

TfReg user manual

TfReg Version 1.0 user manual
Gerald Weber gweberbh@gmail.com
Departamento de Física, Universidade Federal de Minas Gerais, Brazil
January 31, 2013

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/>.



CONTENTS

| | | |
|----------|---|-----------|
| 1 | Introduction | 3 |
| 2 | Installation | 4 |
| 2.1 | Compiling the source files | 4 |
| 2.1.1 | If you don't have a Linux system | 4 |
| 2.1.2 | License | 4 |
| 3 | Usage | 5 |
| 3.1 | Arguments | 5 |
| 3.1.1 | -o=<basename> | 5 |
| 3.1.2 | -reg=<filename> | 5 |
| 3.1.3 | -par=<filename> | 5 |
| 3.1.4 | -data=<filename> | 5 |
| 3.1.5 | -matrix=<new directory> | 5 |
| 3.1.6 | -res=<resulttype> | 5 |
| 3.1.7 | -model=<model acronym> | 6 |
| 3.1.8 | -duplex=<DNA or RNA> | 6 |
| 3.1.9 | -expand=<nearest neighbours> | 6 |
| 3.1.10 | -cutoff=<integer number> | 6 |
| 3.1.11 | -pbc=<0 or 1> | 6 |
| 3.1.12 | -t=<temperature> | 7 |
| 3.1.13 | -int=<range> | 7 |
| 3.1.14 | -ee=<experimental error> | 7 |
| 3.1.15 | -rs=<seed> | 7 |
| 3.1.16 | -mlr=<sequence length> | 7 |
| 3.1.17 | -pm=<prediction method number> | 7 |
| 3.1.18 | -seq=<nucleotide sequence> | 8 |
| 3.1.19 | -cseq=<nucleotide sequence> | 8 |
| 3.1.20 | -salt=<salt concentration> | 8 |
| 4 | Models | 9 |
| 4.1 | Peyrard-Bishop (-model=pb) | 9 |
| 4.2 | The Dauxois variant (-model=dpb) | 9 |
| 4.3 | PB model with added solvent potential (-model=hms) | 10 |
| 4.4 | The Joyeux and Buyukdagli model (-model=jb) | 10 |
| 5 | Parameter specification | 12 |
| 5.1 | Generic * parameters | 12 |
| 5.2 | BP parameters | 12 |
| 5.3 | NN parameters | 12 |
| 5.4 | Parameter precedence | 13 |
| 6 | Result files | 14 |
| 6.1 | .reg regression parameters | 14 |
| 6.2 | .dat melting temperatures and melting index results | 14 |
| 6.3 | .dat average opening if used with -res=averagey | 15 |
| 6.4 | .ver quality of the prediction | 15 |
| 6.5 | .tex convenience L ^A T _E X file | 15 |
| 6.6 | Matrix files if used with -matrix= | 15 |

| | | |
|----------|---|-----------|
| 7 | Examples | 16 |
| 7.1 | Task: given a set of melting temperatures, find the regression parameters | 16 |
| 7.1.1 | Results of Ref. 1 | 16 |
| 7.1.2 | Results of Ref. 2 | 16 |
| 7.2 | Task: prediction of DNA melting temperatures | 17 |
| 7.2.1 | Single sequence example | 17 |
| 7.2.2 | Multiple sequence example | 17 |
| 7.3 | Task: prediction of DNA melting temperatures with different salt concentrations | 18 |
| 7.4 | Task: prediction of RNA melting temperatures | 18 |
| 7.5 | Task: calculating the average opening $\langle y \rangle$ | 18 |
| 7.5.1 | Calculating $\langle y \rangle$ at a given temperature | 19 |
| 7.5.2 | Calculating $\langle y \rangle$ for a range of temperatures | 19 |
| 7.5.3 | Using Perl scripts and parallel processing | 19 |

1 INTRODUCTION

TfReg implements the calculation of Peyrard-Bishop [3] style Hamiltonians to obtain some physical properties of DNA and RNA duplexes. The method uses the transfer matrix technique for the calculation of the classical partition function. Also, TfReg calculates the regression of experimental versus predicted melting temperatures using the equivalent melting index [4].

What will this software do for you? Given a set of experimental melting temperatures and a set of model parameters you will be able to calculate the regression parameters which will allow you predict melting temperatures of any DNA or RNA sequences. Alternatively, you may use one of the calculated regression sets which are provided and start calculating melting temperatures right away. You may chose between four different “flavours” of Hamiltonians if you wish to investigate the effect of different model parameters. If you do have basic programming skills in C++ you should be able to add new types of model Hamiltonians, as long as they fit within the 1D framework of the original Peyrard-Bishop model [3].

Evidently, this is work in progress. I hope to add further programs in the near future as well as increase the number of parameters for other types of oligonucleotide.

I would find it truly helpful indeed if you would let me know if this software is of any use to you. Showing a list of interested users to funding agencies often helps to secure the necessary resources to keeping such projects running. So, please, if you find this software useful let me know and if you use it for your scientific work please cite the appropriate papers which are listed at the end of this manual.

I wish you all the best in using TfReg

Gerald

Belo Horizonte, January 31, 2013

2 INSTALLATION

The easiest way to install is to visit the OpenSuSe Build service and find TfReg for your Linux distribution, this will take care of obtaining the correct libraries which are needed for TfReg to run. Typically, this will install the binary executable file

`/usr/bin/tfreg`

and further files, such as model parameter files, pre-calculated regression parameters and example scripts are to be found in

`/usr/share/TfReg` or `/usr/share/tfreg`

the documentation (which you are reading right now) should be located at

`/usr/share/doc/packages/TfReg` or `/usr/share/doc/packages/tfreg`

If you are unable to find the packages for your specific Linux distribution, please read on.

2.1 Compiling the source files

This software was developed and tested under Opensuse Linux 12.2 and depends on some libraries to function properly:

1. `libboost_filesystem1.49.0` `libboost_system1.49.0` <http://www.boost.org>
2. `gsl` <http://www.gnu.org/software/gsl/>
3. `libblas3` <http://www.netlib.org/lapack/>
4. `libgfortran47` <http://gcc.gnu.org/>

which means that you will need at least these specialized packages in addition to the usual gcc and g++ compiler. After unpacking the tar package you should try to compile using the `make` command. If the compilation is successful, copy the `tfreg` binary file to your installation directory.

2.1.1 If you don't have a Linux system

You may copy an Opensuse Live system on a CD or onto a USB memory stick and reboot your computer with this system. If you do, proceed as you were using a normal Opensuse 12.2 system, following the instructions above. See installation instructions for Live Opensuse at http://en.opensuse.org/Live_USB_stick

2.1.2 License

This software is published under the GNU General Public License version 3 (GPLv3), the complete text of this license can be found in the documentation folder. If you wish to use this software under a different license or wish to make changes to the software without distributing it under the GPLv3 please contact me so that we can arrange for a specific license.

3 USAGE

TfReg take all its program option from the command line, therefore you should invoke the binary `tfreg` with some of the arguments which are described below. Example shell scripts are provided and I recommend you study them as they are the best source to illustrate how to use TfReg. The following section details every available program option for TfReg.

3.1 Arguments

3.1.1 `-o=<basename>`

specifies the output basename, that is, all files which are generated start with `basename`

3.1.2 `-reg=<filename>`

specifies the input regression file name.

3.1.3 `-par=<filename>`

specifies the input parameter file name. This can be a list of files `-par=file1.par,file2.par,file3.par`. In case of multiple specifications of the same model parameter, the last one supersedes previous parameters. For example if `file1.par` has the parameter `AT:Morse.D 0.05` and `file3.par` has `AT:Morse.D 0.03`, the final value for `AT:Morse.D` will be 0.03.

3.1.4 `-data=<filename>`

specifies the input file containing nucleotide sequences and melting temperatures. TfReg was not created with very large datasets in mind, and all sequences are loaded into memory. A safe limit is of the order of 50000 sequences which will require about 1GB of memory to run.

The file format is composed of columns

```
1 temperature
2 ATCAATCATA TAGTTAGTAT 21.3 69 2
3 TTGTAGTCAT AACATCAGTA 24.7 69 2
4 GAAATGAAAG CTTTACTTTC 22.1 69 2
```

data/owczarzy04-69.dat

The first column is the sequence (from 5' → 3') and the second column is the secondary strand (from 3' → 5'). The third column is the melting temperature in °C, the fourth column is the salt concentration (in mM). The last column the species concentration (in μM) which is currently not used. The secondary strand does not necessarily need to be the complementary of the main strand as long as there are parameters for the mismatched pairs. If TfReg fails to find parameters for your sequences it will complain loudly.

3.1.5 `-matrix=<new directory>`

if specified this will create a directory where the calculated matrices will be located. Once calculated these matrices will be reloaded the next time the program is run again. This speeds up the calculations, however if you do change parameters or the model from one run to the next you should not use this option as you will load matrices which are incorrect.

3.1.6 `-res=<resulttype>`

which type of result we wish to obtain.

(default) `-res=regression` given a set of parameters and a set of melting temperatures, calculates the regression parameters.

`-res=prediction` predict temperatures for a given set of parameters and a given regression file.

`-res=nncheck` only checks consistency of the nearest-neighbour decomposition of the nucleotide sequences. This will display and exit immediately without doing any calculations.

`-res=averagey` calculates the average opening $\langle y \rangle$, results are given in Ångström. Note that for short sequences the temperature needs to be unrealistically low, see example 7.5.

3.1.7 `-model=<model acronym>`

Selects the Peyrard-Bishop model which should be used. Each model requires specific parameters which should be passed via the `-par=<filename>` option.

(default) `-model=pb` original Peyrard-Bishop model [3]

`-model=dpb` the anharmonic variant [5]

`-model=hms` PB model with added solvent potential [6]

`-model=jb` finite enthalpy model [2, 7–10]

3.1.8 `-duplextype=<DNA or RNA>`

Selects the type of duplex we should expect, this is important for selecting the base pair complementarity. Note that IUPAC codes cannot be used since we need to know which nucleotide parameters to use, that is, we cannot use N for example since we would not know which potentials to use.

(default) `-duplextype=DNA` , will expect A, C, G and T base pairs and will consider A complementary to T and C complementary to G.

`-duplextype=RNA` , same as for DNA but considers A and U as complementaries [11].

3.1.9 `-expand=<nearest neighbours>`

Selects which nearest neighbours to expand.

(default) `-expand=CG.CG`

3.1.10 `-cutoff=<integer number>`

This controls the truncation P in Eq. (22) of Ref. 12.

(default) `-cutoff=0` no cutoff.

(recommended) `-cutoff=10` a cutoff $P = 10$ gives quite good results and reduces the computational cost.

3.1.11 `-pbc=<0 or 1>`

Controls the type of boundary conditions.

(default) `-pbc=0` open boundary conditions, this is what you normally would have and usually shows end fraying

`-pbc=1` periodic boundary conditions, this would be the case for a circular sequence.

3.1.12 -t=<temperature>

Selects the temperature in kelvin for which the calculation of the matrices is carried out. Please note: this temperature is completely unrelated to the melting temperatures.

(default) -t=370

3.1.13 -int=<range>

Range of integration, this specifies the limits of the integral shown in Eq. (14) of Ref. 12 and the size M of the matrices.

(default) -int=-1:30/100 integrates from $y = -1$ to $y = 30$ Å and uses matrices of size 100

(recommended) -int=-1:200/400 integrates from $y = -1$ to $y = 200$ Å and uses matrices of size 400. This gives very accurate results without too much computational cost.

3.1.14 -ee=<experimental error>

This sets the experimental error (in °C) of the melting temperature set. If given, the data set (provided through the -data option) will be modified by small positive or negative amounts such that the standard deviation is close to the experimental error.

(example) -ee=0.5 will modify the dataset to within 0.5 °C

3.1.15 -rs=<seed>

Sets the seed of the random number generator (C function srand) which is used to modify the dataset. If you use the same seed you will get exactly the same random modifications. Only make sense to be used together with -ee.

3.1.16 -mlr=<sequence length>

Sets the minimal sequence length to be considered in calculating regression parameters.

(example) -mlr=6 only considers sequences with 6 bp and above

3.1.17 -pm=<prediction method number>

Selects the prediction/regression method.

(default) -pm=2 this considers two equations for regression

$$T_p = a_0(N, [\text{Na}^+]) + a_1(N, [\text{Na}^+])\tau, \quad (3.1)$$

there will be one equation for each length N and for each salt concentration $[\text{Na}^+]$ and each coefficient a_k is calculated as

$$a_k(N, [\text{Na}^+]) = b_{0,k}([\text{Na}^+]) + b_{1,k}([\text{Na}^+])N^{1/2}, \quad (3.2)$$

where N is the size of the sequence and τ is the equivalence index,

-pm=3 in addition of the two previous equations, this considers a third equation for the case where there are two or more salt concentrations in the data file

$$b_{j,k}([\text{Na}^+]) = c_{0,j,k} + c_{1,j,k} \log[\text{Na}^+]. \quad (3.3)$$

Note that if there is only one salt concentration $[\text{Na}^+]$ in the data file this reverts automatically to option -pm=2.

-pm=-1 this considers only one regression

$$T_p = a_0 + a_1\tau, \quad (3.4)$$

which is useful if the dataset has only a few sequences, or if all sequences are of the same length or if you have otherwise trouble in getting good linear regression coefficients for Eq. (3.2). The last situation may happen if your model parameters are very far from the optimized values. This option was first introduced for Ref. [11].

3.1.18 -seq=<nucleotide sequence>

Instead of providing a file with your sequences you can give them on the command line. This is useful if you want to see a melting for just one sequence. You should give the main strand from 5' to 3', the complementary sequence will be worked out automatically.

(example) -seq=ACGTTGAATT

3.1.19 -cseq=<nucleotide sequence>

You should provide a 3' → 5' sequence if your sequence is not perfectly complementary, say like in a sequence with nucleotide mismatches.

(example) -cseq=TGCTACTTAA from 3' → 5'

3.1.20 -salt=<salt concentration>

In the case where a -reg file contains several salt concentration, this will select the one that should be used for the calculations. Otherwise, this option is silently ignored. *If you need to calculate melting temperatures for different salt concentrations please see section 7.3.*

4 MODELS

There are several variants of the PB model, each of which requires different model parameters. In this section you will find which models TfReg currently supports and which model parameters are needed. Example files with model parameters are provided and typically you will find these in `/usr/share/TfReg/data` (or `url/usr/share/tfreg/data`) with file extension `.par`. Note that the parameters are specified per nucleotide type and follow a very specific notation described in section 5.

4.1 Peyrard-Bishop (`-model=pb`)

This is the original Peyrard-Bishop model proposed in Ref. [3] which uses a Morse potential for modelling the hydrogen bonds and

$$V_{\text{Morse}}(y_i) = D \left(e^{-y_i/\lambda} - 1 \right)^2, \quad (4.1)$$

and the nearest-neighbour stacking interaction as a harmonic oscillator

$$w_{\text{harm.}}(y_i, y_{i-1}) = \frac{k}{2} (y_i - y_{i-1})^2. \quad (4.2)$$

because of the divergence of the partition function [13] we modified this to

$$w_{\text{harm.}}(y_i, y_{i-1}) = \frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2), \quad (4.3)$$

| Hamiltonian | model parameter | program parameter | units | type |
|--|-----------------|-----------------------------|-------------------|------|
| $D (e^{-y_i/\lambda} - 1)^2$ | D | <code>Morse.D</code> | eV | BP |
| | λ | <code>Morse.lambda</code> | Å | BP |
| $\frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2)$ | k | <code>harmonic.k</code> | eV/Å ² | NN |
| | θ | <code>harmonic.theta</code> | rad | NN |

4.2 The Dauxois variant (`-model=dpb`)

In 1993 Dauxois, Peyrard and Bishop introduced an anharmonicity term to account for sharp transitions in the original PB model [5],

$$w_{\text{an.}}(y_i, y_{i-1}) = \left[1 + \rho e^{-\alpha(y_i + y_{i-1})} \right] w_{\text{harm.}}(y_i, y_{i-1}), \quad (4.4)$$

| Hamiltonian | model parameter | program parameter | units | type |
|---|-----------------|-------------------------------|-------------------|------|
| $D (e^{-y_i/\lambda} - 1)^2$ | D | <code>Morse.D</code> | eV | BP |
| | λ | <code>Morse.lambda</code> | Å | BP |
| $\left[1 + \rho e^{-\alpha(y_i + y_{i-1})} \right]$ | k | <code>harmonic.k</code> | eV/Å ² | NN |
| $\times \frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2)$ | θ | <code>harmonic.theta</code> | rad | NN |
| | ρ | <code>anharmonic.rho</code> | adimensional | NN |
| | α | <code>anharmonic.alpha</code> | adimensional | NN |

data/weber06-1.par

```

1 weber06
2 AT:AU:morse.D      0.05
3 CG:morse.D         0.08
4 AT:AU:morse.lambda 0.33333
5 CG:morse.lambda    0.125
6 *:harmonic.theta   0.01
7 *:harmonic.k        0.025
8 *:anharmonic.alpha  0.35
9 *:anharmonic.rho    2.0

```

4.3 PB model with added solvent potential (-model=hms)

A solvent term was added to the harmonic PB model,

$$V(y_i) = V_{\text{Morse}}(y_i) - f_s D [\tanh(y_i/\lambda_s) + 1], \quad (4.5)$$

| Hamiltonian | model parameter | program parameter | units | type |
|--|-----------------|-------------------|-------------------|------|
| $D (e^{-y_i/\lambda} - 1)^2$ | D | Morse.D | eV | BP |
| | λ | Morse.lambda | Å | BP |
| $-f_s D [\tanh[(y_i + y_e)/\lambda_s] + s]$ | f_s | solvent.f_s | adimensional | BP |
| | λ_s | solvent.lambda | Å | BP |
| | y_e | solvent.eq_sol | Å | BP |
| | s | solvent.sign_sol | adimensional | BP |
| $\frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2)$ | k | harmonic.k | eV/Å ² | NN |
| | θ | harmonic.theta | rad | NN |

data/weber06b-1.par

```

1 weber06b
2 AT:morse.D      0.05
3 CG:morse.D      0.08
4 AT:morse.lambda 0.33333
5 CG:morse.lambda 0.125
6 *:solvent.eq_sol 0.0
7 *:solvent.sign_sol 1.0
8 *:harmonic.theta 0.01
9 *:harmonic.k      0.025
10 *:solvent.lambda 1.0
11 *:solvent.f_s     0.1

```

4.4 The Joyeux and Buyukdagli model (-model=jb)

The model by Joyeux and Buyukdagli [7–10] introduces a finite stacking enthalphy

$$w_{\text{fin.}}(y_i, y_{i-1}) = \frac{\Delta H}{C} [1 - e^{-b(y_i - y_{i-1})^2}] + \frac{K_b}{2} (y_i - y_{i-1})^2. \quad (4.6)$$

data/jb.par

```

1 var_jb_owczarzy04_init2
2 AT:morse.D      0.041
3 CG:morse.D      0.054
4 AT:morse.lambda 0.1667
5 CG:morse.lambda 0.1667
6 *:finite_enthalpy.C 4.0
7 *:finite_enthalpy.DeltaH 0.409
8 *:finite_enthalpy.b 0.80
9 *:finite_enthalpy.kb 4.0e-4

```

| Hamiltonian | model parameter | program parameter | units | type |
|---|-----------------|------------------------|-------------------|------|
| $D (e^{-y_i/\lambda} - 1)^2$ | D | Morse.D | eV | BP |
| | λ | Morse.lambda | Å | BP |
| $\frac{k}{2} (y_i - y_{i-1})^2 s$ | k | harmonic.k | eV/Å ² | NN |
| $\frac{\Delta H}{C} \left[1 - e^{-b(y_i - y_{i-1})^2} \right]$ | ΔH | finite_enthalpy.DeltaH | eV | NN |
| | C | finite_enthalpy.C | adimensional | NN |
| | b | finite_enthalpy.b | adimensional | NN |
| $\frac{K_b}{2} (y_i - y_{i-1})^2$ | K_b | finite_enthalpy.Kb | eV/Å ² | NN |

5 PARAMETER SPECIFICATION

This chapter explains how to write a parameter file.

Let's start with the most simple possible parameter file

```
1 initial
2 *:harmonic.theta 0.01
3 *:harmonic.k 0.025
4 *:morse.D 0.0324083
5 *:morse.lambda 0.350149
```

The first line holds a simple identifier, this identifier can be letters, numbers and some other characters but should contain no spaces. This identifier is used when you specify the `-matrix` option.

The second line starts with an asterisk which means that this parameter should apply to any base pair. In case of a DNA sequence this will apply to AT and CG base-pairs. TA base pairs are considered equivalent to AT, and CG are equivalent to GC, from the perspective of the models considered in this manual.

5.1 Generic * parameters

The generic base * specification can be used when a given parameter should be applied to any base pair (see previous example). However, if the parameter should be applied to a specific base pair you should specify either in BP or NN form (see next section).

5.2 BP parameters

Base-pair parameters are those which do not depend on the context, that is, it is not relevant which are the neighbours of the given base pair. For DNA we use Watson-Crick base pairs AT=TA and CG=GC, but may also use mismatched base pairs such as GT or AA. Switching the nucleotides, eg. GT to TG, makes no difference for BP parameters. The reason for this is that they usually represent properties of the hydrogen bond.

5.3 NN parameters

Nearest neighbour (NN) parameters follow the same conventions as usually found in linear regression models [14, 15]. The convention used is of type AB_CD, where AB is the first base pair and CD the second base pair. Note that DC_BA is equivalent to AB_CD, for example, AT_AT is the same as TA_TA. When specifying a NN sequence always specify in lexical ordering, that is write AT_AT and not TA_TA.

The following example shows all 10 irreducible N parameters for `finite_enthalpy.DeltaH`.

```
Examples of NN parameters
1 AT_AT:finite_enthalpy.DeltaH 0.42085
2 AT_CG:finite_enthalpy.DeltaH 0.416718
3 AT_GC:finite_enthalpy.DeltaH 0.400366
4 AT_TA:finite_enthalpy.DeltaH 0.330006
5 CG_AT:finite_enthalpy.DeltaH 0.445258
6 CG_CG:finite_enthalpy.DeltaH 0.385034
7 CG_GC:finite_enthalpy.DeltaH 0.411902
8 GC_AT:finite_enthalpy.DeltaH 0.44666
9 GC_CG:finite_enthalpy.DeltaH 0.497106
10 TA_AT:finite_enthalpy.DeltaH 0.431154
```

5.4 Parameter precedence

Since the program can read more than one parameter file, the last parameter read is the final value. There may also be several specification in the same file as well. Consider the following example

Precedence example 1

```
1 *:Morse.D 0.3
2 AT:Morse.D 0.5
```

the first line says that all base pairs should a **Morse.D** value of 0.3. The second line however says that AT base pairs should use 0.5. In this case CG base pairs for instance will use 0.3 since nothing different was specified.

However, a generic base pair ***** does not supersedes a specific base pair as in the following example

Precedence example 2

```
1 AT:Morse.D 0.5
2 *:Morse.D 0.3
```

in this case AT base pairs will continue using 0.5, not 0.3. You should understand the generic base pair ***** as: *if nothing else matches, use this value.*

6 RESULT FILES

TfReg produces many result files, here we will make an attempt to describe them in some detail. Some result files can be used again as input file such as the file with extension `.reg`.

6.1 `.reg` regression parameters

The regression parameters will be stored in a file with extension `.reg`

`examples/verify/epl2011-69.reg`

```

1 0.5 0.5 2
2 b
3 69 -86.7785 48.7698 15.3185 -6.02284
4 a
5 69
6 10 -39.4819 30.0373
7 15 -26.5942 25.2957
8 20 -16.6942 21.2883
9 25 -10.8334 18.7827
10 30 -3.51713 16.0349

```

Lets start from line 4 which contains the single letter `a`, this flags the start of all regression parameters which are length and salt concentration dependent. Line 5 holds the first salt concentration 69 mM in this case. The first regression equation starts at line 6 for all sequences of length 10 bp, and the next two numbers are the a_0 and a_1 coefficients of Eq. (3.1), that is

$$T_p(N = 10, [\text{Na}^+] = 69) = -39.4819 + 30.0373\tau, \quad (6.1)$$

At line 2 we see the letter `b` which flags the start of the regression parameters which are length dependent. There will be one line for each salt concentration. At line 3 we see the coefficients for 69 mM corresponding to $b_{0,0}$, $b_{1,0}$, $b_{1,0}$ and $b_{1,1}$ of Eq. (3.2),

$$a_0([\text{Na}^+] = 69) = -86.7785 + 48.7698N^{1/2}, \quad (6.2)$$

$$a_1([\text{Na}^+] = 69) = 15.3185 - 6.02284N^{1/2}. \quad (6.3)$$

these equations are necessary to calculate the a_0 and a_1 coefficients for sequence length which are not in the regression file.

6.2 `.dat` melting temperatures and melting index results

One of the main result files has the extension `.dat` and typically contains the melting temperatures and melting index calculated for each sequence. Here is an example (output of 7.2.2).

`examples/predict-2.dat`

```

1 Main/Complementary alpha salt_concentration species_concentration temperature.measured temperature.adjusted
2 TACTAACATTAATA/ATGATTGTAATTGAT 4 69 2 0 0 34.2215 0 0 0 0 0 0 2 6.006 2.45071
3 ATACTTACTGATTAG/TATGAATGACTAATC 4 69 2 0 0 35.8264 0 0 0 0 0 0 2 6.318 2.51356
4 GTACACTGTCTTATA/CATGTGACAGAATAT 4 69 2 0 0 39.7114 0 0 0 0 0 0 2 7.106 2.66571
5 StdTempDiff 0
6 StdEnthalpyDiff 0
7 StdEntropyDiff 0

```

The first line (which is too long to be shown completely) identifies each column. Importantly, the 7th column is the predicted temperature and the last column is the melting index τ .

6.3 .dat average opening if used with `-res=averagey`

```
examples/open-1.dat
1 0 1.30564
2 1 0.872495
3 2 0.844811
4 3 0.587454
5 4 0.631239
```

Shown are the first few lines of the result from example 7.5. The first column is the base pair position and the second column is the average opening $\langle y \rangle$, results are given in Ångstrom.

6.4 .ver quality of the prediction

The file with extension `.ver` is a short file intended to show how close the predicted melting temperatures are when compared to the experimental melting temperatures. This file only make sense if a data file with the experimental melting temperatures was given.

```
examples/np2009-1-69.ver
1 0.811548 0.771974 115.419 prediction method=2
```

The first column is average deviation of predicted and experimental melting temperatures $\langle \Delta T \rangle$, the second column is the standard deviation between predicted and experimental melting temperatures and the third is χ^2 .

6.5 .tex convenience L^AT_EX file

The file with extension `.tex` provides convenient access to some results as L^AT_EX commands.¹

```
examples/np2009-1-69.tex
1 \Action{ActionIndex}
2 \Model{pb}
3 \SqrDiff{115.419}
4 \AvDiff{0.811548}
5 \StdDiff{0.771974}
6 \StdTempDiff{0}
7 \StdEnthalpyDiff{0}
8 \StdEntropyDiff{0}
```

6.6 Matrix files if used with `-matrix=`

If you specify a matrix directory say, `-matrix=open-1`, the program will create files containing all matrices used in the calculation:

```
Matrix files
1 dna-pb-69-CG_CG-270-A.AT_AT dna-pb-69-CG_CG-270-A.AT_CG dna-pb-69-CG_CG-270-A.AT_GC
2 dna-pb-69-CG_CG-270-A.AT_TA dna-pb-69-CG_CG-270-A.CG_AT dna-pb-69-CG_CG-270-A.CG_CG
3 dna-pb-69-CG_CG-270-A.CG_GC dna-pb-69-CG_CG-270-A.GC_AT dna-pb-69-CG_CG-270-A.GC_CG
4 dna-pb-69-CG_CG-270-A.TA_AT dna-pb-69-CG_CG-270-C.AT_AT dna-pb-69-CG_CG-270-C.AT_CG
5 dna-pb-69-CG_CG-270-C.AT_GC dna-pb-69-CG_CG-270-C.AT_TA dna-pb-69-CG_CG-270-C.CG_AT
6 dna-pb-69-CG_CG-270-C.CG_CG dna-pb-69-CG_CG-270-C.CG_GC dna-pb-69-CG_CG-270-C.GC_AT
7 dna-pb-69-CG_CG-270-C.GC_CG dna-pb-69-CG_CG-270-C.TA_AT dna-pb-69-CG_CG-270-eigenvalues
8 dna-pb-69-CG_CG-270-eigenvectors dna-pb-69-CG_CG-270-Y
```

The first part of the matrix file name `dna-pb-69` is the parameter identifier which you will find in the first line of the file `dna-pb-69.par`. The next field `CG_CG` refers to the base pair used as expansion basis (see option `-expand`). Next comes the temperature at which the calculation was performed (in this case 270 K). The matrix type is given by the last letter or identification and can be either A [Eq. (64) of Ref. 13] or C [Eq. (56) of Ref. 13]. The extensions like `AT_AT` refer to the nearest neighbour base pairs, there will be one for each type of nearest neighbour present in the sequence. The matrices labeled `eigenvalues` and `eigenvectors` are self-explanatory and represent the result of the diagonalization. Finally the matrix ending in `Y` is that of Eq. (60) of Ref. 13.

¹This is no longer in use and may be removed in future versions.

7 EXAMPLES

The example scripts are located in `/usr/share/TfReg/examples` (or `/usr/share/tfreg/examples`). In `/usr/share/TfReg/examples/verify` (or `/usr/share/tfreg/examples/verify`) you will find the output of some of these script which will allow you to check if your installed version of TfReg is working properly.

7.1 Task: given a set of melting temperatures, find the regression parameters

This scenario appears when you have some parameters for the hydrogen bond and perhaps stacking interaction and you wish to know how close these may get to experimental melting temperatures. First you will need to calculate the regression parameters which you may later use to calculate melting temperatures for untested sequences.

What you will need:

1. A set of experimental melting temperatures (for example `data/owczarzy04-69.dat`)
2. A set of parameters (for example `dna_pb_69.par`)

7.1.1 Results of Ref. 1

In this example we take the parameters which were calculated by the minimization procedure, stored in file `dna_pb_X.par` (X being 69, 119, 220, 621 or 1020), and calculate the regression parameters which are going to be stored in file `np2009-X`.

```
8      examples/np2009-1.sh
9      export COMMON="-cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=regression -t=370 -v=1"
10     tfreg ${COMMON} -o=np2009-1-69 \
11     -par=$PREFIX/dna_pb_69.par -data=$PREFIX/owczarzy04-69.dat
12
13     tfreg ${COMMON} -o=np2009-1-119 \
14     -par=$PREFIX/dna_pb_119.par -data=$PREFIX/owczarzy04-119.dat
15
16     tfreg ${COMMON} -o=np2009-1-220 \
17     -par=$PREFIX/dna_pb_220.par -data=$PREFIX/owczarzy04-220.dat
18
19     tfreg ${COMMON} -o=np2009-1-621 \
20     -par=$PREFIX/dna_pb_621.par -data=$PREFIX/owczarzy04-621.dat
21
22     tfreg ${COMMON} -o=np2009-1-1020 \
23     -par=$PREFIX/dna_pb_1020.par -data=$PREFIX/owczarzy04-1020.dat
24
```

7.1.2 Results of Ref. 2

In this example we take the parameters which were calculated by the minimization procedure, stored in file `stat_jb_owczarzy04-69.par`, and calculate the regression parameters which are going to be stored in file `ep12011-69`. The file `var_jb2_owczarzy04_init.par` contains the initial parameters which were used in the minimization procedure, most parameters will be superseded by the contents of `stat_jb_owczarzy04-69.par`.

```
8      examples/ep12011-69.sh
9      tfreg -cutoff=10 -int=-1:200/400 -m=jb -o=ep12011-69 \
10     -par=$PREFIX/dna_jb_69.par -data=$PREFIX/owczarzy04-69.dat \
    -pbc=0 -pm=2 -res=regression -t=370 -v=1
```

7.2 Task: prediction of DNA melting temperatures

You would like to predict melting temperatures of some nucleotide sequence for which you have no experimental data.

What you will need:

1. A set of parameters (for example `data/dna_pb_69.par`)
2. A file with calculated regression parameters (for example `data/np2009-1-69.reg` for salt concentration of 69 mM)

7.2.1 Single sequence example

If you want to predict the melting temperature of just one sequence the easiest is to specify the sequence as a command argument `-seq` as in the following example.

```
7 export COMMON="-cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -t=370 -v=1"
8
9 tfreg ${COMMON} -o=predict-1 \
10 -par=$PREFIX/dna_pb_69.par \
11 -reg=$PREFIX/np2009-1-69.reg -salt=69 \
12 -res=prediction \
13 -seq=ACAGCGAATGGACCTACGTGGCCTT
```

7.2.2 Multiple sequence example

If you want to predict the melting temperature of many sequences it is advisable to edit a simple file like this:

```
1 temperature
2 TACTAACATTA ACTA ATGATTGTAATTGAT 0 69 0
3 ATACTTACTGATTAG TATGAATGACTAATC 0 69 0
4 GTACACTGTCTTATA CATGTGACAGAATAT 0 69 0
```

the 69 refers to the salt concentration at which you want to predict these temperatures. Note that the second column is the complementary sequence of the first column. The way to run this example is as follows:

```
7 export COMMON="-cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=prediction -t=370 -v=1"
8
9 tfreg ${COMMON} -o=predict-2 \
10 -par=$PREFIX/dna_pb_69.par \
11 -reg=$PREFIX/np2009-1-69.reg \
12 -res=prediction \
13 -data=$PREFIX/example2.dat
```

The result file is this:

```
1 Main/Complementary alpha salt_concentration species_concentration temperature.measured temperature.adjusted
2 TACTAACATTA ACTA/ATGATTGTAATTGAT 4 69 2 0 0 34.2215 0 0 0 0 0 0 2 6.006 2.45071
3 ATACTTACTGATTAG/TATGAATGACTAATC 4 69 2 0 0 35.8264 0 0 0 0 0 0 2 6.318 2.51356
4 GTACACTGTCTTATA/CATGTGACAGAATAT 4 69 2 0 0 39.7114 0 0 0 0 0 0 2 7.106 2.66571
5 StdTempDiff 0
6 StdEnthalpyDiff 0
7 StdEntropyDiff 0
```

where the 7th column is the predicted temperature and the last column is the melting index τ .

7.3 Task: prediction of DNA melting temperatures with different salt concentrations

If you need to predict salt concentrations which differ from the ones currently provided you will need to generate a new regression file. To ease this task we provided a Perl script which does this regression for you called `tfreg-salt-regression.pl`. To use it you will need to provide some existing regression files for different salt concentrations which will use the regression equation (3.3). Here is an practical example where we use 5 files containing different salt concentration for which we generate a new regression file for a salt concentration of 50 mM.

```
1  tfreg-salt-regression.pl 50 \
2  "np2009-1-69.reg,np2009-1-119.reg,np2009-1-220.reg,np2009-1-621.reg,np2009-1-1020.reg" \
3  new50.reg
```

the first argument is the new salt concentration, then a list of files comma-separated (the order is unimportant), and the last is the name of the new file. And here is an example script showing how to use the new regression file.

```
7  examples/predict-salt50.sh
8  export COMMON="-cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=3 -t=370 -v=1"
9
10 tfreg ${COMMON} -o=predict-salt50 \
11 -par=$PREFIX/dna_pb_69.par \
12 -reg=$PREFIX/new50.reg -salt=50 \
13 -res=prediction \
14 -seq=ACAGCGAATGGACCTACGTGGCCTT
```

7.4 Task: prediction of RNA melting temperatures

This is very much the same as predicting the melting temperatures for DNA described in the previous sections. The main difference is that you need to use the additional command parameter `-duplextype=RNA`.

What you will need:

1. A set of parameters for RNA (for example `data/rna_pb.par` from Ref. [11])
2. A file with calculated regression parameters (for example `data/reg_pb_xia98-t1.reg` for salt concentration of 1000 mM, this was published as supplementary tables IV and V of Ref. [11])

In the following example we calculate the melting temperatures of all sequences from Ref. 16.

```
8  examples/predict-rna-1.sh
9  export COMMON="-cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=prediction -t=370 -v=1"
10
11 tfreg ${COMMON} -o=predict-rna-1 \
12 -duplextype=RNA \
13 -par=$PREFIX/rna_pb.par \
14 -reg=$PREFIX/reg_pb_xia98-t1.reg \
15 -data=$PREFIX/xia98-t1.dat
```

7.5 Task: calculating the average opening $\langle y \rangle$

Using the option `-res=averagey` you will obtain the average opening $\langle y \rangle$ as a function of nucleotide position (see Ref. 13).

What you will need:

1. A set of parameters (for example `data/dna_pb_69.par`)

7.5.1 Calculating $\langle y \rangle$ at a given temperature

```
examples/open-1.sh
7 export COMMON="-matrix=open-1 -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -t=270 -v=1"
8
9 tfreg ${COMMON} -o=open-1 \
10 -par=$PREFIX/dna_pb_69.par -salt=69 -res=averagey \
11 -seq=ACAGCGAATGGACCTACGTGGCCTT
```

the results will be in the file with extension `.dat`.

7.5.2 Calculating $\langle y \rangle$ for a range of temperatures

Here we use a simple bash loop to loop over the temperatures, each results is stored in a file `open-2-T.dat` where T is the temperature.

```
examples/open-2.sh
7 export COMMON="-cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -v=1 -salt=69 -res=averagey"
8
9 for T in {250..300}
10 do
11 tfreg ${COMMON} -o=open-2-$T -t=$T \
12 -par=$PREFIX/dna_pb_69.par -seq=ACAGCGAATGGACCTACGTGGCCTT
13 done
```

This example takes about 6 min to run (intel i7-2620M 2.70GHz).

7.5.3 Using Perl scripts and parallel processing

Sometimes, doing something more complicated with shell scripts can be very cumbersome. For example, a simple `for` loop with non-integer values can be very complicated. In this case, I would recommend spending a few moments to learn the basics of Perl. In the following example we run the average opening for RNA in increments of 0.5 K. Also each instance is send into background to be processed in a sort of poor man's parallel processing.

```
examples/open-rna-3.pl
1 #!/usr/bin/perl
2 my $prefix='../data';
3 $prefix=$ARGV[0] if (exists $ARGV[0]);
4 my $common='-cutoff=80 -int=-1:30/400 -m=pb -pbc=0 -v=1 -res=averagey';
5
6 for (my $T=280; $T < 320; $T += 0.5)
7 {
8   my $fT=sprintf("%.1f",$T); #here we format $T with 1 decimal place
9   my $com="tfreg $common -o=open-rna-3-$fT -t=$fT -duplextype=RNA -par=$prefix/rna_pb.par";
10   my $seq='-seq=GUGCCCAUUUAGGGUAUAUAUGGCCGAGUGAGCGAGCAGGAUCUCCAUUUUGACCGCAAAUUUGAACG';
11   system("$com $seq" . ' < /dev/null &> open-rna-3-$T.echo &');
12 }
```

BIBLIOGRAPHY

- [1] Weber, G., Essex, J. W., and Neylon, C. Probing the microscopic flexibility of DNA from melting temperatures. *Nature Physics* **5**, 769–773 (2009).
- [2] Weber, G. Finite enthalpy model parameters from DNA melting temperatures. *Europhys. Lett.* **96**, 68001 (2011).
- [3] Peyrard, M. and Bishop, A. R. Statistical mechanics of a nonlinear model for DNA denaturation. *Phys. Rev. Lett.* **62**, 2755–2757 (1989).
- [4] Weber, G., Haslam, N., Whiteford, N., Prügel-Bennett, A., Essex, J. W., and Neylon, C. Thermal equivalence of DNA duplexes without melting temperature calculation. *Nature Physics* **2**, 55–59 (2006).
- [5] Dauxois, T., Peyrard, M., and Bishop, A. R. Entropy-driven DNA denaturation. *Phys. Rev. E* **47**, R44–R47 (1993).
- [6] Weber, G. Sharp DNA denaturation due to solvent interaction. *Europhys. Lett.* **73**, 806–811 (2006).
- [7] Joyeux, M. and Buyukdagli, S. Dynamical model based on finite stacking enthalpies for homogeneous and inhomogeneous DNA thermal denaturation. *Phys. Rev. E* **72**, 051902 (2005).
- [8] Buyukdagli, S. and Joyeux, M. Scaling laws at the phase transition of systems with divergent order parameter and/or internal length: The example of DNA denaturation. *Phys. Rev. E* **73**, 51910 (2006).
- [9] Buyukdagli, S. and Joyeux, M. Theoretical investigation of finite size effects at DNA melting. *Phys. Rev. E* **76**, 021917 (2007).
- [10] Buyukdagli, S. and Joyeux, M. Statistical physics of the melting of inhomogeneous DNA. *Physical Review E* **77**, 031903 (2008).
- [11] Weber, G. Mesoscopic model parametrization of hydrogen bonds and stacking interactions of RNA from melting temperatures. *Nucl. Acids. Res.* **40**, 1 (2012).
- [12] Weber, G., Haslam, N., Essex, J. W., and Neylon, C. Thermal equivalence of DNA duplexes for probe design. *J. Phys. Condens. Matter* **21**, 034106 (2009).
- [13] Zhang, Y.-L., Zheng, W.-M., Liu, J.-X., and Chen, Y. Z. Theory of DNA melting based on the Peyrard-Bishop model. *Phys. Rev. E* **56**, 7100–7115 (1997).
- [14] Breslauer, K. J., Frank, R., Blocker, H., and Marky, L. A. Predicting DNA duplex stability from the base sequence. *Proc. Natl. Acad. Sci. USA* **83**, 3746–3750 (1986).
- [15] SantaLucia, Jr., J. A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics. *Proc. Natl. Acad. Sci. USA* **95**, 1460–1465 (1998).
- [16] Xia, T., SantaLucia, Jr., J., Burkard, M. E., Kierzek, R., Schroeder, S. J., Jiao, X., Cox, C., and Turner, D. H. Thermodynamic parameters for an expanded nearest-neighbor model for formation of RNA duplexes with Watson-Crick base pairs. *Biochem.* **37**, 14719–14735 (1998).