

1) a) From Appendix we have:

$$\begin{cases} p(z) = N(z | \mu, \Sigma) \\ p(x|z) = N(x | Az + b, S) \\ p(x) = N(x | A\mu + b, A\Sigma A^T + S) \\ p(z|x) = N(z | C(A^T S^{-1}(x-b) + \Sigma^{-1}\mu), C) \end{cases}$$

$$\begin{cases} z \sim N(0, 1) \\ p(z) = N(z | \mu, \Sigma) \end{cases} \Rightarrow \boxed{\mu = 0, \Sigma = 1}$$

$$\begin{cases} x|z \sim N(2u, 6^2 I) \\ p(x|z) = N(x | Az + b, S) \end{cases} \Rightarrow \begin{aligned} Az + b &= 2u \\ S &= 6^2 I \end{aligned} \Rightarrow \boxed{\begin{aligned} b &= 0 \\ A &= u \\ S &= 6^2 I \end{aligned}}$$

Therefore we have:

$$p(z) = N(0, 1)$$

$$p(x|z) = N(x | uz, 6^2 I)$$

$$p(x) = N(x | 0, 6^2 I + uu^T)$$

$$p(z|x) = N(z | C(u^T(6^2 I)^{-1}x), C) = N(z | C(\frac{1}{6^2}u^T x), C)$$

$$C = (1 + u^T(6^2 I)^{-1}u)^{-1} = (1 + \frac{1}{6^2}u^T u)^{-1}$$

$$m = E[z|x] = C(\frac{1}{6^2}u^T x) = \frac{\frac{1}{6^2}u^T x}{1 + \frac{1}{6^2}u^T u}$$

$$S = E[z^2|x] = C + E[z|x]^2 = C + m^2$$

$$b) \quad u_{\text{new}} \leftarrow \arg \max_u \frac{1}{N} \sum_{i=1}^N E_{q(z^{(i)})} [\log p(z^{(i)}, x^{(i)})]$$

By expanding out  $\log p(z^{(i)}, x^{(i)})$ :

$$u_{\text{new}} \leftarrow \arg \max_u \frac{1}{N} \sum E_{q(z^{(i)})} [\log p(x^{(i)} | z^{(i)}) + \log p(z^{(i)})]$$

By distributing the expectation:

$$u_{\text{new}} \leftarrow \arg \max_u \frac{1}{N} \sum_{i=1}^N E_{q(z^{(i)})} [\log p(x^{(i)} | z^{(i)})] + E_{q(z^{(i)})} [\log p(z^{(i)})]$$

Gaussian formula with zero mean and unit variance:

$$N(z, \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(z-\mu)^2}{\sigma^2}\right) \Rightarrow N(z, 0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2\right)$$

By substituting  $N(z, 0, 1)$  with  $p(z^{(i)})$ :

$$u_{\text{new}} \leftarrow \arg \max_u \frac{1}{N} \sum_{i=1}^N \left( E_{q(z^{(i)})} [\log p(x^{(i)} | z^{(i)})] + E_{q(z^{(i)})} \left[ \log \left( (2\pi)^{-1/2} \cdot \exp\left(-\frac{1}{2} (z^{(i)})^2\right) \right) \right] \right)$$

$$u_{\text{new}} \leftarrow \arg \max_u \frac{1}{N} \sum_{i=1}^N \left( \underbrace{E_{q(z^{(i)})} [\log p(x^{(i)} | z^{(i)})]}_{(A)} + E_{q(z^{(i)})} \left[ -\frac{1}{2} [\log(2\pi) + (z^{(i)})^2] \right] \right)$$

continue ->

Because we have  $X|Z \sim N(uZ, \sigma^2 I)$

$$\begin{aligned}
 & E_{q(z^{(i)})} \left[ \log \frac{1}{2\pi^{D/2} \sqrt{|\sigma^2 I|}} \exp \left[ (x^{(i)} - z^{(i)}u)^T (\sigma^2 I)^{-1} (x^{(i)} - z^{(i)}u) \right] \right] \\
 &= E_{q(z^{(i)})} \left[ -\frac{D}{2} \log 2\pi - \frac{D}{2} \log \sigma^2 + \frac{1}{\sigma^2} (x^{(i)} - z^{(i)}u)^T (x^{(i)} - z^{(i)}u) \right] \\
 &= E_{q(z^{(i)})} \left[ -\frac{D}{2} \log 2\pi - \frac{D}{2} \log \sigma^2 + \underbrace{\frac{1}{\sigma^2} x^{(i)T} x^{(i)}}_{B^{(i)}} - x^{(i)T} z^{(i)} u - u^T z^{(i)} x^{(i)} \right. \\
 &\quad \left. + u^T z^{(i)} z^{(i)T} u \right]
 \end{aligned}$$

$$\begin{aligned}
 &= E_{q(z^{(i)})} \left[ B^{(i)} - 2z^{(i)} x^{(i)T} u + z^{(i)2} u^T u \right] \\
 &= E_{q(z^{(i)})} [B^{(i)}] - 2x^{(i)T} u E_{q(z^{(i)})} [z^{(i)}] + u^T u E_{q(z^{(i)})} [z^{(i)2}]
 \end{aligned}$$

$$A = E_{q(z^{(i)})} [B^{(i)}] - 2m^{(i)} x^{(i)T} u + s^{(i)} u^T u$$

Therefore by substituting  $A$  in  $u_{new}$ :

$$u_{new} \leftarrow \arg \max_u \frac{1}{N} \sum_{i=1}^N (E_{q(z^{(i)})} [B^{(i)}] - 2m^{(i)} x^{(i)T} u + s^{(i)} u^T u) +$$

$$E_{q(z^{(i)})} \left[ -\frac{1}{2} (\log(2\pi) + z^{(i)})^2 \right])$$

$$\frac{\partial}{\partial u} \sum_{i=1}^N (E_{q(z^{(i)})} [B^{(i)}] - 2m^{(i)} x^{(i)T} u + s^{(i)} u^T u) + E_{q(z^{(i)})} \left[ -\frac{1}{2} \log(2\pi) + z^{(i)2} \right] = 0$$

$$\Rightarrow -2 \sum_{i=1}^N m^{(i)} x^{(i)} + 2 \sum_{i=1}^N s^{(i)} u = 0 \Rightarrow u_{new} = \frac{\sum_{i=1}^N m^{(i)} x^{(i)}}{\sum_{i=1}^N s^{(i)}}$$

2. We want to show that:

$$\|T^\pi Q_1(s,a) - T^\pi Q_2(s,a)\|_\infty \leq \gamma \|Q_1(s',a') - Q_2(s',a')\|_\infty$$

We start with the left part:

$$\max_{s,a} |T^\pi Q_1(s,a) - T^\pi Q_2(s,a)| = \max_{s,a} \left| \left[ r(s,a) + \gamma \sum_{s'} P(s'|a,s) \sum_{a'} \pi(a'|s') Q_1(s',a') \right] - \left[ r(s,a) + \gamma \sum_{s'} P(s'|a,s) \sum_{a'} \pi(a'|s') Q_2(s',a') \right] \right|$$

By factoring:

$$= \left| \gamma \sum_{s'} P(s'|a,s) \left[ \sum_{a'} \pi(a'|s') Q_1(s',a') - \sum_{a'} \pi(a'|s') Q_2(s',a') \right] \right|$$

$$= \left| \underbrace{\gamma \sum_{s'} P(s'|a,s)}_{\text{This is equal to 1 because it's the sum of all probabilities}} \sum_{a'} \pi(a'|s') [Q_1(s',a') - Q_2(s',a')] \right|$$

$$\leq \gamma \max_{s'} \left| \sum_{a'} \pi(a'|s') [Q_1(s',a') - Q_2(s',a')] \right|$$

$$\leq \gamma \max_{s'} \left| \max_{a'} [Q_1(s',a') - Q_2(s',a')] \right|$$

We apply the absolute value to the term inside max, so it is greater than the previous term.

$$\leq \gamma \max_{s',a'} |Q_1(s',a') - Q_2(s',a')|$$

3.

a)

$$r(S, A) = \begin{cases} 1 & \text{if } S = S_1 \\ 2 & \text{if } S = S_2 \end{cases}$$

$$\gamma = 0.9$$

	stay	switch
$S_1$		✓
$S_2$	✓	

$$Q(S_1, \text{stay}) = 1 + \gamma Q(S_1, \text{switch})$$

$$Q(S_1, \text{switch}) = 1 + \gamma Q(S_2, \text{switch})$$

$$Q(S_2, \text{stay}) = 2 + \gamma Q(S_2, \text{stay})$$

$$Q(S_2, \text{switch}) = 2 + \gamma Q(S_1, \text{switch})$$

by solving the 4 equations we get:

$$\begin{cases} Q(S_1, \text{stay}) = 18.1 \\ Q(S_1, \text{switch}) = 19 \\ Q(S_2, \text{stay}) = 20 \\ Q(S_2, \text{switch}) = 19.1 \end{cases}$$

	stay	switch
$S_1$	18.1	19
$S_2$	20	19.1

b)

We want to define function  $Q$  in a way that is equilibrium and also stays at  $S_1$ :

$$Q(S_1, \text{stay}) = 1 + \gamma Q(S_1, \text{stay})$$

$$0.1 Q(S_1, \text{stay}) = 1 \Rightarrow Q(S_1, \text{stay}) = 10$$

If we set the value of  $Q(S_1, \text{switch})$  to a value less than 10, the agent will stay at  $S_1$  and this would be a suboptimal policy.

For example here is a table showing the  $Q$  function:

	stay	switch
$S_1$	10	0
$S_2$	0	0