

HWS ECE498 Spring 2020

Netid:- GOWTHAM4

Name:-

Gowtham Kuntumalla

Problem1:- \rightarrow Jupyter Notebook has been submitted online.

\rightarrow It contains codes for

SVM, DecisionTree, Randomforest classifier and k -fold cross validation results (sklearn)

1) $\begin{matrix} (\mu) & (\sigma) \\ \text{(Mean Accuracy, S.D Accuracy)} \end{matrix}$

	SVM		DT		RF
$C=0.1$	(97.78%, 0.00204)	$d=3$	(97.74%, 0.00288)	$N=5$	(97.73%, 0.00184)
$C=1$	(97.87%, 0.0023)	$d=4$	(97.89%, 0.00222)	$N=11$	(97.782%, 0.00217)
$C=10$	(97.88%, 0.00235)	$d=6$	(97.73%, 0.002898)	$N=13$	(97.79%, 0.00175)

$\begin{matrix} (\mu) & (\sigma) \\ \text{(mean Precision, S.D Precision)} \end{matrix}$

	SVM		DT		RF
$C=0.1$	(94.28%, 0.01996)	$d=3$	(89.74%, 0.03719)	$N=5$	(93.66%, 0.01756)
$C=1$	(94.23%, 0.02005)	$d=4$	(91.91%, 0.0296)	$N=11$	(93.39%, 0.0235)
$C=10$	(94.051%, 0.02081)	$d=6$	(90.54%, 0.03978)	$N=13$	(93.23%, 0.0247)

(mean Recall, s.d Recall)

	SVM (M)	DT (T)	RF
c=0.1	(80.66%, 0.0376)	d=3 (85.29%, 0.0303)	N=5 (80.72%, 0.0354)
c=1	(81.88%, 0.03936)	d=4 (84.62%, 0.0333)	N=11 (81.69%, 0.0378)
c=10	(82.12%, 0.0391)	d=6 (84.26%, 0.0311)	N=13 (81.99%, 0.03426)

2) Label imbalances cause accuracy to be high. Even if most are predicted to of class 0 (not pulsar), 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)

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

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

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

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

$$a_4 = \frac{1}{1+e^{-z_4}} = \frac{1}{1+e^{-\sum_{i=1}^3 w_i \cdot a_i + b_4}}$$

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

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

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

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

$$\hat{y}_i = a_4 = g(z_4) =$$

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

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

similarly

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

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

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

$$\frac{\partial L_i}{\partial b_1} = \frac{\partial L_i}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_4} \frac{\partial z_4}{\partial a_4} \frac{\partial a_4}{\partial z_1} \frac{\partial z_1}{\partial b_1} \rightarrow 1$$

$$\Rightarrow \boxed{\frac{\partial L_i}{\partial b_1} = 2(y_i - \hat{y}_i) [g(z_4)(1-g(z_4))] \omega_4 [g(z_1)(1-g(z_1))] \cdot 1}$$

$$\frac{\partial L_i}{\partial \omega_1} = \frac{\partial L_i}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_4} \frac{\partial z_4}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial \omega_1}$$

$$\Rightarrow \boxed{\frac{\partial L_i}{\partial \omega_1} = 2(y_i - \hat{y}_i) [g(z_4)(1-g(z_4))] \omega_4 [g(z_1)(1-g(z_1))] \eta_i}$$

3) Gradient descent step:-

$$\omega_1' = \omega_1 - \eta \frac{\partial L_i}{\partial \omega_1} \quad (\text{Single } \eta_i)$$

$$\Rightarrow \boxed{\omega_1' = \omega_1 - \eta \times \eta_i}$$

4) Gradient descent multiple samples (n)

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

$$\frac{\partial L}{\partial \omega_1} = \sum_{i=1}^n \frac{\partial L_i}{\partial \omega_1}$$

$$\Rightarrow \boxed{\omega_1' = \omega_1 - \eta \sum_{i=1}^n \frac{\partial L_i}{\partial \omega_1}} \rightarrow \text{given by } \text{for } 9.2.2.$$

5) Sigmoid vs ReLU Activation functions

Ref:- stackexchange, 1262
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 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 $(0, 1) \rightarrow$ 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~~ multiple -ve values in input)