**5.8)**

**a)** Joint distribution can be factorized into

$$p(x,y|\theta) = p(y|x,\theta_2)\, p(x|\theta_1) \qquad \theta = (\theta_1, \theta_2)$$

So, table becomes

| | $y=0$ | $y=1$ |
|---|---|---|
| $x=0$ | $\theta_2(1-\theta_1)$ | $(1-\theta_1)(1-\theta_2)$ |
| $x=1$ | $\theta_1(1-\theta_2)$ | $\theta_1\,\theta_2$ |

Data

| $x$ | $y$ |
|---|---|
| 1 | 1 |
| 1 | 0 |
| 0 | 0 |
| 1 | 0 |
| 1 | 1 |
| 0 | 0 |
| 0 | 1 |

**b)** MLE Solution for likelihood given

$$x = (1,1,0,1,1,0,0) \qquad y = (1,0,0,0,1,0,1) \quad \text{is}$$

$$\hat\theta = \underset{\theta_1,\theta_2}{\arg\max}\ \log \prod_{i=1}^{m} p(x^i/\theta)$$

$$= \underset{\theta}{\arg\max} \cdot \log\Big( (\theta_1 \cdot \theta_2)(\theta_1(1-\theta_2))(\theta_2(1-\theta_1))(\theta_1 \cdot (1-\theta_2))$$
$$\cdot (\theta_1 \cdot \theta_2)(\theta_2 \cdot (1-\theta_1)) \cdot (1-\theta_1)(1-\theta_2)\Big)$$

$$= \underset{\theta}{\arg\max} \cdot \log\ \theta_1^{4} \cdot (1-\theta_1)^3 \cdot \theta_2^{4}(1-\theta_2)^3$$

$$= \underset{\theta}{\arg\max} \cdot \Big(4\log\theta_1 + 3\log(1-\theta_1) + 4\log\theta_2 + 3\log(1-\theta_2)\Big)$$

let call this minimization as $J(\theta)$ and take partial derivatives w.r.t $\theta_1, \theta_2$ and set 0

$$\frac{\partial J}{\partial \theta_1} = \frac{4}{\theta_1} - \frac{3}{1-\theta_1} = 0 \implies 4(1-\theta_1) - 3\theta_1 = 0$$

$$\implies \theta_1 = \tfrac{4}{7}$$

$$\frac{\partial J}{\partial \theta_2} = \frac{4}{\theta_2} - \frac{3}{1-\theta_2} = 0 \implies \theta_2 = \tfrac{4}{7}$$

The marginal likelihood for this data is

$$P(\mathcal{D}/\theta, M) = \prod_{i=1}^{n} P(n_i, y_i/\theta)$$

$$= \left(\frac{4}{7} \cdot \frac{4}{7}\right)\left(\frac{4}{7} \cdot \frac{3}{7}\right)\left(\frac{4}{7} \cdot \frac{3}{7}\right)\left(\frac{3}{7} \cdot \frac{3}{7}\right)\left(\frac{4}{7} \cdot \frac{4}{7}\right)$$

$$\left(\frac{4}{7} \cdot \frac{3}{7}\right)\left(\frac{3}{7} \cdot \frac{3}{7}\right)$$

() MLE with 4 parameters representing all combinations

$$\hat{\theta} = \underset{\theta}{\arg\max} \ \log \prod_{i=1}^{m} P(\mathcal{D}/\theta) \quad s.t \ \sum_{i=0}^{1} \theta_{i,i} = 1$$

$$= \underset{\theta}{\arg\max} \ \log\left(\theta_{1,1} \cdot \theta_{1,0} \cdot \theta_{0,0} \cdot \theta_{1,0} \cdot \theta_{1,1} \cdot \theta_{0,0} \cdot \theta_{0,1}\right)$$

$$s.t \quad \theta_{0,0} + \theta_{1,1} + \theta_{1,0} + \theta_{0,1} = 1$$

$$= \underset{\theta}{\arg\max} \left(2\log\theta_{1,1} + 2\log\theta_{1,0} + 2\log\theta_{0,0} + 1\log\theta_{0,1}\right)$$

lets call this as $J(\theta)$ and do same as above,

lets write this as using lagrangian

$$\frac{dJ}{d\theta_{0,0}} \neq \frac{2}{\theta_{0,0}} \quad \text{multiplier}$$

$$J(\theta) = 2\log\theta_{1,1} + 2\log\theta_{1,0} + 2\log\theta_{0,0} + 1\log\theta_{0,1} + \bar{\lambda}\left(\theta_{0,0} + \theta_{0,1}\right)$$

$$+ \theta_{1,0} + \theta_{1,1} - 1) = 0.$$

$$\frac{\partial J}{\partial \theta_{0,0}} = \frac{2}{\theta_{0,0}} + \lambda = 0 \implies \theta_{0,0} = \frac{2}{\lambda} \to \textcircled{0}$$

Same like this $\theta_{1,0} = \theta_{1,1} = \frac{2}{\lambda}, \theta_{0,1} = \frac{1}{\lambda}$

$$\to \textcircled{0}$$

we can solve $\lambda$ by using $\sum_{i=0}^{1} \theta_{i,i} = 1$

$$\implies \frac{7}{\lambda} = 1 \implies \lambda = 7$$

By substituting $\lambda = 7$ in above

$$\theta_{0,0} = \theta_{1,0} = \theta_{1,1} = \frac{2}{7} \text{ and } \theta_{0,1} = \frac{1}{7}$$

with $4$-parameters in

and marginal likelihood

$$P(\mathcal{D}/\hat{\theta}, M_4) = \prod_{i=1}^{m} P(n_i, y_i/\theta)$$

$$= \frac{2}{7} \cdot \frac{2}{7} \cdot \frac{2}{7} \cdot \frac{1}{7} = \frac{8}{7^4}$$

d) For leave-one-out Cross Validation, we have to remove one data point at every time and see how well the new parameters explain remaining data.

Let's see first for model 1:

$$L(1) = \sum_{i=1}^{m} \log P\left(n_i, y_i / \hat{\theta}(\mathcal{D}_{-i})\right)$$

Here Sample size in $7 \Rightarrow m = 7$

So, for $i=1$, $\hat{\theta}_{\mathcal{D}_{-1}} = \left\{\theta_1 = \frac{3}{6}, \theta_2 = \frac{3}{6}\right\}$ and $P(n_1, y_1/\theta_{\mathcal{D}_{-1}}) = \frac{3}{6} \log \frac{1}{2}$

$$= \log \frac{1}{7}$$

for $i=2$, $\hat{\theta}_{\mathcal{D}_{-2}} = \left[\theta_1 = \frac{3}{6}, \theta_2 = \frac{4}{6}\right]$ and $P(n_2, y_2/\theta_{\mathcal{D}_{-2}}) = \log \frac{3}{6} \cdot \frac{4}{6} = \log \frac{1}{3}$

for $i=3$, $\hat{\theta}_{\mathcal{D}_{-3}} = \left[\theta_1 = \frac{4}{6}, \theta_2 = \frac{3}{6}\right]$ and $P(n_3, y_3/\theta_{\mathcal{D}_{-3}}) = \log \frac{1}{3}$.

for $i=4$, $\hat{\theta}_{\mathcal{D}_{-4}} = \left[\theta_1 = \frac{3}{6}, \theta_2 = \frac{4}{6}\right]$ and $P(n_4, y_4/\theta_{\mathcal{D}-4}) = \log \frac{1}{3}$

for $i=5$, $\hat{\theta}_{\mathcal{D}_{-5}} = \left[\theta_1 = \frac{3}{6}, \theta_2 = \frac{3}{6}\right]$ and $P(n_5, y_5/\theta_{\mathcal{D}-5}) = \log \frac{1}{4}$

for $i=6$, $\hat{\theta}_{\mathcal{D}-6} = \left[\theta_1 = \frac{4}{6}, \theta_2 = \frac{3}{6}\right]$ and $P(n_6, y_4/\theta_{\mathcal{D}-6}) = \log \frac{1}{3}$

for $i=7$, $\hat{\theta}_{D-7} = \left[\theta_1 = \frac{4}{6}, \theta_2 = \frac{2}{6}\right]$ and $p\left(x_7, y_7 / \hat{\theta}_{D-7}\right) =$

$$= \log \frac{4}{6} \cdot \frac{2}{6} = \log \frac{4}{9}$$

the cross validated likelihood for mode $M_2$

so,

$$L(1) = 2\log \frac{1}{4} + 4\log \frac{1}{3} + \log \frac{4}{9}$$

is

let's do same for model 2:

$$L(2) = \sum_{i=1}^{m} \log p\left(x_i, y_i / \hat{\theta}_{D-i}\right)$$

for $i=1$, $\hat{\theta}_{D-1} = \left[\theta_{1,?} = \frac{1}{6}, \theta_{1,0} = \frac{2}{6}, \theta_{0,0} = \frac{2}{6}, \theta_{0,1} = \frac{1}{6}\right]$

and $p\left(x_1, y_1 / \hat{\theta}_{D-1}\right) = \log \frac{1}{6}$

for $i=2$, $\hat{\theta}_{D-1} = \left[\theta_{1,1} = \frac{2}{6}, \theta_{1,0} = \frac{1}{6}, \theta_{0,0} = \frac{2}{6}, \theta_{0,1} = \frac{1}{6}\right]$

and $p\left(x_2, y_2 / \hat{\theta}_{D-2}\right) = \log \frac{1}{6}$

like this for $i=3,4,5,6$ we get $p\left(x_i, y_i / \hat{\theta}_{D-i}\right) = \log \frac{1}{6}$

for $i=7$, $\hat{\theta}_{D-7} = \left[\theta_{1,1} = \frac{2}{6}, \theta_{1,0} = \frac{2}{6}, \theta_{0,0} = \frac{2}{6}, \theta_{0,11} = 0\right]$

and $p\left(x_7, y_7 / \hat{\theta}_{D-7}\right) = \log 0$

$= $ undefined.

So, for Mode $M_4$, if an we saw ~~training~~ examples in prediction that are not there in training samples, solution is undefined. Then is because in mode $M_4$ we overfit the data and set $\theta_{0,1} = 0$, as we didn't see $x = 0, y = 1$ in training samples. we'll choose $M_2$. for this reason.

b1) We have completely random dataset with $N_1$ examples from class 1, and $N_2$ from class 2, (with equal proportions i.e $N_1 = N_2$). Based on this intuition as there in no learning here,

$$P(\hat{y} = \text{class 1}) = \frac{N_1}{N_1 + N_2} = \frac{1}{2}$$

$$P(\hat{y} = \text{class 2}) = \frac{N_2}{N_1 + N_2} = \frac{1}{2}$$

Misclassification rate $E(\text{error}) = 1 \cdot P(\hat{y} \neq \text{correct class})$
$$+ 0 \cdot P(\hat{y} = \text{correct class})$$
$$= \frac{1}{2}$$

The best misclassification rate is $\frac{1}{2}$.

Suppose we use Leave one out - cross validation then

for all $-i$,

if $i$th example belongs to class 1, then prediction for

class 1 $\Rightarrow P\left(\hat{y} = \text{class 1} / i \in \text{class 1}\right) = \frac{N_1 - 1}{N - 1}$

and for class 2 $\Rightarrow P\left(\hat{y} = \text{class 2} / i \in \text{class 1}\right) = \frac{N_2}{N_1 + N_2 - 1}$
$$= \frac{N_2}{N - 1}$$

and suppose if $i$th class belongs to class 2, then

$$P\left(\hat{y} = \text{class 1} / i \in \text{class 2}\right) = \frac{N_1}{N - 1}$$

& $P\left(\hat{y} = \text{class 2} / i \in \text{class 2}\right) = \frac{N_2 - 1}{N - 1}$

Estimated misclassification rate

$i$th example is $P(\hat{y}/i \in class 1) = \dfrac{N_1 - 1}{N - 1}$.

in same way, if $i$th example belongs to class 2 then

$$P(\hat{y}/i \in class 2) = \dfrac{N_2 - 1}{N - 1}$$

then total misclassification rate can be

$$= \dfrac{N_2}{N}\left(1 - \dfrac{N_2 - 1}{N - 1}\right) + \dfrac{N_1}{N}\left(1 - \dfrac{N_1 - 1}{N - 1}\right)$$

$$= \dfrac{N_2 \cdot (N_1 - 1)}{N(N - 1)} + \dfrac{N_1 \cdot (N_2 - 1)}{N(N - 1)}$$

$$= \dfrac{2 N_1 N_2 - N}{N(N - 1)}$$

6.2) $P(y_i/\theta_i) \sim N(\theta_i, \sigma^2)$ and $P_\eta(\theta) \sim N(m_0, J_0^2)$

$$\eta = (m_0, J_0^2)$$

$\sigma^2$ in given 2500

ML-II estimate for $m_0, J_0^2$ is

$$\hat{\eta} = \underset{\eta}{\arg max} \int P(\mathcal{D}/\theta) \cdot P_\eta(\theta) \, d\theta$$

$$= \underset{\eta}{\arg max} \int \cdot e^{-\frac{1}{2} \sum\limits_{i=1}^{m} \frac{(y_i - \theta)^2}{\sigma^2}} \cdot e^{-\frac{1}{2} \cdot \left(\frac{\theta - m_0}{J_0^2}\right)^2} \, d\theta$$

$$= \underset{\eta}{\arg max} \int \cdot e^{-\frac{1}{2} \cdot}$$

let find out the quantity from given data

$$\sum_{i=1}^{m} \frac{(y_i - \theta)^2}{\sigma^2} = \frac{1}{500}\left((1505-\theta)^2 + (15...\right.$$

Based on equations from below

if, $P(\mathcal{D}/\theta) = N(\bar{y}/\theta, \sigma^2/N)$

and $P(\theta) = N(\theta/ m_0, J_0^2)$ then can

$$P(\theta/\mathcal{D}) = N(\theta/ m_N, J_N^2)$$

$$m_N = J_N^2 \left( \frac{\sigma^2(N\bar{y})}{} + J_0 \right)$$

$$\frac{1}{J_N^2} = \frac{1}{J_0^2} + \frac{N}{\sigma^2}$$

and $$m_N = J_N^2 \left( \frac{N\bar{y}}{\sigma^2} + \frac{m_0}{J_0^2} \right)$$

.As we are trying to minimize marginal likelihood

$$P_\eta(\mathcal{D}/\theta) = \int P(\mathcal{D}/\theta) \cdot P_\eta(\theta) \, d\theta$$

$$= \int \cdot N(\theta/ m_N, J_N^2) \, d\theta$$

$$2 \quad \sqrt{2\pi J_N^2}$$

$\Rightarrow$ argmin $P_\eta(\mathcal{D}/\theta) = $ argmin $\sqrt{2\pi J_N^2}$
$\quad \eta \qquad\qquad\qquad\qquad \eta$

$$= \text{argmin} \cdot \sqrt{2\pi} \cdot \left( \frac{1}{J_0^2} + \frac{N}{\sigma^2} \right)^{1/2}$$
$$\eta$$