

Approximate Inference

Phu Sakulwongtana

1 Graphical Model

1.1 Introduction

Definition 1.1. (Types of Graph) There are several kind of graphs that we can use to model the probability distribution: factor graph, undirected graph, and directed graph. Node corresponds to the random variables and the edge in graph indicates statistical dependence between variable.

Definition 1.2. (Dependencies) For the random variable X, Y, V , where we have:

- *Conditional Independence:* $X \perp\!\!\!\perp Y|V$ iff $P(X|Y, V) = P(X|V)$ provided that $P(Y, V) > 0$. We can see that furthermore that:

$$P(X, Y|V) = P(X|Y, V)P(Y|V) = P(X|V)P(Y|V)$$

Please note that, this can generalize the symbol to the sets of random variables as:

$$\mathcal{X} \perp\!\!\!\perp \mathcal{Y}|\mathcal{V} = \{X \perp\!\!\!\perp Y|\mathcal{V} : \forall X \in \mathcal{X}, \forall Y \in \mathcal{Y}\}$$

- *Marginal Independence:* $X \perp\!\!\!\perp Y$ is equivalent to $X \perp\!\!\!\perp Y|\emptyset$ and $P(X, Y) = P(X)P(Y)$

Definition 1.3. (Factor Graph) Factor Graph is a directed graphical representation of the factorized model structure, where each square indicates the factor over the linked variables:

$$P(\mathcal{X}) = \frac{1}{Z} \prod_j f_j(\mathcal{X}_{C_j})$$

where we have the following components:

- $\mathcal{X} = \{X_1, \dots, X_k\}$
- $\mathcal{X}_S = \{X_i : i \in S\}$
- j is index that indicates the factor C_j that contains all indices of variable adjacent to factor j
- f_j is factor function
- Z is normalization constant

The conditional independent is defined by $X \perp\!\!\!\perp Y|\mathcal{V}$ if every path between X and Y contains some $V \in \mathcal{V}$ (this can be shown that).

Remark 1. (Conditional Distribution) Now, if every path between X and Y contains some $V \in \mathcal{V}$, then there exists a factorization. We have the following joint distribution

$$P(X, Y, \mathcal{V}, \dots) = \frac{1}{Z} g_X(X, \mathcal{V}_X, \mathcal{Q}_X) g_Y(Y, \mathcal{V}_Y, \mathcal{Q}_Y) g_R(\mathcal{Q}_R, \mathcal{V}_R)$$

where $\mathcal{V}_X, \mathcal{V}_Y, \mathcal{V}_R \subseteq \mathcal{V}$ and the set containing $\mathcal{Q}_X, \mathcal{Q}_Y, \mathcal{Q}_R$ are disjoint. The conditional is:

$$\begin{aligned} P(X|Y, \mathcal{V}, \dots) &= \frac{P(X, Y, \mathcal{V}, \dots)}{P(Y, \mathcal{V}, \dots)} = \frac{\frac{1}{Z} g_X(X, \mathcal{V}_X, \mathcal{Q}_X) g_Y(Y, \mathcal{V}_Y, \mathcal{Q}_Y) g_R(\mathcal{Q}_R, \mathcal{V}_R)}{\sum_{X'} \frac{1}{Z} g_X(X', \mathcal{V}_X, \mathcal{Q}_X) g_Y(Y, \mathcal{V}_Y, \mathcal{Q}_Y) g_R(\mathcal{Q}_R, \mathcal{V}_R)} \\ &= \frac{g_X(X, \mathcal{V}_X, \mathcal{Q}_X)}{\sum_{X'} g_X(X', \mathcal{V}_X, \mathcal{Q}_X)} \end{aligned}$$

One the RHS doesn't depend on Y as it follows that $X \perp\!\!\!\perp Y | \mathcal{V}$.

Definition 1.4. (Markov Blanket) \mathcal{V} is markov blanket for X iff $X \perp\!\!\!\perp Y | \mathcal{V}$ for all $Y \notin \{X \cup \mathcal{V}\}$

Remark 2. Each variable X is conditionally independent of all non-neighbourhood given its neighbourhood as we have:

$$X \perp\!\!\!\perp Y | \text{ne}(X) \quad \forall Y \notin \{X \cup \text{ne}(X)\}$$

All neighbourhood $\text{ne}(X)$ is markov blanket of X . Please note that it is minimal of such set (markov blanket), which is called *markov boundary*.

Definition 1.5. (Cliques) Cliques is fully connected subgraph, whiel the maximal clique is a clique that isn't contains in the other cliques.

Definition 1.6. (Undirected Graphical Model) The undirected graphical model is a direct representation of conditional independent and nodes are connected iff they are conditionally dependent given all others. The joint probability factors over maximal clique C_j of the graph is given by:

$$P(\mathcal{X}) = \frac{1}{Z} \prod_j f_j(\mathcal{X}_{C_j})$$

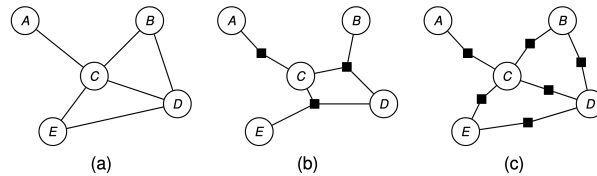
We have the following dependencies properties:

- $X \perp\!\!\!\perp Y | \mathcal{V}$ if every path between X and Y contains some node $V \in \mathcal{V}$
- Each variable X is conditionally independent of all non-neighbour node given its neighbourhood nodes:

$$X \perp\!\!\!\perp Y | \text{ne}(X) \quad Y \in \{X \cup \text{ne}(X)\}$$

And so, the neighbours is a markov blanket.

Remark 3. (Factor Graph vs Undirected Model) Consider 3 difference types of graph, we can see that each nodes has same neighbour:

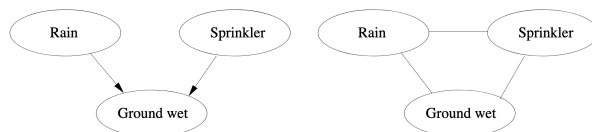


Each graph represents exactly the same conditional independent relationship. However, the maximal factorization differs, suppose we have for each variable K possible values:

- (a) can't distinguish between these (we will adopt the (b) to be safe)
- (b) has 2 three-way factor. This is represented in $\mathcal{O}(K^3)$ -size table.
- (c) has only pairwise factors. This is represented in $\mathcal{O}(K^2)$ -size table.

This means that the factor graphs have richer expressive power than undirected graphical models. But the factors can't be determined by testing for conditional independent.

Remark 4. (Limitation of Undirected Graphical Model) Undirected and Factor graph fails to capture the dependencies as the pair of variables that may be connected because they are some other variable that depends on them, for example:



If the ground is damped, it may suggest that it was rain, but if we see a sprinkler, then this explain away the damp, thus reduce the our belief of the rain into the prior. For example:

$$R \perp\!\!\!\perp S | \emptyset \quad \text{but} \quad R \perp\!\!\!\perp S | G$$

This is where there is difference between marginal and conditional independent.

Definition 1.7. (DAG Graphical Model) A directed acyclic graphical model (DAG) represents a factorization of the joint probability distribution in terms of conditional:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_{\text{pa}(i)})$$

where $\text{pa}(\cdot)$ is the parent of node i . DAG models are called Bayesian network.

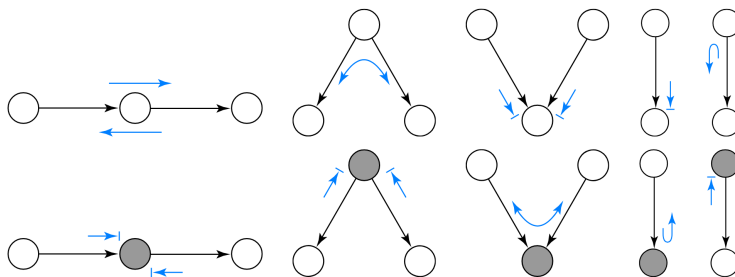
Proposition 1.1. *The conditional Independence between the graph is more complicated than the undirected graph i.e $X \perp\!\!\!\perp Y | \mathcal{V}$. If we consider every undirected path between X and Y , the path is blocked by \mathcal{V} if there is a node V on the path such that:*

- V has convergnece arrows $\rightarrow V \leftarrow$ on the path and neighbour V nor its descendent are in \mathcal{V}
- V doesn't have convergnece arrow $\leftarrow V \rightarrow$ or $\rightarrow V \rightarrow$ and $V \in \mathcal{V}$

If all paths are blocked, then \mathcal{V} is D -separated between X and Y , and so $X \perp\!\!\!\perp Y | \mathcal{V}$. Furthermore, the markov boundary to be:

$$\{\text{pa}(X) \cup \text{ch}(X) \cup \text{pa}(\text{ch}(X))\}$$

Proof. We can see that the conditional independence of the directed graphical model (for example $A \perp\!\!\!\perp B | \mathcal{D}$) can be modeled as the passing of “ball”, in which 2 variables (A, B) aren't independence if there is a way that a ball can be passed between them. We will mark the nodes in \mathcal{V} as shaded. There are 10 simple rules:



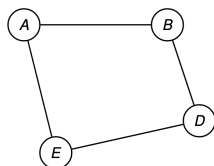
Most of the rules are straightforward to see why it is enforced.

- The first column: Both variables are separated by a middle node, then both of the are independent of each other, thus unable to pass the “ball” to each other, meaning that they are independence given the middle node. (This represents the divergence arrow rule)

- Similar explanation can be done in the second column of the image. (This also represents the divergence arrow rule)
- For the third column (explaining away), we can see that both of the nodes are independence given nothing, however, they becomes dependence once the middle node is shown. This is a reflection of the convergnece arrow rule.
- Finally, the last 2 columns are boundary rule, which is also straightforward to see why it is enforced.

Thus, we have the reason why the rules above are used. □

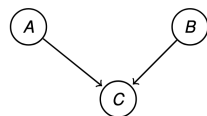
Remark 5. (Differences Between DAG and Factor Graph) There are some types of graphs that DAG can represent its probability distribution, which is:



This is the only graph that that DAG can't represent. This is because there will always be 2 non-adjacent parent sharing the same child, which implies that:

- The variables are dependence in DAG
- But independence in undirected graph.

On the other hand, no undirected or factor graph can represent the following DAG and only these:



This follows from the previous analysis on the explaining away and marginal independence.

Definition 1.8. (Family of Distribution) Each graph \mathcal{G} implies a set of conditional independence statement: $\mathcal{C}(\mathcal{G}) = \{X_i \perp\!\!\!\perp Y_i | \mathcal{Y}\}$. Each set \mathcal{C} defines a family of distribution that satisfies all statement in \mathcal{C} :

$$P_{\mathcal{C}(\mathcal{G})} = \{P(\mathcal{X}) : P(X_i, Y_i | \mathcal{V}) = P(X_i | \mathcal{V}_i)P(Y_i | \mathcal{Y}_i) \text{ for all } X_i \perp\!\!\!\perp Y_i | \mathcal{V}_i \text{ in } \mathcal{C}\}$$

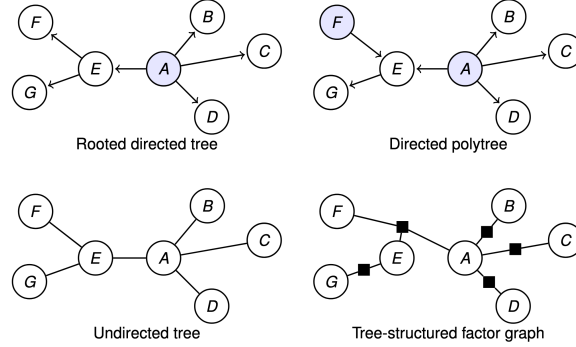
Similarly, we have family distribution in the functional form i.e:

$$P_G = \left\{ P(\mathcal{X}) : \frac{1}{Z} \prod_j f_j(\mathcal{X}_{C_j}) \text{ for some non-negative function } f_j \right\}$$

Remark 6. (Family of Distributions) We can consider the following facts:

- For directed graph: $P_G = P_{\mathcal{C}(\mathcal{G})}$
- For undirected graph: $P_G = P_{\mathcal{C}(\mathcal{G})}$ if all distribution are positive i.e $P(\mathcal{X}) > 0$ for all variable of \mathcal{X}
- Factor graphs are more expressive as for every undirected graph G_1 , there is a factor graph G_2 such that: $P_{G_1} = P_{G_2}$ but not in other direction.

Adding edge implies removing conditional independency statement and thus enlarging family of distribution.
Remark 7. For the next few propositions, we will consider difference kinds of graphical models, which can be shown to interchange with each others:



Proposition 1.2. (Polytree → Tree-Structured Factor Graph) For DAG that has more than one root, we consider:

$$P(\mathcal{X}) = \prod_i P(X_i | X_{pa(i)}) = \prod_i f_i(X_{C_i})$$

where $C_i = i \cup pa(i)$ and $f(X_{C_i}) = P(X_i | X_{pa(i)})$

Proposition 1.3. (Undirected Tree → Factor Graph) Since all undirected tree have maximal clique of size 2, it is equivalent to factor graph with pairwise factor:

$$P(\mathcal{X}) = \frac{1}{Z} \prod_{\text{edge}(i,j)} f_{(i,j)}(X_i, X_j)$$

Proposition 1.4. (Rooted Directed Tree → Undirected Tree) The distribution for single rooted directed tree can be written as a product of pairwise factor of the undirected tree:

$$P(\mathcal{X}) = P(X_r) \prod_{i \neq r} P(X_i | X_{pa(i)}) = \prod_{\text{edge}(i,j)} f_{(i,j)}(X_i, X_j)$$

Proposition 1.5. (Undirected Tree → Rooted Directed Tree) Choose arbitrary node X_r and set it as a root, which we will point every arrow away from it. Compute the conditional in the DAG as:

$$P(\mathcal{X}) = P(X_r) \prod_{i \neq r} P(X_i | X_{pa(i)}) = P(X_r) \prod_{i \neq r} \frac{P(X_i, X_{pa(i)})}{P(X_{pa(i)})} = \frac{\prod_{\text{edge}(ij)} P(X_i, X_j)}{\prod_{\text{nodes}(i)} P(X_i)^{\text{deg}(i)-1}}$$

2 Belief Propagation

Remark 8. We want to calculate the marginal distribution of a single node $P(X_i)$ and on the edge $P(X_i, X_j)$ based on the undirected graphical model, which we would need to use a belief propagation.

Proposition 2.1. (Marginal Distribution) The marginal distribution can be calculated locally as:

$$P(X_i) = \prod_{X_j \in \text{ne}(i)} M_{j \rightarrow i} \quad \text{where} \quad M_{j \rightarrow i} = \sum_{X_j} f_{ij}(X_i, X_j) \prod_{X_k \in \text{ne}(X_j) \setminus \{X_i\}} M_{k \rightarrow j}(X_j)$$

Proof. For each neighbourhood X_j of X_i define a disjoint subtree $T_{j \rightarrow i}$ that split by X_i :

$$\begin{aligned}
P(X_i) &= \sum_{\mathcal{X} \setminus \{X_i\}} P(\mathcal{X}) \propto \sum_{\mathcal{X} \setminus \{X_i\}} \prod_{(i,j) \in \mathcal{E}_T} f_{ij}(X_i, X_j) \\
&= \sum_{\mathcal{X} \setminus \{X_i\}} \prod_{X_j \in \text{ne}(X_i)} f_{ij}(X_i, X_j) \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) \\
&= \prod_{X_j \in \text{ne}(i)} \left(\sum_{\mathcal{X}_{T_{j \rightarrow i}}} f_{ij}(X_i, X_j) \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) \right)
\end{aligned}$$

The last equality comes from the splitting between subtrees, as they are disjoint. $\mathcal{X}_{T_{j \rightarrow i}}$ denotes each edge in the subtree. Let's consider the message:

$$\begin{aligned}
M_{j \rightarrow i} &= \sum_{\mathcal{X}_{T_{j \rightarrow i}}} f_{ij}(X_i, X_j) \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) \\
&= \sum_{X_j} f_{ij}(X_i, X_j) \sum_{\mathcal{X}_{T_{j \rightarrow i}} \setminus \{X_j\}} \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) \\
&= \sum_{X_j} f_{ij}(X_i, X_j) \prod_{X_k \in \text{ne}(X_j) \setminus X_i} M_{k \rightarrow i}(X_j)
\end{aligned}$$

The second equality comes from the fact that the factor of X_i is the root of the tree and for any nodes that connection to the factors. If we consider:

$$\begin{aligned}
\sum_{\mathcal{X}_{T_{j \rightarrow i}} \setminus \{X_j\}} \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) &\propto P_{T_{j \rightarrow i}}(X_j) \\
&\propto \prod_{X_k \in \text{ne}(X_j) \setminus X_i} M_{k \rightarrow i}(X_j)
\end{aligned}$$

This is due to recursive property of the message passing. □

Proposition 2.2. (Pairwise Marginal)

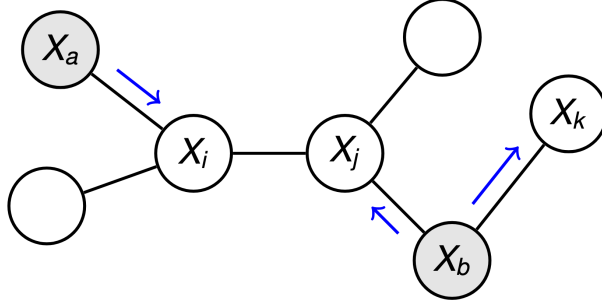
$$P(X_i, X_j) = f_{ij}(X_i, X_j) \prod_{X_k \in \text{ne}(X_j) \setminus \{X_i\}} M_{k \rightarrow j}(X_j) \prod_{X_k \in \text{ne}(X_i) \setminus \{X_j\}} M_{k \rightarrow i}(X_i)$$

Proof. We consider:

$$\begin{aligned}
P(X_i, X_j) &= \sum_{\mathcal{X} \setminus \{X_i, X_j\}} P(\mathcal{X}) \propto \sum_{\mathcal{X} \setminus \{X_i, X_j\}} \prod_{(i,j) \in \mathcal{E}_T} f_{ij}(X_i, X_j) \\
&= \sum_{\mathcal{X} \setminus \{X_i, X_j\}} f_{ij}(X_i, X_j) \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) \prod_{(i',j') \in \mathcal{E}_{T_{i \rightarrow j}}} f_{i'j'}(X_{i'}, X_{j'}) \\
&= f_{ij}(X_i, X_j) \left(\sum_{\mathcal{X}_{T_{j \rightarrow i}}} \prod_{(i',j') \in \mathcal{E}_{T_{j \rightarrow i}}} f_{i'j'}(X_{i'}, X_{j'}) \right) \left(\sum_{\mathcal{X}_{T_{i \rightarrow j}}} \prod_{(i',j') \in \mathcal{E}_{T_{i \rightarrow j}}} f_{i'j'}(X_{i'}, X_{j'}) \right) \\
&= f_{ij}(X_i, X_j) \prod_{X_k \in \text{ne}(X_j) \setminus \{X_i\}} M_{k \rightarrow j}(X_j) \prod_{X_k \in \text{ne}(X_i) \setminus \{X_j\}} M_{k \rightarrow i}(X_i)
\end{aligned}$$

□

Remark 9. (Belief Propagation for Inference) Let's consider the belief propagation for inference on the following graphical model:



To compute the $P(X_i|X_a = a)$ as we have the following message:

$$M_{a \rightarrow i} = f_{ai}(X_a = a, X_i)$$

Please note that computing $P(X_i)$ requires that $M_{a \rightarrow i} = \sum_{X_a} f_{ai}(X_a, X_i)$. For the internal node that partition the graph, like variable like X_b , we have the following message:

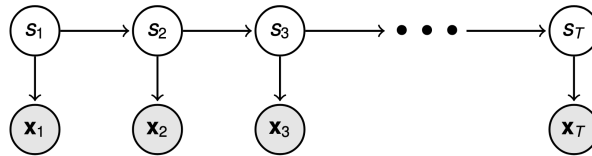
$$M_{b \rightarrow j} = f_{bj}(X_b = b, X_j) \quad M_{b \rightarrow k} = f_{bk}(X_b = b, X_k)$$

Please note that $M_{i \rightarrow j}$ are proportional to likelihood based on any observed variable (within message subtree) and possibly scaled by the prior (depends on factorization):

$$M_{i \rightarrow j}(X_j) \propto P(\mathcal{X}_{T_{i \rightarrow j}} \cap \mathcal{O} | X_j) P(X_j)$$

If we consider the message to the observed node, then we have the likelihood. Keepin all messages unnormlize and any marginal, then the normalizer is the likelihood.

Remark 10. (BP Latent Chain Model) We consider the belief propagation in the latent chain model, which is a rooted directed tree, as we have the following graphical model:



We use the backward-forward algorithm is just belief propagation on this graph:

- $\alpha_t(i) = M_{s_{t-1} \rightarrow s_t}(s_t = i) \propto P(x_{1:t}, s_t)$
- $\beta_t(i) = M_{s_{t+1} \rightarrow s_t}(s_t = i) \propto P(x_{t+1:T} | s_t)$

where we can easily see that:

$$\alpha_t(i)\beta_t(i) = \prod_{j \in \text{ne}(s_t)} M_{j \rightarrow s_t}(s_t = i) \propto P(s_t = i | \mathcal{O})$$

The algorithm like BP extend the power of graphical model beyond just encoding of independent and factorization. A single derivation surves multiple models.

2.1 Junction Tree

Remark 11. (Inference on Graph) The graphical model sometimes, which is represented by an undirected graph isn't a tree as we would like to find the marginal probability of the single value. There are several strategies that we can use:

- Propagate the local message anyway and hope for the best. This is called “loopy belief propagation”, which is an approximation technique.
- Grouping the variable together with multi-variable nodes until the resulting graph, which is a tree as we are going consider this.

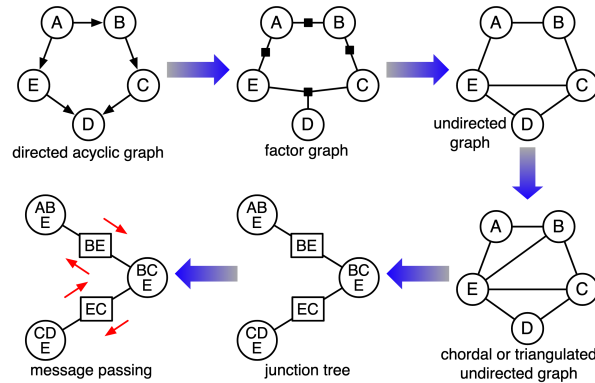
Remark 12. (Transforming the Graph) Consider to transform the graph into one that is easier to handle. As the original graph G encodes a distribution $P(\mathcal{X})$ with a certain factorization or independent structure:

- Transformation from G into an easy to handle G' as it will be valid if $P(\mathcal{X})$ can be represented by G'
- Ensuring this, we need every step of the graph transformation only to remove conditional independence (never adding them).
- Making the family of possible encoding distribution groups grows or stay the same at each step.
- The factor potential on the new graph G' are built from those given on G as to make sure it encodes the same distribution.

Definition 2.1. (Junction Tree) A junction tree is a tree whose node and edges are labelled with set of variables.

- Each node is represented by a cliques, edges are labelled by intersection of cliques called separator.
- The cliques contains all adjacent separator.
- Furthermore, if 2 cliques contain variable X , all cliques and separator on the path between 2 cliques must contain X

Definition 2.2. (Constructing Junction Tree) We consider the following step that transform DAG into a junction tree, which is shown in the figure below:



There are many process that we have to perform.

DAG to Factor Graph: The factor can be see as the conditional distribution of DAG via:

$$P(\mathcal{X}) = \prod_i P(X_i | X_{\text{pa}(i)}) = \prod_i f_i(X_{C_i})$$

where $C_i = i \cup \text{pa}(i)$ and $f_i(X_{C_i}) = P(X_i | X_{\text{pa}(i)})$. Marginal distribution on roots $P(X_r)$ absorbed into an adjacent factor.

Observation in Factor Graph: Usually, inference target a posterior marginal given a set of observed values $P(X_l | \mathcal{O})$; for example, $P(A | D = \text{wet}, C = \text{rain})$. Modify the factor linked to the observed node, or add single factor adjacent to the observed node:

$$f_D(D) = \begin{cases} 1 & \text{if } D = \text{wet} \\ 0 & \text{otherwise} \end{cases} \quad f_C(C) = \begin{cases} 1 & \text{if } C = \text{rain} \\ 0 & \text{otherwise} \end{cases}$$

Factor Graph to Undirected Graph: The transformation from DAG to undirected graph is called moralization, where: Marry all parents of each node by adding an edge to connect them. We also drop arrows on all edges.

Triangulate the Undirected Graph: Want every factor of DAG must be contained within a maximal cliques of the undirected graph. We will have to perform the following modification:

- Replace each factor by an undirected cliques.
- Construct the potential on each maximal clique by multiplying together factor potential that fall on it, and ensure each factor potential only appear once.

To do this, we shall modify the graph as:

- We join the loop into cliques, which can be very inefficient.
- Triangulation is performed to add edges to graph, so that every loop of size ≥ 4 has at least one chord. We will adding it recursively to ensure that the loop ≥ 4 has chords too.
- The graph where every loop of size ≥ 4 has at least one chord is called chordal or triangulated.
- Adding the edge removes conditional independencies, which enlarge the family of distribution.

To find the triangulation is NP-complete problem, so we resort to heuristic called variable elimination. Let's consider the order of elimination as we have:

$$\begin{aligned}
 P(X_{(n)}) &= \sum_{X_{\sigma(n-1)}} \cdots \sum_{X_{\sigma(1)}} P(\mathcal{X}) \\
 &= \frac{1}{Z} \sum_{X_{\sigma(n-1)}} \cdots \sum_{X_{\sigma(1)}} \prod_i f_i(\mathcal{X}_{C_i}) \\
 &= \frac{1}{Z} \sum_{X_{\sigma(n-1)}} \cdots \sum_{X_{\sigma(2)}} \prod_{j:C_j \ni \sigma(2)} f_j(\mathcal{X}_{C_j}) \sum_{X_{\sigma(1)}} \prod_{i:C_i \ni \sigma(1)} f_i(\mathcal{X}_{C_i}) \\
 &= \frac{1}{Z} \sum_{X_{\sigma(n-1)}} \cdots \sum_{X_{\sigma(2)}} \prod_{j:C_j \ni \sigma(2)} f_j(\mathcal{X}_{C_j}) f_{\text{new}}(\mathcal{X}_{\text{new}})
 \end{aligned}$$

Please note that C_{new} is the neighbour of $X_{\sigma(1)}$ and *edges are added to graph connected all nodes in C_{new}* :

- The graph including of all edges would be induced by elimination is chordal.
- Finding a good triangulation depends on finding a good order of elimination $\sigma(1), \dots, \sigma(n)$.
- It is NP-complete to find the best heuristic, as there are 2 ways that we pick the next variable to eliminate as follows:
 - Min-Deficiency Search: Choose variable that induces the fewest new edge.
 - Max-Cardinal Search: Choose node with most previous visited neighbour.

In most experiments, min-deficiency search seem empirically be better.

Chordal Graph to Junction Tree: To build a junction tree, we follows the procedure as we have:

- Find the maximal clique C_1, \dots, C_k of the chordal undirected graph.
- Create a weighted graph, which nodes are labelled by maximal cliques and edges connected each pair of cliques that shares variables.

- Create an edge with size of separator as we find maximal weight spanning tree of weighted graph.

Thus, we have the junction tree.

Remark 13. (Junction Tree Figures) The joint distribution factors over junction tree is:

$$P(\mathcal{X}) = \frac{1}{Z} \prod_i f_i(X_{C_i}) = \dots f_{ABC}(A, B, C) f_{BCD}(B, C, D) \dots$$

This violates the usual undirected tree semantics of factor per edge, and so we add the following constraints:

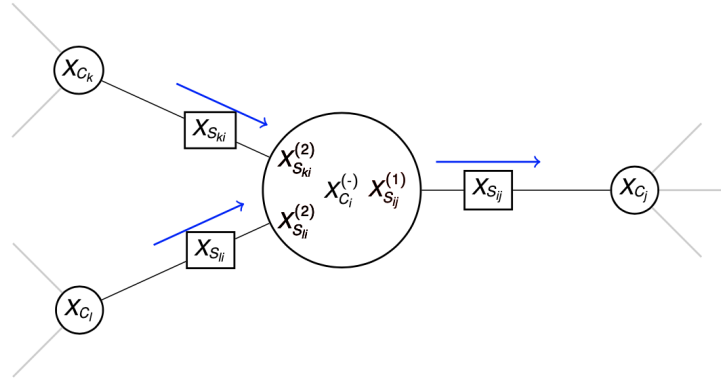
- Introducing the copy of the same variable so that there is no overlaps.
- Adding new delta function that enforce consistency:

$$P(\mathcal{X}) = \dots f_{ABC}(A, B^{(1)}, C^{(1)}) \underbrace{\delta(B^{(1)} - B^{(2)}) \delta(C^{(1)}, C^{(2)})}_{f_{\text{sep}}(B^{(1)}, C^{(1)}, B^{(2)}, C^{(2)})} f_{BCD}(B^{(2)}, C^{(2)}, D) \dots$$

Having a new message passing the junction to be:

- Unshared Variable $X_{C_i}^{(-)} = X_{C_i \setminus \cup S_{ik}}$
- Variable incoming separator $X_{S_{ki}}^{(2)}$ (Same as matching variable $X_{S_{ki}}^{(1)}$ in $k \in \text{ne}(i) \setminus j$)
- Variable outgoing separator $X_{S_{ij}}^{(1)}$ (Same as matching variable $X_{S_{ij}}^{(2)}$ in clique j)
- The variable that appear in more than one separator will need additional copies.

The overall process is shown in the following junction tree figure:



Definition 2.3. (Shafer-Shenopy) We have the following formula of the message:

$$\begin{aligned} M_{i \rightarrow j}(X_{S_{ij}}^{(2)}) &= \sum_{X_{C_i}^{(-)}, \{X_{S_{ki}}^{(2)}\}, S_{S_{ij}}^{(+)}} f_i \left(X_{C_i}^{(-)}, \{X_{S_{ki}}^{(2)}\}, S_{S_{ij}}^{(+)} \right) f_{ij} \left(X_{S_{ij}}^{(1)}, X_{S_{ij}}^{(2)} \right) \prod_{k \in \text{ne}(i) \setminus j} M_{k \rightarrow i}(X_{S_{ij}}^{(2)}) \\ &= \sum_{X_{C_i} \setminus S_{ij}} f_i(X_{C_i}) \prod_{k \in \text{ne}(i) \setminus j} M_{k \rightarrow i}(X_{S_{ij}}) \end{aligned}$$

The marginal distribution on cliques and separator are defined by:

$$P(X_{C_i}) \propto f_i(X_{C_i}) \prod_{k \in \text{ne}(C_i)} M_{k \rightarrow i}(X_{S_{ki}}) \quad P(X_{S_{ij}}) \propto M_{i \rightarrow j}(X_{S_{ij}}) M_{j \rightarrow i}(X_{S_{ij}})$$

Remark 14. (Junction Tree Properties) The running intersection property and tree structure of the junction tree implies that local consistency between cliques and separator marginal guarantee global consistency: If we consider the distribution $q_i(X_{C_i})$ and $r_{ij}(X_{S_{ij}})$ are distribution such that:

$$\sum_{X_{C_i \setminus S_{ij}}} q_i(X_{C_i}) = r_{ij}(X_{S_{ij}})$$

This implies that that joint distribution to be:

$$P(\mathcal{X}) = \frac{\prod_{\text{cliques } i} q_i(X_{C_i})}{\prod_{\text{separator } (ij)} r_{ij}(X_{S_{ij}})}$$

As we have:

$$q_i(X_{C_i}) = \sum_{\mathcal{X} \setminus X_{C_i}} P(\mathcal{X}) \quad r_{ij}(X_{S_{ij}}) = \sum_{\mathcal{X} \setminus X_{S_{ij}}} P(\mathcal{X})$$

Definition 2.4. (Hugin Update) Let's start by initializing the variables to be:

$$q_i(X_{C_i}) \propto f_i(X_{C_i}) \quad r_{ij}(X_{S_{ij}}) \propto 1$$

A Hugin propagation update for $i \rightarrow j$ is given by:

$$r_{ij}^{\text{new}} = \sum_{X_{C_i \setminus S_{ij}}} q_i(X_{C_i}) \quad q_j^{\text{new}}(X_{C_j}) = q_j(X_{C_j}) \frac{r_{ij}^{\text{new}}(X_{S_{ij}})}{r_{ij}(X_{S_{ij}})}$$

Setting the correct marginalization locally for the first update. For the second update, we change the q based on the keeping the joint probability $P(\mathcal{X})$

3 Bayes and Gaussian Processes

3.1 Bayes Method

Remark 15. (Recap of Bayes) Model have a parameter θ_m that specify the probability of data $P(\mathcal{D}|\theta_m, m)$. If the model is known, learning θ_m means finding a posterior or point estimate.

- What if we want to learn the type of model too ?
- Can we combine the model into a single “supermodel” with a composite parameters ?
- We can separate the model selection step: $P(\theta_m, m|\mathcal{D}) = P(\theta_m|m, \mathcal{D})P(m|\mathcal{D})$

Definition 3.1. (Neyman-Pearson Hypothesis Testing) For nested model, starting with simplest model with $m = 1$, comparing null hypothesis m to its alternate $m + 1$ and repeat until $m + 1$ is rejected. Note that this tests only exact when it is asymptotic in data number. Finally, it is conservative as it is asymmetric by design.

Definition 3.2. (Likelihood Validation) Partition data into disjoint training and validation set, where $D = D_{\text{tr}} \cup D_{\text{valid}}$. Then, we choose a model with greatest $P(D_{\text{valid}}|\theta_m^{\text{ML}})$ where:

$$\theta_m^{\text{ML}} = \arg \max_{\theta} P(D_{\text{tr}}|\theta)$$

or if we are able to find $P(D_{\text{valid}}|D_{\text{train}}, m)$.

- This model is consistent and select the most useful model even if they are incorrect.

- It may be biased toward a simpler models and gives a high-variance.
- We can use cross-validation is used with multiple partition and average likelihood.

Definition 3.3. (Baysian Model Selection) We would like to choose a model $P(m|\mathcal{D})$ where consistent as it is a probability principle if true model is in the set.

- However, it might be a problem of assumed prior. Finally, the posterior can weights models for combined prediction.
- A model of class m is a set of distributions parameterized by θ_m . The model implies the prior over parameters as the posterior of the parameter:

$$P(\theta_m|\mathcal{D}, m) = \frac{P(\mathcal{D}|\theta_m, m)P(\theta_m|m)}{P(\mathcal{D}|m)} \quad \text{where}$$

where $P(\mathcal{D}|m)$ is the Baysian evidence for model m , where the ratio is known as Bayes factor:

$$\frac{P(\mathcal{D}|m)}{P(\mathcal{D}|m')} = \frac{P(m|\mathcal{D})}{P(m')} \frac{P(m')}{P(m)}$$

- This is linked to Occam's razor where: the model that are *too complex* can generate many data but they can be unlikely to generate a particular dataset. The model that are *too simple* can't generate the data.

Remark 16. (Conjugate Prior) We will recall the use of exponential model. Suppose, we have $P(\mathcal{D}|\theta_m, m)$ is member of the exponential family, where we can see that:

$$P(\mathcal{D}|\theta_m, m) = \prod_{i=1}^N P(\mathbf{x}_i|\theta_m, m) = \prod_{i=1}^N \exp(\mathbf{s}(\mathbf{x}_i)^T \theta_m - A(\theta_m))$$

Please note that $A(\theta_m)$ is the normalizing factor as it is (recall exponential family of the form in previous course) equal to $\ln(g(\theta_m))$. Consider the prior conjugate to be:

$$P(\theta|m) = \frac{1}{Z(\mathbf{s}_p, n_p)} \exp(\mathbf{s}_p^T \theta_m - n_p A(\theta_m))$$

Then the posterior is equal to:

$$P(\mathcal{D}, \theta_m|m) = \frac{1}{Z(\mathbf{s}_p, p)} \exp\left(\left(\sum_{i=1}^N \mathbf{s}(\mathbf{x}_i) + \mathbf{s}_p\right)^T \theta_m - (N + n_p)A(\theta_m)\right)$$

One can show that:

$$P(\mathcal{D}|m) = \int P(\mathcal{D}, \theta_m|m) d\theta_m = \frac{Z(\sum_i \mathbf{s}(\mathbf{x}_i) + \mathbf{s}_p, N + n_p)}{Z(\mathbf{s}_p, p)}$$

Remark 17. (Laplace Approximation) To find $P(\mathcal{D}|m)$, we will have to perform the following integration:

$$P(\mathcal{D}|m) = \int P(\mathcal{D}, \theta_m|m) d\theta_m$$

as the datasize N grows, θ becomes more concentrated as $P(\mathcal{D}, \theta_m|m) \propto P(\theta_m|\mathcal{D}, m)$ becomes concentrated

on $\boldsymbol{\theta}_m^*$. Let's try to approximate the $P(\mathcal{D}, \boldsymbol{\theta}_m | m)$ to the second order around $\boldsymbol{\theta}_m^*$:

$$\begin{aligned}
\int P(\mathcal{D}, \boldsymbol{\theta}_m | m) d\boldsymbol{\theta}_m &= \int \exp(\log P(\mathcal{D}, \boldsymbol{\theta}_m | m)) d\boldsymbol{\theta}_m \\
&\approx \int \exp \left[\log P(\mathcal{D}, \boldsymbol{\theta}_m^* | m) + \underbrace{\nabla \log P(\mathcal{D}, \boldsymbol{\theta}_m^* | m)}_0 (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^*) \right. \\
&\quad \left. + \frac{1}{2} (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^*)^T \underbrace{\nabla^2 \log P(\mathcal{D}, \boldsymbol{\theta}_m^* | m)}_{-\mathbf{A}} (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^*) \right] d\boldsymbol{\theta}_m \\
&= \int P(\mathcal{D}, \boldsymbol{\theta}_m^* | m) \exp \left[-\frac{1}{2} (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^*)^T \mathbf{A} (\boldsymbol{\theta}_m - \boldsymbol{\theta}_m^*) \right] d\boldsymbol{\theta} \\
&= P(\mathcal{D} | \boldsymbol{\theta}_m^*, m) P(\boldsymbol{\theta}_m^* | m) (2\pi)^{-d/2} |\mathbf{A}|^{-1/2}
\end{aligned}$$

This is approximating the posterior by a Gaussian, where an approximate that is asymmetrically correct.

Definition 3.4. (Bayesian Information Criterion) BIC can be obtained from Laplace approximate:

$$\log P(\mathcal{D} | m) \approx \log P(\boldsymbol{\theta}^* | m) + \log P(\mathcal{D} | \boldsymbol{\theta}_m^*, m) + \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{A}|$$

where we further have:

$$\mathbf{A} = -\nabla^2 \log P(\mathcal{D}, \boldsymbol{\theta}_m^* | m) = -\nabla^2 \log P(\mathcal{D} | \boldsymbol{\theta}^*, m) - \nabla^2 \log P(\boldsymbol{\theta}^* | m)$$

As $N = |\mathcal{D}| \rightarrow \infty$ and $\mathbf{A} \rightarrow N\mathbf{A}_0 + \text{const}$ for a fixed positive definite matrix $\mathbf{A}_0 = \langle -\nabla^2 \log P(\mathbf{x} | \boldsymbol{\theta}^*, m) \rangle$ as we have $\log |N\mathbf{A}_0| = d \log N + \log |\mathbf{A}_0|$. We will retain only the term that grows with N to be:

$$\log P(\mathcal{D} | m) \approx \log P(\mathcal{D} | \boldsymbol{\theta}_m^*, m) - \frac{d}{2} \log N$$

Remark 18. (Properties BIC) Bayesian Information Criterion has the following properties, as we have:

- Quick and Easy to compute. It doesn't depend on the prior as we can use ML to estimate instead of MAP estimate.
- It is related to minimal description length. Given the assumption that in large sample limit, all parameter are well determined. But it negated multiple nodes (permutation of mixture of Gaussian).

Definition 3.5. (Hyperparameter and Evidence Optimization) Need to choose between a family of continuous parameterized models:

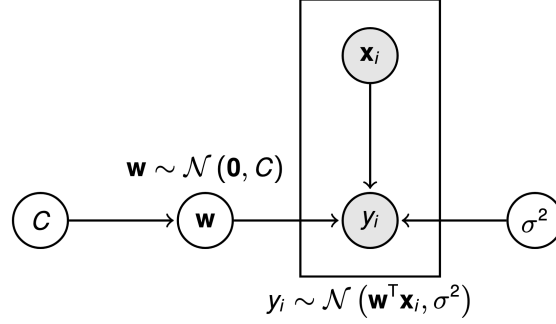
$$P(\mathcal{D} | \boldsymbol{\eta}) = \int P(\mathcal{D} | \boldsymbol{\theta}) P(\boldsymbol{\theta} | \boldsymbol{\eta}) d\boldsymbol{\theta}$$

We can perform an ascending on gradient in: the exact evidence, approximate evidence (Laplace, BIC) or free-energy based on the evidence (variational Bayes). Performing the hyper-prior on $\boldsymbol{\eta}$, which we can sample its posterior via MCMC as:

$$P(\boldsymbol{\eta} | \mathcal{D}) = \frac{P(\mathcal{D} | \boldsymbol{\eta}) P(\boldsymbol{\eta})}{P(\mathcal{D})}$$

3.2 Linear Regression/Gaussian Process

Definition 3.6. (Linear Regression) We consider the following graphical model as:



To find the hyperparameters, we can see that:

$$P(\{y_i\}_{i=1}^N | \{\mathbf{x}_i\}_{i=1}^N, \mathbf{C}, \sigma^2) = \int P(\{y_i\}_{i=1}^N | \{\mathbf{x}_i\}_{i=1}^N, \mathbf{w}, \sigma^2) P(\mathbf{w} | \mathbf{C}) d\mathbf{w}$$

Maximizing the value of \mathbf{C} and σ^2 .

Proposition 3.1. *The update rule for the parameter $\theta = \{\mathbf{C}, \sigma^2\}$ is:*

$$\begin{aligned} \frac{\partial}{\partial \theta} \log \mathcal{E}(\mathbf{C}, \sigma^2) &= \frac{1}{2} \text{Tr} \left[(\mathbf{C} - \Sigma_w - \bar{\mathbf{w}} \bar{\mathbf{w}}^T) \frac{\partial}{\partial \theta} \mathbf{C}^{-1} \right] \\ \frac{\partial}{\partial \sigma^2} \log \mathcal{E}(\mathbf{C}, \sigma^2) &= \frac{1}{\sigma^2} \left(-N + \text{Tr}[\mathbf{I} - \Sigma_w \mathbf{C}^{-1}] + \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}^T \bar{\mathbf{w}})^T (\mathbf{Y} - \mathbf{X}^T \bar{\mathbf{w}}) \right) \end{aligned}$$

Proof. We can see that the posterior of $P(\mathbf{w} | \{y_i\}_{i=1}^N, \{\mathbf{x}_i\}_{i=1}^N, \mathbf{C}, \sigma^2)$. We can see that the posterior on \mathbf{w} is normal with:

$$\Sigma_w = \left(\frac{\mathbf{X} \mathbf{X}^T}{\sigma^2} + \mathbf{C}^{-1} \right)^{-1} \quad \bar{\mathbf{w}} = \Sigma_w \frac{\mathbf{X} \mathbf{Y}}{\sigma^2}$$

This follows from the earlier works. Now, we can consider terms inside the exponential of the joint distributions $P(\mathbf{Y} | \mathbf{X}, \mathbf{w}, \sigma^2) P(\mathbf{w} | \mathbf{C})$:

$$\begin{aligned} & -\frac{1}{2\sigma^2} \mathbf{Y}^T \mathbf{Y} + \frac{1}{\sigma^2} \mathbf{w}^T \mathbf{X} \mathbf{Y} - \frac{1}{2} [\mathbf{w}^T \Sigma_w^{-1} \mathbf{w}] \\ &= -\frac{1}{2\sigma^2} \mathbf{Y}^T \mathbf{Y} - \frac{1}{2} \left\{ \mathbf{w}^T \Sigma_w^{-1} \mathbf{w} - \frac{2}{\sigma^2} \mathbf{w}^T \Sigma_w \Sigma_w^{-1} \mathbf{X} \mathbf{Y} \right\} \\ &= -\frac{1}{2\sigma^2} \mathbf{Y}^T \mathbf{Y} - \frac{1}{2} \left\{ \mathbf{w}^T \Sigma_w^{-1} \mathbf{w} - \frac{2}{\sigma^2} \mathbf{w}^T \Sigma_w \Sigma_w^{-1} \mathbf{X} \mathbf{Y} + \frac{1}{\sigma^4} \mathbf{Y}^T \mathbf{X}^T \Sigma_w^T \Sigma_w^{-1} \Sigma_w^T \mathbf{X} \mathbf{Y} \right\} \\ & \quad + \frac{1}{2\sigma^4} \mathbf{Y}^T \mathbf{X}^T \Sigma_w^T \mathbf{X} \mathbf{Y} \\ &= -\frac{1}{2\sigma^2} \mathbf{Y}^T \mathbf{Y} + \frac{1}{2\sigma^4} \mathbf{Y}^T \mathbf{X}^T \Sigma_w^T \mathbf{X} \mathbf{Y} - \frac{1}{2} \left(\mathbf{w} - \frac{1}{\sigma^2} \Sigma_w \mathbf{X} \mathbf{Y} \right)^T \Sigma_w^{-1} \left(\mathbf{w} - \frac{1}{\sigma^2} \Sigma_w \mathbf{X} \mathbf{Y} \right) \end{aligned}$$

If we consider the quadratic terms, and the integration over it we have:

$$\int \exp \left\{ -\frac{1}{2} \left(\mathbf{w} - \frac{1}{\sigma^2} \Sigma_w \mathbf{X} \mathbf{Y} \right)^T \Sigma_w^{-1} \left(\mathbf{w} - \frac{1}{\sigma^2} \Sigma_w \mathbf{X} \mathbf{Y} \right) \right\} d\mathbf{w}$$

This is the Gaussian integration, which it is a constant that we can ignore. This leads to the evidence to be (the normalizing factor can be easily found):

$$\mathcal{E}(\mathbf{C}, \sigma^2) = \sqrt{\frac{|2\pi \Sigma_w|}{|2\pi \sigma^2 \mathbf{I}| |2\pi \mathbf{C}|}} \exp \left\{ -\frac{1}{2} \mathbf{Y}^T \left(\frac{\mathbf{I}}{\sigma^2} - \frac{\mathbf{X}^T \Sigma_w \mathbf{X}}{\sigma^4} \right) \mathbf{Y} \right\}$$

Now, let's consider the derivative to be, starting with the value of \mathbf{C} , we can see that the relevant values:

$$\begin{aligned}
& \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\theta}} \log |\boldsymbol{\Sigma}_w| - \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\theta}} \log |\mathbf{C}| + \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\theta}} \frac{\mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y}}{\sigma^4} \\
&= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\theta}} |\boldsymbol{\Sigma}_w^{-1}| + \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\theta}} \log |\mathbf{C}^{-1}| + \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\theta}} \frac{\mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y}}{\sigma^4} \\
&= -\frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) + \frac{1}{2} \text{Tr} \left(\mathbf{C} \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) + \frac{1}{2\sigma^4} \text{Tr} \left(\mathbf{X} \mathbf{Y} \mathbf{Y}^T \mathbf{X}^T \frac{\partial \boldsymbol{\Sigma}_w^{-1}}{\partial \boldsymbol{\theta}} \right) \\
&= -\frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) + \frac{1}{2} \text{Tr} \left(\mathbf{C} \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) \\
&\quad - \frac{1}{2\sigma^4} \text{Tr} \left(\mathbf{X} \mathbf{Y} \mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w^{-1} \frac{\partial}{\partial \boldsymbol{\theta}} \left[\frac{\mathbf{X} \mathbf{X}^T}{\sigma^2} + \mathbf{C}^{-1} \right] \boldsymbol{\Sigma}_w \right) \\
&= -\frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) + \frac{1}{2} \text{Tr} \left(\mathbf{C} \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) \\
&\quad - \frac{1}{2\sigma^4} \text{Tr} \left(\mathbf{X} \mathbf{Y} \mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \left[\frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right] \boldsymbol{\Sigma}_w \right) \\
&= -\frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) + \frac{1}{2} \text{Tr} \left(\mathbf{C} \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) \\
&\quad - \frac{1}{2} \text{Tr} \left(\frac{\mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w}{\sigma^2} \left[\frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right] \frac{\boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y}}{\sigma^2} \right) \\
&= -\frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) + \frac{1}{2} \text{Tr} \left(\mathbf{C} \frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right) - \frac{1}{2} \text{Tr} \left(\bar{\mathbf{w}} \bar{\mathbf{w}}^T \left[\frac{\partial \mathbf{C}^{-1}}{\partial \boldsymbol{\theta}} \right] \right)
\end{aligned}$$

Please recall the derivative of the inverse, and we finish the proof. Now consider the derivative of the evidence with respect to $S = \sigma^2$,

$$\underbrace{\frac{1}{2} \frac{\partial}{\partial S} \log |\boldsymbol{\Sigma}_w|}_{\textcircled{1}} - \frac{N}{2} \frac{\partial}{\partial S} \log S - \frac{1}{2} \frac{\partial}{\partial S} \frac{\mathbf{Y}^T \mathbf{Y}}{S} + \underbrace{\frac{1}{2} \frac{\partial}{\partial S} \frac{\mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y}}{S^2}}_{\textcircled{2}}$$

And we have the following derivative for $\textcircled{1}$, as we have:

$$\begin{aligned}
\frac{\partial}{\partial S} \log |\boldsymbol{\Sigma}_w| &= -\frac{\partial}{\partial S} \log |\boldsymbol{\Sigma}_w^{-1}| \\
&= -\text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial \boldsymbol{\Sigma}_w^{-1}}{\partial S} \right) \\
&= -\text{Tr} \left(\boldsymbol{\Sigma}_w \frac{\partial}{\partial S} \left[\frac{\mathbf{X} \mathbf{X}^T}{S} + \mathbf{C}^{-1} \right] \right) \\
&= \frac{1}{S^2} \text{Tr} \left(\boldsymbol{\Sigma}_w \mathbf{X} \mathbf{X}^T \right)
\end{aligned}$$

For for the equation $\textcircled{2}$, we have:

$$\begin{aligned}
\frac{\partial}{\partial S} \frac{\mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y}}{S^2} &= \mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y} \frac{\partial}{\partial S} \frac{1}{S^2} + \frac{1}{S^2} \frac{\partial}{\partial S} \mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y} \\
&= -2 \frac{\mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y}}{S^3} + \frac{1}{S} \mathbf{Y}^T \mathbf{X}^T \boldsymbol{\Sigma}_w \frac{\partial}{\partial S} \left[\frac{\mathbf{X} \mathbf{X}^T}{S} + \mathbf{C}^{-1} \right] \frac{1}{S} \boldsymbol{\Sigma}_w \mathbf{X} \mathbf{Y} \\
&= -2 \mathbf{Y}^T \mathbf{X}^T \bar{\mathbf{w}} + \bar{\mathbf{w}}^T \left(\mathbf{X} \mathbf{X}^T \right) \bar{\mathbf{w}}
\end{aligned}$$

Combining them and we have:

$$\begin{aligned} & \frac{1}{2S} \left(\frac{1}{S} \text{Tr} \left(\boldsymbol{\Sigma}_w \mathbf{X} \mathbf{X}^T \right) - N + \frac{1}{S} \mathbf{Y}^T \mathbf{Y} - \frac{1}{S} 2 \mathbf{Y}^T \mathbf{X}^T \bar{\mathbf{w}} + \frac{1}{S} \bar{\mathbf{w}}^T \boldsymbol{\Sigma}_w \left(\mathbf{X} \mathbf{X}^T \right) \bar{\mathbf{w}} \right) \\ & = \frac{1}{2S} \left(-N + \frac{1}{S} \text{Tr} \left(\boldsymbol{\Sigma}_w \mathbf{X} \mathbf{X}^T \right) + \frac{1}{S} (\mathbf{Y} - \mathbf{X}^T \bar{\mathbf{w}})^T (\mathbf{Y} - \mathbf{X}^T \bar{\mathbf{w}}) \right) \end{aligned}$$

Please note that there are some differences between this and the final proof, which we can't get it but most of them seem correct already \square

Definition 3.7. (Automatic Relevance Determination (ARD)) The most common form of evidence optimization with $\mathbf{C}^{-1} = \text{diag}(\boldsymbol{\alpha})$ and the optimize $\{\alpha_i\}$ setting the gradient to 0, which we now have:

$$\alpha_i^{\text{new}} = \frac{1 - \alpha_i [\boldsymbol{\Sigma}_n]_{ii}}{\bar{w}_i^2} \quad (\alpha^2)^{\text{new}} = \frac{(\mathbf{Y} - \mathbf{X}^T \bar{\mathbf{w}})^T (\mathbf{Y} - \mathbf{X}^T \bar{\mathbf{w}})}{N - \sum_i (1 - [\boldsymbol{\Sigma}_w]_{ii} \alpha_i)}$$

During the optimization, there are 2 possible scenario for $\boldsymbol{\alpha}$:

- $\alpha_i \rightarrow \infty$ where the weight $\alpha_i = 0$
- α_i is finite and so $w_i = \arg \max P(w_i | \mathbf{X}, \mathbf{Y}, \alpha_i)$

This is called Automatic Relevance Determination (ARD), which yields sparse solution that improve ML regression (like LASSO) This Evidence maximization is called ML likelihood or ML-2.

Definition 3.8. (Prediction Averaging) We integrate out the parameter, where we have the density estimation:

$$P(\mathbf{x} | \mathcal{D}, m) = \int P(\mathbf{x} | \boldsymbol{\theta}, m) P(\boldsymbol{\theta} | \mathcal{D}, m) d\boldsymbol{\theta}$$

Or, perform the prediction:

$$P(y | \mathbf{x}, \mathcal{D}, m) = \int P(y | \mathbf{x}, \boldsymbol{\theta}, m) P(\boldsymbol{\theta} | \mathcal{D}, m) d\boldsymbol{\theta}$$

This kind of prediction might resist overfitting with infinitely complex model as this can be called Bayesian non-parametric, which leads to Gaussian Process.

Proposition 3.2. (Prediction) Given a new input vector \mathbf{x} , the predicted output y is given by:

$$y | \mathbf{x} \sim \mathcal{N}(\bar{\mathbf{w}}^T \mathbf{x}, \mathbf{x}^T \boldsymbol{\Sigma}_w \mathbf{x} + \sigma^2)$$

This comes from the linear Gaussian model. The variance $\mathbf{x}^T \boldsymbol{\Sigma}_w \mathbf{x}$ that comes from the uncertainty of \mathbf{w} .

Remark 19. (Observation on Joint Distribution) Consider the joint probability between y_1, \dots, y_N and $\mathbf{x}_1, \dots, \mathbf{x}_N$, and noisy prediction with variance τ^2 , we can see that the mean and the covariances are:

$$\begin{aligned} \mathbb{E}[y_i] &= \mathbb{E}[\mathbf{w}^T \mathbf{x}_i] = 0 \\ \mathbb{E}[(y_i - \bar{y}_i)^2] &= \mathbb{E}[(\mathbf{x}_i^T \mathbf{w})(\mathbf{w}^T \mathbf{x}_i)] + \sigma^2 = \tau^2 \mathbf{x}_i^T \mathbf{x}_i + \sigma^2 \\ \mathbb{E}[(y_i - \bar{y}_i)(y_j - \bar{y}_j)] &= \mathbb{E}[(\mathbf{x}_i^T \mathbf{w})(\mathbf{w}^T \mathbf{x}_j)] = \tau^2 \mathbf{x}_i^T \mathbf{x}_j \end{aligned}$$

So this leads to the following joint distributions:

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \Bigg| \mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{bmatrix} \tau^2 \mathbf{x}_1^T \mathbf{x}_1 + \sigma^2 & \tau^2 \mathbf{x}_1^T \mathbf{x}_2 & \cdots & \tau^2 \mathbf{x}_1^T \mathbf{x}_N \\ \tau^2 \mathbf{x}_2^T \mathbf{x}_1 & \tau^2 \mathbf{x}_2^T \mathbf{x}_2 + \sigma^2 & \cdots & \tau^2 \mathbf{x}_2^T \mathbf{x}_N \\ \vdots & \vdots & \ddots & \vdots \\ \tau^2 \mathbf{x}_N^T \mathbf{x}_1 & \tau^2 \mathbf{x}_N^T \mathbf{x}_2 & \cdots & \tau^2 \mathbf{x}_N^T \mathbf{x}_N + \sigma^2 \end{bmatrix} \right)$$

If we were to consider adding the test input vector \mathbf{x} and test vector y , we have:

$$\begin{bmatrix} \mathbf{Y}^T \\ y \end{bmatrix} \Big| \mathbf{X}, \mathbf{x} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \tau \mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I} & \tau^2 \mathbf{X}^T \mathbf{x} \\ \tau^2 \mathbf{x}^T \mathbf{X} & \tau^2 \mathbf{x}^T \mathbf{x} + \sigma^2 \end{bmatrix} \right) = \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \tilde{\mathbf{K}}_{XX} & \tilde{\mathbf{K}}_{Xx} \\ \tilde{\mathbf{K}}_{xX} & \tilde{\mathbf{K}}_{xx} \end{bmatrix} \right)$$

By the linear Gaussian model, we are given, the following conditional

$$\begin{aligned} y | \mathbf{Y}, \mathbf{X}, \mathbf{x} &\sim \mathcal{N} \left(\tilde{\mathbf{K}}_{xX} \tilde{\mathbf{K}}_{XX}^{-1} \mathbf{Y}^T, \tilde{\mathbf{K}}_{xx} - \tilde{\mathbf{K}}_{xX} \tilde{\mathbf{K}}_{XX}^{-1} \tilde{\mathbf{K}}_{Xx} \right) \\ &\sim \mathcal{N} \left(\tau^2 \mathbf{x}^T \mathbf{X} (\tau^2 \mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{Y}^T, \tau^2 \mathbf{x}^T \mathbf{x} + \sigma^2 - \tau^2 \mathbf{x}^T \mathbf{X} (\tau^2 \mathbf{X}^T \mathbf{X} + \sigma^2)^{-1} \tau^2 \mathbf{X}^T \mathbf{x} \right) \\ &\sim \mathcal{N} \left(\mathbf{x}^T \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X} \mathbf{Y}^T, \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} + \sigma^2 \right) \end{aligned}$$

where $\boldsymbol{\Sigma} = (\frac{1}{\sigma^2} \mathbf{X} \mathbf{X}^T + \frac{1}{\tau^2} \mathbf{I})^{-1}$. The derivation comes from the Woodbury identity.

Remark 20. (Non-Linear Regression) We consider the mapping $\mathbf{x} \mapsto \boldsymbol{\phi}(\mathbf{x})$, the regression:

$$\mathbf{Y} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}_N, \tau^2 \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \sigma^2 \mathbf{I}_N)$$

where i -th column of matrix $\boldsymbol{\Phi}$ is $\boldsymbol{\phi}(\mathbf{x}_i)$. Now we have the following prediction:

$$\mathbf{y} | \mathbf{Y}, \mathbf{X}, \mathbf{x} \sim \mathcal{N}(\tilde{\mathbf{K}}_{xX} \tilde{\mathbf{K}}_{XX}^{-1} \mathbf{Y}^T, \tilde{\mathbf{K}}_{xx} - \tilde{\mathbf{K}}_{xX} \tilde{\mathbf{K}}_{XX}^{-1} \tilde{\mathbf{K}}_{Xx})$$

where we have:

$$\tilde{\mathbf{K}}_{XX} = \tau^2 \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \sigma^2 \mathbf{I} \quad \tilde{\mathbf{K}}_{Xx} = \tau^2 \boldsymbol{\Phi}^T \boldsymbol{\phi}(\mathbf{x}) \quad \tilde{\mathbf{K}}_{xx} = \tau^2 \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}) + \sigma^2$$

Remark 21. (Kernel Vector) Define a covariance kernel function $\tilde{K} : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ if $\mathbf{x}, \mathbf{x}' \in \mathbb{X}$ are 2 inputs vectors with corresponding output y and y' , then we have:

$$\tilde{K}(\mathbf{x}, \mathbf{x}') = \text{Cov}[y, y'] = \mathbb{E}[yy'] - \mathbb{E}[y]\mathbb{E}[y']$$

In non-linear regression problem we have $\tilde{K}(\mathbf{x}, \mathbf{x}') = \tau^2 \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}') + \sigma^2 \delta[\mathbf{x} = \mathbf{x}']$, where the prediction depends on $\tilde{\mathbf{K}}(\mathbf{x}, \mathbf{x}')$ rather than its feature map $\boldsymbol{\phi}(\mathbf{x})$ and so, we have:

$$\mathbf{Y} | \mathbf{X}, \tilde{\mathbf{K}} \sim \mathcal{N}(\mathbf{0}_N, \tilde{\mathbf{K}}_{XX})$$

To perform a prediction, as we have:

$$\mathbf{y} | \mathbf{Y}, \mathbf{X}, \mathbf{x} \sim \mathcal{N}(\tilde{\mathbf{K}}_{xX} \tilde{\mathbf{K}}_{XX}^{-1} \mathbf{Y}^T, \tilde{\mathbf{K}}_{xx} - \tilde{\mathbf{K}}_{xX} \tilde{\mathbf{K}}_{XX}^{-1} \tilde{\mathbf{K}}_{Xx})$$

where $[\tilde{\mathbf{K}}_{XX}]_{ij} = \tilde{K}(\mathbf{x}_i, \mathbf{x}_j)$, $[\tilde{\mathbf{K}}_{Xx}]_i = \tilde{K}(\mathbf{x}_i, \mathbf{X})$ and $\tilde{\mathbf{K}}_{xx} = \tilde{K}(\mathbf{x}, \mathbf{x})$, now we can have the kernel function solely without the use of the feature map.

Definition 3.9. (Gaussian Process) Let $f(\mathbf{x})$ be random variable indexed by \mathbf{x} , then drawn from the whole GP, which is a random function $f : \mathbb{X} \rightarrow \mathbb{R}$ where:

$$f(\cdot) \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot))$$

It is defined such that the finite list of points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. The joint distributions of the function value $f = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^T$ as we have $\mathbf{f} | \mathbf{X}, \mathbf{K} \sim \mathcal{N}(\mathbf{m}, \mathbf{K}_{XX})$ with the following parameters:

$$[\mathbf{K}_{XX}]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) \quad [\mathbf{m}]_i = m(\mathbf{x}_i)$$

Remark 22. (Properties of Gaussian Process + Bayes Theorem) If we assume the random function is drawn from GP prior $f(\cdot) \sim \mathcal{GP}(\mathbf{0}, K(\cdot, \cdot))$. Observation y_i taken to be noisy version of latent function $f(\mathbf{x}_i)$ where we have:

$$\mathbf{y}_i | \mathbf{x}_i, f(\cdot) \sim \mathcal{N}(y(\mathbf{x}_i), \sigma^2)$$

The we have the following quantities:

- The latent f is also Gaussian process:

$$f(\cdot)|\mathbf{X}, \mathbf{Y} \sim \mathcal{GP}\left(\mathbf{K}_{\cdot X}(\mathbf{K}_{XX} + \sigma^2\mathbf{I})^{-1}\mathbf{Y}^T, K(\cdot, \cdot) - \mathbf{K}_{\cdot X}(\mathbf{K}_{XX} + \sigma^2\mathbf{I})^{-1}\mathbf{K}_{X\cdot}\right)$$

- This is given multivariate Gaussian likelihood:

$$P(\mathbf{Y}|\mathbf{X}) = \left|2\pi(\mathbf{K}_{xx} + \sigma^2\mathbf{I})^{-1/2}\right| \exp\left(-\frac{1}{2}\mathbf{Y}(\mathbf{K}_{xx} + \sigma^2\mathbf{I})\mathbf{Y}^T\right)$$

- The posterior of f with observation noise, as we have:

$$y|\mathbf{X}, \mathbf{Y}, \mathbf{x} \sim \mathcal{N}(\mathbb{E}[f(\mathbf{x})|\mathbf{X}, \mathbf{Y}], \text{Var}[f(\mathbf{x})|\mathbf{X}, \mathbf{Y}] + \sigma^2)$$

4 Factored Variation Approximation

Remark 23. (Intractable EM-Method) Sometime E-step or M-step are intractable; therefore, we will have to do some approximation to the problem:

- It might be the case that $P(\mathcal{Z}|\mathcal{X}, \theta)$ and so we might replace M-step with gradient M-step, where every step increase the likelihood.
- In E-step, we should consider parameterized $q = q(\mathcal{Z})$ and take the gradient step in ρ .
- Or we assume some simplification for q usually $q = \prod_i q_i(\mathcal{Z}_i)$ where \mathcal{Z}_i is the partition \mathcal{Z} .

Remark 24. (Restricting Variational Class) This approximation of E-step is the same as choosing q from the limit set \mathcal{Q} such that VE-step is to maximize $\mathcal{F}(q, \theta)$ such that:

$$q^{(k)}(\mathcal{Z}) = \arg \max_{q \in \mathcal{Q}} \mathcal{F}(q, \theta^{(k-1)})$$

The effect of the restricting q to \mathcal{Q} , there is the following effect:

- The free energy is bounded by log-likelihood via Jensen's inequality as $\mathcal{F}(q, \theta) \leq l(\theta^{\text{ML}})$, thus as long as every step increases \mathcal{F} convergence still guarantee.
- Since $P(\mathcal{Z}|\mathcal{X}, \theta^{(k)})$ may not lie in \mathcal{Q} and we no longer saturated the bound after E-step. The likelihood may not increase on each full EM step.

Definition 4.1. (Factored Variational E-Step) We consider the following family of distributions:

$$\mathcal{Q} = \left\{ q \left| q(\mathcal{Z}) = \prod_i q_i(\mathcal{Z}_i) \right. \right\}$$

where \mathcal{Z} are partitioned into $\{\mathcal{Z}_i\}$. We consider the following maximization of the E-step:

$$q_i^{(k)}(\mathcal{Z}_i) = \arg \max_{q_i(\mathcal{Z}_i)} \mathcal{F} \left(q_i(\mathcal{Z}_i) \prod_{j \neq i} q_j(\mathcal{Z}_j), \theta^{(k-1)} \right)$$

Proposition 4.1. *We can show that:*

$$q_i(\mathcal{Z}_i) \propto \exp \left\langle \log P(\mathcal{X}, \mathcal{Z} | \theta^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(\mathcal{Z}_j)}$$

Proof. Given this, the free energy is given by:

$$\begin{aligned} \mathcal{F} \left(q_i(\mathcal{Z}_i) \prod_{j \neq i} q_j(\mathcal{Z}_j), \boldsymbol{\theta}^{(k-1)} \right) &= \left\langle \log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}^{(k-1)}) \right\rangle_{\prod_j q_j(\mathcal{Z}_j)} + H \left[\prod_j q_j(\mathcal{Z}_j) \right] \\ &= \int q_i(\mathcal{Z}_i) \left\langle \log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(\mathcal{Z}_j)} d\mathcal{Z}_i + H[q_i] + \sum_{j \neq i} H[q_j] \end{aligned}$$

We will consider derivative of the Lagrangian with normalizer of q_i , then we have:

$$\frac{\partial}{\partial q_i} \left(\mathcal{F} + \lambda \left(\int q_i - 1 \right) \right) = \left\langle \log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}^{(k-1)}) \right\rangle_{\prod_{j \neq i} q_j(\mathcal{Z}_j)} - \log q_i(\mathcal{Z}_i) - \frac{q_i(\mathcal{Z}_i)}{q_i(\mathcal{Z}_i)} + \lambda = 0$$

This implies that we have the update as defined above. \square

Remark 25. (Observations on Update Rule) This update of factorized model depends on the expected sufficient statistics of q and this means that we don't need the entire distribution and just relevant expectation. If $\mathcal{Z}_i = z_i$ (factorized over all variables) then this technique is called mean field expectation:

- Suppose $P(\mathcal{X}, \mathcal{Z})$ has sufficient statistics that are separated in the latent variables (for example Boltzman machine):

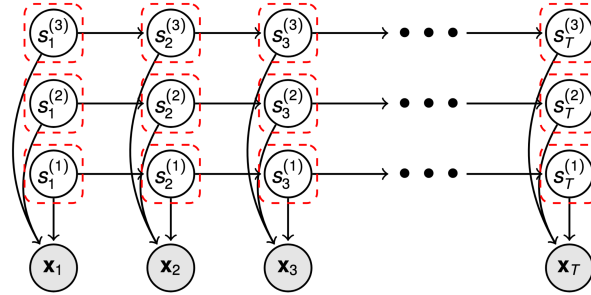
$$P(\mathcal{X}, \mathcal{Z}) = \frac{1}{Z} \exp \left(\sum_{ij} W_{ij} s_i s_j + \sum_i b_i s_i \right)$$

- The expectation with respected to fully factored q over $z_i \in \mathcal{Z}$ gives us:

$$\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{\prod_i q_i} = \sum_{ij} W_{ij} \langle s_i \rangle_{q_i} \langle s_j \rangle_{q_j} + \sum_i b_i \langle s_i \rangle_{q_i}$$

The update each q_i given sufficient statistics of the other. Each variables see the mean field imposed by its neighbour and we update these fields until they all agree.

Remark 26. (Factored Variational Inference in Time Series) Consider the graphical model:



We consider the following inference with the factorization $\prod_{m,t} q_t^m(s_t^m)$, and so we have:

$$\begin{aligned} q_t^{(m)} \left(s_t^{(m)} \right) &\propto \exp \left\langle \log P \left(s_{1:T}^{1:M}, x_{1:T} \right) \right\rangle_{\prod_{-(m,t)} q_{t'}^{m'}(s_{t'}^{m'})} \\ &= \exp \left\langle \sum_{\mu=1}^M \sum_{\tau=1}^T \log P(s_{\tau}^{\mu} | s_{\tau-1}^{\mu}) + \sum_{\tau=1}^T \log P(x_{\tau} | s_{\tau}^{1:M}) \right\rangle_{\prod_{-(m,t)} q_{t'}^{m'}(s_{t'}^{m'})} \end{aligned}$$

We can show that this is proportional to the following:

$$\exp \left[\underbrace{\langle \log P(s_t^m | s_{t-1}^m) \rangle_{q_{t-1}^m} + \langle \log P(x_t | s_t^{1:M}) \rangle_{\prod_{-m} q_{t'}^{m'}}}_{\alpha_t^m(i)} + \underbrace{\langle \log P(s_{t-1}^m | s_t^m) \rangle_{q_{t+1}^m}}_{\beta_t^m(i)} \right]$$

Where we have the following variables (as we have the same message passing algorithm):

$$\alpha_t^m(i) \propto \exp \left[\sum_j \log \Phi_{ji}^m q_{t-1}^m(j) \right] \cdot \exp \left[\langle \log A_i(x_t) \rangle_{q_t^{-m}} \right]$$

$$\beta_t^m(i) \propto \exp \left[\sum_j \log \Phi_{ij}^m q_{t+1}^m(j) \right]$$

The update only depends on the immediate neighbours in chain. The chain couple only through joint output. It has multiple passes and message depends on (approximate) marginal. Finally, the evidence doesn't appear explicitly in backward message.

Remark 27. (Difference Factored Variational Inference) Consider the difference kind of factorization:

$$q(s_{1:T}^{1:M}) = \prod_m q^m(s_{1:T}^m)$$

This gives us the following graphical model:

$$q^{(m)}(s_{1:T}^m) \propto \exp \left(\langle \log p(s_{1:T}^{1:M}, x_{1:T}) \rangle_{\prod_{-m} q^{m'}(s_{1:T}^{m'})} \right)$$

$$= \exp \left[\sum_{t=1}^T \log P(s_t^m | s_{t-1}^m) + \sum_{t=1}^T \langle \log P(x_t | s_t^{1:M}) \rangle_{\prod_{-m} q^{m'}(s_t^{m'})} \right]$$

Please note that this is similar to standard HMM joint with a modified likelihood term, where we cycle through multiple forward-backward pass updating likelihood terms each time.

Remark 28. (Message on Arbitrary Graph) Consider DAG, where we have, the following:

$$P(\mathcal{X}, \mathcal{Z}) = \prod_k P(\mathcal{V}_k | \text{Pa}(\mathcal{V}_k))$$

We let $q(\mathcal{Z}) = \prod_i q_i(\mathcal{Z})$ for disjoint set $\{\mathcal{Z}_i\}$. We have the following VE update:

$$q_i^*(\mathcal{Z}_i) \propto \exp \left(\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{-i}(\mathcal{Z})} \right)$$

Then, we have the following:

$$\log q_i^*(\mathcal{Z}_i) = \left\langle \sum_k P(\mathcal{V}_k | \text{Pa}(\mathcal{V}_k)) \right\rangle_{q_{-i}(\mathcal{Z})} + \text{const.}$$

$$= \sum_{j \in \mathcal{Z}_i} \langle \log P(\mathcal{Z}_j | \text{Pa}(\mathcal{Z}_j)) \rangle_{q_{-i}(\mathcal{Z})} + \sum_{j \in \text{ch}(\mathcal{Z}_i)} \langle \log P(\mathcal{V}_j | \text{Pa}(\mathcal{V}_j)) \rangle_{q_{-i}(\mathcal{Z})} + \text{const.}$$

Each node receives message from its markov boundary: parent/children/parent of children (all neighbours in the corresponding graphs).

Remark 29. (Non-Factored Variational Model) The term variational approximation is used whenever a bound on likelihood is optimized, but not necessary be tight, as we have:

- Closed form update in special case.
- Numerical or sampling based computation of expectation.
- Recognition Network or amortization to approximate variational parameter.
- No free energy based bound on the likelihood.

- We can perform MAP estimate or zero-temperature EM as parametric form of variational inference.

Proposition 4.2. (Variational Bayes) We can consider performing the same for integral over parameter in order to bound the log-marginal likelihood or evidence:

$$\begin{aligned}
\log P(\mathcal{X}|\mathcal{M}) &= \log \iint P(\mathcal{Z}, \mathcal{X}|\theta, \mathcal{M})P(\theta|\mathcal{M}) \, d\mathcal{Z} \, d\theta \\
&= \max_{\mathcal{Q}} \iint Q(\mathcal{Z}, \theta) \log \frac{P(\mathcal{X}, \mathcal{Z}, \theta|\mathcal{M})}{Q(\mathcal{Z}, \theta)} \, d\mathcal{Z} \, d\theta \\
&\geq \max_{\mathcal{Q}_{\mathcal{Z}}, \mathcal{Q}_{\theta}} \iint Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta) \log \frac{P(\mathcal{X}, \mathcal{Z}, \theta|\mathcal{M})}{Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta)} \, d\mathcal{Z} \, d\theta = \max_{\mathcal{Q}_{\mathcal{Z}}, \mathcal{Q}_{\theta}} \mathcal{F}(Q_{\mathcal{Z}}, Q_{\theta})
\end{aligned}$$

The second equality comes from the Jensen's inequality bound being tight. This leads to variational Bayesian EM algorithm with free energy being $\mathcal{F}(Q_{\mathcal{Z}}, Q_{\theta})$:

$$\begin{aligned}
Q_{\mathcal{Z}}^*(\mathcal{Z}) &\propto \exp \langle \log P(\mathcal{Z}, \mathcal{X}|\theta) \rangle_{Q_{\theta}(\theta)} \\
Q_{\theta}^*(\theta) &\propto P(\theta) \exp \langle \log(\mathcal{Z}, \mathcal{X}|\theta) \rangle_{Q_{\mathcal{Z}}(\theta)}
\end{aligned}$$

Proof. The maximizing \mathcal{F} is the same as minimizing KL-divergence between the approximate posterior $Q(\mathcal{Z})Q(\theta)$ and the true posterior $P(\theta, \mathcal{Z}|\mathcal{X})$ as we have:

$$\begin{aligned}
\log P(\mathcal{X}) - \mathcal{F}(Q_{\mathcal{Z}}, Q_{\theta}) &= \iint Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta) \log P(\mathcal{X}) \, d\mathcal{Z} \, d\theta - \iint Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta) \log \frac{P(\mathcal{X}, \mathcal{Z}, \theta|\mathcal{M})}{Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta)} \, d\mathcal{Z} \, d\theta \\
&= \iint Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta) \log \frac{Q_{\mathcal{Z}}(\mathcal{Z})Q_{\theta}(\theta)}{P(\mathcal{Z}, \theta|\mathcal{X}, \mathcal{M})} \, d\mathcal{Z} \, d\theta = \text{KL}(Q\|P)
\end{aligned}$$

□

Remark 30. (Conjugate Exponential) Let's consider the conjugate exponential model on latent variable:

- *Condition 1:* Consider the joint probability over variables in exponential family:

$$P(\mathcal{Z}, \mathcal{X}|\theta) = f(\mathcal{X}, \mathcal{Z})g(\theta) \exp \{ \phi(\theta)^T \mathbf{T}(\mathcal{Z}, \mathcal{X}) \}$$

- *Condition 2:* The joint over parameter is conjugate to this joint probability:

$$P(\theta|\nu, \tau) = h(\nu, \tau)g(\theta)^\nu \exp \{ \phi(\theta)^T \tau \}$$

where ν and τ on the prior and we can see them as number of pseudo-observations and values of pseudo-observations, respectively.

Remark 31. (Conjugate Exponential VB) Given a iid set $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ if the model is conjugate exponential then we have. $Q_{\theta}(\theta)$ is also conjugate, as we have:

$$\begin{aligned}
Q_{\theta}(\theta) &\propto P(\theta) \exp \left\langle \sum_i \log P(\mathcal{Z}_i, \mathcal{X}_i|\theta) \right\rangle_{Q_{\mathcal{Z}}(\mathcal{Z})} \\
&= h(\nu, \tau)g(\theta)^\nu \exp(\phi(\theta)^T \tau) g(\theta)^n \exp \left(\left\langle \log f(\mathcal{X}, \mathcal{Z}) + \sum_{i=1}^n \phi(\theta)^T \mathbf{T}(\mathbf{z}_i, \mathbf{x}_i) \right\rangle_{Q_{\mathcal{Z}}} \right) \\
&\propto h(\tilde{\nu}, \tilde{\tau})g(\theta)^{\nu+n} \exp \left(\phi(\theta)^T \left(\tau + \left\langle \sum_{i=1}^n \mathbf{T}(\mathbf{z}_i, \mathbf{x}_i) \right\rangle_{Q_{\mathcal{Z}}} \right) \right) = h(\tilde{\nu}, \tilde{\tau})g(\theta)^{\tilde{\nu}} \exp(\phi(\theta)^T \tilde{\tau})
\end{aligned}$$

we can set $\tilde{\nu} = \nu + n$ and $\tilde{\tau} = \tau + \sum_i \langle \mathbf{T}(\mathbf{z}_i, \mathbf{x}_i) \rangle_{Q_{\mathcal{Z}}}$

Remark 32. (Factorized EM) We have the following factorized model as we have $Q_{\mathcal{Z}}(\mathcal{Z}) = \prod_{i=1}^n Q_{\mathcal{Z}_i}(\mathcal{Z}_i)$ (over the dataset). It has the same form as regular E-step of regular EM, where we have:

$$\begin{aligned} Q_{\mathcal{Z}_i}(\mathcal{Z}_i) &\propto \exp\left(\langle \log P(\mathbf{z}_i, \mathbf{x}_i) \rangle_{Q_\theta}\right) \\ &\propto f(\mathbf{z}_i, \mathbf{x}_i) \exp\left\{\langle \phi(\boldsymbol{\theta}) \rangle_{Q_\theta}^T \mathbf{T}(\mathbf{z}_i, \mathbf{x}_i)\right\} \\ &= P(\mathbf{z}_i | \mathbf{x}_i, \tilde{\phi}(\boldsymbol{\theta})) \end{aligned}$$

where the natural parameter $\tilde{\phi}(\boldsymbol{\theta}) = \langle \phi(\boldsymbol{\theta}) \rangle_{Q_\theta}$, where the inference is unchanged to regular EM.

Remark 33. (Comparison Between EM for MAP estimation and Variational Bayesian EM) Let's compare the differences between the 2 algorithms:

EM for MAP estimation	Variational Bayesian EM
Goal: maximize $P(\boldsymbol{\theta} \mathcal{X}, m)$ wrt $\boldsymbol{\theta}$	Goal: maximise bound on $P(\mathcal{X} m)$ wrt Q_θ
E Step: compute $Q_{\mathcal{Z}}(\mathcal{Z}) \leftarrow p(\mathcal{Z} \mathcal{X}, \boldsymbol{\theta})$	VB-E Step: compute $Q_{\mathcal{Z}}(\mathcal{Z}) \leftarrow p(\mathcal{Z} \mathcal{X}, \tilde{\phi})$
M Step: $\boldsymbol{\theta} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} \int d\mathcal{Z} Q_{\mathcal{Z}}(\mathcal{Z}) \log P(\mathcal{Z}, \mathcal{X}, \boldsymbol{\theta})$	VB-M Step: $Q_\theta(\boldsymbol{\theta}) \leftarrow \exp \int d\mathcal{Z} Q_{\mathcal{Z}}(\mathcal{Z}) \log P(\mathcal{Z}, \mathcal{X}, \boldsymbol{\theta})$

The following are the feature of VB, as we have:

- It reduces to EM if we set $Q_\theta(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$
- \mathcal{F}_m increase monotonically and interpolate the model complexity penalty.
- We have analytical parameter distribution but we don't have to constraint to Gaussian.
- VB-E step has the same complexity as E-step.
- We can use junction tree, belief propagatio, kalman filter algorithm in VB E-step of VB-EM but we have to use the expected natural parameter $\tilde{\phi}$.

Remark 34. (VB and Model Classification) VB-EM gives us the approximate posterior Q_θ over the model parameter:

- It also gives us the lower bound on the model estimate as we have:

$$\max \mathcal{F}_M(Q_{\mathcal{Z}}, Q_\theta) \leq P(\mathcal{D} | \mathcal{M})$$

- These lower bound can be compared amongsts model to find the right out. However, if we consider continuous domain of model that is specified by the hyperparameter η , then VB free energy depends on parameter

$$\mathcal{F}(Q_{\mathcal{Z}}, Q_\theta, \eta) = \iint Q_{\mathcal{Z}}(\mathcal{Z}) Q_\theta(\boldsymbol{\theta}) \log \frac{P(\mathcal{X}, \mathcal{Z}, \boldsymbol{\theta} | \eta)}{Q_{\mathcal{Z}}(\mathcal{Z}) Q_\theta(\boldsymbol{\theta})} d\mathcal{Z} d\boldsymbol{\theta} \leq P(\mathcal{X} | \eta)$$

Hyperparameter-M step maximizes the current bound with respected to η :

$$\eta \leftarrow \operatorname{argmax}_{\eta} \iint Q_{\mathcal{Z}}(\mathcal{Z}) Q_\theta(\boldsymbol{\theta}) \log P(\mathcal{X}, \mathcal{Z}, \boldsymbol{\theta} | \eta) d\mathcal{Z} d\boldsymbol{\theta}$$

Remark 35. (ARD for Unsupervised Learning) The hyperparameter method to select a useful input regression. This is similar idea with variational Bayes method that can learn a latent dimension. Consider the factor analysis:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{\Lambda}\mathbf{z}, \mathbf{\Psi})$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ with column-wise prior to be $\mathbf{\Lambda}_{:i} \sim \mathcal{N}(\mathbf{0}, \alpha_i^{-1}\mathbf{I})$. Then the VB free-energy is given by:

$$\mathcal{F}\left(Q_{\mathcal{Z}}(\mathcal{Z}), Q_{\Lambda}(\mathbf{\Lambda}), \mathbf{\Psi}, \boldsymbol{\alpha}\right) = \left\langle \log P(\mathcal{X}, \mathcal{Z}|\mathbf{\Lambda}, \mathbf{\Psi}) + \log P(\mathbf{\Lambda}|\boldsymbol{\alpha}) + \log P(\mathbf{\Psi}) \right\rangle_{Q_{\mathcal{Z}}(\mathcal{Z})Q_{\Lambda}(\mathbf{\Lambda})} + \text{const.}$$

We require the hyperparameter optimization require us $\boldsymbol{\alpha} \leftarrow \arg \max \langle \log P(\mathbf{\Lambda}|\boldsymbol{\alpha}) \rangle_{Q_{\Lambda}}$. Now we obtain the following optimization:

- Now, Q_{Λ} is Gaussian with the same form as linear regression but with expected moment of \mathbf{z} instead of the input.
- The optimization with respected to the distribution $\mathbf{\Psi}$ and $\boldsymbol{\alpha}$ will cause some α_i to diverge as in regression ARD.
- Same as selecting relevant latent dimension, effectively learning the dimension of latent variable.

Remark 36. (Problem with GP + Solution) Given the GP prediction, we have:

$$y'|\mathbf{X}, \mathbf{Y}, \mathbf{x} \sim \mathcal{GP}(\mathbf{K}_{xX}(\mathbf{K}_{XX} + \sigma^2\mathbf{I})^{-1}\mathbf{Y}, \mathbf{K}_{xx} - \mathbf{K}_{xX}(\mathbf{K}_{XX} + \sigma^2\mathbf{I})^{-1}\mathbf{K}_{Xx} + \sigma^2)$$

The evidence (for learning the kernel hyperparameter) to be:

$$\log P(\mathbf{Y}|\mathbf{X}) = -\frac{1}{2} \log |2\pi(\mathbf{K}_{XX} + \sigma^2\mathbf{I})|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{Y}^T(\mathbf{K}_{XX} + \sigma^2\mathbf{I})^{-1}\mathbf{Y}\right)$$

Computing this require inverting $N \times N$ matrix with has the time complexity of $\mathcal{O}(K^3)$. We consider the smaller set of possible fictious measurement \mathbf{U} at input \mathbf{Z} such that:

$$P(y'|\mathbf{Z}, \mathbf{U}, \mathbf{x}') \approx P(y'|\mathbf{X}, \mathbf{Y}, \mathbf{x}')$$

Write \mathbf{F} for the smooth GP function that underlie $\mathbf{Y} \sim \mathcal{N}(\mathbf{F}, \sigma^2\mathbf{I})$. Introduce the measurement \mathbf{U} at input \mathbf{Z} , as the likelihood can be written as:

$$\begin{aligned} P(\mathbf{Y}|\mathbf{X}) &= \iint P(\mathbf{Y}, \mathbf{F}, \mathbf{U}|\mathbf{X}, \mathbf{Z}) d\mathbf{F} d\mathbf{U} \\ &= \iint P(\mathbf{Y}|\mathbf{F})P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})P(\mathbf{U}|\mathbf{Z}) d\mathbf{F} d\mathbf{U} \end{aligned}$$

Remark 37. (Free-Energy) As we have \mathbf{F} and \mathbf{U} to be the latent variable, we introduce the variational distribution $q(\mathbf{F}, \mathbf{U})$ to be:

$$\mathcal{F}(q(\mathbf{F}, \mathbf{U}), \theta) = \left\langle \log \frac{P(\mathbf{Y}|\mathbf{F})P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})P(\mathbf{U}|\mathbf{Z})}{q(\mathbf{F}, \mathbf{U})} \right\rangle_{q(\mathbf{F}, \mathbf{U})}$$

Consider the variational distribution of the latent to be $q(\mathbf{F}, \mathbf{U}) = P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})$. We fix $\mathbf{F}|\mathbf{U}$ with reference to \mathbf{U} to make the information about \mathbf{Y} is compressed to $q(\mathbf{U})$, which we have:

$$\begin{aligned} \mathcal{F}(q(\mathbf{F}, \mathbf{U}), \theta, \mathbf{Z}) &= \left\langle \log \frac{P(\mathbf{Y}|\mathbf{F})P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})P(\mathbf{U}|\mathbf{Z})}{P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})} \right\rangle_{P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})} \\ &= \left\langle \log \frac{P(\mathbf{Y}|\mathbf{F})P(\mathbf{U}|\mathbf{Z})}{q(\mathbf{U})} \right\rangle_{P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})} \\ &= \left\langle (\log P(\mathbf{Y}|\mathbf{F}))_{P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})} + \log P(\mathbf{U}|\mathbf{Z}) - q(\mathbf{U}) \right\rangle_{q(\mathbf{U})} \end{aligned}$$

Let's consider the inner expectation, which we can use the Gaussian process results (without noise on \mathbf{F}):

$$\begin{aligned} \langle \log P(\mathbf{Y}|\mathbf{F}) \rangle_{P(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})} &= \left\langle -\frac{1}{2} \log |2\pi\sigma^2 \mathbf{I}| - \frac{1}{2\sigma^2} \text{Tr} [(\mathbf{Y} - \mathbf{F})(\mathbf{Y} - \mathbf{F})^T] \right\rangle_{P(\mathbf{F}|\mathbf{Y})} \\ &= -\frac{1}{2} \log |2\pi\sigma^2 \mathbf{I}| - \frac{1}{2\sigma^2} \text{Tr} \left[\left(\mathbf{Y} - \langle \mathbf{F} \rangle_{P(\mathbf{F}|\mathbf{Y})} \right) \left(\mathbf{Y} - \langle \mathbf{F} \rangle_{P(\mathbf{F}|\mathbf{Y})} \right)^T \right] - \frac{1}{2\sigma^2} \text{Tr} [\Sigma_{\mathbf{F}|\mathbf{Y}}] \\ &= \log \mathcal{N}(\mathbf{Y} | \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{U}, \sigma^2 \mathbf{I}) - \frac{1}{2\sigma^2} \text{Tr} [\mathbf{K}_{XX} - \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX}] \end{aligned}$$

This gives us the following free energy function:

$$\mathcal{F}(q(\mathbf{U}), \boldsymbol{\theta}, \mathbf{Z}) = \left\langle \log \frac{\mathcal{N}(\mathbf{Y} | \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{U}, \sigma^2 \mathbf{I}) P(\mathbf{U} | \mathbf{Z})}{q(\mathbf{U})} \right\rangle_{q(\mathbf{U})} - \frac{1}{2\sigma^2} \text{Tr} [\mathbf{K}_{XX} - \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX}]$$

Now, we can see that the expectation is free energy of PPCA-like model with the latent to be $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{ZZ})$, with the following linear Gaussian model $p(\mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{U}, \sigma^2 \mathbf{I})$, thusm we have the following free energy:

$$\begin{aligned} \mathcal{F}(q^*(\mathbf{U}), \boldsymbol{\theta}, \mathbf{Z}) &= \log \mathcal{N}(\mathbf{Y} | \mathbf{0}, \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} + \sigma^2 \mathbf{I}) - \frac{1}{2\sigma^2} \text{Tr} [\mathbf{K}_{XX} - \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX}] \\ &= \log \mathcal{N}(\mathbf{Y} | \mathbf{0}, \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} + \sigma^2 \mathbf{I}) - \frac{1}{2\sigma^2} \text{Tr} [\mathbf{K}_{XX} - \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX}] \end{aligned}$$

We can now optimize the free energy numerically with respect to $\boldsymbol{\theta}$ and \mathbf{Z} to adjust GP prior and quality of variational approximation.

Remark 38. If \mathbf{X} is unobserved, then assume that $q(\mathbf{X}, \mathbf{F}, \mathbf{U}) = q(\mathbf{X})P(\mathbf{F}, \mathbf{X}, \mathbf{U})q(\mathbf{U})$ then we have the free energy to be:

$$\mathcal{F} = \langle \log P(\mathbf{Y}, \mathbf{F}, \mathbf{U} | \mathbf{X}) + \log P(\mathbf{X}) \rangle_{q(\mathbf{U})q(\mathbf{X})}$$

which can be simplified to tractable components in similar way as above.

5 Expectation Propagation

5.1 Introduction

Remark 39. (Non-Linear State Space) We are given the following non-linear space with the following transition and emission probability as we have:

$$\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{w}_t \quad \mathbf{x}_t = g(\mathbf{z}_t, \mathbf{u}_t) + \mathbf{v}_t$$

where \mathbf{w}_t and \mathbf{v}_t are usually Gaussian, as we can linearized the non-linear function around $\hat{\mathbf{z}}_t^t$:

$$\mathbf{z}_{t+1} \approx \underbrace{f(\hat{\mathbf{z}}_t^t, \mathbf{u}_t)}_{\tilde{\mathbf{B}}_t \mathbf{u}_t} + \underbrace{\frac{\partial f}{\partial \mathbf{z}_t} \Big|_{\hat{\mathbf{z}}_t^t}}_{\tilde{\mathbf{A}}_t} (\mathbf{z}_t - \hat{\mathbf{z}}_t^t) + \mathbf{w}_t \quad \mathbf{x}_t \approx \underbrace{f(\hat{\mathbf{z}}_t^{t-1}, \mathbf{u}_t)}_{\tilde{\mathbf{D}}_t \mathbf{u}_t} + \underbrace{\frac{\partial f}{\partial \mathbf{z}_t} \Big|_{\hat{\mathbf{z}}_t^{t-1}}}_{\tilde{\mathbf{C}}_t} (\mathbf{z}_t - \hat{\mathbf{z}}_t^{t-1}) + \mathbf{w}_t$$

We can run the kalman filter on non-stationary linearized system $(\tilde{\mathbf{A}}_t, \tilde{\mathbf{B}}_t, \tilde{\mathbf{C}}_t, \tilde{\mathbf{D}}_t)$:

- Adaptive approximate non-Gaussian message by Gaussian.
- Local linearization depends on central point of distribution, mening that approximation degrades with more stable uncertainty.
- This might work in the system that is close to linear.

Remark 40. (Other Message Approximation) We have the following message on latent chain:

$$P(\mathbf{z}_t|\mathbf{x}_{1:t}) = \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int P(\mathbf{z}_t|\mathbf{z}_{t-1}) P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1}$$

$$\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) \approx \frac{1}{Z} P(\mathbf{x}_t|\mathbf{z}_t) \int \underbrace{P(\mathbf{z}_t|\mathbf{z}_{t-1})}_{\mathcal{N}(f(\mathbf{z}_{t-1}), \mathbf{Q})} \underbrace{\tilde{P}(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}_{\mathcal{N}(\hat{\mathbf{z}}_{t-1}, \hat{\mathbf{V}}_{t-1})} d\mathbf{z}_{t-1}$$

There are several, way to approximate the integration: linearization at the peak (EKF) is only one approach. We can use the Laplace filter use mode and curvature of integrand. Finally, we can use the sigma-point.

Remark 41. (Why not KL?) We can consider the parametric variational as we have

$$\text{KL} \left[\mathcal{N}(\hat{\mathbf{z}}_t, \hat{\mathbf{V}}) \left\| \int P(\mathbf{z}_t|\mathbf{z}_{t-1}) P(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1} \right. \right]$$

This may be hard to find the closed form solution to this KL divergence and we might have to find the value using Monte-Carlo sampling.

5.2 Expectation Propagation

Proposition 5.1. *Given the $p(\mathbf{x})$, we let q to be exponential family, where we have:*

$$Q(\mathbf{x}) = \frac{\exp(\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta})}{Z(\boldsymbol{\theta})}$$

Now, consider the following minimization of: $Q^ = \arg \min_Q \text{KL} [P(\mathbf{x}) \| Q(\mathbf{x})]$ is solved when:*

$$\langle \mathbf{T}(\mathbf{x}) \rangle_{Q^*} = \langle \mathbf{T}(\mathbf{x}) \rangle_P$$

Or matches the sufficient statistics.

Proof. Consider the KL-divergence to be:

$$\begin{aligned} \arg \min_Q \text{KL} [P(\mathbf{x}) \| Q(\mathbf{x})] &= \arg \min_{\boldsymbol{\theta}} \left[P(\mathbf{x}) \left\| \frac{\exp(\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \right. \right] \\ &= \arg \min_{\boldsymbol{\theta}} - \int P(\mathbf{x}) \log \left(\frac{\exp(\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \right) d\mathbf{x} \\ &= \arg \min_{\boldsymbol{\theta}} \log(Z(\boldsymbol{\theta})) - \int P(\mathbf{x}) [\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta}] d\mathbf{x} \end{aligned}$$

The second equality comes from the fact that $\int P(\mathbf{x}) \log P(\mathbf{x}) d\mathbf{x}$ is constant. Now, we have the following derivative with respected to $\boldsymbol{\theta}$:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}} \left(\log(Z(\boldsymbol{\theta})) - \int P(\mathbf{x}) [\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta}] d\mathbf{x} \right) &= \frac{1}{Z(\boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} \int \exp(\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta}) d\mathbf{x} - \int p(\mathbf{x}) T(\mathbf{x}) d\mathbf{x} \\ &= \frac{1}{Z(\boldsymbol{\theta})} \int \exp(\mathbf{T}(\mathbf{x})^T \boldsymbol{\theta}) T(\mathbf{x}) d\mathbf{x} - \langle T(\mathbf{x}) \rangle_P \\ &= \langle T(\mathbf{x}) \rangle_P - \langle T(\mathbf{x}) \rangle_Q \end{aligned}$$

Setting this to zero as we have the solution. □

Remark 42. (KL-Divergence Solution) It is better for use to perform the following KL minimization:

$$\text{KL} \left[\int P(\mathbf{z}_t | \mathbf{z}_{t-1}) P(\mathbf{z}_{t-1} | \mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1} \left\| \mathcal{N}(\hat{\mathbf{z}}_t, \hat{\mathbf{V}}) \right. \right]$$

As we can consider the expected sufficient statistics of the real distribution, instead of finding the exact version of it. However, to calculate the expected sufficient statistics $\langle \mathbf{T}(\mathbf{x}) \rangle_P$, which may be analytically tractable, however for the high-dimensional integral (we can use various algorithms to calculate the expectation).

Remark 43. (Another Problem of KL-Divergence) Let's consider the KL-divergence of the factored model, as we have:

$$\begin{aligned} \arg \min_{q_i} \text{KL} \left[P(\mathcal{Z} | \mathcal{X}) \left\| \prod_j q_j(\mathcal{Z}_j | \mathcal{X}) \right. \right] &= \arg \min_{q_i} - \int P(\mathcal{Z} | \mathcal{X}) \log \prod_j q_j(\mathcal{Z}_j | \mathcal{X}) d\mathcal{Z} \\ &= \arg \min_{q_i} - \sum_j \int P(\mathcal{Z} | \mathcal{X}) \log q_j(\mathcal{Z}_j | \mathcal{X}) d\mathcal{Z} \\ &= \arg \min_{q_i} - \int P(\mathcal{Z}_i | \mathcal{X}) \log q_i(\mathcal{Z}_i | \mathcal{X}) d\mathcal{Z}_i \\ &= P(\mathcal{Z}_i | \mathcal{X}) \end{aligned}$$

We will have to know the marginal, which is intractable.

Remark 44. (EP Motivation) The posterior distribution in a graphical model is a product of factors:

$$P(\mathcal{Z} | \mathcal{X}) = \frac{P(\mathcal{Z}, \mathcal{X})}{P(\mathcal{X})} = \frac{1}{Z} \prod_i P(\mathcal{Z}_i | \text{Pa}(\mathcal{Z}_i)) \propto \prod_i f_i(\mathcal{Z}_i)$$

where \mathcal{Z}_i isn't necessary disjoint and we call f_i sites, we consider the same factorization, but with approximate site as we have:

$$\min_{\tilde{f}_i} \text{KL} \left[f_i(\mathcal{Z}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j) \left\| \tilde{f}_i(\mathcal{Z}_i) \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j) \right. \right] \iff \min_f \text{KL} \left[f_i(\mathcal{Z}_i) q_{-i}(\mathcal{Z}) \left\| f(\mathcal{Z}_i) q_{-i}(\mathcal{Z}) \right. \right]$$

where we have $q_{-i}(\mathcal{Z}) = \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j)$. This leads to 2 ideas about expectation propagation as:

- *Expectation*: Approximation of factors, where we perform a projection to exponential family, thus requiring a sufficient statistics.
- *Propagation*: Local divergence minimization leads to message passing approach (hence the number of propagation).

Proposition 5.2. (Simpler Update) If we consider the context factor, as $q_{-i}(\mathcal{Z}) = q_{-i}(\mathcal{Z}_i) q_{-i}(\mathcal{Z}_{-i} | \mathcal{Z}_i)$ and we have $\mathcal{Z}_{-i} = \mathcal{Z} \setminus \mathcal{Z}_i$ we can show that:

$$\min_f \text{KL} \left[f_i(\mathcal{Z}_i) q_{-i}(\mathcal{Z}) \left\| f(\mathcal{Z}_i) q_{-i}(\mathcal{Z}) \right. \right] \equiv \min_f \text{KL} \left[f_i(\mathcal{Z}_i) q_{-i}(\mathcal{Z}_i) \left\| f(\mathcal{Z}_i) q_{-i}(\mathcal{Z}_i) \right. \right]$$

Please see that the differences between those 2 between *red* and *blue*. We call $q_{-i}(\mathcal{Z}_i)$

Proof. Consider the following equalities:

$$\begin{aligned}
\min_f \text{KL} \left[f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}) \parallel f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}) \right] &= \max_f \int f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}) \log \left[f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}) \right] d\mathcal{Z} \\
&= \max_f \int f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_{-i}|\mathcal{Z}_i) \log \left[f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_{-i}|\mathcal{Z}_i) \right] d\mathcal{Z}_i d\mathcal{Z}_{-i} \\
&= \max_f \int f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_{-i}|\mathcal{Z}_i) \left[\log f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \right] d\mathcal{Z}_i d\mathcal{Z}_{-i} \\
&= \max_f \int f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \left[\log f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \right] d\mathcal{Z}_i \int q_{-i}(\mathcal{Z}_{-i}|\mathcal{Z}_i) d\mathcal{Z}_{-i} \\
&= \min_f \text{KL} \left[f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \parallel f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \right]
\end{aligned}$$

And so the equality is proven. \square

Definition 5.1. (Expectation Propagation) We consider the following algorithms for expectation propagation, pseudocode:

Algorithm 1 Expectation Propagation Algorithm

- 1: **Input:** $f_1(\mathcal{Z}_1), \dots, f_N(\mathcal{Z}_N)$
- 2: **Initialize:** $\tilde{f}_i(\mathcal{Z}_i) = 1$ for $i \in [n]$ and we have $q(\mathcal{Z}) \propto \prod_i \tilde{f}_i(\mathcal{Z}_i)$
- 3: **while** convergence **do**
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: *Delete:* We do the update to get the cavity

$$q_{-i}(\mathcal{Z}) \leftarrow \frac{q(\mathcal{Z})}{\tilde{f}_i(\mathcal{Z}_i)} = \prod_{j \neq i} \tilde{f}_j(\mathcal{Z}_j)$$

- 6: *Project:* Performing KL-divergence minimization

$$\tilde{f}_i^{\text{new}}(\mathcal{Z}) \leftarrow \arg \min_f \text{KL} \left[f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \parallel f(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i) \right]$$

- 7: *Include:* $q(\mathcal{Z}) \leftarrow \tilde{f}_i^{\text{new}}(\mathcal{Z}_i)q_{-i}(\mathcal{Z})$
 - 8: **end for**
 - 9: **end while**
-

Remark 45. (Calculate the Cavity Distribution) The cavity distribution can be broken down into product of terms for each neighboring cliques:

$$q_{-i}(\mathcal{Z}_i) = \prod_{j \in \text{ne}(i)} M_{j \rightarrow i}(\mathcal{Z}_i \cap \mathcal{Z}_j)$$

The i -th site has been approximated as the message can be passed onto neighboring cliques by normalizing the shared variable. This is belief propagation. Furthermore, the message updates can be scheduled in any order. However, there is no guarantee of convergence.

Remark 46. (Normalizer) As long as approximating class is tractable, normalizer can be computed as needed. Consider the approximation class to be:

$$\tilde{f}_i(\mathcal{Z}_i) \propto \exp \left(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_i - \Phi(\boldsymbol{\theta}_i) \right)$$

This is the same as setting every entries $\boldsymbol{\theta}_i$ to 0 except for the one that is in \mathcal{Z}_i (while finding the sufficient statistics for all latents). This will give us the following probability:

$$q(\mathcal{Z}) \propto \prod_i \tilde{f}_i \propto \exp \left(\mathbf{T}(\mathcal{Z})^T \sum_i \boldsymbol{\theta}_i - \sum_i \Phi(\boldsymbol{\theta}_i) \right)$$

We can re-normalized it as we have $q(\mathcal{Z}) = \exp(\mathbf{T}(\mathcal{Z})^T \sum_i \boldsymbol{\theta}_i - \Phi(\sum_i \boldsymbol{\theta}_i))$

Remark 47. (Computing Likelihood) Consider the unnormalized exponential family, where we have the approximating site to be:

$$\tilde{f}_i = \tilde{C}_i \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_i)$$

We consider the following values

- $\boldsymbol{\theta} = \sum_i \boldsymbol{\theta}_i$, which is the natural parameter of $q(\mathcal{Z})$
- $\boldsymbol{\theta}_{-i} = \sum_{j \neq i} \boldsymbol{\theta}_j$ as the natural parameter of $q_{-i}(\mathcal{Z})$
- Exponential family (tractable) log-normalizer $\Phi(\boldsymbol{\theta})$ of $P(\mathcal{Z})$ to be:

$$\Phi(\boldsymbol{\theta}) = \log \int \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}) d\mathcal{Z}$$

Now, we are interested to find the actual approximation of the normalizer i.e:

$$\log \int \prod_{i=1}^N f_i(\mathcal{Z}_i)$$

To do this, we minimize the unnormalized KL divergence (where we also care about the normalizer, see the additional term):

$$\text{KL}_{\text{un}}[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} + \int (q(\mathbf{x}) - p(\mathbf{x})) d\mathbf{x}$$

Proposition 5.3. (Log Likelihood Approximation) *The approximation of the log-likelihood is given as:*

$$\log \int \prod_{i=1}^N f_i(\mathcal{Z}_i) \approx \log \int \prod_{i=1}^N \tilde{f}_i(\mathcal{Z}_i) = (1 - N)\Phi(\boldsymbol{\theta}) + \sum_{i=1}^N \Phi_i(\boldsymbol{\theta}_{-i})$$

where Φ_i is the log-normalizer of the distribution $\hat{P}_i(\mathcal{Z}) \propto f(\mathcal{Z}) \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta})$

Proof. To follows the unnormalized KL divergence, the first terms simply follows the expected sufficient statistics match. In this case, we have (thinking of it as having KL-divergence equal to zero):

$$\begin{aligned} & \min_{\tilde{C}_i} \text{KL}_{\text{min}} \left[\tilde{C}_i \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_i) \prod_{-i} \tilde{C}_j \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_j) \left\| f_i(\mathcal{Z}_i) \prod_{-i} \tilde{C}_j \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_j) \right. \right] \\ \iff & \int \tilde{C}_i \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_i) \prod_{-i} \tilde{C}_j \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_j) d\mathcal{Z} = \int f_i(\mathcal{Z}_i) \prod_{-i} \tilde{C}_j \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_j) d\mathcal{Z} \\ \iff & \int \tilde{C}_i \prod_i \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_j) d\mathcal{Z} = \int f_i(\mathcal{Z}_i) \prod_{-i} \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_j) d\mathcal{Z} \\ \iff & \tilde{C}_i \exp(\Phi(\boldsymbol{\theta})) = \exp(\Phi_i(\boldsymbol{\theta}_{-i})) \\ \iff & \tilde{C}_i = \exp(\Phi_i(\boldsymbol{\theta}_{-i}) - \Phi(\boldsymbol{\theta})) \end{aligned}$$

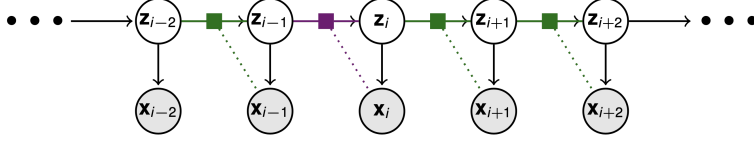
and so, we have the following log-likelihood to be:

$$\begin{aligned} \log \int \prod_{i=1}^N \tilde{f}_i(\mathcal{Z}_i) d\mathcal{Z} &= \log \int \prod_{i=1}^N \tilde{C}_i \exp(\mathbf{T}(\mathcal{Z})^T \boldsymbol{\theta}_i) d\mathcal{Z} \\ &= \log \prod_{i=1}^N \tilde{C}_i \int \exp\left(\mathbf{T}(\mathcal{Z})^T \left[\sum_{i=1}^N \boldsymbol{\theta}_i\right]\right) d\mathcal{Z} \\ &= \Phi(\boldsymbol{\theta}) + \sum_{i=1}^N \log \tilde{C}_i = (1 - N)\Phi(\boldsymbol{\theta}) + \sum_{i=1}^N \Phi_i(\boldsymbol{\theta}_{-i}) \end{aligned}$$

□

5.3 Examples

Remark 48. (Example: Non-Linear State Space Model) Consider the graphical model:



where we have the following factors:

$$P(\mathbf{z}_i | \mathbf{z}_{i-1}) = \phi_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \exp\left(-\frac{\|\mathbf{z}_i - h_s(\mathbf{z}_{i-1})\|^2}{2\sigma^2}\right)$$

$$P(\mathbf{x}_i | \mathbf{z}_i) = \psi_i(\mathbf{z}_i) = \exp\left(-\frac{\|\mathbf{x}_i - h_o(\mathbf{z}_i)\|^2}{2\sigma^2}\right)$$

We can see that $f_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \phi_i(\mathbf{z}_i, \mathbf{z}_{i-1})\psi_i(\mathbf{z}_i)$ as there are non-linear, the inference isn't generally tractable. For the EP, we will consider the approximate $\tilde{f}(\mathbf{z}_i, \mathbf{z}_{i-1})$ to be Gaussian as we have. Consider the cavity distribution to be:

$$q_{-i}(\mathbf{z}_i, \mathbf{z}_{i-1}) = \int_{\mathbf{z}_1, \dots, \mathbf{z}_{i-2}, \mathbf{z}_{i+1}, \dots, \mathbf{z}_n} \prod_{i' \neq i} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1})$$

$$= \underbrace{\int_{\mathbf{z}_1, \dots, \mathbf{z}_{i-2}} \prod_{i' < i} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1})}_{\alpha_{i-1}(\mathbf{z}_{i-1})} \underbrace{\int_{\mathbf{z}_{i+1}, \dots, \mathbf{z}_n} \prod_{i' > i} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1})}_{\beta_i(\mathbf{z}_i)}$$

Please note that α, β is being Gaussian (by default), and so we have the following update rule:

$$\tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \arg \min_{f \in \mathcal{N}} \text{KL} \left[\phi_i(\mathbf{z}_i, \mathbf{z}_{i-1}) \psi_i(\mathbf{z}_i) \alpha_{i-1}(\mathbf{z}_{i-1}) \beta_i(\mathbf{z}_i) \middle\| f(\mathbf{z}_i, \mathbf{z}_{i-1}) \alpha_{i-1}(\mathbf{z}_{i-1}) \beta_i(\mathbf{z}_i) \right]$$

We can consider the following optimization instead as we have:

$$\tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) = \arg \min_{P \in \mathcal{N}} \text{KL} \left[\hat{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) \middle\| P(\mathbf{z}_{i-1}, \mathbf{z}_i) \right]$$

$$\implies \tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \frac{\tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i)}{\alpha_{i-1}(\mathbf{z}_{i-1}) \beta_i(\mathbf{z}_i)}$$

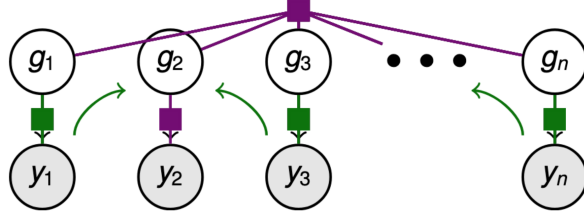
Now, we consider each values $\alpha_i(\mathbf{z}_i)$ and $\beta_{i-1}(\mathbf{z}_{i-1})$ as the propagation step:

$$\alpha_i(\mathbf{z}_i) = \int_{\mathbf{z}_1, \dots, \mathbf{z}_{i-1}} \prod_{i' < i+1} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1}) = \int_{\mathbf{z}_i} \alpha_{i-1}(\mathbf{z}_{i-1}) \tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \frac{1}{\beta_i(\mathbf{z}_i)} \int_{\mathbf{z}_{i-1}} \tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) d\mathbf{z}_{i-1}$$

$$\beta_{i-1}(\mathbf{z}_{i-1}) = \int_{\mathbf{z}_i, \dots, \mathbf{z}_n} \prod_{i' < i+1} \tilde{f}_{i'}(\mathbf{z}_{i'}, \mathbf{z}_{i'-1}) = \int_{\mathbf{z}_i} \beta_i(\mathbf{z}_i) \tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1}) = \frac{1}{\alpha_{i-1}(\mathbf{z}_{i-1})} \int_{\mathbf{z}_i} \tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i) d\mathbf{z}_i$$

The last equality comes from the another definition of $\tilde{f}_i(\mathbf{z}_i, \mathbf{z}_{i-1})$. The update of both α_i and β_i are marginalization of $\tilde{P}(\mathbf{z}_{i-1}, \mathbf{z}_i)$, which is also a Gaussian.

Remark 49. (EP For GP Classification) We can write the GP joint on g_i and y_i as a factor graph:



Consider the following factorization:

$$P(g_1, \dots, g_n, y_1, \dots, y_n) = \underbrace{\mathcal{N}(g_1, \dots, g_n | 0, K)}_{f_0(\mathcal{G})} \prod_i \underbrace{P(y_i | g_i)}_{f_i(g_i)}$$

Some factorization applied to non-Gaussian $P(y_i = 1 | g_i) = 1 / (1 + \exp(-g_i))$. EP approximate non-Gaussian $f_i(g_i)$ by Gaussian to be $\tilde{f}_i(g_i) = \mathcal{N}(\tilde{\mu}_i, \tilde{\psi}_i^2)$. Given $q_{-i}(g_i)$ can be concluded by GP marginalization:

$$q_{-i}(g_i) = \mathcal{N}\left(\Sigma_{i,-i} \Sigma_{-i,-i}^{-1} \tilde{\mu}_{-i}, K_{i,i} - \Sigma_{i,-i} \Sigma_{-i,-i}^{-1} \Sigma_{-i,i}\right)$$

We will consider $\Sigma = \mathbf{K} + \text{diag}(\tilde{\psi}_1^2, \dots, \tilde{\psi}_n^2)$. The update on the Gaussian is based on sufficient statistics matching where $P(g) = f_i(g)q_{-i}(g)$, as now we have the follow:

$$\tilde{f}_{\text{new}}(g_i) = \mathcal{N}\left(\underbrace{\int q_{-i}(g) f_i(g) g \, dg}_{\tilde{\mu}_i^{\text{new}}}, \int q_{-i}(g) f_i(g) g^2 \, dg - (\tilde{\mu}_i^{\text{new}})^2\right) / q_{-i}(g_i)$$

Once approximation site potential have stabilize, we use them to make a prediction from \mathbf{x}' . There are some observations that we made:

- Introduce the prediction point will changes \mathbf{K} , but it will not affect the marginal.
- The unobserved output factor proved no information about g' (the output of \mathbf{x}') and so we don't have to approximate potential \tilde{f}_i
- And so, we have the following prediction:

$$P(y' | \mathbf{x}', \mathcal{D}) = \int P(y' | g') \mathcal{N}\left(g' \mid \mathbf{K}_{x'X} (\mathbf{K}_{XX} + \tilde{\Psi})^{-1} \tilde{\mu}, \mathbf{K}_{x'x'} (\mathbf{K}_{XX} + \tilde{\Psi})^{-1} \mathbf{K}_{Xx'}\right) dg'$$

where $\tilde{\Psi} = \text{diag}(\tilde{\psi}_1^2, \dots, \tilde{\psi}_n^2)$.

5.4 Learning with EP

Remark 50. (Learning and EP) EP can yields the approximate inference posterior. To learn the hyperparameter as we can use:

- Approximate Bayesian inference (like VB): Maybe need to construct a coherent normalizable exponential family on both latent and parameter.
- Approximate EM: As we can maximize $\langle \log P(\mathcal{X}, \mathcal{Z}) \rangle_{q_{\text{EP}}(\mathcal{Z})}$. It is practical but not coherent cost function so no guarantee of convergence even if EP itself converges.
- Direct maximization of EP log-likelihood estimate:
 - Consistent although convergence guarantee still difficult.

– Seem hard to differentiate, but it is simpler than what it looks like.

Proposition 5.4. *We can show that the EP of the moment matching as we now have:*

$$\nabla_\eta \tilde{l} = \sum_{i=1}^N \langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i}$$

where η is the model hyperparameter. It can be computed provided that EP converges.

Proof. Let's consider the derivative of $\log \tilde{C}_i$ as we have:

$$\nabla_\eta \log \tilde{C}_i = \nabla_\eta \Phi_i(\theta_{-i}) - \nabla_\eta \Phi(\theta) = \nabla_\eta \Phi_i(\theta_{-i}) - \boldsymbol{\mu}^T \nabla_\eta \theta$$

Please recall the gradient of the log normalizer in the first equality i.e $\boldsymbol{\mu} = \langle \mathbf{T}(\mathcal{Z}) \rangle_{q(\mathcal{Z})}$. We will have to consider $\nabla_\eta \Phi_i(\theta_{-i})$ as we have (it depends on η via f_i and θ_{-i}):

$$\begin{aligned} \nabla_\eta \Phi_i(\theta_{-i}) &= \frac{1}{\exp(\Phi_i(\theta_{-i}))} \left[\int \exp(\mathbf{T}(\mathcal{Z})^T \theta_{-i}) \nabla_\eta f_i(\mathcal{Z}_i) \, d\mathcal{Z} + \int f_i(\mathcal{Z}_i) \nabla_\eta \exp[\mathbf{T}(\mathcal{Z})^T \theta_{-i}] \, d\mathcal{Z} \right] \\ &= \frac{1}{\exp(\Phi_i(\theta_{-i}))} \int \exp(\mathbf{T}(\mathcal{Z})^T \theta_{-i}) \nabla_\eta f_i(\mathcal{Z}_i) \, d\mathcal{Z} + \int \frac{f_i(\mathcal{Z}_i) \exp(\mathbf{T}(\mathcal{Z})^T \theta_{-i})}{\exp(\Phi_i(\theta_{-i}))} \mathbf{T}(\mathcal{Z})^T \nabla_\eta \theta_{-i} \, d\mathcal{Z} \\ &= \int \frac{f_i(\mathcal{Z}_i) \exp(\mathbf{T}(\mathcal{Z})^T \theta_{-i})}{\exp(\Phi_i(\theta_{-i}))} \nabla_\eta \log f_i(\mathcal{Z}_i) \, d\mathcal{Z} + \int \frac{f_i(\mathcal{Z}_i) \exp(\mathbf{T}(\mathcal{Z})^T \theta_{-i})}{\exp(\Phi_i(\theta_{-i}))} \mathbf{T}(\mathcal{Z})^T \nabla_\eta \theta_{-i} \, d\mathcal{Z} \\ &= \langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i} + \langle \mathbf{T}(\mathcal{Z}) \rangle_{\hat{P}_i}^T \nabla_\eta \theta_{-i} \\ &= \langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i} + \boldsymbol{\mu}^T \nabla_\eta \theta_{-i} \end{aligned}$$

Consider the derivative with respected to the hyperparameter as we have:

$$\begin{aligned} \nabla_\eta \tilde{l} &= \nabla_\eta \Phi(\theta) + \sum_{i=1}^N \nabla_\eta \log \tilde{C}_i \\ &= \boldsymbol{\mu}^T \nabla_\eta \theta + \sum_{i=1}^N \left(\langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i} + \boldsymbol{\mu}^T \nabla_\eta \theta_{-i} - \boldsymbol{\mu}^T \nabla_\eta \theta \right) \\ &= \boldsymbol{\mu}^T \nabla_\eta \left(\sum_{i=1}^N \theta_i + \sum_{i=1}^N (\theta_{-i} - \theta) \right) + \sum_{i=1}^N \left(\langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i} \right) \\ &= \boldsymbol{\mu}^T \nabla_\eta \left(\sum_{i=1}^N (\theta - \theta) \right) + \sum_{i=1}^N \langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i} \\ &= \sum_{i=1}^N \langle \nabla_\eta \log f_i(\mathcal{Z}_i) \rangle_{\hat{P}_i} \end{aligned}$$

Thus complete the proof □

Definition 5.2. (Alpha Divergence) This is a generalization of the KL-divergence to be:

$$D_\alpha[p||q] = \frac{1}{\alpha(1-\alpha)} \int \alpha p(\mathbf{x}) + (1-\alpha)q(\mathbf{x}) - p(\mathbf{x})^\alpha q(\mathbf{x})^{1-\alpha} \, d\mathbf{x}$$

Remark 51. We have the following values for α as we have:

- For $\alpha = -1$:

$$D_{-1}[p||q] = \frac{1}{2} \int \frac{(p(\mathbf{x}) - q(\mathbf{x}))^2}{p(\mathbf{x})} \, d\mathbf{x}$$

- For $\alpha = 1/2$:

$$D_{1/2}[p||q] = 2 \int \left(p(\mathbf{x})^{1/2} - q(\mathbf{x})^{1/2} \right)^2 d\mathbf{x}$$

- For $\alpha = 2$:

$$D_2[p||q] = \frac{1}{2} \int \frac{(p(\mathbf{x}) - q(\mathbf{x}))^2}{q(\mathbf{x})} d\mathbf{x}$$

- Consider the KL-divergence to be:

$$\lim_{\alpha \rightarrow 0} D_\alpha[p||q] = \text{KL}[q||p] \quad \lim_{\alpha \rightarrow 1} D_\alpha[p||q] = \text{KL}[p||q]$$

Remark 52. Local EP minimization gives fixed point update that blends message to power α with previous state approximation:

$$\tilde{f}_{\text{new}} = \arg \min_f \text{KL} \left[f_i(\mathcal{Z}_i)^\alpha \tilde{f}_i(\mathcal{Z})^{1-\alpha} q_{-i}(\mathcal{Z}) \middle\| f(\mathcal{Z}_i) q_{-i}(\mathcal{Z}) \right]$$

Given a small change (like $\alpha < 1$) leads to more stable update and reliable convergence.

6 Belief Propagation: Interpretation

6.1 Introduction

Definition 6.1. (Loopy Propagation) The joint distribution for *any* graph is given by:

$$P(\mathcal{X}) = \frac{1}{Z} \prod_{\text{nodes } i} f_i(\mathbf{x}_i) \prod_{\text{edges } (ij)} f_{ij}(\mathbf{x}_i, \mathbf{x}_j)$$

Message computed recursively with few guarantee of convergence as we have the following message:

$$M_{j \rightarrow i} = \sum_{\mathbf{x}_j} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) f_j(\mathbf{x}_j) \prod_{l \in \text{ne}(j) \setminus \{i\}} M_{l \rightarrow j}(\mathbf{x}_j)$$

The marginal distribution are approximation in general:

$$P(\mathbf{x}_i) \approx b_i(\mathbf{x}_i) \propto f_i(\mathbf{x}_i) \prod_{k \in \text{ne}(i)} M_{k \rightarrow i}(\mathbf{x}_i)$$

$$P(\mathbf{x}_i, \mathbf{x}_j) \approx b_{ij}(\mathbf{x}_i, \mathbf{x}_j) \propto f_{ij}(\mathbf{x}_i, \mathbf{x}_j) f_i(\mathbf{x}_i) f_j(\mathbf{x}_j) \prod_{k \in \text{ne}(i) \setminus \{j\}} M_{k \rightarrow i}(\mathbf{x}_i) \prod_{l \in \text{ne}(j) \setminus \{i\}} M_{l \rightarrow j}(\mathbf{x}_j)$$

Remark 53. (Dealing with Loops) There are various way to deal with loop as we have:

- The belief propagation posterior marginal are approximate on all non-tree because over-counted, but converged approximate are frequently found to be good.
- Converge can be seen in: Tree, Graph with single step, Distribution with weak interaction, Graph with long (and weak) loops, and Gaussian network (variance may also converged).
- Damping, as it is a common approach to enorate of EP:

$$M_{i \rightarrow j}^{\text{new}}(\mathbf{x}_j) = (1 - \alpha) M_{i \rightarrow j}^{\text{old}} + \alpha \sum_{\mathbf{x}_i} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) f_i(\mathbf{x}_i) \prod_{k \in \text{ne}(i) \setminus \{j\}} M_{k \rightarrow i}(\mathbf{x}_i)$$

- Variable can be grouped into cliques to improve accuracy: region graph approximate, cluster variable method, and junction graph.

6.2 Message Based EP

Proposition 6.1. (Loopy BP as Message-Based EP) One can consider the connection between message-based EP and loopy BP, as they are equivalent.

Proof. Consider the approximate pairwise factor \tilde{f}_{ij} as product of messages:

$$f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \approx \tilde{f}_{ij}(\mathbf{x}_i, \mathbf{x}_j) = M_{i \rightarrow j}(\mathbf{x}_j) M_{j \rightarrow i}(\mathbf{x}_i)$$

Consider the approximation of the factorized distribution:

$$\begin{aligned} P(\mathcal{X}) &\approx \frac{1}{Z} \prod_{\text{nodes}(i)} f_i(\mathbf{x}_i) \prod_{\text{edges}(ij)} \tilde{f}_{ij}(\mathbf{x}_i, \mathbf{x}_j) \\ &= \frac{1}{Z} \prod_{\text{nodes}(i)} \left(f_i(\mathbf{x}_i) \prod_{j \in \mathcal{N}(i)} M_{j \rightarrow i}(\mathbf{x}_i) \right) = \prod_{\text{nodes}(i)} b_i(\mathbf{x}_i) \end{aligned}$$

with multiple factors for \mathbf{x}_i , which we consider the update on EP to be:

- *Deletion:* Consider the following $P(\mathcal{X})$ as we have:

$$\begin{aligned} P_{\neg ij}(X_i, X_j) &= \sum_{c \neq i, j} \frac{P(\mathcal{X})}{\tilde{f}(X_i, X_j)} = \sum_{c \neq i, j} \frac{P(\mathcal{X})}{M_{i \rightarrow j}(\mathbf{x}_j) M_{j \rightarrow i}(\mathbf{x}_i)} \\ &= \frac{1}{\tilde{f}(X_i, X_j)} \sum_{c \neq i, j} f_i(\mathbf{x}_i) f_j(\mathbf{x}_j) \prod_{k \in \mathcal{N}(i)} M_{k \rightarrow i}(\mathbf{x}_i) \prod_{l \in \mathcal{N}(j)} M_{l \rightarrow j}(\mathbf{x}_j) \left(\prod_{s \neq i, j} f_s(\mathbf{x}_s) \prod_{t \in \mathcal{N}(s)} M_{t \rightarrow s}(\mathbf{x}_s) \right) \\ &= f_i(\mathbf{x}_i) f_j(\mathbf{x}_j) \prod_{k \in \mathcal{N}(i) \setminus j} M_{k \rightarrow i}(\mathbf{x}_i) \prod_{l \in \mathcal{N}(j) \setminus i} M_{l \rightarrow j}(\mathbf{x}_j) \sum_{c \neq i, j} \left(\prod_{s \neq i, j} f_s(\mathbf{x}_s) \prod_{t \in \mathcal{N}(s)} M_{t \rightarrow s}(\mathbf{x}_s) \right) \\ &= f_i(\mathbf{x}_i) f_j(\mathbf{x}_j) \prod_{k \in \mathcal{N}(i) \setminus j} M_{k \rightarrow i}(\mathbf{x}_i) \prod_{l \in \mathcal{N}(j) \setminus i} M_{l \rightarrow j}(\mathbf{x}_j) \end{aligned}$$

- *Projection:* We consider minimizing the KL-divergence as we have:

$$\{M_{i \rightarrow j}^{\text{new}}, M_{j \rightarrow i}^{\text{new}}\} = \arg \min_{M_{i \rightarrow j}, M_{j \rightarrow i}} \text{KL} \left[f_{ij}(\mathbf{x}_i, \mathbf{x}_j) q_{\neg ij}(\mathbf{x}_i, \mathbf{x}_j) \parallel M_{j \rightarrow i}(\mathbf{x}_i) M_{i \rightarrow j}(\mathbf{x}_j) q_{\neg ij}(\mathbf{x}_i, \mathbf{x}_j) \right]$$

To solve this KL-divergence, this is obvious, as $q_{\neg ij}(\cdot)$ can be factorized and so the minimizer is the marginal between $f_{ij}(\cdot) q_{\neg ij}(\cdot)$, which means that:

$$\begin{aligned} M_{i \rightarrow j}^{\text{new}}(\mathbf{x}_j) q_{\neg ij}(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{\mathbf{x}_j} \left(f_{ij}(\mathbf{x}_i, \mathbf{x}_j) f_j(\mathbf{x}_j) \prod_{l \in \mathcal{N}(j) \setminus i} M_{l \rightarrow j}(\mathbf{x}_j) \right) \underbrace{f_i(\mathbf{x}_i) \prod_{k \in \mathcal{N}(i) \setminus j} M_{k \rightarrow i}(\mathbf{x}_i)}_{q_{\neg ij}(\mathbf{x}_i)} \\ &\implies M_{i \rightarrow j}^{\text{new}}(\mathbf{x}_j) = \sum_{\mathbf{x}_j} \left(f_{ij}(\mathbf{x}_i, \mathbf{x}_j) f_j(\mathbf{x}_j) \prod_{l \in \mathcal{N}(j) \setminus i} M_{l \rightarrow j}(\mathbf{x}_j) \right) \end{aligned}$$

This is the Loopy BP update, and so both of the are equivalent. □

Remark 54. (Comments on the Loopy BP and EP) There are some observation that we can make in the equivalent between loopy BP and EP algorithm:

- Unlike EP, this message based EP doesn't need 2 separate approximate as we have in the normal EP.
- This message based EP is loopy graph can be seen as a more constraint on approximate site and not just exponential family factor but the product of exponential family message.
- On a tree, message forward EP finds the same marginal as standard EP as the messages are calculated the same way. Similarly, the pairwise marginal can be found after converge by compute $\tilde{P}(z_{i-1}, z_i)$
- Factorization still remain valid even when original site lies in the approximation exponential family already, so the loopy BP can be seen as form of EP.
- This doesn't help us with understanding the convergence property of EP.

6.3 Reparameterized on Tree

Remark 55. (Tree-Based Representation) We consider the joint factorization, which can be represented:

$$\begin{aligned}
P(\mathcal{X}) &= \frac{1}{Z} \prod_{\text{nodes}(i)} f_i(\mathbf{x}_i) \prod_{\text{edges}(ij)} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) && \text{(Undirected Tree)} \\
&= P(\mathbf{x}_i) \prod_{i \neq r} P(\mathbf{x}_i | \mathbf{x}_{\text{pa}(i)}) && \text{(Directed Rooted Tree)} \\
&= \prod_{\text{nodes}(i)} P(\mathbf{x}_i) \prod_{\text{edges}(ij)} \frac{P(\mathbf{x}_i, \mathbf{x}_j)}{P(\mathbf{x}_i)P(\mathbf{x}_j)} && \text{(Pairwise Marginal)}
\end{aligned}$$

The last one requires that $\sum_{\mathbf{x}_j} P(\mathbf{x}_i, \mathbf{x}_j) = P(\mathbf{x}_i)$.

- The undirected tree isn't unique as if we multiply the factor $f_{ij}(\mathbf{x}_i, \mathbf{x}_j)$ by $g(\mathbf{x}_i)$ and dividing $f_i(\mathbf{x}_i)$ by the same $g(\mathbf{x}_i)$ doesn't change the distribution.
- BP can be seen as interactive replacement of $f_i(\mathbf{x}_i)$ by local marginal of $p_{ij}(\mathbf{x}_i, \mathbf{x}_j)$ along with corresponding representation of $f_{ij}(\mathbf{x}_i, \mathbf{x}_j)$ (recall the Hugin propagation)
- Converged BP on a tree finds $P(\mathbf{x}_i)$ and $P(\mathbf{x}_i, \mathbf{x}_j)$ allowing up to transform the undirected tree to pairwise marginal.

Remark 56. (Reparameterization in Tree) To consider the tree based reparameterization, we want to transform the representation from undirected tree to pairwise marginal as:

$$\prod_{\text{nodes}(i)} f_i(\mathbf{x}_i) \prod_{\text{edges}(ij)} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \implies \prod_{\text{nodes}(i)} P(\mathbf{x}_i) \prod_{\text{edges}(ij)} \frac{P(\mathbf{x}_i, \mathbf{x}_j)}{P(\mathbf{x}_i)P(\mathbf{x}_j)}$$

We will define the $f_{ij}^0 = f_{ij}$, while the singleton factor to be $f_i^0 = p_i^0 = 1$, we consider the following update: The update is based on the fact that we will act on the factors *as if* it is actually representing the probabilities: We will consider such a procedure on a node that has 2 incoming messages.

- Starting with joint, where if we multiply it by adjacent factors we get $P(\mathbf{x}_i, \mathbf{x}_j)$ i.e

$$p^{(n)}(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{Z_{ij}^{(n)}} f_i^{(n-1)}(\mathbf{x}_i) f_{ij}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_j) f_j^{(n-1)}(\mathbf{x}_j)$$

- Finding the marginal, as we have:

$$f_i^{(n)}(\mathbf{x}_i) = p^{(n)}(\mathbf{x}_i) = \sum_{\mathbf{x}_j} p^{(n)}(\mathbf{x}_i, \mathbf{x}_j) = f_i^{(n-1)}(\mathbf{x}_i) \underbrace{\sum_{\mathbf{x}_j} f_{ij}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_j) f_j^{(n-1)}(\mathbf{x}_j)}_{M_{j \rightarrow i}}$$

- To keep the normalization correctly, we divide the message so that the update on one passing giving us normalized term:

$$f_{ij}^{(n)} = \frac{f_{ij}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_j)}{M_{j \rightarrow i}(\mathbf{x}_j)}$$

- We now consider the next step with the next incoming message from node k to node i :

$$\begin{aligned} p^{(n)}(\mathbf{x}_i, \mathbf{x}_k) &= \frac{1}{Z_{ik}^{(n)}} f_i^{(n)}(\mathbf{x}_i) f_{ik}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_j) f_k^{(n-1)}(\mathbf{x}_j) \\ &= \frac{1}{Z_{ik}^{(n)}} f_i^{(n-1)}(\mathbf{x}_i) M_{j \rightarrow i}(\mathbf{x}_i) f_{ik}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_k) f_k^{(n-1)}(\mathbf{x}_k) \end{aligned}$$

- Finding the singleton factor by marginalization

$$f_i^{(n)}(\mathbf{x}_i) = f_i^{(n-1)}(\mathbf{x}_i) M_{j \rightarrow i}(\mathbf{x}_i) \underbrace{\sum_{\mathbf{x}_k} f_{ik}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_k) f_k^{(n-1)}(\mathbf{x}_k)}_{M_{k \rightarrow i}}$$

- And so, the normalization correction on the joint factor is:

$$f_{ik}^{(n)} = \frac{f_{ik}^{(n-1)}(\mathbf{x}_i, \mathbf{x}_k)}{M_{k \rightarrow i}(\mathbf{x}_i)}$$

We perform this update throughout the tree, which we do it in forward (e.g $i \rightarrow j$) and backward (e.g $j \rightarrow i$) manner, which gives us:

$$\begin{aligned} f_i^{(\infty)}(\mathbf{x}_i) &= \prod_{j \in \text{ne}(i)} M_{j \rightarrow i}(\mathbf{x}_i) = P(\mathbf{x}_i) \\ f_{ij}^{(\infty)}(\mathbf{x}_i, \mathbf{x}_j) &= \frac{f_{ij}(\mathbf{x}_i, \mathbf{x}_j)}{M_{j \rightarrow i}(\mathbf{x}_i) M_{i \rightarrow j}(\mathbf{x}_j)} \\ &= \frac{\prod_{k \in \text{ne}(i) \setminus j} M_{k \rightarrow i} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \prod_{l \in \text{ne}(j) \setminus i} M_{l \rightarrow i}(\mathbf{x}_j)}{\prod_{k \in \text{ne}(i) \setminus j} M_{k \rightarrow i} M_{j \rightarrow i}(\mathbf{x}_i) M_{i \rightarrow j}(\mathbf{x}_j) \prod_{l \in \text{ne}(j) \setminus i} M_{l \rightarrow i}(\mathbf{x}_j)} \\ &= \frac{\prod_{k \in \text{ne}(i) \setminus j} M_{k \rightarrow i} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \prod_{l \in \text{ne}(j) \setminus i} M_{l \rightarrow i}(\mathbf{x}_j)}{\prod_{k \in \text{ne}(i)} M_{k \rightarrow i} \prod_{l \in \text{ne}(j)} M_{l \rightarrow i}(\mathbf{x}_j)} \\ &= \frac{P(\mathbf{x}_i, \mathbf{x}_j)}{P(\mathbf{x}_i) P(\mathbf{x}_j)} \end{aligned}$$

the equation follows from the result from belief propagation. This kind of reparameterization allows us to avoid double counting, which is essentially a book-keeping method, especially the normalizing part (there will be a case where the factor cancel with unnecessary message from singleton factor, as intended).

Remark 57. (Comments on the BP on non-tree) If this converges in a non-tree setting, then we have locally consistent belief i.e:

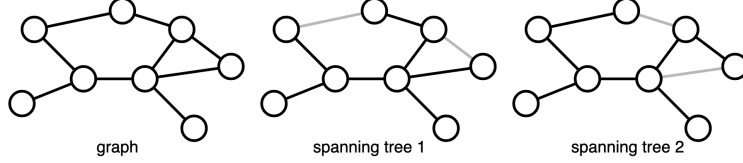
$$p(\mathcal{X}) \propto \prod_i b(\mathbf{x}_i) \prod_{ij} \frac{b(\mathbf{x}_i, \mathbf{x}_j)}{b(\mathbf{x}_i) b(\mathbf{x}_j)} \quad \text{such that} \quad \sum_{\mathbf{x}_j} b(\mathbf{x}_i, \mathbf{x}_j) = b(\mathbf{x}_i)$$

But it doesn't need to be globally consistent:

$$\sum_{\mathcal{X}_{-i}} \left(\prod_i b(\mathbf{x}_i) \prod_{ij} \frac{b(\mathbf{x}_i, \mathbf{x}_j)}{b(\mathbf{x}_i) b(\mathbf{x}_j)} \right) \neq b(\mathbf{x}_i)$$

This kind of marginal is called *pseudo-marginal*.

Remark 58. (Message Schedule Scheme) We consider update the belief on each *subtree* of the graph and passing message on each subtree, looping through all the subtree until converge:



And, we now that the following updates steps:

$$\begin{aligned}
P(\mathcal{X}) &= \frac{1}{Z} \prod_{\text{nodes}(i)} f_i^{(0)}(\mathbf{x}_i) \prod_{\text{edges}(ij)} f_{ij}^{(0)}(\mathbf{x}_i, \mathbf{x}_j) \\
&= \frac{1}{Z} \prod_{\text{nodes}(i) \in T_1} f_i^{(0)}(\mathbf{x}_i) \prod_{\text{edges}(ij) \in T_1} f_{ij}^{(0)}(\mathbf{x}_i, \mathbf{x}_j) \left(\prod_{\text{edges}(ij) \notin T_1} f_{ij}^{(0)}(\mathbf{x}_i, \mathbf{x}_j) \right) \\
&\text{(Update)} \\
&= \frac{1}{Z} \prod_{\text{nodes}(i) \in T_1} f_i^{(1)}(\mathbf{x}_i) \prod_{\text{edges}(ij) \in T_1} f_{ij}^{(1)}(\mathbf{x}_i, \mathbf{x}_j) \left(\prod_{\text{edges}(ij) \notin T_1} f_{ij}^{(1)}(\mathbf{x}_i, \mathbf{x}_j) \right) \\
&\text{(Next Tree)} \\
&= \frac{1}{Z} \prod_{\text{nodes}(i) \in T_2} f_i^{(1)}(\mathbf{x}_i) \prod_{\text{edges}(ij) \in T_2} f_{ij}^{(1)}(\mathbf{x}_i, \mathbf{x}_j) \left(\prod_{\text{edges}(ij) \notin T_2} f_{ij}^{(1)}(\mathbf{x}_i, \mathbf{x}_j) \right) \\
&\dots
\end{aligned}$$

where we have, when we got a new tree,

$$f_i^{(1)}(\mathbf{x}_i) = P^{T_1}(\mathbf{x}_i) \quad f_{ij}^{(1)}(\mathbf{x}_i, \mathbf{x}_j) = \frac{P^{T_1}(\mathbf{x}_i, \mathbf{x}_j)}{P^{T_1}(\mathbf{x}_i)P^{T_2}(\mathbf{x}_j)}$$

If the process converges, suppose it converge to:

$$P(\mathcal{X}) = \frac{1}{Z} \prod_{\text{nodes}(i)} f_i^{(\infty)}(\mathbf{x}_i) \prod_{\text{edges}(ij)} f_{ij}^{(\infty)}(\mathbf{x}_i, \mathbf{x}_j)$$

where for any tree T in the graph, we have:

$$f_i^{(\infty)} = P^T(\mathbf{x}_i) \quad f_{ij}^{(\infty)} = \frac{P^T(\mathbf{x}_i, \mathbf{x}_j)}{P^T(\mathbf{x}_i)P^T(\mathbf{x}_j)}$$

This means that the local marginal of all subtree are consistent with each other, and the pseudo-marginal is valid belief of any of the subtree, as this is stronger constraint.

6.4 Bathe Free Energy

Remark 59. (Introduction to Bathe Free Energy) In reparameterization view, BP solves for marginal belief $b_{ij}(\mathbf{x}_i, \mathbf{x}_j)$ and $b_i(\mathbf{x}_i) = \sum_{\mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j)$ such that:

$$P(\mathcal{X}) \propto \prod_i f_i(\mathbf{x}_i) \prod_{ij} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \propto \prod_i b_i(\mathbf{x}_i) \prod_{ij} \frac{b_{ij}(\mathbf{x}_i, \mathbf{x}_j)}{b_i(\mathbf{x}_i)b_j(\mathbf{x}_j)}$$

Loopy BP is a set of fixed point equation for finding stationary of an objective function called Bathe free energy, which is defined in terms of locally consistent belief (pseudo-marginal) $b_i \geq 0$ and $b_{ij} \geq 0$ such that:

$$\sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) = 1 \quad \sum_{\mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) = b_i(\mathbf{x}_i)$$

Definition 6.2. (Bathe Free Energy) We define it in the form of:

$$\mathcal{F}_{\text{bathe}}(b) = \mathcal{E}_{\text{bathe}}(b) + \mathcal{H}_{\text{bathe}}(b)$$

Both terms are approximated so that it corresponds to variational likelihood terms:

- Bathe average energy is the expected log-joint evaluate as though the pseudomarginal were correct:

$$\mathcal{E}_{\text{bathe}}(b) = \sum_i \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log f_i(\mathbf{x}_i) + \sum_{ij} \sum_{\mathbf{x}_i \mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) \log f_{ij}(\mathbf{x}_i, \mathbf{x}_j)$$

- Bathe entropy is the sum of pseudomarginal entropies corrected for pairwise (pseudo-)interaction, but neglecting higher-order dependence:

$$\begin{aligned} \mathcal{H}_{\text{bathe}}(b) &= \sum_i H[b_i] - \sum_{ij} \text{KL}[b_{ij}|b_i b_j] \\ &= - \sum_i \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log b_i(\mathbf{x}_i) - \sum_{ij} \sum_{\mathbf{x}_i \mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) \log \frac{b_{ij}(\mathbf{x}_i, \mathbf{x}_j)}{b_i(\mathbf{x}_i) b_j(\mathbf{x}_j)} \end{aligned}$$

On tree, both belief and the bathe entropy expression are correct i.e $\mathcal{F}_{\text{bathe}} = \mathcal{F}$. The update rule can be recovered from finding the fixed point.

Proposition 6.2. (Fixed Point for Bathe Free Energy) *The fixed point for Bathe free energy is:*

$$\begin{aligned} b_i(\mathbf{x}_i) &\propto f_i(\mathbf{x}_i) \prod_{j \in \text{ne}(i)} \exp(-\xi_{ij}(\mathbf{x}_i)) \\ b_{ij}(\mathbf{x}_i, \mathbf{x}_j) &\propto f_{ij}(\mathbf{x}_i, \mathbf{x}_j) b_i(\mathbf{x}_i) b_j(\mathbf{x}_j) \exp(\xi_{ij}(\mathbf{x}_i) + \xi_{ij}(\mathbf{x}_j)) \\ \exp(-\xi_{ij}(\mathbf{x}_i)) &\propto \sum_{\mathbf{x}_j} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) f_j(\mathbf{x}_j) \prod_{l \in \text{ne}(j) \setminus i} \exp(-\xi_{lj}(\mathbf{x}_j)) \end{aligned}$$

Proof. We find the Lagrangian with local consistency and normalization, which is given as:

$$\begin{aligned} \mathcal{L} &= \sum_i \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log f_i(\mathbf{x}_i) + \sum_{ij} \sum_{\mathbf{x}_i \mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) \log f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad - \sum_i \sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) \log b_i(\mathbf{x}_i) - \sum_{ij} \sum_{\mathbf{x}_i \mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) \log \frac{b_{ij}(\mathbf{x}_i, \mathbf{x}_j)}{b_i(\mathbf{x}_i) b_j(\mathbf{x}_j)} \\ &\quad + \sum_i \xi_i \left(\sum_{\mathbf{x}_i} b_i(\mathbf{x}_i) - 1 \right) \\ &\quad + \sum_{ij} \left[\sum_{\mathbf{x}_i} \xi_{ij}(\mathbf{x}_i) \left(\sum_{\mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) - b_i(\mathbf{x}_i) \right) + \sum_{\mathbf{x}_j} \xi_{ij}(\mathbf{x}_j) \left(\sum_{\mathbf{x}_i} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) - b_j(\mathbf{x}_j) \right) \right] \end{aligned}$$

Setting the derivate to zero, which gives us the solution:

$$\begin{aligned} \frac{\partial}{\partial b_i(\mathbf{x}_i)} &= \log f_i(\mathbf{x}_i) - \log b_i(\mathbf{x}_i) + \sum_{j \in \text{ne}(j)} \sum_{\mathbf{x}_j} \frac{b_{ij}(\mathbf{x}_i, \mathbf{x}_j)}{b_i(\mathbf{x}_i)} + \xi_i - \sum_{j \in \text{ne}(i)} \xi_{ij}(\mathbf{x}_i) + \text{const.} = 0 \\ \implies b_i(\mathbf{x}_i) &\propto f_i(\mathbf{x}_i) \prod_{j \in \text{ne}(i)} \exp(-\xi_{ij}(\mathbf{x}_i)) \end{aligned}$$

$$\begin{aligned} \frac{\partial f}{\partial b_{ij}(\mathbf{x}_i, \mathbf{x}_j)} &= \log f_{ij}(\mathbf{x}_i, \mathbf{x}_j) - \log b_{ij}(\mathbf{x}_i, \mathbf{x}_j) + \log b_i(\mathbf{x}_i)b_j(\mathbf{x}_j) + \xi_{ij}(\mathbf{x}_i) + \xi_{ji}(\mathbf{x}_j) + \text{const.} = 0 \\ &\implies b_{ij}(\mathbf{x}_i, \mathbf{x}_j) \propto f_{ij}(\mathbf{x}_i, \mathbf{x}_j)b_i(\mathbf{x}_i)b_j(\mathbf{x}_j) \exp(\xi_{ij}(\mathbf{x}_i) + \xi_{ji}(\mathbf{x}_j)) \end{aligned}$$

To solve for $\xi_{ij}(\mathbf{x}_i)$ by enforcing the constant $\sum_{\mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) = b_i(\mathbf{x}_i)$ where we have:

$$\begin{aligned} \sum_{\mathbf{x}_j} b_{ij}(\mathbf{x}_i, \mathbf{x}_j) &\propto \sum_{\mathbf{x}_j} f_{ij}(\mathbf{x}_i, \mathbf{x}_j)b_i(\mathbf{x}_i)b_j(\mathbf{x}_j) \exp(\xi_{ij}(\mathbf{x}_i) + \xi_{ji}(\mathbf{x}_j)) \\ \implies b_i(\mathbf{x}_i) &\propto b_i(\mathbf{x}_i) \exp(\xi_{ij}(\mathbf{x}_i)) \sum_{\mathbf{x}_j} f_{ij}(\mathbf{x}_i, \mathbf{x}_j)b_j(\mathbf{x}_j) \exp(\xi_{ji}(\mathbf{x}_j)) \\ \implies \exp(-\xi_{ij}(\mathbf{x}_i)) &\propto \sum_{\mathbf{x}_j} f_{ij}(\mathbf{x}_i, \mathbf{x}_j)b_j(\mathbf{x}_j) \exp(\xi_{ji}(\mathbf{x}_j)) \\ &= \sum_{\mathbf{x}_j} f_{ij}(\mathbf{x}_i, \mathbf{x}_j)f_j(\mathbf{x}_j) \prod_{l \in \text{ne}(j) \setminus i} \exp(-\xi_{jl}(\mathbf{x}_j)) \end{aligned}$$

□

Remark 60. (Interpretation of Results) Comparing with BP, we have the message to be of the form of $M_{j \rightarrow i}(\mathbf{x}_i) = \exp(-\xi_{ij}(\mathbf{x}_i))$. The fixed point for bathe free energy recovers the message passing rule:

- Stable Fixed point of loopy BP are stationary point of Bathe and local minimum of Bathe free energy.
- For binary attractive network: the Bathe free energy at fixed point of loopy BP provides an upperbound on the log partition function $\log P(\mathbf{Z})$.
- It is useful for learning undirected graphical model as it leads to lower bound on the log-likelihood.
- Belief b_i and b_{ij} in loopy BP are only locally consistent pseudomarginal, not necessary consistent with marginal or the implied joint distribution.
- Bathe free energy accounts for interaction between difference states, while variational free energy that assume independence.
- The log series Plefka expansion of the log-partition Z : the variational energy form the first order while Bathe free energy contains higher term.
- Loopy BP tends to significantly more accurate whenever it converges.

Remark 61. (Extensions and Variations)

- Generalized BP is a group variable together to threat their interaction exactly.
- The algorithm can be derived so that the Bathe free energy at every step and thus guarantee the convergence. Similarly, convex alternative and we will converge to unique global maximum.
- The treatment of loopy Viterbi or max-product algorithm is difference.

7 Exponential families: Convexity, Duality and Free Energies

Definition 7.1. (Log-Partition Function) Consider the exponential family distribution with sufficient statistics $\mathbf{s}(\mathbf{x})$ and natural parameter $\boldsymbol{\theta}$. We have the following probability to be $P(\mathbf{x}|\boldsymbol{\theta}) = \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x}) - \Phi(\boldsymbol{\theta}))$, where $\Phi(\boldsymbol{\theta})$ is the log-partition, where:

$$\Phi(\boldsymbol{\theta}) = \log \int \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x}$$

Proposition 7.1. (Derivative of Log Partitions) We can show that:

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \Phi(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\mathbf{s}(\mathbf{x})] = \boldsymbol{\mu}(\boldsymbol{\theta}) \\ \frac{\partial^2}{\partial \boldsymbol{\theta}^2} \Phi(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\mathbf{s}(\mathbf{x})^2] - \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\mathbf{s}(\mathbf{x})]^2 = \mathbb{V}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\mathbf{s}(\mathbf{x})]\end{aligned}$$

The second derivative is positive semi-definite and so $\Phi(\boldsymbol{\theta})$ is convex in $\boldsymbol{\theta}$.

Proof. The first result is the old results. Now, consider the second derivation as we have:

$$\begin{aligned}\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \Phi(\boldsymbol{\theta}) &= \frac{\partial}{\partial \boldsymbol{\theta}} \left[\exp(-\Phi(\boldsymbol{\theta})) \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \right] \\ &= \left[\frac{\partial}{\partial \boldsymbol{\theta}} \exp(-\Phi(\boldsymbol{\theta})) \right] \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \\ &\quad + \exp(-\Phi(\boldsymbol{\theta})) \left[\frac{\partial}{\partial \boldsymbol{\theta}} \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \right] \\ &= -\exp(-\Phi(\boldsymbol{\theta})) \left[\frac{\partial}{\partial \boldsymbol{\theta}} \Phi(\boldsymbol{\theta}) \right] \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \\ &\quad + \exp(-\Phi(\boldsymbol{\theta})) \left[\frac{\partial}{\partial \boldsymbol{\theta}} \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \right] \\ &= -\exp(-\Phi(\boldsymbol{\theta})) \left[\frac{\partial}{\partial \boldsymbol{\theta}} \Phi(\boldsymbol{\theta}) \right] \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \\ &\quad + \exp(-\Phi(\boldsymbol{\theta})) \left[\int \mathbf{s}(\mathbf{x})^2 \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \right] \\ &= - \left[\exp(-\Phi(\boldsymbol{\theta})) \int \mathbf{s}(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \right]^2 + \exp(-\Phi(\boldsymbol{\theta})) \left[\int \mathbf{s}(\mathbf{x})^2 \exp(\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})) \, d\mathbf{x} \right] \\ &= -\mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\mathbf{s}(\mathbf{x})^2] + \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\mathbf{s}(\mathbf{x})]^2\end{aligned}$$

□

Definition 7.2. (Convex Duality/Conjugate) Given a function $f^*(\mathbf{x})$, we denote the duality/conjugate of the function is:

$$f(\mathbf{y}) = \sup_{\mathbf{x} \in \mathcal{X}} (\mathbf{y}^T \mathbf{x} - f(\mathbf{x}))$$

Definition 7.3. (Negative Entropy of Mean parameter) The negative entropy of the distribution as a function of mean parameter:

$$\Psi(\boldsymbol{\mu}) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})}[\log P(\mathbf{x}|\boldsymbol{\theta})] = \boldsymbol{\theta}^T \mathbb{E}[\mathbf{s}(\mathbf{x})] - \Phi(\boldsymbol{\theta}) = \boldsymbol{\theta}^T \boldsymbol{\mu} - \Phi(\boldsymbol{\theta})$$

Note that the variable $\boldsymbol{\mu}$ is arbitrary:

Lemma 7.1. We can observe that the derivative of negative entropy recovers the natural parameter

$$\frac{d}{d\boldsymbol{\mu}} \Psi(\boldsymbol{\mu}) = \boldsymbol{\theta}$$

This related to dual function.

Proof. We have the following distribution:

$$\begin{aligned}\frac{d}{d\boldsymbol{\mu}} \Psi(\boldsymbol{\mu}) &= \frac{\partial}{\partial \boldsymbol{\mu}} (\boldsymbol{\theta}^T \boldsymbol{\mu} - \Phi(\boldsymbol{\theta})) + \frac{d\boldsymbol{\theta}}{d\boldsymbol{\mu}} \frac{\partial}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}^T \boldsymbol{\mu} - \Phi(\boldsymbol{\theta})) \\ &= \boldsymbol{\theta} + \frac{d\boldsymbol{\theta}}{d\boldsymbol{\mu}} (\boldsymbol{\mu} - \boldsymbol{\mu}) = \boldsymbol{\theta}\end{aligned}$$

□

Proposition 7.2. (Duality Between Partition and Entropy) One can show that partition and entropy are duality of each other:

$$\Psi(\boldsymbol{\mu}) = \sup_{\boldsymbol{\theta}'} \left[(\boldsymbol{\theta}')^T \boldsymbol{\mu} - \Phi(\boldsymbol{\theta}') \right] \quad \Phi(\boldsymbol{\theta}) = \sup_{\boldsymbol{\mu}'} \left[\boldsymbol{\theta}^T \boldsymbol{\mu}' - \Psi(\boldsymbol{\mu}') \right]$$

Proof. Consider the KL-divergence between distribution with natural parameter $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$ as we have:

$$\begin{aligned} \text{KL}[p||q] &= \text{KL} \left[P(\mathbf{x}|\boldsymbol{\theta}) \middle| \middle| P(\mathbf{x}|\boldsymbol{\theta}') \right] \\ &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})} \left[-\log P(\mathbf{x}|\boldsymbol{\theta}') + \log P(\mathbf{x}|\boldsymbol{\theta}) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})} \left[-\log P(\mathbf{x}|\boldsymbol{\theta}') \right] + \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}|\boldsymbol{\theta})} \left[\log P(\mathbf{x}|\boldsymbol{\theta}) \right] \\ &= -(\boldsymbol{\theta}')^T \boldsymbol{\mu} + \Phi(\boldsymbol{\theta}') + \Psi(\boldsymbol{\mu}) \geq 0 \end{aligned}$$

We can see that $\Psi(\boldsymbol{\mu}) \geq (\boldsymbol{\theta}')^T \boldsymbol{\mu} - \Phi(\boldsymbol{\theta}')$. With the minimization of KL to be $\boldsymbol{\theta} = \boldsymbol{\theta}'$, and so we have the first dual formula, while the second formula comes from the dual property (dual of dual is the function itself). \square

Remark 62. (Another Proof) One can consider the sup by finding derivative and set it to zero:

$$\frac{\partial}{\partial \boldsymbol{\theta}'} \left[(\boldsymbol{\theta}')^T \boldsymbol{\mu} - \Phi(\boldsymbol{\theta}') \right] = \boldsymbol{\mu} - \boldsymbol{\mu}(\boldsymbol{\theta}') = 0$$

Setting to zero and we have that $\boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\theta}')$, plugging it back and we get the negative entropy. This is the same as the second formula as:

$$\frac{\partial}{\partial \boldsymbol{\mu}'} \left[\boldsymbol{\theta}^T \boldsymbol{\mu}' - \Psi(\boldsymbol{\mu}') \right] = \boldsymbol{\theta} - \boldsymbol{\theta}(\boldsymbol{\mu}') = 0$$

Note that $\boldsymbol{\theta}(\boldsymbol{\mu})$ doesn't have a full definition but we can consider it from the definition of negative entropy. Thus, complete the proof.

Remark 63. (Free Energy and Dual) We have the following components:

- Joint exponential family distribution of observed \mathbf{x} and \mathbf{z} :

$$P(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) = \exp \left[\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x}, \mathbf{z}) - \Phi_{xz}(\boldsymbol{\theta}) \right]$$

- The posterior \mathbf{z} is in exponential family with clamped sufficient statistics:

$$P(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) = \exp \left[\boldsymbol{\theta}^T \mathbf{s}_z(\mathbf{z}; \mathbf{x}) - \Phi_z(\boldsymbol{\theta}) \right]$$

where we have $s_z(\mathbf{z}; \mathbf{x}) = s_{xz}(\mathbf{x}^{\text{obs}}; \mathbf{z})$

We consider the likelihood (of observed variables) to be as:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= P(\mathbf{x}|\boldsymbol{\theta}) = \int \exp \left[\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x}, \mathbf{z}) - \Phi_{xz}(\boldsymbol{\theta}) \right] d\mathbf{z} \\ &= \int \exp \left[\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x}, \mathbf{z}) \right] \exp \left[-\Phi_{xz}(\boldsymbol{\theta}) \right] d\mathbf{z} \\ &= \exp \left[-\Phi_{xz}(\boldsymbol{\theta}) \right] \int \exp \left[\boldsymbol{\theta}^T \mathbf{s}(\mathbf{z}; \mathbf{x}) \right] d\mathbf{z} \\ &= \exp \left[-\Phi_{xz}(\boldsymbol{\theta}) \right] \exp \left[\Phi_z(\boldsymbol{\theta}) \right] \end{aligned}$$

We can see that the log-likelihood is given as:

$$\begin{aligned}\log \mathcal{L}(\boldsymbol{\theta}) &= l(\boldsymbol{\theta}) = \Phi_z(\boldsymbol{\theta}) - \Phi_{xz}(\boldsymbol{\theta}) \\ &= \sup_{\boldsymbol{\mu}_z} \left[\boldsymbol{\theta}^T \boldsymbol{\mu}_z - \Psi(\boldsymbol{\mu}_z) \right] - \Phi_{xz}(\boldsymbol{\theta}) \\ &= \sup_{\boldsymbol{\mu}_z} \left[\underbrace{\boldsymbol{\theta}^T \boldsymbol{\mu}_z - \Phi_{xz}(\boldsymbol{\theta})}_{\langle \log P(\mathbf{x}, \mathbf{z}) \rangle_q} - \underbrace{\Psi(\boldsymbol{\mu}_z)}_{-H[q]} \right] = \sup_{\boldsymbol{\mu}_z} \mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\mu}_z)\end{aligned}$$

Please note that, we can see $\boldsymbol{\mu}_z = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})}$ for some variational distribution q (that we want to optimize). And so, one can recover the free energy formulation from the duality of entropy and expected sufficient statistics.

Remark 64. (Optimization of Free Energy) We can see that the optimization over the $\boldsymbol{\mu}$ directly, as:

$$\boldsymbol{\mu}_z^* = \arg \max_{\boldsymbol{\mu}_z} \left[\boldsymbol{\theta}^T \boldsymbol{\mu}_z - \Psi(\boldsymbol{\mu}_z) \right]$$

This is concave maximization; however, we have 2 complications to the problem:

- The optimization must be found out feasible means. Interdependence of the sufficient statistics may prevent arbitrary sets of mean sufficient statistics being achieved.
- The feasible means are convex combination of all single configuration of sufficient statistics:

$$\boldsymbol{\mu} = \int \nu(\mathbf{x}) \mathbf{s}(\mathbf{x}) \, d\mathbf{x} \quad \int \mu(\mathbf{x}) \, d\mathbf{x} = 1$$

- Let's consider the Boltzman machine on 2 variables, as we have:

- Consider the Boltzman machine on 2 variables as we have x_1 and x_2 :

$$E = \frac{1}{Z} \exp \left(- \sum_{i < j} w_{ij} x_i x_j + \sum_i \theta_i x_i \right)$$

- Sufficient statistics of the Boltzman machine is: $\mathbf{s}(\mathbf{x}_1, \mathbf{x}_2) = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_1 \mathbf{x}_2]$
- There are only 4 states as we have: $\mathcal{S} = \{[0, 0, 0], [0, 1, 0], [1, 0, 0], [1, 1, 1]\}$
- It is clear that $\boldsymbol{\mu}$ is in the convex hull of \mathcal{S} .
- For discrete distribution, the space of possible mean is bounded exponentially many hyperplanes connected the discrete configuration set called marginal polytrope.
- Even with marginal polytrope to marginal polytrope, evaluating $\Psi(\boldsymbol{\mu})$ is challenging.

Remark 65. (Undirected Tree and Markov Random Field) We can parameterize a discrete MRF:

$$\begin{aligned}P(\mathcal{X}) &= \frac{1}{Z} \prod_i f_i(\mathbf{x}_i) \prod_{ij} f_{ij}(\mathbf{x}_i, \mathbf{x}_j) \\ &= \exp \left(\sum_i \sum_k \theta_i(k) \delta(\mathbf{x}_i = k) + \sum_{ij} \sum_{k,l} \theta_{ij}(k,l) \delta(\mathbf{x}_i = k) \delta(\mathbf{x}_j = l) - \Phi(\boldsymbol{\theta}) \right)\end{aligned}$$

Discrete MRF are always exponential family, with natural and mean parameter:

$$\boldsymbol{\theta} = \left[\theta_i(k) \quad \theta_{ij}(k,l) \right] \quad \boldsymbol{\mu} = \left[P(\mathbf{x}_i = k) \quad P(\mathbf{x}_i = k, \mathbf{x}_j = l) \right]$$

for all i, j, k, l . In particular, the mean parameter are just singleton and pairwise probability tables.

Remark 66. (Entropy of MRF) If the MRF has tree structure T and the negated entropy, can be written in term of single state entropy and mutual information on edges as (note the reparameterization and the tree structure):

$$\begin{aligned}\Psi(\boldsymbol{\mu}_T) &= \mathbb{E}_{\mathbf{x}_i \sim P(\mathbf{x}|\boldsymbol{\theta}_T)} \left[\log \prod_i P(\mathbf{x}_i) \prod_{(ij) \in \text{edge}(T)} \frac{P(\mathbf{x}_i, \mathbf{x}_j)}{P(\mathbf{x}_i)P(\mathbf{x}_j)} \right] \\ &= - \sum_i H(\mathbf{x}_i) + \sum_{(ij) \in \text{edge}(T)} I(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

Remark 67. (Bathe Energy of Free Energy) Let's consider the Bathe free energy, as a relaxation of the free-energy optimization, which is denoted as:

$$\boldsymbol{\mu}_z^* = \arg \max_{\boldsymbol{\mu}_z} \left[\boldsymbol{\theta}^T \boldsymbol{\mu}_z - \Psi(\boldsymbol{\mu}_z) \right]$$

Consider the set \mathcal{M} as the set of feasible means. Now, we have the following relaxation:

- *Feasible Set* Relax the \mathcal{M} to be \mathcal{L} be the set of locally consistent means (all nested means marginalized correctly), instead of being globally estimate.
- *Approximate* $\Psi(\boldsymbol{\mu}_z)$ is the entropy of an arbitrary graph. However, it is hard to evaluate it exactly, and so we consider a result from tree structure but in the final calculation, we consider every edges:

$$\Psi_{\text{Bathe}}(\boldsymbol{\mu}_z) = - \sum_i H(\mathbf{x}_i) + \sum_{(ij) \in \text{edge}(G)} I(\mathbf{x}_i, \mathbf{x}_j)$$

Please note that \mathcal{L} is still a convex set. However, Ψ_{Bathe} isn't convex (Please recall the dual formulation as we can find Φ instead).

Remark 68. (Upperbound on Log Partition) Consider upperbound on $\Phi(\boldsymbol{\theta})$, imagine a set of spanning tree T for the MRF, with corresponding: $\boldsymbol{\theta}_T$ and $\boldsymbol{\mu}_T$ (can be done by padding the entries corresponding to off-tree edges with zero, and so $\boldsymbol{\theta}_T$ have the same size of $\boldsymbol{\theta}$). Consider the following:

- Consider the distribution over the spanning tree β as, we need $\mathbb{E}_\beta[\boldsymbol{\theta}_T] = \boldsymbol{\theta}$. Using the convexity of Φ :

$$\Phi(\boldsymbol{\theta}) = \Phi(\mathbb{E}_\beta[\boldsymbol{\theta}_T]) \leq \mathbb{E}_\beta[\Phi(\boldsymbol{\theta}_T)]$$

- The tighter upperbound can be obtained by minimizing Φ as we have:

$$\Phi(\boldsymbol{\theta}) \leq \inf_{\boldsymbol{\theta}_T: \mathbb{E}_\beta[\boldsymbol{\theta}_T] = \boldsymbol{\theta}} \mathbb{E}_\beta[\Phi(\boldsymbol{\theta}_T)] = \Phi_\beta(\boldsymbol{\theta})$$

- We can maximize the free-energy as we want (together with the Lagrange multiplier) and so, we have, the following optimization problem:

$$\sup_{\boldsymbol{\theta}_T} \inf_{\boldsymbol{\lambda}} \left(\mathbb{E}_\beta[\Phi(\boldsymbol{\theta}_T)] - \boldsymbol{\lambda}^T (\mathbb{E}_\beta[\boldsymbol{\theta}_T] - \boldsymbol{\theta}) \right)$$

Proposition 7.3. (Solving Optimization - Tighter Bound) Given the following optimization problem:

$$\sup_{\boldsymbol{\theta}_T} \inf_{\boldsymbol{\lambda}} \left(\mathbb{E}_\beta[\Phi(\boldsymbol{\theta}_T)] - \boldsymbol{\lambda}^T (\mathbb{E}_\beta[\boldsymbol{\theta}_T] - \boldsymbol{\theta}) \right)$$

This implies that $\boldsymbol{\mu}_T = \Pi_T(\boldsymbol{\lambda})$ as it is Lagrange multiplier corresponding to non-zero element of $\boldsymbol{\theta}$ for each tree (but the $\boldsymbol{\mu}$ stays the same (for all tree), see tree-reparametrisation view of BP). The Lagrange multiplier and $\Phi_\beta(\boldsymbol{\theta})$ can be find by:

$$\Phi_\beta(\boldsymbol{\theta}) = \sup_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \boldsymbol{\theta} + \sum_i H_\lambda(\boldsymbol{\theta}) - \sum_{ij} \beta_{ij} I_\lambda(\mathbf{x}_i, \mathbf{x}_j)$$

where β_{ij} is the probability distribution over tree.

Proof. Consider the derivative of Lagrange multiplier as we have:

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\theta}_T} \sum_T \beta(T) \Phi(\boldsymbol{\theta}_T) - \boldsymbol{\lambda}^T \frac{\partial}{\partial \boldsymbol{\theta}_T} \sum_T \beta(T) \boldsymbol{\theta}_T = 0 \\ \iff & \beta(T) \boldsymbol{\mu}_T - \beta(T) \Pi_T(\boldsymbol{\lambda}) = 0 \\ \iff & \boldsymbol{\mu}_T = \Pi_T(\boldsymbol{\lambda}) \end{aligned}$$

Thus the first point is proven. Now, consider the $\Phi_\beta(\boldsymbol{\beta})$ as we have:

$$\begin{aligned} \Phi_\beta(\boldsymbol{\theta}) &= \sup_{\boldsymbol{\lambda}} \inf_{\boldsymbol{\theta}_T} \left(\mathbb{E}_\beta[\Phi(\boldsymbol{\theta}_T)] - \boldsymbol{\lambda}^T (\mathbb{E}_\beta[\boldsymbol{\theta}_T] - \boldsymbol{\theta}) \right) \\ &= \sup_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \boldsymbol{\theta} + \mathbb{E}_\beta \left[\inf_{\boldsymbol{\theta}_T} \left(\Phi(\boldsymbol{\theta}_T) - \boldsymbol{\theta}_T^T \Pi_T(\boldsymbol{\lambda}) \right) \right] \\ &= \sup_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \boldsymbol{\theta} + \mathbb{E}_\beta \left[-\Psi(\Pi_T(\boldsymbol{\lambda})) \right] \\ &= \sup_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \boldsymbol{\theta} + \mathbb{E}_\beta \left[\sum_i H_\lambda(\mathbf{x}_i) - \sum_{(ij) \in T} I_\lambda(\mathbf{x}_i, \mathbf{x}_j) \right] \\ &= \sup_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \boldsymbol{\theta} + \sum_i H_\lambda(\mathbf{x}_i) - \sum_{(ij) \in T} \beta_{ij} I_\lambda(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

Thus the proposition is proven. \square

Remark 69. (Interpretation of Result) This is the convexification of the Bethe free energy. We can optimization with respect to $\boldsymbol{\lambda}$, which is approximate inference applied to tightest bound on $\Phi(\boldsymbol{\theta})$ for a fixed β . The bound can be tighten by the optimization of β .

Remark 70. (Bethe Free Energy and EP - Setup) A Bethe-like approach also casts EP as variational energy fixed point method. Consider the marginals of a posterior distribution by cliques potential:

$$P(\mathcal{Z}) \propto f_0(\mathcal{Z}) \prod_i f_i(\mathcal{Z})$$

There are some remarks on each of the components:

- All factors are exponential form
- f_0 is tractable exponential family (possible uniform)
- f_i are jointly intractable, but the product can't be marginalized, although individual terms maybe tractable.

We consider the tractable exponential family terms with zero natural parameter together with the tractable sufficient statistics $\tilde{\mathbf{s}}(\mathcal{Z}_i)$ as (can be seen as an approximating sites):

$$P(\mathcal{Z}) = \exp\left(\boldsymbol{\theta}_0^T \mathbf{s}_0(\mathcal{X})\right) \prod_i \exp\left(\boldsymbol{\theta}_i^T \mathbf{s}_i(\mathcal{Z}_i)\right) \exp\left(\mathbf{0}^T \tilde{\mathbf{s}}(\mathcal{Z}_i)\right)$$

Now, we can consider any given natural parameter $\tilde{\boldsymbol{\theta}}$ but it will be initialized with $\mathbf{0}$ (it can change):

$$P(\mathcal{Z}) = \exp\left(\boldsymbol{\theta}_0^T \mathbf{s}_0(\mathcal{X}) + \sum_i \boldsymbol{\theta}_i^T \mathbf{s}_i(\mathcal{Z}_i) + \tilde{\boldsymbol{\theta}}^T \tilde{\mathbf{s}}(\mathcal{Z}_i)\right)$$

The variational dual principle will tell us that the expected sufficient statistics:

$$\boldsymbol{\mu}_0^* = \langle \mathbf{s}_0 \rangle_P \quad \boldsymbol{\mu}_i^* = \langle \mathbf{s}_i(\mathcal{Z}) \rangle_P \quad \tilde{\boldsymbol{\mu}}_i^* = \langle \tilde{\mathbf{s}}_i \rangle_P$$

which are also the maximization of the likelihood, which is given by:

$$\{\boldsymbol{\mu}_0^*, \boldsymbol{\mu}_i^*, \tilde{\boldsymbol{\mu}}_i^*\} = \arg \max_{\{\boldsymbol{\mu}_0, \boldsymbol{\mu}_i, \tilde{\boldsymbol{\mu}}_i\} \in \mathcal{M}} \left[\boldsymbol{\theta}_0^T \boldsymbol{\mu}_0 + \sum_i (\boldsymbol{\theta}_i^T + \mathbf{0}^T \tilde{\boldsymbol{\mu}}_i) - \Phi(\boldsymbol{\mu}_0, \boldsymbol{\mu}_i, \tilde{\boldsymbol{\mu}}_i) \right]$$

Note that the negative entropy terms is evaluated over the full distribution i.e $\tilde{\boldsymbol{\theta}}$ doesn't have to be zero.

Remark 71. (EP Relaxiation) We consider the following relaxation to get the EP:

- We will have to relax the feasible sets \mathcal{M} , as we want a local consistency instead of global, where:
 - We want all edges connecting $\{\boldsymbol{\mu}_0, \tilde{\boldsymbol{\mu}}_i\}$, as $\boldsymbol{\mu}_0$ is the expected sufficient statistics over all \mathcal{Z} , thus we need it, when marginalized, to be consistent with all $\tilde{\boldsymbol{\mu}}_i$
 - We want all pairs $(\boldsymbol{\mu}_i, \tilde{\boldsymbol{\mu}}_i)$ to be consistent also. This is equivalent to projection step (in EP) as we are matching expected sufficient statistics.
- The message entropy should be given by:

$$\Phi_{\text{Bathe}}(\{\boldsymbol{\mu}_0, \tilde{\boldsymbol{\mu}}_i\}) - \sum_i (H[\boldsymbol{\mu}_i, \tilde{\boldsymbol{\mu}}_i] - H[\tilde{\boldsymbol{\mu}}_i])$$

We have all graph structure entropy (without $\boldsymbol{\mu}_i$) together with conditional entropy (which is joint entropy minus by individual entropy)

- By doing this, we have dropped all the intractable terms and use the tractable terms to do the calculation (and talk to each other).
- Free-energy based approximation marginal include μ_i and run reparameterization on a junction graph.
- Direct learning on EP free-energy would use these marginal rather than approximate one and a local normalizer from the integration over $f_i(\mathcal{Z}_i)q_{-i}(\mathcal{Z}_i)$ (see the last part of the EP section.)

8 Variational Method

8.1 Introduction

Remark 72. (Limitations) Our treatment of variational method has emphasis natural choices of variational family often factorized using some functional (exponential) form as a joint. They are mostly restricted to joint exponential family facilitates hieratical and distributional model but not non-linear and non-conjugate.

Remark 73. (Using Unconstraint Optimization) Consider parameteric variational approximation via a constrained family $q(\mathcal{Z}; \boldsymbol{\rho})$. The constrained variational E-step becomes

$$q(\mathcal{Z}) = \arg \max_{q \in \{q(\mathcal{Z}; \boldsymbol{\rho})\}} \mathcal{F}(q(\mathcal{Z}), \boldsymbol{\theta}^{(k-1)}) \implies \boldsymbol{\rho}^{(k)} = \arg \max_{\boldsymbol{\rho}} \mathcal{F}(q(\mathcal{Z}; \boldsymbol{\rho}), \boldsymbol{\theta}^{(k-1)})$$

Remark 74. (Reparameterized Free Energy) We can replace constrained optimization $\mathcal{F}(q, \boldsymbol{\theta})$ with unconstrained optimization of constrained $\mathcal{F}(\boldsymbol{\rho}, \boldsymbol{\theta})$:

$$\mathcal{F}(\boldsymbol{\rho}, \boldsymbol{\theta}) = \left\langle \log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}^{(k-1)}) \right\rangle_{q(\mathcal{Z}; \boldsymbol{\rho})} + H[\boldsymbol{\rho}]$$

We may use the coordinate ascent in $\boldsymbol{\rho}$ and $\boldsymbol{\theta}$ but this no longer necessary, as we have:

- In special case, the expectation of the log-joint under $q(\mathcal{Z}; \boldsymbol{\rho})$ can be expressed in closed form. We can follow $\nabla_{\boldsymbol{\rho}} \mathcal{F}$

- This requires evaluation a high-dimensional expectation with respected to $q(\mathcal{Z}; \boldsymbol{\rho})$ as a function of $\boldsymbol{\rho}$ that isn't simple.
- There are 3 solutions to this problem:
 - “Score Based” gradient estimate and Monte-Carlow.
 - Recognition network training in separate place - not strictly variational.
 - Recognition network training simultaneously with generative model using frozen sample.

Proposition 8.1. *One can show that:*

$$\nabla_{\boldsymbol{\rho}} \mathcal{F}(\boldsymbol{\rho}, \boldsymbol{\theta}) = \left\langle [\nabla_{\boldsymbol{\rho}} \log q(\mathcal{Z}; \boldsymbol{\rho})] \left(\log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}) - q(\mathcal{Z}; \boldsymbol{\rho}) \right) \right\rangle_{q(\mathcal{Z}; \boldsymbol{\rho})}$$

Proof. We consider the following gradient:

$$\begin{aligned} \nabla_{\boldsymbol{\rho}} \mathcal{F}(\boldsymbol{\rho}, \boldsymbol{\theta}) &= \nabla_{\boldsymbol{\rho}} \int q(\mathcal{Z}; \boldsymbol{\rho}) \left[\log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}) - \log q(\mathcal{Z}; \boldsymbol{\rho}) \right] d\mathcal{Z} \\ &= \int [\nabla_{\boldsymbol{\rho}} q(\mathcal{Z}; \boldsymbol{\rho})] \left(\log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}) - \log q(\mathcal{Z}; \boldsymbol{\rho}) \right) + q(\mathcal{Z}; \boldsymbol{\rho}) \nabla_{\boldsymbol{\rho}} [\log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}) - \log q(\mathcal{Z}; \boldsymbol{\rho})] d\mathcal{Z} \end{aligned}$$

We have the following facts:

$$\begin{aligned} \nabla_{\boldsymbol{\rho}} \log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}) &= 0 \\ \int q(\mathcal{Z}; \boldsymbol{\rho}) \nabla_{\boldsymbol{\rho}} q(\mathcal{Z}; \boldsymbol{\rho}) &= \nabla_{\boldsymbol{\rho}} \int q(\mathcal{Z}; \boldsymbol{\rho}) d\mathcal{Z} = 0 \\ \nabla_{\boldsymbol{\rho}} q(\mathcal{Z}; \boldsymbol{\rho}) &= q(\mathcal{Z}; \boldsymbol{\rho}) \nabla_{\boldsymbol{\rho}} \log q(\mathcal{Z}; \boldsymbol{\rho}) \end{aligned}$$

And so we have the solution as required □

Remark 75. (Reducing Variance) To reduce the expectation of gradient, due to the high-dimensional problem, we can evaluate by MC.

- We can reduce by factorization, where $q(\mathcal{Z}) = \prod_i q(\mathcal{Z}_i | \boldsymbol{\rho}_i)$ factor over disjoint cliques.
- Let $\bar{\mathcal{Z}}_i$ be the minimal markov blanket of \mathcal{Z}_i in the joint $P_{\bar{\mathcal{Z}}_i}$ be a product of joint factors that include element of \mathcal{Z}_i and $P_{-\bar{\mathcal{Z}}_i}$

We have the following gradient:

$$\begin{aligned} \nabla_{\boldsymbol{\rho}_i} \mathcal{F}(\{\boldsymbol{\rho}_j\}, \boldsymbol{\theta}) &= \left\langle \left[\nabla_{\boldsymbol{\rho}_i} \sum_j \log q(\mathcal{Z}_j; \boldsymbol{\rho}_j) \right] \left(\log P(\mathcal{X}, \mathcal{Z} | \boldsymbol{\theta}) - \sum_j \log q(\mathcal{Z}_j; \boldsymbol{\rho}_j) \right) \right\rangle_{q(\mathcal{Z})} \\ &= \left\langle [\nabla_{\boldsymbol{\rho}_i} \log q(\mathcal{Z}_i; \boldsymbol{\rho}_i)] (\log P_{\bar{\mathcal{Z}}_i}(\mathcal{X}, \bar{\mathcal{Z}}_i) - \log q(\mathcal{Z}_i; \boldsymbol{\rho}_i)) \right\rangle_{q(\bar{\mathcal{Z}}_i)} \\ &\quad + \left\langle [\nabla_{\boldsymbol{\rho}_i} \log q(\mathcal{Z}_i; \boldsymbol{\rho}_i)] \left(\log P_{-\bar{\mathcal{Z}}_i}(\mathcal{X}, \bar{\mathcal{Z}}_{-i}) - \sum_{j \neq i} \log q(\mathcal{Z}_j; \boldsymbol{\rho}_j) \right) \right\rangle_{q(\mathcal{Z})} \end{aligned}$$

Please note that the second term is propositional to $\langle \nabla_{\boldsymbol{\rho}_i} \log q(\mathcal{Z}_i; \boldsymbol{\rho}_i) \rangle_{q(\mathcal{Z}_i)} = 0$ so we only need to consider the expectation with respected to $q(\bar{\mathcal{Z}}_i)$ which is variational message passing.

Remark 76. (Sampling Methods) We consider the following “black-box” variational approach is as follows:

- Choose a parameteric (factored) variational family $q(\mathcal{Z}) = \prod_i q(\mathcal{Z}_i; \boldsymbol{\theta})$
- Initiate the factors.

- Repeat until convergence:
 - Stochastic VE-step: Sample from $q(\tilde{\mathcal{Z}}_i)$ and estimate expected gradient $\nabla \rho_i \mathcal{F}$, and we update ρ_i along the gradient.
 - Stochastic M-step: Sample from each $q(\tilde{\mathcal{Z}}_i)$ as we can update the corresponding parameter.
- Stochastic may use Robbins Monro to promote convergence.
- Variance of the gradient estimate can also be controlled by MC techniques.

Remark 77. (Batches of Data) We have not distinguish between multi-variate models and iid data instances. As we group, all together in \mathcal{Z} ; for example, large model such as HMM, we often make with multiple data draws and each instance requires a separate variational optimization.

Definition 8.1. (Recognition Model) Suppose we have fixed length vector $\{(\mathbf{x}_i, \mathbf{z}_i)\}$ when \mathbf{z}_i is latent:

- Optimal variational distribution $q^*(\mathbf{z}_i)$ will depends on \mathbf{x}_i
- We want to learn the mapping from $q(\mathbf{z}_i; \boldsymbol{\rho} = f(\mathbf{x}_i; \boldsymbol{\phi}))$
- Now $\boldsymbol{\rho}$ is output of a general function approximate f parameterized by $\boldsymbol{\phi}$, training on map \mathbf{x}_i to the variational parameter of $q(\mathbf{z}_i)$

The mapping function f is called recognition model.

Definition 8.2. (Helmholtz Model) It is a binary sigmoid belief-network with parallel recognition model. There are 2 phase of learning:

- *Wake Phase:* Given current f estimate mean field representation from data:

$$q(\mathbf{z}_i) = \text{Bern}(\hat{\mathbf{z}}_i) \quad \hat{\mathbf{z}}_i = f(\mathbf{x}_i; \boldsymbol{\phi})$$

We will update generative parameter $\boldsymbol{\theta}$ according to $\nabla_{\boldsymbol{\theta}} \mathcal{F}(\{\hat{\mathbf{z}}_i\}; \boldsymbol{\theta})$

- *Sleep Phase:* Sample $\{\mathbf{z}_i, \mathbf{x}_i\}_{i=1}^S$ from the current generative model. Update recognition parameter $f(\mathbf{x}_i)$ toward \mathbf{z}_i :

$$\Delta \boldsymbol{\phi} \propto \sum_{i=1}^S (\mathbf{z}_i - f(\mathbf{x}_i; \boldsymbol{\phi})) \nabla_{\boldsymbol{\phi}} f(\mathbf{x}_i; \boldsymbol{\phi})$$

Please note that this step minimizes:

$$\text{KL} \left[P_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) \parallel q(\mathbf{z}; f(\mathbf{x}; \boldsymbol{\phi})) \right]$$

This is opposite to variational objective. Opposite to variational objective, but may not matter if divergence is small enough

Remark 78. (Comments on Helmholtz Model Evaluation) We have to sample \mathbf{z} from recognition model rather than just evaluate means:

- Expectation in free-energy can be computed directly rather than by means substitution.
- In hieratical model, output of higher recognition layer depends on sample at previous stages, which introduces correlation between sample at difference layer.

Recognition model structure need not necessary exactly echo generative model. Please note that a more general approach is to train f to yields the parameter of exponential family $q(\mathbf{z})$.

Definition 8.3. (Variational Autoencoder) We fuse the wake and sleep phase. We generate the recognition sample used deterministic transformation of external random variable. If f gives marginal μ_i and σ_i for latent z_i and $\epsilon_i^s \sim \mathcal{N}(0, 1)$ then:

$$z_i^s = \mu_i + \sigma_i \epsilon_i^s$$

Given the generative and recognition model can be trained together with gradient descent. Holding ϵ^s fixed:

$$\mathcal{F}_i(\theta, \phi) = \sum_s \log P(\mathbf{x}_i, z_i^s; \theta) - \log q(z_i^s; f(\mathbf{x}_i, \phi))$$

We have the following derivative as:

$$\begin{aligned} \frac{\partial}{\partial \theta} \mathcal{F}_i &= \sum_s \nabla_{\theta} \log P(\mathbf{x}_i, z_i^s; \theta) \\ \frac{\partial}{\partial \phi} \mathcal{F}_i &= \sum_s \frac{\partial}{\partial z_i^s} \left(\log P(\mathbf{x}_i, z_i^s; \theta) - \log q(z_i^s; f(\mathbf{x}_i)) \right) \frac{dz_i^s}{d\phi} + \frac{\partial}{\partial f(\mathbf{x}_i)} \log q(z_i^s; f(\mathbf{x}_i)) \frac{df(\mathbf{x}_i)}{d\phi} \end{aligned}$$

Remark 79. (Observations on Variational Auto-Encoder) We consider the following observations on the training of VAE:

- We start by frozen sample ϵ^s can be redrawn that avoid overfitting.
- It maybe possible to evaluate the entropy and $\log P(z)$ without sampling and reducing variance:
- Differentiable reparameterization are available for number of difference distribution
- The conditional $P(\mathbf{x}|\mathbf{z}, \theta)$ is often implemented as a neural network with additive noise at input.
- In practice, hieratical model appear to be difficult to train.

8.2 Additional Models to VAE

Definition 8.4. (Importance Weigh Free Energy) Consider another interpretation of free energy:

$$\mathcal{F}(q; \theta) = \left\langle \log \frac{P(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right\rangle_q$$

We consider the jensen's inequality on importance sampled estimate:

$$l(\theta) = \log \mathbb{E}_{z \sim q} \left[\frac{P(\mathbf{x}, z)}{q(z)} \right] \leq \mathbb{E}_{z \sim q} \left[\log \frac{P(\mathbf{x}, z)}{q(z)} \right] = \mathcal{F}(q; \theta)$$

Suggest more accurate importance weight as we have:

$$l(\theta) = \log \mathbb{E}_{z_1, \dots, z_k \sim q} \left[\frac{1}{k} \sum_k \frac{P(\mathbf{x}_k, z_k)}{q(z_k)} \right] \geq \mathbb{E}_{z_1, \dots, z_k \sim q} \left[\log \frac{1}{k} \sum_k \frac{P(\mathbf{x}_k, z_k)}{q(z_k)} \right]$$

This allows the tight-bound and reparameterization friendly but as $K \rightarrow \infty$, the signal for learning amortized q grows weaker making VAE learning too slow.

Definition 8.5. (Normalizing Flow) We have the following free energy:

$$\mathcal{F}(q, \theta) = \langle \log P(\mathbf{x}, \mathbf{z}|\theta) \rangle_q - \langle \log q(\mathbf{z}) \rangle_q$$

To evaluate \mathcal{F} , we need to be able to find expectation with respected to q and evaluate the log-density, which usually restrict us to tractable inference families. We consider the followign recognition model $q(\mathbf{z})$ implicitly by:

$$z_0 \sim q(\cdot; \mathbf{x}) \quad \mathbf{z} = f_k(f_{k-1}(\dots f_1(z_0)))$$

where q_0 should be fixed and tractable. And so, we have the following evaluations:

$$\langle F(\mathbf{z}) \rangle_q = \langle F(f_k(f_{k-1}(\dots f_1(\mathbf{z}_0)))) \rangle_{q_0} \quad \log q(\mathbf{z}) = \log q_0(f_1^{-1}(\dots f_{k-1}^{-1}(f_k^{-1}(\mathbf{z})))) - \sum_k \log |\nabla f_k|$$

where $|\nabla f_k|$ being the determinant of the jacobian, where we use the following transformation of variables:

$$\mathbf{z}_k = f_k(\mathbf{z}_{k-1}) \implies q(\mathbf{z}_k) = q(f_k^{-1}(\mathbf{z}_k)) \left| \frac{\partial \mathbf{z}_{k-1}}{\partial \mathbf{z}_k} \right| = q(f_k^{-1}(\mathbf{z}_k)) |\nabla f_k(\mathbf{z}_{k-1})|^{-1}$$

Given a sample $\mathbf{z}_0^i \sim q_0(\cdot; \mathbf{x})$ as we have:

$$\mathcal{F}(q, \boldsymbol{\theta}) \approx \frac{1}{S} \sum_s \log p(\mathbf{x}, f_k(f_{k-1}(\dots f_1(\mathbf{z}_0^i)))) + h[q_0] + \frac{1}{S} \sum_s \sum_k |\nabla f_k(f_{k-1}(\dots f_1(\mathbf{z}_0^i)))|$$

Remark 80. (Special f for Normalizing Flow) Suppose we use the special f and we have:

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^T \mathbf{z} + b) \implies |\nabla f| = |1 + \mathbf{u}^T \boldsymbol{\Psi}(\mathbf{z})| \quad \text{where} \quad \boldsymbol{\Psi}(\mathbf{z}) = h'(\mathbf{x}^T \mathbf{z} + b) \mathbf{w}$$

$$f(\mathbf{z}) = \mathbf{z} + \frac{\beta}{\alpha + |\mathbf{z} - \mathbf{z}_0|} \implies |\nabla f| = [1 + \beta h]^{d-1} [1 + \beta h + \beta h' r] \quad \text{where} \quad r = |\mathbf{z} - \mathbf{z}_0| \quad h = \frac{1}{\alpha + r}$$

Definition 8.6. (DDC Helmholtz Machine) We define q to be unnormalized exponential family with large set of sufficient statistics:

$$q(\mathbf{z}) \propto \exp(\sim_i \eta_i \psi_i(\mathbf{z}))$$

and it is parameterized by mean parameter $\boldsymbol{\mu} = \langle \boldsymbol{\psi}(\mathbf{z}) \rangle$, which we call distributed distribution code (DDC). Train recognition model using a sleep sample as we have:

$$\boldsymbol{\mu} = \langle \boldsymbol{\psi}(\mathbf{z}) \rangle_q = f(\mathbf{x}^*; \boldsymbol{\phi})$$

$$\Delta \boldsymbol{\phi} \propto \sum_s (\boldsymbol{\psi}(\mathbf{z}_s) - f(\mathbf{x}_s; \boldsymbol{\phi})) \nabla_{\boldsymbol{\phi}} f(\mathbf{x}_s; \boldsymbol{\phi})$$

Furthermore, we also learn linear approximation $\nabla \log P(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) \approx \mathbf{A}\boldsymbol{\psi}(\mathbf{z})$, where

$$\mathbf{A} = \left(\sum_s \nabla \log P(\mathbf{x}_s, \mathbf{z}_s|\boldsymbol{\theta}) \boldsymbol{\psi}(\mathbf{z}_s) \right)^T \left(\sum_s \boldsymbol{\psi}(\mathbf{z}_s) \boldsymbol{\psi}(\mathbf{z}_s)^T \right)^{-1}$$

Then we have $\langle \nabla \log P(\mathbf{x}, \mathbf{z}) \rangle_q \approx \mathbf{A} \langle \boldsymbol{\psi}(\mathbf{z}) \rangle \propto f(\mathbf{x}, \boldsymbol{\phi})$ can be generalized into infinite dimension with kernel.

Definition 8.7. (Amortised Learning) We aren't interested in inference. We can short-circled general recognition and compute expectation for learning directly, as we have:

$$\nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{F}(q^*, \boldsymbol{\theta}) = \langle \nabla_{\boldsymbol{\theta}} \log P(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) \rangle_{q^*}$$

We can use the wake-sleep approach:

- Sample $\{\mathbf{x}_s, \mathbf{z}_s\} \sim P(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}^k)$
- Train Regression $\hat{J}_{\boldsymbol{\theta}^k} : \mathbf{x}_k \mapsto \nabla_{\boldsymbol{\theta}} \log P(\mathbf{x}_s, \mathbf{z}_s)|_{\boldsymbol{\theta}^k}$ (Learning the mapping)
- Set $\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \eta \sum_i \hat{J}_{\boldsymbol{\theta}^k}(\mathbf{x}_i)$

Derivate form works for (kernel and GP) regression for which regressor is linear in target. For conditional, exponential family model:

$$\langle \log P(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) \rangle_{q^*} = \langle \boldsymbol{\eta}(\mathbf{z}, \boldsymbol{\theta}) \rangle_{q^*}^T \mathbf{T}(\mathbf{x}) - \langle \boldsymbol{\Phi}(\mathbf{z}, \boldsymbol{\theta}) + \log P(\mathbf{z}|\boldsymbol{\theta}) \rangle_{q^*}$$

and regressor can be trained to function of \mathbf{z} alone, with $T(\mathbf{x})$ evaluated on (wake-phase) data.

Remark 81. (VAE Comments) Much of the VAE and related work has common generative model as:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \mathbf{x} \sim \mathcal{N}(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}), \psi \mathbf{I})$$

where \mathbf{g} is neural network. Let's consider the dimension of \mathbf{z} as we have:

- Overcomplicated: If $\dim(\mathbf{z})$ is large enough, the optimal solution has $\psi \rightarrow 0$ as $\mathbf{q}(\mathbf{z}; \mathbf{x}) \rightarrow \delta(\mathbf{z} - \mathbf{f}(\mathbf{x}; \boldsymbol{\phi}))$. In effect, the generative model learns a flow to transform a model density to the target.
- Oversimplified: If $\dim(\mathbf{z})$ is small as non-linear PCA.

Interesting latent representation are required more structured generative model.

Definition 8.8. (Structured VAE) Consider a model where $P(\mathcal{Z}|\boldsymbol{\theta})$ has tractable joint exponential family potential and:

$$P(\mathcal{X}|\mathcal{Z}, \Gamma) = \prod_i P(\mathbf{x}_i|\mathbf{z}_i, \gamma_i)$$

are interactable. conditional independent observation γ_i might be the same for all i . Consider factored variational inference $q(\mathcal{Z}) = \prod_i q(\mathbf{z}_i)$ with no further constraints:

$$\begin{aligned} \log q_i^*(\mathbf{z}_i) &= \langle \log P(\mathcal{Z}, \mathcal{X}) \rangle_{q_{-i}} + \text{const.} \\ &= \langle \log P(\mathbf{z}_i|\mathcal{Z}_{-i}) + \log P(\mathbf{x}_i|\mathbf{z}_i) \rangle_{q_{-i}} + \text{const.} \\ &= \langle \boldsymbol{\eta}_{-i} \rangle_{q_{-i}}^T \boldsymbol{\psi}_i(\mathbf{z}_i) + \log P(\mathbf{x}_i|\mathbf{z}_i) \end{aligned}$$

Let's consider each variables (exploited the exponential family form of $P(\mathcal{Z})$):

- $\boldsymbol{\psi}_i$ are effective sufficient statistics included log-normalizer of children of DAG.
- $\boldsymbol{\eta}_{-i}$ is function of \mathcal{Z}_{-i}

We will choose the parametric form of $q_i(\mathbf{z}_i) = \exp(\tilde{\boldsymbol{\eta}}^T \boldsymbol{\psi}_i(\mathbf{z}_i) - \Phi_i(\tilde{\boldsymbol{\eta}}_i))$ and so the optimum will have:

$$\log q_i^*(\mathbf{z}_i) = \langle \boldsymbol{\eta}_{-i} \rangle_{q_{-i}}^T \boldsymbol{\psi}_i(\mathbf{z}_i) + \boldsymbol{\rho}(\mathbf{x}_i)^T \boldsymbol{\psi}_i(\mathbf{z}_i)$$

where $\boldsymbol{\rho}(\mathbf{x}_i) = \mathbf{f}_i(\mathbf{x}_i; \mathbf{p}hi)$ recognition function and it might be the same for all i .

Remark 82. (Training of Structured VAE) We consider the free energy:

$$\begin{aligned} \mathcal{F}(\boldsymbol{\theta}, \Gamma, \{\boldsymbol{\phi}_i\}) &= \left\langle \sum_i \log P(\mathbf{x}_i|\mathbf{z}_i, \gamma_i) + \log P(\mathcal{Z}|\boldsymbol{\theta}) \right\rangle_{q(\mathcal{Z}; \boldsymbol{\theta}, \{\boldsymbol{\phi}_i\})} + \sum_i H[q_i] \\ &= \sum_i \underbrace{\langle \log P(\mathbf{x}_i|\mathbf{z}_i, \gamma_i) + \log P(\mathcal{Z}|\boldsymbol{\theta}) \rangle_{q_i(\mathbf{z}_i; \boldsymbol{\theta}, \boldsymbol{\phi}_i)} + H[q_i]}_{\mathcal{F}_i} + \langle \log P(\mathcal{Z}|\boldsymbol{\theta}) \rangle_{q(\mathcal{Z}, \boldsymbol{\theta}, \{\boldsymbol{\phi}_i\})} \end{aligned}$$

Update $\boldsymbol{\theta}$ are just tractable model. To update each $\boldsymbol{\phi}_i$ and γ_i , find $\langle \boldsymbol{\eta}_{-i} \rangle_{q_{-i}}$ to give the prior-like in VAE, and we generated the reparameterization sample $\mathbf{z}_i^s \sim q_i$, then:

$$\begin{aligned} \frac{\partial}{\partial \gamma_i} \mathcal{F}_i &= \sum_s \nabla_{\gamma_i} \log P(\mathbf{x}_i, \mathbf{z}_i^s; \gamma_i) \\ \frac{\partial}{\partial \boldsymbol{\phi}_i} \mathcal{F}_i &= \sum_s \frac{\partial}{\partial \mathbf{z}_i^s} \left(\log P(\mathbf{x}_i, \mathbf{z}_i^s; \gamma_i) - \log q(\mathbf{z}_i^s; \mathbf{f}(\mathbf{x}_i)) \right) \frac{d\mathbf{z}_i^s}{d\boldsymbol{\phi}_i} + \frac{\partial}{\partial \mathbf{f}(\mathbf{x}_i)} \log q(\mathbf{z}_i^s; \mathbf{f}(\mathbf{x}_i)) \frac{d\mathbf{f}(\mathbf{x}_i)}{d\boldsymbol{\phi}_i} \end{aligned}$$

This is like standard VAE.