

1.

(a)(b)

LS 229 PS 3

Chen Cen

1. a. $O^{(i)} = \frac{1}{1 + e^{-z^{(i)}}}$

$$= g(w_0^{[2]} + w_1^{[2]} a_1 + w_2^{[2]} a_2 + w_3^{[2]} a_3)$$

$$= g(w_0^{[2]} + w_1^{[2]} a_1 + w_2^{[2]} (w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2) + w_3^{[2]} a_3)$$

$$= g(w_0^{[2]} + w_1^{[2]} a_1 + w_2^{[2]} (w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2) + w_3^{[2]} a_3)$$

$$\frac{\partial \mathcal{L}}{\partial w_{12}^{[1]}} = \frac{2}{m} \sum_{i=1}^m (O^{(i)} - y^{(i)}) \cdot \frac{\partial O^{(i)}}{\partial w_{12}^{[1]}}, \text{ because } O' = O(1-O)$$

$$= \frac{2}{m} \sum_{i=1}^m (O^{(i)} - y^{(i)}) \cdot O^{(i)} \cdot (1 - O^{(i)}) \cdot w_2^{[2]} \cdot (w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2) \cdot x_1$$

1. b. $O^{(i)} = f(w_0^{[2]} + w_1^{[2]} f(w_{01}^{[1]} + w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2) + w_2^{[2]} f(w_{02}^{[1]} + w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2) + w_3^{[2]} f(w_{03}^{[1]} + w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2))$

$\rightarrow x_1 < 0.5$
 $\rightarrow x_2 < 0.5$
 $\rightarrow x_1 + x_2 > 4$

From the graph, we know: if $x_1 < 0.5$ or $x_2 < 0.5$ or $x_1 + x_2 > 4$, then $y=1$. if $x_1 > 0.5$ and $x_2 > 0.5$ and $x_1 + x_2 \leq 4$, then $y=0$.

So we can let $w_0^{[2]} = -3$, $w_1^{[2]} = w_2^{[2]} = w_3^{[2]} = 1$, let 3 neuron denotes 3 edge,

when $x_1 < 0.5$

let $w_{02}^{[1]} + w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 \geq 0$
 $\Rightarrow w_{12}^{[1]} = -2, w_{22}^{[1]} = 1$
 $w_{22}^{[1]} = 0$

when $x_2 < 0.5$

let $w_{01}^{[1]} + w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 \geq 0$
 $\Rightarrow w_{21}^{[1]} = -2, w_{11}^{[1]} = 1$
 $w_{11}^{[1]} = 0$

when $x_1 + x_2 > 4$

let $w_{03}^{[1]} + w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 \geq 0$
 $\Rightarrow w_{03}^{[1]} = -4, w_{13}^{[1]} = w_{23}^{[1]} = 1$

(c)

$$\begin{aligned} \text{I.C } O^{(i)} &= f(W_0^{(i)} + W_1^{(i)}(W_0^{(1)} + W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2) \\ &\quad + W_2^{(2)}(W_0^{(2)} + W_{12}^{(2)}x_1 + W_{22}^{(2)}x_2) \\ &\quad + W_3^{(3)}(W_0^{(3)} + W_{13}^{(3)}x_1 + W_{23}^{(3)}x_2)) \\ &= f(a + bx_1 + cx_2) \end{aligned}$$

Note from existing dataset, we need

$$a + bx_1 + cx_2 \geq 0 \quad \text{when } \underbrace{x_1 \leq 0.5}_{\textcircled{1}}, \underbrace{x_2 \leq 0.5}_{\textcircled{2}}, \underbrace{x_1 + x_2 \geq 4}_{\textcircled{3}}$$

Note $\textcircled{1}, \textcircled{2}$ needs $O^{(i)}$ to be negatively correlated to x_1, x_2
while $\textcircled{3}$ needs $O^{(i)}$ to be positively correlated to $(x_1 + x_2)$
it's impossible to meet both $\textcircled{1}, \textcircled{2}$ and $\textcircled{3}$ for $a + bx_1 + cx_2$

2.

2. E step: For each set i , set

$$Q_i(z^{(n)}) := P(z^{(n)} | x, \theta)$$

M step: set

$$\theta := \arg \max_{\theta} \left(\sum_i \sum_{z^{(n)}} Q_i(z^{(n)}) \log \frac{P(x^{(n)}, z^{(n)}; \theta)}{Q_i(z^{(n)})} + \log P(\theta) \right)$$

because we set $Q_i(z^{(n)}) = P(z^{(n)} | x; \theta)$, Jensen's inequality holds,

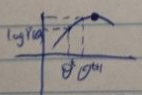
$$\text{so } \ell(\theta^{(t+1)}) \geq \sum_i \sum_{z^{(n)}} Q_i^{(t)}(z^{(n)}) \log \frac{P(x^{(n)}, z^{(n)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(n)})} + \log P(\theta^{(t+1)}) \quad \textcircled{1}$$

$$\geq \sum_i \sum_{z^{(n)}} Q_i^{(t)}(z^{(n)}) \log \frac{P(x^{(n)}, z^{(n)}; \theta^{(t)})}{Q_i^{(t)}(z^{(n)})} + \log P(\theta^{(t+1)}) \quad \textcircled{2}$$

because $\theta^{(t+1)}$ is so chosen that it maximizes the right side,

so $\textcircled{1} \rightarrow \textcircled{2}$ holds

~~because~~ because $\log P(\theta)$ is concave in θ and $\theta^{(t+1)}$ is closer to maximum point, so we have



$$\log P(\theta) \leq \log P(\theta^{(t+1)}) \quad \dots \star$$

Plug \star into $\textcircled{2}$

$$\textcircled{2} \geq \sum_i \sum_{z^{(n)}} Q_i^{(t)}(z^{(n)}) \log \frac{P(x^{(n)}, z^{(n)}; \theta^{(t)})}{Q_i^{(t)}(z^{(n)})} + \log P(\theta^{(t)})$$

$$= \ell(\theta^{(t)}) \quad \text{so we know } \ell(\theta^{(t+1)}) \geq \ell(\theta^{(t)})$$

3.
(a)

3.(a)

$$X^{(p)} = Y^{(p)} + Z^{(p)} + \varepsilon^{(p)}$$

$$\begin{matrix} \mathcal{N}(\mu_p, \sigma_p^2) & \mathcal{N}(\mu_c, \sigma_c^2) & \mathcal{N}(0, \sigma^2) \\ \downarrow & \downarrow & \downarrow \end{matrix}$$

$$E X = E Y + E Z + E \varepsilon$$

$$= \mu_p + \mu_c$$

$$\text{Var}(X) = \text{Var}(Y) + \text{Var}(Z) + \text{Var}(\varepsilon)$$

$$= \sigma_p^2 + \sigma_c^2 + \sigma^2$$

So mean vector of X & Y : (multivariate Gaussian)

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N}(\mu_0, \Sigma_0)$$

$$\mu_0 = \begin{bmatrix} \mu_p + \mu_c \\ \mu_p \end{bmatrix} \quad \Sigma_0 = \begin{bmatrix} \sigma_p^2 + \sigma_c^2 + \sigma^2 & \sigma_p^2 \\ \sigma_p^2 & \sigma_p^2 \end{bmatrix} \quad (1)$$

$$\tilde{\Sigma}_{xy} = \tilde{\Sigma}_{yx}$$

$$= \text{Cov}(X, Y) = \text{Cov}(Y + Z + \varepsilon, Y) = \text{Cov}(Y, Y) + \text{Cov}(Z, Y) + \text{Cov}(\varepsilon, Y)$$

$$= \sigma_p^2$$

$$\text{So } \Sigma_0 = \begin{bmatrix} \sigma_p^2 + \sigma_c^2 + \sigma^2 & \sigma_p^2 \\ \sigma_p^2 & \sigma_p^2 \end{bmatrix}$$

Similarly,

$$\begin{bmatrix} X \\ Z \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_p + \mu_c \\ \mu_c \end{bmatrix}, \begin{bmatrix} \sigma_p^2 + \sigma_c^2 + \sigma^2 & \sigma_c^2 \\ \sigma_c^2 & \sigma_c^2 \end{bmatrix} \right) \quad (2)$$

(b)

3. (b) We want to maximize,

$$\begin{aligned} & \sum \log P(z^{(i)}, y^{(i)} | x^{(i)}, \mu_p, \nu_z, \sigma_p^2, \tau_z^2) \\ &= \sum \log (P(z^{(i)} | x^{(i)}) \cdot P(y^{(i)} | x^{(i)})) \\ &= \sum \log \frac{1}{\sqrt{2\pi}\sigma_p} \exp\left(-\frac{1}{2\sigma_p^2}(y - \mu_p)^2\right) \cdot \frac{1}{\sqrt{2\pi}\tau_z} \exp\left(-\frac{1}{2\tau_z^2}(z - \nu_z)^2\right) \\ &= \sum \left(\log \frac{1}{\sqrt{2\pi}\sigma_p} - \frac{1}{2\sigma_p^2}(y^2 - 2y\mu_p + \mu_p^2) - \frac{1}{2\tau_z^2}(z^2 - 2z\nu_z + \nu_z^2) \right) \\ &\approx \sum \left(\log \frac{1}{\sqrt{2\pi}\sigma_p} - \frac{1}{2\sigma_p^2}(\underbrace{E(y^2)} - 2\underbrace{E(y)} \cdot \mu_p + \mu_p^2) \right. \\ &\quad \left. - \frac{1}{2\tau_z^2}(\underbrace{E(z^2)} - 2\underbrace{E(z)} \nu_z + \nu_z^2) \right) \end{aligned}$$

From ③ we know

$$E(y^2) = \frac{\sigma_p^2(\tau_z^2 + \sigma^2)}{\sigma_p^2 + \tau_z^2 + \sigma^2} \quad E(y) = \mu_p + \frac{\sigma_p^2}{\sigma_p^2 + \tau_z^2 + \sigma^2}(x - \mu_p - \nu_z)$$

$$E(z^2) = \frac{\tau_z^2(\sigma_p^2 + \sigma^2)}{\sigma_p^2 + \tau_z^2 + \sigma^2} \quad E(z) = \nu_z + \frac{\tau_z^2}{\sigma_p^2 + \tau_z^2 + \sigma^2}(x - \mu_p - \nu_z)$$

$$\text{Let } \textcircled{A} = \sum_{i=1}^P -\frac{1}{2\sigma_p^2}(E(y^2) - 2E(y) \cdot \mu_p + \mu_p^2)$$

$$\frac{\partial \textcircled{A}}{\partial \mu_p} = -\frac{1}{2\sigma_p^2} \sum_{i=1}^P (\cancel{E(y^2)} - 2E(y) + 2\mu_p) = -\frac{1}{2\sigma_p^2} \sum_{i=1}^P (-2E(y) + 2\mu_p) = 0$$

$$\Rightarrow \boxed{\mu_p = \frac{1}{P} \sum_{i=1}^P E(y)}$$

$$\textcircled{2A} \\ \text{Let } \frac{\partial \phi_p}{\partial \phi_p^2} = 0 \quad \phi_p^2 = ?$$

$$\text{Let } \textcircled{B} = \sum_{i=1}^R - \frac{1}{2\gamma_i^2} (\bar{\phi}_i^2 - 2E(Z) V_i + V_i^2)$$

$$\textcircled{2B} \\ \text{Let } \frac{\partial \phi}{\partial \gamma_i} = 0 \Rightarrow V_i = \frac{1}{R} \sum_{i=1}^R E(Z)$$

$$\textcircled{2C} \\ \text{Let } \frac{\partial \phi}{\partial \phi^2} = 0 \Rightarrow \phi^2 = ?$$

4.
(a)

4.(a) Note

$$-KL(P||Q) = \sum_x P(x) \log \frac{Q(x)}{P(x)}$$

because \log is a ~~convex~~ ^{concave} function, from Jensen's inequality

$$E(\log(x)) = \sum_x P(x) \log \frac{Q(x)}{P(x)} \leq \log(E(x)) = \log \sum_x P(x) \frac{Q(x)}{P(x)} = \log \sum_x Q(x)$$

$$\text{So } -KL(P||Q) \leq \log \sum_x Q(x)$$

Since $Q(x)$ is probability, $\sum_x Q(x) = 1$

$$\text{So } -KL(P||Q) \leq \log 1 = 0$$

$$\Rightarrow KL(P||Q) \geq 0$$

If $P=Q$, $Q(x)=P(x)$, so $\log \frac{Q(x)}{P(x)} = 1$

$$\Rightarrow KL(P||Q) = \sum_x P(x) \cdot \log 1 = 0$$

(b)

(b). Note $P(x, y) = P(x) \cdot P(y|x)$, $Q(x, y) = Q(x) \cdot Q(y|x)$

So $KL(P(x, y) || Q(x, y))$

$$= KL(P(x) \cdot P(y|x) || Q(x) \cdot Q(y|x))$$

$$= \sum_x (P(x) \sum_y P(y|x) \log \frac{P(x) \cdot P(y|x)}{Q(x) \cdot Q(y|x)})$$

$$= \underbrace{\sum_x (P(x) \sum_y P(y|x) \log \frac{P(x)}{Q(x)})}_{(1)} + \underbrace{\sum_x (P(x) \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)})}_{(2)}$$

Note $P(x) \cdot \sum_y P(y|x) = P(x)$

$$(1) \left\{ \begin{array}{l} \text{So } (1) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = KL(P(x) || Q(x)) \end{array} \right.$$

$$(2) \left\{ \begin{array}{l} \text{Per KL between 2 conditional distributions,} \\ (2) = KL(P(y|x) || Q(y|x)) \end{array} \right.$$

Therefore $KL(P(x, y) || Q(x, y)) = KL(P(x) || Q(x)) +$

$$KL(P(y|x) || Q(y|x))$$

(c)

$$\begin{aligned} \text{(c). } KL(\hat{P} \parallel P_\theta) &= \sum_x \hat{p}(x) \log \frac{\hat{p}(x)}{P_\theta(x)} \\ &= \underbrace{\sum_x \hat{p}(x) \log \hat{p}(x)}_{(1)} - \underbrace{\sum_x \hat{p}(x) \log P_\theta(x)}_{(2)} \end{aligned}$$

we are looking for $\arg \max_{\theta} KL(\hat{P} \parallel P_\theta)$, so θ has no effect on (1)

Since θ will minimize (2), it will maximize $-(2)$

$$\begin{aligned} -(2) &= \sum_x \hat{p}(x) \log P_\theta(x) \\ &= \sum_x \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{x^{(i)}=x\}} \log P_\theta(x) \\ &= \frac{1}{m} \sum_x \left(\underbrace{\sum_{i=1}^m \mathbb{1}_{\{x^{(i)}=x\}}}_{\star} \log P_\theta(x) \right) \end{aligned}$$

For \star , it's \sum of $\log P_\theta(x)$ where $\forall x \neq x^{(i)}$ are filtered

$$\text{so } \star = \log P_\theta(x^{(i)})$$

$$\boxed{-(2) = \frac{1}{m} \sum_x \log P_\theta(x^{(i)})} \rightarrow \text{maximized by } \theta$$

$$\text{so } \arg \max_{\theta} (-(2)) = \arg \max_{\theta} \sum_{i=1}^m \log P_\theta(x^{(i)})$$

5.

(a)

```
def draw_image():
    A = imread('mandrill-large.tiff')
    plt.imshow(A)
    plt.show()
```

(b)

```
class Centroid:
    def __init__(self, image_matrix):
        self.r = random.randint(0, 255)
        self.g = random.randint(0, 255)
        self.b = random.randint(0, 255)
        # a list of index of the pixel assigned to this centroid
        self.assigned_pixels = set()
        self.image_matrix = image_matrix
        (self.height, self.width, c) = image_matrix.shape

    def add_pixel(self, pixel_index):
        self.assigned_pixels.add(pixel_index)

    def remove_pixel(self, pixel_index):
        self.assigned_pixels.remove(pixel_index)

    def contains_pixel(self, pixel_index):
        return pixel_index in self.assigned_pixels

    def update_rgb(self):
        # recalculate this centroids rgb value
        pixel_count = len(self.assigned_pixels)
        if pixel_count == 0:
            return
        red_sum = green_sum = blue_sum = 0
        for index in self.assigned_pixels:
            red_sum += self.get_red_value(index)
            green_sum += self.get_green_value(index)
            blue_sum += self.get_blue_value(index)
        self.r = red_sum / pixel_count
        self.g = green_sum / pixel_count
        self.b = blue_sum / pixel_count

    def distance_from_pixel(self, pixel_index):
        return (self.get_red_value(pixel_index) - self.r) * (self.get_red_value(pixel_index) - self.r) + \
            (self.get_green_value(pixel_index) - self.g) * (self.get_green_value(pixel_index) - self.g) + \
            (self.get_blue_value(pixel_index) - self.b) * (self.get_blue_value(pixel_index) - self.b)

    def get_red_value(self, pixel_index):
        return self.image_matrix[int(pixel_index / self.width), pixel_index % self.width, 0]

    def get_green_value(self, pixel_index):
        return self.image_matrix[int(pixel_index / self.width), pixel_index % self.width, 1]
```

```
def get_blue_value(self, pixel_index):
    return self.image_matrix[int(pixel_index / self.width), pixel_index % self.width, 2]
```

```
def shout(self):
    print("r: " + str(self.r))
    print("g: " + str(self.g))
    print("b: " + str(self.b))
```

```
# return if this pixel is assigned to a different centroid
def assign_pixel_to_centroid(pixel_index, centroids):
    min_cost = sys.maxsize
    new_index = -1
    for i, centroid in enumerate(centroids):
        cost = centroid.distance_from_pixel(pixel_index)
        if cost < min_cost:
            min_cost = cost
            new_index = i
```

```
# remove from previous centroid and add to current centroid
is_updated = False
for i, centroid in enumerate(centroids):
    contains_index = centroid.contains_pixel(pixel_index)
    if i == new_index:
        is_updated = not contains_index
        centroid.add_pixel(pixel_index)
    else:
        if contains_index:
            centroid.remove_pixel(pixel_index)
```

```
return is_updated
```

```
# problem 5.b
def k_means(k, image_matrix):
    centroids = []
    for i in range(k):
        centroids.append(Centroid(image_matrix))
    (height, width, color) = image_matrix.shape
```

```
pixel_updated = 1
iterations = 0
max_iteration = 30
# convergence check is set to no pixel is updated
while pixel_updated > 0 and iterations < max_iteration:
    print("pixel_updated: " + str(pixel_updated))
    iterations += 1
    pixel_updated = 0
    for i in range(height):
        for j in range(width):
            pixel_index = i * width + j
            # count # of pixel updated
            if assign_pixel_to_centroid(pixel_index, centroids):
                pixel_updated += 1
    for centroid in centroids:
        centroid.update_rgb()
    if pixel_updated == 0:
        print("converges in " + str(iterations) + " rounds")
    else:
        print("terminates after " + str(max_iteration) + " rounds")
```



```
return centroids
```

```
def update_picture(centroids, image_matrix):  
    (height, width, c) = image_matrix.shape
```

```
    for centroid in centroids:  
        for pixel_index in centroid.assigned_pixels:  
            y = int(pixel_index / width)  
            x = pixel_index % width  
            image_matrix[x, y, 0] = centroid.r  
            image_matrix[x, y, 1] = centroid.g  
            image_matrix[x, y, 2] = centroid.b
```

```
# problem 5.b
```

```
def run_k_means():  
    # cluster R3 (rgb) to k=16 clusters  
    # a) initialize 16 centroid of R3(0-255, 0-255, 0-255) u = int[16]  
    # b) for each pixel x: repeat until no centroid moves, no new assignment to existing points  
    # i) find which cluster this pixel belongs to by minimize this cost function:  
    #  $J = \sum((x_r - u_r)^2 + (x_g - u_g)^2 + (x_b - u_b)^2)$   
    # ii) update each centroid u by taking the average  
    # let Xs be the pixels assigned to centroid ui  
    # for each ui  
    #  $u_i_r = \sum(Xs_r) / \text{len}(Xs)$   
    #  $u_i_g = \sum(Xs_g) / \text{len}(Xs)$   
    #  $u_i_b = \sum(Xs_b) / \text{len}(Xs)$   
    #  
    #  $Xs_r = \text{sum of } r \text{ value of each } X \text{ in } Xs$ 
```

```
# A = imread('mandrill-small.tiff')
```

```
# sample outputs
```

```
# pixel_updated: 1  
# pixel_updated: 16384  
# pixel_updated: 15263  
# pixel_updated: 8300  
# pixel_updated: 6210  
# pixel_updated: 5682  
# pixel_updated: 4929  
# pixel_updated: 8413  
# pixel_updated: 5804  
# pixel_updated: 4751  
# pixel_updated: 5533  
# pixel_updated: 6313  
# pixel_updated: 5010  
# pixel_updated: 4593  
# pixel_updated: 5795  
# pixel_updated: 6323  
# pixel_updated: 4351  
# pixel_updated: 5171  
# pixel_updated: 5584  
# pixel_updated: 5355  
# pixel_updated: 4202  
# pixel_updated: 4930  
# pixel_updated: 5971  
# pixel_updated: 5217  
# pixel_updated: 4784  
# pixel_updated: 5400  
# pixel_updated: 5999
```

```
# pixel_updated: 4549
# pixel_updated: 4939
# pixel_updated: 5027
# terminates after 30 rounds

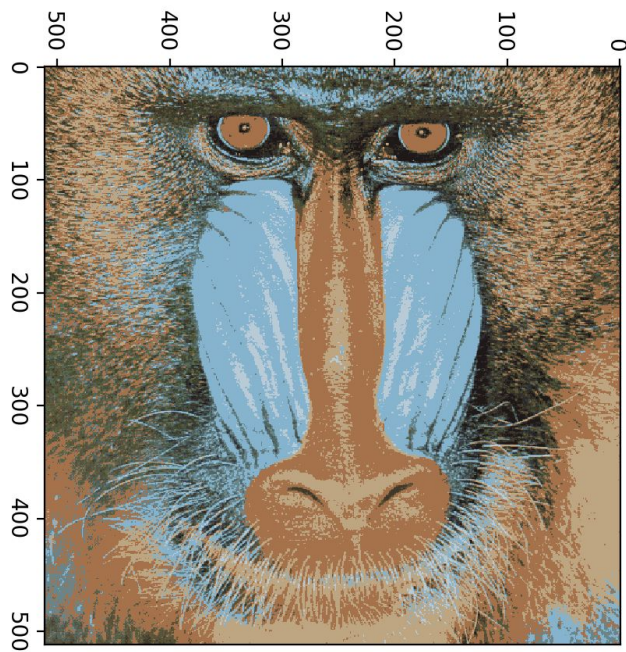
# A = imread('mandrill-large.tiff')
A = imread('mandrill-large.tiff')
A.setflags(write=1)

centroids = k_means(16, A)

update_picture(centroids, A)

plt.imshow(A)
plt.show()
```

(c)




```
#problem 5.c
def update_picture(centroids, image_matrix):
    (height, width, c) = image_matrix.shape

    for centroid in centroids:
        for pixel_index in centroid.assigned_pixels:
            y = int(pixel_index / width)
            x = pixel_index % width
            image_matrix[x, y, 0] = centroid.r
            image_matrix[x, y, 1] = centroid.g
            image_matrix[x, y, 2] = centroid.b
```

(d)

5. (a) (b) (c)

See code and image

(d) Before there are $256 \times 256 \times 256$ colors
 it takes $1 + 1 + 1 \overset{=3}{}$ byte to store

After there are only 16 colors
 why takes 4bit / 0.5 byte to store

So we compressed the factor of $\frac{3}{0.5} = 6$