

eSPC, an Online Data Analysis Platform for Molecular Biophysics

ChiraKit 1.0 User Documentation

June 2024

This document is a work in progress. Some information may be incomplete or subject to change.

Table of Contents

- 1. Import data
 - 1.1. Input file
 - 1.2. Parameters for molar ellipticity / extinction
 - 1.3. Processing
- 2. Analysis
 - 2.1. Creating a dataset
 - 2.2. Thermal unfolding
 - 2.3. Chemical unfolding
 - 2.4. Custom analysis
 - 2.5. Curve fitting and fitting errors
 - 2.6. Comparison of spectra
- 3. Spectra decomposition
 - 3.1. Singular value decomposition (SVD)
 - 3.2. Principal component analysis (PCA)
 - 3.3. Basis spectrum inversion
 - 3.4. Change of basis
 - 3.5. Alternative formulation for the explained variance
- 4. Protein secondary structure estimation
 - 4.1. The Selcon method
 - 4.2. Reference sets
 - 4.3. Secondary structure elements
- Contact details

1. Import data

1.1. Input file

ChiraKit can parse many different types of files:

A) A comma-separated-values (**.csv**) file with a header, and two or three columns. In the case of two columns, the first contains the wavelength and the second the CD signal. In the case of three columns, the first contains the wavelength, the second the sample name, and the third the CD signal. In both cases the first column must be called 'wavelength'.

	Standard	Standard	Standard
1	wavelength	time	CD_signal
2	220	1	10
3	220	2	20
4	220	3	30
5	220	4	40

Figure 1. Example of a valid CSV format.

B) A csv file (**.csv**) with a header and **more than three columns**: 'wavelength', CD signal of the first sample , ..., CD signal of the n-th sample. The first column must be called 'wavelength' or 'xydata'.

	Standard	Standard	Standard	Standard
1	wavelength	CD_signal_A	CD_signal_B	CD_signal_C
2	220	1	2	3
3	221	1	2	3
4	222	1	2	3
5	223	1	2	3

Figure 2. Example of a valid CSV format. Each column represents one CD curve.

C) A csv file (**.csv**) with header lines containing the metadata. The metadata keys and values should be separated by ':'. The CD data should start two lines after a line containing the column heading of 'Wavelength' and should be arranged in columns. The first column has the wavelength data (in nm), and the subsequent columns have the spectral data (e.g., scan repeats). This type of file is exported by the ProData software from AppliedPhotophysics™.

Available Dimensions:		2	
Wavelength	Wavelength: 180nm - 290nm	Step Size: 1nm	
Repeat	2 repeats in set.	-iter option selected	
Available Properties:		1	
CircularDichroism			
Data:			
CircularDichroism			
Wavelength	Repeat		
		0	1
	290	-0.313281	-0.284654
	289	-0.353384	-0.46977
	288	-0.488774	-0.488609
	287	-0.277936	-0.486426
	286	-0.4056	-0.473476

Figure 3. Example of a ProData generated CSV.

D) A PCDDDB file (.pcd, Protein Circular Dichroism Data Bank format) with header lines containing the metadata. The metadata keys and values are organised into two columns, with the second column starting after the 60th character. The CD data starts after the first line containing 'DATA' and 'wavelength', and is arranged in columns. The first column has the wavelength data (in nm), and the subsequent columns have the spectra data (e.g., 1. Wavelength. 2. Final. 3. HT. 4.Smoothed. 5. Avg. Sample. 6. Avg. Baseline.).

Publication Journal	Bioinformatics
Publication Title	A reference database for circi
Publication Volume	22
Publication Pages	1955-1962
Depositor Name	Andrew Miles
Department/School name	Crystallography
University/Institution/Corporation	Birkbeck
Depositor Country	United Kingdom
Name of Principal Investigator (if not depositor)	B.A. Wallace
DATA (1. Wavelength. 2. Final. 3. HT. 4.Smoothed. 5. Avg. Sample. 6. Avg. Baseline.)	
280.0	8.98482E-02 3.84900E+02 8.75000E-02 1.63333E+00 1.50000E+00
279.0	-1.28364E-02 3.84000E+02 -1.25000E-02 1.60000E+00 1.56667E+00

Figure 4. Example of a PCDDDB file. Original source:

<https://pcddb.cryst.bbk.ac.uk/deposit/CD0000003000&dl=1#>

E) A 'generic' file (.gen, CDTool format) with header lines containing the metadata. The metadata keys and values are organised into two columns, separated by tabs. The CD data starts after the first line with more than three tab divided columns. The first column has the wavelength data (in nm), and the subsequent columns have the spectral data (e.g., 'Final Processed Spectrum', 'HT values (from raw spectrum 1)'). Duplicated or missing data is removed.

Smoothing window	5
PDB	1ed9
Swissprot	P00634
Code	Lees, J.G., Wien, F., Miles, A.J. & Wallace, B.A., 2006 Bioinformatics A refere
structure space	22 1955-1962
279.0	5.62969E-02 1.53648E+00 5.62969E-02 0.00000E+00 0.00000E+00 0.00000E+00
278.0	2.21733E-01 1.53303E+00 2.21733E-01 0.00000E+00 0.00000E+00 0.00000E+00
277.0	4.74813E-02 1.52912E+00 4.74813E-02 0.00000E+00 0.00000E+00 0.00000E+00

Figure 5. Example of a 'generic' file. Original source:
<https://pcddb.cryst.bbk.ac.uk/deposit/CD0000002000&dl=1#>

F) A text file (.dat) with header lines containing the metadata. The metadata keys and values should be separated by ':'. The CD data starts two lines after the line containing 'MDCDATA' in it and is arranged in columns. The first column has the wavelength data (in nm), and the subsequent columns have the spectral data. There could be one or more scans separated in blocks. The voltage data is read from the column called 'CD_Dynode'.

```
Averaging Time   : 1.000 seconds
Settling Time    : 0.333 seconds
Multi-Scan Wait  : 1.00 seconds

$DATA
$MDCNAME:Scan_#1
$MDCDATA:1:14:2:3:4:9
  X  CD_Signal  CD_Error  CD_Current_(Abs)  CD_Delta_Absorbance  CD_Dynode  Jacket_Temp.
280.000  0.000  0.047  1.013  0.000  242.6  19.99
279.000 -0.006  0.083  1.013 -0.000  242.9  19.98
278.000  0.044  0.148  1.013  0.000  243.2  19.99
```

Figure 6. Example of a valid text file. Original source:
<https://github.com/samirelanduk/CDtool/blob/master/ftests/files/three-aviv-baseline.dat>

G) Aarhus synchrotron data file (.dx, where x means any number) with header lines containing the metadata. The metadata keys and values are separated into two columns where the second column starts after the 32th character (or the character ':'). The CD data starts after the first line containing 'Lambda' and 'CD', and is arranged in columns. The first column has the wavelength data (in nm), the second column has the spectral data ('CD/mdeg'), the fifth column the voltage data ('Servo_Volts'), and the eighth column the temperature data ('temperature').

```
;Grating / CD slit          LEG / 08.75
; Comments:
; Sample0: Water
; Cell type: AS 121a123
; 330-170nm 1nm 20av 1sc Slit 8.75mm
;Lambda      CD/mdeg      Y_Comp./mdeg      DC_Bias      Servo_Volts      Z_Motor      Beam_current      temperature
330.000      1.4818      0.0095      6.82440      4.37090      0.3505      180.3981      25.00
329.000      1.0437      0.0095      6.82440      4.36770      0.2660      180.3248      25.00
```

Figure 7. Example of a CD data file generated by the AU-SRCD facility at ASTRID2.

H) A csv file (.csv) containing a line with the words 'ORIGIN' and 'JASCO'. The CD data is arranged in columns and the column names are defined between the lines 'XUNITS' and 'FIRSTX'. This type of file is exported by Jasco™ circular dichroism instruments.

XUNITS	NANOMETERS		
YUNITS	CD [mdeg]		
Y2UNITS	HT [V]		
Y3UNITS	ABSORBANCE		
FIRSTX	260		
LASTX	180		
NPOINTS	801		
FIRSTY	-0,28448		
MAXY	0,86205		
MINY	-8,74658		
XYDATA			
260	-0,284475	254,707	0,758667
259,9	-0,279734	254,753	0,759227

Figure 8. Example of a CD data file exported from a Jasco™ circular dichroism instrument.

I) A csv file (**.csv**) with header lines containing the metadata. The metadata keys and values should be separated by ':'. The CD data starts two lines after the line containing 'Wavelength' and 'Temperature' in it and is ordered in columns. The first column has the wavelength data (in nm), and the subsequent columns have the spectral data (one column per temperature). The associated temperature data is available one line before the CD data starts. This type of file is exported by the ProData software from AppliedPhotophysics™.

1.2. Parameters for molar ellipticity / extinction

Each input file can be associated with a certain set of experimental parameters (Box '2. Parameters for molar ellipticity / extinction') . The expected units to input in ChiraKit are:

Molecular weight:	Dalton (g/mol)
Number of residues:	Unitless
Concentration:	mg/ml
Path length:	millimetre

These parameters are useful to compare different samples. Indeed, the CD signal can be normalised based on the sample concentration (C), cell path length (L), and molecular weight (M) or mean residual molecular weight (M_r).

A summary of the CD units is presented below.

Name	Formula	Units
Differential absorbance	$A_l - A_r$	Unitless

Differential milliabsorbance	Absorbance / 1000	Unitless
Molar extinction	Absorbance * M / (C*L)	Lts * mol ⁻¹ * cm ⁻¹
Mean residue molar extinction ($\Delta\epsilon$) ¹	Absorbance * M _r / (C*L)	Lts * mol ⁻¹ * cm ⁻¹
Degrees	Absorbance * 32.98	Degrees (°)
Millidegrees	Degrees / 1000	m°
Molar ellipticity	Degrees * M * 100 / (C*L)	°*cm ² * dmol ⁻¹
Mean residue molar ellipticity ¹	Degrees * M _r * 100 / (C*L)	°*cm ² * dmol ⁻¹

Table 1. Units of CD Measurement. A_l and A_r refer respectively to the left- and right-circularly polarised light. ¹Mean residue molar extinction and mean residue molar ellipticity units are also expressed as Lts * mol⁻¹ * cm⁻¹ * residue⁻¹ and °*cm² * dmol⁻¹ * residue⁻¹, respectively, where ‘residue’ is the total number of residues.

1.3. Processing

Given a selection of spectra (column ‘Spectrum/a 1’ inside the box ‘3. Processing’), we can:

- 1) Subtract another spectrum (column ‘Spectrum 2’)
- 2) Add another spectrum (column ‘Spectrum 2’)
- 3) Smooth them using a Savitzky-Golay filter (window size of 6, 8 or 10 nm)
- 4) Average them
- 5) Average them in batch mode
- 6) Zero them by subtracting the mean signal within a selected wavelength range interval (3, 5, 10, 20 or 40 nm)

The ‘subtract’, ‘sum’, ‘smooth’ and ‘zero’ operations will produce as many spectra as defined in the column ‘Spectrum/a 1’. The ‘average’ operation will generate only one spectrum. The ‘batch average’ operation depends on the selected ‘N’. For instance, if we have six spectra and select ‘N’ equal to 3, we will create two new spectra.

Regarding the ‘subtract’ and ‘sum’ operations, the processing step will as default handle data that was measured with the same wavelength step and wavelength range. If desired, linear interpolation can be activated. In that case, 1) the resulting spectrum will only cover the region where the spectra overlap and 2) a linear interpolation will be performed, implying that the new spectrum will be computed at all the wavelength data points present in both spectra.

Last, the following behaviour is applied for the high tension voltage curves. For the subtract, addition, smooth or zero operations, the newly generated spectra will have the same high tension voltage curve as the spectra selected in the column 'Spectrum/a 1'. For the (batch) average operation, the high tension voltage curves will be averaged.

2. Analysis

2.1. Creating a dataset

When generating a dataset for the thermal /chemical unfolding analysis, custom analysis, or spectra comparison module, the CD signal will be merged based on the wavelength data. In other words, only wavelengths present in all spectra will be used. Non-shared wavelengths will be discarded.

2.2. Thermal unfolding

2.2.1. Equilibrium two-state model

To analyse the CD signal as a function of the temperature, ChiraKit allows fitting an equilibrium two-state model^{1,2}.

Briefly, we assume that the protein only exists in the native (folded) or unfolded state and that there is an equilibrium between these two states given by the unfolding reaction $N \rightleftharpoons U$. The signal is described by the following equation:

$$Signal(T) = \frac{K_u}{1+K_u} (k_U T + b_U) + \frac{1}{1+K_u} (k_N T + b_N) \quad (1)$$

where T is the temperature, k_N , b_N are the slope and intercept of the pre-transition baseline (native), and k_U and b_U are the slope and intercept of the post-transition (unfolded) baseline.

Moreover, the unfolding constant K_u is defined as,

¹ Santoro, M. M., & Bolen, D. W. (1988). Unfolding free energy changes determined by the linear extrapolation method. 1. Unfolding of phenylmethanesulfonyl. alpha.-chymotrypsin using different denaturants. *Biochemistry*, 27(21), 8063-8068.

² Bedouelle, H. (2016). Principles and equations for measuring and interpreting protein stability: From monomer to tetramer. *Biochimie*, 121, 29-37.

$$K_u(T) = e^{-\Delta G / RT} \quad (2)$$

where R is the universal gas constant, and ΔG is the free energy of unfolding given by,

$$\Delta G = \Delta H(1 - 1/T_m) \quad (3)$$

where ΔH is the enthalpy of unfolding and T_m is the melting temperature. This model assumes that the heat capacity ΔC_p unfolding equals zero.

2.2.2. Equilibrium three-state model

To analyse the CD signal as a function of the temperature, ChiraKit allows fitting an equilibrium three-state model³.

Briefly, we assume that the protein exists in the native (folded), intermediate or unfolded state and that there is an equilibrium between these three states given by the reactions $N \rightleftharpoons I$ and $I \rightleftharpoons U$. The signal is described by the following equation:

$$Signal(T) = f_u(k_u T + b_u) + f_n(k_n T + b_n) + f_i(b_i) \quad (4)$$

where T is the temperature, k_n , b_n are the slope and intercept of the pre-transition baseline (native), k_u and b_u are the slope and intercept of the post-transition (unfolded) baseline, and b_i is the baseline of the intermediate state. f_u , f_n and f_i are respectively the unfolded, folded and intermediate fractions, which are calculated as:

$$K_1(T) = e^{(-\Delta H_1(1-\frac{T}{T_1})/RT)} \quad (5)$$

$$K_2(T) = e^{(-\Delta H_2(1-\frac{T}{T_2})/RT)} \quad (6)$$

$$f_n(T) = \frac{1}{1+K_1(T)+K_1(T)K_2(T)} \quad (7)$$

$$f_i(T) = \frac{K_1(T)}{1+K_1(T)+K_1(T)K_2(T)} \quad (8)$$

³ Almeida, Ana V., *et al.* "Structural features and stability of apo-and holo-forms of a simple iron–sulfur protein." *European Biophysics Journal* 50.3 (2021): 561-570.

$$f_u(T) = \frac{K_1(T)K_2(T)}{1+K_1(T)+K_1(T)K_2(T)} \quad (9)$$

where ΔH_1 and ΔH_2 are respectively the enthalpy of unfolding for the first and second transitions, T_1 is the temperature at which f_N equals f_i , and T_2 is the temperature at which f_i equals f_u . R is the gas constant in kcal/mol. The heat capacity of unfolding $C_p(T)$ is assumed to be zero.

2.2.3. Initial parameter estimates

The initial values of the parameters k_u , b_u , k_n and b_n are estimated by fitting the equation of a line to the first or last 12 degrees. For the two-state model, the initial T_m value is estimated based on the maximum (or minimum) of the first derivative in the interval $[(T_{min} + 6) : (T_{max} - 6)]$ where T_{min} and T_{max} are respectively the minimum and maximum measured temperatures. For the three-state model, the initial value of T_1 and T_2 are respectively $T_{min} + 10$ and $T_{max} - 20$. Be sure to check the fitted values, because T_1 is constrained between $[T_{min} + 4, T_{max} - 7]$ and T_2 is constrained between $[T_{min} + 4, T_{max}]$. The initial values of T_1 and T_2 can be overridden by the user to 'guide' the fitting algorithm.

2.2.4. Local and global parameters

When fitting the CD signal against temperature curves, all curves are simultaneously fitted, constrained to share the same values for ΔH_m (or $\Delta H_{1/2}$) and T_m (or $T_{1/2}$). The baseline and slopes are allowed to vary, and users also have the option to set the slopes (k_u , k_n) to zero.

2.2.5. Fitting errors

The standard deviation of all fitted parameters is computed using the square root of diagonal values from the fit parameter covariance matrix reported by `scipy.curve_fit` function. These values are an approximation (**underestimation**) of the real errors. Relative errors are calculated as $100 * (std(\hat{\theta}) / \hat{\theta})$ where $\hat{\theta}$ refers to the estimate of the parameter and $std(\hat{\theta})$ to the standard deviation.

2.3. Chemical unfolding

2.3.1. Equilibrium two-state and the linear extrapolation model

To analyse the CD signal as a function of a chemical denaturant agent concentration, ChiraKit allows fitting an equilibrium two-state model that assumes a linear dependence of stability on the denaturant concentration^{4,5}.

Briefly, the signal is expressed by the following equation:

$$Signal(D) = \frac{K_u}{1+K_u} (k_U D + b_U) + \frac{1}{1+K_u} (k_N D + b_N) \quad (10)$$

where D is the denaturant concentration, k_N , b_N are the slope and intercept of the pre-transition baseline (native state), and k_U and b_U are the slope and intercept of the post-transition (unfolded state) baseline.

Additionally, the unfolding constant K_u is defined as in Equation 2, but the free energy of unfolding depends on the empirical parameters M and $D50$:

$$\Delta G = M(D50 - D) \quad (11)$$

where $D50$ is the denaturant concentration at which 50% of the molecules are folded and M the slope $\frac{\partial \Delta G}{\partial D}$.

2.3.2. Equilibrium three-state and the linear extrapolation model

Briefly, we assume that the protein exists in the native (folded), intermediate or unfolded state and that there is an equilibrium between these three states given by the reactions $N \rightleftharpoons I$ (ΔG_1) and $I \rightleftharpoons U$ (ΔG_2). The signal is described by the following equation:

$$Signal(D) = f_u(k_U D + b_U) + f_n(k_N D + b_N) + f_i(b_i) \quad (12)$$

where D is the denaturant concentration, k_N , b_N are the slope and intercept of the pre-transition baseline (native), k_U and b_U are the slope and intercept of the post-transition (unfolded) baseline, and b_i is the baseline of the intermediate state. f_u , f_n and f_i are respectively the unfolded, folded and intermediate fractions, which are

⁴ Pace, C.N. and Hermans, J. (1975) The stability of globular protein. *CRC Crit. Rev. Biochem.*, **3**, 1–43.

⁵ Myers, J.K., Pace, C.N. and Scholtz, J.M. (1995) Denaturant m values and heat capacity changes: relation to changes in accessible surface areas of protein unfolding. *Protein Sci.*, **4**, 2138–2148.

calculated using Equations 7-9, with the difference that the equilibrium constants depend also on the denaturant concentration:

$$K_1(T, D) = e^{-M_1(D50_1 - D)/RT} \quad (13)$$

$$K_2(T, D) = e^{-M_2(D50_2 - D)/RT} \quad (14)$$

where $D50_1$ is the denaturant concentration at which f_n equals f_i , and $D50_2$ is the denaturant concentration at which f_i equals f_u . $M_{1/2}$ is the slope $\frac{\partial \Delta G_{1/2}}{\partial D}$.

2.3.3. Initial parameter estimates

The initial values of the parameters k_u , b_u , k_n and b_n are estimated by fitting the equation of a line to the first or last data points (2 molar range). For the two-state model, the initial guess for the parameter M is one, and for $D50$ the median of the used concentrations. For the three-state model, the initial value of $D50_1$ and $D50_2$ are respectively $D_{min} + 1$ and $D_{max} - 1$, where D_{min} and D_{max} are respectively the minimum and maximum denaturant agent concentrations. Be sure to check the fitted values, because $D50_1$ is constrained between $[D_{min} + 0.5, D_{max} - 2]$ and $D50_2$ is constrained between $[D_{min} + 0.5, D_{max} - 0.5]$. The initial values of $D50_1$ and $D50_2$ can be overridden by the user to 'guide' the fitting algorithm. M-values are constrained 0 and 20 kcal mol⁻¹ M⁻¹.

2.3.4. Local and global parameters

When fitting the CD signal versus chemical denaturant concentration curves, all curves are simultaneously fitted, constrained to share the same values for $D50$ (or $D50_{1/2}$) and M (or $M_{1/2}$). The baseline and slopes are allowed to vary, and users also have the option to set the slopes (k_u , k_n) to zero.

2.3.5. Fitting errors

See Section '2.2.5. Fitting errors'.

2.4. Custom analysis

2.4.1 The fitting function

The CD signal can also be analysed as a function of any experimental parameter. In this case, the fitting function is given by the user. The following rules are required:

- 1) The given experimental parameter should be included in the fitting function string. For example, if the experimental parameter name is 'T', then using 'T*z' as a fitting function implies that we fit the value of 'z' from the curve of the CD signal against the values of 'T'.
- 2) The function parameters can not contain mathematical characters; e.g., instead of 'dG_h2o', use 'dG_water'.
- 3) To use exponential, logarithm or root-square functions, you need to write 'e^(...)', 'log(...)', or 'sqrt(...)'. The parenthesis are obligatory.
- 4) Parameters with the pattern 'Global', e.g., 'TmGlobal' will be shared across all curves.
- 5) Function parameters with the pattern 'Pos' are constrained to be greater than zero. The opposite holds for the parameters with the pattern 'Neg'.

2.4.2. Initial parameter estimates

The initial estimates for the function parameters are obtained through a log-spaced grid search. The number of combinations to be tested is fixed, but the limits of the search space can be given by the user. The default behaviour is to explore values between 10^{-3} and 10^3 . The combination of initial values resulting in the lowest residual sum of squares is selected.

2.4.3. Local and global parameters

When fitting the CD signal against the experimental parameter values, all curves will be simultaneously fitted. The constraint is that they share the same values for parameters with the pattern 'Global'. As for all other parameters, they can vary independently for each curve.

2.5. Curve fitting and fitting errors

The curve fitting procedure is done by using the non-linear least squares method (`scipy.curve_fit`⁶). The standard deviation is computed using the square root of diagonal values from the fit parameter covariance matrix. For each parameter, the relative errors (in percentage) are computed as the quotient between the parameter standard deviation and the parameter value multiplied by one hundred.

⁶ https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.curve_fit.html

2.6. Comparison of spectra

To comprehend if there are differences between groups of spectra, ChiraKit allows computing averages and the associated standard deviations ('Module 2e. Spectra Comparison', in the app). Moreover, every possible difference spectrum and its associated standard deviation is also computed (through error propagation). ChiraKit provides the euclidean distances between all spectra and a box plot figure with the intra and inter group distributions.

To compare the shapes of spectra and discard the influence of differences in absorbance intensity, the L2 normalisation can be used. This results in the calculation of normalised euclidean distances as defined in the Equation 2 from Oyama, Taiji, *et al* (2022)⁷.

3. Spectra decomposition

To take advantage of all the spectral information, instead of monitoring a single wavelength, the loaded spectra can be decomposed into a set of basis spectra.

The CD values of a given spectrum are decomposed as a linear combination of basis spectra as follows:

$$CD(x) = c_1(x)\phi_1 + c_2(x)\phi_2 + c_3(x)\phi_3 + \dots + c_n(x)\phi_n \quad (15)$$

where each ϕ_i is a basis spectrum and c_i the associated coefficients. The basis spectra are orthogonal to each other and have unit norm. The variable x represents the dimension used for the measurement, such as temperature.

To find the set of basis spectra, we can apply singular value decomposition (SVD) directly on the CD data matrix, or SVD on the centred CD data matrix (equivalent to Principal component analysis)⁸.

3.1. Singular value decomposition (SVD)

SVD consist of factorizing the data matrix M of size $m \times n$ into the matrices U , Σ and V^T such that

⁷ Oyama, Taiji, *et al*. "Performance Comparison of Spectral Distance Calculation Methods." *Applied Spectroscopy* 76.12 (2022): 1482-1493.

⁸ Kim, Donggun, and Kisung You. "PCA, SVD, and Centering of Data." *arXiv preprint arXiv:2307.15213* (2023).

$$M = U\Sigma V^T \quad (16)$$

where U is an $m \times m$ orthogonal matrix, Σ is an $m \times n$ non-negative rectangular diagonal matrix, and V^T is an $n \times n$ orthogonal matrix. U contains the left singular vectors (basis spectra), V^T the right singular vectors and Σ the singular values.

The amount of explained variance associated to each singular value (s_i) is given by the following equation:

$$\text{explained variance}(s_i) = s_i^2 / \sum_{j=1}^n s_j^2 \quad (17)$$

The associated coefficients for the i -th basis spectrum and the j -th acquired spectrum is defined as:

$$\phi_{svd,i} \cdot M_j \quad (18)$$

where M_j is the j -th column (spectrum) of the matrix M .

3.2. Principal component analysis (PCA)

PCA is performed through eigendecomposition of the covariance matrix C :

$$C = VL V^t \quad (19)$$

where V is the matrix containing the eigenvectors (basis spectra) and L is the diagonal matrix with the eigenvalues (sorted in descending order). The amount of explained variance associated to each eigenvalues (λ_i) is given by the following equation:

$$\text{explained variance}(\lambda_i) = \lambda_i / \sum_{j=1}^n \lambda_j \quad (20)$$

It is important to remember that here we refer to the variance of the centred matrix. The associated coefficients for the i -th basis spectrum and the j -th acquired spectrum is defined as:

$$\phi_{pca,i} \bullet M_{meaned,j} \quad (21)$$

where $M_{meaned,j}$ is the j-th column (spectrum) of the centred matrix M .

3.3. Basis spectrum inversion

Any given basis spectrum ϕ_i can be simply inverted. As a result the explained variance will remain constant, the new spectrum ϕ_i' will be equal to $(-\phi_i)$ and the weighting function c_i' will be equal to $(-c_i)$.

3.4. Change of basis

Given a set of two or three basis spectra, we can rotate them to obtain a new set of basis spectra with the following properties: the first new basis spectrum is similar to the first acquired spectrum (e.g., lowest temperature, lowest denaturant concentration, or lowest 'experimental parameter'), the amount of explained variance remains constant.

Case 1) Two basis spectra

The new set of basis spectra is calculated as

$$\begin{bmatrix} | & , & | \\ \phi_1 & , & \phi_2 \\ | & , & | \end{bmatrix} \cdot rotM = \begin{bmatrix} | & , & | \\ \phi'_1 & , & \phi'_2 \\ | & , & | \end{bmatrix} \quad (22)$$

where ϕ_1 and ϕ_2 are the original first and second basis spectra, ϕ'_1 and ϕ'_2 are the new basis spectra, and $rotM$ is the rotation matrix given by

$$rotM = \begin{bmatrix} \cos(\alpha) & \sin(\alpha) \\ -\sin(\alpha) & \cos(\alpha) \end{bmatrix} \quad (23)$$

where α is the angle of a right-triangle between a catheti of length $c_{1,1}$ and the hypotenuse of length $\sqrt{c_{1,1}^2 + c_{2,1}^2}$. $c_{1,1}$ is the associated coefficient for ϕ_1 and the first acquired spectrum. $c_{2,1}$ is the associated coefficient for ϕ_2 and the first acquired spectrum.

Case 2) Three basis spectra

The new set of basis spectra is calculated as

$$\begin{bmatrix} | & , & | & , & | \\ \phi_1 & , & \phi_2 & , & \phi_3 \\ | & , & | & , & | \end{bmatrix} \cdot rotZ \cdot rotY = \begin{bmatrix} | & , & | & , & | \\ \phi'_1 & , & \phi'_2 & , & \phi'_3 \\ | & , & | & , & | \end{bmatrix} \quad (24)$$

where ϕ_1 , ϕ_2 and ϕ_3 are the original first, second and third basis spectra, ϕ'_1 , ϕ'_2 and ϕ'_3 are the new basis spectra, $rotZ$ is the rotation matrix given by

$$rotZ = \begin{bmatrix} \cos(\gamma) & \sin(\gamma) & 0 \\ -\sin(\gamma) & \cos(\gamma) & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (25)$$

where γ is the angle of a right-triangle between a catheti of length $c_{1,1}$ and the hypotenuse of length $\sqrt{c_{1,1}^2 + c_{2,1}^2}$. Here, $c_{i,1}$ refers to the associated coefficient for i -th original basis spectrum and the first acquired spectrum. Then, $rotY$ is the rotation matrix given by

$$rotY = \begin{bmatrix} \cos(\beta) & 0 & \sin(\beta) \\ 0 & 1 & 0 \\ -\sin(\beta) & 0 & \cos(\beta) \end{bmatrix} \quad (26)$$

where β is the angle of a right-triangle between a catheti of length $\sqrt{c_{1,1}^2 + c_{2,1}^2}$ and the hypotenuse of length $\sqrt{c_{1,1}^2 + c_{2,1}^2 + c_{3,1}^2}$.

3.5. Alternative formulation for the explained variance

Given a set of orthogonal vectors (basis spectra) and the associated coefficients, the amount of explained variance for a given vector (basis spectrum) is calculated as:

$$\frac{\|c_i\|^2}{\|\phi_i\|^2} / \|M\|^2 \quad (27)$$

where ϕ_i is the i -th basis spectrum, c_i the associated coefficients, and M the data matrix. In the case of PCA, M is replaced by M_{meaned} .

4. Protein secondary structure estimation

Proteins containing different amounts of secondary structure elements, such as alpha helices, beta sheets, and random coils, produce different CD spectra⁹. Given a certain query spectrum and a reference set, the Selcon algorithm can be used to calculate the secondary structure content.¹⁰

4.1. The Selcon method

The Selcon algorithm can be divided into three versions (1, 2 and 3). ChiraKit will try to run Selcon 3 and if it fails it will run Selcon2.

Selcon version 1

Briefly, we start with a set of linear equations that relate the CD spectra with the secondary structure content:

$$F = XC \quad (28)$$

where C is the $m \times n$ matrix containing the CD spectra in delta epsilon (mean residue molar extinction) units, F is the $l \times n$ matrix containing the known secondary structure elements, and X is the $l \times m$ matrix that allows the transition between the space of secondary structure elements F and the space of CD spectra C . Here, n , l and m represent, respectively, the number of proteins, secondary structure elements and wavelengths.

Let c be the query spectrum, an initial guess for the secondary structure components is done based on the most similar spectrum from the reference set, and we build the set of equations

$$F_2 = XC_2 \quad (29)$$

where C_2 now includes c , has $(n + 1)$ columns, and is ordered according to the root mean square (RMS) distance to c . F_2 is also rearranged accordingly. The Equation

⁹ Greenfield, Norma J. "Using circular dichroism spectra to estimate protein secondary structure." Nature protocols 1.6 (2006): 2876-2890.

¹⁰ Sreerama, Narasimha, Sergei Yu Venyaminov, and Robert W. Woody. "Estimation of protein secondary structure from circular dichroism spectra: inclusion of denatured proteins with native proteins in the analysis." Analytical biochemistry 287.2 (2000): 243-251.

is solved for X and then the secondary structure f (of the protein to be analysed) is calculated as

$$f = Xc \quad (30)$$

To solve the Equation regarding F_2 , C_2 is factorized using the SVD algorithm. Indeed, by varying the number of proteins from the reference set (from 2 to n) and the number of relevant singular values (from 1 to 7), several solutions are obtained for f . Those satisfying the following criteria are considered: 1) the sum of the secondary structure elements' fractions lies between 95% and 105%, and 2) each secondary structure element fraction is greater than - 2.5 %. If no acceptable solutions are found during this step, the rules are relaxed by adding a tolerance of 1 % to the total sum of fractions, and of 0.5 % to the individual fractions.

For each set of valid solutions that share a certain number of proteins from the reference set but vary in the number of relevant singular values, only the solution whose sum of secondary structure element fractions is closest to 100% is kept. Let N_p be the number of valid solutions, we will have at most $2 \leq N_p \leq n$ solutions.

The N_p solutions are averaged to obtain f_1 , which replaces the initial guess in the matrix F_2 as the second approximation. Finally, the same process is repeated until the RMS difference between two successive solutions is less than 0.0025 delta epsilon units.

Selcon version 2

The Selcon Version 2 method starts with the same algorithm as Selcon Version 1. After self-consistency is achieved, the final N_p set of valid solutions are filtered by removing those where the RMS between the reconstructed spectrum and the query spectrum c is larger than 0.25 delta epsilon units. If the filter does not leave any solution, the tolerance is relaxed using 0.01 delta epsilon units steps until at least one solution is found. Finally, the mean solution of the filtered solutions is reported.

Selcon version 3

The Selcon Version 3 method follows the same procedure as Selcon Version 2 with the difference that one more rule is applied. The N_p set of valid solutions is further filtered based on the amount of alpha helix.

Let h_i be the alpha-helix content fraction of the i -th solution, h_q the predicted alpha-helix content fraction for the query spectrum (after the Selcon 2 method); and h_{min} , h_{ave} and h_{max} the minimum, average, and maximum fraction of alpha-helix content, respectively, considering all the reference proteins required to obtain the N_p

solutions. Then, for the i -th solution to be valid, it needs to satisfy any of these criteria:

$$\text{Case 1)} \quad h_q > 0.65 \ \& \ h_i > 0.65 \quad (31)$$

$$\text{Case 2)} \quad 0.25 \leq h_q \leq 0.65 \ \& \ \frac{h_q + h_{max}}{2} - 0.03 \leq h_i \leq \frac{h_q + h_{max}}{2} + 0.03 \quad (32)$$

$$\text{Case 3)} \quad 0.15 \leq h_q \leq 0.25 \ \& \ \frac{h_q + h_{ave}}{2} - 0.03 \leq h_i \leq \frac{h_q + h_{ave}}{2} + 0.03 \quad (33)$$

$$\text{Case 4)} \quad h_q < 0.15 \ \& \ \frac{h_q + h_{min}}{2} - 0.03 \leq h_i \leq \frac{h_q + h_{min}}{2} + 0.03 \quad (34)$$

4.2. Reference sets

By default, two reference sets are provided: AU_SP175 and 2) AU_SM180. These sets contain CD spectra from the same group of proteins as the original SP175¹¹ and SMP180¹² sets. AU_SP175 contains the CD spectra of 71 soluble proteins measured down to 175 nm. AU_SM180 contains the AU_SP175 set and 57 additional CD spectra, including membrane proteins. The lower limit of the AU_SM180 set is 180 nm. In both cases, the data sets have 1 nm steps and start at 240 nm.

The number of secondary structure elements to be detected depends on the reference set (or the lower wavelength limit). Six secondary structure elements for values lower than 180 nm (AU_SP175), and four secondary structure elements for values between 180 and 190 nm (AU_SM180). **No calculation is possible if the query spectrum does not go to at least 190 nm.**

To select the reference data, the CD spectra were downloaded from the Protein Circular Dichroism Data Bank (PCDDDB)¹³. In case that many versions were available, the latest one was chosen (September 2023). In two cases we did not use the PCDDDB data directly: The Human Serum Albumin CD spectrum was replaced with an in-house measured spectrum. The Lysozyme spectrum CD0000045100 was zeroed again and averaged with the spectrum CD0000045000.

Custom reference sets can be imported too. The requisite is to import the corresponding matrices F (secondary structure elements) and C (reference spectra).

¹¹ Lees, Jonathan G., et al. "A reference database for circular dichroism spectroscopy covering fold and secondary structure space." *Bioinformatics* 22.16 (2006): 1955-1962.

¹² Abdul-Gader, Ali, Andrew John Miles, and Bonnie A. Wallace. "A reference dataset for the analyses of membrane protein secondary structures and transmembrane residues using circular dichroism spectroscopy." *Bioinformatics* 27.12 (2011): 1630-1636.

¹³ <https://pcddb.cryst.bbk.ac.uk/>

C is arranged column-wise (one spectrum per column) and the CD data goes from the highest wavelength to the lowest wavelength. F should have the same columns as C, and as many rows as different secondary structure elements.

4.3. Secondary structure elements

All secondary structures elements come from the DSSP routine¹⁴. Following the method of Sreerama, Narasimha et al., 1999, the secondary structure elements are further categorised as regular or distorted alpha-helix, regular or distorted beta-sheet, turns and 'other'. Alpha helix are the sum of fractions H and G. Beta sheets are the fraction E. Turns are fraction T. All other fractions (I, B, S, C) are assigned to 'other'. The alpha helix fraction is divided into two groups: Alpha_D (distorted) which are the two first and two last residues in a helix and the rest are Alpha_R (regular). If the helix is less than 4 residues all go into the Alpha_D fraction. The same division is made for beta sheets, with the difference that only the first and the last residue of a beta strand goes into the distorted fraction Beta_D. The rest goes into the regular fraction Beta_R.

For a given spectrum with an unknown associated secondary structure, if the lower wavelength range value is below 180 nm, the six secondary structure elements can be detected. If the lower limit is between 180 and 190 nm (inclusive), four different fractions will be estimated: Alpha, Beta, Turns and Other (Alpha = Alpha_R + Alpha_D, Beta = Beta_R + Beta_D).

Contact details

For further assistance, please contact us:

 spc@embl-hamburg.de

 EMBL (c/o DESY), Notkestrasse 85, Build. 25a, 22607 Hamburg, Germany

¹⁴ Kabsch, Wolfgang, and Christian Sander. "Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features." *Biopolymers: Original Research on Biomolecules* 22.12 (1983): 2577-2637.