### pre\_process

September 18, 2024

## 1 Preprocess Required for Generating Train Data:

```
[1]: import urllib.request
                    import time
                    import sys
                    import getopt
                    import pandas as pd
                    import numpy as np
                    import pickle
[11]: %run ../utils.ipynb
   [2]: embSize = 200
                    ftrain='../data/EUADR_target_disease.csv'
                    # Replace with path of word embdding file
                    \#wefile = \#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin\#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA\_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/Admin/GDA_backup/Dataset/embeddings/PubMed-and-PMC-w2v.bin#/mnt/Admin/GDA_backup/Dataset/embeddings/PubMed-and-pmc-admin/GDA_backup/Dataset/embeddings/PubMed-and-pmc-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Dataset/embeddings/PubMed-admin/GDA_backup/Datas
                    wefile = "../support/PubMed-and-PMC-w2v.bin"
                    random_seed=1331
[11]: import pandas as pd
                    from tabulate import tabulate
                    ftrain = '../data/EUADR_target_disease.csv'
                    with open(ftrain, 'r', encoding='latin1') as file:
                                 first_line = file.readline()
                                 print(first_line)
                    import pandas as pd
                    ftrain = '../data/EUADR_target_disease.csv'
                    # Specify the delimiter for tabs
                    df = pd.read_csv(ftrain, encoding='latin1', sep='\t')
                    # Display the first 10 rows of the DataFrame as a table
                    print(tabulate(df.head(4), headers='keys', tablefmt='grid'))
```

"ASSOCIATION\_TYPE" "PMID" "NUM\_SENTENCE" "ENTITY1\_TEXT" "ENTITY1\_INI"

| "ENTITY2_TYPE   | "ENTITY1_TYPE" "ENTITY2_TEXT" "ENTITY2_INI" "ENTITY2_END" "SENTENCE"  |
|---|---|
| +   |   |
| ASSOCI<br>ENTITY1_INI  <br>ENTITY2_INI  <br> <br>+===+=====   | ATION_TYPE   PMID   NUM_SENTENCE   ENTITY1_TEXT   ENTITY1_END   ENTITY1_TYPE   ENTITY2_TEXT   ENTITY2_END   ENTITY2_TYPE   SENTENCE   |
|   | 17241106   16   LRP5   23   Genes & Molecular Sequences   osteoporosis   92   Diseases & Disorders   Our work supported LRP5 genetic cossible susceptibility factors for osteoporosis and fractures in              |
| +   | +   |
| 1   SA<br>20  <br>108  <br>its strongly<br>osteoporosis<br>++ | 28   SNP & Sequence variations   osteoporosis   120   Diseases & Disorders   Especially, the SNP rs491347 and associated SNPs (e.g., rs1784235) could be important to human   |
| 2   SA<br>69  <br>108  <br>its strongly<br>osteoporosis<br>++ | 17241106   17   rs1784235   78   SNP & Sequence variations   osteoporosis   120   Diseases & Disorders   Especially, the SNP rs491347 and associated SNPs (e.g., rs1784235) could be important to human phenotypes. |
|   | 18697826   0   fetal haemoglobin  |

```
110 | Genes & Molecular Sequences | beta-thalassaemia |
                 149 | Diseases & Disorders | The HBS1L-MYB intergenic region on
    132 l
    chromosome 6q23 is a quantitative trait locus controlling fetal haemoglobin
    level in carriers of beta-thalassaemia.
    ______
    ______
[19]: import pandas as pd
    from tabulate import tabulate
    # Load the CSV file
    ftrain = '../data/EUADR_target_disease.csv'
    df = pd.read_csv(ftrain, encoding='latin1', sep='\t')
     # Display the first line to understand the structure
    with open(ftrain, 'r', encoding='latin1') as file:
        first_line = file.readline()
        print("First line of the file:", first_line)
     # Display the column names to identify the one that corresponds to the entity_
     ⇔association type
    print("Column names:", df.columns)
    # Check for unique association types
    unique_association_types = df['ASSOCIATION_TYPE'].unique()
    print("Unique association types:", unique_association_types)
     # Filter the DataFrame to exclude rows where the association type is NaN (if \Box
     \rightarrowapplicable)
    df_filtered = df[df['ASSOCIATION_TYPE'].notna()]
    # Display the filtered rows
    print(tabulate(df_filtered.head(15), headers='keys', tablefmt='grid'))
                                            "PMID" "NUM_SENTENCE"
    First line of the file: "ASSOCIATION_TYPE"
    "ENTITY1_TEXT"
                 "ENTITY1 INI"
                              "ENTITY1_END"
                                            "ENTITY1_TYPE" "ENTITY2_TEXT"
    "ENTITY2 INI"
                 "ENTITY2 END"
                              "ENTITY2_TYPE" "SENTENCE"
    Column names: Index(['ASSOCIATION_TYPE', 'PMID', 'NUM_SENTENCE', 'ENTITY1_TEXT',
          'ENTITY1_INI', 'ENTITY1_END', 'ENTITY1_TYPE', 'ENTITY2_TEXT',
          'ENTITY2_INI', 'ENTITY2_END', 'ENTITY2_TYPE', 'SENTENCE'],
         dtype='object')
    Unique association types: ['SA' 'FA' 'PA' nan]
    ______
```

93 I

| +   |
|---|
|   |
| '<br>   |
|   |
|   |
|   |
|   |
| +   ASSOCIATION_TYPE   PMID   NUM_SENTENCE   ENTITY1_TEXT   ENTITY1_INI   ENTITY1_END   ENTITY1_TYPE   ENTITY2_TEXT   ENTITY2_INI   ENTITY2_END   ENTITY2_TYPE   SENTENCE |
| +===+======+=====+=====+=====+=========   |
| =====+====+====++=====++=====++====++====   |
|   |
| =======+=====+=====+=====+=====+=====+====  |
| =====+=================================   |
|   |
|   |
| =======================================   |
|   |
| ======================================  |
| 80   92   Diseases & Disorders   Our work supported LRP   |
| genetic variants as possible susceptibility factors for osteoporosis and  |
|   |
| fractures in humans.  |
|   |
| ++  |
|   |
|   |
| +   |
|   |
|   |
|   |
|   |
|   |
|   |
| <del>-</del>  |
| 1   SA  |
| 20   28   SNP & Sequence variations   osteoporosis  |
| 108   120   Diseases & Disorders   Especially, the SNP  |
| rs491347 and its strongly associated SNPs (e.g., rs1784235) could be important  |
|   |
| to human osteoporosis phenotypes.   |
|   |
| ++  |
| ++++  |
|   |
|   |
| <del>-</del>  |
|   |
|   |
|   |

|                                       |                | +                       |                                 |
|---------------------------------------|----------------|-------------------------|---------------------------------|
| 2   SA                                | I              | 17241106                | 17   rs1784235                  |
| 1                                     | 69             | 78   SNP & Sequence v   | variations   osteoporosis       |
| 1 1                                   | .08 I          | <del>-</del>            | rders   Especially, the SNP     |
| rs491347 and                          | its strongly a |                         | rs1784235) could be important   |
|                                       | oporosis pheno |                         | isiro4250) could be important   |
| i i i i i i i i i i i i i i i i i i i | oporosis pheno | ctypes.                 |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         | +                               |
|                                       |                |                         |                                 |
|                                       |                | ·                       |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                | +                       |                                 |
| 3   FA                                | I              | 18697826                | 0   fetal haemoglobin           |
| I                                     | 93             |                         | _                               |
| +1-7                                  | 93             | 110   Genes & Molecula  | <del>-</del>                    |
| thalassaemia                          |                | 132                     |                                 |
|                                       |                |                         | romosome 6q23 is a quantitative |
|                                       | _              | al haemoglobin level in | n carriers of beta-             |
| thalassaemia.                         |                |                         |                                 |
|                                       |                |                         |                                 |
| +                                     | +              |                         | +                               |
|                                       | +              | +                       | +                               |
|                                       |                | +                       |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                | +                       |                                 |
| 4   PA                                |                | 18697826                | 1   HbF                         |
|                                       | 19             | 22   Genes & Molecula   | ar Sequences   HBB disorders    |
| 1                                     | 64             | 77   Diseases & Dison   | rders   Fetal haemoglobin (HbF) |
| level modifie                         | s the clinical | severity of HBB disord  | ders.                           |
| 1                                     |                | ·                       |                                 |
| ·<br>+                                | +              |                         | +                               |
|                                       |                |                         | ·<br>+                          |
|                                       |                |                         | ·<br>                           |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                |                         |                                 |
|                                       |                | +                       |                                 |
| 5   PA                                | I              | 18697826                | 1   Fetal haemoglobin           |

| 1  |   | 77   Diseases & Disorder severity of HBB disorder      |  |
|--|---|--|--|
|  |   |  | +  |
|  |   |  |  |
| 6   FA<br> <br>thalassaemia<br>Disorders | 4  <br>The HBS1L-MYB i<br>controlling fet | 132  | Sequences   beta-<br>149   Diseases &<br>mosome 6q23 is a quantitative |
|  |   |  | +  |
|  |   |  |  |
|  |   |  |  |
|  |   |  |  |
| 7   FA<br> <br>thalassaemia              | <br>  46  <br>                            | ·  | 149   Diseases &   |
|  | controlling fet                           | ntergenic region on chron<br>al haemoglobin level in o | nosome 6q23 is a quantitative<br>carriers of beta-                     |
|  | +   |  | +  |
| +  |   |  |  |
|  |   |  |  |
| 8   FA<br>                               |   |  | 7   HbF<br>Sequences   HBB disorders<br>ers   Functional studies to    |

| unravel the biological significance of this region in regulating HbF production is clearly indicated, which may lead to new strategies to modify the disease course of severe HBB disorders. |
|--|
| ++   |
| ++++   |
| ·  |
|  |
|  |
|  |
| 9   PA   |
| ++   |
|  |
|  |
|  |
|  |
|  |
| 10   FA  |
|  |
| 11   PA  |

| (p=0.028), IL-1R C TNFalpha G -308 all significantly in th polymorphic variati role in susceptibil + | pst1 1970 allele (pele (p=0.0002) and e patients versus rons of these pro-ir ity of Iranian mult | quency of IL-1a p=0.0001) and C GG genotype (p normal subjects iflammatory cyt ciple sclerosis | 19ha TT -889 genotype C genotype (p=0.00006) =0.000001) decreased These results suggest okines may play an importants. | t that portant |
|--|--|--|--|----------------|
| 12   PA<br>  0  <br>sclerosis<br>Disorders   IL-1, I<br>multiple sclerosis.<br> <br>++               | 18322311  <br>  4   Gene<br> <br>  | es & Molecular 69   gene polymorphi  | 0   IL-1  Sequences   multiple  87   Diseases & sms in Iranian patient   | kts with       |
|  |  |  |  |                |
| 13   PA the GNAS1   oropharyngeal Diseases & Disorder hypopharyngeal squa polymorphism of the   ++   | 18347176   127   s   Overall and rel mous cell carcinoma GNAS1 gene.                             | 158   SNP & Se<br> <br>  | 0   T393C polymorphism quence variations   37   50   ival in oropharyngeal d with genotypes of T3                      | and<br>393C    |
|  |  | +<br>+   | <br>0   T393C polymorphism   | <br>n of       |

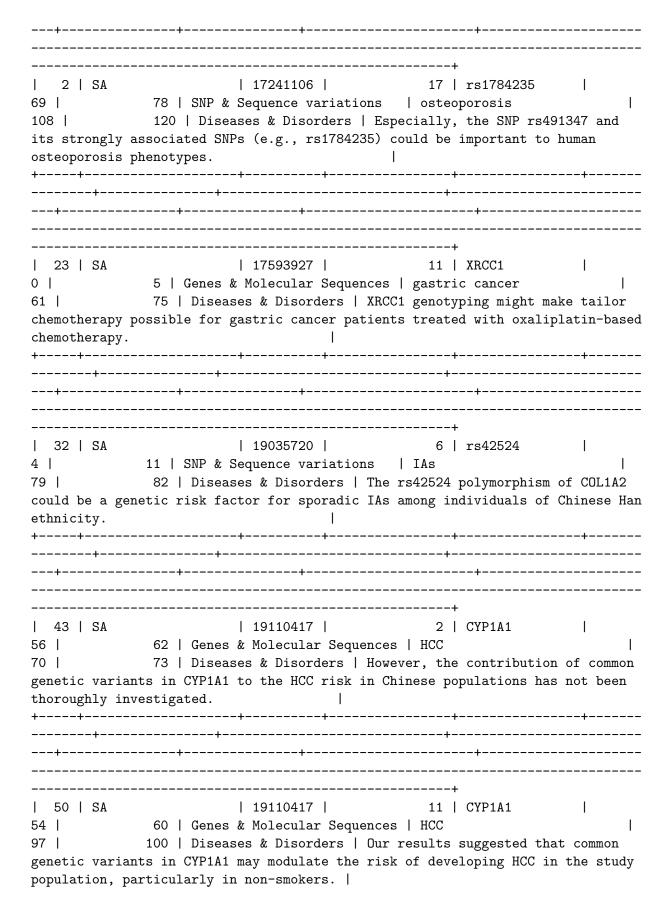
```
the GNAS1 |
                       127 l
                                    158 | SNP & Sequence variations
    hypopharyngeal squamous cell carcinoma |
                                                55 I
                                                             93 I
    Diseases & Disorders | Overall and relapse-free survival in oropharyngeal and
    hypopharyngeal squamous cell carcinoma are associated with genotypes of T393C
    polymorphism of the GNAS1 gene.
                -----
    ______
            ______
    _____+___
    ______
[16]: import pandas as pd
     from tabulate import tabulate
     # Load the CSV file
     ftrain = '../data/EUADR_target_disease.csv'
     df = pd.read_csv(ftrain, encoding='latin1', sep='\t')
     # Display the first line to understand the structure
     with open(ftrain, 'r', encoding='latin1') as file:
        first line = file.readline()
        print("First line of the file:", first_line)
     # Display the column names to identify the one that corresponds to the entity_
     →association type
     print("Column names:", df.columns)
     # Get unique association types and their counts
     association_type_counts = df['ASSOCIATION_TYPE'].value_counts()
     print("Association types and their counts:")
     print(tabulate(association_type_counts.reset_index(), headers=['Association_u
      →Type', 'Count'], tablefmt='grid'))
     # If you want to see a sample of rows for each association type
     # Display the first 10 rows for each unique association type
     for assoc_type in df['ASSOCIATION_TYPE'].unique():
        print(f"\nSample rows for association type '{assoc_type}':")
        sample df = df[df['ASSOCIATION TYPE'] == assoc type].head(10)
        print(tabulate(sample_df, headers='keys', tablefmt='grid'))
```

First line of the file: "ASSOCIATION\_TYPE" "PMID" "NUM\_SENTENCE"

"ENTITY1\_TEXT" "ENTITY1\_INI" "ENTITY1\_END" "ENTITY1\_TYPE" "ENTITY2\_TEXT"

"ENTITY2\_INI" "ENTITY2\_END" "ENTITY2\_TYPE" "SENTENCE"

```
Column names: Index(['ASSOCIATION_TYPE', 'PMID', 'NUM_SENTENCE', 'ENTITY1_TEXT',
   'ENTITY1_INI', 'ENTITY1_END', 'ENTITY1_TYPE', 'ENTITY2_TEXT',
   'ENTITY2_INI', 'ENTITY2_END', 'ENTITY2_TYPE', 'SENTENCE'],
   dtype='object')
Association types and their counts:
+---+
  | Association Type | Count |
+===+=======+
I O I PA
+---+
1 | FA
+---+
Sample rows for association type 'SA':
____+___
___+____
______
    ______
   | ASSOCIATION TYPE |
                PMID |
                     NUM SENTENCE | ENTITY1 TEXT
ENTITY1_INI | ENTITY1_END | ENTITY1_TYPE
                             | ENTITY2_TEXT
 ENTITY2_INI | ENTITY2_END | ENTITY2_TYPE
                           | SENTENCE
_____+__+___+____+___+___+____
______
_____+
             | 17241106 |
                          16 | LRP5
        23 | Genes & Molecular Sequences | osteoporosis
19 l
        92 | Diseases & Disorders | Our work supported LRP5 genetic
variants as possible susceptibility factors for osteoporosis and fractures in
humans.
-------
| 1 | SA
             | 17241106 |
                          17 | rs491347
20 |
        28 | SNP & Sequence variations | osteoporosis
        120 | Diseases & Disorders | Especially, the SNP rs491347 and
108
its strongly associated SNPs (e.g., rs1784235) could be important to human
osteoporosis phenotypes.
.____+___
```



| +              |   | +                 | -+       |
|----------------|---|-------------------|----------|
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   | +                 |          |
| 86   SA        | 18708184   1                                | NRG3              | 1        |
| 51 l           | 55   Genes & Molecular Sequences   schi     | zonhrenia         | · 1      |
| •              |   | -                 |          |
| 65             | 78   Diseases & Disorders   The study in    | -                 | possible |
| association of | f NRG3 gene and schizophrenia in a Han Chi  | nese population.  |          |
|                |   |                   |          |
| +              |   | +                 | -+       |
|                |   |                   |          |
|                |   |                   |          |
| ·              |   | •                 |          |
|                |   |                   |          |
|                |   | +                 |          |
| 98   SA        | 19098911   6                                | NLRP3             | 1        |
| 31             | 36   Genes & Molecular Sequences   Croh     | n's disease       | 1        |
| 130 l          | 145   Diseases & Disorders   These result   |                   | the .    |
| •              |   |                   | one      |
| _              | is also implicated in the susceptibility of | r more common     |          |
|                | diseases such as Crohn's disease.           |                   |          |
| +              |   | +                 | -+       |
|                | +   |                   |          |
| +              | +   | +                 |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
| 102   SA       |   | MDR1              | ı        |
| 32             | 36   Genes & Molecular Sequences   infla    | ammatory bowel d: | isease   |
| 93             | 119   Diseases & Disorders   Therefore,     | the mutations of  | the MDR1 |
| gene are thou  | ght to be related with the pathogenesis of  | inflammatory box  | wel      |
| disease.       | <br>  |                   |          |
|                | ·<br>                                       |                   |          |
|                |   |                   |          |
|                | +   |                   |          |
| +              | +   | +                 |          |
|                |   |                   |          |
|                |   | +                 |          |
|                |   |                   |          |
| C              | i-tion town ITAL.                           |                   |          |
| -              | or association type 'FA':                   |                   |          |
|                | +   |                   |          |
|                | -+  |                   |          |
| +              |   | +                 |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
|                |   |                   |          |
|                | +   |                   |          |
| ASSOCIA        | ATION_TYPE   PMID   NUM_SENTENCE            | ENTITY1 TEXT      | 1        |
|                | ENTITY1_END   ENTITY1_TYPE                  | ENTITY2_TEXT      | Г        |
|                |   | , -MIIIIA_IDA.    | •        |

| 1  | ENTITY2_END   ENTITY2_TYPE   SENTENCE  |   |
|--|--|---|
| =======================================  | :======+==+====+===++====++====++====++====  | = |
|  |  |   |
|  |  |   |
|  |  |   |
| _  | 18697826   0   fetal haemoglobin   110   Genes & Molecular Sequences   beta-thalassaemia   149   Diseases & Disorders   The HBS1L-MYB intergenic region on is a quantitative trait locus controlling fetal haemoglobin ers of beta-thalassaemia. |   |
|  | ·  |   |
| +  | +  | _ |
|  |  |   |
|  |  |   |
| 6   FA<br>4  <br>132  <br>chromosome 6q2 |  |   |
| _  | ers of beta-thalassaemia.  |   |
| ++                                       |  | _ |
|  | +  |   |
|  |  |   |
|  |  | - |
| _  | 18697826   0   6q23   50   Genes & Molecular Sequences   beta-thalassaemia   149   Diseases & Disorders   The HBS1L-MYB intergenic region on is a quantitative trait locus controlling fetal haemoglobiners of beta-thalassaemia.                |   |
| ++                                       | +++++  | - |

| <del></del>  |
|--|
|  |
|  |
|  |
|  |
|  |
|  |
| +  |
|  |
| · · · · · · · · · · · · · · · · · · ·  |
| 87   90   Genes & Molecular Sequences   HBB disorders                            |
| 196   209   Diseases & Disorders   Functional studies to unravel the             |
| biological significance of this region in regulating HbF production is clearly   |
| indicated, which may lead to new strategies to modify the disease course of      |
| severe HBB disorders.  |
|  |
| ı<br>++  |
|  |
|  |
| +++++++  |
|  |
|  |
|  |
|  |
|  |
|  |
| +  |
| 10   FA  |
| 35   52   SNP & Sequence variations   multiple sclerosis                         |
| 406   424   Diseases & Disorders   On the other hand the frequency of            |
| IL-1alpha TT -889 genotype (p=0.028), IL-1R C pst1 1970 allele (p=0.0001) and CC |
|  |
| genotype (p=0.00006), TNFalpha G -308 allele (p=0.0002) and GG genotype          |
| (p=0.000001) decreased significantly in the patients versus normal               |
| subjects. These results suggest that polymorphic variations of these pro-        |
| inflammatory cytokines may play an important role in susceptibility of Iranian   |
| multiple sclerosis patients.   |
|  |
|  |
|  |
| +++++  |
|  |
|  |
|  |
|  |
|  |
| ·  |
|  |
| 15   FA  |
| 4   9   SNP & Sequence variations   hypopharyngeal cancer                        |
| 130   151   Diseases & Disorders   The T393C SNP could be considered             |
| as a genetic marker to predict the clinical course of patients suffering from    |
| oropharyngeal and hypopharyngeal cancer.   |
| <br>   |
|  |

| ++   |
|--|
|  |
| +  |
|  |
|  |
|  |
|  |
|  |
|  |
| +  |
| 17   FA   18347176   3   T393C   |
|  |
| 8   33   SNP & Sequence variations   oropharyngeal   |
| 21   134   Diseases & Disorders   The prognostic value of the T393C  |
| NP was evaluated in an unselected series of patients treated with curative   |
| ntent for oropharyngeal and hypopharyngeal squamous cell carcinomas, including   |
| ll tumor stages with different therapeutic regimens.   |
|  |
| +  |
|  |
|  |
| ++++++   |
|  |
|  |
|  |
|  |
|  |
|  |
|  |
|  |
| 18   FA   17430902   3   Wnt   |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   |
| 18   FA   17430902   3   Wnt   |
| 18   FA         17430902         3   Wnt                 9         72   Genes & Molecular Sequences   colon cancer                 5         37   Diseases & Disorders   However, the majority of colon            |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA   17430902   3   Wnt   9   72   Genes & Molecular Sequences   colon cancer   5   37   Diseases & Disorders   However, the majority of colon ancer cells have deregulation of the Wnt/beta-catenin pathway. |
| 18   FA  |

|   | +  |   |
|---|--|---|
| 28   FA<br>5  <br>89  <br>etween the CYP      | 19097922   1   CYP4F2<br>51   Genes & Molecular Sequences   MI<br>141   Diseases & Disorders   This study assessed associated associated associated (MI), using a haplotic control of the co |   |
|   | tudy of 234 MI patients and 248 controls genotyped for symorphisms (rs3093105, rs3093135, rs1558139, rs2108622,  | 5 single                                  |
|   |  |   |
|   | +  |   |
|   | +  |   |
|   |  |   |
|   |  |   |
|   |  |   |
|   |  |   |
|   |  |   |
|   | +  |   |
| <del>-</del>                                  | +++++  | -+  |
|   |  |   |
|   |  |   |
|   |  |   |
|   |  |   |
| <br>  ASSOCIAT                                | +  | Y2 TEXT                                   |
|   | I   ENTITY2_END   ENTITY2_TYPE   SENTENCE  | _   |
|   |  |   |
| ===+=====                                     |  |   |
| ===+======                                    |  | =+=====                                   |
| ===+===================================       | +  | =+=====                                   |
| ===+======<br>=====+======<br>====+========   | ======+====+====+=====+=====+=====+=====   | =+=====<br>========<br>=======            |
| ===+======<br>=====+======<br>====+========== | +++++  | =+=====<br>========<br>=======            |
| +   |  | =+=====<br>============================== |
|   | ++++   | -+<br><br>                                |
|   |  | -+  |

| 64                          | 77   Diseases & Disorders   Fetal haemoglobin (HbF)   |
|-----------------------------|---|
| level modifies the clinical | severity of HBB disorders.                            |
|                             |   |
|                             |   |
|                             | +   |
|                             | · · · · · · · · · · · · · · · · · · ·                 |
|                             |   |
|                             |   |
|                             |   |
|                             |   |
|                             | +   |
| 5   PA                      | 18697826   1   Fetal haemoglobin                      |
| 0                           | 17   Genes & Molecular Sequences   HBB disorders      |
| 64                          | 77   Diseases & Disorders   Fetal haemoglobin (HbF)   |
| level modifies the clinical | severity of HBB disorders.                            |
| 1                           |   |
|                             | ++  |
|                             |   |
|                             | +   |
| •                           |   |
|                             |   |
|                             |   |
|                             |   |
|                             |   |
|                             | 18322311   0   TNFalpha gene                          |
| polymorphisms               | 16   43   Genes & Molecular Sequences                 |
| multiple sclerosis          | 69   87   |
| <del>-</del>                | , IL-1R and TNFalpha gene polymorphisms in Iranian    |
| patients with multiple scle |   |
| 1                           |   |
| ++                          |   |
|                             | ++  |
|                             | +   |
| +                           |   |
|                             |   |
|                             |   |
|                             |   |
|                             |   |
|                             | +   |
|                             | 18322311   7   TNFalpha G -308 allele                 |
| 138                         | 160   Genes & Molecular Sequences   multiple          |
| sclerosis                   | 406   424   Diseases &                                |
|                             | nd the frequency of IL-1alpha TT -889 genotype        |
| -                           | 0 allele (p=0.0001) and CC genotype (p=0.00006),      |
| <del>-</del>                | .0002) and GG genotype (p=0.000001) decreased         |
| significantly in the patien | ts versus normal subjects. These results suggest that |

| polymorphic variations of these pro-inflammatory cytokines may play an important process of these pro-inflammatory cytokines may play an important play in susceptibility of Iranian multiple sclerosis patients. |                  |
|---|------------------|
| 12   PA   |                  |
|   |                  |
| 13   PA   | of<br>and<br>93C |
|   |                  |
| 14   PA   | and              |

| ++   |
|--|
|  |
|  |
|  |
|  |
|  |
|  |
| +  |
| 16   PA  |
| 4   9   SNP & Sequence variations   oropharyngeal  |
| 112   125   Diseases & Disorders   The T393C SNP could be  |
| considered as a genetic marker to predict the clinical course of patients  |
| suffering from oropharyngeal and hypopharyngeal cancer.  |
| ++   |
|  |
|  |
| +  |
|  |
|  |
|  |
|  |
| +<br>  20   PA   |
| 20   PA  |
| 57   69   Diseases & Disorders   A recent study showed   |
| that LPA-mediated proliferation of colon cancer cells requires activation of   |
| beta-catenin.  |
|  |
| ++   |
| +  |
|  |
| <del>t</del>   |
|  |
|  |
|  |
| +  |
| 21   PA  |
| *3/*3   29   57   Genes & Molecular Sequences  |
| leukocytopenia   211   225   |
| Diseases & Disorders   On multivariate analysis the CYP3A5 A6986G genotype *3/*3   |
| (OR 8.205, 95% CI 1.616-41.667, p = 0.011) and smaller number of treatment   |
| cycles (OR 0.156, 95% CI 0.037-0.659, p = 0.011) were independent factors for leukocytopenia (grade 3 or greater) throughout the period of chemotherapy. |
| reakocy copenia (grade o or greater) unroughout the period or chemotherapy.  |

```
Sample rows for association type 'nan':
   +-----
   --+----+
   | ASSOCIATION_TYPE
               | PMID | NUM_SENTENCE
                                | ENTITY1_TEXT
                                          | ENTITY1_INI
   | ENTITY1_END | ENTITY1_TYPE | ENTITY2_TEXT
                                 | ENTITY2_INI
                                          | ENTITY2_END
   | ENTITY2_TYPE
            | SENTENCE
   +-----
   +-----
   --+----+
[18]: import pandas as pd
   # Load the CSV file
   ftrain = '../data/EUADR_target_disease.csv'
   df = pd.read_csv(ftrain, encoding='latin1', sep='\t')
   # Display the column names to identify the one that corresponds to the entity \Box
    →association type
   print("Column names:", df.columns)
   # Check if 'SA' exists in the 'ASSOCIATION_TYPE' column
   association_types = df['ASSOCIATION_TYPE'].unique()
   if 'NA' in association_types:
      print("Association type 'SA' is present in the data.")
   else:
      print("Association type 'SA' is not present in the data.")
   Column names: Index(['ASSOCIATION_TYPE', 'PMID', 'NUM_SENTENCE', 'ENTITY1_TEXT',
       'ENTITY1_INI', 'ENTITY1_END', 'ENTITY1_TYPE', 'ENTITY2_TEXT',
       'ENTITY2_INI', 'ENTITY2_END', 'ENTITY2_TYPE', 'SENTENCE'],
       dtype='object')
   Association type 'SA' is not present in the data.
```

### 2 Read Data

```
[13]: from tabulate import tabulate
      Tr_sent_contents, Tr_entity1_list, Tr_entity2_list, Tr_sent_lables = utils.
       ⇒dataRead_befree_EUADR(ftrain)
      # Convert the lists to a DataFrame
      df = pd.DataFrame({
          'Sent Contents': Tr_sent_contents,
          'Entity 1': Tr_entity1_list,
          'Entity 2': Tr entity2 list,
          'Sent Labels': Tr_sent_lables
      })
      # Display the first 10 rows of the DataFrame as a table
      print(tabulate(df.head(4), headers='keys', tablefmt='grid'))
      Tr_word_list, Tr_d1_list, Tr_d2_list = utils.
       →get_wordList_and_distances_befree(Tr_sent_contents)
      df = pd.DataFrame({
          'word list': Tr_word_list,
          'd1 list': Tr_d1_list,
          'd2 list': Tr_d2_list,
      })
      # Display the first 10 rows of the DataFrame as a table
      print(tabulate(df.head(4), headers='keys', tablefmt='grid'))
     print ("train_size", len(Tr_word_list))
```

Input File Reading
train\_size 355

### 2.1 Prepare Lable Matrix

```
[17]: # Y : is positive association
# N: is negative association
label_dict = {'FA':0, 'NA':0,'PA':1,'SA':1}

Y_t = mapLabelToId_befree_EUADR(Tr_sent_lables, label_dict)
Y_train = np.zeros((len(Y_t), 2))
for i in range(len(Y_t)):
    Y_train[i][Y_t[i]] = 1.0
```

# 3 Generate Word and Position Embedding Vectors

#### 3.0.1 Word Embedding

```
Found 1355 unique words (10062 in total) word dictonary length 1355
Reading word vectors
Loaded 1356 pretrained embeddings.
number of unknown word in word embedding 814
W_train 355
word_vectors 1355
```

#### 3.0.2 Position Embedding

```
[19]: d1_dict = makeDistanceList([Tr_d1_list])
d2_dict = makeDistanceList([Tr_d2_list])
d1_train = mapWordToId_list(Tr_d1_list, d1_dict)
d2_train = mapWordToId_list(Tr_d2_list, d2_dict)
```

#### 3.0.3 Pad Embdding Vectors

```
sentMax 102
W_train 355
d1_train 355
d2_train 355
```

# 4 Save Prepared Data as Pickle File

```
[21]: with open('../data/pickles/befree_EUADR_2class_PubMed-and-PMC-w2v.pickle',
       pickle.dump(W_train, handle)
         pickle.dump(d1_train, handle)
         pickle.dump(d2_train, handle)
         pickle.dump(Y_train, handle)
         pickle.dump(Tr_word_list, handle)
         pickle.dump(word_vectors, handle)
         pickle.dump(word_dict, handle)
         pickle.dump(d1_dict, handle)
         pickle.dump(d2_dict, handle)
         pickle.dump(label_dict, handle)
         pickle.dump(sentMax, handle)
 []:
 []:
 []:
```