# CS 229, Fall 2018
# Problem Set #3 Solutions:
## Deep Learning & Unsupervised learning

## 1. A Simple Neural Network

a) Gradient w.r.t $w_{1,2}^{[1]}$.

Denote the hidden layer output as $o$, the final output $\tilde{y}$. Forward propagation:

$$z^{[1]} = W^{[1]^T} x^{(i)}$$
$$o = g(z^{[1]})$$
$$z^{[2]} = W^{[2]^T} h$$
$$\tilde{y}^{(i)} = g(z^{[2]})$$

Cost function:

$$\ell = \frac{1}{m} \Sigma_{i=1}^m (\tilde{y}^{(i)} - y^{(i)})^2$$

Derivative of $W^{(2)}$ ($W^{(2)} \in R^3$, $h \in R^3$, $\tilde{y}$ and $y$ are scalars).

$$\frac{\partial \ell}{\partial W^{[2]}} = \frac{\partial}{\partial W^{[2]}} (\tilde{y}^{(i)} - y^{(i)})^2$$
$$= 2(\tilde{y}^{(i)} - y^{(i)}) \frac{\partial}{\partial W^{[2]}} (g(z^{[2]}))$$
$$= 2(\tilde{y}^{(i)} - y^{(i)}) g(z^{[2]}) (1 - g(z^{[2]})) o$$
$$= 2(\tilde{y}^{(i)} - y^{(i)}) \tilde{y}^{(i)} (1 - \tilde{y}^{(i)}) o$$

Derivative of $W^{[1]} \in R^{2 \times 3}$, $\frac{\partial \ell}{\partial W^{[1]}}$ has the same dimension.

$$\frac{\partial \ell}{\partial W^{[1]}} = \frac{\partial \ell}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial o} \frac{\partial o}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial W^{[1]}}$$
$$= 2(\tilde{y} - y) \tilde{y} (1 - \tilde{y}) W^{[2]} o (1 - o) \circ x$$

Here $\circ$ means outer product. Add subscription and sum up

$$\frac{\partial \ell}{\partial W_{1,2}^{[1]}} = \frac{1}{m} \sum_{i=1}^m 2(\tilde{y}^{(i)} - y^{(i)}) \tilde{y}^{(i)} (1 - \tilde{y}^{(i)}) W^{[2]} o^{(i)} (1 - o^{(i)}) \circ x^{(i)}$$

The update rule: $W_{1,2}^{[1]} = W_{1,2}^{[1]} - \alpha * \frac{\partial \ell}{\partial W_{1,2}^{[1]}}$

b) In the dataset plot, find three points clock wise which can form a triangle to separate the dataset in two:

$$a^{(1)} = (0.5, 0.5)$$
$$a^{(2)} = (0.5, 3.5)$$

$$a^{(3)} = (3.5, 0.5)$$

Subtract with the next to form three vectors, and add intercept in front:

$$v^{(1)} = (0, -3.0)$$
$$v^{(2)} = (-3.0, 3.0)$$
$$v^{(3)} = (3.0, 0)$$

Equations for the above three vectors. For $v^{(1)}$ and $v^{(2)}$:

$$\frac{x_1 - a_1^{(2)}}{v_1^{(1)}} = \frac{x_2 - a_2^{(2)}}{v_2^{(1)}}$$

Which can re rewritten as:

$$v_2^{(1)} x_1 - v_1^{(1)} x_2 + \left( v_1^{(1)} a_2^{(2)} - v_2^{(1)} a_1^{(2)} \right) = 0$$

This gives us the first component of matrix $W^{[1]}$:

$$W_1^{[1]} = \begin{pmatrix} v_2^{(1)} \\ -v_1^{(1)} \\ v_1^{(1)} a_2^{(2)} - v_2^{(1)} a_1^{(2)} \end{pmatrix} = \begin{pmatrix} -3.0 \\ 0 \\ 1.5 \end{pmatrix}$$

Similarly, for $v^{(2)}$ and $v^{(3)}$ we have:

$$W_2^{[1]} = \begin{pmatrix} v_2^{(2)} \\ -v_1^{(2)} \\ v_1^{(2)} a_2^{(3)} - v_2^{(2)} a_1^{(3)} \end{pmatrix} = \begin{pmatrix} 3.0 \\ 3.0 \\ -12.0 \end{pmatrix}$$

And for $v^{(3)}$ and $v^{(1)}$ we have:

$$W_3^{[1]} = \begin{pmatrix} v_2^{(3)} \\ -v_1^{(3)} \\ v_1^{(3)} a_2^{(1)} - v_2^{(3)} a_1^{(1)} \end{pmatrix} = \begin{pmatrix} 0 \\ -3.0 \\ 1.5 \end{pmatrix}$$

Stack them to from a matrix $W^{[1]}$:

$$W^{[1]} = \begin{pmatrix} -3.0 & 3.0 & 0 \\ 0 & 3.0 & -3.0 \\ 1.5 & -12.0 & 1.5 \end{pmatrix}$$

Explanation: given any point x, if its projection on the normal of vector $v^{(i)}$ is positive then it is on the right hand of $v^{(i)}$, otherwise on the left. If it lies on the right of all the vectors, it is inside the triangle. (note the intercept is the last component of W which can be easily adjusted to be the first component)

Output layer matrix can be:

$$W^{[2]} = \begin{pmatrix} -1 \\ -1 \\ -1 \\ 2.5 \end{pmatrix}$$

Because the input can only be one of the eight values of h:

$$h = \begin{pmatrix} 0/1 \\ 0/1 \\ 0/1 \\ 1 \end{pmatrix}$$

This guarantees that only for dataset inside the triangle could the product $W^{[2]}h$ be negative, which will then make f(x) = 0. All the other cases will have $W^{[2]}h \geq 0$, which makes f(x) = 1.

c) The activation function for $h_1, h_2, h_3$ is changed to $f(x) = x$, provide a set of weights that makes it achieve 100% accuracy.

Answer:
If the activation function is just f(x)=x, the NN is degraded to logistic regression. It is impossible to achieve 100% accuracy since the feature mappings doesn't include high order components. The decision boundary is going to be a straight line and will have misclassification.

# 2. KL divergence

## 2 KL divergence and Maximum Likelihood

## (a) Nonnegativity

Prove the following:

$$\forall P, Q, D_{KL}(P \parallel Q) \geq 0$$

And

$$D_{KL}(P\|Q) = 0 \iff P = Q$$

PROOF 1:

$$D_{KL}(P \parallel Q) = \Sigma_x P(x) log \frac{P(x)}{Q(x)}$$
$$= -\Sigma_x P(x) log \frac{Q(x)}{P(x)}$$
$$>= -log(\Sigma_x P(x) \frac{Q(x)}{P(x)})$$
$$= -log(\Sigma_x Q(x))$$
$$= -log(1) = 0$$

PROOF 2:

$$D_{KL}(P \parallel Q) = 0 \iff P = Q$$

a. If $P = Q$, $D_{KL}(P \parallel Q) = \sum_{x \in X} P \log \frac{P}{Q} = \sum_{x \in X} P \log(1) = 0$

b. If $D_{KL}(P \parallel Q) = 0$, given $-\log x$ is strictly convex, then

$$E[-log(\frac{P}{Q})] \geq -log(E[\frac{P}{Q}]) = 0$$
$$= log(E[\frac{Q}{P}])$$
$$= log(\sum[P\frac{Q}{P}])$$
$$= log(1)$$
$$= 0$$

The equality holds iff $\frac{Q}{P}$ is a constant with probability 1, given the fact that both $P$ and $Q$ are pdf, we can only have $P = Q$.

## (b) Chain rule for KL divergence

Prove that:

$$D_{KL}(P(X,Y) \parallel Q(X,Y)) = D_{KL}(P(X) \parallel Q(X)) + D_{KL}(P(Y|X) \parallel Q(Y|X))$$

PROOF:

$$D_{KL}(P(X) \parallel Q(X)) + D_{KL}(P(Y|X) \parallel Q(Y|X)) = \sum_x P(x) \log \frac{P(x)}{Q(x)} + \sum_x P(x)(\sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)})$$
$$= \sum_x P(x)\left( \log \frac{P(x)}{Q(x)} + \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right)$$
$$= \sum_x P(x)\left( \sum_y P(y|x) \log \frac{P(x)}{Q(x)} + \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right)$$
$$= \sum_x P(x)\left( \sum_y P(y|x)\left( \log \frac{P(x)}{Q(x)} + \log \frac{P(y|x)}{Q(y|x)} \right) \right)$$
$$= \sum_x \left( \sum_y P(y|x)P(x)(\log \frac{P(x)P(y|x)}{Q(x)Q(y|x}) \right)$$
$$= \sum_x \sum_y P(x,y) \log \frac{P(x,y)}{Q(x,y)}$$
$$= D_{KL}(P(X,Y) \parallel Q(X,Y))$$

## (c) KL and maximum likelihood

Prove that

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

$$
\begin{aligned}
D_{KL}(\hat{P} \parallel P_\theta) &= \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P_\theta(x)} \\
&= -\sum_x \hat{P}(x) \log \frac{P_\theta(x)}{\hat{P}(x)} \\
&= -\sum_x \frac{1}{m} \sum_{i=1}^{m} 1\{x^{(i)} = x\} \log \frac{P_\theta(x)}{\frac{1}{m}\sum_{i=1}^{m} 1\{x^{(i)} = x\}} \\
&= -\frac{1}{m} \sum_{i=1}^{m} \log \frac{P_\theta(x^{(i)})}{\frac{1}{m}\sum_{i=1}^{m}} \\
&= -\frac{1}{m} \sum_{i=1}^{m} \log P_\theta(x^{(i)}) \\
&= -\frac{1}{m} \text{log-likelihood}
\end{aligned}
$$

Which implies that

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

# 3. KL divergence, Fisher Information, Natural gradient

**(For simplicity, $\nabla_{\theta'}\log p(y;\theta')|\theta' = \theta$ is simply written as $\nabla_\theta \log p(y;\theta)$, $E_{y\sim p(y;\theta)}$ is sometimes simply written as $E_\theta$ as it is in some Statistics textbooks.)**

PS3.3 [For simplicity, $\nabla_{\theta'}\log p(y;\theta')|\theta'=\theta$ is sometimes written as $\nabla_\theta \log p(y;\theta)$ directly]

(a) Suppose the score function is
$$s(y;\theta) = \nabla_\theta \log p(y;\theta);$$
prove that:
$$E_{y\sim p(y;\theta)}[s(y;\theta')|\theta'=\theta'] = 0$$

PROOF:
$$E_{y\sim p(y;\theta)}[s(y;\theta')|\theta'=\theta] = \int p(y;\theta)s(y;\theta)\,dy$$
$$= \int p(y;\theta)\nabla_\theta \log p(y;\theta)\,dy = \int p(y;\theta)\cdot\frac{1}{p(y;\theta)}\cdot\frac{\partial p(y;\theta)}{\partial\theta}\,dy$$
$$= \frac{\partial}{\partial\theta}\int p(y;\theta)\,dy = \frac{\partial}{\partial\theta}(1) = 0. \quad \text{Done.}$$

(b) Fisher Information:
$$I(\theta) = \text{Cov}_{y\sim p(y;\theta)}[\nabla_{\theta'}\log p(y;\theta')|\theta'=\theta]$$
Show that
$$I(\theta) = E_{y\sim p(y;\theta)}[\nabla_\theta\log p(y;\theta')\nabla_\theta\log p(y;\theta')^T|\theta'=\theta]$$

PROOF: By definition of Covariance, Given a vector $X = \nabla_\theta \log p(y;\theta)$
$$I(\theta) = \text{Cov}_{y\sim p(y;\theta)}[\nabla_{\theta'}\log p(y;\theta')|\theta'=\theta]$$
$$= E_{y\sim p(y;\theta)}[\nabla_{\theta'}\log p(y;\theta')\nabla_\theta\log p(y;\theta')^T|\theta'=\theta]$$
$$- E_{y\sim p(y;\theta)}[\nabla_\theta\log p(y;\theta')]\cdot E_{y\sim p(y;\theta)}[\nabla_\theta\log p(y;\theta')|\theta'=\theta]$$
$$= E_{y\sim p(y;\theta)}[\nabla_\theta\log p(y;\theta')\nabla_\theta\log p(y;\theta')^T|\theta'=\theta]$$

The last part is from the fact that $E_{y\sim p(y;\theta)}[s(y;\theta)]=0$

$$\text{Cov}(X) = \begin{bmatrix} E(x_1^2)-E^2(x_1) & E(x_1x_2)-E(x_1)E(x_2) & \cdots & E(x_1x_n)-E(x_1)E(x_n) \\ E(x_2x_1)-E(x_2)E(x_1) & E(x_2^2)-E^2(x_2) & \cdots & E(x_2x_n)-E(x_2)E(x_n) \\ E(x_3x_1)-E(x_3)E(x_1) & & & \\ \cdots & & & \cdots \\ E(x_nx_1)-E(x_n)E(x_1) & E(x_nx_2)-E(x_n)E(x_2) & \cdots & E(x_nx_n)-E^2(x_n) \end{bmatrix}$$

From (a) we know that $E(x_i) = 0$; So :

$$Cov(X) = \begin{bmatrix} E(x_1^2) & E(x_1x_2) & \cdots & E(x_1x_n) \\ E(x_2x_1) & E(x_2^2) & \cdots & E(x_2x_n) \\ \cdots & & \ddots & \\ E(x_nx_1) & & & E(x_n^2) \end{bmatrix} = E \begin{bmatrix} x_1^2 & x_1x_2 & \cdots & x_1x_n \\ x_2x_1 & x_2^2 & \cdots & x_2x_n \\ \cdots & & & \cdots \\ x_nx_1 & & & x_n^2 \end{bmatrix}$$

$$= E[XX^T] \quad \text{replace } X = \nabla_{\theta'} \log P(y; \theta')\big|_{\theta'=\theta}, \quad \text{Proved.}$$

(c). Show that:

$$I(\theta) = E_{y \sim p(y;\theta)}\left[ -\nabla_{\theta'}^2 \log p(y;\theta')\big|_{\theta'=\theta}\right]$$

PROOF:

$$\nabla_\theta^2 \log p(y;\theta) = \nabla_\theta (\nabla_\theta \log p(y;\theta))$$
$$= \nabla_\theta (\nabla_\theta p(y;\theta)/p(y;\theta))$$
$$= \nabla_\theta^2 p(y;\theta)/p(y;\theta) - (\nabla_\theta p(y;\theta)/p(y;\theta))(\nabla_\theta p(y;\theta)/p(y;\theta))^T$$
$$= \nabla_\theta^2 p(y;\theta)/p(y;\theta) - (\nabla_\theta \log p(y;\theta))(\nabla_\theta \log p(y;\theta))^T$$

Sum over $y$ on both sides:

$$E_{y \sim p(y;\theta)}\left[\nabla_\theta^2 \log p(y;\theta)\right] = E_\theta\left[\nabla_\theta^2 p(y;\theta)/p(y;\theta)\right]$$
$$- E_\theta\left[\nabla_\theta \log p(y;\theta)\nabla_\theta \log p(y;\theta)^T\right]$$

Consider the first part of the right hand side:

$$E_{y \sim p(y;\theta)}\left[\nabla_\theta^2 p(y;\theta)/p(y;\theta)\right] = \int p(y;\theta) \cdot \frac{\nabla_\theta^2 p(y;\theta)}{p(y;\theta)} dy$$
$$= \int \nabla_\theta^2 p(y;\theta) dy = \nabla_\theta^2 \int p(y;\theta) dy = \nabla_\theta^2 (1) = 0;$$

The second part of the right hand side:

$$E_\theta\left[\nabla_\theta \log p(y;\theta) \nabla_\theta \log p(y;\theta)^T\right] = I(\theta)$$

So:

$$E_\theta\left[\nabla_\theta^2 \log p(y;\theta)\right] = -I(\theta)$$
$$I(\theta) = E_\theta\left[-\nabla_\theta^2 \log p(y;\theta)\right] = E_{y \sim p(y;\theta)}\left[-\nabla_{\theta'}^2 \log p(y;\theta')\big|_{\theta'=\theta}\right]$$

Done!

(d) $D_{KL}(P_\theta \| P_{\theta+d}) \approx \frac{1}{2} d^T I(\theta) d$

PROOF:

$D_{KL}(P_\theta \| P_{\theta+d}) = E_\theta[\log P_\theta - \log P_{\theta+d}] = E_\theta[\log P_\theta] - E_\theta[\log P_{\theta+d}]$

Consider $\log P_{\theta+d}$ using Taylor's expansion:

$\log P_{\theta+d} \approx \log P_\theta + d \nabla_\theta \log P_\theta + \frac{1}{2} d^T (\nabla_\theta^2 \log P_\theta) d$

Sum over $y \sim p(y;\theta)$ for both sides:

$E_\theta[\log P_{\theta+d}] \approx E_\theta[\log P_\theta] + E_\theta[d \nabla_\theta \log P_\theta] + E_\theta[\frac{1}{2} d^T (\nabla_\theta^2 \log P_\theta) d]$

$= E_\theta[\log P_\theta] + d E_\theta[\nabla_\theta \log P_\theta] + \frac{1}{2} d^T E_\theta[\nabla_\theta^2 \log P_\theta] d$

$= E_\theta[\log P_\theta] - \frac{1}{2} d^T I(\theta) d$

Which implies:

$D_{KL}(P_\theta \| P_{\theta+d}) \approx \frac{1}{2} d^T I(\theta) d.$    done!

(e).    $l(\theta) = \log p(y;\theta)$ ;    $\theta \in R^n$ ;

From (d) above we have:    $D_{KL}(P_\theta \| P_{\theta+d}) \approx \frac{1}{2} d^T I(\theta) d$ ;

The Taylor expansion for $l(\theta+d) \approx l(\theta) + d^T \nabla_\theta l(\theta)$ ;

the Original problem:

$d^* = \arg\max_d l(\theta+d)$ , s.t. $D_{KL}(P_\theta \| P_{\theta+d}) = c$ ,

Replace both Taylor expansions:

$d^* = \arg\max_d l(\theta) + d^T \nabla_\theta l(\theta)$ , s.t. $D_{KL}(P_\theta \| P_{\theta+d}) = \frac{1}{2} d^T I(\theta) d = c$

→ the Lagrangian is:

$L(d,\lambda) = l(\theta) + d^T \nabla_\theta l(\theta) - \lambda \cdot (\frac{1}{2} d^T I(\theta) d - c)$

→ Derivatives of linear equations:

$\frac{\partial L}{\partial d} = \nabla_\theta l(\theta) - \lambda I(\theta) d = 0$ ;   ⓐ

$\frac{\partial L}{\partial \lambda} = c - \frac{1}{2} d^T I(\theta) d = 0$   ;   ⓑ

→ From ⓐ :   $d = \lambda^{-1} I(\theta)^{-1} \nabla_\theta l(\theta)$ ;

→ Plug $d$ to ⓑ:

$c = \frac{1}{2} (\lambda^{-1} I(\theta)^{-1} \nabla_\theta l(\theta))^T I(\theta) (\lambda^{-1} I(\theta)^{-1} \nabla_\theta l(\theta))$

$= \frac{1}{2\lambda^2} \nabla_\theta l(\theta)^T I(\theta)^{-T} I(\theta) I(\theta)^{-1} \nabla_\theta l(\theta)$

$= \frac{1}{2\lambda^2} \nabla_\theta l(\theta)^T I(\theta)^{-1} \nabla_\theta l(\theta)$

$\Rightarrow \lambda = \frac{1}{\sqrt{2c}} \sqrt{\nabla_\theta l(\theta)^T I(\theta)^{-1} \nabla_\theta l(\theta)}$

Remove the '-T' part: use
(AB)^T = B^T A^T and the fact
that the result is a scalar

⊙ → plug $\lambda$ above back to (a), get $d^*$:

$$d^* = \sqrt{c} \cdot \left(\nabla_\theta l(\theta)^T I(\theta)^{-1} \nabla_\theta l(\theta)\right)^{-\frac{1}{2}} I(\theta)^{-1} \nabla_\theta l(\theta)$$

the dimension of $d^*$ is easy to know: $d^* \in \mathbb{R}^n$;

(f). Relation to Newton's method.

Compare the two updates (compare the $\Delta$s for both)

natural gradient : $d = \lambda^{-1} \cdot I(\theta)^{-1} \nabla_\theta l(\theta)$;

Newton's method : $\Delta \hat{\theta} = H^{-1} \nabla_\theta l(\theta)$

the Hessian $H$ is just given by:

$$H = \nabla_\theta^2 l(\theta);$$

~~$= Var(y; \eta) xx^T$, this is already proved in PS1.4.(c);
Also in PS1.4.(c) we know:
$\nabla_\theta l(\theta)$~~

in ~~PS1.4.(c)~~ PS1.4.(c) we already proved:

$\nabla_\theta^2 l(\theta) = Var(y; \eta) xx^T$ this is with natural variable $\eta = \theta^T x$; if ~~both~~ we simple use $\eta = \theta$, we have:

$\nabla_\theta^2 l(\theta) = Var(y; \theta)$ which is just the same as $I(\theta)$;

$\lambda^{-1}$ is simply a scalar. the direction is the same for both, because:

$I(\theta)^{-1} = H^{-1}$;

# 4. Semi-supervised EM

4. Semi-supervised EM

(a): PROOF.

For ease of hand writing, we $l(\theta)$ as $l_{semi-sup}(\theta)$.

Jensen's inequality tells us:

$$l(\theta) = \sum_{i=1}^{m} \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) + \alpha\left(\sum_{i=1}^{\tilde{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta)\right)$$

$$= \sum_{i=1}^{m} \sum_{z^{(i)}} \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} + \alpha\left(\sum_{i=1}^{\tilde{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta)\right)$$

$$\geq \sum_{i=1}^{m} \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} + \alpha\left(\sum_{i=1}^{\tilde{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta)\right)$$

it is already shown in lecture note that the tight equality only happenes when:

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta), \quad \text{suppose or}$$

suppose $\theta^{(t)}$ is the current tight point, it means:

$$l(\theta^{(t)}) = \sum_{i=1}^{m} \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} + \alpha\left(\sum_{i=1}^{\tilde{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t)})\right)$$

For parameter $\theta^{(t+1)}$ will have

$$l(\theta^{(t+1)}) \geq \sum_{i=1}^{m} \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} + \alpha\left(\sum_{i=1}^{\tilde{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t+1)})\right)$$

because the right part of the inequality is a lower bound of $l(\theta)$;

Consider the right part again, because $\theta^{(t+1)}$ is an arg max $(g)$ $(g$ here denotes the above right part). we again have:

$$g \geq \sum_{i=1}^{m} \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} + \alpha\left(\sum_{i=1}^{\tilde{m}} \log P(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta^{(t)})\right)$$

$$= l(\theta^{(t)})$$

So, $l(\theta^{(t+1)}) \geq l(\theta^{(t)})$. Done!

## 3. Semi-supervised GMM

(b). $z^{(i)}$'s are the latent randoms which can take value $\{1, 2, \dots k\}$.

The updated values are:

$$w_j^{(i)} = Q_i(z^{(i)} = j) = p(z^{(i)} = j \mid x^{(i)}; \phi, \mu_j, \mu_\ell, \Sigma_j, \Sigma_\ell)$$

$$= \frac{p(x^{(i)} \mid z^{(i)} = j; \mu_j, \Sigma_j) \, p(z^{(i)} = j; \phi)}{\sum_{\ell=1}^{R} p(x^{(i)} \mid z^{(i)} = \ell; \mu_\ell, \Sigma_\ell) \, p(z^{(i)} = \ell; \phi)}$$

$$= \frac{\frac{1}{\sqrt{2\pi} |\Sigma_j|} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1}(x^{(i)} - \mu_j)\right) \cdot \phi_j}{\sum_{\ell=1}^{R} \frac{1}{\sqrt{2\pi} |\Sigma_\ell|} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_\ell)^T \Sigma_\ell^{-1}(x^{(i)} - \mu_\ell)\right) \phi_\ell}$$

$$= \frac{\frac{1}{\sqrt{2\pi} |\Sigma_j|} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1}(x^{(i)} - \mu_j)\right) \phi_j}{\sum_{\ell=1}^{R} \frac{1}{\sqrt{2\pi} |\Sigma_\ell|} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_\ell)^T \Sigma_\ell^{-1}(x^{(i)} - \mu_\ell)\right) \phi_\ell}$$

( The parameters are : $\mu_1, \mu_2, \dots \mu_R, \Sigma_1, \Sigma_2, \dots \Sigma_R, \phi_1, \phi_2, \dots \phi_R$ )

> Corrections: in (b)
> $\sqrt{2\pi}$ should be $(2\pi)^{\frac{n}{2}}$
> $|\Sigma_\ell|$, $|\Sigma_j|$ should be $|\Sigma_\ell|^{\frac{1}{2}}$, $|\Sigma_j|^{\frac{1}{2}}$

(c) in M-step, just plug $w_j^{(i)}$ to $\ell(\theta)$ and take derivatives w.r.t $\mu$, $\Sigma$ and $\phi$. The first part is the same as it is in lecture notes. and we also need to plus the supervised part:

$$\ell(\theta) = \sum_{i=1}^{m} \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} + \alpha \left( \sum_{i=1}^{\tilde{m}} \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \tilde{\mu}, \tilde{\phi}, \tilde{\Sigma}) \right)$$

The above parameters are without subscripts, expand and add subscript:

$$\ell(\theta) = \sum_{i=1}^{m} \sum_{j=1}^{R} w_j^{(i)} \log \left( \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1}(x^{(i)} - \mu_j)\right) \phi_j \Big/ w_j^{(i)} \right) +$$

$$\alpha \left( \sum_{i=1}^{\tilde{m}} \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \phi_{\tilde{z}^{(i)}}, \mu_{\tilde{z}^{(i)}}, \Sigma_{\tilde{z}^{(i)}}) \right)$$

$$= \sum_{i=1}^{m} \sum_{j=1}^{R} w_j^{(i)} \log \left( \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1}(x^{(i)} - \mu_j)\right) \phi_j \Big/ w_j^{(i)} \right) +$$

$$\alpha \left( \sum_{i=1}^{\tilde{m}} \log \left( \frac{1}{(2\pi)^{n/2} |\Sigma_{\tilde{z}^{(i)}}|^{1/2}} \exp\left(-\frac{1}{2}(\tilde{x}^{(i)} - \mu_{\tilde{z}^{(i)}})^T \Sigma_{\tilde{z}^{(i)}}^{-1}(\tilde{x}^{(i)} - \mu_{\tilde{z}^{(i)}})\right) \phi_{\tilde{z}^{(i)}} \right) \right)$$

take the derivative w.r.t $\mu_\ell$, the first part is just the same as lecture notes. and the second part can be written out similar to the first part:

$$\nabla_{\mu_\ell} = \sum_{i=1}^{m} w_\ell^{(i)} \left( \Sigma_\ell^{-1} x^{(i)} - \Sigma_\ell^{-1} \mu_\ell \right) + \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} \left( \Sigma_{\tilde{z}^{(i)}}^{-1} \tilde{x}^{(i)} - \Sigma_{\tilde{z}^{(i)}}^{-1} \mu_{\tilde{z}^{(i)}} \right)$$

Set it to $\vec{0}$ and left multiply with $\Sigma_\ell$ give us:

> Correction: should have an '$\alpha$' in front of the 2nd term

$$\sum_{i=1}^{m} w_\ell^{(i)}(x^{(i)} - \mu_\ell) + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} (\tilde{x}^{(i)} - \mu_\ell) = \vec{0}$$

$$\mu_\ell = \frac{\sum_{i=1}^{m} w_\ell^{(i)} x^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} \tilde{x}^{(i)}}{\sum_{i=1}^{m} w_\ell^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\}}$$

~~take derivative w.r.t $\phi_j$:~~

Now derive the update rule for $\phi_j$. Group all terms that depend on $\phi_j$, we have the following:

$$\sum_{i=1}^{m}\sum_{j=1}^{k} W_j^{(i)} \log \phi_j + \text{\rule{2cm}{0.4cm}} \alpha \sum_{i=1}^{\tilde{m}} \log \phi_{\tilde{z}^{(i)}}$$

Similar to lecture notes, the Lagrangian is:

$$L(\phi) = \sum_{i=1}^{m}\sum_{j=1}^{k} W_j^{(i)} \log \phi_j + \alpha \sum_{i=1}^{\tilde{m}} \log \phi_{\tilde{z}^{(i)}} + \beta\left(\sum_{j=1}^{k}\phi_j - 1\right)$$

take derivative w.r.t $\phi_j$:

$$\frac{\partial}{\partial \phi_j} = \sum_{i=1}^{m} \frac{W_j^{(i)}}{\phi_j} + \alpha \sum_{i=1}^{\tilde{m}} \frac{1\{\tilde{z}^{(i)}=j\}}{\phi_j} + \beta = 0 \implies$$

$$\phi_j = \frac{\sum_{i=1}^{m} W_j^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} 1(\tilde{z}^{(i)}=j)}{-\beta}$$

We the fact that $\sum_{j=1}^{k}\phi_j = 1$, its easy to come to:

$$-\beta = \sum_{i=1}^{m}\sum_{j=1}^{k} W_j^{(i)} + \alpha \sum_{i=1}^{m} 1\{\tilde{z}^{(i)} \text{ in } \{1,2,\dots k\}\}$$

$$= m + \alpha\tilde{m},$$

therefore:

$$\phi_j = \frac{\sum_{i=1}^{m} W_j^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)}=j\}}{m + \alpha\tilde{m}}$$

Finally lets do $\Sigma_\ell$. Go to next page. →

$\rightarrow$

terms related to $\Sigma_j$s:

$$l(\theta) = \sum_{i=1}^{m} \sum_{j=1}^{R} W_j^{(i)} \left( -\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) - \log \left( (2\pi)^{N/2} |\Sigma_j|^{1/2} \right) \right)$$

$$+ \alpha \left( \sum_{i=1}^{\tilde{m}} \left( -\frac{1}{2}(\tilde{x}^{(i)} - \mu_{\tilde{z}^{(i)}})^T \Sigma_{\tilde{z}^{(i)}}^{-1} (\tilde{x}^{(i)} - \mu_{\tilde{z}^{(i)}}) - \log \left( (2\pi)^{N/2} |\Sigma_{\tilde{z}^{(i)}}|^{1/2} \right) \right) \right)$$

$$= -\sum_{i=1}^{m} \sum_{j=1}^{R} W_j^{(i)} \left( \frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1}(x^{(i)} - \mu_j) + \frac{1}{2}\log(|\Sigma_j|) \right)$$

$$- \alpha \left( \sum_{i=1}^{\tilde{m}} \left( \frac{1}{2}(\tilde{x}^{(i)} - \mu_{\tilde{z}^{(i)}})^T \Sigma_{\tilde{z}^{(i)}}^{-1}(\tilde{x}^{(i)} - \mu_{\tilde{z}^{(i)}}) + \frac{1}{2}\log(|\Sigma_{\tilde{z}^{(i)}}|) \right) \right)$$

take derivative w.r.t $\Sigma_\ell$; it's still an $n \times n$ matrix:

$$\nabla_{\Sigma_\ell} = \frac{1}{2} \sum_{i=1}^{m} \left( (x^{(i)} - \mu_\ell)(x^{(i)} - \mu_\ell)^T \Sigma_\ell^{-2} - \Sigma_\ell^{-1} \right) W_\ell^{(i)}$$

$$+ \frac{1}{2}\alpha \left( \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} (\tilde{x}^{(i)} - \mu_\ell)(\tilde{x}^{(i)} - \mu_\ell)^T \Sigma_\ell^{-2} - \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} \Sigma_\ell^{-1} \right)$$

$$= \frac{1}{2} \left( \sum_{i=1}^{m} W_\ell^{(i)}(x^{(i)} - \mu_\ell)(x^{(i)} - \mu_\ell)^T \Sigma_\ell^{-2} - \sum_{i=1}^{m} W_\ell^{(i)} \Sigma_\ell^{-1} \right)$$

$$+ \frac{1}{2}\alpha \left( \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\}(\tilde{x}^{(i)} - \mu_\ell)(\tilde{x}^{(i)} - \mu_\ell)^T \Sigma_\ell^{-2} - \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} \Sigma_\ell^{-1} \right)$$

Set it to $0$ and right multiply with $\Sigma_\ell^2$:

$$\sum_{i=1}^{m} W_\ell^{(i)}(x^{(i)} - \mu_\ell)(x^{(i)} - \mu_\ell)^T + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\}(\tilde{x}^{(i)} - \mu_\ell)(\tilde{x}^{(i)} - \mu_\ell)^T$$

$$= \left( \sum_{i=1}^{m} W_\ell^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\} \right) \Sigma_\ell$$

$$\Rightarrow \Sigma_\ell = \frac{\sum_{i=1}^{m} W_\ell^{(i)}(x^{(i)} - \mu_\ell)(x^{(i)} - \mu_\ell)^T + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\}(\tilde{x}^{(i)} - \mu_\ell)(\tilde{x}^{(i)} - \mu_\ell)^T}{\sum_{i=1}^{m} W_\ell^{(i)} + \alpha \sum_{i=1}^{\tilde{m}} 1\{\tilde{z}^{(i)} = \ell\}}$$
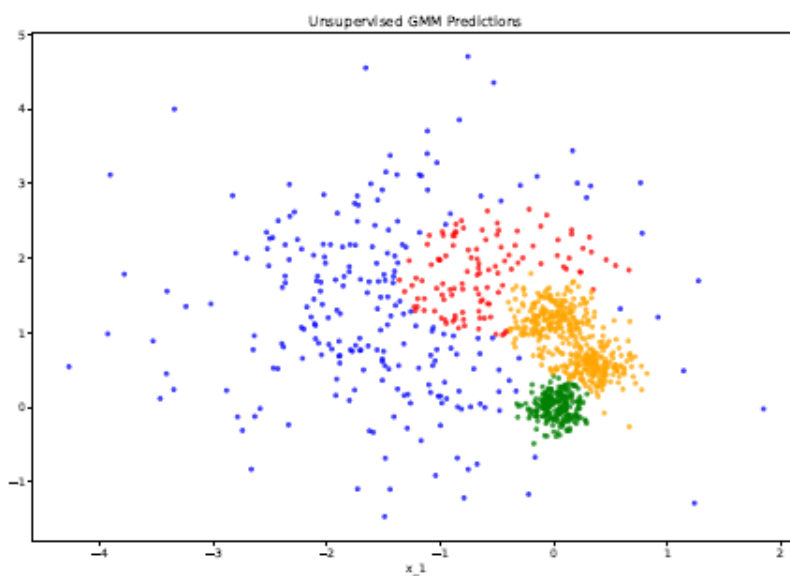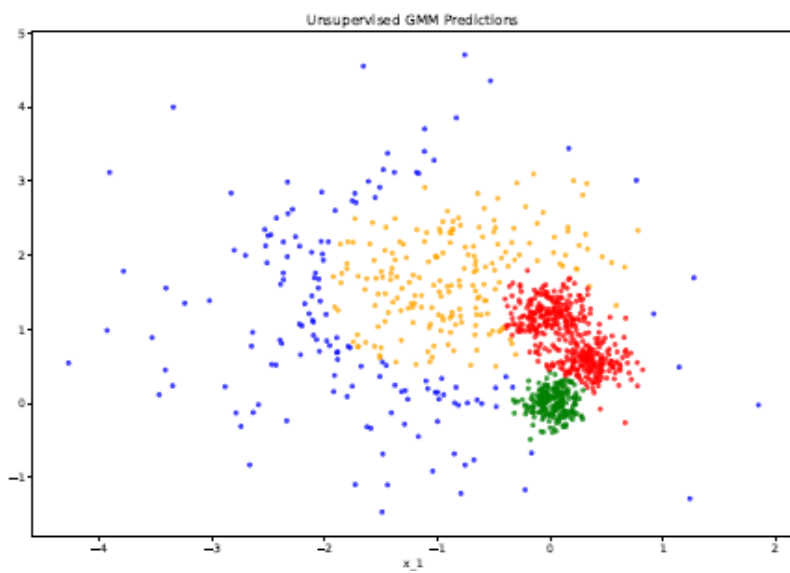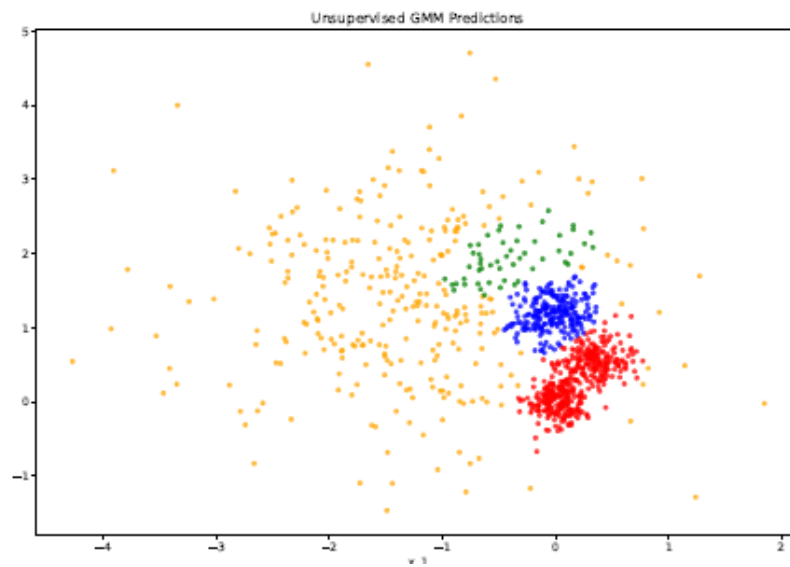
the denominator is a scalar, the numerator is an $n \times n$ matrix. Done!

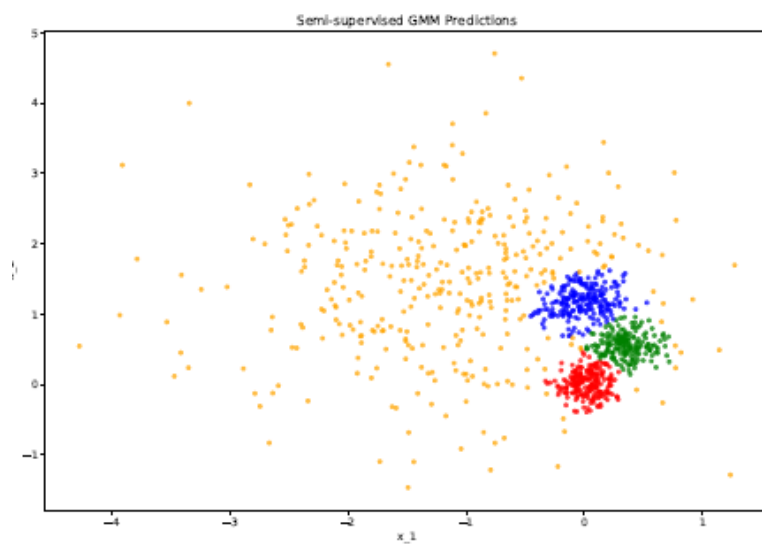(d) plot.
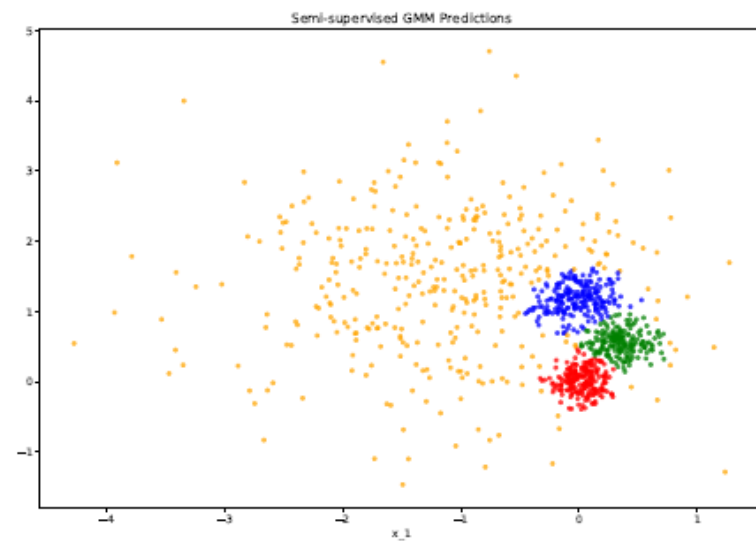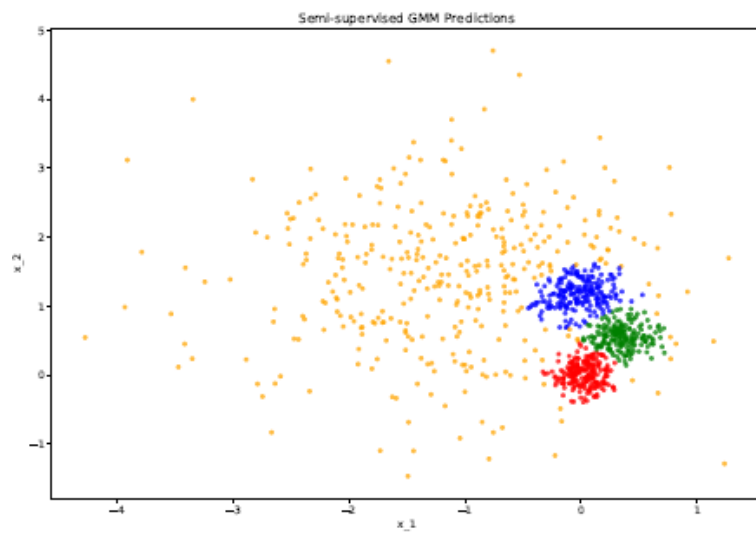
(e) plot.

d) Unsupervised GMM predictions, , run three times, not stable



Unsupervised GMM Predictions



Unsupervised GMM Predictions



Unsupervised GMM Predictions

e) Semi-supervised GMM predictions, run three times, stable

f) Comparison of Unsupervised and Semi-supervised EM

i) Number of iterations taken to converge
Answer:
Unsupervised EM takes 162,118,117 iterations to converge. Three runs are shown below:
Unsupervised EM: converged in 162 iterations, log-likelihood=-1831.349043
Unsupervised EM: converged in 118 iterations, log-likelihood=-1835.448461
Unsupervised EM: converged in 117 iterations, log-likelihood=-1835.444311

Semi-supervised EM only takes 15,31,31 iterations to converge.
Semi-supervised EM: converged in **15 iterations**, **log-likelihood=-2344.767485**, ll_sup=-552.478311, ll_unsup=-1792.289174
Semi-supervised EM: converged in **31 iterations**, **log-likelihood=-2344.614540**, ll_sup=-553.009998, ll_unsup=-1791.604542
Semi-supervised EM: converged in **31 iterations**, **log-likelihood=-2344.613954**, ll_sup=-553.011754, ll_unsup=-1791.602200

ii) Stability
Answer:
As can be seen in the above plot, Unsupervised EM is very sensitive to initialization and the class assignment changes dramatically with different random initializations.
Semi-supervised EM is very stable and class assignments do not change much.

iii) Overall quality of assignments
Answer:
Overall the Semi-supervised EM has a much better quality compared to Unsupervised EM.

# 5. K-means for compression.

a) K-means compression implementation.

```python
from matplotlib.image import imread
import matplotlib.pyplot as plt
import numpy as np


# reduce to n colors
def kmeans_compress(X, n = 16):
    X = X.reshape(-1, 3)
    centroids = X[np.random.choice(np.arange(X.shape[0]), size=n, replace=False)]
    centers = np.zeros([n, 3])
    #
    while not np.array_equal(centers, centroids):
        centers = centroids
        norms = [np.sqrt(np.sum((centers[i] - X) ** 2, axis=1)) for
            i in range(0, len(centers))]
        idxs = np.stack(norms).argmin(axis=0)
        #recalc centroids
        centroids = np.array([np.mean(X[np.where(idxs == i,True,False)],
            axis = 0).round() for i in range(n)])
    #you need the data type conversion, otherwise you are in big trouble!!
    return np.uint8(centroids), np.uint8(idxs)

def show_img(X, ax, s):
    ax.imshow(X)
    ax.set_title(s)

def main():
    fig, axes = plt.subplots(2, 1, figsize=(12, 4))
    axb1, axb2 = axes.ravel()

    A1 = imread('..\data\peppers-small.tiff')
    centroids,idxs = kmeans_compress(A1)
    B1 = np.array([centroids[idxs[i]] for i in range(0, len(idxs))])
    B1 = B1.reshape(A1.shape)
    show_img(B1, axb1,'compressed peppers-small')

    A2 = imread('..\data\peppers-large.tiff')
    centroids,idxs = kmeans_compress(A2)
    B2 = np.array([centroids[idxs[i]] for i in range(0, len(idxs))])
    B2 = B2.reshape(A2.shape)
    show_img(B2, axb2, 'compressed peppers-small')

    plt.show()

if __name__ == '__main__':
    main()
```
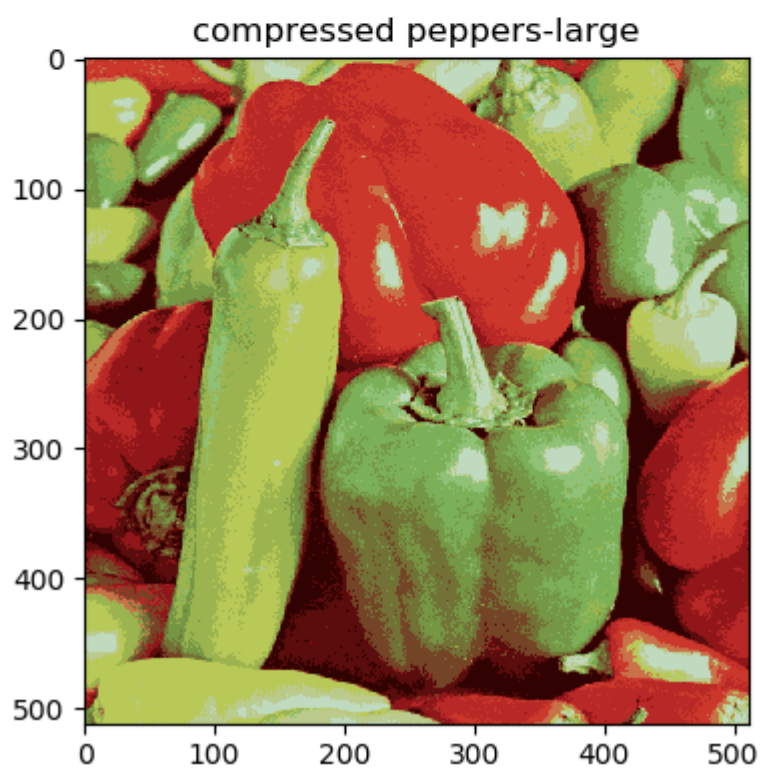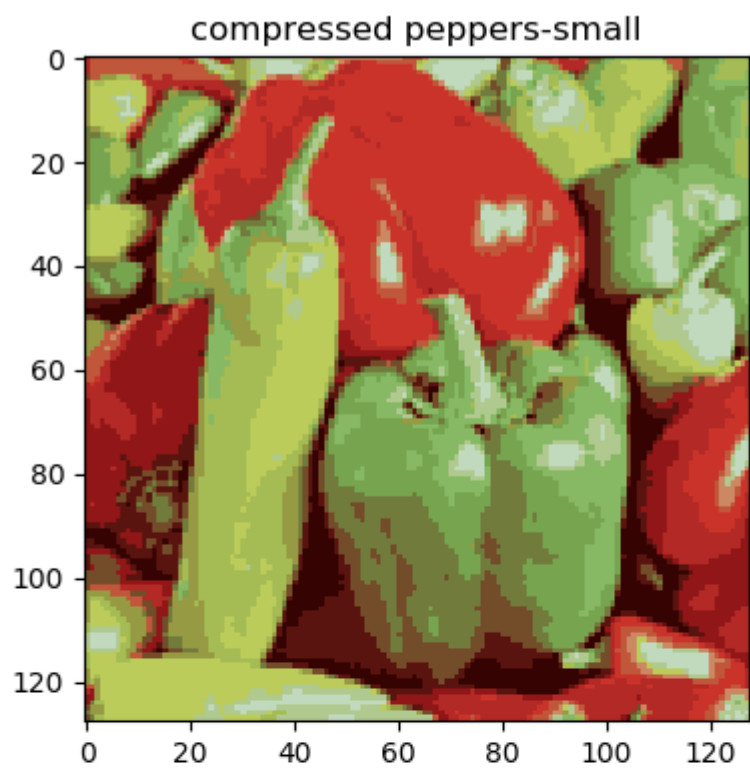
Compressed images:



compressed peppers-small



compressed peppers-large

b) K-means compression Factor

Answer:

To store the original image, each pixel requires 24bit for the RGB value (255 x 255 x 255 colors);

To store the compressed image, each pixel will only need 16 colors, which need only 4bit storage.

The compression factor is thus 24:4 = 6:1.