

2) a) Taking into account that Beta distribution is normalized, it implies $\left(\int_0^1 \text{Beta}(\mu|\alpha, \beta) d\mu\right) = 1$ where

$$\text{Beta}(\mu|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1-\mu)^{\beta-1}. \text{ We'll prove the}$$

normalization of Dirichlet distribution using induction principle.

Since Beta distribution is a special case of Dirichlet for $M=2$, and considering Beta distribution is normalized, we can assume that Dirichlet distribution is normalized for $(M-1)$ variables, and will guarantee to show that it is also normalized for M variables.

Consider Dirichlet distribution over M variables. As we have

$$\sum_{k=1}^M \mu_k = 1 \Rightarrow \text{We write } p_M(\mu_1, \dots, \mu_{M-1}) = D_M \prod_{k=1}^{M-1} d\mu_k \left(1 - \sum_{j=1}^{M-1} \mu_j\right)^{dM-1}$$

for Dirichlet

When integrating over μ_{M-1} , lower limit of integration = 0, upper limit of integration = $1 - \sum_{j=1}^{M-2} \mu_j$. Therefore, we get

$$p_{M-1}(\mu_1, \dots, \mu_{M-2}) = \int_0^{1 - \sum_{j=1}^{M-2} \mu_j} p_M(\mu_1, \dots, \mu_{M-1}) d\mu_{M-1} = \\ = C_M \prod_{k=1}^{M-2} d\mu_k \int_0^{1 - \sum_{j=1}^{M-2} \mu_j} \mu_{M-1}^{dM-1-1} \left(1 - \sum_{j=1}^{M-1} \mu_j\right)^{dM-1} d\mu_{M-1}$$

changing integration variables from μ_{M-1} to t where

$$t = \frac{\mu_{M-1}}{1 - \sum_{j=1}^{M-2} \mu_j} \Rightarrow \mu_{M-1} = t(1 - \sum_{j=1}^{M-2} \mu_j) \text{ and}$$

$$1 - \sum_{j=1}^{M-2} \mu_j \quad (\mu_{M-1} \in [0, 1 - \sum_{j=1}^{M-2} \mu_j] \Rightarrow t \in [0, 1])$$

$$P_{M-1}(\mu_1, \dots, \mu_{M-2}) = D_M \prod_{k=1}^{M-2} \mu_k^{\alpha_{k-1}} \left(\frac{1}{1 - \sum_{j=1}^{M-2} \mu_j} \right)^{\alpha_{M-1} + \alpha_{M-2}}$$

$\int_0^1 t^{\alpha_{M-1}-1} (1-t)^{\alpha_{M-2}} dt$. Since Beta distribution is normalized,

$$\int_0^1 t^{\alpha_{M-1}-1} (1-t)^{\alpha_{M-2}} dt = \frac{\Gamma(\alpha_1) \Gamma(\alpha_2)}{\Gamma(\alpha_1 + \alpha_2)} =$$

$$= \frac{\Gamma(\alpha_{M-1}) \Gamma(\alpha_M)}{\Gamma(\alpha_{M-1} + \alpha_M)}$$

$$= D_M \prod_{k=1}^{M-2} \mu_k^{\alpha_{k-1}} \left(\frac{1}{1 - \sum_{j=1}^{M-2} \mu_j} \right)^{\alpha_{M-1} + \alpha_{M-2}} \cdot \frac{\Gamma(\alpha_{M-1}) \Gamma(\alpha_M)}{\Gamma(\alpha_{M-1} + \alpha_M)}$$

RHS is normalized Dirichlet distribution of $M-1$ variables,

where coefficients $\sim d_1, \dots, d_{M-2}, d_{M-1} + d_M$ (integrating over past groups)

$$\text{and from the formula, } \text{Dir}(\mu | \alpha) = \frac{\Gamma(d_0)}{\Gamma(d_1) \dots \Gamma(d_K)} \prod_{k=1}^K \mu_k^{d_{k-1}}$$

$$\text{where } d_0 = \sum_{k=1}^K d_k$$

By determining normalization coefficient from previous formula and \star , we obtain

$$\mathcal{D}_M = \frac{\Gamma(\alpha_1 + \dots + \alpha_M)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_{M-2}) \Gamma(\alpha_{M-1} + \alpha_M)} \cdot \frac{\Gamma(\alpha_{M-1} + \alpha_M)}{\Gamma(\alpha_{M-1}) \Gamma(\alpha_M)} =$$

$$= \frac{\Gamma(\alpha_1 + \dots + \alpha_M)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_{M-2}) \Gamma(\alpha_{M-1}) \Gamma(\alpha_M)}$$

as desired ✓

$\Rightarrow \text{Dir}(\alpha) \text{ is normalized}$

b) Based on the definition of Expectation and formula

$$\text{Dir}(\mu | \alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \alpha_k^{d_k-1} \quad \text{where } \alpha_0 = \sum_{k=1}^K \alpha_k \Rightarrow$$

$$E[\mu_j] = \int \mu_j \text{Dir}(\mu | \alpha) d\mu = \int \mu_j \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \alpha_k^{d_k-1} d\mu =$$

$$= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \int \mu_j \prod_{k=1}^K \alpha_k^{d_k-1} d\mu = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \cdot \frac{\Gamma(\alpha_1) \dots \Gamma(\alpha_{j-1})}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \cdot$$

$$\cdot \frac{\Gamma(\alpha_{j+1}) \dots \Gamma(\alpha_K)}{\Gamma(\alpha_{j+1}) \dots \Gamma(\alpha_K)} \cdot \frac{\Gamma(\alpha_j+1)}{\Gamma(\alpha_j+1)}$$

$$= \frac{\Gamma(\alpha_0) \Gamma(\alpha_j+1)}{\Gamma(\alpha_j) \Gamma(\alpha_0+1)} = \frac{\Gamma(\alpha_0) \alpha_j \Gamma(\alpha_j)}{\Gamma(\alpha_j) \alpha_0 \Gamma(\alpha_0)} = \frac{\alpha_j}{\alpha_0} \Rightarrow E[\mu_j] = \frac{\alpha_j}{\alpha_0}$$

For variance, we calculate $E[\mu_j^2] \Rightarrow E[\mu_j^2] = \int \mu_j^2 \text{Dir}(\mu | \alpha) d\mu$

$$= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \Gamma(\alpha_2) \dots \Gamma(\alpha_K)} \int \mu_j^2 \prod_{k=1}^K \alpha_k^{d_k-1} d\mu =$$

$$= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \cdot \frac{\Gamma(\alpha_1) \dots \Gamma(\alpha_{j-1}) \Gamma(\alpha_j+2) \Gamma(\alpha_{j+1}) \dots \Gamma(\alpha_K)}{\Gamma(\alpha_0+2) \Gamma(\alpha_0+1) \dots \Gamma(\alpha_0+1)}$$

$$= \frac{\Gamma(\alpha_0) \Gamma(\alpha_j + 2)}{\Gamma(\alpha_j) \Gamma(\alpha_0 + 2)} = \frac{\Gamma(\alpha_0)(\alpha_j + 1)\alpha_j \Gamma(\alpha_j)}{\Gamma(\alpha_j)(\alpha_0 + 1)\alpha_0 \Gamma(\alpha_0)} = \frac{\alpha_j(\alpha_j + 1)}{\alpha_0(\alpha_0 + 1)}$$

$$E[\mu_j^2] = \frac{\alpha_j(\alpha_j + 1)}{\alpha_0(\alpha_0 + 1)} \text{ and from previous result, } E[\mu_j] = \frac{\alpha_j}{\alpha_0}$$

$$\Rightarrow \text{Var}[\mu_j] = E[\mu_j^2] - E[\mu_j]^2 = \frac{\alpha_j(\alpha_j + 1)}{\alpha_0(\alpha_0 + 1)} - \frac{\alpha_j^2}{\alpha_0^2} = \frac{\alpha_j^2}{\alpha_0^2}$$

$$= \frac{\alpha_j(\alpha_j + 1)\alpha_0 - \alpha_j^2(\alpha_0 + 1)}{\alpha_0^2(\alpha_0 + 1)} = \frac{\alpha_j\alpha_0 - \alpha_j^2}{\alpha_0^2(\alpha_0 + 1)} \quad \boxed{\text{Var}[\mu_j] = \frac{\alpha_j(\alpha_0 - \alpha_j)}{\alpha_0^2(\alpha_0 + 1)}}$$

Since $\text{Cov}(x_i, x_j) = E[(x_i - E[x_i])(x_j - E[x_j])] = E[x_i x_j] - E[x_i]E[x_j] = E[x_i x_j] - \frac{\alpha_i}{\alpha_0} \frac{\alpha_j}{\alpha_0}$ from previous results

$$E[x_i x_j] = \int_{\Delta^{K-1}} x_i x_j \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} x_1^{\alpha_1-1} x_K^{\alpha_K-1} dX =$$

$$= \int_{\Delta^{K-1}} \frac{\Gamma(\alpha_0+2)}{\alpha_0(\alpha_0+1)} = (1+\frac{1}{\alpha_0}) \frac{\alpha_0+1}{\alpha_0} = (1+\frac{1}{\alpha_0}) x_1^{\alpha_1-1} x_2^{\alpha_2-1} \dots x_K^{\alpha_K-1} dX$$

$$= \frac{\alpha_i \alpha_j}{\alpha_0(\alpha_0+1)} \int_{\Delta^{K-1}} \frac{\Gamma(\alpha_0+2)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_j+1) \dots \Gamma(\alpha_i+1) \dots \Gamma(\alpha_K)} x_1^{\alpha_1-1} x_j^{\alpha_j-1} \dots x_i^{\alpha_i-1} x_K^{\alpha_K-1} dX$$

$$= \frac{\alpha_i \alpha_j}{\alpha_0(\alpha_0+1)} \text{ since Dirichlet distribution is normalized, where } \Gamma(\alpha_0+2) = \alpha_0(\alpha_0+1) \Gamma(\alpha_0) \text{ and } \Gamma(\alpha_i+1) = \alpha_i \Gamma(\alpha_i)$$

$$\frac{\partial}{\partial \mu} \left(-\frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) \right) = -\frac{1}{2} (x_j - \mu)^T (\Sigma^{-1} + (\Sigma^{-1})^T) \cdot \frac{\partial}{\partial \mu} (x_j - \mu) =$$

$$-\frac{1}{2} (x_j - \mu)^T (\Sigma^{-1} + (\Sigma^{-1})^T) = \frac{1}{2} (x_j - \mu)^T (2\Sigma^{-1}) =$$

$= (x_j - \mu)^T \Sigma^{-1}$, considering the fact $\Sigma \sim$ positive definite
such that $\Sigma^{-1} = ((\Sigma^{-1})^T)$ and independent of μ
 $(\Sigma \text{-positive definite} \Rightarrow \Sigma^{-1} \text{-positive definite})$

Since $\frac{\partial \log(l)}{\partial \mu} = 0$ for finding MLE of $\mu \Rightarrow$

$$\sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} = 0, \text{ and } x_j - \mu \sim \text{vector column} \Rightarrow$$

$$\sum_{j=1}^n (x_j - \mu) = 0, \quad \boxed{\text{MLE} = \bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}} \quad (1)$$

b) It's well-known that posterior distribution
 $P(\mu | X) \propto P(X|\mu) P(\mu)$ or equivalently,

$$N(\mu | X, \Sigma) \propto N(X|\mu, \Sigma) N(\mu | \mu_0, \Sigma_0)$$

where considering $N(X|\mu, \Sigma) = \prod_{j=1}^n N(x_j | \mu, \Sigma)$

$\Rightarrow N(X|\mu, \Sigma) N(\mu | \mu_0, \Sigma_0)$ becomes identical to

$$\text{Thus, } \left[E[x_i x_j] = \frac{\alpha_i \alpha_j}{\alpha_0(\alpha_0+1)} \right] \Rightarrow \text{Cor}(x_i, x_j) = E[x_i x_j] - \frac{\alpha_i \alpha_j}{\alpha_0^2} =$$

$$= \frac{\alpha_i \alpha_j}{\alpha_0(\alpha_0+1)} - \frac{\alpha_i \alpha_j}{\alpha_0^2} = \frac{\alpha_i \alpha_j (\alpha_0 - \alpha_0 - 1)}{\alpha_0^2(\alpha_0+1)} = \frac{-\alpha_i \alpha_j}{\alpha_0^2(\alpha_0+1)} \Rightarrow \boxed{\text{Cor}(x_i, x_j) = -\frac{\alpha_i \alpha_j}{\alpha_0^2(\alpha_0+1)}}$$

3) a) We consider the likelihood function where product of n random variables X_1, X_2, \dots, X_n is taken

$$\mathcal{L}(X_1, X_2, \dots, X_n | \mu) = N(X_1 | \mu, \Sigma) N(X_2 | \mu, \Sigma) \dots N(X_n | \mu, \Sigma)$$

$$N(X_n | \mu, \Sigma) = \left(\frac{1}{(2\pi)^{k/2}} \right)^n \cdot \left(\frac{1}{|\Sigma|^{1/2}} \right)^n \cdot \exp \left(-\frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) \right)$$

To compute MLE, we use $\text{Pog}(\mathcal{L}(X_1, \dots, X_n | \mu)) \Rightarrow$

$$\text{Pog}(\mu) = -\frac{nk}{2} \text{Pog}(2\pi) - \frac{n}{2} \text{Pog}(|\Sigma|) - \frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)$$

Since Σ and Π do not depend on μ , we'll take derivative of $\text{Pog}(\mu)$ wrt μ to obtain μ_{MLE} . In order to do that,

observe from lecture notes that $\frac{\partial}{\partial \mu} (x^T A x) = x^T (A + A^T)$

where $A \sim$ independent of X

$$\text{Thus, } \frac{\partial}{\partial \mu} \left(-\frac{nk}{2} \text{Pog}(2\pi) \right) = \frac{\partial}{\partial \mu} \left(-\frac{n}{2} \text{Pog}(|\Sigma|) \right) = 0 \text{ and}$$

$$\text{as } \frac{\partial}{\partial \mu} (x^T A x) = x^T (A + A^T) \frac{\partial x}{\partial \mu} \text{ (where } A \sim \text{indep. of } \mu \text{ and } x \sim \text{dependent of } \mu \text{)}$$

$$\left(\frac{1}{(2\pi)^{k/2}} \right)^{n+1} \cdot \left(\frac{1}{|\Sigma|^{1/2}} \right)^{n+1} \cdot \left(\frac{1}{|\Sigma_0|^{1/2}} \right)^{n+1} \cdot \exp \left(-\frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) \right)$$

From lecture notes, we know $\mu_{MAP} = \arg \max_{\mu} P(\mu | X)$
and using the same derivative/Pop strategy from past problem

$$\text{Pop} (P(\mu | X)) = \text{Pop} (P(X|\mu) P(\mu)) = -\frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)$$

$$-\frac{1}{2} (\mu - \mu_0)^T \Sigma_0^{-1} (\mu - \mu_0) - \frac{n+1}{2} \text{Pop} (|\Sigma_0|) - \frac{n+1}{2} \text{Pop} (|\Sigma|)$$

$$-\frac{(n+1)k}{2} \text{Pop} 2\pi \Rightarrow \partial (\text{Pop} \cdot P(\mu | X)) = 0 \text{ should satisfy}$$

for obtaining MAP estimator of μ

$$= 0 \text{ where } \partial \left(-\frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu) \right) = \begin{aligned} & -\frac{n+1}{2} \text{Pop} |\Sigma_0| - \\ & -\frac{n+1}{2} \text{Pop} (|\Sigma|) - \\ & -\frac{(n+1)k}{2} \text{Pop} 2\pi \end{aligned}$$

$$= \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} \quad \text{had been proved in previous section}$$

$$\text{along with } \partial \left(-\frac{1}{2} (\mu - \mu_0)^T \Sigma_0^{-1} (\mu - \mu_0) \right) = (\mu - \mu_0)^T (-I) \Sigma^{-1}$$

$$\Rightarrow (A + B)^T = A^T + B^T$$

Using this formula, becomes true for LH8 operations,

$$\sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} = (\mu - \mu_0)^T \Sigma_0^{-1}$$

$$\sum_{j=1}^n (x_j^T - \mu^T) \Sigma^{-1} = \left(\left(\sum_{j=1}^n x_j \right)^T - n\mu^T \right) \Sigma^{-1} = \left(\sum_{j=1}^n x_j \right)^T \Sigma - n\mu^T \Sigma$$

$$(\mu - \mu_0)^T \Sigma_0^{-1} = \mu^T \Sigma_0^{-1} - \mu_0^T \Sigma_0^{-1} = (\sum_{j=1}^n x_j)^T \Sigma^{-1} - n \mu^T \Sigma^{-1} =$$

$$= \mu^T \Sigma_0^{-1} - \mu_0^T \Sigma_0^{-1} = \mu^T (\Sigma_0^{-1} + n \Sigma^{-1}) = (\sum_{j=1}^n x_j)^T \Sigma^{-1} + \mu_0^T \Sigma_0^{-1}$$

Considering $n \Sigma^{-1} + \Sigma_0^{-1}$ is invertible, we find that

$$\mu^T = \left(\sum_{j=1}^n x_j^T \Sigma^{-1} + (\mu_0^T \Sigma_0^{-1}) \right) (n \Sigma^{-1} + \Sigma_0^{-1})^{-1} \text{ where } (AB)^{-1} = B^T A^{-1} \Rightarrow$$

$$\mu_{MAP} = \left((n \Sigma^{-1} + \Sigma_0^{-1})^{-1} \right)^T \left((\mu_0 - \sum_{j=1}^n x_j)^T \Sigma^{-1} + \mu_0^T \Sigma_0^{-1} \right)$$

		True Label		Truth Table (confusion matrix)
		Positive	Negative	
Prediction	Positive	TP	FP	score > x \Rightarrow pred is positive
	Negative	FN	TN	score \leq x \Rightarrow pred is negative

a) $X=15 \Rightarrow TP=8$ and $FP=5$
 $FN=2$ and $TN=35$

Sensitivity $\approx \frac{8}{10} = \frac{4}{5}$ and Specificity $\approx \frac{35}{40} = \frac{7}{8}$

b) Precision $\approx \frac{TP}{TP+FP}$ and Recall $\approx \frac{TP}{TP+FN}$

$$F_1 \approx \frac{2TP}{2TP+FP+FN} = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}$$

$$X=5 \Rightarrow TP=10 \quad FP=84 \quad FN=0 \quad TN=6 \quad F_1 = \frac{20}{84} = \frac{10}{42} \quad \checkmark$$

$$\begin{array}{l} \underline{X=10 \Rightarrow TP=8, FN=1, FP=14, TN=26} \\ \quad \left| \begin{array}{l} F_1 = \frac{18}{33} = \frac{6}{11} \end{array} \right. \end{array}$$

$$\underline{X=15 \Rightarrow TP=8, FN=2, FP=5, TN=35} \quad \left| \begin{array}{l} F_1 = \frac{16}{23} \end{array} \right. \quad \checkmark$$

$$\underline{X=20 \Rightarrow TP=5, FN=5, FP=1, TN=16} \quad \left| \begin{array}{l} F_1 = \frac{10}{16} = \frac{5}{8} \end{array} \right. \quad \checkmark$$

Finding the best value for X means

to choose highest precision and recall \rightarrow highest F_1 -score

$$\frac{16}{28} > \frac{10}{27}, \quad \frac{16}{28} > \frac{6}{11}, \quad \frac{16}{28} > \frac{5}{8} \Rightarrow \text{highest } F_1\text{-score } \approx \frac{16}{23}$$

and thus, it is achieved at $\boxed{X=15} \star \checkmark \oplus$

$$3) A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}, \quad B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \dots & b_{mn} \end{bmatrix}$$

$$a) (AB)_{ij} = \sum_{k=1}^n a_{ik} b_{kj} \text{ and } \text{tr}(A|B) = \sum_{i=1}^n (AB)_{ii} = \sum_{i=1}^n a_{ii} b_{ii}$$

$$= \sum_{i=1}^n \sum_{k=1}^n a_{ik} b_{ki} \Rightarrow \text{According to definition in the problem}$$

$$y \in \mathbb{R}, \quad \mathbf{x} \in \mathbb{R}^{p \times n} \Rightarrow \frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix}^T \quad \text{tr}(AB) \in \mathbb{R}$$

We use numerator-Payout notation,

$$\frac{\partial \text{tr}(AB)}{\partial A_{ij}} = \begin{bmatrix} \frac{\partial \text{tr}(AB)}{\partial a_{11}} & \frac{\partial \text{tr}(AB)}{\partial a_{12}} & \dots & \frac{\partial \text{tr}(AB)}{\partial a_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \text{tr}(AB)}{\partial a_{m1}} & \frac{\partial \text{tr}(AB)}{\partial a_{m2}} & \dots & \frac{\partial \text{tr}(AB)}{\partial a_{mn}} \end{bmatrix}$$

and using the $\text{tr}(AB)$ value from previous

$$\left[\frac{\partial \operatorname{tr}(AB)}{\partial A} \right]_{ij} = \frac{\partial \operatorname{tr}(AB)}{\partial a_{ji}} = \frac{\partial}{\partial a_{ji}} \left(\sum_{i=1}^n \sum_{k=1}^n a_{ik} b_{ki} \right) = \sum_{k=1}^n b_{ki}$$

In the summation $\sum_{i=1}^n \sum_{k=1}^n a_{ik} b_{ki}$, term containing a_{ji} coefficient

(dot product of j^{th} row of A and i^{th} column of B) will only be $a_{ji} \times b_{ij} = \frac{\partial}{\partial a_{ji}} \left(\sum_{i=1}^n \sum_{k=1}^n a_{ik} b_{ki} \right) = \frac{\partial}{\partial a_{ji}} (\alpha_{ji} \times \beta_{ij}) = \frac{\partial \operatorname{tr}(AB)}{\partial A}$ according to indices rule.

$$= \beta_{ij}, \text{ thus } \left[\frac{\partial \operatorname{tr}(AB)}{\partial A} \right]_{ij} = \beta_{ij} \Rightarrow \boxed{\frac{\partial \operatorname{tr}(AB)}{\partial A} = \beta}$$

$$\text{where } \beta = I \Rightarrow \boxed{\frac{\partial \operatorname{tr}(AB)}{\partial A} = I} = \boxed{\frac{\partial \operatorname{tr}(A)}{\partial A} = I}$$

$$\text{c) } \Rightarrow \text{We have to prove } \frac{\partial}{\partial A} \log |A| = A^{-1} \text{ where } |A| > 0$$

$$\text{According to chain rule, } \frac{\partial}{\partial A} \log |A| = \frac{1}{|A|} \cdot \frac{\partial |A|}{\partial A}$$

$$\text{It's a well-known fact that } A^{-1} = \frac{1}{|A|} \cdot [\operatorname{adj}(A)]. \text{ Thus,}$$

$$\text{We have to prove } \frac{\partial |A|}{\partial A} = \operatorname{adj}(A)$$

According to the definition, if we denote C_{ij} by $(-1)^{i+j} \det \begin{pmatrix} \text{matrix after removal of } i^{\text{th}} \text{ row} \\ \text{and } j^{\text{th}} \text{ column} \end{pmatrix}$, then adjugate is the transpose of its cofactor matrix.

$$C = \begin{bmatrix} C_{11} & C_{12} & \dots & C_{1n} \\ C_{21} & C_{22} & \dots & C_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ C_{n1} & C_{n2} & \dots & C_{nn} \end{bmatrix} \Rightarrow \text{adj}(A) = \begin{bmatrix} C_{11} & C_{21} & \dots & C_{n1} \\ C_{12} & C_{22} & \dots & C_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ C_{1n} & C_{2n} & \dots & C_{nn} \end{bmatrix}$$

It's also well-known rule that $\det(A) = |A| = \sum_{k=1}^n a_{kk} C_{kk}$ along the k^{th} row. According to the definition of cofactor matrix, neither any of a_{ik} nor C_{jk} are dependent on a_{ii} , where the only exception is a_{ji} for $k=i$.
 $|A| = \sum_{j=1}^n a_{ij} C_{ij}$ and $\frac{\partial |A|}{\partial a_{ij}} = \frac{\partial}{\partial a_{ij}} (a_{ij} C_{ij}) = C_{ij}$
 $[\text{adj}(A)]_{ij} = C_{ji}$ from the construction ($\text{adj}(A)$ is C_{ij} i^{th} row, j^{th} column of A)

Taking into account numerator-Payout notation \Rightarrow
 $\left[\frac{\partial |A|}{\partial A} \right]_{ij} = \frac{\partial |A|}{\partial a_{ij}}$, since $\frac{\partial |A|}{\partial A} = \begin{bmatrix} \frac{\partial |A|}{\partial a_{11}} & \dots & \frac{\partial |A|}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial |A|}{\partial a_{n1}} & \dots & \frac{\partial |A|}{\partial a_{nn}} \end{bmatrix}$
 As we showed $\frac{\partial |A|}{\partial a_{ij}} = C_{ji} = \left[\frac{\partial |A|}{\partial A} \right]_{ij} = C_{ji}$ (i^{th} row, j^{th} col of $\frac{\partial |A|}{\partial A}$ is C_{ji})

From the definition, $[\text{adj}(A)]_{ij} = C_{ji}$, concluding that

$$\boxed{\frac{\partial |A|}{\partial A} = \text{adj}(A)} \Rightarrow \text{Therefore, claim is satisfied and so, } \boxed{\frac{\partial |A|}{\partial A} = A^{-1}, \text{ where } \det(A) = |A| \neq 0}$$

4) a) From famous Jensen's Inequality, it is evident that
for real convex function φ , $\varphi(a_1x_1 + \dots + a_m x_m) \leq a_1\varphi(x_1) + \dots + a_m\varphi(x_m)$ where $a_1 + \dots + a_m = 1$ and $a_i > 0$

$$\text{From the definition, } H(x) = -E[\log P(x)] = E\left[\log\left(\frac{1}{P(x)}\right)\right] \\ = \sum_{i=1}^m p_i \log\left(\frac{1}{p_i}\right) = -\sum_{i=1}^m p_i \log p_i = \sum_{i=1}^m p_i (-\log p_i)$$

Take $\varphi(x) = \log\left(\frac{1}{x}\right)$, where $\varphi'(x) = -\frac{1}{x^2} < 0$, $\varphi''(x) = \frac{2}{x^3} > 0 \Rightarrow$

$$\sum_{i=1}^m p_i \log\left(\frac{1}{p_i}\right) = \sum_{i=1}^m p_i \varphi(p_i) \geq \varphi\left(\sum_{i=1}^m p_i\right) = \log\left(\frac{1}{\sum_{i=1}^m p_i}\right)$$

$$\sum_{i=1}^m p_i \log(p_i) = \sum_{i=1}^m p_i \log\left(\frac{1}{1/p_i}\right) = \sum_{i=1}^m p_i \varphi\left(\frac{1}{p_i}\right) \geq$$

$$\geq \varphi\left(1 + 1 + \dots + 1\right) = \varphi(m) = \log\left(\frac{1}{m}\right) = -\log m \Rightarrow \text{that}$$

$$\sum_{i=1}^m p_i \log(p_i) \geq -\log m \Rightarrow -\sum_{i=1}^m p_i \log(p_i) \leq \log m \text{ or j48}$$

$$\log m \geq \sum_{i=1}^m p_i (-\log p_i) = \sum_{i=1}^m p_i \log\left(\frac{1}{p_i}\right) = -E[\log P(x)] = H(x)$$

In conclusion, $H(x) \leq \log m \quad \boxed{\checkmark}$

Note: Discrete p_i 's are greater than 0 with
 $P(X=x_i) = p_i$ for $1 \leq i \leq m$

B) From previously solved problem 3, it's easy to observe that univariate Gaussian distributions are

$$N(m, \sigma^2) = \frac{1}{(2\pi)^{1/2}} \cdot \frac{1}{|\sigma|} \exp\left(-\frac{1}{2} \frac{(x-m)^2}{\sigma^2}\right)$$

$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi)^{1/2}} \cdot \frac{1}{|\sigma|} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$$

Therefore, $D_{KL}(N(\mu, \sigma^2) || N(m, \sigma^2)) =$

$$= E\left[\log\left(\frac{N(x|\mu, \sigma^2)}{N(x|m, \sigma^2)}\right)\right] = E\left[\log(N(x|\mu, \sigma^2))\right]$$

$$- E\left[\log(N(x|m, \sigma^2))\right] = E\left[\log\left(\frac{|\sigma|}{|m|} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} - \frac{(x-m)^2}{\sigma^2}\right)\right)\right]$$

$$D_{KL}(N(\mu, \sigma^2) || N(m, \sigma^2)) =$$

$$= E\left[\log\left(\frac{|\sigma|}{|m|} \exp\left(\frac{1}{2} \left(\frac{(x-m)^2}{\sigma^2} - \frac{(x-\mu)^2}{\sigma^2}\right)\right)\right)\right] =$$

$$= E\left[\log \frac{|\sigma|}{|m|}\right] + \frac{1}{2} E\left[\frac{(x-m)^2}{\sigma^2} - \frac{(x-\mu)^2}{\sigma^2}\right] =$$

$$= \log \frac{|\sigma|}{|m|} + \frac{1}{2\sigma^2} E[(x-m)^2] - \frac{1}{2\sigma^2} E[(x-\mu)^2]$$

$$\text{Var}(x) = E[x^2] - \mu^2 \Rightarrow E[x^2] = \mu^2 + \sigma^2, E[x^2 - 2mx + m^2] =$$

$$= \mu^2 + \sigma^2 - 2m\mu + m^2, \text{ similarly } E[(x-\mu)^2] = \mu^2 + \sigma^2 - 2\mu^2 + \mu^2 = \sigma^2 \text{ after plugging } m=\mu \text{ in general form}$$

$$\text{So, answer is } \log \frac{|s|}{|G|} + \frac{1}{2s^2} (u^2 + G^2 - 2m^2 u + m^2) -$$

$$-\frac{1}{2s^2} \cdot \cancel{s^2} = \boxed{\log \frac{|s|}{|G|} - \frac{1}{2} + \frac{1}{2s^2} (u^2 + G^2 - 2m^2 u + m^2)}$$

+