

Boltzmann Machines

Дата: 5 октября 2011

Содержание

Ликбез	1
Экспоненциальное семейство распределений	1
Схема Гиббса генерации выборки из распределения	2
Вариационный подход	3
Оценка нормировочной константы распределения с помощью схемы Гиббса	4
Boltzmann Machine	6
Restricted Boltzmann Machine	8
Deep Boltzmann Machine	9
Использование Deep Boltzmann Machine	11
Виды переменных в Deep Boltzmann Machine	14
Примеры применения Deep Boltzmann Machine	15
Список литературы	17

Ликбез: экспоненциальное семейство распределений

Распределение вероятностей $p(\mathbf{x})$ принадлежит экспоненциальному семейству распределений, если оно может быть представлено как

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{u}(\mathbf{x})),$$

где $\boldsymbol{\theta}$ – набор параметров распределения, количество компонент вектора $\mathbf{u}(\mathbf{x})$ совпадает с размерностью $\boldsymbol{\theta}$, $h(\cdot)$ – некоторая функция, а $Z(\boldsymbol{\theta}) = \int h(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{u}(\mathbf{x})) d\mathbf{x}$ – нормировочная константа распределения.

Многие стандартные вероятностные распределения принадлежат экспоненциальному семейству, например, нормальное, гамма, бета, Бернулли, Дирихле и многие другие. Соответствие между параметрами этих распределений и компонентами $\boldsymbol{\theta}$, $\mathbf{u}(\mathbf{x})$ в экспоненциальном представлении показано в таблице ниже.

Распределение	Плотность	$\mathbf{u}(\mathbf{x})$	$\boldsymbol{\theta}$
Бернулли	$q^x(1-q)^{1-x}$	x	$\log \frac{q}{1-q}$
Мультиномиальное	$\prod_k \mu_k^{x_k}$	$[x_1, \dots, x_{K-1}]$	$\theta_i = \log \frac{\mu_i}{1 - \sum_j \mu_j}$
Нормальное	$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$	$[x, x^2]$	$[-\frac{1}{2\sigma}, \frac{\mu}{\sigma^2}]$
Гамма	$\frac{b^a}{\Gamma(a)} x^{a-1} \exp(-bx)$	$[\log x, x]$	$[a-1, -b]$
Бета	$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1}(1-x)^{b-1}$	$[\log(x), \log(1-x)]$	$[a-1, b-1]$
Пуассона	$\exp(-\lambda) \frac{\lambda^x}{x!}$	$[x, \log \Gamma(x+1)]$	$[k, -1]$

Легко показать, что для экспоненциального семейства градиент логарифма правдоподобия по параметрам $\boldsymbol{\theta}$ может быть вычислен как

$$\nabla_{\boldsymbol{\theta}} \log p(\hat{\mathbf{x}}|\boldsymbol{\theta}) = \mathbf{u}(\hat{\mathbf{x}}) - \mathbb{E}\mathbf{u}(\mathbf{x}),$$

где математическое ожидание берется по распределению $p(\mathbf{x}|\boldsymbol{\theta})$. Предположим далее, что у нас есть выборка $\hat{X} = \{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N\}$ независимых объектов из распределения $p(\mathbf{x}|\boldsymbol{\theta})$. Тогда усредненный градиент логарифма правдоподобия по параметрам может быть вычислен как

$$\frac{1}{N} \nabla_{\boldsymbol{\theta}} \log p(\hat{X}|\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\hat{\mathbf{x}}_n) - \mathbb{E}\mathbf{u}(\mathbf{x}) = \mathbb{E}_{data} \mathbf{u}(\mathbf{x}) - \mathbb{E}_{model} \mathbf{u}(\mathbf{x}),$$

где $p_{data}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \hat{\mathbf{x}}_n)$ – выборочная плотность распределения, а $p_{model}(\mathbf{x}) = p(\mathbf{x}|\boldsymbol{\theta})$.

Рассмотрим теперь экспоненциальное семейство для случая наблюдаемых переменных \mathbf{x} и ненаблюдаемых переменных \mathbf{t} :

$$p(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}, \mathbf{t}) \exp(\boldsymbol{\theta}^T \mathbf{u}(\mathbf{x}, \mathbf{t})).$$

Пусть имеется выборка $\hat{X} = \{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N\}$ и нас интересует градиент логарифма неполного правдоподобия по параметрам:

$$\frac{1}{N} \nabla_{\boldsymbol{\theta}} \log p(\hat{X}|\boldsymbol{\theta}) = \mathbb{E}_{data} \mathbf{u}(\mathbf{x}, \mathbf{t}) - \mathbb{E}_{model} \mathbf{u}(\mathbf{x}, \mathbf{t}), \quad (1)$$

где $p_{data}(\mathbf{t}, \mathbf{x}) = p(\mathbf{t}|\mathbf{x}, \boldsymbol{\theta}) p_{data}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N p(\mathbf{t}|\mathbf{x}, \boldsymbol{\theta}) \delta(\mathbf{x} - \hat{\mathbf{x}}_n)$, а $p_{model}(\mathbf{x}, \mathbf{t}) = p(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta})$.

Ликбез: схема Гиббса генерации выборки из распределения

Рассмотрим вероятностное распределение $p(\mathbf{t})$. Схема Гиббса позволяет сгенерировать выборку из этого многомерного распределения:

$$\mathbf{t}_1, \dots, \mathbf{t}_N \sim p(\mathbf{t}).$$

Данная выборка не является набором независимых величин, но при этом может быть использована для оценки вероятностных интегралов вида

$$\mathbb{E}_{\mathbf{t}} f(\mathbf{t}) = \int f(\mathbf{t}) p(\mathbf{t}) d\mathbf{t} \simeq \frac{1}{N} \sum_{n=1}^N f(\mathbf{t}_n). \quad (2)$$

Рассмотрим шаг генерации по схеме Гиббса. Пусть на шаге n сгенерирована конфигурация $\mathbf{t}^n = \{t_1^n, \dots, t_d^n\}$. Тогда генерация следующей точки выборки \mathbf{t}^{n+1} происходит следующим образом:

$$\begin{aligned} t_1^{n+1} &\sim p(t_1 | t_2^n, t_3^n, \dots, t_d^n), \\ t_2^{n+1} &\sim p(t_2 | t_1^{n+1}, t_3^n, t_4^n, \dots, t_d^n), \\ t_3^{n+1} &\sim p(t_3 | t_1^{n+1}, t_2^{n+1}, t_4^n, \dots, t_d^n), \\ &\dots \\ t_d^{n+1} &\sim p(t_d | t_1^{n+1}, t_2^{n+1}, \dots, t_{d-1}^{n+1}). \end{aligned} \tag{3}$$

Здесь через $p(t_i | \mathbf{t}_{\setminus i})$ обозначено условное одномерное распределение значений i -ой компоненты при условии всех остальных. Таким образом, согласно схеме Гиббса генерация выборки из многомерного распределения заменяется на итерационную генерацию точек из одномерных распределений. По аналогии с методами одномерной оптимизации генерация выборки из одномерного распределения является существенно более простой задачей, чем генерация выборки из многомерного распределения.

В схеме Гиббса (3) следует обратить внимание на следующие моменты:

- При реализации схемы Гиббса на практике часто допускается следующая ошибка: вместо шага

$$t_p^{n+1} \sim p(t_p | t_1^{n+1}, \dots, t_{p-1}^{n+1}, t_{p+1}^n, \dots, t_d^n)$$

делается шаг

$$t_p^{n+1} \sim p(t_p | t_1^n, \dots, t_{p-1}^n, t_{p+1}^n, \dots, t_d^n),$$

т.е. в условие подставляются значения компонент только с предыдущей итерации. При таком подходе интересующее нас распределение $p(\mathbf{t})$ не будет инвариантным относительно введенной марковской цепи и, следовательно, генерируемая выборка не будет выборкой из распределения $p(\mathbf{t})$.

- Порядок генерации компонент t_1, \dots, t_d может быть произвольным, но не должен меняться на разных итерациях. При изменении порядка генерации компонент на итерациях соответствующая марковская цепь перестает быть однородной, а, следовательно, необходимо отдельно доказывать эргодичность такой цепи (т.е. сходимости к единственному инвариантному распределению из любого начального распределения). В однородном случае для эргодичности достаточно, чтобы все условные вероятности были строго больше нуля.
- В схеме Гиббса необязательно рассматривать только одномерные условные распределения. Здесь достаточно разбить множество компонент вектора \mathbf{t} на набор непересекающихся подмножеств и для каждого такого подмножества уметь генерировать выборку из условного распределения $p(\mathbf{t}_i | \mathbf{t}_{\setminus i})$, где \mathbf{t}_i – множество компонент из i -го подмножества.

Ликбез: вариационный подход

Пусть имеется вероятностная модель $p(\mathbf{x}, \mathbf{t})$ с наблюдаемыми и ненаблюдаемыми переменными. Пусть также имеется некоторое произвольное распределение $q(\mathbf{t})$. Тогда:

$$\log p(\mathbf{t}) = \underbrace{\int \log \frac{p(\mathbf{x}, \mathbf{t})}{q(\mathbf{t})} q(\mathbf{t}) d\mathbf{t}}_{\mathcal{L}(q)} + \text{KL}(q || p(\mathbf{t} | \mathbf{x})). \tag{4}$$

Здесь через KL обозначена дивергенция Кульбака-Лейблера между двумя вероятностными распределениями. Эта дивергенция является мерой расстояния между двумя распределениями, т.к. она всегда неотрицательна и равна нулю тогда и только тогда, когда распределения тождественно совпадают.

На равенство (4) можно смотреть с нескольких точек зрения. Допустим, что апостериорное распределение $p(\mathbf{t}|\mathbf{x})$ не поддается вычислению (не вычисляется нормировочная константа), и мы хотим найти приближение $q(\mathbf{t})$ для распределения $p(\mathbf{t}|\mathbf{x})$. Будем искать это приближение путем минимизации KL-дивергенции между распределением $q(\mathbf{t})$ и $p(\mathbf{t}|\mathbf{x})$ в некотором семействе распределений \mathcal{Q} . Тогда из равенства выше следует, что

$$\text{KL}(q||p(\mathbf{t}|\mathbf{x})) \rightarrow \min_{q \in \mathcal{Q}} \Leftrightarrow \mathcal{L}(q) \rightarrow \max_{q \in \mathcal{Q}}.$$

Теперь задача минимизации, которая зависит от недоступного для вычисления распределения $p(\mathbf{t}|\mathbf{x})$, сведена к задаче максимизации функционала $\mathcal{L}(q)$, который зависит от известного полного совместного распределения модели $p(\mathbf{x}, \mathbf{t})$.

С другой точки зрения, нас может интересовать значение маргинального распределения $p(\mathbf{x})$ (нормировочной константы для распределения $p(\mathbf{t}|\mathbf{x})$). Так как $\text{KL}(q||p(\mathbf{t}|\mathbf{x})) \geq 0$, то значение функционала $\mathcal{L}(q)$ является нижней границей для $\log p(\mathbf{x})$. Таким образом, решая задачу максимизации функционала $\mathcal{L}(q)$ в некотором семействе распределений \mathcal{Q} , мы одновременно получаем аналитическое приближение апостериорного распределения $p(\mathbf{t}|\mathbf{x})$ и нижнюю границу для обоснованности $\log p(\mathbf{x})$.

Рассмотрим решение задачи максимизации $\mathcal{L}(q)$ в семействе т.н. факторизованных распределений:

$$q(\mathbf{t}) = \prod_{j=1}^J q_j(\mathbf{t}_j).$$

Здесь множество переменных \mathbf{t} разбито на непересекающиеся подмножества \mathbf{t}_j , причем $\cup_j \mathbf{t}_j = \mathbf{t}$, а $q_j(\mathbf{t}_j)$ – произвольное распределение в пространстве переменных \mathbf{t}_j . Таким образом, мы приходим к следующей задаче оптимизации:

$$\mathcal{L}(q) = \int \log \frac{p(\mathbf{x}, \mathbf{t})}{\prod_j q_j(\mathbf{t}_j)} \prod_j q_j(\mathbf{t}_j) d\mathbf{t}_j \rightarrow \max_{q_1, \dots, q_J}.$$

Рассмотрим решение этой задачи с помощью покоординатного подъема, т.е. зафиксируем все компоненты распределения q , кроме q_i , и рассмотрим оптимизацию $\mathcal{L}(q)$ по отдельной компоненте $q_i(\mathbf{t}_i)$. Оказывается, что решение такой (вариационной) задачи оптимизации можно получить аналитически:

$$q_i^{new}(\mathbf{t}_i) = \frac{\exp \left(\int \log p(\mathbf{x}, \mathbf{t}) \prod_{j \neq i} q_j^{old}(\mathbf{t}_j) d\mathbf{t}_j \right)}{\int \exp \left(\int \log p(\mathbf{x}, \mathbf{t}) \prod_{j \neq i} q_j^{old}(\mathbf{t}_j) d\mathbf{t}_j \right) d\mathbf{t}_i}. \quad (5)$$

Эта формула является основным результатом вариационного подхода. Заметим, что оптимальное распределение $q_i(\mathbf{t}_i)$ зависит от всех остальных распределений $q_j(\mathbf{t}_j)$ для $j \neq i$. Поэтому при применении вариационного подхода возникает итерационная схема, в которой последовательно пересчитываются отдельные компоненты факторизованного распределения $q(\mathbf{t})$. При этом на каждом шаге итерации происходит монотонное увеличение нижней границы $\mathcal{L}(q)$ для обоснованности $\log p(\mathbf{x})$, а итерационная оптимизация проводится до сходимости по значению $\mathcal{L}(q)$. Заметим также, что в построениях выше не накладывалось никаких ограничений на семейство распределений $q(\mathbf{t})$, кроме факторизации.

Алгоритм 1 Оценка отношения нормировочных констант двух распределений с помощью схемы Гиббса

Вход: Ненормированные распределения $\tilde{p}_A(\mathbf{t})$ и $\tilde{p}_B(\mathbf{t})$.

Выход: Оценка отношения нормировочных констант Z_B/Z_A .

- 1: Построить последовательность распределений $\tilde{p}_k(\mathbf{t}) = [\tilde{p}_A(\mathbf{t})]^{1-\beta_k} [\tilde{p}_B(\mathbf{t})]^{\beta_k}$ для набора значений $0 = \beta_1 < \beta_2 < \dots < \beta_{K-1} < \beta_K = 1$.
 - 2: для $m = 1, \dots, M$
 - 3: Сгенерировать \mathbf{t}^1 из распределения $p_1(\mathbf{t})$;
 - 4: Сделать шаг по схеме Гиббса $\mathbf{t}^2 \sim T_2(\mathbf{t}, \mathbf{t}^1)$ генерации из распределения $p_2(\mathbf{t})$;
 - 5: ...
 - 6: Сделать шаг по схеме Гиббса $\mathbf{t}^{K-1} \sim T_{K-1}(\mathbf{t}, \mathbf{t}^{K-2})$ генерации из распределения $p_{K-1}(\mathbf{t})$;
 - 7: $w_m = \prod_{k=1}^{K-1} (\tilde{p}_{k+1}(\mathbf{t}^k) / \tilde{p}_k(\mathbf{t}^k))$;
 - 8: $Z_B/Z_A \simeq \frac{1}{M} \sum_{m=1}^M w_m$.
-

Ликбез: оценка нормировочной константы распределения с помощью схемы Гиббса

Предположим, что у нас имеется вероятностное распределение, известное с точностью до нормировочной константы $p(\mathbf{t}) = \tilde{p}(\mathbf{t})/Z$. Как было отмечено выше, схема Гиббса позволяет сгенерировать выборку из этого распределения $\mathbf{t}^1, \dots, \mathbf{t}^N$, которая затем может быть использована для оценки статистики распределения $f(\mathbf{t})$ по формуле (2).

Рассмотрим задачу оценки нормировочной константы распределения Z . Эта нормировочная константа играет роль обоснованности модели и может быть использована для сравнения различных вероятностных моделей между собой, а также для оценки вероятности $p(\mathbf{t})$ для тестовых объектов. Нормировочная константа является «нулевой статистикой» распределения и поэтому не может быть оценена с помощью формулы (2).

Предположим, что у нас имеется два распределения $p_A(\mathbf{t}) = \tilde{p}_A(\mathbf{t})/Z_A$ и $p_B(\mathbf{t}) = \tilde{p}_B(\mathbf{t})/Z_B$. Тогда отношение двух нормировочных констант можно оценить по следующей схеме:

$$\frac{Z_B}{Z_A} = \frac{\int \tilde{p}_B(\mathbf{t}) d\mathbf{t}}{Z_A} = \int \frac{\tilde{p}_B(\mathbf{t})}{Z_A} d\mathbf{t} = \int \frac{\tilde{p}_B(\mathbf{t})}{\tilde{p}_A(\mathbf{t})} p_A(\mathbf{t}) d\mathbf{t} \simeq \frac{1}{M} \sum_{m=1}^M \frac{\tilde{p}_B(\mathbf{t}^m)}{\tilde{p}_A(\mathbf{t}^m)}. \quad (6)$$

Здесь $\mathbf{t}^1, \dots, \mathbf{t}^M$ – выборка из распределения $p_A(\mathbf{t})$, которую можно сгенерировать, например, по схеме Гиббса. Если у распределения $p_A(\mathbf{t})$ нормировочная константа известна, то тогда мы можем оценить абсолютное значение Z_B .

К сожалению, схема (6) применима только в случае, когда распределение $p_A(\mathbf{t})$ является хорошим приближением для $p_B(\mathbf{t})$. На практике поиск хорошего аналитического приближения для интересующего нас распределения $p_B(\mathbf{t})$ может оказаться очень трудной задачей. В этом случае можно построить серию промежуточных распределений $p_A = p_1, p_2, \dots, p_{K-1}, p_K = p_B$ и оценить нормировочную константу Z_B из соотношения:

$$\frac{Z_B}{Z_A} = \frac{Z_K}{Z_1} = \frac{Z_2}{Z_1} \frac{Z_3}{Z_2} \dots \frac{Z_{K-1}}{Z_{K-2}} \frac{Z_K}{Z_{K-1}}.$$

Здесь каждое отношение Z_{k+1}/Z_k оценивается по схеме (6) путем генерации выборки из распределения $p_k(\mathbf{t})$, а серия промежуточных распределений строится как $\tilde{p}_k(\mathbf{t}) = [\tilde{p}_A(\mathbf{t})]^{1-\beta_k} [\tilde{p}_B(\mathbf{t})]^{\beta_k}$ для некоторого набора значений $0 = \beta_1 < \beta_2 < \dots < \beta_{K-1} < \beta_K = 1$.

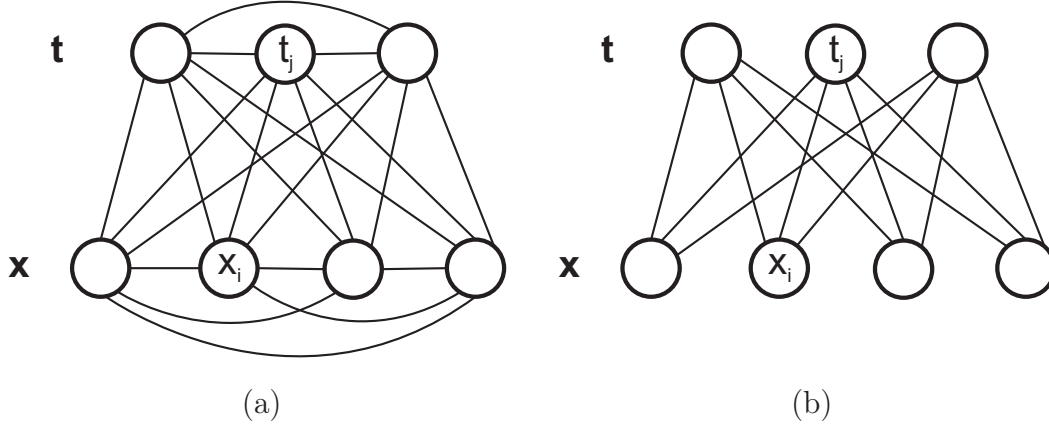


Рис. 1: Граф марковской сети для Boltzmann Machine (a) и для Restricted Boltzmann Machine (b).

Предположим, что для каждого промежуточного распределения $p_k(\mathbf{t}) = \tilde{p}_k(\mathbf{t})/Z_k$, известного с точностью до нормировочной константы, мы можем применить схему Гиббса генерации выборки из этого распределения. Обозначим один шаг такой схемы Гиббса через $\mathbf{t}_{next} \sim T_k(\mathbf{t}, \mathbf{t}_{pred})$. Тогда итоговую схему оценки нормировочной константы Z_B можно представить как Алгоритм 1.

Boltzmann Machine

Модель Boltzmann Machine (BM) [1] представляет собой марковскую сеть с двумя слоями (см. рис. 1,a). В нижнем слое стоят наблюдаемые переменные \mathbf{x} , в верхнем слое – ненаблюдаемые переменные \mathbf{t} . Граф сети является полностью связным. Все переменные являются бинарными $x_i \in \{0, 1\}, t_j \in \{0, 1\}$. Энергию введенной марковской сети можно записать как

$$E(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta}) = -\mathbf{t}^T W \mathbf{x} - \frac{1}{2} \mathbf{x}^T L \mathbf{x} - \frac{1}{2} \mathbf{t}^T J \mathbf{t}.$$

Здесь $\boldsymbol{\theta} = \{W, L, J\}$, $L_{ii} = J_{jj} = 0$. Как правило, в энергию $E(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta})$ добавляются также унарные потенциалы $-\mathbf{x}^T \mathbf{a} - \mathbf{t}^T \mathbf{b}$. Однако, для упрощения дальнейших выкладок унарные потенциалы будут опускаться. Энергия $E(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta})$ задает вероятностное распределение

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \sum_{\mathbf{t}} \exp(-E(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta})),$$

где $Z(\boldsymbol{\theta}) = \int \exp(-E(\mathbf{x}, \mathbf{t}|\boldsymbol{\theta})) d\mathbf{x} d\mathbf{t}$ – нормировочная константа распределения.

Рассмотрим задачу настройки параметров BM по выборке $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ с помощью метода максимального правдоподобия:

$$p(X|\boldsymbol{\theta}) \rightarrow \max_{\boldsymbol{\theta}}.$$

Распределение $p(\mathbf{x}|\boldsymbol{\theta})$ принадлежит экспоненциальному семейству распределений. Поэтому мы

Алгоритм 2 Алгоритм обучения ВМ

Вход: Выборка $\mathbf{x}_1, \dots, \mathbf{x}_N$

Выход: Параметры ВМ $\boldsymbol{\theta} = \{W, L, J\}$

- 1: Инициализация параметров $\boldsymbol{\theta}_0$ случайным образом.
- 2: для $iter = 1, \dots, \#iter$
- 3: для $n = 1, \dots, N$
- 4: Сгенерировать выборку $\mathbf{t}_1, \dots, \mathbf{t}_M$ из распределения $p(\mathbf{t}|\mathbf{x}_n)$ по схеме Гиббса;
- 5: Оценить достаточные статистики $\mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{x}_n^T) = \frac{1}{M} \sum_{m=1}^M \mathbf{t}_m \mathbf{x}_n^T$, $\mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{t}^T)$;
- 6: Сгенерировать выборку $(\mathbf{t}_1, \mathbf{x}_1), \dots, (\mathbf{t}_M, \mathbf{x}_M)$ из распределения $p(\mathbf{t}, \mathbf{x})$ по схеме Гиббса;
- 7: Оценить достаточные статистики $\mathbb{E}_{\mathbf{t}, \mathbf{x}}(\mathbf{t}\mathbf{x}^T) = \frac{1}{M} \sum_{m=1}^M \mathbf{t}_m \mathbf{x}_m^T$, $\mathbb{E}_{\mathbf{t}, \mathbf{x}}(\mathbf{t}\mathbf{t}^T)$, $\mathbb{E}_{\mathbf{t}, \mathbf{x}}(\mathbf{x}\mathbf{x}^T)$;
- 8: Сделать шаг по градиенту:

$$\begin{aligned} W_{iter} &= W_{iter-1} + \alpha_{iter}^W \left(\sum_{n=1}^N \mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{x}_n^T)/N - \mathbb{E}_{\mathbf{t}, \mathbf{x}}(\mathbf{t}\mathbf{x}^T) \right), \\ J_{iter} &= J_{iter-1} + \alpha_{iter}^J \left(\sum_{n=1}^N \mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{t}^T)/N - \mathbb{E}_{\mathbf{t}, \mathbf{x}}(\mathbf{t}\mathbf{t}^T) \right), \\ L_{iter} &= L_{iter-1} + \alpha_{iter}^L \left(\sum_{n=1}^N (\mathbf{x}_n \mathbf{x}_n^T)/N - \mathbb{E}_{\mathbf{t}, \mathbf{x}}(\mathbf{x}\mathbf{x}^T) \right). \end{aligned}$$

можем воспользоваться общим результатом (1) и записать:

$$\begin{aligned} \frac{1}{N} \nabla_W \log p(X|\boldsymbol{\theta}) &= \mathbb{E}_{data}(\mathbf{t}\mathbf{x}^T) - \mathbb{E}_{model}(\mathbf{t}\mathbf{x}^T), \\ \frac{1}{N} \nabla_J \log p(X|\boldsymbol{\theta}) &= \mathbb{E}_{data}(\mathbf{t}\mathbf{t}^T) - \mathbb{E}_{model}(\mathbf{t}\mathbf{t}^T), \\ \frac{1}{N} \nabla_L \log p(X|\boldsymbol{\theta}) &= \mathbb{E}_{data}(\mathbf{x}\mathbf{x}^T) - \mathbb{E}_{model}(\mathbf{x}\mathbf{x}^T). \end{aligned}$$

Следовательно, задачу максимизации правдоподобия в модели ВМ можно решать с помощью градиентного подъема, в котором на каждом шаге необходимо оценивать значение градиента путем оценки статистик по данным и по модели. Эти статистики можно оценить с помощью схемы Гиббса. Для применения схемы Гиббса необходимо уметь генерировать выборку из одномерных условных распределений вида $p(x_i|\mathbf{t}, \mathbf{x}_{\setminus i})$ и $p(t_j|\mathbf{x}, \mathbf{t}_{\setminus j})$. Эти распределения для ВМ могут быть найдены аналитически:

$$\begin{aligned} p(x_i = 1|\mathbf{t}, \mathbf{x}_{\setminus i}) &= \frac{1}{1 + \exp(-\sum_j W_{ij}t_j - \sum_k L_{ki}x_i)}, \\ p(t_j = 1|\mathbf{x}, \mathbf{t}_{\setminus j}) &= \frac{1}{1 + \exp(-\sum_i W_{ij}x_i - \sum_k J_{kj}t_k)}. \end{aligned}$$

Рассмотренная схема обучения ВМ представлена в алгоритме 2. Данный алгоритм был предложен в начале 80-х годов в работе [2]. К сожалению, на практике данный алгоритм не используется, т.к. он требует очень большого времени работы.

Значительное улучшение рассмотренного алгоритма было предложено в начале 90-х годов в работе [3] (см. Алгоритм 3). В этом варианте для оценки статистик по каждому из распределений предлагается запускать несколько т.н. продолжающихся (persistent) цепей Гиббса. При этом на каждой итерации делается всего один или несколько шагов по каждой из схем Гиббса, а шаг движения по градиенту выбирается маленьким. Здесь предполагается, что при небольшом изменении параметров $\boldsymbol{\theta}$ условные распределения $p(\mathbf{t}|\mathbf{x}_n, \boldsymbol{\theta})$ и совместное распределение

Алгоритм 3 Улучшенный алгоритм обучения ВМ с помощью продолжающихся цепей Гиббса

Вход: Выборка $\mathbf{x}_1, \dots, \mathbf{x}_N$

Выход: Параметры ВМ $\theta = \{W, L, J\}$

- 1: Инициализация параметров θ_0 случайным образом;
- 2: Инициализация M общих цепей Гиббса и M цепей Гиббса для каждого объекта обучения $(\mathbf{t}_0^{(0),1}, \mathbf{x}_0^1), \dots, (\mathbf{t}_0^{(0),M}, \mathbf{x}_0^M), \mathbf{t}_0^{(1),1}, \dots, \mathbf{t}_0^{(N),M}$;
- 3: для $iter = 1, \dots, \#iter$
- 4: для $n = 1, \dots, N$
- 5: Сделать шаг по схеме Гиббса для каждой из M цепей для объекта n и получить $\mathbf{t}_{iter}^{(n),1}, \dots, \mathbf{t}_{iter}^{(n),M}$;
- 6: Оценить достаточные статистики $\mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{x}_n^T) = \frac{1}{M} \sum_{m=1}^M \mathbf{t}_{iter}^{(n),m} \mathbf{x}_n^T$, $\mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{t}^T)$;
- 7: Сделать шаг по схеме Гиббса для каждой из M общих цепей и получить $(\mathbf{t}_{iter}^{(0),1}, \mathbf{x}_{iter}^1), \dots, (\mathbf{t}_{iter}^{(0),M}, \mathbf{x}_{iter}^M)$;
- 8: Оценить достаточные статистики $\mathbb{E}_{\mathbf{t},\mathbf{x}}(\mathbf{t}\mathbf{x}^T) = \frac{1}{M} \sum_{m=1}^M \mathbf{t}_{iter}^{(0),m} (\mathbf{x}_{iter}^m)^T$, $\mathbb{E}_{\mathbf{t},\mathbf{x}}(\mathbf{t}\mathbf{t}^T)$, $\mathbb{E}_{\mathbf{t},\mathbf{x}}(\mathbf{x}\mathbf{x}^T)$;
- 9: Сделать шаг по градиенту:

$$W_{iter} = W_{iter-1} + \alpha_{iter}^W \left(\sum_{n=1}^N \mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{x}_n^T) / N - \mathbb{E}_{\mathbf{t},\mathbf{x}}(\mathbf{t}\mathbf{x}^T) \right),$$

$$J_{iter} = J_{iter-1} + \alpha_{iter}^J \left(\sum_{n=1}^N \mathbb{E}_{\mathbf{t}|\mathbf{x}_n}(\mathbf{t}\mathbf{t}^T) / N - \mathbb{E}_{\mathbf{t},\mathbf{x}}(\mathbf{t}\mathbf{t}^T) \right),$$

$$L_{iter} = L_{iter-1} + \alpha_{iter}^L \left(\sum_{n=1}^N (\mathbf{x}_n \mathbf{x}_n^T) / N - \mathbb{E}_{\mathbf{t},\mathbf{x}}(\mathbf{x}\mathbf{x}^T) \right).$$

$p(\mathbf{t}, \mathbf{x}|\theta)$ незначительно отличаются от инвариантных распределений цепей Гиббса. Поэтому достаточно сделать лишь небольшое число шагов по каждой из цепей Гиббса, чтобы достигнуть новых инвариантных распределений.

Алгоритм 3 используется на практике для обучения небольших ВМ. При переходе к ВМ, которые представляют практический интерес во многих задачах, описанный алгоритм работает слишком долго. Дальнейшие попытки улучшения алгоритма обучения ВМ связаны с рассмотрением упрощенной модели – Restricted Boltzmann Machine.

Restricted Boltzmann Machine

Модель Restricted Boltzmann Machine [4] полностью повторяет модель Boltzmann Machine за тем исключением, что в графе марковской сети отсутствуют связи на одном уровне (см. рис. 1,b). Энергия и вероятностное распределение такой сети могут быть записаны как

$$E(\mathbf{x}, \mathbf{t}|W) = -\mathbf{t}^T W \mathbf{x}, \quad p(\mathbf{x}|\mathbf{t}, W) = \frac{1}{Z(W)} \sum_{\mathbf{x}} \exp(-E(\mathbf{x}, \mathbf{t}|W)).$$

Это вероятностное распределение принадлежит экспоненциальному семейству. Поэтому задачу обучения RBM по данным можно решать с помощью градиентного подъема, где градиент по-прежнему вычисляется как

$$\frac{1}{N} \nabla_W \log p(X|W) = \mathbb{E}_{data}(\mathbf{t}\mathbf{x}^T) - \mathbb{E}_{model}(\mathbf{t}\mathbf{x}^T).$$

Ключевое отличие RBM от ВМ состоит в том, что скрытые переменные \mathbf{t} являются условно независимыми при известных \mathbf{x} и наоборот: переменные \mathbf{x} являются условно независимыми

Алгоритм 4 Алгоритм обучения RBM по схеме 1-Contrastive Divergence

Вход: Выборка $\mathbf{x}_1, \dots, \mathbf{x}_N$

Выход: Параметры RBM W

- 1: Инициализация параметров W случайным образом;
 - 2: для $iter = 1, \dots, \#iter$
 - 3: для $n = 1, \dots, N$
 - 4: $\mathbf{t}_n \sim p(\mathbf{t}|\mathbf{x}_n)$;
 - 5: $\mathbf{x}_n^{pred} \sim p(\mathbf{x}|\mathbf{t}_n)$;
 - 6: $\mathbf{t}_n^{pred} \sim p(\mathbf{t}|\mathbf{x}_n^{pred})$;
 - 7: Оценить необходимые достаточные статистики $\mathbb{E}_{data}(\mathbf{t}\mathbf{x}^T) = \frac{1}{N} \sum_{n=1}^N \mathbf{t}_n \mathbf{x}_n^T$, $\mathbb{E}_{model}(\mathbf{t}\mathbf{x}^T) = \frac{1}{N} \sum_{n=1}^N \mathbf{t}_n^{pred} (\mathbf{x}_n^{pred})^T$;
 - 8: Сделать шаг по градиенту: $W_{iter} = W_{iter-1} + \alpha_{iter}(\mathbb{E}_{data}(\mathbf{t}\mathbf{x}^T) - \mathbb{E}_{model}(\mathbf{t}\mathbf{x}^T))$.
-

при известных \mathbf{t} . Соответствующие вероятностные распределения могут быть вычислены аналитически:

$$p(\mathbf{t}|\mathbf{x}, W) = \prod_j p(t_j|\mathbf{x}, W), \quad p(t_j = 1|\mathbf{x}, W) = \frac{1}{1 + \exp(-\sum_i W_{ij}x_i)},$$
$$p(\mathbf{x}|\mathbf{t}, W) = \prod_i p(x_i|\mathbf{t}, W), \quad p(x_i = 1|\mathbf{t}, W) = \frac{1}{1 + \exp(-\sum_j W_{ij}t_j)}.$$

Таким образом, генерация выборки из распределений $p(\mathbf{t}|\mathbf{x})$ и $p(\mathbf{x}|\mathbf{t})$ в модели RBM не требует привлечения затратной по времени схемы Гиббса и может быть осуществлена непосредственно. Это обстоятельство позволяет предложить эффективный алгоритм обучения RBM, представленный как Алгоритм 4. В этом алгоритме для оценки статистик по данным и по модели на каждой итерации осуществляется всего один шаг генерации. При таком подходе оценка обеих статистик будет обладать большой дисперсией. Однако, оцениваемый градиент представляет собой разность двух статистик. По утверждению авторов алгоритма 4, шум, который появляется при оценке обеих статистик, имеет в значительной степени одну и ту же природу и взаимно уничтожается при вычитании статистик. Тем не менее, для более точной сходимости к оптимуму, авторы метода рекомендуют увеличивать число шагов генерации на каждом шаге итерационного процесса по мере приближения к оптимуму.

Несмотря на наличие эффективного алгоритма обучения, модель RBM редко используется сама по себе в практических задачах, т.к. является слишком простой и не позволяет описывать данные со сложными взаимосвязями. Вместо этого на практике используется модель Deep Boltzmann Machine, которая, с одной стороны, является значительным усложнением RBM для моделирования сложных экспериментальных данных, а, с другой стороны, имеет относительно эффективный алгоритм обучения.

Deep Boltzmann Machine

Модель Deep Boltzmann Machine [1] представляет собой марковскую сеть, граф которой показан на рис. 2, справа. В этом графе есть несколько слоев, нижний слой соответствует наблюдаемым переменным \mathbf{x} , остальные – ненаблюдаемым. Связи в графе есть только между переменными соседних слоев. Энергия модели DBM с тремя слоями скрытых переменных выглядит следую-

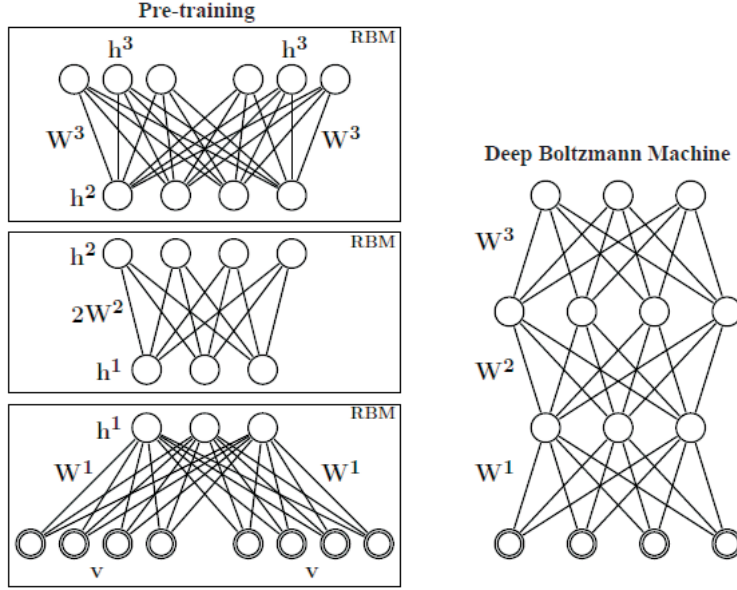


Рис. 2: Слева: схема алгоритма предобучения для модели DBM, справа: граф марковской сети для DBM.

щим образом:

$$E(\mathbf{x}, \mathbf{t}^1, \mathbf{t}^2, \mathbf{t}^3 | W^1, W^2, W^3) = -\mathbf{x}^T W^1 \mathbf{t}^1 - (\mathbf{t}^1)^T W^2 \mathbf{t}^2 - (\mathbf{t}^2)^T W^3 \mathbf{t}^3. \quad (7)$$

Как и раньше, распределение вероятностей для введенной марковской сети принадлежит экспоненциальному семейству:

$$p(\mathbf{x} | W^1, W^2, W^3) = \frac{1}{Z(W^1, W^2, W^3)} \sum_{\mathbf{t}^1, \mathbf{t}^2, \mathbf{t}^3} \exp(-E(\mathbf{x}, \mathbf{t}^1, \mathbf{t}^2, \mathbf{t}^3 | W^1, W^2, W^3)).$$

Эффективный алгоритм обучения DBM [1] предлагает два основных новшества по сравнению с алгоритмом 3: 1) использование специальной процедуры предобучения параметров сети W с помощью многократного вызова процедуры обучения RBM, описанной выше, и 2) использование вариационного подхода вместо схем Гиббса для оценки достаточных статистик по данным.

Рассмотрим процедуру предобучения параметров DBM W (см. рис. 2). Веса первого слоя W^1 находятся путем обучения RBM, в которой нижний слой соответствует наблюдаемым переменным \mathbf{x} в DBM, а верхний слой соответствует ненаблюдаемым переменным первого уровня \mathbf{t}^1 . Затем значения \mathbf{t}^1 генерируются из обученной модели RBM $p(\mathbf{t} | \mathbf{x}, W^1)$ и используются в качестве наблюдаемых переменных в RBM для обучения весов следующего уровня W^2 . Процедура повторяется вплоть до верхнего уровня. Одним из недостатков описанной процедуры является тот факт, что при обучении очередной RBM никак не учитывается влияние верхних уровней скрытых переменных на текущий слой скрытых переменных. Для частичной компенсации данного обстоятельства авторы метода предобучения предлагают следующую модификацию. В RBM для первого слоя наблюдаемые переменные дублируются, а обучение весов W^1 происходит при ограничении равенства весов у соответствующих связей между скрытыми и наблюдаемыми переменными. Аналогично, в RBM для последнего слоя дублируются скрытые

Алгоритм 5 Предобучение модели DBM

Вход: Выборка $\mathbf{x}_1, \dots, \mathbf{x}_N$.

Выход: Начальное приближение для весов DBM W^1, \dots, W^L .

Продублировать узлы \mathbf{x} на нижнем слое и обучить RBM (настроить W^1) для первого скрытого слоя при условии равенства весов W^1 для дублированных узлов;

Зафиксировать W^1 и сгенерировать выборку \mathbf{t}^1 из распределения $p(\mathbf{t}|\mathbf{x}, 2W^1)$;

Обучить RBM для второго слоя с наблюдаемыми данными \mathbf{t}^1 с весами $2W^2$;

...

При обучении RBM на последнем слое продублировать скрытые переменные и приравнять веса продублированных узлов;

Использовать веса $\{W^1, \dots, W^L\}$ в качестве начального приближения для весов DBM.

переменные и обучение происходит при ограничении равенства весов у соответствующих связей между переменными. Результирующие веса для промежуточных слоев устанавливаются как половина от весов, полученных в результате обучения набора RBM. Итоговая схема предобучения DBM представлена в алгоритме 5.

Рассмотрим применение вариационного подхода для оценки достаточных статистик по данным при обучении DBM. Для этого рассмотрим приближение для апостериорного распределения $p(\mathbf{t}^1, \mathbf{t}^2, \mathbf{t}^3|\mathbf{x}, W^1, W^2, W^3)$ в полностью факторизованном семействе распределений:

$$q(\mathbf{t}^1, \mathbf{t}^2, \mathbf{t}^3|\boldsymbol{\mu}^1, \boldsymbol{\mu}^2, \boldsymbol{\mu}^3) = \prod_j q(t_j^1|\mu_j^1) \prod_k q(t_k^2|\mu_k^2) \prod_m q(t_m^3|\mu_m^3).$$

Здесь $\boldsymbol{\mu}^1, \boldsymbol{\mu}^2, \boldsymbol{\mu}^3$ — параметры распределения q , μ_l^i имеет смысл вероятности того, что $t_l^i = 1$. Наилучшее приближение q в выбранном семействе может быть найдено путем максимизации нижней границы для неполного правдоподобия $\log p(\mathbf{x}|W)$:

$$\log p(\mathbf{x}|W) \geq \mathcal{L}(q) = \mathbf{x}^T W^1 \boldsymbol{\mu}^1 + (\boldsymbol{\mu}^1)^T W^2 \boldsymbol{\mu}^2 + (\boldsymbol{\mu}^2)^T W^3 \boldsymbol{\mu}^3 - \log Z(W) + \mathcal{H}(q) \rightarrow \max_{\boldsymbol{\mu}^1, \boldsymbol{\mu}^2, \boldsymbol{\mu}^3}. \quad (8)$$

Здесь через $\mathcal{H}(q)$ обозначена энтропия распределения q . Задача оптимизации (8) может быть решена с помощью итерационного процесса (5), который в данном случае выглядит как

$$\begin{aligned} \mu_j^1 &= \frac{1}{1 + \exp(-\sum_i W_{ij}^1 x_i - \sum_k W_{jk}^2 \mu_k^2)}, \\ \mu_k^2 &= \frac{1}{1 + \exp(-\sum_j W_{jk}^2 \mu_j^1 - \sum_m W_{km}^3 \mu_m^3)}, \\ \mu_m^3 &= \frac{1}{1 + \exp(-\sum_k W_{km}^3 \mu_k^2)}. \end{aligned}$$

После того, как параметры $\boldsymbol{\mu}$ определены, новые значения весов W можно найти путем движения вдоль градиента нижней границы (8)

$$\begin{aligned} \nabla_{W^1} \mathcal{L}(q) &= \mathbf{x}(\boldsymbol{\mu}^1)^T - \nabla_{W^1} \log Z(W), \\ \nabla_{W^2} \mathcal{L}(q) &= \boldsymbol{\mu}^1(\boldsymbol{\mu}^2)^T - \nabla_{W^2} \log Z(W), \\ \nabla_{W^3} \mathcal{L}(q) &= \boldsymbol{\mu}^2(\boldsymbol{\mu}^3)^T - \nabla_{W^3} \log Z(W). \end{aligned}$$

Здесь, как и раньше, градиенты нормировочной константы $Z(W)$ могут быть оценены с помощью продолжающихся цепей Гиббса по схеме алгоритма 3.

Объединяя вместе все вышесказанное, можно составить итоговую схему обучения DBM (см. Алгоритм 6).

Алгоритм 6 Обучение модели DBM

Вход: Выборка $\mathbf{x}_1, \dots, \mathbf{x}_N$.

Выход: Параметры DBM W^1, \dots, W^L .

- 1: Предобучить параметры DBM W^1, \dots, W^L по алгоритму 5;
- 2: Инициализировать M цепей Гиббса случайным образом $(\mathbf{x}_0^1, \mathbf{t}_0^{1,1}, \dots, \mathbf{t}_0^{L,1}), \dots, (\mathbf{x}_0^M, \mathbf{t}_0^{1,M}, \dots, \mathbf{t}_0^{L,M})$;
- 3: **для** $iter = 1, \dots, \#iter$
- 4: **для** $n = 1, \dots, N$
- 5: Сделать вариационное приближение для распределения $p(\mathbf{t}^1, \dots, \mathbf{t}^L | \mathbf{x}_n)$ и найти его параметры $\{\boldsymbol{\mu}_n^1, \dots, \boldsymbol{\mu}_n^L\}$;
- 6: **для** $m = 1, \dots, M$
- 7: Сделать шаг по схеме Гиббса и получить $(\mathbf{x}_{iter}^m, \mathbf{t}_{iter}^{1,m}, \dots, \mathbf{t}_{iter}^{L,m})$;
- 8: Сделать шаг по градиенту:

$$\begin{aligned} W_{iter}^1 &= W_{iter-1}^1 + \alpha_{iter} \left[\left(\sum_{n=1}^N \mathbf{x}_n (\boldsymbol{\mu}_n^1)^T \right) / N - \left(\sum_{m=1}^M \mathbf{x}_{iter}^m (\mathbf{t}_{iter}^{1,m})^T \right) / M \right], \\ W_{iter}^2 &= W_{iter-1}^2 + \alpha_{iter} \left[\left(\sum_{n=1}^N \boldsymbol{\mu}_n^1 (\boldsymbol{\mu}_n^2)^T \right) / N - \left(\sum_{m=1}^M \mathbf{t}_{iter}^{1,m} (\mathbf{t}_{iter}^{2,m})^T \right) / M \right], \\ &\dots \\ W_{iter}^L &= W_{iter-1}^L + \alpha_{iter} \left[\left(\sum_{n=1}^N \boldsymbol{\mu}_n^{L-1} (\boldsymbol{\mu}_n^L)^T \right) / N - \left(\sum_{m=1}^M \mathbf{t}_{iter}^{L-1,m} (\mathbf{t}_{iter}^{L,m})^T \right) / M \right]. \end{aligned}$$

Использование Deep Boltzmann Machine

Deep Boltzmann Machine как порождающая модель

Рассмотрим вероятностную модель DBM как способ восстановления плотности в пространстве \mathbf{x} по данным. В этом случае нам необходимо уметь оценивать значение плотности $p(\mathbf{x}|W)$ для обучающих и/или тестовых объектов, что в свою очередь требует оценки нормировочной константы распределения $Z(W)$. Для ее оценивания воспользуемся методом построения промежуточных распределений, описанный выше.

Пусть модель DBM имеет два слоя скрытых переменных. Следовательно, энергия системы определяется формулой (7). Суммируя по \mathbf{t}^1 и \mathbf{t}^3 , получим ненормированную плотность в пространстве $(\mathbf{x}, \mathbf{t}^2)$

$$\tilde{p}(\mathbf{x}, \mathbf{t}^2 | W) = \sum_{\mathbf{t}^1, \mathbf{t}^3} \tilde{p}(\mathbf{x}, \mathbf{t}^1, \mathbf{t}^2, \mathbf{t}^3 | W) = \prod_j \left(1 + \exp \left(\sum_i W_{ij}^1 x_i + \sum_k W_{jk}^2 t_k^2 \right) \right) \prod_m \left(1 + \exp \left(\sum_k W_{km}^3 t_k^2 \right) \right).$$

Для DBM с произвольным числом слоев мы всегда можем аналитически просуммировать по всем четным или нечетным слоям. Для оценки нормировочной константы применим метод промежуточных распределений между равномерным распределением и интересующим нас распределением $p(\mathbf{x}, \mathbf{t}^2 | W)$:

$$p_k(\mathbf{x}, \mathbf{t}^2 | W) \propto [p_{const}(\mathbf{x}, \mathbf{t}^2)]^{1-\beta_k} [p(\mathbf{x}, \mathbf{t}^2 | W)]^{\beta_k},$$

где $p_{const}(\mathbf{x}, \mathbf{t}^2) = 1/2^{N_1+N_2}$, N_1, N_2 – число узлов на первом и втором слое DBM соответственно. Легко показать, что схема Гиббса для генерации выборки из распределения p_k реализуется с

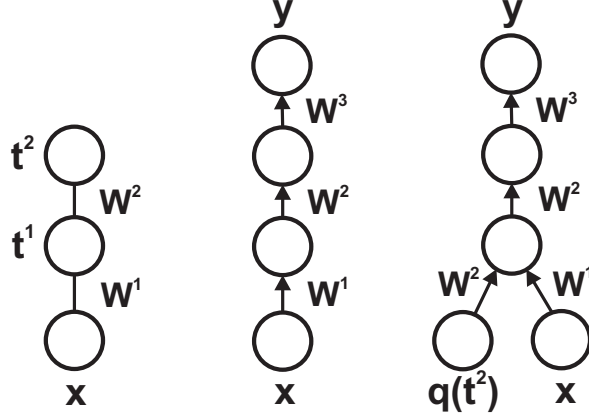


Рис. 3: Слева: модель DBM с двумя слоями скрытых переменных, Центр: преобразование DBM к многослойному перцептрону, Справа: преобразование DBM к многослойному перцептрону, в котором дополнительно на вход подается вектор признаков, полученный с помощью DBM.

помощью следующих формул:

$$p_k(x_i = 1 | \mathbf{t}^2, \mathbf{x}_{\setminus i}) = \frac{1}{1 + \prod_j \left(\frac{1 + \exp(\sum_{i_1 \neq i} W_{i_1 j}^2 x_{i_1} + \sum_k W_{jk}^2 t_k^2)}{1 + \exp(W_{ij}^2 + \sum_{i_1 \neq i} W_{i_1 j}^2 x_{i_1} + \sum_k W_{jk}^2 t_k^2)} \right)^{\beta_k}},$$

$$p_k(t_k^2 = 1 | \mathbf{x}, \mathbf{t}_{\setminus k}^2) = \frac{1}{1 + \left[\prod_j \frac{1 + \exp(\sum_j W_{ij}^1 x_i + \sum_{k_1 \neq k} W_{jk_1}^2 t_{k_1}^2)}{1 + \exp(\sum_j W_{ij}^1 x_i + \sum_{k_1 \neq k} W_{jk_1}^2 t_{k_1}^2 + W_{jk}^2)} \prod_m \frac{1 + \exp(\sum_{k_1 \neq k} W_{k_1 m}^3 t_{k_1}^2)}{1 + \exp(\sum_{k_1 \neq k} W_{k_1 m}^3 t_{k_1}^2 + W_{km})} \right]^{\beta_k}}.$$

После получения оценки для нормировочной константы $Z(W)$ значение плотности для тестового объекта \mathbf{x}_* можно оценить с помощью вариационного подхода:

$$\log p(\mathbf{x}_* | W) \geq - \sum_{t^1, t^2, t^3} q(t^1, t^2, t^3 | \mu) E(\mathbf{x}_*, t^1, t^2, t^3 | W) + \mathcal{H}(q) - \log Z(W).$$

Здесь q – факторизованное приближение для апостериорного распределения $p(t^1, t^2, t^3 | \mathbf{x}, W)$ из вариационного подхода.

Deep Boltzmann Machine как дискриминативная модель

Модель DBM можно использовать различными способами для решения задач классификации и регрессии.

Первый способ предполагает использовать скрытые признаки, полученные с помощью DBM, в качестве новых признаков объектов, для которых запускается отдельный независимый метод классификации или регрессии. Пусть у нас имеется обученная модель DBM с двумя слоями скрытых переменных (см. рис. 3, слева). Тогда новые признаки для объекта \mathbf{x} представляют собой вероятности $q_k(t_k^2 = 1)$ из вариационной оценки апостериорного распределения $p(\mathbf{t}^2 | \mathbf{x}, W)$. Модификацией этой идеи является преобразование модели DBM в модель многослойного перцептрона вместе с добавлением на последнем слое линейного (для задачи регрессии) или логистического/мультиномиального (для задачи классификации) узла (см. рис. 3, центр). Построенный таким образом многослойный перцептрон обучается с помощью метода обратного

распространения ошибки, в котором в качестве начального приближения для весов сети используются обученные значения W для DBM. Можно также добавить признаки $q_k(t_k^2 = 1)$, обученные с помощью DBM, в набор входов многослойного персептрона (см. рис. 3, справа).

Описанный способ применения DBM как дискриминативной модели позволяет решать задачи классификации и регрессии с частично размеченной обучающей выборкой (semi-supervised learning). Действительно, на этапе обучения DBM информация о метках объектов не используется. Фактически с помощью DBM восстанавливается плотность $p(\mathbf{x})$, которая соответствует осмысленным входным объектам. Затем производится обучение многослойного персептрона по полностью размеченной выборке, которое приводит лишь к небольшой коррекции весов сети.

Второй способ решения задач классификации с помощью DBM предполагает использование байесовского классификатора, где для моделирования плотности каждого из классов используется своя модель DBM. В отличие от предыдущего способа, здесь необходимо обучать отдельную DBM для каждого класса. При этом нельзя решать задачу классификации с частично размеченными выборками.

При третьем способе метка класса добавляется в качестве входа в модель DBM. На этапе распознавания тестового объекта к признакам тестового объекта добавляется по очереди каждая из возможных меток классов и вычисляется энергия $E(\mathbf{x}|W)$ такой конфигурации. Метка класса, соответствующая минимальной энергии, выдается в качестве результата распознавания. В отличие от второго способа, здесь производится обучение одной общей модели DBM, но также нельзя решать задачи классификации с частично размеченными выборками.

Виды переменных в Deep Boltzmann Machine

К настоящему моменту рассматривались модели DBM только с бинарными наблюдаемыми и скрытыми переменными. Однако, на практике существует необходимость работы с переменными других типов, в частности, с непрерывными и K -значными переменными.

K -значные переменные являются удобными для моделирования метки класса, например, при третьем способе использования DBM как дискриминативной модели. K -значные переменные добавляются в DBM как K бинарных переменных x_1, \dots, x_K без внутренних связей с дополнительным ограничением, что одна и только одна из этих переменных может принимать значение единицы. Для применения схемы Гиббса необходимо уметь вычислять условное распределение на каждую из переменных x_k при условии всех ненаблюдаемых переменных \mathbf{t} . Данное условное распределение определяется мультиномиальной функцией:

$$p(x_k|\mathbf{t}, W) = \frac{\exp(\sum_j W_{kj}t_j)}{\sum_{m=1}^K \exp(\sum_j W_{mj}t_j)}.$$

Рассмотрим случай непрерывных наблюдаемых переменных \mathbf{x} и бинарных скрытых переменных \mathbf{t} в модели RBM. В этом случае энергия системы записывается как

$$E(\mathbf{x}, \mathbf{t}|W, \mathbf{a}, \mathbf{b}) = \sum_i \frac{(x_i - a_i)^2}{2} - \sum_j b_j t_j - \sum_{i,j} W_{ij} x_i t_j.$$

Таким образом, для непрерывных переменных добавляются квадратичные унарные потенциалы вместо линейных для случая бинарных переменных. Условные одномерные распределения в

такой модели вычисляются как

$$p(x_i|\mathbf{t}, W, \mathbf{a}, \mathbf{b}) = \mathcal{N}(x_i|a_i + \sum_j W_{ij}t_j, 1),$$

$$p(t_j = 1|\mathbf{x}, W, \mathbf{a}, \mathbf{b}) = \frac{1}{1 + \exp(-b_j - \sum_i W_{ij}x_i)}.$$

Наконец, в случае непрерывных наблюдаемых и скрытых переменных энергия RBM записывается как

$$E(\mathbf{x}, \mathbf{t}|W, \mathbf{a}, \mathbf{b}) = \sum_i \frac{(x_i - a_i)^2}{2} + \sum_j \frac{(t_j - b_j)^2}{2} - \sum_{i,j} W_{ij}x_it_j.$$

Одномерные условные распределения в такой модели вычисляются как

$$p(x_i|\mathbf{t}, W, \mathbf{a}, \mathbf{b}) = \mathcal{N}(x_i|a_i + \sum_j W_{ij}t_j, 1),$$

$$p(t_j|\mathbf{x}, W, \mathbf{a}, \mathbf{b}) = \mathcal{N}(t_j|b_j + \sum_i W_{ij}x_i, 1).$$

Примеры применения Deep Boltzmann Machine

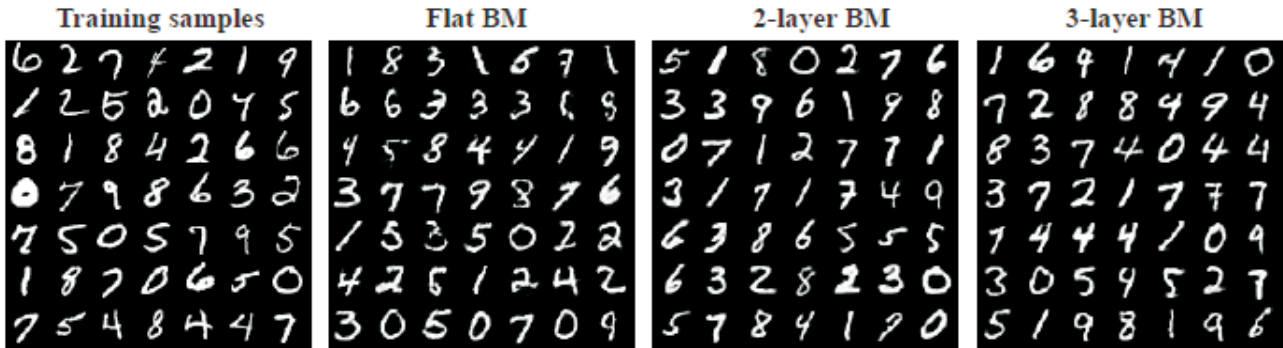


Рис. 4: Примеры сгенерированных объектов в задаче MNIST для различных моделей. Слева: примеры изображений из обучающей выборки, Центр слева: примеры сгенерированных изображений из модели BM, Центр справа: примеры сгенерированных изображений из модели DBM с двумя слоями скрытых переменных, Справа: примеры сгенерированных изображений из модели DBM с тремя слоями скрытых переменных.

Рассмотрим известную задачу распознавания рукописных цифр MNIST¹. В этой задаче объектами являются изображения размера 28×28 , каждое из которых содержит написание одной из цифр 0–9 (см. рис. 4,слева). Всего имеется 60000 изображений в обучающей выборке и 10000 изображений в тестовой выборке. При применении BM в данной задаче каждое изображение интерпретировалось как бинарный вектор длины 784, где значение каждого элемента вектора генерировалось с вероятностью, пропорциональной интенсивности пиксела.

Рассмотрим применение трех моделей Boltzmann Machine для задачи MNIST: 1) двухслойная полносвязная BM с 784 наблюдаемыми переменными и 500 скрытыми переменными, 2)

¹<http://yann.lecun.com/exdb/mnist/>

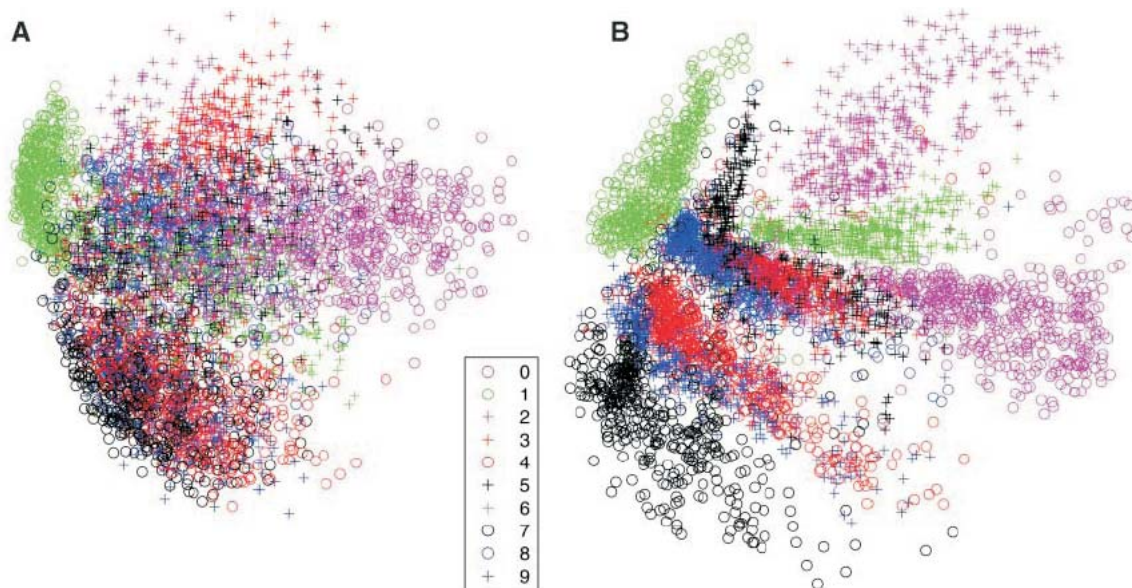


Рис. 5: Визуализация выборки MNIST с помощью метода главных компонент (А) и с помощью модели DBM с конфигурацией 784-1000-500-250-2.

DBM с двумя слоями скрытых переменных 784-500-1000 и 3) DBM с тремя слоями скрытых переменных 784-500-500-1000. На рис. 4 показаны примеры сгенерированных изображений из обученных моделей. Как видно из этого рисунка, генерируемые изображения очень похожи на изображения рукописных цифр. При этом модели DBM обучаются существенно быстрее, чем полносвязная модель ВМ. Дополнительное дискриминативное обучение моделей DBM в этой задаче дает процент ошибок для тестовой выборки 0.95% для DBM с двумя слоями и 1.01% для DBM с тремя слоями скрытых переменных. По утверждению авторов статьи [1] это наилучший процент, которого удалось добиться человечеству для задачи MNIST на сегодняшний день.

Рассмотрим также применение модели DBM с конфигурацией 784-1000-500-250-2 в качестве метода уменьшения размерности в задаче MNIST. На рис. 5 показана проекция обучающей выборки на первые две главные компоненты (слева) и найденное в помощью DBM двухпризнаковое описание для обучающей выборки (справа). Как видно из этого рисунка, модель DBM обеспечивает значительно лучшую разделяемость по классам по сравнению с методом главных компонент.

В заключение данного раздела рассмотрим применение DBM для задачи NORB. Исходные данные в этой задаче представляют собой картинки размера 96×96 , содержащие изображения трехмерных игрушек, снятых в различных условиях освещенности и с различных углов зрения (см. рис. 6, слева). При этом данные объекты разбиты на 5 общих классов: машины, грузовики, самолеты, животные и люди. Задача состоит в предсказании общих классов для тестовых изображений.

Для решения данной задачи применялась модель DBM с конфигурацией 9216-4000-4000-4000, где наблюдаемые переменные являлись непрерывными величинами, значения которых соответствовали уровню интенсивности пикселей. На рис. 6, справа показаны примеры сгенерированных изображений из обученной модели DBM. Дополнительное дискриминативное обучение данной модели DBM дает уровень ошибок 10.8% на тестовой выборке. По утверждению авторов статьи [1] этот уровень ошибок является наилучшим из опубликованных для дан-

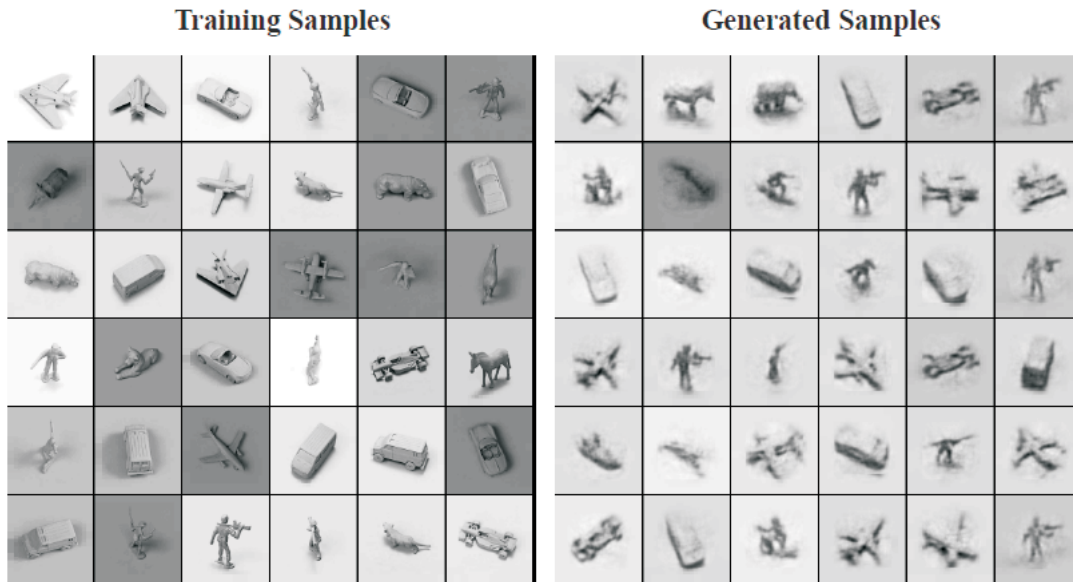


Рис. 6: Слева: примеры изображений в задаче NORB, справа: примеры сгенерированных изображений из обученной модели DBM.



Рис. 7: Иллюстрация к задаче восстановления изображений. Сверху: примеры изображений из тестовой выборки, середина: частично-стертые изображения, подаваемые на вход модели DBM, снизу: восстановленные изображения с помощью обученной модели DBM.

ной задачи на настоящий момент. Кроме того, модель DBM может обучаться на частично-классифицированной обучающей совокупности. При добавлении в обучение большого количества изображений, полученных из исходных простыми преобразованиями пикселей и различными размытиями, позволяет снизить уровень ошибок до 7.2%.

Наконец, рассмотрим применение модели DBM для решения задачи восстановления изображений по данным NORB. На рис. 7 показаны примеры тестовых изображений (сверху), их частично-стертые варианты, подаваемые на вход обученной модели DBM (середина) и результат восстановления (снизу). В частности, здесь модель сумела восстановить стертые части изображений даже для объектов ковбоя, которые ни разу не встречались на этапе обучения.

Список литературы

- [1] R. Salakhutdinov, G. Hinton. An Efficient Learning Procedure for Deep Boltzmann Machines // MIT Technical Report MIT-CSAIL-TR-2010-037, 2010.
- [2] G. Hinton, T. Sejnowski. Optimal perceptual inference // CVPR, 1983.
- [3] R. Neal. Connectionist learning of belief networks // Artificial Intelligence, V. 56, No. 1, 1992, pp. 71–113.
- [4] P. Smolensky. Information processing in dynamical systems: Foundations of harmony theory // Parallel Distributed Computing, V. 1, 1986.