

# Assignment 3

(Due: Nov 30, 6am)

**Question 1:** Let  $\mathbf{Q} \in \mathbb{R}^{N \times d}$  denote a set of  $N$  query vectors, which attend to  $M$  key and value vectors, denoted by matrices  $\mathbf{K} \in \mathbb{R}^{M \times d}$  and  $\mathbf{V} \in \mathbb{R}^{M \times c}$  respectively. For a query vector at position  $n$ , the softmax attention function computes the following quantity:

$$\text{Attn}(\mathbf{q}_n, \mathbf{K}, \mathbf{V}) = \sum_{m=1}^M \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_m)}{\sum_{m'=1}^M \exp(\mathbf{q}_n^\top \mathbf{k}_{m'})} \mathbf{v}_m^\top := \mathbf{V}^\top \text{softmax}(\mathbf{K} \mathbf{q}_n), \quad (1)$$

which is an average of the set of value vectors  $\mathbf{V}$  weighted by normalized similarity between different queries and keys.

**Question 1a:** Please briefly explain what is the time and space complexity for the attention computation from query  $\mathbf{Q}$  to  $\mathbf{K}, \mathbf{V}$ , using the big  $O$  notation.

**Question 2:** Consider a probabilistic context-free grammar with the following rules (assume that S is the start symbol):

$S \rightarrow NP VP$	1.0
$VP \rightarrow V_t NP$	0.7
$VP \rightarrow VP PP$	0.3
$NP \rightarrow DT NN$	0.8
$NP \rightarrow NP PP$	0.2
$PP \rightarrow IN NP$	1.0
$Vi \rightarrow \text{sleeps}$	1.0
$V_t \rightarrow \text{saw}$	1.0
$NN \rightarrow \text{man}$	0.1
$NN \rightarrow \text{woman}$	0.1
$NN \rightarrow \text{telescope}$	0.3
$NN \rightarrow \text{dog}$	0.5
$DT \rightarrow \text{the}$	1.0
$IN \rightarrow \text{with}$	0.6
$IN \rightarrow \text{in}$	0.4

**Question 2a:** What's the most likely parse tree for the following sentence under this PCFG? Show CYK chart you developed below.

*the man saw the woman with the dog*

**Question 2b:** What's the (marginal) probability of the following sentence under this PCFG?

*the man saw the woman with the dog*

**Question 3:** A trigram language model is also often referred to as a second-order Markov language model. It has the following form:

$$P(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i \mid X_{i-2} = x_{i-2}, X_{i-1} = x_{i-1})$$

**Question 3a:** Could you briefly explain the advantages and disadvantages of a high-order Markov language model when comparing with the second-order one?

**Question 3b:** Could you give some examples in English where English grammar suggests that the second-order Markov assumption is clearly violated.