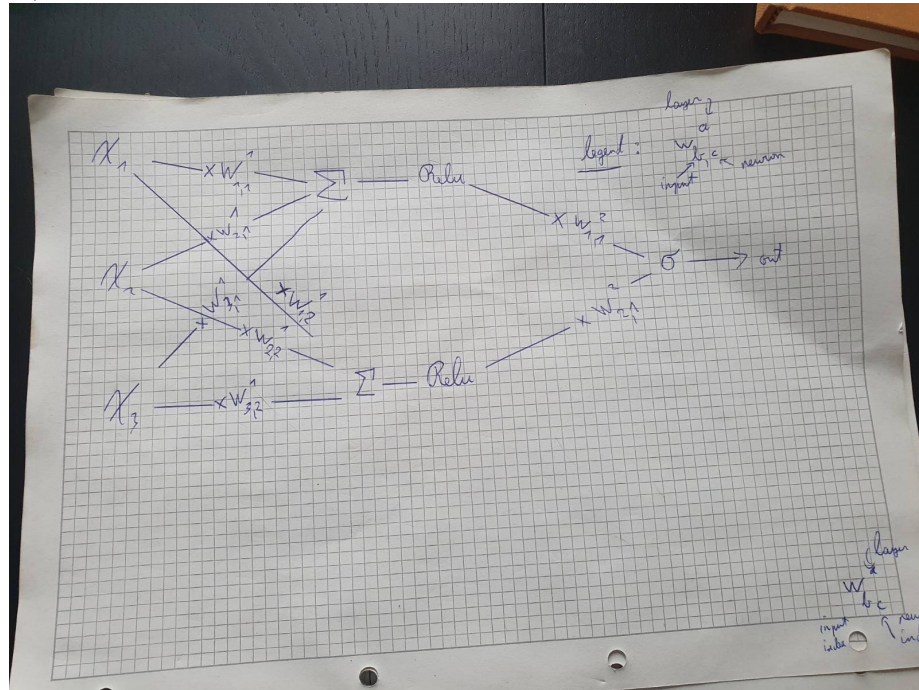


NLP assignments #1 and # 2

Andrius Buinovskij - andriusb - 18-940-270

1 Assignment #1

Q1 a)



Q1 b) i)

Well since all inputs are 1 and all weights are 1 then

$$s_1^1 = \sum_{i=1}^3 x_i \cdot w_{i,1}^1 = 3 \quad (1)$$

$$s_2^1 = \sum_{i=1}^3 x_i \cdot w_{i,2}^1 = 3 \quad (2)$$

Where s_j^i is the sum input to the j 'th neuron in the i 'th layer, and s^i is a vector whose length is equal to the number of neurons in the i 'th layer, with the

input being the 0'th layer. So \mathbf{s}^1 is of length 2 since 1'st layer has 2 neurons.

Let \mathbf{n}_j^i be the output of the j 'th neuron in the i 'th layer, then of course \mathbf{n}^i is a vector of length equal to the number of neurons in the i 'th layer:

$$n_j^i = ReLU(\mathbf{s}_j^i) \quad (3)$$

So in our case We get

$$\mathbf{n}_1^1 = ReLU(\mathbf{s}_1^1) = ReLU(3) = 3 \quad (4)$$

$$\mathbf{n}_2^1 = ReLU(\mathbf{s}_2^1) = ReLU(3) = 3 \quad (5)$$

Same steps in the next layer:

$$\mathbf{s}_1^2 = 3 \cdot 1 + 3 \cdot 1 = 6 \quad (6)$$

And now We pass this through a sigmoid for our output instead of a $ReLU$ so We get

$$out = \sigma(\mathbf{s}_1^2) = \frac{1}{1 + e^{-6}} = 0.99752737684 \quad (7)$$

Q1 b) ii-iv)

Alright let's just roll all of these questions into one and do an iteration of backprop.

$$L_{BCE} = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})) \quad (8)$$

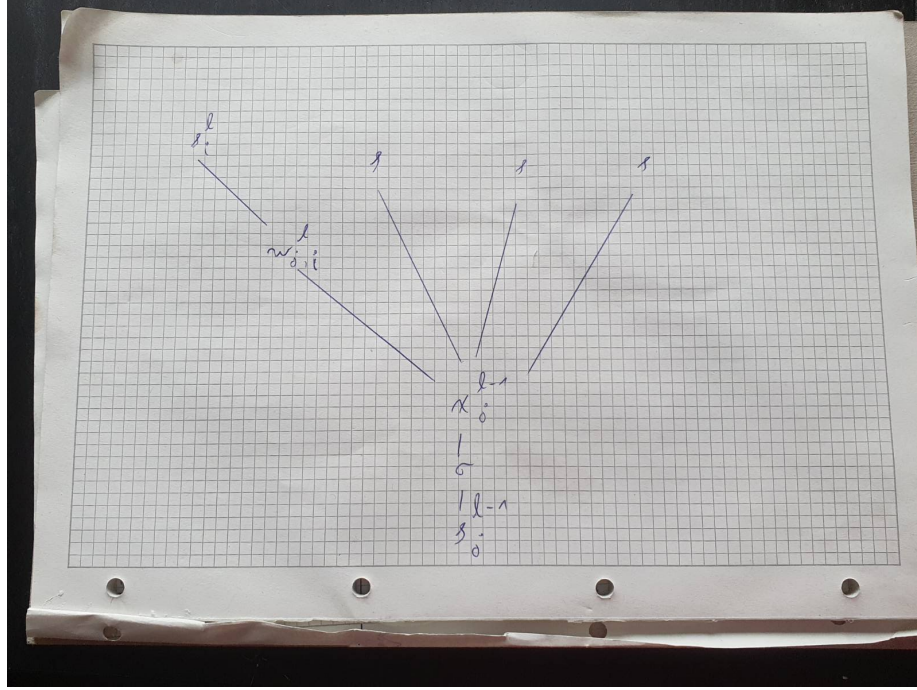
$$= -(0 \cdot \log(0.99752737684) + (1 - 0) \log(1 - 0.99752737684)) \quad (9)$$

$$= -(\log(0.00247262316)) \quad (10)$$

$$= 2.60684206722 \quad (11)$$

The forward pass is trivial.

Here is a pictorial view of the terms about to be defined:



Let \mathbf{x}_i^l be the output of the i 'th neuron in the l 'th layer.

Let $\mathbf{w}_{i,j}^l$ be the weight belonging to j 'th neuron multiplying the i 'th input in the l 'th layer. $\mathbf{w}_{i,j}^l$ is then simply the weight vector associated with the j 'th neuron in the l 'th layer.

Let \mathbf{s}_j^l be $(\mathbf{x}^{l-1})^\top \mathbf{w}_{\cdot,j}^l$, the weighted sum of inputs to the j 'th neuron in the l 'th layer.

We can then say that $x_i^l = \sigma(\mathbf{s}_j^l)$, where σ is the nonlinearity of choice.

Then define

$$\delta_j^l = \frac{\partial L_{BCE}}{\partial \mathbf{s}_j^l} \quad (12)$$

And finally

$$\delta_j^{l-1} = \sum_{i=1}^{d^l} \frac{\partial L_{BCE}}{\partial \mathbf{s}_i^l} \frac{\partial \mathbf{s}_i^l}{\partial \mathbf{x}_j^{l-1}} \frac{\partial \mathbf{x}_j^{l-1}}{\partial \mathbf{s}_j^{l-1}} \quad (13)$$

$$= \sum_{i=1}^{d^l} \delta_i^l \frac{\partial \mathbf{s}_i^l}{\partial \mathbf{x}_j^{l-1}} \frac{\partial \mathbf{x}_j^{l-1}}{\partial \mathbf{s}_j^{l-1}} \quad (14)$$

Where d^l is the number of neurons in layer l . So now We have a recursive definition which uses dynamic programming. This could be further nuanced by expressing things in matrix notation, but it's good enough for present purposes.

So now taking the derivative of the loss w.r.t. first sum:

$$\frac{\partial L_{BCE}}{\partial \mathbf{s}_1^2} = \frac{\partial L_{BCE}}{\partial \hat{y}} \frac{\partial \hat{y}}{\mathbf{s}_1^2} \quad (15)$$

Calculating terms individually:

$$\frac{\partial L_{BCE}}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} - (y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})) \quad (16)$$

$$= -y \frac{1}{\hat{y}} - (1 - y) \frac{\partial}{\partial \hat{y}} \log(1 - \hat{y}) \quad (17)$$

$$= -y \frac{1}{\hat{y}} - (1 - y) \frac{-1}{1 - \hat{y}} \quad (18)$$

$$= -y \frac{1}{\hat{y}} + \frac{1 - y}{1 - \hat{y}} \quad (19)$$

And the second bit in equation 15:

$$\frac{\partial \hat{y}}{\mathbf{s}_1^2} = \frac{\partial}{\mathbf{s}_1^2} \sigma(\mathbf{s}_1^2) \quad (20)$$

$$= \sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \quad (21)$$

So then We have

$$\delta_1^2 = \frac{\partial L_{BCE}}{\partial \mathbf{s}_1^2} = \frac{\partial L_{BCE}}{\partial \hat{y}} \frac{\partial \hat{y}}{\mathbf{s}_1^2} \quad (22)$$

$$= \left(-y \frac{1}{\hat{y}} + \frac{1 - y}{1 - \hat{y}} \right) \left(\sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \right) \quad (23)$$

Alright. We only have one more of these to figure out:

$$\delta_j^1 = \sum_{i=1}^{d^2} \delta_i^2 \frac{\partial \mathbf{s}_i^2}{\partial \mathbf{x}_j^1} \frac{\partial \mathbf{x}_j^1}{\partial \mathbf{s}_j^1} \quad (24)$$

But since $d^2 = 1$, i.e. there is only one neuron in the output layer, We have

$$\delta_j^1 = \sum_{i=1}^{d^2} \delta_i^2 \frac{\partial \mathbf{s}_i^2}{\partial \mathbf{x}_j^1} \frac{\partial \mathbf{x}_j^1}{\partial \mathbf{s}_j^1} \quad (25)$$

$$= \delta_1^2 \frac{\partial \mathbf{s}_1^2}{\partial \mathbf{x}_j^1} \frac{\partial \mathbf{x}_j^1}{\partial \mathbf{s}_j^1} \quad (26)$$

$$= \delta_1^2 \mathbf{w}_j^2 \sigma(\mathbf{s}_j^1) (1 - \sigma(\mathbf{s}_j^1)) \quad (27)$$

Now for the gradient update We will of course need

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^l} \quad (28)$$

But this can be easily derived since

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^l} = \frac{\partial L_{BCE}}{\partial \mathbf{s}_j^l} \frac{\partial \mathbf{s}_j^l}{\partial \mathbf{w}_{i,j}^l} \quad (29)$$

$$= \delta_j^l \mathbf{x}_i^{l-1} \quad (30)$$

Now We do a forward pass and We know that $\mathbf{s}_i^1 = 3$ and $\mathbf{s}_1^2 = 6$. Plugging in values We then get

$$\delta_1^2 = \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \right) \quad (31)$$

$$= \left(\frac{1}{1 - 0.99752737684} \right) \left(\sigma(6) \cdot (1 - \sigma(6)) \right) \quad (32)$$

$$= \left(\frac{1}{1 - 0.99752737684} \right) \left(0.99752737684 \cdot (1 - 0.99752737684) \right) \quad (33)$$

$$= 404.428792942 \cdot 0.00246650929 \quad (34)$$

$$= 0.99752737493 \quad (35)$$

Similarly We have

$$\delta_j^1 = \delta_1^2 \mathbf{w}_j^2 \sigma(\mathbf{s}_j^1) (1 - \sigma(\mathbf{s}_j^1)) \quad (36)$$

$$\delta_1^1 = \delta_1^2 \mathbf{w}_1^2 \sigma(\mathbf{s}_1^1) (1 - \sigma(\mathbf{s}_1^1)) \quad (37)$$

$$= 0.997527374931 \cdot \sigma(3) (1 - \sigma(3)) \quad (38)$$

$$= 0.997527374931 \cdot 0.95257412682 \cdot 0.04742587317 \quad (39)$$

$$= 0.04506495478 \quad (40)$$

δ_2^1 is equal to δ_1^1 .

Now that We have the deltas We can use

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^l} = \delta_j^l \mathbf{x}_i^{l-1} \quad (41)$$

For weights in the first layer this means

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^1} = \delta_j^1 \mathbf{x}_i^0 \quad (42)$$

$$= 0.04506495478 \cdot 1 \quad (43)$$

$$= 0.04506495478 \quad (44)$$

Where the 0'th layer is the input layer.

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^2} = \delta_j^2 \mathbf{x}_i^1 \quad (45)$$

$$= 0.99752737493 \cdot 3 \quad (46)$$

$$= 2.99258212479 \quad (47)$$

All that is left is to perform a step for each weight. All weights in the first layer simplify to

$$w_{i,j}^1 = 1 - 0.1 \cdot 0.04506495478 \quad (48)$$

$$= 0.99549350452 \quad (49)$$

and similarly for the second layer:

$$w_{i,j}^2 = 1 - 0.1 \cdot 2.99258212479 \quad (50)$$

$$= 0.70074178752 \quad (51)$$

Sidenote

So when actually using this in the implementation of Q2, We get the following simplification:

$$\delta_1^2 = \frac{\partial L_{BCE}}{\partial \mathbf{s}_1^2} = \frac{\partial L_{BCE}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{s}_1^2} \quad (52)$$

$$= \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \right) \quad (53)$$

$$= \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\hat{y} \cdot (1 - \hat{y}) \right) \quad (54)$$

$$= -\hat{y} \cdot (1 - \hat{y}) \cdot \frac{y}{\hat{y}} + \hat{y} \cdot (1 - \hat{y}) \cdot \frac{1-y}{1-\hat{y}} \quad (55)$$

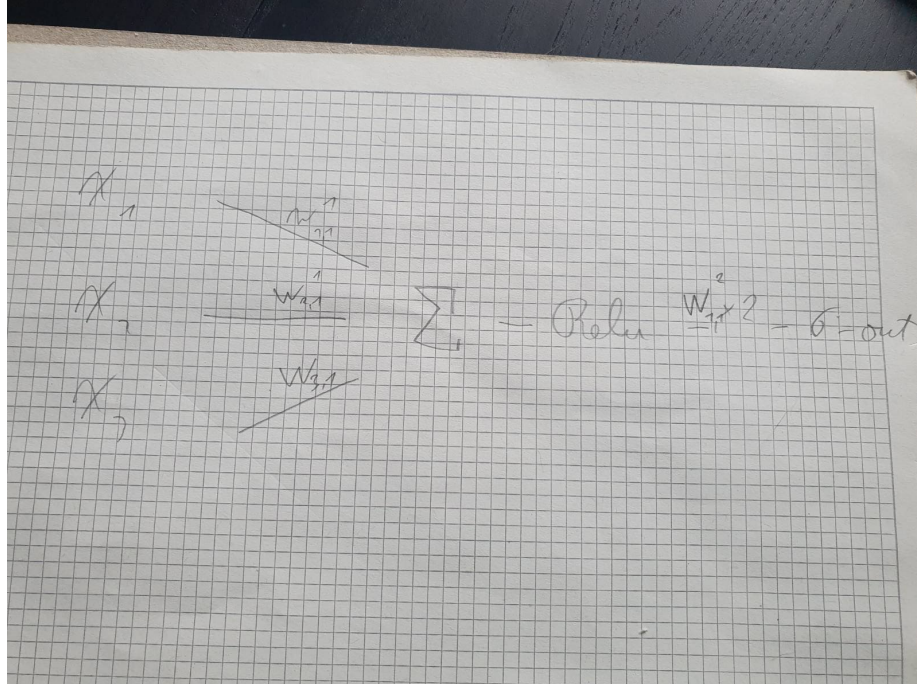
$$= -(1 - \hat{y}) \cdot y + \hat{y} \cdot (1 - y) \quad (56)$$

Q1 c))

Due to the fact that the weights were initialized as a uniform constant, I think that the second neuron in the hidden layer is redundant.

I don't think We can get rid of the hidden layer all together due to the ReLU
 - if We simply passed the output of the sum to the sigmoid, the sigmoid could receive negative terms, which is not possible under the current architecture. As it is, the input to the sigmoid is bounded to be above zero regardless of the weights.

The below is certainly more efficient and has fewer parameters:



The reason the $\times 2$ is there is to make the computation identical regardless of what the weights are initialized to, since currently the source of the redundancy is that the two neurons in the first hidden layer would output identical values.

As for other benefits, numerical stability maybe?

I don't think that this is the intended solution haha, but I don't think You can get rid of the hidden layer all together due to the ReLU as mentioned above, so "_(ツ)_/"

Q2

Running the script outputs the relevant comparisons.

The scikit implementation performs slightly better, but the two are very comparable accuracy-wise e.g. 0.8925 0.8765 out-of-sample accuracy in the third test case.

I should note that my implementation is almost an order of magnitude slower: prediction time sklearn: 0:00:00.000876 v.s. mine: 0:00:00.003137

Finally, as can be seen by the output of the script, the sklearn implementation prefers smaller absolute values of the weights, where as my implementation does not seem shy about using larger values. Sklearn is superior in this regard.

Q4 a) i)

We're going to use linearity of expectation.

Let $I_i, i \in \{1 \dots |V|\}$ be the indicator variable for whether or not the word i appears in Mary's sampling.

Since draws are with replacement and the distribution is uniform, for a token *not* to appear, all n draws must have chosen another token.

The probability of not choosing token i in a particular draw is of course $(|V| - 1)/|V|$. The probability of not choosing token i over n samples is then

$$\left(\frac{|V| - 1}{|V|}\right)^n \quad (57)$$

However, We are interested in the opposite event - We want to know the probability of *not* not choosing token i , which is given by

$$1 - \left(\frac{|V| - 1}{|V|}\right)^n \quad (58)$$

By linearity of expectation and the fact that the expected value of an indicator is simply the probability of it being 1, We get

$$E\left(\sum_{i=1}^{|V|} I_i\right) = \sum_{i=1}^{|V|} E(I_i) = \sum_{i=1}^{|V|} P(I_i = 1) = |V| \cdot \left(1 - \left(\frac{|V| - 1}{|V|}\right)^n\right) \quad (59)$$

Q4 a) ii)

Given n draws, what is the probability that all the words will appear in the sample?

All of my attempts have failed, and a whole lot of searching leads me to conclude that the probability is:

$$\frac{|V|!}{|V|^n} S_2(n - 1, |V| - 1) \quad (60)$$

Where $S_2(n - 1, |V| - 1)$ is the Stirling number of the second kind.

Q4 b) i)

Alright, so, let X be the number of samples before "work" and "hard" bigram appears.

Let ps (partial success) mean that the word "work" has been sampled. Then by law of total expectation We can say

$$E(X) = (1 + E(X|ps)) \cdot P(ps) + (1 + E(X|ps')) \cdot P(ps') \quad (61)$$

Where We the 1's in there come from the fact that a single sample has been taken.

Now, given partial success We either fully succeed (fs) or fail and go back to the beginning, so:

$$\begin{aligned} E(X|ps) &= (1 + E(X|ps, fs)) \cdot P(fs|ps) + (1 + E(X|ps, fs')) \cdot P(fs'|ps) \quad (62) \\ &= (1) \cdot \frac{1}{|V|} + (1 + E(X)) \cdot \frac{|V| - 1}{|V|} \quad (63) \end{aligned}$$

Where given full success, expected number of draws to get a success is zero. Given a failure We're back to $E(X)$.

Of course, given ps' , the compliment of partial success i.e. partial failure, $E(X|ps') = E(X)$.

Putting it all together We get:

$$\begin{aligned} E(X) &= (1 + E(X|ps)) \cdot P(ps) + (1 + E(X|ps')) \cdot P(ps') \quad (64) \\ &= \left(1 + \frac{1}{|V|} + (1 + E(X)) \cdot \frac{|V| - 1}{|V|}\right) \cdot \frac{1}{|V|} + (1 + E(X)) \cdot \frac{|V| - 1}{|V|} \quad (65) \\ &= \frac{1}{|V|} + \frac{1}{|V|^2} + (1 + E(X)) \cdot \frac{|V| - 1}{|V|^2} + (1 + E(X)) \cdot \frac{|V| - 1}{|V|} \quad (66) \\ &= \frac{1}{|V|} + \frac{1}{|V|^2} + \frac{|V| - 1}{|V|^2} + \frac{(|V| - 1)E(X)}{|V|^2} + \frac{|V| - 1}{|V|} + \frac{E(X)(|V| - 1)}{|V|} \quad (67) \end{aligned}$$

Now We just need to simplify this mess.

$$E(X) = \frac{1}{|V|} + \frac{1}{|V|^2} + \frac{|V| - 1}{|V|^2} + \frac{(|V| - 1)E(X)}{|V|^2} + \frac{|V| - 1}{|V|} + \frac{E(X)(|V| - 1)}{|V|} \quad (68)$$

$$E(X) - \frac{(|V| - 1)E(X)}{|V|^2} - \frac{E(X)(|V| - 1)}{|V|} = \frac{1}{|V|} + \frac{1}{|V|^2} + \frac{|V| - 1}{|V|^2} + \frac{|V| - 1}{|V|} \quad (69)$$

$$E(X) \left(1 - \frac{(|V| - 1)}{|V|^2} - \frac{(|V| - 1)}{|V|}\right) = \frac{1}{|V|} + \frac{1}{|V|^2} + \frac{|V| - 1}{|V|^2} + \frac{|V| - 1}{|V|} \quad (70)$$

$$E(X) = \frac{\frac{1}{|V|} + \frac{1}{|V|^2} + \frac{|V| - 1}{|V|^2} + \frac{|V| - 1}{|V|}}{1 - \frac{(|V| - 1)}{|V|^2} - \frac{(|V| - 1)}{|V|}} \quad (71)$$

Beautiful.

For a sanity check, of You plug $|V| = 2$ into that equation, You get 6, which is the well known result of expected number of tosses to get 2 heads in a row when tossing a fair coin, so it seems to work at least for that case.

For an even further sanity check, I ran simulations for $|V| = 1 \dots 10$, i.e. perform this experiment 100,000 times and check on average how many samples

it took to receive two of the same word in a row given a uniform distribution, and for $|V| = 1 \dots 10$, it's correct.

Q4 b) ii)

Well, the probability of "work" not appearing in n draws is

$$\left(\frac{|V| - 1}{|V|}\right)^n \quad (72)$$

And of course the probability of it appearing is

$$1 - \left(\frac{|V| - 1}{|V|}\right)^n \quad (73)$$

And Mary wants this probability to be above 0.95:

$$1 - \left(\frac{|V| - 1}{|V|}\right)^n \geq 0.95 \quad (74)$$

$$-\left(\frac{|V| - 1}{|V|}\right)^n \geq -0.05 \quad (75)$$

$$\left(\frac{|V| - 1}{|V|}\right)^n \leq 0.05 \quad (76)$$

$$\log\left(\left(\frac{|V| - 1}{|V|}\right)^n\right) \leq \log(0.05) \quad (77)$$

$$n \cdot \log\left(\frac{|V| - 1}{|V|}\right) \leq \log(0.05) \quad (78)$$

$$n \geq \frac{\log(0.05)}{\log\left(\frac{|V| - 1}{|V|}\right)} \quad (79)$$

$$(80)$$

Where the last sign flips since \log of a quantity less than 1 is negative.

Q4 c) i)

The answer here is identical to the answer to question **b)i)**, by the fact that the tokens are all equally likely.

Q4 c) ii)

Okay so We have

$$h_t^0 = f(w_0 x_{t-1} + w_1 x_t + w_2 h_{t-1}^0 + b_0) \quad (81)$$

So, that's our first hidden state. f is some arbitrary non-linearity (or maybe linearity, whatever, a function).

w_0 multiplies previous input, w_1 multiplies current input, w_2 multiplies the previous hidden state and We have a bias term.

The output then is

$$y_t = g(w_3 h_t^0 + b_1) \quad (82)$$

So w_3 is multiplying the hidden state, there's a bias term and another function.

The goal is to output 1 if x_{t-1} and x_t are the same word.

Okay so assume that the two words are equal and the bias is 0:

$$h_t^0 = f(w_0 x_{t-1} + w_1 x_t + w_2 h_{t-1}^0 + b_0) \quad (83)$$

$$= f(w_0 x_{t-1} + w_1 x_{t-1}) \quad (84)$$

$$= f(x_{t-1}(w_0 + w_1)) \quad (85)$$

So clearly one way to go is to let w_0 be -1 and w_1 be 1, so that when they are added, We get zero in the case of equality. Then f can just be like a check as to whether the input is 0, and if it is output 1, otherwise output 0.

g is then simply the identity function and b_1 is also 0.

So, $w_0 = -1$, $w_1 = 1$, $w_2 = 0$, $w_3 = 1$, all biases are zero, f is a boolean check for whether the input is 0 and g is the identity function.

Q4 c) iii)

So We simply use the second hidden layer as a counter. $w_4, w_5 = 1$, $b_1, b_2 = 0$ and g and h are still the identity function, and f is the boolean 0 check.

Q4 c) iv)

I think the non-uniform unigram language model has the greater probability of drawing two of the same token in a row.

This isn't a proof, but the reasoning is as follows: let p_i be the probability of drawing the i 'th token. Now suppose We alter the distribution such that $p'_i = p_i - x$. Now, We'll have to add that x to some other p_j to preserve a distribution, and since the sample are identical, sampling p_j twice in a row will be more likely than if We had just left the probabilities alone.

Q5 a) i)

For the person correcting this: below is a re-derivation of semiring thing. If You'd like to skip it, scroll down until You see **stop scrolling**. I left the work in since, well, I did it, and so if something went wrong in the derivation, You can at least see what led to the error.

Well first let's re-derive this mess.

$$\sum_{\mathbf{t} \in \mathcal{T}^N} \exp \left\{ \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, \mathbf{w}) \right\} \quad (86)$$

So let's start with inside the bracket comes from

$$\text{score}(\mathbf{t}, \mathbf{w}) = \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, \mathbf{w}) \quad (87)$$

So this is our simplifying assumption, and I suppose it's also what makes this a "Conditional Random Field". \mathbf{t} stands for "tag" and it's the, well, tagging of the word sequence \mathbf{w} . So We're saying that the score for the tag does not have to be calculated all in one go, but rather can be done in parts, where each part only depends on it's predecessor.

N is the length of the sentence by the way, so the length of both \mathbf{t} and \mathbf{w} .

\mathcal{T} is the set of all possible tags for any particular word. Since there are N words, the total number of possible tags for \mathbf{w} is then \mathcal{T}^N .

So then the expression is simply calculating the normalizing constant, since it's summing over all possible tags $\mathbf{t} \in \mathcal{T}^N$.

So, up next We get

$$\sum_{\mathbf{t} \in \mathcal{T}^N} \exp \left\{ \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, \mathbf{w}) \right\} \quad (88)$$

$$= \sum_{\mathbf{t}_{1:n} \in \mathcal{T}^N} \prod_{n=1}^N \exp \{ \text{score}(\langle \mathbf{t}_{n-1}, \mathbf{t}_n \rangle, \mathbf{w}) \} \quad (89)$$

Okay so what have We done here. For one, \mathbf{t} now has a subscript. Sure.

We've converted the inner sum to a product, which is fine - before We were taking the exponent of a sum, but product of a bunch of exponents will sum the exponents so no worries there. We then still sum over every possible tag. Sure.

$$\sum_{\mathbf{t} \in \mathcal{T}^N} \exp \left\{ \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, \mathbf{w}) \right\} \quad (90)$$

$$= \sum_{\mathbf{t}_{1:n} \in \mathcal{T}^N} \prod_{n=1}^N \exp \{ \text{score}(\langle \mathbf{t}_{n-1}, \mathbf{t}_n \rangle, \mathbf{w}) \} \quad (91)$$

$$= \sum_{\mathbf{t}_{1:N-1} \in \mathcal{T}^{N-1}} \sum_{t_N \in \mathcal{T}} \prod_{n=1}^N \exp \{ \text{score}(\langle \mathbf{t}_{n-1}, \mathbf{t}_n \rangle, \mathbf{w}) \} \quad (92)$$

So, the first change is that the subscript for \mathbf{t} now goes to $N - 1$. So We're just taking out the last tag in the sequence of tags \mathbf{t}_N .

Since that last tag could be anything, We sum over all the possible tags, so $t_N \in \mathcal{T}$. There are $|\mathcal{T}|$ choices.

I mean really You could rewrite the first sum as N sums each over $|\mathcal{T}|$ terms each. So We're just splitting one out. No worries.

$$\sum_{\mathbf{t} \in \mathcal{T}^N} \exp \left\{ \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, \mathbf{w}) \right\} \quad (93)$$

$$= \sum_{\mathbf{t}_{1:n} \in \mathcal{T}^N} \prod_{n=1}^N \exp \{ \text{score}(\langle \mathbf{t}_{n-1}, \mathbf{t}_n \rangle, \mathbf{w}) \} \quad (94)$$

$$= \sum_{\mathbf{t}_{1:N-1} \in \mathcal{T}^{N-1}} \sum_{t_N \in \mathcal{T}} \prod_{n=1}^N \exp \{ \text{score}(\langle \mathbf{t}_{n-1}, \mathbf{t}_n \rangle, \mathbf{w}) \} \quad (95)$$

$$= \sum_{\mathbf{t}_{1:N-1} \in \mathcal{T}^{N-1}} \prod_{n=1}^{N-1} \exp \{ \text{score}(\langle \mathbf{t}_{n-1}, \mathbf{t}_n \rangle, \mathbf{w}) \} \times \sum_{t_N \in \mathcal{T}} \exp \{ \text{score}(\langle \mathbf{t}_{N-1}, \mathbf{t}_N \rangle, \mathbf{w}) \} \quad (96)$$

Okay so if I had to explain this in words, how would I do it?

Well, everything hinges on two parts - the first being that the score can be decomposed into partial computations. If there is no substructure, no simplification is possible.

The second I suppose is the distributive property I suppose. The key bit is that You can take stuff out of the product and sum it, and if You take the product of the stuff and the sum it all still works. Or rather, You have some operator B . B is iterating over all possible versions of the tags. Then You also have some operator A , and You can take stuff out of B , combine it with A and take output of A and stick it back into B and it all still works. Now You're reducing the exponential number of things You had to deal with. Sure.

Another way of saying it is that A has to give You the ability to group and it has to work with B via distributativity.

The viterbi algorithm just keeps track of max possible value for each node going backwards using multiplication and max.

Vertices V are the tags, edges E are pairs of tags, and the weight on each tag would be $\text{score}(a, b)$. These should really be directed in our case but it ought to work regardless.

stop scrolling

Data: Graph, source node

Result: Shortest path from source to all nodes

```
prq = initialize priority queue with source node with distance 0;
//scores is a dictionary of tuples such that
//scores[a] = (b, c) where c is the cumulative score of a
//and c is the node We took to arrive at a with score c
scores = empty dictionary;
//initialize scores to include source node scores[source] = (source, 1);
while prq not empty do
    //popMin retrieves and removes minimum entry in prq
    //the minimum is decided by the associated score
    //see prq.add below
    currentNode = prq.popMax() ;
    for neighbour in neighbours(currentNode) do
        //scores[neighbour](1) accesses the tuple for
        //neighbour and (1) accesses the cumulative score
        tempScore = scores[neighbour](1) + score(neighbour,
            currentNode)
        if tempScore > score[neighbour] then
            scores[neighbour] = (currentNode, tempScore)
            //tempScore is what is used to rank the options
            //and prq returns the neighbour entry
            //when popMax is called
            //loops will not enter prq due to increasing score
            //and anyway the graph We care about is acyclic
            prq.add(neighbour, tempScore)
        end
    end
end
```

Algorithm 1: Dijkstra's algorithm using a priority queue.

To get the best path, one simply picks the node they would like to find the path from, find the node in the scores list, look at it's predecessor, then look at the predecessor of the predecessor and so on.

Q5 a) ii)

Well the algorithm will visit every node once, so it's $\mathcal{O}(|V|)$ for that.

It'll also examine every edge twice - let a, b , be nodes, then it will look at the edge (a, b) first and sometime later at the edge (b, a) , so $2|E|$ there or $\mathcal{O}(|E|)$. In those cases it will also consult the priority queue prq , and the complexity of

that in the best case¹ is $\log(|V|)$, so in total We get $\mathcal{O}(|E| \cdot \log(|V|) + |V|)$.

In contrast, Viterbi visits each edge and node once, so We have $\mathcal{O}(|V| + |E|)$ time complexity. This obviously compares favourably with Dijkstra.

Now, since Dijkstra explores more promising venues first, it is possible to construct examples where, despite being worse asymptotically, Dijkstra will beat Viterbi.

Likewise one could construct examples where Viterbi will beat Dijkstra. Neither is strictly better than the other, but Viterbi is better on average.

Q5 b) i)

Well first let's identify the semiring.

A are strictly positive real numbers of course.

\oplus is the max operation.

\otimes is the \times operation, i.e. just simple multiplication. This works due to the fact that We're exponentiating, so taking the product leads to the sum in the exponent We want.

$\hat{0}, \hat{1}$ are 0 and 1 respectively.

In semiring notation We then get an algorithm that looks basically nothing like the first version:

Data: Graph G , score function, source node s

Result: Shortest path from source to all nodes

prq = initialize max queue with source node and score of 1;

scores = dictionary initialized with source node and value of 1;

while Q not empty **do**

 currentNode = prq.popMax();

for each neighbour of currentNode **do**

 tempScore = score[neighbour] \oplus (scores(currentNode) \otimes score(neighbour));

if scores[neighbour](1) < tempScore **then**

 //scores[neighbour]= parent, new best score

 scores[neighbour] =(currentNode, tempScore);

 //priority decided based on score obviously

 prq.add(neighbour, tempScore);

end

end

end

Algorithm 2: Dijkstra's algorithm using a priority queue in semiring notation.

It's a little different since I don't like initializing the queue with the entire graph at the beginning. Instead, the queue is initialized with just the source and entries are added as You go. This could be dicey in cyclic settings, but We don't have to worry about cycles in our case. Don't think I'd use this with a cyclic graph, since cycles would just continue increasing the score and get continuously added.

¹With Fibonacci heaps.

To get the path one traverses the parents in scores, same as before - each node has a parent, and to get the path You take the parent of the goal node, then the parent of the parent and so on.

Q5 b) ii)

I think just using $\oplus = \min$ operation should work, and of course the priority queue would now use popMin instead of popMax, and also the if condition changes to $>$. I think this would not work in general due to cycles, but our graph is acyclic.

Q5 b) iii)

Well, I think the $\otimes = \min$, since given a path We'll want to take the minimum out of all the edge values.

Could You use $\oplus = \max$? I think that's it. Across an entire path We want to take the minimum to find the minimal value, but when choosing amongst a number of possible edges to take, We'll choose the maximum.

2 Assignment # 2

Q1 a)

$$S \rightarrow NP, VP \quad (1)$$

$$NP \rightarrow Det, N \mid NP, PP \mid "I" \mid "glasses" \quad (2)$$

$$VP \rightarrow VP, PP \mid V, NP \quad (3)$$

$$Det \rightarrow "a" \mid "an" \quad (4)$$

$$N \rightarrow "man" \mid "pencil" \mid "ball" \mid "umbrella" \quad (5)$$

$$PP \rightarrow P, NP \quad (6)$$

$$V \rightarrow "draw" \mid "hit" \quad (7)$$

$$P \rightarrow "with" \quad (8)$$

Q1 b)

Warning: advanced calculations ahead (instead of probabilities We just do counts first):

$$S \rightarrow NP, VP(4) \quad (9)$$

$$NP(14) \rightarrow Det, N(7) \mid NP, PP(2) \mid "I"(4) \mid "glasses"(1) \quad (10)$$

$$VP(6) \rightarrow VP, PP(2) \mid V, NP(4) \quad (11)$$

$$Det(7) \rightarrow "a"(6) \mid "an"(1) \quad (12)$$

$$N(7) \rightarrow "man"(4) \mid "pencil"(1) \mid "ball"(1) \mid "umbrella"(1) \quad (13)$$

$$PP(4) \rightarrow P, NP(4) \quad (14)$$

$$V(4) \rightarrow "draw"(2) \mid "hit"(2) \quad (15)$$

$$P(4) \rightarrow "with"(4) \quad (16)$$

Convert to probabilities (rounded, for precise figures just divide counts on the right by total counts on the left):

$$S \rightarrow NP, VP(1.0) \quad (17)$$

$$NP \rightarrow Det, N(0.5) \mid NP, PP(0.14) \mid "I"(0.29) \mid "glasses"(0.7) \quad (18)$$

$$VP \rightarrow VP, PP(2 = 0.33) \mid V, NP(0.66) \quad (19)$$

$$Det \rightarrow "a"(0.86) \mid "an"(0.14) \quad (20)$$

$$N \rightarrow "man"(0.57) \mid "pencil"(0.14) \mid "ball"(0.14) \mid "umbrella"(0.14) \quad (21)$$

$$PP \rightarrow P, NP(1) \quad (22)$$

$$V \rightarrow "draw"(0.5) \mid "hit"(0.5) \quad (23)$$

$$P \rightarrow "with"(1) \quad (24)$$

Q1 c)

Our goal is to have different likelihood of expansions on noun phrase depending on whether the expansion is going to be an object or a subject.

I suppose a logical way of going about it would be to replace the noun phrase non-terminal with two other non terminals - object phrase (*OP*) and subject phrase (*SP*).

$$S \rightarrow NP, VP(1.0) \quad (25)$$

$$NP \rightarrow SP \mid OP \quad (26)$$

$$SP \rightarrow Det, N(0.5) \mid NP, PP(0.14) \mid "I"(0.29) \mid "glasses"(0.7) \quad (27)$$

$$OP \rightarrow Det, N(0.5) \mid NP, PP(0.14) \mid "I"(0.29) \mid "glasses"(0.7) \quad (28)$$

$$VP(\rightarrow VP, PP(2 = 0.33) \mid V, NP(0.66) \quad (29)$$

$$Det \rightarrow "a"(0.86(\mid "an"(0.14) \quad (30)$$

$$N \rightarrow "man"(0.57) \mid "pencil"(0.14) \mid "ball"(0.14) \mid "umbrella"(0.14) \quad (31)$$

$$PP \rightarrow P, NP(1) \quad (32)$$

$$V \rightarrow "draw"(0.5) \mid "hit"(0.5) \quad (33)$$

$$P \rightarrow "with"(1) \quad (34)$$

Here are the counts as before:

$$S \rightarrow NP, VP(1.0) \quad (35)$$

$$NP(14) \rightarrow SP(4) \mid OP(10) \quad (36)$$

$$SP(4) \rightarrow Det, N(0) \mid NP, PP(0) \mid "I"(4) \mid "glasses"(0) \quad (37)$$

$$OP(10) \rightarrow Det, N(7) \mid NP, PP(2) \mid "I"(0) \mid "glasses"(1) \quad (38)$$

$$VP(\rightarrow VP, PP(2 = 0.33) \mid V, NP(0.66) \quad (39)$$

$$Det \rightarrow "a"(0.86(\mid "an"(0.14) \quad (40)$$

$$N \rightarrow "man"(0.57) \mid "pencil"(0.14) \mid "ball"(0.14) \mid "umbrella"(0.14) \quad (41)$$

$$PP \rightarrow P, NP(1) \quad (42)$$

$$V \rightarrow "draw"(0.5) \mid "hit"(0.5) \quad (43)$$

$$P \rightarrow "with"(1) \quad (44)$$

Finally converting to probabilities:

$$S \rightarrow NP, VP(1.0) \quad (45)$$

$$NP \rightarrow SP(0.29) \mid OP(0.71) \quad (46)$$

$$SP(4) \rightarrow Det, N(0) \mid NP, PP(0) \mid "I"(1) \mid "glasses"(0) \quad (47)$$

$$OP(10) \rightarrow Det, N(0.7) \mid NP, PP(0.2) \mid "I"(0) \mid "glasses"(0.1) \quad (48)$$

$$VP(\rightarrow VP, PP(2 = 0.33) \mid V, NP(0.66) \quad (49)$$

$$Det \rightarrow "a"(0.86) \mid "an"(0.14) \quad (50)$$

$$N \rightarrow "man"(0.57) \mid "pencil"(0.14) \mid "ball"(0.14) \mid "umbrella"(0.14) \quad (51)$$

$$PP \rightarrow P, NP(1) \quad (52)$$

$$V \rightarrow "draw"(0.5) \mid "hit"(0.5) \quad (53)$$

$$P \rightarrow "with"(1) \quad (54)$$

Well this at least shows that subject phrases are likely to be "I", and object phrases are unlikely to be "I", so it's capturing some of what We wanted to.

There is a problem here of course - nothing is stopping anyone from using *SP* instead of *OP* wherever they like, so the onus is on the user of the grammar to use the subject non-terminal for subjects and the object non-terminal for objects.

Q1 d)

$$S \rightarrow NP, VP^+ \quad (55)$$

$$NP \rightarrow Det, N^+ \mid NP^+, PP \mid "I" \mid "glasses" \quad (56)$$

$$VP \rightarrow VP^+, PP \mid V^+, NP \quad (57)$$

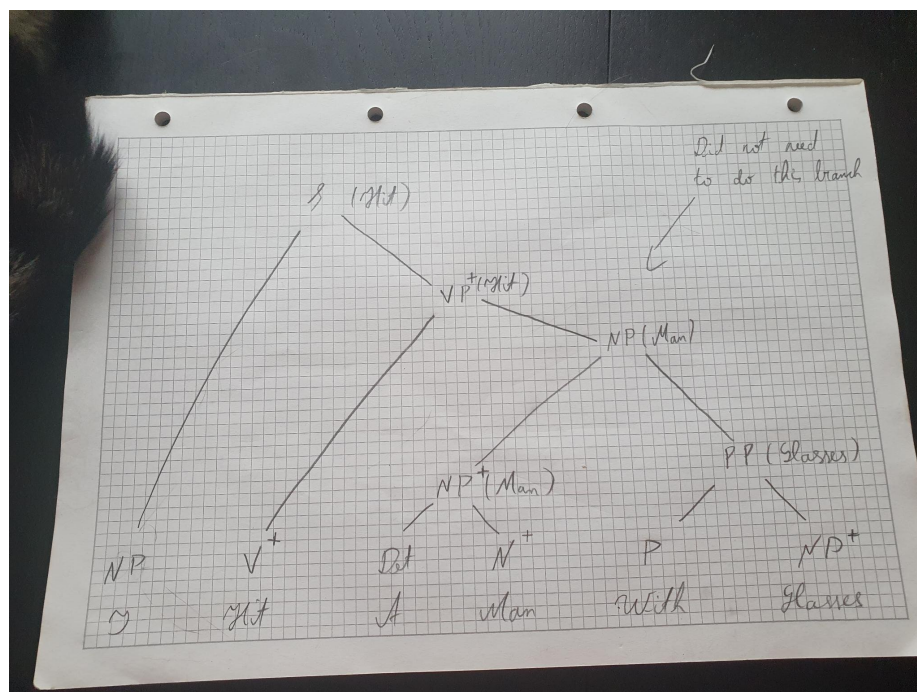
$$Det \rightarrow "a" \mid "an" \quad (58)$$

$$N \rightarrow "man" \mid "pencil" \mid "ball" \mid "umbrella" \quad (59)$$

$$PP \rightarrow P, NP^+ \quad (60)$$

$$V \rightarrow "draw" \mid "hit" \quad (61)$$

$$P \rightarrow "with" \quad (62)$$



Q2 a)

So, it is my understanding that for the projective dependency tree, to obtain the score We would simply take all dependency pairs and multiply the scores corresponding to those dependencies (taking care to pay attention to the direction of the dependency etc.)

To obtain the score of a lexicalized probabilistic context free grammar tree, one also simply multiplies the score for each "rule" applied, except this rule is also lexicalized so You have very many parameters indeed (page 15 from - click me-.)

With all that out of the way, I would construct a weighted lexicalized CFG as follows:

1. Add root productions $R \rightarrow X_j, j \in 1 \dots M$.
2. Add self-productions $X_j \rightarrow j$.
3. For each production $\psi(i \rightarrow j)$, add bidirectional productions to the language as follows:

$$X_i \rightarrow X_i^+, X_j \quad (63)$$

$$X_i \rightarrow X_j, X_i^+ \quad (64)$$

And of course the weight for that transformation corresponds to the value of $\psi(i \rightarrow j)$ does not tell you anything about the ordering of words within a sentence.

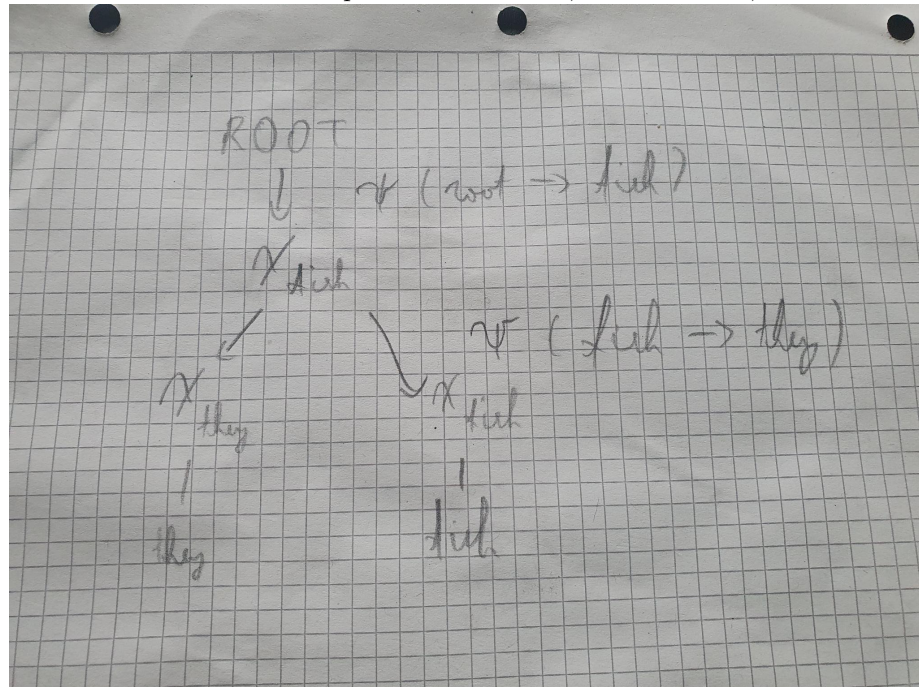
The reason for adding bidirectional rules is that just having access to the set of weights generated by

And that's all she wrote.

Q2 b)

So the dependency tree would simply be $ROOT \rightarrow fish$, $fish \rightarrow they$, and the score would be $\psi(ROOT, fish) \cdot \psi(fish, they)$.

In our CFG we would end up with $ROOT \rightarrow S$, then $S \rightarrow NP$,



Q4 a)

number	sample strings	accepted	weight
1	educational is this not	NO	
2	is this assignment educational	NO	
3	not educational is not educational	YES	12
4	this assignment is not educational	YES	7
5	is this assignment educational	NO	
6	this assignment course is educational	NO	
7	is this assignment not educational	YES	15
8	this assignment not	YES	8
9	this course assignment is not educational	YES	12
10	this course is not not educational	YES	21
11	not educational is this	YES	10
12	course assignment is not educational	YES	10
13	not this assignment is educational	NO	
14	not not not educational	YES	25
15	is this course assignment not educational	YES	18
16	course assignment is this	YES	8
17	this course is interesting	NO	
18	this course assignment not educational	YES	14

Table 1: Some strings from $\mathcal{V}_{\geq 2, \leq 6}$

Q4 b)

Cost matrices on left, backtracking matrices on the right.

"inf" means an infinite cost, i.e. no path exists, and similarly "N" in the backtracking matrix means a backtrack is not possible.

Iteration 0

-	A	B	C	D	E	F	-	A	B	C	D	E	F
A	inf	1	2	inf	inf	inf	A	-	B	C	-	-	-
B	inf	inf	inf	1	5	inf	B	-	-	-	D	E	-
C	inf	2	1	2	3	inf	C	-	B	C	D	E	-
D	4	inf	inf	1	2	4	D	A	-	-	D	E	F
E	inf	inf	inf	inf	10	2	E	-	-	-	-	E	F
F	inf	inf	inf	3	inf	inf	F	-	-	-	D	-	-

Iteration 1

-	A	B	C	D	E	F
A	inf	1	2	inf	inf	inf
B	inf	inf	inf	1	5	inf
C	inf	2	1	2	3	inf
D	4	5	6	1	2	4
E	inf	inf	inf	inf	10	2
F	inf	inf	inf	3	inf	inf

-	A	B	C	D	E	F
A	-	B	C	-	-	-
B	-	-	-	D	E	-
C	-	B	C	D	E	-
D	A	A	A	D	E	F
E	-	-	-	-	E	F
F	-	-	-	D	-	-

Iteration 2

-	A	B	C	D	E	F
A	inf	1	2	2	6	inf
B	inf	inf	inf	1	5	inf
C	inf	2	1	2	3	inf
D	4	5	6	1	2	4
E	inf	inf	inf	inf	10	2
F	inf	inf	inf	3	inf	inf

-	A	B	C	D	E	F
A	-	B	C	B	B	-
B	-	-	-	D	E	-
C	-	B	C	D	E	-
D	A	A	A	D	E	F
E	-	-	-	-	E	F
F	-	-	-	D	-	-

Iteration 3

-	A	B	C	D	E	F
A	inf	1	2	2	5	inf
B	inf	inf	inf	1	5	inf
C	inf	2	1	2	3	inf
D	4	5	6	1	2	4
E	inf	inf	inf	inf	10	2
F	inf	inf	inf	3	inf	inf

-	A	B	C	D	E	F
A	-	B	C	B	C	-
B	-	-	-	D	E	-
C	-	B	C	D	E	-
D	A	A	A	D	E	F
E	-	-	-	-	E	F
F	-	-	-	D	-	-

Iteration 4

-	A	B	C	D	E	F
A	6	1	2	2	4	6
B	5	6	7	1	3	5
C	6	2	1	2	3	6
D	4	5	6	1	2	4
E	inf	inf	inf	inf	10	2
F	7	8	9	3	5	7

-	A	B	C	D	E	F
A	B	B	C	B	B	B
B	D	D	D	D	D	D
C	D	B	C	D	E	D
D	A	A	A	D	E	F
E	-	-	-	-	E	F
F	D	D	D	D	D	D

Iteration 5

-	A	B	C	D	E	F	-	A	B	C	D	E	F
A	6	1	2	2	4	6	A	B	B	C	B	B	B
B	5	6	7	1	3	5	B	D	D	D	D	D	D
C	6	2	1	2	3	5	C	D	B	C	D	E	E
D	4	5	6	1	2	4	D	A	A	A	D	E	F
E	inf	inf	inf	inf	10	2	E	-	-	-	-	E	F
F	7	8	9	3	5	7	F	D	D	D	D	D	D

Final 6th iteration

-	A	B	C	D	E	F	-	A	B	C	D	E	F
A	6	1	2	2	4	6	A	B	B	C	B	B	B
B	5	6	7	1	3	5	B	D	D	D	D	D	D
C	6	2	1	2	3	5	C	D	B	C	D	E	E
D	4	5	6	1	2	4	D	A	A	A	D	E	F
E	9	10	11	5	7	2	E	F	F	F	F	F	F
F	7	8	9	3	5	7	F	D	D	D	D	D	D

Q4 c)

The number of iterations is bound by $\mathcal{O}(N^3)$, where N is the number of nodes.

The algorithm terminated after 6 iterations, so after looking at all of the nodes.

Q4 d)

The time complexity is $\mathcal{O}(N^3)$ and for the method I used the space complexity is N^2 . Strictly speaking it's more like $4N^2$ since I was keeping new and old versions of matrices around in my implementation, but this can be reduced. Asymptotically, it's simply $\mathcal{O}(N^2)$ for space complexity and $\mathcal{O}(N^3)$ for time complexity.

Backtracking takes $\mathcal{O}(n)$ and is achieved as follows: let starting node be u and ending node be v . Now take the v 'th column of the backtracking matrix, and begin at the u 'th entry. The value of the u 'th entry is the next node in the shortest path. Then simply follow the entries until the desired destination is reached.

This makes sense intuitively - it may be the case that to go from u to v one has to traverse all the nodes, but no node should be traversed more than once since if it were, then there is a cycle in the path that could be cut out.