

Машинное обучение и нейросетевые модели

Лекция 7. Точный инференс дискретных переменных

Лектор: Кравченя Павел Дмитриевич

Волгоград 2025

План лекции

1. Подходы к точному инференсу дискретных латентных переменных.
 2. Процедура устранения переменной.
 3. Графы факторов и байесовские сети.
 4. Сообщения между факторами и переменными.
 5. Передача сообщений в общем виде.
 6. Алгоритм суммирования-произведения. Пример.
 7. Алгоритм Junction Tree.
 8. Случай plate-графов. Алгоритм устранения тензорной переменной. Пример.
 9. Особенности реализации алгоритма в Pyro.
 10. Механика работы алгоритма устранения тензорной переменной в Pyro.
-

- В случае дискретных латентных переменных интеграл в формуле апостериорного распределения превращается в сумму, которую можно точно посчитать.
- Однако, если количество дискретных переменных и число значений, которые они могут принимать, большое, то вычисление этой суммы является вычислительно сложным и требующим огромного объёма памяти.
- Ситуация может осложняться наличием в модели как дискретных, так и непрерывных переменных.
- Часто вместо непосредственного суммирования функции по всем значениям переменных используют другой подход.

- Рассмотрим в качестве примера следующую вероятностную модель:

$$p(a, b, c, d) = p(a | b)p(b | c)p(c | d)p(d).$$



- Рассчитаем вероятность $p(a = 0)$. Согласно правилу суммирования:

$$p(a = 0) = \sum_{b,c,d} p(a = 0, b, c, d) = \sum_{b,c,d} p(a = 0 | b)p(b | c)p(c | d)p(d).$$

- Можно вычислять эти суммы, последовательно складывая все слагаемые с конкретными значениями b, c, d . Однако, более эффективный способ – вынести слагаемые за знаки сумм там, где это возможно:

$$p(a = 0) = \sum_b p(a = 0 | b) \sum_c p(b | c) \underbrace{\sum_d p(c | d)p(d)}_{\gamma_d(c)}.$$

- Переменная $\gamma_d(c)$ называется потенциалом. Она содержит всю информацию о переменной d .
- Выполняя аналогичные операции по другим переменным, получаем:

$$p(a = 0) = \sum_b p(a = 0 | b) \underbrace{\sum_c p(b | c) \gamma_d(c)}_{\gamma_c(b)} = \sum_b p(a = 0 | b) \gamma_c(b)$$

- Проведённая процедура называется **устранением переменной** (*variable elimination*), т.к. каждый раз, когда суммируются состояния переменной, она исключается из распределения.
- Устранение переменной можно рассматривать как передачу сообщения (информации) соседнему узлу графа.

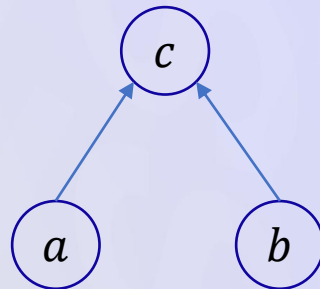
- Часто вместо *байесовской сети* вероятностная модель представляется в виде графа факторов, что порой может оказаться удобнее для расчетов.
- Любое совместное распределение $p(\mathbf{X})$, $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$, может быть представлено как произведение факторов f_i , $f_i(x_i) > 0$:

$$p(\mathbf{X}) = \prod_i f_i(x_i).$$

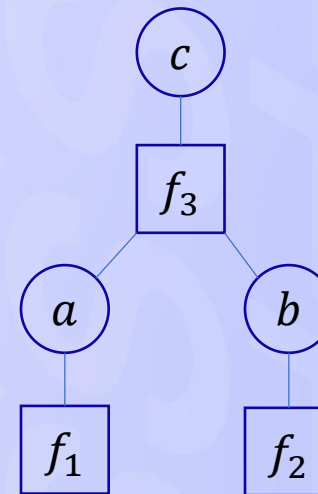
- При этом один фактор (некоторая функция переменных x_i) ставится в соответствие одной клике (*clique*): подмножеству переменных $\{x_1, x_2, \dots, x_n\}$, в котором каждая переменная соединена со всеми другими.
- В *графе факторов* вершины, соответствующие факторам, обозначаются квадратом, а вершины, соответствующие переменным, – кружком.

- Граф факторов является двудольным графом.
- Пример: вероятностная модель, представленная байесовской сетью и соответствующим ей графом факторов.

$$p(a, b, c) = p(a)p(b)p(c | a, b).$$



$$p(a, b, c) = f_1(a)f_2(b)f_3(a, b, c).$$



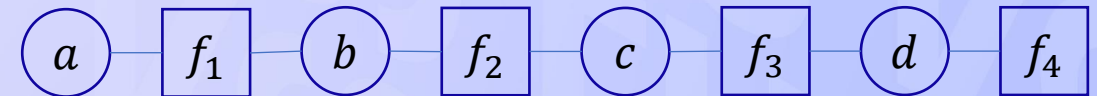
- Рассмотрим графическую модель, не содержащую ветвлений. Например:

$$p(a, b, c, d) = p(a | b)p(b | c)p(c | d)p(d).$$



- Представим её в виде *графа факторов*:

$$p(a, b, c, d) = f_1(a, b)f_2(b, c)f_3(c, d)f_4(d).$$



- Вычислим маргинальное распределение:

$$p(a, b, c) = \sum_d f_1(a, b)f_2(b, c)f_3(c, d)f_4(d) = f_1(a, b)f_2(b, c) \underbrace{\sum_d f_3(c, d)f_4(d)}_{\mu_{d \rightarrow c}(c)}.$$

- В данном выражении $\mu_{d \rightarrow c}(c)$ определяет сообщение от узла d к узлу c .

- Аналогично:
$$p(a, b) = f_1(a, b) \underbrace{\sum_c f_2(b, c) \mu_{d \rightarrow c}(c)}_{\mu_{c \rightarrow b}(b)}.$$

Сообщения между факторами и переменными

- Рассмотрим более сложный случай:

$$p(a, b, c, d, e) = p(a | b)p(b | c, d)p(c)p(e | d)p(d).$$

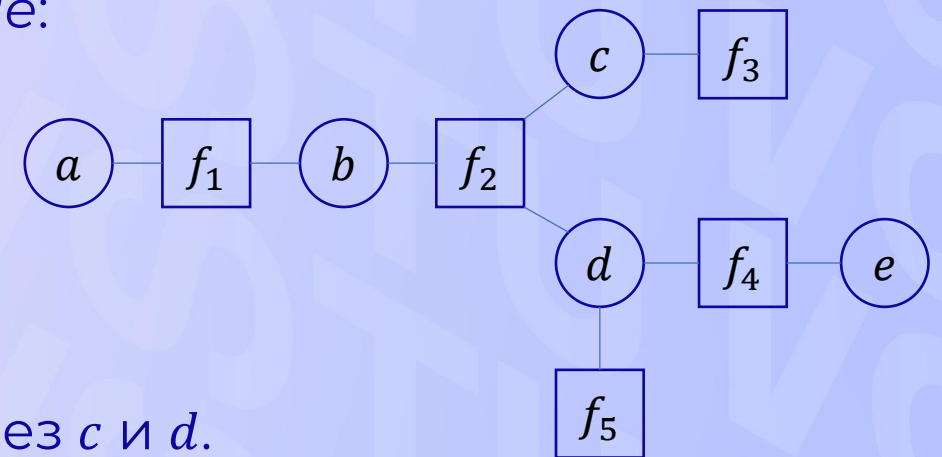
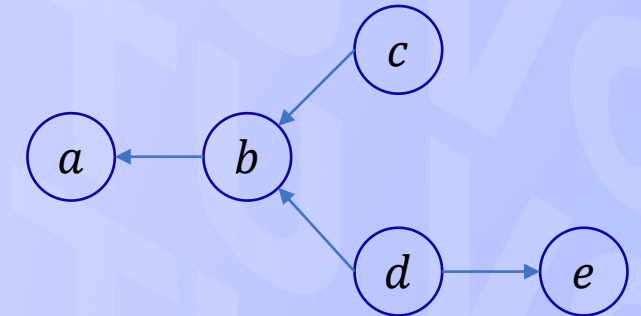
- Или в нотации факторов:

$$p(a, b, c, d, e) = f_1(a, b)f_2(b, c, d)f_3(c)f_4(d, e)f_5(d).$$

- Вычислим маргинальное распределение:

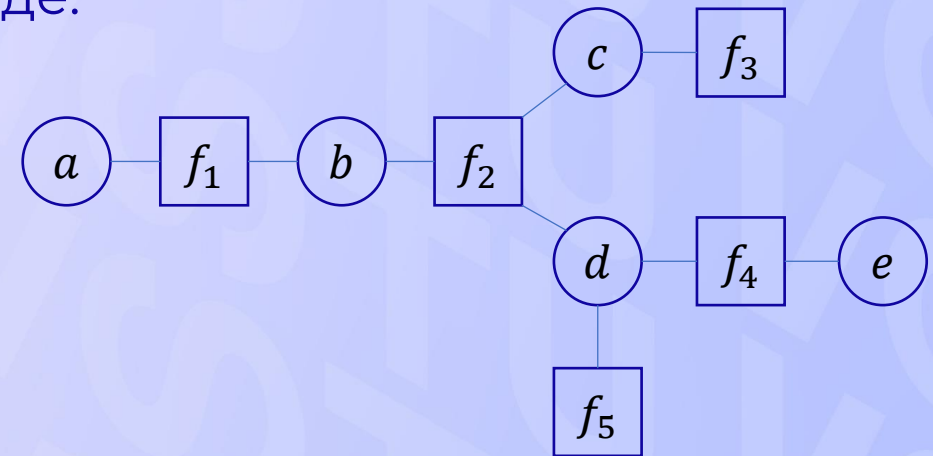
$$p(a, b) = f_1(a, b) \underbrace{\sum_{c, d} f_2(b, c, d)f_3(c)f_5(d) \sum_e f_4(d, e)}_{\mu_{f_2 \rightarrow b}(b)}.$$

- Здесь $\mu_{f_2 \rightarrow b}(b)$ обозначает сообщение от фактора к переменной. Оно состоит из сообщений, полученных из двух ветвей через c и d .



- Выражение $\mu_{f_2 \rightarrow b}(b)$ можно записать в виде:

$$\mu_{f_2 \rightarrow b}(b) = \sum_{c,d} f_2(b, c, d) \underbrace{\underbrace{f_3(c)}_{\mu_{c \rightarrow f_2}(c)} \underbrace{f_5(d) \sum_e f_4(d, e)}_{\mu_{d \rightarrow f_2}(d)}}.$$



- Аналогично:

$$\mu_{d \rightarrow f_2}(d) = \underbrace{f_5(d)}_{\mu_{f_5 \rightarrow d}(d)} \underbrace{\sum_e f_4(d, e)}_{\mu_{f_4 \rightarrow d}(d)}.$$

- Таким же способом можно посчитать маргинальное распределение $p(a)$:

$$p(a) = \underbrace{\sum_b f_1(a, b) \mu_{f_2 \rightarrow b}(b)}_{\mu_{f_1 \rightarrow a}(a)},$$

$$\mu_{f_1 \rightarrow a}(a) = \sum_b f_1(a, b) \underbrace{\mu_{f_2 \rightarrow b}(b)}_{\mu_{b \rightarrow f_1}(b)}.$$

- Можно увидеть, что сообщение от фактора к переменной формируется с помощью суммирования произведений входящих сообщений от переменной к фактору.
- Аналогично, сообщение от переменной к фактору определяется произведением входящих сообщений от фактора к переменной.
- Удобство такого подхода заключается в возможности повторного использования рассчитанных сообщений для оценки других маргинальных распределений.
- Например, маргинальное распределение $p(b)$ может быть выражено так:

$$p(b) = \underbrace{\sum_a f_1(a, b)}_{\mu_{f_1 \rightarrow b}(b)} \mu_{f_2 \rightarrow b}(b).$$

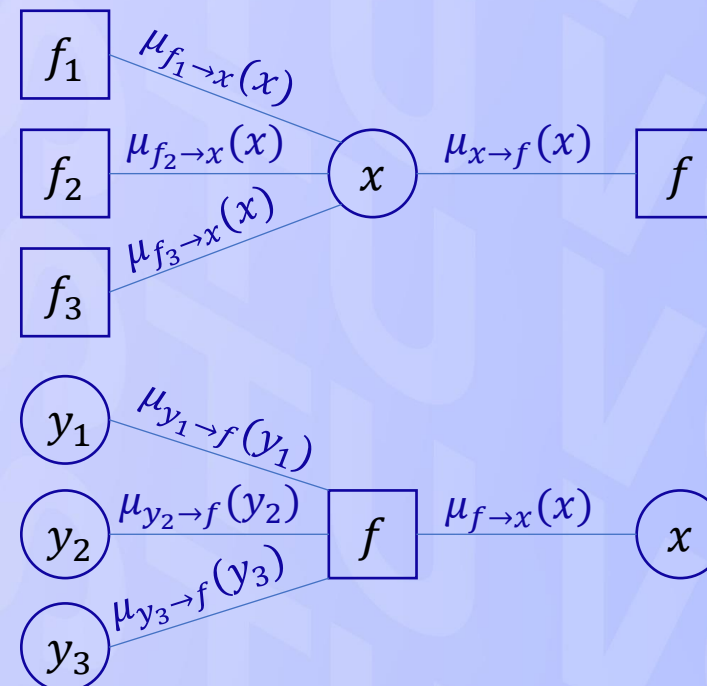
- Таким образом, передачу сообщений в общем виде можно записать так:
 - Сообщение от переменной к фактору:

$$\mu_{x \rightarrow f}(x) = \prod_{g \in \text{neighbors}(x) \setminus f} \mu_{g \rightarrow x}(x),$$

- Сообщение от фактора к переменной:

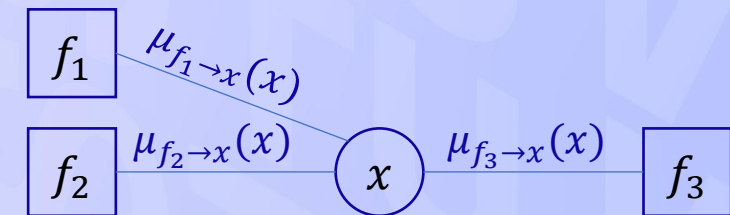
$$\mu_{f \rightarrow x}(x) = \sum_{\mathcal{X}_f \setminus x} f(\mathcal{X}_f) \prod_{y \in \text{neighbors}(f) \setminus x} \mu_{y \rightarrow f}(y),$$

Суммирование выполняется по всем значениям переменных $\mathcal{X}_f \setminus x$.



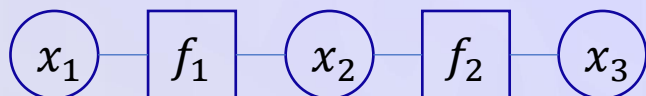
- Сообщения от факторов в листьях графа инициализируются фактором.
- Сообщения от переменных в листьях графа устанавливаются в единицу.
- Маргинализация переменной в этом случае может быть рассчитана так:

$$p(x) = \frac{1}{Z} \prod_{f \in \text{neighbors}(x)} \mu_{f \rightarrow x}(x).$$



- Данная процедура называется алгоритмом суммирования-произведения
- Сообщения в алгоритме суммирования-произведения рассчитываются в определённом порядке. Он часто называется планировкой сообщений (*message schedule*). Сообщения, от которых зависят другие сообщения, рассчитываются перед ними.

- Рассмотрим следующий факторный граф с бинарными переменными x_1 , x_2 , x_3 и факторами $f_1(x_1, x_2)$ и $f_2(x_2, x_3)$, заданными, например, таблично:



$$p(x_1, x_2, x_3) = \frac{1}{Z} \cdot f_1(x_1, x_2) \cdot f_2(x_2, x_3).$$

- Требуется вычислить маргинальное распределение $p(x_2)$.
- Инициализация: $\mu_{x_1 \rightarrow f_1}(x_1) = 1$, $\mu_{x_3 \rightarrow f_2}(x_3) = 1$.
- Вычисление сообщений от факторов f_1 и f_2 к переменной x_2 :

$$\mu_{f_1 \rightarrow x_2}(x_2) = \sum_{x_1} f_1(x_1, x_2) \mu_{x_1 \rightarrow f_1}(x_1) = \sum_{x_1} f_1(x_1, x_2) \cdot 1 = \sum_{x_1} f_1(x_1, x_2);$$

$$\mu_{f_2 \rightarrow x_2}(x_2) = \sum_{x_3} f_2(x_2, x_3) \mu_{x_3 \rightarrow f_2}(x_3) = \sum_{x_3} f_2(x_2, x_3) \cdot 1 = \sum_{x_3} f_2(x_2, x_3).$$

- Маргинальное распределение $p(x_2)$ пропорционально произведению всех входящих сообщений в x_2 :

$$p(x_2) \propto \mu_{f_1 \rightarrow x_2}(x_2) \cdot \mu_{f_2 \rightarrow x_2}(x_2).$$

- С учетом полученных сообщений:

$$p(x_2) \propto \left(\sum_{x_1} f_1(x_1, x_2) \right) \cdot \left(\sum_{x_3} f_2(x_2, x_3) \right).$$

- Для получения точного распределения $p(x_2)$ потребуется нормализовать результат по всем значениям переменной x_2 :

$$p(x_2) = \frac{(\sum_{x_1} f_1(x_1, x_2)) \cdot (\sum_{x_3} f_2(x_2, x_3))}{\sum_{x_2} (\sum_{x_1} f_1(x_1, x_2)) \cdot (\sum_{x_3} f_2(x_2, x_3))}.$$

- Факторные графы представляют собой совместное распределение вероятностей в виде произведения факторов.
- Если факторный граф достаточно сложный, то прямое вычисление маргинальных или условных вероятностей может быть вычислительно сложным из-за большого числа переменных и их зависимостей.
- Эту проблему решает **Junction tree**, преобразуя граф в древовидную структуру, где вычисления становятся *локальными и эффективными*.
- Junction tree – это дерево, построенное на основе исходного графа факторов, и содержащее:
 - ✓ вершины (кластеры переменных, которые связаны между собой в исходном графе);
 - ✓ рёбра (сепараторы, содержащие переменные, общие для обоих кластеров).

1. Построение морального графа (moral graph).

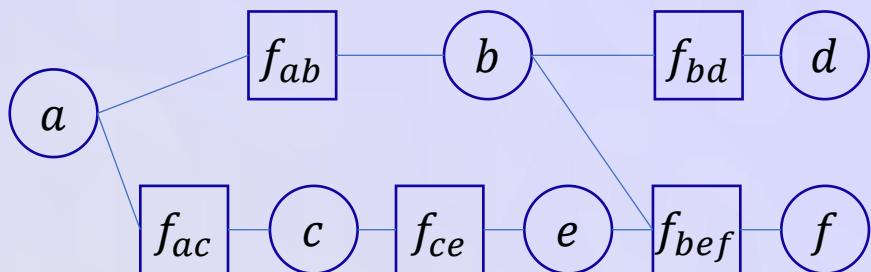
- ✓ Морализованная копия направленного ациклического графа образуется добавлением рёбер между всеми парами узлов, которые имеют общих детей, а затем преобразования всех рёбер в графе в неориентированные.
 - Для каждого фактора нужно соединить все переменные, которые входят в этот фактор. При этом образуется клика (полный подграф).
 - А затем нужно удалить все факторы, оставив только переменные и ребра между ними.

2. Выполнение триангуляции графа.

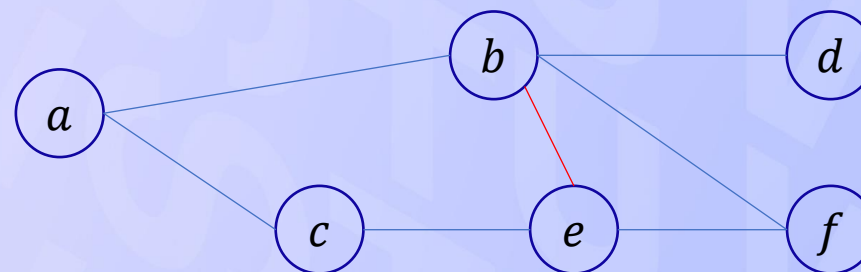
- ✓ Граф называется хордальным, если каждый из его длинных циклов (имеющих четыре ребра и более) имеет хорду. Для построения **junction tree** граф должен быть хордальным. Процесс добавления ребер для устранения длинных циклов называется триангуляцией.
- ✓ Для выбранного порядка устранения переменных (variables elimination):
 - для каждой устраняемой переменной нужно соединить ребром всех её соседей в текущем графе, чтобы избежать длинных циклов. Добавляемые рёбра остаются в графе.
 - Этот процесс повторяется для всех узлов в графе.

3. Определение максимальных клик (clique) в графе.
 - ✓ Максимальная клика — это полный подграф, который не может быть расширен добавлением новых вершин без потери свойства полноты.
4. Объединение максимальных клик в граф кластеров.
 - ✓ Нужно построить взвешенный граф из определённых максимальных клик, в котором:
 - каждая максимальная клика (кластер) соответствует одной вершине;
 - каждое ребро между вершинами графа (кликами) имеет вес, равный количеству общих переменных в вершинах.
5. Определение максимального остовного дерева в графе (**Junction Tree**).
 - ✓ Максимальное остовное дерево — это дерево, включающее все вершины взвешенного неориентированного графа, максимизирующее общий вес рёбер.
6. Распределение факторов.
 - ✓ Каждый фактор из исходного факторного графа должен быть назначен одному из кластеров (вершин, клик) **junction tree**, который содержит все переменные этого фактора.

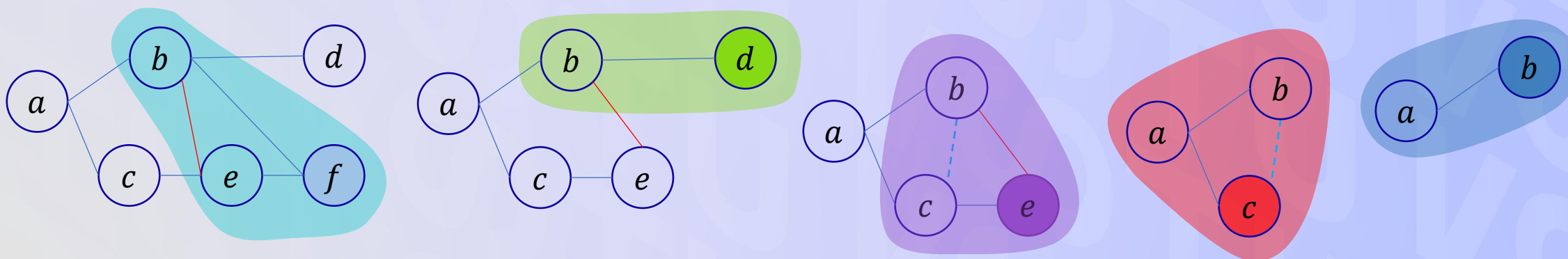
Пример построения Junction Tree



Исходный факторный граф



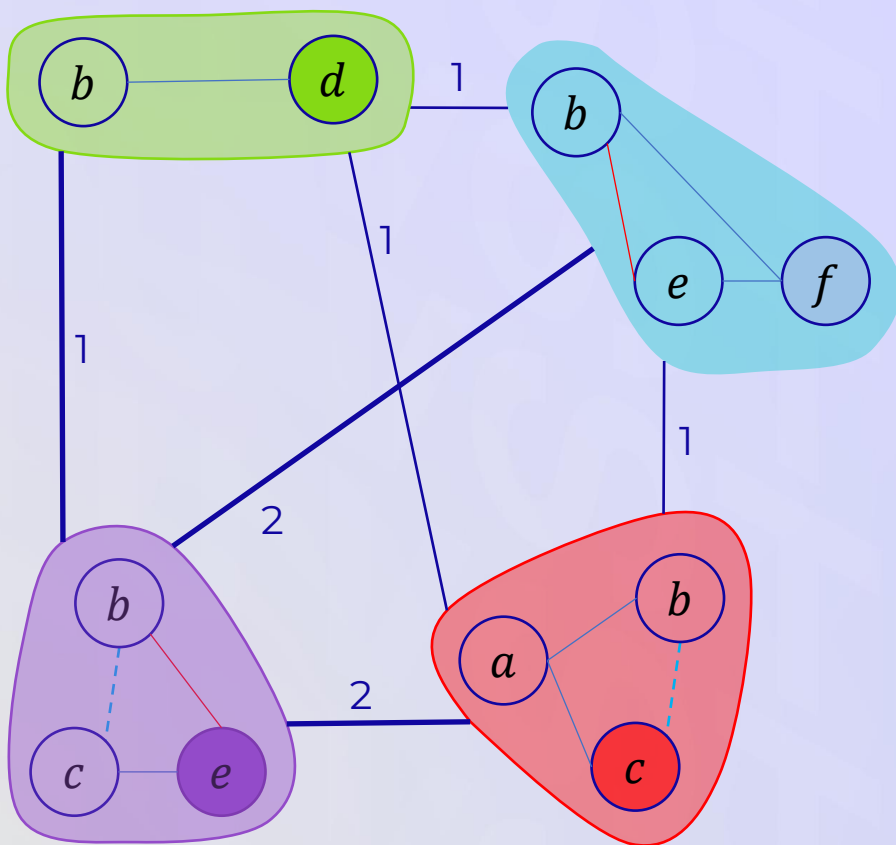
Морализованная версия графа



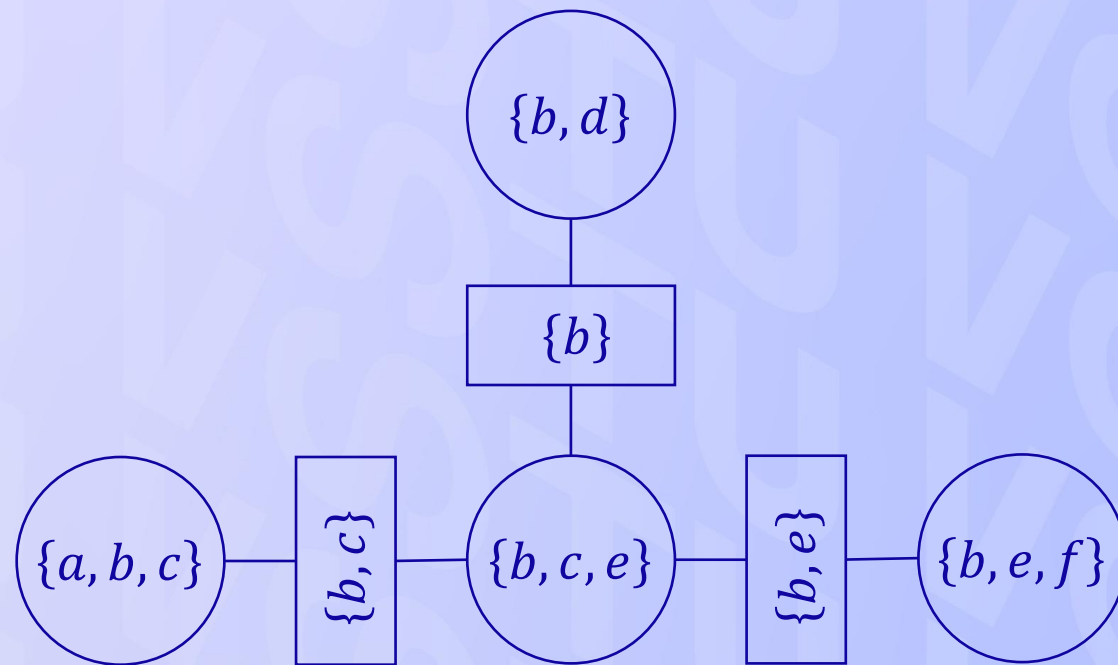
Триангуляция графа в последовательности $\{f, d, e, c, b, a\}$



Пример построения Junction Tree



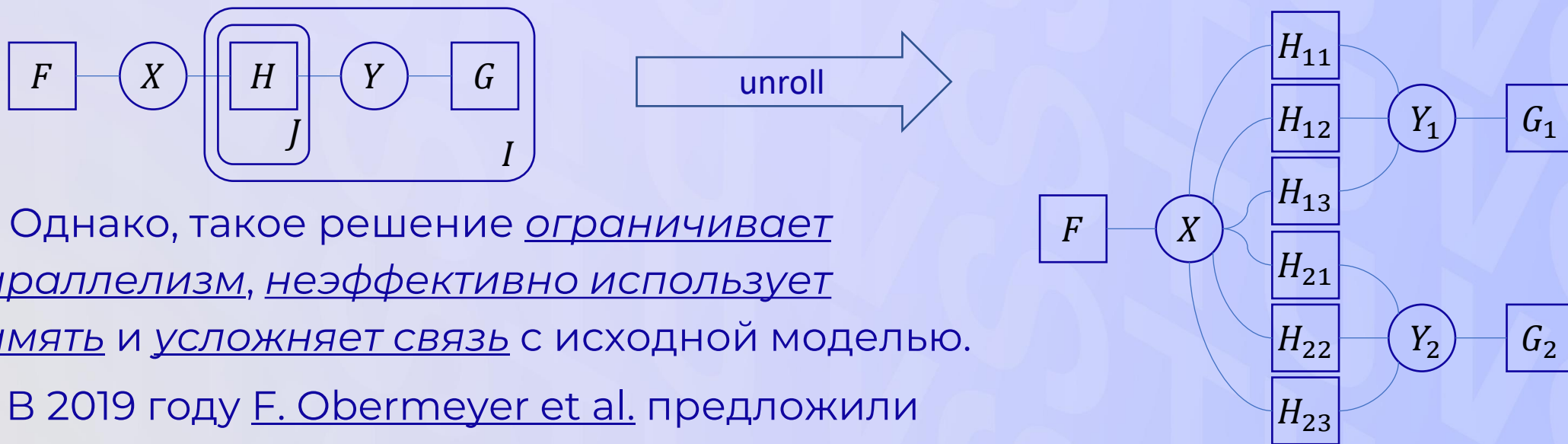
Граф кластеров



Junction Tree

- Основная идея алгоритма Junction Tree заключается в преобразовании факторного графа в дерево кластеров, на котором эффективно выполняется алгоритм устранения переменных.
- Затем в Junction Tree просто выполняется передача сообщений.
- Однако, построение Junction Tree неоднозначно. Например, всегда можно найти тривиальное Junction Tree с одной вершиной, содержащей все переменные исходного графа. Но данное дерево бесполезно.
- Оптимальными можно считать деревья, которые делают кластеры настолько маленькими и модульными, насколько это возможно.
- Но нахождение оптимального дерева является NP-сложной задачей.

- А как выполнять процедуру variable elimination для plated-графов?
- Можно реализовать алгоритм суммирования-произведения на графе факторов с plate-нотацией, применив его к развернутой версии графа:



- Однако, такое решение ограничивает параллелизм, неэффективно использует память и усложняет связь с исходной моделью.
- В 2019 году F. Obermeyer et al. предложили новый алгоритм: Tensor variable elimination.

Предложенный алгоритм в псевдокоде записывается следующим образом:

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow$ leaf plate set в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.\text{append}(\text{Product}(f, L))$

else:

$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.

Вычисляет сильно
связанные компоненты
двудольного графа

Выполняет устранение
переменных в группе
структурно идентичных
факторных графов

Выполняет поэлементное
произведение факторов по
одному или нескольким
индексам в plate

Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow$ *leaf plate set* в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

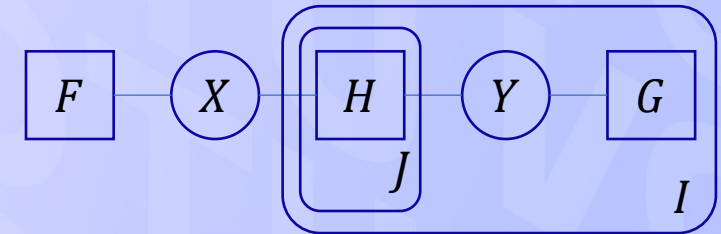
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

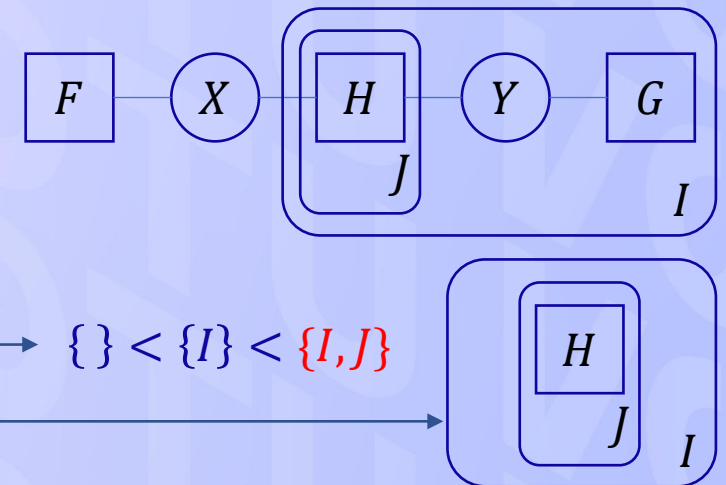
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$; —

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.\text{append}(\text{Product}(f, L))$

else:

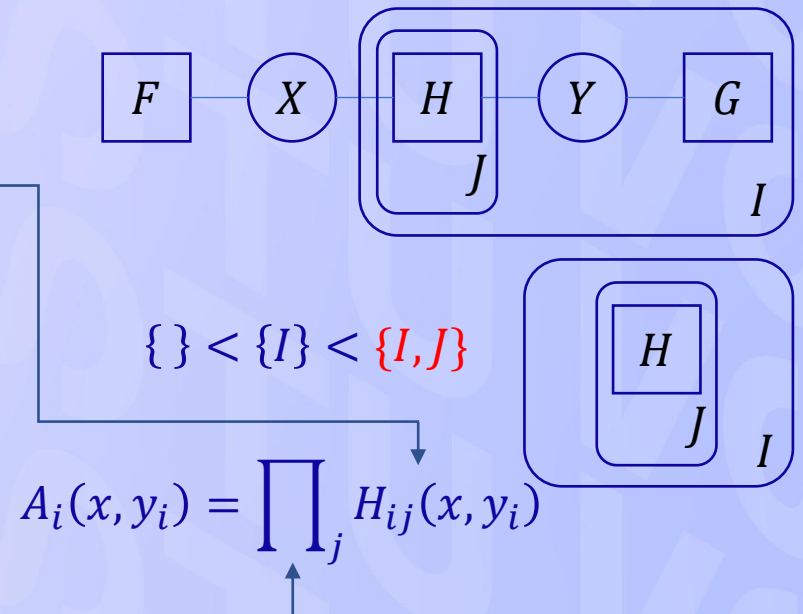
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$; —

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

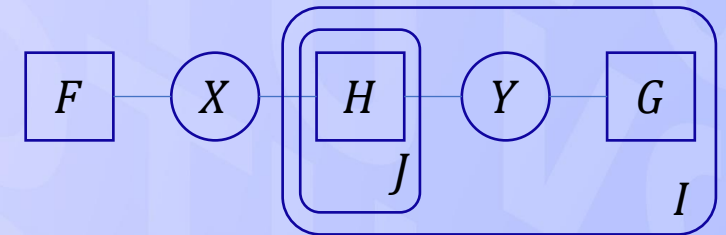
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

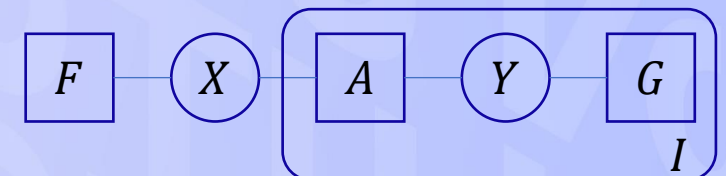
вставить f' в G ;

return $\text{SumProduct}(S, \{ \})$.



$$\{\} < \{I\} < \{I, J\}$$

$$A_i(x, y_i) = \prod_j H_{ij}(x, y_i)$$



Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow$ *leaf plate set* в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

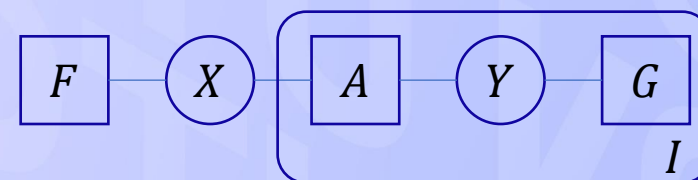
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set в } G \text{ максимального размера};$

$G_L \leftarrow \text{подграф графа } G \text{ в } L;$

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C);$

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.\text{append}(\text{Product}(f, L))$

else:

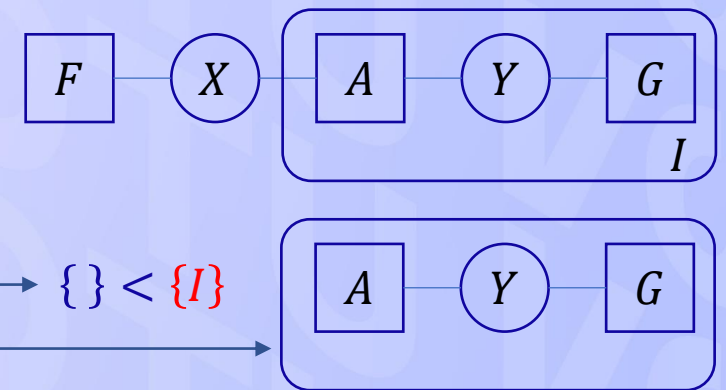
$L' \leftarrow \text{plate всех переменных } f \text{ в } G;$

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L');$

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$; —

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

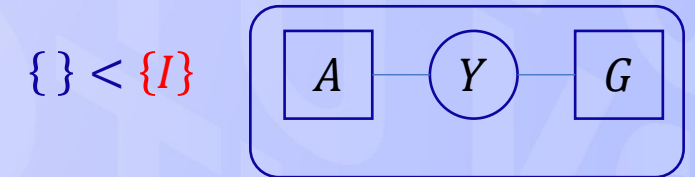
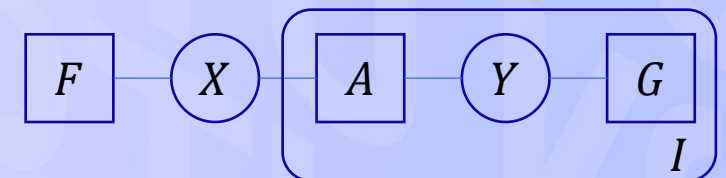
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$; —

вставить f' в G ;

return $\text{SumProduct}(S, \{ \})$.



$$B(x) = \prod_i \sum_{y_i} A_i(x, y_i) \cdot G_i(y_i)$$

Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

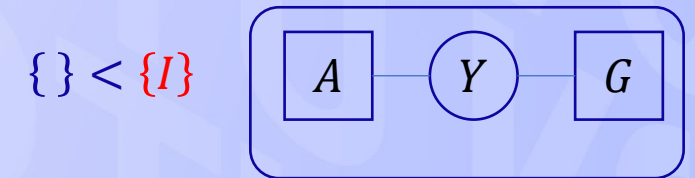
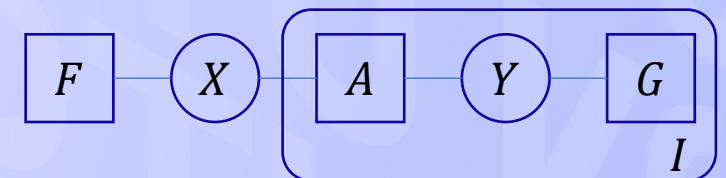
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

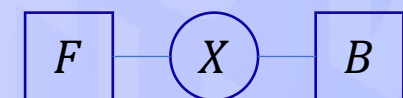
$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{ \})$.



$$B(x) = \prod_i \sum_{y_i} A_i(x, y_i) \cdot G_i(y_i)$$



Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow$ *leaf plate set* в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

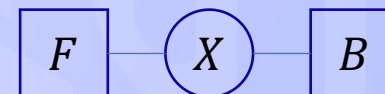
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



Алгоритм устранения тензорной переменной (Tensor variable elimination)

$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

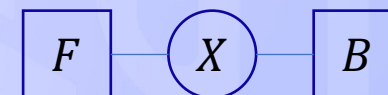
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



$\{\}$



$S \leftarrow \text{list}([])$.

while (есть факторы в графе G):

$L \leftarrow \text{leaf plate set}$ в G максимального размера;

$G_L \leftarrow$ подграф графа G в L ;

for подграф G_C in **Partition**(G_L):

$f \leftarrow \text{SumProduct}(G_C)$;

Удалить подграф G_C из G

If $V_f = \emptyset$ (множество оставшихся в f переменных):

$S.append(\text{Product}(f, L))$

else:

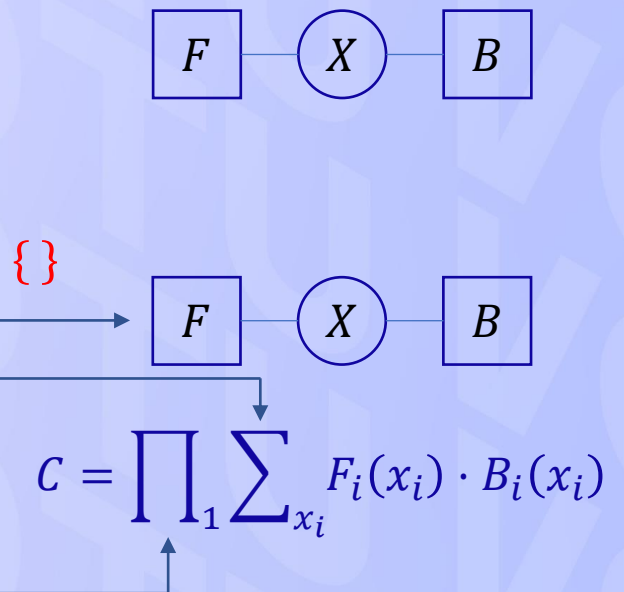
$L' \leftarrow$ plate всех переменных f в G ;

if $L' = L$ **then error** ("Intractable!")

$f' \leftarrow \text{Product}(f, L - L')$;

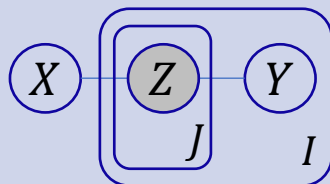
вставить f' в G ;

return $\text{SumProduct}(S, \{\})$.



```
 $S \leftarrow \text{list}([]).$   
while (есть факторы в графе  $G$ ):  
     $L \leftarrow \text{leaf plate set}$  в  $G$  максимального размера;  
     $G_L \leftarrow$  подграф графа  $G$  в  $L$ ;  
    for подграф  $G_C$  in Partition( $G_L$ ):  
         $f \leftarrow \text{SumProduct}(G_C);$   
        Удалить подграф  $G_C$  из  $G$   
        If  $V_f = \emptyset$  (множество оставшихся в  $f$  переменных):  
             $S.append(\text{Product}(f, L))$   
        else:  
             $L' \leftarrow$  plate всех переменных  $f$  в  $G$ ;  
            if  $L' = L$  then error ("Intractable!")  
             $f' \leftarrow \text{Product}(f, L - L');$   
            вставить  $f'$  в  $G$ ;  
return  $\text{SumProduct}(S, \{\}).$ 
```

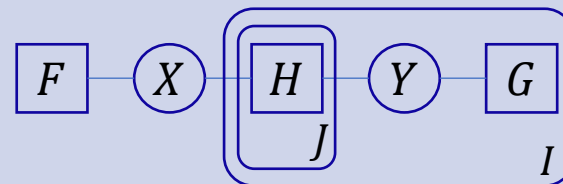
Высокоуровневый интерфейс для
определения модели:



@pyro.infer.config_enumerate

```
def model(z):
    I, J = z.shape
    x = pyro.sample("x", Bernoulli(Px))
    with pyro.plate("I", I):
        y = pyro.sample("y", Bernoulli(Py))
        with pyro.plate("J", J):
            pyro.sample(
                "z", Bernoulli(Pz[x,y]), obs=z)
```

Низкоуровневый интерфейс для
определения графа факторов:



```
pyro.ops.contract.einsum(
    "x,iy,ijxy->",
    F, G, H,
    plates="ij"
)
```

- Декоратор *config_enumerate* предписывает выполнить устранение всех дискретных переменных.
- Если требуется выполнить процедуру устранения только для некоторых дискретных переменных, то при их объявлении в метод *sample()* нужно передать параметр: *infer={"enumerate": "parallel"}* или *infer={"enumerate": "sequential"}*.
- Все дискретные переменные, участвующие в процедуре устранения, должны отсутствовать в вариационном распределении (*guide*). Если *guide* формируется автоматически, устраняемые переменные необходимо указать явно при создании класса *guide*:

```
guide = AutoNormal(poutine.block(model, hide=["x", "data"]))
```

- Стратегия устранения переменной может использоваться совместно с методами *HMC, NUTS, SVI*.
- Для работы с процедурой устранения переменной необходимо использовать варианты классов *ELBO* с приставкой “*Enum_*”.
- Для получения доступа к предсказанным в ходе вариационного инференса дискретным переменным необходимо использовать обработчик *infer_discrete*:

```
serving_model = infer_discrete(model, first_available_dim=-1)
```

- Работа процедуры устранения переменной базируется на работе с формой тензоров и использует правила *broadcasting*.

- Тензор, содержащий количество элементов некоторого тензора *вдоль* каждого его измерения, называется формой тензора (*shape*).
- Тензор в PyTorch имеет только один атрибут формы (*.shape*):

```
x = torch.Tensor([[1, 2], [3, 4], [5, 6]])  
assert x.shape == torch.Size([3, 2])
```

- Распределение имеет два атрибута формы (*.batch_shape* и *.event_shape*)

```
y = dist.MultivariateNormal(torch.zeros([2, 3]), diag(torch.ones(3)))  
assert (y.batch_shape, y.event_shape) == (torch.Size([2]), torch.Size([3]))
```

- *Индексы* в *.batch_shape* определяют условно-независимые переменные, а в *.event_shape* – зависимые (один семпл из распределения). Их комбинация определяет полную форму.

- В Pyro можно изменять размерности и формы распределений.
- Для получения батча семплов можно к распределению применить метод `expand()` или воспользоваться нотацией `plate`. Данные батчи условно независимы. Приведенные способы изменяют `batch_shape` распределения.
- Некоторую размерность распределения можно объявить зависимой. Для этого используется метод распределения `to_event(n)`, в который передается количество размерностей в `batch_shape` справа, которые должны интерпретироваться как зависимые. После вызова метода они добавляются к размерностям в `event_shape` слева.

```
d = Bernoulli(torch.tensor([0.1, 0.2, 0.3, 0.4])).expand([3, 4]).to_event(1)
```

- Предполагать зависимость переменных всегда безопасно.

- Главная идея перечисления (*enumeration*) заключается в интерпретации операторов `sample()` для дискретных переменных как полного перечисления, а не как *случайной выборки*. В дальнейшем, другие алгоритмы вывода могут суммировать перечисляемые значения.
- Для этого для каждой *устраняемой* переменной в тензорах, зависящих от нее, вводится новая размерность, значения тензоров в которой зависят от *конкретного значения* устраняемой переменной.
- Данная размерность должна быть введена “левее” уже существующих размерностей, что контролируется параметрами *first_available_dim* и *max_plate_nesting*, передаваемыми в функции во время вывода.
- Для сложной индексации по устраняемой переменной применяют Vindex

Демонстрация практических примеров

Заключение

1. Рассмотрели подходы к точному инференсу дискретных латентных переменных и разобрали особенности этих подходов.
 2. Познакомились с процедурой устранения переменной.
 3. Ввели понятия графа факторов и сравнили его с байесовской сетью.
 4. Поговорили про механизм сообщений между факторами и переменными, сформулировали метод суммирования-произведения для маргинализации.
 5. Познакомились с алгоритмом Junction Tree.
 6. Поговорили про особенности процедуры устранения переменной в plate-графах, разобрали в псевдокоде и на примере алгоритм устранения тензорной переменной, обсудили его реализацию в Pyro.
 7. Поговорили об особенностях реализации алгоритма в Pyro.
-

Спасибо за внимание!

Волгоград 2025
