

1) Recall from Gaussian mixture distribution that it can be written as $p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$. Let's also introduce K-dimensional binary random variable \mathbf{z} in which particular element $z_k = 1$ and all others are equal to zero.

Values z_k satisfy $\forall k \in \{0, 1\}$ and $\sum_k z_k = 1$. After this, we define marginal distribution over \mathbf{z} in terms of mixing coefficients π_k where $P(z_k=1) = \pi_k$ such that $\sum_k \pi_k = 1$ with $0 \leq \pi_k \leq 1$.

Since \mathbf{z} uses 1-of-K representation, we can also rewrite this distribution as $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \Rightarrow p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} = \prod_{k=1}^K \pi_k$. By summing up joint distribution over all possible states of \mathbf{z}

$$p(x|\mathbf{z}) = \prod_{k=1}^K N(x | \mu_k, \Sigma_k)$$

$$p(x) = \sum_{\mathbf{z}} p(\mathbf{z}) p(x|\mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

Another quantity $\gamma(z_k)$ to denote $P(z_k=1|x)$ where from Bayes

$$\gamma(z_k) = P(z_k=1|x) = \frac{P(z_k=1) p(x|z_k=1)}{\sum_{j=1}^K p(z_j=1) p(x|z_j=1)} = \frac{\pi_k N(x | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x | \mu_j, \Sigma_j)}$$

We will view π_k as prior probability of $z_k=1$, and quantity $\gamma(z_k)$ as corresponding posterior probability once we have observed x . In fact, $\gamma(z_k)$ is viewed as the responsibility that component k takes for explaining the observation x .

In general, EM Algorithm can be regarded as a soft-version of K-means for Gaussian mixture model. The EM algorithm uses responsibilities to make a soft assignment of each data point to one of the clusters. When ϵ is fixed, responsibility of data point i assigning to cluster K is given by the following relation:

$$r_i^{(K)} = \frac{\exp\left(-\frac{1}{2\sigma^2} \|y_i - \mu_K\|^2\right)}{\sum_{p=1}^K \exp\left(-\frac{1}{2\sigma^2} \|y_i - \mu_p\|^2\right)}$$

It's observable fact that as $\epsilon \rightarrow 0$, $r_i^{(K)} \rightarrow 1$ for the cluster K that is closest to y_i and $r_i^{(K)} \rightarrow 0$ for other clusters, indicating how K-means algorithm is recovered as

More rigorous: Notice $r_i^{(K)} = \frac{\exp\left(-\frac{1}{2\sigma^2} \|y_i - \mu_K\|^2\right)}{\sum_{p=1}^K \exp\left(-\frac{1}{2\sigma^2} \|y_i - \mu_p\|^2\right)}$

$$\sum_{p=1}^K \exp\left(-\frac{1}{2\sigma^2} (\|y_i - \mu_p\|^2 - \|y_i - \mu_K\|^2)\right)$$

$$\exp\left(-\frac{1}{2\sigma^2} (\|y_i - \mu_p\|^2 - \|y_i - \mu_K\|^2)\right) \rightarrow \begin{cases} 1, & \|y_i - \mu_p\|^2 = \|y_i - \mu_K\|^2 \\ 0, & \|y_i - \mu_p\|^2 > \|y_i - \mu_K\|^2 \\ \infty, & \|y_i - \mu_p\|^2 < \|y_i - \mu_K\|^2 \end{cases}$$

Thus, it means if i is the unique closest cluster to y_i , then the denominator of $r_i^{(1)}$ will have all terms tending to zero, except for 1 term tending to 1, so the limit is $\frac{1}{1+0^+} = 1$ while the denominator of $r_i^{(k)}$ for $k \neq i$ will have at least 1 term ending to ∞ , making the limit 0 instead.

It's also worthwhile to mention that if there are multiple "closest clusters" to a point y_i , then the limiting probability will be uniform among these clusters.

Now, coming to the more rigorous approach where we initially defined $f(z_{hk})$ as $\pi_k N(x_h | \mu_k, \Sigma_k)$

$$\frac{\sum_{j=1}^K \pi_j N(x_h | \mu_j, \Sigma_j)}{\sum_{j=1}^K \pi_j N(x_h | \mu_j, \Sigma_j)}$$

since we are considering GMM with covariance matrices defined as Σ (where without loss of generality, we assume it's variance parameter)

in which the normal distribution becomes

$$p(x | \mu_k, \Sigma_k) = \frac{1}{\sqrt{(2\pi)^M \det \Sigma}} e^{-\frac{1}{2\Sigma} \|x - \mu_k\|^2}$$

where $f(z_{hk})$ is accompanied in the formula as

$\pi_k e$

$$\frac{\sum_{j=1}^K \frac{-\|x_h - \mu_j\|^2}{2\Sigma}}{\sum_{j=1}^K \pi_j e}$$

z_{hk}

e

Now, if we also define

$$f_{hk} = \begin{cases} 1, & \text{if } k = \arg \min_j \|x_h - \mu_j\|^2 \\ 0, & \text{in the other case} \end{cases}$$

Taking limit $\Sigma \rightarrow 0$, we'll showcase how responsibilities $f(z_{hk})$ for data pt x_h will go to unity except for specific term

Denote the distance by δ_k where $\delta_k = \|x_n - \mu_k\|^2$

If we substitute those notations into the original formula, one can

get that $\pi_k e^{-\frac{\|x_n - \mu_k\|^2}{2\varepsilon}} = \pi_k e^{-\frac{\delta_k}{2\varepsilon}}$

$$\frac{\sum_{j=1}^K e^{-\frac{\|x_n - \mu_j\|^2}{2\varepsilon}}}{\sum_{j=1}^K \pi_j e^{-\frac{\|x_n - \mu_j\|^2}{2\varepsilon}}} = \frac{\sum_{j=1}^K e^{-\frac{\delta_j}{2\varepsilon}}}{\sum_{j=1}^K \pi_j e^{-\frac{\delta_j}{2\varepsilon}}}$$

If one defines the minimum value of all δ_i 's, let us denote it as P : $\min(\delta_1, \delta_2, \dots, \delta_K) = P$.

then $\pi_k e^{-\frac{\delta_k}{2\varepsilon}} = \pi_k e^{-\frac{P-\delta_k}{2\varepsilon}}$

$$\frac{\sum_{j=1}^K e^{-\frac{P-\delta_j}{2\varepsilon}}}{\sum_{j=1}^K \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}}} = \frac{\sum_{j=1}^K e^{-\frac{P-\delta_j}{2\varepsilon}}}{\sum_{j=1}^K \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}}}$$

it's clear that $P \leq \delta_i$ for all i , and if $P = \delta_t$ for some t ,

then taking k different from those t indices yield that

$$\lim_{\varepsilon \rightarrow 0} \frac{\pi_k e^{-\frac{P-\delta_k}{2\varepsilon}}}{\sum_{j=1}^K \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}}} = \lim_{\varepsilon \rightarrow 0} \frac{\pi_k e^{-\frac{P-\delta_k}{2\varepsilon}}}{\sum_{j=1}^t \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}} + \sum_{j \neq t} \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}}} \quad \text{with}$$

$$P \neq \delta_k \text{ giving } \lim_{\varepsilon \rightarrow 0} e^{-\frac{P-\delta_k}{2\varepsilon}} = \lim_{\varepsilon \rightarrow 0} \frac{1}{e^{\frac{P-\delta_k}{2\varepsilon}}} = \lim_{\varepsilon \rightarrow 0} \frac{1}{\infty} = 0 \quad \checkmark$$

However, choosing the values k in which δ_k becomes the minimum (set of values t)

$$\text{then, } \lim_{\varepsilon \rightarrow 0} \frac{\pi_t e^{-\frac{P-\delta_t}{2\varepsilon}}}{\sum_{j=1}^K \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}}} = \lim_{\varepsilon \rightarrow 0} \frac{\pi_t}{\pi_t + \sum_{j \neq t} \pi_j e^{-\frac{P-\delta_j}{2\varepsilon}}} = 1 \quad \checkmark$$

which concludes that Pointing case is equivalent to K-means clustering.

3) We first prove that $(I + UVV^T)^{-1} = I - U(V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1}$ where $I \in \mathbb{R}^{n \times n}$, $V \in \mathbb{R}^{m \times n}$, $U \in \mathbb{R}^{m \times m}$.

The formula could be shown by checking that $I + UVV^T$ times the desired inverse on the RHS gives identity.

$$(I + UVV^T) \left[I - U(V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} \right] =$$

$$= I + UVV^T I^{-1} - U(V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} - \\ - UVV^T I^{-1} U (V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} = [I + UVV^T I^{-1}]$$

$$- \{ U(V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} + UVV^T I^{-1} U (V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} \}$$

$$= \{ I + UVV^T I^{-1} \} - \left[U(V^{-1} + V^T I^{-1} U)^{-1} + UVV^T I^{-1} U (V^{-1} + V^T I^{-1} U)^{-1} \right]$$

$$= \{ I + UVV^T I^{-1} \} - \left[U(V^{-1} + V^T I^{-1} U)^{-1} + UVV^T I^{-1} U (V^{-1} + V^T I^{-1} U)^{-1} \right] V^T I^{-1}$$

$$= \{ I + UVV^T I^{-1} \} - \left[U + UVV^T I^{-1} U (V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} \right] =$$

$$= \{ I + UVV^T I^{-1} \} - \left[U + UVV^T I^{-1} U (V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} \right] (V^{-1} + V^T I^{-1} U)^{-1} =$$

$$= I + UVV^T I^{-1} - U(V^{-1} + V^T I^{-1} U) (V^{-1} + V^T I^{-1} U)^{-1} V^T I^{-1} =$$

$$= I + UVV^T I^{-1} - UVV^T I^{-1} = I \text{ since the squiggle lines are multiplied into identity matrix } I = \boxed{I}.$$

Now, we'll prove the first task of given question. The key factor is to consider that A, B -symmetric matrices since they are covariance matrices, with size $A, B \in \mathbb{R}^{D \times D}$ and if we take into account $a, b \in \mathbb{R}^D$, then column n -vector c and covariance matrix $C \in \mathbb{R}^{D \times D}$ become in dimension $(\mathbb{R}^{D \times 1})$ and $\mathbb{R}^{D \times D}$, respectively. ($c \in \mathbb{R}^D, C \in \mathbb{R}^{D \times D}$)

Initially, we apply the given formulas for Gaussian distributions and check whether they satisfy or not.

$$N(x|a, A) = \frac{1}{\sqrt{(2\pi)^D \det(A)}} e^{-\frac{1}{2}(x-a)^T A^{-1} (x-a)}$$

$$N(x|b, B) = \frac{1}{\sqrt{(2\pi)^D \det(B)}} e^{-\frac{1}{2}(x-b)^T B^{-1} (x-b)}$$

$$N(x|c, C) = \frac{1}{\sqrt{(2\pi)^D \det(C)}} e^{-\frac{1}{2}(x-c)^T C^{-1} (x-c)}$$

$$N(x|a, A) N(x|b, B) = \frac{1}{\sqrt{(2\pi)^D \det(A+B)}} e^{-\frac{1}{2}(x-a)^T (A+B)^{-1} (x-b)}$$

$$N(x|a, A) N(x|b, B) = \frac{1}{\sqrt{(2\pi)^D \det(A+B)}} e^{-\frac{1}{2}(x-a)^T (A+B)^{-1} (x-b)}$$

Since the $(2\pi)^D$ terms are canceled out whenever we multiply. It's sufficient to show two properties: $\det(A)\det(B)=\det(A+B)$.

$$\text{and } -\frac{1}{2}(x-a)^T A^{-1} (x-a) + -\frac{1}{2}(x-b)^T B^{-1} (x-b) = -\frac{1}{2}(x-a)^T (A+B)^{-1} (x-a) + -\frac{1}{2}(x-b)^T (A+B)^{-1} (x-b)$$

Claim: $\det(A) \det(B) = \det(C) \det(A+B)$

Proof: $C = (A^{-1} + B^{-1})^{-1} \Rightarrow C^{-1} = A^{-1} + B^{-1}$ where we know that
 $\det(C^{-1}) = \frac{1}{\det(C)} = \det(A^{-1} + B^{-1})$, $\det(C) = \frac{1}{\det(A^{-1} + B^{-1})}$

$A \otimes (A^{-1} + B^{-1}) = A \otimes A^{-1} + A$, It's enough to show $\det(A) \det(B) =$
 $= \frac{1}{\det(A^{-1} + B^{-1})} \det(A+B)$ or $\det(A+B) = \det(A^{-1} + B^{-1}) \det(A) \det(B)$

$A(A^{-1} + B^{-1})B = AA^{-1}B + AB^{-1}B = B + A$, meaning that
 $\det(A(A^{-1} + B^{-1})B) = \det(A+B)$. From the multiplication rule
we get $\det(A+B) = \det(A) \det(A^{-1} + B^{-1}) \det(B)$ ✓ as desired

Claim: $-\frac{1}{2}(x-a)^T A^{-1}(x-a) + \frac{-1}{2}(x-b)^T B^{-1}(x-b) = -\frac{1}{2}(a-b)^T (A+B)^{-1} \cdot (a-b)$

$+ -\frac{1}{2}(x-c)^T C^{-1}(x-c)$

Proof: It's enough to show $(x-a)^T A^{-1}(x-a) + (x-b)^T B^{-1}(x-b) =$

$= (a-b)^T (A+B)^{-1} (a-b) + (x-c)^T C^{-1}(x-c)$

$(x^T - a^T) A^{-1}(x-a) = x^T A^{-1} x - a^T A^{-1} x - x^T A^{-1} a + a^T A^{-1} a$

$(a-b)^T (A+B)^{-1} (a-b) = (a^T - b^T) (A+B)^{-1} (a-b) = a^T (A+B)^{-1} a -$

$- b^T (A+B)^{-1} a - a^T (A+B)^{-1} b + b^T (A+B)^{-1} b$. Plugging those

values yield $x^T A^{-1} x - a^T A^{-1} x - x^T A^{-1} a + a^T A^{-1} a +$

$+ x^T B^{-1} x - b^T B^{-1} x - x^T B^{-1} b + b^T B^{-1} b$

$$\begin{aligned}
&= \alpha^T (A + \beta)^{-1} q - \beta^T (A + \beta)^{-1} q - \alpha^T (A + \beta)^{-1} b + \beta^T (A + \beta)^{-1} b + \\
&+ x^T C^{-1} x - c^T C^{-1} x - x^T C^{-1} c + c^T C^{-1} c = \\
&= \alpha^T (A + \beta)^{-1} q - \beta^T (A + \beta)^{-1} q - \alpha^T (A + \beta)^{-1} b + \beta^T (A + \beta)^{-1} b + \\
&+ x^T A^{-1} x + x^T \beta^{-1} x - (A^{-1} \alpha + \beta^{-1} b)^T C^T C^{-1} x - \\
&- x^T C^{-1} C (A^{-1} \alpha + \beta^{-1} b) + (A^{-1} \alpha + \beta^{-1} b)^T C^T C^{-1} c = \\
&= \alpha^T (A + \beta)^{-1} q - \beta^T (A + \beta)^{-1} q - \alpha^T (A + \beta)^{-1} b + \beta^T (A + \beta)^{-1} b + \\
&+ x^T \cancel{A^{-1} x} + x^T \cancel{\beta^{-1} x} - (A^{-1} \alpha + \beta^{-1} b)^T x - x^T (A^{-1} \alpha + \beta^{-1} b) + \\
&+ (A^{-1} \alpha + \beta^{-1} b)^T C (A^{-1} \alpha + \beta^{-1} b) \quad \text{It's enough to show} \\
&- \cancel{\alpha^T A^{-1} x} - \cancel{x^T A^{-1} q} + \alpha^T A^{-1} q - \cancel{\beta^T \beta^{-1} x} - \cancel{x^T \beta^{-1} b} + \beta^T \beta^{-1} b \\
&\stackrel{?}{=} \alpha^T (A + \beta)^{-1} q - \beta^T (A + \beta)^{-1} q - \alpha^T (A + \beta)^{-1} b + \beta^T (A + \beta)^{-1} b \\
&- (A^{-1} \alpha + \beta^{-1} b)^T x - x^T (A^{-1} \alpha + \beta^{-1} b) + (A^{-1} \alpha + \beta^{-1} b)^T C \cdot \\
&(A^{-1} \alpha + \beta^{-1} b)^T = (A^{-1} \alpha)^T + (\beta^{-1} b)^T = \alpha^T (A^{-1})^T + \beta^T (\beta^{-1})^T \\
&\text{Since } (A^{-1})^T = (A^T)^{-1} = A^{-1} \text{ as } A, \beta \text{ symmetric} \Rightarrow \\
&- \cancel{\alpha^T A^{-1} x} + \alpha^T A^{-1} q - \cancel{\beta^T \beta^{-1} x} + \beta^T \beta^{-1} b = \alpha^T (A + \beta)^{-1} q - \beta^T (A + \beta)^{-1} q \\
&- \cancel{\alpha^T (A + \beta)^{-1} b} + \beta^T (A + \beta)^{-1} b - \cancel{\alpha^T A^{-1} x} - \cancel{\beta^T \beta^{-1} x} + \\
&+ (A^{-1} \alpha + \beta^{-1} b)^T (A^{-1} + \beta^{-1})^{-1} (A^{-1} \alpha + \beta^{-1} b), \text{ so it's left} \\
&\text{as follows}
\end{aligned}$$

$$a^T A^{-1} a + \beta^T \beta^{-1} b = a^T (A + \beta)^{-1} a - \beta^T (A + \beta)^{-1} \beta - a^T (A + \beta)^{-1} \beta + \beta^T (A + \beta)^{-1} b + \underbrace{(a^T A^{-1} + \beta^T \beta^{-1}) (A^{-1} + \beta^{-1})^{-1}}_{(A^{-1} + \beta^{-1})^{-1}} (A^{-1} a + \beta^{-1} b)$$

Since we proved that $A(A^{-1} + \beta^{-1})\beta = \beta + A$ and in the similar case, $\beta(A^{-1} + \beta^{-1})A = A + \beta$ taking inverses of both sides

$$(A + \beta)^{-1} = \beta^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} = A^{-1} (A^{-1} + \beta^{-1})^{-1} \beta^{-1}$$

Computing the straight line $\rightarrow a^T A^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} a +$

$$+ \beta^T \beta^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} a + a^T A (A + \beta)^{-1} \beta^{-1} b +$$

$$+ \beta^T \beta^{-1} (A^{-1} + \beta^{-1})^{-1} \beta^{-1} b = a^T A^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} a +$$

$$+ \beta^T \beta^{-1} (A^{-1} + \beta^{-1})^{-1} \beta^{-1} b + \beta^T (A + \beta)^{-1} a + a^T (A + \beta)^{-1} b$$

$$a^T A^{-1} a + \beta^T \beta^{-1} b = a^T (A + \beta)^{-1} a - \cancel{\beta^T (A + \beta)^{-1} a} - \cancel{a^T (A + \beta)^{-1} b} +$$

$$+ \beta^T (A + \beta)^{-1} b + \cancel{\beta^T (A + \beta)^{-1} a} + \cancel{a^T (A + \beta)^{-1} b} +$$

$$+ a^T A^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} a + \beta^T \beta^{-1} (A^{-1} + \beta^{-1})^{-1} \beta^{-1} b$$

$$a^T A^{-1} a + \beta^T \beta^{-1} b = a^T (A + \beta)^{-1} a + \beta^T (A + \beta)^{-1} b +$$

$$+ a^T A^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} a + \beta^T \beta^{-1} (A^{-1} + \beta^{-1})^{-1} \beta^{-1} b$$

$$(CD)^{-1} = D^{-1} C^{-1} \Rightarrow (A^{-1} + \beta^{-1})^{-1} A^{-1} = (A(A^{-1} + \beta^{-1}))^{-1} = (I + A\beta^{-1})^{-1}$$

$$\text{Since } (A + \beta)^{-1} = \beta^{-1} (A^{-1} + \beta^{-1})^{-1} A^{-1} \Rightarrow \beta(A + \beta)^{-1} A =$$

$$= (A^{-1} + \beta^{-1})^{-1}, \text{ where } (A^{-1} + \beta^{-1})^{-1} = \beta(A + \beta)^{-1} \text{ and similarly}$$

$(A^{-1} + \beta^{-1})^{-1} \beta^{-1}$ can be computed by $\rightarrow (A + \beta)^{-1} = A^{-1}(A^{-1} + \beta^{-1})^{-1} \beta^{-1}$
 $A(A + \beta)^{-1} \beta = (A^{-1} + \beta^{-1})^{-1}$ with $(A^{-1} + \beta^{-1})^{-1} \beta^{-1} = A(A + \beta)^{-1}$

$a^T A^{-1} a + b^T \beta^{-1} b = a^T (A + \beta)^{-1} a + b^T (A + \beta)^{-1} b +$ second expression if
 $+ a^T A^{-1} \beta (A + \beta)^{-1} a + b^T \beta^{-1} A (A + \beta)^{-1} b$ where the
 $[a^T (A + \beta)^{-1} a + a^T A^{-1} \beta (A + \beta)^{-1} a] + [b^T (A + \beta)^{-1} b + b^T \beta^{-1} A (A + \beta)^{-1} b]$
 $= (a^T + a^T A^{-1} \beta) (A + \beta)^{-1} a + (b^T + b^T \beta^{-1} A) (A + \beta)^{-1} b =$
 $= a^T (I + A^{-1} \beta) (A + \beta)^{-1} a + b^T (I + \beta^{-1} A) (A + \beta)^{-1} b.$ So

$a^T \boxed{A^{-1} a + b^T \beta^{-1} b} = a^T \boxed{(I + A^{-1} \beta) (A + \beta)^{-1} a + b^T \boxed{(I + \beta^{-1} A) (A + \beta)^{-1} b}}$
Claim: $A^{-1} = (I + A^{-1} \beta) (A + \beta)^{-1}$ and $\beta^{-1} = (I + \beta^{-1} A) (A + \beta)^{-1}$
Proof: $A \left((I + A^{-1} \beta) (A + \beta)^{-1} \right) = (A + \beta) (A + \beta)^{-1} = I$, therefore
 $\boxed{A^{-1} = (I + A^{-1} \beta) (A + \beta)^{-1}}$ Similarly, $\boxed{\beta \left((I + \beta^{-1} A) (A + \beta)^{-1} \right)} =$
 $= (\beta (I + \beta^{-1} A)) (A + \beta)^{-1} = (\beta + A) (\beta + A)^{-1} = I \Rightarrow$ Therefore,
 $\boxed{\beta^{-1} = (I + \beta^{-1} A) (A + \beta)^{-1}}$ ✓ Following the above results, it's
 indeed holds. Therefore, the topmost claim would also satisfy.
 In conclusion, we obtain that $N(x|0, A) N(x|\beta, \beta) =$
 holds correctly $\boxed{\checkmark} \bullet \boxed{\oplus}$ $= f^{-1} N(x|c, C)$ $\boxed{\checkmark}$

4) a) In order to solve the problem, we'll prove the general result:

Claim: Let $p(x) = N(x|\mu, \Sigma)$ and $p(y|x) = N(y|Ax+b, L)$

then $p(y) = N(y|A\mu+b, L + A\Sigma A^T)$ and

$p(x|y) = N(x|\Sigma^{-1}(y-b) + \Sigma^{-1}\mu, \Sigma)$ where Σ is defined as
 $\Sigma = (\Sigma^{-1} + A^T L^{-1} A)$

Proof: $y|x \sim N(Ax+b, L)$ and $x \sim N(\mu, \Sigma)$

$$\begin{aligned} y^T L^{-1} y &= 2y^T L^{-1} (Ax+b) + (Ax+b)^T L^{-1} (Ax+b) + \\ &+ x^T \Sigma^{-1} x - 2x^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu = \\ &= y^T L^{-1} y - 2y^T L^{-1} Ax - 2y^T L^{-1} b + \underline{x^T A^T L^{-1} A x} + 2x^T A^T L^{-1} b + \\ &+ b^T L^{-1} b + \underline{x^T \Sigma^{-1} x - 2x^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu}. \end{aligned}$$

The key idea is to make quadratic forms out of the x terms to integrate them out: By grouping with straight line terms and squirly terms, with the fact that $(y^T L^{-1} Ax)^T = x^T A^T L^{-1} y$ due to the $y^T L^{-1} Ax$ being scalar and Σ -symmetric \Rightarrow L^{-1} -symmetric

$$\begin{aligned} x^T (\Sigma^{-1} + A^T L^{-1} A) x - 2x^T (\Sigma^{-1} \mu + A^T L^{-1} y - A^T L^{-1} b) + \\ + \mu^T \Sigma^{-1} \mu + b^T L^{-1} b \end{aligned}$$

Let $V = [A^T L^{-1} A + \Sigma^{-1}]$

$$U = [\Sigma^{-1} \mu + A^T L^{-1} y - A^T L^{-1} b] \quad C = \mu^T \Sigma^{-1} \mu + b^T L^{-1} b$$

$$h = VU$$

x is Gaussian and integrates to its normalizing constant. The term $C - K$, $K = u^T V u$ factors into the y exponential giving,

$$y^T L^{-1} y - 2y^T L^{-1} b + C - K \text{ and } u^T V u = \left[\Delta \mu + A^T L^{-1} y - A^T L^{-1} b \right]^T \left[\Delta \mu + A^T L^{-1} y - A^T L^{-1} b \right]$$

$$= \mu^T \Delta^{-1} V \Delta \mu + y^T L^{-1} A V A^T L^{-1} y + b^T L^{-1} A V A^T L^{-1} b$$

$$+ 2y^T L^{-1} A V \Delta \mu - 2\mu^T \Delta V A^T L^{-1} b - 2y^T L^{-1} A V A^T L^{-1} b$$

In together, this becomes the following:

$$y^T \left[L^{-1} - L^{-1} A V A^T L^{-1} \right] y - 2y^T \left[L^{-1} b + L^{-1} A V \Delta \mu - L^{-1} A V A^T L^{-1} b \right] + \mu^T \left[\Delta^{-1} - \Delta V \Delta^{-1} \right] \mu + b^T \left[L^{-1} - L^{-1} A V A^T L^{-1} \right] b - 2b^T L^{-1} A V \Delta \mu$$

which simplifies in the following style:

$$y^T R^{-1} y - 2y^T g + b^T (R^{-1} b + \mu)^T \left[\Delta^{-1} - \Delta V \Delta^{-1} \right] \mu - 2b^T L^{-1} A V \Delta^{-1} \mu$$

where we substituted

$$R^{-1} = L^{-1} - L^{-1} A V A^T L^{-1} = \left[L + A \Delta A^T \right]^{-1}$$

$$g = L^{-1} b + L^{-1} A V \Delta^{-1} \mu - L^{-1} A V A^T L^{-1} b$$

Moreover, $Rg = Ah + b$:

$$Rg = R \left[L^{-1}b + L^{-1}AV\Delta^{-1}y - L^{-1}AVATL^{-1}b \right] =$$

$$= R \left[(L^{-1} - L^{-1}AVATL^{-1})b + L^{-1}AV\Delta^{-1}y \right] =$$

$$= R(L^{-1} - L^{-1}AVATL^{-1})b + RL^{-1}AV\Delta^{-1}y$$

Since $R^{-1} = L^{-1} - L^{-1}AVATL^{-1}$ ($\Rightarrow Rg = RL^{-1}AV\Delta^{-1}y + b$)

After cancelling out the b term on both sides, it's sufficient to solve for the other term, equating the desired coefficient:

$$RL^{-1}AV\Delta^{-1}y = [I + A\Delta A^T]L^{-1}AV\Delta^{-1}y = Ah$$

$$= A + A\Delta A^T L^{-1}A = A[\Delta^{-1} + A^T L^{-1}A]$$

(using the definition of $V \rightsquigarrow [ATL^{-1}A + \Delta^{-1}]^{-1}$)

$$I + \Delta A^T L^{-1}A = I + A^T L^{-1}A \Delta A^T = A^T A$$

$$\Delta A^T L^{-1}A = \Delta A^T L^{-1}A \checkmark$$

And, therefore,

$$Rg = Ah + b \text{ with } R^{-1}Ah = g - R^{-1}b = L^{-1}b + L^{-1}AV\Delta^{-1}y -$$

$$- (L^{-1} - L^{-1}AVATL^{-1})b = L^{-1}b + L^{-1}AV\Delta^{-1}y - L^{-1}AVATL^{-1}b - L^{-1}b$$

$$+ L^{-1}AVATL^{-1}b = L^{-1}AV\Delta^{-1}y$$

Now, simplifying the previous equation, we get the following equation for y :

$$y^T R^{-1} y - 2y^T R^{-1} (Ay + b) + P^T R^{-1} b + \mu^T \left[\Delta^{-1} - \Delta^{-1} V \Delta^{-1} \right] \mu$$

$$- 2\mu^T L^{-1} A V \Delta^{-1} \mu$$

Taking into account the identity proved in problem 3, it's easily seen that

$$\Delta^{-1} - \Delta^{-1} I \left(F^{-1} + T^T \Delta^{-1} T \right)^{-1} I^T \Delta^{-1} = (\Delta + T^T F T)^{-1}$$

$$\text{So, } \Delta^{-1} - \Delta^{-1} \left(F^{-1} + T^T \Delta^{-1} T \right)^{-1} T^T \Delta^{-1} = (\Delta + F)^{-1}$$

$$\text{Substitute } F^{-1} = A^T L^{-1} A \Rightarrow V = [F^{-1} + \Delta^{-1}]^{-1} \text{ and}$$

$$(\Delta + F)^{-1} = \Delta^{-1} - \Delta^{-1} V \Delta^{-1} \text{ where } R = L^{-1} - L^T A V A^T L^{-1}$$

$$A^T R^{-1} A = A^T (L^{-1} - L^T A V A^T L^{-1}) A \text{ and } V = [A^T L^{-1} A + \Delta]$$

$$V = [F^{-1} + \Delta^{-1}]^{-1} \Rightarrow \text{Our claim is to prove that}$$

$$A^T R^{-1} A = \Delta^{-1} - \Delta^{-1} V \Delta^{-1} \Rightarrow A^T (L^{-1} - L^T A V A^T L^{-1}) A =$$

$$= A^T L^{-1} A - A^T L^{-1} A V A^T L^{-1} A = \Delta^{-1} - \Delta^{-1} (A^T L^{-1} A + \Delta^{-1}) A$$

$$A^T L^{-1} A - A^T L^{-1} A [A^T L^{-1} A + \Delta^{-1}] A^T L^{-1} A = \Delta^{-1} - \Delta^{-1}$$

$$- \Delta^{-1} (A^T L^{-1} A + \Delta^{-1}) \Delta^{-1} = 0 \quad (\Delta V \Delta^{-1} = 0)$$

$$(A^T L^{-1} A + \Delta^{-1}) (A^T L^{-1} A - A^T L^{-1} A [A^T L^{-1} A + \Delta^{-1}]) = 0$$

$$F^{-1} = A^T L^{-1} A \Rightarrow F^{-1} - F^{-1} (F^{-1} + \Delta^{-1}) F^{-1} = \Delta^{-1} - \Delta^{-1} (F^{-1} + \Delta^{-1})$$

Claim: $A^T R^{-1} A = \Delta^{-1} - \Delta^{-1} V \Delta^{-1}$

Since $F = A^T L^{-1} A$ and $V = [A^T L^{-1} A + \Delta^{-1}] = [F + \Delta^{-1}]^{-1}$

$$A^T R^{-1} A = A^T [L^{-1} - L^{-1} V A V A^T L^{-1}] A = A^T L^{-1} A -$$

$$= A^T [L^{-1} - L^{-1} A [F + \Delta^{-1}]^{-1} A^T L^{-1}] A = A^T L^{-1} A -$$

$$- A^T L^{-1} A [F + \Delta^{-1}] A^T L^{-1} A = F^{-1} - F^{-1} [F + \Delta^{-1}]^{-1} F^{-1}$$

and $\Delta^{-1} - \Delta^{-1} [F + \Delta^{-1}]$ is the second term \Rightarrow

$$\text{Since we proved that } (\Delta + F) = \Delta^{-1} - \Delta^{-1} (F + \Delta^{-1}) \Delta^{-1}$$

(inspired from problem 3) \Rightarrow it's enough to show $F^{-1} - F^{-1} [F + \Delta^{-1}] F^{-1}$

is equal to $(\Delta + F)^{-1} \Rightarrow (\Delta + F) (F^{-1} - F^{-1} [F + \Delta^{-1}]^{-1} F^{-1}) =$

$$= \cancel{\Delta F^{-1} + I - \Delta F^{-1} [F + \Delta^{-1}]^{-1} F^{-1}} - \cancel{[F + \Delta^{-1}]^{-1} F^{-1} I} = I$$

$$[F^{-1} + \Delta^{-1}]^{-1} [\cancel{\Delta F^{-1} - \Delta F^{-1} [F + \Delta^{-1}]^{-1} F^{-1}} - \cancel{[F + \Delta^{-1}]^{-1} F^{-1} F}]$$

$$(\Delta - \Delta F^{-1} [F + \Delta^{-1}]^{-1} - [F^{-1} + \Delta^{-1}]^{-1}) F^{-1} = 0$$

$$\Delta - (\Delta F^{-1} [F + \Delta^{-1}]^{-1} - [F^{-1} + \Delta^{-1}]^{-1}) = 0$$

$$(\Delta - \Delta F^{-1} [F + \Delta^{-1}]^{-1} - [F^{-1} + \Delta^{-1}]^{-1}) (F^{-1} + \Delta^{-1}) = 0$$

$\Delta F^{-1} + (I - \Delta F^{-1} - I) = 0$ which is true. \checkmark In conclusion,
claim is true $\Rightarrow A^T R^{-1} A = \Delta^{-1} - \Delta^{-1} V \Delta^{-1}$ \checkmark

Since we proved $A^T R^{-1} A = \Delta^{-1} - \Delta^{-1} V \Delta^{-1}$ where in the other claim we also had $(\Delta + F)^{-1} = \Delta^{-1} - \Delta^{-1} V \Delta^{-1}$ (using problem 8 result) we conclude $\boxed{A^T R^{-1} A = (\Delta + F)^{-1}}$ \Rightarrow Taking inverse, it yields

$$\begin{aligned}
 A^{-1} R(A^{-1}) &= \Delta + F, \text{ since } R = L + A \Delta A^T \Rightarrow A \Delta A^T = \\
 &= R - L \text{ where } F^{-1} = A^T L^{-1} A, A (A^{-1} R(A^{-1})) A^T = \\
 &= A(\Delta + F) A^T = (AA^{-1}) R(A^{-1}) A^T = R[(A^{-1})^T A^T] = \\
 &= R = (A \Delta A^T + AF A^T), \text{ with } F^{-1} = A^T L^{-1} A \text{ indicates} \\
 &F = A^{-1} L (A^{-1})^T \text{ and } AF A^T = AA^{-1} L (A^{-1})^T A^T = \\
 &= (AA^{-1}) L (A^{-1})^T A^T = L[(A^{-1})^T A^T] = [\Rightarrow AF A^T = L, \text{ so}]
 \end{aligned}$$

$\boxed{R = A \Delta A^T + L}$ So, the key expression is in the red star ★

turn to be $\boxed{\Delta^{-1} - \Delta^{-1} V \Delta^{-1} = A^T R^{-1} A}$ and

$$\begin{aligned}
 R^{-1} A &= (L^{-1} - L^{-1} A V A^T L^{-1}) A = L^{-1} A - L^{-1} A V A^T L^{-1} A = \\
 &= L^{-1} A - L^{-1} A [A^T L^{-1} A + \Delta^{-1}]^{-1} A^T L^{-1} A = L^{-1} A V \Delta^{-1}
 \end{aligned}$$

It's sufficient to show $I - \boxed{[A^T L^{-1} A + \Delta^{-1}]^{-1} A^T L^{-1} A = V \Delta^{-1}}$
 $I - V A^T L^{-1} A = V \Delta^{-1} \Rightarrow I = V (A^T L^{-1} A + \Delta^{-1})$ which is obviously true \checkmark

In the end, we receive $R^{-1}A = L^{-1}AV\Delta^{-1}$. Now, coming back

into the equation red star

$$y^T R^{-1} y - 2y^T R^{-1}(Ay + b) + b^T R^{-1} b + \mu^T [L^{-1} - L^{-1} V \Delta^{-1}] y \\ - 2b^T [L^{-1} A V \Delta^{-1}] \mu$$

if we substitute the new results, we obtain

$$y^T R^{-1} y - 2y^T R^{-1}(Ay + b) + b^T R^{-1} b + \mu^T A^T R^{-1} A y - 2b^T R^{-1} A y$$

which gives the result that y is distributed normally

$$y \sim N(Ay + b, R) \quad \text{Considering } R = A \Delta A^T + L \Rightarrow$$

$$y \sim N(Ay + b, A \Delta A^T + L)$$

Initially, in the ~~red star~~ equation at the beginning of solution problem 4, that long equation given for y can be simplified when conditioning on X since everything not involving X can be summed into the constant of proportionality. This will leave only

$$-2y^T L^{-1} Ax + x^T A^T L^{-1} Ax + 2x^T A^T L^{-1} b + x^T \Delta^{-1} X - 2x^T \Delta^{-1} y$$

$$= -2x^T A^T L^{-1} y + x^T A^T L^{-1} Ax + 2x^T A^T L^{-1} b + x^T \Delta^{-1} X - 2x^T \Delta^{-1} y$$

After this, refactoring this term will yield

$$x^T [A^T L^{-1} A + \Delta^{-1}] X - 2x^T [A^T L^{-1} (y - b) + \Delta^{-1} y] + K$$

where K-constant (specifically, $y^T L^{-1} y - 2y^T L^{-1} b + b^T L^{-1} b + \mu^T \Delta^{-1} \mu$)

If we define $A^T L^{-1} A + \Sigma^{-1} = \Sigma^{-1}$, then it's easily seen that

$$x|y \sim N\left(\Sigma [A^T L^{-1}(y - b) + \Sigma^{-1} \mu], \Sigma\right)$$

Finally, we proved that $p(y) = N(y | A\mu + b, L + A\Sigma A^T)$
 and $\checkmark p(x|y) = N(x | \Sigma [A^T L^{-1}(y - b) + \Sigma^{-1} \mu], \Sigma)$
 where Σ is defined

Sidenote: As an easily understandable logic for proving $p(y)$, we can mention the following ideas as help: (Idea is to just present bigger picture)
 what happens

Let's write the r.v.'s X, y as $X = \mu + \epsilon_x$ and $y = AX + b + \epsilon_y$ where $\epsilon_x \sim N(0, \Sigma)$ and $\epsilon_y \sim N(0, L)$. Plugging in X in the second equation gives us $y = A\mu + b + A\epsilon_x + \epsilon_y$. This is a linear combination of normal distributed random variables and as such itself normal distributed with expectation $A\mu + b$ and covariance matrix $A\Sigma A^T + L$ since $\text{var}(AX) = A \text{var}(X) A^T$. From this, $p(y) = N(y | A\mu + b, A\Sigma A^T + L)$ (precise calculations were indicated in the main solution, this is just main logic). For the 2nd part, $p(x|y) \propto p(y|x)p(x)$ from Bayes, therefore we see $\propto \exp(-(y - Ax - b)^T L^{-1} (y - Ax - b) + (x - \mu)^T \Sigma^{-1} (x - \mu))$. Using rigorous calculations and factoring out x , we obtain proportional to $\exp(-(x - (\Sigma^{-1} + A^T L^{-1} A)^{-1} A^T L^{-1} (y - b))^T (\Sigma^{-1} + A^T L^{-1} A) \cdot (x - (\Sigma^{-1} + A^T L^{-1} A)^{-1} A^T L^{-1} (y - b)))$ which is proportional to given norm.

Combining the logic presented in sidenote as a big picture and how detailed rigorous proofs were presented before that, we successfully proved that $p(x) = N(x | \mu, \Sigma)$ and $p(y|x) = N(y | Ax + b, L) \Rightarrow$

$p(y) = N(y | Ay + b, L + A\Sigma A^T)$ and

$p(x|y) = N(x | \sum [A^T L^{-1} (y - b) + \Sigma^{-1} \mu], \Sigma)$ where

$$\Sigma = (L^{-1} + A^T L^{-1} A)^{-1}$$

Plugging those variables into original problem - where it's explicitly said

$$p(\exists) = N(\exists | 0, I) \Rightarrow \mu = 0, \Sigma = I$$

$$p(x|\exists) = N(x | W\exists + \mu, 6^2 I) \Rightarrow A = W, b = \mu, L = 6^2 I$$

It's evident that $L^{-1} = \frac{1}{6^2} I$ and $\Sigma^{-1} = I \Rightarrow$

$$p(x) = N(x | W \cdot 0 + \mu, 6^2 I + W I W^T) = N(x | \mu, 6^2 I + W W^T)$$

$$p(x) = N(x | \mu, 6^2 I + W W^T)$$

$$p(\exists | x) = N(\exists | \sum [W^T \cdot \frac{1}{6^2} I (y - \mu) + I \cdot 0], \Sigma)$$

$$\text{where } \Sigma = (I + W^T \cdot \frac{1}{6^2} I W) = (I + \frac{1}{6^2} W^T W)^{-1}$$

$$p(\exists | X) = N(\exists | \sum W^T \frac{1}{6^2} (X - \mu), \Sigma)$$

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

6) Using the first part from problem 4, we successfully proved that
 $p(x) = \mathcal{N}(x | \mu, \sigma^2 I + WW^T)$ Let us denote $C = \sigma^2 I + WW^T$

where we consider the determination of the model parameters using maximum likelihood. Given a data set $\mathcal{D} = \{x_i \in \mathbb{R}^d : 1 \leq i \leq N\}$ corresponding log likelihood function is given by

$$\ln p(x | \mu, W, \sigma^2) = \sum_{n=1}^N \ln p(x_n | W, \mu, \sigma^2) \text{ where } X = \{x_n\}$$

clear that

So, using the formula for multivariate Gaussian distribution, it is

$$p(x_i | W, \mu, \sigma^2) = \frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{\det(\text{covariance matrix})}} \exp \left[-\frac{1}{2} (x_i - \mu)^T C^{-1} (x_i - \mu) \right]$$

Since we have $p(x) = \mathcal{N}(x | \mu, C)$ where $C \sim \text{covariance matrix} = \sigma^2 I + WW^T$

$$p(x_i | W, \mu, \sigma^2) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|C|^{1/2}} \exp \left[-\frac{1}{2} (x_i - \mu)^T C^{-1} (x_i - \mu) \right]$$

$$\begin{aligned} \Rightarrow \ln p(x | \mu, W, \sigma^2) &= \sum_{i=1}^N \left[-\frac{d}{2} \ln(2\pi) + \frac{1}{2} \ln |C| - \frac{1}{2} (x_i - \mu)^T C^{-1} (x_i - \mu) \right] \\ &= -\frac{Nd}{2} \ln(2\pi) - \frac{N}{2} \ln |C| - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T C^{-1} (x_i - \mu) \end{aligned}$$

Setting the derivative of the log likelihood w.r.t μ equal to zero gives the expected result μ : maximum likelihood estimator for μ

Our goal is to evaluate $\frac{\partial}{\partial \mu} \ln p(x | \mu, W, \sigma^2)$ in which

the terms $-\frac{Nd}{2} \ln(2\pi)$ and $-\frac{N}{2} \ln |C|$ don't play any role \Rightarrow
 We just need to compute $\frac{\partial}{\partial \mu} \left(-\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T C^{-1} (x_i - \mu) \right)$

$$\text{Assume } \mathbf{x}_i - \boldsymbol{\mu} = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iD} \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_D \end{bmatrix} = \begin{bmatrix} x_{i1} - \mu_1 \\ x_{i2} - \mu_2 \\ \vdots \\ x_{iD} - \mu_D \end{bmatrix}, \quad C = \begin{bmatrix} C_{11} & \dots & C_{1D} \\ \vdots & \ddots & \vdots \\ C_{D1} & \dots & C_{DD} \end{bmatrix}$$

$$(\mathbf{x}_i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = \begin{bmatrix} x_{i1} - \mu_1 & x_{i2} - \mu_2 & \dots & x_{iD} - \mu_D \end{bmatrix} \begin{bmatrix} C_{11} & \dots & C_{1D} \\ \vdots & \ddots & \vdots \\ C_{D1} & \dots & C_{DD} \end{bmatrix} \begin{bmatrix} x_{i1} - \mu_1 \\ x_{i2} - \mu_2 \\ \vdots \\ x_{iD} - \mu_D \end{bmatrix}$$

$$= \sum_{t=1}^D \sum_{k=1}^D (x_{it} - \mu_t) (x_{ik} - \mu_k) C_{kt}, \text{ where } C = G^T T + W W^T \text{ is symmetric as it's covariance matrix}$$

$$C^T = G^T T + W W^T = C \Rightarrow C_{kt} = C_{tk} \quad \text{Taking derivative of } (\mathbf{x}_i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

$$\text{wrt each term } \mu_j \Rightarrow \frac{\partial}{\partial \mu_j} (\mathbf{x}_i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = \frac{\partial}{\partial \mu_j} (x_{ij} - \mu_j)^2 C_{jj}$$

$$+ \frac{\partial}{\partial \mu_j} \left[\sum_{\substack{k=1 \\ k \neq j}}^D (x_{ij} - \mu_j) (x_{ik} - \mu_k) C_{kj} + \sum_{\substack{t=1 \\ t \neq j}}^D (x_{it} - \mu_t) (x_{ij} - \mu_j) C_{it} \right]$$

$$= -2(x_{ij} - \mu_j) C_{jj} + \frac{\partial}{\partial \mu_j} \left[2 \sum_{\substack{k=1 \\ k \neq j}}^D (x_{ik} - \mu_k) (x_{ij} - \mu_j) C_{kj} \right] =$$

$$= -2(x_{ij} - \mu_j) C_{jj} + -2 \sum_{\substack{k=1 \\ k \neq j}}^D (x_{ik} - \mu_k) (C_{kj}) = -2 \sum_{k=1}^D (x_{ik} - \mu_k) C_{kj}$$

$$\text{So, } \frac{\partial}{\partial \mu_j} (\mathbf{x}_i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = -2 \sum_{k=1}^D (x_{ik} - \mu_k) C_{kj} \text{ since } C \text{-symmetric}$$

Fifthly, if we consider $C^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = \begin{bmatrix} C_{11} & \dots & C_{1D} \\ \vdots & \ddots & \vdots \\ C_{D1} & \dots & C_{DD} \end{bmatrix} \begin{bmatrix} x_{i1} - \mu_1 \\ x_{i2} - \mu_2 \\ \vdots \\ x_{iD} - \mu_D \end{bmatrix}$ and thus,

$$C^{-1}(x_i - \mu) = \begin{bmatrix} C_{11}(x_{i1} - \mu_1) + \dots + C_{10}(x_{i0} - \mu_0) \\ \vdots \\ C_{01}(x_{i1} - \mu_1) + \dots + C_{00}(x_{i0} - \mu_0) \end{bmatrix}$$

where the j^{th}
entry is equivalent
to $\sum_{k=1}^0 C_{jk}(x_{ik} - \mu_k)$

which yields that $\frac{\partial}{\partial \mu_j} (x_i - \mu)^T C^{-1}(x_i - \mu) = -2 \sum_{k=1}^0 C_{jk}(x_{ik} - \mu_k)$

and $[C^{-1}(x_i - \mu)]_j = \sum_{k=1}^0 C_{jk}(x_{ik} - \mu_k)$, therefore,

$$\frac{\partial}{\partial \mu} \left(-\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T C^{-1}(x_i - \mu) \right) = \sum_{i=1}^N C^{-1}(x_i - \mu) \text{ due to}$$

$$\frac{\partial}{\partial \mu_j} \left(-\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T C^{-1}(x_i - \mu) \right) = -\frac{1}{2} \sum_{i=1}^N (-2) \cdot \sum_{k=1}^0 C_{jk}(x_{ik} - \mu_k)$$

$$= \sum_{i=1}^N \sum_{k=1}^0 C_{jk}(x_{ik} - \mu_k)$$

where j^{th} entry of $C^{-1}(x_i - \mu)$ is
derivative of log likelihood

equal to $\sum_{k=1}^0 C_{jk}(x_{ik} - \mu_k) \rightsquigarrow$ wrt μ is given by
 $\sum_{i=1}^N C^{-1}(x_i - \mu)$

Setting this derivative to zero $\Rightarrow C^{-1} \left(\sum_{i=1}^N x_i - n\mu \right) = 0$ where

$C = G^2 T + V_r X^T$ with nonzero elements inside entries of matrix C^{-1}

Thus, $\mu = \frac{\sum_{i=1}^N x_i}{n} \rightsquigarrow$ goin for
maximum likelihood estimate

$\hat{\mu}_{ML}$

↳ mean of observed set of data points

If we differentiate $\sum_{i=1}^N C^{-1}(x_i - \mu)$ once again, we obtain $-\sum_{i=1}^N C = -NC^{-1}$

A quadratic form $-y^T C^{-1} y$ has unique maximum at $y=0$, if C^{-1} is p.d.

C^{-1} is clearly p.d. since $C = G^2 I + VV^T \Rightarrow y^T C y = y^T G^2 I y + y^T VV^T y = G^2 y^T y + (V^T y)^T V^T y > 0$ whenever $y \neq 0$

As we know inverse of p.d. matrix is also p.d. $\Rightarrow -y^T C^{-1} y < 0$

except for $y=0$. Therefore, [maximum Likelihood Solution MLE] indicates unique maximum

$$\Rightarrow MLE = \bar{x} \quad \checkmark$$

c) As we showed in the previous problem, $\ln p(X|\mu, V, G^2) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |C| - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T C^{-1} (x_i - \mu)$ where

$C = G^2 I + VV^T$, in order to find the maximum Likelihood estimator of V , it's required to set the derivative of log-likelihood, wrt V as zero:

$$= -N \frac{\partial}{\partial V} \log \det(C) - \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial V} (x_i - \mu)^T C^{-1} (x_i - \mu)$$

where we use the fact $\frac{\partial}{\partial A} (y^T A^{-1} y) = -A^{-1} y y^T A^{-1}$ (from class notes)

$$\frac{\partial}{\partial A} \det(A) = \det(A) (A^{-1})_{ji} \quad \text{Note that } \frac{\partial C}{\partial V} = \frac{\partial}{\partial V} (VV^T)$$

$\therefore \frac{\partial}{\partial A} \det(A) = \det(A) (A^{-1})^T$ The key idea is to take derivatives wrt C , and then chain rule

$\therefore \frac{\partial}{\partial A} \log \det(A) = (A^{-1})^T$ Consequently,

$$\frac{\partial}{\partial X} \ln p(x|\mu, \Sigma, \sigma^2) = -\frac{N}{2} \frac{\partial}{\partial C} \log \det(C) \cdot \frac{\partial C}{\partial X} -$$

$$-\frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial C} \cdot (x_i - \mu)^T C^{-1} (x_i - \mu) \cdot \frac{\partial C}{\partial X}$$

where C -covariance matrix (symmetric)
a symmetric matrix.

Thus, $\frac{\partial}{\partial C} \log \det(C) = \frac{1}{\det(C)} C^{-1} = C^{-1}$ where C^{-1} is

As mentioned, $\frac{\partial}{\partial A} (y^T A^{-1} y) = - (A^{-1})^T y y^T (A^{-1})^T$ where

$$\frac{\partial}{\partial C} ((x_i - \mu)^T C^{-1} (x_i - \mu)) = -(C^{-1})^T (x_i - \mu) (x_i - \mu)^T (C^{-1})^T$$

$$(C^{-1})^T = (C^T)^{-1} = C^{-1} \Rightarrow \frac{\partial}{\partial X} \ln p(x|\mu, \Sigma, \sigma^2) = -\frac{N}{2} C \frac{\partial C}{\partial X} =$$

$$-\frac{1}{2} \sum_{i=1}^N -C^{-1} (x_i - \mu) (x_i - \mu)^T C \cdot \frac{\partial C}{\partial X}, \text{ where } C = \sigma^2 I + X X^T$$

It's easily seen that $\boxed{\frac{\partial C}{\partial X} = 2X}$ and if we set derivative into

$$0 = -\frac{N}{2} C \cdot 2X + \frac{1}{2} \sum_{i=1}^N C^{-1} (x_i - \mu) (x_i - \mu)^T C^{-1} \cdot 2X$$

$$NC^{-1}X = \sum_{i=1}^N C^{-1} (x_i - \mu) (x_i - \mu)^T C^{-1} X \text{ where } C = \sigma^2 I + X X^T$$

$$\sum_{i=1}^N (x_i - \mu) (x_i - \mu)^T = N\delta, \quad \sum_{i=1}^N C^{-1} (x_i - \mu) (x_i - \mu)^T = C^{-1} N\delta = NC^{-1}\delta \text{ where}$$

$$NC^{-1}X = NC^{-1}\delta C^{-1}X, \text{ since } N\text{-scalar} \Rightarrow \boxed{C^{-1}X = C^{-1}\delta C^{-1}X}$$

Considering C -covariance matrix (positive semi-definite matrix), it turns out to be symmetric and C^{-1} is symmetric where $\det(C^{-1})$ is nonzero \Rightarrow we can eliminate C^{-1} from $C^{-1}(V - \delta C^{-1}W) = 0$ and obtain $V = \delta C^{-1}W = \delta(G^2 I + VV^T)^{-1}W$. In conclusion

$$\boxed{\delta(G^2 I + VV^T)^{-1}W = V} \quad \text{as desired } \boxed{V} \bullet \boxed{F}$$

2) It's important to define some crucial definitions and characteristics for Latent class analysis. We assume that data points are generated by different Gaussian distributions as $p(x_i) = \sum_{k=1}^K \pi_k p(x_i | \mu_k)$

where data points are $\{x_1, x_2, \dots, x_N\}$ and

$\sum_{k=1}^K \pi_k = 1$ is satisfied. After that, we define a latent random

variable $f_i = (f_{i1}, f_{i2}, \dots, f_{iK})$ with $f_{ik} \in \{0, 1\}$, $\sum_{k=1}^K f_{ik} = 1$ and

$p(f_{ik}=1) = \pi_k$, with $1 \leq k \leq K$. In alternative manner, we obtain $p(f_i) = \prod_{k=1}^K \pi_k^{f_{ik}} (\pi_1, \pi_2, \dots, \text{or } \pi_K)$

We also notice that $p(x_i | f_{ik}=1) = p(x_i | \mu_k)$, and

$p(x_i | f_i) = \prod_{k=1}^K p(x_i | \mu_k, \pi_k)^{f_{ik}}$

$p(x_i, f_i) = p(x_i | f_i) p(f_i) = \prod_{k=1}^K \pi_k^{f_{ik}} p(x_i | \mu_k, \pi_k)^{f_{ik}}$

So, $p(x_i) = \sum_{f_i} p(x_i, f_i) = \sum_{f_i} \pi_k p(x_i | \mu_k)$

($f_i = (f_{i1}, f_{i2}, \dots, f_{iK})$ independently generated)

Now, we consider the joint probability of all data points

$x = \{x_1, x_2, \dots, x_N\}$, $f = \{f_1, f_2, \dots, f_N\} \Rightarrow p(x, f) = \prod_{i=1}^N p(x_i | f_i) p(f_i)$

$= \prod_{i=1}^N \prod_{k=1}^K \pi_k^{f_{ik}} p(x_i | \mu_k)^{f_{ik}}$ where variables $f_i \sim$ variables

$\pi \sim$ assignment probabilities. It's also important to get $p(f_{ik}=1 | x_i)$

Also we have to notice that $p(\mathbb{I}_{ik=1} | x_i) = \frac{p(\mathbb{I}_{ik=1}, x_i)}{p(x_i)} =$

$$= \frac{p(\mathbb{I}_{ik=1}, x_i)}{\sum_{p=1}^K p(x_i, \mathbb{I}_{ip=1})} = \frac{p(x_i | \mathbb{I}_{ik=1}) p(\mathbb{I}_{ik=1})}{\sum_{p=1}^K p(x_i | \mathbb{I}_{ip=1}) p(\mathbb{I}_{ip=1})}$$

$$= \frac{\pi_k p(x_i | \mu_k)}{\sum_{p=1}^K \pi_p p(x_i | \mu_p)}$$

using the current parameters,
we compute the E-step ✓

Regarding the M-step, we have to update parameters μ_k and π_k

Considering $\sum_{k=1}^K \pi_k = 1$ and Lagrange multiplier λ , we formulate

$$\mathcal{L}(\theta) := \log p(x | \theta) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) = J(\theta) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \pi} = 0 \Rightarrow \sum_{i=1}^N \frac{p(x_i | \mu_k)}{\sum_{p=1}^K \pi_p p(x_i | \mu_p)} + \lambda = 0, \quad 1 \leq k \leq K$$

$$\text{In the end, we get } \frac{1}{\pi_k} \sum_{i=1}^N \frac{\pi_k p(x_i | \mu_k)}{\sum_{p=1}^K \pi_p p(x_i | \mu_p)} + \lambda = 0$$

$$\text{Recall that } p(\mathbb{I}_{ik=1} | x_i) = \frac{\pi_k p(x_i | \mu_k)}{\sum_{p=1}^K \pi_p p(x_i | \mu_p)} \text{ holds true}$$

Therefore, it follows that $\bar{\pi}_k = \frac{1}{\lambda} \sum_{i=1}^N p(z_{ik}=1 | x_i)$

Considering that $\sum_k \bar{\pi}_k = 1$, λ satisfies the following relation

$$\lambda = - \sum_{k=1}^K \sum_{i=1}^N p(z_{ik}=1 | x_i) = - \sum_{i=1}^N \sum_{k=1}^K p(z_{ik}=1 | x_i) = -N$$

In conclusion, we obtain $\bar{\pi}_k = \frac{1}{N} \sum_{i=1}^N p(z_{ik}=1 | x_i)$

updating the parameters

for $\bar{\pi}_k$ during M-step

techniques applied in problem 1

For updating the μ_k parameters, we'll use the terminologies and
(defining responsibilities $\delta(z_{ik})$, and similar approaches)

Inertia for K-means clustering with all four features is 97.225

