

Prof. Ryan Cotterell

Simon Wachter: Assignment 2

siwachte@ethz.ch, 19-920-198

27/11/2022 - 16:46h

Question 1:

a) Prove that the expectation semiring satisfies the semiring axioms:

- $(\mathbb{R} \times \mathbb{R}, \oplus, \mathbf{0})$ must be a commutative monoid with identity element $\mathbf{0}$:

$$(\langle x, y \rangle \oplus \langle x', y' \rangle) \oplus \langle x'', y'' \rangle = \langle x + x', y + y' \rangle \oplus \langle x'', y'' \rangle \quad (1)$$

$$= \langle x + x' + x'', y + y' + y'' \rangle \quad (2)$$

$$= \langle x, y \rangle \oplus \langle x' + x'', y' + y'' \rangle \quad (3)$$

$$= \langle x, y \rangle \oplus (\langle x', y' \rangle \oplus \langle x'', y'' \rangle) \quad (4)$$

$$\mathbf{0} + \langle x, y \rangle = \langle 0, 0 \rangle \oplus \langle x, y \rangle \quad (5)$$

$$= \langle 0 + x, 0 + y \rangle \quad (6)$$

$$= \langle x, y \rangle \quad (7)$$

$$= \langle x + 0, y + 0 \rangle \quad (8)$$

$$= \langle x, y \rangle + \mathbf{0} \quad (9)$$

$$\langle x, y \rangle + \langle x', y' \rangle = \langle x + x', y + y' \rangle \quad (10)$$

$$= \langle x' + x, y' + y \rangle \quad (11)$$

$$= \langle x', y' \rangle + \langle x, y \rangle \quad (12)$$

- $(\mathbb{R} \times \mathbb{R}, \otimes, \mathbf{1})$ must be a monoid with identity element $\mathbf{1}$:

$$(\langle x, y \rangle \otimes \langle x', y' \rangle) \otimes \langle x'', y'' \rangle = \langle x \cdot x', x \cdot y' + y \cdot x' \rangle \otimes \langle x'', y'' \rangle \quad (13)$$

$$= \langle x \cdot x' \cdot x'', x \cdot x' \cdot y'' + (x \cdot y' + y \cdot x') \cdot x'' \rangle \quad (14)$$

$$= \langle x \cdot x' \cdot x'', x \cdot x' \cdot y'' + x \cdot y' \cdot x'' + y \cdot x' \cdot x'' \rangle \quad (15)$$

$$= \langle x, y \rangle \otimes \langle x' \cdot x'', x' \cdot y'' + y' \cdot x'' \rangle \quad (16)$$

$$= \langle x, y \rangle \otimes (\langle x', y' \rangle \otimes \langle x'', y'' \rangle) \quad (17)$$

$$\mathbf{1} \otimes \langle x, y \rangle = \langle 1, 0 \rangle \otimes \langle x, y \rangle \quad (18)$$

$$= \langle 1 \cdot x, 1 \cdot y \rangle \quad (19)$$

$$= \langle x, y \rangle \quad (20)$$

$$= \langle x \cdot 1, y \cdot 1 \rangle \quad (21)$$

$$= \langle x, y \rangle \otimes \mathbf{1} \quad (22)$$

- Multiplication left and right distributes over addition:

$$\langle x, y \rangle \otimes (\langle x', y' \rangle \oplus \langle x'', y'' \rangle) = \langle x, y \rangle \otimes \langle x' + x'', y' + y'' \rangle \quad (23)$$

$$= \langle x \cdot x' + x \cdot x'', x \cdot y' + x \cdot y'' + y \cdot x' + y \cdot x'' \rangle \quad (24)$$

$$= \langle x \cdot x', x \cdot y' + y \cdot x' \rangle \oplus \langle x \cdot x'', x \cdot y'' + y \cdot x'' \rangle \quad (25)$$

$$= (\langle x, y \rangle \otimes \langle x', y' \rangle) \oplus (\langle x, y \rangle \otimes \langle x'', y'' \rangle) \quad (26)$$

$$(\langle x, y \rangle \oplus \langle x', y' \rangle) \otimes \langle x'', y'' \rangle = \langle x + x', y + y' \rangle \otimes \langle x'', y'' \rangle \quad (27)$$

$$= \langle x \cdot x'' + x' \cdot x'', x \cdot y'' + x' \cdot y'' + y \cdot x'' + y' \cdot x'' \rangle \quad (28)$$

$$= \langle x \cdot x'', x \cdot y'' + y \cdot x'' \rangle \oplus \langle x' \cdot x'', x' \cdot y'' + y' \cdot x'' \rangle \quad (29)$$

$$= (\langle x, y \rangle \otimes \langle x'', y'' \rangle) \oplus (\langle x', y' \rangle \otimes \langle x'', y'' \rangle) \quad (30)$$

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

$$\mathbf{0} \otimes \langle x, y \rangle = \langle 0, 0 \rangle \otimes \langle x, y \rangle \quad (31)$$

$$= \langle 0 \cdot x, 0 \cdot y \rangle \quad (32)$$

$$= \langle 0, 0 \rangle \quad (33)$$

$$= \mathbf{0} \quad (34)$$

$$= \langle 0, 0 \rangle \quad (35)$$

$$= \langle x \cdot 0, y \cdot 0 \rangle \quad (36)$$

$$= \langle x, y \rangle \otimes \langle 0, 0 \rangle \quad (37)$$

$$= \langle x, y \rangle \otimes \mathbf{0} \quad (38)$$

b) Our initial graph looks like Fig. 1.

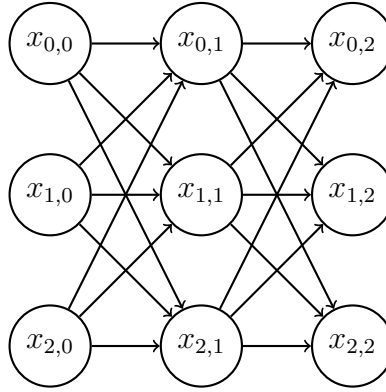


Figure 1: The initial graph

Where the columns represent the words in \mathbf{w} and the rows represent different tags. We use the algorithm from the script:

Algorithm 1: Forward pass

```

1  $\beta(\mathbf{w}, t_0) = 1$ 
2 for  $n = 1 \rightarrow N$  do
3    $\beta(\mathbf{w}, t_n) = \sum_{t_{n-1} \in \mathcal{T}} \exp(\text{score}_\theta(\langle t_{n-1}, t_n \rangle, \mathbf{w})) \otimes \beta(\mathbf{w}, t_{n-1})$ 
4 end
```

When we now lift the CRF into the expectation semiring, the forward propagation algorithm changes to:

Algorithm 2: Forward pass

```

1  $\beta(\mathbf{w}, t_0) = \langle 1, 0 \rangle$ 
2 for  $n = 1 \rightarrow N$  do
3    $\mid \beta(\mathbf{w}, t_n) = \oplus_{t_{n-1} \in \mathcal{T}} \langle w, -w \log w \rangle \otimes \beta(\mathbf{w}, t_{n-1})$ 
4 end

```

Where $w = \exp(\text{score}_\theta(\langle t_n, t_{n+1} \rangle, \mathbf{w}))$.

The output of the forward algorithm lifted into the semiring will yield:

$$\bigoplus_{t_{1:N} \in T^N} \bigotimes_{n=1}^N \langle w, -w \log w \rangle \quad (39)$$

We want to show that the result of the forward propagation lifted in the semiring is the same as the unnormalized Entropy:

$$H_u(T_w) = - \sum_{\mathbf{t} \in \mathcal{T}^N} \exp(\text{score}_\theta(\mathbf{t}, \mathbf{w})) \text{score}_\theta(\mathbf{t}, \mathbf{w}) \quad (40)$$

$$(41)$$

We show this by induction. Starting with the base case where $N = 1$:

$$\bigoplus_{t_1 \in T^1} \bigotimes_{n=1}^1 \langle w, -w \log w \rangle = \bigoplus_{t_1 \in T^1} \langle \exp(\text{score}_\theta(\langle t_0, t_1 \rangle, \mathbf{w})), \quad (42)$$

$$- \exp(\text{score}_\theta(\langle t_0, t_1 \rangle, \mathbf{w})) \log(\exp(\text{score}_\theta(\langle t_0, t_1 \rangle, \mathbf{w}))) \rangle \quad (43)$$

$$= \bigoplus_{t \in T} \left\langle \exp(\text{score}_\theta(t, \mathbf{w})), \right. \quad (44)$$

$$\left. - \exp(\text{score}_\theta(t, \mathbf{w})) \log(\exp(\text{score}_\theta(t, \mathbf{w}))) \right\rangle$$

$$= \left\langle \sum_{t \in T} \exp(\text{score}_\theta(t, \mathbf{w})), \right. \quad (45)$$

$$\left. - \sum_{t \in T} \exp(\text{score}_\theta(t, \mathbf{w})) \log(\exp(\text{score}_\theta(t, \mathbf{w}))) \right\rangle$$

Our induction hypothesis is the following:

$$\bigoplus_{t_{1:i} \in T^i} \bigotimes_{n=1}^i \langle w, -w \log w \rangle = \left\langle \sum_{t \in T^i} \exp(\text{score}_\theta(t, \mathbf{w})), - \sum_{t \in T^i} \exp(\text{score}_\theta(t, \mathbf{w})) \log(\exp(\text{score}_\theta(t, \mathbf{w}))) \right\rangle \quad (46)$$

Meaning we assume that $\beta(\mathbf{w}, t_i)$ corresponds to the unnormalized entropy of all sequences of length i .

Now we proceed with the induction step, where $i \rightarrow i + 1$:

$$\bigoplus_{t_{1:N} \in T^n} \bigotimes_{n=1}^N \langle w, -w \log w \rangle \quad (47)$$

$$= \bigoplus_{t_{1:N-1} \in T^{N-1}} \bigoplus_{t_N \in T} \bigotimes_{n=1}^N \langle w, -w \log w \rangle \quad (48)$$

$$= \bigoplus_{t_{1:N-1} \in T^{N-1}} \bigotimes_{n=1}^{N-1} \langle \exp(\text{score}_\theta(\langle t_{n-1}, t_n \rangle, \mathbf{w})), -\exp(\text{score}_\theta(\langle t_{n-1}, t_n \rangle, \mathbf{w})) \log(\exp(\text{score}_\theta(\langle t_{n-1}, t_n \rangle, \mathbf{w}))) \rangle \\ \otimes \bigoplus_{t_N \in T} \langle \exp(\text{score}_\theta(\langle t_{N-1}, t_N \rangle, \mathbf{w})), -\exp(\text{score}_\theta(\langle t_{N-1}, t_N \rangle, \mathbf{w})) \log(\exp(\text{score}_\theta(\langle t_{N-1}, t_N \rangle, \mathbf{w}))) \rangle \rangle \quad (49)$$

$$= \bigoplus_{t_1 \in T} \langle \exp(\text{score}_\theta(\langle t_0, t_1 \rangle, \mathbf{w})), -\exp(\text{score}_\theta(\langle t_0, t_1 \rangle, \mathbf{w})) \log(\exp(\text{score}_\theta(\langle t_0, t_1 \rangle, \mathbf{w}))) \rangle \\ \otimes \bigoplus_{t_2 \in T} \langle \exp(\text{score}_\theta(\langle t_1, t_2 \rangle, \mathbf{w})), -\exp(\text{score}_\theta(\langle t_1, t_2 \rangle, \mathbf{w})) \log(\exp(\text{score}_\theta(\langle t_1, t_2 \rangle, \mathbf{w}))) \rangle \\ \otimes \dots \\ \otimes \bigoplus_{t_N \in T} \langle \exp(\text{score}_\theta(\langle t_{N-1}, t_N \rangle, \mathbf{w})), -\exp(\text{score}_\theta(\langle t_{N-1}, t_N \rangle, \mathbf{w})) \log(\exp(\text{score}_\theta(\langle t_{N-1}, t_N \rangle, \mathbf{w}))) \rangle \rangle \\ = \left\langle \sum_{\mathbf{t} \in \mathcal{T}^N} \exp(\text{score}_\theta(\mathbf{t}, \mathbf{w})), -\sum_{\mathbf{t} \in \mathcal{T}^N} \exp(\text{score}_\theta(\mathbf{t}, \mathbf{w})) \text{score}_\theta(\mathbf{t}, \mathbf{w}) \right\rangle \quad (50)$$

c) We want to prove:

$$H(T_w) = Z(\mathbf{w})^{-1} H_U(T) + \log(Z(\mathbf{w})) \quad (51)$$

$$H(T_w) = - \sum_{t \in T^N} p(t \mid w) \cdot \log(p(t \mid w)) \quad (\text{def. H}) \quad (52)$$

$$= - \sum_{t \in T^N} \frac{\exp(\text{score}_\theta(t, \mathbf{w}))}{Z(\mathbf{w})} \log \left(\frac{\exp(\text{score}_\theta(t, \mathbf{w}))}{\sum_{t' \in T^N} \exp(\text{score}_\theta(t', \mathbf{w}))} \right) \quad (\text{def. p}) \quad (53)$$

$$= - \sum_{t \in T^N} \frac{\exp(\text{score}_\theta(t, \mathbf{w}))}{Z(\mathbf{w})} \log \left(\frac{\exp(\text{score}_\theta(t, \mathbf{w}))}{Z(\mathbf{w})} \right) \quad (54)$$

$$= - \sum_{t \in T^N} \frac{\exp(\text{score}_\theta(t, \mathbf{w}))}{Z(\mathbf{w})} (\text{score}_\theta(t, \mathbf{w}) - \log Z(\mathbf{w})) \quad (55)$$

$$= - \sum_{t \in T^N} \frac{\exp(\text{score}_\theta(t, \mathbf{w})) \text{score}_\theta(t, \mathbf{w}) - \exp(\text{score}_\theta(t, \mathbf{w})) \log Z(\mathbf{w})}{Z(\mathbf{w})} \quad (56)$$

$$= - \sum_{t \in T^N} \frac{\exp(\text{score}_\theta(t, \mathbf{w})) \text{score}_\theta(t, \mathbf{w})}{Z(\mathbf{w})} + \sum_{t \in T^N} \frac{\exp(\text{score}_\theta(t, \mathbf{w})) \log Z(\mathbf{w})}{Z(\mathbf{w})} \quad (57)$$

$$= H_U(T_w) Z(\mathbf{w})^{-1} + \frac{\log(Z(\mathbf{w}))}{Z(\mathbf{w})} \sum_{t \in T^N} \exp(\text{score}_\theta(t, \mathbf{w})) \quad (58)$$

$$= H_U(T_w) Z(\mathbf{w})^{-1} + \log(Z(\mathbf{w})) \quad (59)$$

$$(60)$$