

Лабораторная работа № 6

Алгоритмы обучения СММ с дискретным пространством наблюдений

Цель работы

Научится оценивать параметры СММ с дискретным пространством наблюдений

Указания к работе

Для получения описания исследуемого процесса или объекта в виде СММ по имеющимся наблюдаемым последовательностям необходимо оценить параметры этой модели. Для этого решается задача обучения, состоящая в подборе параметров модели λ так, чтобы она правильно распознавала последовательность мультинаблюдений $O^* = \{O^1, O^2, \dots, O^K\}$, где K – это число наблюдаемых последовательностей.

Необходимо выбрать один из способов обучения:

1) распознать эти последовательности наблюдений, сравнить результаты распознавания $\tilde{Q}^1, \tilde{Q}^2, \dots, \tilde{Q}^K$, где \tilde{Q}^k – найденная последовательность скрытых состояний (с использованием, например, алгоритма Витерби), $k = \overline{1, K}$ с правильными ответами Q^1, Q^2, \dots, Q^K , вычислить в каком-либо смысле среднюю ошибку и минимизировать ее, варьируя λ ;

2) распознать последовательность мультинаблюдений и максимизировать функцию правдоподобия наблюдения последовательности O в предположении, что последовательность скрытых состояний найдена правильно. Значит необходимо

максимизировать $\prod_{k=1}^K P(O^k | \tilde{Q}^k)$, варьируя λ ;

3) максимизировать функцию правдоподобия наблюдений, т. е. максимизировать вероятность $L(O^* | \lambda) = \prod_{k=1}^K P(O^k | \lambda)$, варьируя параметры модели λ .

Заметим, что в качестве оптимизационного критерия можно использовать не только максимум правдоподобия, но и максимум взаимной информации или минимизировать различающую информацию.

Первый способ обучения – это обучение с учителем. Его можно проводить, например, методом градиентного спуска, и он хорош всем, кроме своей трудоемкости. Остальные два способа – это обучение без учителя, хотя во втором способе можно использовать учителя. Чаще всего применяется третий способ обучения (иногда с последующим дообучением другими способами), поскольку для него известен быстрый алгоритм. Это алгоритм, в общей ситуации называемый ЕМ (ЕМ – expectation maximization; максимизация ожидания) или, применительно к СММ, алгоритмом Баума-Велша. Данный алгоритм является итеративным и сходится, вообще говоря, не к глобальному максимуму правдоподобия, а к локальному. Далее будет более подробно рассмотрен третий способ обучения.

Определим вероятность того, что последовательность $Q = \{q_1, q_2, \dots, q_T\}$ скрытых состояний порождает последовательность $O = \{o_1, o_2, \dots, o_T\}$ наблюдений:

$$P(O|Q, \lambda) = P(o_1, o_2, \dots, o_T | q_1, q_2, \dots, q_T) = \prod_{t=1}^T P(o_t | q_t, \lambda) = \prod_{t=1}^T b_{q_t}(o_t).$$

Вероятность появления последовательности $Q = \{q_1, q_2, \dots, q_T\}$ вычисляется как:

$$P(Q | \lambda) = P(q_1, q_2, \dots, q_T) = \pi_{q_1} \prod_{t=1}^{T-1} a_{q_t q_{t+1}}.$$

Обе эти формулы следуют непосредственно из условий независимости в определении скрытой марковской модели, которые приведены ниже.

1 условие. Последовательность $Q = \{q_1, q_2, \dots, q_T\}$ случайных величин со значениями в S удовлетворяет следующему условию (свойство марковости – см. **Ошибка! Источник ссылки не найден.**):

$$P(q_t = s_{i_t} | q_{t-1} = s_{i_{t-1}}, \dots, q_1 = s_{i_1}) = P(q_t = s_{i_t} | q_{t-1} = s_{i_{t-1}}), \forall t, i_1, \dots, i_t.$$

2 условие. Каждой последовательности скрытых состояний $\{s_{i_1}, s_{i_2}, \dots, s_{i_T}\}$ ставится в соответствие последовательность случайных величин $O = \{o_1, o_2, \dots, o_T\}$.

При этом выполняется следующее условие:

$$P(o_t | o_1, \dots, o_{t-1}, o_{t+1}, \dots, o_T, q_1 = s_{i_1}, \dots, q_T = s_{i_T}) = P(o_t | q_t = s_{i_t}),$$

$\forall t, i_1, i_2, \dots, i_T$. Таким образом, вероятность появления некоторого наблюдения зависит только от того, в каком состоянии находится скрытый случайный процесс в данный момент.

Вероятность наблюдения последовательности O , порожденной последовательностью Q , равна:

$$P(O, Q | \lambda) = P(O | Q, \lambda) P(Q | \lambda) = \pi_{q_1} \prod_{t=1}^{T-1} a_{q_t q_{t+1}} \prod_{t=1}^T b_{q_t}(o_t).$$

Вероятность наблюдения последовательности O без каких-либо условий, по определению, равна:

$$P(O | \lambda) = \sum_{q_1, q_2, \dots, q_T} P(O, Q | \lambda) = \sum_{q_1, q_2, \dots, q_T} \left(\pi_{q_1} \prod_{t=1}^{T-1} a_{q_t q_{t+1}} \prod_{t=1}^T b_{q_t}(o_t) \right).$$

Заметим, что вычислительная сложность этой формулы – порядка $2TN^T$ операций. Для эффективного вычисления этой вероятности используют алгоритм

прямого-обратного прохода (forward-backward algorithm). Данный алгоритм основывается на методах динамического программирования.

Определим *forward-вероятность* следующим образом:

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, q_t = s_i | \lambda) = \sum_{q_1, q_2, \dots, q_t | q_t = s_i} P(o_1, o_2, \dots, o_t, q_1, q_2, \dots, q_t | \lambda),$$

т. е. это вероятность того, что данная последовательность наблюдений $\{o_1, o_2, \dots, o_t\}$ будет сгенерирована моделью λ и эта модель находится в состоянии s_i .

Для вычисления α необходимо провести следующие шаги.

1) инициализация:

$$\alpha_1(i) = \pi_i b_i(o_1), \quad i = \overline{1, N};$$

2) индукция:

$$\alpha_{t+1}(i) = b_i(o_{t+1}) \left[\sum_{j=1}^N \alpha_t(j) a_{ji} \right], \quad i = \overline{1, N}, \quad t = \overline{1, T-1};$$

3) завершение:

$$P(O | \lambda) = \sum_{i=1}^N \alpha_T(i).$$

Таким способом вероятность $P(O | \lambda)$, а значит и все вероятности $\alpha_t(i)$ при $i = \overline{1, N}$, $t = \overline{1, T}$, можно вычислить за порядка $3Tn'$ операций, где множитель n' – это количество ненулевых элементов матрицы вероятностей переходов.

Backward-вероятность:

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | q_t = s_i, \lambda),$$

т. е. это вероятность того, что последовательность наблюдений $\{o_{t+1}, o_{t+2}, \dots, o_T\}$, данная состоянием s_i , будет сгенерирована моделью λ .

Для вычисления β необходимо провести следующие шаги.

1) инициализация:

$$\beta_T(i) = 1, \quad i = \overline{1, N};$$

2) индукция:

$$\beta_t(i) = \sum_{j=1}^N \beta_{t+1}(j) b_j(o_{t+1}) a_{ij}, \quad i = \overline{1, N}, \quad t = \overline{1, T-1};$$

3) завершение:

$$P(O | \lambda) = \sum_{i=1}^N \alpha_1(i) \beta_1(i).$$

Кроме того, заметим:

$$P(O | \lambda) = \sum_{i=1}^N \alpha_t(i) \beta_t(i) \quad \forall t = \overline{1, T}.$$

В данных лабораторных работах, как мы уже отмечали, обучение СММ будем вести через максимизацию функции правдоподобия $L(O^* | \lambda)$, используя метод Баум-Велша. Фактически необходимо подобрать последовательность скрытых состояний Q к последовательности наблюдений O , т.е. решить задачу с недостающими (пропущенными) данными. Эту проблему решают при помощи ЕМ-алгоритма, который ориентирован на поиск максимума функции правдоподобия по параметрам ненаблюдаемой функции распределения для множества наблюдений, где данные неполны или имеются пропуски. В такой модели недостающие данные – это переменные, указывающие, из какого компонента смеси извлечен элемент данных.

Заметим, что вычисление значения функции правдоподобия проводится по формуле:

$$\begin{aligned} L(O^* | \lambda) &= \prod_{k=1}^K P(O^k | \lambda) = \prod_{k=1}^K \sum_{q_1, q_2, \dots, q_{T^k}} P(O^k, Q | \lambda) = \\ &= \prod_{k=1}^K \sum_{q_1, q_2, \dots, q_{T^k}} \left(\pi_{q_1} \prod_{t=1}^{T^k-1} a_{q_t q_{t+1}} \prod_{t=1}^{T^k} b_{q_t}(o_t^k) \right). \end{aligned}$$

Поиск параметров модели λ , обеспечивающих максимум правдоподобия $L(O^* | \lambda)$, традиционно формулируется как поиск точки минимума «функции ошибки» $E(O^* | \lambda) = -\ln L(O^* | \lambda)$. Зафиксируем обучающий набор O^* из K последовательностей наблюдений O^k длины T^k , $k = \overline{1, K}$ и начальный набор параметров модели λ^0 (любое допустимое значение). На каждом шаге итерации функция ошибки мажорируется функцией более простого вида, минимум которой единственен и находится явно:

$$\begin{aligned}
E(O^*, \lambda) - E(O^*, \lambda^0) &= - \sum_{k=1}^K \ln \frac{P(O^k | \lambda)}{P(O^k | \lambda^0)} = - \sum_{k=1}^K \ln \frac{P(O^k, Q | \lambda)}{P(O^k | \lambda^0)} = \\
&= - \sum_{k=1}^K \ln \sum_{q_1, q_2, \dots, q_{T^k}} \frac{P(O^k, Q | \lambda)}{P(O^k | \lambda^0)} = \\
&= - \sum_{k=1}^K \ln \sum_{q_1, q_2, \dots, q_{T^k}} \frac{P(Q | O^k, \lambda^0)}{P(Q | O^k, \lambda^0)} \frac{P(O^k, Q | \lambda)}{P(O^k | \lambda^0)} = \\
&= - \sum_{k=1}^K \ln \sum_{q_1, q_2, \dots, q_{T^k}} P(Q | O^k, \lambda^0) \frac{P(O^k, Q | \lambda)}{P(O^k | \lambda^0)} \leq \\
&\leq - \sum_{k=1}^K \sum_{q_1, q_2, \dots, q_{T^k}} P(Q | O^k, \lambda^0) \ln \frac{P(O^k, Q | \lambda)}{P(O^k | \lambda^0)} = \\
&= - \sum_{k=1}^K \sum_{q_1, q_2, \dots, q_{T^k}} P(Q | O^k, \lambda^0) \ln P(O^k, Q | \lambda) + \\
&+ \sum_{k=1}^K \sum_{q_1, q_2, \dots, q_{T^k}} P(Q | O^k, \lambda^0) \ln P(O^k | \lambda^0).
\end{aligned}$$

Неравенство в этой цепочке преобразований следует из выпуклости функции E и того, что условные вероятности неотрицательны и их сумма по всем возможным последовательностям Q длины T^k равна 1.

Обозначим первое слагаемое в этой формуле следующим образом:

$$G(O^*, \lambda^0, \lambda) = - \sum_{k=1}^K \sum_{q_1, q_2, \dots, q_{T^k}} P(Q | O^k, \lambda^0) \ln P(O^k, Q | \lambda).$$

Тогда полученное неравенство означает, что:

$$E(O^*, \lambda) \leq G(O^*, \lambda^0, \lambda) - G(O^*, \lambda^0, \lambda^0) + E(O^*, \lambda^0),$$

т. е. правая часть мажорируется функцией $E(O^*, \lambda)$ и совпадает с ней в точке λ^0 . Значит, если удастся найти точку λ глобального минимума правой части или, что то же самое, точку минимума $G(O^*, \lambda^0, \cdot)$ по всей области определения, то будет выполняться

неравенство $E(O^*, \lambda) \leq E(O^*, \lambda^0)$. Вычисление производных от функции $G(O^*, \lambda^0, \lambda)$ по параметрам показывает, что критическая точка $G(O^*, \lambda^0, \lambda)$ всегда единственна и является точкой минимума. Введем ряд обозначений, которые будут использованы в формулах для оценок параметров СММ.

Определим вероятность того, что в текущий момент t скрытый процесс СММ находится в состоянии s_i :

$$\gamma_t(i) = P(q_t = s_i | O, \lambda) = \frac{\alpha_t(i)\beta_t(i)}{P(O | \lambda)}, \quad i = \overline{1, N}, \quad t = \overline{1, T-1}.$$

Введем вероятность того, что в текущий момент t скрытый процесс СММ находится в состоянии s_i и в последующий момент он перейдет в состояние s_j :

$$\begin{aligned} \xi_t(i, j) &= P(q_t = s_i, q_{t+1} = s_j | O, \lambda) = \\ &= \frac{\alpha_t(i)a_{ij}b_j(o_{t+1})\beta_{t+1}(j)}{P(O | \lambda)}, \quad i, j = \overline{1, N}, \quad t = \overline{1, T-1}. \end{aligned}$$

На рисунке 1 поясняется, каким образом вычисляется совместная вероятность пребывания модели в момент времени t в состоянии s_i и в момент времени $t+1$ – в состоянии s_j , т. е. вероятность $\xi_t(i, j)$.

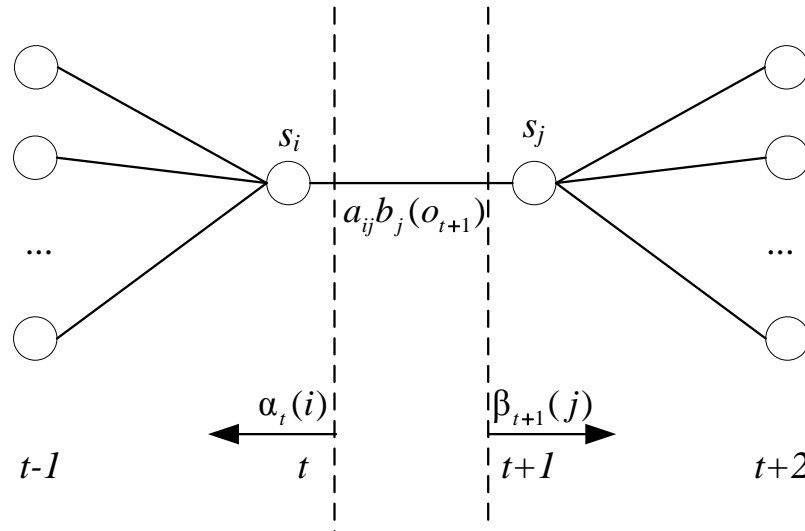


Рисунок 1 – Связь между соседними скрытыми состояниями s_i и s_j

Заметим, что $\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$.

Введем следующую вероятность того, что в текущий момент t скрытый процесс СММ находится в состоянии s_i и при этом m -ая компонента смеси порождает текущее наблюдение из наблюдаемой последовательности:

$$\gamma_t(i, m) = P(q_t = s_i, \mathcal{G}_{it} = m | O, \lambda),$$

где \mathcal{G}_{it} – случайная переменная, показывающая номер компоненты смеси в момент времени t для состояния i . Тогда:

$$\gamma_t(i, m) = \gamma_t(i) \frac{\tau_{im} g(o_t; \Theta_{im})}{\sum_{\tilde{m}=1}^{M_i} \tau_{i\tilde{m}} g(o_t; \Theta_{i\tilde{m}})}.$$

Для СММ с непрерывной вероятностью описания наблюдаемых состояний минимум функции $Q(O^*, \lambda^0, \lambda)$ достигается в точке λ^* со следующими координатами:

$$\pi_i^* = \frac{1}{K} \sum_{k=1}^K \gamma_1^{(k)}(i) = \frac{1}{K} \sum_{k=1}^K \frac{\alpha_1(i) \beta_1(i)}{P(O^k | \lambda)},$$

$$a_{ij}^* = \frac{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \xi_t^{(k)}(i, j)}{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i)} = \frac{\sum_{k=1}^K \frac{1}{P(O^k | \lambda)} \sum_{t=1}^{T^k-1} \alpha_t^{(k)}(i) a_{ij} b_j(o_{t+1}^{(k)}) \beta_{t+1}^{(k)}(j)}{\sum_{k=1}^K \frac{1}{P(O^k | \lambda)} \sum_{t=1}^{T^k-1} \alpha_t^{(k)}(i) \beta_t^{(k)}(i)},$$

$$\tau_{im}^* = \frac{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i, m)}{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i)},$$

$$\mu_{im}^* = \frac{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i, m) o_t^k}{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i, m)},$$

$$\sigma_{im}^{2*} = \frac{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i,m)(o_t^{(k)} - \mu_{im}^*)(o_t^{(k)} - \mu_{im}^*)^T}{\sum_{k=1}^K \sum_{t=1}^{T^k-1} \gamma_t^{(k)}(i,m)}.$$

С учетом введенных обозначений для СММ с дискретным пространством наблюдений новое приближение оценки элементов матрицы B будет находиться следующим образом:

$$b_{im}^* = \frac{\sum_{k=1}^K \sum_{\substack{npu \\ o_t^{(k)}=v_m}}^{T^k} \gamma_t^{(k)}(i)}{\sum_{k=1}^K \sum_{t=1}^{T^k} \gamma_t^{(k)}(i)}.$$

Непосредственно алгоритм Баум-Велша состоит из трех основных этапов.

1 этап. Forward-Backward алгоритм.

2 этап. Перевычисление параметров модели λ по формулам.

3 этап. Повторение 1-ого и 2-ого этапов, пока не будет достигнут порог сходимости ε , т. е. пока выполняется условие:

$$|L(O^* | \lambda)^{iter-1} - L(O^* | \lambda)^{iter}| > \varepsilon.$$

Начальные значения параметров A и π модели можно задавать произвольно, учитывая при этом вероятностные нормировки. Алгоритм обучения всегда сходится, при этом почти всегда к точке локального, а не глобального максимума функции правдоподобия $L(O^* | \lambda)$.

Во всех вышеописанных формулах вероятности умножаются друг на друга: т. е. числа, не превышающие 1 и имеющие типичные значения, обратные количеству состояний, умножаются в количестве, пропорциональном длине последовательности. Для длинных (длины порядка 100 элементов и более) последовательностей эти произведения меньше минимальных аппаратно реализуемых чисел типичных компьютеров. Таким образом, нужно либо программно реализовывать неограниченную точность вычислений, что связано с временными затратами, либо как-то масштабировать все промежуточные результаты, чтобы они не стремились к нулю. Методы масштабирования, почти не замедляющие обучение, хорошо известны.

Задание

1. Найти оценки параметров модели по данным, смоделированным в предыдущей лабораторной работе
2. Отобразить результаты оценивания в таблице:

Начальное приближение параметров модели	Оценки параметров модели	Достигнутая точность по параметрам (ρ_A, ρ_B)	Кол-во итераций ($iter$)	Достигнутая точность по значению невязки функции правдоподобия: $ \ln L(O \lambda)^{iter} - \ln L(O \lambda)^{iter-1} $

Начальное приближение параметров модели выбирать как минимум три раза:

- 1) Близкими к истинным параметрам
- 2) Равными истинным параметрам
- 3) Далекими от истинных параметров

Контрольные вопросы

1. Какие существуют подходы к обучению СММ?
2. Описание алгоритма Баум-Велша.
3. Что такое прямые вероятность? Как они используются в алгоритме Баум-Велша?
4. Что такое обратные вероятность? Как они используются в алгоритме Баум-Велша?
5. Чем отличается алгоритм Баум-Велша для СММ с дискретным пространством наблюдений от алгоритма для СММ с непрерывным пространством наблюдений