

# SciDBoss Documentation

## 1. Quasar sample

I used a testing sample of 5000 quasars from the DR12 catalog (<http://www.sdss.org/dr12/algorithms/boss-dr12-quasar-catalog/>). These quasars were required to have  $2.15 < z < 5$  (using  $Z_{VI}$ ) and  $BAL\_FLAG\_VI = 0$ ; that is, they were required to have no visually apparent broad absorption line troughs. From this set I randomly selected 5000 quasars and downloaded the spec files for each of these quasars, from [http://data.sdss3.org/sas/dr12/boss/spectro/redux/%d/spectra/%04d/spec-%04d-%05d-%04d.fits%\(rerun,plate,plate,mjd\)](http://data.sdss3.org/sas/dr12/boss/spectro/redux/%d/spectra/%04d/spec-%04d-%05d-%04d.fits%(rerun,plate,plate,mjd)) (see <https://www.sdss3.org/dr9/spectro/pipeline.php> for a description, under "per object spec files"). For each quasar I used the coadded data combining each exposure. These data are stored in the fits file *boss\_test\_sample.fits* and in the directory *boss\_test\_sample*. The function loading all the data is *load\_all\_data* within *lya\_functions.py*. I also mask the data by rejecting any pixel with the  $AND\_MASK > 0$  (see <http://www.sdss3.org/dr12/algorithms/bitmasks/#SPPIXMASK> for a description of the pixel maskbits; the  $AND\_MASK$  means that the maskbit was set for any of the exposures within the coadd, not necessarily all). Since the SDSS sky subtraction is often inadequate around bright night sky lines, I additionally mask 6 Å on either side of the night sky lines at 5577 Å, 6300 Å, 6363 Å, and the ISM sodium D line at 5890 Å. The function doing this masking is *mask\_data* in *lya\_functions.py*.

## 2. Finding $\langle F(z) \rangle$

As a first step, I fit a model of the following form to each quasar spectrum:

$$f(\lambda) = A_{qso} \bar{C}(\lambda_r) < F(z_{abs}) > \quad (1)$$

where  $\bar{C}(\lambda_r)$  is the mean quasar continuum,  $A_{qso}$  is the amplitude of each individual quasar, and  $< F(z_{abs}) >$  is the mean transmitted flux fraction through the Ly $\alpha$  forest. This model assumes that all quasars have the same shape and does not consider fluctuations from the mean transmission.

Rather than fitting this model in one go I proceed iteratively. In the first iteration I take the mean continuum to be 1:

$$\bar{C}_0(\lambda_r) = 1 \quad (2)$$

The mean continuum is defined over a specified grid; for now, the grid starts at 600 Å, ends at 3600 Å, and has spacing of 4 Å between points. I use the Faucher-Giguere 2008 model as a rough estimate for  $< F(z) >$  (the analytic fit cited in the Font-Ribera DLA paper):

$$< F(z)_0 > = \exp[-0.0018(1+z)^{3.92}] \quad (3)$$

$\langle F(z) \rangle$  is defined in a grid where each point has  $\Delta z = 0.03$ , and spans a total width of 1.5 in redshift. The beginning of  $\langle F(z) \rangle$  is the lowest measurable absorption redshift: the blue end SDSS spectroscopic limit, 3566.97 Å, divided by Ly $\alpha$ , 1215.24 Å, minus 1.

Last, I find the initial quasar amplitude by finding the weighted mean of the quasar flux between 1275 and 1285 Å.

$$A_{qso,0} = \frac{\sum f_i w_i}{\sum w_i} \quad (4)$$

$w_i$  is the weight of each pixel, calculated from its total inverse variance:

$$w_i = (\sigma_{pipe}^2 + f_i^2 \sigma_i^2)^{-1} \quad (5)$$

where  $\sigma_{pipe}^2$  is the variance outputted by the SDSS pipeline (from the spec file) and  $\sigma_i^2$  is the intrinsic variance, a step function

$$\sigma_i^2 = \begin{cases} 0.1 & \text{if } \lambda_r < 1215.24 \\ 0.01 & \text{if } \lambda_r > 1215.24 \end{cases} \quad (6)$$

Note that  $\sigma_i^2$  is dimensionless. We need  $\sigma_i^2$  for two reasons. First, there is intrinsic variance associated with the transmission fraction in the Lyman $\alpha$  forest (due to cosmic structure) so even a high S/N quasar should have some floor to its variance in the Lyman $\alpha$  forest region. Second, we need an intrinsic variance to prevent high S/N quasars from being weighted too heavily; if we had no intrinsic variance, then the mean continuum would be heavily biased towards high S/N quasars. Note also the units of the weight: since the weights are inverse-variance weights, their units are flux $^{-2}$ .

Now I begin to iterate. First I compute  $\bar{C}$ :

$$\bar{C}(\lambda_r)_i = \frac{\sum_j w_j \frac{f_j}{A_{qso,i-1} \langle F(z) \rangle_{i-1,j}}}{\sum_j w_j} \quad (7)$$

The weights are equal to the total inverse variance of the quantity being averaged

$$w = \left( \frac{\sigma_{pipe}^2}{A_{qso,i-1}^2 \langle F(z) \rangle_{i-1}^2} + \bar{C}_{i-1} \sigma_i^2 \right)^{-1} \quad (8)$$

Note that I use the  $(i-1)$ th estimates of  $A_{qso}$ ,  $\langle F(z) \rangle$  and  $\bar{C}$ . For any pixel with  $z_{abs} > z_{max} + \Delta z/2$ , I have no estimate of  $\langle F(z) \rangle_{i-1}$ , so I use the Faucher-Giguere estimate at these redshifts for all iterations.

Next I compute  $A_{qso}$ . Here I use a chisq minimization that is formally identical to the weighted average but is easier to implement if I switch to using a multi-parameter model for the continuum. Specifically I only fit to the region red of Ly $\alpha$ , and I use weighted least

squares to properly account for the different variances between different pixels; the weights are given by:

$$w = \sqrt{\left( \frac{\sigma_{pipe}^2}{\langle F(z) \rangle_{i-1}^2} + A_{qso,i-1}^2 \bar{C}_i^2 \sigma_i^2 \right)^{-1}} \quad (9)$$

in other words the inverse of the standard deviation of each pixel, in units of flux<sup>-1</sup>. Note also that  $\langle F(z) \rangle_{i-1}$  will always be equal to 1 since we only fit red of Ly $\alpha$ . Written in matrix form we then solve

$$W \vec{T} A_{qso} = W \vec{f} \quad (10)$$

where  $W$  is the diagonal matrix of the weights,  $\vec{T}$  is the template ( $\bar{C}(\lambda_r)_i$  re-binned to the quasar wavelengths),  $\vec{f}$  is the measured fluxes, and  $A_{qso}$  is the scalar amplitude that we're fitting for. This is a least squares problem of the form  $\vec{a}X = \vec{b}$  so it can be solved quickly by already-existing libraries to find  $A_{qso}$ .

Now I find  $\langle F(z) \rangle$ . I consider only the region between 1041 Å and 1185 Å rest frame. I assign each pixel to the nearest  $z_{abs}$  and calculate  $\langle F(z) \rangle$  in the following way:

$$\langle F(z) \rangle_i = \frac{\sum_j w_j \frac{f_j}{A_{qso,i} \bar{C}_i}}{\sum_j w_j} \quad (11)$$

where now  $w_j$  is given by

$$w_j = \left( \frac{\sigma_{pipe}^2}{(A_{qso,i} \bar{C}_i)^2} + \langle F(z) \rangle_{i-1}^2 \sigma_i^2 \right)^{-1} \quad (12)$$

I also calculate the error on each measurement of  $\langle F(z) \rangle$ :

$$\sigma_F^2 = \frac{1}{\sum w_i} \quad (13)$$

This follows from the variance of a weighted mean (see appendix). Also,  $\langle F(z) \rangle$  is fixed in order to break the degeneracy between  $\langle F(z) \rangle$  and  $\bar{C}$  in the Ly $\alpha$  forest region. I fix the mean of  $\langle F(z) \rangle$  between  $z = 2.2$  and  $z = 2.6$  to equal the mean of Faucher-Giguere's measurement of  $\langle F(z) \rangle$  between 2.2 and 2.6.

To measure convergence, I calculate a "chisq" statistic:

$$\chi^2 = \sum \left( \frac{\langle F(z) \rangle_i - \langle F(z) \rangle_{i-1}}{\sigma_{Fi}} \right)^2 \quad (14)$$

where the sum runs over the redshift bins, and  $i$  refers to the  $i$ th iteration. Convergence is reached when this statistic is less than the number of redshift bins.

### 3. Appendix: variance of a weighted mean

Consider the estimation of the mean and standard error for a set of data  $x_i$  with weights  $w_i$ . These weights can be interpreted as non-integer “counts” of each data point. For instance, a weight of 1.3 means that a particular datapoint is counted 1.3 times.

Take the probability of obtaining a particular measurement  $x_i$ . Assume this distribution is Gaussian:

$$P_i(\mu') \propto \exp \left[ -\frac{1}{2} \left( \frac{x_i - \mu'}{\sigma} \right)^2 \right] \quad (15)$$

The probability of observing the entire data set of  $N$  observations is given by

$$P(\mu') = \prod_{i=1}^N P_i(\mu')^{w_i} \quad (16)$$

We see that the interpretation of weights as effective counts means that we must raise each probability to the  $w_i$  power. We thus have

$$P(\mu') \propto \exp \left[ -\frac{1}{2} \sum w_i \left( \frac{x_i - \mu'}{\sigma} \right)^2 \right] \quad (17)$$

According to the method of maximum likelihood, the most probable value for the mean is the value that maximizes the probability in equation 17. This is equivalent to minimizing the argument of the exponential:

$$X = -\frac{1}{2} \sum w_i \left( \frac{x_i - \mu'}{\sigma} \right)^2 \quad (18)$$

The minimization condition is

$$\frac{\partial X}{\partial \mu'} = 0 \quad (19)$$

yielding (after algebraic simplification)

$$\sum w_i (x_i - \mu') = 0 \quad (20)$$

or

$$\mu' = \frac{\sum w_i x_i}{\sum w_i} \quad (21)$$

This checks with our intuition.

The standard error is given by

$$\sigma_\mu^2 = \sum \left[ \sigma_i^2 \left( \frac{\partial \mu}{\partial x_i} \right)^2 \right] \quad (22)$$

using error propagation. We have

$$\frac{\partial \mu}{\partial x_i} = \frac{\partial}{\partial x_i} \left( \frac{\sum w_i x_i}{\sum w_i} \right) = \frac{w_i}{\sum w_i} \quad (23)$$

Thus

$$\sigma_\mu^2 = \sum [\sigma_i^2 \left( \frac{w_i}{\sum w_i} \right)^2] \quad (24)$$

In our case the weights are equal to the inverse of the variance. Therefore

$$\sigma_\mu^2 = \left( \frac{1}{\sum w_i} \right)^2 \sum \frac{w_i^2}{w_i} = \frac{1}{\sum w_i} \quad (25)$$