# Spectral Analysis of High Resolution X-ray Binary Data

Michael Nowak, mnowak@space.mit.edu

X-ray Astronomy School; Aug. 1-5, 2011

## Introduction

This exercise takes a look at X-ray binary observations using the *Chandra* high energy transmission gratings spectrometer (HETGS). Specifically, we look at *Chandra*-HETG observations of GRO J1655−40, ObsID 5461 (Miller et al., 2008, ApJ, 680, 1359).

Whereas some fraction of this exercise can be performed in any of the major analysis systems (`XSPEC`, `Sherpa`, `ISIS`, `SPEX`), you will probably get the furthest using `ISIS` or `Sherpa`. The reasons for this are two-fold: 1) I'll be suggesting that you "bin and combine data on the fly" (i.e., grouping data *during analysis* and co-adding data *during analysis*), and 2) end off suggesting programming to perform Monte Carlo simulations of error bars. Those going the `XSPEC` route will have to create sets of files *before analysis* that represent each of the sets of binnings/data combinations (and work much harder to do the Monte Carlo simulations).

For those trying `ISIS`, you can find a helpful set of "startup files" ( files named `.isisrc*`) located at:

    `http://space.mit.edu/home/mnowak/isis_vs_xspec/download.html`

If these files are placed in your home directory, and the path variable in the main `.isisrc` file is edited to point to your home directory, then these will automatically be loaded when you start `ISIS`. They provide plotting functions a little nicer than the intrinsic functions in `ISIS`, as well as a number of helper functions.

## Obtaining the Data

Open up a web browser and go to *TGCat*, the *Transmission Gratings Catalog*, which can be accessed at:

    `http://tgcat.mit.edu`

There you can browse, plot, and download spectral products for all publicly available *Chandra* gratings observations.

Search for data associated with ObsID 5461, *or* search for the object GRO J1655−40, and select the data from the above ObsID. Before downloading the products, try plotting the data with different combinations of signal-to-noise criteria for the binning, and in different energy/wavelength ranges, and in different units.

You can download the data products as "Type 1" spectral files and/or "Type 2" spectral files (explained below). It is perhaps easiest to download both kinds. Be sure to get also the spectral response products for the data.

In what follows, I will be assuming the use of the "Type 2" products.

## Loading Gratings Data

There are two sets of *Chandra* transmission gratings: the High Energy Grating (HEG) and the Medium Energy Grating (MEG), each of which disperses in two directions away from the aimpoint (the negative

and positive dispersion orders). Furthermore, at any location along the dispersed spectra, one finds multiple dispersion orders corresponding to wavelengths $\lambda$, $\lambda/2$, $\lambda/3$, .... These orders are separated from one another using the energy resolution of the CCD. The standard spectral extraction routines typically create spectra for the first three orders of each set of gratings in each direction. That is, one extracts twelve spectra: HEG -3,-2,-1, 1, 2, 3, and MEG -3, -2, -1, 1, 2, 3. Rather than create twelve separate spectral files, all twelve spectra are often stored in a *single* FITS file, referred to as a "Type 2 PHA" file. For the case of the HETG spectra, the 12 spectra are stored in the order listed above. The advantage is that there is one file with all the associated spectra. The disadvantage is that there aren't standard protocols for storing the information about the names of the associated ARF and RMF files for each spectrum.

Reading such a PHA2 file is not a problem for any of the analysis packages. One just has to make sure to also read the proper ARF and RMF files, and then associate them with the correct spectra. Here we will work with just the first order spectra (both positive and negative dispersion orders), since they contain the vast majority of counts, and are also the best calibrated of the spectra.

1. Read the HEG $\pm1$ and MEG $\pm1$ spectra from the PHA2 file (the third, fourth, ninth, and tenth spectra in the PHA2 file). Also read the ARF and RMF files, and then associate them with the data.

2. Plot the data (as counts/bin, then as counts/Å/sec). Since the gratings disperse *linearly* in wavelength, and the spectra have constant width wavelength bins, first plot the spectra vs. Å. The useful range of the *Chandra* HETG is $\approx$ 1.5–30 Å. Try plotting different wavelength ranges. How do each of the individual gratings compare to one another? Which has the most counts at low energies/long wavelengths? Which has the most counts at high energies/short wavelengths? Which has the greatest spectral resolution?

For `Sherpa` users, a useful thread to use as a reference is:
    http://asc.harvard.edu/sherpa/threads/grating/
while `ISIS` users can follow the tutorial at:
    http://space.mit.edu/home/mnowak/isis_vs_xspec/

## Combining Data

To combine or not combine data? In principle, if one uses the proper statistical tests, there isn't any real advantage to combining data. However, combining data might allow one to raise the counts/bin sufficiently to use $\chi^2$ statistics, it might serve the purposes of "averaging over" systematic deviations from one observation to another (or in this case, among the four different dispersed spectra), and it reduces the computational time. (The model is evaluated once, by default, rather than four times.) Combined data might also be easier to plot and visualize. (In `ISIS`, however, one can combine the data in a plot without having to combine the data for a fit.)

For purposes of this exercise, add together the gratings spectra. This can be done via a tool based approach:
    http://cxc.harvard.edu/ciao/threads/add_grating_spectra/
Note, however, in the above procedure one can only add the two MEG spectra, or the two HEG spectra, but not all four.

If using `ISIS`, one can combine all four datasets; however, before adding the spectra they must be placed on the same spectral grid. The HEG data have twice the spectral resolution of the MEG (i.e., over a given wavelength interval there are two HEG bins for ever one MEG bin). The `match_dataset_grids` function in `ISIS` can be used to match the HEG data grids to the MEG data grids. All the spectra and responses can then be added together using the `combine_datasets` command.

3. Create a combined set of spectra (either the two MEG spectra, or all four spectra), and group this new spectrum to a minimum signal-to-noise of 5 and a minimum of two channels per bin ($\approx$ half width half maximum resolution of the MEG, and full width half maximum for the regridded HEG data).

The response matrices for the HETG are *almost* diagonal, so using "flux corrected" data is "less dangerous" than normal. One should never *fit* with flux corrected data; however, it is sometimes useful for visualization. (However, *always*, revert back to "detector space" plots to check your data and your fits!)

4. Plot flux corrected spectra, and then zoom in on the 12–16 Å region. Note the features that you see here. The goal will be to fit some of these.

## Fitting an Edge and Lines

You should notice an absorption edge in these data, as well as several prominent absorption lines. Do a *local* fit to describe these features. That is, do not attempt to describe the global spectrum, rather try to describe the location and depth of the edge, as well as the location of the absorption lines.

5. Restrict the range of the noticed data to 13.5–15 Å. Feel free now to switch into keV units (and maybe flux units for the y-axis) for the plots. Start with a really simple local continuum model – a powerlaw – fit this, and look at the ratio residuals.

The ratio residuals should give you an idea as to the depth of the absorption edge. (The fractional residual at the edge will be close to the value of the optical depth.)

6. Add an edge to the model and fit the data. In general, when attempting to fit high resolution features in such data, it's best to restrict the locations and widths of the model components, to prevent them from wandering off, or becoming broad and fitting continuum features instead.

The presence of narrow features embedded in a broader, noisy continuum makes fitting these data a good candidate for fit methods other than the usual default "Levenberg-Marquardt" methods, even if these other methods are usually slower. So long as one doesn't have too many bins, and isn't attempting to fit many, many lines, it won't slow you down too much in this case, and it might help find a global minimum.

7. The fit should have improved and you should have obtained a first estimate of the edge parameters; however, there are clearly absorption lines present in the data. Incorporate the most prominent one into the model by subtracting a `gaussian` function. Fit the data and plot your results. Again, constrain the `gaussian` parameters to help the fit from becoming "lost", and to keep the `gaussian` from becoming broad and fitting continuum features instead.

8. The fit should have improved; however, additional absorption lines should remain, including a possible line very close to the edge. Incorporate the next three most prominent residual features by subtracting further gaussians from the model, and fit.

9. Run an error bar search on all the parameters, and save the final fit results to a file. Also save the fit statistic information.

Note that when subtracting a gaussian, it's completely possible for the summed model to *become locally negative*, which then goes unnoticed after the model is smeared by the detector response. (The forward folding doesn't care that the model has gone negative - it's just a vector of numbers related to a function that

is being minimized.) *There have been spectral analyses published in the literature where this has occurred.* So, be careful, and double check your work, and make sure you are in a regime where a gaussian line is an acceptable approximation for the absorption lines. Use a more sophisticated model, such as a Voit profile, if warranted and required by your data!

The expected location of the Neon edge is $14.295 \pm 0.003$ Å. The Neon II 1s→2p line is expected at $14.608 \pm 0.002$, and the Neon III 1s→2p is expected at $14.508 \pm 0.002$. (For many X-ray lines of ionized species, *Chandra* HETG observations have provided better determinations of their positions than either theoretical calculation or laboratory measurements!) How close do your values come to the above? Do your results argue for the edge and line being intrinsic to the black hole system, or due to absorption by the interstellar medium?

## Monte Carlo Simulations

10. The next most prominent residual occurs at $\approx 859$ eV. Is this another significant absorption line? Add one in, fit the data, and run the error bars. Plot and save your results, and save the fit statistic for use in the next step.

The results of the above error bar search suggest that this fifth line is indeed significant – it's 90% confidence value lower limit for the line flux is well above zero. But should we believe that? At what point do we start worrying that we have just fit a random noise fluctuation with a narrow gaussian? (Narrow gaussians will probably describe well any noise fluctuation that's only a few bins wide.) Here is where simulations can be very useful.

The idea is that we take the model parameters from our fit with only four lines, simulate data of the same exposure as our real data, use the same grouping/noticing criteria, then fit these fake data with the five line model. We then store the difference in $\chi^2$ values, and repeat many, many times. We then histogram our results, and see how many times the simulated data (which we *know* has only four lines) yields an improvement in $\chi^2$ as large as the one we found with the real data. (Those who closely followed the statistics lectures will already note some objections to even this scenario. We discuss some of these below.)

To obtain the most meaningful results for such simulations *we need to replicate our analysis procedures as closely as possible* (and ideally our analysis procedure should be one that is well-defined and quantitative). In this case that would mean that we fit, and then run the error bar search to guarantee that we have found the best fit. That's going to be very time consuming. As a "first cut" compromise, run the fits but not the error bars. (Before publishing the results, one would likely increase the fidelity of the simulations.)

10. Write and run a script to evaluate such Monte Carlo simulations. You first have to delete the real data, then create fake datasets with the four line model. Combine, group, and notice these fake data as before. Store the $\chi^2$ value. Load the five line model, and then fit the data. Store this $\chi^2$ value. Repeat many times. (More than 1000 might be prohibitively long depending upon the speed of your computer. Those with slower computers might want to start with 300 trials.) Histogram the results. How many simulations reach or exceed the $\Delta\chi^2$ value that was found with the real data?

Here's the big, obvious objection to the above. When adding the fifth line, we added it to *exactly* the same region as we did for the real data. However, we *chose* that location based upon the fact that it was the largest remaining residual in the spectrum. If it had been *some other* location with that large of a residual, we would have chosen that instead. Therefore, we really should write script to repeat that procedure. First, find the largest remaining residual, then look for a line in a limited bandpass around that residual. That procedure undoubtedly would increase the number of simulations with as large $\chi^2$ changes.

In fact, one might argue that we should just run through all possible independent wavelength regions, and try adding a line. We might expect that we have $\approx 40$ such regions given the energy range we allowed for the fifth line. Given these considerations, how would you expect our Monte Carlo-derived significances to change? What changes would you make to the simulation script? How many trials would you run? (Exercises for the reader!)