

Prof. Ryan Cotterell

Entiol Liko: Assignment 2

enliko@student.ethz.ch

1 Theory

1.1 Question 1

1.1.1 Point a

$(\mathbb{R} \times \mathbb{R}, +)$ is a commutative monoid with identity element $\langle 0, 0 \rangle$:

1. $\langle a, b \rangle + \langle c, d \rangle = \langle a + c, b + d \rangle = \langle c, d \rangle + \langle a, b \rangle$
2. $\langle a, b \rangle + \langle 0, 0 \rangle = \langle a, b \rangle = \langle 0, 0 \rangle + \langle a, b \rangle$
3. $(\langle a, b \rangle + \langle c, d \rangle) + \langle e, f \rangle = \langle a + c, b + d \rangle + \langle e, f \rangle = \langle a + c + e, b + d + f \rangle$
 $\langle a, b \rangle + (\langle c, d \rangle + \langle e, f \rangle) = \langle a, b \rangle + \langle c + e, d + f \rangle = \langle a + c + e, b + d + f \rangle$

$(\mathbb{R} \times \mathbb{R}, \times)$ is a monoid with identity element $\langle 1, 0 \rangle$:

1. $\langle a, b \rangle \times \langle 1, 0 \rangle = \langle a \times 1, a \times 0 + b \times 1 \rangle = \langle a, b \rangle = \langle 1, 0 \rangle + \langle a, b \rangle$
2. $(\langle a, b \rangle \times \langle c, d \rangle) \times \langle e, f \rangle = \langle ac, ad + bc \rangle \times \langle e, f \rangle = \langle ace, acf + ade + bce \rangle$
 $\langle a, b \rangle \times (\langle c, d \rangle \times \langle e, f \rangle) = \langle a, b \rangle \times \langle ce, cf + de \rangle = \langle ace, acf + ade + bce \rangle$

Multiplication left and right distributes over addition:

1. $\langle a, b \rangle \times (\langle c, d \rangle + \langle e, f \rangle) = \langle a, b \rangle \times \langle c + e, d + f \rangle = \langle ac + ae, ad + af + bc + be \rangle$
 $(\langle a, b \rangle \times \langle c, d \rangle) + (\langle a, b \rangle \times \langle e, f \rangle) = \langle ac, ad + bc \rangle + \langle ae, af + be \rangle = \langle ac + ae, ad + af + bc + be \rangle$
2. $(\langle a, b \rangle + \langle c, d \rangle) \times \langle e, f \rangle = \langle a + c, b + d \rangle \times \langle e, f \rangle = \langle ae + ce, af + fc + be + de \rangle$
 $(\langle a, b \rangle \times \langle e, f \rangle) + (\langle c, d \rangle \times \langle e, f \rangle) = \langle ae, af + be \rangle + \langle ce, cf + de \rangle = \langle ae + ce, af + fc + be + de \rangle$

Multiplication by 0 annihilates $\mathbb{R} \times \mathbb{R}$:

1. $\langle a, b \rangle \times \langle 0, 0 \rangle = \langle a0, a0 + b0 \rangle = \langle a0, a0 + b0 \rangle = \langle 0, 0 \rangle \times \langle a, b \rangle$

1.1.2 Point b

We will start by defining some symbols so we don't have to write the whole notation. Let $\omega(t_n^i, t_m^j) := \exp(\text{score}_\theta(\langle t_n^i, t_m^j \rangle, \omega))$ and $\alpha(t_n^i, t_m^j) := \text{score}_\theta(\langle t_n^i, t_m^j \rangle, \omega)$

We can think of having a work of length N and of a set of tags of length T . For the last but one layer we have:

$$\beta(t_{N-1}^i, N-1) = \langle \omega(t_{N-1}^i, t_N^0), -\omega(t_{N-1}^i, t_N^0) \alpha(t_{N-1}^i, t_N^0) \rangle$$

We can calculate the same thing for the last to second layer.

$$\begin{aligned} \beta(t_{N-2}^i, n) &= \sum_{j=0}^T < \omega(t_{N-2}^i, t_{N-1}^j), -\omega(t_{N-2}^i, t_{N-1}^j) \alpha(t_{N-2}^i, t_{N-1}^j) > \times \\ &< \omega(t_{N-1}^j, t_N^0), -\omega(t_{N-1}^j, t_N^0) \alpha(t_{N-1}^j, t_N^0) > = \dots = \\ &< \sum_{j=0}^T \exp(\alpha(t_{N-2}^i, t_{N-1}^j) + \alpha(t_{N-1}^j, t_N^0)) - \exp(\alpha(t_{N-2}^i, t_{N-1}^j) + \alpha(t_{N-1}^j, t_N^0)) (\alpha(t_{N-2}^i, t_{N-1}^j) + \\ &\alpha(t_{N-1}^j, t_N^0)) > \end{aligned}$$

It can be easily seen that if we continue the procedure the second component of the result is $H_U(T_\omega)$, what we are looking for.

1.1.3 Point c

$$\begin{aligned} H(T_\omega) &= -\sum_{t \in T} p(t|\omega) \log(p(t|\omega)) = -\sum_{t \in T} \frac{\exp(\text{score}_\theta(t, \omega))}{Z(\omega)} (\text{score}_\theta(t, \omega) - \log(Z(\omega))) = \\ &\frac{\sum_{t \in T} \exp(\text{score}_\theta(t, \omega)) \text{score}_\theta(t, \omega)}{Z(\omega)} + \log(Z(\omega)) = \frac{H_U(T_\omega)}{Z(\omega)} + \log(Z(\omega)) \end{aligned}$$

1.1.4 Point d

Since we can compute $H_U(T_\omega)$ and Z_ω in $\mathcal{O}(N|T|^2)$ we can calculate $H(T_\omega)$ in $\mathcal{O}(N|T|^2)$. We can compute the gradient of $H_U(T_\omega)$ in $\mathcal{O}(N|T|^2)$ if we assume that we can calculate the gradient of the score function with respect to the parameters θ . In fact we could reuse the backlog algorithm but with a different semiring. For example to calculate the gradient of Z we can use the expectation semiring where the elements are vectors where each component is a tuple $\langle w, \text{derivate of } w \text{ with respect to the } i\text{-th parameter} \rangle$.

1.2 Question 2

1.2.1 Point a

We could do a proof by induction. We can clearly see that the algorithm is true when we have a single node, the BOT node. We assume that the algorithm is true and correct up to the last popped node and we have to prove that it holds for the next popped node. We can do a proof by contradiction and assume that there exists another path that passes by a node which is not popped that connects BOT with the current popped node with a better score than the current path. If we prove that this cannot be true then we have finished the proof. The proof of this contradiction is pretty straightforward since it is very similar to the proof that the Dijkstra algorithm for shortest path is correct. In essence we show that since we have negative weights each node that is added to the path lowers the score and so if there exist a more optimal path then in that path there is a node with a higher score than the score of the node we just popped (because of negative weights which are specified by the ring) but that is a contradiction so there can't exist a more optimal path.

1.2.2 Point b

I think that we can easily see that both algorithms return the same result by noting that the $+$ and \max operations are commutative. This means that it doesn't matter if we start

from the left or from the right. In the end we calculate, for each node, the score of the best part-of-speech tagging for that node and that score is unique so we get the same result. From the result of point a we can say that the algorithm works for every one, in the sense that it calculates the best tagging from the BOT to that node and since the path is unique it returns the same result as Viterbi.

1.2.3 Point c

We assume that the insert and pop operations for the PriorityQueue have a time complexity equal to $\mathcal{O}(\log V)$, where V is the number of elements in the queue. The time complexity of Dijkstra is $\mathcal{O}(N|T|^2 \log V)$, where N is the number of words in the sentence, T the number of possible tags and V the number of elements in the priority queue. We have to do a push for every "connection" between nodes. There are $N \cdot T$ nodes and T connections for each node. That gives us the this time complexity. In the case of the Viterbi algorithm runs in $\mathcal{O}(N|T|^2)$. The benefit of using Dijkstra over Viterbi is that Viterbi visits every single edge but Dijkstra may reach the end, so best tagging sequence, without necessarily having visited every edge so even if the time complexity is greater in reality it may finish sooner.

1.2.4 Point d

I think that this operation is not adequate for what we are trying to do. We are searching for the best tagging but the problem is that the using this operation gives us a value which is closest to the smallest value not the biggest.

1.2.5 Point e

In our setting the best path represents the path with the highest score. What we need is a semiring for which the proof of point a is still true. We could write that we want that if $a = a \oplus b \Rightarrow a = a \oplus (b \otimes c)$, where \oplus is *max* or *min*. What we are basically saying is that if one node has a larger value than another node the path that passes from the second node and end with the first node will have a lower score than the optimal path from BOT to the first node. This way we guarantee that the proof of point a is correct. This is the reason the other semirings specified in the question do not work. Note: The semiring $\langle R_{\geq 0}, \min, +, +\infty, 0 \rangle$ used with Dijkstra returns the path with the lowest score.

2 Practice

For the practical part of the algorithms were vectorized using a for loop and the implementation was done using the pseudo code found in the course notes.