

ALGORITHMEN UND DATENSTRUKTUREN

ÜBUNG 14: EM-ALGORITHMUS

Eric Kunze
eric.kunze@tu-dresden.de

TU Dresden, 26.01.2022

Aufgabe 1

WAHRSCHEINLICHKEITSTHEORIE

Wir betrachten ein Zufallsexperiment (X, p) mit

- ► Ergebnismenge X und
- einer Funktion $p: X \to [0,1]$ mit $\sum_{x \in X} p(x) = 1$ (Wahrscheinlichkeitsverteilung von X)

Die Menge aller Wahrscheinlichkeitsverteilungen über X sei $\mathcal{M}(X)$. Jede Teilmenge $\mathcal{M} \subseteq \mathcal{M}(X)$ heißt **Wahrscheinlichkeitsmodell**.

Ein Wahrscheinlichkeitsmodell \mathcal{M} heißt **beschränkt**, falls $\mathcal{M} \neq \mathcal{M}(X)$; andernfalls unbeschränkt.

Führen wir nun zwei Zufallsexperimente nacheinander aus und nehmen dabei an, dass die beiden Experimente unabhängig voneinander sind. Folge das erste Experiment einer Verteilung $p_1 \in \mathcal{M}(X_1)$ und das zweite Experiment einer Verteilung $p_2 \in \mathcal{M}(X_2)$, dann ist $p_1 \times p_2 \in \mathcal{M}(X_1 \times X_2)$ eine Verteilung auf der Ergebnismenge $X_1 \times X_2$ unseres zweistufigen Experiments:

$$(p_1 \times p_2)(a,b) = p_1(a) * p_2(b).$$

"Einzelwahrscheinlichkeiten multiplizieren / erste Pfadregel"

KORPORA UND KORPUSWAHRSCHEINLICHKEITEN

Oftmals wissen wir aber die zugrundeliegende Verteilung nicht, sondern können lediglich die Ergebnisse des Experiments wahrnehmen. Zählen wir diese Beobachtungen, dann nennen wir das einen X-Korpus modelliert durch eine Funktion $h:X\to\mathbb{R}^{\geq 0}$. Man definiert die Korpuswahrscheinlichkeit / Likelihood von h unter einer Verteilung p als

$$L(h,p) = \prod_{x \in X} p(x)^{h(x)}.$$

Nun kennen wir aber die Verteilung *p* nicht und müssen sie daher aus den beobachteten Daten schätzen. Dies macht der **Maximum-Likelihood-Schätzer** (MLE)

$$mle(h, \mathcal{M}) = \underset{p \in \mathcal{M}}{arg \, max} \, L(h, p).$$

Solange das Modell unbeschränkt gewählt wird, d.h. es werden alle Verteilungen über *X* zugelassen, dann wird der MLE zur relativen Häufigkeit von *h*.

UNVOLLSTÄNDIGE DATEN

Bisher sind wir davon ausgegangen, dass die Daten stets vollständig waren, d.h. wir konnten jedes Ergebnis beobachten. In der Realität können aber oftmals nur Gruppen von Ergebnissen beobachtet werden; z.B. gewinne oder verliere ich bei einem Spiel. Wir wissen aber nicht, welches Ergebnis genau erzielt wurde.

Sei Y die Menge der Beobachtungen. Die Beobachtungsfunktion yield: $X \rightarrow Y$ ordnet jedem Ergebnis seine Beobachtung zu. Die Umkehrabbildung ordnet dann jeder Beobachtung eine Menge von möglichen Ergebnissen zu, die zu dieser Beobachtung führen, d.h.

$$A: Y \to \mathcal{P}(X)$$
 mit $A(y) = \{x \in X : yield(x) = y\}$.

Diese Funktion heißt Analysator.

Sei h ein Y-Korpus, d.h. h zählt Beobachtungen (nicht Ergebnisse). Die **Korpuswahrscheinlichkeit** / **Likelihood** von h unter einer Verteilung p ist

$$L(h,p) = \prod_{y \in Y} \left(\sum_{x \in A(y)} p(x) \right)^{h(y)}.$$

Der MLE bleibt wie er war: $mle(h, \mathcal{M}) = arg \max_{p \in \mathcal{M}} L(h, p)$.

AUFGABE 1

Bestimmen Sie für die folgenden Szenarien die Menge X der Ergebnisse und die Menge Y der Beobachtungen. Bestimmen Sie außerdem den Analysator.

- (a) Werfen zweier unabhängiger Münzen. Sie können nur beobachten, ob beide Münzen dieselbe oder verschiedene Seiten zeigen.
- (b) Werfen zweier Würfel, wobei Sie nur die Summe der Augenzahlen beobachten.
- (c) Zwei Spieler spielen Schere-Stein-Papier. Sie beobachten lediglich, welcher Spieler gewonnen hat bzw. ob das Spiel unentschieden ausging.

AUFGABE 1

(a) zweimaliger Münzwurf – Beobachtung der Gleichheit

$$X = \{K, Z\} \times \{K, Z\}$$
 und $Y = \{gleich, ungleich\}$

Der Analysator ordnet jeder Beobachtung $y \in Y$ die Menge der Ergebnisse aus X zu, die zur Beobachtung y führen, also

$$A(\mathsf{gleich}) = \{(K, K), (Z, Z)\}$$
$$A(\mathsf{ungleich}) = \{(K, Z), (Z, K)\}.$$

(b) zweimaliger Würfelwurf - Beobachtung der Augensumme

$$X = \{1, \dots, 6\} \times \{1, \dots, 6\} \quad \text{und} \quad Y = \{2, \dots, 12\}$$
 Analysator: $A(x) = \{(i, j) \in X : i + j = x\}$, d.h. konkret
$$A(2) = \{(1, 1)\}$$

$$A(3) = \{(1, 2), (2, 1)\}$$

$$A(4) = \{(1, 3), (3, 1), (2, 2)\}$$

$$\vdots$$

 $A(12) = \{(6,6)\}$

AUFGABE 1

(c) Schere, Stein, Papier - Beobachtung des Gewinners

```
X = \{Schere, Stein, Papier\}^2

Y = \{Spieler1, Spieler2, Unentschieden\}
```

Analysator:

```
\begin{split} A(\mathsf{Spieler1}) &= \{(\mathsf{Schere}, \mathsf{Papier}), (\mathsf{Stein}, \mathsf{Schere}), (\mathsf{Papier}, \mathsf{Stein})\} \\ A(\mathsf{Spieler2}) &= \{(\mathsf{Papier}, \mathsf{Schere}), (\mathsf{Schere}, \mathsf{Stein}), (\mathsf{Stein}, \mathsf{Papier})\} \\ A(\mathsf{Unentschieden}) &= \{(\mathsf{Papier}, \mathsf{Papier}), (\mathsf{Stein}, \mathsf{Stein}), (\mathsf{Schere}, \mathsf{Schere})\} \end{split}
```

Aufgabe 2

PROBLEM: ERGEBNISSE VS. BEOBACHTUNGEN

Gegeben: ein Y-Korpus von Beobachtungen

Gesucht: eine Verteilung $p \in \mathcal{M} \subseteq \mathcal{M}(X)$ auf *Ergebnissen*

Problem: Wir wollen die Verteilung der *Ergebnisse* schätzen. Um den MLE nutzen zu können, brauchen wir dafür einen X-Korpus. Jedoch ist uns nur ein Y-Korpus gegeben, da wir nur *Beobachtungen* wahrnehmen können.

Ausweg: Erweiterung des Y-Korpus h zu einem X-Korpus h_1 auf vollständigen Daten

$$h_i(x) = h(\text{yield}(x)) \cdot \frac{p_{i-1}(x)}{\sum_{x' \in A(\text{yield}(x))} p_{i-1}(x')}$$
 für alle $x \in X$

Dazu benötigen wir eine gewisse Vorkenntnis mit der Verteilung p_{i-1} , die wir aus dem vorherigen Iterationsschritt bzw. einer initialen Vermutung bekommen.

Dies ist der E-Schritt des EM-Algorithmus.

EM-ALGORITHMUS

Input: Analysator $A: Y \to \mathcal{P}(X)$, Modell $\mathcal{M} \subseteq \mathcal{M}(X)$,

- Y-Korpus h von Beobachtungen,
- ▶ initiale Wahrscheinlichkeitsverteilung p₀

Eine Iteration des Algorithmus besteht aus den folgenden beiden Schritten:

E-Schritt Expectation

Bestimmte die versteckten Eigenschaften mithilfe der Parameter aus der vorherigen Iteration.

$$h_i(x) = h(\mathsf{yield}(x)) \cdot \frac{p_{i-1}(x)}{\sum_{x' \in A(\mathsf{yield}(x))} p_{i-1}(x')}$$

M-Schritt Maximization

Bestimmte die neuen Parameter mithilfe der vollständigen Eigenschaften aus dem E-Schritt.

$$p_i = \underset{p \in \mathcal{M}}{\operatorname{arg\,max}} L(h_i, p)$$

EIN WEITERES PROBLEM?

Problem: Im M-Schritt bleibt ein MLE zu berechnen, der für beschränkte Modelle $\mathcal{M} \neq \mathcal{M}(X)$ nicht leicht zu bekommen ist. Eine Vereinfachung ist aber für mehrstufige Zufallsversuche möglich, die als *unabhängig* vorausgesetzt werden.

Haben wir zwei unabhängige Teilversuche mit Ergebnismengen X_1 und X_2 , dann soll die Verteilung die Unabhängigkeit widerspiegeln und wir verwenden daher das Modell

$$\mathcal{M} = \left\{ p^1 \times p^2 : p^1 \in \mathcal{M}(X_1), p^2 \in \mathcal{M}(X_2) \right\} \neq \mathcal{M}(X_1 \times X_2),$$

der MLE wird also nicht zur relativen Häufigkeit auf $X_1 \times X_2$.

Ausweg: Die Unabhängigkeit rettet uns: wir erhalten den MLE bzgl. \mathcal{M} als unabhängiges Produkt der relativen Häufigkeiten *auf den Teilexperimenten*. Dazu ist ein $(X_1 \times X_2)$ -Korpus h nicht unbedingt hilfreich. Wir brauchen vielmehr die Teilkorpora h^1 auf X_1 und h^2 auf X_2 . Dann gilt

$$\mathsf{mle}(h,\mathcal{M}) = \mathsf{rfe}(h^1) \times \mathsf{rfe}(h^2) \neq \mathsf{rfe}(h).$$

Problem: Wie bekommen wir h^1 und h^2 ? — **Ausweg:** Marginalisierung

MARGINALISIERUNG

Wir betrachten die zwei Ergebnismengen X_1 und X_2 . Das Modell sei gegeben durch das unabhängige Produkt der Modelle auf X_1 und der Modelle auf X_2 , d.h. $\mathcal{M} = \left\{ p^1 \times p^2 : p^1 \in \mathcal{M}(X_1), p^2 \in \mathcal{M}(X_2) \right\}$. Weiter sei h ein $X_1 \times X_2$ -Korpus. Die Teilkorpora h^1 auf X_1 und h^2 auf X_2 erhalten wir durch **Marginalisierung**

$$h^{1}(x_{1}) = \sum_{x_{2} \in X_{2}} h(x_{1}, x_{2})$$
 für alle $x_{1} \in X_{1}$
 $h^{2}(x_{2}) = \sum_{x_{1} \in X_{1}} h(x_{1}, x_{2})$ für alle $x_{2} \in X_{2}$

Die Summen entsprechen dabei gerade Zeilen- bzw. Spaltensummen, wenn man h in einer Tabelle notiert.

$X_1 \setminus X_2$	α		ω	
а	$h(a, \alpha)$		$h(a,\omega)$	$h^1(a)$
:	:	٠.	÷	:
z	$h(z, \alpha)$		$h(z,\omega)$	$h^1(z)$
	$h^2(\alpha)$		$h^2(\omega)$	

AUFGABE 2 — TEIL (A)

Das Spiel wird gewonnen, wenn beide Münzen auf der gleichen Seite landen.

Damit ist der Analysator A: {Gewinn, keinGewinn} $\rightarrow \mathcal{P}(X)$ gegeben durch

$$A(\mathsf{Gewinn}) = \{x \in X : \mathsf{yield}(x) = \mathsf{Gewinn}\}$$

$$= \{(K, K), (Z, Z)\}$$

$$A(\mathsf{keinGewinn}) = \{x \in X : \mathsf{yield}(x) = \mathsf{keinGewinn}\}$$

$$= \{(K, Z), (Z, K), (R, K), (R, Z)\}$$

AUFGABE 2 — TEIL (B)

Wir können nur die Beobachtungen Gewinn und keinGewinn feststellen.

Wir spielen das Spiel 24 Mal und gewinnen 6 Mal. Gesucht ist nun der Y-Korpus h, d.h. wie oft beobachten wir die Ereignisse Gewinn und keinGewinn.

$$h(Gewinn) = 6$$
 $h(keinGewinn) = 18$

AUFGABE 2 — TEIL (C)

Gegeben ist nun eine initiale Wahrscheinlichkeitsverteilung $q_0=q_0^1\times q_0^2$ über den vollständigen Daten mit

$$q_0^1(K) = \frac{2}{5}$$

$$q_0^2(K) = \frac{1}{3}$$

$$q_0^1(R) = \frac{1}{5}$$

$$\Rightarrow q_0^1(Z) = 1 - q_0^1(K) - q_0^1(R) = \frac{2}{5}$$

$$q_0^2(Z) = 1 - q_0^1(K) = \frac{2}{3}$$

Mit dem unabhängigen Produkt erhalten wir

$$\begin{split} q_0(K,K) &= q_0^1(K) \cdot q_0^2(K) = \frac{2}{15} & q_0(K,Z) &= q_0^1(K) \cdot q_0^2(Z) = \frac{4}{15} \\ q_0(Z,K) &= q_0^1(Z) \cdot q_0^2(K) = \frac{2}{15} & q_0(Z,Z) &= q_0^1(Z) \cdot q_0^2(Z) = \frac{4}{15} \\ q_0(R,K) &= q_0^1(R) \cdot q_0^2(K) = \frac{1}{15} & q_0(R,Z) &= q_0^1(R) \cdot q_0^2(Z) = \frac{2}{15} \end{split}$$

AUFGABE 2 — TEIL (C)

E-Schritt: Erweiterung von h auf h_1 mit folgender Formel:

$$h_1(x) = h(\text{yield}(x)) \cdot \frac{q_0(x)}{\sum\limits_{x' \in A(\text{yield}(x))} q_0(x')}$$

Damit ergibt sich dann zum Beispiel für das Ergebnis (K, K)

$$h_{1}(K,K) = h(Gewinn) \cdot \frac{q_{0}(K,K)}{\sum\limits_{x' \in \{(K,K),(Z,Z)\}} q_{0}(x')}$$

$$= h(Gewinn) \cdot \frac{q_{0}(K,K)}{q_{0}(K,K) + q_{0}(Z,Z)}$$

$$= 6 \cdot \frac{\frac{2}{15}}{\frac{2}{15} + \frac{4}{15}}$$

$$= 2$$

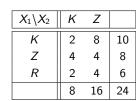
Mit gleicher Rechnung erhalten wir für die restlichen Ereignisse

$$h_1(Z,K) = 4$$
 $h_1(R,K) = 2$
 $h_1(K,Z) = 8$ $h_1(Z,Z) = 4$ $h_1(R,Z) = 4$

AUFGABE 2 — TEIL (D)

M-Schritt: Bestimmung der Teilkorpora h_1^1 bzw. h_1^2 durch *Marginalisierung*:

$X_1 \setminus X_2$	K	Ζ	
K	$h_1(K,K)$	$h_1(K,Z)$	$h_1^1(K)$
Ζ	$h_1(Z,K)$	$h_1(Z,Z)$	$h_1^1(Z)$
R	$h_1(R,K)$	$h_1(R,Z)$	$h_1^1(R)$
	$h_1^2(K)$	$h_1^2(Z)$	



AUFGABE 2 — TEIL (E)

Nun bestimmen wir noch die relativen Häufigkeiten mit der Formel

$$\mathsf{rfe}(h)(x) \coloneqq \frac{h(x)}{|h|} \quad \mathsf{mit} \quad |h| \coloneqq \sum_{x \in X} h(x)$$

Wenden wir dies nun auf h_1 und h_2 an, so erhalten wir

$$\begin{split} q_1^1(K) &= \mathsf{rfe}(h_1^1)(K) = \frac{h_1^1(K)}{h_1^1(K) + h_1^1(Z) + h_1^1(R)} = \frac{10}{24} = \frac{5}{12} \\ q_1^1(Z) &= \mathsf{rfe}(h_1^1)(Z) = \frac{h_1^1(Z)}{h_1(K) + h_1^1(Z) + h_1^1(R)} = \frac{8}{24} = \frac{1}{3} \\ q_1^1(R) &= \mathsf{rfe}(h_1^1)(R) = \frac{h_1^1(R)}{h_1^1(K) + h_1^1(Z) + h_1^1(R)} = \frac{6}{24} = \frac{1}{4} \end{split}$$

und

$$q_1^2(K) = \text{rfe}(h_1^2)(K) = \frac{h_1^2(K)}{h_1^2(K) + h_1^2(Z)} = \frac{8}{24} = \frac{1}{3}$$
$$q_1^2(Z) = \text{rfe}(h_1^2)(Z) = \frac{h_1^2(Z)}{h_1^2(K) + h_1^2(Z)} = \frac{16}{24} = \frac{2}{3}$$