Humidity} | r) = \frac{(0.0202)^2}{1} = 0.00041$$

$$MI(\text{Wind} | r) = \frac{(0.971)^2}{1} = 0.943$$

Therefore, Wind should be the child of rain because it has the largest marginal information gain.

Therefore, the cost-effective decision tree should look like the following



4. (a) I think Random Forest is the best choice for this problem because of following reasons. (1) the data is mixed, which means this dataset has either numeric variable or categorical, and random forest algorithm is good at deal with this situation. (2) the task is to perform accurate diagnosis of patients, and random forest is suitable for this kind of classification. (3) Additional task is that "to obtain insight regarding the relationships between features for different diseases.", and random forest learning algorithm has an advantage that it provides the feature importance after training, therefore, we can obtain the insight from feature importance using random forest. Above three reasons are why I prefer random forests algorithm for this problem.

### Random Forests.

(b) I think the random forests algorithm is also applicable for this problem because of following reasons. (1) Since the dataset is very large the random forest algorithm would divide the very large dataset into several subsets, therefore, we don't have to worry about the size problem. (2) our task is predicting users' interest. Thus, we can predict the user's interest based on the decision trees created by random forest. (3) Variables could be mixed, and Random forests algorithm is suitable for processing mixed data. Overall, the choice of algorithm should depending on the real situation. But, ideally, I think random forests algorithm is the best choice for this problem.

### Random Forests.

(C) I think Decision Tree is the best choice for this case. First, the problem says "they can automate the decision to approve or deny a loan based on a simple rule," and the key word is simple. Among the four algorithms, perceptron algorithm and naive bayes algorithm require mathematical inference, so they do not meet the "simple rule" requirement. Thus, we have to compare decision tree and random forest. Another requirement of this problem is the decision-making process be transparent, and one of disadvantage of random forest is that this algorithm is very difficult to interpret the resulting model because of its randomness. For instance, when a client asks us why we deny his request, we can show the model of decision tree algorithm to him, and tell him which part caused his request to be rejected. If we are facing the same problem using random forest, we wouldn't be able to quickly find out which part was causing this result because random forest is a collection of many decision trees, and then we don't have an easy way to figure out which tree made this prediction. All in all, the "simple rule" requirement and the "decision-making process be transparent" requirement lead to decision tree algorithm is the best choice for this problem.

