

Name: Sina

Last name: Heidari

Student number: 400464201

Course: Natural Language Processing

## Importing Libraries

```
In [1]: from os import walk
import os
from parsivar import Normalizer
from parsivar import Tokenizer
from parsivar import FindStems
from hazm import Lemmatizer
import re
```

If you don't have hazm and parsivar installed, run the following command inside the working directory:

```
pip install -r requirements.txt
```

## Reading the poem files into a dictionary

```
In [2]: # derives current working folder path from getcwd function
poems_path = os.getcwd() + '/Poems'

# create a dic for storing the poems
poems = {}

# walk through all of the files inside Poems directory
for root, dirs, files in walk(poems_path):
    # iterate through all poem files
    for name in files:
        # join cwd path and file name and open the file
        to_read_file = open(os.path.join(root, name))
        # split file name with '.' and choose first field
        file_key = name.split('.')[0]
        # store in dic
        poems[file_key] = {'txt' : to_read_file.read(), 'tokenized': []}
        to_read_file.close()

# Ex: open file with key: 4844
print(poems['4844']['txt'])
```

هله ای شاه میپچان سر و دستار مرو  
هله ای ماه که نغزت رخ و رخسار مرو  
در همه روی زمین چشم و دل باز که راست  
مکن آزار مکن جانب اغیار مرو  
میر از یار مبر خانه اسرار مسوز  
گل و گلزار مکن جانب هر خار مرو

مکن ای یار ستیزه دغل و جنگ مجوی  
هله آن بار برفتی مکن این بار مرو  
بنده و چاکر و پرورده و مولای توایم  
ای دل و دین و حیات خوش ناچار مرو  
هله سرنای توام مست نواهای توام  
مشکن چنگ طرب را مشکل تار مرو  
هله مخمور چه نالی بر مخمور دگر  
پهلوی خم بنشین از بر خمار مرو  
هله جان بخش بیا ای صدقات تو حیات  
به از این خیر نباشد بجز این کار مرو  
خاتم حسن و جمالی هله ای یوسف دهر  
سوی مکاری اخوان ستمکار مرو  
هله دیدار مهل برگزین فکر و خیال  
از عیان سر مکشان در پی آثار مرو  
هله موسی زمان گرد برآر از دریا  
دل فرعون مجو جانب انکار مرو  
هله عیسی قران صحت رنجور گران  
از برای دو سه ترسا سوی زنا مرو  
هله ای شاهد جان خواجه جانهای شهان  
شیوه کن لب بگز و غنچه افشار مرو  
هله صدیق زمانی به تو ختم است وفا  
جز سوی احمد بگزیده مختار مرو  
جبرئیل کرمی سدره مقام و وطن  
همچو مرغان زمین بر سر شخسار مرو  
تو یقین دار که بی تو نفسی جان نرید  
در احسان بگشا و پس دیوار مرو  
همه رندان و حریفان و بتان جمع شدند  
وقت کار است بیا کار کن از کار مرو  
هله باقی غزل را ز شهنشاه بجوی  
همگی گوش شو اکنون سوی گفتار مرو

## Normalize each poem document

In [3]:

```
normalizer = Normalizer()
for key in poems:
    poems[key] = normalizer.normalize(poems[key]['txt'])
    # replace half spaces and new line characters with single space
    poems[key] = {'txt': re.sub('|\n| ', ' ', poems[key]), 'tokenized': []}
```

In [4]:

```
print(poems['4844'])
```

{'txt': 'هله ای شاه میپچان سر و دستار مرو هله ای ماه که نفزت رخ و رخسار مرو در همه روی زمین چشم و دل باز که راست مکن' آزار مکن جانب اغیار مرو مبر از یار مبر خانه اسرار مسوز گل و گلزار مکن جانب هر خار مرو مکن ای یار ستیزه دغل و جنگ مجوی هله آن بار برفتی مکن این بار مرو بنده و چاکر و پرورده و مولای توایم ای دل و دین و حیات خوش ناچار مرو هله سرنای توام مست نواهای توام مشکن چنگ طرب را مشکل تار مرو هله مخمور چه نالی بر مخمور دگر پهلوی خم بنشین از بر خمار مرو هله جان بخش بیا ای صدقات تو حیات به از این خیر نباشد بجز این کار مرو خاتم حسن و جمالی هله ای یوسف دهر سوی مکاری اخوان ستمکار مرو هله دیدار مهل برمگزین ف کر و خیال از عیان سر مکشان در پی آثار مرو هله موسی زمان گرد برآر از دریا دل فرعون مجو جانب انکار مرو هله عیسی قران صحت رن جور گران از برای دو سه ترسا سوی زنار مرو هله ای شاهد جان خواجه جان های شهان شیوه کن لب بگز و غنچه افشار مرو هله صدیق زما نی به تو ختم است وفا جز سوی احمد بگزیده مختار مرو جبرئیل کرمی سدره مقام و وطنت همچو مرغان زمین بر سر شخسار مرو تو یقین دا ر که بی تو نفسی جان نرید در احسان بگشا و پس دیوار مرو همه رندان و حریفان و بتان جمع شدند وقت کار است بیا کار کن از کار مرو 'tokenized': []}

## Tokenize each poem document

Tokenized poems are stored in 'tokenized' field of every dictionary element (which is a dictionary itself).

In [5]:

```
tokenizer = Tokenizer()
for key in poems:
    poems[key]['tokenized'] = tokenizer.tokenize_words(poems[key]['txt'])
```

In [6]:

```
print(poems['4844']['tokenized'])
```

هله ای شاه میپچان سر و دستار مرو هله ای ماه که نفزت رخ و رخسار مرو در همه روی زمین چشم و دل باز که راست مکن آزار مکن جانب اغیار مرو مبر از یار مبر خانه اسرار مسوز گل و گلزار مکن جانب هر خار مرو مکن ای یار ستیزه دغل و جنگ مجوی هله آن بار برفتی مکن ای ن بار مرو بنده و چاکر و پرورده و مولای توایم ای دل و دین و حیات خوش ناچار مرو هله سرنای توام مست نواهای توام مشکن چنگ طرب را مشکل تار مرو هله مخمور چه نالی بر مخمور دگر پهلوی خم بنشین از بر خمار مرو هله جان بخش بیا ای صدقات تو حیات به از این خیر نباشد بجز این کار مرو خاتم حسن و جمالی هله ای یوسف دهر سوی مکاری اخوان ستمکار مرو هله دیدار مهل برمگزین ف آثار مرو هله موسی زمان گرد برآر از دریا دل فرعون مجو جانب انکار مرو هله عیسی قران صحت رنجور گران از برای دو سه ترسا سوی زنار مرو هله ای شاهد جان خواجه جان های شهان شیوه کن لب بگز و غنچه افشار مرو هله صدیق زما نی به تو ختم است وفا جز سوی احمد بگزیده مختار مرو جبرئیل کرمی سدره مقام و وطنت همچو مرغان زمین بر سر شخسار مرو تو یقین دا

ن', 'دار', 'که', 'بی', 'تو', 'نفسی', 'جان', 'نزید', 'در', 'احسان', 'بگشا', 'و', 'پس', 'دیوار', 'مرو', 'همه',  
 'رندان', 'و', 'حریفان', 'و', 'بتان', 'جمع', 'شدند', 'وقت', 'کار', 'است', 'بیا', 'کار', 'کن', 'از', 'کار', 'مر  
 [و', 'هله', 'باقی', 'غزل', 'را', 'ز', 'شهنشاه', 'بجوی', 'همگی', 'گوش', 'شو', 'اکنون', 'سوی', 'گفتار', 'مرو

## Removing stop-words

In [7]:

```
# read stop-words file
stop_file = open(os.getcwd() + '/Stopwords/Stopwords.txt', mode='r')
stop_file = stop_file.read()
# split them into a list of words
stop_file_split = stop_file.split('\n')
to_read_file.close()
tokenizer = Tokenizer()
# for every key inside poem dictionary
for key in poems:
    # a list comprehension to only choose words that are not included in our stop-words file (ig
    poems[key]['tokenized'] = [token for token in poems[key]['tokenized'] if token not in stop_f
```

In [8]:

```
print(poems['4844']['tokenized'])
```

هله', 'شاه', 'مپیچان', 'دستار', 'مرو', 'هله', 'ماه', 'نفزت', 'رخ', 'رخسار', 'مرو', 'روی', 'زمین', 'چشم', 'د'  
 'باز', 'راست', 'مکن', 'آزار', 'مکن', 'جانب', 'اغیار', 'مرو', 'میر', 'یار', 'میر', 'خانه', 'اسرار', 'مسوز',  
 'گل', 'گلزار', 'مکن', 'جانب', 'خار', 'مرو', 'مکن', 'یار', 'ستیزه', 'دغل', 'جنگ', 'مجوی', 'هله', 'بار', 'برفت  
 ی', 'مکن', 'بار', 'مرو', 'بنده', 'چاکر', 'پرورده', 'مولای', 'توایم', 'دل', 'دین', 'حیات', 'خوش', 'ناچار', 'مرو',  
 'هله', 'سرنای', 'توام', 'مست', 'نواهای', 'توام', 'مشکن', 'چنگ', 'طرب', 'مسکل', 'تار', 'مرو', 'هله', 'مخمور',  
 'نالی', 'مخمور', 'دگر', 'پهلوی', 'خم', 'بنشین', 'خمار', 'مرو', 'هله', 'جان', 'بخش', 'بیا', 'صدقات', 'حیات', 'خی  
 ر', 'نباشد', 'بجز', 'کار', 'مرو', 'خاتم', 'حسن', 'جمالی', 'هله', 'یوسف', 'دهر', 'مکاری', 'اخوان', 'ستمکار', 'مر  
 و', 'هله', 'دیدار', 'مهل', 'برمگزین', 'فکر', 'خیال', 'عیان', 'مکشان', 'بی', 'آثار', 'مرو', 'هله', 'موسی', 'زمان',  
 'گرد', 'برآر', 'دریا', 'دل', 'فرعون', 'مجو', 'جانب', 'انکار', 'مرو', 'هله', 'عیسی', 'قران', 'صحت', 'رنجور',  
 'گران', 'برای', 'دو', 'سه', 'ترسا', 'زنار', 'مرو', 'هله', 'شاهد', 'جان', 'خواجه', 'جان', 'های', 'شهان', 'شیو  
 ه', 'لب', 'بگر', 'غنبه', 'افشار', 'مرو', 'هله', 'صدیق', 'زمانی', 'ختم', 'وفا', 'جز', 'احمد', 'بگریده', 'مختار',  
 'مرو', 'جبرئیل', 'کرمی', 'سدره', 'مقام', 'وطن', 'همچو', 'مرغان', 'زمین', 'شخسار', 'مرو', 'یقین', 'دار', 'بی',  
 'نفسی', 'جان', 'نزید', 'احسان', 'بگشا', 'دیوار', 'مرو', 'رندان', 'حریفان', 'بتان', 'جمع', 'شدند', 'وقت', 'کار', 'ب  
 یا', 'کار', 'کار', 'مرو', 'هله', 'باقی', 'غزل', 'شهنشاه', 'بجوی', 'همگی', 'گوش', 'شو', 'اکنون', 'گفتار', 'مر  
 و']

## Lemmatizing tokenized words

In [9]:

```
lemmatizer = Lemmatizer()
# A list comprehension to lemmetize each word inside the tokenized list of words corresponding t
example_lemmetized = [lemmatizer.lemmatize(token) for token in poems['4844']['tokenized']]
```

```
In [10]: print(example_lemmetized)
```

هله, 'شاه', 'مپیجان', 'دستار', 'مرو', 'هله', 'ماه', 'نغز', 'رخ', 'رخسار', 'مرو', 'روی', 'زمین', 'چشم', 'د', 'ل', 'باز', 'راست', 'مکن', 'آزار', 'مکن', 'جانب', 'اغیار', 'مرو', 'میر', 'یار', 'میر', 'خانه', 'اسرار', 'مسوز', 'گل', 'گلزار', 'مکن', 'جانب', 'خار', 'مرو', 'مکن', 'یار', 'ستیزه', 'دغل', 'جنگ', 'مجوی', 'هله', 'بار', 'برفت', 'ی', 'مکن', 'بار', 'مرو', 'بنده', 'چاکر', 'پرورده', 'مولا', 'توایم', 'دل', 'دین', 'حیات', 'خوش', 'ناچار', 'مرو', 'هله', 'سرنا', 'توایم', 'مست', 'نوا', 'توایم', 'مشکن', 'جنگ', 'طرب', 'مسکل', 'تار', 'مرو', 'هله', 'مخمور', 'نال', 'مخمور', 'دگر', 'پهلوی', 'خم', 'نشست#نشین', 'خمار', 'مرو', 'هله', 'جان', 'بخش', 'آمد#آ', 'صدقات', 'حیات', 'خ', 'یر', 'بود#باش', 'بجز', 'کار', 'مرو', 'خاتم', 'حسن', 'جمال', 'هله', 'یوسف', 'دهر', 'مکار', 'اخوان', 'ستمکار', 'م', 'رو', 'هله', 'دیدار', 'مهل', 'برمگزین', 'فکر', 'خیال', 'عیان', 'مکشان', 'پی', 'آثار', 'مرو', 'هله', 'موسی', 'زما', 'ن', 'گرد', 'برآر', 'دریا', 'دل', 'فرعون', 'مجو', 'جانب', 'انکار', 'مرو', 'هله', 'عیسی', 'قران', 'صحت', 'رنجو', 'ر', 'گران', 'برای', 'دو', 'سه', 'ترسا', 'زنار', 'مرو', 'هله', 'شاهد', 'جان', 'خواجه', 'جان', 'های', 'شه', 'شی', 'وه', 'لب', 'بگز', 'غیبه', 'افشار', 'مرو', 'هله', 'صدیق', 'زمان', 'ختم', 'وفا', 'جز', 'احمد', 'بگزیده', 'مختار', 'مرو', 'جبرئیل', 'کرمی', 'سدره', 'مقام', 'وطن', 'همجو', 'مرغ', 'زمین', 'شخسار', 'مرو', 'یقین', 'دار', 'پی', 'ن', 'فس', 'جان', 'نزید', 'احسان', 'گشود#گشا', 'دیوار', 'مرو', 'رند', 'حریف', 'بتان', 'جمع', 'شد#شو', 'وقت', 'کار', 'آمد#آ', 'کار', 'کار', 'مرو', 'هله', 'باقی', 'غزل', 'شهنشاه', 'جست#جو', 'همگی', 'گوش', 'شو', 'اکنون', 'گفتا', 'ر', 'مرو']

## Stemming tokenized words

```
In [11]: my_stemmer = FindStems()
# A list comprehension to find stems of each word inside the tokenized list of words correspondi
example_stems = [my_stemmer.convert_to_stem(token) for token in poems['4844']['tokenized']]
```

```
In [12]: print(example_stems)
```

هله, 'شاه', 'مپیجان', 'دستار', 'مرو', 'هله', 'ماه', 'نغز', 'رخ', 'رخسار', 'مرو', 'روی', 'زمین', 'چشم', 'د', 'ل', 'باز', 'راست', 'مکن', 'آزار', 'مکن', 'جانب', 'اغیار', 'مرو', 'میر', 'یار', 'میر', 'خانه', 'اسرار', 'مسوز', 'گل', 'گلزار', 'مکن', 'جانب', 'خار', 'مرو', 'مکن', 'یار', 'ستیزه', 'دغل', 'جنگ', 'مجوی', 'هله', 'بار', 'برفت', 'ی', 'مکن', 'بار', 'مرو', 'بنده', 'چاکر', 'پرورده', 'مولا', 'توایم', 'دل', 'دین', 'حیات', 'خوش', 'ناچار', 'مرو', 'هله', 'سرنا', 'توایم', 'مست', 'نوا', 'توایم', 'مشکن', 'جنگ', 'طرب', 'مسکل', 'تار', 'مرو', 'هله', 'مخمور', 'نا', 'ل', 'مخمور', 'دگر', 'پهلوی', 'خم', 'نشست&نشین', 'خمار', 'مرو', 'هله', 'جان', 'بخش', 'بیا', 'صدقات', 'حیات', 'خیر', 'بود&باش', 'بجز', 'کار', 'مرو', 'خاتم', 'حسن', 'جمال', 'هله', 'یوسف', 'دهر', 'مکاری', 'اخوان', 'ستمکا', 'ر', 'مرو', 'هله', 'دیدار', 'مهل', 'برمگزین', 'فکر', 'خیال', 'عیان', 'مک', 'پی', 'اثر', 'مرو', 'هله', 'موسی', 'زمان', 'گرد', 'برآر', 'دریا', 'دل', 'فرعون', 'مجو', 'جانب', 'انکار', 'مرو', 'هله', 'عیسی', 'قران', 'صحت', 'رن', 'جور', 'گران', 'برای', 'دو', 'سه', 'ترسا', 'زنار', 'مرو', 'هله', 'شاهد', 'جان', 'خواجه', 'جان', 'های', 'شهان', 'شیوه', 'لب', 'گزید&گز', 'غیبه', 'افشار', 'مرو', 'هله', 'صدیق', 'زمانی', 'ختم', 'وفا', 'جز', 'احمد', 'بگزیده', 'مختار', 'مرو', 'جبرئیل', 'کرمی', 'سدره', 'مقام', 'وطن', 'همجو', 'مرغ', 'زمین', 'شخسار', 'مرو', 'یقین', 'دار',

'بی', 'نفس', 'جان', 'نزدید', 'احسان', 'گشود&گشا', 'دیوار', 'مرو', 'رند', 'حریف', 'بت', 'جمع', 'شد&شد', 'وقت',  
'کار', 'ییا', 'کار', 'کار', 'مرو', 'هله', 'باقی', 'غزل', 'شهنشاه', 'جست&جو', 'همگی', 'گوش', 'شو', 'اکنون', 'گ  
[فتار', 'مرو

