

## Homework Set 3 of CS229 class.

3 a) Done.

b) ~~Previously each pixel was represented by  $8 \times 3 = 24$  bits. Now, after compression, only 4 bits are required. This gives a compression ratio of  $24/4 = 6$ .~~

4(a) Regarding the derivation of the E-M steps:

~~The~~ E-step is obviously right per definition: we're just estimating the pdf of  $z$ , to softly classify the datapoints.

M-step:  $\log p(x^i, z^i; \theta) = \log Q(z^i) \frac{p(x^i, z^i; \theta)}{Q(z^i)} \geq Q(z^i) \log \frac{p(x^i, z^i; \theta)}{Q(z^i)}$  per Jensen's inequality, since "log" is a concave function. ~~thus~~ Continuing, maximizing  $L_{\text{sup}}$  is straight forward relative to  $L_{\text{lower}}$ , so for  $L_{\text{sup}}$  we don't need to bother with any lower-bounds.

$$\Rightarrow \hat{\theta}^{(t+1)} := \underset{\theta^{(t)}}{\operatorname{argmax}} \left[ \sum_{i=1}^n \left( \sum_{z^i} Q_i^*(z^i; \theta) \log \frac{p(x^i, z^i; \theta)}{Q_i(z^i; \theta)} \right) + \lambda \sum_{i=1}^n \log p(x^i, z^i; \theta) \right]$$

And then to the E-step again, and back...

Have this after  
 $t+1$  E-steps,  $t$  M-steps.

$$\begin{aligned}
 4(a2) \quad l_{\text{semi-sup}}(\Theta^{t+1}) &= \sum_{i=1}^n \log \sum_{z_i} p(x_i, z_i; \Theta^{t+1}) + l_{\text{sup}}(\Theta^{t+1}) \\
 &= \sum_{i=1}^n \log \sum_{z_i} Q(z_i) \frac{p(x_i, z_i; \Theta^{t+1})}{Q(z_i)} + l_{\text{sup}}(\Theta^{t+1}) \\
 &\xrightarrow{\text{Jensen's + log is concave}} \sum_{i=1}^n \sum_{z_i} Q(z_i) \log \frac{p(x_i, z_i; \Theta^t)}{Q(z_i)} + l_{\text{sup}}(\Theta^t) \\
 &\xrightarrow[\text{from argmaxing } l(\Theta^t) \text{ wrt } \Theta^t]{} \sum_{i=1}^n \sum_{z_i} Q(z_i) \log \frac{p(x_i, z_i; \Theta^t)}{Q(z_i)} + l_{\text{sup}}(\Theta^t) \\
 &\xrightarrow[\text{since } Q(z_i) \text{ was set to } p(z_i | x_i; \Theta^t)]{} = \sum_{i=1}^n \log \sum_{z_i} p(x_i, z_i; \Theta^t) + l_{\text{sup}}(\Theta^t) = l_{\text{semi-sup}}(\Theta^t). \text{ qed.}
 \end{aligned}$$

Making  $\sum_{z_i} Q(z_i) \log \frac{p(x_i, z_i; \Theta^t)}{Q(z_i)} = \log p(x_i; \Theta^t) = \log \sum_{z_i} p(x_i, z_i; \Theta^t)$

4(b) Clearly, all  $z_i$  for  $i = \{1, \dots, n\}$  (unlabelled cases) need to be estimated. The  $\hat{z}_i$ 's are all known.

$$Q_i^t(z_i) = p(z_i | x_i; \Theta^t) = \frac{p(x_i, z_i; \Theta^t)}{p(x_i; \Theta^t)} = \frac{p(x_i | z_i; \Theta^t) p(z_i; \Theta^t)}{\sum_{z_i} p(x_i | z_i; \Theta^t) p(z_i; \Theta^t)}$$

$$\begin{aligned}
 &\Leftrightarrow \\
 p(z_i=j | x_i; \Theta^t) &= \frac{p(x_i | z_i=j; \Theta^t) p(z_i=j; \Theta)}{\sum_{j=1}^k p(x_i | z_i=j; \Theta^t) p(z_i=j; \Theta)} = \cancel{\frac{p(x_i | z_i=j; \Theta^t) p(z_i=j; \Theta)}{\sum_{j=1}^k p(x_i | z_i=j; \Theta^t) p(z_i=j; \Theta)}}
 \end{aligned}$$

$$\frac{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}{Q} \cdot \varphi_j / (\sqrt{2\pi}^d |\Sigma_j|^{1/2})$$

$$\sum_{i=1}^k \tilde{e}^{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)} \cdot \varphi_i / (\sqrt{2\pi}^d |\Sigma_i|^{1/2})$$

C)  $\Theta = \{\mu, \Sigma, \varphi\}$  needs to be re-estimated, while  $Q^t$  is kept constant (since otherwise we'd be ~~iterating~~ directly, not iteratively, finding the maximum of  $l_{\text{semi-sup}}$ ).

$$\begin{aligned}
 \Rightarrow \nabla_\mu l(\Theta) &= \nabla_\mu \left[ \sum_{i=1}^n \sum_{z_i} Q_i^t \log \frac{p(x_i, z_i; \Theta)}{Q_i^t} + \alpha \sum_{i=n+1}^m \log (p(x_i^m, z_i^m; \Theta)) \right] \\
 &= \nabla_\mu \left[ \sum_{i=1}^n \sum_{z_i} Q_i^t \log \frac{p(x_i^i | z_i^i; \Theta)}{p(z_i^i; \Theta)} + \alpha \sum_{i=n+1}^m \log (p(x_i^m | z_i^m; \Theta) p(z_i^m; \Theta)) \right] \\
 &= \Omega
 \end{aligned}$$

$$= \nabla_{\mu} \left[ \sum_{i=1}^n \sum_{z^i} Q_i^t \log \frac{-\frac{1}{2} (x^i - \mu_z^i)^T \Sigma_z^{-1} (x^i - \mu_z^i)}{\sqrt{2\pi^d} |\Sigma_z|^{1/2}} + \alpha \sum_{i=n}^{\tilde{n}} \log \frac{\frac{1}{2} (x^i - \mu_{\tilde{z}}^i)^T \Sigma_{\tilde{z}}^{-1} (x^i - \mu_{\tilde{z}}^i)}{\sqrt{2\pi^d} |\Sigma_{\tilde{z}}|^{1/2}} P(z^i; \theta) \right]$$

$$\Leftrightarrow O = \nabla_{\mu} \left[ \sum_{i=1}^n \sum_{z^i} Q_i^t (x^i - \mu_{z^i})^T \Sigma_z^{-1} (x^i - \mu_{z^i}) + \alpha \sum_{i=n}^{\tilde{n}} \underbrace{P(z^i; \theta)}_{\text{multinomial dist.}} \cdot (x^i - \mu_{z^i})^T \Sigma_z^{-1} (x^i - \mu_{z^i}) \right]$$

$$\Leftrightarrow O = \nabla_{\mu_L} \left[ \sum_{i=1}^n Q_i^t (z^i = l; \theta) \cdot (x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L) + \alpha \sum_{i=n}^{\tilde{n}} \delta(z^i = l) \cdot (x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L) \right]$$

$$\Leftrightarrow O = \left[ \sum_{i=1}^n Q_i^t (z^i = l) \sum_L (x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L) + \alpha \sum_{i=n}^{\tilde{n}} \delta(z^i = l) \sum_L (x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L) \right]$$

(distr. mult. by  $\alpha / \sum_L$ )

$$\Leftrightarrow \mu_L \left[ \sum_{i=1}^n Q_i^t (z^i = l) + \alpha \sum_{i=n}^{\tilde{n}} \delta(z^i = l) \right] = \sum_{i=1}^n Q_i^t (z^i = l) x^i + \sum_{i=n}^{\tilde{n}} \delta(z^i = l) x^i$$

$$\Leftrightarrow \mu_L = \frac{\sum_{i=1}^n Q_i^t (z^i = l) x^i + \sum_{i=n}^{\tilde{n}} \delta(z^i = l) x^i}{\sum_{i=1}^n Q_i^t (z^i = l) + \sum_{i=n}^{\tilde{n}} \alpha \delta(z^i = l)} = \frac{\sum_{i=1}^n w^t(z^i = l) x^i}{\sum_{i=1}^n w^t(z^i = l)}$$

with  $w^t(z^i = l)$

~~$\nabla_{\mu} \left[ \sum_L \sum_{i=1}^n Q_i^t (z^i = l; \theta) (x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L) \right]$~~

~~$\nabla_{\mu} \left[ \sum_L \sum_{i=1}^n \omega_i^t \log \frac{\frac{1}{2} (x^i - \mu_z^i)^T \Sigma_z^{-1} (x^i - \mu_z^i)}{\sqrt{2\pi^d} |\Sigma_z|^{1/2}} \right]$~~

with  $w^t(z^i = l) = \begin{cases} Q_i^t(z^i = l), & i \leq n, \\ \alpha \delta(z^i = l), & i > \tilde{n} \end{cases}$

~~$\nabla_{\Sigma} \left[ \sum_{i=1}^n \sum_{z^i} \omega_i^t \log \frac{\frac{1}{2} (x^i - \mu_z^i)^T \Sigma_z^{-1} (x^i - \mu_z^i)}{\sqrt{2\pi^d} |\Sigma_z|^{1/2}} \right] = 0$~~

~~$\nabla_{\Sigma_L} \left[ \sum_{i=1}^n \omega_i^t (z^i = l) \left[ \frac{(x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L)}{\sqrt{2\pi^d} |\Sigma_L|^{1/2}} \right] \right] = 0$~~

~~$= \sum_{i=1}^n \omega_i^t (z^i = l) \left[ \frac{(x^i - \mu_L)^T \Sigma_L^{-1} (x^i - \mu_L)}{\sqrt{2\pi^d} |\Sigma_L|^{1/2}} - \frac{1}{\sqrt{2\pi^d} |\Sigma_L|^{1/2}} \right]$~~

~~$= \sum_{i=1}^n \omega_i^t (z^i = l) \left( -\sum_L (x^i - \mu_L)^T \Sigma_L^{-1} + \frac{|\Sigma_L|^{-1}}{|\Sigma_L|} \right)$~~

$$= \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) \left( -(\mathbf{x}^i - \mu_e) (\mathbf{x}^i - \mu_e)^T + \Sigma_e \right)$$

$$\Leftrightarrow \underline{\sum_e} = \frac{\sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) (\mathbf{x}^i - \mu_e) (\mathbf{x}^i - \mu_e)^T}{\sum_{i=1}^{n+\tilde{n}} w^t(z^i=l)}$$

Finally, for  $\phi$ : Here we need to add a Lagrange multiplier to  $l(\theta)$  that enforces  $\sum_{i=1}^k \phi_i = 1$ .

$$\Rightarrow \nabla_{\phi} [l(\theta) + \lambda \left( \sum_{i=1}^k \phi_i - 1 \right)] = 0 \Leftrightarrow \nabla_{\phi_e} \left[ \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) \log(p(x^i | z^i=l; \theta)) \right] + \lambda \left[ \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) \right] = 0$$

$$\# \# - \lambda \left[ \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) \right] = 0$$

$$\Leftrightarrow \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) \frac{1}{\phi_e} = \lambda \Leftrightarrow \frac{1}{\lambda} \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l) = \underline{\phi_e} . \text{ Now,}$$

$$\sum_{e=1}^k \phi_e = 1 = \frac{1}{\lambda} \sum_{i=1}^{n+\tilde{n}} \sum_{e=1}^k w^t(z^i=l) \Rightarrow \underline{\lambda} = \underline{\sum_i \sum_e w^t(z^i=l)}$$

$$\Rightarrow \underline{\phi_e} = \frac{\sum_{i=1}^{n+\tilde{n}} w^t(z^i=l)}{\sum_{e=1}^k \sum_{i=1}^{n+\tilde{n}} w^t(z^i=l)} = \frac{\sum_{i=1}^{n+\tilde{n}} w^t(z^i=l)}{h + \alpha \tilde{n}}$$

Note that  $w_i^t$  is a conditional probability  $p(z^i | x^i; \theta^t)$  that is computed during E-Step for  $i=1, 2, \dots, n$ , and set to some delta-distribution.

4 d) Done.

4 e) Done.

4 f)

i) Convergence is much faster for the semi-supervised one.

ii) Stability is ~~poor~~ for semi-sup., the unsup. one is <sup>rather</sup> ~~more~~ unstable.

iii) Overall Quality is much better in the semi-supervised case.

5) The  $u$  which minimises  $\|x - u\|^2$  is  $f_u(x) = \frac{x^T u}{\|u\|} u$  given  $u = \underline{\underline{u}}$ .

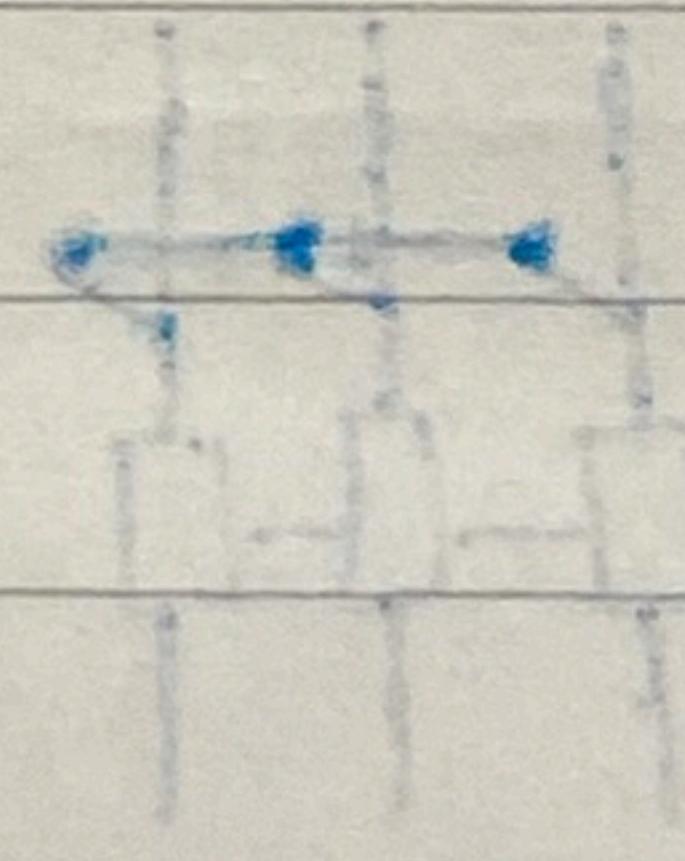
Plugging this result into the latter cost function gives:

$$\begin{aligned}
 L &= \|x^i - \underbrace{(u^T x^i) u}_{= u^T x^i} u\|_2^2 = ((I - uu^T)x^i)^T (I - uu^T)x^i = x^{i^T} (I - uu^T)x^i \\
 &= x^{i^T} (I - uu^T)x^i = x^{i^T} x^i - u^T x^i x^{i^T} u = \text{tr} \sum_i - u^T \sum_i u \\
 &\stackrel{\text{sum of } X \text{ from } \Sigma_i = x^T x^i}{=} \text{tr} (U S^2 U^T) - u^T U S^2 U^T u = \text{tr} S^2 - u^T [S, U, S_2 U_2, \dots] [S U_1, S_2 U_2, \dots] u \\
 &= \text{tr} S^2 - (S_1^2 (u^T U_1)^2 + S_2^2 (u^T U_2)^2 + \dots) = S_1^2 + S_2^2 + \dots
 \end{aligned}$$

because the minimum of  $L$  is found when  $u^T U_i = 1 \Rightarrow \underline{\underline{u}} = U_1$ .

Therefore, the first PC of the data corresponds to ~~the~~ the axis of maximum variation.

OA



## Homework Problem 6

a) If  $g(s) = \frac{e^{-\frac{s^2}{2}}}{\sqrt{\pi}}$ , then  $\log g' = -s^2/2 + \dots$ , which in turn implies

$$\begin{aligned} l(W) &= \sum_{i=1}^n \left( \log |W| - \frac{1}{2} \sum_{j=1}^d (w_j^T x_i)^2 \right) = n \log |W| - \frac{1}{2} \sum_{i=1}^n x^T W^T W x_i \\ &= n \log |W| - \frac{1}{2} \text{tr}(x^T W^T W x). \end{aligned}$$

$$\Rightarrow \frac{\partial l}{\partial W} = \frac{n}{|W|} W W^T - \frac{1}{2} x^T W x x^T = n W^T - W x x^T = 0$$

$\Rightarrow n I = W^T W x x^T \Rightarrow W^T W = \underline{n(x x^T)}^{-1}$ . This expression shows that there is no way to compute a unique  $W$ , because if  $W=W_0$  is a solution, then  $W=AW_0$  for any unitary  $A$  would also be a solution. ~~inherent in the missing matrix multiplication~~ ~~because no solution is possible. Thus the source terms cannot be randomly distributed like this for a solution to be found.~~

There is therefore a rotational invariance that exists for  $W$  which for Gaussian <sup>indeed</sup> source terms which precludes a unique solution to the problem.

$$b) l(W) = n \log |W| + \sum_{i,j} \log \frac{1}{2} e^{-|s_j^i|} \sim n \log |W| + \sum_{i,j} -|w_j^T x_i|$$

$$= n \log |W| - \sum_i |W x_i|. \Rightarrow \frac{\partial l}{\partial W} = \frac{n}{|W|} W^T - \sum_i \text{sgn}(W x_i) x^i = 0$$

$$\Rightarrow W := W + \alpha \frac{\partial l}{\partial W} = W + \alpha (n W^T - \sum_i \text{sgn}(W x_i) x^i) \text{ in batch gradient descent, and}$$

$$W := W + \alpha (W^T - \text{sgn}(W x_i) x^i) \text{ for single-update (stochastic) gradient descent.}$$

c) Ok.

CS229 Problem set #3

1.  $a \in \{L, R\}$ ,  $S = \{x, \dot{x}, \theta, \dot{\theta}\}$ ,  $\gamma = 0.995$  are all known beforehand. Reward  $R = R(S)$  and  $P_{sa}$  need to be learned.  
Eg through:

$$P_{sa}(s') = \frac{\#_{s,a \rightarrow s'}}{\#_{s,a}} , \quad R(S) = \frac{\sum_i R_i \#_{i \in S}}{\sum_i \#_{i \in S}}$$

Then Bellman's equations are to be solved to find the policy  $\pi^*$  which optimally selects actions at states:

$$\pi^*: V^*(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

$$\Rightarrow \tilde{\pi}_{(s)}^* = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

8:20

Answers to questions:

- 286 in my case
- Done.
- Seed "1" gives convergence at 299, seed 2 @ 135, seed 3 @ 377.  
 $\Rightarrow$  Algorithm is quite sensitive to initial condition, likely due to the nature of the underlying physical system (chaotic?).