

Bioinformatics III

Eigth Assignment

Alexander Flohr (2549738)

Andrea Kupitz (2550260)

June 18, 2018

Exercise 8.1: Data Preprocessing

Listing 1: Listing of source code

```
(a) import csv
import numpy
import math
from scipy import stats

5 class DataMatrix:
    def __init__(self, file_path):
        """
        :param file_path: path to the input matrix file
        """
10     self.file_path = file_path
    # TODO: define and initialise the class fields you need for your implementation
    # read the matrix in the input file, remove rows with empty values and merge duplicate
    self.cols, self.rows = self.read_data()

15     def read_data(self):
        """
        Reads data from a given matrix file, where the first line gives the names of the columns
        gives the names of the rows. Removes rows with empty or non-numerical values and merges
        name into one.
        """
20     with open(self.file_path) as tsvfile:
        reader = csv.reader(tsvfile, delimiter='\\t')
        cols = next(reader)
        fields = len(cols)
        rows = dict()
25     for row in reader:
        if len(row) == fields:
            floats = []
            nan = False
            for i in range(1, fields):
30                 try:
                    f = float(row[i])
                    if math.isnan(f):
                        nan = True
                except ValueError:
                    nan = True
                    floats.append(f)
            if row[0] not in rows and not nan:
                rows[row[0]] = [floats]
40     elif not nan:
        rows[row[0]].append(floats)
    for k, v in rows.items():
        if len(v) > 1:
            value = []
            for i in range(0, len(v[0])):
45                 mean = []
```

```

                    for j in range(0, len(v)):
                        mean.append(v[j][i])
                        value.append(numpy.mean(mean))
50                 rows[k] = [value]
            return cols, rows

    def get_rows(self):
        """
55         :return: dictionary with keys = row names, values = list of row values
        """
        rows = dict()
        for k, v in self.rows.items():
            rows[k] = v[0]
60        return rows

    def get_columns(self):
        """
65         :return: dictionary with keys = column names, values = list of column values
        """
        dic = dict()
        for k, v in self.rows.items():
            for i, col in enumerate(self.cols):
                if i != 0:
70                     if col not in dic:
                         dic[col] = []
                        else:
                            dic[col].append(v[0][i-1])
75        return dic

    def not_normal_distributed(self, alpha, rows):
        """
        Uses the Shapiro-Wilk test to compute all rows (or columns) that are not normally distributed
        :param alpha: significance threshold
80         :param rows: True if the Shapiro-Wilk p-values should be computed for the rows, False if for columns
        :return: dictionary with keys = row/columns names, values = Shapiro-Wilk p-value
        """
        dic = dict()
        values = self.get_rows() if rows else self.get_columns()
85        for k, v in values.items():
            W, p = stats.shapiro(v)
            if p < alpha:
                dic[k] = p
        return dic

90    def to_tsv(self, file_path):
        """
        Writes the processed matrix into a tab-separated file, with the same column order as the
        the rows in lexicographical order.
95         :param file_path: path to the output file
        """
        with open(file_path, 'w') as f:
            values = [value for (key, value) in sorted(self.get_rows().items())]
            keys = sorted(self.get_rows().keys())
100           for i in range(0, len(values)):
                s = keys[i] + '\t'
                for j in values[i]:
                    s += str(j) + '\t'
                s += '\n'
105           f.write(s)

```

Listing 1 shows the source code.

Listing 2: Listing of source code

```

(b) from data_matrix import DataMatrix
    from network import CorrelationNetwork
    from correlation import CorrelationMatrix
    from cluster import CorrelationClustering

```

```
5  def exercise_1():
    DM1 = DataMatrix('expression.tsv')
    DM1.to_tsv('Flohr-Kupitz-expression.tsv')
    print(len(DM1.not_normal_distributed(0.05, True)))
10  DM2 = DataMatrix('methylation.tsv')
    DM2.to_tsv('Flohr-Kupitz-methylation.tsv')
    print(len(DM2.not_normal_distributed(0.05, True)))

    def exercise_3(threshold=0.75):
15      CN1 = CorrelationNetwork(CorrelationMatrix(DataMatrix('expression.tsv'), 'Pearson', True),
                                CN1.to_sif('Flohr-Kupitz-expression-network-pearson.sif'))
        CN1 = CorrelationNetwork(CorrelationMatrix(DataMatrix('expression.tsv'), 'Spearman', True),
                                CN1.to_sif('Flohr-Kupitz-expression-network-spearman.sif'))
        CN1 = CorrelationNetwork(CorrelationMatrix(DataMatrix('expression.tsv'), 'Kendall', True),
                                CN1.to_sif('Flohr-Kupitz-expression-network-kendall.sif'))
20      CN1 = CorrelationNetwork(CorrelationMatrix(DataMatrix('methylation.tsv'), 'Pearson', True),
                                CN1.to_sif('Flohr-Kupitz-methylation-network-pearson.sif'))
        CN1 = CorrelationNetwork(CorrelationMatrix(DataMatrix('methylation.tsv'), 'Spearman', True),
                                CN1.to_sif('Flohr-Kupitz-methylation-network-spearman.sif'))
25      CN1 = CorrelationNetwork(CorrelationMatrix(DataMatrix('methylation.tsv'), 'Kendall', True),
                                CN1.to_sif('Flohr-Kupitz-methylation-network-kendall.sif'))

    def exercise_4():
        # TODO
30      pass

# only execute the following if this module is the entry point of the program, not when it is imported
    if __name__ == '__main__':
        exercise_1()
35      exercise_3()
        exercise_4()
```

Listing 2 shows the source code.

The number of genes and samples whose data does not follow a normal distribution for the expression data is 74 and for the methylation data is 66.

Exercise 8.2: Correlation Measures

All subtasks are implemented in `correlation.py`, see listing 3

Listing 3: Listing of source code

```
0 from itertools import combinations
  from copy import copy
  import math

  def computeMean(vec):
      """
      5   Compute the mean for a list
      """
      temp = 0
      for i in vec:
          temp += i
      10  return float(temp)/float(len(vec))

  def initializeIndices(x_s):
      """
      15  Create list of indices for values, where equal values
          in x_s get equal averaged idices
      """
      x_r = list(range(0, len(x_s)))
      # add an terminate-sign to end the algorithm
      20  x_r.append(-1)
      i = 0
      for j in range(1, len(x_s)):
          if x_s[i] != x_s[j]:
              index = computeMean(x_r[i:j])
              25  while i < j:
                  x_r[i] = index
                  i += 1
          # return all but the last termination-sign
          return x_r[0:len(x_s)]
      30

  def bubble(x_s):
      """
      Recursive bubble sorting with "rebubbling"
      the indices with respect to multi-occurring
      35  values
      """
      for i in range(1, len(x_s)):
          if x_s[i-1] < x_s[i]:
              # swap list elements
              40  temp = x_s[i-1]
                  x_s[i-1] = x_s[i]
                  x_s[i] = temp
              # next bubble
              x_r = bubble(x_s)
              45  # reverse bubbleing of the indices
                  temp = x_r[i-1]
                  x_r[i-1] = x_r[i]
                  x_r[i] = temp
              return x_r
      50  # called only in the deepest recursive step
          # here the indices get initialized
          return initializeIndices(x_s)

  def rank(x):
      """
      55  :param x: a list of values
          :return: ranking of the input list
      """
      # sort the values of x by bubble sort
      60  # original indices are kept in a second list
```

```
    return bubble(copy(x))

def pearson_correlation(x, y):
    """
65     :param x: a list of values
        :param y: a list of values
        :return: Pearson correlation coefficient of X and Y
        """
    mean_x = computeMean(x)
70    mean_y = computeMean(y)

    # covariance
    upper = 0
    for i in range(0, len(x)):
75        upper += (x[i] - mean_x) * (y[i] - mean_y)
    lower_left = 0

    # standard deviation of x
    for i in range(0, len(x)):
80        lower_left += (x[i] - mean_x) * (x[i] - mean_x)
    lower_left = math.sqrt(lower_left)

    # standard deviation of y
    lower_right = 0
85    for i in range(0, len(y)):
        lower_right += (y[i] - mean_y) * (y[i] - mean_y)
    lower_right = math.sqrt(lower_right)

    # pearson-correlation formula
90    return (upper/(lower_left*lower_right))

def spearman_correlation(x, y):
    """
95     :param x: a list of values
        :param y: a list of values
        :return: Spearman correlation coefficient of X and Y
        """
    rank_x = rank(x)
    rank_y = rank(y)
100    return pearson_correlation(rank_x, rank_y)

def kendall_correlation(x, y):
    """
105     :param x: a list of values
        :param y: a list of values
        :return: Kendall-B correlation coefficient of X and Y
        """

    # processing part (1)
110    rank_x = rank(x)
    rank_y = rank(y)

    # processing part (2)
    pairing = []
115    for i in range(0, len(x)):
        pairing.append((rank_x[i], rank_y[i]))

    # processing part (3)
    concordant = 0
120    discordant = 0
    tied_x = 0
    tied_y = 0
    for i in range(0, len(pairing)):
        for j in range(0, len(pairing)):
125            if i != j:
                (a,b) = pairing[i]
                (c,d) = pairing[j]
```

```

        if (a < b and c < d) or (a > b and c > d):
            concordant += 1
130         elif (a < b and c > d) or (a > b and c < d):
            discordant += 1
        elif (a == b and c != d):
            tied_x += 1
        elif (a != b and c == d):
135             tied_y += 1

    # processing part (4)
    score1 = float(concordant - discordant)
    score2 = float(concordant + discordant + tied_x)
140    score3 = float(concordant + discordant + tied_y)
    if score1*score2*score3 == 0:
        return 0
    else:
        return score1/math.sqrt(score2 * score3)
145

class CorrelationMatrix(dict):
    """
    This class behaves like a dictionary, where the correlation between two elements 1 and 2 is acc
    cor_matrix[(element_1, element_2)] or cor_matrix[(element_2, element_1)] since the matrix is sym
    It also stores the row (or column) names of the input DataMatrix.
150    """

    def __init__(self, data_matrix, method, rows):
        """
        :param data_matrix: a DataMatrix (see data_matrix.py)
155        :param method: string specifying the correlation method, must be 'Pearson', 'Spearman' or 'Kendall'
        :param rows: True if the correlation matrix should be constructed for the rows, False if for columns
        """
        # initialise the dictionary
        super().__init__(self)

160        # if rows = True, then compute the correlation matrix for the row data
        if rows:
            data = data_matrix.get_rows()
        # if rows = False, then compute the correlation matrix for the column data
165        else:
            data = data_matrix.get_columns()

        # sorted list of row names (or column names) in the input data matrix
        self.names = list(sorted(data.keys()))

170        # compute the correlation between all pairs of rows (or columns)
        for name_1, name_2 in combinations(data.keys(), 2):
            # use the specified correlation method
            if method == 'Pearson':
175                correlation = pearson_correlation(data[name_1], data[name_2])
            elif method == 'Spearman':
                correlation = spearman_correlation(data[name_1], data[name_2])
            elif method == 'Kendall':
                correlation = kendall_correlation(data[name_1], data[name_2])
180            else:
                raise ValueError('The correlation method not supported must be either Pearson, Spearman or Kendall')

            # add the correlation symmetrically
            self[(name_1, name_2)] = correlation
185            self[(name_2, name_1)] = correlation
```

Exercise 8.3: Gene Co-Expression Networks

Listing 4: Listing of source code

```
(a) from Node import Node

class CorrelationNetwork:
    def __init__(self, correlation_matrix, threshold):
        """
5         Constructs a co-expression network from a correlation matrix by adding edges between nodes
        correlation bigger than the given threshold.
        :param correlation_matrix: a CorrelationMatrix (see correlation.py)
        :param threshold: a float between 0 and 1
        """
10        self.cor = correlation_matrix
        self.nodes = {}
        for entry, cor in correlation_matrix.items():
            node1 = Node(entry[0])
            node2 = Node(entry[1])
15            self.appendNode(node1)
            self.appendNode(node2)
            if abs(cor) > threshold:
                node1.addLinkTo(node2)
                node2.addLinkTo(node1)
20
        def appendNode(self, node):
            """
            Appends node to network
            """
25            if node.id not in self.nodes:
                self.nodes[node.id] = node

        def to_sif(self, file_path):
            """
30            Write the network into a simple interaction file (SIF).
            Column 0: label of the source node
            Column 1: interaction type
            Columns 2+: label of target node(s)
            :param file_path: path to the output file
            """
35            with open(file_path, 'w') as f:
                values = [value for (key, value) in sorted(self.nodes.items())]
                keys = sorted(self.nodes.keys())
                written = []
40                s = ''
                for position, value in enumerate(keys):
                    for node2id in values[position].nodelist:
                        if node2id not in written: #and node2id != value:
                            s = str(value) + '\t'
                            s += str(round(self.cor[(value, node2id.id)], 2)) + '\t'
45                            s += str(node2id) + '\n'
                            written.append(node2id)
                f.write(s)
```

Listing 4 shows the source code.

(b) Listing 2 shows the source code.

(c)

Exercise 8.4: Hierarchical Clustering

- (a)
- (b)
- (c)