

Московский государственный университет имени М. В. Ломоносова

Факультет Вычислительной Математики и Кибернетики

Кафедра Математических Методов Прогнозирования

ДИПЛОМНАЯ РАБОТА

Вычислительные методы приближенного подсчета нормировочной константы марковского случайного поля

Выполнил:

студент 517 группы

Новиков Александр Витальевич

Научный руководитель:

к.ф.-м.н., доцент

Ветров Дмитрий Петрович

Москва, 2015

Содержание

1	Введение	5
2	Обзор литературы	7
3	Обозначения и формальная постановка задачи	8
4	ТТ-разложение	9
5	Тензоры марковского случайного поля	12
5.1	ТТ-формат для тензора энергии	13
5.2	ТТ-формат для тензора вероятности	15
6	Тензорный метод	15
6.1	Алгоритм	16
6.2	Анализ алгоритма 1	18
6.3	Доказательство теоремы 2	20
6.4	Доказательство следствия 1	21
7	Метод n-окрестностей	22
7.1	Ограничения подхода	24
8	Метод «Суммировать всё»	25
8.1	Эффективные формулы для модели Изинга	27
8.2	Эффективные формулы для модели RBM	28
9	Метод, основанный на формуле Тейлора	29
10	Вычислительные эксперименты	29
10.1	Модель Изинга	30
10.2	ТТ-ранги энергии и вероятности	31
10.3	Обучение модели RBM	31
10.4	Обсуждение и выводы	33
11	Заключение	34

Аннотация

В данной работе предлагаются три метода приближенного подсчёта нормировочной константы марковского случайного поля. Первый метод основан на разложении Tensor Train, которое позволяет компактно хранить тензор и эффективно применять к нему операции линейной алгебры. Свойства этого разложения используются для оценки нормировочной константы и получения теоретических гарантий на точность работы метода. Второй метод ускоряет и обобщает недавно предложенный метод n -окрестностей до применимости к произвольным марковским случайным полям. Третий метод основан на разложении показательной функции в определении нормировочной константы в ряд Тейлора. Предложенные методы сравниваются друг с другом и с доступными аналогами на моделях Изинга и на задаче обучения ограниченной машины Больцмана по выборке MNIST.

1 Введение

Вероятностные графические модели являются удобным инструментом для решения задач в таких областях, как компьютерное зрение, машинное обучение, анализ социальных сетей и т. д. Одно из их основных достоинств — это возможность компактно задать ненормированное распределение на сотни переменных в виде произведения факторов низких порядков. Это позволяет, с одной стороны, учесть сложные зависимости в данных, а с другой стороны точно или приближенно вычислять статистики полученного распределения.

Многие возникающие задачи становятся сложными, если фактор-граф содержит циклы. В работе рассматриваются ненаправленные графические модели (марковские случайные поля, MRF), и поднимается актуальная для них задача оценки нормировочной константы распределения по параметрам модели. Решение этой задачи особенно важно при настройке параметров графической модели по тренировочной выборке (обучение по методу максимального правдоподобия [24]). В последние годы для этой задачи предложено множество различных методов, но, тем не менее, её пока рано называть полностью решенной.

Совместное распределение дискретных переменных можно задать через многомерный массив (тензор) ненормированной вероятности. Тогда подсчет нормировочной константы сводится к суммированию всех элементов тензора. Для хранения и обработки тензора в явном виде требуется экспоненциальное количество памяти и вычислительных ресурсов. Тензорные разложения позволяют компактно хранить тензор и эффективно совершать над ним различные операции.

В дипломной работе рассматривается недавно предложенное И. Оселедцом [18] разложение тензорного произведения (ТТ-разложение или ТТ-формат). Алгоритмы ТТ-разложения хранят тензор в специальном формате, который позволяет эффективно оперировать над тензорами. Эффективность ТТ-разложения определяется величиной ТТ-рангов тензора. В работе показано, что тензор энергии (минус логарифма вероятности) марковского случайного поля обладает низкими ТТ-рангами, а тензор вероятности — высокими ТТ-рангами. Для подсчета нормировочной константы в условиях высоких ТТ-рангов тензора вероятности предложен метод, основанный на факторизации графической модели. Факторы графической модели переводятся

в ТТ-формат и комбинируются таким образом, чтобы оценить нормировочную константу, не строя сам тензор вероятности.

Так же в работе рассматривается предложенный в 2014 году метод n -окрестностей [10], придуманный для оценки нормировочной константы (или, как её называют в физической литературе, «статистической суммы») модели Изинга¹. Основная идея метода n -окрестностей – разбить множество значений энергии на подмножества и приблизить эмпирическую функцию распределения энергий в каждом подмножестве с помощью нормального распределения. В дипломной работе предлагается отказаться от рассмотрения отдельных подмножеств и рассматривать множество всех энергий целиком. Полученные формулы позволяют применять такой метод к произвольному MRF (тогда как метод n -окрестностей предложен только для модели Изинга), а использование полного множества энергий в ряде случаев позволяет получить более точные результаты.

Формулы, выведенные для обобщения метода n -окрестностей, используются в третьем предложенном методе, основанном на разложении показательной функции в определении нормировочной константы в ряд Тейлора.

Основной вклад настоящей работы заключается в следующем:

- Предложен алгоритм оценки нормировочной константы через ТТ-разложение факторов графической модели, который не строит ТТ-представление тензора ненормированного совместного распределения.
- Приведены теоретические гарантии на точность оценки нормировочной константы.
- Предложено обобщение метода n -окрестностей, применимое к произвольным марковским случайным полям и обладающее более высокой скоростью работы.
- Предложен простой и быстрый метод оценки нормировочной константы через разложение в ряд Тейлора.

¹Модель Изинга – графическая модель, используемая в статистической физики для моделирования магнитных свойств вещества.

Настоящая дипломная работа организована следующим образом. В разделе 4 приведён краткий обзор тензорного разложения Tensor Train, в разделе 5 поясняется его связь с MRF, а в разделе 6 описывается предлагаемый метод оценки нормировочной константы, основанный на разложении Tensor Train. В разделе 7 приведён краткий обзор недавно предложенного в физической литературе метода для подсчета нормировочной константы модели Изинга, а в разделе 8 описывается предлагаемое обобщение и ускорение данного метода, которое называется методом «Суммировать всё». В разделе 9 предлагается простой и быстрый метод для оценки нормировочной константы, который основан на разложении показательной функции в определении ненормированного распределения Гиббса в ряд Тейлора. Раздел 10 содержит результаты экспериментального сравнения предлагаемых подходов друг с другом и с лучшими аналогами.

2 Обзор литературы

Обзор можно разделить на три основных направления: тензорные разложения, позволяющие компактно представить тензор, использование тензорных методов в графических моделях и различные алгоритмы оценки нормировочной константы MRF.

В отличие от классических тензорных форматов (канонический формат [1] и формат Таккера [22]), для ТТ-разложения существуют устойчивые алгоритмы построения, а само ТТ-разложение не подвержено проклятию размерности. Альтернативой ему является иерархический формат Таккера [5, 3], который также позволяет устойчиво искать компактное представление тензора. ТТ-формат эквивалентен иерархическому формату Таккера с линейным деревом размерностей, что является преимуществом с алгоритмической точки зрения: большинство алгоритмов ТТ-формата значительно проще аналогов для иерархического формата Таккера.

В статьях по графическим моделям тензорные подходы часто упоминаются применительно к восстановлению структуры модели. Джернит и др.; Мария Иштева и др. принимают локальные решения, основываясь на свойствах четырехмерных тен-

зоров [8, 7]. Сонг и др. используют иерархические тензорные разложения для восстановления параметров графической модели [21].

Можно выделить следующие группы методов подсчета нормировочной константы: методы, основанные на генерации выборки (например Annealed Importance Sampling [15, 4]), методы передачи сообщений (Loopy Belief Propagation [11] и его многочисленные модификации), методы, основанные на минимизации KL-дивергенции (Mean Field [24] и Expectation Propagation [13]), методы декомпозиции графов (Tree-Reweighted Message-passing [23]) и методы, основанные на минимизации энергии (randomized MAP-predictors [6]).

3 Обозначения и формальная постановка задачи

В данной статье часто будет использоваться понятие многомерного массива вещественных чисел. Одномерные массивы будем называть *векторами*, двумерные — *матрицами*, а массивы большей размерности будем называть *тензорами*. Для обозначения векторов будут использоваться маленькие жирные буквы (\mathbf{a}), а для матриц и тензоров большие жирные буквы (\mathbf{A}).

Будем рассматривать все массивы как функции их индексов: $a(i) = a_i$, $A(x_1, x_2)$, $A(\mathbf{x}) = A(x_1, \dots, x_n)$, где n — это размерность тензора \mathbf{A} .

Символом $\|\cdot\|_F$ будем обозначать норму Фробениуса матриц и её обобщение на случай тензоров:

$$\|\mathbf{A}\|_F = \sqrt{\sum_{x_1, \dots, x_n} A^2(x_1, \dots, x_n)}.$$

Под $\|\cdot\|_2$ будем понимать евклидову норму векторов и спектральную норму матриц: $\|\mathbf{A}\|_2 = \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}$.

В работе будут использоваться несколько различных матричных произведений: произведение Адамара (символ “ \odot ”), произведение Кронекера (символ “ \otimes ”) и обычное матричное произведение (символ “ \cdot ”).

Введём теперь формально определение марковского случайного поля. Рассмотрим гиперграф $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ с конечным множеством вершин \mathcal{V} и гиперребра \mathcal{E} . Пусть все вершины пронумерованы от 1 до n , а все гиперребра от 1 до m .

Сопоставим каждой вершине $i = 1, \dots, n$ переменную x_i принимающую значения из множества \mathcal{X}_i мощности d_i : $|\mathcal{X}_i| = d_i$ (например $\mathcal{X}_i = \{1, \dots, d_i\}$). Обозначим символом \mathcal{X} совместную область определения переменных x_1, \dots, x_n : $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_n$. Обозначим множество вершин, входящих в гиперребро $\ell = 1, \dots, m$, через \mathbf{x}^ℓ . Каждому гиперребру $\ell = 1, \dots, m$ сопоставим вещественнозначную функцию Θ_ℓ , определенную на совместной области определения переменных из \mathbf{x}^ℓ . Функцию Θ_ℓ будем называть *потенциалом*.

Функцией энергии марковского случайного поля (Markov random field, MRF) заданного на гиперграфе \mathcal{G} , назовем сумму всех потенциалов: $E(\mathbf{x}) = \sum_{\ell=1}^m \Theta_\ell(\mathbf{x}^\ell)$.

Экспонента от минус энергии задает *ненормированное распределение Гиббса*: $\hat{P}(\mathbf{x}) = \exp(-E(\mathbf{x}))$. Для обозначения нормировочной константы будет использоваться символ Z :

$$Z = \sum_{\mathbf{x}} \hat{P}(\mathbf{x}). \quad (1)$$

В дипломной работе рассматривается следующая задача: по заданному MRF (т. е. гиперграфу и функции энергии) требуется оценить величину логарифма нормировочной константы $\ln(Z)$.

4 ТТ-разложение

Будем говорить, что n -мерный тензор \mathbf{A} представлен в ТТ-формате, если для всех размерностей $i = 1, \dots, n$ и всех значений индексов по этой размерности $x_i = 1, \dots, d_i$ ($d = \max_{i=1, \dots, n} d_i$) существуют матрицы $G_i^{\mathbf{A}}[x_i]$, такие, что каждый элемент тензора \mathbf{A} представим в виде произведения матриц:

$$A(x_1, \dots, x_n) = \mathbf{G}_1^{\mathbf{A}}[x_1] \mathbf{G}_2^{\mathbf{A}}[x_2] \dots \mathbf{G}_n^{\mathbf{A}}[x_n]. \quad (2)$$

При этом все матрицы относящиеся к одной и той же размерности i должны иметь одинаковые размеры $r_{i-1}(\mathbf{A}) \times r_i(\mathbf{A})$. Чтобы результат матричного произведения (2) был равен числу, положим $r_0(\mathbf{A}) = r_n(\mathbf{A}) = 1$. Последовательность $\{r_i(\mathbf{A})\}_{i=0}^n$ будем называть *ТТ-рангами* тензора \mathbf{A} , а максимальный элемент последовательности — *максимальным ТТ-рангом* тензора \mathbf{A} : $r(\mathbf{A}) = \max_{i=0, \dots, n} r_i(\mathbf{A})$. Набор матриц $\mathbf{G}_i^{\mathbf{A}}$, со-

ответствующих одному измерению, называется *ТТ-ядром* тензора \mathbf{A} . Представление тензора в ТТ-формате будем называть *ТТ-разложением* или *ТТ-представлением*.

Для любого n -мерного тензора \mathbf{A} существует ТТ-представление с максимальным ТТ-рангом $r(\mathbf{A}) \leq d^{\frac{n}{2}}$ (см. теорему 2.1 Оселедца [18]). Отметим, что ТТ-представление тензора не единственно.

Для обозначения (α_{i-1}, α_i) -ого элемента матрицы $\mathbf{G}_i^{\mathbf{A}}[x_i]$ будет использоваться символ $G_i^{\mathbf{A}}[x_i](\alpha_{i-1}, \alpha_i)$. Пользуясь определением произведения матриц, можно переписать формулу (2) через элементы ТТ-ядер:

$$A(\mathbf{x}) = \sum_{\alpha_0, \dots, \alpha_n} G_1^{\mathbf{A}}[x_1](\alpha_0, \alpha_1) \dots G_n^{\mathbf{A}}[x_n](\alpha_{n-1}, \alpha_n). \quad (3)$$

Для хранения всех элементов тензора \mathbf{A} требуется $\prod_{i=1}^n d_i$ ячеек памяти, тогда как хранение \mathbf{A} в ТТ-формате требует $\sum_{i=1}^n d_i r_{i-1}(\mathbf{A}) r_i(\mathbf{A})$. Таким образом, ТТ-представление тензора с низкими ТТ-рангами существенно компактнее перечисления его элементов.

Существуют две различные постановки задачи перевода тензора в ТТ-формат: точное ТТ-представление (алгоритм ТТ-SVD [18], применимый для небольших тензоров), и построение приближенного ТТ-представления по небольшому подмножеству элементов тензора. Наилучшим из алгоритмов второго класса в настоящее время является метод AMEn-cross [2].

Одним из основных достоинств ТТ-формата является возможность эффективно применять различные операции к тензорам в ТТ-формате: умножение тензора на константу, добавление константы к тензору, поточечное сложение и умножение тензоров (результат этих операций — это тензор в ТТ-формате с возросшими ТТ-рангами); подсчет глобальных характеристик тензора, таких как сумма всех элементов или норма Фробениуса. Обзор операций, которые используются в этой работе, приведен в таблице 1 (детальный обзор проведен Оселедцом [18]).

Применение операций к ТТ-тензорам увеличивает ТТ-ранги даже в том случае, когда существует низкоранговое ТТ-представление результата. Чтобы контролировать рост ТТ-рангов существует операция ТТ-округления. По ТТ-представлению тензора \mathbf{A} и относительной точности $\varepsilon \geq 0$ операция ТТ-округления $\text{round}(\mathbf{A}, \varepsilon)$ найдет тензор $\hat{\mathbf{A}}$ в ТТ-формате который, во-первых, достаточно близок к тензору \mathbf{A} :

Таблица 1: Операции, которые можно эффективно выполнять над тензорами в ТТ-формате. Для каждой операции указана её вычислительная сложность и ТТ-ранг результата для ситуаций, когда результат является тензором в ТТ-формате.

ОПЕРАЦИЯ	РАНГ РЕЗУЛЬТАТА	ВЫЧИСЛИТЕЛЬНАЯ СЛОЖНОСТЬ
$\mathbf{C} = \mathbf{A} \cdot \text{const}$	$r(\mathbf{C}) = r(\mathbf{A})$	$O(d r(\mathbf{A}))$
$\mathbf{C} = \mathbf{A} + \text{const}$	$r(\mathbf{C}) = r(\mathbf{A}) + 1$	$O(nd r^2(\mathbf{A}))$
$\mathbf{C} = \mathbf{A} + \mathbf{B}$	$r(\mathbf{C}) \leq r(\mathbf{A}) + r(\mathbf{B})$	$O(nd (r(\mathbf{A}) + r(\mathbf{B}))^2)$
$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	$r(\mathbf{C}) \leq r(\mathbf{A}) r(\mathbf{B})$	$O(nd r^2(\mathbf{A}) r^2(\mathbf{B}))$
$\mathbf{c} = \mathbf{M}\mathbf{b}$	$r(\mathbf{c}) \leq r(\mathbf{M}) r(\mathbf{b})$	$O(nd^2 r^2(\mathbf{M}) r^2(\mathbf{b}))$
sum \mathbf{A}	–	$O(nd r^2(\mathbf{A}))$
$\ \mathbf{A}\ _F$	–	$O(nd r^3(\mathbf{A}))$
$\mathbf{C} = \text{round}(\mathbf{A}, \varepsilon)$	$r(\mathbf{C}) \leq r(\mathbf{A})$	$O(nd r^3(\mathbf{A}))$

$\|\mathbf{A} - \widehat{\mathbf{A}}\|_F \leq \varepsilon \|\mathbf{A}\|_F$ и, во-вторых, обладает минимальными ТТ-рангами среди всех тензоров \mathbf{B} : $\|\mathbf{A} - \mathbf{B}\|_F \leq \frac{\varepsilon}{\sqrt{n-1}} \|\mathbf{A}\|_F$. Наличие операции ТТ-округления позволяет применять последовательность операций к тензорам (например, округляя результат после применения каждой операции), контролируя рост ТТ-рангов.

Для повышения эффективности работы с векторами и матрицами специальным образом вводятся понятия ТТ-формата вектора и ТТ-формата матрицы. Пусть существует отображение между индексами вектора $\mathbf{b} \in \mathbb{R}^q$ и n -мерными векторами $\mathbf{y} = (y_1, \dots, y_n)^2$. ТТ-представлением вектора \mathbf{b} называется ТТ-представление тензора $\mathbf{B}(y_1, \dots, y_n)$, содержащего все элементы \mathbf{b} .

Рассмотрим пример. Пусть вектор \mathbf{b} содержит 18 элементов. Его можно интерпретировать как (например) трехмерный тензор \mathbf{B} :

$$\begin{aligned}
 1 &\leftrightarrow (1, 1, 1) & B(1, 1, 1) &= b(1), \\
 2 &\leftrightarrow (1, 1, 2) & B(1, 1, 2) &= b(2), \\
 \dots & & \dots & \\
 18 &\leftrightarrow (2, 3, 3) & B(2, 3, 3) &= b(18).
 \end{aligned}$$

² Количество элементов в векторе \mathbf{b} равно $q = \prod_{i=1}^n d_i$.

Здесь составной индекс $\mathbf{y} = (y_1, y_2, y_3)$ нумерует элементы вектора \mathbf{b} .

Определим понятие ТТ-формата для матриц. Пусть существует отображение между индексами строк и столбцов матрицы \mathbf{M} в n -мерные вектора \mathbf{x} и \mathbf{y} соответственно. Переупорядочим размерности и представим получившийся тензор в ТТ-формате:

$$M((x_1, y_1), \dots, (x_n, y_n)) = \mathbf{G}_1^M[x_1, y_1] \dots \mathbf{G}_n^M[x_n, y_n],$$

где \mathbf{G}_i^M , $i = 1, \dots, n$ — это ТТ-ядра, а $\mathbf{G}_i^M[x_i, y_i]$ — матрицы. *ТТ-представлением матрицы \mathbf{M}* будем называть ТТ-представление тензора \mathbf{M} . Отметим, что матрица в ТТ-формате не обязана быть квадратной, т. к. x_i и y_i могут принимать разное число возможных значений

Для матрицы \mathbf{M} и вектора \mathbf{b} представленных в ТТ-формате можно эффективно вычислять произведение $\mathbf{c} = \mathbf{M}\mathbf{b}$ (если соответствующие размерности совпадают). Результатом этой операции является вектор \mathbf{c} в ТТ-формате с рангами равными произведению рангов \mathbf{M} и \mathbf{b} : $r_i(\mathbf{c}) = r_i(\mathbf{M}) r_i(\mathbf{b})$.

Наличие специального определения ТТ-формата для векторов и матриц позволяет применять операции линейной алгебры к задачам большого размера. Например, для поиска минимального элемента тензора можно вытянуть его в диагональную матрицу и применить приближенный метод поиска минимальных собственных значений, основанный на алгоритме DMRG [9].

5 Тензоры марковского случайного поля

Как энергия, так и ненормированная вероятность являются n -мерными тензорами. Потенциалы и факторы можно интерпретировать как n -мерные тензоры, если добавить в них переменные, от которых они зависят несущественно: $\Theta_\ell(\mathbf{x}^\ell) = \Theta_\ell(\mathbf{x})$. Тензоры энергии и вероятности можно определить следующим образом: $\mathbf{E} = \sum_{\ell=1}^m \Theta_\ell$ и $\hat{\mathbf{P}} = \odot_{\ell=1}^m \Psi_\ell$.

Традиционно для компактного представления энергии используется набор потенциалов MRF. ТТ-разложение является альтернативным способом компактного представления тензора энергии. Эффективность применения ТТ-формата для тензора энергии и тензора вероятности обсуждается в разделах 5.1 и 5.2.

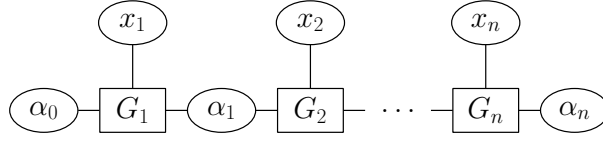


Рис. 1: Графическая модель, полученная в результате ТТ-разложения тензора вероятности $P(x_1, \dots, x_n)$.

ТТ-представление тензора совместного распределения $\mathbf{P} = \widehat{\mathbf{P}}/Z$ обладает особой интерпретацией. Переменные α_i , $i = 0, \dots, n$, (см. (3)) можно интерпретировать как скрытые переменные в графической модели типа цепочка (рис. 1). Маргинализация по ним дает исходную вероятность:

$$P(\mathbf{x}) = \sum_{\boldsymbol{\alpha}} P(\mathbf{x}, \boldsymbol{\alpha}), \quad (4)$$

где $P(\mathbf{x}, \boldsymbol{\alpha})$ — совместное распределение вероятностей на векторы \mathbf{x} и $\boldsymbol{\alpha}$. Каждая добавленная переменная α_i может принимать значения от 1 до соответствующего ТТ-ранга $r_i(\mathbf{P})$.

5.1 ТТ-формат для тензора энергии

В этом разделе предложен алгоритм для представление тензора энергии \mathbf{E} в ТТ-формате. Предлагаемый алгоритм существенно превосходит АМEn-cross по скорости и точности.

Алгоритм состоит из трех основных шагов.

Шаг 1. Найдем ТТ-представление для каждого потенциала $\Theta_\ell(\mathbf{x}^\ell)$, который будем интерпретировать как тензор зависящий только от переменных \mathbf{x}^ℓ . Обычно потенциалы существенно зависят лишь от небольшого числа переменных и можно искать их ТТ-представление с помощью алгоритма ТТ-SVD.

Шаг 2. Добавим в каждый тензор Θ_ℓ несущественные переменные $\mathbf{x} \setminus \mathbf{x}^\ell$, сделав его n -мерным тензором.

Пусть тензор Θ_ℓ существенно зависит от p переменных: $\mathbf{x}^\ell = (x_{i_1}, \dots, x_{i_p})$, $i_1 < i_2 < \dots < i_p$. После первого шага алгоритма получим следующее ТТ-пред-

ставление $\Theta_\ell(\mathbf{x}^\ell)$:

$$\Theta_\ell(x_{i_1}, \dots, x_{i_p}) = \overline{\mathbf{G}}_1^{\Theta_\ell}[x_{i_1}] \dots \overline{\mathbf{G}}_p^{\Theta_\ell}[x_{i_p}], \quad (5)$$

где $\overline{\mathbf{G}}_k^{\Theta_\ell}[x_{i_k}]$ — матрица размера $\bar{\Gamma}_{k-1} \times \bar{\Gamma}_k$. Чтобы добавить несущественные размерности, можно, например, положить недостающие ядра равными набору единичных матриц подходящего размера³:

$$\begin{aligned} \mathbf{G}_1^{\Theta_\ell}[x_1] &\equiv \dots \equiv \mathbf{G}_{i_1-1}^{\Theta_\ell}[x_{i_1-1}] \equiv \mathbf{I}_{\bar{\Gamma}_0} = \mathbf{I}_1, \\ \mathbf{G}_{i_1+1}^{\Theta_\ell}[x_{i_1+1}] &\equiv \dots \equiv \mathbf{G}_{i_2-1}^{\Theta_\ell}[x_{i_2-1}] \equiv \mathbf{I}_{\bar{\Gamma}_1}, \\ &\dots \\ \mathbf{G}_{i_p+1}^{\Theta_\ell}[x_{i_p+1}] &\equiv \dots \equiv \mathbf{G}_n^{\Theta_\ell}[x_n] \equiv \mathbf{I}_{\bar{\Gamma}_p} = \mathbf{I}_1, \end{aligned}$$

где $\mathbf{G}_k^{\Theta_\ell}[x_{i_k}] = \overline{\mathbf{G}}_k^{\Theta_\ell}[x_{i_k}]$. Полученные таким образом ГТ-ядра задают ГТ-формат тензора $\Theta_\ell(\mathbf{x})$. Отметим, что в процессе данного преобразования максимальный ГТ-ранг Θ_ℓ не изменяется.

Шаг 3. Просуммируем полученные на шаге 2 ГТ-тензоры Θ_ℓ , чтобы получить тензор энергии:

$$\mathbf{E} = \sum_{\ell=1}^m \Theta_\ell. \quad (6)$$

Теорема 1 содержит верхнюю оценку на максимальный ГТ-ранг тензора, построенного описанным алгоритмом.

Теорема 1. *Если порядок каждого потенциала не превосходит p , то предложенный алгоритм строит ГТ-представление тензора энергии \mathbf{E} с максимальным ГТ-рангом, ограниченным сверху следующим образом:*

$$r(\mathbf{E}) \leq d^{\frac{p}{2}} \cdot m. \quad (7)$$

Теорема 1 показывает, что можно построить точное низкоранговое ГТ-представление тензора энергии. Отметим, что после шага 3 к тензору энергии можно применить процедуру ГТ-округления, чтобы уменьшить ГТ-ранги, если это возможно.

³ \mathbf{I}_k используется для обозначения единичной матрицы размера $k \times k$.

Доказательство теоремы 1. Оценим, какими ТТ-рангами обладают факторы после шага 1 алгоритма. Из теоремы 2.1 [18] следует, что ТТ-ранги тензора не превосходят $d^{\frac{n}{2}}$, где n — это размерность тензора, а d — максимальное число возможных значений, которые может принимать каждая переменная x_i . По условию теоремы, порядок каждого потенциала не превосходит p , а значит после шага 1 верна следующая оценка:

$$r(\Theta_\ell) \leq d^{\frac{p}{2}}. \quad (8)$$

На шаге 2 максимальный ТТ-ранг потенциалов не возрастает, а значит после шага 2 неравенство (8) всё ещё выполнено.

При суммировании тензоров на шаге 3, ТТ-ранги растут аддитивно. Чтобы завершить доказательство, осталось вспомнить, что всего есть m потенциалов. \square

5.2 ТТ-формат для тензора вероятности

Алгоритм из раздела 5.1 легко адаптировать для построения ТТ-представления тензора ненормированной вероятности $\hat{\mathbf{P}}$.

На шагах 1 и 2 алгоритма вместо потенциалов Θ_ℓ следует работать с факторами Ψ_ℓ . На шаге 3 вместо суммы нужно вычислить поэлементное произведение:

$$\hat{\mathbf{P}} = \bigodot_{\ell=1}^m \Psi_\ell. \quad (9)$$

Данный алгоритм построит точное ТТ-представление тензора $\hat{\mathbf{P}}$. Тем не менее, ТТ-ранги $\hat{\mathbf{P}}$ экспоненциально зависят от количества вершин (рис. 2), что делает ТТ-представление тензора вероятности неприменимым на больших задачах. На рис. 2 изображены ранги точного ТТ-представления тензора энергии и вероятности для MRF разного размера. ТТ-ранги тензора вероятности остаются экспоненциально большими и после округления $\hat{\mathbf{P}}$ с точностью $\varepsilon = 10^{-8}$. Таким образом, точного низкорангового ТТ-представления $\hat{\mathbf{P}}$ не существует.

6 Тензорный метод

Естественный подход к подсчету нормировочной константы состоит в представлении всего тензора ненормированной вероятности $\hat{\mathbf{P}}$ в ТТ-формате и в последующем

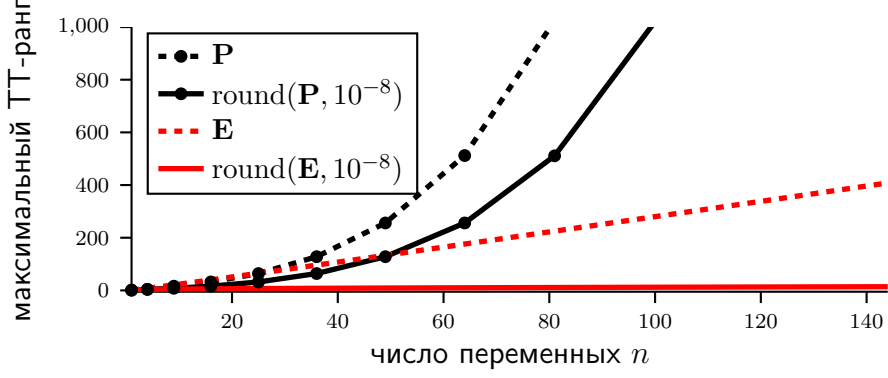


Рис. 2: Максимальный ТТ-ранг тензора энергии \mathbf{E} и тензора ненормированной вероятности $\hat{\mathbf{P}}$ для модели Изинга с температурой равной 10 и весами парных потенциалов равными 1. Детали см. в разделе 10.2.

суммировании всех его элементов. На практике тензор $\hat{\mathbf{P}}$ обладает экспоненциально большими ТТ-рангами, и работа с его ТТ-представлением становится неэффективной. С другой стороны, каждый отдельный фактор графической модели можно точно представить в ТТ-формате. В этом разделе предложен алгоритм подсчета нормировочной константы, который работает с ТТ-представлением отдельных факторов Ψ_ℓ , не строя ТТ-представление всего тензора $\hat{\mathbf{P}}$.

6.1 Алгоритм

Пусть все факторы графической модели уже представлены в ТТ-формате (см. раздел 5.2):

$$\Psi_\ell(\mathbf{x}^\ell) = \Psi_\ell(\mathbf{x}) = \mathbf{G}_1^{\Psi_\ell}[x_1] \dots \mathbf{G}_n^{\Psi_\ell}[x_n]. \quad (10)$$

Далее символ $G_i^\ell[x_i](\alpha_{i-1}^\ell, \alpha_i^\ell)$ будет использоваться как сокращенное обозначение для $G_i^{\Psi_\ell}[x_i](\alpha_{i-1}^{\Psi_\ell}, \alpha_i^{\Psi_\ell})$.

По определению, нормировочная константа Z вычисляется как сумма значений ненормированного распределения на всех конфигурациях:

$$Z = \sum_{\mathbf{x}} \prod_{\ell=1}^m \underbrace{\Psi_\ell(\mathbf{x})}_{\in \mathbb{R}} = \sum_{x_1, \dots, x_n} \bigotimes_{\ell=1}^m (G_1^\ell[x_1] \dots G_n^\ell[x_n]).$$

Второе равенство выполнено, т. к. кронекерово произведение чисел (матриц размера 1×1) эквивалентно обычному произведению.

Пользуясь свойством смешанного произведения $\mathbf{AC} \otimes \mathbf{BD} = (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D})$, преобразуем выражение для Z :

$$Z = \sum_{x_1, \dots, x_n} (\mathbf{G}_1^1[x_1] \otimes \dots \otimes \mathbf{G}_1^m[x_1]) \dots (\mathbf{G}_n^1[x_n] \otimes \dots \otimes \mathbf{G}_n^m[x_n]).$$

Обозначим кронекерово произведение матриц $\mathbf{G}_i^\ell[x_i]$ через $\mathbf{A}_i[x_i]$:

$$\mathbf{A}_i[x_i] = \mathbf{G}_i^1[x_i] \otimes \dots \otimes \mathbf{G}_i^m[x_i].$$

Для любого значения x_i матрица $\mathbf{A}_i[x_i]$ имеет размеры $(r_{i-1}(\Psi_1) \dots r_{i-1}(\Psi_m)) \times (r_i(\Psi_1) \dots r_i(\Psi_m))$. Значение её элементов выражается через элементы матриц $\mathbf{G}_i^\ell[x_i]$:

$$A_i[x_i](\alpha_{i-1}^1, \dots, \alpha_{i-1}^m; \alpha_i^1, \dots, \alpha_i^m) = G_i^1[x_i](\alpha_{i-1}^1, \alpha_i^1) \dots G_i^m[x_i](\alpha_{i-1}^m, \alpha_i^m).$$

Таким образом, матрица $\mathbf{A}_i[x_i]$ представлена в ТТ-формате, а её ТТ-ранг равен одному (т.к. $G_i^\ell[x_i](\alpha_{i-1}^\ell, \alpha_i^\ell)$ — это матрица размера 1×1).

Представим нормировочную константу Z в виде произведения n матриц:

$$Z = \sum_{x_1, \dots, x_n} \mathbf{A}_1[x_1] \dots \mathbf{A}_n[x_n] = \left(\sum_{x_1} \mathbf{A}_1[x_1] \right) \dots \left(\sum_{x_n} \mathbf{A}_n[x_n] \right) = \mathbf{B}_1 \dots \mathbf{B}_n,$$

где

$$\mathbf{B}_i = \sum_{x_i=1}^{d_i} \mathbf{A}_i[x_i].$$

ТТ-представление матрицы \mathbf{B}_i можно получить просуммировав d_i матриц в ТТ-формате. Все матрицы \mathbf{B}_i можно построить и держать в оперативной памяти, т.к. ТТ-ранги \mathbf{B}_i не превосходят d_i .

Матрицы \mathbf{B}_1 и \mathbf{B}_n являются вектором-строкой и вектором-столбцом соответственно, а значит результат произведения $\mathbf{B}_1 \dots \mathbf{B}_n$ — это число. Построив ТТ-матрицы \mathbf{B}_i , их можно перемножать, округляя результат после каждого умножения (см. алгоритм 1). Параметр ε контролирует баланс между точностью и скоростью работы алгоритма.

Помимо нормировочной константы Z , предложенный метод также позволяет найти маргинальные распределения на переменные графической модели. Ненормированное маргинальное распределение $\hat{P}_i(x_i)$ вычисляется следующим образом: $\hat{P}_i(x_i) = \mathbf{B}_1 \dots \mathbf{B}_{i-1} \mathbf{A}_i[x_i] \mathbf{B}_{i+1} \dots \mathbf{B}_n$. При этом произведения $\mathbf{B}_1 \dots \mathbf{B}_{i-1}$ и $\mathbf{B}_{i+1} \dots \mathbf{B}_n$, $i =$

Алгоритм 1 Подсчет нормировочной константы Z

Вход: факторы Ψ_1, \dots, Ψ_m , точность округления ε

Выход: оценка нормировочной константы \widehat{Z}

- 1: для $\ell := 1$ до m
 - 2: Найти ГТ-ядра $\mathbf{G}_1^\ell, \dots, \mathbf{G}_n^\ell$ для тензора Ψ_ℓ
 - 3: Инициализировать $\mathbf{f}_{n+1} := \mathbf{1}$
 - 4: для $i := n$ до 1
 - 5: Инициализировать $\mathbf{B}_i := 0$
 - 6: для $x_i := 1$ до d_i
 - 7: Построить ГТ-матрицу $\mathbf{A}_i[x_i] = \bigotimes_{\ell=1}^m \mathbf{G}_i^\ell[x_i]$
 - 8: $\mathbf{B}_i := \mathbf{B}_i + \mathbf{A}_i[x_i]$
 - 9: $\overline{\mathbf{f}}_i := \mathbf{B}_i \cdot \mathbf{f}_{i+1}$
 - 10: $\mathbf{f}_i := \text{round}(\overline{\mathbf{f}}_i, \varepsilon)$
 - 11: $\widehat{Z} := \mathbf{f}_1$
-

$1, \dots, n$ можно предрассчитать за $2(n-1)$ умножений ГТ-матриц. Дополнительно рассчитав все произведения вида $\mathbf{B}_i \dots \mathbf{B}_j$, $1 \leq i < j \leq n$, можно вычислять маргинальное распределение на любое подмножество переменных.

6.2 Анализ алгоритма 1

В этом разделе предоставлены теоретические гарантии точности оценки нормировочной константы алгоритмом 1.

Обозначим оценку произведения ГТ-матриц $\{\mathbf{B}_j\}_{j=i}^n$ за \mathbf{f}_i , $\mathbf{f}_n = \mathbf{B}_n$, $\widehat{Z} = \mathbf{f}_1$. Умножая ГТ-матрицу на ГТ-вектор и применяя ГТ-округление, получаем \mathbf{f}_i

$$\mathbf{f}_i = \text{round}(\mathbf{B}_i \mathbf{f}_{i+1}, \varepsilon),$$

где точность ГТ-округления контролирует относительную точность по евклидовой норме⁴:

$$\|\mathbf{B}_i \mathbf{f}_{i+1} - \mathbf{f}_i\|_2 \leq \varepsilon \|\mathbf{B}_i \mathbf{f}_{i+1}\|_2. \quad (11)$$

⁴Алгоритм ГТ-округления контролирует относительную точность тензоров по норме Фробениуса, но для векторов норма Фробениуса совпадает с L_2 -нормой.

Основной результат представлен в теореме 2, затем приведено следствие, которое легче интерпретировать.

Теорема 2. Для любого набора факторов Ψ_1, \dots, Ψ_m и любого значения параметра точности округления $\varepsilon \geq 0$ абсолютная ошибка оценки нормировочной константы, посчитанной алгоритмом 1, не превосходит:

$$\begin{aligned} |Z - \hat{Z}| &\leq \|\mathbf{B}_1\|_2 \cdots \|\mathbf{B}_{n-2}\|_2 \cdot \|\mathbf{B}_{n-1} \mathbf{f}_n - \mathbf{f}_{n-1}\|_2 + \\ &+ \|\mathbf{B}_1\|_2 \cdots \|\mathbf{B}_{n-3}\|_2 \cdot \|\mathbf{B}_{n-2} \mathbf{f}_{n-1} - \mathbf{f}_{n-2}\|_2 + \dots + \\ &+ \|\mathbf{B}_1 \mathbf{f}_2 - \mathbf{f}_1\|_2 \end{aligned} \quad (12)$$

Следствие 1. Для любого набора факторов Ψ_1, \dots, Ψ_m и любого значения параметра точности округления $\varepsilon \geq 0$ абсолютная ошибка оценки нормировочной константы, посчитанной алгоритмом 1, не превосходит:

$$|Z - \hat{Z}| \leq \|\mathbf{B}_1\|_2 \cdots \|\mathbf{B}_n\|_2 ((1 + \varepsilon)^{n-1} - 1) \quad (13)$$

Оценка из следствия менее точная, но позволяет по требуемой итоговой точности выбрать достаточный ε .

Чтобы пользоваться результатами теоремы 2, необходимо вычислять 2-норму векторов и матриц в ТТ-формате. 2-норма вектора совпадает с нормой Фробениуса соответствующего тензора, поэтому значения $\|\mathbf{B}_i \mathbf{f}_{i+1} - \mathbf{f}_i\|_2$ легко вычисляются. Хотя подсчет 2-нормы матриц в ТТ-формате является вычислительно сложной задачей, 2-норму можно оценить сверху с помощью нормы Фробениуса или с помощью эмпирически более точной оценки, использующей структуру матрицы \mathbf{B}_i :

$$\begin{aligned} \|\mathbf{B}_i\|_2 &= \left\| \sum_{x_i} \mathbf{G}_i^1[x_i] \otimes \dots \otimes \mathbf{G}_i^m[x_i] \right\|_2 \leq \sum_{x_i} \|\mathbf{G}_i^1[x_i] \otimes \dots \otimes \mathbf{G}_i^m[x_i]\|_2 = \\ &= \sum_{x_i} \|\mathbf{G}_i^1[x_i]\|_2 \cdots \|\mathbf{G}_i^m[x_i]\|_2 = U_i. \end{aligned} \quad (14)$$

Здесь используется равенство: $\|\mathbf{A} \otimes \mathbf{B}\|_2 = \|\mathbf{A}\|_2 \|\mathbf{B}\|_2$.

На рис. 3 представлены результаты сравнения величины нормы Фробениуса $\|\mathbf{B}_i\|_F$, спектральной нормы $\|\mathbf{B}_i\|_2$ и её верхней оценки U_i . Значения разных норм указаны для всех индексов $i = 1, \dots, n$ фиксированной модели Изинга.

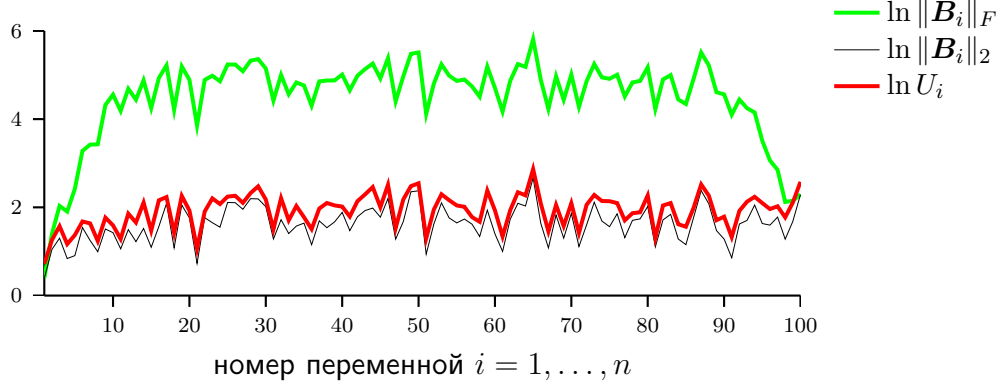


Рис. 3: Сравнение Фробениусовой и спектральной нормы матрицы \mathbf{B}_i с верхней оценкой U_i . График построен для модели Изинга с решеткой размера 10×10 , в которой коэффициенты унарных и парных потенциалов сгенерированы из равномерного распределения на $[-1, 1]$, а температура равна 1.

6.3 Доказательство теоремы 2

Лемма 1. В условиях теоремы 2 для всех индексов $i = 1, \dots, n-1$ верно следующее неравенство:

$$\begin{aligned} \|\mathbf{B}_i \dots \mathbf{B}_n - \mathbf{f}_i\|_2 &\leq \|\mathbf{B}_i\|_2 \dots \|\mathbf{B}_{n-2}\|_2 \cdot \|\mathbf{B}_{n-1}\mathbf{f}_n - \mathbf{f}_{n-1}\|_2 + \\ &+ \|\mathbf{B}_i\|_2 \dots \|\mathbf{B}_{n-3}\|_2 \cdot \|\mathbf{B}_{n-2}\mathbf{f}_{n-1} - \mathbf{f}_{n-2}\|_2 + \dots + \|\mathbf{B}_i\mathbf{f}_{i+1} - \mathbf{f}_i\|_2. \end{aligned} \quad (15)$$

Доказательство. Проведем доказательство по индукции.

В качестве базы индукции рассмотрим $i = n-1$:

$$\|\mathbf{B}_{n-1}\mathbf{B}_n - \mathbf{f}_{n-1}\|_2 = \|\mathbf{B}_{n-1}\mathbf{f}_n - \mathbf{f}_{n-1}\|_2.$$

(т.к. $\mathbf{f}_n = \mathbf{B}_n$ по построению.)

Предположим теперь, что (15) верно $\forall i = j+1, \dots, n-1$. Тогда для $i = j$ получаем

$$\begin{aligned} \|\mathbf{B}_j \dots \mathbf{B}_n - \mathbf{f}_j\|_2 &= \\ &= \|(\mathbf{B}_j \dots \mathbf{B}_n - \mathbf{B}_j\mathbf{f}_{j+1}) + (\mathbf{B}_j\mathbf{f}_{j+1} - \mathbf{f}_j)\|_2 \leq \\ &\leq \|\mathbf{B}_j\|_2 \|\mathbf{B}_{j+1} \dots \mathbf{B}_n - \mathbf{f}_{j+1}\|_2 + \|\mathbf{B}_j\mathbf{f}_{j+1} - \mathbf{f}_j\|_2 \leq \\ &\leq \|\mathbf{B}_j\|_2 (\|\mathbf{B}_{j+1}\|_2 \dots \|\mathbf{B}_{n-2}\|_2 \cdot \|\mathbf{B}_{n-1}\mathbf{f}_n - \mathbf{f}_{n-1}\|_2 + \\ &+ \|\mathbf{B}_{j+1}\|_2 \dots \|\mathbf{B}_{n-3}\|_2 \cdot \|\mathbf{B}_{n-2}\mathbf{f}_{n-1} - \mathbf{f}_{n-2}\|_2 + \\ &+ \dots + \|\mathbf{B}_{j+1}\mathbf{f}_{j+2} - \mathbf{f}_{j+1}\|_2) + \|\mathbf{B}_j\mathbf{f}_{j+1} - \mathbf{f}_j\|_2. \end{aligned}$$

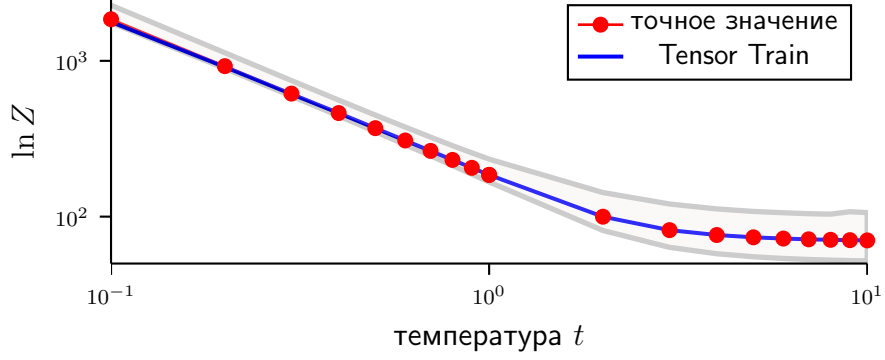


Рис. 4: Доверительный интервал на значение логарифма нормировочной константы $\ln Z$, полученный из теоремы 2 и неравенства (14). Детали см. в разделе 10.1.

Таким образом, (15) верно для $i = j$. □

Доказательство теоремы 2. По построению, $Z = \mathbf{B}_1 \dots \mathbf{B}_n$, а $\widehat{Z} = \mathbf{f}_1$. Таким образом, выполнены равенства

$$|Z - \widehat{Z}| = |\mathbf{B}_1 \dots \mathbf{B}_n - \mathbf{f}_1| = \|\mathbf{B}_1 \dots \mathbf{B}_n - \mathbf{f}_1\|_2.$$

Здесь используется тот факт, что $\mathbf{B}_1 \dots \mathbf{B}_n$ и \mathbf{f}_1 — числа, а для чисел модуль совпадает с L_2 -нормой одноэлементного вектора. Для завершения доказательства осталось применить лемму 1 к полученному выражению. □

6.4 Доказательство следствия 1

Лемма 2. В условиях следствия 1 для всех индексов $i = 1, \dots, n$ верно следующее неравенство:

$$\|\mathbf{f}_i\|_2 \leq \|\mathbf{B}_i\|_2 \dots \|\mathbf{B}_n\|_2 (1 + \varepsilon)^{n-i}. \quad (16)$$

Доказательство. Проведем доказательство по индукции.

Для $i = n$ утверждение следует из определения \mathbf{f}_n .

Пусть (16) верно $\forall j = i + 1, \dots, n$. Тогда из (11) следует, что

$$\begin{aligned} \|\mathbf{f}_j\|_2 &= \|\mathbf{f}_j - \mathbf{B}_j \mathbf{f}_{j+1} + \mathbf{B}_j \mathbf{f}_{j+1}\|_2 \leq \\ &\leq \varepsilon \|\mathbf{B}_j\|_2 \|\mathbf{f}_{j+1}\|_2 + \|\mathbf{B}_j\|_2 \|\mathbf{f}_{j+1}\|_2 = \\ &= \|\mathbf{B}_j\|_2 \|\mathbf{f}_{j+1}\|_2 (1 + \varepsilon). \end{aligned}$$

Чтобы завершить доказательство, остаётся воспользоваться предположением индукции для $\|\mathbf{f}_{j+1}\|_2$. \square

Лемма 3. В условиях следствия 1 для всех индексов $i = 1, \dots, n-1$ верно следующее неравенство:

$$\|\mathbf{B}_i \mathbf{f}_{i+1} - \mathbf{f}_i\|_2 \leq \|\mathbf{B}_i\|_2 \dots \|\mathbf{B}_n\|_2 \varepsilon (1 + \varepsilon)^{n-i-1}. \quad (17)$$

Доказательство. Утверждение следует из леммы 2 и неравенства

$$\|\mathbf{B}_i \mathbf{f}_{i+1} - \mathbf{f}_i\|_2 \leq \varepsilon \|\mathbf{B}_i \mathbf{f}_{i+1}\|_2 \leq \varepsilon \|\mathbf{B}_i\|_2 \|\mathbf{f}_{i+1}\|_2,$$

которое в свою очередь следует из (11). \square

Доказательство следствия 1. Применяя лемму 3 к неравенству (12), получаем

$$\begin{aligned} |Z - \widehat{Z}| &\leq \|\mathbf{B}_1\|_2 \dots \|\mathbf{B}_n\|_2 \varepsilon + \|\mathbf{B}_1\|_2 \dots \|\mathbf{B}_n\|_2 \varepsilon (1 + \varepsilon) + \dots + \\ &+ \|\mathbf{B}_1\|_2 \dots \|\mathbf{B}_n\|_2 \varepsilon (1 + \varepsilon)^{n-2} = \|\mathbf{B}_1\|_2 \dots \|\mathbf{B}_n\|_2 \varepsilon (1 + (1 + \varepsilon) + \dots + (1 + \varepsilon)^{n-2}). \end{aligned}$$

Пользуясь формулой суммы геометрической прогрессии $1 + (1 + \varepsilon) + \dots + (1 + \varepsilon)^{n-2} = ((1 + \varepsilon)^{n-1} - 1)/\varepsilon$, получаем (13). \square

7 Метод n-окрестностей

В этом разделе будет приведен обзор метода n-окрестностей [10]. Пусть задана модель Изинга, т. е. графическая модель, чья функция энергии следующим образом выражается через параметры $\mathbf{T} \in \mathbb{R}^{n \times n}$ и $\mathbf{h} \in \mathbb{R}^{n \times 1}$:

$$E(\mathbf{x}) = -\frac{1}{2} \sum_{i,j} T_{ij} x_i x_j - \sum_{i=1}^n h_i x_i, \quad x_i \in \{-1, 1\} \quad \forall i = 1 \dots n, \quad (18)$$

где \mathbf{x} – это вектор переменных модели, \mathbf{T} – симметричная матрица параметров с нулями на диагонали и \mathbf{h} – n -мерный вектор параметров. Положим \mathbf{x}_0 – конфигурацию минимума энергии, т. е. $\mathbf{x}_0 = \operatorname{argmin} E(\mathbf{x})$. Запишем определение нормировочной константы:

$$Z = \sum_{\mathbf{x}} \exp(-E(\mathbf{x})) = \sum_{k=0}^n |\mathcal{X}_{\text{layer}}^k| \sum_{\mathbf{x} \in \mathcal{X}_{\text{layer}}^k} \frac{1}{|\mathcal{X}_{\text{layer}}^k|} \exp(-E(\mathbf{x})),$$

где $\mathcal{X}_{\text{layer}}^k$ – это множество, состоящее из всех наборов \mathbf{x} , которые отличаются от \mathbf{x}_0 ровно в k позициях. Множество $\mathcal{X}_{\text{layer}}^k$ называется k -ым слоем. Легко показать, что мощность k -ого слоя равна биномиальному коэффициенту: $|\mathcal{X}_{\text{layer}}^k| = \binom{n}{k}$.

Рассмотрим равномерное распределение вероятности на множестве $\mathcal{X}_{\text{layer}}^k$:

$$p_k(\mathbf{x}) = \begin{cases} \frac{1}{|\mathcal{X}_{\text{layer}}^k|}, & \text{если } \mathbf{x} \in \mathcal{X}_{\text{layer}}^k \\ 0, & \text{иначе.} \end{cases}$$

Тогда:

$$\mathbb{E}_{p_k(\mathbf{x})} \exp(-E(\mathbf{x})) = \sum_{\mathbf{x} \in \mathcal{X}_{\text{layer}}^k} \frac{1}{|\mathcal{X}_{\text{layer}}^k|} \exp(-E(\mathbf{x})) \quad (19)$$

Обращаем внимание, что компоненты случайной величины $\mathbf{x} \sim p_k$ не являются независимыми (например последняя компонента вектора $\mathbf{x} \in \mathcal{X}_{\text{layer}}^k$ однозначно определяется по первым $n - 1$).

Рассмотрим теперь распределение на значения энергии: $q_k(e) = \frac{|\{\mathbf{x} \in \mathcal{X}_{\text{layer}}^k, E(\mathbf{x})=e\}|}{|\mathcal{X}_{\text{layer}}^k|}$.

Тогда:

$$\mathbb{E}_{q_k(E)} \exp(-E) = \sum_{E \in \{E(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_{\text{layer}}^k\}} q_k(E) \exp(-E) = \sum_{\mathbf{x} \in \mathcal{X}_{\text{layer}}^k} \frac{1}{|\mathcal{X}_{\text{layer}}^k|} \exp(-E(\mathbf{x}))$$

Случайная величина E – это сумма большого числа слагаемых вида $T_{ij}x_i x_j$, большинство из которых являются независимыми. Приближим $q_k(E)$ нормальным распределением⁵. Тогда:

$$\sum_{\mathbf{x} \in \mathcal{X}_{\text{layer}}^k} \frac{1}{|\mathcal{X}_{\text{layer}}^k|} \exp(-E(\mathbf{x})) = \sum_{E \in \{E(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_{\text{layer}}^k\}} q_k(E) \exp(-E) \approx \int_{E_{\min}}^{E_{\max}} f_k(E) \exp(-E) dE,$$

где $f_k(E)$ – плотность нормального распределения: $E \sim \mathcal{N}(E \mid \mu_k, \sigma_k^2)$, E_{\min} и E_{\max} – минимальное и максимальное значения функции энергии соответственно. Этот прием является ключевым для метода n -окрестностей.

Для поиска параметров нормального распределения $f_k(E)$, наилучшим образом приближающего распределение $q_k(E)$, приравняем моменты: $\mu_k = \mathbb{E}_{q_k} E$, $\sigma_k^2 = \mathbb{D}_{q_k} E$. Эти величины рассчитываются аналитически [10]. Пример выражения для среднего:

$$\mu_k = -\frac{\mathbf{x}_0^T \mathbf{T} \mathbf{x}_0}{2} \frac{(n - 2k)^2 - n}{n(n - 1)} - \left(1 - \frac{2k}{n}\right) \mathbf{x}_0^T \mathbf{h}$$

⁵Отметим, что $E \sim q_k$ – это дискретная случайная величина, и нормальное распределение является довольно грубым приближением для неё.

В ходе дальнейших преобразований⁶ приходим к следующему приближенному выражению для нормировочной константы:

$$\begin{aligned}
 Z &= \sum_{k=0}^n \binom{n}{k} \sum_{\mathbf{x} \in \mathcal{X}_{\text{layer}}^k} \frac{1}{|\mathcal{X}_{\text{layer}}^k|} \exp(-E(\mathbf{x})) \approx \\
 &\approx n \frac{1}{\sqrt{2\pi}} \int_0^1 \exp(nf(x)) (\Phi(B_x) - \Phi(A_x)) dx,
 \end{aligned} \tag{20}$$

где

$$\begin{aligned}
 f(x) &= -x \ln(x) - (1-x) \ln(1-x) - \frac{\mu_k}{n} + \frac{1}{2} \frac{\sigma_k^2}{n}, \\
 A_x &= \sqrt{n} \left(\frac{E_{\min} - \mu}{\sigma} + \frac{\sigma}{n} \right), \\
 B_x &= \sqrt{n} \left(\frac{E_{\max} - \mu}{\sigma} + \frac{\sigma}{n} \right).
 \end{aligned}$$

7.1 Ограничения подхода

Рассмотрим особенности метода, ограничивающие его применение в области графических моделей:

1. Ограниченность модели. В большинстве прикладных задач возникающие марковские случайные поля существенно выходят за рамки модели Изинга (многоклассовые модели, присутствие потенциалов высоких порядков, модели RBM, и т. п.).
2. Нарушение предположений метода. Т.к. энергия – это сумма *зависимых* слагаемых, то предположение о нормальности возникающего распределения нарушается на практике.
3. Скорость работы. Необходимость поиска минимума и максимума энергии, а так же численного взятия интеграла (20) ограничивает применимость метода к обучению MRF.

В разделе 8 предлагается изменение метода, направленное на устранение описанных недостатков.

⁶Основными преобразованиями являются: применение формулы Стирлинга, приближение внешней суммы интегралом и выделение полного квадрата под экспонентой.

8 Метод «Суммировать всё»

В этом разделе описывается предлагаемый метод, основанный на методе p -окрестностей и получивший название «Суммировать всё».

Рассмотрим произвольное марковское случайное поле. Предположим, что эмпирическая плотность распределения энергий E близка к нормальному распределению:

$$\forall y \quad P(E < y) = \frac{\sum_{\mathbf{x}} [E(\mathbf{x}) < y]}{|\mathcal{X}|} \approx \int_{-\infty}^y \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(t-\mu)^2}{2\sigma^2}\right)}_{f(t | \mu, \sigma^2)} dt,$$

где $|\mathcal{X}|$ – это мощность совместной области определения всех переменных x_1, \dots, x_n : $|\mathcal{X}| = \prod_{i=1}^n d_i$.

Заменяем сумму в определении нормировочной константы на интеграл по все числовой оси:

$$\ln Z \approx \ln \left(|\mathcal{X}| \int_{-\infty}^{\infty} f(E | \mu, \sigma^2) \exp(-E) dE \right) = -\mu + \frac{\sigma^2}{2} + \ln(|\mathcal{X}|). \quad (21)$$

Найдём аналитические выражения для μ и σ^2 , рассчитанные по методу максимального правдоподобия для множества всех энергий $\{E(\mathbf{x}) | \mathbf{x} \in \mathcal{X}\}$:

$$\begin{aligned} \mu &= \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}} E(\mathbf{x}) = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}} \sum_{\ell} \Theta_{\ell}(\mathbf{x}^{\ell}) \\ &= \frac{1}{|\mathcal{X}|} \sum_{\ell} \sum_{\mathbf{x}} \Theta_{\ell}(\mathbf{x}^{\ell}) = \frac{1}{|\mathcal{X}|} \sum_{\ell} \frac{|\mathcal{X}|}{|\mathcal{X}^{\ell}|} \sum_{\mathbf{x}^{\ell}} \Theta_{\ell}(\mathbf{x}^{\ell}) \\ &= \sum_{\ell} \frac{1}{|\mathcal{X}^{\ell}|} \sum_{\mathbf{x}^{\ell}} \Theta_{\ell}(\mathbf{x}^{\ell}), \end{aligned} \quad (22)$$

где $|\mathcal{X}^{\ell}|$ – это мощность совместной области определения тех переменных x_i , от которых существенно зависит ℓ -ый потенциал.

Т.е. μ – это по сути взвешенная сумма всех значений всех потенциалов. Асимптотическая сложность расчёта данной величины линейна по количеству параметров марковского случайного поля.

Получим теперь формулу для параметра дисперсии σ^2 :

$$\begin{aligned} \sigma^2 &= \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}} (E(\mathbf{x}))^2 - \mu^2 = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}} \left(\sum_{\ell} \Theta_{\ell}(\mathbf{x}^{\ell}) \right)^2 - \mu^2 \\ &= \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}} \sum_{\ell_1, \ell_2} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) - \mu^2 \\ &= \frac{1}{|\mathcal{X}|} \sum_{\ell_1, \ell_2} \sum_{\mathbf{x}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) - \sum_{\ell_1, \ell_2} \left(\frac{1}{|\mathcal{X}^{\ell_1}|} \sum_{\mathbf{x}^{\ell_1}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \right) \left(\frac{1}{|\mathcal{X}^{\ell_2}|} \sum_{\mathbf{x}^{\ell_2}} \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) \right) \end{aligned}$$

Для любой пары потенциалов ℓ_1, ℓ_2 введем следующие обозначения:

$$\begin{aligned}\mathbf{x}_{1,2}^{\ell_1} &= \mathbf{x}^{\ell_1} \setminus \mathbf{x}^{\ell_2}, \\ \mathbf{x}_{1,2}^{\ell_2} &= \mathbf{x}^{\ell_2} \setminus \mathbf{x}^{\ell_1}, \\ \mathbf{x}_{1,2}^{\cup} &= \mathbf{x}^{\ell_2} \cup \mathbf{x}^{\ell_1}, \\ \mathbf{x}_{1,2}^{\cap} &= \mathbf{x}^{\ell_2} \cap \mathbf{x}^{\ell_1}.\end{aligned}$$

Так же обозначим совместную область определения переменных $\mathbf{x}_{1,2}^{\cup}$ за $\mathcal{X}_{1,2}^{\cup}$, а совместную область определения переменных $\mathbf{x}_{1,2}^{\cap}$ за $\mathcal{X}_{1,2}^{\cap}$. Введём дополнительное обозначение для суммы всех значений потенциала: $\tau^{\ell} = \sum_{\mathbf{x}^{\ell}} \Theta_{\ell}(\mathbf{x}^{\ell})$.

Тогда

$$\sigma^2 = \sum_{\ell_1, \ell_2} \left(\frac{1}{|\mathcal{X}_{1,2}^{\cup}|} \sum_{\mathbf{x}_{1,2}^{\ell_1}, \mathbf{x}_{1,2}^{\ell_2}, \mathbf{x}_{1,2}^{\cap}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) - \frac{\tau^{\ell_1}}{|\mathcal{X}^{\ell_1}|} \frac{\tau^{\ell_2}}{|\mathcal{X}^{\ell_2}|} \right). \quad (23)$$

Рассмотрим пару потенциалов ℓ_1, ℓ_2 таких, что они существенно зависят от не пересекающихся групп переменных: $\mathbf{x}_{1,2}^{\cap} = \emptyset$. Тогда $|\mathcal{X}_{1,2}^{\cup}| = |\mathcal{X}^{\ell_1}| \cdot |\mathcal{X}^{\ell_2}|$, и

$$\sum_{\mathbf{x}_{1,2}^{\ell_1}, \mathbf{x}_{1,2}^{\ell_2}, \mathbf{x}_{1,2}^{\cap}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) = \left(\sum_{\mathbf{x}^{\ell_1}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \right) \left(\sum_{\mathbf{x}^{\ell_2}} \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) \right),$$

т. е. соответствующее слагаемое обнуляется. Таким образом, можно переписать выражение (23) добавив туда условие $\mathbf{x}_{1,2}^{\cap} \neq \emptyset$:

$$\begin{aligned}\sigma^2 &= \sum_{\substack{\ell_1, \ell_2: \\ \mathbf{x}_{1,2}^{\cap} \neq \emptyset}} \left(\frac{1}{|\mathcal{X}_{1,2}^{\cup}|} \sum_{\mathbf{x}_{1,2}^{\ell_1}, \mathbf{x}_{1,2}^{\ell_2}, \mathbf{x}_{1,2}^{\cap}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) - \frac{\tau^{\ell_1}}{|\mathcal{X}^{\ell_1}|} \frac{\tau^{\ell_2}}{|\mathcal{X}^{\ell_2}|} \right) \\ &= \sum_{\substack{\ell_1, \ell_2: \\ \mathbf{x}_{1,2}^{\cap} \neq \emptyset}} \frac{1}{|\mathcal{X}_{1,2}^{\cup}|} \left(\gamma^{\ell_1, \ell_2} - \frac{1}{|\mathcal{X}_{1,2}^{\cap}|} \tau^{\ell_1} \tau^{\ell_2} \right),\end{aligned} \quad (24)$$

где

$$\gamma^{\ell_1, \ell_2} = \sum_{\mathbf{x}_{1,2}^{\cap}} \left(\sum_{\mathbf{x}_{1,2}^{\ell_1}} \Theta_{\ell_1}(\mathbf{x}^{\ell_1}) \right) \left(\sum_{\mathbf{x}_{1,2}^{\ell_2}} \Theta_{\ell_2}(\mathbf{x}^{\ell_2}) \right).$$

Асимптотическая сложность вычисления σ^2 равна $O(\sum_{\ell=1}^m |\mathcal{X}^{\ell}|^2)$. В случае, когда все переменные принимают d значений ($d_i = d, i = 1, \dots, n$), а все потенциалы существенно зависят от p переменных (в реальных задачах чаще всего $p \leq 3$), то оценка упрощается до $O(md^{2p})$.

Замечание: чтобы получить несмещенную выборочную оценку дисперсии, нужно поделить сумму квадратов отклонений на размер выборки минус один: $\overline{\sigma^2} = \frac{|\mathcal{X}|}{|\mathcal{X}|-1}\sigma^2$.

Отметим, что в отличие от большинства методов оценки нормировочной константы, точность данного метода возрастает при увеличении размера модели, т. к. увеличивается размер выборки $|\{E(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}\}|$, по которой вычисляются параметры нормального распределения.

8.1 Эффективные формулы для модели Изинга

Выведем формулы для μ и σ^2 для частного случая – когда рассматриваемое марковское случайное поле является моделью Изинга (см. раздел 7). В этом случае каждая переменная принимает два возможных значения $x_i \in \{-1, 1\}$, $i = 1, \dots, n$, а MRF состоит из двух видов потенциалов: унарных $\Theta_i(x_i) = -h_i x_i$, $i = 1, \dots, n$ и парных $\Theta_{i,j}(x_i, x_j) = -\frac{1}{2}T_{ij}x_i x_j$, $i, j = 1, \dots, n$. Среднее значение всех потенциалов равно нулю:

$$\begin{aligned}\tau^i &= \sum_{x_i} \Theta_i(x_i) = -h_i + h_i = 0, \\ \tau^{i,j} &= \sum_{x_i, x_j} \Theta_{i,j}(x_i, x_j) = -2T_{ij} + 2T_{ij} = 0.\end{aligned}$$

Подставляя средние значения потенциалов в формулу (22) получим, что $\mu = 0$.

Для подсчёта дисперсии воспользуемся формулой (24). Рассмотрим отдельно все пары потенциалов зависящие от пересекающихся подмножеств переменных.

Рассмотрим парный потенциал $\Theta_{i,j}(x_i, x_j)$ и унарный потенциал $\Theta_i(x_i)$. Для них множество общих переменных $\mathbf{x}_{1,2}^\cap = \{x_i\}$, множество $\mathbf{x}_{1,2}^{\ell_1} = \{x_j\}$, множество $\mathbf{x}_{1,2}^{\ell_2} = \emptyset$, мощность совместной области определения $|\mathcal{X}_{1,2}^\cup|$ равна 4.

$$\gamma^{(ij),i} = \sum_{x_i} \underbrace{\left(\sum_{x_j} \frac{1}{2} x_i x_j T_{ij} \right)}_0 x_i h_i = 0.$$

Рассмотрим парные потенциалы $\Theta_{i,j}(x_i, x_j)$ и $\Theta_{i,k}(x_i, x_k)$, где $j \neq k$. Для них множество общих переменных $\mathbf{x}_{1,2}^\cap = \{x_i\}$, множество $\mathbf{x}_{1,2}^{\ell_1} = \{x_j\}$, множество $\mathbf{x}_{1,2}^{\ell_2} = \{x_k\}$.

$$\gamma^{(ij),(ik)} = \sum_{x_i} \underbrace{\left(\sum_{x_j} -\frac{1}{2} x_i x_j T_{ij} \right)}_0 \underbrace{\left(\sum_{x_k} -\frac{1}{2} x_i x_k T_{ik} \right)}_0 = 0.$$

Рассмотрим парные потенциалы $\Theta_{i,j}(x_i, x_j)$ и $\Theta_{i,j}(x_i, x_j)$. Для них множество общих переменных $\mathbf{x}_{1,2}^\cap = \{x_i, x_j\}$, множество $\mathbf{x}_{1,2}^{\ell_1} = \emptyset$, множество $\mathbf{x}_{1,2}^{\ell_2} = \emptyset$, мощность совместной области определения $|\mathcal{X}_{1,2}^\cup|$ равна 4.

$$\gamma^{(ij),(ij)} = \sum_{x_i, x_j} \frac{1}{2} x_i x_j T_{ij} \frac{1}{2} x_i x_j T_{ij} = T_{ij}^2.$$

Аналогично рассмотрим парные потенциалы $\Theta_{i,j}(x_i, x_j)$ и $\Theta_{j,i}(x_j, x_i)$: $\gamma^{(ij),(ji)} = T_{ij}^2$.

Рассмотрим унарные потенциалы $\Theta_i(x_i)$ и $\Theta_i(x_i)$. Для них множество общих переменных $\mathbf{x}_{1,2}^\cap = \{x_i\}$, множество $\mathbf{x}_{1,2}^{\ell_1} = \emptyset$, множество $\mathbf{x}_{1,2}^{\ell_2} = \emptyset$, мощность совместной области определения $|\mathcal{X}_{1,2}^\cup|$ равна 2.

$$\gamma^{i,i} = \sum_{x_i} x_i h_i x_i h_i = 2h_i^2.$$

Подставив выражения для γ^{ℓ_1, ℓ_2} в формулу (24) получим:

$$\sigma^2 = \frac{1}{2} \sum_{i,j} T_{ij}^2 + \sum_i h_i^2 = \frac{1}{2} \text{Tr}(\mathbf{T}^\top \mathbf{T}) + \mathbf{h}^\top \mathbf{h}, \quad (25)$$

где под $\text{Tr}(\mathbf{A})$ понимается след матрицы \mathbf{A} .

8.2 Эффективные формулы для модели RBM

Модель ограниченной машины Больцмана (Restricted Boltzmann Machines, RBM) [20] – это модель с бинарными переменными $x_i \in \{0, 1\}$, которые разбиваются на две группы: $\mathbf{x} = (v_1, \dots, v_q, h_1, \dots, h_f)$, где q – это число *наблюдаемых переменных*, а f – число *латентных переменных*, $q + f = n$. Функция энергии RBM задаётся следующим образом:

$$E(\mathbf{x}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{h}^\top \mathbf{W} \mathbf{v},$$

где \mathbf{b} и \mathbf{c} – это q -мерный и f -мерный вектора параметров соответственно, а \mathbf{W} – это матрица параметров размера $f \times q$.

Применяя формулы (22) и (24) к потенциалам RBM получим:

$$\begin{aligned} \mu &= -\frac{1}{4} \mathbf{1}^\top \mathbf{W} \mathbf{1} - \frac{1}{2} \mathbf{c}^\top \mathbf{1} - \frac{1}{2} \mathbf{b}^\top \mathbf{1}, \\ \sigma^2 &= \frac{1}{4} (\mathbf{c}^\top \mathbf{W} \mathbf{1} + \mathbf{1}^\top \mathbf{W} \mathbf{b}) + \frac{1}{16} (\mathbf{1}^\top \mathbf{W} \mathbf{W}^\top \mathbf{1} + \mathbf{1}^\top \mathbf{W}^\top \mathbf{W} \mathbf{1}) + \frac{1}{16} \text{Tr}(\mathbf{W}^\top \mathbf{W}) + \frac{1}{4} (\mathbf{b}^\top \mathbf{b} + \mathbf{c}^\top \mathbf{c}). \end{aligned}$$

9 Метод, основанный на формуле Тейлора

В этом разделе предлагается метод, не делающий предположения о нормальности распределения энергий.

Заменим в определении нормировочной константы (1) показательные функции на их разложение в ряд Тейлора до второго порядка вокруг точки E_0 :

$$\begin{aligned} Z &= \sum_{\mathbf{x}} \exp(-E(\mathbf{x})) \approx \\ &\approx \sum_{\mathbf{x}} \left(\exp(-E_0) - \exp(-E_0)(E(\mathbf{x}) - E_0) + \exp(-E_0)\frac{1}{2}(E(\mathbf{x}) - E_0)^2 \right) \\ \ln Z &\approx -E_0 + \underbrace{\ln \left(\sum_{\mathbf{x}} 1 - (E(\mathbf{x}) - E_0) + \frac{1}{2}(E(\mathbf{x}) - E_0)^2 \right)}_{t(E_0)} \end{aligned}$$

Воспользовавшись определением μ (22) и σ^2 (24) можно упростить данное выражение:

$$\ln Z \approx t(E_0) = -E_0 + \ln \left(1 - \mu(1 + E_0) + \frac{1}{2}\mu^2 + \frac{1}{2}\sigma^2 + E_0 + \frac{1}{2}E_0^2 \right) + \ln(|\mathcal{X}|)$$

Отметим, что оценка логарифма нормировочной константы по метод «Суммировать всё» (21) является оценкой $t(\mu)$ через разложение в ряд Тейлора:

$$t(\mu) = -\mu + \ln \left(1 + \frac{1}{2}\sigma^2 \right) + \ln(|\mathcal{X}|) \approx \underbrace{-\mu + \frac{1}{2}\sigma^2 + \ln(|\mathcal{X}|)}_{\text{оценка по методу «Суммировать всё»}}$$

10 Вычислительные эксперименты

В этом разделе проводится экспериментальное сравнение метода n-окрестностей, предлагаемых методов «Суммировать всё», «Тейлор» и «Tensor Train», и альтернатив из библиотеки LibDAI [14]: Belief Propagation (BP) [11] и метод Mean Field (MF) [24]. Так же проводится сравнение с методом Annealed Importance Sampling (AIS) [15] — представителем методов MCMC.

Во всех экспериментах используется неоптимизированная реализация предложенных методов в среде Matlab⁷. Для операций связанных с ТТ-форматом используется ТТ-Toolbox⁸, реализованный в среде Matlab.

⁷<https://github.com/bihaqo/TT-MRF>

⁸<http://spring.inm.ras.ru/osel/download/tt22.zip>

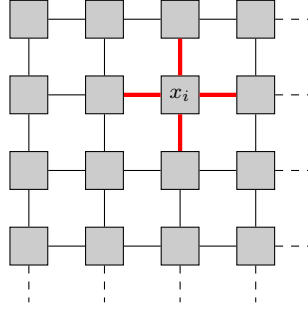


Рис. 5: 4-х связная решетка, задающая систему соседства: $(i, j) \in \mathcal{E}$ тогда и только тогда, когда переменные x_i и x_j соединены ребром.

10.1 Модель Изинга

Добавим параметр *температуры* $t \in \mathbb{R}_+$ в модель Изинга (см. раздел 7):

$$E(\mathbf{x}) = \frac{1}{t} \left(-\frac{1}{2} \sum_{i,j} T_{ij} x_i x_j - \sum_{i=1}^N h_i x_i \right).$$

Для каждого значения температуры t сгенерировано 50 моделей Изинга размера 10×10 со унарными весами сгенерированными из равномерного распределения $h_i \sim U[-1, 1]$, $i = 1, \dots, n$ и парными весами равными 1 на 4-х связной решетке (рис. 5):

$$T_{ij} = \begin{cases} 1, & \text{если } (i, j) \in \mathcal{E} \\ 0, & \text{иначе.} \end{cases}$$

Размер 10×10 выбран достаточно большим, чтобы расчёт нормировочной константы по определению был невыполнимым (т.к. требует суммирования по $2^{10 \times 10}$ слагаемым), но достаточно маленьким для того, чтобы можно было посчитать ошибку методов с помощью метода Junction Tree [24]. На рис. 6 указана медиана, а так же верхняя и нижняя квартили абсолютной ошибки оценки логарифма нормировочной константы. Медианное время работы методов указано в таблице 2.

На рис. 4 изображен доверительный интервал на значение логарифма нормировочной константы, полученный путем оценивания результата теоремы 2 неравенством (14).

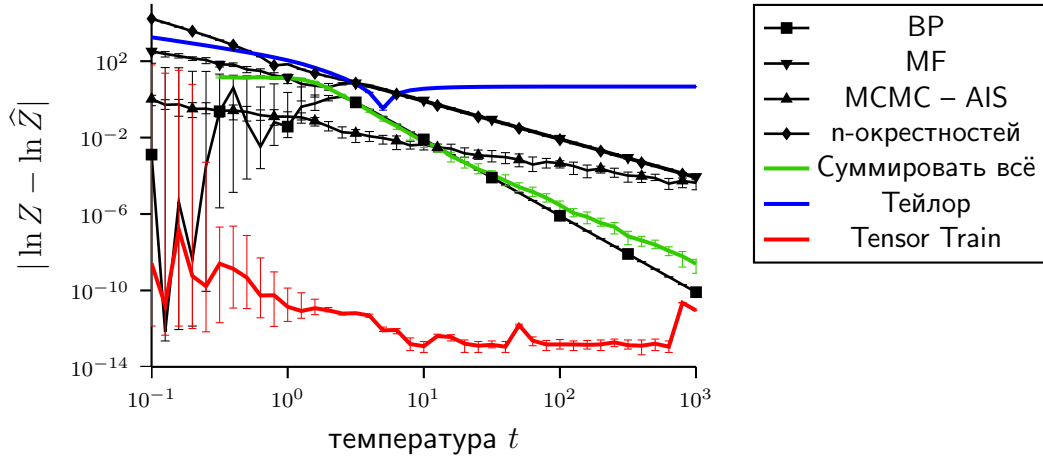


Рис. 6: Абсолютная ошибка подсчета логарифма нормировочной константы в зависимости от температуры (чем меньше, тем лучше).

10.2 ТТ-ранги энергии и вероятности

В данном эксперименте проверяется утверждение из разделов 5.1 и 5.2: при увеличении количества переменных ТТ-ранги тензора энергии растут медленно, а ТТ-ранги тензора ненормированной вероятности — экспоненциально быстро (рис. 2). Рассматривается последовательность моделей Изинга с 4-х связной решеткой возрастающих размеров: от 1×1 до 12×12 . Унарные веса h_i сгенерированы из равномерного распределения $U[-1, 1]$, все парные веса T_{ij} равны 1, температура t равна 10. Для моделей Изинга указан максимальный ТТ-ранг энергии и ненормированной вероятности при точном ТТ-представлении и после ТТ-округления с точностью 10^{-8} . Для точного ТТ-представления тензора вероятности модели Изинга с решеткой размера 11×11 не хватает 8GB оперативной памяти.

10.3 Обучение модели RBM

Исследуем точность рассматриваемых оценок логарифма нормировочной константы в процессе обучения модели RBM (см. раздел 8.2). В качестве обучающей выборки использовалась база данных MNIST [12] — черно-белые изображения рукописных цифр размера 28×28 . Из 784 признаков (пикселей) были выбраны 250 с

Метод	Изинг, время (с)	RBM, время (с)
BP	0.005	509
MF	0.004	0.1147
MCMC – AIS	60	–
n-окрестностей	0.1505	–
«Суммировать всё»	0.00043	0.0041
«Тейлор»	0.0005	0.0047
«Tensor Train»	32	596

Таблица 2: Медианное время работы методов оценки нормировочной константы модели Изинга размера 10×10 и модели RBM с 250 наблюдаемыми и 12 латентными переменными.

наибольшей дисперсией ⁹. Рассматривалась модель RBM с 12 латентными переменными (число 12 выбрано достаточно маленьким, чтобы был применим метод Junction Tree). Для обучения по методу максимального правдоподобия использовался метод наискорейшего подъёма:

$$\ln L(\mathbf{w}) = \sum_{\ell} \Theta_{\ell}(\mathbf{T}_{\text{train}}; \mathbf{X}_{\text{train}}, \mathbf{w}) - \ln Z(\mathbf{X}_{\text{train}}, \mathbf{w})$$

$$\alpha_e = \arg \max_{\alpha} \ln L(\mathbf{w}^e + \alpha \nabla \ln L(\mathbf{w}^e))$$

$$\mathbf{w}^{e+1} = \mathbf{w}^e + \alpha_e \nabla \ln L(\mathbf{w}^e)$$

Величина $\sum_{\ell} \Theta_{\ell}(\mathbf{T}_{\text{train}}; \mathbf{X}_{\text{train}}, \mathbf{w})$ и её градиент могут быть эффективно рассчитаны по определению, а для подсчёта логарифма нормировочной константы $\ln Z$ и её градиента использовался (точный) метод Junction Tree.

Для каждой из моделей RBM, заданных параметрами \mathbf{w}^e , $e = 1 \dots Q$, вычислены оценки логарифма правдоподобия $\ln \hat{L}(\mathbf{w}^e)$ и приведены на рис. 7 и 8. Медианное время работы методов приведено в таблице 2.

⁹Признаки, равные 0 (1) на большинстве объектов обучающей выборки, обладают априорной вероятностью близкой к 1 и вызывают численные неустойчивости.

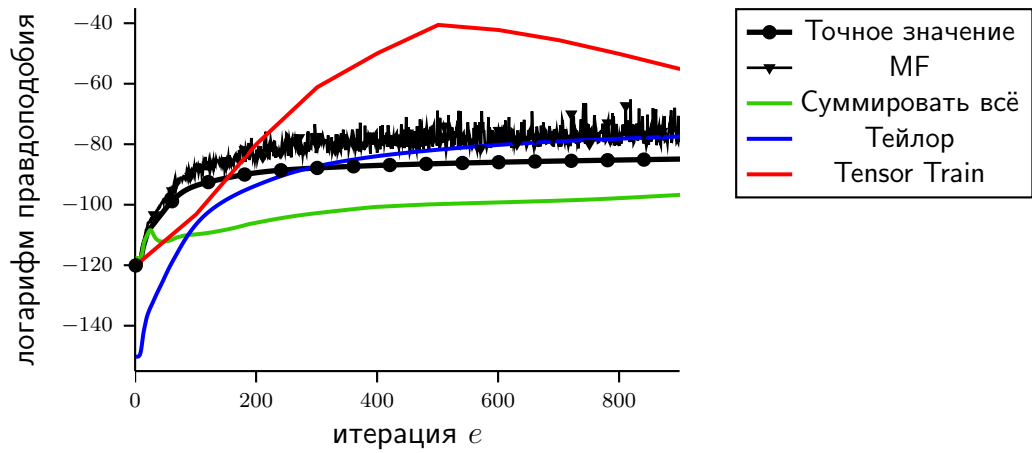


Рис. 7: Логарифм правдоподобия и его оценки на различных итерациях обучения модели RBM.

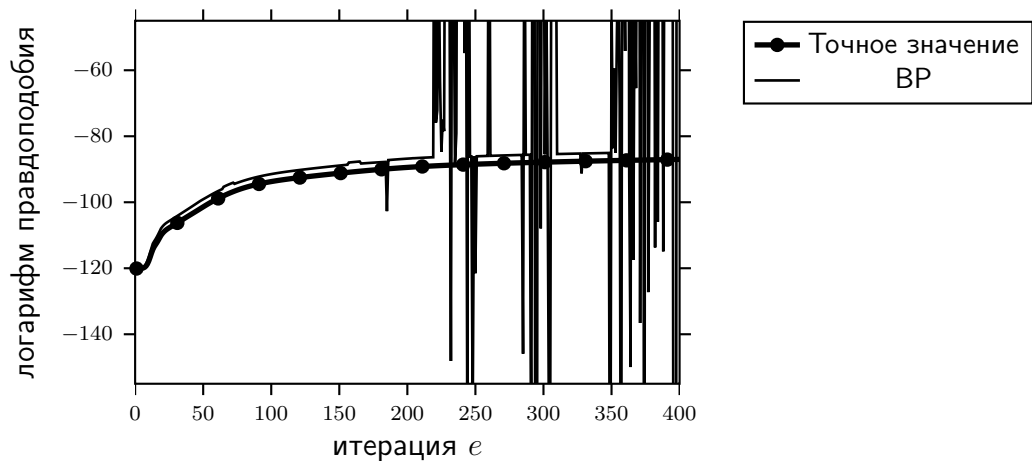


Рис. 8: Логарифм правдоподобия и его оценка методом Belief Propagation на различных итерациях обучения модели RBM.

10.4 Обсуждение и выводы

Предложенные методы «Суммировать всё» и метод, основанный на разложении в ряд Тейлора, на порядки опережают все аналоги по скорости работы (табл. 2). При этом метод «Суммировать всё» работает точнее метода n -окрестностей (рис. 6). В эксперименте по оценке логарифма нормировочной константы модели Изинга со случайными унарными потенциалами, метод, основанный на разложении Tensor Train, на порядки опережает все аналоги по точности работы (рис. 6). На первых итерациях обучения модели RBM (пока веса модели близки к случайной инициализации) метод «Tensor Train» так же показывает высокое качество работы, но по мере обу-

чения точность падает. В эксперименте по обучению RBM из всех рассмотренных методов только методы Mean Field и «Тейлор» показывают качественно верный результат (т. е. форма графика оценки повторяет форму графика точного логарифма правдоподобия), при этом метод Mean Field значительно менее устойчив по сравнению с предложенными методами (рис. 7). Метод Belief Propagation по мере обучения расходится (рис. 8).

11 Заключение

В настоящей дипломной работе получены следующие основные результаты:

- Предложен алгоритм оценки нормировочной константы через ТТ-разложение факторов графической модели, который не строит ТТ-представление тензора ненормированного совместного распределения;
- Приведены теоретические гарантии на точность оценки нормировочной константы;
- Найден практически важный класс задач, на котором тензорный метод опережает все доступные аналоги по точности работы;
- Предложено обобщение метода n -окрестностей, применимое к произвольным марковским случайным полям и обладающее более высокой скоростью работы;
- Предложен простой и быстрый метод оценки нормировочной константы через разложение в ряд Тейлора, который даёт качественно верную оценку на задаче обучения RBM.

По теме работы сделано 3 доклада на международных конференциях [17, 16, 19] и опубликована статья в журнале ВАК [25].

Список литературы

- [1] J. D. Carroll and J. J. Chang. Analysis of individual differences in multidimensional scaling via n-way generalization of Eckart-Young decomposition. *Psychometrika*, 35:283–319, 1970.
- [2] S. V. Dolgov and D. V. Savostyanov. Alternating minimal energy methods for linear systems in higher dimensions. Part II: Faster algorithm and application to nonsymmetric systems. arXiv preprint 1304.1222, 2013.
- [3] L. Grasedyck. Hierarchical singular value decomposition of tensors. *SIAM J. Matrix Anal. Appl.*, 31:2029–2054, 2010.
- [4] R. Grosse, C. Maddison, and R. Salakhutdinov. Annealing between distributions by averaging moments. In *Neural Information Processing Systems (NIPS)*, 2013.
- [5] W. Hackbusch and S. Kühn. A new scheme for the tensor representation. *J. Fourier Anal. Appl.*, 15:706–722, 2009.
- [6] T. Hazan, S. Maji, and T. Jaakkola. On sampling from the Gibbs distribution with random maximum a posteriori perturbations. In *Neural Information Processing Systems (NIPS)*, 2013.
- [7] M. Ishteva, H. Park, , and L. Song. Unfolding latent tree structures using 4th order tensors. In *International Conference on Machine Learning (ICML)*, 2013.
- [8] Y. Jernite, Y. Halpern, and D. Sontag. Discovering hidden variables in noisy-or networks using quartet tests. In *Neural Information Processing Systems (NIPS)*, 2013.
- [9] B. N. Khoromskij and I. V. Oseledets. DMRG+QTT approach to computation of the ground state for the molecular Schrödinger operator. Preprint 69, MPI MIS, Leipzig, 2010.
- [10] Boris Kryzhanovsky and Leonid Litinskii. Approximate method of free energy calculation for spin system with arbitrary connection matrix. *arXiv preprint arXiv:1410.6696*, 2014.

- [11] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Transactions on Information Theory*, 47(2):498–519, 2001.
- [12] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [13] T. Minka and Y. Qi. Tree-structured approximations by expectation propagation. In *Neural Information Processing Systems (NIPS)*, 2004.
- [14] J. M. Mooij. libDAI: A free and open source C++ library for discrete approximate inference in graphical models. *Journal of Machine Learning Research*, 11:2169–2173, August 2010.
- [15] R. Neal. Annealed importance sampling. *Statistics and Computing*, 11:125–139, 2001.
- [16] Alexander Novikov, Anton Rodomanov, Anton Osokin, and Dmitry Vetrov. Computationally efficient methods for map-inference and partition function estimation in mrf in tt format. In *SIAM Conference on Imaging Science*, 2014.
- [17] Alexander Novikov, Anton Rodomanov, Anton Osokin, and Dmitry Vetrov. Putting mrfs on a tensor train. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 811–819, 2014.
- [18] I. V. Oseledets. Tensor-Train decomposition. *SIAM J. Scientific Computing*, 33(5):2295–2317, 2011.
- [19] Anton Rodomanov, Alexander Novikov, Anton Osokin, and Dmitry Vetrov. Low-rank approximation of energies in markov random fields and their representation in tt-format. In *SIAM Conference on Imaging Science*, 2014.
- [20] Paul Smolensky. Information processing in dynamical systems: Foundations of harmony theory. 1986.
- [21] L. Song, M. Ishteva, A. Parikh, E. Xing, and H. Park. Hierarchical tensor decomposition of latent tree graphical models. In *International Conference on Machine Learning (ICML)*, 2013.

- [22] L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.
- [23] M. J. Wainwright, T. Jaakkola, and A. S. Willsky. A new class of upper bounds on the log partition function. *IEEE Transactions on Information Theory*, 51(7):2313–2335, 2005.
- [24] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.
- [25] Александр Новиков, Антон Родоманов, Антон Осокин, and Дмитрий Ветров. Тензорный поезд в марковском случайном поле. *Журнал «Интеллектуальные системы. Теория и приложения»*, pages 293–318, 2014.