Q

Q.1. For logistic regression:-

$$LL(\theta) = \sum_{i=1}^{m} \log P(y^{(i)}|x^{(i)};\theta), \quad \text{where } P(y=1|x^{(i)};\theta) = \frac{1}{1+e^{-\theta^T x^{(i)}}}$$

$$\Rightarrow LL(\theta) = \sum_{i=1}^{m}\left[ y^{(i)} \log \frac{1}{1+e^{-\theta^T x^{(i)}}} + [1-y^{(i)}] \log\left(1- \frac{1}{1+e^{-\theta^T x^{(i)}}}\right)\right]$$

we know that

$$\frac{d g(z)}{dz} = g(z)(1-g(z)) \quad \text{where } g(z) \text{ is the sigmoid function.}$$

$$\Rightarrow \frac{d}{d\theta_j} LL(\theta) = \sum_{i=1}^{m}\left[ \frac{y^{(i)}}{g(\theta^T x^{(i)})} g(\theta^T x^{(i)}) \cdot (1- g(\theta^T x^{(i)})) \, x_j^{(i)} + \frac{(1-y^{(i)})(-1) g(\theta^T x^{(i)})(1-g(\theta^T x^{(i)}))}{(1-g(\theta^T x^{(i)}))} (x_j^{(i)})\right]$$

$$= \sum_{i=1}^{m}\left[ y^{(i)} - y^{(i)} g(\theta^T x^{(i)}) + (1-y^{(i)})(-1) g(\theta^T x^{(i)})\right] x_j^{(i)}$$

$$= \sum_{i=1}^{m}\left[ y^{(i)} - g(\theta^T x^{(i)})\right] x_j^{(i)}$$

$$\Rightarrow \frac{d}{d\theta_k d\theta_j} LL(\theta) = \sum_{i=1}^{m} (-1) g(\theta^T x^{(i)}) [1-g(\theta^T x^{(i)})] x_j^{(i)} x_k^{(i)}$$

Now, $H_{jk} = \dfrac{d \, LL(\theta)}{d\theta_k d\theta_j} = \displaystyle\sum_{i=1}^{m}(-1) g(\theta^T x^{(i)}) [1-g(\theta^T x^{(i)})] x_j^{(i)} x_k^{(i)}$

$$Z^T H Z = \sum_{j=1}^{n} \sum_{k=1}^{n} Z_j H_{jk} Z_k$$

$$\Rightarrow Z^T H Z = \sum_{j,k=1,1}^{n,n} Z_j \left[ \sum_{i=1}^{m}(-1) g(\theta^T x^{(i)})[1-g(\theta^T x^{(i)})] x_j^{(i)} x_k^{(i)}\right] Z_k$$

$$= \sum_{i=1}^{m}\left[ \underbrace{\sum_{j} Z_j x_j^{(i)} \sum_{k} Z_k x_k^{(i)}}_{\text{identical}}\right] (g(\theta^T x^{(i)})(1-g(\theta^T x^{(i)})))$$

$$= \sum_{i=1}^{m}(-1) Z^T x^{(i)} [Z^T x^{(i)}]^T \; g(\theta^T x^{(i)})(1-g(\theta^T x^{(i)})) \leq 0$$

since each term inside the sum is +ve (or nm-n$_e$ positive)

$\Rightarrow$ $Z^T H Z \cancel{=0} \leq 0$

$\Rightarrow$ H is +ve semi-definite

$\Rightarrow$ $LL(\theta)$ is a concave function $\odot$

(since its ~~se~~ matrix of corresponding second order derivatives is +ve semi-definite).

---

**Q.2.**

Procedure:-

First Divide the training set $T_r$ into two subsets. Train & validation. Let us say the split be $80|20$.

Then, we first "train" the model on 80 examples & test/validate on 20 w for various values of $2$ $\tau$.

te_best $= \infty$; $\tau_{best} = \cancel{+}0$;

~~For $2\tau$ in 2range 2(0, $\tau_{max}$) 2~~

~~for($\tau=0$; $\tau < \tau_{max}$; $\tau \pm \tau + S$)~~ → increment

Let $T_v$ set of examples in validation set.

For $\tau$ in range $(0, \tau_{max})$ with increments $\tau_\Delta$ {

$\hat{T}_r = T_r - T_v$; $te = 0$;

For $(x^{(i)}, y^{(i)}) \in T_v$ {

$M^{\#} = $ learn LWR $(\hat{T}_r, x^{(i)}, \tau)$

$\hat{y}^{(i)} = $ pre $M(x^{(i)})$;

$te = te + Err(y^{(i)}, \hat{y}^{(i)})$;

if $(te < te\_best)$ {

$\tau_{best} = \tau$;

$te = te\_best$;

}

}

Compute error of model on each example

Compute total error

update $\tau$-best if required

// if total error is best till // now, update // te-best & $\tau$-best

Next, compute the error on test set.

```
te = 0;
For ∀ (x⁽ᵘ⁾, y⁽ᵘ⁾) ∈ Te {
    M = learnLWF(Tr, x⁽ᵘ⁾, Zbest);
    ŷ⁽ᵘ⁾ = M(x⁽ᵘ⁾);
    te = te + Err(ŷ⁽ᵘ⁾, y⁽ᵘ⁾);
}
return te;   } Return the error on test set.
```

3

---

**Q.3.**

In gradient de ascent,

$$g_j = \pm \frac{d}{d\theta_j} LL(\theta)$$

$$LL(\theta) = \sum_{i=1}^{m} \log(P(y^{(i)}, x^{(i)}; \theta))$$

$$\Rightarrow g_j = \frac{d}{d\theta_j} \sum_{i=1}^{m} \log[P(y^{(i)}; x^{(i)}; \theta)] = \sum_{i=1}^{m} \frac{d}{d\theta_j} [\log P(y^{(i)}; x^{(i)}; \theta)]$$

Now, Consider SGD.

Define $LL(\theta)^{(i)} = \log(P(y^{(i)}, x^{(i)}; \theta)]$

$$g_j^{(i)} = \frac{d}{d\theta_j} \log P(y^{(i)}, x^{(i)}; \theta)$$

Now $E[g_j^{(i)}] = $ ...

$$E[g_j^{(i)}] = \sum_{i=1}^{m} \frac{1}{m} \frac{\partial}{\partial \theta_j} [\log p(y^{(i)}, x^{(i)}; \theta)]$$

Prob. of seeing $i$th example
(sample uniformly at random from training set)

$$= \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial \theta_j} \log [p(y^{(i)}, x^{(i)}; \theta)]$$

$$= \frac{1}{m} g_j$$

$$\Rightarrow g_j = m * E[g_j^{(i)}] \qquad \forall j \quad \text{(state we do)}$$

Thus holds $\forall j$.

Hence, proved.

---

**Q4.** Assumptions:- Yes, we did make i.i.d. assumption over the training set. Because of this assumption we can write:-

$$LL(\theta) = \log \#(p(y^{(1)}, \ldots y^{(m)}, x^{(1)} \ldots x^{(m)}; \theta))$$

$$\quad \textcircled{1} \text{ since examples are independent}$$

$$= \log \prod_{i=1}^{m} [p(y^{(i)}, x^{(i)}; \theta)]$$

$$= \sum_{i=1}^{m} \log [p(y^{(i)}, x^{(i)}; \theta)] \qquad z$$

& the expression for $g_j = \frac{\partial}{\partial \theta_j} LL(\theta) = \sum_{i=1}^{m} \frac{\partial [\log p(y^{(i)}, x^{(i)}; \theta)]}{\partial \theta_j}$ follows. This won't hold if $(x^{(i)}, y^{(i)})$ were not i.i.d.

3)

Normal Distribution

$$P(x^{(i)} | y^{(i)} = 1; \, \theta, \mu_1, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)}{2}}$$

Decision boundary :-

$$P(y^{(i)} = 1 | x^{(i)}; \theta) = 0.5$$

$$= \frac{P(x^{(i)} | y^{(i)} = 1) \; P(y^{(i)} = 1)}{P(x^{(i)}) \longrightarrow}$$

$$\hookrightarrow P(x^{(i)} | y^{(i)} = 1) P(y^{(i)} = 1)$$
$$+ P(x^{(i)} | y^{(i)} = 0) P(y^{(i)} = 0)$$

$$\geq \frac{1}{P(x^{(i)} | y^{(i)} = 0) P(y^{(i)} = 0) + P(x^{(i)} | y^{(i)} = 1) P(y^{(i)} = 1)} \cdot \frac{P(x^{(i)} | y^{(i)} = 1) P(y^{(i)} = 1)}{} = \frac{1}{2}$$

$$\Rightarrow \frac{1}{1 + \frac{P(x^{(i)} | y^{(i)} = 0) \; P(y^{(i)} = 0)}{P(x^{(i)} | y^{(i)} = 1) \; P(y^{(i)} = 1)}} = \frac{1}{2}$$

$$\Rightarrow \frac{2}{1} = 1 + \frac{P(x^{(i)} | y^{(i)} = 0) \; P(y^{(i)} = 0)}{P(x^{(i)} | y^{(i)} = 1) \; P(y^{(i)} = 1)}$$

Taking log

$$0 = \log \left[ P(x^{(i)} | y^{(i)} = 0) \; P(y^{(i)} = 0) \right]$$
$$- \log \left[ P(x^{(i)} | y^{(i)} = 1) \; P(y^{(i)} = 1) \right]$$

2) $\log P(x^{(i)}|y^{(i)}=0) + \log P(y^{(i)}=0)$

$$= \log P(x^{(i)}|y^{(i)}=1) + \log P(y^{(i)}=1)$$

$\Rightarrow \log \dfrac{1}{(2\pi)^{n/2}|\Sigma_0|^{1/2}} + (-1)\dfrac{(x-\mu_0)^T \Sigma_0^{-1}(x-\mu_0)}{2} + \log \phi$

$$= \log \dfrac{1}{(2\pi)^{n/2}|\Sigma_1|^{1/2}} - \dfrac{(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)}{2}$$

$$+ \log(1-\phi)$$

$\Rightarrow$

$$\frac{1}{2}\left[ (x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1) - (x-\mu_0)^T \Sigma_0^{-1}(x-\mu_0)\right]$$

$$= \log \dfrac{(2\pi)^{n/2}|\Sigma_0|^{1/2}}{(2\pi)^{n/2}|\Sigma_1|^{1/2}} + \log\left[\dfrac{(1-\phi)}{\phi}\right]$$

$\frac{\ }{\ }$ ~~$x^T[\Sigma_1^{-1}-\Sigma_0^{-1}]x$~~ ~~$+ x^T\Sigma_1^{-1}\mu_1 - x^T\Sigma_1^{-1}x$~~

$x^T \Sigma_1^{-1} x + x^T \Sigma_1^{-1}\mu_1 - x\mu_1^T\Sigma_1^{-1}x + \mu_1^T\Sigma_1^{-1}\mu_1$

$-\left[ x^T\Sigma_0^{-1}x + x^T\Sigma_0^{-1}\mu_0 - \mu_0^T\Sigma_0^{-1}x + \mu_0^T\Sigma_0^{-1}\mu_0\right]$

Note $\Sigma_1^{-1} = (\Sigma_1^{-1})^T$ & $\Sigma_0^{-1} = (\Sigma_0^{-1})^T$ $\quad = \log \dfrac{|\Sigma_0|^{1/2}}{|\Sigma_1|^{1/2}} * \left[\dfrac{1-\phi}{\phi}\right]$

$\Rightarrow x^T[\Sigma_1^{-1}-\Sigma_0^{-1}]x + 2\left[x^T\Sigma_1^{-1}\mu_1 - x^T\Sigma_0^{-1}\mu_0\right]$

$$+ \mu_1^T\Sigma_1^{-1}\mu_1 - \mu_0^T\Sigma_0^{-1}\mu_0$$

$$= \log \dfrac{|\Sigma_0|^{1/2}(1-\phi)}{|\Sigma_1|^{1/2}\phi}$$

$\Rightarrow x^T(\Sigma_1^{-1}-\Sigma_0^{-1})x + 2x^T\left[\Sigma_1^{-1}\mu_1 - \Sigma_0^{-1}\mu_0\right] = \log\left[\dfrac{|\Sigma_0|^{1/2}(1-\phi)}{|\Sigma_1|^{1/2}\phi}\right] + --$

Looking for quadratic terms.

coefficients of $x_i y_j = 0 \; \forall \; j \neq k$.

Since $\hat{\Sigma}_1 \& \hat{\Sigma}_0$ both are diagonal

coefficients of $x_j^2$ is given as:-

$$(\hat{\Sigma}_1^{-1})_{jj} - (\hat{\Sigma}_0^{-1})_{jj}$$

$$= \left[ \frac{1}{(\hat{\Sigma}_1)_{jj}} - \frac{1}{(\hat{\Sigma}_0)_{jj}} \right]$$

(b). When $\hat{\Sigma}_0 = \hat{\Sigma}_1$. Quadratic terms vanish.

So, Decision boundary is simply a linear function of $x$. (Hyper plane).

(c) → when $x \in \mathbb{R}^2$.

(a) → Decision boundary is of form.

$$a x_1^2 + a' x_1 + b x_2^2 + b' x_2 + c = 0$$

⇒ Equation of an ellipse or hyperbola whose principal axes are aligned with $x$-$y$ axis.

**(2) Q.5:**

$$P(y|x; \lambda) = \frac{e^{-\lambda}\lambda^k}{k!} \qquad y=k \qquad \lambda = \theta^T x$$

for $y = k$

$$\Rightarrow \log P(y^{(i)}|x^{(i)}; \theta) = \log e^{-\lambda} + \log \lambda^{y^{(i)}} - \log k!$$

$$= \log e^{-\lambda} + y^{(i)} \log \lambda - \log y!$$

$$= -\lambda + y^{(i)} \log \lambda - \log y!$$

substituting $\lambda = \theta^T x^{(i)}$

$$= -e^{\theta^T x^{(i)}} + y^{(i)} \log[e^{\theta^T x^{(i)}}] - \log y!$$

$$\Rightarrow LL(\theta) = \sum_{i=1}^{m} \log P(y^{(i)}|x^{(i)}; \theta)$$

$$\Rightarrow = \sum_{i=1}^{m} -e^{\theta^T x^{(i)}} + y^{(i)} \log[e^{\theta^T x^{(i)}}] - \log y!$$

$$= \sum_{i=1}^{m} -e^{\theta^T x^{(i)}} + y^{(i)} \theta^T x^{(i)} - \log y!$$

$$\frac{d}{d\theta_j} LL(\theta) = \frac{d}{d\theta_j} \sum_{i=1}^{m} - e^{\theta^T x^{(i)}} .$$

$$= \sum_{i=1}^{m} \left[ y^{(i)} - e^{\theta^T x^{(i)}} \right] x_j^{(i)}$$

$$\Rightarrow \frac{d}{d\theta_j} LL(\theta) = \sum_{i=1}^{m} \left[ y^{(i)} - h_\theta(x^{(i)}) \right] x_j^{(i)}$$

where $\boxed{h_\theta(x^{(i)}) = e^{\theta^T x^{(i)}} = e^{x^T \lambda}}$

Hence, $E[y^{(i)}|x^{(i)}; \theta] = \lambda = e^{\theta^T x^{(i)}} = h_\theta(x^{(i)})$   mean of distribution   $P(y|x; \lambda)$
Hence, proved