

BAYESIAN NETWORKS (DGM)

Conditional Independence: $X \perp Y | Z$; $P(X, Y | Z) = P(X | Z)P(Y | Z), P(X, Y, Z) = P(X | Z)P(Y | Z)$

Product rule: $P(A, B | C) = \frac{P(A | B, C)}{P(C)} = \frac{P(A | B)P(B | C)}{P(C)} = P(A | B)P(B | C)$

Markov Anytime: Any random variable is locally dependent on its parent nodes. $X \perp (X_{nonDesc} \setminus X_{parent}) | X_{parent}$ (useful to find all CIs in DGM), so $P(x_1, x_2, \dots, x_N) = \prod_{i=1}^N P(x_i | x_{\pi_i})$.

Graph Separation: From the three canonical 3-node graphs, a path is said to be *blocked* / d-separated (conditional independence) if it contains a node satisfying either:

- the arrows meet 'head-to-tail' or 'tail-to-tail' at that node, and the node is in conditioning set C;
- the arrows meet head-to-head at that node, and neither the node nor any of its descendants is in C.

Bayes Ball Algorithm: This is a reachability-based procedure: (1) shade the nodes in the conditioning set C; (2) place a ball at each node in the query set A; (3) let the balls traverse the graph according to the d-separation rules (separation means no passing). If none of the balls reach any node in B, then $A \perp B | C$; otherwise $A \not\perp B | C$. The procedure can be implemented using a breadth-first search.

Markov Blanket: For DGM, the Markov blanket of a node X_j comprises the set of parents, children, and co-parents of X_j : $\{X_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N\} = P(x_j | MB(x_j))$. For UGM, the Markov blanket of a node X_j is its direct neighbour.

MARKOV RANDOM FIELDS (UGM)

Rules of Conditional Independence: $X | Y$ means $P(X, Y) = P(X)P(Y)$

- Symmetry:** $X \perp Y | Z \Rightarrow Y \perp X | Z$.
- Decomposition:** $X \perp \{Y, W\} | Z \Rightarrow X \perp Y | Z$ and $X \perp W | Z$.
- Weak Union:** $X \perp \{Y, W\} | Z \Rightarrow X \perp Y | \{W, Z\}$ and $X \perp W | \{Y, Z\}$.
- Contraction:** $X \perp Y | \{W, Z\}$ and $X \perp W | Z \Rightarrow X \perp \{Y, W\} | Z$.
- Intersection:** $X \perp Y | \{W, Z\}$ and $X \perp W | \{Y, Z\} \Rightarrow X \perp \{Y, W\} | Z$.

I-MAP: All independencies probability distribution P or graph G.

Conditional independence for UGMs: (1) global Markov property: Given node sets A, B, C, we have $X_A \perp X_B | X_C$ if C separates A from B in graph G. (2) local Markov property: A node X_S is conditionally independent of all others given its Markov blanket: $X_S \perp V \setminus \{mb(X_S), X_S\} | mb(X_S)$. (3) pairwise Markov property: Two nodes X_S and X_T are conditionally independent given the rest nodes if no edge connects them: $X_S \perp X_T | V \setminus \{X_S, X_T\}$, where $\mathcal{E}_{ST} = \emptyset$.

Hammersley-Clifford Theorem: A positive distribution $p(y) > 0$ satisfies the CI properties of an undirected graph G if p can be written as a product of factors, one per maximal clique: $p(y | \theta) = \prod_{C \in \mathcal{C}} \prod_{c \in C} \psi_c(y_C)$

Moralization: A DGM can be converted into a UGM by marrying the unmarried parents of a node.

Conditional Random Fields: A CRF is an MRF where all clique potentials are conditioned on all observed variables x : $p(y | x, w) = \frac{1}{Z(x, w)} \prod_c \psi_c(y_C | x, w)$. MRF models $P(x, y)$, while CRF models $P(y | x)$.

VARIABLE ELIMINATION AND BELIEF PROPAGATION

Reconstituted Graph: Elimination means removing a node from the graph and connecting all its remaining neighbors. The original and newly created edges are recorded in the reconstituted graph.

Variable Elimination

- Goal:** Compute $p(x_F | x_E)$ via Variable Elimination on a directed graph.
- Initialize(G, F):** choose elimination order I with F last; for each node X_i , place its CPT $p(x_i | x_{\pi_i})$ into the active list.
- Evidence(E):** for each evidence variable $i \in E$, add $\delta(x_i, \bar{x}_i)$ to the active list.
- Update(G) (marginalization):** for each variable i in ordering I:

 - collect all active factors involving x_i and multiply potentials to get $\phi_i(x_i)$
 - Example: $\phi_2(x_2, x_3) = p(x_2 | x_3) \phi_1(x_2)$
 - marginalize: $m_2(x_2) = \sum_{x_3} \phi_2(x_2, x_3) = \sum_{x_3} p(x_2 | x_3) m_6(x_2, x_5)$
 - replace old factors with m_i in the active list

- Normalize(F):** obtain the desired conditional distribution $p(x_F | x_E) = \frac{\phi_F(x_F)}{\sum_{x_F} \phi_F(x_F)}$.

Treewidth: one less than the smallest achievable cardinality of the largest clique over all possible elimination orderings. The maximum clique will influence the computation complexity.

Sum-Product Algorithm

- Initialize:** For each node $i \in V$, set $\psi_i^E(x_i) = \begin{cases} \psi_i(x_i) \delta(x_i, \bar{x}_i), & i \in E, \\ \psi_i(x_i), & i \notin E. \end{cases}$
- Choose a root r , perform inward pass (collect): For each edge $(j \rightarrow i)$ toward r send $m_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_j^E(x_j) \psi_{j \rightarrow i}(x_j) \prod_{k \in N(j) \setminus i} m_{k \rightarrow j}(x_j)$.
- Outward pass (distribute): For each edge $(i \rightarrow j)$ away from r send $m_{i \rightarrow j}(x_j) = \sum_{x_i} \psi_i^E(x_i) \psi_{i \rightarrow j}(x_i) \prod_{k \in N(i) \setminus j} m_{k \rightarrow i}(x_i)$.
- Marginal:** For any node i , $p(x_i | E) \propto \psi_i^E(x_i) \prod_{j \in N(i)} m_{j \rightarrow i}(x_i)$.

FACTOR GRAPH AND JUNCTION TREE ALGORITHM

Joint probability of factor graph: $p(x) = \prod_S f_S(x_S)$. The method for constructing factor graphs is based on the original joint distribution. Taking DGMs as an example, if it's $P(X_i)$, there's an independently connected factor X_i ; if it's a conditional probability, this factor is between multiple variables.

Messages from leaf node: (1) message from leaf variable \rightarrow factor: $v_{IS}(x_i) = 1$; (2) message from leaf factor \rightarrow variable: $M_{SI}(x_i) = f_S(x_i)$

Sum-Product of graph factor

- Evidence incorporation: $\psi_i^E(x_i) = \psi_i(x_i) \delta(x_i, \bar{x}_i)$ if $i \in E$, else $\psi_i^E(x_i) = \psi_i(x_i)$.
- Collect (leaf \rightarrow root): $M_{SI}(x_i) = \sum_{x_S \in N(i) \setminus i} f_S(x_S) \prod_{j \in N(S) \setminus i} v_{JS}(x_j)$, $v_{IS}(x_i) = \prod_{j \in N(i) \setminus i} M_{SI}(x_j)$.
- Distribute (root \rightarrow leaves): $v_{JS}(x_j) = \prod_{i \in N(j) \setminus j} M_{SI}(x_i)$, $M_{SI}(x_i) = \sum_{x_S \in N(i) \setminus i} f_S(x_S) \prod_{j \in N(S) \setminus i} v_{JS}(x_j)$.

Iterative Proportional Fitting (IPF) for Tabular MRFs

- Initialize clique potentials $\psi_c(y_c) = 1$ for all $c = 1, \dots, C$.
- for $c = 1$ to C do
- Model marginal: $\psi_c(y_c | \psi) = \sum_{y \in \mathcal{Y}} p(y | \psi)$
- Empirical marginal: $\hat{\psi}_c(y_c) = \rho_{emp}(y_c) = \frac{1}{N} \sum_{i=1}^N \delta(y_{ci}, y_c)$
- IPF multiplicative update: $\psi_c(y_c) \leftarrow \psi_c(y_c) \frac{\hat{\psi}_c(y_c)}{\psi_c(y_c | \psi)}$
- end for

MIXTURE MODEL AND EM ALGORITHM

Mixture Models: Probabilistic models formed by taking linear combinations of more basic distributions.

Marginal: $p(x_i) \propto \psi_i^E(x_i) \prod_{s \in N(i)} \mu_{si}(x_i)$.

MAP-Elimination Algorithm

- Input: graph $G(\mathcal{V}, \mathcal{E})$, evidence set E
- Initialize: choose elimination order I; for each $X_i \in \mathcal{V}$ place $p(x_i | x_{\pi_i})$ on active list
- Evidence: for each $i \in E$ place $\delta(x_i, \bar{x}_i)$ on active list
- Update (maximization): for each $i \in I$:
 - collect all active factors involving x_i and remove them, then multiply them: $\phi_i^{\max}(x_i) = \prod \phi_i$
 - maximize out x_i : $m_i^{\max}(x_i) = \max_{x_i} \phi_i^{\max}(x_i)$ and place $m_i^{\max}(x_i)$ back on the active list
- Maximum: final scalar on active list is $\mu_p^E(x)$

Max-Product for trees

- Input: tree $\mathcal{T}(\mathcal{V}, \mathcal{E})$, evidence set E
- Evidence: for each $i \in \mathcal{V}$ set $\psi_i^E(x_i) = \psi(x_i) \mathbf{1}[i \notin E] + \psi(x_i) \delta(x_i, \bar{x}_i) \mathbf{1}[i \in E]$
- ChooseRoot(\mathcal{T})
- Inward pass (collect): for $s \in N(f)$ do Collect(f, s), where Collect(i, j) : for $k \in N(j) \setminus \{i\}$ do Collect(j, k): $\mu_{j \rightarrow i}(x_j) = \max_{x_j} (\psi_j^E(x_j) \psi_{j \rightarrow i}(x_j) \prod_{k \in N(j) \setminus \{i\}} \mu_{k \rightarrow j}(x_j))$; $\delta_{ji}(x_j) = \arg \max_{x_j} (\psi_j^E(x_j) \psi_{j \rightarrow i}(x_j) \prod_{k \in N(j) \setminus \{i\}} \mu_{k \rightarrow j}(x_j))$
- Root MAP at X_f : $x_f^* = \arg \max_{x_f} (\psi_f^E(x_f) \prod_{e \in f} \mu_{e \rightarrow f}(x_f))$
- Outward pass (distribute): for $s \in N(f)$ do Distribute(f, s), where Distribute(i, j): set $x_j^* = \delta_{ji}(x_j^*)$; for $k \in N(j) \setminus \{i\}$ do Distribute(j, k)
- Output: MAP configuration $\{x_i^*\}_{i \in \mathcal{V}}$

HIDDEN MARKOV MODELS

Joint probability: $p(z_1, \dots, z_T) = p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \prod_{t=1}^T p(x_t | z_t)$

Transition matrix: Properties of the state transition matrix $A \in \mathbb{R}^{K \times K}$: (1) $A_{jk} = p(z_{nk} = 1 | z_{n-1}, j = 1, \dots, K) \leq 1$, $\sum_k A_{jk} = 1$, (3) there are $K(K - 1)$ independent parameters.

Pairwise MRF: $p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i, j) \in E} \psi_{ij}(x_i, x_j)$

Viterbi Algorithm: The maximal probability of the joint distribution $p(x_1, \dots, x_N, z_1, \dots, z_N)$ is given by the max of $\omega_{zN}(z_N)$ at the root node, $\max_{z_1, \dots, z_N} p(x_1, \dots, x_N, z_1, \dots, z_N) = \max_{z_N} \omega_{zN}(z_N)$. $\omega(z_1) = \ln p(z_1) + \ln p(x_1 | z_1)$, $\omega(z_{n+1}) = \ln p(x_{n+1} | z_{n+1}) + \max_{z_n} \{\ln p(z_{n+1} | z_n) + \omega(z_n)\}$.

MONTE CARLO INFERENCE

Monte Carlo Principle: Draw samples $x^{(i)} \sim p(x)$ i.i.d.; define empirical measure $P_N(x) = \frac{1}{N} \sum_{i=1}^N \delta(x^{(i)})$. For any function f , define $I(f) = \int f(x) p(x) dx$ and estimator $I_N(f) = \frac{1}{N} \sum_{i=1}^N f(x^{(i)})$. By law of large numbers, $I_N(f) \xrightarrow{a.s.} I(f)$. If $\text{Var}[f(x)] = \sigma_f^2 < \infty$, then $\text{Var}[I_N(f)] = \sigma_f^2 / N$. By central limit theorem, $\sqrt{N}(I_N(f) - I(f)) \xrightarrow{d} \mathcal{N}(0, \sigma_f^2)$.

Summary: By independently sampling from the target distribution $p(x)$ to construct the empirical measure P_N and using the corresponding sample mean $I_N(f)$ to estimate the integral $I(f)$, one can rely on the law of large numbers and the central limit theorem as $N \rightarrow \infty$ to ensure consistency and asymptotic normality, thereby achieving a reliable numerical approximation of complex integrals or expectations.

Rejection Sampling: $p(z) = \frac{1}{Z} p(z)$, but $p(z)$ is difficult to sample so we need a common distribution $q(z)$. Sample $z^{(i)} \sim q(z)$ and $u \sim U(0, 1)$, if $u < \frac{p(z^{(i)})}{k q(z^{(i)})}$, then accept $z^{(i)}$.

Importance Sampling: To estimate expectations under a target distribution $p(z)$, we draw samples from an easier proposal distribution $q(z)$ and correct the mismatch using importance weights. Using the identity $E_p[f(z)] = \int f(z) \frac{p(z)}{q(z)} dz$, the Monte Carlo estimator becomes $\frac{1}{L} \sum_{l=1}^L \frac{p(z_l)}{q(z_l)} f(z_l)$ with $z_l \sim q(z)$. When only unnormalized densities $\tilde{p}(z)$ and $\tilde{q}(z)$ are available, self-normalized importance sampling is used: $E_p[f(z)] \approx \sum_{l=1}^L w_l f(z_l)$, where $w_l = \frac{\tilde{p}(z_l) / \tilde{q}(z_l)}{\sum_k \tilde{p}(z_k) / \tilde{q}(z_k)}$.

Ancestral Sampling: Given a joint distribution factorized as $p(z) = \prod_{k=1}^M p(z_k | pa_k)$ in a Bayesian network, we can generate exact samples by drawing each variable in topological order. For root nodes, sample directly from $p(z_1)$; for each subsequent node, sample $z_i \sim p(z_i | pa_i)$ using already-sampled parent values. Repeating this for $i = 1, \dots, M$ yields a complete sample from the joint distribution $p(z)$.

MCMC (Metropolis-Hastings):

- Initialize $x^{(0)}$
- for $i = 0$ to $N - 1$ do
 - Sample $u \sim U[0, 1]$, $x' \sim q(x^{(i)} | x^{(i)})$ \triangleright draw acceptance threshold and proposal sample
 - If $u < \mathcal{A}(x', x^{(i)}) = \min \left\{ 1, \frac{\tilde{p}(x') / q(x'|x^{(i)})}{\tilde{p}(x^{(i)}) / q(x'|x^{(i)})} \right\}$, $x^{(i+1)} = x'$
- end for

Gibbs Sampling:

- Initialize $\{x_i : i = 1, \dots, M\}$
- for $\tau = 1, \dots, T$ do
 - Sample $x_1^{(\tau+1)} \sim p(x_1 | x_2^{(\tau)}, x_3^{(\tau)}, \dots, x_M^{(\tau)})$
 - Compute model expectation: $\hat{\phi} = \frac{1}{S} \sum_{s=1}^S \psi_s^{(\tau)}$.
 - Compute per-sample gradients in minibatch: $g_{i,k} = \phi_i(y_i) - \hat{\phi}$, $i \in B$.
 - Minibatch gradient: $g_k = \frac{1}{B} \sum_{i \in B} g_{i,k}$.
 - Parameter update: $\theta_{k+1} = \theta_k - \eta g_k$.
 - Increment iteration: $k \leftarrow k + 1$. Then, decrease step size η .

VARIANCE INFERENCE

Mean Field Approximation: Given a joint probability, the goal of MFA is to find a $q(z)$ for $p(z|x)$, which satisfies $q(z) = \prod_i q_i(z_i)$.

Mean Field Approximation

- Goal: Derive $q_j^*(Z_j)$ given Joint $p(X, Z)$, q_j^* is normalized q_j
- $\ln q_j^*(Z_j) = \mathbb{E}_{i \neq j} [\ln p(X, Z)] + \text{const}$
- while Not Converged do
 - Log-joint: $\ln p = \sum_i \ln p(\text{factors})$ (Expand all terms)
 - Drop irrelevant terms: Keep terms with Z_j ; others \rightarrow const
 - Expectation: Replace neighbors $Z_{i \neq j} \rightarrow \mathbb{E}[Z_i]$
 - Match & update: $-AZ_j^2 + BZ_j \rightarrow \mathcal{N}(B/2A, 1/2A)$, $(A - 1) \ln Z_j - BZ_j \rightarrow \text{Gam}(A, B)$, $\sum Z_k \ln \pi_k \rightarrow \text{Cat}(\pi)$

Forward vs Backward KL-Divergence: KL is not symmetrical, minimize $KL(q || p) \neq KL(p || q)$. Reverse KL: $KL(q || p) = \sum_z q(z) \ln \frac{q(z)}{p(z)}$, if $p(z) = 0$, $q(z) = 0 \Rightarrow q$ under-estimate p . Forward KL: $KL(p || q) = \sum_z p(z) \ln \frac{p(z)}{q(z)}$, if $p(z) > 0$, $q(z) > 0 \Rightarrow q$ over-estimate p . Choose the Backward KL: $KL(q || p) = \sum_z q(z) \ln \frac{q(z)}{p(z)}$. Tractable lower bound. 2. Statically more sensible.

Example 1: The Univariate Gaussian Observed: $D = \{x_1, \dots, x_N\}$, Goal: infer posterior. Likelihood: $p(D | \mu, \sigma^2) = \left(\frac{\tau}{2\pi} \right)^N \exp(-\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2)$, conjugate prior: $p(\mu | \tau) = N(\mu | \mu_0, (\lambda_0 \tau)^{-1})$, $p(\tau) = \text{Gam}(\tau | a_0, b_0)$, factorized: $q(\mu, \tau) = q(\mu | \mu_0, \tau)$, optimal solutions: $q_{\mu}^*(\mu) = N(\mu | \mu_N, \lambda_N^{-1})$, $\mu_N = \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N}$, $b_N = (\lambda_0 + N) \bar{x} + b_0$, $\lambda_N = (\lambda_0 + N) \tau$, $q_{\tau}^*(\tau) = \text{Gam}(\tau | a_N, b_N)$, $a_N = a_0 + \frac{N^2}{2}$, $b_N = b_0 + \frac{1}{2} \mu_N (\lambda_N \bar{x}^2 + \lambda_0 \mu_0^2 - \lambda_0 \mu_0 \bar{x}^2)$. Solutions are coupled \Rightarrow iteration approach: (1) Initial guess for \bar{x} , recompute $q_{\mu}^*(\mu)$. (2) Use revised $q_{\mu}^*(\mu)$ to extract the moments $E[\mu]$ and $E[\mu^2]$, use these to re-compute $q_{\tau}^*(\tau)$. (3) Use revised $q_{\tau}^*(\tau)$ to extract the moment $E[\tau]$, and use this to re-compute $q_{\mu}^*(\mu)$. (4) Repeat until convergence.

Example 3: Ising Model $p(x) = p(x) p(y | x)$, prior: $p(x) = \frac{1}{Z} \exp(E_0(x))$, $E_0(x) = -\sum_i \sum_{j \in \text{nbr}_i} W_{ij} x_i x_j$. Likelihood: $p(y | x) = \prod_i p(y_i | x_i) = \exp(\sum_i L_i(x_i))$, $\ln p(x, y) = \sum_i \sum_{j \in \text{nbr}_i} W_{ij} x_i y_j + L_i(x_i) - \ln Z_0$, optimized factor: $\log q(x_j) = E_{x_j} [\log p(x, y)] + \text{const}$, thus $q(x_j) \propto \exp(\sum_i \sum_{j \in \text{nbr}_i} W_{ij} x_i y_j + L_i(x_i))$, $\mu_j = \sum_i x_j y_j q_j(x_j)$.

MAP INFERENCE AND AUGMENTING PATHS

$E(w) = \sum_n U_n(w_n) + \sum_{(m,n)} P_{mn}(w_m, w_n)$

Cut and capacity: A cut is a node partition (S, T) on \mathcal{G} such that $S \cap T = \emptyset$ and $S \cup T = \mathcal{V}$. Capacity $C(S, T)$ is the sum of weights of edges leaving S.

MAP: Differentiating cases and first identify possible 2^n solutions (having different cuts), the solution with minimum $E(w)$ is the min-cut. The cut must sever the path from S to T. Whether a pixel belongs to the S or T depends on how it was cut. If the S was cut off, then it belongs to the T.

Converting MRF: In some cases, you may need to perform a graph transformation first, converting it into a standard form based on the number of observed variables. Considering S is label '1' and T is label '0', so we regard 'right' direction as $1 \rightarrow 0$.

Augment: Iterate using the shortest path first.

EXAMPLE 1: EM

We are given two biased coins, Coin A and Coin B, each with an unknown probability of landing heads. In a series of trials, we randomly choose one of the two coins and flip it a single time. However, we do not know which coin was chosen for each flip. Let us denote the random variables of choosing each coin as a 1-of-2 representation $Z_k \in \{0, 1\}$, where $Z_0 = 0 :=$ Coin A and $Z_1 = 1 :=$ Coin B. The outcome of a coin flip is $X \in \{0, 1\}$, where $X = 0 :=$ Head and $X = 1 :=$ Tail. We denote the probability of choosing Coin B as π , and the probability of getting a tail for Coin A and B as μ_0 and μ_1 respectively.

a) Given N i.i.d. observations of the coin flips, draw the Bayesian network and write the joint probability distribution that represents the coin choosing and flipping process.

The probability of choosing each coin is $p(Z) = \prod_k \pi_k^{Z_k}$, where $\pi_0 = 1 - \pi$ and $\pi_1 = \pi$.

The likelihood of outcome X given the coin is $p(X | \mu_k) = \mu_k^{X_k} (1 - \mu_k)^{1-X_k}$.

Thus, the joint distribution is $p(X, Z) = p(Z)p(X | Z) = \prod_k \pi_k^{Z_k} p(X | \mu_k)^{1-Z_k}$.

b) Using the EM algorithm, derive the expressions of μ_0 and μ_1 after the Maximization step.

The EM objective is $Q(\theta, \theta^{old}) = \sum_Z p(Z | X, Z, \theta^{old}) \ln p(X, Z | \theta)$.

Let responsibilities be $\gamma(Z_k) = p(Z_k = 1 | X, \theta^{old})$. Then, $Q(\theta, \theta^{old}) = \sum_k \ln \pi_k \gamma(Z_k) + \ln p(X_k | \mu_k)$.

Taking derivative w.r.t. μ_k : $\frac{\partial Q}{\partial \mu_k} = \sum_k \gamma(Z_k) \left(\frac{\pi_0}{\mu_k} - \frac{1 - \pi_0}{1 - \mu_k} \right) = 0$.

Solving gives $\mu_k = \frac{\sum_k \gamma(Z_k) \pi_0}{\sum_k \gamma(Z_k) \pi_0}$, where $N_k = \sum_i \gamma(Z_k)$, then $\mu_k = \frac{1}{N_k} \sum_k \gamma(Z_k) \pi_0$.

c) The results of these flips are recorded as a sequence of heads (H) and tails (T): (H, T, H) . Taking equal probability of choosing either Coin A or B and the initial values of $\mu_0 = 0.6$ and $\mu_1 = 0.5$, perform one iteration of the EM algorithm to update μ_0 and μ_1 .

The likelihood is $p(X | \mu_k) = \mu_k^{X_k} (1 - \mu_k)^{1-X_k}$ and responsibilities: $\gamma(Z_k) = \frac{0.5(1-\mu_k^{old})}{0.5(1-\mu_k^{old}) + 0.5(1-\mu_1^{old})} = \frac{1-0.6}{(1-0.6)+(1-0.5)} = 0.44$.

For $X_2 = 1$: $\gamma(Z_2, 0) = \frac{0.5(1-\mu_0^{old})}{0.5(1-\mu_0^{old}) + 0.5(1-\mu_1^{old})} = \frac{(1-0.6)(1-0.5)}{(1-0.6)+(1-0.5)} = 0.454$.

For $X_3 = 0$: $\gamma(Z_3, 0) = 0.444$.

Total responsibility: $N = 0.444 + 0.545 + 0.444 = 1.433$.

Parameter update: $\mu_0 = \frac{1}{N} \sum_k \gamma(Z_k) \pi_0 = \frac{0.545}{1.433} = 0.380$.

EXAMPLE 2: EM

Figure 9.1 shows a Bayesian network with both binary and continuous state latent random variables, $i, Z \in \{0, 1\}$ and $T \in \mathbb{R}$. In addition, $X = 0.5$ is the observed random variable. The maximum log-likelihood of T : $\arg \max_T \log p(T | X)$, can be obtained from the Expectation-Maximization (EM) algorithm. The EM algorithm iterates between the Expectation step that evaluates the expected complete data log-likelihood with respect to $p(Z | X, T, \theta^{old})$. T^{old} is the value of T from the previous iteration of the EM algorithm. $\{\lambda = 0.1, w_a = 0.5, w_{a1} = 0.5, w_b = 0.8, w_{b1} = 0.2, \tau_a = 1.0, \tau_b = 1.2, U = 0.6\}$ are known hyperparameters of the following distributions: $p(Z) = \lambda Z(1 - \lambda)^{1-Z}$, $p(X | T, Z) = \mathcal{N}(X | w_a + w_{a1} T, \tau_a)^2$, $p(T | w_b + w_{b1} T, \tau_b) = \mathcal{N}(T | w_b + w_{b1} T, \tau_b)^2$, $p(U | T) = \sqrt{\frac{\pi}{2\pi}} \exp(-0.5\pi(T - w_b - w_{b1} T)^2)$, $p(U) = U$.

a. Derive the expression for the posterior $p(Z | X, T, \theta^{old})$ from the Bayesian Network.

Derive the expression for T that maximizes the expected complete data log-likelihood with respect to $p(Z|X, T^{old})$.

c. Given the initial value of $T = 2.0$, find the value of T in the next EM iteration.

Answer:

Joint distribution: $p(X, Z, T) = p(T)p(Z)p(X|T, Z) = p(Z)p(X|T, Z)$.

$$p(Z|X, T^{old}) = \frac{p(Z)p(X|T, Z)}{\sum_Z p(Z)p(X|T, Z)} = \frac{p(Z)p(X|T, Z)}{\sum_Z p(Z)p(X|w_{a0} + w_{a1}T, \tau_a)^Z \cdot p(X|w_{b0} + w_{b1}T, \tau_b)^{(1-Z)}}.$$

$$= \frac{\lambda^Z(1-\lambda)^{(1-Z)} \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a)^Z \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)^{(1-Z)}}{\sum_Z \lambda^Z(1-\lambda)^{(1-Z)} \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a)^Z \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)^{(1-Z)}}.$$

$$= \frac{\lambda^Z(1-\lambda)^{(1-Z)} \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a)^Z \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)^{(1-Z)}}{\lambda \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a) + (1-\lambda) \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)}$$

Let's define: $\gamma(Z) = p(Z=0|X, T^{old}) = \frac{(1-\lambda) \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)}{\lambda \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a) + (1-\lambda) \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)}$

$$\gamma(Z=0) = p(Z=0|X, T^{old}) = \frac{(1-\lambda) \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a)}{\lambda \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a) + (1-\lambda) \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b)}.$$

$$b. Q = \sum_Z p(Z|X, T^{old}) \ln p(X, Z|T) = \sum_Z \gamma(Z) \ln p(X|T, Z)$$

$$= \gamma(Z=0) \ln p(Z=0|X|T, Z=0) + \gamma(Z=1) \ln p(Z=1|X|T, Z=1)$$

$$= \gamma(Z=0) \ln \{ (1-\lambda) \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b) \} + \gamma(Z=1) \ln \{ \lambda \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a) \}$$

$$= \gamma(Z=0) \frac{\partial}{\partial T} \{ \ln (1-\lambda) \cdot \mathcal{N}(X|w_{b0} + w_{b1}T, \tau_b) \} + \gamma(Z=1) \frac{\partial}{\partial T} \{ \ln \lambda + \ln \lambda \cdot \mathcal{N}(X|w_{a0} + w_{a1}T, \tau_a) \}$$

$$= \frac{X(\gamma(Z=0)w_{b1}\tau_b + \gamma(Z=1)w_{a1}\tau_a) - \gamma(Z=0)w_{b0}w_{b1}\tau_b - \gamma(Z=1)w_{a0}w_{a1}\tau_a}{\gamma(Z=0)\tau_b^2 + \gamma(Z=1)\tau_a^2}$$

c. Just use $T = 2.0$ to compute terms in T .

EXAMPLE 3: VARIANCE INFERENCE

Given a dataset $\mathcal{D} = \{x_1, \dots, x_N\}$ of observed values of a random variable X , which are assumed to be drawn independently from a univariate Gaussian distribution: $p(x|\mu, \tau) = \mathcal{N}(x|\mu, \tau^{-1})$ parameterized by the mean μ and precision τ . The prior distributions for μ and τ are given by a univariate Gaussian distribution: $p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1})$ and a Gamma distribution: $p(\tau) = \text{Gam}(\tau|a_0, b_0)$, respectively. $\{\mu_0, \lambda_0, a_0, b_0\}$ are hyperparameters of the prior distributions. (i) Using variation inference and mean field approximation, derive the distributions $q(\mu)$ and $q(\tau)$ that best approximate the conditional distributions $p(\mu|\mathcal{D})$ and $p(\tau|\mathcal{D})$.

(ii) Derive the expressions for the expected mean μ and precision τ under the approximated distributions $q(\mu)$ and $q(\tau)$.

(iii) Derive the expressions for the mean μ and precision τ that maximize the log-likelihood.

(iv) Comment on the difference*** between the mean μ and precision τ obtained from variational inference and maximum log-likelihood.

1. Gamma distribution: $\text{Gam}(\tau|a, b) = \frac{1}{\Gamma(a)} b^a \tau^{a-1} \exp(-b\tau)$, where $\Gamma(a)$ is the Gamma function;

2. Gaussian distribution: $p(x|\mu, \tau^{-1}) = \left(\frac{\tau}{2\pi}\right)^{1/2} \exp\left(-\frac{\tau}{2}(x-\mu)^2\right)$;

3. Completing the square: $ax^2 + bx + c = 0 \Rightarrow a(x+d)^2 + e = 0$, where $d = \frac{b}{2a}$, $e = c - \frac{b^2}{4a}$.

Solution

i) Joint probability distribution: $p(D, \mu, \tau) = p(D|\mu, \tau)p(\mu|\tau)p(\tau)$

$= \prod_n \mathcal{N}(x_n|\mu, \tau^{-1}) \cdot \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \cdot \text{Gam}(\tau|a_0, b_0)$ Mean field approximation: $q(\mu, \tau) = q(\mu)p(\tau)$ The optimal factors $q(\mu)$ and $q(\tau)$ are obtained from variational inference: $\ln q^*(Z_j) = \ln \int \ln p(D, \mu, \tau) d\mu d\tau + \text{const}$.

We get: $\ln q^*(\mu) = \mathbb{E}_\tau[\ln p(D, \mu, \tau)] + \text{const}$

$= \mathbb{E}_\tau[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) + \ln \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1})] + \text{const}$

$= \mathbb{E}_\tau[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) + \ln \left(\frac{\lambda_0\tau}{2\pi} \right)^{1/2} \exp\left(-\frac{\lambda_0\tau}{2} (\mu - \mu_0)^2\right)] + \text{const}$

$= \mathbb{E}_\tau \left[\sum_n \left(-\frac{\tau}{2} (x_n - \mu)^2 - \frac{\lambda_0\tau}{2} (\mu - \mu_0)^2 \right) \right] + \text{const} = -\frac{\mathbb{E}_\tau[\tau]}{2} \left[\sum_n (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right] + \text{const}$

Let $a = \lambda_0 + N$, $b = 2(\lambda_0\mu_0 + \sum_n x_n)$ and $c = (\lambda_0\mu_0^2 + \sum_n x_n^2)$, we further let $d = \frac{b}{2a}$ and $e = c - \frac{b^2}{4a^2}$, we get: $\ln q^*(\mu) = -\frac{\mathbb{E}_\tau[\tau]}{2} [a(\mu + d)^2 + e] + \text{const}$, where e can be put into the constant terms since it is independent of μ .

Thus, we get: $\ln q^*(\mu) = -\frac{(\lambda_0+N)}{2} \left[\left(\mu - \frac{\lambda_0\mu_0 + \sum_n x_n}{\lambda_0+N} \right)^2 \right] + \text{const}$. Taking exponential on both sides and normalize, we have: $\hat{q}^*(\mu) = \exp\left(-\frac{(\lambda_0+N)}{2} \mathbb{E}_\tau[\tau] \left[\left(\mu - \frac{\lambda_0\mu_0 + \sum_n x_n}{\lambda_0+N} \right)^2 \right]\right)$, which is a Gaussian distribution, i.e., $\hat{q}^*(\mu) = \mathcal{N}\left(\mu \mid \frac{\lambda_0\mu_0 + \sum_n x_n}{\lambda_0+N}, (\lambda_0+N)\mathbb{E}_\tau[\tau]^{-1}\right)$.

ii) Expected mean: $\mathbb{E}_{q(\mu)}[\mu] = \int q(\mu) \mu d\mu$. Expected precision: $\mathbb{E}_{q(\tau)}[\tau] = f(\tau) \tau d\tau$

iii) Maximum log-likelihood: $\arg \max_\mu \ln p(D|\mu, \tau) = \arg \max_\mu \left\{ \ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right\} = \arg \max_\mu \left\{ -\frac{\tau}{2} \sum_n (x_n - \mu)^2 \right\}$

Taking partial derivatives w.r.t. μ and equate to 0, we get: $\tau(\sum_n x_n - \mu N) = 0 \Rightarrow \mu = \frac{\sum_n x_n}{N}$

$\arg \max_\tau \left\{ \ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right\} = \arg \max_\tau \left\{ \frac{N}{2} \ln \left(\frac{\tau}{2\pi} \right) - \frac{\tau}{2} \sum_n (x_n - \mu)^2 \right\} = \frac{N}{2\tau} - \frac{\sum_n (x_n - \mu)^2}{2\tau}$

Taking partial derivatives w.r.t. τ and equate to 0, we get: $\left(\frac{N}{2}\right) \frac{2\pi}{(\frac{\tau}{2\pi})^2} - \frac{1}{2} \sum_n (x_n - \mu)^2 = 0 \Rightarrow \tau = \frac{N}{\sum_n (x_n - \mu)^2}$

iv) Solution obtained from variational inference is an approximation and often under-estimate. However, maximum log-likelihood ignores prior and would not be accurate if number of observations is low.

EXAMPLE 4: VARIANCE INFERENCE

Figure 13.1 shows a three-node undirected graphical model, where $X_i \in \mathbb{R}_{\geq 0}$, $\psi(X_1, X_2) = \exp(-\alpha X_1 X_2)$, $\psi(X_t) = \exp(-\beta X_t)$, and $\alpha = 0.5$ and $\beta = 2.5$ are constants.

(a) Using variational inference, find the expressions of the expectation of X_1, X_2 , and X_3 under $q(X_1)$,

$q(X_2)$, and $q(X_3)$, respectively, where

$$q(X_1, X_2, X_3) = q(X_1)q(X_2)q(X_3)$$

is the mean-field approximation of the posterior distribution $p(X_1, X_2, X_3)$.

(b) Taking the initial expected values of X_2 and X_3 under $q(X_2)$ and $q(X_3)$ to be 2.0 and 1.0, respectively, find the mean-field approximation $q(X_1, X_2, X_3)$ after one iteration.

$$1. \int \exp\{kx\} dx = \frac{1}{k} \exp(kx) + \text{const}, \quad k \neq 0.$$

$$2. \int x \exp\{kx\} dx = \left(\frac{kx-1}{k^2} \right) \exp(kx), \quad k \neq 0.$$

(a) **Joint probability.**

$$p(X) = \frac{1}{k} \prod_{i < j} \psi(X_i, X_j) \prod_k \psi(X_k) = \frac{1}{k} \exp[-\alpha(X_1 X_2 + X_1 X_3 + X_2 X_3) - \beta(X_1 + X_2 + X_3)].$$

(b) **Mean-field update of X_1 .** $\ln q^*(X_1) = \mathbb{E}_q[X_1] = \frac{1}{k} \cdot 1$, where $\lambda_1 = \alpha(\mu_2 + \mu_3) + \beta$, $\mu_1 = \mathbb{E}[X_1] = \frac{1}{k}$.

Numerical parameters: $\alpha = 0.5$, $\beta = 0.25$. Initial: $\mu_2 = 0.2$, $\mu_3 = 0.1$.

1. **Update X_1 :** $\lambda_1 = 0.5(0.2) + 0.5(0.1) + 0.25 = 0.4$, $\mu_1 = \frac{1}{0.4} = 2.5$.

2. **Update X_2 (use new μ_1):** $\lambda_2 = \alpha\mu_1 + \alpha\mu_3 + \beta = 0.5(2.5) + 0.5(0.1) + 0.25 = 1.55$, $\mu_2 = \frac{1.55}{0.5} \approx 3.14$.

3. **Update X_3 (use new μ_1, μ_2):** $\lambda_3 = \alpha\mu_1 + \alpha\mu_2 + \beta = 0.5(2.5) + 0.5(0.64) + 0.25 = 0.695$.

$$q(X) = q^*(X_1)q^*(X_2)q^*(X_3) \propto \exp[-0.4X_1 - 1.55X_2 - 0.695X_3].$$

EXAMPLE 5: HMM

a. Figure 8.1 shows a homogeneous hidden Markov Model (HMM) over three time steps. The latent random variables are Y_1, Y_2, Y_3 , where $Y_n \in \{0, 1, 2\}$, and the observed random variables are X_1, X_2, X_3 , where $X_n \in \mathbb{R}$. The prior probability of the random variable Y_1 is $p(Y_1 | \pi) = \prod_k \pi_k^{y_k^1}$, where $\pi = \{0.2, 0.5, 0.3\}$. Furthermore, the transition probability is given by:

$$p(Y_n | Y_{n-1}, A) = \prod_k \prod_j \frac{\pi_{j-1}^{y_{n-1}^j}}{\pi_j^{y_n^j}}, \quad A = \begin{bmatrix} 0.2 & \alpha & \beta \\ 0.1 & 0.6 & 0.3 \\ 0.4 & 0.5 & 0.1 \end{bmatrix},$$

and the emission probabilities of the respective observed random variables X_n are:

Given that the minimum probability of the joint distribution $p(Y_1, Y_2, Y_3, X_1, X_2, X_3)$ is 0.000216 and occurs at $Y_1 = 0, Y_2 = 1, Y_3 = 0$, find the unknown values α and β in the transition probability.

b. Figure 8.2 shows an undirected graphic model with six random variables $X_1, X_2, X_3, X_4, X_5, X_6$, where $X_i \in \{0, 1, 2\}$. The potential $\psi(X_i, X_j)$ between any pair of nodes X_i and X_j , where $i < j$, is given in Table 2.2. Given $X_1 = 0, X_3 = 1$ and $X_5 = 2$, find the states of X_2, X_4, X_6 that maximize the joint distribution $p(X_1, X_2, X_3, X_4, X_5, X_6)$.

Answer: a. Joint probability: $p(Y_1) = p(Y_1) \prod_{n=2} p(Y_n | Y_{n-1}) \prod_{n=1} p(X_n | Y_n)$

$$= \min(Y_1) \prod_{n=2} p(Y_n | Y_{n-1}) p(X_n | Y_n) = \min(Y_1) \prod_{n=2} \min(Y_2) p(Y_2 | Y_1) p(X_n | Y_n)$$

$$= \min(Y_3) \prod_{n=2} \min(Y_3) p(Y_3 | Y_2) p(X_n | Y_n) = \min(Y_3) \prod_{n=2} \min(Y_1) p(Y_1 | Y_n) p(X_n | Y_n)$$

$$= \min(Y_0) = 0.000216$$

Given that the minimum probability equals 0.000216 and occurs at $Y_1 = 0, Y_2 = 1, Y_3 = 0$, this implies:

$$\min(Y_0) = 0.000216$$

Since each row of the transition matrix sums to one, we have $0.2 + \alpha + \beta = 1 \Rightarrow \beta = 0.5$.

b. Joint probability: $p(X) = \frac{1}{k} \psi(X_1 = 0, X_2 = 1) \psi(X_1 = 0, X_3 = 1) \psi(X_2, X_3 = 1) \psi(X_2 = 1, X_5 = 2) \psi(X_3 = 1, X_5 = 2) \psi(X_4 = 1, X_5 = 2) \psi(X_4 = 1, X_6 = 2) \psi(X_5 = 1, X_6 = 2) \psi(X_5 = 2, X_6 = 2)$

$$= \max_{X_2} \max_{X_4} \max_{X_5} p(X) = \max_{X_2} \psi(X_1 = 0, X_2) \psi(X_2, X_3 = 1) \psi(X_2, X_5 = 2) \max_{X_4} \psi(X_1 = 0, X_4) \psi(X_4, X_5 = 2) \max_{X_5} \psi(X_1 = 1, X_5) \psi(X_5 = 2, X_6) = \max\{2, 4, 72\} = 72 \quad (X_6 = 2)$$

$$\max_{X_2} \psi(X_1 = 0, X_2) = \max\{2(2), 4(6), 8(9)\} = \max\{6, 24, 72\} = 72 \quad (X_6 = 2)$$

$$\max_{X_4} \psi(X_1 = 0, X_4) = \max\{1(7), 5(8), 7(9)\} = \max\{7, 40, 63\} = 63 \quad (X_4 = 2)$$

$$\max_{X_5} \psi(X_1 = 0, X_5) = \max\{1(5), 7(5), 4(8), 7(6)(9)\} = \max\{15, 160, 378\} = 378 \quad (X_2 = 2)$$

Let $a = \lambda_0 + N$, $b = 2(\lambda_0\mu_0 + \sum_n x_n)$ and $c = (\lambda_0\mu_0^2 + \sum_n x_n^2)$, we further let $d = \frac{b}{2a}$ and $e = c - \frac{b^2}{4a^2}$, we get: $\ln q^*(\mu) = -\frac{\mathbb{E}_\tau[\tau]}{2} [a(\mu + d)^2 + e] + \text{const}$, where e is independent of μ .

Thus, we get: $\ln q^*(\mu) = -\frac{(\lambda_0+N)}{2} \left[\left(\mu - \frac{\lambda_0\mu_0 + \sum_n x_n}{\lambda_0+N} \right)^2 \right] + \text{const}$.

We get: $\ln q^*(\mu) = \mathbb{E}_\tau[\ln p(D, \mu, \tau)] + \text{const}$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) + \ln \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right)^{N/2} \exp\left(-\frac{\tau}{2} \sum_n (x_n - \mu)^2\right) \right) \right] + \text{const}$$

$$= \mathbb{E}_\tau \left[\ln \left(\left(\frac{\tau}{2\pi} \right$$