



## Künstliche Intelligenz: Grundlagen und Anwendungen

Albayrak, Fricke (AOT) – Oppen, Thiel (KI)

Wintersemester 2016 / 2017

### 6. Aufgabenblatt

Abgabetermin: 18.01.2017

#### Aufgabe 1 – Hidden Markov-Prozess (50%)

Hidden-Markov-Modelle werden in der Bioinformatik zur Analyse von DNA-Sequenzen eingesetzt. Eine Anwendung ist das Auffinden von *CpG-Inseln* in der beobachteten Sequenz  $Y_t \in \{a, c, g, t\}$ . Der verborgene Zustand  $X_t \in \{w, f\}$  gibt an, ob das aktuelle Nukleotid zu einer CpG-Insel gehört ( $X_t = w$ ) oder nicht ( $X_t = f$ ). Für die Wahrscheinlichkeiten der einzelnen Nukleotide und die Übergangswahrscheinlichkeiten der Zustände gilt in einem stark vereinfachten Modell:

| $x$                      | $w$ | $f$ | $y$                  | $a$ | $c$ | $g$ | $t$ |
|--------------------------|-----|-----|----------------------|-----|-----|-----|-----|
| $P(X_{t+1} = x X_t = w)$ | 0.7 | 0.3 | $P(Y_t = y X_t = w)$ | 0.2 | 0.3 | 0.3 | 0.2 |
| $P(X_{t+1} = x X_t = f)$ | 0.2 | 0.8 | $P(Y_t = y X_t = f)$ | 0.3 | 0.2 | 0.2 | 0.3 |

Als Anfangsbedingung wird  $P(X_0 = w) = 0.5$  angenommen.

- Wie wahrscheinlich ist es, eine CpG-Insel der Länge  $k$  zu finden? Geben Sie  $P(X_1 = \dots = X_k = w, X_{k+1} = f|X_0 = f)$  für  $k \geq 1$  an!
- Sie wollen effizient CpG-Inseln in einer DNA-Sequenz finden und berechnen hierzu die Wahrscheinlichkeit  $p_t = P(X_t = w|Y_1, \dots, Y_t)$  aus  $p_{t-1}$  und  $Y_t$ . Wie sieht ein solcher Filter-Schritt für die Beobachtung  $Y_t = g$  aus?
- Wie hoch ist die Wahrscheinlichkeit  $P(X_t|Y_1 = c, Y_2 = g)$  für eine CpG-Insel an den Positionen  $t = 3$  und  $t = 4$ , wenn Sie nur die ersten zwei Nukleotide der DNA-Sequenz kennen?
- Wie hoch ist die Wahrscheinlichkeit  $P(X_1 = w|Y_1 = c, Y_2 = g)$ , dass bereits das erste Nukleotid der DNA-Sequenz  $Y_1 = c, Y_2 = g, \dots$  zu einer CpG-Insel gehört?
- Verwenden Sie den Viterbi-Algorithmus, um die wahrscheinlichste Folge von  $X_t$  für die DNA-Sequenz  $Y_1 = a, Y_2 = c, Y_3 = g, Y_4 = t$  zu finden!

**Aufgabe 2 – Hidden-Markov-Modell****(50%)**

Eine neuentdeckte Chamäleonart nutzt ihre Hautfarbe, um komplexe Botschaften zu kommunizieren. Wir unterscheiden zwischen einer Folge von Segmenten  $x_1, x_2, x_3, \dots$  und der tatsächlich beobachteten Folge von Farben  $y_1, y_2, y_3, \dots$

| $x_i$   | $x_{i+1}$ | $P(x_{i+1} x_i)$ |
|---------|-----------|------------------|
| Beginn  | Feind     | 0.4              |
| Beginn  | Nahrung   | 0.6              |
| Feind   | Ort       | 1.0              |
| Nahrung | Ort       | 0.2              |
| Nahrung | Menge     | 0.8              |
| Ort     | Beginn    | 0.3              |
| Ort     | Ende      | 0.7              |
| Menge   | Ort       | 0.3              |
| Menge   | Beginn    | 0.2              |
| Menge   | Ende      | 0.5              |

| $x_i$   | $y_i$   | $P(y_i x_i)$ |
|---------|---------|--------------|
| Beginn  | weiß    | 1.0          |
| Feind   | rot     | 0.6          |
| Feind   | blau    | 0.4          |
| Nahrung | rot     | 0.7          |
| Nahrung | grün    | 0.3          |
| Ort     | blau    | 0.8          |
| Ort     | orange  | 0.2          |
| Menge   | blau    | 0.1          |
| Menge   | grün    | 0.9          |
| Ende    | schwarz | 1.0          |

Die Markovkette beginnt immer mit  $x_1 = \text{Beginn}$  und endet mit  $x_k = \text{Ende}$ . Alle nicht angegebenen Wahrscheinlichkeiten  $P(x_{i+1}|x_i)$  und  $P(y_i|x_i)$  sind Null. Sie können die Segmenttypen mit großen und die Farben mit kleinen Anfangsbuchstaben abkürzen, um Platz zu sparen.

- Stellen Sie das Modell für die Ausdrücke in einem Übergangsdigramm graphisch dar! Sie brauchen keine Wahrscheinlichkeiten einzutragen.
- Mit welcher Wahrscheinlichkeit tritt die Farbfolge „weiß-rot-blau-orange-schwarz“ in diesem Modell auf?
- Sie beobachten die Folge „weiß-rot-orange-schwarz“. Ist es wahrscheinlicher, dass es um Nahrung oder Feinde geht? Wie sicher ist dies?
- Wie wahrscheinlich ist eine Nachricht aus 4 Segmenten?