# Reducing noise in protein multialignments

Marina Herrera Sarrias

January 6, 2019

# Abstract

In this research a noise-reduction method is implemented in protein multialigments to evaluate its impact on phylogenic inference.The data used in this project is a reduced data set from the original data used by the creators of TrimAI. To test the data the program fastprot was used to obtain the distance matrices and fnj to infer a phylogenic tree for each alignment. Dendropy has been used to measure the symmetric distance between the inferred tree and the reference tree.

# Contents

# 1 Introduction

The implemented noise-reduction method evaluates a multialigment column as noisy if there are more than 50% indels, at least 50% of the amino acids are unique and non of the amino acids appear more than twice. Each column takes the form of a sequence of amino acids given a position of the aligned sequences.

The reduced data set used in this project is composed by six different directories, each directory contains a reference tree and 300 alignments created by evolving sequences along the reference tree. Each pair of directories present an average amount of mutations of 0.5, 1.0 and 2.0 per sequence site. Also, per each mutation rate there is one symmetric reference tree and one asymmetric reference tree.

For validating the effect of reducing reducing

# 2 Results and Discussion

# 3 Methods and Materials

# 4 run_program.py

```python
1  import sys
2  import glob
3  import os
4  import json
5  import subprocess
6  from collections import Counter
7
8  from Bio.Seq import Seq
9  from Bio import SeqIO, AlignIO, Align, SeqRecord
10 import dendropy
11 from dendropy.calculate import treecompare
12 from Bio import Phylo
13
14 """
15 This script will perform all the program's required tasks.
16 It will first create a list of all the data directories where
       we will reduce noise,
17 new directories for storing the reduced alignments and an empty
        dictionary
18 with key 'main directory name' (e.g asymmetric_0.5) and which
       will
19 later contain a nested dictionary:
20 {
21     'filename': [
22     original_alignment_symmetric_distance,
23     noise_reduced_alignment_symmetric_distance,
24     noise_reduction_ratio_
25     ]
26     ...
27 }
28  with 'filename' e.g: s001.align.1.msl
29 and original_alignment_symmetric_distance,
      noise_reduced_alignment_symmetric_distance,
      noise_reduction_ratio_:
30     distance of the reference tree with the original alignment
          tree and
31     distance of the reference tree with the noise reduced
          alignment tree
32     ratio between the difference of number of columns between
          the original and reduced alignment and the number of
          rows
33     of the original alignment.
```

```python
34
35  The program will start parsing through the raw data , it will
        first read the reference tree ,
36  set the keys of the dictionary and write the noise reduced
        alignment file using the
37  perform_noise_reduction function , afterwards it will generate
        inferred tree files , for both :
38  the raw data and the noise reduced data using the function
        computing_and_writing_alignment_tree .
39  Lastly it will compute the distance between the reference tree
        and both inferred tree cases
40  (with and without noise reduction ). The file names (key),
        distances and ratios (key values) will be nested
41  in the dictionary .
42
43  The program will write a json file in the results folder
        containing the dictionary .
44  """
45
46
47  def noise_filter ( column ):
48      """
49      noise_filter is a function that determines whether an
            alignment column is noisy or not .
50      :param column : multiple sequence alignment column .
51      :return : boolean .
52      """
53      indel = '-'
54      amino_counter = Counter ( column )
55      if column . count ( indel ) > len ( column ) / 2:
56          return True
57      if len ([ key for key , value in amino_counter . items () if
            value == 1]) >= len ( column ) / 2:
58          return True
59      if max ( amino_counter . values ()) <= 2:
60          return True
61      return False
62
63
64  def reduce_noise ( multiple_seq_alignment , _alignment_path ):
65      """
66      reduce_noise is a function that filters all noisy columns .
67      :param multiple_seq_alignment : MultipleSeqAlignment .
68      :param _alignment_path : msa path .
69      :return : a MultipleSeqAlignment .
```

```python
70          """
71          noise_free_column_list = []
72          new_aligned_list = []
73          alignment_filename_ = _alignment_path.split('/')[-1]
74
75          for i in range(multiple_seq_alignment.get_alignment_length
                ()):
76              if not noise_filter(multiple_seq_alignment[:, i]):
77                  noise_free_column_list.append(
78                      multiple_seq_alignment[:, i])
78          for i in range(len(multiple_seq_alignment)):
79              new_aligned_list.append(''.join([x[i] for x in
                    noise_free_column_list]))
80          record_list_ = Align.MultipleSeqAlignment([
81              SeqRecord.SeqRecord(Seq(new_seq), id=record.id, name=
                    record.name, description=record.description)
82              for new_seq, record in zip(new_aligned_list,
                    multiple_seq_alignment)
83          ])
84          noise_reduction_ratio_ = (multiple_seq_alignment.
                get_alignment_length()
85                          - record_list_.get_alignment_length()) /
                            multiple_seq_alignment.
                            get_alignment_length()
86
87          if multiple_seq_alignment.get_alignment_length() ==
                record_list_.get_alignment_length():
88              print(f'{alignment_filename_}'+ '>> WARNING: no noise
                    reduction')
89          if 0 < record_list_.get_alignment_length()  <
                multiple_seq_alignment.get_alignment_length() / 2:
90              print(f'{alignment_filename_}'+ '>> WARNING: more than
                    0.5 of the sequence is noise')
91          return record_list_, noise_reduction_ratio_
92
93
94
95  def perform_noise_reduction(alignment_path_):
96          """
97          perform_noise_reduction is a function that takes the
                parameter alignment path and applies
98          the reduce_noise function.
99          :param alignment_path_: msa path.
100         :return: a MultipleSeqAlignment.
101         """
```

```python
102    with open(alignment_path_, mode='r') as aligned_file:
103        my_sequence_recorded = AlignIO.read(aligned_file, '
               fasta')
104    return reduce_noise(my_sequence_recorded, alignment_path_)
105
106
107 def computing_and_writing_alignment_tree(msa_filename,
        tree_outfile_):
108    """
109    computing_and_writing_alignment_tree is a function that
           pipes two processes,
110    it first creates a distance matrix for each msa using
           fastprot,
111    an then infers an alignment tree using fnj. The function
           writes the inferred
112    tree in a file per msa.
113    :param msa_filename: msa path.
114    :param tree_outfile_: inferred tree path.
115    :return: file containing the inferred tree.
116    """
117    args= "cat " + msa_filename + " | fastprot | fnj -O newick
           -o " + tree_outfile_
118    child_noise_reduced = subprocess.Popen(args, shell=True)
119    child_noise_reduced.wait()
120
121
122 def compare_trees(ref_tree, inferred_tree,
        is_bipartitions_updated=False):
123    """
124    compare_trees is a function that computes the simmetric
           difference between
125    the reference tree and the inferred tree.
126    :param ref_tree: reference tree read using dendropy.
127    :param inferred_tree: inferred tree read using dendropy.
128    :param is_bipartitions_updated: recalculates bipartitions.
129    :return: (int) symmetric distance between the reference and
            inferred tree.
130    """
131    sd = treecompare.symmetric_difference(ref_tree,
           inferred_tree)
132    return sd
133
134
135 def compute_distance_between_trees(inferred_tree_file_path,
        reference_tree_):
```

```
136          """
137          compute_distance_between_trees is a function that reads
                 inferred trees using
138          dendropy and computes the distance using the function
                 compare_trees
139          :param inferred_tree_file_path: path
140          :param reference_tree_: reference tree read using dendropy.
141          :return: (int) symmetric distance between the reference and
                 inferred tree.
142          """
143          with open(inferred_tree_file_path, mode='r') as
                 reduced_noise_tree_file:
144              reduced_noise_tree_str = ''.join(list(
                     reduced_noise_tree_file))
145              reduced_noise_tree = dendropy.Tree.get_from_string(
146                  reduced_noise_tree_str,
147                  schema="newick",
148                  taxon_namespace=tns
149              )
150          return compare_trees(reference_tree_, reduced_noise_tree)
151
152
153  if __name__ == '__main__':
154      args_in = sys.argv[1:]
155      if len(args_in) > 0:
156          print('              ')
157          sys.exit('WARNING: run_program do not require any input
                 .')
158      print('processing data...')
159
160      original_dir = './data/raw_data'
161      reduced_dir = './data/noise_reduced_data'
162      result_dir = './results'
163      directories = []
164      tns = dendropy.TaxonNamespace()
165
166      for folder in glob.glob(original_dir +'/*'):
167          if len(glob.glob(original_dir + '/*')) == 0:
168              sys.exit('empty data directory, please verify your
                     data is on the right directory')
169
170          sub_folder_name = folder.split('/')[-1]
171          directories.append(sub_folder_name)
172
```

```python
173        # this will contain all the symmetric distances between
               inferred trees
174        # and reference trees, for all different data
               subdirectories
175        compare_trees_dictionary = dict()
176
177        for directory in directories:  # directories = [
               asymmetric_0.5, asymmetric_1.0, ...]
178            # create asymmetric_0.5, asymmetric_1.0, ...
                   subdirectories into the reduced_data directory
179            # and getting the reference tree for each subdirectory
                   of alignments.
180            new_folder_path = os.path.join(reduced_dir, directory)
181            os.makedirs(new_folder_path, exist_ok=True)
182            original_dir_path = os.path.join(original_dir,
                   directory)
183            if len(glob.glob(original_dir_path + '/*')) == 0:
184                sys.exit('empty data sub directory, please verify
                       your data is on the right directory')
185
186            reference_tree_path = glob.glob(original_dir_path +'/*.
                   tree')[0]
187            with open(reference_tree_path, mode='r') as
                   ref_tree_file:
188                ref_tree_str = ''.join(list(ref_tree_file))
189                reference_tree = dendropy.Tree.get_from_string(
190                    ref_tree_str,
191                    schema="newick",
192                    taxon_namespace=tns
193            )
194            folder_dict_key_name = original_dir_path.split('/')[-1]
195
196            compare_trees_dictionary[folder_dict_key_name] = dict()
197
198            for alignment_path in glob.glob(original_dir_path + '
                   /*.msl'):
199                if os.stat(alignment_path).st_size == 0:
200                    sys.exit(f'{alignment_path}'+ '>> ERROR: empty
                           msl file.')
201                # performing noise reduction and writing the
                       alignment
202                # into the corresponding noise_reduction directory
203                record_list, noise_reduction_ratio =
                       perform_noise_reduction(alignment_path)
204                if record_list.get_alignment_length() == 0:
```

```python
                sys.exit(f'{alignment_path}'+ '>> WARNING: all
                    columns are noisy.')
            reduced_alignment_name = alignment_path.split('/')
                [-1]
            reduced_filename_out = os.path.join(new_folder_path
                , reduced_alignment_name)
            AlignIO.write(record_list, reduced_filename_out, "
                fasta")

            # computing and writing noise reduced alignment
                trees
            tree_outfile_reduced = reduced_filename_out[:-3] +
                'tree'
            computing_and_writing_alignment_tree(
                reduced_filename_out, tree_outfile_reduced)
            if os.stat(tree_outfile_reduced).st_size == 0:
                tree_file_name = tree_outfile_reduced.split('/'
                    )[-3:]
                sys.exit(f'{tree_outfile_reduced}' + '>> ERROR:
                    empty tree file.')

            # computing and writing original alignment trees
            alignment_name = alignment_path.split('/')[-1]
            filename_out = os.path.join(original_dir_path,
                alignment_name)
            tree_outfile = filename_out[:-3] + 'tree'
            computing_and_writing_alignment_tree(alignment_path
                , tree_outfile)
            if os.stat(tree_outfile).st_size == 0:
                tree_file_name = tree_outfile.split('/')[-3:]
                sys.exit(f'{tree_outfile}' + '>> ERROR: empty
                    tree file.')

            # computing distance between ref tree and inferred
                trees
            noise_reduced_distance =
                compute_distance_between_trees(
                tree_outfile_reduced, reference_tree)
            original_distance = compute_distance_between_trees(
                tree_outfile, reference_tree)

            alignment_key_name = alignment_name[:-3]
            compare_trees_dictionary[folder_dict_key_name][
                alignment_key_name] = (
```

```
232                    original_distance , noise_reduced_distance ,
                          round ( noise_reduction_ratio , 2)
233              )
234
235     distance_results_path = os.path.join( result_dir , '
           distance_result_dict ')
236     with open( distance_results_path , 'w ') as result_dir_file :
237         json.dump ( compare_trees_dictionary , result_dir_file )
238
239     print ('your data has been processed !')
```