## 2.2. QUALITY INDICATORS

If the objective is to estimate electron-density distributions, exclusion of large numbers of weak intensity terms beyond a limiting resolution is unlikely to be of any significant impact except, perhaps, a favourable one in computation time. If a sparse population of more intense reflections is also excluded, the effect may also be positive by reducing aberrations that might interfere with interpretation. If the objective is refinement, the benefit is less clear, except possibly in computation time. The cost for that may be the exclusion of a few more intense terms that provide positive guidance to refinement. Another positive effect of a sharply defined limiting resolution might be improved estimation of the spherical interference function where it is needed. It appears that, in practice, the consensus is that, in applying a cutoff, the exclusion of a few intense terms is at worst of negligible impact overall.

Because many imposed limits that exclude data are dependent on estimates of standard uncertainty, procedures for estimation of standard uncertainty need to be considered. As suggested earlier, a diffraction data set without a standard uncertainty for each intensity measurement is certainly the exception in current practice. The methods used by data-processing programs to estimate these standard uncertainties may be difficult, even impossible, to discern, but at core they must all be based on counting statistics. With the continuing trend toward data sets of high multiplicity (also referred to as redundancy), however, estimates of standard uncertainties are available from distributions of replicate measurements about means. Both of these estimates have value and express impacts of somewhat different sources of error. It would seem therefore that, where applicable, the best way to accommodate both would be to calculate weighted average values of individual intensities, where the weights are derived from the standard uncertainties from data processing, and then to estimate the standard uncertainties for the weighted averages by application of standard propagation-of-error procedures. This, finally, is the third area we highlight here that deserves focused discussion within the MX community with the objective of defining consensus best practices.

Little remains to be included in these introductory preliminaries. It should be understood that most of these indicators of quality may be cast in terms of either structure factors or intensities. We draw little attention to this as we define individual indicators, except that, where one of the two possibilities seems to dominate in usage, we tend to define that form alone. It is also true that many of these indicators have counterparts in which individual reflection terms are weighted. In current practice, these forms find little use in application to biological macromolecules and we ignore them here. Finally, many of these indicators may be expressed as fractions, that is, as numbers between one and zero. They are also often expressed in the percentage form, that is, as numbers between one hundred and zero. While we confine our use to the former here, we make no statement of preference.

The remainder of this chapter consists primarily of individual sections that reflect the various steps from the crystallographic data-collection experiment to the refinement of the final model. Relevant quality indicators and definitions are given for each of the steps and the most commonly used are collected in Table 2.2.11.1.

### 2.2.2. Quality indicators for diffraction data

Once useful crystalline samples have been obtained, the collection of X-ray diffraction data is the next (and the last) experi-

mental step in a structure determination. Although the greatest care may be taken to collect data of as high quality as possible, there remain circumstances and influences that limit the quality of the data. Over time, many indicators have been defined to describe various aspects of diffraction data quality. The most important ones are discussed here.

**Nominal resolution, $d_{min}$.** The resolution of a diffraction data set describes the extent of measurable data and is calculated by Bragg's law [equation (2.2.2.1)] based on the maximum Bragg angle $2\theta$ included in the data set for a given data-collection wavelength $\lambda$. As discussed above, the nominal resolution is a limit set by the experimenters and is well known to be prone to subjective judgment. A number of suggestions have been made to reduce the subjectivity associated with this limit. One defines the limit as the resolution within which the intensities of a fraction of the unique reflections, for example 70%, are above a threshold, for example zero or three times their standard uncertainties. Another suggested limiting criterion, discussed further below, recommends that the nominal resolution be set as the midpoint of the resolution range of the shell at which the mean signal-to-noise ratio falls below 2.

$$d_{min} = \lambda/2\sin\theta. \qquad (2.2.2.1)$$

**True resolution, $d_{true}$.** The true resolution of a diffraction data set is defined as the minimum distance between two objects in a crystal that permits their images in the resultant electron-density map to be resolved. Often, $d_{true}$ is approximated as $d_{min}$.

To illustrate this crucial distance, represent two equivalent atoms by equal overlapping Gaussians. One might then consider that the distance between them that just permits distinguishing them as individual atoms might be the distance at which the electron-density value at the midpoint between the atoms drops to a value just below that at the positions of the atoms. For a normal distribution, this distance is $2\sigma$, twice the standard deviation of the distribution.

Another perspective is provided by the realization that, when a Fourier synthesis is terminated at a resolution cutoff $d_{min}$, successive spheres of negative and positive density of decreasing amplitude surround the maxima of positive density at atomic positions in that synthesis. It has been shown that the distance from the centre of such a maximum to the first zero is $0.715d_{min}$ (James, 1948), which is a useful estimate of the limiting distance between distinguishable features in the electron-density map or $d_{true}$. Similar estimates are used in other areas, notably for defining resolution in astronomy. A more recent re-evaluation suggests that a limit of $0.917d_{min}$ is a better value, especially when the effects of form factors and atomic displacement parameters are considered (Stenkamp & Jensen, 1984). Add to that the effects of errors in experimental amplitudes and derived phases, and the approximation of $d_{true}$ as $d_{min}$ seems quite reasonable.

**Optical resolution, $d_{opt}$.** The optical resolution $d_{opt}$ is calculated from the standard deviation of a Gaussian fitted to the origin peak of the Patterson function ($\sigma_{Patt}$) of the diffraction data set and the standard deviation of another Gaussian fitted to the origin peak of the spherical interference function ($\sigma_{sph}$). This definition is based on Vaguine *et al.* (1999) and is implemented in the program *SFCHECK*. The optical resolution is intended to account for uncertainties in the data, atomic displacement factors, effects of crystal quality and series-termination effects by means of a propagation-of-error-like approach (Blundell & Johnson, 1976; Vaguine *et al.*, 1999). It has been suggested that $d_{opt}$ is a better approximation of $d_{true}$ than

$d_{\mathrm{min}}$ (Weiss, 2001).

$$d_{\mathrm{opt}} = [2(\sigma_{\mathrm{Patt}}^2 + \sigma_{\mathrm{sph}}^2)]^{1/2}. \qquad (2.2.2.2)$$

**Completeness, *C*.** The completeness $C$ of a diffraction data set is defined as the fraction of the unique reflections in a given space group to a given nominal resolution $d_{\mathrm{min}}$ that have been measured at least once during data collection. $C$ may be given assuming that Friedel symmetry is either applied or not. In the latter case, $C$ is also referred to as anomalous completeness. In the program *SCALA* (Evans, 2006), the anomalous completeness is defined based on acentric reflections only.

**Effective resolution, $d_{\mathrm{eff}}$.** Since any missing reflection of a data set leads to a deterioration of the model parameters (Hirshfeld & Rabinovich, 1973; Arnberg *et al.*, 1979), an effective resolution may be defined based on the nominal resolution $d_{\mathrm{min}}$ and the cube root of the completeness $C$ of the data set.

$$d_{\mathrm{eff}} = d_{\mathrm{min}} C^{-1/3}. \qquad (2.2.2.3)$$

**Multiplicity (or redundancy), *N*.** The multiplicity or redundancy $N$ of a diffraction data set defines on average how many times a reflection *hkl* has been observed during the data-collection experiment including symmetry mates and replicate measurements. $N$ may be given assuming that Friedel symmetry is either applied or not.

**Merging *R* factor, $R_{\mathrm{merge}}$.** The merging $R$ factor of a diffraction data set describes the spread of the individual intensity measurements $I_i$ of a reflection *hkl* around the mean intensity $\langle I(hkl)\rangle$ of this reflection. Sometimes $R_{\mathrm{merge}}$ is also referred to as $R$ factor (observed) (in the program *XDS*; Kabsch, 1988, 1993, 2010; Chapter 11.6), as $R_{\mathrm{sym}}$ or as $R_{\mathrm{linear}}$ (in the program *SCALEPACK*; Otwinowski & Minor, 1997). In fractional form, this is

$$R_{\mathrm{merge}} = \sum_{hkl}\sum_i |I_i(hkl) - \langle I(hkl)\rangle| / \sum_{hkl}\sum_i I_i(hkl), \qquad (2.2.2.4)$$

where $\langle I(hkl)\rangle$ is the mean of the several individual measurements $I_i(hkl)$ of the intensity of reflection *hkl*. The sums $\sum_{hkl}$ and $\sum_i$ run over all observed unique reflections *hkl* and over all individual observations *i* of a given reflection *hkl*. It should be noted that alternative definitions of $R_{\mathrm{merge}}$ exist. In one, $I_i(hkl)$ in the denominator is replaced by $\langle I(hkl)\rangle$, thereby producing an expression that is formally equivalent to the one above. In another, $I_i(hkl)$ in the denominator is replaced by $|I_i(hkl)|$ with the suggestion that the denominator is thereby prevented from becoming negative or zero, even in the case of many negative-intensity observations. One should note, however, the counter-intuitive side effect: artificial damping of $R_{\mathrm{merge}}$ values, that is, reducing expected higher $R_{\mathrm{merge}}$ values of data sets with more weak reflections.

The usefulness of $R_{\mathrm{merge}}$ as a quality indicator for diffraction data is limited because it is dependent on the multiplicity of a data set (Diederichs & Karplus, 1997*a,b*; Weiss & Hilgenfeld, 1997; Weiss, 2001). The higher the multiplicity of a data set, the higher its $R_{\mathrm{merge}}$ will be, although, based on statistics, the better determined the averaged intensity values should be. Despite these shortcomings, $R_{\mathrm{merge}}$ is still widely used today.

**Redundancy-independent merging *R* factor, $R_{\mathrm{r.i.m.}}$ or $R_{\mathrm{meas}}$.** The redundancy-independent merging $R$ factor $R_{\mathrm{r.i.m.}}$ or $R_{\mathrm{meas}}$ describes the precision of the individual intensity measurements $I_i$, independent of how often a given reflection has been measured. Because of its independence of the redundancy (hence its name), it has been proposed that $R_{\mathrm{r.i.m.}}$ or $R_{\mathrm{meas}}$ should be used as a substitute for the conventional $R_{\mathrm{merge}}$ (Diederichs & Karplus,

1997*a,b*; Weiss & Hilgenfeld, 1997; Weiss, 2001). In fractional form, this is

$$\begin{aligned} R_{\mathrm{r.i.m.}} = \sum_{hkl}\big\{N(hkl)/[N(hkl)-1]\big\}^{1/2} \\ \times \sum_i |I_i(hkl) - \langle I(hkl)\rangle| / \sum_{hkl}\sum_i I_i(hkl), \end{aligned}$$
$$(2.2.2.5)$$

where $\langle I(hkl)\rangle$ is the mean of the $N(hkl)$ individual measurements $I_i(hkl)$ of the intensity of reflection *hkl*. As for $R_{\mathrm{merge}}$, the sums $\sum_{hkl}$ and $\sum_i$ run over all observed unique reflections *hkl* and over all individual observations *i* of a given reflection *hkl*.

**Precision-indicating merging *R* factor, $R_{\mathrm{p.i.m.}}$.** The precision-indicating merging $R$ factor $R_{\mathrm{p.i.m.}}$ describes the precision of the averaged intensity measurements $\langle I(hkl)\rangle$ (Weiss, 2001). In fractional form, this is

$$\begin{aligned} R_{\mathrm{p.i.m.}} = \sum_{hkl}\big\{1/[N(hkl)-1]\big\}^{1/2} \\ \times \sum_i |I_i(hkl) - \langle I(hkl)\rangle| / \sum_{hkl}\sum_i I_i(hkl), \end{aligned}$$
$$(2.2.2.6)$$

where $\langle I(hkl)\rangle$ is the mean of the $N(hkl)$ individual measurements $I_i(hkl)$ of the intensity of reflection *hkl*. As with $R_{\mathrm{merge}}$ and $R_{\mathrm{r.i.m.}}$ or $R_{\mathrm{meas}}$, the sums $\sum_{hkl}$ and $\sum_i$ run over all observed unique reflections *hkl* and over all individual observations *i* of a given reflection *hkl*.

***R* factor of merged intensities or amplitudes, $R_{\mathrm{mrgd}\text{-}I}$ and $R_{\mathrm{mrgd}\text{-}F}$.** An alternative precision-indicating merging $R$ factor, called $R_{\mathrm{mrgd}}$, is defined as the $R$ factor between two or more data sets or between two subsets of a data set created by randomly apportioning the individual intensity measurements between the two subsets (Diederichs & Karplus, 1997*a,b*). $R_{\mathrm{mrgd}}$ can be calculated for intensities ($R_{\mathrm{mrgd}\text{-}I}$) or structure-factor amplitudes ($R_{\mathrm{mrgd}\text{-}F}$). The latter quantity was suggested to present a lower limit for the crystallographic $R$ factor of a model against the observed data (Diederichs & Karplus, 1997*a,b*). In fractional form

$$R_{\mathrm{mrgd}\text{-}I} = 2\sum_{hkl} |\langle I_1(hkl)\rangle - \langle I_2(hkl)\rangle| / \sum_{hkl}\langle I_1(hkl)\rangle + \langle I_2(hkl)\rangle,$$
$$(2.2.2.7)$$

where $\langle I_1(hkl)\rangle$ and $\langle I_2(hkl)\rangle$ are the mean intensity values for the individual observations of the reflections *hkl*, which have been partitioned into the two subsets 1 and 2. The sums $\sum_{hkl}$ run over all observed unique reflections. $R_{\mathrm{mrgd}\text{-}I}$ is related to $R_{\mathrm{p.i.m.}}$ by a constant factor ($R_{\mathrm{mrgd}\text{-}I} = 2^{1/2} R_{\mathrm{p.i.m.}}$).

$R_{\mathrm{mrgd}\text{-}F}$ is defined analogously to $R_{\mathrm{mrgd}\text{-}I}$ (Diederichs & Karplus, 1997*a,b*). In the equation, only the intensities are replaced by structure-factor amplitudes.

$$R_{\mathrm{mrgd}\text{-}F} = 2\sum_{hkl} |\langle F_1(hkl)\rangle - \langle F_2(hkl)\rangle| / \sum_{hkl}\langle F_1(hkl)\rangle + \langle F_2(hkl)\rangle.$$
$$(2.2.2.8)$$

In order to cope with negative-intensity observations, pseudo-amplitudes had to be introduced just for the purpose of calculating $R_{\mathrm{mrgd}\text{-}F}$ ($F = I^{1/2}$ if $I \geq 0$ and $F = -|I|^{1/2}$ if $I < 0$).

*Note.* The approach of comparing randomly partitioned subsets of a given data set is used for a variety of quality indicators. While there is potential for variation in these indicators from one partitioning of the data set to another, an average of several random partitionings should be expected to give a useful estimate. There is also potential for subjectivity,

but the principal value of these indicators is to assist the experimenter in proper analysis and they are less often applied to compare experiments from different laboratories and are seldom published.

**Pooled coefficient of variation, PCV.** The pooled coefficient of variation PCV is the ratio of the sum of the standard deviations to the sum of the reflection intensities (Diederichs & Karplus, 1997a,b). PCV is related to $R_{meas}$ or $R_{r.i.m.}$ by the factor $(\pi/2)^{1/2}$. In fractional form, this is

$$\text{PCV} = \frac{\sum_{hkl} \left( \{1/[N(hkl)-1]\}^{1/2} \sum_i |I_i(hkl) - \langle I(hkl)\rangle|^2 \right)^{1/2}}{\sum_{hkl} \sum_i I_i(hkl)}, \tag{2.2.2.9}$$

where $\langle I(hkl)\rangle$ is the mean of the $N(hkl)$ individual measurements $I_i(hkl)$ of the intensity of reflection $hkl$.

**Mean signal-to-noise ratio, $\langle I/\sigma(I)\rangle$.** The signal-to-noise ratio $I_i/\sigma(I_i)$ of an individual intensity measurement describes the statistical significance of a measured intensity. As a measure of the overall quality of a data set, the mean signal-to-noise ratio for all reflections is useful as an indication of the robustness of the data, that is, the average intensity as a multiple of the standard uncertainty. In addition, as mentioned above, the mean signal-to-noise ratio for all reflections within the outer resolution shell can be used to define the nominal resolution of a data set. For the data set as a whole or for a resolution shell of that data set, the mean signal-to-noise ratio, $\langle I/\sigma(I)\rangle$, is the sum of the signal-to-noise ratios of all individual reflections $hkl$ within resolution limits divided by the number of individual reflections $hkl$ within those resolution limits.

In principle there are two ways to define a mean signal-to-noise ratio of a data set (or a given resolution shell). The two ways yield different quantities, although, unfortunately, they are both called the mean signal-to-noise ratio. They differ in the manner in which mean signal-to-noise ratios of individual reflections $hkl$ are calculated.

(i) $\langle I(hkl)\rangle/\sigma[I(hkl)]$. The mean signal-to-noise ratio of individual reflections $hkl$ may be calculated as the ratio of the mean intensity $\langle I(hkl)\rangle$ and the r.m.s. scatter of $I_i(hkl)$ about that mean. This is a measure of the average significance of individual observations, but it does not take into account the multiplicity or redundancy of the measurements.

$$\langle I(hkl)\rangle/\sigma[I(hkl)]$$
$$= \langle I(hkl)\rangle/[(1/N)\sum_i |I_i(hkl) - \langle I(hkl)\rangle|^2]^{1/2}. \tag{2.2.2.10}$$

In the program *SCALA* (Evans, 2006), this value is reported as I/sigma.

(ii) $\langle I(hkl)\rangle/\sigma\langle I(hkl)\rangle$. $\langle I(hkl)\rangle$ is the average over all observations of the reflection $hkl$, and is sometimes weighted. $\sigma\langle I(hkl)\rangle$ is the propagation-of-error combination of standard uncertainties assigned at data processing for the individual measurements $I_i(hkl)$, that is, a modification of equation (2.2.2.10) in which the term $|I_i(hkl) - \langle I(hkl)\rangle|$ in the denominator is replaced by $\sigma_i(hkl)$, the experimental standard uncertainty for the measurement $I_i(hkl)$.

$$\langle I(hkl)\rangle/\sigma\langle I(hkl)\rangle = \langle I(hkl)\rangle/[(1/N)\sum_i \sigma_i(hkl)^2]^{1/2}. \tag{2.2.2.11}$$

An error model[1] is often applied in the denominator here to scale to the r.m.s. scatter in (i) above. In the program *SCALA* (Evans, 2006), this value is reported as Mn(I)/sd.

Both methods of defining the mean signal-to-noise ratio for the reflection $hkl$ have merit. As suggested for individual intensities in Section 2.2.1, perhaps the best approach would be to calculate weighted averages and weighted standard uncertainties of the $I(hkl)$ where weights are the experimental standard uncertainties $\sigma_i(hkl)$ for individual measurements $I_i(hkl)$.

**Highest possible signal-to-noise ratio, $I/\sigma(I)_{asymptotic}$.** A relatively recent addition to the collection of diffraction-data quality indicators is the highest possible signal-to-noise ratio of a data set $I/\sigma(I)_{asymptotic}$ or ISa (Diederichs, 2010). ISa is calculated from the parameters of the error model used for inflating the standard deviations of the reflections with an intensity-dependent term.[1] Since ISa is practically independent of counting statistics, it was suggested to be a good measure of instrument errors manifesting themselves in the data set, provided the crystal is close to ideal and radiation damage is negligible. Data sets with ISa values of 25 or greater are considered to be very good and amenable to straightforward structure determination, while data sets exhibiting ISa values of 15 or less are considered marginal at best. The calculation of ISa is implemented in *XDS* versions of December 2009 or later (Kabsch, 2010).

**Anomalous $R$ factor, $R_{anom}$.** The anomalous $R$ factor $R_{anom}$ describes the sum of the differences in intensities of Friedel-related reflections $(hkl)$ and $(\bar{h}\bar{k}\bar{l})$ relative to the sum of their mean intensities. In fractional form, this is

$$R_{anom} = \sum_{hkl} |I(hkl) - I(\bar{h}\bar{k}\bar{l})|/\sum_{hkl} \langle I(hkl)\rangle, \tag{2.2.2.12}$$

where, in this case, $\langle I(hkl)\rangle$ is the mean intensity of the Friedel mates of the reflections $hkl$, or $\frac{1}{2}[I(hkl) + I(\bar{h}\bar{k}\bar{l})]$. Here, the sums $\sum_{hkl}$ run over all unique reflections with one of the indices, typically $h$, greater than zero ($h > 0$) for which both Friedel mates have been observed at least once.

The ratio of $R_{anom}$ to $R_{p.i.m.}$ has been proposed as a possible indicator for the strength of the anomalous signal (Panjikar & Tucker, 2002).

**Anomalous correlation coefficient, $CC_{anom}$.** The anomalous correlation coefficient $CC_{anom}$ quantifies the linear dependence of observed anomalous differences in two diffraction data sets. These can be data sets, for example, collected at two different wavelengths in a MAD experiment. In cases where only one data set is available, two randomly partitioned half data sets can be created for comparison.

*Note.* The correlation coefficient referred to here and elsewhere in this chapter is invariably the Pearson linear correlation coefficient (Rodgers & Nicewander, 1988):

$$CC = \sum(x - \langle x\rangle)(y - \langle y\rangle)\Big/\Big[\sum(x - \langle x\rangle)^2 \sum(y - \langle y\rangle)^2\Big]^{1/2}, \tag{2.2.2.13}$$

with $x$ and $y$ being, in this case, the anomalous differences $[I(hkl) - I(\bar{h}\bar{k}\bar{l})]$ or $[F(hkl) - F(\bar{h}\bar{k}\bar{l})]$ in the two data sets, $\langle x\rangle$ and $\langle y\rangle$ are their averages, and the summations are over all reflections

---

[1] Typically, the experimental standard uncertainties are modified by various correction factors in order to account for errors in the data that arise from other sources than counting statistics. The inflation factors are applied in various ways in different programs. In *SCALA* (Collaborative Computational Project, Number 4, 1994; Evans, 2006) they are called SDFAC, SDB and SDADD. In *SCALEPACK* (Otwinowski & Minor, 1997), they are called ERROR SCALE FACTOR and ESTIMATED ERROR. In *D*TREK* (Pflugrath, 1999), they are called $E_{mul}$ and $E_{add}$. In *XDS* (Kabsch, 1988, 1993, 2010; Chapter 11.6) they are called a and b.

*hkl* for which observations exist in both data sets across the entire resolution range or within a particular resolution shell. $CC_{anom}$ is a reliable indicator of the strength of the anomalous signal. Values above 0.30 are considered good.

**R.m.s. correlation ratio**. This is another statistic based on randomly partitioned data sets, which is calculated by the program *SCALA* (Evans, 2006; Collaborative Computational Project, Number 4, 1994). It is an analysis of the scatterplot of $[I_1(hkl) - I_1(\bar{h}\bar{k}\bar{l})]$ *versus* $[I_2(hkl) - I_2(\bar{h}\bar{k}\bar{l})]$, where the subscripts 1 and 2 identify the two half data sets. The analysis assumes that the correlation is ideally 1.0. The r.m.s. correlation ratio is defined as the ratio of the r.m.s. widths of the scatterplot distribution along the diagonal and perpendicular to the diagonal. This statistic seems to be more robust than $CC_{anom}$ to the presence of outliers. It cannot, however, be applied to analysing the correlations between different data sets.

**Mean anomalous signal-to-noise ratio, $\langle d''/\sigma(d'') \rangle$**. The anomalous signal-to-noise ratio of an individual reflection measurement $d''(hkl)/\sigma[d''(hkl)]$ is defined as the ratio of the observed anomalous intensity difference $d'' = |I(hkl) - I(\bar{h}\bar{k}\bar{l})|$ and the corresponding estimated standard uncertainty $\sigma(d'')$ in the measurement of this anomalous difference. The average of the anomalous signal-to-noise ratios for all reflections within a certain resolution range is used as an indicator of utility for phasing. A value of $(2/\pi)^{1/2} \simeq 0.8$ for mean $d''/\sigma(d'')$ of a resolution shell, for example, is taken to indicate that no anomalous signal is present (G. Sheldrick & G. Bunkoczi, personal communication).

**Decay $R$ factor, $R_d$**. The decay $R$ factor $R_d$ is defined as a pairwise $R$ factor based on the intensities of symmetry-related reflections occurring on different diffraction images (Diederichs, 2006). An increase in $R_d$ as a function of difference in image-collection times is a good indicator of radiation damage occurring during data collection. In fractional form, this is

$$R_d = 2 \sum_{hkl} \sum_{m-n} |I_m(hkl) - I_n(hkl)| / \sum_{hkl} \sum_i [I_m(hkl) + I_n(hkl)], \quad (2.2.2.14)$$

where $I_m(hkl)$ and $I_n(hkl)$ are the intensities of the reflection *hkl* occurring on images $m$ and $n$. The only program in which this is currently implemented is *XDSSTAT*.

**Wilson-plot $B$ factor, $B_{Wilson}$**. A Wilson plot (Wilson, 1949) is a plot for a contiguous series of resolution shells of the logarithm of the mean intensity in a given resolution shell divided by the sum of the squared atomic form factors for all atoms in the unit cell evaluated at the mean of the resolution limits of the shell. From a least-squares fit of a straight line to the linear part of the Wilson plot, the $B$ factor $B_{Wilson}$ can be derived. Typically, data of lower than 4.5 Å resolution are excluded from the fit. The more meaningful determinations of $B_{Wilson}$ come from Wilson plots that are linear all the way to the nominal resolution $d_{min}$ and minimize the occurrence of spikes due to ice rings.

$$\ln(\langle I_{obs} \rangle / \sum_i f_i^2) = -\ln K_{Wilson} - 2B_{Wilson}/d^2, \quad (2.2.2.15)$$

where $\langle I_{obs}(hkl) \rangle$ is the mean over the intensities of all observed reflections *hkl* in a given resolution shell. The sum $\sum_i$ runs over all atoms in the structure. The parameter $d$ is the midpoint of the resolution shell over which $I_{obs}$ has been averaged. $K_{Wilson}$ is an absolute scale factor.

### 2.2.3. Comparing different diffraction data sets

In MX, there are many instances when two or more diffraction data sets have to be compared to each other. An important decision which has to be made is, for instance, whether data sets collected from different crystals are isomorphous enough so that they can be merged together. Another task is the comparison of a native data set and a heavy-atom derivative. The question here is how much of the observed difference is due to non-isomorphism and how much is due to isomorphous replacement.

**Scaling $R$ factor, $R_{scale}$**. The scaling $R$ factor $R_{scale}$ between two data sets is defined as the difference in the structure-factor amplitudes of one data set relative to the structure-factor amplitudes of the other data set after the two data sets have been scaled to each other.

$$R_{scale} = \sum_{hkl} |F_1 - F_2| / \sum_{hkl} |F_1|. \quad (2.2.3.1)$$

However, the above formula is asymmetric with respect to data sets 1 and 2. An alternative symmetric formula is

$$R_{scale} = 2 \sum_{hkl} |F_1 - F_2| / \sum_{hkl} (F_1 + F_2). \quad (2.2.3.2)$$

$R_{scale}$ is also often given based on intensities rather than structure-factor amplitudes.

**Gradient from a normal probability analysis**. For each reflection *hkl* in two data sets that have been scaled to each other, the quantity $\Delta(real)$ is computed and compared with its expectation value $\Delta(expected)$ assuming a normal distribution of errors (Howell & Smith, 1992). If $\Delta(real)$ is plotted against $\Delta(expected)$, random differences between the two data sets give rise to a slope of 1.0. Significant differences result in slopes significantly larger than 1.0. Such an analysis is implemented in the program *SCALEIT* (Collaborative Computational Project, Number 4, 1994).

$$\Delta(real) = (F_1 - F_2) / [\sigma(F_1)^2 + \sigma(F_2)^2]^{1/2}. \quad (2.2.3.3)$$

### 2.2.4. Quality indicators for substructure determination

The first step of a crystal structure determination after processing and analysing the diffraction data is to use the observed isomorphous and/or anomalous differences in order to determine the isomorphous or anomalous substructure. Since such substructures typically consist of rather few atoms, which are relatively far apart, direct-methods approaches have been adopted from small-molecule crystallography and have largely superseded pure Patterson-based methods. The most widely used computer programs for substructure determination are *SHELXD* (Schneider & Sheldrick, 2002; Sheldrick, 2008), *SnB* (Weeks *et al.*, 1993, 1994) and *HySS* (Grosse-Kunstleve & Adams, 2003).

**Correlation coefficient CC(all)**. CC(all) measures the correlation of the $E$ values (analogous to normalized structure factors; Hauptman & Karle, 1953) derived from the observed isomorphous or anomalous differences and those calculated from the substructure model. CC(all) is another application of the Pearson linear correlation coefficient [see equation (2.2.2.13)] and is calculated for all reflections. In the equation, $x$ and $y$ are the $E$ values derived from the observed and calculated differences. CC(all) is typically used to assess the successful determination of the isomorphous or anomalous substructure in a macromolecule. A value of $CC(all) \geq 0.3$ often indicates that the substructure has been correctly identified (Schneider & Sheldrick, 2002).

**Correlation coefficient CC(weak)**. Like CC(all), CC(weak) is the Pearson linear correlation coefficient [see equation (2.2.2.13)]

**references**