# Supplementary Information

# MPF-BML: A standalone GUI-based package for maximum entropy model inference

Ahmed A. Quadeer[1], Matthew R. McKay[1,2], John P. Barton[3], Raymond H. Y. Louie[4,5,*]

[1]Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Hong Kong, China, [2]Department of Chemical and Biological Engineering, The Hong Kong University of Science and Technology, Hong Kong, China, [3]Department of Physics and Astronomy, University of California, Riverside, USA, [4]The Kirby Institute, University of New South Wales, Sydney, Australia, [5]School of Medical Sciences, University of New South Wales, Sydney, Australia.

*To whom correspondence should be addressed.
Contact: rlouie@kirby.unsw.edu.au

# Table of contents

# Supplementary Text

## Text S1. Inputs

Besides the mandatory input, optional inputs are: (i) the maximum fraction of gaps above which the variable is removed (default value set to 0.5); (ii) an optional sub-sample number to train MPF (used if MPF is running slow); (iii) an optional input file in csv/xls format comprising of sample weights to reduce sampling bias (each sample is weighted according to a maximum fraction threshold of similar sequences as default); and (iv) an optional file comprising model parameters to initialize the method.

For (iii), reweighting is performed based on the provided threshold $0 < x < 1$ as follows (Morcos 2011): Two sequences with Hamming distance less than $xN$ (where $N$ is the length of protein) are considered to carry almost the same information, and vice versa. Thus, for each sequence $m$, a weight $w_m = 1/z_m$ is calculated, where $z_m$ are all the sequences in the MSA having a hamming distance less than $xN$ with the sequence $m$. A threshold value of 0.1 is used as default if the user does not provide any weight file.

For (iv), the user can provide an optional input maximum entropy model to initialize the MPF-BML method. This model may be obtained from some other inference method or from a previous run of the MPF method (the model obtained after running MPF is automatically saved once the MPF method converges) for which the user is interested in running either BML method alone or both MPF (again) and BML methods with a different set of parameters. For the latter, the user needs to check the provided "Run only BML step on user-provided parameters" option in the "Input information" panel.

## Text S2. MPF-BML parameters

Here, we provide details of all parameters available in the "Parameters" panel of the MPF-BML GUI.

### Step 1: Mutant combining

**phi_opt:** Mutant combining factor which indicates the ratio of entropy with grouping of states to the entropy without grouping. See (Louie *et al.*, 2018) for details. A value of 0 indicates each site is coarse-grained to only two states, a value of 1 indicates all the states from the original MSA are kept.

### Step 2: MPF

**L1 reg:** L1 regularization parameter. Higher value indicates more couplings will be set to zero.

**L2 reg:** L2 regularization parameter. Higher value indicates high couplings values will be more substantially reduced to zero.

**Grad. tol:** If the maximum gradient of the MPF function falls below this tolerance level, the MPF algorithm will terminate. Smaller value will require more time for the algorithm to converge but will converge with more accuracy.

**Func. tol:** If the absolute difference between the MPF function in the previous and current iteration falls below this tolerance level, the MPF algorithm terminates. Smaller value will require more time for the algorithm to converge but will converge with more accuracy.

**Max iter:** Maximum number of iterations before the MPF algorithm terminates.

### Step 3: BML

**Max iter:** Maximum number of iterations before the BML algorithm terminates.

**Max eps:** Maximum epsilon value for the individual frequencies, pairwise frequencies and connected correlations. The epsilon value is chosen to balance between underfitting and overfitting the individual frequencies, pairwise frequencies and connected correlations. Higher value will mean less accuracy but faster convergence. See (Louie *et al.*, 2018) for details.

**Param. tol:** If the sum absolute difference in parameter values in the previous and current iteration falls below this tolerance level, the BML algorithm terminates. Smaller value will require more time for the algorithm to converge but will converge with more accuracy.

**Grad. tol:** If the sum absolute difference in the gradient of the parameter values in the previous and current iteration falls below this tolerance level, the BML algorithm terminates. Smaller value will require more time for the algorithm to converge but will converge with more accuracy.

**Thinning:** Keep every x samples in the MCMC algorithm, where x is the thinning parameter. Higher value will mean less samples, thus less accuracy but faster speed.

**Burn-in:** Remove the first x samples in the MCMC algorithm, where x is the burn-in parameter. Higher value will mean the MCMC algorithm will produce more accurate output, but slower speed.

**Cores:** Number of cores to use in the MCMC algorithm. If set to a value greater than or equal to the maximum number of cores ($N_c$) available in the system, the application resets it to $N_c - 1$ to avoid using all available resources.

**MCMC samples:** Number of samples in the MCMC algorithm before thinning and burn-in. Higher value will mean the MCMC algorithm will produce more accurate output, but slower speed.

The following are RPROP specific parameters, See (Riedmiller and Braun, 1993) for more details.

**Gamma h:** Initial weight update value for field (h) parameters. Set larger for larger step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**Gamma J:** Initial weight update value for coupling (J) parameters. Set larger for larger step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**h pos:** Weight increase factor for the fields with positive gradient. Set larger for larger step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**h neg:** Weight increase factor for the fields with negative gradient. Set larger for larger step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**J pos:** Weight increase factor for the couplings with positive gradient. Set larger for larger step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**J neg:** Weight increase factor for the couplings with negative gradient. Set larger for larger step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**h delta max:** The maximum weight update per iteration for the fields. Set larger for larger possible step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**h delta min:** The minimum weight update per iteration for the fields. Set larger for larger possible step sizes. Smaller value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**J delta max:** The maximum weight update per iteration for the couplings. Set larger for larger possible step sizes. Higher value will mean initial faster convergence of algorithm, but converged values may be less accurate.

**J delta min:** The minimum weight update per iteration for the couplings. Set larger for larger possible step sizes. Smaller value will mean initial faster convergence of algorithm, but converged values may be less accurate.

## Text S3. Contact prediction using the inferred couplings

A tab-delimited text file is also provided as an output with pairs of positions arranged in descending order according to their computed Frobenius norm, a metric representative of the pairs in contact. This metric is calculated using the inferred couplings (Cocco *et al.*, 2018) and the higher the value of this metric, the higher is the chance of this pair to be in contact, and vice versa.

# Supplementary Tables

**Table S1.** Summary of the input datasets analyzed using MPF-BML.

| Dataset | Type of data | Format | Weights available | Number of samples | Number of variables | Number of parameters | Associated figure |
|---------|--------------|--------|-------------------|-------------------|---------------------|----------------------|-------------------|
| HCV E2 surface glycoprotein | Amino acid sequence data | FASTA | Yes | 3,363 | 352 | ~ $2 \times 10^6$ | Fig. 1 |
| Erdos-Renyi random graphs (ER05) | Synthetic categorical data | Microsoft Excel Open XML | No | 10,000 | 50 | ~ $5 \times 10^4$ | Fig. S1 |
| HIV p7 nucleocapsid protein | Amino acid sequence data | FASTA | No | 4,131 | 71 | ~ $1 \times 10^5$ | Fig. S2 |
| Trypsin inhibitor protein family (PF00014) | Amino acid sequence data | FASTA | No | 4,915 | 53 | ~ $5 \times 10^4$ | Fig. S3 |
| Breast cancer data | Somatic mutations (binary data) | Microsoft Excel Open XML | No | 2,327 | 100 | ~ $5 \times 10^3$ | Fig. S4 |
| Wisconsin breast cancer data (UCI repository) | Categorical data | Microsoft Excel Open XML | No | 699 | 9 | ~ $3 \times 10^3$ | Fig. S5 |
| Chess data (UCI repository) | Categorical data | Microsoft Excel Open XML | No | 3,196 | 36 | ~ $5 \times 10^2$ | Fig. S6 |

**Table S2.** Format of the inferred maximum entropy model parameters using MPF-BML. Values are shown for a synthetic test data set. The first row and column show the change of configuration at a position with A/1/B denoting a change from configuration A to configuration B at position 1. The diagonal entries of the matrix represent the inferred fields while the non-diagonal entries represent the inferred couplings. These parameters are saved automatically as a tab delimited file.

| | M/1/- | R/2/K | R/2/T | V/3/A | K/4/M | K/4/T | G/5/E | I/6/T | I/6/M | R/7/K | R/7/M | K/8/R | N/9/S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| M/1/- | 6.43877 | -0.31164 | 0.487343 | -0.8045 | 0.325733 | -0.14218 | -0.27445 | 0.150707 | 0.668197 | 0.068006 | 0.28129 | 0.27983 | 0.071754 |
| R/2/K | -0.31164 | 4.690157 | 0.004776 | 0.246884 | 0.266135 | -1.31717 | -0.72883 | 0.50299 | -0.35844 | -0.53223 | 1.163576 | -0.03381 | 0.365269 |
| R/2/T | 0.487343 | 0.004776 | -1.0441 | -0.12154 | 0.953398 | 0.620634 | 0.588852 | 0.322001 | -0.88575 | 0.614748 | 0.12424 | 0.870705 | 0.503503 |
| V/3/A | -0.8045 | 0.246884 | -0.12154 | 4.172256 | 0.283196 | 0.787504 | -0.00353 | -0.26528 | 0.033055 | -0.44571 | -1.2515 | -0.23769 | 1.322479 |
| K/4/M | 0.325733 | 0.266135 | 0.953398 | 0.283196 | 2.792618 | -2.1E-05 | -0.95053 | 0.333316 | 1.192085 | 0.744931 | -0.27076 | -0.62612 | 0.34005 |
| K/4/T | -0.14218 | -1.31717 | 0.620634 | 0.787504 | -2.1E-05 | 3.632102 | -0.1222 | 0.800195 | 0.250115 | 0.937377 | -0.27037 | -0.09487 | 1.162473 |
| G/5/E | -0.27445 | -0.72883 | 0.588852 | -0.00353 | -0.95053 | -0.1222 | 2.732045 | -0.03059 | -1.84851 | 0.655304 | 0.064181 | -0.39161 | -0.14443 |
| I/6/T | 0.150707 | 0.50299 | 0.322001 | -0.26528 | 0.333316 | 0.800195 | -0.03059 | 3.744127 | -5.5E-05 | 0.274624 | 0.329657 | 0.040394 | -0.69886 |
| I/6/M | 0.668197 | -0.35844 | -0.88575 | 0.033055 | 1.192085 | 0.250115 | -1.84851 | -5.5E-05 | 3.82687 | 1.352347 | -0.40026 | 1.040857 | 0.853013 |
| R/7/K | 0.068006 | -0.53223 | 0.614748 | -0.44571 | 0.744931 | 0.937377 | 0.655304 | 0.274624 | 1.352347 | 5.206358 | -1E-04 | -0.8084 | -0.01596 |
| R/7/M | 0.28129 | 1.163576 | 0.12424 | -1.2515 | -0.27076 | -0.27037 | 0.064181 | 0.329657 | -0.40026 | -1E-04 | 5.098788 | 0.930253 | -0.61156 |
| K/8/R | 0.27983 | -0.03381 | 0.870705 | -0.23769 | -0.62612 | -0.09487 | -0.39161 | 0.040394 | 1.040857 | -0.8084 | 0.930253 | 4.569314 | -0.99103 |
| N/9/S | 0.071754 | 0.365269 | 0.503503 | 1.322479 | 0.34005 | 1.162473 | -0.14443 | -0.69886 | 0.853013 | -0.01596 | -0.61156 | -0.99103 | 5.057654 |

# Supplementary Figures

**MPF-BML**

**Input information**

**Input data**

Load data

Gap threshold | 0.5

No. of samples to use for sub-sampling [Optional]

**Sample weights**

1. Similarity-based weighting (default) | 2. Predefined weighting

Threshold | 0.1 | Load sample weights

**Initialization [Optional]** | Load initial parameters

☑ Run only BML step on user-provided parameters

**Parameters**

**Step 1: Mutant combining**

phi_opt

**Step 2: MPF**

| L1 reg | 10 | L2 reg | 30 |
| Grad tol | 1e-4 | Func tol | 1e-4 |
| Max iter | 1e4 | | |

**Step 3: BML**

| Max iter | 1000 | Max eps | 1.5 |
| Param tol | 1e-5 | Grad tol | 1e-5 |
| Thinning | 3e3 | Burn-in | 1e4 |
| Cores | 4 | MCMC samples | 1e7 |
| Gamma h | 1e-4 | Gamma J | 1e-4 |
| h pos | 1.05 | h neg | 0.95 |
| J pos | 1.05 | J neg | 0.95 |
| h delta max | 1e-8 | h delta min | 1e-8 |
| J delta max | 1e-4 | J delta min | 1e-8 |

**Run MPF-BML** | Stop

**Compute energies**

Load data for computing energies

**Processing information**

```
Input data: ER05-configurations.xlsx.........................
Init. params. file: ER05-configurations_J_MPF.csv..............
-----------------------------------------------------------
Checking input data and parameters.........................
-----------------------------------------------------------
Computing sequence weights using threshold 0.10..............
Initial model parameters provided..........................
Skipping MPF and running only BML, as specified...............
Warning! Max cores available = 2; using cores = 1.............
All tests cleared!.........................................
Samples = 10000; Effective samples after re-weighting = 10000.0
Positions = 50; Informative positions = 50...................
-----------------------------------------------------------
Creating output directory..................................
-----------------------------------------------------------
/Users/ahmed/MATLAB/HKUST/MPF-BML-master/ER05-configurations...
_277_output/...............................................
-----------------------------------------------------------
Step 1: Mutant combining...................................
-----------------------------------------------------------
Inferred phi_opt = 1.00....................................
-----------------------------------------------------------
Step 2: MPF................................................
-----------------------------------------------------------
Skipping MPF. Using provided initial model parameters.........
Plotting MPF-based model verification results.................
-----------------------------------------------------------
Step 3: BML................................................
-----------------------------------------------------------
Iter   Single ε   Double ε    Conn. ε.......................
 990  1.7588344  1.2747432  1.2613410.......................
 995  1.6863893  1.2741051  1.2629095.......................
1000  1.8567414  1.2710147  1.2531179.......................
Terminated: Max. BML iterations reached.....................
Plotting MPF-BML-based model verification results............
Saving inferred model paramters "J_MPF_BML"..................
-----------------------------------------------------------
```

**Model verification**

**Model: MPF**

Individual frequencies, ε = 3.84

Pair-wise frequencies, ε = 1.20

Connected corrs., ε = 1.04

r = 0.963, p = 3.49e-22

**Model: MPF-BML**

Individual frequencies, ε = 1.73

Pair-wise frequencies, ε = 1.26

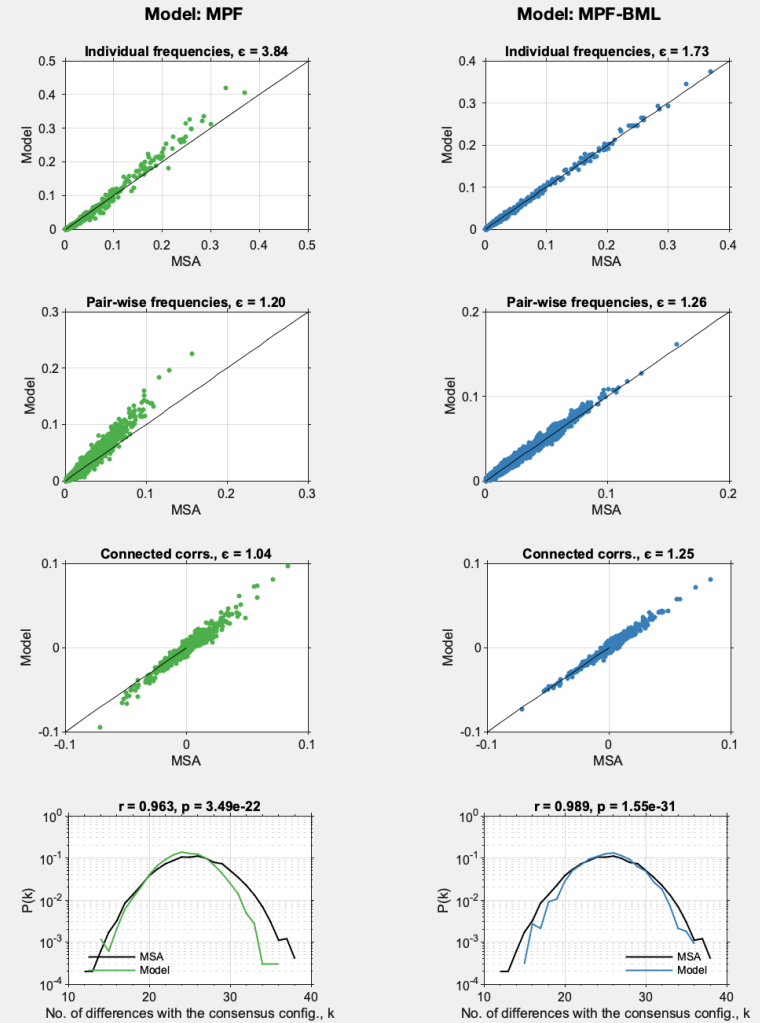Connected corrs., ε = 1.25

r = 0.989, p = 1.55e-31

**Figure S1.** Model inferred using the MPF-BML package for Erdos-Renyi random graphs (ER05), analyzed in (Barton *et al.*, 2016). Here, the model is initialized by a set of parameters previously inferred using the MPF method.
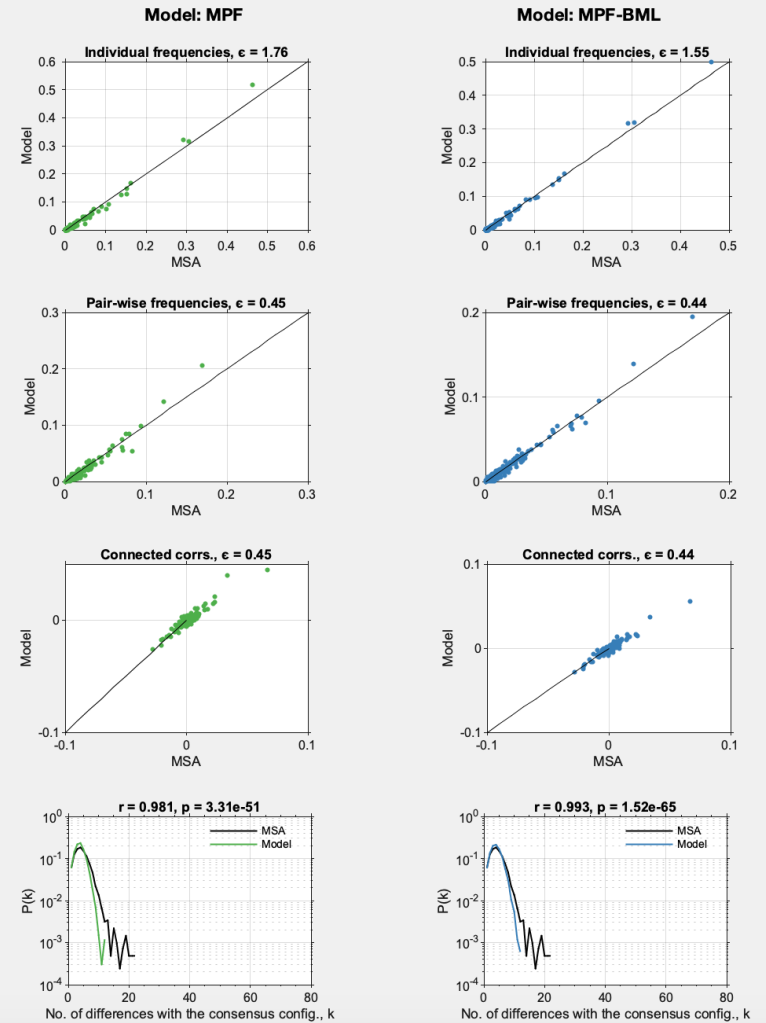
**Figure S2.** Model inferred using the MPF-BML package for HIV p7 nucleocapsid protein, analyzed in (Barton *et al.*, 2016).

# MPF-BML

## Input information

**Input data**

[Load data]

Gap threshold: 0.5

No. of samples to use for sub-sampling [Optional]: [ ]

**Sample weights**

1. Similarity-based weighting (default)    |    2. Predefined weighting

Threshold: 0.1    |    [Load sample weights]

**Initialization [Optional]**    [Load initial parameters]

☐ Run only BML step on user-provided parameters

## Parameters

**Step 1: Mutant combining**

phi_opt: [ ]

**Step 2: MPF**

| L1 reg | 0.01 | L Grad tol | 0.5 |
| Grad tol | 1e-4 | Func tol | 1e-4 |
| Max iter | 1e4 | | |

**Step 3: BML**

| Max iter | 1000 | Max eps | 1 |
| Param tol | 1e-5 | Grad tol | 1e-5 |
| Thinning | 3e3 | Burn-in | 1e4 |
| Cores | 4 | MCMC samples | 1e7 |
| Gamma h | 1e-4 | Gamma J | 1e-4 |
| h pos | 1.05 | h neg | 0.95 |
| J pos | 1.05 | J neg | 0.95 |
| h delta max | 1e-8 | h delta min | 1e-8 |
| J delta max | 1e-4 | J delta min | 1e-8 |

[**Run MPF-BML**]    [Stop]

## Compute energies

[Load data for computing energies]

## Processing information

```
Input data: PF00014-alignment.fasta...........................
------------------------------------------------------------
------------------------------------------------------------
Checking input data and parameters..........................
------------------------------------------------------------
Computing sequence weights using threshold 0.10.............
Warning! Max cores available = 2; using cores = 1...........
All tests cleared!..........................................
Samples = 4915; Effective samples after re-weighting = 2589.66.
Positions = 53; Informative positions = 53..................
------------------------------------------------------------
Creating output directory...................................
------------------------------------------------------------
/Users/ahmed/MATLAB/HKUST/MPF-BML-master/PF00014-alignment_4...
22_output/..................................................
------------------------------------------------------------
Step 1: Mutant combining....................................
------------------------------------------------------------
Inferred phi_opt = 1.00.....................................
------------------------------------------------------------
Step 2: MPF.................................................
------------------------------------------------------------
Iter  fEvals    stepLen      fVal      optCond     nnz.........
   1      2  2.308e-08  1.391e+06  3.626e+03   835396.........
Terminated: Progress in parameters or objective below Func Tol.
Saving inferred model parameters "J_MPF"....................
Plotting MPF-based model verification results...............
------------------------------------------------------------
Step 3: BML.................................................
------------------------------------------------------------
Iter    Single ε    Double ε    Conn. ε.....................
 990  0.9667344  1.0815723  1.0903464......................
 995  0.9372664  1.0774342  1.0959069......................
1000  0.9657123  1.0811389  1.0951105......................
Terminated: Max. BML iterations reached.....................
Plotting MPF-BML-based model verification results...........
Saving inferred model paramters "J_MPF_BML".................
------------------------------------------------------------
```

## Model verification

**Model: MPF**

Individual frequencies, ε = 1.80

X: 0.41341
Y: 0.31682
Config: P21

Pair-wise frequencies, ε = 1.17

Connected corrs., ε = 1.32

r = 0.630, p = 7.90e-06

**Model: MPF-BML**

Individual frequencies, ε = 0.89

Pair-wise frequencies, ε = 1.07
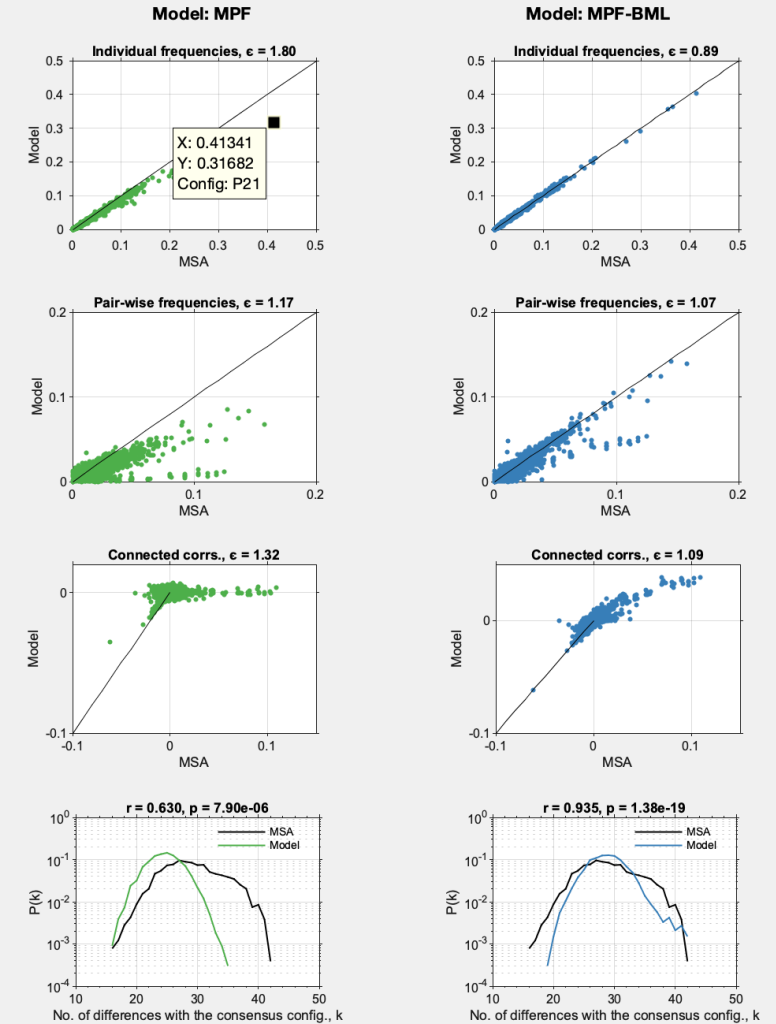
Connected corrs., ε = 1.09

r = 0.935, p = 1.38e-19

**Figure S3.** Model inferred using the MPF-BML package for Trypsin inhibitor protein family (PF00014), analyzed in (Barton *et al.*, 2016).
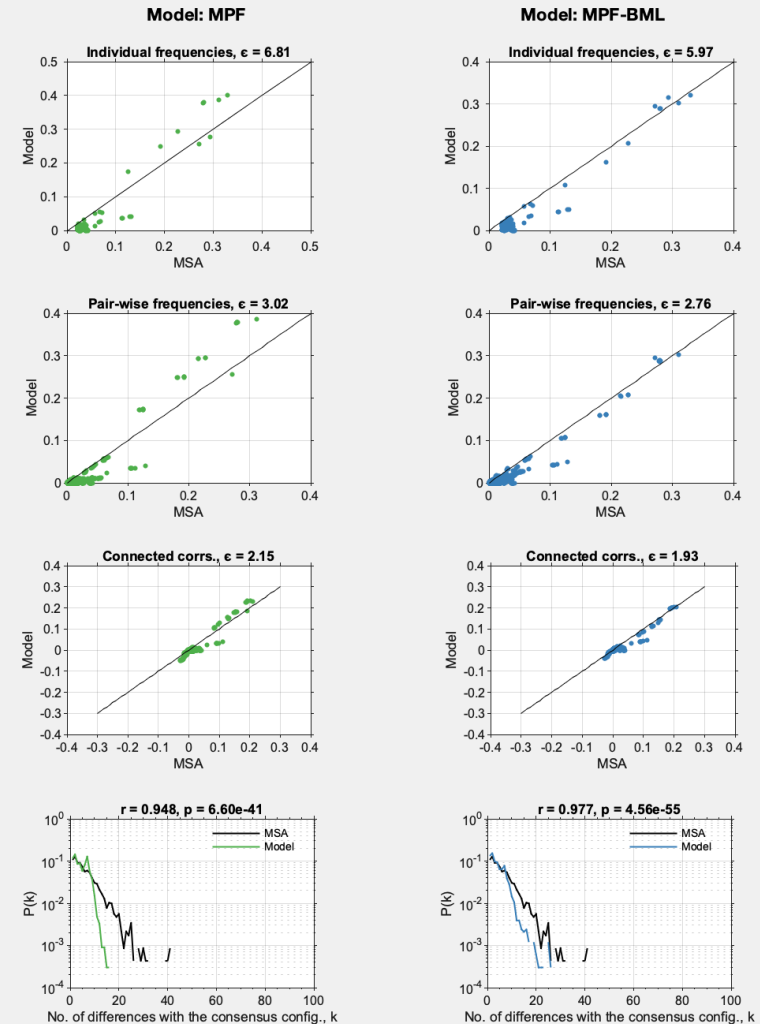
**Figure S4.** Model inferred using the MPF-BML package for whole genome somatic mutation (100 most highly expressed genes) breast cancer data (Bamford *et al.*, 2004) obtained from http://cancer.sanger.ac.uk/.

# MPF-BML

## Input information

**Input data**

Load data

Gap threshold: 0.5

No. of samples to use for sub-sampling [Optional]: ____

**Sample weights**

1. Similarity-based weighting (default)  |  2. Predefined weighting

Threshold: 0  |  Load sample weights

**Initialization [Optional]**  Load initial parameters

☐ Run only BML step on user-provided parameters

## Parameters

**Step 1: Mutant combining**

phi_opt: ____

**Step 2: MPF**

| | | | |
|---|---|---|---|
| L1 reg | 0.01 | L2 reg | 0.02 |
| Grad tol | 1e-4 | Func tol | 1e-4 |
| Max iter | 1e4 | | |

**Step 3: BML**

| | | | |
|---|---|---|---|
| Max iter | 1000 | Max eps | 1.0 |
| Param tol | 1e-5 | Grad tol | 1e-5 |
| Thinning | 3e3 | Burn-in | 1e4 |
| Cores | 4 | MCMC samples | 1e7 |
| Gamma h | 1e-4 | Gamma J | 1e-4 |
| h pos | 1.05 | h neg | 0.95 |
| J pos | 1.05 | J neg | 0.95 |
| h delta max | 1e-8 | h delta min | 1e-8 |
| J delta max | 1e-4 | J delta min | 1e-8 |

**Run MPF-BML**    Stop

## Compute energies

Load data for computing energies

## Processing information

```
Input data: breast_cancer_wisconsin_data_processed.csv.........
Release the "Stop" button and re-run MPF-BML!.................
----------------------------------------------------------------
Checking input data and parameters...........................
----------------------------------------------------------------
Computing sequence weights using threshold 0.00..............
Warning! Max cores available = 2; using cores = 1.............
All tests cleared!...........................................
Samples = 699; Effective samples after re-weighting = 699.00...
Positions = 9; Informative positions = 9.....................
----------------------------------------------------------------
Creating output directory....................................
----------------------------------------------------------------
/Users/ahmed/MATLAB/HKUST/MPF-BML-master/breast_cancer_wisco...
nsin_data_processed_352_output/..............................
----------------------------------------------------------------
Step 1: Mutant combining.....................................
----------------------------------------------------------------
Inferred phi_opt = 1.00......................................
----------------------------------------------------------------
Step 2: MPF..................................................
----------------------------------------------------------------
Iter   fEvals    stepLen      fVal     optCond     nnz.........
270      271   1.000e+00  9.987e+03  1.150e+01    6394.........
280      281   1.000e+00  9.969e+03  1.145e+01    6393.........
290      291   1.000e+00  9.951e+03  1.143e+01    6393.........
300      301   1.000e+00  9.932e+03  1.139e+01    6392.........
Terminated: Directional derivative below Func Tol.............
Saving inferred model parameters "J_MPF"......................
Plotting MPF-based model verification results.................
----------------------------------------------------------------
Step 3: BML..................................................
----------------------------------------------------------------
Iter    Single є    Double є     Conn. є.......................
285   1.2747017   1.3861288    1.3478115.......................
290   3.0805567   3.1426328    2.8622453.......................
295   1.5145814   1.6819056    1.6184130.......................
300   1.2800425   1.3979934    1.3579532.......................
304   0.8184159   0.9727669    0.9768026.......................
Terminated: Average є < 1.00 ................................
Plotting MPF-BML-based model verification results.............
Saving inferred model paramters "J_MPF_BML"...................
----------------------------------------------------------------
```

## Model verification

**Model: MPF**

- Individual frequencies, є = 10.97
- Pair-wise frequencies, є = 10.58
- Connected corrs., є = 7.98
- r = 0.594, p = 9.19e-02 (MSA, Model)

**Model: MPF-BML**

- Individual frequencies, є = 1.37
- Pair-wise frequencies, є = 1.49
- Connected corrs., є = 1.43
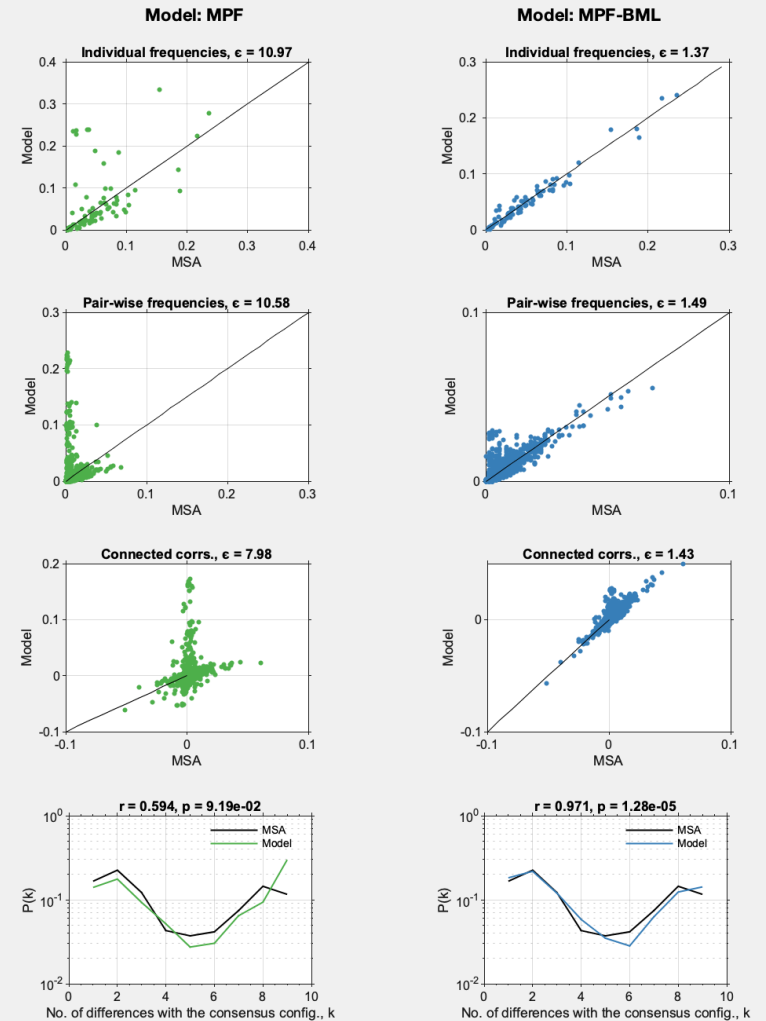- r = 0.971, p = 1.28e-05 (MSA, Model)

**Figure S5.** Model inferred using the MPF-BML package for the classical Wisconsin breast cancer data (Wolberg and Mangasarian, 1990) obtained from the UCI repository (Asuncion and Newman, 2007).
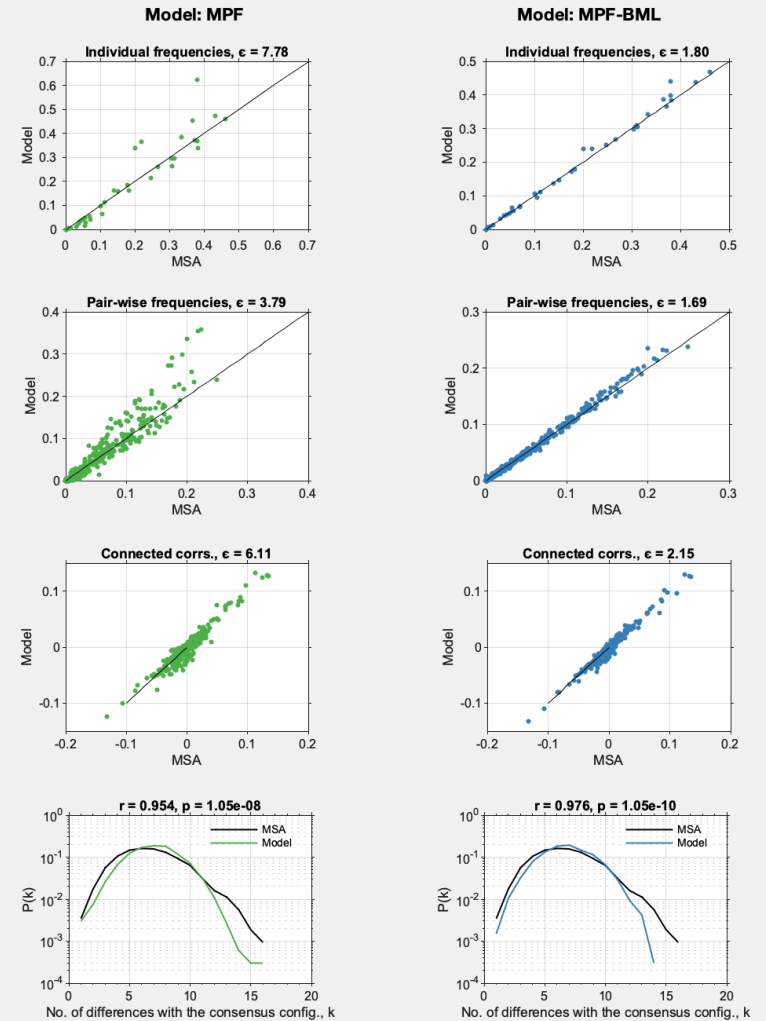
**Figure S6.** Model inferred using the MPF-BML package for the classical chess data set obtained from the UCI repository (Asuncion and Newman, 2007).

# Supplementary References

Asuncion,A. and Newman,D.J. (2007) UCI Machine Learning Repository
[http://www.ics.uci.edu/~mlearn/MLRepository.html]. *Univ. California, Irvine, Sch. Inf. Comput. Sci.*

Bamford,S. *et al.* (2004) The COSMIC (Catalogue of Somatic Mutations in Cancer) database and website. *Br. J. Cancer*, **91**, 355–358.

Barton,J.P. *et al.* (2016) ACE: adaptive cluster expansion for maximum entropy graphical model inference. *Bioinformatics*, **32**, 3089–3097.

Cocco,S. *et al.* (2018) Inverse statistical physics of protein sequences: a key issues review. *Reports Prog. Phys.*, **81**, 032601.

Louie,R.H.Y. *et al.* (2018) Fitness landscape of the human immunodeficiency virus envelope protein that is targeted by antibodies. *Proc. Natl. Acad. Sci.*, **115**, E564–E573.

Riedmiller,M. and Braun,H. (1993) A direct adaptive method for faster backpropagation learning: the RPROP algorithm. *IEEE Int. Conf. Neural Networks*, **1**, 586–591.

Wolberg,W.H. and Mangasarian,O.L. (1990) Multisurface method of pattern separation for medical diagnosis applied to breast cytology. *Proc. Natl. Acad. Sci.*, **87**, 9193–9196.