

**Ασκήσεις μελέτης της ενότητας Β2 (γλωσσικά μοντέλα  $n$ -γραμμάτων, διόρθωση και πρόβλεψη κειμένου)**

1. Επιβεβαιώστε τους υπολογισμούς της διαφάνειας «Υπολογισμός απόστασης Levenshtein». Δοκιμάστε και άλλα ζεύγη λέξεων. Υλοποιήστε τον αλγόριθμο. Υπόδειξη: Χρειάζονται απλά δύο φωλιασμένοι βρόχοι επανάληψης.

2. (α) Χρησιμοποιώντας τις παρακάτω προτάσεις ως (μικροσκοπικό) σώμα κειμένων εκπαίδευσης:

<start> he plays football

<start> he plays cricket

<start> she enjoys good football

<start> she plays good music

<start> he prays to god

<start> please buy me the other ball

<start> he pleases the other players by playing good football

εκτιμήστε τις πιθανότητες  $P(t_1^4)$  που θα επέστρεφε ένα γλωσσικό μοντέλο διγραμμάτων με εξομάλυνση Laplace για κάθε μία από τις δύο παρακάτω προτάσεις  $t_1^4$ :

$t_1^4$ : <start> he please god football

$t_1^4$ : <start> he plays good football

Υποθέστε ότι το λεξιλόγιο  $V$  περιέχει όλες τις λέξεις του σώματος κειμένων (εξαιρώντας το <start>), οπότε  $|V| = 21$ . Δείξτε λεπτομερώς τους υπολογισμούς σας, χωρίς να εκτελέσετε τις τελικές αριθμητικές πράξεις.

Απάντηση:

Το γλωσσικό μοντέλο διγραμμάτων εκτιμά την  $P(\langle \text{start} \rangle, \text{he}, \text{please}, \text{god}, \text{football})$  ως εξής:

$$P(\text{he} | \langle \text{start} \rangle) P(\text{please} | \text{he}) P(\text{god} | \text{please}) P(\text{football} | \text{god}) =$$

$$(4+1)/(7+21) (0+1)/(4+21) (0+1)/(1+21) (0+1)/(1+21)$$

όπου χρησιμοποιήσαμε εξομάλυνση Laplace.

Ομοίως για την  $P(\langle \text{start} \rangle, \text{he}, \text{plays}, \text{good}, \text{football})$ . (Κάντε τους υπολογισμούς μόνοι σας.)

Σημείωση: Στην πράξη αποφεύγουμε διαδοχικούς πολλαπλασιασμούς πιθανοτήτων, γιατί συχνά οδηγούν σε πολύ μικρούς αριθμούς που δεν μπορούν να παρασταθούν με ακρίβεια στον υπολογιστή. Γι' αυτό υπολογίζουμε συνήθως τον λογάριθμο της πιθανότητας μιας ακολουθίας λέξεων, δηλαδή θα υπολογίζαμε το  $\log P(\langle \text{start} \rangle, \text{he}, \text{please}, \text{god}, \text{football})$ , αντί του  $P(\langle \text{start} \rangle, \text{he}, \text{please}, \text{god}, \text{football})$ , οπότε θα καταλήγαμε στο παρακάτω άθροισμα τεσσάρων λογαρίθμων, αντί του παραπάνω γινομένου τεσσάρων πιθανοτήτων.

$$\log[(4+1)/(7+21)] + \log[(0+1)/(4+21)] + \log[(0+1)/(1+21)] + \log[(0+1)/(1+21)]$$

(β) Υποθέστε ότι ένας χρήστης έγραψε στο πληκτρολόγιο του κινητού του την παρακάτω ακολουθία λέξεων  $w_1^4$ :

$w_1^4$ : <start> he pls gd fball

Εκτιμήστε τις πιθανότητες  $P(t_1^4 | w_1^k)$  των δύο υποθέσεων (ακολουθιών λέξεων που ίσως ήθελε να γράψει)  $t_1^4$  του σκέλους (α), χρησιμοποιώντας ένα μοντέλο θορυβώδους καναλιού (βλ. σχετικές διαφάνειες) και το γλωσσικό μοντέλο διγραμμάτων του σκέλους (α). Θεωρήστε ότι  $P(w_i | t_i) \cong \frac{1}{LD(w_i, t_i) + 1}$ , όπου  $LD(w_i, t_i)$  η απόσταση Levenshtein από τη λέξη  $w_i$  στην  $t_i$ . Δείξτε λεπτομερώς τους υπολογισμούς σας, χωρίς να εκτελέσετε τις τελικές αριθμητικές πράξεις και χωρίς να υπολογίσετε τις αποστάσεις Levenshtein.

Απάντηση: Χρησιμοποιώντας το θορυβώδες κανάλι των διαφανειών της διάλεξης, έχουμε:

$$P(t_1^4 | w_1^4) = P(t_1^4) P(w_1^4 | t_1^4) / P(w_1^4)$$

Για  $t_1^4 = \text{<start> he please god football}$ :

$P(\text{<start>, he, please, god, football})$

$$P(\text{<start> he pls gd fball} | \text{<start>, he, please, god, football}) / P(w_1^4)$$

Η πιθανότητα  $P(\text{<start>, he, please, god, football})$  εκτιμάται από το γλωσσικό μοντέλο, όπως στο σκέλος (α).

Χρησιμοποιώντας την προσέγγιση  $P(w_i | t_i) \cong \frac{1}{LD(w_i, t_i) + 1}$  της εκφώνησης (βλ. και διαφάνειες), η πιθανότητα  $P(\text{<start> he pls gd fball} | \text{<start>, he, please, god, football})$  γίνεται:

$$P(\text{he, he}) P(\text{pls, please}) P(\text{gd, god}) P(\text{fball, football}) =$$

$$1/(LD(\text{he, he})+1) \cdot 1/(LD(\text{pls, please})+1) \cdot 1/(LD(\text{gd, god})+1) \cdot 1/(LD(\text{fball, football})+1)$$

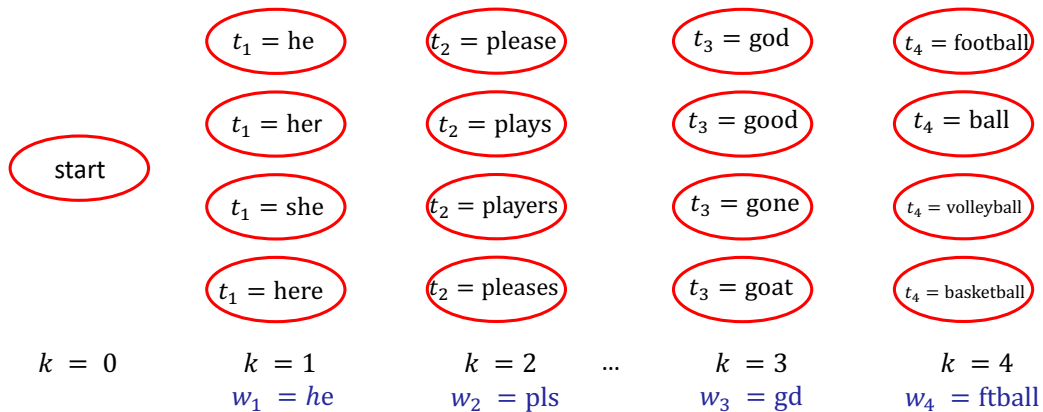
Η πιθανότητα  $P(w_1^4)$  δεν χρειάζεται να εκτιμηθεί, γιατί είναι ίδια και για τις δύο υποθέσεις  $t_1^4 = \text{<start> he please god football}$  και  $t_1^4 = \text{<start> he plays good football}$ .

Ομοίως εκτιμούμε την πιθανότητα  $P(t_1^4 | w_1^4)$  για την υπόθεση  $t_1^4 = \text{<start> he plays good football}$ . (Κάντε τους υπολογισμούς μόνοι σας.) Επιλέγουμε τελικά την υπόθεση  $t_1^4$  με το μεγαλύτερο  $P(t_1^4 | w_1^4)$ .

Σημείωση: Και πάλι στην πράξη θα υπολογίζαμε τον λογάριθμο κάθε γινομένου πιθανοτήτων, οπότε θα καταλήγαμε σε αθροίσματα λογαρίθμων πιθανοτήτων, αντί γινόμενα πιθανοτήτων.

(γ) Εξηγήστε αναλυτικά πώς θα γινόταν η αποκωδικοποίηση με beam search (διαφάνειες «Beam search decoder»), αν ο χρήστης γράψει στο πληκτρολόγιο την ακολουθία λέξεων  $w_1^4$  του σκέλους (β). Χρησιμοποιούμε πάλι το γλωσσικό μοντέλο διγραμμάτων λέξεων του σκέλους (α), εκπαιδευμένο στο μικροσκοπικό σώμα εκπαίδευσης εκείνου του σκέλους, μαζί με το μοντέλο θορυβώδους καναλιού του σκέλους (β). Θεωρήστε ότι το πλέγμα (lattice) αναζήτησης είναι το ακόλουθο, δηλαδή περιλαμβάνει 4 κοντινές (κατά απόσταση διόρθωσης) υποψήφια σωστές λέξεις (του λεξικού), για κάθε λέξη  $w_i$  που έχει γράψει ο χρήστης. Σε κάθε βήμα του beam search, κρατάμε τα  $b = 2$  καλύτερα μονοπάτια.

# Beam search decoder



Για  $k = 1$ , τα υποψήφια μονοπάτια  $t_1^k$  είναι τα ακόλουθα. Δίπλα στο καθένα, υπολογίζουμε την ποσότητα  $P(t_1^k)P(w_1^k|t_1^k)$ . (Στη πραγματικότητα θα υπολογίζαμε τον λογάριθμο αυτής της ποσότητας.) Τα δύο καλύτερα μονοπάτια σημειώνονται με «\*\*».

$\langle \text{start, he} \rangle$ :  $P(\text{he}|\text{start}) P(\text{he}|\text{he}) = (4+1)/(7+21) 1/(0+1) = 5/28 = 0.179$  \*\*  
 $\langle \text{start, her} \rangle$ :  $P(\text{her}|\text{start}) P(\text{he}|\text{her}) = (0+1)/(7+21) 1/(1+1) = 1/28 1/2 = 0.018$   
 $\langle \text{start, she} \rangle$ :  $P(\text{she}|\text{start}) P(\text{he}|\text{she}) = (2+1)/(7+21) 1/(1+1) = 3/28 1/2 = 0.054$  \*\*  
 $\langle \text{start, here} \rangle$ :  $P(\text{here}|\text{start}) P(\text{he}|\text{here}) = (0+1)/(7+21) 1/(2+1) = 1/28 1/3 = 0.012$

Για  $k = 2$ , τα υποψήφια μονοπάτια  $t_1^k$  είναι τα ακόλουθα. Δίπλα στο καθένα, υπολογίζουμε πάλι την ποσότητα  $P(t_1^k)P(w_1^k|t_1^k)$ . Συμπληρώστε μόνοι σας τους υπολογισμούς. Σημειώστε πάλι τα δύο καλύτερα μονοπάτια με «\*\*».

$\langle \text{start, he, please} \rangle$ :  $P(\text{he}|\text{start}) P(\text{he}|\text{he}) P(\text{please}|\text{he}) P(\text{pls}|\text{please}) = \dots$   
 $\langle \text{start, he, plays} \rangle$ :  $P(\text{he}|\text{start}) P(\text{he}|\text{he}) P(\text{plays}|\text{he}) P(\text{pls}|\text{plays}) = \dots$   
 $\langle \text{start, he, players} \rangle$ :  $P(\text{he}|\text{start}) P(\text{he}|\text{he}) P(\text{players}|\text{he}) P(\text{pls}|\text{players}) = \dots$   
 $\langle \text{start, he, pleases} \rangle$ :  $P(\text{he}|\text{start}) P(\text{he}|\text{he}) P(\text{pleases}|\text{he}) P(\text{pls}|\text{pleases}) = \dots$   
 $\langle \text{start, she, please} \rangle$ :  $P(\text{she}|\text{start}) P(\text{he}|\text{she}) P(\text{please}|\text{she}) P(\text{pls}|\text{please}) = \dots$   
 $\langle \text{start, she, plays} \rangle$ :  $P(\text{she}|\text{start}) P(\text{he}|\text{she}) P(\text{plays}|\text{she}) P(\text{pls}|\text{plays}) = \dots$   
 $\langle \text{start, she, players} \rangle$ :  $P(\text{she}|\text{start}) P(\text{he}|\text{she}) P(\text{players}|\text{she}) P(\text{pls}|\text{players}) = \dots$   
 $\langle \text{start, she, pleases} \rangle$ :  $P(\text{she}|\text{start}) P(\text{he}|\text{she}) P(\text{pleases}|\text{she}) P(\text{pls}|\text{pleases}) = \dots$

Για  $k = 3$ , τα υποψήφια μονοπάτια  $t_1^k$  είναι τα ακόλουθα. Δίπλα στο καθένα, υπολογίζουμε πάλι την ποσότητα  $P(t_1^k)P(w_1^k|t_1^k)$ . Συμπληρώστε μόνοι σας τους υπολογισμούς. Σημειώστε πάλι τα δύο καλύτερα μονοπάτια με «\*\*».

Συμπληρώστε και τους υπολογισμούς για  $k = 4$ . Τελικά ποια υπόθεση  $\langle \text{start, } t_1, t_2, t_3, t_4 \rangle$  επιλέγεται;