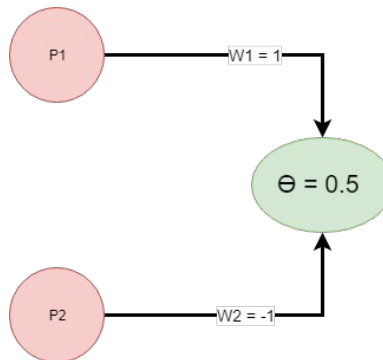


# Problem Set 1

[smatthews39@gatech.edu](mailto:smatthews39@gatech.edu)

Sean Matthews

**2a)** Design a two-input perceptron that implements the boolean function  $A \wedge \neg B$

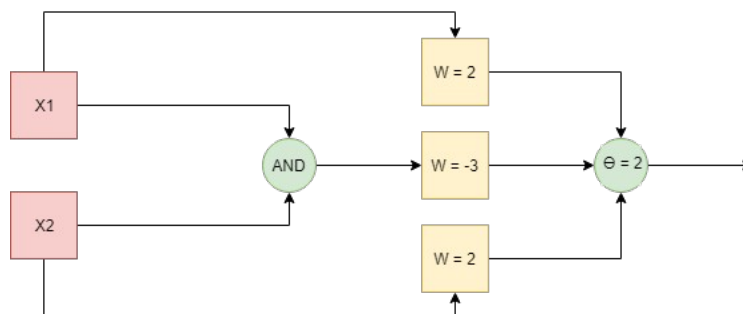


**2b)** Design a two-layer network of perceptrons that implements  $A \oplus B$  ( $\oplus$  is XOR).

$A \oplus B$  - A XOR B

$$A \text{ XOR } B = (A \wedge \neg B) \vee (B \wedge \neg A)$$

X1	X2	AND	OR	XOR
0	0	0	0	0
0	1	0	1	1
1	0	0	1	1
1	1	1	1	0



3)

Derivation problem set 1

Thursday, February 14, 2019 4:07 PM

$$0 = w_0 + w_1 x_1 + w_1 x_1^2 + \dots + w_n x_n + w_n x_n^2$$

gradient of  $\nabla E$  of error  $E$  is a vector

$$\nabla E(w) = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

remember...  $\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$

$$\Delta \vec{w} = -\eta \nabla E(\vec{w})$$

outputs:

$o_d$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \Rightarrow \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} (t_d - o_d)^2$$

$$\frac{1}{2} \sum_{d \in D} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$\sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial w_i} \left( t_d - (w_0 + w_1 x_1 + w_1 x_1^2 + \dots + w_n x_n + w_n x_n^2) \right)$$

$$\sum_{d \in D} (t_d - o_d) \left( - (x_1 + x_1^2 + \dots + x_n + x_n^2) \right)$$

$$= \sum_{d \in D} (t_d - o_d) (-x_{id} - x_{id}^2)$$

4) A decision tree can be used for regression by replacing information gain with standard deviation reduction. Standard deviation is used for creating the tree, and the other parameters are Coefficient of Deviation which is used to tell the model when to stop branching, and Average, which is the value within the leaf nodes.

Standard deviation reduction is based on the decrease in standard deviation after a dataset is split on an attribute. The construction of a decision tree is based on finding the attribute that has the most alike branches. This would signify that it returns the highest standard deviation reduction.

The process of using Decision trees for reduction are below:

1. Calculate the Standard Deviation of the target
2. Split the data on different attributes and calculate the STD for each branch. Subtract that STD from the STD calculated before the split and the result is the STD reduction
3. The LARGEST STD reduction is the chosen decision node
4. Split the data by the selected attribute, and recursively run the data for the next level on the decision tree

On each leaf, a squared error function is equal to using the mean of the class for each instance, therefore, the mean of Y is used as the expected value at each leaf:

$$\sum (Y_i - h_i)^2 = \sum (\text{mean}(Y) - h_i)^2$$

```
1 import pandas as pd
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import mean_squared_error
5
6
7 data = pd.read_csv("housing.csv")
8 X = data.iloc[:,0:13]
9 y = data.iloc[:,13]
10 training_size= 60
11
12 print("Training Model and testing with: ", 100-training_size, "%")
13 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0, test_size=100-training_size)
14 regr_1 = DecisionTreeRegressor(random_state=0, max_depth=4)
15 regr_1.fit(X_train, y_train)
16 y_1 = regr_1.predict(X_test)
17 print("Error with Max Depth 4: ", mean_squared_error(y_1, y_test))
18
19 regr_2 = DecisionTreeRegressor(random_state=0, max_depth=5)
20 regr_2.fit(X_train, y_train)
21 y_2 = regr_2.predict(X_test)
22 print("Error with Max Depth 5: ", mean_squared_error(y_2, y_test))
23
24 regr_3 = DecisionTreeRegressor(random_state=0, max_depth=6)
25 regr_3.fit(X_train, y_train)
26 y_3 = regr_3.predict(X_test)
27 print("Error with Max Depth 6: ", mean_squared_error(y_3, y_test))
28 |
```

```
C:\Users\Sean\PycharmProjects\DTRegression\venv\Scripts\python.exe C:/Users/Sean/PycharmProjects/DTRegression/housing_data.py
Training Model and testing with: 40 %
Error with Max Depth 4: 7.326168811753211
Error with Max Depth 5: 7.234786420281537
Error with Max Depth 6: 7.617227360554128

Process finished with exit code 0
```

5)

ID3 (training examples, class attributes, attributes):

Loop:

- $A \leftarrow$  best attribute
- Assign A as decision attribute for node
- For Each Value of A create a descendant of node
- Sort training examples to leaves
- If examples perfectly classified, STOP
- Else iterate over leaves

Lazy DT (unclassified instance, training examples)

- Select a test, X (the one that maximizes the decreases in entropy)
- Let x be the value of the test on the unclassified instance
- Assign the set of instances with  $X = x$  to training examples
- Iterate (unclassified instance, training examples)

This Lazy DT takes advantage of the additional information given by a single instance and fixes some problems that occur in the regular DT algorithms such as replication and fragmentation. Fragmentation is where subsequent tests on some attributes have lower significance due being based on fewer instances. Lazy decision trees select better tests with respect to a particular instance and it avoids splitting on unnecessary attributes which produce a more accurate classifier and shorter paths. Another advantage of lazy DTs is they never branch on a missing value in the instance

Lazy learning requires significantly less time during training than eager methods, but they require more time to classify each query instance. They also require large space requirements to store the entire training dataset.

6) I would use kNN. Decision trees would have a hard time determining the classification for those which are on the line  $y=x$ . They are both nonlinear classifiers, but I think kNN, would be able to determine the classification

pretty well whether it was on the (+) side or the (-), as well as on the line  $y=x$ . This is true for a dense amount of data as well as few samples.

7)

1. An origin-centered circle (2D)

$$VC(H) = 3$$

Figure 7.3 in the book represents a 2D space, and shows  $2^3$  hypothesis shattering the 3 points. For each of the dichotomies, there must exist a hypothesis that covers it. In the case where the 3 points are colinear, there isn't a way to find  $2^3$  linear surfaces that shatter, however, the VC dimension definition states if there are *ANY* set of instances of size  $d$  that can be shattered, then  $VC(H) \geq d$

However after rereading the questions... I'm picturing that the circle cannot be moved, it can just be grown (considering it is origin-centered). If the circle can't be moved, then I believe it would be  $VC(H)=2$ , I may not be understanding the questions then.

2. I would also say for a sphere it is  $VC(H) = 3$

I can't come up with a solution where there are 4 points in 3D space that can be shattered. Therefore it is still the same as a 2D space of  $VC(H)=3$