

SAINTq: scoring protein-protein interactions in affinity purification - mass spectrometry experiments with fragment or peptide intensity data

Guo Ci Teo and Hyungwon Choi

February 27, 2016

1 Compilation

Run “make” to compile the software. g++ version 4.4 or above is required to compile the source code. When the compilation finishes, the executable will appear in the “bin” directory.

2 Input files

Two input files will be used by the program for the analysis. The first is a table of intensity measurements, and the second is a parameter file for the input data information and optional parameters for inclusion criteria in scoring. Example input files can be found in the “examples” directory (MSPLIT/PeakView data from the SAINTq manuscript, Teo *et al*, 2015).

2.1 Parameter file

In the parameter file, the following information should be provided by the user.

- `input_filename`: the name of the file containing intensity data
- `normalize_control`: normalize control intensities by multiplying a constant to all control intensities so that the average observed test intensities is equal to the average observed control intensities. “true” or “false”.
- `input_level`: the type of intensity data (valid entries are “protein”, “peptide” or “fragment” and they are case sensitive).

- 24 • `protein_colname`: the header in the column of protein names in the intensity
25 table
- 26 • `pep_colname`: the header of the column of peptide names (valid for fragment
27 or peptide intensities only).
- 28 • `frag_colname`: the header of the column of fragment names (valid for fragment
29 intensity data only).
- 30 • `compress_n_ctrl`: the number of control baits used in calculations, with priority
31 for baits with greater intensities. Setting this number to a large number makes
32 the program use all available control data (recommended in cases with at most
33 several controls).
- 34 • `compress_n_rep`: the number of test bait replicates used for scoring, with pri-
35 ority given to the baits with higher probability scores. If this number is greater
36 than or equal to the number of available replicates, then the scores will use the
37 data from all replicates. Otherwise, the highest scoring replicate scores will be
38 averaged to yield the final probability score.
- 39 • `best_prop_pep`: the proportion of peptides to be used for protein score calcu-
40 lation. Default is 0.5 (50%).
- 41 • `min_n_pep`: a minimum number of peptide intensities to be used for protein
42 score calculation. This sets a lower bound on the number of peptides used when
43 the “`best_prop_pep`” parameter selects too few peptides. Default is 3.
- 44 • `best_prop_frag`: the proportion of fragments to be used for peptide score cal-
45 culation. Default is 0.5 (50%).
- 46 • `min_n_frag`: a minimum number of fragment intensities to be used for peptide
47 score calculation. This sets a lower bound on the number of fragments used
48 when the “`best_prop_frag`” parameter selects too few fragments. Default is 3.

49 Note that `compress_n_ctrl`, `compress_n_rep`, `min_n_pep`, `best_prop_pep`, `min_n_frag`
50 and `best_prop_frag` control the number of intensity measurements utilized for the
51 scoring, with priority for higher intensities. The program will not execute when unused
52 options are set. For example, a protein level data should not have `pep_colname` and
53 a peptide level data should not have `frag_colname`. Option lines can be commented
54 out by putting a pound symbol (`#`) in the beginning of the line.

55 2.2 Intensity table file

56 This file should be a tab-separated values file, with the first three lines describing the
57 bait and experimental information. The third line is a row of column names containing
58 names of protein, peptide, fragment columns and the bait replicate IP names. Directly
59 above the bait IP names, in the second line, are the corresponding bait names. The first
60 line indicates whether each bait is a test protein (T) or control (C). Missing intensities
61 can be represented with "0" or an empty character. Bait replicate IP names of the
62 same bait names must be placed next to one another.

63 3 Running the program

64 In a command line prompt, execute the program with the parameter filename:

```
65 $ ./saintq input_parameter_file
```

66 4 Output file format

67 The output file is a tab-separated value file, with each row corresponding to a bait-prey
68 interaction pair.

- 69 • The first two columns "Bait" and "Prey" are bait and prey proteins, respectively.
- 70 • #Rep: the number of bait replicates.
- 71 • #Pep: the number of peptides in each prey protein. This appears only for
72 peptide and fragment level inputs.
- 73 • #Frag: the numbers of fragments available for each peptide in each prey protein.
74 The numbers are separated with a bar "|" character. This appears only for
75 fragment level inputs.
- 76 • AvgP: probability score for a bait-prey interaction.
- 77 • BFDR: the Bayesian false discovery rate (FDR).

78 5 Copyrights

79 Copyright (C) 2015 Guo Ci Teo <ci@nus.edu.sg> and Hyungwon Choi <hyung_won_choi@nuhs.edu.sg>
80 National University of Singapore.

81 SAINTq is free software: you can redistribute it and/or modify it under the terms of
82 the GNU General Public License as published by the Free Software Foundation, either
83 version 3 of the License, or (at your option) any later version.

84 SAINTq is distributed in the hope that it will be useful, but WITHOUT ANY
85 WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS
86 FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

87 You should have received a copy of the GNU General Public License along with
88 SAINTq. If not, see <<http://www.gnu.org/licenses/>>.