Structure factor statistics and likelihood

Randy J Read





Principle of maximum likelihood

- Best model is most consistent with data
- Measure consistency by probabilities
- Optimise model by adjusting parameters in probability distribution
- Crystallographic likelihood is based on probability distributions of structure factors
 - univariate for maps, molecular replacement, refinement
 - multivariate for experimental phasing and other advanced applications

Structure factors with random atoms

- Assume atoms randomly scattered relative to Bragg planes
- Random walk in complex plane



Wilson distribution

+

The Central Limit Theorem

- Probability distribution of a sum of independent random variables tends to be Gaussian
 - regardless of distributions of variables in sum
- Conditions:
 - sufficient number of independent random variables
 - none may dominate the distribution
- Centroid (mean) of Gaussian is sum of centroids
- Variance of Gaussian is sum of variances

Wilson distribution for space group P1

- Apply central limit theorem to real and imaginary parts of structure factor separately
 - sums of real and imaginary atomic contributions

Derivation of Wilson distribution: centroids

$$\mathbf{F} = \sum_{j=1}^{Natom} f_j \exp(2\pi i \mathbf{h} \cdot \mathbf{x}_j) = \sum_{j=1}^{Natom} f_j \cos(2\pi \mathbf{h} \cdot \mathbf{x}_j) + i \sum_{j=1}^{Natom} f_j \sin(2\pi \mathbf{h} \cdot \mathbf{x}_j)$$
$$= A + iB$$
$$\langle \mathbf{F} \rangle = \langle A + iB \rangle = \langle A \rangle + i \langle B \rangle$$
$$\langle A \rangle = \sum_{j=1}^{Natom} \langle f_j \cos(2\pi \mathbf{h} \cdot \mathbf{x}_j) \rangle$$
$$= 0, \text{ (assume random position relative to Bragg planes)}$$
$$\langle B \rangle = \sum_{j=1}^{Natom} \langle f_j \sin(2\pi h \cdot x_j) \rangle = 0$$
$$\langle \mathbf{F} \rangle = \mathbf{0}$$

Derivation of Wilson distribution: variances

$$\left\langle \left| \mathbf{F} - \mathbf{0} \right|^2 \right\rangle = \left\langle A^2 + B^2 \right\rangle = \left\langle A^2 \right\rangle + \left\langle B^2 \right\rangle$$

$$\left\langle A^2 \right\rangle = \sum_{j=1}^{Natom} \left\langle \left(f_j \cos\left(2\pi h \cdot x_j\right) \right)^2 \right\rangle, \text{ (assume atoms uncorrelated)}$$

$$= \frac{1}{2} \sum_{j=1}^{Natom} f_j^2, \text{ (assume random position relative to Bragg planes)}$$

$$= \sum_N / 2$$

$$\left\langle B^2 \right\rangle = \sum_{j=1}^{Natom} \left\langle \left(f_j \sin\left(2\pi h \cdot x_j\right) \right)^2 \right\rangle = \sum_N / 2$$

$$\left\langle |F|^2 \right\rangle = \sum_N$$

Derivation of Wilson distribution: joint distribution of *A* and *B*

 $p(A) = \frac{1}{\sqrt{\pi \Sigma_N}} \exp\left(-\frac{A^2}{\Sigma_N}\right), \text{ (Gaussian with mean of 0, variance of } \Sigma_N / 2)$ $p(B) = \frac{1}{\sqrt{\pi \Sigma_N}} \exp\left(-\frac{B^2}{\Sigma_N}\right)$ $p(\mathbf{F}) = p(A, B) = \frac{1}{\pi \Sigma_N} \exp\left(-\frac{A^2 + B^2}{\Sigma_N}\right) = \frac{1}{\pi \Sigma_N} \exp\left(-\frac{|\mathbf{F}|^2}{\Sigma_N}\right)$

Alternative derivation of Wilson distribution

- Complex normal distribution
 - Gaussian for complex numbers
 - joint distribution of real and imaginary parts
 - Central Limit Theorem also applies

$$\langle \mathbf{F} \rangle = \sum_{j=1}^{Natom} \langle f_j \exp(2\pi i \mathbf{h} \cdot \mathbf{x}_j) \rangle = \mathbf{0}$$



$$\left\langle \left| \mathbf{F} \right|^{2} \right\rangle = \left\langle \mathbf{F} \mathbf{F}^{*} \right\rangle = \sum_{j=1}^{Natom} \left\langle f_{j} \exp\left(2\pi i \mathbf{h} \cdot \mathbf{x}_{j}\right) f_{j} \exp\left(-2\pi i \mathbf{h} \cdot \mathbf{x}_{j}\right) \right\rangle = \sum_{j=1}^{Natom} f_{j}^{2} = \Sigma_{N}$$
$$p(\mathbf{F}) = \frac{1}{\pi \Sigma_{N}} \exp\left(-\frac{\left| \mathbf{F} \right|^{2}}{\Sigma_{N}}\right), \text{ (complex Gaussian with mean of 0, variance of } \Sigma_{N} \right)$$

Sim distribution: known and unknown atoms



Sim distribution for amplitudes

- Likelihood function is the probability of the observations
 - but only the intensity (or amplitude) is measured
 - phase component has to be eliminated
- Change variables from real and imaginary to amplitude and phase
- Integrate over all possible values of (unknown) phase

Sim distribution for amplitudes

$$p(\mathbf{F}_{N}) = \frac{1}{\pi \Sigma_{Q}} \exp\left(-\frac{|\mathbf{F}_{N} - \mathbf{F}_{P}|^{2}}{\Sigma_{Q}}\right), \, dA \, dB = F \, dF \, d\alpha$$

$$p(F_{N}, \alpha_{N}) = \frac{F_{N}}{\pi \Sigma_{Q}} \exp\left(-\frac{F_{N}^{2} + F_{P}^{2} - 2F_{N}F_{P}\cos(\alpha_{N} - \alpha_{P})}{\Sigma_{Q}}\right)$$
Jacobian
$$p(F_{N}) = \int_{0}^{2\pi} p(F_{N}, \alpha_{N}) d\alpha$$

$$= \frac{F_{N}}{\pi \Sigma_{Q}} \exp\left(-\frac{F_{N}^{2} + F_{P}^{2}}{\Sigma_{Q}}\right)\int_{0}^{2\pi} \exp\left(\frac{2F_{N}F_{P}\cos(\alpha_{N} - \alpha_{P})}{\Sigma_{Q}}\right) d\alpha$$

$$= \frac{2F_{N}}{\Sigma_{Q}} \exp\left(-\frac{F_{N}^{2} + F_{P}^{2}}{\Sigma_{Q}}\right)I_{0}\left(\frac{2F_{N}F_{P}}{\Sigma_{Q}}\right)$$

Effect of atomic errors (or differences)

- Atomic errors give "boomerang" distribution of possible atomic contributions
- Portion of atomic contribution is correct



Effect of atomic errors (or differences)

 Work out centroid and variance for the contribution from a single atom



Structure factor with coordinate errors

- Luzzati (1952)
 - all atoms subject to same errors
- Complex normal distribution





Srinivasan and colleagues: σ_A

- Effects of errors and incompleteness are equivalent in terms of E-values
 - variance of error in E-value is $1-\sigma_A^2$
 - conservation of scattering power



Advanced applications of likelihood

- Refinement and MR involve pairs of structure factors with one observation: \mathbf{F}_{O} and \mathbf{F}_{C}
- Other applications involve larger collections
 - MIR: \mathbf{F}_{P} , \mathbf{F}_{PH1} , \mathbf{H}_{1} , \mathbf{F}_{PH2} , \mathbf{H}_{2} , ...
 - SAD: F⁺, F⁻, H⁺, H⁻
- Need joint distributions of collections of structure factors

Likelihood for more structure factors

- σ_A can be interpreted as complex covariance between normalised structure factors
- New applications based on multivariate complex normal distribution
 - difficulty is integrating out more than one phase!
 - can always isolate one phase by factoring probability distribution
 - see paper on SAD likelihood target

The normal (Gaussian) distribution

Gaussian distribution for one variable

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(x-\langle x \rangle)^2}{2\sigma^2}\right)$$
$$= \frac{1}{(2\pi\Sigma)^{1/2}} \exp\left[-\frac{1}{2}(x-\langle x \rangle)\Sigma^{-1}(x-\langle x \rangle)\right]$$

Multivariate normal distribution

$$\mathbf{p}(\mathbf{x}) = \frac{1}{|2\pi\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \langle \mathbf{x} \rangle)^T \Sigma^{-1} (\mathbf{x} - \langle \mathbf{x} \rangle)\right], \text{ where }$$

elements of Σ given by $\sigma_{ij} = \langle (x_i - \langle x_i \rangle) (x_j - \langle x_j \rangle) \rangle$

Multivariate complex normal distribution

- Complex normal $p(\mathbf{z}_{1}) = \frac{1}{\pi \Sigma} \exp\left[-\frac{|\mathbf{z}_{1} - \langle \mathbf{z}_{1} \rangle|^{2}}{\Sigma}\right]$ $= \frac{1}{\pi \Sigma} \exp\left[-\left(\mathbf{z}_{1} - \langle \mathbf{z}_{1} \rangle\right)^{*} \Sigma^{-1}\left(\mathbf{z}_{1} - \langle \mathbf{z}_{1} \rangle\right)\right]$ Re
- Multivariate complex normal distribution
 - Hermitian covariance matrix

$$\mathbf{p}(\mathbf{z}) = \frac{1}{|\mathbf{\pi}\mathbf{\Sigma}|} \exp\left[-\left(\mathbf{z} - \langle \mathbf{z} \rangle\right)^{H} \mathbf{\Sigma}^{-1}\left(\mathbf{z} - \langle \mathbf{z} \rangle\right)\right], \text{ where }$$

elements of
$$\Sigma$$
 given by $\sigma_{ij} = \langle (\mathbf{z}_i - \langle \mathbf{z}_i \rangle) (\mathbf{z}_j - \langle \mathbf{z}_j \rangle)^* \rangle$

Deriving conditional Gaussian probability

- Start with large joint distribution $p(x_1, x_2, \dots, x_m, y_1, y_2, \dots, y_n)$
- Fix known or model (y) terms
 - partition covariance matrix $\begin{vmatrix} \text{data} - \text{data} & \text{data} - \text{model} \\ \text{data} - \text{model} & \text{model} - \text{model} \end{vmatrix} = \begin{vmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{vmatrix}$ update covariances and expected values for remaining variables $\Sigma'_{11} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$ $\begin{vmatrix} \langle x_n \rangle \\ \langle x_m \rangle \end{vmatrix} = \Sigma_{12} \Sigma_{22}^{-1} \begin{bmatrix} y_1 \\ \vdots \\ \langle x_n \rangle \\ y_n \end{vmatrix}$

$$\Sigma_{11}' = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Re-deriving Srinivasan distribution

- Start from joint distribution of \mathbf{E}_O and \mathbf{E}_C
- Fix E_C, manipulate covariance matrix
 - turn joint distribution into conditional

$$\begin{bmatrix} \langle \mathbf{E}_{o} \mathbf{E}_{o}^{*} \rangle & \langle \mathbf{E}_{o} \mathbf{E}_{c}^{*} \rangle \\ \langle \mathbf{E}_{o} \mathbf{E}_{c}^{*} \rangle^{*} & \langle \mathbf{E}_{c} \mathbf{E}_{c}^{*} \rangle \end{bmatrix} = \begin{bmatrix} 1 & \sigma_{A} \\ \sigma_{A} & 1 \end{bmatrix} \longrightarrow \begin{cases} \langle \mathbf{E}_{o} \rangle_{\mathbf{E}_{c}} = \sigma_{A} \mathbf{E}_{C} \\ \operatorname{var}(\langle \mathbf{E}_{o} \rangle_{\mathbf{E}_{c}}) = 1 - \sigma_{A}^{2} \end{cases}$$

Molecular replacement with an ensemble

- Start from joint distribution of structure factors, including true and all models
- Conditional distribution turns collection of models into a statistically-weighted ensembleaverage structure factor

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{N} & D_{1}\boldsymbol{\Sigma}_{P_{1}} & \cdots & D_{n}\boldsymbol{\Sigma}_{P_{n}} \\ D_{1}\boldsymbol{\Sigma}_{P_{1}} & \boldsymbol{\Sigma}_{P_{1}} & \cdots & \left\langle \mathbf{F}_{C_{1}}\mathbf{F}_{C_{n}}^{*} \right\rangle \\ \vdots & \vdots & \ddots & \vdots \\ D_{n}\boldsymbol{\Sigma}_{P_{n}} & \left\langle \mathbf{F}_{C_{1}}^{*}\mathbf{F}_{C_{n}} \right\rangle & \cdots & \boldsymbol{\Sigma}_{P_{n}} \end{bmatrix}$$

SAD: probabilities for Friedel pair

- Start from joint distribution of true and calculated structure factors
- Use standard manipulations to get conditional probability: $p(\mathbf{F}_{o}^{+}, \mathbf{F}_{o}^{-}, \mathbf{H}^{+}, \mathbf{H}^{-}) \rightarrow p(\mathbf{F}_{o}^{+}, \mathbf{F}_{o}^{-}; \mathbf{H}^{+}, \mathbf{H}^{-})$ $\left[\begin{array}{c|c} \Sigma_{N} & \langle \mathbf{F}_{o}^{+}\mathbf{F}_{o}^{-} \rangle & \langle \mathbf{F}_{o}^{+}\mathbf{H}^{+*} \rangle & \langle \mathbf{F}_{o}^{+}\mathbf{H}^{-} \rangle \\ \langle \mathbf{F}_{o}^{+}\mathbf{F}_{o}^{-} \rangle^{*} & \Sigma_{N} & \langle \mathbf{F}_{o}^{-}\mathbf{H}^{+} \rangle^{*} & \langle \mathbf{F}_{o}^{-*}\mathbf{H}^{-} \rangle \\ \hline \langle \mathbf{F}_{o}^{+*}\mathbf{H}^{+} \rangle & \langle \mathbf{F}_{o}^{-}\mathbf{H}^{+} \rangle & \Sigma_{H} & \langle \mathbf{H}^{+}\mathbf{H}^{-} \rangle \\ \hline \langle \mathbf{F}_{o}^{+}\mathbf{H}^{-} \rangle^{*} & \langle \mathbf{F}_{o}^{-}\mathbf{H}^{-*} \rangle & \langle \mathbf{H}^{+}\mathbf{H}^{-} \rangle^{*} & \Sigma_{H} \end{array}\right]$



Dealing with translational NCS

- Diffraction from copies in different orientations is uncorrelated
 - zero covariances



Dealing with translational NCS

- Diffraction from copies in different orientations is uncorrelated
 - zero covariances
- Diffraction from copies in the same orientation is correlated
 - covariances are modulated
- Add covariances to get
 expected intensity



Effect of rotation on translational NCS

• Rotation parallel to diffraction vector has no effect



Effect of rotation on translational NCS

- Rotation parallel to diffraction vector has no effect
- Rotation around other axes reduces correlation
- Random coordinate differences between copies reduce correlation



Accounting for translational NCS

- Model effect of translation combined with small rotation and random differences between copies
- Also model vs. data covariances





Other potential applications

- SIR likelihood target
 - SIR phasing
 - joint refinement of native and liganded structures
- Fast translation function for SAD target
- Understanding of solvent flattening

Acknowledgements

- Phaser: Airlie McCoy, Laurent Storoni, Gábor Bunkóczi, Rob Oeffner
- Hyp-1: Mariusz Jaskolski, Joanna Sliwiak, Zbyszek Dauter

