

Assignment 1

Ans 1:

$$A = \begin{bmatrix} \vec{a}_{11} & \vec{a}_{12} & \dots & \vec{a}_{1m} \\ \vec{a}_{21} & \vec{a}_{22} & \dots & \vec{a}_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \vec{a}_{n1} & \vec{a}_{n2} & \dots & \vec{a}_{nm} \end{bmatrix}_{n \times m}$$

$$B = \begin{bmatrix} \vec{b}_{11} & \vec{b}_{12} & \dots & \vec{b}_{1n} \\ \vec{b}_{21} & \vec{b}_{22} & \dots & \vec{b}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \vec{b}_{m1} & \vec{b}_{m2} & \dots & \vec{b}_{mn} \end{bmatrix}_{m \times n}$$

$$AB = \begin{bmatrix} \sum_{i=1}^m \vec{a}_{1i} \vec{b}_{i1} & \sum_{i=1}^m \vec{a}_{1i} \vec{b}_{i2} & \dots & \sum_{i=1}^m \vec{a}_{1i} \vec{b}_{in} \\ \sum_{i=1}^m \vec{a}_{2i} \vec{b}_{i1} & \sum_{i=1}^m \vec{a}_{2i} \vec{b}_{i2} & \dots & \sum_{i=1}^m \vec{a}_{2i} \vec{b}_{in} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^m \vec{a}_{ni} \vec{b}_{i1} & \sum_{i=1}^m \vec{a}_{ni} \vec{b}_{i2} & \dots & \sum_{i=1}^m \vec{a}_{ni} \vec{b}_{in} \end{bmatrix}$$

Now

$$\text{tr}(AB) = \sum_{i=1}^m \vec{a}_{1i} \vec{b}_{i1} + \sum_{i=1}^m \vec{a}_{2i} \vec{b}_{i2} + \dots + \sum_{i=1}^m \vec{a}_{ni} \vec{b}_{in}$$

P.T.O.

So,

$$\frac{2 \operatorname{tr}(AB)}{2 \mathbf{1}^T \mathbf{1}} = \sum_{i=1}^m b_{e_i} + \sum_{i=1}^m b_{i_2} + \dots + \sum_{i=1}^m b_{i_n}$$
$$= \mathbf{b}^T \mathbf{1}$$

i.e.

$$\boxed{\frac{2 \operatorname{tr}(AB)}{2 \mathbf{1}^T \mathbf{1}} = \mathbf{B}^T}$$

Ans 2

Binary classes ($K=2$)

Suppose we have two classes c_1 and c_2 and the training data sets given as

$$(x_1, t_1) (x_2, t_2) \dots (x_n, t_n)$$

So, the posteriors of class c_1

$$p(c_1/x) = \frac{p(x, c_1)}{p(x)}$$

$$= \frac{p(x/c_1) p(c_1)}{\sum_{i=1}^n p(x_i)} \quad \text{[Crossed out]}$$

$$= \frac{p(x/c_1) \cdot p(c_1)}{\sum_{i=1}^2 p(x, c_i)}$$

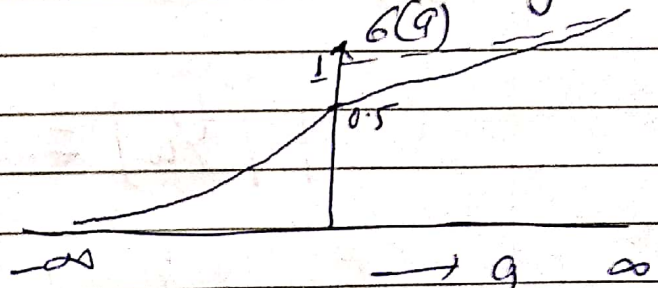
$$= \frac{p(x/c_1) \cdot p(c_1)}{p(x/c_1) p(c_1) + p(x/c_2) p(c_2)}$$

Say,

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

where $a = \ln \left[\frac{p(x/c_1) p(c_1)}{p(x/c_2) \cdot p(c_2)} \right]$

and $\sigma(a)$ is also called as sigmoid funⁿ.



$$\Rightarrow p(c_1/x) = \sigma(w^T x + w_0) \quad \text{--- (1)}$$

Assuming the probabilities are sigmoid funⁿ

So, $p(c_2/x) = 1 - p(c_1/x)$

So, if we know the weights, then we can get the probabilities of $p(c_1/x)$ or $p(c_2/x)$

RT-0

Now,

when $(w^T x + w_0) > 0$ then the
 x belongs to C_1 (say)
and when $(w^T x + w_0) < 0$ then
 x belongs to C_2 (say).

Also when

$w^T x + w_0 = 0$ this means
 x is ~~right~~ on the decision
surface.

eqⁿ (1) can also be written as,

$$y_n = \sigma(w^T x_n + w_0) = p(C_1/x_n)$$

So, as per likelihood

$$p(t/w) = \prod_{n=1}^N (y_n)^{t_n} (1-y_n)^{1-t_n} \quad \text{--- (1)}$$

i.e. if $t_n = 1$ then maximize y_n
 $t_n = 0$ then maximize $(1-y_n)$

For maximum likelihood,

$$w^* = \underset{w}{\operatorname{argmax}} (p(t/w)) \quad \text{--- (ii)}$$

//_

Since we know that the "log" of likelihood is also the likelihood.

So,

$$w^* = \underset{w}{\operatorname{argmax}} \left(\ln p(t/w) \right)$$

$$= \underset{w}{\operatorname{argmax}} \left\{ \ln \left[\prod_{n=1}^N y_n^{t_n} (1-y_n)^{(1-t_n)} \right] \right\}$$

$$= \underset{w}{\operatorname{argmax}} \left\{ \sum_{n=1}^N t_n \ln(y_n) + \sum_{n=1}^N (1-t_n) \ln(1-y_n) \right\}$$

$$= \underset{w}{\operatorname{argmax}} \left\{ \sum_{n=1}^N t_n \ln[\sigma(w^T x_n + w_0)] + \sum_{n=1}^N (1-t_n) \ln[1 - \sigma(w^T x_n + w_0)] \right\}$$

This is called L [loss]

Also called

Now,

$$\frac{\partial \sigma(q)}{\partial q} = \frac{\partial}{\partial q} \left[\frac{1}{1+e^{-q}} \right] \quad \text{log likelihood loss.}$$

$$= \frac{\partial}{\partial q} \left[\frac{e^q}{1+e^q} \right]$$

$$= \frac{e^{-q}}{(1+e^{-q})^2}$$

P.F.O.

$$\Rightarrow \left[\frac{\partial \theta(q)}{\partial q} = [1 - \theta(q)] \theta(q) \right] \text{--- (p)} \quad \text{--- 1/1}$$

$$\Rightarrow \frac{\partial L}{\partial w} = \sum_{n=1}^N t_n \frac{\partial (\ln y_n)}{\partial w} + \sum_{n=1}^N (1-t_n) \frac{\partial (\ln(1-y_n))}{\partial w}$$

Since $\frac{\partial \log x}{\partial x} = \frac{1}{x}$

So,

$$\frac{\partial L}{\partial w} = \sum_{n=1}^N t_n \cdot \frac{1}{y_n} + \sum_{n=1}^N (1-t_n) \cdot \frac{1}{1-y_n}$$

$$\frac{\partial L}{\partial w} = \sum_{n=1}^N t_n \cdot \frac{1}{y_n} \cdot (1-y_n) y_n \cdot x_n$$

$$+ \sum_{n=1}^N (1-t_n) \cdot \left(\frac{-1}{1-y_n} \right) y_n (1-y_n) \cdot x_n$$

$$= \sum_{n=1}^N \left[t_n (1-y_n) + (1-t_n) (-y_n) \right] x_n$$

$$\Rightarrow \left[\frac{\partial L}{\partial w} = \sum_{n=1}^N (t_n - y_n) x_n \right] \text{--- (v)}$$

Now, As per gradient ascent method,
for moving from one point to another
we have a rule

$$w^{t+1} = w^t + \eta \left. \frac{\partial L}{\partial w} \right|_{w=w^t}$$

So, putting eqⁿ (v)th value in above
eqⁿ we get,

$$w^{t+1} = w^t + \eta \left[\sum_{n=1}^N (t_n - y_n) x_n \right]_{w=w^t}$$

Multi-class ($K > 2$)

$$\begin{bmatrix} \vdots \\ p(c_k/x_k) \\ \vdots \end{bmatrix}_{K \times 1} = \frac{p(x/c_k) p(c_k)}{\sum_{j=1}^K p(x/c_j) p(c_j)}$$

$$= \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)} \quad \text{--- (7)}$$

where.

$$\exp(a_k) = p(x/c_k) \cdot p(c_k)$$

$$\exp(a_j) = p(x/c_j) \cdot p(c_j)$$

And

$$a_j = \ln \left[p(x/c_j) \cdot p(c_j) \right]$$

$$\underline{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix}_{k \times 1}$$

eqⁿ (1) is called the softmax of "a".

Logistic regression for k-classes

$$a_{nk} = W_k^T \underline{x}_n + W_{k0}$$

$$\text{where } \underline{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix}_{k \times 1}$$

$$\text{Output } \underline{y} = \text{softmax}(\underline{a})$$

//_

So, for training data

$(x_n, t_n)_{n=1}^N$ Assuming t_n is
one hot encoded.

then the likelihood model is defined as,

$$p(\mathcal{T}|\omega) = \prod_{n=1}^N \prod_{k=1}^K (y_{nk})^{t_{nk}}$$

Maximum likelihood model is,

$$\tilde{\omega}^* = \underset{\tilde{\omega}}{\operatorname{argmax}} \left[\prod_{n=1}^N \prod_{k=1}^K (y_{nk})^{t_{nk}} \right]$$

where $\tilde{\omega} = \begin{bmatrix} \tilde{\omega}_1^T \\ \tilde{\omega}_2^T \\ \vdots \\ \tilde{\omega}_K^T \end{bmatrix}$ $\omega = \begin{bmatrix} \omega \\ 1 \end{bmatrix}$

Since we know that the "log" of likelihood
is also a likelihood hence,

$$\tilde{\omega}^* = \underset{\tilde{\omega}}{\operatorname{argmax}} \left[\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln(y_{nk}) \right]$$

P.T.P.

Now, $\frac{\partial \text{softmax}(a)}{\partial a} = \frac{\partial}{\partial a} \frac{e^a}{\sum_{j=1}^K e^{a_j}}$

$$= \frac{\partial}{\partial a} \left\{ \frac{e^a}{e^{a_1} + e^{a_2} + e^{a_3} + \dots + e^{a_K}} \right\}$$

Using quotient rule

$$\frac{\partial \text{softmax}(a)}{\partial a} = \frac{e^a \times \sum_{j=1}^K e^{a_j} - (e^a)^2}{\left(\sum_{j=1}^K e^{a_j} \right)^2}$$

$$= \frac{e^a}{\sum_{j=1}^K e^{a_j}} - \frac{(e^a)^2}{\left(\sum_{j=1}^K e^{a_j} \right)^2}$$

$$= \text{softmax}(a) - \text{softmax}^2(a)$$

$$= \text{softmax}(a) [1 - \text{softmax}(a)]$$

Using this

$$w^* = \underset{\tilde{w}}{\text{argmax}} \left\{ \sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{\partial}{\partial w} (\ln(Y_{nk})) \right\}$$

$$\Rightarrow \tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \left[\sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \frac{1}{y_{nk}} \cdot \frac{\partial \mathcal{L}}{\partial \tilde{w}} \right]$$

$$\Rightarrow \tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \left[\sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \frac{1}{y_{nk}} \cdot (1 - y_{nk}) \cdot y_{nk} \right]$$

$$\tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \left[\sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot (1 - y_{nk}) \right]$$

This is $\frac{\partial \mathcal{L}}{\partial \tilde{w}}$.

So, as per gradient ascent method, for moving from one point to another we have a rule,

$$w^{t+1} = w^t + \eta \frac{\partial \mathcal{L}}{\partial w} \Big|_{w=w^t}$$

$$\text{So, } w^{t+1} = w^t + \eta \left[\sum_{n=1}^N \sum_{k=1}^K t_{nk} (1 - y_{nk}) \right]$$