

Lab 1: Gaia, RR Lyrae stars, and Galactic Dust

Astro 128 / 256 (UC Berkeley)

Assigned: Thursday February 6, 2025

Checkpoints: Thursday, February 13 2025; Tuesday, February 25, 2024; Tuesday March 4, 2025

Final Write Up Due: Thursday March 6, 2025

In this lab, we'll explore the sample of RR Lyrae stars observed by Gaia. The final goal is to build a map of the distribution of dust in the Milky Way. We'll make use of the fact that RR Lyrae stars are "standardizable candles" – that is, they follow a period-luminosity trend (the [Leavitt Law](#)) that can be used to infer intrinsic luminosity using measurements of brightness and periodicity. This, in turn, can be used to infer distances to each star.

As in Lab 0, the catalogs we'll need can all be found on the Gaia archive. To see which catalogs are available, go to the [Gaia archive](#), click "search", and click on the "Advanced (ADQL)" tab on the upper left. Available catalogs are listed on the left side of the page in nested drop-down menus. We will use both catalogs from DR3 and eDR3.

Technical Components

- Periodograms
- Fourier decomposition
- Databases; SQL
- Python web queries
- Linear and nonlinear optimization
- Markov Chain Monte Carlo (MCMC)
- Bayesian modeling
- Data visualization

Broad Summary of Steps

You should expect the lab to proceed in the following order:

- You will need to write ADL queries to download RR Lyrae stars from the Gaia archive
- The Gaia archive provides time-averaged magnitude, periods, etc. but of varying quality. Your next task will be to confirm the reliability of these values. Your tasks will be to estimate the periods and mean magnitudes of RR Lyrae on your own, and compare to the reported Gaia values.
- We now need to establish a period-relationship for the RR Lyrae. You will need to write an ADL query to get parallaxes to nearby RR Lyrae in order to determine distances. You'll then explore the reliability of the reported distances.

- To determine the period luminosity (or absolute magnitude) relationship, your next task will be to use a Markov chain Monte Carlo (MCMC) approach to fitting. First, we'll ask you to write and test your own, simple MCMC routine. And then you'll compare this to results from a package called `pymc`. Please note that this step can be quite time consuming.
- Finally, you will create a dust map of the Milky Way. To do this, you will apply your measured P-L relationship to the full RR Lyrae catalog from the Gaia archive. You'll compute the "color excess", which tells you about extinction along the line-of-sight, and make a 2D map of dust in the Milky Way.
- Note that we include a number of links to useful concepts (e.g., error propagation, MCMC tutorials) that we strongly encourage you to read.

Preamble

Note: you likely already did some or all of Problem 0 for Lab 0.

- (a) For this lab, you will have to write a lot of ADQL queries. If you run into difficulties, you can consult the [Gaia ADQL cookbook](#). This should contain everything you need for the lab. We particularly recommend reading through the "Query timeouts" section.

This tutorial is for your own edification, so spend as much or as little time as you need on it.

- (b) If you haven't already, make yourself an account on the Gaia Archive by going to the [Gaia archive](#) and clicking "sign in" and "register new user" in the upper right-hand corner. You don't need an account to query the Gaia catalogs, but having an account will allow you to save your previous queries and to upload and query your own tables.
- (c) Download [Astroquery](#). (Already done for you if you are working on Datahub). The `astroquery.utils.tap.core.TapPlus` utility will allow you to combine ADQL queries with Python code. *Note: For Lab 0, some people used a different package, `astroquery.gaia.Gaia`. Feel free to use whatever packages you like to submit queries, but be aware that this one seems to time out faster than the `TapPlus` package.*
- (d) As in Lab 0, you should develop a system for saving/caching the results of web queries (i.e., ADQL queries, light curves, anything else that requires an internet connection) locally. This is a safeguard in case, say, the Gaia archive crashes an hour before the lab is due. After Lab 0, you know this is not just a hypothetical! That is, once you've run a cell or query once, you should be able to run it again without a web connection. You'll also want to implement a system for overwriting the cache when you improve your queries.

Assignment

1. The `gaiadr3.vari_rrlyrae` catalog contains the results of RR Lyrae classification done by the Gaia data processing and analysis consortium for more than 10^5 RR Lyrae stars. Write an ADQL query to download the first 100 rows of the table for which a fundamental pulsation frequency ("pf") has been measured and more than 40 clean epochs were obtained in the G-band. Submit your query using Astroquery and display the first 10 rows of the thus-obtained catalog.

To find the meaning of individual columns, you can search for them in the [Gaia DR3 data model](#). To learn about the meaning of RR Lyrae-specific parameters and how they were derived, see [Clementini et al. 2016](#) and [Clementini et al. 2022](#). Table 12 lists the relevant columns for the data you're downloading from the Gaia archive.

2. The table you downloaded above contains the results of fitting the light curves of RR Lyrae stars, but it does not contain the raw light curves. The raw light curves can be accessed as described in the [Datalink and light curves](#) tutorial.

Download the light curves (i.e., G-band magnitude vs. time, with magnitude uncertainties) for the 100 RR Lyrae stars in your table. This is most easily accomplished with a Python utility for fetching web urls. `astroquery` is one good option. Plot one of the light curves.

3. Estimate the period and mean G-band magnitude of the 100 light curves you downloaded above. A simple estimate of the period can be obtained using a Lomb-Scargle (L-S) periodogram. An implementation is available in `astropy.timeseries.LombScargle`. Plot the periodogram for one light curve, marking your estimate of the period on the plot.

When estimating mean magnitudes, note that magnitudes are a logarithmic quantity, so simply taking the mean of the measured value is *not* correct. Note that your plot should include error bars on the absolute magnitudes. These error bars are a combination of the flux errors and the [zero point errors](#). More information on error propagation can be found [here](#) and [here](#). *Hint: You can use the range of plausible periods for RR Lyrae to set the minimum and maximum frequency in your L-S analysis.*

4. Compare the periods you computed from the 100 light curves to the values reported in the `vari_rrlyrae` catalog. Comment on your results.
5. Suppose we want to predict how bright an RR Lyrae star will be at some time in the future. To do this, we need to know a functional form, $m_G(t)$, which we can evaluate at some future time t . Stellar atmospheres are complicated, so predicting a closed-form expression for $m_G(t)$ from first principles is hard. However, since we expect the fluctuations to be periodic, we can use some results from Fourier analysis to get a good estimate of $m_G(t)$.

Any periodic, smooth function $f(t)$ can be described by a sum of sines and cosines:

$$f(t) = A_0 + \sum_{k=1}^K [a_k \sin(k\omega t) + b_k \cos(k\omega t)], \quad (1)$$

where A_0 , a_k , and b_k are to-be-determined constants (a_k and b_k are arrays that are K terms long), and $\omega = 2\pi/P$ is the angular frequency. This equation is exact in the limit of $K \rightarrow \infty$. The reason it is useful for us is that for typical periodic functions, it is reasonably accurate even for small K (say, $K = 5$ or 10).

In the most general case, determining all the $2K + 2$ free parameters needed to represent $f(t)$ is a complex nonlinear optimization problem. However, if ω is already known (say, from a periodogram), the problem becomes much simpler.

Show that if ω is known, then the problem of determining the remaining $2K + 1$ free parameters can be re-cast as a linear algebra problem; i.e., $y = \mathbf{X}\beta$. Here y is an array of the measured fluxes, \mathbf{X} is a matrix that can be constructed from known quantities, and β is an array of unknowns. What are \mathbf{X} and β ? (write out the terms).

6. **Note: Past students have indicated this is a particularly challenging problem** For the star with Gaia DR3 id = 4659713442253931776, determine the series representations for $K = 1, 3, 5, 7$, and 9 . Plot the phased light curve, the series representation on a fine grid, and the residuals between the data and series representation for each K .
7. **Note: Past students have indicated this is a particularly challenging problem** How many terms is enough? Use cross-validation to find the optimal K . Designate 20% of the observed points as the cross-validation set. For K ranging from 1 to 25, calculate χ^2/N_{data} for both the training data (which enters y) and the cross-validation data (which does not). Plot this quantity as a function of K .

Discuss what an appropriate value of K for this data set might be. It maybe useful to use a log-scale on the y-axis.

8. Using a value of K chosen based on (7), predict the expected magnitude of the star exactly 10 days after the last observed data point (in the units returned from the Gaia archive). Plot the light curve, showing the last few days of Gaia data points, the extrapolation over the next 12 days, and indicating your best estimate for the magnitude exactly 10 days after the last data point.
9. When you calculated the average magnitude in part (3), you probably used some sort of mean in flux space of the observed data points. This is not ideal because the observations are not necessarily evenly sampled in phase. A more accurate value can be obtained by calculating the mean magnitude based on the mean flux when averaging over one pulsation period. Now estimate the mean magnitude using your Fourier model for each of the 100 light curves. Once again, the averaging should be done in flux space, not magnitude space. Make a plot comparing your estimates of the mean G-band magnitude – both from this estimate and from part (3) – to the estimate in the Gaia catalog. Also plot the residuals, and comment on your results.

Required checkpoint 1, due Thursday February 13 by 1:59pm: (a) the plot produced in part (6), and (b) the plot produced in part (8). Submit this via gradescope. It should be a pdf, named Firstname_Lastname_lab1_cp1.pdf

10. All the RR Lyrae you have looked at so far are in the variability class “RRab”. (This is a consequence of our requirement in part (1) that there be a measured “pf”). There are, however, other observationally-motivated classes of RR Lyrae, the most common of which is “RRc”.

Let’s compare the light curve shapes of RR Lyrae in the these two classes. Following a similar procedure to the one in part (2), download light curves for the top 3 RR Lyrae in the `vari_rrlyrae` catalog that have `best_classification = “RRc”`, mean G-band magnitudes brighter than 15, and more than 80 clean epochs in the G-band. Using a suitable number of terms, compute Fourier expansions for these three light curves. Plot the phased light curves and overplot their Fourier models.

Repeat this for the top 3 RR Lyrae with `best_classification = “RRab”`. Compare the 6 phased light curves and models, using some plotting scheme that makes it easy to see and compare their shapes. Comment on the difference in light curve shape between the two classes. Do some reading (published papers, Wikipedia, etc) about the difference between the classes.

11. Are the 6 light curves you plotted well-described by a single period, or is there evidence of intrinsic scatter? Read about deviations from simple periodicity in RR Lyrae in Netzel et al. 2018 and discuss your findings.
12. Now that you’ve looked at the light curves for a few objects, we’ll assume for the rest of the problem that the periods reported in the Gaia catalog are reliable. We’ll now use the data in the `vari_rrlyrae` table to infer the RR Lyrae period-luminosity relation in different bandpasses.

First, we’ll need to use Gaia distances to estimate the absolute magnitude of RR Lyraes. This will only work if there isn’t a lot of dust between us and the RR Lyraes. Explain why this is.

It turns out that most of the dust in the Milky Way is in the disk, at low Galactic latitude. Write an ADQL query to select RR Lyrae stars that (a) have accurately measured distances, with parallax errors of less than 20%, (b) are above or below the disk, with $|b| > 30$ degrees, where b is Galactic latitude, and (c) are relatively nearby, with distances less than 4 kpc. To do this, you’ll need to join the `gaiadr3.vari_rrlyrae` and `gaiadr3.gaia_source` catalogs. You should find about 500 objects.

13. Plot the distribution of targets you obtained in Galactic coordinates. Verify that your ADQL query has removed stars in the Galactic disk.

14. Plot period vs. absolute G-band magnitude for all stars returned by your query. Include error bars on your plots.
15. You should see that the majority of the stars have similar absolute magnitude, but a non-negligible fraction of them scatter far off the median relation. This is mostly due to incorrectly measured parallaxes. Apply the quality cuts in Equations C1 and C2 of [Lindgren et al. 2018](#) to the sample. Then plot the period-luminosity relation again. Has the scatter decreased?
16. Most of the “bad” objects should have been removed by the above cut, but a few will remain. You can crudely remove these outliers by removing objects with absolute G magnitudes greater than some threshold of your choice. You can also try more sophisticated outlier rejection, if you wish. Plot the period-absolute magnitude relation, including error bars on absolute magnitude due to distance uncertainties, from the resulting cleaned sample.
17. **(this problem will be started as an in-class group coding activity. Results from your implementation should be lab report).** We are going to use Markov chain Monte Carlo (MCMC) techniques to fit models to data. This technique is quite powerful and frequently used in cutting-edge research, in part because it can yield reliable uncertainty estimates. However, it is also likely to be new to you and it is not intuitive when learning it for the first time.
 To develop an understanding for how Markov chain Monte Carlo (MCMC) routines work, you’re going to start by coding up your own Metropolis-Hastings (M-H) MCMC sampler using a Gaussian proposal distribution. To verify that your sampler works, pretend you have a single measurement of a number, $x = 1$, with error $\sigma = 0.1$. Draw 10,000 MCMC samples assuming a one-dimensional Gaussian likelihood, $p(\mu) = 1/\sqrt{2\pi\sigma^2} \exp\left[-(\mu - x)^2 / (2\sigma^2)\right]$. Compare a normalized histogram of these samples to the analytic likelihood $p(\mu)$. Choose a step size such that the acceptance fraction is of order 0.5. To qualitatively demonstrate convergence, plot μ and $\ln P$ versus step number.
18. **Note: Past students have indicated this is a particularly challenging problem.** Fit a line to the period vs. absolute magnitude relation in the G band, using DR3 data. Assume a model $M_G = a \times \log[P/\text{day}] + b$, where a and b are free parameters. Allow for intrinsic scatter in your relation. That is, fit for a positive constant σ_{scatter} , such that if there were no measurement uncertainties, values of M_G at a given P would follow a Gaussian distribution with variance $\sigma_{\text{scatter}}^2$. Do this in several steps:

- (i) **(optional extra credit – whether you do this or not, some of the terminology and definitions will be useful for (ii) and (iii))** Once you are satisfied that your sampler works, use it to constrain the values of a , b , and σ_{scatter} in the RR Lyrae period-luminosity relation. You will need to write down a likelihood function, $p(\vec{d}|a, b, \sigma_{\text{scatter}})$, where \vec{d} is the data, and a prior, $p(a, b, \sigma_{\text{scatter}})$. State clearly what priors you are assuming. Use Bayes’ theorem to calculate the posterior probability distribution, $p(a, b, \sigma_{\text{scatter}}|\vec{d})$, which will include an unknown multiplicative factor $p(\vec{d})$, the “evidence”. Use your sampler to draw samples from the posterior distribution (you don’t need to know the evidence to do this). You may need to modify it slightly to account for the fact that the function you are sampling from is now three-dimensional rather than one-dimensional. Once again, tune the step size so that the acceptance fraction is about 0.5, and make diagnostic plots to verify that the sampler has converged. Use the [corner](#) package to visualize constraints on the posterior for the fit.

Plot 50 random, independent samples from the posterior over the data. Does the spread between samples as a function of period seem consistent with what you’d expect given the data? Explain.

- (ii) Repeat your fit to the period-luminosity relation, using the same likelihood and priors as in (i). But instead of using your M-H sampler, use the [no-U-turn](#) Hamiltonian Monte Carlo sampler

provided in [pymc](#). Discuss briefly what this sampler is and what some of the advantages are of using it over Metropolis-Hastings. Add your likelihood function to your pymc model as a “potential” term.

- (iii) Repeat step (ii), but instead of enrolling your likelihood function into the pymc model explicitly, simply tell pymc your model: namely, that you expect values of M_G at a given P to have a mean value $\mu = a \log(P/\text{day}) + b$, and to follow a normal distribution with mean μ and variance $\sigma_{\text{scatter}}^2 + \sigma_i^2$, where σ_i are the measurement uncertainties.

On a single corner plot, compare the posterior constraints you obtained in parts (i) and (iii) [, and (i) if doing the extra credit]. If everything worked, they should be basically identical. When we do MCMC fits in the future, you will usually be able to skip right to (iii). We hope (i) and (ii) have provided some insight into what is happening under the hood.

Required checkpoint 2, due Tuesday February 25 by 1:59pm: the corner plot produced at the end of part (18.iii), where you compare posterior constraints for the period-luminosity relation from two (or three, if doing extra credit) sampling methods. Submit this via gradescope.

It should be a pdf, named `Firstname_Lastname_lab1_cp2.pdf`

19. **(optional for extra credit in AY128)** Many of the RR Lyrae identified by Gaia were also observed by the WISE survey, which observed the whole sky in the near-infrared. Cross-match your sample of clean RR-Lyrae stars with WISE. The catalogs `gaiadr3.allwise_best_neighbour` and `gaiadr3.allwise_neighborhood`, available on Gaia archive, will be useful.
20. **(optional for extra credit in AY128)** Repeat step (17), now using the WISE “W2” magnitude rather than the G band. You can simply use the magnitude reported in the Wise catalog; you don’t need to average any light curves. You can skip directly to fitting method (iv).
21. **(optional for extra credit in AY128)** Comment on the differences between your inferred period-luminosity relations in the optical and in the near-infrared. In which band is the period-luminosity relationship steeper?
22. Compare your derived period-luminosity relations to results in the literature, which can be found in [Beaton et al. 2018](#) or [Klein & Bloom 2014](#). You can compare your G-band relation to their V-band relation. If there are systematic differences between your results and the literature, what might account for them?
23. Following a similar procedure to the one in step (17), derive a period-color relation for RR Lyrae stars in the Gaia bands. That is, a relation between $\log(\text{period})$ and the Gaia $G_{BP} - G_{RP}$ color. You will need to calculate color uncertainties using the uncertainties in BP and RP flux. You don’t need to use your own sampler, but can skip right to the pymc version.
24. Write an ADQL query to download the full Gaia RR Lyrae catalog (i.e., no longer excluding stars with imprecise parallaxes or low galactic latitude) and cross-match it with the Gaia source catalog.
25. Calculate the *color excess*, $E(G_{BP} - G_{RP}) = (G_{BP} - G_{RP})_{\text{observed}} - (G_{BP} - G_{RP})_{\text{intrinsic}}$, for all RR Lyraes in the catalog. From this, calculate A_G , the G-band extinction, for each star. You may assume

$$R_G \equiv \frac{A_G}{E(G_{BP} - G_{RP})} = 2.0.$$

26. Compare your calculated A_G values to the “G_absorption” value provided in the RR Lyrae catalog.
27. Plot a 2-d map of $E(G_{BP} - G_{RP})$ as a function of Galactic longitude and latitude, using an Aitoff projection. A simple way to do this is to plot each RR Lyrae as a point in a scatter plot with semi-transparent points, coloring each point by its A_G .

When you first make the plot, you will find that some large-scale structure is clearly visible, but also that there are some points for which the reddening looks clearly wrong; i.e., the color excess in one point is very different from the adjacent points. Construct appropriate quality cuts to remove these objects. You may want to cut on photometric signal-to-noise and/or BP/RP excess. To ensure reasonable spatial coverage, you should not have any less than $\sim 60,000$ stars (and you may have more). Be sure to clearly explain your cuts and the effects they have on the resulting map.

You will notice that the distribution of RR Lyrae stars in the catalog is not at all uniform. Why is this?

Required checkpoint 3, due Tuesday March 4 by 1:59pm: the cleaned reddening map produced in part (27). Submit this via gradescope.

It should be a pdf, named `Firstname_Lastname_lab1_cp3.pdf`

28. The most widely used Galactic dust map is the “SFD” map produced by [Schlegel, Finkbeiner, and Davis 1998](#). Incidentally, all three authors have strong ties to UC Berkeley.

Compare the attenuation map you produced above to the SFD map. You can query SFD using the “[dustmaps](#)” Python package. You only need the SFD map; don’t worry about installing the larger maps (some of which are several GB).

Plot the SFD optical reddening map, $E(B - V)$, sampled at the same positions as your RR Lyrae map. Does the general structure of your map agree with SFD? What about the small-scale details?

29. Do you *expect* your map to look exactly like SFD? Hint: think about how the SFD map is constructed. What are the differences between what it measured and what your map measures?

Next Steps

Your full lab write up is due on Thursday 3/6/25 by 11:59pm. It should be a full lab report including discussion of all plots and sections highlighted in magenta, named `Firstname_Lastname_lab1.pdf`.