

Ограниченная машина Больцмана (Restricted Boltzmann machine)

Введение

Ограниченная машина Больцмана (RBM) — это генеративная стохастическая искусственная нейронная сеть, которая моделирует распределение вероятности входных данных.

Как следует из названия, RBM являются разновидностью машин Больцмана с ограничением. Ограничение состоит в том, что нейроны данной сети образуют неориентированный двудольный граф без петель. Данный факт означает, что множество нейронов разбивается на две группы (видимые и скрытые нейроны), причем связана любая пара узлов, принадлежащих разным группам, а узлы одной группы связей не имеют. Такая нейронная сеть содержит симметричные синаптические соединения, т.е. видимые и скрытые нейроны взаимодействуют в обоих направлениях. Указанное ограничение позволяет использовать более эффективные алгоритмы обучения по сравнению с теми, которые доступны для общего класса машин Больцмана. Кроме того, особенностью ограниченных машин Больцмана является возможность их обучения без учителя, т.е. не требуется разметка тренировочных данных, на которых настраиваются веса нейронной сети.

Обозначения и вывод формул

Введем следующие обозначения.

- Состояния нейронов видимого слоя. Видимые нейроны являются **стохастическими**, т.е. могут находиться в двух состояниях "включен" или "выключен". Принимают значения 0 или 1. 1 означает, что нейрон активирован ("включен"), 0 - не активирован ("выключен").

$$v_1, v_2, \dots, v_{N_v}; v_i \in \{0, 1\}, i = \overline{1, N_v}, v \in R^{N_v}$$

- Состояния нейронов скрытого слоя. Скрытые нейроны по аналогии с видимыми являются стохастическими и принимают значения 0 или 1.

$$h_1, h_2, \dots, h_{N_h}; h_i \in \{0, 1\}, i = \overline{1, N_h}, h \in R^{N_h}$$

- Матрица связей между нейронами видимого и скрытого слоя:

$$W \in R^{N_v \times N_h}$$

- Вес связи между i -ым нейроном видимого слоя и j -ым нейроном скрытого слоя:

$$w_{ij}, i = \overline{1, N_v}, j = \overline{1, N_h}$$

- Смещение видимого нейрона:

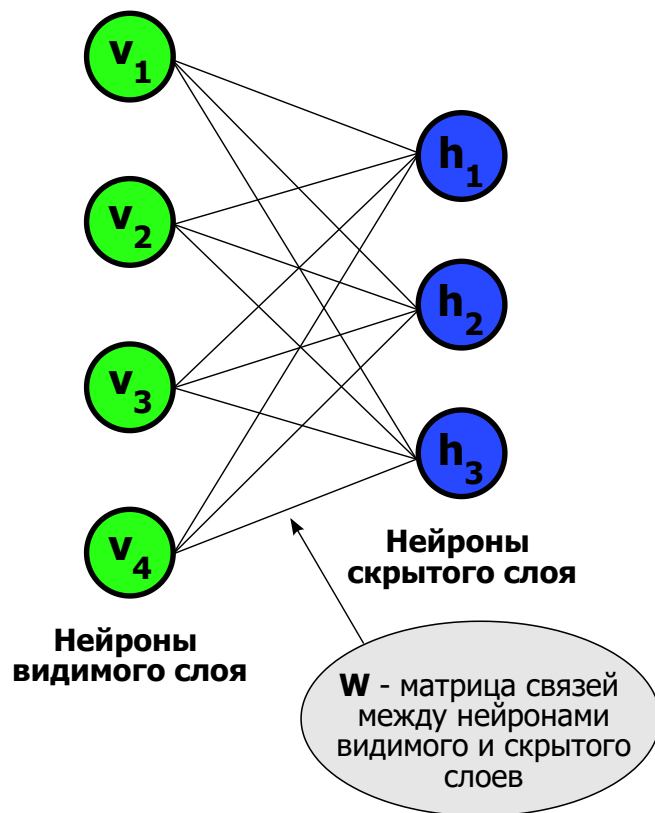
$$a_i, i = \overline{1, N_v}$$

- Смещение скрытого нейрона:

$$b_j, j = \overline{1, N_h}$$

Примечания:

1. Состояние нейрона может принимать значения 1 и -1, 1 означает, что нейрон активирован ("включен"), -1 - не активирован ("выключен").
2. Вектора смещений можно внести в матрицу связей между видимыми и скрытыми нейронами, если ввести фиктивный нейрон, который всегда находится в активированном состоянии. В результате получается равносильная модель RBM.
3. В общем случае состояния нейронов могут принимать вещественные значения.



Энергия и вероятности

- Введем понятие энергии для ограниченной машины Больцмана:

$$E(v, h) = - \sum_{i=1}^{N_v} a_i v_i - \sum_{j=1}^{N_h} b_j h_j - \sum_{i=1}^{N_v} \sum_{j=1}^{N_h} w_{ij} v_i h_j = -a^\top v - b^\top h - v^\top W h$$

Примечания:

1. Если вектора смещений внесены в матрицу связей между видимыми и скрытыми нейронами, то функция энергии имеет упрощенный вид.

$$E(v, h) = - \sum_{i=0}^{N_v} \sum_{j=0}^{N_h} w_{ij} v_i h_j = -v^\top W h$$

2. Если вектора смещений внесены в матрицу связей между видимыми и скрытыми нейронами, и состояния нейронов принимают значения 1 и -1 (1 соответствует состоянию "включен", -1 - состоянию "выключен"), то функция энергии отличается от предыдущей наличием дополнительного коэффициента.

$$E(v, h) = -\frac{1}{2} \sum_{i=0}^{N_v} \sum_{j=0}^{N_h} w_{ij} v_i h_j = -\frac{1}{2} v^\top W h$$

- Моделируемая совместная плотность вероятности представляет собой распределение Гиббса.

$$p(v, h) = \frac{1}{Z} e^{-E(v, h)},$$

где Z — это статсумма следующего вида (суммирование по всевозможным векторам v и h , L и S - количество образов видимых и скрытых векторов соответственно):

$$Z = \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} e^{-E(v^{(l)}, h^{(s)})}$$

- Полная вероятность вектора v (суммирование по всевозможным векторам h):

$$p(v) = \frac{1}{Z} \sum_{s=0}^{S-1} e^{-E(v, h^{(s)})}$$

Вычисление условных вероятностей

- Вероятность того, что при данном входном векторе v активируется скрытый нейрон $h_k = 1$:

$$\begin{aligned} p(h_k = 1|v) &= \frac{p(v, h_k = 1)}{p(v)} = \frac{p(v, h_k = 1)}{p(v, h_k = 0) + p(v, h_k = 1)} = \frac{e^{-E(v, h_k=1)}}{e^{-E(v, h_k=1)} + e^{-E(v, h_k=0)}} \\ &= \frac{1}{1 + e^{E(v, h_k=1) - E(v, h_k=0)}} = \frac{1}{1 + e^{-b_k - \sum_{i=1}^{N_v} w_{ik} v_i}} = \sigma(b_k + \sum_{i=1}^{N_v} w_{ik} v_i), \end{aligned}$$

где $\sigma(z) = \frac{1}{1+e^{-z}}$ - сигмоидальная функция.

- Вероятность того, что при данном скрытом векторе h активируется видимый нейрон $v_k = 1$:

$$\begin{aligned} p(v_k = 1|h) &= \frac{p(v_k = 1, h)}{p(h)} = \frac{e^{-E(v_k=1, h)}}{e^{-E(v_k=1, h)} + e^{-E(v_k=0, h)}} = \frac{1}{1 + e^{E(v_k=1, h) - E(v_k=0, h)}} \\ &= \frac{1}{1 + e^{-a_k - \sum_{j=1}^{N_h} w_{kj} h_j}} = \sigma(a_k + \sum_{j=1}^{N_h} w_{kj} h_j) \end{aligned}$$

- Поскольку при фиксированном входном векторе v состояния скрытых нейронов h_j независимы друг от друга (как и при фиксированном h состояния видимых нейронов v_i), то условные вероятности для векторов имеют вид:

$$p(h|v) = \prod_{j=1}^{N_h} p(h_j|v)$$

$$p(v|h) = \prod_{i=1}^{N_v} p(v_i|h)$$

Цель обучения

Состояния видимых нейронов можно наблюдать, поскольку видимые нейроны соответствуют входным данным, но состояния скрытых нейронов напрямую увидеть нельзя. Несмотря на это, опираясь на наблюдаемые состояния, можно сделать вероятностный вывод относительно скрытых состояний. После обучения модели можно делать вероятностные выводы (при помощи теоремы Байеса) о том, какими являются видимые нейроны, опираясь на скрытые. Таким образом, **цель** - обучить модель таким образом (настроить ее параметры), чтобы полученный в результате работы RBM вектор состояний видимых нейронов был близок к входному состоянию, из которого он был восстановлен. Множество весов RBM реализует совершенную модель, если оно приводит к такому же распределению вероятности состояния видимых элементов, что и подача входных векторов.

Данную цель можно достигнуть путем максимизации вероятности $p(v)$. Применяя метод максимального правдоподобия, вместо максимизации вероятности $p(v)$, можно максимизировать натуральный логарифм от нее $p(v)$. Максимум $p(v)$ соответствует максимуму $\ln p(v)$, поскольку при переходе от функции x к $\ln x$ монотонность не изменяется.

Максимизация вероятности $p(v)$

Найдем производную от натурального логарифма $p(v)$ по весу w_{ij} :

$$\begin{aligned}\frac{\partial \log(p(v))}{\partial w_{ij}} &= \frac{1}{p(v)} \frac{\partial p(v)}{\partial w_{ij}} = \frac{1}{p(v)} \frac{\partial}{\partial w_{ij}} \left(\frac{\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})}}{Z} \right) \\ &= \frac{1}{p(v)} \frac{\frac{\partial}{\partial w_{ij}} \left(\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})} \right) Z - \sum_{s=0}^{S-1} e^{-E(v, h^{(s)})} \frac{\partial Z}{\partial w_{ij}}}{Z^2},\end{aligned}$$

где

$$p(v) = \frac{1}{Z} \sum_{s=0}^{S-1} e^{-E(v, h^{(s)})},$$

$$Z = \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} e^{-E(v^{(l)}, h^{(s)})},$$

$$E(v, h) = - \sum_{i=1}^{N_v} a_i v_i - \sum_{j=1}^{N_h} b_j h_j - \sum_{i=1}^{N_v} \sum_{j=1}^{N_h} w_{ij} v_i h_j.$$

Вычислим производные:

$$\frac{\partial}{\partial w_{ij}} \left(\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})} \right) = \sum_{s=0}^{S-1} v_i h_j^{(s)} e^{-E(v, h^{(s)})}$$

$$\frac{\partial Z}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \left(\sum_{l=0}^{L-1} \sum_{s=0}^{S-1} e^{-E(v^{(l)}, h^{(s)})} \right) = \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} e^{-E(v^{(l)}, h^{(s)})}$$

Подставим полученные производные в формулу производной от логарифма плотности распределения вероятности входных данных.

$$\frac{\partial \ln(p(v))}{\partial w_{ij}} = \frac{1}{p(v)} \frac{\frac{\partial}{\partial w_{ij}} \left(\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})} \right) Z - \sum_{s=0}^{S-1} e^{-E(v, h^{(s)})} \frac{\partial Z}{\partial w_{ij}}}{Z^2} =$$

$$\begin{aligned}
&= \frac{Z}{\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})}} \frac{Z(\sum_{s=0}^{S-1} v_i h_j^{(s)} e^{-E(v, h^{(s)})}) - (\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})})(\sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} e^{-E(v^{(l)}, h^{(s)})})}{Z^2} = \\
&= \frac{\sum_{s=0}^{S-1} v_i h_j^{(s)} e^{-E(v, h^{(s)})}}{\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})}} - \frac{\sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} e^{-E(v^{(l)}, h^{(s)})}}{Z} = \\
&= \left(\frac{Z}{\sum_{s=0}^{S-1} e^{-E(v, h^{(s)})}} \right) \left(\sum_{s=0}^{S-1} v_i h_j^{(s)} \frac{e^{-E(v, h^{(s)})}}{Z} \right) - \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} \frac{e^{-E(v^{(l)}, h^{(s)})}}{Z} = \\
&= \sum_{s=0}^{S-1} v_i h_j^{(s)} \frac{p(v, h^{(s)})}{p(v)} - \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} p(v^{(l)}, h^{(s)}) = \\
&= \sum_{s=0}^{S-1} v_i h_j^{(s)} p(h^{(s)}|v) - \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} p(v^{(l)}, h^{(s)})
\end{aligned}$$

Таким образом, получаем формулу для вычисления производных функции ошибки по весам сети.

$$\frac{\partial \ln(p(v))}{\partial w_{ij}} = \sum_{s=0}^{S-1} v_i h_j^{(s)} p(h^{(s)}|v) - \sum_{l=0}^{L-1} \sum_{s=0}^{S-1} v_i^{(l)} h_j^{(s)} p(v^{(l)}, h^{(s)}) = E_{data}[v_i h_j] - E_{model}[v_i h_j],$$

где $E(\cdot)$ - математическое ожидание.

Аналогично, можно получить формулы для производных по смещениям:

$$\frac{\partial \ln(p(v))}{\partial a_i} = E_{data}[v_i] - E_{model}[v_i], \quad \frac{\partial \ln(p(v))}{\partial b_j} = E_{data}[h_j] - E_{model}[h_j]$$

Правила обновления параметров модели

Введем параметр η (learning rate), отвечающий за скорость обучения модели (он влияет на то, насколько сильно изменяются веса при обучении). Тогда получаем следующие правила обновления параметров модели (весов и смещений):

$$\Delta w_{ij} = \eta(E_{data}[v_i h_j] - E_{model}[v_i h_j])$$

$$\Delta a_i = \eta(E_{data}[v_i] - E_{model}[v_i])$$

$$\Delta b_j = \eta(E_{data}[h_j] - E_{model}[h_j])$$

Алгоритм обучения Contrastive Divergence (CD-k)

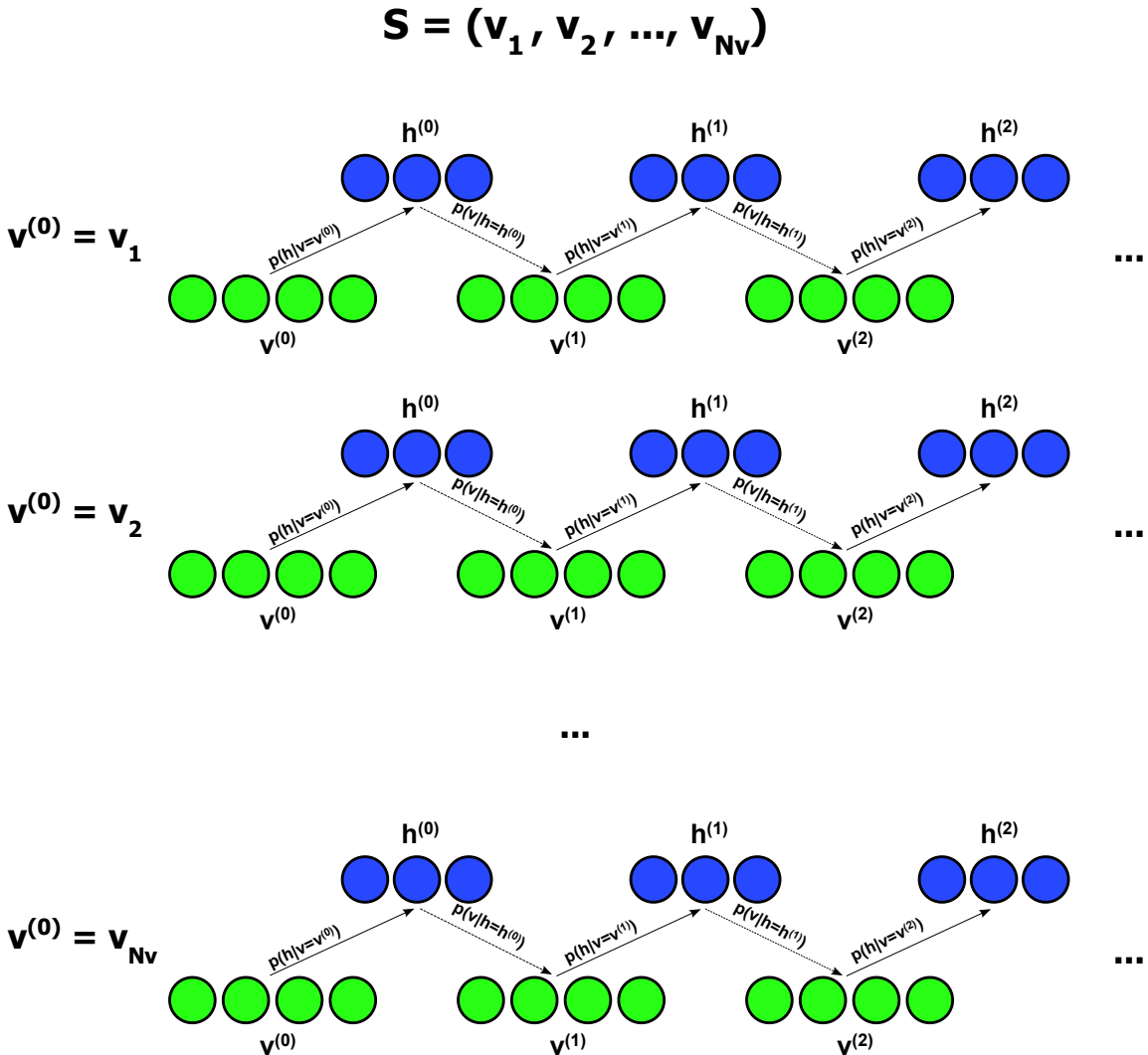
Алгоритм обучения имеет название **алгоритма контрастной дивергенции** (Contrastive Divergence, CD-k). Алгоритм основывается на процедуре **сэмплирования Гиббса** (Gibbs sampling) - процедура получения значения случайной величины на основании ее плотности распределения. Общая схема алгоритма имеет вид:

1. Состояние видимых нейронов x приравнивается к входному образу сети $v^0 = v^{(0)}$ ($v^{(0)}$ - элемент тренировочных данных).
2. Реализуется цикл по индексу t , состоящий из k итераций (k - параметр алгоритма). При этом осуществляется сбор статистики для вычисления математических ожиданий в формулах коррекции весов.

- A. Выводятся вероятности состояний скрытого слоя из распределения $p(h|v = v^t)$, из полученного распределения сэмплируется вектор h^t .
- B. Выводятся вероятности состояний видимого слоя из распределения $p(v|h = h^{t-1})$, из полученного распределения сэмплируется вектор v^t .

1. Выполняется коррекция весов сети и переход к следующему входному образу (элементу тренировочной выборки).

Иллюстрация алгоритма:



Псевдокод алгоритма приведен ниже ([ссылка на источник](#)).

Algorithm 1. k -step contrastive divergence

Input: RBM $(V_1, \dots, V_m, H_1, \dots, H_n)$, training batch S **Output:** gradient approximation Δw_{ij} , Δb_j and Δc_i for $i = 1, \dots, n$,
 $j = 1, \dots, m$

```
1 init  $\Delta w_{ij} = \Delta b_j = \Delta c_i = 0$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
2 forall the  $v \in S$  do
3    $v^{(0)} \leftarrow v$ 
4   for  $t = 0, \dots, k - 1$  do
5     for  $i = 1, \dots, n$  do sample  $h_i^{(t)} \sim p(h_i | v^{(t)})$ 
6     for  $j = 1, \dots, m$  do sample  $v_j^{(t+1)} \sim p(v_j | h^{(t)})$ 
7   for  $i = 1, \dots, n, j = 1, \dots, m$  do
8      $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | v^{(0)}) \cdot v_j^{(0)} - p(H_i = 1 | v^{(k)}) \cdot v_j^{(k)}$ 
9      $\Delta b_j \leftarrow \Delta b_j + v_j^{(0)} - v_j^{(k)}$ 
10     $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1 | v^{(0)}) - p(H_i = 1 | v^{(k)})$ 
```

Сэмплирование значений случайной величины из распределения

Пусть в результате вычисления сигмоидальной функции была получена вероятность 0.6723. Округлим её до двух знаков дробной части и получим число 0.67. То есть:

$$p(h|v_0) = \sigma(b + W^\top v)_0 = 0.6723 \approx 0.67$$

Теперь можно поступить следующим образом:

1. Заведём массив A из 100 элементов, в котором первые $0.67 * 100 = 67$ элементов будут единицами, а остальные $100 - 67 = 33$ элемента — нулями.
2. Сгенерируем псевдослучайное число следующего вида:

$$\xi \sim U[1, 100],$$

где $U[1, 100]$ — непрерывное равномерное распределение целых чисел на отрезке $[1, 100]$.

3. Будем считать, что (с учетом того, что индексы в массиве начинаются с нуля):

$$h_0 = A[\xi - 1]$$

Можно, например, округлять вероятность до трех знаков дробной части, и заводить массив на 1000 элементов, но это более затратно.