

HWS ECE498 Spring 2020
 Netid:- GOWTHAM4

Problem1:- Jupyter Notebook has been submitted online.

It contains codes for
 SVM, DecisionTree, Randomforest classifiers
 and 5-fold cross validation results (sklearn)

1) (Mean Accuracy, S.D Accuracy)

	<u>SVM</u>	<u>DT</u>	<u>RF</u>
SVM C=0.1	d=3		N=5
C=1	d=4		N=11
DT C=10	d=6		N=13
RF			

(mean Precision, S.D Precision)

	<u>SVM</u>	<u>DT</u>	<u>RF</u>
C=0.1	d=3		N=5
C=1	d=4		N=11
C=10	d=6		N=13

(mean Recall, s.d Recall)

	<u>SVM</u>	<u>DT</u>	<u>R F</u>
$c=0.1$	$d=3$		$N=5$
$c=1$	$d=4$		$N=11$
$c=10$	$d=6$		$N=13$

2) Label imbalances cause accuracy to be high. Even if most are predicted to of class 0 (not pos/neg), accuracy will be $> 90\%$. Recall/Precision take TP into account

See Jupyter Notebook

3) Best classifier: SVM $c=10$ we used F1-Score weighted \rightarrow takes care of label imbalance

Problem 2:- Activation function
 $f(x) = \frac{1}{1+e^{-x}}$ (sigmoid function)

$$g'(x) = g(x)(1-g(x))$$

1)

$$z_1 = w_1 \cdot x_1 + b_1$$

Vector
(dot products)

$$z_4 = w_4 \cdot a_1 + w_5 \cdot a_2 + w_6 \cdot a_3 + b_4$$

$$a_1 = \frac{1}{1+e^{-z_1}}$$

$$a_4 = \frac{1}{1+e^{-z_4}}$$

$$\text{Output} = \hat{y} = a_4$$

$$a_4 = \frac{1}{1+e^{-z_4}} = \frac{1}{1+e^{-\sum_{j=1}^3 w_j \cdot a_j + b_4}}$$

$$\hat{y} = a_4 = \frac{1}{1+e^{-\left[w_4 \cdot \left(\frac{1}{1+e^{-w_1 x_1 + b_1}} \right) + w_5 \cdot \left(\frac{1}{1+e^{-w_2 x_2 + b_2}} \right) + w_6 \cdot \left(\frac{1}{1+e^{-w_3 x_3 + b_3}} \right) + b_4 \right]}}$$

Note:- $g(x) = \frac{1}{1+e^{-x}}$ \therefore Sigmoid function

2) loss, $L_i = \left(\overset{\text{Target label}}{\uparrow} y_i - \overset{\text{Predicted label}}{\uparrow} \hat{y}_i \right)^2$

$$\frac{\partial L_i}{\partial b_4} = 2(y_i - \hat{y}_i) \frac{\partial \hat{y}_i}{\partial b_4}$$

$$\hat{y}_i = a_4 = g(z_4) = \frac{1}{1+e^{-z_4}}$$

$$\Rightarrow \frac{\partial \hat{y}_i}{\partial b_4} = g(z_4)(1-g(z_4)) \cdot \frac{\partial z_4}{\partial b_4}$$

$$\Rightarrow \frac{\partial L_i}{\partial b_4} = 2(y_i - \hat{y}_i) g(z_4)(1-g(z_4))$$

Similarly

$$\frac{\partial L_i}{\partial w_4} = \frac{\partial \hat{y}_i}{\partial w_4} 2(y_i - \hat{y}_i)$$

$$\text{also, } \frac{\partial \hat{y}_i}{\partial w_4} = g(z_4)(1-g(z_4)) \frac{\partial z_4}{\partial w_4}$$

$$\Rightarrow \frac{\partial L_i}{\partial w_4} = 2(y_i - \hat{y}_i) g(z_4)(1-g(z_4)) a_1$$

$$\frac{\partial L_i}{\partial b_1} = \frac{\partial L_i}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_4} \frac{\partial z_4}{\partial a_4} \frac{\partial a_4}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

$$\Rightarrow \boxed{\frac{\partial L_i}{\partial b_1} = 2(y_i - \hat{y}_i) [g(z_4)(1-g(z_4))] w_4 [g(z_1)(1-g(z_1))] \eta_i}$$

$$\frac{\partial L_i}{\partial w_1} = \frac{\partial L_i}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_4} \frac{\partial z_4}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$\Rightarrow \boxed{\frac{\partial L_i}{\partial w_1} = 2(y_i - \hat{y}_i) [g(z_4)(1-g(z_4))] w_4 [g(z_1)(1-g(z_1))] \eta_i}$$

3) Gradient descent step:-

$$w_1' = w_1 - \eta \frac{\partial L_i}{\partial w_1} \quad (\text{single } \eta_i)$$

$$\Rightarrow \boxed{w_1' = w_1 - \eta \Delta / \chi}$$

4) Gradient descent multiple samples (n)

$$L = \frac{1}{n} \sum_{i=1}^n L_i = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\frac{\partial L}{\partial w_1} = \sum_{i=1}^n \frac{\partial L_i}{\partial w_1}$$

$$\Rightarrow \boxed{w_1' = w_1 - \eta \sum_{i=1}^n \frac{\partial L_i}{\partial w_1}} \quad \begin{array}{l} \text{given by} \\ \text{per } 9.2.2. \end{array}$$

5) Sigmoid vs ReLU Activation functions

Ref: stackexchange, 126238
wikipedia entries

Advantages of ReLU over Sigmoid

1) Gradient is easier to calculate and is 1 when $a > 0$
if $\text{ReLU}(a) = \max(0, a)$

2) Sparsity, when $a \leq 0$
the resulting computation is sparse because ~~gradient~~ function outputs zero.

★ Sparsity functions better in many cases than dense representations (common in sigmoid)

Advantages of Sigmoid over ReLU

1) ^{Sigmoid} Activation produces output in the range $(0, 1)$ → easy to store. doesn't blow up output of neurons.

2) It is better than ReLU in some cases where large number of neurons die (ex: multiple -ve values in input)