

NLP – Assignment 1

Group 4

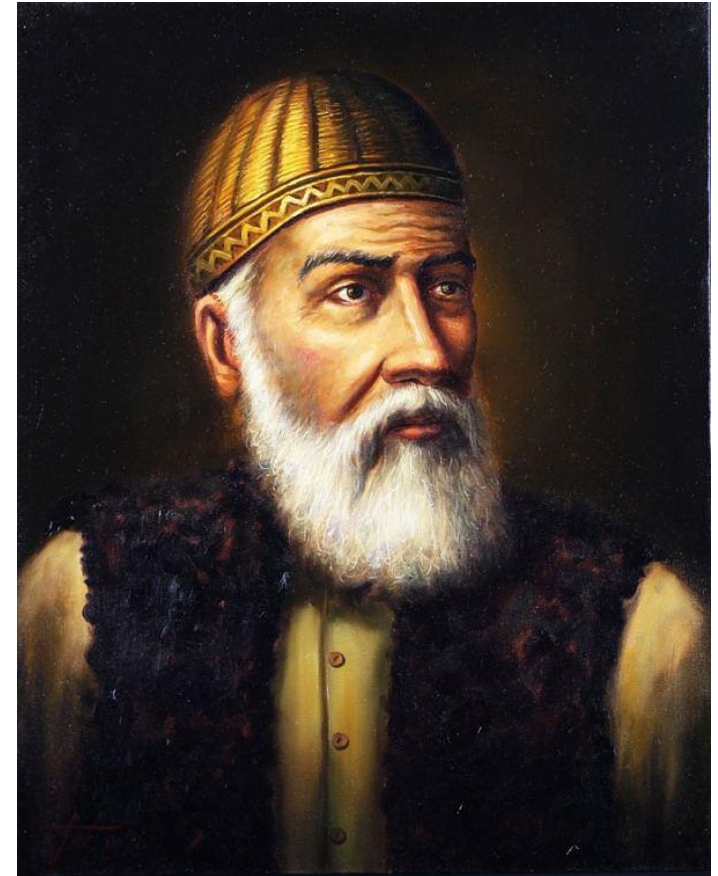
Dataset – Fuzuli's poems

“Qəd ənarəl-eşqə-lil-üşşaqi minhacəl hüdä!

Saliki-rahi-həqiqət eşqə eylər iqtida.

Eşqdir ol nəş'eyi-kamil kim, ondandır müdam

Meydə təşviri-hərarət, neydə tə'siri-səda”...



Füzuli

Data description

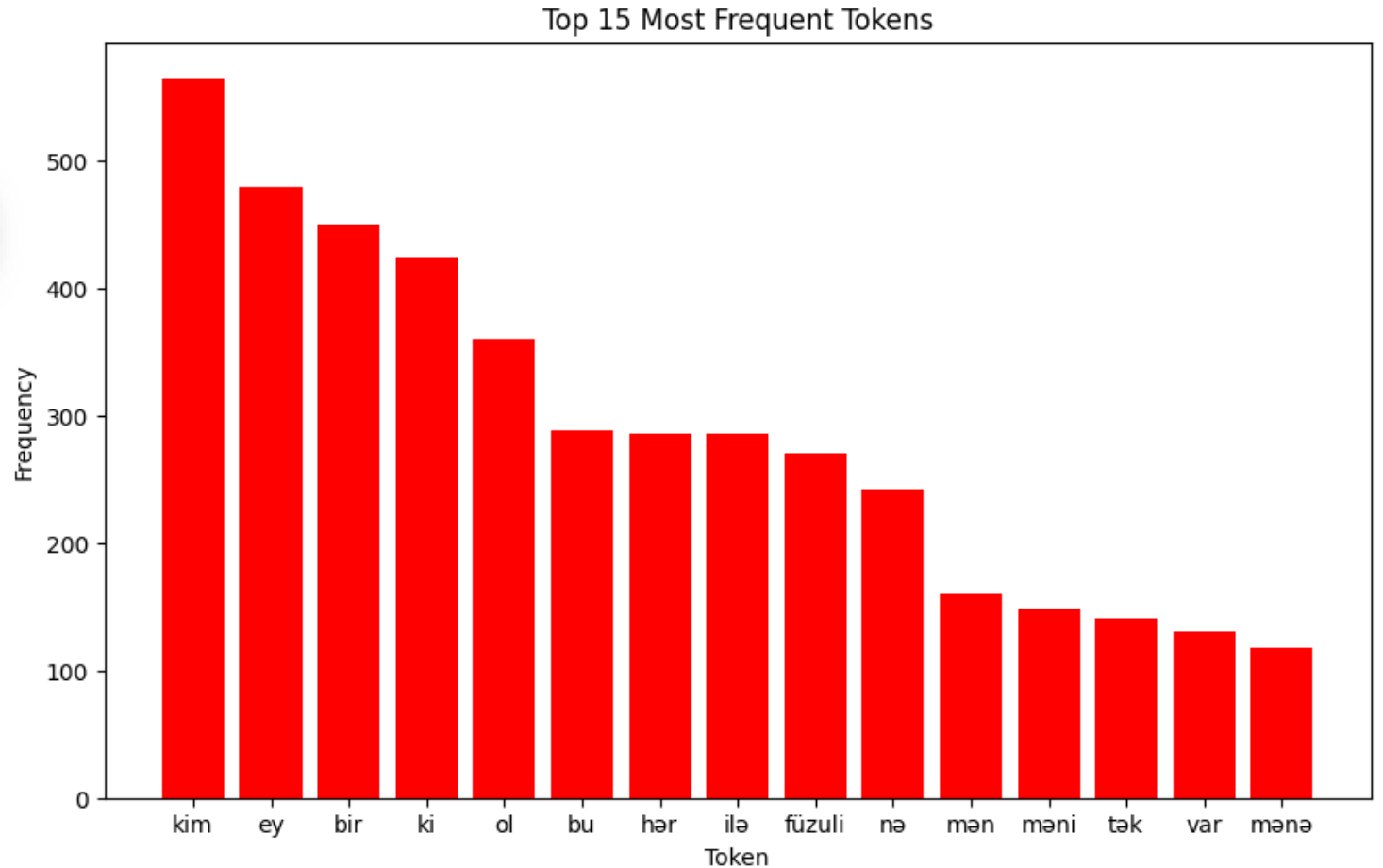
- Motivation : complex word forms, possible usage in linguistic research
- Situation: poetry written in 16 century
- Language: classical Azerbaijani
- Speaker: Demographics: written by Məhəmməd Füzuli
- Collection process: Data was sourced from open presidential library. Conversion from pdf to plain text
- Annotation Process: no annotation
- Distribution: Original text is public domain. The source for the project is © “ŞƏRQ-QƏRB”, 2005.

Whitespace Tokenization

```
r"\p{L}+(?:[-'']\p{L}+)*"
```

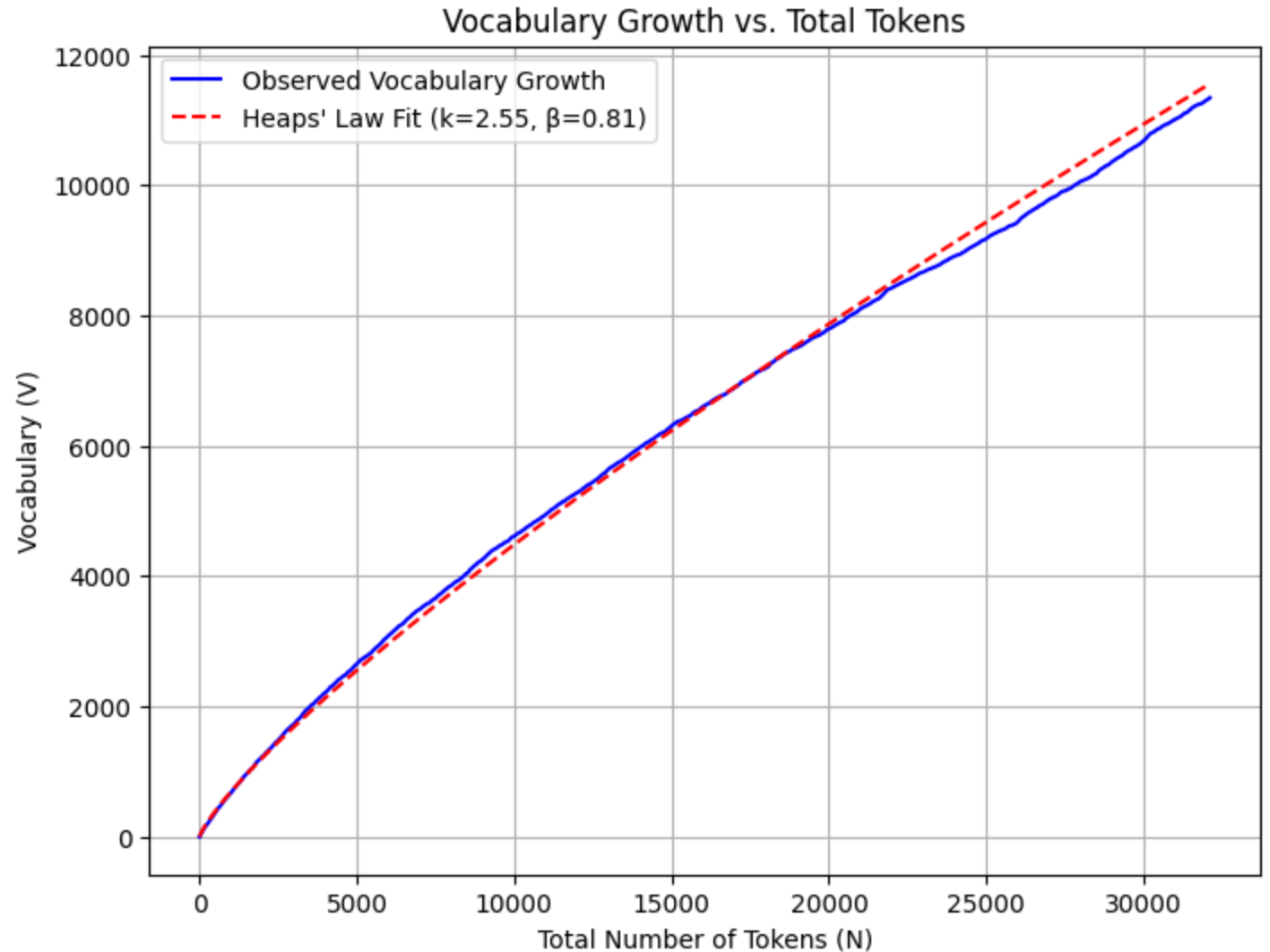
Tokenizer regex

- 32,000 Tokens
- 11,000 Unique tokens



Heaps' Law

- $k = 2.55$
- $\beta = 0.81$

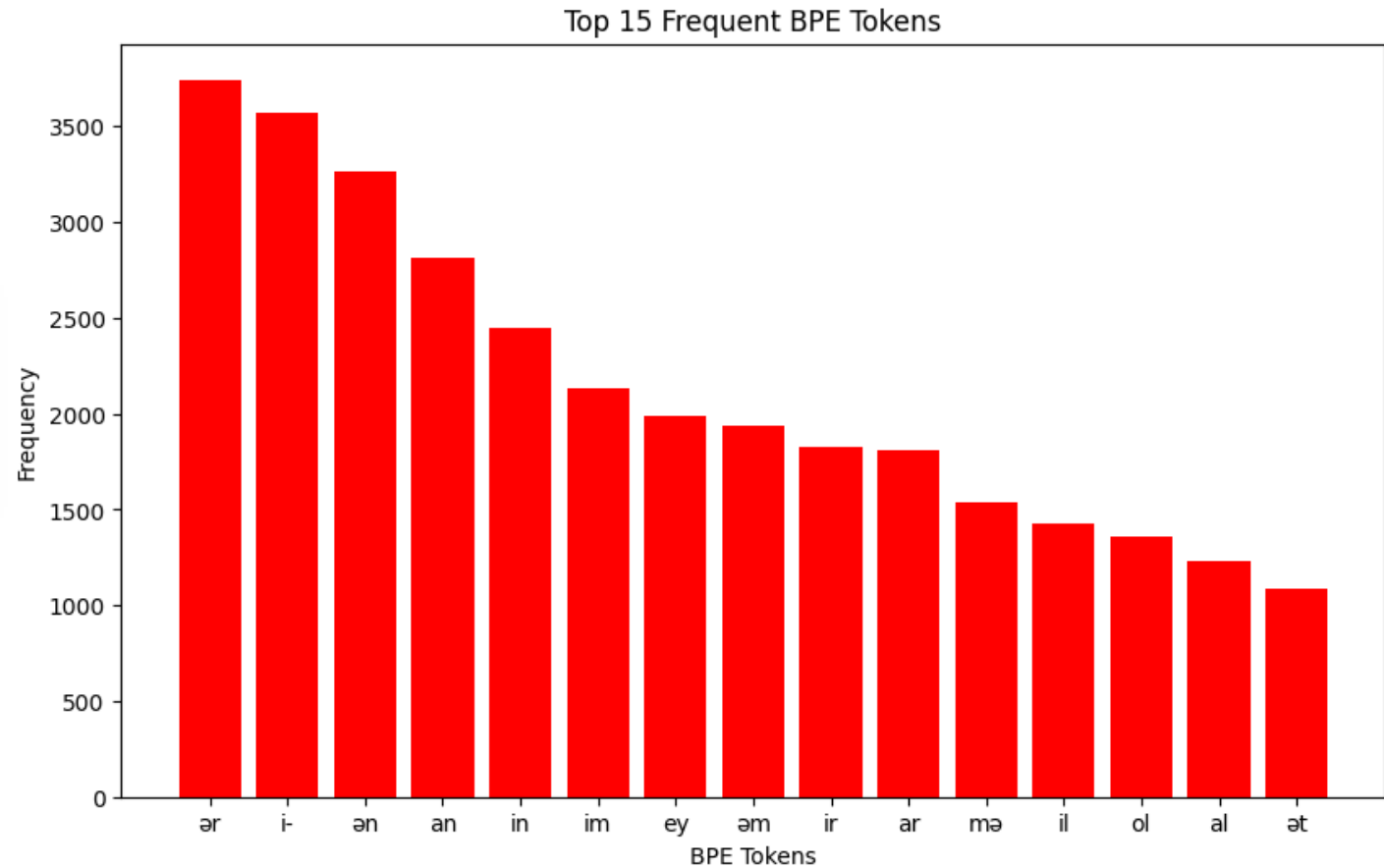


BPE Tokenization

- 1000 merges

```
çox  lə z z ət li  ş or ba dır ! 11 tokens  
t e z - t e z  me n yu mu z a  q ola q  olur . 18 tokens  
bu y ur un ,  siz  de  da dın a  bax ın . 13 tokens
```

Result from a test input



Sentence Segmentation

```
sample_text = "Bakıdan əlimə çatan qadın jurnallarının birində Kələkoş sözü diqqətimi cəlb etdi. " \
"Bu, elə reseptini axtardığım Kələkoş idi! Amma bu dəfı adın yanında resept də var idi! " \
"Beləliklə sirr dolu Kələkoşu bişirmək fürsətini əldə etmiş oldum. Çox ləzzətli şorbadır! " \
"Tez-tez menyumuza qolaq olur. Buyurun, siz de dadına baxın."
```



```
1) ba k ı dan ə l im ə ç at an q ad ın j ur nal lar ın ın bir ind ə k ə l ə k o ş s ö z ü di q q ə t im i c ə l b et di . 35 tokens
2) bu , e l ə r e s e p t ini ax tar d ı ğ ım k ə l ə k o ş idi ! 23 tokens
3) am ma bu də f ı ad ın yan ında r e s e p t də var idi ! 19 tokens
4) b el ə l ik l ə sir r d ol u k ə l ə k o ş u bi ş ir m ə k f ür s ə t ini ə l də et mi ş ol du m . 30 tokens
5) çox l ə z z ə t li ş or ba d ır ! 11 tokens
6) t e z - t e z me n y u mu z a q ol a q ol ur . 18 tokens
7) bu y ur un , siz de da d ın a b ax ın . 13 tokens

Total number of tokens: 149
```

Spelling checker

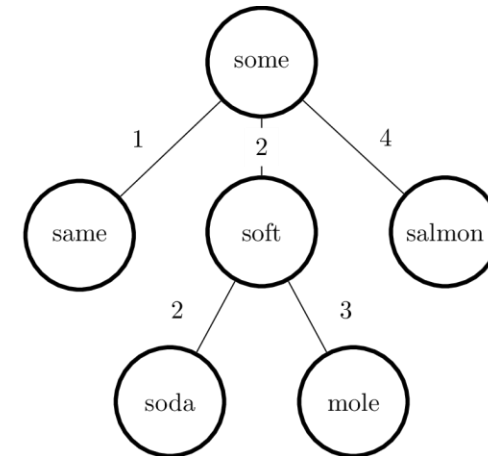
- Levenstein Distance

$$D[i, j] = \begin{cases} 0 & \text{if } i = 0 \text{ and } j = 0 \\ j & \text{if } i = 0 \\ i & \text{if } j = 0 \\ \min \begin{cases} D[i-1, j] + 1 & \text{(deletion)} \\ D[i, j-1] + 1 & \text{(insertion)} \\ D[i-1, j-1] + \text{cost}(s_i, t_j) & \text{(substitution)} \end{cases} & \text{otherwise} \end{cases}$$

where

$$\text{cost}(s_i, t_j) = \begin{cases} 0 & \text{if } s_i = t_j \\ 1 & \text{if } s_i \neq t_j \end{cases}$$

- BK-tree



Example

- Candidates sorted according to Levenshtein distance in ascending order
- Search word: lə'li-nabın
- Results: lə'li-nabın; lə'li-nabin; lə'li-nab
- Search word: ölmaz
- Results: ölməz; olmaz; almaz; olmaq; açmaz; qılmaz; urmaz; olman; olma; alma; qomaz; namaz; bulmaz; salmaz; ölmək; qımaz; qalmaz; almaq; nəmaz

Confusion matrix

- Based on keyboard layout
- Probabilities depend on Euclidean distance

```
keyboard_az = {  
  
    'q': (0, 0), 'ü': (1, 0), 'e': (2, 0), 'r': (3, 0),  
    't': (4, 0), 'y': (5, 0), 'u': (6, 0), 'i': (7, 0),  
    'o': (8, 0), 'p': (9, 0), 'ö': (10, 0), 'ğ': (11, 0),  
    '-': (12, 0),  
  
    'a': (0.5, 1), 's': (1.5, 1), 'd': (2.5, 1), 'f': (3.5, 1),  
    'g': (4.5, 1), 'h': (5.5, 1), 'j': (6.5, 1), 'k': (7.5, 1),  
    'l': (8.5, 1), 'ñ': (9.5, 1), 'ə': (10.5, 1),  
    '"': (11.5, 1),  
  
    'z': (1, 2), 'x': (2, 2), 'c': (3, 2), 'v': (4, 2),  
    'b': (5, 2), 'n': (6, 2), 'm': (7, 2), 'ç': (8, 2),  
    'ş': (9, 2)  
}
```

Weighted Levenstein Distance

$$D[i, j] = \begin{cases} 0 & \text{if } i = 0 \text{ and } j = 0 \\ j & \text{if } i = 0 \\ i & \text{if } j = 0 \\ \min \begin{cases} D[i-1, j] + 1 & \text{(deletion)} \\ D[i, j-1] + 1 & \text{(insertion)} \\ D[i-1, j-1] + w(s_i, t_j) & \text{(substitution)} \end{cases} & \text{otherwise} \end{cases}$$

Example

- Candidates sorted according to Weighted Levenstein distance in ascending order
- Search word: læ'li-nabın
- Results: læ'li-nabın; læ'li-nabin; læ'li-nab
- Search word: ölmaz
- Results: olmaz; olma; ölməz; almaz

