

## Electronic Supplementary Information

### Bayesian inference assessment of protein secondary structure analysis using circular dichroism data – how much structural information is contained in protein circular dichroism spectra?

Simon E.F. Spencer and Alison Rodger

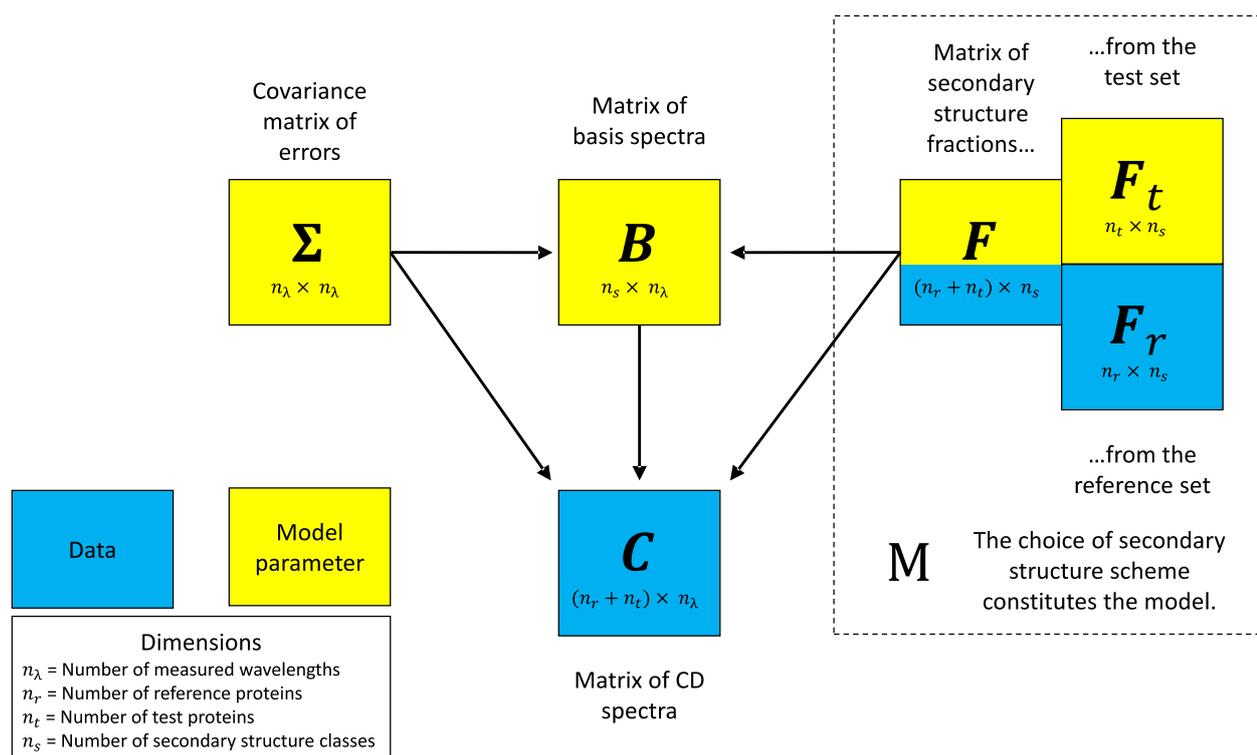


Fig S1: Conditional independence diagram showing the relationships between the model parameters and the data for our Bayesian model.

<sup>a</sup> Department of Statistics, University of Warwick, Coventry, UK.

[s.e.f.spencer@warwick.ac.uk](mailto:s.e.f.spencer@warwick.ac.uk)

<sup>b</sup> Author contribution: J. Franco worked on this project in its early stages, but he is unaware that this paper is being submitted for publication as we have no contact details for him. Therefore, he does not take any responsibility for the contents.

<sup>c</sup> Department of Molecular Sciences, Macquarie University, NSW, 2109, Australia. [alison.rodger@mq.edu.au](mailto:alison.rodger@mq.edu.au); here.

Electronic Supplementary Information (ESI) available: [details of any supplementary information available should be included here]. See DOI: 10.1039/x0xx00000x

Table S1. Cross-validation results for the SP175 proteins with 3 secondary structure classes from DSSP:  $\alpha$ -helix, -sheet and Other. Results for competing approaches (SEL-MAT3, PCR, PLS, NN and SVM, taken from reference<sup>31</sup>). The best per-forming approach for each measure is given in bold.

	$\alpha$ -helix			$\beta$ -sheet			Other		
	$\delta$	$\zeta$	$r$	$\delta$	$\zeta$	$r$	$\delta$	$\zeta$	$r$
Bayesian	0.061	3.48	0.959	0.127	1.25	0.770	0.137	0.77	0.652
SELMAT3	0.063	3.35	0.957	0.083	1.90	0.862	0.078	1.35	0.701
PCR	0.057	3.70	0.966	0.069	2.29	0.906	0.066	1.61	0.796
PLS	<b>0.053</b>	<b>3.98</b>	<b>0.971</b>	0.073	2.17	0.895	0.068	1.56	0.781
NN	0.055	3.84	0.968	<b>0.067</b>	<b>2.36</b>	<b>0.912</b>	<b>0.062</b>	<b>1.71</b>	<b>0.816</b>
SVM	0.057	3.70	0.966	0.069	2.29	0.908	0.066	1.61	0.792

Table S2. Cross-validation results for the SP175 proteins with the SELCON secondary structure scheme (see Section 2.5). Results for competing approaches (SELMAT3 and PCR, see Section 3.2) taken from Lees et al. (2006a). The best performing approach for each measure is given in bold.

Structure	Bayesian			SELMAT3			PLS		
	$\delta$	$\zeta$	$r$	$\delta$	$\zeta$	$r$	$\delta$	$\zeta$	$r$
Regular helix	0.091	1.73	0.836	0.048	3.28	0.956	<b>0.040</b>	<b>3.94</b>	<b>0.971</b>
Distorted helix	0.129	0.46	0.043	<b>0.035</b>	<b>1.70</b>	<b>0.809</b>	0.036	1.66	0.791
Regular $\beta$ -strand	0.090	1.31	0.695	0.073	1.62	0.792	<b>0.063</b>	<b>1.88</b>	<b>0.853</b>
Distorted $\beta$ -strand	0.281	0.17	-0.081	<b>0.020</b>	<b>2.41</b>	<b>0.913</b>	0.023	2.10	0.889
Turn	0.201	0.27	0.098	<b>0.052</b>	<b>1.04</b>	0.325	<b>0.052</b>	<b>1.04</b>	<b>0.332</b>
Other	0.169	0.43	0.278	<b>0.050</b>	<b>1.45</b>	0.717	<b>0.050</b>	<b>1.45</b>	<b>0.720</b>

Table S3. Cross validation results for the SP175 proteins with the BeStSel secondary structure scheme for different minimum wavelengths. The classification scheme for each number of classes was chosen to be the scheme with the largest marginal likelihood. Values of the normalised measure zeta above one indicate that the predictions are more accurate than a random draw from the reference set. Values below one have been highlighted in orange.

Minimum wavelength (nm) BeStSel classes		$\zeta$ values					
		175	180	185	190	195	200
3 classes	Helix1	3.57	3.55	3.50	3.37	3.45	3.52
	Anti1+Anti2+Parallel	2.02	2.00	2.03	2.02	1.97	1.86
	Helix2+Anti3+Turn+Others	1.84	1.83	1.86	1.90	1.86	1.72
4 classes	Helix1	4.08	4.16	3.99	4.19	4.23	4.14
	Anti1+Anti2	1.37	1.31	1.33	1.30	1.35	1.51
	Parallel	0.54	0.52	0.43	0.50	0.58	0.69
	Helix2+Anti3+Turn+Others	1.32	1.29	1.11	1.33	1.46	1.38
5 classes	Helix1	3.47	3.44	3.49	3.46	3.59	3.67
	Anti1+Anti2	1.42	1.35	1.36	1.35	1.40	1.69
	Helix2+Anti3	0.32	0.32	0.34	0.37	0.34	0.33
	Parallel	0.51	0.50	0.53	0.50	0.65	0.86
	Turn+Others	0.43	0.45	0.48	0.50	0.50	0.52