

CS148 HW3

1 A) Entropy = $-\left[p(0) \times \log p(0) + p(1) \times \log p(1)\right]$

$p(0) = 3/8 \quad p(1) = 5/8$

Entropy = $-\left[\left(\frac{3}{8} \times \log\left(\frac{3}{8}\right)\right) + \left(\frac{5}{8} \times \log\left(\frac{5}{8}\right)\right)\right]$

Entropy = 0.954

B) CHEST PAIN Information Gain

YES - E = 0

NO - E = 0.811

Gain = 0.954 - 0.4055

$\left(\frac{1}{2} \times 0\right) + \left(\frac{1}{2} \times 0.811\right) = 0.4055$

= 0.5485

MALE Information Gain

YES - Entropy = 1

NO - E = 0.811

Gain = 0.954 - 0.9055

$\left(\frac{1}{2} \times 1\right) + \left(\frac{1}{2} \times 0.811\right) = 0.9055$

= 0.0485

SMOKES Information Gain

YES - Entropy = 0.722

NO - E = 0.918

Gain = 0.954 - 0.7955

$\left(\frac{5}{8} \times 0.722\right) + \left(\frac{3}{8} \times 0.918\right) = 0.7955$

= 0.1585

EXERCISES Information Gain

YES - E = 0.971

NO - E = 0

Gain = 0.954 - 0.6069

$\left(\frac{5}{8} \times 0.971\right) + \left(\frac{3}{8} \times 0\right) = 0.6069$

= 0.3471

I would split on CHEST PAIN, since it has highest information gain. (0.5485)

C) ① Stop splitting into more regions if all instances in a region belong in the same class

② Stop splitting if all of the information gains are negative/less than a certain threshold

③ Stop splitting if the number of instances in a region falls below the pre-defined threshold

④ Stop splitting if the total number of regions exceeds a pre-defined threshold

⑤ Stop splitting if the points in a region are independent of the predictors

D) Standardizing or normalizing shouldn't be necessary for decision trees

as such transformations will not affect how a decision tree is created.

We look at Entropy/Information gain, neither of which are affected by standardization/normalization.

E) Decision trees are generally robust to outliers. We split based on observed

entropy or Gini Index. Neither of these measurements should really be affected by outliers, meaning that the decision trees can still work fine with a few outliers.

$$W = 0$$

- 2 a)
- X_1 $W = y_1 X_1$
 - X_2 $W = y_1 X_1$
 - X_3 $W = y_1 X_1 + y_3 X_3$
 - X_4 $W = y_1 X_1 + y_3 X_3 + y_4 X_4$
 - X_5 $W = y_1 X_1 + y_3 X_3 + y_4 X_4$

After the training epoch, $W = y_1 X_1 + y_3 X_3 + y_4 X_4$

- b) Based on my answer from part a, I got that $W = y_1 X_1 + y_3 X_3 + y_4 X_4$
 $W = \langle 1, 1, 0 \rangle - \langle 1, 1, -3 \rangle + \langle 1, 3, -1 \rangle$ with the values from the table
 $W = \langle 1, 3, 2 \rangle$

Does $y_i = \text{sgn}(W^T x_i)$? If yes, then the prediction is correct (+1)
 If no, then the prediction is wrong (-1)

$$(110) \begin{pmatrix} 1 \\ 3 \\ 2 \end{pmatrix} = 4$$

$$y_1 = \text{sgn}(4)$$

$$+1 = 1 \quad \checkmark$$

Model predicts $X_1 = 1$, which is correct.

- c) The activation function for this Perceptron model is called sgn ,
 whereas in class the professor said that traditionally a
 Step function is used (for the logistic regression model).
 If the perceptron in this problem used the sigmoid activation
 function, it would basically be the same as the model
 discussed during class.

3. a) For the hidden layers, we could use Sigmoid, tanh, ReLU (and its variations), Softplus, or Swish. Sigmoid outputs a value between 0 and 1, which makes sense for binary classification. ReLU can regulate the input by eliminating negative inputs and acting as a linear function for non-negative inputs. The variations of ReLU have the same basic idea but have minor advantages (ie Exponential ReLU is differentiable at all points, tanh is a compressed version of sigmoid, so it follows the same idea but has outputs between -1 and 1 instead. This ^{and can be a bit faster} would still work fine for this binary classification. Softplus and Swish are similar to ReLU so they would also work fine, but they are typically used for large neural networks. All of the activation functions listed above could work in the hidden layer(s) of a binary classification neural network; it really depends on the context of the problem. Personally, I would try using ReLU or sigmoid first, then possibly experiment with other functions.

For the output layer, we want outputs to be binary for the sake of binary classification. Activation functions such as sigmoid or logistic would probably work best since they output from 0 to 1. A modified version of tanh that outputs from 0 to 1 may be used as well. Personally, I would try using sigmoid for a binary classification problem and then experiment with other activation functions if the results are unsatisfactory.

b) $x_1 = 2$ $x_2 = -3$ $y = 1$

$$w_{11} = 0.9 \quad w_{12} = 0.4 \quad w_{21} = -1.5 \quad w_{22} = -0.7 \quad w_{31} = -0.2 \quad w_{32} = 1.6$$

$$w_{10} = 0 \quad w_{20} = 0 \quad w_{30} = 0$$

$$\text{Neuron 1: } w_1 = (x_1 \cdot w_{11}) + (x_2 \cdot w_{12})$$

$$= (2 \cdot 0.9) + (-3 \cdot 0.4) = 0.6 \rightarrow \text{ReLU} \rightarrow 0.6$$

$$\text{Neuron 2: } w_2 = (x_1 \cdot w_{21}) + (x_2 \cdot w_{22})$$

$$= (2 \cdot -1.5) + (-3 \cdot -0.7) = -0.9 \rightarrow \text{sigmoid} \rightarrow 0.289$$

$$y = (w_1 \cdot w_{31}) + (w_2 \cdot w_{32}) = (0.6 \cdot -0.2) + (0.289 \cdot 1.6)$$

$$y = 0.34 \quad \text{Output} = 0.34$$

c) Binary Cross-Entropy Loss Function

$$L(W; X, Y) = -\sum [y_i \log p_i + (1-y_i) \log (1-p_i)]$$

$$L = -[(1)(\log 0.34) + (1-1)(\log 0.66)] = 1.5564$$

Basically: $L = -\ln(\hat{y})$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{1}{\hat{y}} = -2.94$$

$$d) L = -\ln \left(\underbrace{-0.2}_{W_{s1}} \cdot \underbrace{\max(0, 0.9 \cdot 2 - 3W_{12})}_{\text{Neuron 1 (ReLU)}} + \underbrace{1.6}_{W_{s2}} \cdot \underbrace{\frac{1}{1 + e^{-(-1.5 \cdot 2 + -0.7 \cdot -5)}}}_{\text{Neuron 2 (Sigmoid)}} \right)$$

$$L = -\ln(-0.2 \cdot \max(0, 1.8 - 3W_{12}) + 0.4625)$$

If $W_{12} < 0.6$: $L = -\ln(0.6x + 0.1025)$

$W_{12} \geq 0.6$: $L = -\ln(0.4625)$

From part b we see that $W_{12} = 0.4$

$$L = -\ln(0.6x + 0.1025)$$

$$\frac{\partial L}{\partial W_{12}} = \frac{1}{0.6W_{12} + 0.1025} \cdot 0.6$$

$$-1.75$$

e) The network has a total of 9 parameters (6 weights + 3 bias terms).

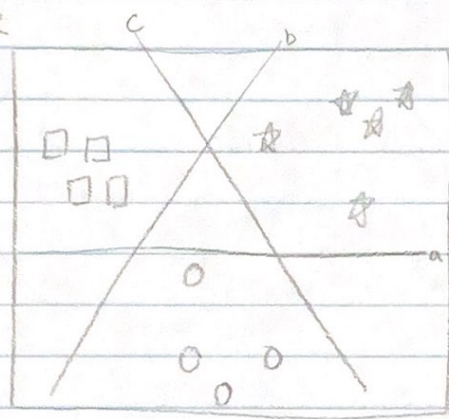
4. a) In logistic regression, the number of parameters is :

$$\# \text{ of features} + 1 (\text{bias term}) = 20 + 1 = 21 \text{ parameters}$$

$$\text{Total number of parameters} = 5 (\text{classes}) \times 21 = 105$$

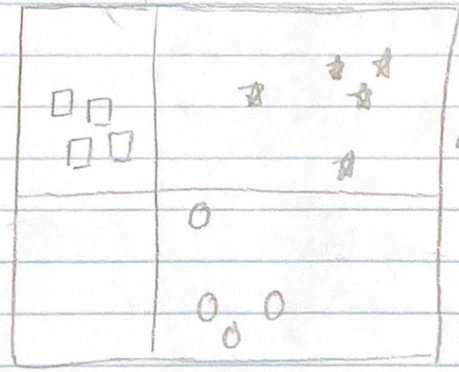
We need 105 total parameters for a multi-class classification problem with 5 classes and 20 features.

b) OVR



- a - separates circles from squares + stars
- b - separates squares from circles + stars
- c - separates stars from circles + squares

Multinomial



Use the square as the reference group
 Line 1: predict $y \neq \text{square}$ from $y = \text{square}$
 Line 2: predict $y = \text{circle}$ from $y \neq \text{square}$

square vs star
 and
 square vs circles

Note : I realize my sketches of the multi-class classification aren't perfect, but hopefully you get the general idea.

6

3 > Neural Network
4ReLU - piecewise
straight line
tanh - much
smoother

5. (1) C

Decision Trees
makestraight,
hard
boundaries

Decision Tree shouldn't generate a curve like A or B. In general,
Decision trees should have very clear-cut boundaries

2) D

Random Forest should look similar to decision tree, but the
boundaries are a bit less clear-cut as many decision
trees are averaged out to get the final model

3) B

Neural
Network
is a
curve

ReLU is more of a piecewise straight line, and this is
reflected in the decision boundary of B

4) A

tanh is a much smoother function than ReLU, and this is
reflected in the decision boundary.