

Maschinelles Lernen 2

Sommersemester 2010

Blatt 10

Abgabe bis Montag, 21. Juni 2010, 13:00 Uhr bei Dr. Konrad Rieck (rieck@cs.tu-berlin.de)

In der Vorlesung wurden maschinelle Lernmethoden in der Bioinformatik anhand des Beispiels *mGene* besprochen. Die Gensuchmaschine *mGene* funktioniert nach einem 2-Schichten Prinzip. In der ersten Schicht werden biologische Signale durch Support-Vector-Maschinen (SVM) gelernt und in der zweiten Schicht werden die Ausgaben der SVMs als Eingaben zum Strukturlernen verwendet, um eine genaue Segmentierung in “nicht-Gen”, “Gen” (Intron, Exon, usw.) zu erhalten. In der Übung sollen der Spektrum und der Weighted Degree Kernel—zwei Kerne aus der Gruppe der String-Kerne—weiter vertieft werden, da sie eine entscheidende Rolle beim Lernen der ersten Schicht in *mGene* spielen.

1. (10 Punkte) Der Weighted Degree Kern ist wie folgt definiert (Beispiel siehe Bild):

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^K \beta_k \sum_{l=1}^{L-k+1} \mathbf{I}(\mathbf{u}_{k,l}(\mathbf{x}_i) = \mathbf{u}_{k,l}(\mathbf{x}_j)). \quad (1)$$

Dabei ist $\mathbf{u}_{k,l}(\mathbf{x})$ ein String der Länge k welcher an Position l der Sequenz \mathbf{x} anfängt und $\mathbf{I}(\cdot)$ die Indikatorfunktion, welche 1 ist wenn ihr Argument wahr ist, sonst 0. Als Gewichtung wird $\beta_k = 2^{\frac{K-k+1}{K(K+1)}}$ gewählt.

```

x AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTGAG
#1-mers .|.|.|||.|.|||.|.|||.|||.|||.|||.|||.|||.|||.|||.|||.|||.
#2-mers .....||.....|.....|.....|.....|.....|.....|.....|.....|.....
#3-mers .....|.....|.....|.....|.....|.....|.....|.....|.....|.....
x' TACCTAATTATGAAATTAATTTTCTGTGCTGATGGAACGGAGAAGTC
    
```

Zeige für den Spektrum und den Weighted Degree (WD) Kern, dass es sich um positiv-definite Mercerkerne handelt.

2. (20 Punkte) Alternativ kann man den WDK auch in der “Block”-Formulierung betrachten. Hierbei werden zunächst längste gemeinsame Teilstrings identifiziert, und dann entsprechende Gewichte für die Länge dieser Teilstrings aufaddiert (im Bild dargestellt).

$k(\mathbf{s}_1, \mathbf{s}_2) = w_7 + w_1 + w_2 + w_2 + w_3$

```

S1 → AGTCAGATAGAGGACATCAGTAGACAGATTAAA →
      ||| ||| ||| |||
S2 → TTATAGATAGACAAAGACATCAGTAGACTTATT →
    
```

Berechne die Gewichte w_L , die den in Aufgabe 1 definierten WDK mit Gewichten β_k ergeben. Was ist zu beachten, wenn die Teilstrings länger werden als K ?

Für Fragen zum Übungsblatte bitte in der Google Group <http://groups.google.com/group/mikiobraun-lehre> registrieren und die Frage an die Mailingliste stellen.