

Projekt 2 - Metody klasyfikacji i redukcji wymiaru

Klaudia Weigel

1 Ogólna motywacja

Niech $\{\mathbf{x}_i\}$, $i = 1, \dots, n$ będzie zbiorem niezależnych obserwacji o takim samym rozkładzie, oraz niech θ oznacza nieznany parametr/parametry tego rozkładu. Możemy wyestymować ten parametr za pomocą estymatora największej wiarygodności:

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \prod_{i=1}^n p_{\theta}(\mathbf{x}_i)$$

Funkcja wiarygodności, może okazać się trudna lub niemożliwa do wyznaczenia. W takim przypadku możemy skorzystać z algorytmu EM (Expectation Maximization).

2 Specyfikacja problemu

Mamy dane ciągi $\mathbf{x}_1, \dots, \mathbf{x}_k$ określone: $\mathbf{x}_1 = (x_{11}, \dots, x_{1w})$, gdzie każdy element $x_{1i} \in \{A, C, G, T\} \equiv \{1, 2, 3, 4\}$. Każdy z ciągów jest realizacją zmiennej losowej $X = (X_1, \dots, X_w)$.

Rozważmy następujące rozkłady:

1. Dla $\boldsymbol{\theta}_j = (\theta_{1,j}, \dots, \theta_{4,j})^T$, określamy:

$$P(X_j = a) = \theta_{a,j}, \quad a = 1, \dots, 4$$

Zatem $\theta_{a,j}$ jest prawdopodobieństwem, że na j -tej pozycji znajduje się litera a .
Czyli prawdopodobieństwo otrzymania ciągu $\mathbf{x}_1 = (x_{11}, \dots, x_{1w})$ jest równe:

$$P(\mathbf{x}_1; \boldsymbol{\theta}) = \prod_{i=1}^w \theta_{x_{1j}, j}$$

2. Dla $\boldsymbol{\theta}^b = (\theta_1^b, \theta_2^b, \theta_3^b, \theta_4^b)$:

$$P(X_j = a) = \theta_a^b, \quad a = 1, \dots, 4$$

Kolejne litery losowane są niezależnie, bez względu na pozycję.
Prawdopodobieństwo otrzymania ciągu $\mathbf{x}_1 = (x_{11}, \dots, x_{1w})$ jest równe:

$$P(\mathbf{x}_1; \boldsymbol{\theta}^b) = \prod_{i=1}^w \theta_{x_{1j}}^b$$

Założmy teraz, że losowo wybieramy liczbę 0 albo 1, gdzie prawdopodobieństwo wylosowania 1 to α , a prawdopodobieństwo wylosowania 0 to $1 - \alpha$. Czyli, oznaczając przez Z_i zmienną losową, reprezentującą wynik i -tego rzutu:

$$P(Z_i = 1) = \alpha, \quad P(Z_i = 0) = 1 - \alpha, \quad i = 1, \dots, k$$

Na podstawie tego, jaka wypadła liczba wybieramy z jakiego rozkładu losujemy \mathbf{x}_i . Niech $Z_i = z_i$, wtedy jeśli $z_i = 1$, to \mathbf{x}_i losujemy z rozkładu opisanego w punkcie 1 powyżej, a jeśli $z_i = 0$, to \mathbf{x}_i losujemy z rozkładu opisanego w 2.

Chcemy wyestymować parametry $\boldsymbol{\theta}$ oraz $\boldsymbol{\theta}^b$. Gdyby $\mathbf{z} = (z_1, \dots, z_k)$ było znane, moglibyśmy bezpośrednio wyznaczyć estymator największej wiarygodności ze wzoru:

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^b) = \prod_{i=1}^k z_i(P(\mathbf{x}_i; \boldsymbol{\theta}) + (1 - z_i)P(\mathbf{x}_i; \boldsymbol{\theta}^b))$$

Zakładamy jednak, że nie znamy \mathbf{z} , do estymacji $\boldsymbol{\theta}$ oraz $\boldsymbol{\theta}^b$ posłuży nam algorytm EM.

3 Algorytm EM

3.1 Ogólny opis algorytmu EM

Algorytm EM iteracyjnie szuka maksimum funkcji wiarygodności, w celu wyznaczenia estymatorów nieznanymi parametrów rozkładu. Jego stosowanie jest przydatne w przypadkach kiedy bezpośrednia estymacja jest trudna bądź niemożliwa, na przykład kiedy w modelu probabilistycznym występują ukryte zmienne (jak w naszym przypadku, gdzie ukrytymi zmiennymi są Z_1, \dots, Z_k). Algorytm EM składa się z dwóch kroków:

1. Krok E(xpectation)

Znaleźć dolne ograniczenie dla logarytmu funkcji wiarygodności dla obecnych parametrów.

2. Krok M(aximization)

Zoptymalizować dolną granicę i odpowiednio uaktualnić parametry.

Czyli przy założeniu, że nasze obserwacje to $\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_n$, $\mathbf{x}_i \in \mathbb{R}^d$, realizacje zmiennych ukrytych $Z_i \in \{1, \dots, M\}$ to $\mathbf{z} = (z_1, \dots, z_n)$, a nieznanne parametry to $\boldsymbol{\theta}$, chcemy znaleźć:

$$\mathcal{L}(\boldsymbol{\theta}) = \log P(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta}) = \sum_{i=1}^n \log \sum_{z_i} P(\mathbf{x}_i, z_i; \boldsymbol{\theta})$$

Założmy, że dla każdego $i \in \{1, \dots, n\}$, Q_i jest pewnym rozkładem zmiennej losowej Z_i ($Q_i(z) = P(Z_i = z)$, $\sum_z Q_i(z) = 1$, $Q_i(z) \geq 0$). Wtedy:

$$\sum_{i=1}^n \log \sum_{z_i} P(\mathbf{x}_i, z_i; \boldsymbol{\theta}) = \sum_{i=1}^n \log \sum_{z_i} Q_i(z_i) \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{Q_i(z_i)} \geq \sum_{i=1}^n \sum_{z_i} Q_i(z_i) \log \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{Q_i(z_i)}$$

Ostatnia nierówność wynika z nierówności Jensena, która mówi, że jeśli funkcja jest wklęsła to $f(E[X]) \geq E[f(X)]$ (równość otrzymujemy kiedy X jest stałą). Ponieważ logarytm jest funkcją wklęsłą, a wartość:

$$\sum_{z_i} Q_i(z_i) \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{Q_i(z_i)}$$

jest wartością oczekiwaną funkcji $h(z_i) = \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{Q_i(z_i)}$ zmiennej losowej Z_i z rozkładu Q_i . Teraz chcemy dobrać tak Q_i aby ograniczenie było jak najlepsze, czyli zamiast nierówności otrzymać równość. Czyli musi być spełnione:

$$\frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{Q_i(z_i)} = c$$

Gdzie c jest stałą. Wystarczy zatem, że wybierzemy Q_i proporcjonalne do $P(\mathbf{x}_i, z_i; \boldsymbol{\theta})$. Stąd i z faktu, że $\sum_z Q_i(z) = 1$, możemy zatem określić Q_i jako:

$$Q_i(z_i) = \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{\sum_z P(\mathbf{x}_i, z; \boldsymbol{\theta})} \quad (1)$$

$$= \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{P(\mathbf{x}_i; \boldsymbol{\theta})} \quad (2)$$

$$= P(z_i | \mathbf{x}_i; \boldsymbol{\theta}) \quad (3)$$

$$= P(Z_i = z_i | \mathbf{x}_i; \boldsymbol{\theta}) \quad (4)$$

Czyli tak dobrana funkcja Q_i daje nam dolną granicę logarytmu funkcji wiarygodności. W kroku E, obliczymy Q_i dla obecnych parametrów. W kroku M zoptymalizujemy dolną granicę ze względu na $\boldsymbol{\theta}$:

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \sum_{z_i} Q_i(z_i) \log \frac{P(\mathbf{x}_i, z_i; \boldsymbol{\theta})}{Q_i(z_i)}$$

Ponieważ szukamy maksimum tylko po $\boldsymbol{\theta}$, to powyższy wzór można zapisać:

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \sum_{z_i} Q_i(z_i) \log P(\mathbf{x}_i, z_i; \boldsymbol{\theta})$$

Czyli algorytm EM wygląda następująco:

Algorithm 1 EM Algorithm

Require: Observations $\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_n$, $\mathbf{x}_i \in \mathbb{R}^d$

- 1: Set: $t = 0$. Initialize $\boldsymbol{\theta}^{(0)}$
 - 2: **E(xpectation step)**. Compute for each i $Q_i(z_i) = P(z_i|\mathbf{x}_i; \boldsymbol{\theta})$
 - 3: **M(aximization step)**.
 - 4: Compute $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \sum_{i=1}^n \sum_{z_i} Q_i(z_i) \log P(\mathbf{x}_i, z_i; \boldsymbol{\theta})$
 - 5: Find $\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$
 - 6: **Repeat until converged**
-

3.2 Algorytm EM dla danych

Mając ogólny opis algorytmu EM, musimy teraz wyprowadzić wykorzystywane w nim wzory dla problemu opisanego w sekcji 2. Mamy więc:

$$Q_i^{(t)}(0) = P(Z_i = 0|\mathbf{x}_i; \boldsymbol{\Theta}^{(t)}) \quad (5)$$

$$= \frac{P(Z_i = 0)P(\mathbf{x}_i|Z_i = 0; \boldsymbol{\Theta}^{(t)})}{P(\mathbf{x}_i; \boldsymbol{\Theta}^{(t)})} \quad (6)$$

$$= \frac{P(Z_i = 0)P(\mathbf{x}_i|Z_i = 0; \boldsymbol{\Theta}^{(t)})}{\alpha P(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}) + (1 - \alpha)P(\mathbf{x}_i; \boldsymbol{\theta}^{b,(t)})} \quad (7)$$

$$\frac{(1 - \alpha) \prod_{i=1}^w \theta_{x_{ij}}^{b,(t)}}{\alpha \prod_{j=1}^w \theta_{x_{1j},j}^{(t)} + (1 - \alpha) \prod_{j=1}^w \theta_{x_{ij}}^{b,(t)}} \quad (8)$$

$$Q_i^{(t)}(1) = 1 - Q_i^{(t)}(0) \quad (9)$$

W kroku M musimy wyznaczyć:

$$\boldsymbol{\Theta}^{(t+1)} = \arg \max_{\boldsymbol{\Theta}} \sum_{i=1}^k \sum_{z_i} Q_i^{(t)}(z_i) \log P(\mathbf{x}_i, z_i; \boldsymbol{\Theta}) \quad (10)$$

$$= \arg \max_{\boldsymbol{\Theta}} \sum_{i=1}^k \left[Q_i^{(t)}(0)(1 - \alpha) \sum_{j=1}^w \log \theta_{x_{ij}}^b + Q_i^{(t)}(1)\alpha \sum_{j=1}^w \log \theta_{x_{ij},j} \right] \quad (11)$$

$$= \arg \max_{\boldsymbol{\Theta}} Q_1(\boldsymbol{\theta}^b) + Q_2(\boldsymbol{\theta}) \quad (12)$$

Ponieważ:

$$\boldsymbol{\Theta}^{(t+1)} = (\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\theta}^{b,(t+1)})$$

$$\boldsymbol{\theta}^{b,(t+1)} = \arg \max_{\boldsymbol{\theta}^b} \mathcal{Q}(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{(t+1)})$$

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{(t+1)})$$

Czyli szukamy maksimum $Q_1(\boldsymbol{\theta}^b)$ ze względu na $\boldsymbol{\theta}^b$, oraz $Q_2(\boldsymbol{\theta})$ ze względu na $\boldsymbol{\theta}$.
Dla $Q_1(\boldsymbol{\theta}^b)$:

$$\begin{aligned} & \text{maximize } Q_1(\boldsymbol{\theta}^b) \\ & \text{subject to } \theta_1^b + \theta_2^b + \theta_3^b + \theta_4^b = 1 \end{aligned}$$

Po rozwiązaniu metodą mnożników Lagrange'a otrzymujemy:

$$\theta_1^b = \frac{\sum_{i=1}^k Q_i^{(t)}(0) |\{j : x_{ij} = 1\}|}{w \sum_{i=1}^k Q_i^{(t)}(0)} \quad (13)$$

$$\theta_2^b = \frac{\sum_{i=1}^k Q_i^{(t)}(0) |\{j : x_{ij} = 2\}|}{w \sum_{i=1}^k Q_i^{(t)}(0)} \quad (14)$$

$$\theta_3^b = \frac{\sum_{i=1}^k Q_i^{(t)}(0) |\{j : x_{ij} = 3\}|}{w \sum_{i=1}^k Q_i^{(t)}(0)} \quad (15)$$

$$\theta_4^b = \frac{\sum_{i=1}^k Q_i^{(t)}(0) |\{j : x_{ij} = 4\}|}{w \sum_{i=1}^k Q_i^{(t)}(0)} \quad (16)$$

Postępując analogicznie dla $Q_2(\boldsymbol{\theta})$, otrzymujemy:

$$\theta_{t_1 t_2} = \frac{\sum_{i=1, x_{it_2}=t_1}^k Q_i^{(t)}(1)}{\sum_{i=1}^k Q_i^{(t)}(1)} \quad (17)$$

4 Implementacja w Pythonie

Program został napisany w Pythonie wersji 3.7.1.

4.1 Generowanie danych

Na podstawie danych $\boldsymbol{\theta}$, $\boldsymbol{\theta}^b$, w , k oraz α możemy wysymulować macierz \mathbf{X} , tak jak zostało to opisane w specyfikacji problemu. Program generujący dane należy uruchomić poleceniem:

```
python3 NAZWA_PROGRAMU --params param_file.json --output generated_data.json
```

Gdzie `params_file` to plik zawierający dane, a `generated_data.json` to plik, do którego zostanie zapisana wygenerowana macierz.

4.2 Estymacja parametrów

Program do estymacji parametrów $\boldsymbol{\theta}$ i $\boldsymbol{\theta}^b$ uruchamiany jest poleceniem:

```
python3 NAZWA_PROGRAMU --input generated_data.json -- output estimated_params.json
```

4.3 Działanie algorytmu

Algorytm EM wymaga początkowej inicjalizacji szukanych parametrów. W tym przypadku macierze zostały zainicjalizowane losowo. Jakość wyestymowanych parametrów została obliczona za pomocą MSE (Mean Squared Error). Przykładowe wyniki otrzymane dla różnych w i k oraz przybliżenia dla θ^b :

$w = 50, k = 100, \alpha = 0.3, \text{MSE Theta} = 0.003657, \text{MSE ThetaB} = 0.000477$

Original ThetaB = [0.25 0.25 0.25 0.25]

Estimated ThetaB = [0.26344315 0.2782388 0.22371683 0.23460122]

$w = 50, k = 120, \alpha = 0.4, \text{MSE Theta} = 0.008031, \text{MSE ThetaB} = 0.005872$

Original ThetaB = [0.39759036 0.03212851 0.37349398 0.19678715]

Estimated ThetaB = [0.32628343 0.13392961 0.29670589 0.24308106]

$w = 50, k = 500, \alpha = 0.3, \text{MSE Theta} = 0.003166, \text{MSE ThetaB} = 0.009603$

Original ThetaB = [0.15533981 0.22815534 0.48058252 0.13592233]

Estimated ThetaB = [0.20108117 0.24352806 0.31923179 0.23615899]

Iteracje algorytmu wykonujemy dopóki różnice w logarytmie funkcji wiarygodności stana się nieznaczne.