STA414 : Partial Problem Set Solutions.

## PS, Q4

$$l_x(\theta) = \log f(x|\theta)$$

a) $f(x,z|\theta) = f(x|\theta)f(z|x,\theta)$

$\Rightarrow f(x|\theta) = f(x,z|\theta) / f(z|x,\theta)$   (Bayes')

$$\boxed{\Rightarrow l_x(\theta) = \log f(x|\theta) = \log f(x,z|\theta) - \log f(z|x,\theta)}$$

b) $Q(\theta, \theta^t) \equiv E_{z|x,\theta^t}\left[\log f(x,z|\theta)\right]$ as defined in class.

Note: $\theta^t$, not $\theta$

Note: $\theta$, not $\theta^t$

$$H(\theta, \theta^t) \equiv E_{z|x,\theta^t}\left(\log f(z|x,\theta)\right)$$

c) i) The M-step maximizes $Q(\theta, \theta^t)$ w.r.t. $\theta$. So $Q(\theta^{t+1}, \theta^t) \geq Q(\theta, \theta^t) \ \forall \theta$; in particular, $Q(\theta^{t+1}, \theta^t) \geq Q(\theta^t, \theta^t)$

ii) $H(\theta^{t+1}, \theta^t) - H(\theta^t, \theta^t)$

$= E_{z|x,\theta^t}\left(\log f(z|x,\theta^{t+1}) - \log f(z|x,\theta^t)\right)$

$= E_{z|x,\theta^t}\left[\log \dfrac{f(z|x,\theta^{t+1})}{f(z|x,\theta^t)}\right]$

$\leq \log E_{z|x,\theta^t}\left(\dfrac{f(z|x,\theta^{t+1})}{f(z|x,\theta^t)}\right)$   (Jensen)

$= \log \int \dfrac{f(z|x,\theta^{t+1})}{f(z|x,\theta^t)} f(z|x,\theta^t) dz$

$= \log \int f(z|x,\theta^{t+1}) dz = \log(1) = 0$

P5,Q5.

The notation here is very confusing.

$$\{ \underline{X}_n \}_{n=1}^N \quad - \text{ SAMPLE.}$$ Each $\underline{X}_n$ is one datapoint

$$\underline{X}_n = (X_{n_1}, \ldots X_{n_D})$$ Each $X_n$ is composed of $D$

more → (possibly dependent) Bernoulli

trials. (INDEPENDENT)

-eg, I flip a coin $D$ times, and get a binary vector.
The $i^{th}$ flip had $P(X_{ni} = 1) = \mu_i$

The probability of getting any particular
sequence of results — any particular binary
sequence — is

$i^{th}$ flip is heads        $i^{th}$ flip is tails.

$$P(\underline{X}_n | \mu, \ldots \mu_D) = \prod_{i=1}^{D} \mu_i^{X_i} (1-\mu_i)^{1-X_i}$$

Observed     parameters     $P(i^{th}$ flip is     $P(i^{th}$ flip
data.                        heads)              is tails)

Now, we have $K > 1$ component distributions, each
with their own set of parameters $\{ \mu_k \}_{k=1}^K$, with

$$\mu_k = (\mu_{k_1}, \ldots \mu_{k_D})$$

Define, for each $\underline{X}_n$, $\underline{z}_n = (0, \ldots, 1, \ldots 0)$ as the
length-$K$ binary (latent) vector indicating which group
$X_n$ came from. If $P(X_n$ came from group $k) = \pi_k$, then

$$P(\underline{z}_n | \pi) = \prod_{k=1}^{K} \pi_k^{z_{nk}}$$

Finally, denote the actual ~~density~~ distribution for $X_n$ coming from group $K$ as $\cancel{P(X_n|Z_n)}$ $P(X_n|\mu_k)$. So the distribution of $X_n$ conditional on $Z_n$ is

$$P(X_n|Z_n) = \prod_{k=1}^{K} P(X_n|\mu_k)^{Z_{nk}}$$

which just equals $P(X_n|\mu_k)$ for the one and only $K$ for which $Z_{nk} = 1$.

a) $P(X_n|\mu,\pi) = \cancel{\sum_{Z} P(X_n|Z_n)}$

$\qquad = \sum_{Z} P(X_n, Z_n)$ ← Sum is taken over all $K$ possible $Z_n$ vectors.

$\qquad = \sum_{Z} P(X_n|Z_n) P(Z_n)$

$\qquad = \sum_{Z} \left[ \prod_{k=1}^{K} P(X_n|\mu_k)^{Z_{nk}} \times \prod_{k=1}^{K} \pi_k^{Z_{nk}} \right]$

$\qquad = \sum_{Z} \prod_{k=1}^{K} \left( \pi_k P(X_n|\mu_k) \right)^{Z_{nk}}$

**Key step!** $\begin{cases} = \pi_1 P(X_n|\mu_1) + \pi_2 P(X_n|\mu_2) + \cdots + \pi_K P(X_n|\mu_K) \\ \qquad\quad \uparrow \qquad\qquad\qquad \uparrow \qquad\qquad\qquad\qquad \uparrow \\ \quad Z=(1,0,\cdots 0) \quad Z=(0,1,\cdots,0) \qquad Z=(0,0,\cdots 1) \end{cases}$

$$\boxed{= \sum_{k=1}^{K} \pi_k P(X_n|\mu_k).}$$

b) $P(X) = \prod_{n=1}^{N} P(X_n)$

$\qquad = \prod_{n=1}^{N} \sum_{k=1}^{K} \pi_k P(X_n|\mu_k)$

$$\boxed{\Rightarrow \ell(\mu,\pi|X) = \log P(X) = \sum_{n=1}^{N} \log\left( \sum_{k=1}^{K} \pi_k P(X_n|\mu_k) \right)}$$

c) As stated in the problem, the complete-cluster dist$^n$ for each $x_n$ (ie its dist$^n$ if we knew $z_n$) is

$$P(x_n | \mu_k) = \prod_{i=1}^{D} \mu_i^{x_{ni}} (1-\mu_i)^{1-x_{ni}}$$

So $P(x_n, z_n) = P(x_n | z_n) P(z_n)$

$$= \prod_{k=1}^{K} \left( \prod_{i=1}^{D} \mu_i^{x_{ni}} (1-\mu_i)^{1-x_{ni}} \right)^{z_{nk}} \prod_{k=1}^{K} \pi_k^{z_{nk}}$$

and $P(X, Z) = \prod_{n=1}^{N} P(x_n, z_n)$

$$= \prod_{n=1}^{N} \prod_{k=1}^{K} \left( \pi_k \prod_{i=1}^{D} \mu_i^{x_{ni}} (1-\mu_i)^{1-x_{ni}} \right)^{z_{nk}}$$

which gives $\ell(\mu, \pi | X, Z) = \log P(X, Z)$

$$\boxed{= \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \left( \log \pi_k + \sum_{i=1}^{D} \left( x_{ni} \log \mu_i + (1-x_{ni}) \log(1-\mu_i) \right) \right)}$$

d) Dist$^n$ of $z_{nk} | X, \mu^+, \pi^+$ ...

$$P(z_{nk} = 1 | X) = \frac{P(X, 1)}{P(X, 1)} \qquad \text{Notation is hard!}$$

$$P(z_{nk} = 1 | x_n) = P(z_n = (0, \cdots \overset{k}{1} \cdots 0) | x_n)$$

$$= \frac{P(x_n, (0, \cdots, 1, \cdots 0))}{P(x_n)}$$

$$\boxed{= \frac{\pi_k P(x_n | \mu_k)}{\sum_{g=1}^{K} \pi_g P(x_n | \mu_g)}}$$

The question did ask for $P(z_{nk} | X)$, not $x_n$ ... but the x's are independent, so others don't affect the answer.

Now because complete-data log-likelihood is linear in $z_{nk}$, $Q(\theta, \theta^t)$ is obtained by plugging in $\hat{z}$ for $z$.

e) These are just the weighted MLE's for the Bernoulli dist$^t$ — you may attempt the calculus yourselves.

## P6Q2

c) Recall that a covariance matrix for a vector-valued random variable satisfies

$$\Sigma_{ij} = \text{Cov}(X_i, X_j)$$

so

$$\Sigma_{ii} = \text{Var}(X_i)$$

Hence $S_{ii}$ is a sample estimate of $\text{Var}(X_i)$. The sum of the variances of all the $X_i$ then is

$$\sum_{i=1}^{p} S_{ii} = \text{tr}(S) \quad (\text{trace})$$

But it is known for any matrix $A$ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_D$ that

$$\text{tr}(A) = \sum_{d=1}^{D} \lambda_d.$$

It is also true that the variance of the $d^{th}$ principal component is $\lambda_d$, the corresponding eigenvalue of $S$. Putting it all together gives

$$\sum_{d=1}^{D} \lambda_d / \text{tr}(S)$$

as the prop$^n$ of total variance explained by the

## P7 Q5

a) A saturated model has $\hat{P}_n = t_n$.

For a binomial* distribution, the distribution function of $t_n$ is

$$P(t_n = x) = \hat{P}_n^x (1-P_n)^{1-x}$$

Binom(1,p) is Bern(p)

So the likelihood for $P$ is

$$L(P) = \prod_{n=1}^{N} P(t_n|P)$$
$$= \prod_{n=1}^{N} P_n^{t_n} (1-P_n)^{1-t_n}$$

and the log-likelihood,

$$\ell(P) = \sum_{n=1}^{N} t_n \log P_n + (1-t_n)\log(1-P_n)$$

For the saturated model, this gives

$$\ell_{sat}(P) = \sum_{n=1}^{N} t_n \log t_n + (1-t_n)\log(1-t_n)$$

But $t_n \in \{0,1\}$, so the above is not defined unless we fudge it and say $0\log 0 = 0$. In that case, every term in the log-likelihood is

$$0\log 0 + 1\log 1 = 0.$$

b) Just plug $\bar{t}$ into the above formula. (see below)

c) $D = 2\left(\ell_{sat}(P) - \ell_{model}(P)\right)$

$\uparrow$
$= 0$.

$$= -2\ell_{model}(P)$$
$$= -2\sum_{n=1}^{N} t_n \log \hat{P}_n + (1-t_n)\log(1-\hat{P}_n)$$

d) $AIC = -2l_{model}(p) + 2d$ , $d = \#$ parameters

$BIC = -2l_{model}(p) + d \log N$.

Forward Stepwise selection:

I) Start with the null model, or the smallest model you would be willing to accept.

II) Do:
- Add each available feature separately into the model, and calculate the AIC/BIC for the resulting (more complex) model

- choose the feature that yields the largest decrease in AIC/BIC, and add it to the model permanently

- repeat.

UNTIL:
- AIC/BIC stops decreasing OR
- AIC/BIC " " "too much" OR
- All of the features have been added.

## P7Q3

Let $y = g(x)$. Note that
$$G_i(y) = P(Y_i < y)$$
$$= P(g(T_i) < g(x))$$
$$= P(T_i < x)$$
$$= F_i(x)$$

although this isn't quite enough to prove the statement; it is suggestive though.

We can write $F_1(x) - F_0(x) = \int_{-\infty}^{x}(f_1(s) - f_0(s))ds$.

Now,
$$G_1(x) - G_0(x) = \int_{-\infty}^{x}(g_1(s) - g_0(s))ds$$

$$= \int_{-\infty}^{g^{-1}(x)}\left(f_1(g^{-1}(s)) - f_0(g^{-1}(s))\right)g^{-1\,'}(s)ds.$$

$$= \int_{-\infty}^{u}(f_1(u) - f_0(u))du.$$

$$= F_1(u) - F_0(u)$$

$$\implies \max_x |G_1(x) - G_0(x)| = \max_u |F_1(u) - F_0(u)|,$$

so $KS = KS^{\dagger}$.