## PROBLEM 1

2 (a)

$$I(X,Y) = H(X) - H(X|Y)$$

We need to prove that

$$I(X,Y) = KL\left( p(x,y) \| p(x)p(y) \right)$$

$$= -\sum_x \sum_y p(x,y) \log \frac{p(x)p(y)}{p(x,y)} \quad \text{(from definition)}$$

First, let us write

$H(X|Y) = 8$

$$H(X|Y=\varnothing) = -\sum_{i=1}^{n} P(X=i|Y=\varnothing) \log P(X=i|Y=\varnothing)$$

$$H(X|Y) = \sum_\varnothing P(Y=\varnothing) H(X|Y=\varnothing)$$

$$= -\sum_\varnothing \sum_{i=1}^{n} P(Y=\varnothing) P(X=i|Y=\varnothing) \log P(X=i|Y=\varnothing)$$

$$= -\sum_y \sum_x p(y) p(x|y) \log P(x|y)$$

$$= -\sum_y \sum_x p(x,y) \log p(x|y)$$

$$= -\sum_y \sum_x p(x,y) \log \frac{p(x,y)}{p(y)}.$$

Now, we can write

$$H(X) = -\sum_x p(x) \log p(x)$$

$$= -\sum_x \sum_y p(x,y) \log p(x).$$

② 

Therefore,

$I(x,y) = H(x) - H(x/y)$

$$= -\sum_x \sum_y p(x,y) \log p(x) + \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{p(y)}$$

$$= -\sum_x \sum_y p(x,y) \log \frac{p(x)p(y)}{p(x,y)}$$

$$= KL\left( P(x,y) \,||\, P(x)P(y) \right)$$

Hence, proved.

(2) The definition is

$$I(x,y) = -\sum_x \sum_y p(x,y) \log \frac{p(x)p(y)}{p(x,y)}$$

We have $I(x,y) = 0$ when

$p(x,y)$ ... then x and y are ...

1) $\log \frac{p(x)p(y)}{p(x,y)} = 0 \Rightarrow \frac{p(x)p(y)}{p(x,y)} = 1$

$$p(x,y) = p(x)\,p(y)$$

This is the condition for independence of x and y.

Thus, $I(x,y)$ is zero when x and y are independent. Given y, there is no new information about x and y, thus $I(x,y) = 0$.

2) $p(x,y) = 0$ on x and y are disjoint events.

2) $H(x) = -\int p(x) \ln p(x) \, dx$

we can write $\ln p(x)$ as:

$$\ln p(x) = \ln \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp - \frac{(x-\mu)^2}{2\sigma^2} \right]$$

$$= \ln \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \exp \frac{(x-\mu)^2}{2\sigma^2}$$

$$= -\ln(\sqrt{2\pi}\sigma) - \frac{(x-\mu)^2}{2\sigma^2}$$

$$H(x) = -\int_{-\infty}^{\infty} \left[ -\ln \sqrt{2\pi}\sigma - \frac{(x-\mu)^2}{2\sigma^2} \right] p(x) \, dx$$

$$= \int_{-\infty}^{\infty} \ln \sqrt{2\pi}\sigma \, p(x) \, dx + \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{2\sigma^2} p(x) \, dx$$

$$= \ln \sqrt{2\pi}\sigma + \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{2\sigma^2} p(x) \, dx \qquad \left( \text{since} \int_{-\infty}^{\infty} p(x) \, dx = 1 \right)$$

To compute the second integral term,

$$\int_{-\infty}^{\infty} \frac{(x-\mu)^2}{2\sigma^2} p(x) \, dx = \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{2\sigma^2} \frac{1}{\sqrt{2\pi}\sigma} \exp^{-}\frac{(x-\mu)^2}{2\sigma^2} \cdot dx$$

Write $x_1 = \frac{x-\mu}{\sqrt{2}\sigma}$.  $dx_1 = \frac{dx}{\sqrt{2}\sigma}$.  when $x \to \infty$, $x_1 \to \infty$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $x \to -\infty$, $x_1 \to -\infty$

Then we can write this term as:

$$\int_{-\infty}^{\infty} x_1^2 \frac{1}{\sqrt{2\pi}\sigma} \exp(-x_1^2) \, dx = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} x_1^2 \exp(-x_1^2) \frac{dx}{\sqrt{2}\sigma}$$

$$= \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} x_1^2 \exp(-x_1^2) \, dx_1$$

Integrating by parts:

(3)

Integrating by parts

$$= \frac{1}{\sqrt{\pi}} x_1 \int_{-\infty}^{\infty} x_1 e^{-x_1^2} dx_1 - \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} 1 \cdot \frac{1}{2} (-e^{-x_1^2}) dx_1$$

$$= \frac{x_1}{\sqrt{\pi}} \left[ \frac{e^{-x_1^2}}{2} \right]_{-\infty}^{\infty} + \frac{1}{\sqrt{\pi}} \int \frac{1}{2} e^{-x_1^2} dx_1$$

$$= 0 + \frac{1}{\sqrt{\pi}} \frac{1}{2} \int_{-\infty}^{0} e^{-x_1^2} dx_1$$

$$= \frac{1}{2} \cdot \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2}\sigma} \exp\left(\frac{-(x-\mu)^2}{2\sigma^2}\right)$$

$$= \frac{1}{2} \cdot \frac{1}{\sqrt{2\sigma}} \int_{\sigma}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(x-\mu)^2}{2\sigma^2}\right)$$

$$= \frac{1}{2} \cdot 1 = \frac{1}{2}.$$

Hence,

$$\boxed{H(X) = \ln(\sqrt{2\pi}\sigma) + \frac{1}{2}}$$ ✓

This is same as

$$H(X) = \frac{1}{2} \ln\left(\sqrt{2\pi}\,\sigma\right)^2 + \frac{1}{2}$$

$$= \frac{1}{2}\left( \ln(2\pi\sigma^2) + 1 \right).$$

④

PROBLEM 2

1. $p(\text{disease}) = 0.01$

   $p(\sim\text{disease}) = 0.99$

   $p(+\text{ve} \mid \text{disease}) = 0.95$

   $p(+\text{ve} \mid \sim\text{disease}) = 0.05$

(a) $P(+\text{ve}) = p(+\text{ve}, \text{disease}) + p(+\text{ve}, \sim\text{disease})$

   $= p(+\text{ve} \mid \text{disease})\, p(\text{disease})$

   $\quad + p(+\text{ve} \mid \sim\text{disease})\, p(\sim\text{disease})$

   $= 0.95 \times 0.01 + 0.05 \times 0.99$

   $= 0.059$ ✓

(b) $P(\text{disease} \mid +\text{ve}) = \dfrac{P(\text{disease}, +\text{ve})}{P(+\text{ve})}$

   $= \dfrac{0.95 \times 0.01}{0.059}$

   $= 0.161.$ ✓

2) $\lambda_{MLE} = \underset{\lambda}{\text{argmax}}\; P(x_1 x_2 \cdots x_n \mid \lambda).$

Since the random variables are drawn independently,

$$P(x_1 x_2 \cdots x_n \mid \lambda) = P(x_1 \mid \lambda)\, P(x_2 \mid \lambda) \cdots P(x_n \mid \lambda).$$

$$\ln P(x_1 x_2 \cdots x_n \mid \lambda) = \sum_{i=1}^{n} P(x_i \mid \lambda)$$

⑤

$\dfrac{30}{30}$

$$\frac{d}{d\lambda} \ln P(x_i/\lambda)$$

$$\frac{d}{d\lambda} \ln P(x_i) = x_i \ln \lambda - \lambda - \ln x_i!$$

$$\frac{d}{d\lambda} \ln P(x_i/\lambda) = x_i \cdot \frac{1}{\lambda} - 1$$

$$\therefore \frac{d}{d\lambda} \ln P(x_1, x_2 \cdots x_n /\lambda)$$

$$= \frac{d}{d\lambda} \sum_{i=1}^{n} \ln P(x_i/\lambda)$$

$$= \sum_{i=1}^{n} \frac{x_i}{\lambda} - \sum_{i=1}^{n} 1$$

$$= \frac{\sum x_i}{\lambda} - n \quad \checkmark$$

for maximum $\lambda$, $\frac{d}{d\lambda} \ln P(x_1, x_2 \cdots x_n/\lambda) = 0$

$$\therefore \frac{\sum x_i}{\lambda_{MLE}} - n = 0$$

$$\underline{\lambda_{MLE} = \hat{\lambda}_{MLE} = \frac{\sum x_i}{n}} \quad \checkmark$$

To show that it is an unbiased estimate, we write

$$E(\hat{\lambda}) = E\left(\frac{\sum x_i}{n}\right) = \frac{1}{n} E\left(\sum x_i\right)$$

$$= \frac{1}{n} \sum E(x_i) \qquad \text{(from linearity}$$
$$\text{of expectation)}$$

$$= \frac{1}{n} \cdot n\lambda$$

$$= \lambda. \quad \checkmark$$

(b)

$$\cancel{p(\lambda/x) \propto p(x/\lambda)}$$

$$p(\lambda / x_1, x_2 \cdots x_n) \propto p(x_1, x_2 \cdots x_n \mid \lambda) \, p(\lambda)$$

$$p(x_1, x_2 \cdots x_n \mid \lambda) = \prod_{i=1}^{n} p(x_i \mid \lambda)$$

$$= \frac{\lambda^{\Sigma x_i} \, e^{-n\lambda}}{\prod_{i=1}^{n} x_i !}$$

$$p(\lambda) = p(\lambda \mid \alpha, \beta) = \frac{\beta^{\alpha} \, \lambda^{\alpha-1} \, e^{-\beta\lambda}}{\Gamma(\alpha)}.$$

$$\therefore \ p(x_1, x_2 \cdots x_n \mid \lambda) \propto \frac{\lambda^{\Sigma x_i} \, e^{-n\lambda}}{\prod (x_i !)} \cdot \frac{\beta^{\alpha} \, \lambda^{\alpha-1} \, e^{-\beta\lambda}}{\Gamma(\alpha)}$$

$$= \frac{\lambda^{\Sigma x_i + \alpha - 1} \, e^{-(\beta\lambda + n\lambda)}}{\prod (x_i !) \, \Gamma(\alpha)} \cdot \checkmark$$

(c) To compute MAP, we take

$$\underset{\lambda}{\text{argmax}} \ p(x_1, x_2 \cdots x_n \mid \lambda).$$

$$\ln P(x_1, x_2 \cdots x_n \mid \lambda) \propto (\Sigma x_i + \alpha - 1) \ln \lambda - \lambda(\beta + n)$$
$$- \Sigma \ln x_i ! - \ln \Gamma(\alpha)$$

⑦

$$\frac{d}{d\lambda} P(x_1, x_2 \cdots x_n | \lambda) = \frac{(\sum x_i + \alpha - 1)}{\lambda} - (\beta + n) = 0$$

$$\therefore \lambda_{MAP} = \frac{\sum x_i + \alpha - 1}{\beta + n} \checkmark$$

PROBLEM 3

3.1> Attached the code added to entropy.c and prune-dt.c $\boxed{6/20}$

3.2> 1) For the fully-grown tree: $\boxed{\dfrac{0}{20}}$

   Tree size: depth = 9

       Nodes = 768

  Accuracy on training set: 90.3% ✗

  Accuracy on testing set: 87.1%

 2) For post-pruning with top-down approach:

   Tree size: depth = 7

       Nodes = 184

  Accuracy on training set: 89.5% ✗

  Accuracy on testing set: 88.3%

 3) For post-pruning with bottom-up approach

   Tree size: depth = 9

       Nodes = 512

  Accuracy on training set: 89.9% ✗

  Accuracy on test set: 89.0%

→ For the case where the tree is not pruned, the size is highest. This also has the highest training accuracy, and lowest test accuracy, because it overfits the training data.

→ Both the prunings lead to an increase in test accuracy.

→ With the top down approach, if the accuracy increases by making the current node a leaf, we remove the subtree belonging the node.

In the bottom up approach, we first check the subtree of the current node. Thus, in the bottom up approach, we might decide it to retain the node, if, after pruning its subtree, we have already attained a high accuracy, (and hence pruning current node does not increase accuracy).

→ Thus, probability of current node being pruned in bottom up approach is lower. Hence tree size is higher and tree depth is higher for bottom up approach.

→ also, in bottom up approach, instead of removing the whole subtree of current node, we look for best possible subtree.
(Top down approach is a special case of this, where removing the whole subtree gives best accuracy).

Thus, bottom up approach takes into account more combinations of nodes, and gives better training and test accuracies. ✗

3.3 >

| Epsilon | No. of nodes | |
|---|---|---|
| | Bottom up | top down |
| 0.001 → | 635 | 90 |
| 0.005 → | 512 | 184 |
| 0.01 → | 752 | 752 |
| 0.05 → | 768 | 768 |

✗

→ In general, as epsilon increases, probability that node will get pruned decreases, Hence less nodes are pruned, and tree size gets higher.

→ Only exception to this is $\epsilon = 0.001$ for bottom-up approach. Here, subtree is checked first. Although it is more probable for a subtree node to get pruned, that leads to the node a accuracy changes, and then, removing the current node, does not increase accuracy in many cases. So, nodes removed is not very high. ⑪