

MSARC – implementation details

Michał Modzelewski

Norbert Dojer

The MSARC implementation is written in the Python language in the version 2.7 dialect. It makes heavy use of the `numpy` library for array operations, as well as `scipy.sparse` for storing the posterior probability matrices in compressed row format. The BioPython library is used for sequence input and output. A customized version of the `probA` program is used for the pairwise stochastic alignment of sequences. The `probA` implementation of the partition function was chosen because the available source code provides several substitution matrices to work with, and solves the problem of numerical overflow when aligning large sequences. The `probA` tool is exposed to Python as an extension module written using Cython. The consistency transformation, the main loop of the Fiduccia-Mattheyses partitioner, the gains computation methods, and the maximal expected accuracy alignment algorithm are also written as Cython extension modules for speed. Code listings below are simplified for legibility.

1 Module `balibase`

The `balibase` module is a thin wrapper around the `bali_score` program from BALIBASE 3.0, which is expected to be in the same directory as the calling script. It is used by the `balitest` and `balitest-compare` programs.

2 Module `clusteralign`

This module contains the main implementation of the MSARC method conceptually split into several files.

`alignment.py`

Contains the definition of the `Alignment` class which describes and can construct an alignment. This class is the main access point of the MSARC algorithm. An instance of `Alignment` is created with a string containing the name of a file containing the aligned sequences in *fasta* format, and keyword arguments containing the algorithm parameters. An instance of `Alignment` exposes some public methods and properties.

alignment A property for accessing the alignment as a `MultipleSeqAlignment` instance. This can be calculated on access from the sequences.

value A property that returns the sum of pairwise posterior probabilities value of the alignment.

```
1 @property
2 def value(self):
3     # self._columns is a list of Graph instances
4     return sum(c.value for c in self._columns)
```

output A method that prints the alignment to the screen in any format supported by the BioPython library or the *msf* format (default).

write A method that writes the alignment to the named file in any format supported by the BioPython library or the *msf* format (default).

read A method that reads an alignment from the named file.

Accessing the alignment property causes an alignment to be calculated if the object does not yet contain an alignment. To this end the alignment property getter method calls some private methods.

`_relax` A method that performs the consistency transformation on the posterior probabilities the desired number of times. The actual consistency transformation is delegated to the Graph class.

`_align` A method that recursively partitions the alignment graph (in the form of an instance of the Graph class) until a list of columns is produced.

```

1 def _align(self, graph_or_starts, ends=None):
2     """Align the sequences."""
3     if isinstance(graph_or_starts, Graph):
4         graph = graph_or_starts
5     else:
6         graph = self._graph.slice(graph_or_starts, ends)
7     graphs = [graph]
8     columns = []
9     while len(graphs) > 0:
10        g = graphs.pop()
11        # Graph depicts a single column, add to columns
12        if g.is_column:
13            columns.append(g)
14            continue
15        # Iteratively partition the graph
16        l, r = g.split()
17        graphs.append(r)
18        graphs.append(l)
19    return columns

```

`_refine` A method that performs iterative refinement on the alignment.

graph.py

The definition of the Graph class which is an instance of an alignment graph. An instance of Graph is created with a list of SeqRecord objects describing the aligned sequences, and keyword arguments containing the algorithm parameters. The Graph class knows how to partition itself, return a slice of itself, perform an iteration of the consistency transformation, and coarsen and uncoarsen itself for the purpose of multilevel partitioning. It populates, on demand, the values of the weights of the edges using the clusteralign.proba module. Some public properties and methods of the Graph class:

`size` A property returning the total number of nodes in the graph.

`num_sequences` A property returning the number of sequences.

`lengths` A property that returns an array of lengths of the sequences.

`max_length` A property that returns the length of the longest sequence.

`value` A property returning the total *value* of the graph. Only inter-sequence edges are counted, because only those can remain uncut after alignment.

```

1 @property
2 def value(self):
3     return sum((m.sum() for m in self._pair_pp.flat if m is not None))

```

`starts` A property that returns an array of starting points of the sub-sequences within the sequences.

`ends` A property that returns an array of end points of the sub-sequences within the sequences.

`is_column` A property that returns whether the graph represents a valid *column*.

```

1 @property
2 def is_column(self):
3     return self.max_length <= 1

```

nonzero A property returning the number of non-zero length sequences.

```
1 @property
2 def nonzero(self):
3     return (self.lengths > 0).sum()
```

partition A property that returns the partition of the graph as an array of partition points within the sequences.

relax A method that runs an iteration of the consistency transformation.

```
1 from numpy import empty_like
2 from clusteralign.matrix import prune
3 from clusteralign._graph import relax_XZ_ZY, relax_ZX_ZY, relax_XZ_YZ
4
5 def relax(self, weighted_transform=False):
6     pair_pp = self._pair_pp
7     self._pair_pp = new_pair_pp = empty_like(pair_pp)
8     cut = self._params['cut']
9     # for every pair of sequences
10    n = self.num_sequences
11    weights=[[1]*n for i in xrange(n)]
12    if weighted_transform:
13        for x in xrange(n):
14            for y in xrange(x + 1, n):
15                weights[x][y]=pair_pp[x,y].sum()/min(pair_pp[x,y].shape)
16    for x in xrange(n):
17        for y in xrange(x + 1, n):
18            # P/(z==x) + P/(z==y)
19            sum_w = weights[x][y]*2
20            new_pp_xy = pair_pp[x, y] * sum_w
21            # P/(z!=x&&z!=y)
22            for z in xrange(n):
23                if z is x or z is y:
24                    continue
25                if z < x:
26                    relax_ZX_ZY(pair_pp[z, x], pair_pp[z, y], new_pp_xy, weights[z][x]*weights[z][y]
27                    sum_w += weights[z][x]*weights[z][y]
28                elif x < z < y:
29                    relax_XZ_ZY(pair_pp[x, z], pair_pp[z, y], new_pp_xy, weights[x][z]*weights[z][y]
30                    sum_w += weights[x][z]*weights[z][y]
31                else:
32                    relax_XZ_YZ(pair_pp[x, z], pair_pp[y, z], new_pp_xy, weights[x][z]*weights[y][z]
33                    sum_w += weights[x][z]*weights[y][z]
34            # normalize
35            new_pp_xy.data /= sum_w
36            # cut-off
37            prune(new_pp_xy, cut)
38            new_pair_pp[x, y] = new_pp_xy
```

slice A method that returns a sub-graph corresponding to sub-sequences with the given start and end points.

split A method that splits the graph into to sub-graphs along the given partition.

```
1 def split(self, parts=None):
2     if parts is None:
3         parts = self.partition
4     return (self.slice(self.starts, parts),
5             self.slice(parts, self.ends))
```

matrix.py

Some helper functions that operate on numpy arrays.

prune A function that zeros all values in a matrix below a given threshold.

```

1 from scipy.sparse import isspmatrix_csr
2
3 def prune(array, cut):
4     assert isspmatrix_csr(array)
5     array.data[array.data < cut] = 0.
6     array.eliminate_zeros()

```

msf.py

Unlike BioPerl, the Python library BioPython does not support the *msf* multiple sequence alignment format. Since this is the preferred alignment format of the BALIBASE benchmark, a custom implementation was required. This module provides functions for reading and writing such files.

options.py

This module implements an options parser for the alignment programs using Python's built-in argparse library.

partition.py

The Graph class delegates the heavy lifting of searching for an optimal partition to the GraphPartitioner class implemented in this module. A GraphPartitioner instance is created with an instance of Graph that is to be partitioned. The public method partition is exposed that takes an initial partition as an argument and returns a new partition optimized by an implementation of the Fiduccia-Mattheyses algorithm. Helper properties and methods implement the details of the process.

partition A method that finds a balanced minimum-cut partition of the alignment graph by iterating the Fiduccia-Mattheyses algorithm until convergence.

```

1 def partition(self, parts):
2     self._parts = parts
3     while True:
4         best_parts = self._testmoves()
5         if best_parts is None:
6             break
7         self._move([(i, best_parts[i]) for i in xrange(self._numgroups)])
8     return self._parts

```

_testmoves Iteratively selects nodes to move between the two sides of the partition, and calculates the gain.

```

1 from clusteralign._partition import bestmove
2
3 def _testmoves(self):
4     parts = self._parts.copy()
5
6     gain = 0.
7     bestgain = None if min(self._balance_weight) < self._minbalance else gain
8     bestparts = parts
9
10    allowed = self._new_allowed(parts)
11    while True:
12        if self._balance_weight[0] < self._minbalance:
13            direction = 1
14        elif self._balance_weight[1] < self._minbalance:
15            direction = -1
16        else:
17            direction = 0
18
19        move = bestmove(self._parts, self._gains, allowed, direction)
20        if move is None:
21            break
22        group, partition, delta_gain = move

```

```

23         gain += delta_gain
24
25         self._movepartition(group, partition)
26         allowed[group][partition] = 0
27
28         if min(self._balance_weight) < self._minbalance:
29             continue
30         if bestgain is None or gain > bestgain:
31             bestparts = self._parts.copy()
32             bestgain = gain
33
34         if bestparts is parts:
35             self._parts = parts
36             return None
37         return bestparts

```

_move Performs a series of moves.

_movepartition Moves the partition boundary and recalculates balance and gains.

```

1 def _movepartition(self, group, partition):
2     if partition == self._parts[group]:
3         #no change
4         return
5
6     if self._gains is not None:
7         self._fixgains(group, partition)
8
9     p_group = self._parts[group]
10    max_partition, min_partition = max(partition, p_group), min(partition, p_group)
11    dweights = self._weights_cumsum[group][max_partition - 1]
12    if min_partition > 0:
13        dweights -= self._weights_cumsum[group][min_partition - 1]
14
15    side = self._getside(group, partition)
16    other = 1 if side == 0 else 0
17    self._balance_weight[side] -= dweights
18    self._balance_weight[other] += dweights
19
20    self._parts[group] = partition

```

_edges_sum A property that returns the edge weight matrix.

_calcgain_all Calculates the gains of moving the partition point to any other point in the sequence.

```

1 def _calcgain_all(self, group):
2     l_group = self._lengths[group]
3     gains = zeros(l_group + 1, dtype=float)
4     if not l_group:
5         return gains
6
7     for i in xrange(self._numgroups):
8         if i == group:
9             continue
10
11        p_i = self._parts[i] # the current point of partition
12        l_i = self._lengths[i]
13
14        edges = self._edges[group, i]
15        if p_i > 0:
16            gains[1:] += edges[:, :p_i].sum(axis=1)
17        if p_i < l_i:
18            gains[1:] -= edges[:, p_i:].sum(axis=1)
19
20        # divide by 2 here so we don't have to double the difference in
21        # _fixgain_all, which is called more frequently
22        return gains / 2.

```

`_calcgains` Calculate the gains for all the sequences.

```

1 def _calcgains(self):
2     g = empty(self._numgroups, dtype=object)
3     for i in xrange(self._numgroups):
4         g[i] = self._calcgain_all(i)
5     return g

```

`_fixgain_all` Adjusts the gains for the sequence after the partition has been moved in another sequence.

```

1 def _fixgain_all(self, group, move_group, move_partition):
2     if move_group == group:
3         raise ValueError
4     l_group = self._lengths[group]
5     p_move_group = self._parts[move_group]
6     gains = zeros(l_group + 1, dtype=float)
7     if not l_group or move_partition == p_move_group:
8         return gains
9
10    edges = self._edges[group, move_group]
11    if move_partition > p_move_group:
12        gains[1:] += edges[:, p_move_group:move_partition].sum(axis=1)
13    else: # less
14        gains[1:] -= edges[:, move_partition:p_move_group].sum(axis=1)
15
16    return gains

```

`_fixgains` Adjusts the gains for all the sequences.

```

1 def _fixgains(self, group, partition):
2     p_group = self._parts[group]
3     if partition == p_group:
4         #no change
5         return
6     g = self._gains
7     for i in xrange(self._numgroups):
8         if i == group or not self._lengths[i]:
9             continue
10        g[i] += self._fixgain_all(i, group, partition)
11    g[group] -= g[group][partition]

```

probA.py

This module defines functions for calculating the pairwise posterior probabilities of sequence residue alignment. The actual calculations are delegated to the `clusteralign._probA` module.

`align` A function that returns the pairwise posterior probabilities matrix for the given pair of sequences.

`pairwise` A function that returns a matrix of pairwise posterior probabilities matrices for the sequences provided.

sequence.py

A module implementing a subclass of BioPython's `Seq` class called `EmptySeq` that is an empty sequence. This is required when partitioning the alignment graph leaves a slice of one sequence that is of zero length.

_graph.pyx

A Cython extension module that implements efficient helper functions for the `clusteralign.graph` module.

`relax_XZ_ZY` A function that computes a sparse matrix multiplication and accumulates the result in a third sparse matrix. Only results that are non-zero in the third matrix are stored. Sparse matrices are stored in index sorted compressed row format.

```

1 def relax_XZ_ZY(matXZ, matZY, matXY, weight):
2     for x in range(matXY.shape[0]):
3         startXY = matXY.indptr[x]
4         endXY = matXY.indptr[x + 1]
5         for z in range(matXZ.indptr[x], matXZ.indptr[x + 1]):
6             startZY = matZY.indptr[maxXZ.indices[z]]
7             endZY = matZY.indptr[maxXZ.indices[z] + 1]
8             valXZ = matXZ.data[z]
9             iterXY = startXY
10            indexZY = matZY.indices[startZY]
11            indexXY = matXY.indices[iterXY]
12            while startZY < endZY and iterXY < endXY:
13                if indexZY < indexXY:
14                    startZY += 1
15                    indexZY = maxZY.indices[startZY]
16                elif indexZY > indexXY:
17                    iterXY += 1
18                    indexXY = matXY.indices[iterXY]
19                else: # equal
20                    matXY.data[iterXY] += valXZ * matZY.data[startZY] * weight
21                    startZY += 1
22                    iterXY += 1
23                    indexZY = maxZY.indices[startZY]
24                    indexXY = matXY.indices[iterXY]

```

relax_ZX_ZY Analogous to the previous function but for the case when the first matrix is transposed.

relax_XZ_YZ Analogous to the previous function but for the case when the second matrix is transposed.

__partition.pyx

A companion module for the `clusteralign.partition` module that efficiently implements long running code as an extension module for Python written in Cython.

bestmove A function that returns the best allowed move for the Fiduccia-Mattheyses algorithm.

```

1 def bestmove(parts, gains, allowed, direction):
2     for i in range(num_groups):
3         l = allowed[i].shape[0]
4         if not l:
5             continue
6         if direction < 0:
7             start, end = 0, parts[i]
8         elif direction > 0:
9             start, end = parts[i], l
10        else:
11            start, end = 0, l
12            group_max_partition = -1
13            group_max_gain = 0.
14            gain = 0.
15            curgain = 0.
16            for j in range(min(start, curpart), max(end, curpart + 1)):
17                gain += gains_item[j]
18                if j == curpart:
19                    curgain = gain
20                if not allowed[i][j]:
21                    continue
22                if group_max_partition == -1 or gain > group_max_gain:
23                    group_max_partition = j
24                    group_max_gain = gain
25            if group_max_partition >= 0:
26                group_max_gain -= curgain
27                if max_partition == -1 or group_max_gain > max_gain:
28                    max_group = i
29                    max_partition = group_max_partition
30                    max_gain = group_max_gain

```

```

31
32     if max_partition == -1:
33         return None
34     return max_group, max_partition, max_gain

```

`_ppalign.pyx`

An extension module for Python written in Cython that implements the dynamic programming algorithm for computing the maximum expected accuracy alignment of a posterior probability matrix. This is used during iterative refinement by the `clusteralign.alignment` module.

align_probabilities A function that calculates the maximum expected accuracy alignment path for the supplied posterior probability matrix, and returns mappings of matches and spaces for the two sequences.

`_probA.pyx`

An extension module for Python written in Cython that wraps and exposes the customized `probA` tool.

align A function that returns the pairwise posterior probabilities matrix for the given pair of sequences and parameters, as calculated by the customized `probA`.

3 Module `clusteralign.tests`

Unit tests for the different parts of the `clusteralign` module, based on the Python `unittest` framework. The test suite can be run using a compatible test runner, for example using the `nose` test runner by running the `nosetests` executable in the `msarc` directory.

```
$ nosetests
```

4 Customized `probA`

The `probA` source was modified to ease interaction with the Python code and optimize the computations. Some small modifications were needed to allow the construction of consecutive alignments with different parameters (specifically the temperature parameter). The array of posterior probabilities is allocated as a contiguous block so that it can be exposed directly as a `numpy` array. The original `probA` source code ordered the input sequences by length, which was not required and disabled. The code for computing the forwards and backwards partition functions, and the posterior probabilities matrix was changed to minimize the amount of multiplications and exponent calculations performed.

Furthermore, the `probA` program has a notion of a canonical unambiguous alignment¹ that is detrimental to the multiple sequence alignment problem. Consider the following alignments

<pre> A - - - X X X X B A Y Y Y - - - B </pre>	and	<pre> A X X X X - - - B A - - - - Y Y Y B </pre>
--	-----	--

The `probA` definition of a canonical alignment excludes the second alternative which leads to this modified

¹Mückstein, U., Hofacker, I.L., Stadler, P.F.: *Stochastic pairwise alignments*. *Bioinformatics* 18 Suppl 2, S153–S160 (2002)

recursion for the Gotoh algorithm

$$M_{i,j} = \max \begin{cases} M_{i-1,j-1} + s(a_i, b_j) \\ E_{i-1,j-1} + s(a_i, b_j) \\ F_{i-1,j-1} + s(a_i, b_j) \end{cases} \quad (1)$$

$$E_{i,j} = \max \begin{cases} M_{i,j-1} + g_o \\ E_{i,j-1} + g_{ext} \end{cases} \quad (2)$$

$$F_{i,j} = \max \begin{cases} M_{i-1,j} + g_o \\ E_{i-1,j} + g_o \\ F_{i-1,j} + g_{ext} \end{cases} \quad (3)$$

And this modified recursion for the partition function

$$Z_{i,j}^M = (Z_{i-1,j-1}^M + Z_{i-1,j-1}^E + Z_{i-1,j-1}^F) e^{\beta s(a_i, b_j)} \quad (4)$$

$$Z_{i,j}^E = Z_{i,j-1}^M e^{\beta g_o} + Z_{i,j-1}^E e^{\beta g_{ext}} \quad (5)$$

$$Z_{i,j}^F = (Z_{i-1,j}^M + Z_{i-1,j}^E) e^{\beta g_o} + Z_{i-1,j}^F e^{\beta g_{ext}} \quad (6)$$

$$Z_{i,j} = Z_{i,j}^M + Z_{i,j}^E + Z_{i,j}^F \quad (7)$$

The probA source code was modified to use a symmetrical recursion. Either the full recursion by replacing equation 5 with

$$Z_{i,j}^E = (Z_{i,j-1}^M + Z_{i,j-1}^F) e^{\beta g_o} + Z_{i,j-1}^E e^{\beta g_{ext}} \quad (8)$$

and equation 2 with

$$E_{i,j} = \max \begin{cases} M_{i,j-1} + g_o \\ E_{i,j-1} + g_{ext} \\ F_{i,j-1} + g_o \end{cases} \quad (9)$$

or the restricted recursion by replacing equation 6 with

$$Z_{i,j}^F = Z_{i-1,j}^M e^{\beta g_o} + Z_{i-1,j}^F e^{\beta g_{ext}} \quad (10)$$

and equation 3 with the symmetrical pair of equation 2

$$F_{i,j} = \max \begin{cases} M_{i-1,j} + g_o \\ F_{i-1,j} + g_{ext} \end{cases} \quad (11)$$

Source modifications in *diff* format can be found in the patches sub-directory of the included probA source tree.