

Name: Yashwanth. M

Reg no: 18GAEC9077

Class: VII Sem CSE

Assignment II

1.

- a) Write Gradient descent Algorithm to train a linear unit along with the derivation of gradient descent rule.

The direction of steepest can be found by computing the derivative of E with respect to each component of the \vec{w} . The vector derivative is called the gradient of E with respect to \vec{w} , written as

$$\nabla E[\vec{w}] = \begin{bmatrix} \frac{\partial E}{\partial w_0} & \frac{\partial E}{\partial w_1} & \dots & \frac{\partial E}{\partial w_n} \end{bmatrix}$$

The gradient specifies the direction of steepest increase of E , the training rule for gradient descent is,

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

where

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

$\rightarrow \eta$ is positive learning rate.

This training rule can also be written in its component form

$$w_i \leftarrow w_i + \Delta w_i$$

where,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} \rightarrow \textcircled{1}$$

Calculate the gradient at each step. The vector of $\frac{\partial E}{\partial w_i}$ derivatives that form the gradient can be

obtained by differentiating E from $E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d)$$

$$\frac{\partial E}{\partial w_i} = \sum_d (t_d - o_d) (-x_{id}) \rightarrow \textcircled{2}$$

Substituting $\textcircled{2}$ in $\textcircled{1}$ yields the weight update rule for gradient descent

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{id}$$

Algorithm

Each training example is a pair of the form (\vec{x}, t) , where \vec{x} is the vector of input values, and t is the target output value, η is the learning rate

* Initialize each w_i to some small random value

* Until the termination condition is met, Do

→ Initialize each Δw_i to zero

→ For each (\vec{x}, t) in training examples, Do

→ Input the instance \vec{x} to the unit and compute the output o

→ For each linear unit weight w_i , Do

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i$$

→ → → For each linear unit weight w_i , Do

$$w_i \leftarrow w_i + \Delta w_i$$

b) Write derivation of Backpropagation rule considering unit j as output unit and unit i as hidden unit.

For each training example d every weight w_{ji} is updated by adding to it Δw_{ji}

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} \rightarrow \textcircled{1}$$

where, E_d is the error on training example d ,

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k \in \text{output}} (t_k - o_k)^2$$

t_k is target value of unit k

o_k is output value of unit k

x_{ji} = the i^{th} input to unit j

w_{ji} = the weight associated with the i^{th} input to unit j

$\text{net}_j = \sum_i w_{ji} x_{ji}$

o_j = output

t_j = target

σ = sigmoid f^n

Downstream(j) = set of units whose immediate inputs include the output of unit j

derive an expression for $\frac{\partial E_d}{\partial w_{ji}}$

Using chain rule,

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$= \frac{\partial E_d}{\partial \text{net}_j} x_{ji} \rightarrow (2)$$

Derive a convenient expression for $\frac{\partial E_d}{\partial \text{net}_j}$

Case 1:- Training rule for O/p unit weights.

$$\frac{\partial E_d}{\partial \text{net}_j} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \rightarrow (3)$$

First term, in (3), $\frac{\partial E}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \in \text{output}} (t_k - o_k)^2$

The derivative of $\frac{\partial}{\partial o_j} (t_k - o_k)^2$ is 0 for k except $k=j$

$$\begin{aligned} \frac{\partial E_d}{\partial o_j} &= \frac{\partial}{\partial o_j} \frac{1}{2} (t_j - o_j)^2 \\ &= \frac{1}{2} 2(t_j - o_j) \frac{\partial (t_j - o_j)}{\partial o_j} \\ &= -(t_j - o_j) \rightarrow (4) \end{aligned}$$

Second term, in (3), since $o_j = \sigma(\text{net}_j)$

$$\frac{\partial o_j}{\partial \text{net}_j} = \sigma(\text{net}_j) (1 - \sigma(\text{net}_j))$$

$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial \sigma(\text{net}_j)}{\partial \text{net}_j}$$

$$= o_j (1 - o_j) \rightarrow (5)$$

(4) & (5) in (3)

$$\frac{\partial E_d}{\partial \text{net}_j} = -(t_j - o_j) o_j (1 - o_j)$$

and combining with ① and ②

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} = \eta (t_j - o_j) o_j (1 - o_j) (x_{ji}) \rightarrow \textcircled{7}$$

Case 2: - Training rule for Hidden Unit weights.

- In this case, j is hidden unit
- the derivation of training rule w_{ji} must take into account the indirect ways in which w_{ji} can influence the network outputs.
- we will find it useful to refer to the set of all units immediately downstream of unit j
- net_j can influence the network outputs only through the units in $Downstream(j)$

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in Downstream(j)} \frac{\partial E_d}{\partial net_k} \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in Downstream(j)} -\delta_k \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in Downstream(j)} -\delta_k \frac{\partial net_k}{\partial o_j} \frac{\partial o_j}{\partial net_j}$$

$$= \sum_{k \in Downstream(j)} -\delta_k w_{kj} \frac{\partial o_j}{\partial net_j}$$

$$= \sum_{k \in Downstream(j)} -\delta_k w_{kj} o_j (1 - o_j) \rightarrow \textcircled{8}$$

$$\delta_j = o_j(1-o_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}$$

$$\Delta w_{ji} = \eta \delta_j x_{ji}$$

2

- a) Derive derivation to show maximum likelihood hypothesis has least squared error hypotheses.

Consider the problem of learning a continuous valued target function such as neural network learning, linear regression, and polynomial curve fitting

Using the definition of hML we have

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h)$$

Assuming training examples are mutually independent given h , we can write $P(D|h)$ as the product of the various $\{d_i|h\}$

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m p(d_i|h)$$

Given the noise ϵ_i obeys a normal distribution with 0 mean and unknown variance σ^2 , each d_i must also obey a Normal distribution around the true target value $f(x_i)$.

Because we are writing the expression for $P(D|h)$, we assume h is the correct description of f .

$$\text{Hence } \mu = f(x_i) = h(x_i)$$

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - \mu)^2}$$

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$$

Maximise the less complicated logarithm, which is justified because of the monotonicity of function p

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

The first term in this expression is a constant independent of h , and can therefore be discarded

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

Maximising this negative quantity is equivalent to minimize +ve - || -

$$h_{ML} = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

discard constant

$$h_{ML} = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

Thus, above equation shows that the maximum likelihood hypothesis H_{ML} is the one that minimises the sum of the squared errors between the observed training rules d_i and hypothesis predictions $h(x_i)$

- b. Classify the test data and {Red, SUV, Domestic} using Naive Bayes classifier for the data show in Table

Color	Type	Origin	Stolen
Red	Sports	Domestic	Yes
Red	Sports	Domestic	No
Red	Sports	Domestic	Yes
Yellow	Sports	Domestic	No
Yellow	Sports	Imported	Yes
Yellow	SUV	Imported	No
Yellow	SUV	Imported	Yes
Yellow	SUV	Domestic	No
Red	SUV	Imported	No
Red	Sports	Imported	Yes

$$V_{nb} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) \prod P(a_i | v_j) \rightarrow (1)$$

$$P(a_i | v_j) = \frac{n_c + mp}{n + m} \rightarrow (2)$$

n = no. of training examples for which $v = v_j$

n_c = no. of examples for $v = v_j$ & $a = a_i$

p = a priori estimate for $P(a_i | v_j)$

m = the equivalent sample size
chandra's

We want to classify { Red, Domestic, SUV }

We need to calculate

$P(\text{Red} | \text{Yes})$, $P(\text{SUV} | \text{Yes})$, $P(\text{Domestic} | \text{Yes})$

$P(\text{Red} | \text{No})$, $P(\text{SUV} | \text{No})$, and $P(\text{Domestic} | \text{No})$

and multiply them by $P(\text{Yes})$ and $P(\text{No})$ respectively.

~~Yes:~~

~~Red:~~

~~$n = 5$~~

~~$n_c = 3$~~

~~$p = 0.5$~~

~~$m = 3$~~

~~SUV:~~

~~$n = 5$~~

~~$n_c = 1$~~

color

Stolen		
	$P(\text{Yes})$	$P(\text{No})$
Red	$3/5$	$2/5$
Yellow	$2/5$	$3/5$

Type

Sports	$4/5$	$2/5$
SUV	$1/4$	$3/4$

Origin	Stolen	
	$P(Y_u)$	$P(N_u)$
Domestic	$2/5$	$3/5$
Imported	$3/5$	$2/5$

$$P(Y_u | X) = P(\text{Red} | Y_u) * P(\text{SUV} | Y_u) * P(\text{Domestic} | Y_u) * P(Y_u)$$

$$= \frac{3}{5} * \frac{1}{4} * \frac{2}{5} * \frac{1}{2}$$

$$= 0.03$$

$$P(N_u | X) = P(\text{Red} | N_u) * P(\text{SUV} | N_u) * P(\text{Domestic} | N_u) * P(N_u)$$

$$= \frac{2}{5} * \frac{3}{4} * \frac{3}{5} * \frac{1}{5}$$

$$= 0.036$$

Since $0.036 > 0.03$, it is classified as 'No'.

3) a) Define

i) Sample Error

ii) True Error

iii) Confidence intervals for discrete valued hypothesis.

Sample Error:

The Sample Error of a hypothesis with respect to some sample S of instances drawn from X is the fraction of S that it misclassifies

Definition: The Sample Error ($\text{error}_S(h)$) of hypothesis h with respect to target function f and data sample S is

$$\text{error}_S(h) \equiv \frac{1}{n} \sum_{x \in S} \delta(f(x), h(x))$$

where n is the no. of examples in S , and the quantity $\delta(f(x), h(x))$ is 1 if $f(x) \neq h(x)$ and 0 otherwise.

True Error

The true error ($\text{error}_D(h)$) of hypothesis h with respect to target function f and distribution D , is the probability that h will misclassify an instance drawn at random according to D

$$\text{error}_D(h) \equiv \Pr_{x \in D} [f(x) \neq h(x)]$$

Confidence Intervals for Discrete - Valued Hypothesis

Suppose we wish to estimate the true error for some discrete valued hypothesis h , based on its observed sample error over a sample S , where

→ The sample S contains n examples drawn independent of one another, and independent of h , according to the probability distribution D

→ $n \geq 30$

→ Hypothesis h commits r errors over these n examples (i.e. $\text{error}_S(h) = r/n$)

Under these conditions, statistical theory allows to make the following assertions.

1. Given no other info, the most probable value of $\text{error}_D(h)$ is $\text{error}_S(h)$
2. With approx 95% probability, the true $\text{error}_D(h)$ lies in the interval

$$\text{error}_S(h) \pm 1.96 \sqrt{\frac{\text{error}_S(h)(1 - \text{error}_S(h))}{n}}$$

b) Last year, five randomly selected students took a math aptitude test before they began their statistics course. The Statistics Department has three questions.

- i) What linear regression equation best predicts statistics performance, based their aptitude scores?

ii) if a student made on 80 on the aptitude test, what grade would we expect her to make in statistics?

iii) How well does the regression equation fit the data?

Student	x_i	y_i
1	95	85
2	85	95
3	80	70
4	70	65
5	60	70

i) $\hat{y} = b_0 + b_1 x_i$

$$b_0 = \frac{(\sum y)(\sum x^2) - \sum(x) \sum xy}{n(\sum x^2) - (\sum x)^2}$$

$$b_1 = \frac{n(\sum xy) - \sum x \sum y}{n(\sum x^2) - (\sum x)^2}$$

x	y	x^2	xy
95	85	9025	8075
85	95	7225	8075
80	70	6400	5600
70	65	4900	4550
60	70	3600	4200
390	385	31150	30500

$$b_0 = \frac{385 \times 31150}{3650} = 26.780$$

$$b_1 = \frac{2350}{3656} = 0.64383$$

$$\therefore \hat{y} = 26.780 + x \times 0.64383$$

i) $x = 80$

$$\hat{y} = 26.780 + 80 \times 0.64383$$

$$\hat{y} = 78.28$$

iii)

$$E = \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2$$

$$E = 33.78$$