

1. Introducing Galaxies

Galaxies are a slippery topic in astronomy at present. Galaxies are much less well understood than say stars; certainly our understanding is changing (and hopefully improving) noticeably each year. The standard texts/references are *Galactic Astronomy* by Binney and Merrifield, and *Galactic Dynamics* by Binney and Tremaine. *The Physical Universe* by Shu is more elementary, but very insightful and always repays reading.

The reason galaxies are difficult to understand is that they are made of three very different kinds of things. There are stars of course, but there's also the interstellar medium (which produces stars, and is in turn fed by dying stars), and dark matter (about which we know very little, except that it's there). And these three all influence each other. We'll study each of these, and to a small extent how they influence each other. Some galaxies (more of them in earlier epochs) have 'active nuclei' which can vastly outshine the starlight, but we won't go into that—we'll confine ourselves to 'normal' galaxies.

There are three broad categories of galaxies:

DISC GALAXIES

These have masses of $10^6 M_\odot$ to $10^{12} M_\odot$. The discs brightness tend to be roughly exponential, i.e.,

$$I(R) = I_0 \exp[-R/R_0] \quad (1.1)$$

I_0 is $\sim 10^2 L_\odot \text{pc}^{-2}$. The scale radius R_0 is $\simeq 4 \text{kpc}$ for the Milky Way. The visible component is $\simeq 95\%$ stars (dominated by F and G stars for giant spirals), and the rest dust and gas. The more gas-rich discs have spiral arms, and arms are regions of high gas density that tend to form stars; clumps of nascent stars are observed as H II regions. Disc galaxies have bulges which appear to be much the same as small ellipticals. All disc galaxies seem to be embedded in much larger dark halos; the ratio of total mass to visible stellar mass is $\simeq 5$, but we don't really have a good mass estimate for any disc galaxy.

ELLIPTICAL GALAXIES

These have masses from $10^6 M_\odot$ to $10^{12+} M_\odot$. There are various functional forms around for fitting the surface brightness, of which the best known is the de Vaucouleurs model

$$I(R) = I_0 \exp \left[-(R/R_0)^{1/4} \right]. \quad (1.2)$$

with $I_0 \sim 10^5 L_\odot \text{pc}^{-2}$ for giant ellipticals. (To fit to observations, one typically un-squashes the ellipses to circles first. Also, the functional forms are only fitted to observations over the restricted range in which $I(R)$ is measurable. So don't be surprised to see very different looking functional forms being fit to the same data.) The visible component is almost entirely stars (dominated by K giants for giant ellipticals), but there appears to be dark matter in a proportion similar to disc galaxies. Ellipticals of masses $\lesssim 10^{11} M_\odot$ rotate as fast as you'd expect from their flattening; giant ellipticals rotate much slower, and tend to be triaxial—more on this later.

At the small end of ellipticals, we might put the globular clusters, even though they occur inside galaxies rather than in isolation. These are clusters of masses from $10^4 M_\odot$ to $10^{6.5} M_\odot$, consisting exclusively of very old stars.

2 Introducing Galaxies

IRREGULARS

Everything else! They tend to have strong emission lines, and their starlight is dominated by B,A and F types. Basically, they look like they've just been shaken up and are responding by forming stars.

PROBLEM 1.1: Instead of functional forms for the surface brightness $I(R)$, people sometimes pick a functional form for the 3D density $\rho(r)$. These are related by the projection

$$I(R) = 2 \int_R^\infty \frac{r\rho(r) dr}{\sqrt{r^2 - R^2}}$$

which is easily worked out numerically if not analytically.

A popular example are the Dehnen models:

$$\rho(r) = \frac{q}{4\pi} \frac{r^q}{r^3(1+r)^{q+1}}$$

where q is an adjustable parameter. Here the normalization is chosen so that ρ integrates to unity. The special case of $q = 1$ (called the Jaffe model) is particularly important because it is found to fit the observed $I(R)$ of ellipticals at least as well as de Vaucouleurs' profile.

What is the potential of a mass distribution with a Jaffe $\rho(r)$? [10]

The Dehnen models have an interesting limit as $q \rightarrow 0$. What is it? [5]

EXAMPLE [The fundamental plane for ellipticals] If we assume that all ellipticals have the same constant mass to light ratio and the same form for the mass distribution (only scalable) then $M \propto I_0 R_0^2$, where I_0 is a characteristic surface brightness and R_0 a characteristic radius. The virial theorem implies $M \propto R_0 \sigma_0^2$ where σ_0 is a characteristic velocity dispersion (if we assume dispersion dominates rotation). So under these assumptions we'd expect

$$R_0 I_0 \sigma_0^{-2} = \text{constant.} \quad (1.3)$$

Observationally, ellipticals are found to satisfy

$$R_0 I_0^{0.9} \sigma_0^{-1.4} = \text{constant} \quad (1.4)$$

to within observational uncertainties. In the space of $\log R_0, \log I_0, \log \sigma_0$, equation (1.4) is of course a plane, and it is called the fundamental plane. Deviation from the virial prediction presumably has something to do with varying mass to light, but nobody seems to have much idea of why it's a very good correlation in practice.

In diffuse dwarf ellipticals, $I(R)$ falls off faster than in giant ellipticals or compact dwarf ellipticals, so $M \propto I_0 R_0^2$, wouldn't have the same proportionality factor. And observationally, diffuse dwarf ellipticals don't lie on the fundamental plane. \square

PROBLEM 1.2: Suppose some category of galaxies has $I(R) = I_0 f(R/R_0)$ with all galaxies having the same I_0 and function f but different galaxies having different R_0 . If the mass to light is constant everywhere then show that

$$L \propto v^4$$

where L is the total luminosity and v is a characteristic velocity. [10]

For spirals, the $L \propto v^4$ relates the total light to the disc rotation velocity (as measured in radio or infrared), and is called the Tully-Fisher relation. In ellipticals (with v identified with the velocity dispersion) it is called the Faber-Jackson relation. Tully-Fisher is important in distance scale work.

HUBBLE TYPES

On the whole, galaxy classification probably shouldn't be taken as seriously as stellar classification, because there isn't (yet) a clear physical interpretation of what the gradations mean. But some physical properties do clearly correlate with the so-called Hubble types, so it's worth learning about these at least.

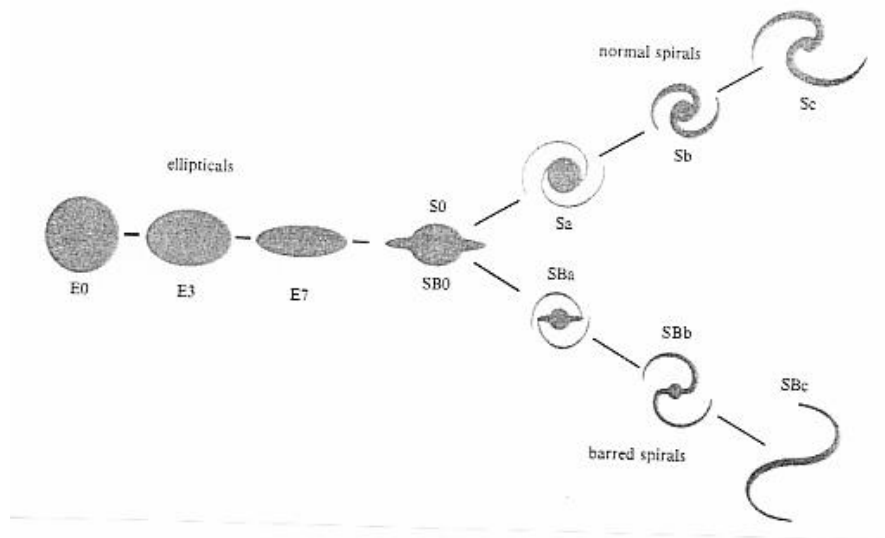


Figure 1.1: The tuning fork diagram of Hubble types.

Figure 1.1 shows the Hubble types. Ellipticals go on the left, labelled as E_n , where $n = 10(1 - \langle \text{axis ratio} \rangle)$. Then the lenticulars or disc galaxies without spiral arms: S0 and SB0. Then spirals with increasingly spaced arms, Sa etc. if unbarred, SBa etc. if barred.

The left ones are called early types, and the right ones late types. People once thought this represented an evolutionary sequence, but that's long been obsolete. (Our current understanding is that, if anything, galaxies tend to evolve towards early types.) But the old names are still used.

We never see ellipticals flatter than about E7. The reason (as indicated by simulations and normal mode analyses) seems to be that a stellar system any flatter is unstable to buckling, and will eventually settle into something rounder.

Note that bulges get smaller as spiral arms get more widely spaced. Theory for spiral density waves predicts that the spacing between arms is proportional to the disc's mass density.

HANDWAVING DYNAMICS

Stars are so compact on the scale of a galaxy that a stellar system behaves like a collisionless fluid (except in the cores of galaxies and globular clusters), resembling a plasma in some respects. Gas and dust are collisional. This leads to two very important differences between stellar and gas dynamics in a galaxy.

- 1) Gas tends to settle into discs, but stars don't.
- 2) Gravity must be balanced by motion in stellar and gas dynamics, but in equilibrium gas must follow closed orbits (and in the same sense), but stars in general don't. Two streams of stars can go through each other and hardly notice, but two streams of gas will shock (and probably form stars). You could have a disc of stars with no net rotation (just reverse the directions of motion of some stars), but not so with a disc of gas. People sometimes speak of 'rotation-support' and 'pressure-support' balancing self-gravity. Pressure support refers to the high velocity dispersion (compensating for low net-rotation) that comes from reversing stellar motions; this stellar dynamical pressure needn't be isotropic. Observationally gas dispersions are never more than $\simeq 10$ km/sec while stellar dispersions can easily be $\simeq 300$ km/sec.

We can start putting together a general picture now. (The rest of this paragraph varies from well-accepted to controversial to wildly speculative, so don't take it too seriously.) Primordial gas will tend to form rotating discs. Differential rotation in the discs will cause spiral density waves, enhancing density along spiral arms and preferentially forming stars. A bulge-less stellar disc is actually unstable to buckling, and produces a bulge with part of its mass. (That's what simulations indicate.) A bulge formed this way will be rotationally supported like the disc that gave rise to it. Meanwhile the disc will continue to form stars, so disc stars will tend to be younger than a bulge stars. Discs that have turned almost all their gas into stars will have stellar discs, but no spiral arms. Now, a disc galaxy can be disrupted by the gravitational influence of another galaxy. It can be a merger of two or more galaxies, or the tidal disruption of a single galaxy; both tending to disrupt discs and produce irregulars with much star formation, then ellipticals. Disruptions of single galaxies will tend to produce rotationally supported ellipticals; but for mergers the angular momentum vectors will tend to cancel, producing pressure support. So we might expect giant ellipticals to be pressure supported. But even a completely gas-free elliptical will generate gas from its dying stars. This second-generation gas will of course settle into discs, and there we might see spiral arms all over again. . . And all this while, dark matter (whatever it is) will be finding gravitational potential wells in the neighbourhood of galaxies and form halos (sort of like polarization clouds) around them.

Note, by the way, that all galaxies appear to have *some* stars $\sim 10^{10}$ yr old. Evidently galaxies all formed fairly early, though they have merged or been otherwise disrupted much more recently.

To end this introductory chapter, let's look at a picture that says rather a lot—it's a very deep photograph of the Sombrero galaxy: Figure 1.2. (You may have across a gorgeous colour poster of this galaxy.) Is it an elliptical with a large embedded disc or a spiral or lenticular with an extra large bulge? But in Figure 1.2 the main galaxy is just an inset within a much larger dark halo. And what is that diffuse fan to the NE and the loop to the SW? Almost certainly traces of past encounters with other galaxies.

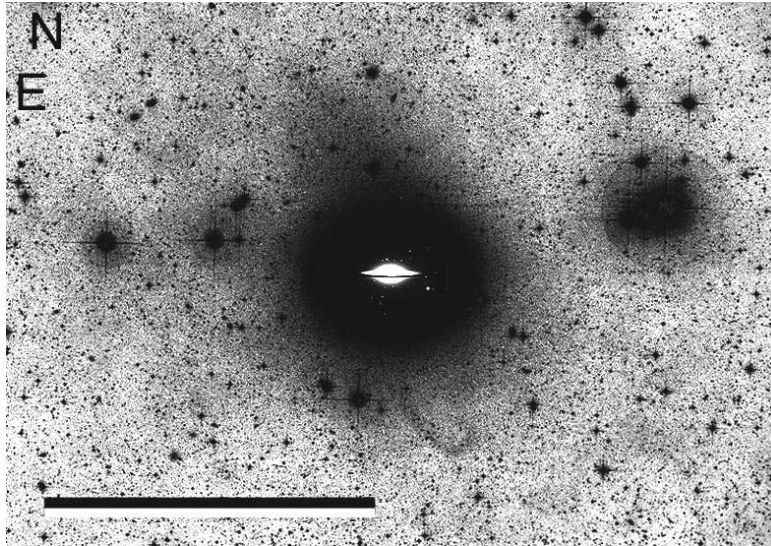


Figure 1.2: A recent deep photograph by David Malin of the Sombrero galaxy (aka M 104 and NGC 4594) with a ‘normal’ image inset to the same scale. The scale bar is 30’.

2. Stellar Dynamics

A system of stars behaves like a fluid, but one with unusual properties. In a normal fluid two-body interactions are crucial in the dynamics, but stellar encounters are very rare. Instead stellar dynamics is mostly governed by interaction of individual stars with the mean gravitational field of all the other stars.

THE VIRIAL THEOREM

Before going into the main material on stellar dynamics, it is worth deriving this basic result. It states for any system of particles bound by an inverse-square force law, the time-averaged kinetic energy (say $\langle T \rangle$) and the time-averaged potential energy (say $\langle V \rangle$) satisfy

$$2 \langle T \rangle + \langle V \rangle = 0. \quad (2.1)$$

To prove this, consider the quantity

$$F = \sum_i m_i \dot{\mathbf{x}}_i \cdot \mathbf{x}_i \quad (2.2)$$

where m_i are the masses. Clearly

$$\frac{dF}{dt} = 2T + \sum_i m_i \ddot{\mathbf{x}}_i \cdot \mathbf{x}_i. \quad (2.3)$$

If F is bounded then the long-time average $\langle dF/dt \rangle$ will vanish. Thus

$$2 \langle T \rangle + \sum_i m_i \langle \ddot{\mathbf{x}}_i \cdot \mathbf{x}_i \rangle = 0. \quad (2.4)$$

If the system is gravitationally bound, we have

$$2 \langle T \rangle - G \sum_{ij} m_i m_j \left\langle \frac{(\mathbf{x}_i - \mathbf{x}_j)}{|\mathbf{x}_i - \mathbf{x}_j|^3} \cdot \mathbf{x}_i \right\rangle = 0. \quad (2.5)$$

Interchanging the dummy indices in the second term and adding, we have

$$2 \langle T \rangle - \frac{1}{2} G \sum_{ij} m_i m_j \left\langle \frac{1}{|\mathbf{x}_i - \mathbf{x}_j|} \right\rangle = 0. \quad (2.6)$$

But the second term is now just minus the total potential energy, which proves the result (2.1).

The virial theorem provides an easy way to make rough estimates of masses, because velocity measurements can give $\langle T \rangle$. But it is prudent to consider virial mass estimates as order-of-magnitude only, because (i) generally one can measure only line-of-sight velocities, and getting $T = \frac{1}{2} \sum_i m_i \dot{\mathbf{x}}_i^2$ from there requires more assumptions (e.g. isotropy of the velocity distribution); and (ii) the systems involved may not be in a steady state, in which case of course the virial theorem does not apply—clusters of galaxies are particularly likely to be quite far from a steady state.

TWO IMPORTANT TIME SCALES

Consider a stellar system of size R , having N stars each of mass m ; the stars are distributed roughly homogeneously, with v being a typical velocity, and the system is in dynamical equilibrium. Then from the virial theorem

$$v^2 \simeq NGm/R. \quad (2.7)$$

The crossing time (sometimes called dynamical time)

$$T_{\text{cross}} = \frac{R}{v} \simeq \sqrt{\frac{R^3}{NGm}} \quad \text{or} \quad \frac{1}{\sqrt{G\rho}}. \quad (2.8)$$

The relaxation time is how long it takes for a star's velocity to be changed significantly changed from two-body interactions. To estimate this, consider first one encounter, with a star going past another with impact parameter b . The change δv in the star's velocity due to this encounter is

$$\delta v = Gmb \int_{-\infty}^{\infty} \frac{dt}{(b^2 + v^2 t^2)^{\frac{3}{2}}} = \frac{2Gm}{bv}. \quad (2.9)$$

(Note that this will be perpendicular to the direction of motion.) Next we consider all the encounters in one crossing time with impact parameters in the range $(b, b + db)$. There are $2Nb db/R^2$ of these, since the surface density of stars is $N/(\pi R^2)$. The δv 's due these encounters will tend to cancel, so we add their squares and then integrate over b to get the total change in v^2 over one crossing time:

$$\Delta v^2(T_{\text{cross}}) = \int_{b_{\text{min}}}^R \left(\frac{2Gm}{bv}\right)^2 \frac{2N}{R^2} b db = 8N \left(\frac{Gm}{Rv}\right)^2 \ln\left(\frac{R}{b_{\text{min}}}\right). \quad (2.10)$$

The relaxation time T_{relax} is the time needed for $\Delta v^2 \simeq v^2$. Thus

$$T_{\text{relax}} = \frac{v^2}{\Delta v^2(T_{\text{cross}})} \times T_{\text{cross}} = \frac{1}{8N \ln(R/b_{\text{min}})} \frac{(Rv)^3}{(Gm)^2}. \quad (2.11)$$

It's easier to remember T_{relax} in crossing times. Taking $R/b_{\text{min}} \simeq N$ and then using equation (2.7) to eliminate R , we get

$$\frac{T_{\text{relax}}}{T_{\text{cross}}} \simeq \frac{N}{8 \ln N}. \quad (2.12)$$

Galaxies are $\lesssim 10^3 T_{\text{cross}}$ old and have $\gtrsim 10^6$ stars, so stellar encounters have negligible dynamical effect. In globular clusters, which may have $\sim 10^6$ stars and be $\sim 10^5$ crossing times old, stellar encounters start to become important, and in the cores of globular clusters two-body relaxation is very important.

PROBLEM 2.1: The v and m dependences of the relaxation time can actually be extracted by a back of the envelope calculation.

Consider N stars of mass m each in a box of side R , and let these stars be fixed. Then send another star through this box with speed v . How long does it take for the star to pass near enough to another star that kinetic and two-body potential energies are equal? (Order of magnitude only.) [15]

PROBLEM 2.2: Most researchers doing N -body simulations study the dynamics of galaxies, but some study the dynamics of globular clusters. The latter group of people would seem to have an easier job, because they can easily afford as many particles as there are stars, and they don't have to worry about gas dynamics. So you'd think that globular cluster dynamics would have been cleaned up by now. But in fact, globular cluster dynamics has *not* been cleaned up, and plenty of difficult research remains to be done. This problem is to work out why.

Consider a globular cluster and a galaxy, both $\sim 10^{10}$ yr old. The globular cluster has size ~ 100 pc and $\sim 10^6$ stars with typical velocity 50 km s^{-1} . The galaxy has ~ 10 kpc and $\sim 10^{10}$ stars with typical velocity 200 km s^{-1} . Now let's say both of these are simulated using 10^6 particles. Can you see two reasons why the globular cluster simulation will be more difficult? [10]

THE COLLISIONLESS BOLTZMANN EQUATION

In the absence of two-body relaxation, stars move under the total gravitational field of all other stars. This field depends only on location in space and we can express it by the potential $\Phi(\mathbf{x})$. Thus the motion of any star is given by Hamilton's equations

$$\frac{d\mathbf{x}}{dt} = \frac{\partial H}{\partial \mathbf{p}}, \quad \frac{d\mathbf{p}}{dt} = -\frac{\partial H}{\partial \mathbf{x}}, \quad (2.13)$$

with Hamiltonian

$$H = \frac{p^2}{2m} + \Phi(\mathbf{x}). \quad (2.14)$$

If you haven't met Hamiltonian mechanics before, not to worry: you can easily verify that equations (2.14) and (2.13) give the usual Newtonian equations; but remember the form of equations (2.14).¹ It's very useful to consider the density of stars in 6-dimensional 'phase' space (\mathbf{x}, \mathbf{p}) ; that density is called the distribution function and denoted by f .

Since stars are conserved, f must satisfy a continuity equation:

$$\frac{\partial f}{\partial t} + \frac{\partial}{\partial \mathbf{x}} \cdot \left(f \frac{d\mathbf{x}}{dt} \right) + \frac{\partial}{\partial \mathbf{p}} \cdot \left(f \frac{d\mathbf{p}}{dt} \right) = 0. \quad (2.15)$$

Substituting from Hamilton's equations gives

$$\frac{\partial f}{\partial t} + \frac{d\mathbf{x}}{dt} \cdot \frac{\partial f}{\partial \mathbf{x}} + \frac{d\mathbf{p}}{dt} \cdot \frac{\partial f}{\partial \mathbf{p}} \equiv \frac{df}{dt} = 0. \quad (2.16)$$

In Hamiltonian dynamics, (2.16) is known as Liouville's theorem, but in stellar dynamics it's usually called the collisionless Boltzmann equation. Physically, it means that if you move with a star, the phase space density around you stays constant. As the sun moves inwards in the Galaxy, the stellar density around it will increase, but at the same time the spread of stellar velocities around it will increase so as to keep phase space density constant.

¹ Hamiltonian dynamics is a beautiful subject in itself, and helps understand the relations—and differences—between classical mechanics and optics, quantum mechanics, and quantum field theory.

The collisionless Boltzmann equation, and the Poisson equation (which is the gravitational analogue of Gauss's law in electrostatics) together constitute the basic equations of stellar dynamics:

$$\frac{df}{dt} = 0, \quad \nabla^2 \Phi(\mathbf{x}) = 4\pi G \rho(\mathbf{x}). \quad (2.17)$$

EXAMPLE [*N*-body simulations] You have probably come across *N*-body simulations of stars in galaxies. The particles in a galaxy simulation do not correspond to stars. They cannot, they have too few particles (10^5 to maybe 10^8 particles max, versus maybe 10^{12} stars in the galaxies being modelled). The appropriate interpretation of simulation particles is as Monte-Carlo samplers of f . Simulation particles have to made collisionless artificially (since there are comparatively few of them, the two-body relaxation time will be correspondingly shorter). The standard way of doing this is to replace the $1/r$ gravitational potential by $(r^2 + a^2)^{-\frac{1}{2}}$, which amounts to smearing out the mass on the 'softening length' scale a .

N-body simulations are widely used now to study the evolution of galaxies, and a trendy research area at present is to incorporate gas dynamics in them. \square

Though f is a density in phase space, the full form of the collisionless Boltzmann equation doesn't have to be written in terms of \mathbf{x} and \mathbf{p} . We can express df/dt in any set of six variables in phase space.

EXAMPLE [Cylindrical coordinates] In terms of cylindrical coordinates R, ϕ, z and velocities v_R, v_ϕ, v_z we have

$$\frac{\partial f}{\partial t} + \dot{R} \frac{\partial f}{\partial R} + \dot{\phi} \frac{\partial f}{\partial \phi} + \dot{z} \frac{\partial f}{\partial z} + \dot{v}_R \frac{\partial f}{\partial v_R} + \dot{v}_\phi \frac{\partial f}{\partial v_\phi} + \dot{v}_z \frac{\partial f}{\partial v_z} = 0. \quad (2.18)$$

To eliminate the dots we use the standard relations for velocity and acceleration components. We have

$$\begin{aligned} \dot{R} &= v_R, & \dot{\phi} &= \frac{v_\phi}{R}, & \dot{z} &= v_z \\ \dot{v}_R &= -\frac{\partial \Phi}{\partial R} + v_\phi^2, & \dot{v}_\phi &= -\frac{1}{R} \frac{\partial \Phi}{\partial \phi} - \frac{v_R v_\phi}{R}, & \dot{v}_z &= -\frac{\partial \Phi}{\partial z}, \end{aligned} \quad (2.19)$$

where we have noted substituted $-\nabla \Phi$ for the acceleration. \square

You should remember that f is always taken to be a density in six-dimensional phase space, even in situations where it is a function of fewer variables. For example, if f happens to be a function of energy alone, it is not the same as the density in energy space.

ORBITS

The trajectories of individual stars (sometimes just called orbits) is in general highly chaotic. This can be so even if there is no collective motion at all (f in equilibrium). Actually, it's not difficult to appreciate why. Think about making bread, the baker's dough being a sort of fluid. Dough is incompressible, but that doesn't prevent you stretching it in one direction and shrinking it in others, and then folding it back. So while the dough keeps much the same shape, initially nearby particles within it can be dispersed to widely different parts of it, through the repeated stretching and folding. The same stretching and folding operation can take place in phase space. In fact it appears that phase space is typically riddled with regions where f gets stretched in one directions while being shrunk in others. Thus nearby orbits tend to diverge, and the divergence is exponential in time, which is the technical definition of chaos in dynamical systems. Simulations suggest that the e -folding time of the divergence is comparable to T_{cross} , and gets shorter the more particles there are.

In some special situations, there is no chaos, and the system is said to be 'integrable'. If the dynamics is confined to one real-space dimension (hence two phase-space dimensions) then no stretching-and-folding can happen, and orbits are regular. So in a spherical system all orbits are regular. In addition, there are certain potentials (usually referred to as Stäckel potentials) where the dynamics decouples into three effectively one-dimensional systems; so if some equilibrium f generates a Stäckel potential, the orbits will stay chaos-free. Also, small perturbations of non-chaotic systems tend to produce only small regions of chaos,² and orbits may be well described through perturbation theory.

In integrable systems there are significant simplifications. Each orbit is (i) confined to a three-dimensional toroidal subspace of six-dimensional phase space, and (ii) fills its torus evenly.³ Phase space itself is filled by nested orbit-carrying tori—they have to be nested, since orbits can't cross in phase space. Therefore the time-average of each orbit is completely specified once we have specified which torus it is on; this takes three numbers for each orbit, and these are called 'isolating integrals'—they are constants for each orbit of course. Think of the isolating integrals as a coordinate system that parametrizes orbital tori; transformations to a different set of isolating integrals is like a coordinate transformation.

If isolating integrals exist, then any f that depends only on them will automatically satisfy the collisionless Boltzmann equation. Conversely, since orbits fill their tori evenly, any equilibrium f cannot depend on location *on* the tori, it can only depend on the tori themselves, i.e., on the isolating integrals. This result is known as Jeans' theorem.

You should be wary of Jeans' theorem, especially when people tacitly assume it, because as we saw, it assumes that the system is integrable, which is in general not the case.

² If you ever come across the 'KAM theorem', that's basically it.

³ These two statements are important results from Hamiltonian dynamical systems which we won't try to prove here. But the statements that follow in this section are straightforward consequences of (i) and (ii).

SPHERICAL SYSTEMS

In spherical systems Jeans' theorem does apply, so f can depend on (at most) three integrals of motion. The simplest case is for f to be a function of energy $E = \frac{1}{2}v^2 + \Phi$ only. (Since we are considering bound systems, $f = 0$ for $E < 0$ always.) To find an equilibrium solution, we only have to satisfy Poisson's equation.

We'll take $G = 1$ for this section, to simplify the expressions a bit. Poisson's equation is now

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi}{dr} \right) = (4\pi)^2 \int_0^{\sqrt{-2\Phi}} f\left(\frac{1}{2}v^2 + \Phi\right) v^2 dv. \quad (2.20a)$$

We can also replace the integral over v by an integral over E :

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi}{dr} \right) = (4\pi)^2 \sqrt{2} \int_{\Phi}^0 \sqrt{E - \Phi} f(E) dE \quad (2.20b)$$

In (2.20a) we take $f(E)$ as given and try to solve for Φ and hence $\rho(r)$; this is a nonlinear differential equation. In (2.20b) we take Φ as given, and try to solve for $f(E)$; this is a linear integral equation.

There are $f(E)$ models in the literature, and you can always concoct a new one by picking some $\rho(r)$, computing $\Phi(r)$ and then solving (2.20b) numerically. Note that the velocity distribution is isotropic for any $f(E)$. If f depends on other integrals of motion, say angular momentum L or its z component, or both—thus $f(E, L^2, L_z)$ —then the velocity distribution will be anisotropic, and there are many examples of these around too.

EXAMPLE [Two spherical isotropic distribution functions] The Plummer model has

$$\Phi(r) = -(r^2 + a^2)^{-\frac{1}{2}}, \quad \rho(r) = -\frac{3a^2}{4\pi} \Phi^5, \quad (2.21)$$

and the distribution function

$$f(E) = \frac{24\sqrt{2}a^2}{7\pi^3} (-E)^{\frac{7}{2}}. \quad (2.22)$$

can be verified by inserting in (2.20a). Because of the simple functional forms, the Plummer model is occasionally useful for doing rough calculations, but the r^{-5} density profile is much steeper than elliptical galaxies are observed to have.

The isothermal sphere is defined by analogy with a Maxwell-Boltzmann gas, as

$$f(E) = \frac{\rho_0}{(2\pi\sigma^2)^{\frac{3}{2}}} \exp\left(-\frac{E}{\sigma^2}\right) = \frac{\rho_0}{(2\pi\sigma^2)^{\frac{3}{2}}} \exp\left(-\frac{v^2 + \Phi}{\sigma^2}\right), \quad (2.23)$$

and σ^2 is like a temperature. Integrating over velocities gives

$$\rho = \rho_0 \exp\left(-\frac{\Phi}{\sigma^2}\right). \quad (2.24)$$

Using this, Poisson's equation becomes

$$\frac{d}{dr} \left(r^2 \frac{d \ln \rho}{dr} \right) = -\frac{4\pi}{\sigma^2} r^2 \rho, \quad (2.25)$$

for which the solution is

$$\rho(r) = \frac{\sigma^2}{2\pi r^2}, \quad (2.26)$$

or $\sigma^2/(2\pi Gr^2)$ if we put back the G . The isothermal sphere has infinite mass! (A side effect of this is that the boundary condition $\Phi(\infty) = 0$ cannot be used, which is why we needed the redundant-looking constant ρ_0 in (2.24) and (2.23).) Nevertheless, it is often used as a model, with some large- r truncation assumed, for the dark halos of disc galaxies. \square

The same $\rho(r)$ can be produced by many different f , all having different velocity distributions.

THE JEANS EQUATIONS

Phase space quantities are hard to measure. Much more often we have information only about averages, e.g., bulk velocities and velocity dispersions. So it is useful to derive equations for the quantities

$$\begin{aligned}\rho &= \int f d^3\mathbf{v}, \\ \rho \langle v_i \rangle &= \int v_i f d^3\mathbf{v}, \\ \rho \sigma_{ij} &= \int (v_i - \langle v_i \rangle)(v_j - \langle v_j \rangle) f d^3\mathbf{v}.\end{aligned}\tag{2.27}$$

by taking moments of the collisionless Boltzmann equation (expressed in the cartesian variables x_i and v_i).

Consider first the zeroth moment

$$\int \left(\frac{\partial f}{\partial t} + v_i \frac{\partial f}{\partial x_i} - \frac{\partial \Phi}{\partial x_i} \frac{\partial f}{\partial v_i} \right) d^3\mathbf{v}.\tag{2.28}$$

If we integrate the last term by parts (equivalently, apply the divergence theorem) and assume that f and its derivatives vanish for large enough \mathbf{v} , the term vanishes. In the middle term we can take the gradient outside the integral. This gives us

$$\frac{\partial \rho}{\partial t} + \frac{\partial \rho \langle v_i \rangle}{\partial x_i} = 0,\tag{2.29}$$

which is a continuity equation.

Now we consider the first moment

$$\int \left(v_i \frac{\partial f}{\partial t} + v_i v_j \frac{\partial f}{\partial x_j} - \frac{\partial \Phi}{\partial x_j} v_i \frac{\partial f}{\partial v_j} \right) d^3\mathbf{v}.\tag{2.30}$$

Again, we integrate the last term by parts, and since

$$\int v_i \frac{\partial f}{\partial v_j} d^3\mathbf{v} = - \int \delta_{ij} f d^3\mathbf{v},$$

we get

$$\frac{\partial \rho \langle v_i \rangle}{\partial t} + \frac{\partial \rho \langle v_i v_j \rangle}{\partial x_j} = -\rho \frac{\partial \Phi}{\partial x_i}.\tag{2.31}$$

From this we subtract $\langle v_i \rangle$ times the continuity equation, and then substitute

$$\langle v_i v_j \rangle = \sigma_{ij} + \langle v_i \rangle \langle v_j \rangle,$$

to get

$$\rho \frac{\partial \langle v_i \rangle}{\partial t} + \rho \langle v_j \rangle \frac{\partial \langle v_i \rangle}{\partial x_j} = -\rho \frac{\partial \Phi}{\partial x_i} - \frac{\partial \rho \sigma_{ij}}{\partial x_j},\tag{2.32}$$

which is the same as⁴

$$\frac{d\langle \mathbf{v} \rangle}{dt} = -\nabla\Phi - \frac{1}{\rho}\nabla \cdot (\rho\boldsymbol{\sigma}). \quad (2.33)$$

Finally we have an equation that reminds us of ordinary fluid dynamics but also shows us why a stellar fluid is different. An ordinary fluid has

$$\frac{d\langle \mathbf{v} \rangle}{dt} = -\nabla\Phi - \frac{p}{\rho} + \text{viscous terms.} \quad (2.34)$$

where the pressure p arises because of the high rate of molecular encounters, which also leads to the equation of state, and p is isotropic. In a stellar fluid $\nabla \cdot (\rho\boldsymbol{\sigma})$ behaves like a pressure, but it is anisotropic. A related fact is that in the flow of an ordinary fluid the particle paths and streamlines coincide, whereas stellar orbits and the streamlines $\langle \mathbf{v} \rangle$ do not generally coincide.

EXAMPLE [Useful forms of the hydrodynamic equation] In a steady state axisymmetric system like a disc galaxy we use cylindrical coordinates, and then $\partial/\partial t = \partial/\partial\phi = 0$. Neglecting $\langle v_R \rangle$ and $\langle v_z \rangle$ (which is realistic), we have

$$\begin{aligned} \frac{1}{R} \frac{\partial}{\partial R} (R\rho\sigma_{RR}) + \frac{\partial}{\partial z} (\rho\sigma_{Rz}) - \frac{\rho}{R} \left(\langle v_\phi \rangle^2 + \sigma_{\phi\phi} \right) &= -\rho \frac{\partial\Phi}{\partial R}, \\ \frac{\partial}{\partial z} (\rho\sigma_{zz}) + \frac{1}{R} \frac{\partial}{\partial R} (R\rho\sigma_{Rz}) &= -\rho \frac{\partial\Phi}{\partial z}. \end{aligned} \quad (2.35)$$

For a steady state spherical system, we use spherical polar coordinates, so $\partial/\partial\theta = \partial/\partial\phi = 0$, and then neglect $\langle v_r \rangle$ and $\langle v_\theta \rangle$. Then we have

$$\frac{d}{dr} (\rho\sigma_{rr}) + \frac{\rho}{r} \left[2\sigma_{rr} - \left(\sigma_{\theta\theta} + \sigma_{\phi\phi} + \langle v_\phi \rangle^2 \right) \right] = -\rho \frac{d\Phi}{dr}. \quad (2.36)$$

These forms are quite useful. Note that Φ is the total potential but $\rho, \langle \mathbf{v} \rangle, \langle \boldsymbol{\sigma} \rangle$ could be for any subpopulation.

As a simple test to see if this apparatus really does work, let us make a crude model of the Milky way halo. We take $\Phi = v_0^2 \ln r$, assume $\boldsymbol{\sigma}$ is constant and isotropic with all diagonal components = σ^2 (say). Then we say $\rho \propto r^{-n}$ and $\langle v_\phi \rangle = 0$. This gives $\sigma = v_0/\sqrt{n}$. For the Milky Way halo, $\rho \propto r^{-3.5}$, v_0 as measured from gas on circular orbits is 220 km/sec, and rotation is negligible. So we expect $\sigma \simeq 120$ km/sec. And it is. \square

PROBLEM 2.3: An E0 galaxy has a total density distribution

$$\rho_{\text{tot}}(r) = \frac{\rho_0}{1 + r^2/a^2}.$$

Show that the enclosed mass $M(r) \propto r^3$ for $r \ll a$ and $M(r) \propto r$ for $r \gg a$. [3]

Now take a population of massless test particles in the potential of this galaxy. Assume that this population is spherical, non-rotating, isothermal and isotropic, with velocity dispersion σ in each velocity component. What is the radial density distribution of this test particle population? [8]

At large r the test particle distribution simplifies and its form depends on a dimensionless number. Give a physical interpretation of this number. What is the condition that the density distributions of the test particle population and the galaxy itself have similar forms at large r ? [7]

⁴ Note that d/dt is not $\partial/\partial t$, but

$$\frac{d}{dt} \equiv \frac{\partial}{\partial t} + \mathbf{v} \cdot \nabla$$

sometimes called the convective derivative; also sometimes written as D/Dt to emphasize that it's not $\partial/\partial t$.

3. The Interstellar Medium

The interstellar medium (ISM) is a mixture of the primordial gas left over from galaxy formation and the material spewed out by dying stars. It is only a few percent of a galaxy's mass, and very very diffuse ($\lesssim 10^3$ atoms cm^{-3}). But it is very important because it is the stuff that forms stars. It is also the site of varied physical processes that make it observable and fascinate the people studying them.

GAS

Under laboratory conditions, spectral lines with low transition probabilities are 'forbidden' because the excited states get collisionally de-excited before they can radiate. In the ISM, collisional times are typically much longer than the lifetimes of excited states with only forbidden transitions. So forbidden lines are observable from the ISM, and in fact they can dominate the spectrum.

Cold gas emits only in radio, and the most important ISM line of all is the 21 cm line of atomic hydrogen (HI). It comes from the hyperfine split ground state of the hydrogen atom (split because of the coupling of the nuclear and electron spins). The spin flip transition itself cannot be observed in a laboratory, but the split ground state shows up in the hyperfine splitting of the Lyman lines. HI is observed in both emission or in absorption against a background continuum source. One of the uses of HI observations is to measure rotation velocities of gas. Molecular hydrogen (H_2) has no radio lines, which is unfortunate, since it prevents the coldest and densest parts of the ISM being absorbed directly. What saves the situation somewhat is that CO has strong lines at 1.3 mm and 2.6 mm from transition between rotational states, and CO gets used as a tracer of H_2 .

Hot gas is readily observed in optical. An important kind of object are HII regions, which partially ionized hydrogen surrounding a very hot young star or stars (O or B). Hot stars produce a large flux of ultraviolet photons, and any Lyman continuum photons (i.e., $\lambda < 912 \text{ \AA}$) will photoionize hydrogen. The ionized hydrogen then recombines. But it doesn't have to recombine into the ground state, it can recombine into an excited state and then radiatively decay after that. This process produces a huge variety of observable lines and continuums, of Lyman, Balmer and on through the infrared and into radio. Of each series, the longest wavelength (or α) line will be the strongest, because the transition rate from principal quantum number n is strongest to $n - 1$. Atoms in HII regions can also be collisionally excited. Atomic hydrogen has no levels accessible at collision energies characteristic of HII regions ($T \sim 10^4 \text{ K}$) but NII, OII, SII, OIII, NeIII all do. The [OIII] lines at 4959 \AA and 5007 \AA are particularly prominent.

A planetary nebula is like a compact HII region, except that it surrounds the exposed core of a highly evolved star rather than a hot young star. Because of their bright emission lines and compactness, planetary nebulae can be detected from much greater distances than individual ordinary stars; they are used as sort of tracers of stars.

The photoionization and recombination process in HII regions and planetary nebulae produces, by a happy accident, one Balmer photon for each Lyman continuum photon from the hot star, so the UV flux can be measured by observing an optical spectrum. The reason is basically that the gas is opaque to Lyman photons and transparent to other photons, since almost the H atoms are in the ground state. A Lyman

continuum photon initially from the star will get absorbed by a hydrogen atom, producing a free electron. This electron will then be captured into some bound state. If it gets captured to the ground state we are back where we started (with a ground state atom and a Lyman continuum photon), so consider the case where the electron is captured to some $n > 1$ state. Such a capture releases a free-bound continuum photon which then escapes, and leaves an excited state which wants to decay to $n = 1$. If it decays to $n = 1$ bypassing $n = 2$, it will just produce a Lyman photon which will get almost certainly get absorbed again. Only if it decays to some $n > 1$ will a photon escape. In other words, if the decay bypasses $n = 2$ it almost always gets another chance to decay to $n = 2$ and produce a Balmer photon that escapes. The Ly α photons produced by the final decay from $n = 2$ to $n = 1$ random-walk through the gas as they get absorbed and re-emitted again and again. The total Balmer photon flux thus equals the Lyman continuum photon flux. One can then place the source star in an optical-UV colour magnitude diagram, and determine a colour temperature which is called the Zanstra temperature in this context.

H II regions and planetary nebulae also produce thermal continuum radiation. The process that produces this is free-free emission: free electrons in the H II can interact with protons without recombination, and the acceleration of the electrons in this process produces radiation. (Electrons can interact with other electrons in similar fashion as well, but this produces no radiation because the net electric dipole moment doesn't change.) The resulting spectrum is not blackbody because the gas is transparent to free-free photons. In fact the spectrum is quite flat at radio frequencies—this is the same thing as saying that the time scale for free-free encounters is $\ll 1/\nu$ for radio frequency ν .

When an interstellar gas cloud is seen in front of a continuum source, it produces an enormous variety of absorption lines and bands, by no means all of them well understood. Perhaps the most puzzling ones are the so-called diffuse interstellar bands in the infrared; apparently these are similar to what you get if you take bacteria out of the river at Cardiff and stick them in a spectrograph, which led to some interesting speculations some years ago. . .

EXAMPLE [Cold interstellar CN] Here is a really cute (and slightly poignant) example of what interstellar absorption lines can do for you. Like most heteronuclear molecules, CN has rotational modes which produce radio lines. The radio lines can be observed directly, but more interesting are the optical lines that have been split because of these rotational modes. Observations of cold CN against background stars reveal, through the relative widths of the split optical lines, the relative populations of the rotational modes, and hence the temperature of the CN. The temperature turns out to be 2.73 K, i.e., these cold clouds are in thermal equilibrium with the microwave background. The temperature of interstellar space was first estimated as $\simeq 3$ K in 1941, well before the Big Bang predictions of 1948 and later, but nobody made the connection at the time. \square

In the highest density H II regions ($\sim 10^8 \text{ cm}^{-3}$), either very near a young star, or in a planetary-nebula-like system near the evolved star, population inversion between certain states becomes possible. The overpopulated excited state then decays by stimulated emission, i.e., it becomes a maser. An artificial maser or laser uses a cavity with reflecting walls to mimic an enormous system, but in an astrophysical maser the enormous system is available for free; so an astrophysical maser is not directed perpendicular to some mirrors but shines in all directions. But as in an artificial maser,

the emission is coherent (hence polarized), with very narrow lines and high intensity. Masers from OH and H₂O are known. Their high intensity and relatively small size makes masers very useful as kinematic tracers.

Finally, we'll just briefly mention synchrotron radiation, which you'll cover in more detail in the high-energy astrophysics part. It's a broad-band non-thermal radiation emitted by electrons gyrating relativistically in a magnetic field, and can be observed in both optical and radio. The photons are emitted in the instantaneous direction of electron motion and polarized perpendicular to the magnetic field. The really spectacular sources of synchrotron emission are systems with jets (young stellar objects with bipolar outflows, or active galactic nuclei). It is synchrotron emission that lights up the great lobes of radio galaxies.

DUST

Interstellar dust consists of particles of silicates or carbon compounds; the largest are $\simeq 0.5 \mu\text{m}$ with $\sim 10^4$ atoms, but some appear to have $\lesssim 10^2$ atoms and thus might be thought of as large molecules. Their nastiest property is that they absorb and scatter light, and the observational effect of these two are known as extinction. (Extinction in magnitudes is denoted as A_V for V -band and so on.) Extinction gets less severe for $\lambda \gtrsim 1 \mu\text{m}$ as the wavelength gets much longer than the grains, but it is worse for blue than red light. Hence objects are said to be 'reddened' by interstellar extinction. Grains are transparent to X-rays, though. From our location, extinction is worst along the Milky Way disc, and the Galactic Centre is completely opaque to optical observations.

However, extinction by dust does one very useful thing for optical astronomers. Spinning dust grains tend to align with their long axes perpendicular to the local magnetic field. They thus preferentially block light perpendicular to the magnetic field. Thus the observed polarization will tend to be parallel to the magnetic field. Hence polarization measurements of starlight reveal the direction of the magnetic field (or at least the sky-projection of the direction).

Dust also reflects light, with some polarization. This is observable as reflection nebulae, where the stars cannot be seen (at least in optical) but faint diffuse starlight can be seen as reflected by dust.

Light absorbed by dust will be reradiated as a blackbody-ish spectrum. Such a spectrum is observed (from space, by IRAS) as diffuse emission superimposed on a reflected starlight spectrum, but the associated temperature is extremely high— $\sim 10^3$ K. The interpretation is that some dust grains are so small (< 100 atoms) that a single ultraviolet photon packs enough energy to heat them to $\sim 10^3$ K, after which these 'stochastically heated' grains cool again by radiating, mostly in the infrared. This process may be part of the explanation for the correlation between infrared and radio continuum luminosities of galaxies (e.g., at 0.1 mm and 6 cm), which seems to be independent of galaxy type. The idea is that ultraviolet photons from the formation of massive stars cause stochastic heating of dust grains, which then reradiate them to give the infrared luminosity. The supernovae resulting from the same stellar populations produce relativistic electrons which produce the radio continuum as synchrotron emission.

CHEMICAL ENRICHMENT

The birth and death of stars and what that does to the interstellar medium is a large and very important subject. We won't be able to do it any sort of justice, but just for a sampler let's discuss the effect on chemical evolution of the ISM.

Consider a region of a galaxy, small enough to be fairly homogeneous, but large enough to contain a good sample of stars. Suppose at time t , the total mass of this region is M_{total} , M_{stars} in stars and M_{gas} in gas; also say M_{metal} is the part of M_{gas} in metals. Thus the metallicity of the gas is

$$Z \equiv \frac{M_{\text{metal}}}{M_{\text{gas}}}. \quad (3.1)$$

Now we consider the effect of forming some new stars over some time δt . This time is longer than the time massive stars spend on the main sequence, so the newly formed massive stars are supposed to have already gone supernova and spewed some more metals into the ISM. Let δM_{stars} be the change of stellar (or stellar remnant) mass, and let the metal mass contributed to the ISM by this generation of stars be $p\delta M_{\text{stars}}$ (p is known as the 'yield' and we will take it to be constant). We want to find the time evolution of Z , from

$$\delta Z = \delta \left(\frac{M_{\text{metal}}}{M_{\text{gas}}} \right) = \frac{\delta M_{\text{metal}} - Z\delta M_{\text{gas}}}{M_{\text{gas}}} \quad (3.2)$$

We will assume that the system starts with only gas and at $Z = 0$.

The simplest approximation is the 'closed box model', where gas and stars neither enter nor leave this region of the galaxy. Then

$$\delta M_{\text{metal}} = p\delta M_{\text{stars}} - Z\delta M_{\text{stars}} = (p - Z)\delta M_{\text{stars}} \quad (3.3)$$

and

$$0 = \delta M_{\text{total}} = \delta M_{\text{stars}} + \delta M_{\text{gas}}. \quad (3.4)$$

Inserting these in equation (3.2) gives

$$\delta Z = -p \frac{\delta M_{\text{gas}}}{M_{\text{gas}}}, \quad (3.5)$$

whence

$$Z = -p \ln \left(\frac{M_{\text{gas}}(t)}{M_{\text{gas}}(0)} \right). \quad (3.6)$$

In other words,

$$Z = -p \ln(\text{gas fraction}).$$

Magellanic irregulars fit this reasonably well, and p is estimated to be $\simeq 0.0025$. In spiral galaxies, the gas fraction in the disc increases as we go outwards, and Z is observed to decrease, though perhaps more steeply than this crude model predicts.

The closed box model can also be used to calculate the distribution of stellar metallicities, because the metallicity of each star approximately indicates Z when that star was formed. If we take all the stars now with metallicities less than some Z_1 , the

sum of their masses equals $M_{\text{stars}}(t)$ for the t when Z equalled Z_1 . To get $M_{\text{stars}}(t)$ we rewrite (3.5) as

$$\delta Z = \frac{p\delta M_{\text{stars}}}{M_{\text{gas}}(0) - M_{\text{stars}}(t)} \quad (3.7)$$

which gives

$$M_{\text{stars}}(t) = \left(1 - e^{-Z/p}\right) M_{\text{gas}}(0). \quad (3.8)$$

This gives a tolerably good fit for metal-poor globular clusters. But it fails badly for the solar neighbourhood: the most metal rich stars have $Z \simeq Z_{\odot} \simeq 0.02$, and (3.8) predicts that $\sim 50\%$ of solar neighbourhood stars will have $Z \leq \frac{1}{4}Z_{\odot}$; in fact only about 2% do. This is known as the ‘G-dwarf problem’.

The G-dwarf problem indicates that the closed-box model is an oversimplification, and that loss and/or accretion of material into a star-forming region needs to be considered.

PROBLEM 3.1: In this problem we consider a ‘leaky-box’ model, which simulates the effect of shocks from supernovae and winds from young massive stars by making gas leave the formerly closed box at a rate proportional to the star formation rate:

$$\delta M_{\text{total}} = -c\delta M_{\text{stars}}.$$

Use this to work out $M_{\text{gas}}(t)$ in terms of $M_{\text{total}}(0)$ and $M_{\text{stars}}(t)$. Now modify the closed-box relation between δM_{metal} and δM_{stars} by adding an appropriate leaking term. [6]

Use these two expressions to derive

$$\delta Z = \frac{p\delta M_{\text{stars}}}{M_{\text{total}}(0) - (1+c)M_{\text{stars}}}. \quad [2]$$

This expression shows that the leaky box model won’t solve the G-dwarf problem? Why? [5]

If we allow the box to accrete gas, that does make metal poor stars rarer.

PROBLEM 3.2: In this problem we consider the ‘accreting-box’ model, another modification of the closed-box model, this time allowing for metal-free gas to be accreted into the system.

From the assumption that no metal enters or leaves the region, relate δM_{metal} and δM_{stars} . Allowing for (metal-free) gas accretion, relate δM_{stars} to δM_{total} and δM_{gas} . Use the above to show that

$$\delta Z = \frac{(p-Z)\delta M_{\text{total}} - p\delta M_{\text{gas}}}{M_{\text{gas}}}. \quad [8]$$

This equation can be solved exactly with some awkwardness, but for us it’s enough to consider the simplest case whether the gas accretion rate equals the star formation rate, so M_{gas} stays constant. For this simple case show that Z asymptotes to p . [6]

Can you argue physically why we should expect such behaviour for stellar metallicities in this case? [5]

In fact this model predicts that $\simeq 3\%$ of solar neighbourhood stars will have $Z \leq \frac{1}{4}Z_{\odot}$.

4. Rotation Curves

Gas and young stars will move on nearly closed orbits, and if the underlying potential is axisymmetric these will be nearly circular. So if you measure the bulk velocity v (of gas or young stars, *not* old stars) at any place on a galactic disc, you've measured $R(\partial\Phi/\partial R)$; and if you measure $v(R)$ —the ‘rotation curve’—you have information on the mass distribution.

When people first starting measuring rotation curves (c. 1970), it quickly became clear that the mass in disc galaxies doesn't follow the visible disc. Disc galaxies generically have rotation curves that are fairly flat to as far out as they can be measured (several scale radii). The simplest interpretation of a flat rotation curve is that enclosed mass $M(r) \propto r$, or $\rho(r) \propto 1/r^2$, a ‘dark halo’. The deep picture of M104 in part 1 of these notes suggests that dark halos are not entirely dark, but as yet nobody knows really knows how far they extend. And there is no good estimate of the total mass of any disc galaxy. This is what makes disc rotation curves very important.

However, one needs to be a little careful about interpreting flat rotation curves. The maximum contribution to the rotation curve from an e^{-R/R_0} disc is not (as we might naively expect) around R_0 but around $2.5R_0$. Adding the effect of a bulge can easily give a fairly flat rotation curve to $4R_0$ without a dark halo. To be confident about the dark halo, one needs to have the rotation curve for $\gtrsim 5R_0$. In practice, that means HI measurements; optical rotation curves don't go out far enough to say anything about dark halos.

The rest of part 4 is a more detailed working out of the previous paragraph. It follows an elegant derivation and explanation due to A.J. Kalnajs.

The potential from a disc with surface density $\Sigma(R)$ is

$$\Phi(R) = -G \int_0^\infty R' \Sigma(R') dR' \int_0^{2\pi} \frac{d\phi}{\sqrt{R^2 + R'^2 - 2RR' \cos \phi}}. \quad (4.1)$$

To make this tractable, let us first define¹

$$2\pi L(u) \equiv \int_0^{2\pi} \frac{d\phi}{\sqrt{1 + u^2 - 2u \cos \phi}} = 1 + \left(\frac{1}{2}u\right)^2 + \left(\frac{\frac{1}{2} \frac{3}{2} u^2}{2!}\right)^2 + \dots \quad (u < 1). \quad (4.2)$$

Then

$$\Phi(R) = -2\pi G \int_0^R \Sigma(R') \left(\frac{R'}{R}\right) L\left(\frac{R'}{R}\right) dR' - 2\pi G \int_R^\infty \Sigma(R') L\left(\frac{R}{R'}\right) dR', \quad (4.3)$$

and hence

$$\begin{aligned} v^2(R) = R \frac{\partial\Phi}{\partial R} = & 2\pi G \int_0^R \left[\left(\frac{R'}{R}\right) L\left(\frac{R'}{R}\right) + \left(\frac{R'}{R}\right)^2 L'\left(\frac{R'}{R}\right) \right] \Sigma(R') dR' \\ & - 2\pi G \int_R^\infty \left(\frac{R}{R'}\right) L'\left(\frac{R}{R'}\right) \Sigma(R') dR'. \end{aligned} \quad (4.4)$$

¹ If you really want to know where that came from, look up any musty old celestial mechanics book under ‘Laplace coefficients’.

(Here L' means a derivative!) The important thing to take away with you is not the algebraic mess but the form of the relation, which is

$$v^2(R) = 2\pi G \int_0^\infty K\left(\frac{R}{R'}\right) \Sigma(R') dR'. \quad (4.5)$$

Changing variables to

$$x = \ln R, \quad y = \ln R',$$

we can write this as a convolution

$$v^2(R) = 2\pi G \int_{-\infty}^\infty K(e^{x-y}) R' \Sigma(R') dy. \quad (4.6)$$

The kernel $K(R/R')$ is in Figure 4.1.

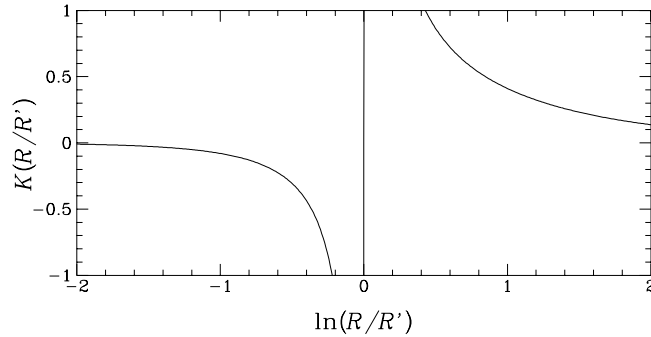


Figure 4.1: The kernel $K(R/R')$. Observe that the $R > R'$ part tends to have higher absolute value than the $R < R'$ part.

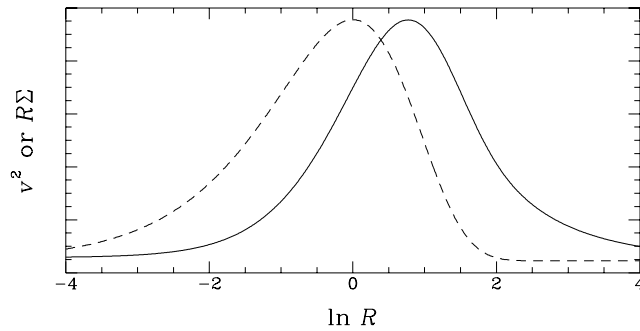


Figure 4.2: The dashed curve is $R\Sigma(R)$ for an exponential disc with $\Sigma \propto e^{-R}$ and the solid curve is $v^2(R)$. Note that R is measured in disc scale lengths, but the vertical scales are arbitrary.

Figure 4.2 shows $R\Sigma(R)$ and v^2 for an exponential disc, but the general shapes aren't very sensitive to whether $\Sigma(R)$ is precisely exponential. The important qualitative fact is that whatever $R\Sigma(R)$ does, v^2 does roughly the same, but expanded by a factor of $\simeq e$.

The distinctive shape of the $v^2(\ln R)$ curve for realistic discs makes it very easy to recognize *non*-disc mass. Figure 4.3, following Kalnajs, shows the rotation curves you

get by adding either a bulge or a dark halo. (Actually this figure fakes the bulge/halo contribution by adding a smaller/larger disc; but if you properly add spherical mass distributions for disc/halo, the result is very similar.) Kalnajs' point is that a bulge+disc rotation curve has a similar shape to a disc+halo rotation curve—only the scale is different. So when examining a flat(-ish) rotation curve, you must ask what the disc scale radius is.

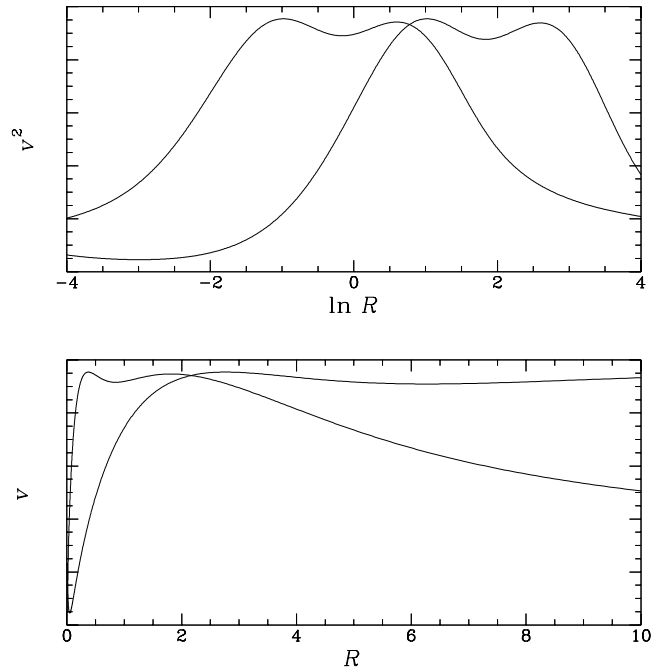


Figure 4.3: Plots of v^2 against $\ln R$ (upper panel) or v against R (lower panel) For one curve in each panel, a second exponential disc with mass and scale radius both scaled down by $e^2 \simeq 7.39$ has been added (to mimic a bulge); for the other curve a second exponential disc with mass and scale radius both scaled up by $e^2 \simeq 7.39$ has been added (to mimic a dark halo).

PROBLEM 4.1: Express the integral equation (4.3) relating $\Phi(R)$ and $\Sigma(R)$ as a convolution in $\ln R$. [10]

The convolution kernel differs from $K(R/R')$ of course, and in a particularly interesting way in the $R/R' \ll 1$ limit. Can you explain this difference using a physical argument? [10]

5. Gravitational Lensing

Gravitational lensing is about how the appearance of distant bright objects is altered by the gravity of foreground mass. Being a purely gravitational effect makes lensing astrophysically important as a probe of dark matter.

This part is more detailed than it needs to be. Only the section on microlensing in the Milky Way is really syllabus material. The rest you should consider as relevant background material plus general interest.

Photons are affected by a gravitational field, but not in the same way as massive particles are. For the details we need general relativity, but fortunately, for astrophysical applications we only need to take over a few simple results. The most important is that if a light ray passes by a mass M with impact parameter R ($\gg GM/c^2$ and \gg the size of the mass), it gets deflected by an angular amount

$$\alpha = \frac{4GM}{c^2 R}. \quad (5.1)$$

In contrast, a massive body at high speed v gets deflected by $\alpha = 2GM/(v^2 R)$.

THE LENSING EQUATION

To make (5.1) useful we need two approximations, both very good in almost all astrophysical situations:

- (i) The deflector is much smaller than the distances to the observer and the object being viewed (the ‘source’);
- (ii) The deflections are always very small, so we can freely use $\sin \alpha = \alpha$, and also we can get the total deflection from a mass distribution by integrating (5.1).

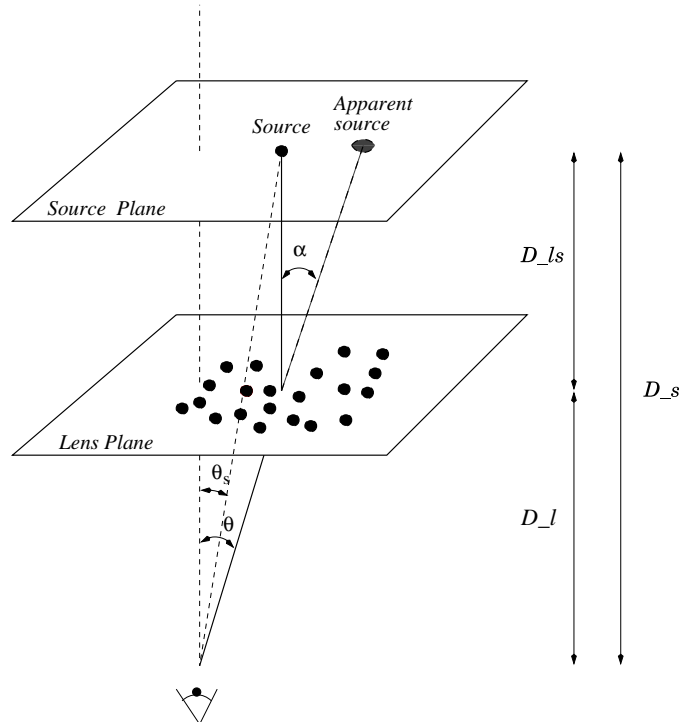


Figure 5.1: Definitions of D_L , D_S , D_{LS} , θ , θ_S , and α .

Accordingly, let us consider a situation as in Figure 5.1: observer is viewing a source at distance D_S , with a lens (a mass screen) intervening at distance D_L ; D_{LS} is the distance from lens to the source.¹ We'll use angular coordinates for the transverse position.² Thus, θ_S is the position of the source, θ is its observed position after being deflected—note that these are two-dimensional angles. Let $\Sigma(\theta)$ be the lens's density surface mass density (as in solar masses per steradian). Let $\alpha(\theta)$ be the deflection angle. Then, comparing vectors in the source plane, we get

$$D_S\theta = D_S\theta_S + D_{LS}\alpha. \quad (5.2)$$

(By convention,³ α is directed outwards from the deflecting mass rather than towards it.) Using (5.1) to get α in terms of Σ , we get

$$\theta = \theta_S + \frac{D_{LS}}{D_S}\alpha(\theta), \quad \alpha(\theta) = \frac{4G}{c^2 D_L} \int \frac{\Sigma(\theta')(\theta - \theta') d^2\theta'}{|\theta - \theta'|^2}. \quad (5.3)$$

This is known as the lens equation. It gives θ_S as an explicit function of θ , but θ as an implicit function of θ_S . Moreover, $\theta(\theta_S)$ need not be single-valued, so sources can be multiply imaged.

THE ARRIVAL TIME SURFACE

It's possible to work entirely with the form (5.3), but there's a much more intuitive reformulation, which we'll now derive.

We start by noting that the lens equation (5.3) amounts to equating a gradient to zero:

$$\begin{aligned} \nabla T &= 0, \quad T = \frac{1}{2}T_0(\theta - \theta_S)^2 - \Psi(\theta), \\ \Psi(\theta) &= \frac{4G}{c^3} \int \Sigma(\theta') \ln |\theta - \theta'| d^2\theta', \quad T_0 = \frac{D_L D_S}{c D_{LS}}. \end{aligned} \quad (5.4)$$

The two terms in T express the change in light travel time for an arbitrary deflection:⁴ the first term is what we would get from geometrical considerations alone; the second term is an extra time delay caused by the gravitational field.⁵ The requirement that T be stationary is just Fermat's principle.

Next we consider a point mass M , which happens to be precisely between us and a point source. In other words $\theta_S = 0$ and $\Sigma(\theta) = M\delta(\theta)$. Then the lens equation is solved by $\theta = \theta_E$, with

$$\theta_E^2 = \frac{4GM}{c^2} \frac{D_{LS}}{D_L D_S}, \quad R_E^2 = \frac{4GM}{c^2} \frac{D_L D_{LS}}{D_S}. \quad (5.5)$$

Here R_E is just the non-angular form of θ_E —it is called the Einstein radius. The image will consist of a ring of angular radius θ_E , called the Einstein ring.

¹ On galactic scales D_L, D_S, D_{LS} are ordinary distances, but on cosmological scales they must be understood as angular diameter distances, and $D_S \neq D_L + D_{LS}$. The reason for this complication is that the universe will have expanded substantially over the light travel time.

² Later on, we'll use $\theta_r, \theta_x, \theta_y$ as coordinates rather than r, x, y , to remind us that these are angles on the sky, not distances.

³ The astrophysical convention being that you first think how a rational person would do it, and then you *change the sign*.

⁴ In cosmology both terms need to be multiplied by $(1 + z_L)$.

⁵ The gravitational time delay can be derived directly from general relativity, independently of (5.1), and is known as the Shapiro time delay. Radio astronomers can measure it directly.

PROBLEM 5.1: For very distant sources (i.e., $D_S \gg D_L$) we can write

$$\theta_E = (\dots) \times \sqrt{M/D_L}.$$

Find (...) in arcsec, if M is measured in solar masses and D_L in parsecs. [4]

By a Gauss's-law type argument, for any circular mass distribution $\Sigma(\theta_r)$, $\Psi(\theta_r)$ and $\alpha(\theta)$ will be influenced only by interior mass. So we'll get the same images for any circular distribution of the mass M , *provided it fits within an Einstein radius*. Bodies that fit within their own Einstein radius are said to be 'compact'. But the Einstein radius depends on where the source and observer are:

$$R_E \sim (\text{Schwarzschild radius} \times D_L)^{\frac{1}{2}}.$$

This sort of means that the further away you look, the easier it gets to see examples of gravitational lensing. It's a surprising fact at first, but it's really just the gravitational analogue of a familiar fact about glass lenses—to get the maximum effect from a lens you have to be near the focal plane, if you're too near the lens doesn't have much effect.

For given D_L, D_S , to get a compact object you have to pack a mass (in projection) into a circle of radius θ_E ; but the area of the circle is proportional to the mass. So clearly there has to be a critical density, say Σ_{crit} , such that if $\Sigma \geq \Sigma_{\text{crit}}$ somewhere then there is a compact (sub)-object. Working out the algebra we easily get

$$\Sigma_{\text{crit}} = \frac{D_L D_S}{D_{LS}} \frac{c^2}{4\pi G}. \quad (5.6)$$

Using this we can write (5.4) more concisely as

$$\begin{aligned} \nabla T &= 0, \quad T = T_0 \left[\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_S)^2 - \psi(\boldsymbol{\theta}) \right] \\ \psi(\boldsymbol{\theta}) &= \frac{1}{\pi} \int \kappa(\boldsymbol{\theta}') \ln |\boldsymbol{\theta} - \boldsymbol{\theta}'| d^2 \boldsymbol{\theta}', \end{aligned} \quad (5.7)$$

where κ is the projected mass density in units of the critical density. From the second line of 5.7 it should be evident that ψ satisfies a two-dimensional Poisson equation

$$\nabla^2 \psi = 2\kappa. \quad (5.8)$$

The fact there is a critical density, and that it depends on distances, has important astrophysical consequences. For example, a galaxy as a whole (a smooth distribution of $\sim 10^{12} M_\odot$ on a scale of $\sim 10^5$ pc) is not compact to lensing for $D_L \lesssim 10^9$ pc—cosmological distances. But clumps within the galaxy may be compact at much smaller distances. In particular, a star is compact to lensing at distances of even $\lesssim 1$ pc.

The surface $T(\boldsymbol{\theta})$ is known as the time delay surface or the arrival time surface. Wherever the arrival time is stationary (i.e., the surface as a maximum, minimum, or saddle point) there'll be constructive interference, and an image. This is Fermat's principle. Furthermore, the less the curvature of the surface at the images, the more magnified the image will be. We'll formalize this in the next section.

Try to visualize the arrival time surface. The geometrical part is a parabola with a minimum at $\boldsymbol{\theta}_S$. Having mass in the lens pushes up the surface variously. If $\kappa(\boldsymbol{\theta}) > 1$

anywhere, there will be a maximum somewhere near there, hence another image. There must be a third image too, because to have a minimum and a maximum in a surface you must have a saddle point somewhere. In fact

$$\text{maxima} + \text{minima} = \text{saddle points} + 1. \quad (5.9)$$

This is a really a statement about geometry that should be intuitively clear, though a formal proof is difficult.

A good way of gaining some intuition about the arrival time surface is to take a transparency with a blank piece of paper behind it and look at the reflections of a light bulb. Notice how images merge and split, and how you get grotesquely stretched images just as they do. Deep images of rich clusters of galaxies show just these effects!

MAGNIFICATION

By magnification we mean: how much does the image move when we move the source? It should be clear that this magnification can't be a scalar, because an image doesn't in general move in the same direction as the source. In fact the magnification is a tensor. We'll denote it by M (A for 'amplification' is also used). Formalizing our definition, we have

$$M^{-1} = \frac{\partial \boldsymbol{\theta}_s}{\partial \boldsymbol{\theta}} = \frac{\partial^2}{\partial \boldsymbol{\theta}^2} T(\boldsymbol{\theta}). \quad (5.10)$$

In cartesian coordinates

$$M^{-1} = \begin{pmatrix} 1 - \frac{\partial^2 \psi}{\partial \theta_x^2} & \frac{\partial^2 \psi}{\partial \theta_x \partial \theta_y} \\ \frac{\partial^2 \psi}{\partial \theta_y \partial \theta_x} & 1 - \frac{\partial^2 \psi}{\partial \theta_y^2} \end{pmatrix}. \quad (5.11)$$

Notice that M^{-1} is basically taking the curvature of the arrival time surface.

It is helpful to write M^{-1} in terms of its eigenvalues, and the usual form is like

$$M^{-1} = (1 - \kappa) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - \gamma \begin{pmatrix} \cos 2\phi & \sin 2\phi \\ \sin 2\phi & -\cos 2\phi \end{pmatrix}. \quad (5.12)$$

The eigenvalues are of course $1 - \kappa \pm \gamma$. The first term in (5.12) is the trace part—and comparing equations (5.11) and (5.8) shows that it must be κ —while the second term is traceless. The κ thing produces an isotropic expansion or contraction, while the γ thing produces a stretching in the ϕ direction and a shrinking in the perpendicular direction; κ is known as 'convergence' and γ as 'shear'.

The determinant of M can be thought of as a scalar magnification.

$$|M| = [(1 - \kappa)^2 + \gamma^2]^{-1}. \quad (5.13)$$

The places where one of the eigenvalues of M^{-1} becomes zero (and in consequence $|M|$ is infinite) are in general curves and are known as critical curves. When critical curves are mapped onto the source plane through the lens equation, they give caustics; a source lying on a caustic gets infinitely magnified.

EXAMPLE [Point mass and isothermal lenses] For a point mass, the lens equation is

$$\theta_{Sx} = \theta_x - \frac{\theta_x}{\theta_r^2} \theta_E^2, \quad \theta_{Sy} = \theta_y - \frac{\theta_y}{\theta_r^2} \theta_E^2,$$

and this gives

$$M^{-1} = \begin{pmatrix} 1 - \left(\frac{1}{\theta_r^2} + 2 \frac{\theta_x^2}{\theta_r^4} \right) \theta_E^2 & 2 \frac{\theta_x \theta_y}{\theta_r^4} \theta_E^2 \\ 2 \frac{\theta_x \theta_y}{\theta_r^4} \theta_E^2 & 1 - \left(\frac{1}{\theta_r^2} + 2 \frac{\theta_y^2}{\theta_r^4} \right) \theta_E^2 \end{pmatrix}.$$

Taking the determinant and simplifying, we get

$$|M|^{-1} = 1 - \frac{\theta_E^4}{\theta_r^4}. \quad (5.14)$$

For a circular mass distribution $\Sigma \propto \theta_r^{-1}$ (known as the ‘isothermal lens’, because it is just the $\rho \propto 1/r^2$ isothermal sphere in projection) the lens equation is

$$\theta_{Sx} = \theta_x - \frac{\theta_x}{\theta_r} \theta_E^2, \quad \theta_{Sy} = \theta_y - \frac{\theta_y}{\theta_r} \theta_E^2,$$

and gives

$$M^{-1} = \begin{pmatrix} 1 - \left(\frac{1}{\theta_r} + \frac{\theta_x^2}{\theta_r^3} \right) \theta_E^2 & \frac{\theta_x \theta_y}{\theta_r^3} \theta_E^2 \\ \frac{\theta_x \theta_y}{\theta_r^3} \theta_E^2 & 1 - \left(\frac{1}{\theta_r} + \frac{\theta_y^2}{\theta_r^3} \right) \theta_E^2 \end{pmatrix}.$$

And from this we get

$$|M|^{-1} = 1 - \frac{\theta_E}{\theta_r}. \quad (5.15)$$

It’s shorter in polar coordinates, but tensor components in polar coordinates can get confusing. \square

Magnification in lensing conserves surface brightness. We can prove this in a rather fun way. Let us consider the axial direction as a formal time variable t ; then light rays can be thought of as trajectories. Now allow observers to be at arbitrary transverse position (say \mathbf{w} —two dimensional) and arbitrary t . Then $\boldsymbol{\theta}$ as observed at (\mathbf{w}, t) is just the local $d\mathbf{w}/dt$ for the corresponding light ray, up to a constant factor. This means we can make a formal analogy with Hamiltonian formulation of stellar dynamics, with $\boldsymbol{\theta}$ (up to a constant) playing the role of the momentum, \mathbf{w} playing the role of the coordinates, and $\psi(\mathbf{w}, t)$ replacing the Newtonian potential. The phase space density f is the density of photons in $(\mathbf{w}, \boldsymbol{\theta})$ space, or the number of photons per unit solid angle on the sky per unit telescope area, i.e., the surface brightness. The collisionless Boltzmann equation applies (as it does for any Hamiltonian system) and it tells us that surface brightness is conserved along trajectories! Surface brightness must be conserved by the act of placing the lens there too—think of surface brightness before and after going through the lens. QED. We must be careful, though, to understand ‘along the trajectories’ correctly. It means we must always be looking at photons from the same source, so if the image is moved in the sky by lensing we must follow it when we measure surface brightness.

In other words, lensing changes the apparent sizes (and shapes) of objects, but without altering their surface brightness.

PROBLEM 5.2: For unresolved sources, we don't observe the surface brightness but only the luminosity, say L . In a survey of objects with luminosity function $f(L)$ to a limit of L_{\min} , the number of objects detected will be

$$\int_{L_{\min}}^{\infty} f(L) dL \times \langle \text{area of survey} \rangle.$$

Now suppose there is a foreground lens in the survey area with uniform scalar magnification $|M|$. This will increase the effective luminosity limit of the survey to $L_{\min}/|M|$, and hence change the number of objects detected. The changed number of objects is not, however,

$$\int_{L_{\min}/|M|}^{\infty} f(L) dL \times \langle \text{area of survey} \rangle.$$

Correct this formula.

[10]

This effect is known as 'amplification bias'.

That's more than enough theory, let's discuss real systems a little.

MULTIPLE-IMAGE QSOS

These happen when a foreground galaxy is within $\lesssim \theta_E$ (in projection) of a QSO, and produces two or four images with arcsecond order separations. Two-image systems have a minimum and a saddle point, while four-image systems have two minima and two saddle points. In both cases there's a maximum too, at the bottom of the galaxy's potential well; but since that is also generally the densest part of the galaxy, κ is very high and $|M|$ nearly vanishes, so these central images are too faint to detect.

Multiple-image QSOS are of great astrophysical interest, and two things make them so.

The first is that since QSOS are often very time-variable and the different images have different arrival times, the images will show the same time-variability, but with offsets. These offsets are simply the differences in $T(\boldsymbol{\theta})$ between different images. (So far they have been explicitly measured for two lenses.) Provided we know (or can model) $\kappa(\boldsymbol{\theta})$, the measured time offsets tell us T_0 , and hence H_0 . Basically it's this: normally we can only measure dimensionless things (image separations, relative magnifications) in lenses systems; but if we succeed in measuring a quantity that has a scale (the time delays) that tells us the scale of the universe (H_0). In practice, there is considerable uncertainty about the distribution of mass in the lensing galaxies, and this translates into an uncertainty in the inferred H_0 that is much larger than errors in the time delays. Maybe this problem can be overcome, maybe not...

The second thing has to do with the extremely small size of QSOS in optical continuum. Now the $\kappa(\boldsymbol{\theta})$ of a galaxy isn't perfectly smooth, it becomes granular on the scale of individual stars. This produces a very complicated network of critical lines (in the lens plane), and a corresponding complicated network of caustics in the source plane (like the pattern at the bottom of a swimming pool). The optical continuum emitting regions of QSOS are small enough to fit between the caustics, but the line emitting regions straddle several caustics. As proper motions move the caustic network, the continuum region will sometimes cross a caustic, and show a sudden change in brightness; the time taken for the brightness to change is the time it takes to cross the

caustic. This is the phenomenon of QSO microlensing: continuum shows it but lines don't. (It's just the gravitational version of stars twinkling and planets not twinkling.) This has been observed, and modelling the caustic network and putting in plausible values for the proper motion leads to an estimate of the intrinsic size of the continuum regions of QSOs. It's very small ~ 100 AU.

GALAXY CLUSTERS

Galaxy clusters are generally not in dynamical equilibrium (there haven't been enough crossing times since they formed). Their mass distributions and ψ potentials are thus warped in more complicated ways than for single galaxies. They are also much bigger on the sky and thus have many more background objects (faint blue galaxies) to lens.

The transparency with a paper behind it and several lightbulbs overhead is a good simulacrum of lensing by a cluster. Rich clusters show many highly stretched images of background galaxies, and these are known as arcs. A deep HST image of Abell 2218 shows over a hundred arcs, including seven multiple image systems.

An arc is close to a zero eigenvalue of M^{-1} , and is stretched along the corresponding eigenvalue. Thus each arc provides some sort of constraint on the ψ of the cluster.

Clusters also show weak lensing. That's when the eigenvalues $1 - \kappa \pm \gamma$ are too close to unity to show up as arcs, but if many galaxies in the same region are examined then statistically a stretching is measurable. The statistical stretching measures the ratio of the two eigenvalues, and thus $\gamma/(1 - \kappa)$.

Several groups have been reconstructing cluster mass profiles from information provided by multiple-images, arcs, and weak lensing.

MICROLENSING IN THE MILKY WAY

One possibility for the dark matter in the Milky Way halo is that it consists of brown dwarfs, compact objects below the hydrogen burning threshold of $0.08M_{\odot}$. Such objects would act as point lenses. A point lens has two images, at

$$\theta = \frac{1}{2} \left(\theta_S \pm \sqrt{\theta_S^2 + 4\theta_E^2} \right). \quad (5.16)$$

(There is formally a third image at $\theta = 0$, i.e., at the lens itself, but for a point mass this image has zero magnification.) The image separation for a $\sim M_0$ lens at distances of ~ 10 kpc is < 1 mas, far too small to resolve. What will be observed is a brightening equal to the combined magnification of both images. Using the result 5.14 for $|M|$ for a point lens, and adding the absolute values of $|M|$ at the two image positions, we get

$$M_{\text{tot}} = \frac{u^2 + 2}{u(u^2 + 4)^{\frac{1}{2}}}, \quad u = \frac{\theta_S}{\theta_E}. \quad (5.17)$$

Now because of stellar motions, θ_S will change by an amount θ_E over times of order a month, so microlensing in the Milky Way can be observed by monitoring light curves. If the background source star has impact parameter b and velocity v (projected onto the lens plane) with respect to the lens, then

$$u = \frac{(b^2 + v^2 t^2)^{\frac{1}{2}}}{D_L \theta_E}. \quad (5.18)$$

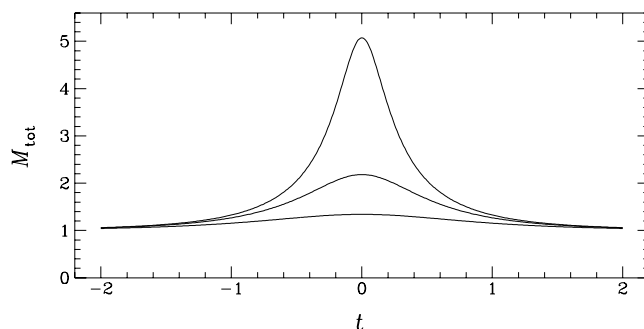


Figure 5.2: Light curves for impact parameters of R_E (lowest), $0.5R_E$ and $0.2R_E$. The unit of time is how long it takes the source to move a distance R_E .

Inserting (5.18) into (5.17) gives us $M_{\text{tot}}(0)$, i.e., the light curve, plotted for three different b in Figure 5.2. The height of a measured light curve immediately gives R_E/b , and the width gives R_E/v .

Though trying to resolve the images in microlensing seems hopeless with foreseeable technology, there are some prospects for tracking the moving double image indirectly. By combining the positions and magnifications of the two images, we have for the centroid

$$\theta_{\text{cen}} = \frac{u(3+u^2)}{2+u^2}\theta_E. \quad (5.19)$$

Such microlensing events are rare, because θ_S has to be $\lesssim \theta_E$ for significant magnification. People speak of an optical depth τ to microlensing in a field. This is the probability of a star being (in projection) within θ_E of a foreground lens, at any given time. From equation (5.18) it amounts to the probability of $M_{\text{tot}} \geq 2/\sqrt{5} = 1.34$. It's just the covering factor of discs of radius θ_E (Einstein rings) from all lenses between us and the stars in the field.⁶ The source stars might be bright stars in the Large Magellanic Cloud (LMC) and the lenses very faint stars or brown dwarfs in the Milky Way halo.

Using equation (5.5) for R_E and considering the total area covered by the Einstein rings of lenses at distances between D_L and $D_L + dD_L$ in a patch of sky, and then integrating over D_L , we have

$$\tau = \frac{4\pi G}{c^2 D_S} \int_0^{D_S} D_L D_{LS} \rho(D_L) dD_L. \quad (5.20)$$

PROBLEM 5.3: Derive the formula (5.20) for the microlensing optical depth. [10]

Imagine an observer at radius $r = 1$ in an isothermal sphere made of machos, looking outwards (i.e., towards the anti-centre) at sources at radius $r = a$, and monitoring for microlensing. Show that τ for this observer will be

$$\tau = 2 \frac{\sigma^2}{c^2} \left[\frac{a+1}{a-1} \ln a - 2 \right]. \quad [6]$$

$$\left[\int_1^a x^{-2}(x-1)(a-x) dx = (1+a) \ln a - 2(a-1). \right]$$

⁶ So optical depth is a bit of a misnomer.

The really nice thing about the formula (5.20) is that it doesn't depend on the mass distribution of the lenses, as long as each mass fits within its own Einstein radius (diffuse gas clouds don't count, nor does any kind of diffuse dark matter). So τ estimated from light curve monitoring could be used to make inferences about ρ .

How large is τ through the Galactic halo? To estimate that, we need an estimate for ρ . Now the Milky Way rotation curve suggests an isothermal halo, $\rho = \sigma^2/(2\pi Gr^2)$, with $\sigma \sim 200$ km/sec. If we then say that r will be of order the D factors in (5.20), we get

$$\tau \sim \frac{\sigma^2}{c^2}, \quad \text{or} \quad \tau \sim 10^{-7} \text{ to } 10^{-6}. \quad (5.21)$$

So to have any hope of detecting such microlensing events, it is necessary to monitor the light curves of millions of stars. Four such surveys have been started up in the last two years, observing fields in the LMC and/or the Milky Way bulge. (The bulge surveys don't go through the halo of course, but through part of the Milky Way disc.)⁷ The current estimates for τ are $\sim 10^{-7}$ towards the LMC and $\simeq 3 \times 10^{-6}$ towards the bulge. How much of the lensing mass is in brown dwarfs as distinct from faint stars, and whether the lensing mass alone can account for rotation curve data are not yet clear. Meanwhile, the huge number of variable stars discovered by these surveys are revolutionizing that field of study.

⁷ An estimate of τ from a survey will include a correction for the detection efficiency. Surveys have to be very wary of spurious detections; hence any light curve possibly contaminated by stellar variability has to be discarded for microlensing purposes. Detection efficiencies are of order 30%.

6. The Milky Way

The Milky Way is, as far as we know, a typical disc galaxy. Figure 6.1 is a cartoon to remind you of the different components of the Milky Way. The luminous parts are mostly a disc of Pop I stars and a bulge of older Pop II stars. We live in the disc, about 8.5 kpc from the centre. Apart from stars, the disc also has clusters of young stars and H II regions, and dust and gas; the gas is mostly observed as an H I layer which flares at large radii. There is some evidence that there are two spiral arms in the disc (the dust makes it hard to tell). The bulge has a bar, though the dimensions of it are unclear. But the most massive part is the halo; there are some old stars (and globular clusters of very old stars) in the halo, but most of the halo is dark matter of unknown composition.

That is not all: there are also the small companion galaxies. The best known of these are the the Large and Small Magellanic Clouds (LMC and SMC) which are $\simeq 50$ kpc away; these may or may not be part of a trail of debris known as the Magellanic Stream. Then there is the Sagittarius Dwarf galaxy which seems be to being sucked into the Milky Way halo now.

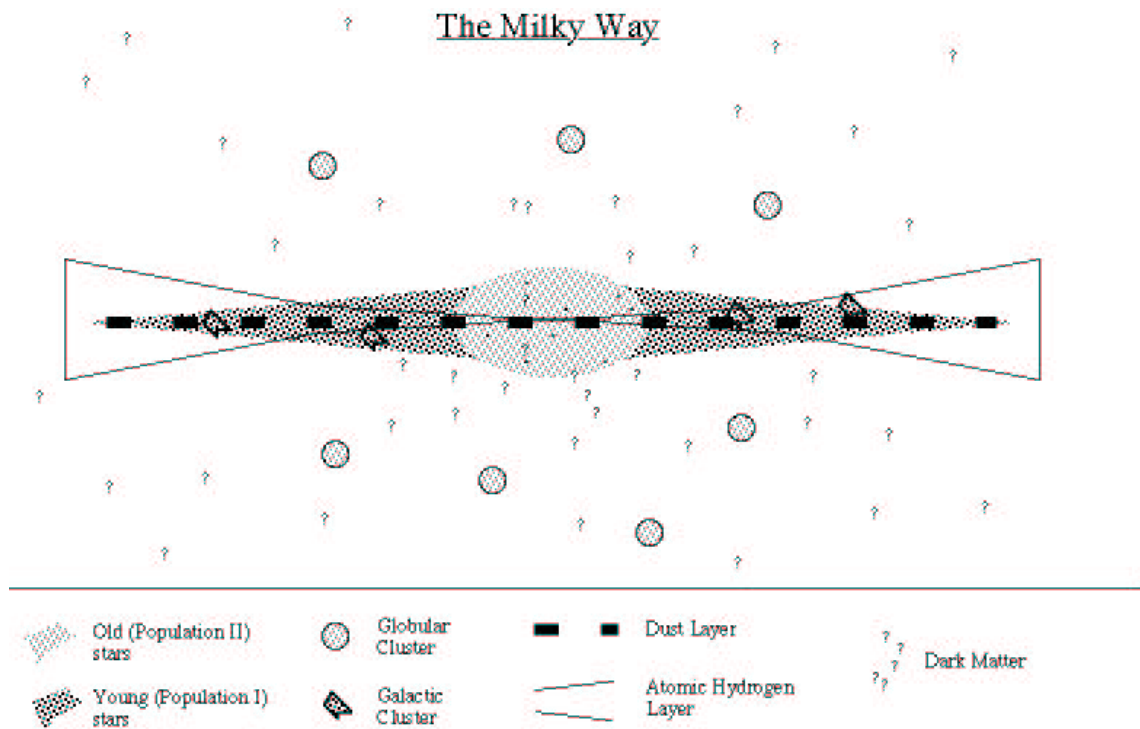


Figure 6.1: Cartoon of the Milky Way (by Mike Merrifield).

THE MASS OF THE MILKY WAY

While there are good estimates of enclosed mass of the Milky Way within different radii, it is not known where the halo of the Milky Way finally fades out (or even if the size of the halo is a very meaningful concept). So the only way to get at the *total* mass of the Milky Way is to observe its effect on other galaxies. The simplest but most robust of these comes from an analysis of the mutual dynamics of the Milky Way and M31 (Andromeda); it is known as the timing argument.

The observational inputs are (i) M31 is $\simeq 750$ kpc away, and (ii) the Milky Way and M31 are approaching at $\simeq 120$ km/sec. A simple approximation for their dynamics is to suppose that they started out at the same point with initial recessional velocities from the Big Bang, and have since turned around because of mutual gravity. This is not strictly true of course, because galaxies had not already formed at the Big Bang; however it is thought that galaxies (at least galaxies like these) formed early in the history of the universe, so the approximation may be acceptable. Writing l for the distance and M for the combined mass in the Milky Way and in M31, the equation of motion for the reduced Keplerian one-body problem is

$$\frac{d^2l}{dt^2} = -\frac{GM}{l^2}. \quad (6.1)$$

In considering a Keplerian problem without perturbation we are, of course, assuming that the gravity from Local Group dwarfs and the cosmological tidal field is negligible; but as there are no other large galaxies within a few Mpc this seems a fair approximation. It is not obvious how to solve this nonlinear equation, but fortunately the solution is known and easy to verify; it is most conveniently expressed in parametric form, as

$$\begin{aligned} t &= \tau_0(\eta - \sin \eta), \\ l &= (GM\tau_0^2)^{\frac{1}{3}}(1 - \cos \eta). \end{aligned} \quad (6.2)$$

Here τ_0 is an integration constant, the other integration constant has been eliminated by the boundary condition $l(t=0) = 0$. Now consider the dimensionless quantity

$$\left(\frac{t_0}{l_0}\right) \left(\frac{dl}{dt}\right)_{t_0} = \frac{\sin \eta_0(\eta_0 - \sin \eta_0)}{(1 - \cos \eta_0)^2}, \quad (6.3)$$

where the subscripts in t_0 and so on refer to the current time, as conventional in cosmology. Inserting the observed values for l_0 and $(dl/dt)_{t_0}$ and a plausible value of 15 Gyr for t_0 (the age of the universe), we get -2.4 for the left hand side. The solution for η_0 to give the same value on the right hand side is 4.3. Inserting these values in (6.2) we get $\tau_0 = 2.9$ Gyr and¹

$$M \simeq 4 \times 10^{12} M_{\odot}. \quad (6.4)$$

From its luminosity and rotation curve, M31 appears to have of order twice the mass of the Milky Way. This implies that the mass of Milky Way exceeds $10^{12} M_{\odot}$. Estimates for the mass of the luminous part of the Milky Way range from 0.05 to $0.12 \times 10^{12} M_{\odot}$.

¹ It is useful to remember G in useful astrophysical units as $4.98 \times 10^{-15} M_{\odot}^{-1} \text{pc}^3 \text{yr}^{-2}$.

PROBLEM 6.1: The mass estimate from the timing argument involves solving a difficult differential equation and then an algebraic equation. But one can do a back-of-the-envelope version of the calculation using just dimensional analysis.

Show that the inputs (a) the Universe is ~ 10 Gyr old and the Milky Way and M31 formed early, (b) M31 is turning round about now, (c) M31 is ~ 1 Mpc away, and (d) $G = 4.98 \times 10^{-15} M_{\odot}^{-1} \text{pc}^3 \text{yr}^{-2}$ imply that the combined mass of the Milky Way and M31 is $M \simeq 2 \times 10^{12} M_{\odot}$. [10]

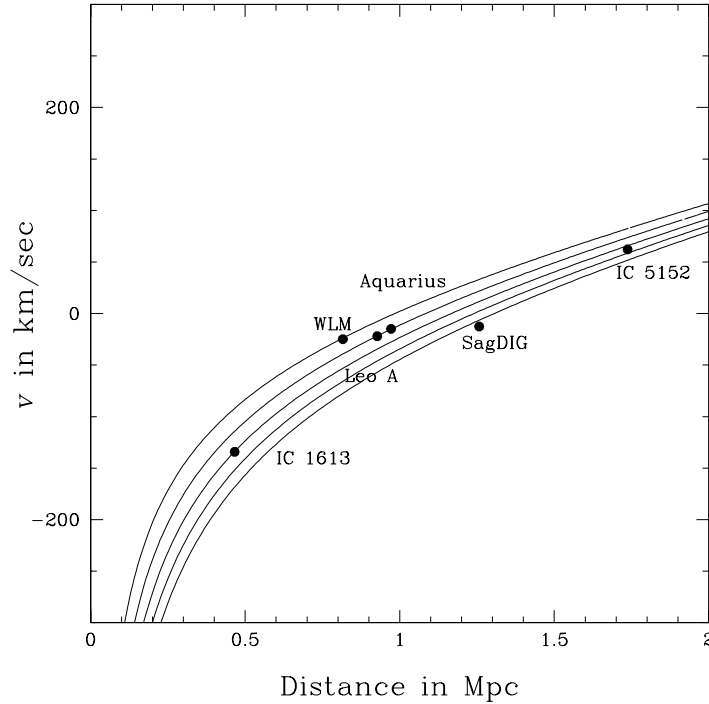


Figure 6.2: Distances and velocities of six Local Group dwarf galaxies, and predictions for different values of GM/τ_0 (by Alan Whiting).

The timing argument can be applied not only to Andromeda, but also to Local Group dwarf galaxies (which have much less mass and behave just as tracers). Figure 6.2 shows plots l against dl/dt for some Local Group dwarfs, along with the predictions of the timing argument for different values of GM/τ_0 .

$$l = (GM\tau_0^2)^{\frac{1}{3}} (1 - \cos \eta). \tag{6.5}$$

$$\frac{dl}{dt} = \left(\frac{GM}{\tau_0}\right)^{\frac{1}{3}} \frac{\sin \eta}{1 - \cos \eta}$$

PROBLEM 6.2: If we replace sin and cos in (6.5) with sinh and cosh, the result still satisfies the differential equation (6.1). Verify this, and explain how it relevant to Figure 6.2. [25]

THE SOLAR NEIGHBOURHOOD

The Milky Way is a differentially rotating system. The local standard of rest (LSR) is a system located at the sun and moving with the local circular velocity (which is $\simeq 200$ km/sec). The sun has its own peculiar motion of $\simeq 13$ km/sec with respect to the LSR.

The rotation velocity and its derivative at the solar position are traditionally expressed in terms of Oort's constants:

$$\begin{aligned} A &= \frac{1}{2} \left(\frac{v_\phi}{R} - \frac{\partial v_\phi}{\partial R} \right), \\ B &= -\frac{1}{2} \left(\frac{v_\phi}{R} + \frac{\partial v_\phi}{\partial R} \right). \end{aligned} \tag{6.6}$$

The reason is that A vanishes for solid body rotation, and can be measured from line of sight velocity data without proper motions (which in the past were hard to measure). But now that we have accurate proper motions from Hipparcos, and hence (combining with ground-based line-of-sight velocities) three-dimensional stellar velocities in the solar neighbourhood, A and B are less important.

If you take the average (three-dimensional) velocity and dispersions of any class of stars in the solar neighbourhood, then $\langle v_R \rangle$ and $\langle v_z \rangle$ turn out to be nearly zero, while $\langle v_\phi \rangle$ is such that $\langle v_\phi \rangle - v_{\text{LSR}}$ is negative and $\propto \sigma_{RR}$. This is known as the 'asymmetric drift' and is nothing but our old friend rotational support versus pressure support. Young stars are almost entirely supported by $\langle v_\phi \rangle$, like the gas that produced them. Older stars pick up increasing amounts of pressure support in the form of σ_{RR} ; they then need less v_ϕ to support them, and thus tend to lag behind the LSR. The linear relation can be derived from the Jeans equations, but we won't go through that because you've probably had enough of Jeans equations for now...

When examined in detail using Hipparcos proper motions, the velocity structure in the solar neighbourhood is more complicated than anyone expected. Figure 6.3 shows a reconstruction of the stellar (u, v) (i.e., radial and tangential velocity) distribution in the solar neighbourhood for stars in different ranges of the main sequence.² Notice the clumps in the velocity distribution which appear for stars of all ages. (And these are clumps only in velocity space, not in real space.) The idea that there are groups of stars at similar velocities is itself not new—it actually dates from the early proper motion measurements of nearly a century ago. But these 'streams' have generally been interpreted as groups of stars which formed in the same complex and were later stretched in real space over several galactic orbits. The surprising new finding is that the 'streams' are seen for stars of all ages, which indicates a dynamical origin; they seem to be wanting to tell us something interesting about Milky Way dynamics, but as yet we don't know what.

² The Schwarzschild ellipsoid and its vertex deviation that you may find in textbooks should now be considered obsolete—they are essentially the result of washing out the structure in Figure 6.3.

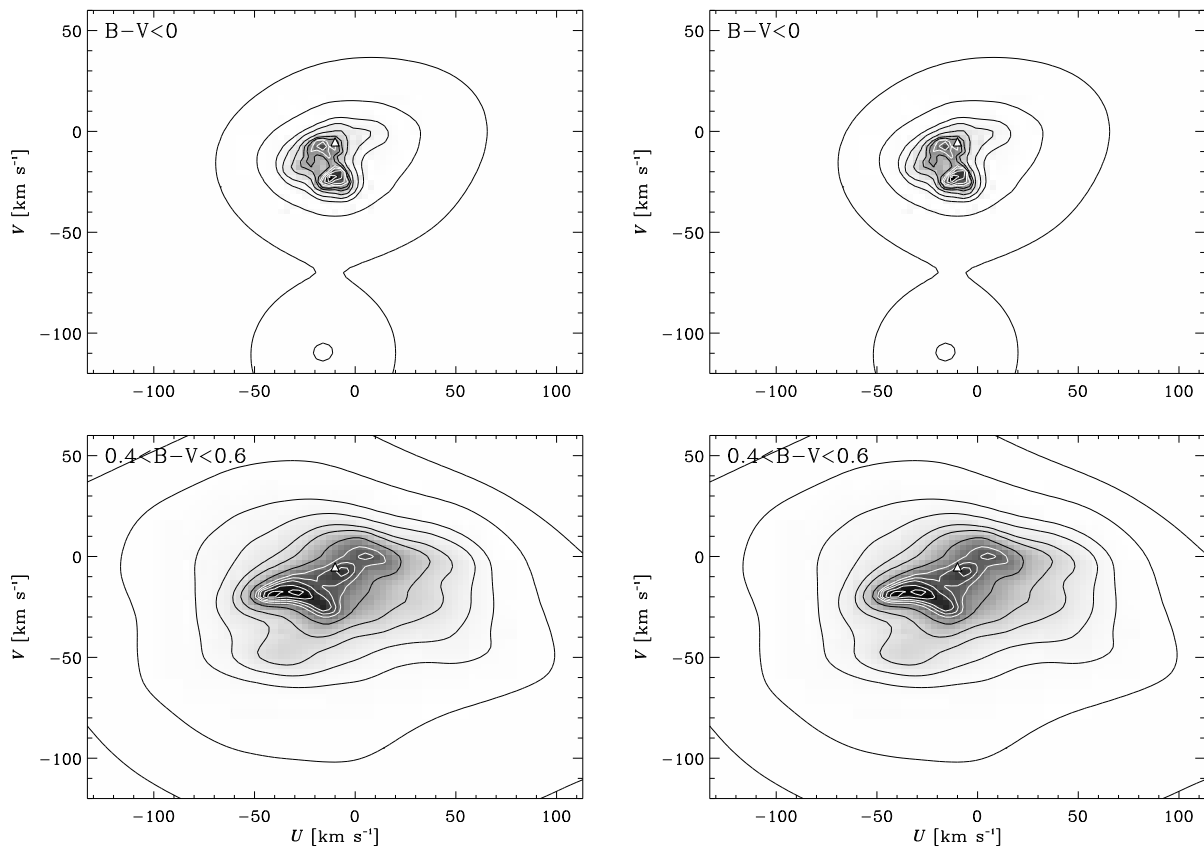


Figure 6.3: Distribution of radial (u) and tangential (v) velocities of main sequence stars in the solar neighbourhood, recently reconstructed from Hipparcos proper motions by Walter Dehnen (1998). The upper left panel is for the youngest (and bluest) stars; these are estimated to be < 0.4 Gyr old. The upper right panel is for stars younger than 2 Gyr, and the lower left panel is for stars younger than 8 Gyr. The lower right panel shows the combined distribution for all main sequence stars. The sun is at $(0, 0)$ and the LSR is marked by a triangle.

THE BAR

There is little doubt now that the Milky Way bulge is triaxial—there is a (rotating) bar with the positive l side nearer to us and moving away. The evidence for this was at first indirect, as the following. Consider gas in the ring, which must move on closed orbits. If it moved on circular orbits in the disc, and we measured its Galactic longitude l and line of sight velocity v , then all the gas at positive l would have one sign for v and similarly all the gas at negative l would have the opposite sign for v . In fact gas at positive l is seen with both signs for v , and likewise at negative l . So the gas orbits must be non-circular, and hence the gravitational potential must be non-circular in the disc. This suggests a bar and indeed the observed gas kinematics is well fitted by a bar.

The features of a bar can in fact be seen in an infrared map of the bulge, if you know what to look for. Figure 6.4 shows a bar in the plane, and its effect on an l, b map.

- (i) The side nearer to us is brighter. Contours of constant surface brightness are further apart in both l and b on the nearer size.
- (ii) Very near the centre, the *further* side appears brighter, so the brightest spot is slightly to the further size of $l = 0$. The reason is that on the further side our

line of sight passes through a greater depth of bar material, which more than compensates for it being slightly further.

The features (i) can be discerned in many different data sets; the feature (ii) is harder to find, it just about shows up in the COBE maps of the bulge.

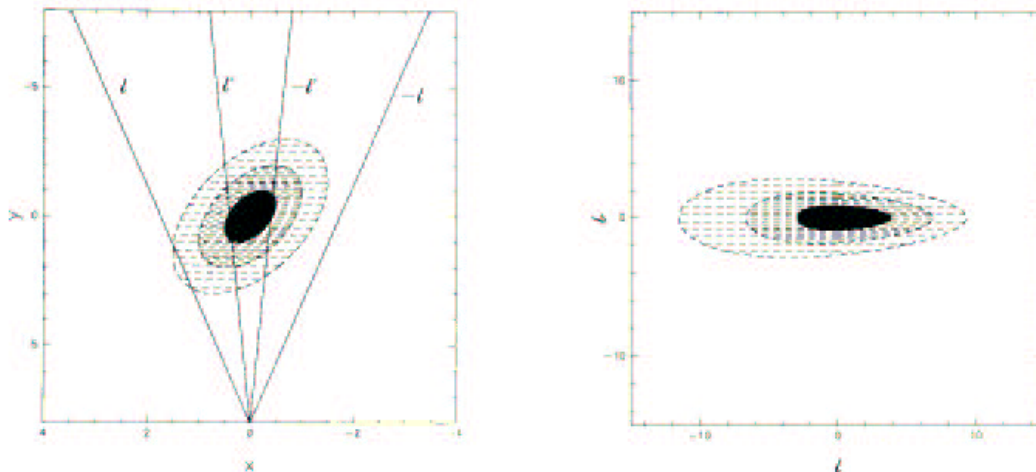


Figure 6.4: Schematic of the bar in the Milky Way Bulge, viewed from the North Galactic pole (left), and from the Sun (right). (From Blitz and Spergel, *ApJ*, 1991. The right panel uses minus the usual convention for l .)

THE SAGITTARIUS DWARF

We'll end our discussion of galaxies with the Sagittarius Dwarf. It may seem amazing that this fairly substantial companion galaxy of the Milky Way remained undiscovered till 1993; the reason is that it's behind the bulge, and thus has the densest part of the Milky Way in the foreground as camouflage. We don't know yet how large the Sagittarius Dwarf is, because it can't be spotted against the foreground in an image. A lower limit on its size comes indirectly from microlensing surveys, because they detect RR Lyraes in their fields. Figure 6.5 shows its rough extent.

The Sagittarius Dwarf is presumably being tidally stretched as it falls into the Milky Way halo; that would explain its being long and thin. Has the Milky Way eaten many such galaxies in the past?

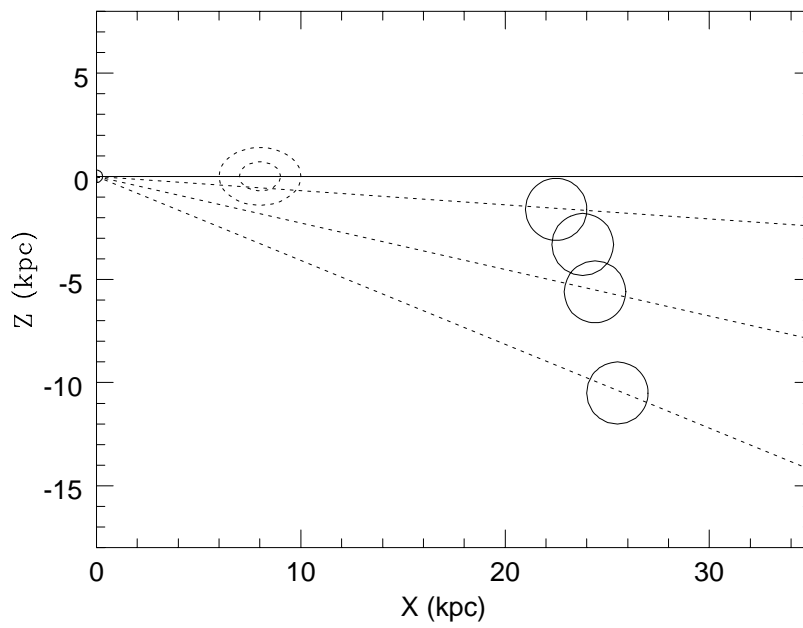


Figure 6.5: A partial map of the Sagittarius dwarf galaxy, from RR Lyraes. We are at $(0,0)$, the ellipses around $(8.5,0)$ represent the bulge, and the four circles indicate the four microlensing survey fields where the RR Lyraes were found. (From Minniti et al. 1997.)