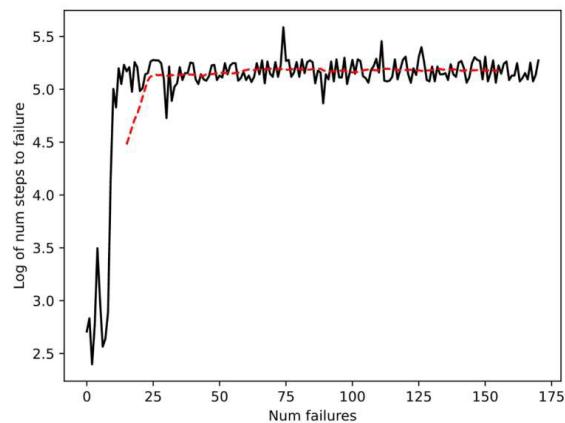


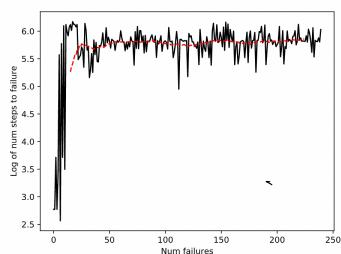
Question 1

i) It takes about 175 trials to converge

ii)



seed 1 — 240-



iii) As seed index changes to 1, 2, 3, the model has different convergance time than index at 0. In general, the time for the model to converge is longer at index = 1, 2, 3. Since the random generator is completely depending on the random seed, it directly influences the initial value vector. Hence, we can imply that the initial value vector can affect the performance of our reinforced learning model.

## Question 2

(a) Given  $P(x) > 0, Q(x) > 0 \quad \forall P, Q.$

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \left[ \frac{P(x)}{Q(x)} \right]$$

let  $f(x) = -\log(x)$  be the strictly convex function

Since  $X$  is a RV. we can apply the Jensen's Inequality:

$$\begin{aligned} D_{KL}(P||Q) &= \mathbb{E}_{x \sim P} \left[ -\log \frac{Q(x)}{P(x)} \right] \\ &\geq -\log \mathbb{E}_{x \sim P} \left[ \frac{Q(x)}{P(x)} \right] \\ &= -\log \sum_{x \in X} P(x) \cdot \frac{Q(x)}{P(x)} \\ &= -\log \sum_{x \in X} Q(x) \quad [Q(x) \text{ is distribution of } X] \\ &= -\log 1 \\ &= 0. \end{aligned}$$

Hence,  $D_{KL}(P||Q) \geq 0 \quad \forall P, Q.$

" $\Rightarrow$ "

$$\text{if } P = Q. \text{ we have } D_{KL}(P||Q) = \mathbb{E}_{x \sim P} \left[ -\log \frac{Q(x)}{P(x)} \right] \geq -\log \mathbb{E}_{x \sim P} \left[ \frac{Q(x)}{P(x)} \right]$$

RHS = 0 as proved in part i.

$$\text{LHS} = \mathbb{E}_{x \sim P} \left[ -\log 1 \right] = \sum_{x \in X} P(x) \cdot 0 = 0 = \text{RHS}.$$

proved!

" $\Leftarrow$ "

$$\text{Suppose } D_{KL}(P||Q) = 0 = \mathbb{E}_{x \sim P} \left[ -\log \frac{Q(x)}{P(x)} \right]$$

$$\text{because we know } -\log \mathbb{E}_{x \sim P} \left[ \frac{Q(x)}{P(x)} \right] = 0$$

$$\text{Hence } E[f(x)] = f(E[x])$$

which implies that  $X = E[X]$  with probability 1.

So, the distribution  $P, Q$  over  $X$  must be the same.

$$P = Q.$$



$$\begin{aligned}
 (b) D_{KL}(P(X,Y) || Q(X,Y)) &= \sum_x \sum_y P(x,y) \log \frac{P(x,y)}{Q(x,y)} \\
 &= \sum_x \sum_y P(x,y) \log \frac{P(y|x)P(x)}{Q(y|x)Q(x)} \\
 &= \sum_x \sum_y P(x,y) \log \frac{P(x)}{Q(x)} + \sum_x \sum_y P(x,y) \log \frac{P(Y|x)}{Q(Y|x)}
 \end{aligned}$$

The first part :

$$\begin{aligned}
 \sum_x \sum_y P(x,y) \log \frac{P(x)}{Q(x)} &= \sum_x \sum_y P(x)P(y|x) \log \frac{P(x)}{Q(x)} \\
 &= \sum_x P(x) \sum_y P(y|x) \log \frac{P(x)}{Q(x)} \quad \text{since } y \text{ is irrelevant} \\
 &= \sum_x P(x) \log \frac{P(x)}{Q(x)} = D_{KL}(P(X) || Q(X))
 \end{aligned}$$

The second part :

$$\begin{aligned}
 \sum_x \sum_y P(x,y) \log \frac{P(Y|x)}{Q(Y|x)} &= \sum_x \sum_y P(y) P(x|y) \log \frac{P(Y|x)}{Q(Y|x)} \\
 &= \sum_y P(y) \left( \sum_x P(x|y) \log \frac{P(Y|x)}{Q(Y|x)} \right) \\
 &= D_{KL}(P(X|Y) || Q(X|Y))
 \end{aligned}$$

Hence, the chain rule is proved.

(c) Expand  $D_{KL}(\hat{P} \parallel P_\theta)$

$$\begin{aligned} &= \sum_{x \in \mathcal{X}} \hat{P}(x) \log \frac{\hat{P}(x)}{P_\theta(x)} \\ &= \sum_{x \in \mathcal{X}} \hat{P}(x) \log \hat{P}(x) - \sum_{x \in \mathcal{X}} \hat{P}(x) \log P_\theta(x) \end{aligned}$$

Given  $\hat{P}(x)$  is a uniform distribution over  $x \in \mathcal{X}$ .

so the first half of the equation is a constant term.

$$\begin{aligned} &\underset{\theta}{\operatorname{argmin}} D_{KL}(\hat{P} \parallel P_\theta) \\ &= \underset{\theta}{\operatorname{argmin}} - \sum_{x \in \mathcal{X}} \hat{P}(x) \log P_\theta(x) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{x \in \mathcal{X}} \hat{P}(x) \log P_\theta(x) \quad \text{since } \hat{P}(x) \text{ is constant} \\ &\quad \text{for } \forall x \in \mathcal{X}. \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{x \in \mathcal{X}} \log P_\theta(x) \end{aligned}$$

Say the  $x \in \mathcal{X}$  is defined as  $\{x_i; i=1, \dots, n\}$ .

We thus show the equivalence.

Question : 3

part a :

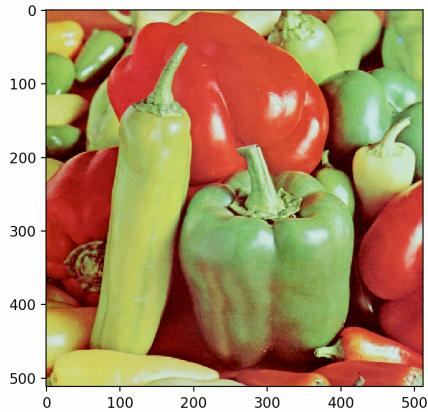


fig. large peppers (original)

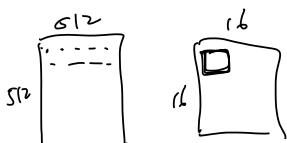


fig. large peppers (generated)

$3 \times 8$

part b :

originally, each pixel in the image is stored using 24-bit ( $\log_2(256)=8$ ;  $3 \times 8 = 24$ )

After compressing using only 16 colors, each pixel only takes 4( $\log_2(16)=4$ ) bits.

so the image compressed by a factor  $\frac{4}{24} = \frac{1}{6}$ .

## Question 4

Part a.

We want to prove the algorithm converges

$\Leftrightarrow$  proving  $\ell_{\text{semi-sup}}(\theta)$  is monotonic increasing

i.e. let  $\theta^{(t)}$  be the parameter at the end of  $t$  EM-step, we show  $\ell_{\text{semi-sup}}(\theta^{(t+1)}) \geq \ell_{\text{semi-sup}}(\theta^{(t)})$

$$\begin{aligned}
 \ell_{\text{semi-sup}} \theta^{(t+1)} &= \ell_{\text{unsup}} \theta^{(t+1)} + \alpha \ell_{\text{sup}} \theta^{(t+1)} \\
 &= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta^{(t+1)}) + \alpha \ell_{\text{sup}} \theta^{(t+1)} \\
 &= \sum_{i=1}^n \log p(x^{(i)}; \theta^{(t+1)}) + \alpha \ell_{\text{sup}} \theta^{(t+1)} \\
 &\geq \sum_{i=1}^n \text{ELBO}(x^{(i)}, \theta_i^{(t)}, \theta^{(t+1)}) + \alpha \ell_{\text{sup}} \theta^{(t+1)} \quad [\text{by Jensen's Inequality}] \\
 &\geq \sum_{i=1}^n \text{ELBO}(x^{(i)}, \theta_i^{(t)}, \theta^{(t+1)}) + \alpha \ell_{\text{sup}} \theta^{(t)} \quad [\text{The E-step}] \\
 &= \ell_{\text{unsup}}(\theta^{(t)}) + \alpha \ell_{\text{sup}}(\theta^{(t)}) \\
 &= \ell_{\text{semi-sup}}(\theta^{(t)})
 \end{aligned}$$

Since  $\theta^{(t)}$  is the end of the last EM step. it satisfies that

$$\ell_{\text{unsup}}(\theta^{(t)}) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, \theta_i^{(t)}, \theta^{(t+1)})$$

for  $\forall t$ , we show that  $\ell_{\text{semi-sup}}(\theta^{(t+1)}) \geq \ell_{\text{sup}}(\theta^{(t)})$

$\{\theta\}$  converges.

Part b:

We need to reestimate  $\bar{z}^{(i)}$  for  $i = 1, \dots, n$

$$\begin{aligned}
 \text{Calculate } w_j^{(i)} &= Q_i(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}, \phi, \mu, \Sigma) \\
 &= \frac{P(x^{(i)} | z^{(i)} = j; \mu, \Sigma) P(z^{(i)} = j; \phi)}{\sum_{l=1}^K P(x^{(i)} | z^{(i)} = l; \mu, \Sigma) P(z^{(i)} = l; \phi)} \\
 &= \frac{\frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)) \cdot \phi_j}{\sum_{l=1}^K \frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}} \exp(-\frac{1}{2}(x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l)) \cdot \phi_l} \quad (\text{according to note 22})
 \end{aligned}$$

Part c: Will reestimate  $\phi_j, \mu_j, \Sigma_j$  for  $j = 1, \dots, K$

Say  $Q_i(z^{(i)} = j)$  denotes the prob of  $z^{(i)}$  taking the value  $j$  under dist  $Q_i$ .

$$\text{Let } w_j^{(i)} = Q_i(z^{(i)} = j) \text{ and } \tilde{w}_j^{(i)} = \alpha \cdot \{z^{(i)} = j\}$$



The quantity we want to maximize is:

$$\begin{aligned}
 &\sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} + \sum_{i=1}^n \log (P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)) \\
 &= \sum_{i=1}^n \sum_{l=1}^K w_l^{(i)} \log \frac{\frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}} \exp(-\frac{1}{2}(x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l)) \cdot \phi_l}{w_l^{(i)}} + \sum_{i=1}^n \sum_{l=1}^K \tilde{w}_l^{(i)} \log \frac{\frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}} \exp(-\frac{1}{2}(x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l)) \cdot \phi_l}{\tilde{w}_l^{(i)}}
 \end{aligned}$$

① Maximize w.r.t.  $\mu$ . We attain the gradient:

$$\begin{aligned}
 &-\nabla_{\mu_l} \sum_{i=1}^n w_l^{(i)} \frac{1}{2} (x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l) + \sum_{i=1}^n \tilde{w}_l^{(i)} \frac{1}{2} (x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l) \\
 &= \frac{1}{2} \sum_{i=1}^n w_l^{(i)} \nabla_{\mu_l} 2\mu_l^T \Sigma_l^{-1} x^{(i)} - \mu_l^T \Sigma_l^{-1} \mu_l + \left[ \frac{1}{2} \sum_{i=1}^n w_l^{(i)} \nabla_{\mu_l} 2\mu_l^T \Sigma_l^{-1} x^{(i)} - \mu_l^T \Sigma_l^{-1} \mu_l \right] \\
 &= \sum_{i=1}^n w_l^{(i)} (\sum_l^{-1} x^{(i)} - \sum_l^{-1} \mu_l) + \sum_{i=1}^n \tilde{w}_l^{(i)} (\sum_l^{-1} x^{(i)} - \sum_l^{-1} \mu_l) \\
 &= \sum_l^{-1} \left( \sum_{i=1}^n w_l^{(i)} x^{(i)} + \sum_{i=1}^n \tilde{w}_l^{(i)} x^{(i)} \right) - \sum_l^{-1} \left( \sum_{i=1}^n w_l^{(i)} + \sum_{i=1}^n \tilde{w}_l^{(i)} \right) \mu_l = 0.
 \end{aligned}$$

$$80 \quad M_i = \frac{\sum_{j=1}^n w_j^{(i)} x^{(i)} + \sum_{j=1}^m \tilde{w}_j^{(i)} \tilde{x}^{(i)}}{\sum_{j=1}^n w_j^{(i)} + \sum_{j=1}^m \tilde{w}_j^{(i)}}$$

② Maximize w.r.t.  $\phi$

construct the lagrangian

$$\mathcal{L}(\phi) = \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \phi_j + \sum_{i=1}^m \sum_{j=1}^k \tilde{w}_j^{(i)} \log \phi_j + \beta (\sum_{j=1}^k \phi_j - 1)$$

$$\frac{\partial \mathcal{L}(\phi)}{\partial \phi_j} = \sum_{i=1}^n w_j^{(i)} \frac{1}{\phi_j} + \sum_{i=1}^m \tilde{w}_j^{(i)} \frac{1}{\phi_j} + \beta = 0.$$

$$\phi_j = \frac{\sum_{i=1}^n w_j^{(i)} + \sum_{i=1}^m \tilde{w}_j^{(i)}}{-\beta}$$

using the constraint  $\sum_j \phi_j = 1$ , we easily find that

$$-\beta = \sum_{j=1}^k \sum_{i=1}^n w_j^{(i)} + \sum_{j=1}^k \sum_{i=1}^m \tilde{w}_j^{(i)} = n + \alpha m$$

following the Lecture note, we have

$$\phi_j = \frac{\sum_{i=1}^n w_j^{(i)} + \sum_{i=1}^m \tilde{w}_j^{(i)}}{n + \alpha m}$$

③ Maximize w.r.t.  $\Sigma$ , the gradient is:

$$\begin{aligned} & \nabla_{\Sigma_j} \left( \sum_{i=1}^n w_j^{(i)} (\log |\Sigma_j| - (x^{(i)} - w_j)^T \Sigma_j^{-1} (x^{(i)} - w_j)) + \left( \sum_{i=1}^m \tilde{w}_j^{(i)} (\log |\Sigma_j| - (x^{(i)} - w_j)^T \Sigma_j^{-1} (x^{(i)} - w_j)) \right) \right) \\ &= (\sum_{i=1}^n w_j^{(i)} + \sum_{i=1}^m \tilde{w}_j^{(i)}) \Sigma_j - (\sum_{i=1}^n w_j^{(i)} (x^{(i)} - w_j) (x^{(i)} - w_j)^T) + \sum_{i=1}^m \tilde{w}_j^{(i)} ((x^{(i)} - w_j) (x^{(i)} - w_j)^T) = 0 \end{aligned}$$

$$\text{so, } \Sigma_j = \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - w_j) (x^{(i)} - w_j)^T + \sum_{i=1}^m \tilde{w}_j^{(i)} ((x^{(i)} - w_j) (x^{(i)} - w_j)^T)}{\sum_{i=1}^n w_j^{(i)} + \sum_{i=1}^m \tilde{w}_j^{(i)}}$$

part d:

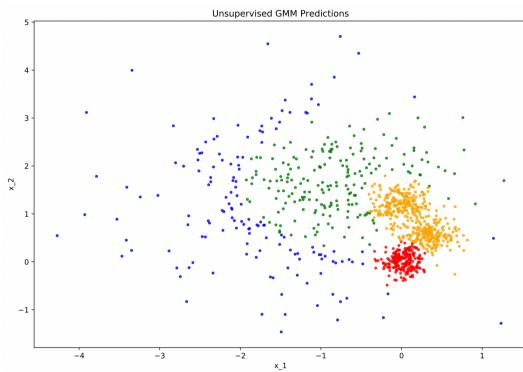


fig. unsupervised GMM

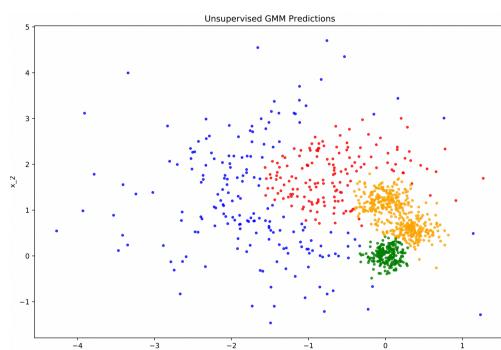


fig. unsupervised GMM

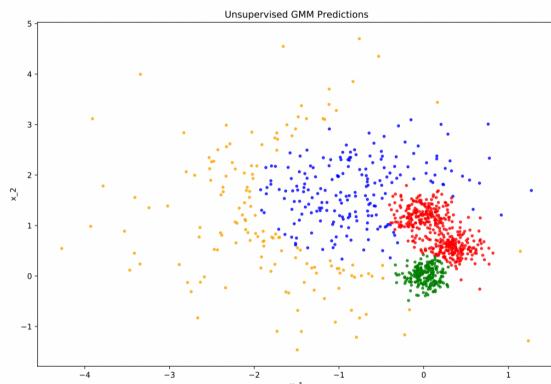


fig. unsupervised GMM

part e:

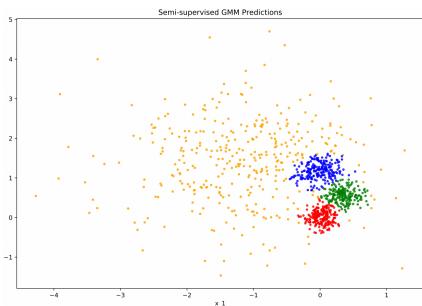
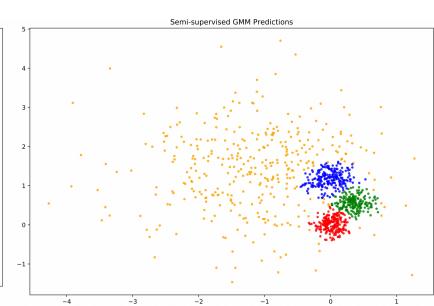
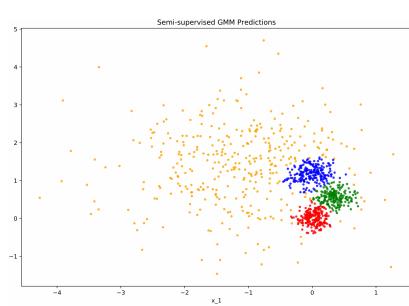


fig. semi-supervised GMM



f.

- i) The semi-supervised model has a faster convergence rate. The semi-supervised model takes at average 25 iterations to converge, whereas unsupervised model takes at average 157 iterations
- ii) with random initialization, unsupervised EM is less stable, which is obvious in 3 graphs attached at part d.  
semi-supervised EM is more stable.
- iii) Overall, semi-supervised EM outperforms unsupervised EM. Since semi-supervised EM has additional matching pairs, it provides more accurate identification of clusters under random initialization and when the distribution of the training data is not quite obvious.

## Question 5

Show first principle component:

$$\begin{aligned}
 & \underset{\|u\|=1}{\operatorname{arg \min}} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2 \\
 &= \underset{\|u\|=1}{\operatorname{arg \max}} \sum_{i=1}^n \|x^{(i)}\|^2 - \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2 \\
 &= \underset{\|u\|=1}{\operatorname{arg \max}} \sum_{i=1}^n \left( \|x^{(i)}\|^2 - \|x^{(i)} - f_u(x^{(i)})\|_2^2 \right) \\
 &= \underset{\|u\|=1}{\operatorname{arg \min}} \frac{1}{n} \sum_{i=1}^n \|f_u(x^{(i)})\|_2^2 \\
 &= \underset{\|u\|=1}{\operatorname{arg \min}} \frac{1}{n} \sum_{i=1}^n ((x^{(i)T} u)^2) \\
 &= \underset{\|u\|=1}{\operatorname{arg \min}} u^T \left( \frac{1}{n} \sum_{i=1}^n (x^{(i)T} x^{(i)}) \right) u
 \end{aligned}$$

Here we see that minimizing  $\sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2$  is equivalent to finding the max variance of projections onto the direction:  $u: \frac{1}{n} \sum_{i=1}^n \|f_u(x^{(i)})\|_2^2$ , which is what first principle component of PCA is about: Maximizing the variance.

## Question 6

part a:

Maximizing  $\ell(w)$ :

$$\begin{aligned}\ell(w) &= \sum_{i=1}^n \log P_x(x^{(i)}) \\ &= \sum_{i=1}^n \log P_s(wx^{(i)}|w) \\ &= \sum_{i=1}^n \log \left( \frac{1}{(2\pi)^{d/2}} \exp \left( -\frac{1}{2} (wx^{(i)})^\top (wx^{(i)}) \right) |w| \right) \\ &= \sum_{i=1}^n \left( -\frac{d}{2} \log(2\pi) - \frac{1}{2} x^{(i)\top} w^\top w x^{(i)} + \log|w| \right)\end{aligned}$$

Then we take the gradient w.r.t.  $w$

$$\begin{aligned}\nabla_w \ell(w) &= \sum_{i=1}^n \left( -\frac{1}{2} \nabla_w x^{(i)\top} w^\top w x^{(i)} + \nabla_w \log|w| \right) \\ &= \sum_{i=1}^n \left( -W x^{(i)} x^{(i)\top} + (W^{-1})^\top \right) \\ &= -W \left( \sum_{i=1}^n x^{(i)} x^{(i)\top} \right) + n (W^{-1})^\top \\ &= -W X^\top X + n (W^{-1})^\top = 0\end{aligned}$$

$$\text{So } W^\top W = \frac{1}{n} (X^\top X)^{-1}$$

Assuming RHS is invertible:

if there exists a valid solution  $w$ , then any  $w' = W \times U$  where  $U$  is an arbitrary orthogonal matrix,  $w'$  can be a valid solution too. This leads to an ambiguity, where  $U$  cannot be determined by data  $X$ .

So, ICA cannot recover the original.

Part b:

let  $x^{(i)}$  be an arbitrary example, the likelihood func is:

$$\begin{aligned} l_i(w) &= \sum_{j=1}^d \log p_s(w_j^T x^{(i)}) + \log |w| \\ &= \sum_{j=1}^d \log \left( \frac{1}{2} \exp(-|w_j^T x^{(i)}|) \right) + \log |w| \\ &= -d \log(2) - \sum_{j=1}^d |w_j^T x^{(i)}| + \log |w| \end{aligned}$$

Take the derivative:

$$\begin{aligned} \nabla_w l_i(w) &= -\sum_{j=1}^d \nabla_w |w_j^T x^{(i)}| + \nabla_w \log |w| \\ &= -\sum_{j=1}^d \text{sign}(w_j^T x^{(i)}) \begin{bmatrix} x^{(i)\top}_{\text{row}-j} \\ \vdots \\ 0 \end{bmatrix} + (w^{-1})^T \\ &= \begin{bmatrix} \text{sign}(w_1^T x^{(i)}) \\ \vdots \\ \text{sign}(w_d^T x^{(i)}) \end{bmatrix} x^{(i)\top} + (w^{-1})^T \end{aligned}$$

Take the SGA update rule to the above case, we get:

$$w := w + \alpha \left( - \begin{bmatrix} \text{sign}(w_1^T x^{(i)}) \\ \vdots \\ \text{sign}(w_d^T x^{(i)}) \end{bmatrix} x^{(i)\top} + (w^{-1})^T \right)$$

part c:

$\hat{w}$

[ 52.8352532 16.79619701 19.94171825 -10.19846303 -20.89757762 ]
[ -9.9292747 -0.97875614 -4.67786427 8.04377382 1.7865852 ]
[ 8.31096507 -7.47675728 19.31500349 15.17429591 -14.32612384 ]
[ -14.66742843 -26.64517989 2.44081559 21.38210464 -8.4207738 ]
[ -0.26929644 18.37414675 9.31198649 9.10287095 30.59463426 ] ]

List. The unmixing matrix with Laplace