**a)** $KL(P\|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$

Prove $KL(P\|Q)$ is non negative.

**c)**

$KL(P\|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$.

$= -\sum_x P(x) \log \frac{Q(x)}{P(x)}$.

$= -E\left[\log \frac{Q}{P}\right]$

$\geq -\log\left(E\left[\frac{Q}{P}\right]\right)$ $\qquad$ { Jensen's Inequality.

$= -\log\left(\sum_{x \in X} P(x) \frac{Q(x)}{P(x)}\right)$

$= 0$.

The inequality is introduced due to the application of Jensen's Inequality & the concavity of log.

**b)** Ans $\quad KL(P(x,y)\|Q(x,y)) = KL(P(x)\|Q(x)) + KL(P(y|x)\|Q(y|x))$

**Proof**

$KL(P(x,y)\|Q(x,y)) = \sum_x \sum_y P(x,y) \log \frac{P(x,y)}{Q(x,y)}$

$= \sum_x \sum_y P(x,y) \log\left[\frac{P(x) P(y|x)}{Q(x) Q(y|x)}\right]$ $\qquad$ { By conditional probability

$= \sum_x \sum_y P(x,y) \log \frac{P(x)}{Q(x)} + \sum_x \sum_y P(x,y) \log \frac{P(y|x)}{Q(y|x) Q(y|x)}$

$= KL(P(x)\|Q(x))$

$= \sum_x \sum_y P(x,y) \log \frac{P(x)}{Q(x)} + \sum_x \sum_y P(x,y) \log \frac{P(y|x)}{Q(y|x)}$.

$= \sum_x \sum_y P(x,y) \log \frac{P(x)}{Q(x)} + \sum_x P(x) \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)}$

$= KL(P(x)\|Q(x)) + KL(P(y|x)\|Q(y|x))$

Hence proved.

i) For $\hat{P}(x) = 1/m \sum_{i=1}^{m} 1\{x^{(i)} = x\}$, for family of distributions $P_\theta$.

Prove that

$$\arg\min_\theta KL(\hat{P} \| P_\theta) = \arg\max_\theta \sum_{i=1}^{m} \log P_\theta(x^{(i)})$$

This indicates that finding the maximum likelihood estimate for the parameter $\theta$ is equivalent to finding $P_\theta$ with minimal KL divergence from $\hat{P}$.

**Ans.** By KL divergence.

$$KL(P\|Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx \quad ①$$

$\longrightarrow$ Let $\hat{P}(x)$ be the empirical distribution.

$$\hat{P}(x) = 1/m \sum_{i=1}^{m} 1\{x^{(i)} = x\}.$$

$$KL[\hat{P}(x) \| P(x|\theta)] = \int \hat{P}(x) \log \frac{\hat{P}(x)}{P(x|\theta)} dx$$

$$= -H(\hat{P}) - \int \hat{P}(x) \log[P(x|\theta)] dx \qquad \text{where } H(\hat{P}) = -\int \hat{P}(x) \log \hat{P}(x) dx$$

$$\qquad \qquad \qquad \qquad ②$$

From ② it follows that

$$\arg\min_\theta KL[\hat{P}(x) \| P(x|\theta)] P(x|\theta)] = \arg\max_\theta \left( \log P(x|\theta) \right)_{\hat{P}}$$

where $(\cdots)_{\hat{P}}$ represents expectation with respect to the distribution of $\hat{P}$.

using ① in ② RHS.

$$\left( \log P(x|\theta) \right)_{\hat{P}} = \frac{1}{m} \int \sum_{t=1}^{m} \hat{P}(x) \log P(x|\theta) dx.$$

$$= \frac{1}{m} \sum_{t=1}^{m} \log P(x_t|\theta)$$

Apart from the scaling factor $1/m$, this is a log-likelihood function. ⊠.

2) EM for MAP estimation

$$l(\theta) = \sum_{i=1}^{m} \log P(x^{(i)}|\theta) + \log P(\theta)$$

$$= \sum_{i=1}^{m} \log \sum_{z} P(x^i, z^i|\theta) + \log P(\theta)$$

$$= \sum_{i=1}^{m} \log \sum_{z} q_i(z^i) \frac{P(x^i, z^i|\theta)}{q_i(z_i)} + \log P(\theta)$$

$$\geq \sum_{i=1}^{m} \sum_{z^i} q_i(z^i) \log \left( \frac{P(x^i, z^{(i)}|\theta)}{q_i(z^i)} \right) + \log P(\theta)$$

The above equality holds when

$$\frac{P(x^i, z^i|\theta)}{q_i(z^i)} = c \quad \& \text{ since } \sum_z q_i(z^i) = 1,$$

$$q_i(z^i) = \frac{c \, P(x^i, z^i|\theta)}{P(x^i|\theta)}$$

EM for MAP is following:
Repeat until convergence
{ (E-step) For each $i$    $q_i(z^i) = P(z^i|x^i, \theta)$

(M-step) set

$$\theta = \arg\max_\theta \sum_{i=1}^{m} \sum_{z^i} q_i(z^i) \log \frac{P(x^i, z^i|\theta)}{q_i(z^i)} + \log P(\theta)$$
}

Prove that $l(\theta) = \sum_{i=1}^{m} \log P(x^i|\theta) + \log P(\theta)$ monotonically increases with each iteration. This is to just prove that $l(\theta^i) \leq l(\theta^{i+1})$

Firstly, given $q_i(z^i) := p(z^i|x^i, \theta)$ the following equality holds.

$$l(\theta^t) = \sum_{i=1}^{m} \sum_{z^i} q_i(z^i) \log \frac{(P(x^i, z^i|\theta^t)}{q_i(z^i)} + \log P(\theta^t) = L(\theta^t)$$

with respect to the lower bound;

$$L(\theta) = \sum_{i=1}^{m} \sum_{z^i} q_i(z^i) \log \frac{P(x^i, z^i|\theta)}{q_i(z^i)} + \log P(\theta)$$

The next iteration is done by explicitly choosing $\theta^{t+1}$ w.r.t the posterior $q \cdot q_i^{(t)}$ to maximize $L(\theta)$ which means

$$L(\theta^t) \leq L(\theta^{t+1})$$

w.r.t the loss, the equality

$$l(\theta^{tr}) = L(\theta^{t+1})$$

$$= \sum_{i=1}^{m} \log P(x^i | \theta^{t+1}) + \log P(\theta^{t+1})$$

holds only for

$$q_{t+1}(z^i) = P(z^i | x^i, \theta^{t+1})$$

which means

$$q_{t+1}(z^{(i)}) = P(z^i | x^i, \theta^{t+1})$$

$$l(\theta^{t+1}) = L(\theta^{t+1}) \geq - L(\theta^t) \cdot l(\theta^t)$$