

Michaelmas term, 2022

C7.5: General relativity 1

University of Oxford: Part C

Chris Couzens

Office S1.46

Mathematical Institute, University of Oxford

Andrew Wiles Building, Radcliffe Observatory Quarter
Woodstock Road
Oxford, OX2 6GG

Disclaimer: There are almost certainly typos in the notes, if something does not look correct or needs further explanation please let me know.

Please send any comments/corrections to Christopher.Couzens@maths.ox.ac.uk

Recommended books and resources

There are a large variety of good textbooks and lecture notes on general relativity. This course borrows from a number of them, in various different places. An assortment of textbooks that have been used in writing these notes are:

- Wald, General Relativity

A very thorough introduction to the subject.

- Weinberg, Gravitation and cosmology
- Carroll, An introduction to general relativity, spacetime and geometry.

Aimed more at particle physicists. We will follow this in the cosmology section and borrow bits for elsewhere.

- Hartle, Gravity, an introduction to Einstein's general relativity
- Misner, Thorne and Wheeler, Gravitation

It is a very big book.

- Nakahara, Geometry, Topology and Physics

An excellent book for learning about geometry and topology and will be useful for the differential geometry section of the notes.

There are also a number of useful lecture notes online. In particular:

- Joe Kier's lecture notes from 2020
- David Tong's lecture notes
- Sean Carroll's lecture notes
- Harvey Reall's lecture notes

Contents

1	Introduction	7
2	Newtonian gravity and Special relativity	8
2.1	Special relativity	8
2.2	Newtonian gravity	18
2.3	Equivalence Principles	20
2.4	Some worked examples	21
2.5	Problems with Newtonian gravity and why we need GR	28
3	Differential Geometry	29
3.1	Manifolds	30
3.2	Calculus on manifolds	38
3.2.1	Tangent Vectors	39
3.2.2	One-forms	40
3.2.3	Tensors	41
3.2.4	Tensor fields	41
3.2.5	Induced maps	42
3.3	Flows and Lie derivatives	43
3.3.1	One-parameter group of transformations	44
3.3.2	Lie Derivatives	45
3.4	Differential forms	49
3.4.1	Exterior product	51
3.4.2	Exterior derivative	52
3.4.3	Interior product	55
3.4.4	Integration	57
4	Riemannian geometry	61
4.1	The metric	61
4.1.1	Riemannian metric	62
4.1.2	Lorentzian manifolds	62
4.1.3	Why is the metric useful	63
4.2	Connections and curvature	66
4.3	Torsion and curvature	69

4.3.1	Levi–Civita connection	70
4.4	Parallel transport and geodesics	72
4.4.1	Geodesics	73
4.4.2	Normal coordinates	74
4.4.3	Path dependence: Curvature and Torsion	76
4.4.4	Geodesic deviation	79
4.5	Riemann tensor and its symmetries	81
4.5.1	Ricci and Einstein tensors	82
5	Einstein’s equations	82
5.1	The Einstein–Hilbert action	83
5.1.1	Newton’s constant	86
5.1.2	Cosmological constant	86
5.1.3	Higher derivative terms	87
5.1.4	Diffeomorphisms	87
5.1.5	Coupling to matter	89
5.2	Newtonian gravity as a limit	91
6	Schwarzschild solution	95
6.1	The Schwarzschild black hole	95
6.1.1	Birkhoff’s theorem	97
6.1.2	Geodesics	101
6.2	Schwarzschild solution as a black hole	112
6.2.1	Kruskal spacetime	118
7	Cosmology	121
7.1	FRW metric	121
7.1.1	Cosmological red-shift	124
7.2	The Friedmann equations	125
7.2.1	Equation of state	126
7.2.2	Deriving the Friedmann equations	127
7.3	Cosmological solutions	128
7.3.1	Solutions with $k = 0$	129
7.3.2	Solutions with $\Lambda = 0$	131
7.3.3	The Big Bang	131

Conventions

- We will use the god-given signature convention of mostly plus $(-, +, +, +)$. This may differ with the convention you have used in other courses, especially field theory courses. This convention is preferable when thinking about geometry as it gives positive spatial distances. For quantum field theory the other convention is preferable since it ensures that energies and frequencies are positive. You may map between the two conventions through *Wick rotation*, essentially allowing the coordinates to become complex.
- Spacetime indices will be taken to be greek letters from the middle of the alphabet: μ, ν, ρ, \dots and run over $0, 1, 2, 3$. Latin indices i, j, k, \dots run over the spatial directions and take values $1, 2, 3$.
- We employ Einstein summation convention, repeated indices are summed over, unless otherwise stated.
- We work in units where the speed of light c is set to 1. Occasionally it is instructive to reintroduce c which can be done by dimensional analysis.
- The Minkowski metric will be denoted by $\eta_{\mu\nu} = \text{diagonal}(-1, 1, 1, 1)_{\mu\nu}$.
- After introducing curvature we will take the metric to be $g_{\mu\nu}$ and the determinant will be $\det(g_{\mu\nu}) \equiv g$.

Useful formulae

- The Lagrangian for the geodesic equation of a massive test particle is

$$\mathcal{L}\left(\frac{dx^\mu}{d\lambda}, x^\mu\right) = \sqrt{-g_{\mu\nu}(x)\frac{dx^\mu}{d\lambda}\frac{dx^\nu}{d\lambda}},$$

with λ an arbitrary parameter along the worldline.

- The geodesic equation for a massive particle is

$$\frac{d^2x^\mu}{d\tau^2} + \Gamma^\mu_{\nu\rho}\frac{dx^\nu}{d\tau}\frac{dx^\rho}{d\tau} = 0, \quad g_{\mu\nu}(x)\frac{dx^\nu}{d\tau}\frac{dx^\rho}{d\tau} = -1,$$

where τ is the proper time. For light, the first equation takes the same form just replacing τ with an affine parameter. The second is modified by $-1 \rightarrow 0$.

- The Christoffel symbols (Levi–Civita connection) are

$$\Gamma^\mu_{\nu\rho} = \frac{1}{2}g^{\mu\sigma}\left(\partial_\nu g_{\sigma\rho} + \partial_\rho g_{\sigma\nu} - \partial_\sigma g_{\nu\rho}\right).$$

- The Riemann tensor is

$$R^\mu_{\nu\rho\sigma} = \partial_\rho \Gamma^\mu_{\nu\sigma} - \partial_\sigma \Gamma^\mu_{\nu\rho} + \Gamma^\mu_{\rho\lambda} \Gamma^\lambda_{\nu\sigma} - \Gamma^\mu_{\sigma\lambda} \Gamma^\lambda_{\nu\rho}.$$

- Symmetries

$$R_{\mu\nu\rho\sigma} = -R_{\mu\nu\sigma\rho},$$

$$R_{\mu\nu\rho\sigma} = R_{\sigma\rho\mu\nu}.$$

- Bianchi identity 1

$$R^\mu_{\nu\rho\sigma} + R^\mu_{\rho\sigma\nu} + R^\mu_{\sigma\nu\rho} = 0.$$

- Bianchi Identity 2

$$\nabla_\mu R^\sigma_{\lambda\nu\rho} + \nabla_\nu R^\sigma_{\lambda\rho\mu} + \nabla_\rho R^\sigma_{\lambda\mu\nu} = 0.$$

- Ricci tensor

$$R_{\mu\nu} = R^\rho_{\mu\rho\nu}$$

- Ricci scalar

$$R = R_{\mu\nu} g^{\mu\nu}.$$

- Einstein tensor

$$G^{\mu\nu} = R^{\mu\nu} - \frac{1}{2}Rg^{\mu\nu}.$$

- Einstein–Hilbert action plus cosmological constant,

$$S = \frac{1}{16\pi G} \int d^4x \sqrt{-g} (R + \Lambda).$$

- Under a variation $g_{\mu\nu} \rightarrow g_{\mu\nu} + \delta g_{\mu\nu}$ we have

$$\delta g^{\mu\nu} = -g^{\mu\rho} g^{\nu\sigma} \delta g_{\rho\sigma},$$

$$\delta g = gg^{\mu\nu} \delta g_{\mu\nu},$$

$$\delta R_{\mu\nu} = \nabla_\rho \delta \Gamma^\rho_{\mu\nu} - \nabla_\mu \delta \Gamma^\rho_{\rho\nu}.$$

1 Introduction

Gravity is one of the four¹ fundamental forces alongside electromagnetism, the strong nuclear force and the weak nuclear force. Of these forces gravity is by far the weakest force, the ratio of the gravitational force to electric force acting on an electron is 10^{-36} .² Despite this gravity plays a dominant role in shaping the large scale structure of the universe, this is because the strong and weak forces have a very short range, while, though electromagnetism is a long range force it is both attractive and repulsive and for bodies of macroscopic dimensions the repulsion of like charges is approximately balanced by the attraction of oppositely charged. On the other hand gravity is only an attractive force, thus for sufficiently large bodies the gravitational field of the sum of all its constituents adds up to become the dominant force.

The leading candidate for a theory of gravity for some time was Newton's theory of gravitation. This however is a non-relativistic theory of gravity and therefore is incompatible with special relativity: it is not invariant under Lorentz transformations. One can see this by thinking about what would happen if the sun suddenly disappeared. For 8 minutes, the time it takes for light to travel from the sun to Earth, we would be completely oblivious. This is because special relativity tells us that no signal can travel faster than light: the Earth must continue on its orbit for these 8 minutes, after which, it is flung out of the solar system leading to almost certain death for all life on Earth. However, Newton's theory of gravity acts instantaneously, we would be flung out of the solar system immediately. In Newton's theory, the force on one mass depends on the location of the other mass at the same time.

Einstein's breakthrough lead to a conceptual revolution in the way that we view spacetime. The fact that objects with the same initial conditions travel along the same curve, independent of their mass, hints that the curve that is followed is a property of the geometry of spacetime rather than a force acting on the body. General relativity (GR) understands gravity as the curvature of spacetime and the trajectories within spacetime as geodesics on this curved space. Or as John Wheeler once said, "*Mass tells space how to curve, while curved space tells matter how to move*".

The aim of this course is to introduce you to General relativity and by the end of it to allow you to perform calculations. Among other topics we will see how gravity bends light, the corrections to the motion of the planets and some cosmology. This is a large topic and

¹One should probably add *currently known to physics* at this point.

²You can see this very clearly by holding two magnets together, gravity is not strong enough to pull one magnet to the floor.

we will therefore omit many interesting directions, but this will lay the foundation for further study and for the follow up course general relativity II.

The notes are organised as follows. We begin by reviewing special relativity and Newtonian gravity in section 2. To understand general relativity properly we need to understand the underlying geometry of spacetime. This requires knowledge of the sophisticated tools of *differential geometry* to describe curved spacetime. With these new tools we are finally in a position to introduce Einstein's equations and physics in curved spacetime. The Schwarzschild solution is the go to solution of general relativity and we will use it as a testing ground for studying many interesting topics in GR including black holes, the motion of the planets and the bending of light. We will also see what GR has to say about the large scale structure of the universe with a trek through the world of cosmology.

2 Newtonian gravity and Special relativity

2.1 Special relativity

We begin with a whirlwind exploration of special relativity. This section is by no means meant to be an introduction to special relativity, more a refresher on the subject and to emphasise the pertinent points. Excellent texts for a more detailed treatment are []. [Add refs](#)

By the end of the 18th century two areas of physics that were in conflict had emerged: Newtonian mechanics and Electromagnetism. Newtonian mechanics has a notion of absolute time with the equations of motion are invariant under Galilean coordinate transformations. The transformation law between two inertial frames moving at a uniform speed v in the x direction is

$$(t', x', y', z') = (t, x - vt, y, z). \quad (2.1)$$

Galilean transformations imply that the speed of light should changes in different inertial frames moving with respect to each other. This is incompatible with Maxwell's equations describing electromagnetism where the speed of light is fixed. A resolution to this problem was proposed by conjecturing a preferred frame, the frame of the physical medium in which light propagates, called the *Ether*. The speed of light in any other rest frame would then be modified by the Newtonian addition of velocities. An experiment by Michelson and Morley in 1887 to detect the Ether failed, the Newtonian law of addition of velocities was not correct, either Newtonian mechanics or Maxwell's equations required modification.

Einstein gave the resolution to this problem in 1905 with his theory of special relativity. The principle of special relativity states that the laws of nature are invariant under Lorentz

transformations, a group of spacetime coordinate transformations. In particular the speed of light is the same in any reference frame and requires an abandonment of the notion of absolute time. Events which are simultaneous in one reference frame need not be in another frame (see problem sheet 1 for an example to work through).

Lorentz transformations A Lorentz transformation is a transformation from one spacetime coordinate system $x^\mu = (ct, x, y, z)$ to another x'^μ ,

$$x'^\mu = \Lambda^\mu_\nu x^\nu, \quad (2.2)$$

where Λ is a constant matrix which satisfies

$$\Lambda^\mu_\rho \Lambda^\nu_\sigma \eta_{\mu\nu} = \eta_{\rho\sigma}. \quad (2.3)$$

The matrix η is the famed Minkowski metric

$$\eta_{\mu\nu} = \text{diag}(-1, 1, 1, 1)_{\mu\nu} \equiv \begin{pmatrix} -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}_{\mu\nu}. \quad (2.4)$$

The Lorentz group is denoted $O(1, 3)$, the numbers signify the signature of the space. We could also add in constant shifts of the coordinates, $x'^\mu = \Lambda^\mu_\nu x^\nu + a^\mu$, with a^μ a constant four-vector. This would enhance the Lorentz group to the Poincaré group. Here we will focus only on the Lorentz group. The fundamental property that distinguishes the Lorentz group is that it leaves the line element (sometimes also called length element, invariant interval)

$$ds^2 = \eta_{\mu\nu} dx^\mu dx^\nu, \quad (2.5)$$

invariant.³ Here, d stands for a small displacement, you also see δ and Δ to mean the same thing.

Aside: The group described above is sometimes called the *homogeneous Lorentz group*. It admits a proper subgroup defined by imposing

$$\Lambda^0_0 \geq 1, \quad \det \Lambda = 1. \quad (2.6)$$

The proper subgroup restricts to all transformations which can be smoothly joined to the identity. The *improper* Lorentz transformations involve either space inversion $\det \Lambda = -1$, $\Lambda^0_0 \geq 1$, or time reversal $\det \Lambda = 1$, $\Lambda^0_0 \leq 1$. Space and time inversions are known not to be exact symmetries of nature and therefore when we say Lorentz transformation what we really mean is the proper Lorentz transformations.

³One can show that the Lorentz transformations are the only non-singular coordinate transformations that leave ds^2 invariant. Here non-singular means that both $x'(x)$ and $x(x')$ are well behaved differential functions and thus $\frac{\partial x^\mu}{\partial x'^\nu}$ has an inverse. When we consider $ds^2 = 0$ there is an enhancement of the symmetry group.

The proper Lorentz transformations have a further subgroup consisting of rotations taking the form:

$$\Lambda^0_0 = 1, \quad \Lambda^0_i = \Lambda^i_0 = 0, \quad \Lambda^i_j = R_{ij}, \quad (2.7)$$

with R an $\text{SO}(3)$ matrix: $RR^T = 1, \det R = 1$. The remaining transformations are known as *boosts* which mix the space and time directions. You may think of the boosts as rotations between space and time. Examples of the two types of transformation are⁴

$$\Lambda^{\text{Rotation}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos \theta & \sin \theta & 0 \\ 0 & -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad \Lambda^{\text{Boost}} = \begin{pmatrix} \cosh \phi & -\sinh \phi & 0 & 0 \\ -\sinh \phi & \cosh \phi & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}. \quad (2.8)$$

The first is a rotation in the x, y directions and the second is a boost in the x direction. The rotation parameter is compact $\theta \in [0, 2\pi)$ while the boost parameter, known as the *rapidity* is non-compact $\phi \in (-\infty, \infty)$. Altogether the Lorentz group has six parameters, split evenly between boost and rotations. Rotations commute amongst themselves but do not commute with boosts, thus it the Lorentz group is non-abelian. *Exercise: Compute the addition of the rapidity under two successive boosts along the x axis.*

The interpretation of the rotations is clear from our understanding of Galilean symmetries but what is the interpretation of the boosts? This corresponds to changing coordinates to that of a moving frame which travels at a constant velocity. The transformed coordinates under such a boost are

$$t' = t \cosh \phi - x \sinh \phi, \quad x' = -t \sinh \phi + x \cosh \phi. \quad (2.9)$$

The point $x' = 0$ is then moving, as viewed from the original frame, with velocity

$$v = \frac{x}{t} = \tanh \phi. \quad (2.10)$$

Motivated by this it is useful to replace $\phi = \text{arctanh } v$ in the transformations to obtain

$$\begin{aligned} t' &= \gamma(t - vx), & \text{with } \gamma = (1 - v^2)^{-1/2}. \\ x' &= \gamma(x - vt), \end{aligned} \quad (2.11)$$

Applying these transformations leads to time dilation, length contraction, and other phenomena. In problem sheet 0 you will review some of these problems.

A useful way of visualising the causal structure of spacetime is the spacetime diagram. We begin by portraying the original t, x directions as axes, suppressing the y and z directions.

⁴Note that we have implicitly taken the proper Lorentz group.

Under a boost, (2.9), the x' axis is given by $t = x \tanh \phi$ and the t' axis is given by $x = t \tanh \phi$. The boost rotates the space and time axes into each other, with the angle between them seemingly closing. This is a Euclidean view-point however, the axis remain orthogonal in the Lorentzian sense. The paths corresponding to the motion of light in the diagram are the $x = \pm t$ lines. The paths defined by $t' = \pm x'$ are precisely the same as the $x = \pm t$ lines (*Exercise: check that this is correct*). A set of points which are all connected to a single event by straight lines moving at the speed of light is called a *light cone*, and is invariant under Lorentz transformations. Light cones are divided into the future and the past. The set of all points inside the future and past light cones of a point p are called *time-like separated* from p . Those outside of the light cones are *space-like separated* while those lying on the cone are called *lightlike* or *null separated*. The interval between time-like separated points is negative, for space-like separated it is positive and for light-like/null it is vanishing.

To probe the structure of Minkowski space it is necessary to introduce the concepts of vectors and tensors. We will give a full treatment of this subject later in section 3 introducing only the necessary notation for the moment. You may be used to thinking of a vector as something stretching from one point to another and which can be freely moved around. In relativity this is no longer true and so we must be more careful by what we mean by a vector. To each point p in spacetime we associate the set of all possible vectors located at that point. This is known as the *tangent space* at p , and denoted as T_p . A vector is a perfectly well-defined object geometric object, so too is a *vector field* which is a set of vectors with exactly one at each point in spacetime. The set of all the tangent spaces T_p of a manifold⁵ M is known as the *tangent bundle* $T(M)$. It is often useful to decompose vectors into components in terms of some basis. Recall that a *basis* is a set of vectors which both spans the vector space and is linearly independent. There are an infinite number of possible bases, but each will have the same number of basis elements, the dimension of the manifold. Let us imagine that at every point we set up a basis with four vectors \hat{e}_μ .

A standard example of a vector in spacetime is the tangent to a curve. We can specify a curve by specifying coordinates in terms of a parameter, $x^\mu(\lambda)$. The tangent vector has components

$$V^\mu = \frac{dx^\mu(\lambda)}{d\lambda} . \quad (2.12)$$

The full vector is then

$$V = V^\mu \hat{e}_\mu . \quad (2.13)$$

⁵We will define a manifold later in section ??.

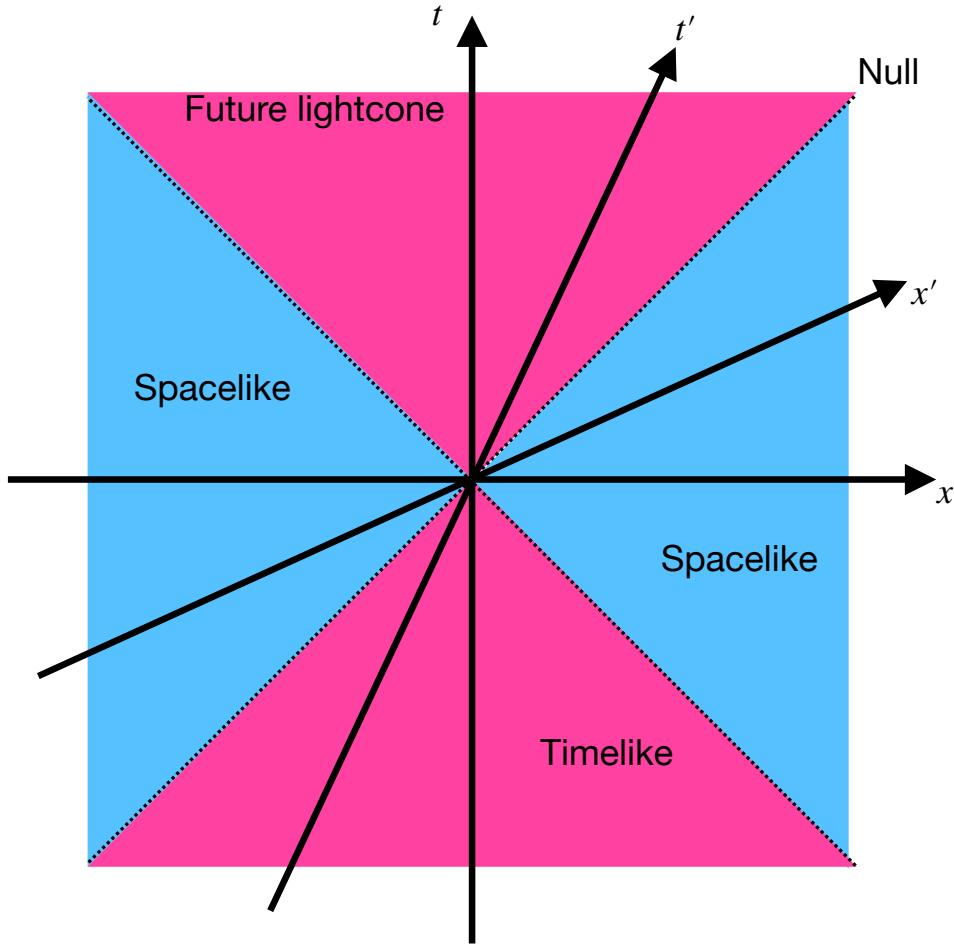


Figure 1: The lightcone diagram. The pink areas are time-like separated from the point at the centre, while points in the blue area are space-like separated. The dotted lines are light-like separated.

Under a Lorentz transformation the coordinates transform according to (2.2), and from this we may deduce the transformation of the components of the four-vector V^μ ,

$$V^\mu \rightarrow V'^\mu = \Lambda^\mu{}_\nu V^\nu, \quad (2.14)$$

when the coordinate system is transformed as in (2.2). Since the vector itself does not change under Lorentz transformations, and the parametrisation with λ is unaltered it follows that the basis vectors transform according to

$$\hat{e}_\mu = \Lambda^\nu{}_\mu \hat{e}'_\nu. \quad (2.15)$$

This is just multiplication by the inverse of the Lorentz transformation which transforms the

coordinates, therefore

$$\hat{e}'_\mu = \Lambda_\mu^\nu \hat{e}_\nu . \quad (2.16)$$

To summarise, we have introduced a set of coordinates labelled by upper indices which transform in a certain way under Lorentz transformations. We then considered vector components with upper indices which transformed in the same way as the coordinates. The basis vectors associated with the coordinate system transformed via the inverse matrix and were labelled by a lower index. These transformations leave invariant the vector, that is summing over the vector components with the basis vectors.

Once we have a vector space we can define and associated vector space known as the *dual vector space*. It is usually denoted with an asterisk, so that the dual vector space of the Tangent space T_p is T_p^* . The dual space is the space of all linear maps from the original vector space to the real numbers, so that if $\omega \in T_p^*$ then

$$\omega(aV + bW) = a\omega(V) + b\omega(W) \in \mathbb{R} , \quad (2.17)$$

for $V, W \in T_p$ and $a, b \in \mathbb{R}$. It follows that T_p^* is a vector space itself. We may introduce a basis of dual vectors $\hat{\theta}^\mu$ by fixing

$$\hat{\theta}^\mu(\hat{e}_\nu) = \delta_\nu^\mu . \quad (2.18)$$

Every dual vector can be written in components in terms of this basis as

$$\omega = \omega_\mu \hat{\theta}^\mu . \quad (2.19)$$

We will usually simply write ω_μ for the entire dual vector, and similarly write V^μ for the vector. Typically one refers to the elements of T_p as *contravariant four vectors* and elements of T_p^* as *covariant vectors*, or even *one-forms*, (a name that will make more sense after we have introduced differential geometry in section 3). The set of all cotangent spaces over M is called the *cotangent bundle* $T^*(M)$. The action of a dual vector field on a vector field is no longer a single number but a *scalar*, depending on the spacetime position. A scalar has no indices and is left invariant under Lorentz transformations.

The component notation is useful when considering the action of a dual vector on a vector:

$$\omega(V) = \omega_\mu V^\nu \hat{\theta}^\mu(\hat{e}_\nu) = \omega_\mu V^\nu \delta_\nu^\mu = \omega_\mu V^\mu . \quad (2.20)$$

The scalar product of a contravariant and covariant vector, which is invariant under Lorentz transformations

$$\omega'_\mu V'^\mu = \Lambda_\mu^\rho \Lambda^{\mu\sigma} \omega_\rho V^\sigma = \omega_\mu V^\mu , \quad (2.21)$$

in other words it is a *scalar*. It is from here that we can obtain the transformation of the dual vector: a *covariant* four-vector is a quantity which transforms as

$$\omega_\mu \rightarrow \omega'_\mu = \Lambda_\mu^\nu \omega_\nu , \quad (2.22)$$

where

$$\Lambda_\mu^\nu \equiv \eta_{\mu\rho} \eta^{\nu\sigma} \Lambda^\rho_\sigma , \quad (2.23)$$

with $\eta^{\mu\nu}$ the inverse of $\eta_{\mu\nu}$, which are numerically the same.⁶

To every contravariant vector we may associate a covariant vector by

$$\omega_\mu = \eta_{\mu\nu} V^\nu , \quad (2.24)$$

and vice-versa. We see that the Minkowski metric raises and lowers the indices of four-vectors.

One may extend the notion of a vector to a *tensor*. A tensor of type (rank) (k, l) , is a multilinear map from a collection of dual vectors and vectors to \mathbb{R} :

$$T : T_p^* \times \dots \times T_p^* \times T_p \times \dots \times T_p \rightarrow \mathbb{R} , \quad (2.25)$$

for example a scalar is a tensor of rank $(0,0)$, a vector a rank $(0, 1)$ tensor and a contravariant vector of rank $(1, 0)$. The space of all tensors of a fixed rank (k, l) forms a vector space. To construct a basis for this space it is useful to define the *tensor product* \otimes . If T is a (k, l) -tensor and S a (m, n) -tensor then $T \otimes S$ is a $(k + m, l + n)$ tensor defined to be

$$\begin{aligned} T \otimes S &(\omega^{(1)}, \dots, \omega^{(k)}, \dots, \omega^{(k+m)}, V^{(1)}, \dots, V^{(l)}, \dots, V^{(l+n)}) \\ &= T(\omega^{(1)}, \dots, \omega^{(k)}, V^{(1)}, \dots, V^{(l)}) S(\omega^{(k+1)}, \dots, \omega^{(k+m)}, V^{(l+1)}, \dots, V^{(l+n)}). \end{aligned} \quad (2.26)$$

As with vectors we will let T be a tensor of rank (k, l) , then under a Lorentz transformation it transforms as

$$T'^{\mu_1 \dots \mu_k}_{\nu_1 \dots \nu_l} = \Lambda^{\mu_1}_{\mu'_1} \dots \Lambda^{\mu_k}_{\mu'_k} \Lambda^{\nu'_1}_{\nu_1} \dots \Lambda^{\nu'_l}_{\nu_l} T^{\mu'_1 \dots \mu'_k}_{\nu'_1 \dots \nu'_l} . \quad (2.27)$$

One can use tensors to construct additional tensors either by taking linear combinations of tensors with the same upper and lower indices, direct products, contraction, or differentiation. The order of the indices of a tensor matters.

Note that because of the map between contravariant and covariant vectors via the Minkowski metric we can define an *inner product* on two vectors as

$$\eta(V, W) = \eta_{\mu\nu} V^\mu W^\nu . \quad (2.28)$$

⁶Using the properties of the Lorentz transformation it is not hard to show that Λ_μ^ν is the inverse of Λ^μ_ν .

Two vectors whose inner product vanishes are called *orthogonal*. Since it is a scalar the dot product is left invariant under Lorentz transformations and therefore orthogonality is basis and frame independent. We can define the *norm* of a vector to be the inner product with itself. Unlike in Euclidean geometry this is not positive definite, instead

$$\text{if } \eta_{\mu\nu} V^\mu V^\nu \text{ is } \begin{cases} < 0, & V^\mu \text{ is timelike,} \\ = 0, & V^\mu \text{ is lightlike or null,} \\ > 0, & V^\mu \text{ is spacelike,} \end{cases} \quad (2.29)$$

This is the more mathematical definition of these concepts from our earlier discussion.

Some tensors that will appear regularly are: the metric which is a $(0, 2)$ tensor, with the inverse being a $(2, 0)$ tensor, the *Kronecker delta* δ_ν^μ which is a $(1, 1)$ tensor, and finally the Levi–Civita tensor which is a $(0, 4)$ tensor. Not only can the metric be used to raise and lower indices of a tensor, it can also be used to contract indices. *Contraction* takes a (k, l) tensor to a $(k - 1, l - 1)$ tensor by

$$T^{\mu\nu\rho}_{\mu\sigma} \equiv S^{\nu\rho}_{\sigma}. \quad (2.30)$$

So far we have been very good and everything we have defined applies equally well for curved spacetime. We will now start to introduce some technology which requires modification when going to curved space. This will be given in the differential geometry section 3. Let us see how physics works in Minkowski space.

Let us start with a worldline, this is a curve in spacetime: $\gamma : [0, 1] \rightarrow \mathbb{R}^{1,3}$. Usually we will think of such a curve in inertial coordinates such that $\gamma(\lambda) = x^\mu(\lambda)$. The tangent vector to the curve at the point p with coordinate $x^\mu(\lambda_0)$ is

$$v^\mu|_p := \frac{d}{d\lambda} x^\mu(\lambda)|_{\lambda=\lambda_0}. \quad (2.31)$$

Note that the tangent depends on the parametrisation of the curve. An object of interest is the norm of the tangent vector as this characterises the path: if the tangent is timelike/null/spacelike for some value of λ we say that the path is timelike/null/spacelike at that point. Note that the sign of the norm of a tangent vector is quite a natural thing to classify the different tangent vectors by. The interval between two points on the other hand is not: it depends on the specific choice of path. In flat spacetime we think of straight lines between points and this is unique. When the manifold is curved this is no longer true.

A more natural object is the *line element* we introduced earlier in (2.5). Since ds^2 need not be positive we should not just take the square root and integrate along a curve, it depends

on the type of path. For spacelike paths we define the *path length*

$$\Delta s = \int_{\lambda_i}^{\lambda_f} \sqrt{\eta_{\mu\nu} \frac{dx^\mu}{d\lambda} \frac{dx^\nu}{d\lambda}} d\lambda, \quad (2.32)$$

where the integral is over the path start and end points. For null paths the integral is zero so we need not define anything. For timelike paths we define the *proper time*

$$\Delta\tau = \int_{\lambda_i}^{\lambda_f} \sqrt{-\eta_{\mu\nu} \frac{dx^\mu}{d\lambda} \frac{dx^\nu}{d\lambda}} d\lambda, \quad (2.33)$$

which is positive. The proper time is useful because of the *clock postulate*.

Clock Postulate *An accurate clock moving along a timelike worldline measures the proper time along the worldline.*

This point of view makes the “twin paradox” and similar puzzles clear. Two worldlines which have two intersections at different events will have proper times which measure their respective proper times, however these numbers in general will be different since the paths are different.

Note that the proper time is a convenient choice for parametrising a curve since it satisfies

$$\eta_{\mu\nu} v^\mu v^\nu = \eta_{\mu\nu} \frac{d}{d\tau} x^\mu(\tau) \frac{d}{d\tau} x^\nu(\tau) = -1. \quad (2.34)$$

Such a parameter can always be found along a timelike curve and is unique up to the start point of the curve. For spacelike curves the proper length gives a similarly useful parametrisation since

$$\eta_{\mu\nu} v^\mu v^\nu = \eta_{\mu\nu} \frac{d}{ds} x^\mu(s) \frac{d}{ds} x^\nu(s) = 1. \quad (2.35)$$

There is no such analogue along a null curve.

Massive paths Let us now consider the worldlines of massive particles. We will use the proper time as the parameter along the path with the path starting at $\tau = 0$. The tangent vector is known as the *four-velocity* U^μ :

$$U^\mu = \frac{dx^\mu}{d\tau}. \quad (2.36)$$

This is automatically normalised, $\eta_{\mu\nu} U^\mu U^\nu = -1$ since we parametrised the curve using the proper time. We may define the *energy-momentum four-vector* as

$$p^\mu = m U^\mu, \quad (2.37)$$

with m the mass of the particle. The mass is a fixed quantity independent of inertial frame, this is what you may have been used to calling the rest mass. The energy is simply p^0 , and as

one component of a four-vector is not invariant under Lorentz transformations. Note that in the particles rest frame we have $p^0 = m$ (recall $c = 1$) and so this is the celebrated $E = mc^2$. Note that the energy in the rest frame is the norm of the energy momentum four vector. In a general frame we have

$$E^2 - p^i p_i = m^2, \quad (2.38)$$

which is the full version of Einstein's famous formula.

We now want the special relativity version of Newton's second law. The requirement that it be tensorial puts some stringent constraints on the possible form, we must introduce a force four-vector f^μ satisfying

$$f^\mu = m \frac{d^2}{d\tau^2} x^\mu(\tau) = \frac{d}{d\tau} p^\mu(\tau). \quad (2.39)$$

For electromagnetism and the Lorentz force law ($f = q(E + v \times B)$) we find

$$f^\mu = q U^\nu F_\nu{}^\mu, \quad (2.40)$$

where F is the field strength of the electromagnetism gauge field.

Although p^μ provides a complete description of the energy and momentum of a particle for extended systems it is necessary to go further and define the *energy-momentum tensor*, or *stress tensor*, $T^{\mu\nu}$. This is a symmetric $(2,0)$ tensor which tells us all we need to know about the energy like aspects of a system: energy density, pressure, stress etc.. Consider a *fluid*. This is a continuum of matter described macroscopic quantities such as temperature, pressure, entropy, viscosity, etc. We will work with *perfect fluids* which are completely characterised by their pressure and density. This in particular means that they are isotropic (same in every direction) in the rest frame.

To understand this let us first consider *dust*. This is a collection of particles which are at rest with respect to each other, as a perfect fluid they have zero pressure. Since all the particles have an equal velocity in any fixed inertial frame we can imagine a four-velocity field $U^\mu(x)$ defined over all spacetime. We can define the *number-flux four-vector*

$$N^\mu = n U^\mu, \quad (2.41)$$

where n is the number density of the particles as measured in their rest frame. Then N^0 is the number density of particles as measured in any other frame, while N^i is the flux of particles in the i 'th direction. Let us imagine each of the particles have the same mass m . Then in the rest frame the energy density of the dust is given by

$$\rho = nm. \quad (2.42)$$

This completely specifies the dust, however this only measures the energy density in the rest frame, how do we measure it in other frames? Notice that both n and m are 0-components of four-vectors in their rest frame: $N^\mu = (m, 0, 0, 0)$ and $p^\mu = (m, 0, 0, 0)$. Therefore ρ is the $\mu = 0, \nu = 0$ component of the tensor $p \otimes N$ as measured in the rest frame. We are therefore lead to define the energy momentum tensor for dust

$$T_{\text{dust}}^{\mu\nu} = p^\mu N^\nu = nmU^\mu U^\nu = \rho U^\mu U^\nu, \quad (2.43)$$

where ρ is the energy density as measured in the rest frame.

We can now consider other perfect fluids. The key point is the isotropic in the rest frame property which implies that the energy momentum tensor must take a diagonal form in the rest frame, since there cannot be a net flux of momentum in an orthogonal direction. Moreover the spacelike components must all be equal $T^{11} = T^{22} = T^{33}$, there are only two independent components. We will take the two independent parameters to be the energy density ρ and the pressure p (note that p is also used for momentum but will always come with a superscript or subscript). In the rest frame the energy momentum tensor takes the form

$$T^{\mu\nu} = \begin{pmatrix} \rho & 0 & 0 & 0 \\ 0 & p & 0 & 0 \\ 0 & 0 & p & 0 \\ 0 & 0 & 0 & p \end{pmatrix}. \quad (2.44)$$

We want a formula which is good in any frame and therefore we want to write this in terms of tensors. For dust we had $T^{\mu\nu} = \rho U^\mu U^\nu$, so we may guess that there should be $(\rho + p)U^\mu U^\nu$, which gives $\rho + p$ in the 00 component and zero elsewhere in the rest frame. To include the remainder we should find something which is of the form $p\text{diag}(-1, 1, 1, 1)$ this is of course given by the Minkowski metric! The general form of the energy momentum tensor for a perfect fluid is

$$T^{\mu\nu} = (\rho + p)U^\mu U^\nu + p\eta^{\mu\nu}. \quad (2.45)$$

This will be important when we consider the cosmology section of the course.

2.2 Newtonian gravity

We can cast Newtonian gravity in terms of a field theory. The force acting on a particle of mass m is

$$F = -m\nabla\Phi(t, \vec{x}), \quad (2.46)$$

where the gravitational field $\Phi(t, \vec{x})$ is determined by the surrounding matter distribution $\rho(t, \vec{x})$,

$$\nabla^2 \Phi(t, \vec{x}) = 4\pi G \rho(t, \vec{x}), \quad (2.47)$$

where G is Newton's constant with approximate value

$$G \sim 6.67 \times 10^{-11} \text{m}^3 \text{kg}^{-1} \text{s}^{-2}. \quad (2.48)$$

This is simply a rewriting into field theory language of the inverse square law of Newton. For example if there is a mass M concentrated at a single point at $(t, \vec{0})$, then the mass density is

$$\rho(t, \vec{x}) = M \delta^{(3)}(\vec{x}), \quad (2.49)$$

which gives the gravitational field,

$$\Phi(\vec{x}) = -\frac{GM}{r}, \quad r^2 = \vec{x} \cdot \vec{x}. \quad (2.50)$$

This can be extended to more complicated matter distributions, either summing up contributions from the location of point-like particles or more generally by using the Greens function for the Laplacian and the mass density

$$\Phi(\vec{x}) = - \int d^3x' \frac{G\rho(\vec{x}')}{|\vec{x} - \vec{x}'|}. \quad (2.51)$$

Exercise: Newton's theorem

Newton's theorem states that the gravitational field outside of a spherically symmetric mass distribution depends only on its total mass. Show this by using (2.46), (2.47) and Gauss' theorem.

We can now insert the gravitational force law into Newton's second law of motion $F = ma$. At this point one should ask oneself whether the inertial mass appearing in Newton's second law is the same as the one appearing in the gravitational force law (2.46), there is no reason that they need to be the same. Application of Newton's second law gives

$$\vec{a} = -\frac{m_G}{m_i} \nabla \Phi, \quad (2.52)$$

with \vec{a} the acceleration. Starting with Galileo, Christaan Huygens all the way to more recent experimental data has shown that $m_i = m_G$ to an accuracy of 10^{-13} . This is known as the *weak equivalence principle*. In the Newtonian theory this appears as an isolated unexplained fact, however it is this experimental fact that underlies general relativity. Since all bodies with the same initial conditions fall along the same curve regardless of their composition, we can interpret that curve to be a property of the geometry of the spacetime not of a force acting on the body.

2.3 Equivalence Principles

The *Weak equivalence principle* was one of the starting points for the development of GR. It is motivated by thought experiments using Newtonian gravity. The exact equality of $m_i = m_g$ is one version of the weak equivalence principle. Newtonian gravity gives no explanation for why this should be true. A theory of gravity should be able to explain this. Another way to formulate the weak equivalence principle is

The trajectory of a freely falling test body depends only on its initial position and initial velocity and is independent of the composition of the body.

A consequence of the weak equivalence principle is that it is not possible to tell the difference between constant acceleration and a constant gravitational field. Suppose that you are in a closed box and consider the two situations 1) you are on earth, 2) you are in a spaceship undergoing constant acceleration. Within Newtonian mechanics there is no local experiment that you can perform which distinguishes the two.⁷ Another version of this is 1) the box is in free fall 2) you are floating in deep space. Again there is not local experiment that you can conduct to tell the difference. If the two situations can't be distinguished why do we describe them so differently?

This motivated the Einstein equivalence principle:

1) The weak equivalence principle is valid and 2) In a local inertial frame the results of all non-gravitational experiments will be indistinguishable from the results of the same experiments in an inertial frame in Minkowski spacetime.

The weak equivalence principle implies that 2) is valid for test bodies. The fact that test bodies which include ordinary matter which is held together by the three other forces, gives evidence that the electromagnetic and nuclear forces also obey 2).

Implications The Einstein equivalence principle implies that light is bent in a gravitational field. Consider a uniform gravitational field and a freely falling lab. Inside the lab the Einstein equivalence principle says that light rays must move on straight lines. But a straight line with respect to the lab corresponds to a curved path with respect to a frame at rest relative to the Earth. The effect is small but

⁷One of the important words is *local*. You can use tidal forces to distinguish between the two. Roughly if you drop two masses on Earth they will ever so slightly come together because the direction gravity acts on them is slightly different, they are pulled to the centre of the Earth. On a spaceship this is not the case and they fall down never getting closer together. This however is a non-local experiment, you need to watch the masses fall for a while and for a distance.

2.4 Some worked examples

Proper time along an accelerated worldline We treat the planets as being at rest relative to each other in this question.

Leia begins at rest on the planet Polis Massa and sets off in a spaceship to visit a distant planet called Alderaan. Alderaan is at rest relative to Polis Massa and is a proper distance D away. Leia's spaceship accelerates during the journey at a constant rate α ,

$$\eta_{\mu\nu}a^\mu a^\nu = \alpha^2, \quad (2.53)$$

where a^μ is the four-acceleration of Leia. We want to answer two questions: 1) what path does Leia take in terms of coordinates centred on Polis Massa? 2) How much time passes, from Leia's point of view until she reaches Alderaan?

We can choose coordinates (t, x, y, z) where the worldline of Polis Massa is simply $(t, 0, 0, 0)$ and the worldline of Alderaan is $(t, D, 0, 0)$ (recall that the two planets are at rest relative to each other). Leia's world line is then of the form

$$(t(\tau), x(\tau), 0, 0), \quad (2.54)$$

where τ is the proper time along Leia's worldline. Since we have parametrised Leia's worldline by the proper time we have

$$-\dot{t}(\tau)^2 + \dot{x}(\tau)^2 = -1 \quad \bullet \equiv \frac{d\bullet}{d\tau}. \quad (2.55)$$

Leia's acceleration is therefore,

$$a = (\ddot{t}(\tau), \ddot{x}(\tau), 0, 0) = \left(\frac{\dot{x}(\tau)\ddot{x}(\tau)}{\sqrt{1+\dot{x}(\tau)^2}}, \ddot{x}(\tau), 0, 0 \right), \quad (2.56)$$

where for the second equality we have used (2.55) to eliminate $\ddot{t}(\tau)$. Since Leia's acceleration is constant, (2.53), we have

$$\alpha^2 = \frac{\ddot{x}(\tau)^2}{1+\dot{x}(\tau)^2}. \quad (2.57)$$

We have that $\dot{x}(\tau) > 0$ and therefore the solution for $\dot{x}(\tau)$ is

$$\dot{x}(\tau) = \sinh(\alpha\tau + \beta), \quad (2.58)$$

with β a constant of integration. Since Leia began at rest on Polis Massa, we take $\beta = 0$. Integrating again and using that Leia begins at Polis Massa at $\tau = 0$, i.e. $x(0) = 0$, we have

$$x(\tau) = \frac{1}{\alpha} (\cosh(\alpha\tau) - 1). \quad (2.59)$$

Inserting this into (2.55), solving for $t(\tau)$ and imposing $t(0) = 0$ we find

$$t(\tau) = \frac{1}{\alpha} \sinh(\alpha\tau). \quad (2.60)$$

Leia reaches Alderaan when

$$\tau = \operatorname{arccosh}(1 + \alpha D), \quad (2.61)$$

If αD is large then $\tau \sim \frac{1}{\alpha} \log(\alpha D)$ and therefore no matter how large D is, for a sufficiently large acceleration Leia can reach Alderaan in a "reasonable" proper time. On the other hand, when Leia reaches Alderaan

$$t = \sqrt{D^2 + \frac{2D}{\alpha}}, \quad (2.62)$$

and therefore no matter how large α is it always takes at least a time of D (recall $c = 1$) to reach Alderaan as viewed from Polis Massa.

Null curves in Minkowski space By now we have all seen that a straight line is a null curve in Minkowski space but are there more? Note that we are not asking about geodesics. Consider the curve, given in inertial coordinates, by

$$x^\mu = (\lambda, \sin \lambda, \cos \lambda, 0). \quad (2.63)$$

The tangent to the vector is

$$v^\mu = \frac{dx^\mu}{d\lambda} = (1, \cos \lambda, -\sin \lambda, 0), \quad (2.64)$$

and has norm

$$v^\mu v^\nu \eta_{\mu\nu} = -1 + \cos^2 \lambda + \sin^2 \lambda = 0. \quad (2.65)$$

This is a null curve that is not straight, it is not a geodesic however.

Ladders and barns Barry and Paul Chuckle have been employed by Albert E. to put a ladder in a barn, a simple feat you would imagine but these are the Chuckle brothers and nothing is simple with them. Albert E. stands outside the barn, and tells Barry and Paul to run very quickly at a constant speed in a straight line through the barn carrying the ladder. The barn has doors at the front and back, and two apprentices (Jimmy and Brian) stand at either door ready to close or open them. Initially the front door is open and the back door is closed. The proper length of the ladder is l , while the proper length of the barn is b with $b < l$.

Albert E. claims that if Barry and Paul run fast enough, and that there is no slacking, then both doors of the barn can be temporarily closed with both Barry, Paul and the ladder inside the barn. One of the apprentices can then open the back door again so that Barry, Paul and the ladder can pass through the barn safely. The brothers are stumped, “oh dear, oh dear” says Barry, “the ladder is bigger than the barn, it will never work”. To put their minds at rest show that the ladder will fit in a chosen reference frame.

Let us work in inertial coordinates where the barn is at rest, which corresponds to Albert E.’s point of view. In these coordinates the front of the barn is at $x^\mu = (\lambda, 0, 0, 0)$ while the back of the barn is at $x^\mu = (\lambda, b, 0, 0)$.

The worldline of the front of the ladder in this reference frame is $x^\mu = (\lambda, v\lambda, 0, 0)$, where v is the velocity of the ladder. We have chosen coordinates so that the front of the ladder enters the barn at $\lambda = 0$. The back of the ladder follows the worldline $x^\mu = (\lambda, \lambda v - L, 0, 0)$ for some L which is not l !

First we must work out what L is in terms of l . We could of course perform a Lorentz transformation to switch to the rest frame of the ladder, the proper length of the ladder is then the coordinate length in this frame. We will use an alternative approach, staying in the original coordinate frame. How can we measure a length? Well we can define it to be half the proper time along the worldline at one end of the ladder between the emission and reception of a light signal which bounces off the other end of the body. The worldlines of the points making up the ladder are given by $x^\mu = (\lambda, \lambda v - \beta, 0, 0)$ where $\beta \in [0, L]$ and their tangent vectors are

$$\frac{dx^\mu}{d\lambda} = (1, v, 0, 0). \quad (2.66)$$

We now want to find a spacelike straight line orthogonal to this tangent vector. Such a worldline is given by $n^\mu = (-v\tilde{\lambda}, -\tilde{\lambda}, 0, 0)$. This curve meets the front of the ladder at $\tilde{\lambda} = 0$ and the back of the ladder at $\tilde{\lambda} = L(1 - v^2)^{-1}$. We want to calculate the proper length of this curve with $\tilde{\lambda} \in [0, \frac{L}{1-v^2}]$. To do so we should parametrise the curve by the proper length. The norm of the tangent of the above vector is $\eta_{\mu\nu}\dot{n}^\mu(\tilde{\lambda})\dot{n}^\nu(\tilde{\lambda}) = 1 - v^2$. Then the proper length is

$$s = \tilde{\lambda}\sqrt{1 - v^2}. \quad (2.67)$$

The ladder then has proper length

$$l = s\Big|_{\tilde{\lambda}=\frac{L}{1-v^2}} = \frac{L}{\sqrt{1 - v^2}}. \quad (2.68)$$

The entire ladder can fit into the barn from Albert E.'s perspective if $b \geq L$ and therefore the Chuckle brothers must run at a speed of

$$v \geq \sqrt{1 - \frac{b^2}{l^2}}. \quad (2.69)$$

Both doors of the barn can be closed if: the front of the ladder is still in the barn, $b > v\lambda$ and the back of the ladder is in the barn $\lambda v - L > 0$. Since $t = \lambda$ the ladder is in the barn for

$$\frac{l\sqrt{1-v^2}}{v} \leq t \leq \frac{b}{v}. \quad (2.70)$$

We see that Albert E. sees the ladder fully inside the barn with the doors closed.

Now consider what happens from the Chuckle brother's perspective. We can do a Lorentz transformation to coordinates in which they are at rest:

$$(t', x', y', z') = (\gamma t - \gamma vx, \gamma x - \gamma vt, y, z), \quad \gamma = \frac{1}{\sqrt{1-v^2}}. \quad (2.71)$$

In the Chuckle brother's coordinates the barn follows the worldline $(\lambda, -v\lambda, 0, 0)$, while the back of the barn follows the worldline $(\lambda, -v\lambda + b\sqrt{1-v^2}, 0, 0)$.

The front door can be closed when the front of the barn passes the back of the ladder, so $-v\lambda < -l$ and therefore the front door of the barn is closed for $t' > \frac{l}{v}$.

The back door must open when the front of the ladder is about to go through it. So it is closed until $b\sqrt{1-v^2} - vt' = 0$ and therefore the back door is closed for $t' \in [0, \frac{b}{v\gamma}]$. In summary we have

$$\begin{cases} \text{Front door closed} & \frac{l}{v} \leq t', \\ \text{Back door closed} & 0 \leq t' \leq \frac{b}{v\gamma}. \end{cases} \quad (2.72)$$

Since the ladder is longer than the barn $l > b$ and $\gamma \geq 1$ it follows that there is no time for which both doors are closed from the point of view of the Chuckle brothers. The entire ladder never fits into the barn from their perspective. The two view-points are depicted in figure 2.

Having been convinced by your arguments the brothers were off with a “to me, to you”.⁸

Planetary orbits in Newtonian mechanics Let us now consider the orbits of the planets in Newtonian mechanics. Let us set up a coordinate system where the massive body of mass M is at $r = 0$ and the planet of mass m is a distance r from that point. The Lagrangian describing the system is

$$\mathcal{L} = \frac{m}{2}(\dot{x}(t)^2 + \dot{y}(t)^2 + \dot{z}(t)^2) - V(r), \quad V(r) = -\frac{mMG}{r}. \quad (2.73)$$

⁸ChuckleVision was a British children's comedy tv show following the antics of the Chuckle brothers Barry and Paul. Carrying a ladder was a common theme.

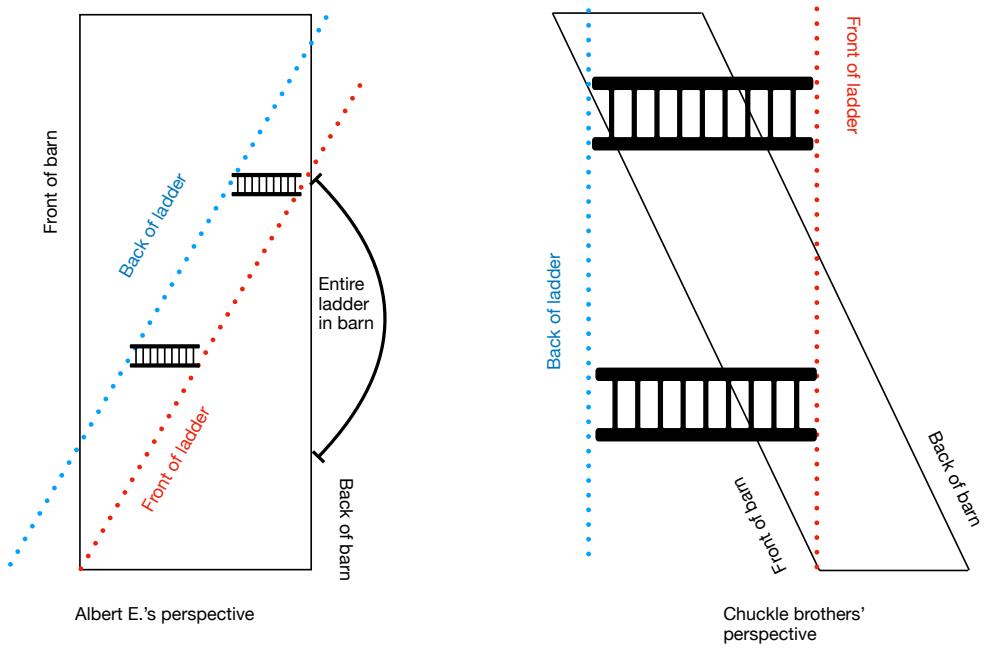


Figure 2: The two different perspectives of the ladder and barn. On the left from the perspective of Albert E., a stationary observer in the rest frame of the barn. On the right from the perspective of the Chuckle brothers carrying the ladder.

To make this more tractable it is useful to change coordinates to polar coordinates rather than Cartesian coordinates:

$$x = r \sin \theta \cos \phi, \quad y = r \sin \theta \sin \phi, \quad z = r \cos \theta. \quad (2.74)$$

The Lagrangian becomes

$$\mathcal{L} = \frac{m}{2} \left(\dot{r}^2 + r^2 (\dot{\theta}^2 + \sin^2 \theta \dot{\phi}^2) \right) + \frac{mMG}{r}. \quad (2.75)$$

We can now compute the equations of motion via the Euler–Lagrange equations: we find

$$\begin{aligned} \ddot{r} - r(\dot{\theta}^2 + \sin^2 \theta \dot{\phi}^2) + \frac{MG}{r^2} &= 0, \\ \frac{d}{dt}(r^2 \dot{\theta}) - r^2 \sin \theta \cos \theta \dot{\phi}^2 &= 0, \\ \frac{d}{dt}(r^2 \sin^2 \theta \dot{\phi}) &= 0. \end{aligned} \quad (2.76)$$

First let us consider the $\dot{\theta}$ equation. If we kick the particle off in the $\theta = \frac{\pi}{2}$ plane with $\dot{\theta} = 0$ then it will remain in that plane. We will make this choice from now on. The coordinate ϕ

is an ignorable coordinate since it does not appear explicitly in the Lagrangian. Recall that for every ignorable coordinate there is an associated conserved charge, in this case it will be the angular momentum. We may define

$$l = r^2 \dot{\phi}, \quad (2.77)$$

which is conserved. We have now solved the last two equations of (2.76) and only the first remains. Then we have

$$\ddot{r} - \frac{l^2}{r^3} + \frac{MG}{r^2} = 0. \quad (2.78)$$

To proceed further it is useful to note that there is one more conserved quantity, the Energy of the system. This follows since the Lagrangian is explicitly time independent, thus

$$Em = \frac{m}{2} \left(\dot{r}^2 + r^2 (\dot{\theta}^2 + \sin^2 \theta \dot{\phi}^2) \right) + V(r), \quad (2.79)$$

is conserved. We can now substitute $\dot{\theta}$ and $\dot{\phi}$ into this final condition to obtain an equation for \dot{r} only:

$$E = \frac{1}{2} \dot{r}^2 + \frac{l^2}{2r^2} - \frac{MG}{r} \equiv \frac{1}{2} \dot{r}^2 + V_N(r). \quad (2.80)$$

We can now study the orbits by looking at the Newtonian potential. At large distances the attractive $-r^{-1}$ dominates, while the angular momentum prohibits the particle from getting too close to the origin, see figure 3.

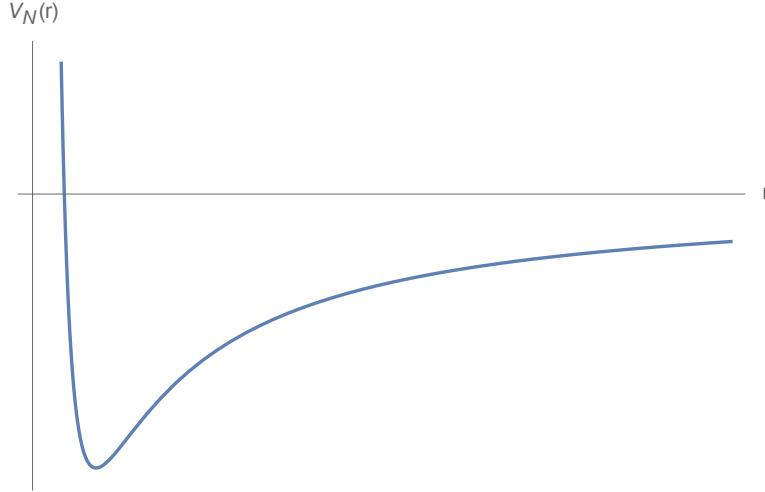


Figure 3: A representative example of the Newtonian potential.

The potential has a minimum when

$$V'(r_*) - \frac{MG}{r_*^2} - \frac{l^2}{r_*^3} = 0 \quad r_* = \frac{l^2}{MG}. \quad (2.81)$$

The planet can happily sit at $r = r_*$ for all time on a circular orbit, note that $E < 0$ in this case. The planet could also oscillate back and forth around the minima. This happens when $E < 0$ so that the planet cannot escape off to infinity. This describes an orbit where the distance to the massive body varies, as you may expect this is the usual elliptic orbit. For $E > 0$ the motion describes a flyby, the planet gets close to the massive body, never reaching it, before being flung off to infinity. Clearly such a planet would be dead and inhospitable for life.

So far we have discussed the radial motion of the planet, this does not tell us about the full motion however. We consider the orbit trajectory, the flyby motion is not so interesting for us. We now need to solve the angular momentum equations (2.77). To solve the coupled equations we start by employing a change of coordinates

$$u = r^{-1}, \quad (2.82)$$

and then view this as a function of ϕ . This works nicely because

$$\dot{u} = \frac{du}{d\phi} \dot{\phi} = lu^2 \frac{du}{d\phi}, \quad (2.83)$$

where we have used (2.77). We have

$$\dot{r} = -\frac{1}{u^2} \dot{u} = -l \frac{du}{d\phi}. \quad (2.84)$$

The conservation of energy equation (2.80) becomes

$$\left(\frac{du}{d\phi}\right)^2 + \left(u - \frac{GM}{l^2}\right)^2 = \frac{2E}{l^2} + \frac{G^2 M^2}{l^4}. \quad (2.85)$$

This turns out to be straightforward to solve, the solution is

$$u(\phi) = \frac{GM}{l^2} (1 + e \cos \phi). \quad (2.86)$$

In the original radial coordinate we have

$$r(\phi) = \frac{l^2}{GM} \frac{1}{1 + e \cos \phi}. \quad (2.87)$$

This is the equation for a conic section with the eccentricity given by

$$e = \sqrt{1 + \frac{2El^2}{G^2 M^2}}. \quad (2.88)$$

The shape of the orbit depends on the eccentricity. Motion with $E > 0$ is not in a bounded orbit, tracing out a hyperbola for $e > 1$ and a parabola for $e = 1$. Objects in orbit have

$e < 1$ with elliptical orbits. An important thing to note about this solution is that the orbit does not *precess*, its closest approach to the origin, known as the *perihelion*⁹ is always at the same point it does not precess, nor does the furthest point of the orbit, the *aphelion*. This disagrees with observations of Mercury's orbit and is the first observational discrepancy of Newtonian gravity.¹⁰

2.5 Problems with Newtonian gravity and why we need GR

Newton's theory of gravitation is successful in explaining the motions of the moon and planets. Some irregularities in the orbit of Uranus remained unexplained until the irregularities were used independently by John Couch Adams and Jean Joseph Le Verrier in 1846, to predict the existence and position of Neptune. There were still issues with predictions from Newtonian gravity and experimental data however. The precession of the perihelion of Mercury was one such problem. It was shown to be out by $43''/\text{century}$ ¹¹, recall that in the section above we showed that the perihelion does not precess in Newtonian gravity. We will see later how GR corrects this. A more obvious (and mathematical) problem arose after Einstein's work on special relativity in 1905. Newtonian gravity is incompatible with special relativity. A body can, in principle, be accelerated to a speed greater than the speed of light. Moreover, effects are instantaneous in Newtonian gravity clearly this is not allowed in special relativity where the speed of light gives an upper bound on the transfer of information.

Despite Newtonian gravities' failings it is sufficient for studying a large range of phenomena. To understand when a relativistic theory is needed let us consider a circular orbit around a star of mass M . The speed of the planet is easily computed by equating the centripetal force with the gravitational force giving,

$$\frac{v^2}{r} = \frac{GM}{r^2}. \quad (2.89)$$

Relativistic effects become important when $v \sim c$ and therefore the dimensionless parameter which governs corrections to Newtonian gravity is

$$\frac{GM}{rc^2}. \quad (2.90)$$

⁹Strictly this is for the closest approach to the sun. *Helios* is the word for the sun in greek, while *peri* means around.

¹⁰To perform a more accurate computation one should also take into account the effect of the gravitation fields of the other planets. This is notoriously difficult since one has to study a multi-body problem. Instead what one can do is imagine that the other planets to consider form a shell of mass along their orbit. One can then evaluate the force due to these. This approximation works if one considers the problem over a long enough time. Since planets closer to the sun have quicker orbits over a long enough time this approximation will give a reasonable result.

¹¹The " stands for arcseconds, with 3600 arcseconds(=3600") in a degree.

There is a convenient length scale which one can construct from a mass and the fundamental constants known as the *Schwarzschild radius*,¹²

$$R_s = \frac{2GM}{c^2}. \quad (2.91)$$

Relativistic corrections to gravity are then necessary when $R_s \sim 2r$. By this measure the earth is not a relativistic system $R_s \sim 10^{-2}m$ and the corrections on the surface of the Earth are of the order 10^{-8} . For satellites in orbit this is even smaller $\sim 10^{-9}$ however for GPS satellites clocks with such high precision are needed that this effect can be seen and if GR was not taken into account would fail very soon. The sun has $R_s \sim 3\text{km}$ and for Mercury the corrections are of order 10^{-7} , clearly very small but over a century the precession of Mercury's perihelion adds up to the previously quoted $43''$.

General relativity is the theory that replaces both Newtonian gravity and special relativity. However, general relativity is not the final theory of gravity, one eventually needs a theory of quantum gravity. General relativity breaks down for very extreme phenomena where quantum effects become important, e.g. the Big Bang and inside black holes. If one views gravity as a classical field theory and attempts to quantise it one finds that it is perturbatively non-renormalizable. Essentially this means that to obtain sensible observable results we must absorb infinities in computations by introducing new parameters. For a renormalizable theory we need to introduce only a finite number of these new parameters but for a non-renormalizable theory we need to introduce an infinite number, rendering the theory unable to give meaningful predictions. A candidate theory for quantum gravity, but no means the only candidate, is string theory. We should emphasise that a theory of quantum gravity is only needed for these extreme phenomena, general relativity is still worth learning and using.

3 Differential Geometry

Gravity is geometry and to properly understand general relativity we need to be able to understand curved spacetime. This is the language of differential geometry. Our discussion will not be all encompassing, there will be both topics and proofs that we omit. Instead we will build up all the necessary mathematical structure, in a logical order, that we will need to understand general relativity. As we proceed many of the objects that we will introduce may already be familiar to you, they will however take a different guise in places.

This section closely follows the excellent book by Nakahara: *Geometry, Topology and physics*.

¹²We will see this appear later when we consider the Schwarzschild solution in section 6.1.

3.1 Manifolds

Definition Let X be any set and $\mathcal{T} = \{U_i | i \in I\}$ denote a certain collection of subsets of X . The pair (X, \mathcal{T}) is called a *topological space* if \mathcal{T} satisfies

1. Both the set X and the empty set \emptyset are open subsets: $M \in \mathcal{T}$ and $\emptyset \in \mathcal{T}$.
2. If \mathcal{T} is any, possibly infinite, sub-collection of I , then the family $\{U_j | j \in J\}$ satisfies $\cup_{j \in J} U_j \in \mathcal{T}$.
3. If K is any finite sub-collection of I then the set $\{U_k | k \in K\}$ satisfies $\cap_{k \in K} U_k \in \mathcal{T}$.

Sometimes X alone is called a topological space, i.e. without associating to it a topology. The sets U_i are called *open sets* (we may sometimes refer to them as coordinate patches, the reason why will become obvious later) and \mathcal{T} gives a *topology* to X .

Examples

- a) If X is a set and \mathcal{T} a collection of all subsets of X then this is a topological space, and is known as the *discrete topology*.
- b) Let X be a set and take $\mathcal{T} = \{\emptyset, X\}$. This is then a topological space and the topology is known as the *trivial topology*. While the discrete topology is too stringent, this topology is too trivial.
- c) Take $X = \mathbb{R}$. All open subsets (a, b) (a, b may be $\mp\infty$ respectively) and their unions define a topology known as the *usual topology*.

Exercise: Consider the usual topology on \mathbb{R} and show that if we allow for an infinite number of open sets in condition 3 for the definition of a topological space, then the usual topology reduces to the discrete topology.

A *metric* $d : X \times X \rightarrow \mathbb{R}$ is a function that for any $x, y, z \in X$ satisfies:

1. $d(x, y) = d(y, x)$,
2. $d(x, y) \geq 0$ with equality iff $x = y$,
3. $d(x, y) + d(y, z) \geq d(x, z)$.

If X is endowed with a metric then X is made a topological space whose open sets are given by open discs

$$U_\epsilon(x) = \{y \in X | d(x, y) < \epsilon\}, \quad (3.1)$$

and all possible unions. The topology \mathcal{T} is called the *metric topology* determined by d .

Definition: Suppose \mathcal{T} gives a topology to X . Then N is a neighbourhood of the point $x \in X$ if N is a subset of X and N contains at least one open set U_i which contains x . Note that there is no requirement for N to be open, in the case where it is open it is called an *open neighbourhood*.

Definition: A topological space (X, \mathcal{T}) is a *Hausdorff space* if for an arbitrary pair of distinct points $x, y \in X$, there always exists neighbourhoods U_x and U_y such that $U_x \cap U_y = \emptyset$.

Example Let $X = \{A, B, C, D\}$ define the sets

$$U_0 = \emptyset, \quad U_1 = \{A\}, \quad U_2 = \{A, B\}, \quad U_3 = \{A, B, C, D\}. \quad (3.2)$$

Then the topology $\mathcal{T} = \{U_0, U_1, U_2, U_3\}$ makes X a topological space but it is not Hausdorff. First note that both the empty set and X are in the topology \mathcal{T} , satisfying 1 of the definition of a topological space. Note that the union of these sets is within \mathcal{T} thereby satisfying the second requirement. Finally the intersection of any of the sets is within \mathcal{T} and therefore it is a topological space. To see why it is not Hausdorff it suffices to show that we can pick two points which have no open sets in which one of the points is in and that the intersection of these open sets is not the empty set. There are a few choices we could take but an obvious one is C, D . They both appear in only one open set and therefore the space cannot be Hausdorff.

Most examples in physics that one encounters are Hausdorff spaces. We will assume this is the case in this course since the property protects us against funky things happening.

Definition Let X and Y be topological space. A map $f : X \rightarrow Y$ is *continuous* if the inverse image of an open set in Y is an open set in X . Note that a continuous function does *not* need to map an open set in X to an open set in Y , $f(x) = x^2$ is an example of a continuous function that would fail this requirement.

Definition Let (X, \mathcal{T}) be a topological space. A subset A of X is *closed* if its complement $X - A \in \mathcal{T}$ in X is an open set. The *closure* of the subset A is the smallest closed set that contains A and is denoted by \bar{A} . The *interior* of A is the largest open subset of A and is denoted by A° . The *boundary* $b(A)$ of A is the complement of A° in \bar{A} : $b(A) = \bar{A} - A^\circ$. An open set is always disjoint from its boundary while a closed set always contains its boundary.

To make this a little more clear let us consider a concrete example.

Example: Let us consider \mathbb{R}^2 with the metric topology and let A be the open set $\{(x, y) \in \mathbb{R}^2 | x^2 + y^2 < 1\}$. Then the closure of A is

$$\bar{A} = \{(x, y) \in \mathbb{R}^2 | x^2 + y^2 \leq 1\}. \quad (3.3)$$

The interior of A is itself $A^\circ = A$. The boundary is then the complement of A° in \bar{A} , thus

$$b(A) = \{(x, y) \in \mathbb{R}^2 \mid x^2 + y^2 = 1\}. \quad (3.4)$$

It therefore agrees with our usual understanding of these concepts.

Definition Let (X, \mathcal{T}) be a topological space. A family $\{A_i\}$ of subsets of X is called a *covering* of X if

$$\bigcup_{i \in I} A_i = X. \quad (3.5)$$

If all the A_i happen to be the open sets of the topology \mathcal{T} then the covering is called an *open covering*.

Definition Consider a set X and all possible coverings of X . The set X is *compact* if for every open covering $\{U_i \mid i \in I\}$ there exists a *finite* subset J of I such that $\{U_j \mid j \in J\}$ is also a covering of X .

Theorem Let X be a subset of \mathbb{R}^n , then X is compact iff it is closed and bounded.

Definition

- i) A topological space X is *connected* if it cannot be written as $X = X_1 \cup X_2$ where X_1 and X_2 are both open and $X_1 \cap X_2 = \emptyset$. Otherwise X is called *disconnected*.
- ii) A topological space is called *arcwise connected* if for any points $x, y \in X$ there exists a continuous map $f : [0, 1] \rightarrow X$ such that $f(0) = x$ and $f(1) = y$. Only in a few pathological cases is arcwise connectedness not equivalent to connectedness.
- iii) A *loop* in a topological space X is a continuous map $f : [0, 1] \rightarrow X$ such that $f(0) = f(1)$. If every loop in X can be continuously shrunk to a point, X is called *simply connected*.

Some simple examples are:

- $\mathbb{R}^2 - \mathbb{R}$ is not arcwise connected.
- $\mathbb{R}^2 - \{0\}$ is arcwise connected but not simply connected.
- $\mathbb{R}^3 - \{0\}$ is arcwise connected and simply connected.
- The n -dimensional torus is arcwise connected but not simply connected.

The main purpose of topology is to classify spaces. Suppose we have several figures, we want to be able to say which are equal and which are different, and probably more fundamentally what does being equal or different mean. In topology two figures are equivalent if it is

possible to deform them continuously into each other. We therefore construct an equivalence relation under which geometrical objects are classified according to whether it is possible to deform one into the other. Of course these are just words and we should define this more mathematically. To wit let us define

Definition Let X_1 and X_2 be two topological spaces. A map $f : X_1 \rightarrow X_2$ is a *homeomorphism* if it is continuous and has an inverse $f^{-1} : X_2 \rightarrow X_1$ which is also continuous. If there exists a homeomorphism between X_1 and X_2 we say that X_1 and X_2 are *homeomorphic* to each other.

The classic example of two homeomorphic spaces are a donut and a coffee mug.

One would like a quick way to understand whether two spaces are homeomorphic to each other. Even today we cannot fully characterise the equivalence classes between spaces. One modest statement that we can make is that if two spaces have different *topological invariants* then they are not homeomorphic to each other. A topological invariant is conserved under homeomorphisms. It may be a number such as the number of connected components of the space, an algebraic structure such as a group or a ring which can be constructed from the space, or something like connectedness, compactness or the Hausdorff property. If we knew the complete class of topological invariants we could specify the equivalence classes easily, however so far we only know a partial list. As such even if all the known topological invariants of two spaces coincide these spaces may still not be homeomorphic.

We are now finally in a position to define a *manifold*. An n -dimensional manifold is a space which looks locally like \mathbb{R}^n . Globally it need not be \mathbb{R}^n but we may glue local patches, each of which look like \mathbb{R}^n together to get the full global space. A manifold is then homeomorphic to \mathbb{R}^n locally. The local homeomorphism allows us to give each point on the manifold a set of n numbers called local *coordinates*. If the manifold is not homeomorphic to \mathbb{R}^n then we need to cover it in more than one patch, and so we need to introduce several local coordinates. We will require that the transition functions between these coordinates on the overlapping region are *smooth*. In this way we can develop the usual notion of calculus on a manifold. Topology is based on continuity, while manifolds is based on smoothness. With that let us begin with our definitions again.

Definition M is an n -dimensional *differentiable manifold* if it satisfies:

1. M is a Hausdorff topological space,
2. M is provided with a family of pairs $\{(U_i, \varphi_i)\}$;
3. $\{U_i\}$ is a family of open sets which covers M : $\cup_i U_i = M$.

4. φ_i is a homeomorphism from U_i onto an open subset U'_i of \mathbb{R}^n ,
5. Given U_i and U_j such that $U_i \cap U_j \neq \emptyset$, then the map $\psi_{ij} = \varphi_i \circ \varphi_j^{-1}$ from $\varphi_j(U_i \cap U_j)$ to $\varphi_i(U_i \cap U_j)$ is infinitely differentiable. ψ_{ij} is known as a *transition function*.

In figure 4 we have represented (well copied the image from Nakahara) the ideas above.

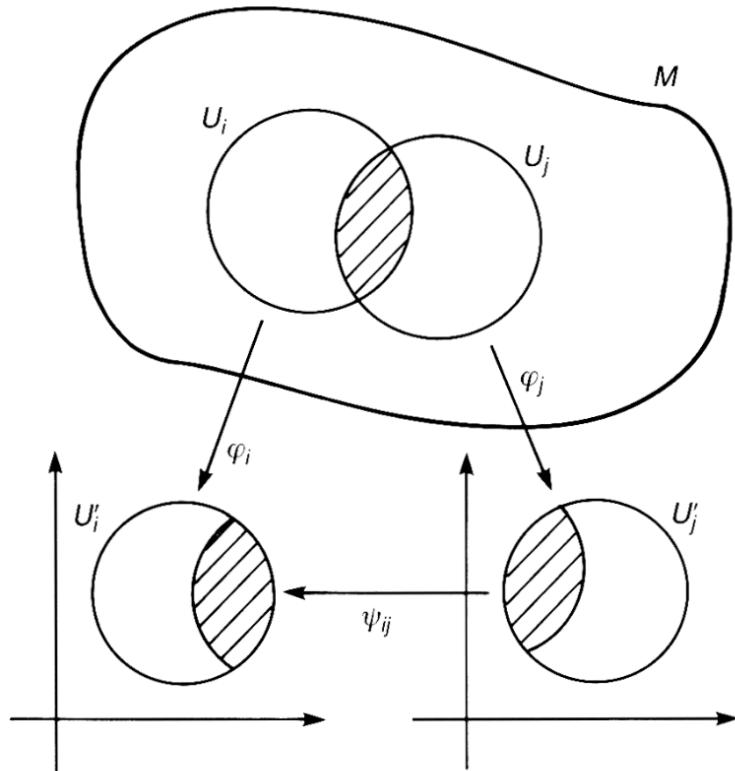


Figure 4: Here we see the manifold M and two coordinate charts. The homeomorphisms φ_i maps U_i onto an open set of $U'_i \subset \mathbb{R}^n$ providing coordinates for the point $p \in M$. If $U_i \cap U_j \neq \emptyset$ the transition functions from one coordinate system to another is smooth.

The pair (U_i, φ_i) are called a *chart* and the collection of charts is called an *atlas*. The subsets U_i are called the *coordinate neighbourhood* while the φ_i is called the *coordinate function*, or simply the *coordinate*. The homeomorphism φ_i is represented by n functions $\{x^1(p), \dots, x^n(p)\}$, with this set $\{x^\mu(p)\}$ also called the *coordinate*. A point $p \in M$ exists independently of its coordinates, however we will often be sloppy and denote the point p through its coordinates.

If U_i and U_j overlap, two coordinate systems are assigned to the same point in $U_i \cap U_j$. Axiom 5 asserts that the transition function from one coordinate system to another be smooth C^∞ . One may be alarmed by this but there is no reason for trepidation, it is analogous to labelling a point by Euclidean coordinates and polar coordinates. The map φ_i assigns n coordinates values x^μ , ($1 \leq \mu \leq n$) to a point $p \in U_i \cap U_j$, while φ_j assigns coordinates y^μ to the same point. The transition function from y to x , $x^\mu = x^\mu(y)$ is given by n functions of n variables, and is the explicit form of the map $\psi_{ji} = \varphi_j \circ \varphi_i^{-1}$. The differentiability in the definition is then in the usual sense we are familiar from calculus. All this leads to use being able to move over M however we may choose with the coordinates varying in a smooth way.

If the union of two atlases $\{(U_i \varphi_i)\}$ and $\{(V_j \psi_j)\}$ is again an atlas, then these two atlases are said to be *compatible*. The compatibility is an equivalence relation. This equivalence class is called the *differentiable structure*. Mutually compatible atlases define the same differentiable structure on M .

Let us briefly comment on manifolds with a boundary. We have assumed that the coordinate neighbourhood U_i is homeomorphic to an open set of \mathbb{R}^n . In some cases this is too restrictive. If a topological space M is covered by a family of open sets $\{U_i\}$ each of which is homeomorphic to an open set $H^n \equiv \{(x^1, \dots, x^n) \in \mathbb{R}^n | x^n \geq 0\}$, M is said to be a manifold with boundary. The analogous plot of figure 4 for the manifold with a boundary is given in figure 5.

The set of points which are mapped to points with $x^n = 0$ is called the *boundary* of M and is denoted by ∂M . The coordinates on ∂M are given by $n - 1$ numbers $(x^1, \dots, x^{n-1}, 0)$. We now need to be careful when we define smoothness on the overlaps. The map $\psi_{ij} : \varphi_j(U_i \cap U_j) \rightarrow \varphi_i(U_i \cap U_j)$ is defined on an open set of H^n in general, and ψ_{ij} is said to be smooth if it is C^∞ in an open set of \mathbb{R}^n which contains $\varphi_j(U_i \cap U_j)$.

Examples

- \mathbb{R}^n is a differentiable manifold trivially. A single chart covers the whole space and we take φ to be the identity map.
- Let $n = 1$ and let us impose connectedness. Then there are two choices, either \mathbb{R} or the circle S^1 . Let us work out an atlas for S^1 . For concreteness let us embed the circle in \mathbb{R}^2 via $x^2 + y^2 = 1$. We will need at least two charts. We can take them as in figure 6.

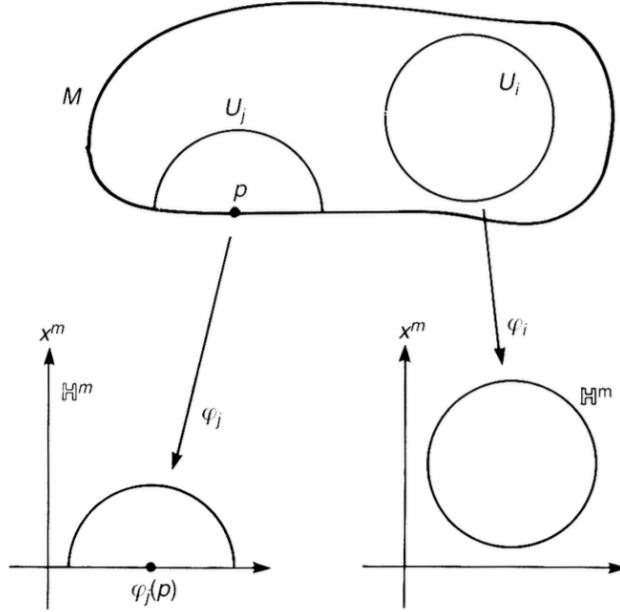


Figure 5: A manifold with a boundary. The point p is on the boundary. Note the subtle difference, for a manifold without a boundary the left figure would be extended below $x^n=0$.

Define $\varphi_1^{-1} : (0, 2\pi) \rightarrow S^1$ by¹³

$$\varphi_1^{-1} : \theta \rightarrow (\cos \theta, \sin \theta), \quad (3.6)$$

whose image is $S^1 - \{(1, 0)\}$. Similarly define $\varphi_2^{-1} : (-\pi, \pi) \rightarrow S^1$ by

$$\varphi_2^{-1} : \theta \rightarrow (\cos \theta, \sin \theta), \quad (3.7)$$

whose image is $S^1 - \{(-1, 0)\}$. Clearly both φ_i^{-1} are invertible and all the maps are continuous, thus the φ_i 's are homeomorphisms. The transition functions seem trivial for this example but one must be careful to end up in the correct domain. The two charts overlap on the upper and lower hemispheres and therefore we have

$$\varphi_2(\varphi_1^{-1}(\theta)) = \begin{cases} \theta & \text{if } \theta \in (0, \pi) \\ \theta - 2\pi & \text{if } \theta \in (\pi, 2\pi) \end{cases}. \quad (3.8)$$

The transition function isn't defined at $\theta = 0$ or $\theta = \pi$, nonetheless it is smooth on each of the two overlapping open sets as required.

¹³Until now we would just have taken the range to be $\theta \in [0, 2\pi)$ and been happy with this. However this does not meet our requirement of being a chart since it is not an open set. This would present problems later when we try to differentiate anything at $\theta = 0$. Recall that the derivative requires us to be able to take limits from both sides, and since there is nothing smaller than 0 we are stuck.

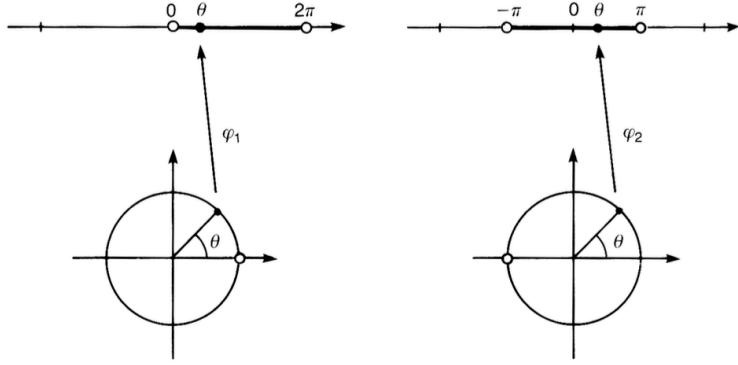


Figure 6: Two charts on S^1 .

- Let us consider a slightly less trivial example, the n -dimensional sphere S^n . We may realise it by embedding it in \mathbb{R}^{n+1} . (Note that embedding it in a higher-dimensional space is just for convenience and not a necessary requirement for being a manifold, in fact some n -dimensional spaces cannot be embedded in \mathbb{R}^{n+1} , for example hyperbolic space.)

We can realise the n -dimensional sphere S^n in \mathbb{R}^{n+1} as

$$\sum_{i=0}^n (x^i)^2 = 1. \quad (3.9)$$

We can introduce coordinate neighbourhoods

$$\begin{aligned} U_{i+} &\equiv \{(x^0, x^1, \dots, x^n) \in S^n | x^i > 0\}, \\ U_{i-} &\equiv \{(x^0, x^1, \dots, x^n) \in S^n | x^i < 0\}. \end{aligned} \quad (3.10)$$

Next define the coordinate map $\varphi_{i+} : U_{i+} \rightarrow \mathbb{R}^n$ to be

$$\varphi_{i+}(x^0, \dots, x^n) = (x^0, \dots, x^{i-1}, x^{i+1}, \dots, x^n), \quad (3.11)$$

and $\varphi_{i-} : U_{i-} \rightarrow \mathbb{R}^n$ to be

$$\varphi_{i-}(x^0, \dots, x^{i-1}, x^{i+1}, \dots, x^n). \quad (3.12)$$

Note that the domains of φ_{i+} and φ_{i-} are different and they have no overlap. Instead they are the projections of the hemispheres $U_{i\pm}$ to the plane $x^i = 0$. The transition functions can be obtained simply from the above maps. As an example let us take S^2 , then we have six coordinate neighbourhoods: $U_{x\pm}, U_{y\pm}, U_{z\pm}$. The transition function $\psi_{(y-)(x+)}$ is given by

$$\psi_{(y-)(x+)} : (y, z) \rightarrow \left(\sqrt{1 - y^2 - z^2}, z\right). \quad (3.13)$$

This is infinitely differentiable on $U_{x+} \cap U_{y-}$.

We have seen that to describe n -dimensional spheres we need more than one chart. The need to deal with multiple charts arises when we consider manifolds of non-trivial topology. When we come to discuss general relativity we will care a lot about changing coordinates and the limitations of the coordinate systems. In almost all situations that we will consider a single set of coordinates generally covers enough of the space to tell us everything we need to know. However as one progresses in physics, topology becomes more important. We will not see much of this but you should see this in some of your other physics/mathematics courses.

3.2 Calculus on manifolds

The reason why differentiable manifolds are useful is because it allows us to use the usual calculus we have developed on \mathbb{R}^n . Smoothness of the transition functions implies that the calculus is independent of the chosen coordinates.

Differentiable maps Let $f : M \rightarrow N$ be a map from an m -dimensional manifold M to an n -dimensional manifold N . A point $p \in M$ is mapped to a point $f(p) \in N$. We may take a chart (U, φ) on M and a chart (V, ψ) in N where for all $p \in U$, $f(p) \in V$. Then f has the following coordinate presentation:

$$\psi \circ f \circ \varphi^{-1} : \mathbb{R}^m \rightarrow \mathbb{R}^n. \quad (3.14)$$

If we write $\varphi(p) = \{x^\mu\}$ and $\psi(f(p)) = \{y^\alpha\}$ then, $\psi \circ f \circ \varphi^{-1}$ is just the usual vector-valued function $y = \psi \circ f \circ \varphi^{-1}(x)$ of m variables. Sometimes it is useful to abuse notation and write $y = f(x)$ or $y^\alpha = f^\alpha(x^\mu)$ when we know the coordinate systems on M and N that are in use.

Definition We say that a function $f : M \rightarrow \mathbb{R}$ is *smooth* if the map $f \circ \varphi^{-1} : U \rightarrow \mathbb{R}$ is smooth for all charts. We let the set of all small functions on M be denoted by $\mathcal{F}(M)$.

Definition We say that a map $f : M \rightarrow N$ between two manifolds is smooth if the map $\psi \circ f \circ \varphi^{-1} : U \rightarrow V$ is smooth for all charts $\varphi : M \rightarrow \mathbb{R}^m$ and $\psi : N \rightarrow \mathbb{R}^n$. If $y = \psi \circ f \circ \varphi^{-1}(x)$ is C^∞ then we say that f is *differentiable* at p . This is actually independent of the coordinate system.

Definition Let $f : M \rightarrow N$ be a homeomorphism and ψ and φ coordinate functions. If $\psi \circ f \circ \varphi^{-1}$ is invertible, f is called a *diffeomorphism* and M is said to be *diffeomorphic* to N and vice-versa. This is denoted by $M \equiv N$.

Since the map is invertible it follows that if $M \equiv N$ then $\dim M = \dim N$. Homeomorphisms classify spaces according to whether it is possible to deform one space into another

continuously. *Diffeomorphisms* classify spaces into equivalence classes according to whether it is possible to deform one space into the other *smoothly*. As such a diffeomorphism is stronger than a homeomorphism, it requires that both the map and its inverse are smooth. Two diffeomorphic manifolds are viewed as the same manifold.

3.2.1 Tangent Vectors

Having defined maps on a manifold we can define other objects on the manifold. The elementary notion of a vector no longer works: where is the origin, what is a straight arrow. On a manifold a vector is defined to be a *tangent vector* to a curve in M .

To define a tangent vector we need a curve $\gamma : (a, b) \rightarrow M$ and a function $f : M \rightarrow \mathbb{R}$. Let $t \in (a, b)$ then we can define the tangent vector at $\gamma(0)$ as a directional derivative of a function $f(\gamma(t))$ along the curve $\gamma(t)$ at $t = 0$. The rate of change of $f(\gamma(t))$ at $t = 0$ along the curve is

$$\frac{df(\gamma(t))}{dt} \Big|_{t=0}. \quad (3.15)$$

In terms of local coordinates this becomes

$$\frac{\partial f}{\partial x^\mu} \frac{dx^\mu(\gamma(t))}{dt} \Big|_{t=0}. \quad (3.16)$$

Notice the abuse of notation, the first term should really be

$$\frac{\partial(f \circ \varphi^{-1}(x))}{\partial x^\mu}, \quad (3.17)$$

we will persist with this abuse of notation. This is then equivalent to applying the differential operator

$$X = X^\mu \left(\frac{\partial}{\partial x^\mu} \right), \quad X^\mu = \frac{dx^\mu(\gamma(t))}{dt} \Big|_{t=0} \quad (3.18)$$

then

$$\frac{df(\gamma(t))}{dt} \Big|_{t=0} = X^\mu \frac{\partial f}{\partial x^\mu} \equiv X[f]. \quad (3.19)$$

We define X as the tangent vector to M at $p = \gamma(0)$ along the direction given by the curve $\gamma(t)$.

One can define an equivalence class of curves on M . If two curves $\gamma_1(t)$ and $\gamma_2(t)$ satisfy

$$(i) \quad \gamma_1(0) = \gamma_2(0) = p,$$

$$(ii) \quad \frac{dx^\mu(\gamma_1(t))}{dt} \Big|_{t=0} = \frac{dx^\mu(\gamma_2(t))}{dt} \Big|_{t=0},$$

then $\gamma_1(t)$ and $\gamma_2(t)$ yield the same differential operator X at p . This allows us to define the equivalence relation between curves at the point p , $\gamma_1(t) \sim \gamma_2(t)$. We identify the *tangent vector* X with the *equivalence class of curves*

$$[\gamma(t)] = \left\{ \tilde{\gamma}(t) \middle| \gamma(0) = \tilde{\gamma}(0) \text{ and } \frac{dx^\mu(\gamma(t))}{dt} \Big|_{t=0} = \frac{dx^\mu(\tilde{\gamma}(t))}{dt} \Big|_{t=0} \right\} \quad (3.20)$$

rather than a particular representative of the curve.

All the equivalence classes of curves at a point $p \in M$, i.e. all the tangent vectors at p , form a vector space called the *tangent space* of M at p , $T_p(M)$. We can take a basis of vectors for $T_p(M)$ to be $e_\mu = \frac{\partial}{\partial x^\mu}$. It follows that $\dim T_p(M) = \dim(M)$. The basis $\{e_\mu\}$ is called the *coordinate basis*. If a vector $X \in T_p(M)$ is written as $X = X^\mu e_\mu$ the numbers X^μ are called the components of X with respect to the basis $\{e_\mu\}$. The vector X exists without specifying a choice of coordinates, it is just simpler to assign coordinates and work with this. The coordinate independence of the vector allows us to understand how the components of the vector must transform. Let $p \in U_i \cap U_j$ and let $x = \varphi_i(p)$ and $y = \varphi_j(p)$. Then we have two expressions for the vector X in two different coordinate bases:

$$X = X^\mu \frac{\partial}{\partial x^\mu} = \tilde{X}^\mu \frac{\partial}{\partial y^\mu}, \quad (3.21)$$

therefore

$$\tilde{X}^\mu = X^\nu \frac{\partial y^\mu}{\partial x^\nu}. \quad (3.22)$$

Note that for two distinct points p and q the tangent spaces $T_p(M)$ and $T_q(M)$ are different. We cannot add vectors from one to a vector in the other. In fact even to compare the vectors in $T_p(M)$ with the vectors in $T_q(M)$ we need to introduce the notion of *parallel transport*.

3.2.2 One-forms

Since $T_p(M)$ is a vector space, there exists a dual vector space to $T_p(M)$ whose element is a linear function from $T_p(M) \rightarrow \mathbb{R}$. The dual space is called the *cotangent space* at p , and is denoted by $T_p^*(M)$. An element $\omega : T_p(M) \rightarrow \mathbb{R}$ of $T_p^*(M)$ is called a *dual vector/cotangent vector* or in the context of differential forms a *one-form*. The simplest example of a one-form is the differential df for a smooth function f on M . The action of a vector V on f is $V[f] = V^\mu \frac{\partial f}{\partial x^\mu} \in \mathbb{R}$. The action of $df \in T_p^*(M)$ on $V \in T_p(M)$ is defined by

$$\langle df, V \rangle \equiv V[f] = V^\mu \frac{\partial f}{\partial x^\mu} \in \mathbb{R}. \quad (3.23)$$

This is then \mathbb{R} -linear in both V and f . In terms of the coordinate basis we have

$$df = \frac{\partial f}{\partial x^\mu} dx^\mu, \quad (3.24)$$

and it is natural to regard $\{\mathrm{d}x^\mu\}$ as a basis of $T_p^*(M)$. This is a dual basis since

$$\left\langle \mathrm{d}x^\mu, \frac{\partial}{\partial x^\nu} \right\rangle = \frac{\partial x^\mu}{\partial x^\nu} = \delta_\nu^\mu. \quad (3.25)$$

We can then write an arbitrary one-form as

$$\omega = \omega_\mu \mathrm{d}x^\mu. \quad (3.26)$$

If we take a vector V and a one-form ω we may define the *inner product* $\langle \cdot, \cdot \rangle : T_p^*(M) \times T_p(M) \rightarrow \mathbb{R}$ to be

$$\langle \omega, V \rangle = \omega_\mu V^\nu \left\langle \mathrm{d}x^\mu, \frac{\partial}{\partial x^\nu} \right\rangle = \omega_\mu V^\nu \delta_\nu^\mu = \omega_\mu V^\mu. \quad (3.27)$$

The inner product is defined between a vector and a covector. Since ω is defined without reference to any coordinate system for a point $p \in U_i \cap U_j$ we have

$$\omega = \omega_\mu \mathrm{d}x^\mu = \tilde{\omega}_\mu \mathrm{d}y^\mu, \quad (3.28)$$

with x and y as before. Then we have

$$\tilde{\omega}_\nu = \omega_\mu \frac{\partial x^\mu}{\partial y^\nu}. \quad (3.29)$$

3.2.3 Tensors

We can now define a *tensor* of type (q, r) to be a multilinear object which maps q elements of $T_p^*(M)$ and r elements of $T_p(M)$ to \mathbb{R} . Let $\mathcal{T}_p^{(q,r)}(M)$ denote the set of (q, r) tensors at $p \in M$. An element of $\mathcal{T}^{(q,r)}(M)$ can be written in terms of the bases described above as

$$T = T^{\mu_1 \dots \mu_q}_{\nu_1 \dots \nu_r} \frac{\partial}{\partial x^{\mu_1}} \dots \frac{\partial}{\partial x^{\mu_q}} \mathrm{d}x^{\nu_1} \dots \mathrm{d}x^{\nu_r}. \quad (3.30)$$

T is a linear function

$$T : \otimes^q T_p^*(M) \otimes^r T_p(M) \rightarrow \mathbb{R}. \quad (3.31)$$

Let $V_i = V_i^\mu \frac{\partial}{\partial x^\mu}$ with $1 \leq i \leq r$ and $\omega_j = \omega_{j\mu} \mathrm{d}x^\mu$ with $1 \leq j \leq q$ then the action of T is

$$T(\omega_1, \dots, \omega_q; V_1, \dots, V_r) = T^{\mu_1 \dots \mu_q}_{\nu_1 \dots \nu_r} \omega_{1\mu_1} \dots \omega_{q\mu_q} V_1^{\mu_1} \dots V_r^{\mu_r}. \quad (3.32)$$

3.2.4 Tensor fields

So far we have defined vectors, one-forms and tensors at a particular point $p \in M$. We want to be able to smoothly assign such an object to every point of M . For a vector we call such an object a *vector field*. In other words if V is a vector field then for every $f \in \mathcal{F}(M)$ then $V[f] \in \mathcal{F}(M)$. We will denote the set of all vector fields on M as $\mathcal{X}(M)$. A vector field X at $p \in M$ is denoted by $X|_p$ which is an element of $T_p(M)$. Similarly we may define a *tensor field* of type (q, r) by a smooth assignment of an element of $\mathcal{T}_{r,p}^q(M)$ at each point $p \in M$. The set of tensor fields of type (q, r) on M is denoted by $\mathcal{T}_r^q(M)$.

3.2.5 Induced maps

A smooth map $f : M \rightarrow N$ naturally induces a map f_* called the *differential map* or *push-forward*,

$$f_* : T_p(M) \rightarrow T_{f(p)}(N). \quad (3.33)$$

The explicit form of f_* is obtained by the definition of the directional derivative along a

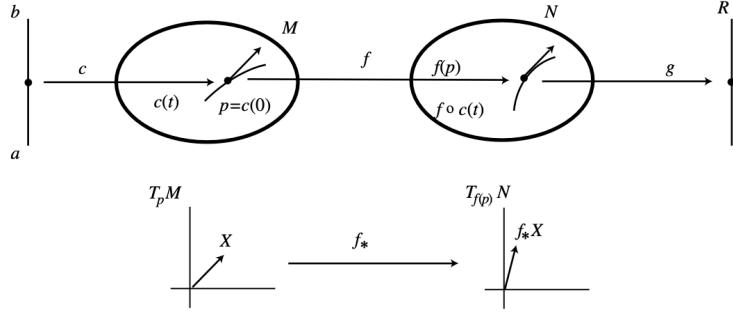


Figure 7: A map $f : M \rightarrow N$ induces the differential map $f_* : T_p(M) \rightarrow T_{f(p)}(N)$. Note that the mapping is performed by mapping the curve $c(t)$ between the two manifolds using the map f .

curve. Let $g \in \mathcal{F}(N)$ then $g \circ f \in \mathcal{F}(M)$. A vector $V \in T_p(M)$ acts on $g \circ f$ to give a number $V[g \circ f]$. We can now define $f_* V \in T_{f(p)}(N)$ by

$$(f_* V)[g] \equiv V[g \circ f]. \quad (3.34)$$

We can write this more explicitly by introducing coordinates. Let us introduce the charts (U, φ) on M and (V, ψ) on N , then

$$(f_* V)[g \circ \psi^{-1}(y)] = V[g \circ f \circ \varphi^{-1}(x)], \quad (3.35)$$

where $x = \varphi(p)$ and $y = \psi(f(p))$. Let $V = V^\mu \frac{\partial}{\partial x^\mu}$ and $f_* V = W^\alpha \frac{\partial}{\partial y^\alpha}$, then in components it reads

$$W^\alpha \frac{\partial}{\partial y^\alpha} [g \circ \psi^{-1}(y)] = V^\mu \frac{\partial}{\partial x^\mu} [g \circ f \circ \varphi^{-1}(x)]. \quad (3.36)$$

If we take the function $g = y^\alpha$, i.e. we map the point in N to the α 'th component of the coordinate, then we find

$$W^\alpha = V^\mu \frac{\partial}{\partial x^\mu} y^\alpha(x). \quad (3.37)$$

This is nothing but the Jacobian of the map $f : M \rightarrow N$. This can be extended to tensors of type $(q, 0)$.

Example Let (x^1, x^2) and (y^1, y^2, y^3) be coordinates on M and N respectively, and let $V = a\frac{\partial}{\partial x^1} + b\frac{\partial}{\partial x^2}$. Take the map $f : M \rightarrow N$ whose coordinate representation is

$$y = (x^1, x^2, \sqrt{1 - (x^1)^2 - (x^2)^2}). \quad (3.38)$$

Then

$$f_* V = V^\mu \frac{\partial y^\alpha}{\partial x^\mu} \frac{\partial}{\partial y^\alpha} = a \frac{\partial}{\partial y^1} + b \frac{\partial}{\partial y^2} - \left(a \frac{y^1}{y^3} + b \frac{y^2}{y^3} \right) \frac{\partial}{\partial y^3}. \quad (3.39)$$

A map f also induces a map between cotangent space

$$f^* : T_{f(p)}^*(N) \rightarrow T_p^*(M), \quad (3.40)$$

which is called the *pull-back*. If we take $V \in T_p(M)$ and $\omega \in T_{f(p)}^*(N)$ then the pull-back of ω by f^* is defined to be

$$\langle f^*\omega, V \rangle = \langle \omega, f_* V \rangle. \quad (3.41)$$

In components we have

$$(f^*\omega)_\mu = \omega_\alpha \frac{\partial y^\alpha}{\partial x^\mu}. \quad (3.42)$$

The pull-back can be extended to tensors of type $(0, r)$.

3.3 Flows and Lie derivatives

Let X be a vector field on M . An *integral curve* $x(t)$ of X is a curve in M whose tangent vector at $x(t)$ is $X|_x$. Given a chart (U, φ) , this means that

$$\frac{dx^\mu(t)}{dt} = X^\mu(x(t)), \quad (3.43)$$

where $x^\mu(t)$ is the μ 'th component of $\varphi(x(t))$ and $X = X^\mu \frac{\partial}{\partial x^\mu}$. As always we have very much abused notation, using x to denote a point in M as well as its coordinates. Finding an Integral curve is equivalent to solving the ODE with initial conditions $x^\mu(0) = x^\mu$. The existence and uniqueness theorems for ODEs implies that there is always a unique solution, at least locally, with the given initial data.

Let $\sigma(t, x_0)$ be an integral curve of X which passes through the point x_0 at $t = 0$, and denote the coordinate by $\sigma^\mu(t, x_0)$. The flow equation becomes

$$\frac{d}{dt} \sigma^\mu(t, x_0) = X^\mu(\sigma(t, x_0)), \quad (3.44)$$

with the initial condition

$$\sigma^\mu(0, x_0) = x_0^\mu. \quad (3.45)$$

The map $\sigma : \mathbb{R} \times M \rightarrow M$ is called a *flow* generated by $X \in \mathcal{X}(M)$. A flow satisfies the rule

$$\sigma(t, \sigma^\mu(s, x_0)) = \sigma(t + s, x_0), \quad (3.46)$$

for any $s, t \in \mathbb{R}$, such that both sides make sense. This follows from the uniqueness of the ODE with fixed initial condition.

Theorem For any point $x \in M$, there exists a differentiable map $\sigma : \mathbb{R} \times M \rightarrow M$ such that

- (i) $\sigma(0, x) = x$,
- (ii) $t \mapsto \sigma(t, x)$ is a solution of (3.44) and (3.45),
- (iii) $\sigma(t, \sigma^\mu(s, x)) = \sigma(t + s, x)$

note that the initial point is denoted by x to emphasise that σ is a map $\mathbb{R} \times M \rightarrow M$.

We may imagine a flow as a steady stream flow. If a particle is observed at a point x at $t = 0$ it will be found at $\sigma(t, x)$ at later time t .

Example Let $M = \mathbb{R}^2$ and let $X((x, y)) = -y \frac{\partial}{\partial x} + x \frac{\partial}{\partial y}$ be a vector field in M . Then

$$\sigma(t, (x, y)) = (x \cos t - y \sin t, x \sin t + y \cos t), \quad (3.47)$$

is a flow generated by X . The flow through (x, y) is a circle whose centre is at the origin. Clearly $\sigma(t, (x, y)) = (x, y)$ if $t = 2\pi n, n \in \mathbb{Z}$. If $(x, y) = (0, 0)$, the flow stays at $(0, 0)$.

3.3.1 One-parameter group of transformations

For fixed $t \in \mathbb{R}$ a flow $\sigma(t, x)$ is a diffeomorphism from M to M which we denote by $\sigma_t : M \rightarrow M$. This map is made into a commutative group by the following rules (*Exercise: Check this*):

1. $\sigma_t(\sigma_s(x)) = \sigma_{t+s}(x)$ i.e. $\sigma_t \circ \sigma_s = \sigma_{t+s}$,
2. σ_0 = identity map (unit element),
3. $\sigma_{-t} = (\sigma_t)^{-1}$.

This group is called the *one-parameter group of transformations*. Locally the group looks like the additive group \mathbb{R} , although they may not be isomorphic globally. For example in

the example above (see equation (3.47)) we had that $\sigma_{2\pi n+t} = \sigma_t$ and we find that the one-parameter group is isomorphic to $\text{SO}(2)$ the multiplicative group of 2×2 real matrices of the form;

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \quad (3.48)$$

or $\text{U}(1)$ the multiplicative group of complex numbers of unit modulus $e^{i\theta}$.

We can consider an infinitesimal transformation and see where it maps the point x . Using (3.44) and (3.45) we find

$$\sigma_\epsilon^\mu(x) = \sigma^\mu(\epsilon, x) = x^\mu + \epsilon X^\mu(x). \quad (3.49)$$

The vector field X in this context is called the *infinitesimal generator* of the transformation σ_t .

Given a vector field X the corresponding flow σ is often referred to as the *exponentiation* of X and is denoted by

$$\sigma^\mu(t, x) = \exp(tX)x^\mu. \quad (3.50)$$

To see why this is so, let us take a parameter t and evaluate the coordinate of a point which is separated from the initial point $x = \sigma(0, x)$ by the parameter distance t along the flow σ . The coordinate corresponding to the point $\sigma(t, x)$ is

$$\begin{aligned} \sigma^\mu(t, x) &= x^\mu + t \frac{d}{ds} \sigma^\mu(s, x) \Big|_{s=0} + \frac{t^2}{2!} \left(\frac{d}{ds} \right)^2 \sigma^\mu(s, x) \Big|_{s=0} + \dots \\ &= \left[1 + t \frac{d}{ds} + \frac{t^2}{2!} \left(\frac{d}{ds} \right)^2 + \dots \right] \sigma^\mu(s, x) \Big|_{s=0} \\ &\equiv \exp \left(t \frac{d}{ds} \right) \sigma^\mu(s, x) \Big|_{s=0}. \end{aligned} \quad (3.51)$$

The last expression can also be written as $\sigma^\mu(t, x) = \exp(tX)x^\mu$ as in the definition above. Then the flow satisfies the following exponential properties

$$\begin{aligned} \sigma(0, x) &= x = \exp(0X)x, \\ \frac{\sigma(t, x)}{dt} &= X \exp(tX)x = \frac{d}{dt} \left(\exp(tX)x \right), \\ \sigma(t, \sigma(s, x)) &= \sigma(t, \exp(sX)x) = \exp(tX) \exp(sX)x = \exp((t+s)X)x = \sigma(t+s, x). \end{aligned} \quad (3.52)$$

3.3.2 Lie Derivatives

Let $\sigma(t, x)$ and $\tau(t, x)$ be two flows generated by the vector fields X and Y respectively:

$$\frac{d\sigma^\mu(s, x)}{ds} = X^\mu(\sigma(s, x)), \quad \frac{d\tau^\mu(t, x)}{dt} = Y^\mu(\tau(t, x)). \quad (3.53)$$

Let us evaluate the change of the vector field Y along $\sigma(s, x)$. To do this we need to compare the vector Y at a point x with Y at a nearby point $x' = \sigma_\epsilon(x)$, see figure 8. We cannot simply take the difference between the components of Y at the two points since they belong to different tangent spaces: $T_x(M)$ and $T_{\sigma_\epsilon(x)}(M)$, and so the difference between the two vectors is ill-defined. To define a sensible derivative, we first map $Y|_{\sigma_\epsilon(x)}$ to $T_x(M)$ by using the push-forward $(\sigma_{-\epsilon})_* : T_{\sigma_\epsilon(x)}(M) \rightarrow T_x(M)$, after which the two vectors are in the same tangent space and we can take the difference between them. The *Lie derivative* of a vector field Y along the flow σ of the vector field X is defined by

$$\mathcal{L}_X Y = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left[(\sigma_{-\epsilon})_* Y|_{\sigma_\epsilon(x)} - Y|_x \right]. \quad (3.54)$$

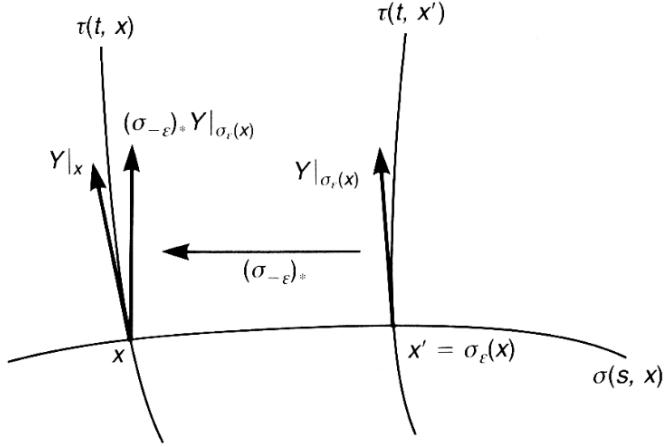


Figure 8: To compare a vector $Y|_x$ with the vector $Y|_{\sigma_\epsilon(x)}$ the latter must be transported back to x by the differential map $(\sigma_{-\epsilon})_*$, that is we use the push-forward.

By writing this in components we may obtain another expression for the Lie derivative of a vector field. Let (U, φ) be a chart with the coordinates x and let $X = X^\mu \frac{\partial}{\partial x^\mu}$ and $Y = Y^\mu \frac{\partial}{\partial x^\mu}$ be vector fields defined on U . Then $\sigma_\epsilon(x)$ has the coordinates $x^\mu + \epsilon X^\mu(x)$ and

$$\begin{aligned} Y|_{\sigma_\epsilon(x)} &= Y^\mu(x^\nu + \epsilon X^\nu(x)) e_\mu|_{x+\epsilon X} \\ &\simeq \left[Y^\mu(x) + \epsilon X^\nu(x) \partial_\nu Y^\mu(x) \right] e_\mu|_{x+\epsilon X}, \end{aligned} \quad (3.55)$$

with $e_\mu = \frac{\partial}{\partial x^\mu} \equiv \partial_\mu$. Mapping this vector at $\sigma_\epsilon(x)$ to x using $(\sigma_{-\epsilon}(x))_*$ we obtain

$$(\sigma_{-\epsilon}(x))_* Y|_{\sigma_\epsilon(x)} = \left[Y^\mu(x) + \epsilon X^\lambda(x) \partial_\lambda Y^\mu(x) \right] \partial_\mu(x^\nu - \epsilon X^\nu(x)) e_\nu|_x$$

$$\begin{aligned}
&= \left[Y^\mu(x) + \epsilon X^\lambda(x) \partial_\lambda Y^\mu(x) \right] \left[\delta_\mu^\nu - \epsilon \partial_\mu X^\nu(x) \right] e_\nu|_x \\
&= Y^\mu(x) e_\mu|_x + \epsilon \left[X^\mu(x) \partial_\mu Y^\nu(x) - Y^\mu(x) \partial_\mu X^\nu(x) \right] e_\nu|_x + O(\epsilon^2),
\end{aligned} \tag{3.56}$$

and therefore we find

$$\mathcal{L}_X Y = (X^\mu \partial_\mu Y^\nu - Y^\mu \partial_\mu X^\nu) e_\nu. \tag{3.57}$$

This motivates the introduction of the *Lie bracket*, [,]. For vector fields X, Y on M we have

$$[X, Y]f = X[Y[f]] - Y[X[f]], \tag{3.58}$$

for all $f \in \mathcal{F}(M)$. In components $[X, Y]$ reads

$$(X^\mu \partial_\mu Y^\nu - Y^\mu \partial_\mu X^\nu) e_\nu. \tag{3.59}$$

Then the Lie derivative of Y along X is

$$\mathcal{L}_X Y = [X, Y]. \tag{3.60}$$

Exercise: Show that the Lie bracket does define a vector field. In addition show that it satisfies the following properties:

1. Bilinearity

$$\begin{aligned}
[X, c_1 Y_1 + c_2 Y_2] &= c_1 [X, Y_1] + c_2 [X, Y_2], \\
[c_1 X_1 + c_2 X_2, Y] &= c_1 [X_1, Y] + c_2 [X_2, Y],
\end{aligned} \tag{3.61}$$

for any constants c_1 and c_2 .

2. Skew symmetry

$$[X, Y] = -[Y, X]. \tag{3.62}$$

3. Jacobi Identity

$$[[X, Y], Z] + [[Z, X], Y] + [[Y, Z], X] = 0, \tag{3.63}$$

4. For X, Y vector fields and f a smooth function on M then

$$\begin{aligned}
\mathcal{L}_{fX} Y &= f[X, Y] - Y[f]X, \\
\mathcal{L}_X(fY) &= f[X, Y] + X[f]Y
\end{aligned} \tag{3.64}$$

5. For $f : M \rightarrow N$ then

$$f_*[X, Y] = [f_*X, f_*Y]. \tag{3.65}$$

Geometrically the Lie bracket shows the non-commutativity of two flows. Let us take the flows $\sigma(s, x)$ and $\tau(t, x)$ generated by X and Y respectively. If we first move a small parameter distance ϵ along the flow σ and then by δ along the second flow τ we end up at a point whose coordinates are

$$\begin{aligned}\tau^\mu(\delta, \sigma(\epsilon, x)) &\simeq \tau^\mu(\delta + x^\nu + \epsilon X^\nu(x)) \\ &\simeq x^\mu + \epsilon X^\mu(x) + \delta Y^\mu(x^\nu + \epsilon X^\nu) \\ &\simeq x^\mu + \epsilon X^\mu(x) + \delta Y^\mu(x) + \epsilon \delta X^\nu(x) \partial_\nu Y^\mu(x).\end{aligned}\tag{3.66}$$

If we instead first move along τ and then move along σ we find

$$\sigma^\mu(\epsilon, \tau(\delta, x)) \simeq x^\mu + \delta Y^\mu(x) + \epsilon X^\mu(x) + \epsilon \delta Y^\nu(x) \partial_\nu X^\mu(x).\tag{3.67}$$

The difference between the two points is proportional to the Lie bracket

$$\tau^\mu(\delta, \sigma(\epsilon, x)) - \sigma^\mu(\epsilon, \tau(\delta, x)) = \epsilon \delta [X, Y]^\mu.\tag{3.68}$$

The Lie bracket measures the failure of the parallelogram in figure 9 to close. It is easy to see that

$$\mathcal{L}_X Y = [X, Y] = 0 \iff \sigma(s, \tau(t, x)) = \tau(t, \sigma(s, x)).\tag{3.69}$$

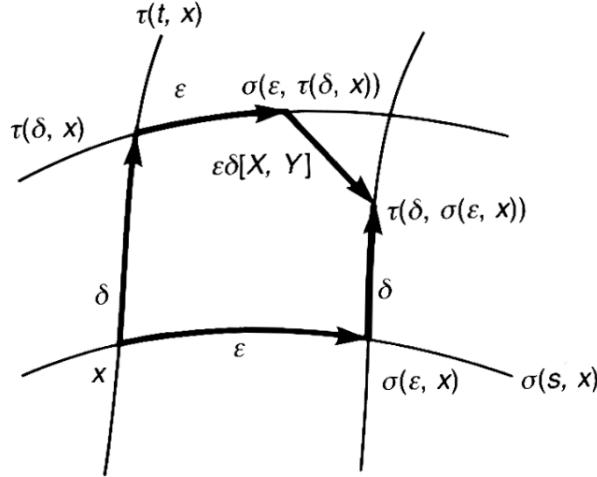


Figure 9: Moving first along the flow σ and then the flow τ or first along τ and then along σ we find that we may not end up at the same point. The difference is measured by the failure of the Lie bracket to vanish.

We may also define the Lie derivative of a one-form $\omega \in \Omega^1(M)$ along X . This time we need to use the pull-back, then the Lie derivative of the one-form ω is

$$\mathcal{L}_X\omega \equiv \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left[(\sigma_\epsilon)^* \omega|_{\sigma_\epsilon(x)} - \omega|_x \right], \quad (3.70)$$

where $\omega|_x \in T_x^*(M)$ is ω at x . Introducing coordinates such that $\omega = \omega_\mu dx^\mu$, then we have

$$(\sigma_\epsilon)^* \omega|_{\sigma_\epsilon(x)} = \omega_\mu(x) dx^\mu + \epsilon [X^\nu(x) \partial_\nu \omega_\mu(x) + \partial_\mu X^\nu(x) \omega_\nu(x)] dx^\mu, \quad (3.71)$$

which leads to

$$\mathcal{L}_X\omega = (X^\nu \partial_\nu \omega_\mu + \partial_\mu X^\nu \omega_\nu) dx^\mu. \quad (3.72)$$

This remains a one-form, that is $\mathcal{L}_X\omega \in T_x^*(M)$ since it is the difference of two one-forms at the same point.

This may also be extended to functions f on M . Then

$$\begin{aligned} \mathcal{L}_X f &\equiv \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [f(\sigma_\epsilon(x)) - f(x)] \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [f(x^\mu + \epsilon X^\mu(x)) - f(x^\mu)] \\ &= X^\mu(x) \frac{\partial f}{\partial x^\mu} = X[f], \end{aligned} \quad (3.73)$$

which is just the usual directional derivative of f along X .

To extend this to more general tensors we need the following proposition:

Proposition: The Lie derivative satisfies

$$\mathcal{L}_X(t_1 + t_2) = \mathcal{L}_X t_1 + \mathcal{L}_X t_2, \quad (3.74)$$

where t_1 and t_2 are tensor fields of the same type. Moreover

$$\mathcal{L}_X(t_1 \otimes t_2) = (\mathcal{L}_X t_1) \otimes t_2 + t_1 \otimes (\mathcal{L}_X t_2), \quad (3.75)$$

with t_1 and t_2 tensors of arbitrary type this time.

3.4 Differential forms

Not all tensors are created equally, some will play a more prominent role than others. One class of interesting tensors are the p -forms. To define them we must first introduce some additional notation. The symmetry operation on a tensor $\omega \in \mathcal{T}_p^{(0,r)}(M)$ is defined by

$$P\omega(V_1, \dots, V_r) \equiv \omega(V_{P(1)}, \dots, V_{P(r)}), \quad (3.76)$$

with the $V_i \in T_p(M)$, and P an element of the symmetric group of order r . Let us take the coordinate basis, $\{e_\mu = \frac{\partial}{\partial x^\mu}\}$, then the components of ω in this basis are

$$\omega(e_{\mu_1}, \dots, e_{\mu_r}) = \omega_{\mu_1 \dots \mu_r}. \quad (3.77)$$

It follows that the components of $P\omega$ are

$$P\omega(e_{\mu_1}, \dots, e_{\mu_r}) = \omega_{\mu_{P(1)} \dots \mu_{P(r)}}. \quad (3.78)$$

For a general tensor of type (q, r) the symmetry operations are defined for the q and r indices separately. For $\omega \in \mathcal{T}_p^{(0,r)}(M)$ the *symmetrizer* S is defined by

$$S\omega = \frac{1}{r!} \sum_{P \in S_r} P\omega, \quad (3.79)$$

while the anti-symmetrizer A is defined to be

$$A\omega = \frac{1}{r!} \sum_{P \in S_r} \text{sgn}(P)P\omega, \quad (3.80)$$

with $\text{sgn}(P) = +1$ for an even permutation and -1 for an odd permutation. $S\omega$ is totally symmetric so that $PS\omega = S\omega$ for any $P \in S_r$ while $A\omega$ is totally anti-symmetric so that $A\omega = \text{sgn}(P)A\omega$.

Definition A *differential form* of order r , or more succinctly an r -form, is a totally anti-symmetric tensor of type $(0, r)$.

The *Wedge product* \wedge of r one-forms is defined to be the totally anti-symmetric tensor product of the one-forms

$$dx^{\mu_1} \wedge dx^{\mu_2} \wedge \dots \wedge dx^{\mu_r} \equiv \sum_{P \in S_r} \text{sgn}(P) dx^{\mu_{P(1)}} \otimes dx^{\mu_{P(2)}} \otimes \dots \otimes dx^{\mu_{P(r)}}. \quad (3.81)$$

Thus

$$dx^\mu \wedge dx^\nu = dx^\mu \otimes dx^\nu - dx^\nu \otimes dx^\mu. \quad (3.82)$$

The wedge product satisfies the following conditions

- $dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r} = 0$ if some index is repeated.
- $dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r} = \text{sgn}(P) dx^{\mu_{P(1)}} \wedge \dots \wedge dx^{\mu_{P(r)}}.$
- $dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r}$ is linear in each dx^μ .

We will denote the vector space of r -forms at the point $p \in M$ by $\Omega_p^r(M)$, a basis is provided by the set of all wedge products in (3.81). We can then expand an element of $\Omega_p^r(M)$ as

$$\omega = \frac{1}{r!} \omega_{\mu_1 \dots \mu_r} dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r}, \quad (3.83)$$

where $\omega_{\mu_1 \dots \mu_r}$ are taken to be totally anti-symmetric. Since there are $\binom{m}{r}$ choices of the set $\{\mu_1, \dots, \mu_r\}$ out of $(1, 2, \dots, m)$ the dimension of the vector space $\Omega_p^{(r)}(M)$ is

$$\binom{m}{r} = \frac{m!}{r!(m-r)!}. \quad (3.84)$$

We take $\Omega_p^0(M) = \mathbb{R}$ and it is obvious that $\Omega_p^1(M) = T_p^*(M)$. Also since we are anti-symmetrising it follows that if r exceeds $m = \dim(M)$ then it vanishes identically. Moreover since $\binom{m}{r} = \binom{m}{m-r}$ it follows that $\dim \Omega_p^r(M) = \dim \Omega_p^{m-r}(M)$. Since $\Omega_p^r(M)$ is a vector space it is isomorphic to $\Omega_p^{(r-m)}(M)$.¹⁴

3.4.1 Exterior product

We may define the *exterior product* to be the map $\wedge : \Omega_p^q(M) \times \Omega_p^r(M) \rightarrow \Omega_p^{q+r}(M)$. Its action follows by trivial extension of the wedge product defined above. Let $\omega \in \Omega_p^q(M)$ and $\xi \in \Omega_p^r(M)$ be an q -form and and r -form respectively. The action of the $(q+r)$ -form $\omega \wedge \xi$ on $q+r$ vectors V_i is

$$(\omega \wedge \xi)(V_1, \dots, V_{q+r}) = \frac{1}{q!r!} \sum_{P \in S_{q+r}} \text{sgn}(P) \omega(V_{P(1)}, \dots, V_{P(q)}) \xi(V_{P(q+1)}, \dots, V_{P(q+r)}). \quad (3.85)$$

It follows that if $q+r > m$ then $\omega \wedge \xi$ vanishes. With this product we can define and algebra

$$\Omega_p^*(M) \equiv \Omega_p^0(M) \oplus \Omega_p^1(M) \oplus \dots \oplus \Omega_p^m(M). \quad (3.86)$$

Exercise: From the properties of the wedge product show that for $\xi \in \Omega_p^q(M)$, $\eta \in \Omega_p^r(M)$ and $\omega \in \Omega_p^s(M)$ that

$$\begin{aligned} \xi \wedge \xi &= 0 && \text{if } q \text{ odd,} \\ \xi \wedge \eta &= (-1)^{qr} \eta \wedge \xi, \\ (\xi \wedge \eta) \wedge \omega &= \xi \wedge (\eta \wedge \omega). \end{aligned} \quad (3.87)$$

We may assign an r -form smoothly at each point on a manifold M . We denote the space of smooth r -forms on M by $\Omega^r(M)$, and take $\Omega^0(M) = \mathcal{F}(M)$ to be the space of smooth functions.

¹⁴When the manifold is equipped with a metric the isomorphism is provided by the Hodge star operation \star .

3.4.2 Exterior derivative

Definition The *exterior derivative* d_r is a map $\Omega^r(M) \rightarrow \Omega^{r+1}(M)$, whose action on an r -form

$$\omega = \frac{1}{r!} \omega_{\mu_1 \dots \mu_r} dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r}, \quad (3.88)$$

is

$$d_r \omega = \frac{1}{r!} \left(\frac{\partial}{\partial x^\nu} \omega_{\mu_1 \dots \mu_r} \right) dx^\nu \wedge dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r}. \quad (3.89)$$

It is common to drop the r subscript and simply write d . The wedge product automatically anti-symmetrises the coefficient so it is indeed a $(r+1)$ -form that we obtain. It follows that for $\xi \in \Omega_p^q(M)$, $\eta \in \Omega_p^r(M)$ we have

$$d(\xi \wedge \eta) = d\xi \wedge \eta + (-1)^q \xi \wedge d\eta. \quad (3.90)$$

Example: Let us take \mathbb{R}^3 with coordinates (x, y, z) . The generic r -forms are

$$\begin{aligned} \omega_0 &= f(x, y, z), \\ \omega_1 &= \omega_x(x, y, z)dx + \omega_y(x, y, z)dy + \omega_z(x, y, z)dz, \\ \omega_2 &= \omega_{xy}(x, y, z)dx \wedge dy + \omega_{yz}(x, y, z)dy \wedge dz + \omega_{zx}(x, y, z)dz \wedge dx, \\ \omega_3 &= \omega_{xyz}(x, y, z)dx \wedge dy \wedge dz. \end{aligned} \quad (3.91)$$

The exterior derivative of these forms is

$$\begin{aligned} d\omega_0 &= \frac{\partial}{\partial x} f(x, y, z)dx + \frac{\partial}{\partial y} f(x, y, z)dy + \frac{\partial}{\partial z} f(x, y, z)dz, \\ d\omega_1 &= \left(\frac{\partial}{\partial x} \omega_y - \frac{\partial}{\partial y} \omega_x \right) dx \wedge dy + \left(\frac{\partial}{\partial y} \omega_z - \frac{\partial}{\partial z} \omega_y \right) dy \wedge dz + \left(\frac{\partial}{\partial z} \omega_x - \frac{\partial}{\partial x} \omega_z \right) dz \wedge dx, \\ d\omega_2 &= \left(\frac{\partial}{\partial x} \omega_{yz} + \frac{\partial}{\partial y} \omega_{zx} + \frac{\partial}{\partial z} \omega_{xy} \right) dx \wedge dy \wedge dz, \\ d\omega_3 &= 0. \end{aligned} \quad (3.92)$$

In the usual 3d vector calculus you may identify these as ‘grad’ for d acting on the scalar, ‘curl’ for the one-form and the ‘divergence’ for the two-form.

We have used coordinates to give the definition of the exterior derivative above, we may also write it in coordinate free notation. For an r -form, $\omega \in \Omega^r(M)$ we have

$$\begin{aligned} d\omega(X_1, \dots, X_{r+1}) &= \sum_{i=1}^r (-1)^{i+1} X_i \omega(X_1, \dots, \hat{X}_i, \dots, X_{r+1}) \\ &\quad + \sum_{i < j} (-1)^{i+j} \omega([X_i, X_j], X_1, \dots, \hat{X}_i, \dots, \hat{X}_j, \dots, X_{r+1}), \end{aligned} \quad (3.93)$$

where the hats denote that this term should be removed.

From either the coordinate free expression (3.93) or the one using the coordinates in (3.89), we can prove the important result that

$$d^2 = 0, \quad (d_{r+1}dr = 0). \quad (3.94)$$

Using the coordinate form (3.89) we find

$$d^2\omega = \frac{1}{r!} \frac{\partial^2}{\partial x^\nu \partial x^\sigma} \omega_{\mu_1 \dots \mu_r} dx^\nu \wedge dx^\sigma \wedge dx^{\mu_1} \wedge \dots \wedge dx^{\mu_r}. \quad (3.95)$$

Using that the derivative term is symmetric in $\nu\sigma$ while the wedge product is anti-symmetric in these indices it follows that this vanishes.

It then follows that an exact form is always closed, though the converse need not be true. The failure of a closed form to be exact tells us interesting information about the topology of the underlying manifold. The exterior derivative induces the sequence

$$0 \xrightarrow{i} \Omega^0(M) \xrightarrow{d_0} \Omega^1(M) \xrightarrow{d_1} \dots \xrightarrow{d_{m-1}} \Omega^m(M) \xrightarrow{d_m} 0, \quad (3.96)$$

with i the inclusion map. This is known as the *de Rahm complex*. If we let the set of all closed r -forms on M be $Z^r(M)$, so that for $d_r : \Omega^r(M) \rightarrow \Omega^{r+1}(M)$, $\ker(d_r) = Z^r(M)$. Moreover, let us also denote the set of all exact r -forms to be $B^r(M)$, i.e. the B^r is the image of $\Omega^{r-1}(M)$ under $d^{r-1} : \Omega^{r-1}(M) \rightarrow \Omega^r(M)$. Then the r th de-Rahm cohomology group is defined to be

$$H^r(M) = Z^r(M)/B^r(M). \quad (3.97)$$

This is the dual space of the *homology group*, though we will not have time to consider this. The cohomology groups tell us important information about a manifold, the dimension of them are topological invariants. Let $b^r = \dim(H^r(M))$, these are the *Betti numbers* of the manifold and are always finite. For a connected manifold one always has $b_0 = 1$, these are just the constant functions. The higher Betti numbers are non-zero when the manifold has some interesting topology. The Euler characteristic of a manifold is

$$\chi(M) = \sum_{r=0}^m (-1)^r b^r(M). \quad (3.98)$$

We have seen that every exact form is closed, however not every closed form is exact, instead we have:

Theorem Poincaré's lemma If a coordinate neighbourhood U of a manifold M is contractible

to a point $p \in M$, any closed form on U is also exact. In particular on $M = \mathbb{R}^m$, closed implies exact.

Since we have been mapping manifold to \mathbb{R}^m this says that for a general manifold any closed form is locally exact. That is if ω is a closed r -form, then in any neighbourhood $U \subset M$ it is always possible to find $\eta \in \Omega^{r-1}(M)$ such that $\omega = d\eta$ on U . Since we cannot generally cover the manifold with a single patch, it may not be possible to find such an η everywhere on M , hence we say that the form is only *locally* exact.

Let us consider some examples. First consider $M = \mathbb{R}$. We can take a one-form $\omega = f(x)dx$. This is trivially closed since it is a top form, it is also exact since we can write

$$g(x) = \int_0^x dx' f(x') , \quad (3.99)$$

such that $\omega = dg(x)$.

Now consider a circle, S^1 . We can view this as the phase $e^{i\theta} \in \mathbb{C}$ and can introduce the one-form $\omega = d\theta$. Clearly this is once again closed since it is a top form. By the way we have written this it makes it seem that this is once again an exact form however the caveat is that θ is *not* a good coordinate everywhere on S^1 , it is not single valued. As such θ is not a good smooth function and so is not a zero-form. Hence it is closed but not exact.

Next consider $M = \mathbb{R}^2$. The Poincaré lemma ensures that all closed forms are exact. This changes if we remove a point, consider $\mathbb{R}^2 - \{0,0\}$ and the one-form

$$\omega = -\frac{y}{x^2 + y^2} dx + \frac{x}{x^2 + y^2} dy . \quad (3.100)$$

This is not a smooth one-form on \mathbb{R}^2 , the problematic point is the origin. However, on $\mathbb{R}^2 - \{0,0\}$ it is smooth since we no longer have the problem area at the origin. It is not difficult to see that ω is closed, but is it exact. If such a smooth function exists such that $\omega = df$ then the function f must satisfy:

$$\frac{\partial f}{\partial x} = -\frac{y}{x^2 + y^2} , \quad \frac{\partial f}{\partial y} = \frac{x}{x^2 + y^2} . \quad (3.101)$$

The solution is

$$f(x,y) = \arctan\left(\frac{y}{x}\right) + \text{constant} , \quad (3.102)$$

so have we found an exact form. The answer is no, this is not a smooth function everywhere on $\mathbb{R}^2 - \{0,0\}$, and so ω is *not* exact. Removing a point makes a big difference: closed no longer implies exact! A similar story holds for \mathbb{R}^3 and this is how magnetic monopoles sneak back into physics despite being forbidden by Maxwell's equations.

3.4.3 Interior product

We can now go from $\Omega^r(M) \rightarrow \Omega^{r+1}(M)$, what about the other way around? To do this we have to define the *Interior product*. Let X be a vector field and $\omega \in \Omega^r(M)$ then

$$i_X \omega(X_1, \dots, X_{r-1}) \equiv \omega(X, X_1, \dots, X_{r-1}). \quad (3.103)$$

If we introduce coordinates: $X = X^\mu \frac{\partial}{\partial x^\mu}$ then

$$i_X \omega = \frac{1}{(r-1)!} X^\nu \omega_{\nu \mu_1 \dots \mu_{r-1}} dx^{\mu_2} \wedge \dots \wedge dx^{\mu_r}. \quad (3.104)$$

Example: Let us take \mathbb{R}^3 again with coordinates (x, y, z) , and the usual coordinate basis, then we have

$$i_{e_x}(dx \wedge dy) = dy, \quad i_{e_x}(dy \wedge dz) = 0, \quad i_{e_x}(dz \wedge dx) = -dz. \quad (3.105)$$

The interior product and exterior derivative combine beautifully to give a simple way of computing the Lie derivative of a form along the vector X :

$$\mathcal{L}_X \omega = (d i_X + i_X d)\omega, \quad (3.106)$$

this applies for any r -form.

The interior product satisfies (*Exercise:* show this)

$$\begin{aligned} i_X^2 &= 0, \\ i_X(\omega \wedge \eta) &= i_X \omega \wedge \eta + (-1)^r \omega \wedge i_X \eta, \\ i_{[X,Y]} \omega &= X(i_Y \omega) - Y(i_X \omega), \\ \mathcal{L}_X i_X \omega &= i_X \mathcal{L}_X \omega. \end{aligned} \quad (3.107)$$

Hamiltonian mechanics in differential geometry We can now combine some of the differential geometry we have learnt so far to reformulate classical Hamiltonian mechanics. Recall that in classical mechanics the phase space is a manifold M parametrised by coordinates (q^i, p_j) where q^i are the positions of particles and p_j the momenta. Note that M is even dimensional. The Hamiltonian $H(q, p)$ is a function on M and Hamilton's equations are

$$\dot{q}^i = \frac{\partial H}{\partial p_i}, \quad \text{and} \quad \dot{p}_i = -\frac{\partial H}{\partial q^i}. \quad (3.108)$$

Phase space comes equipped with the *Poisson bracket*, defined on a pair of functions f, g to act as

$$\{f, g\} = \frac{\partial f}{\partial q^j} \frac{\partial g}{\partial p_j} - \frac{\partial f}{\partial p_j} \frac{\partial g}{\partial q^j}. \quad (3.109)$$

The time evolution of a function is

$$\dot{f} = \{f, H\}, \quad (3.110)$$

with H the Hamiltonian. For $f = q^i$ and $f = p_i$ one obtains Hamilton's equations.

Underlying this structure are forms. The key idea behind this is to convert the scalar function H into a vector field X_H on M . Particles will then follow trajectories which are the integral curves generated by X_H . To convert the scalar into a vector we introduce the *symplectic two-form* ω . This is a two-form which is closed $d\omega = 0$ and is non-degenerate, $\omega \wedge \omega \wedge \dots \wedge \omega \neq 0$. A manifold equipped with such a two-form is called a *symplectic manifold*.

Any two-form provides a map $\omega : T_p(M) \rightarrow T_p^*(M)$, since given a vector field X we can simply take the inner product with ω to obtain a one-form, $i_X \omega$. For our purposes we want to go in the opposite direction, we want to convert a scalar function into a vector field. This is possible if the map $\omega : T_p(M) \rightarrow T_p^*(M)$ is an isomorphism. This is equivalent to ω being non-degenerate. In this case we can define a vector field X_H via

$$i_{X_H} \omega = -dH. \quad (3.111)$$

In coordinate notation we have

$$X_H^\mu \omega_{\mu\nu} = -\partial_\nu H. \quad (3.112)$$

If we take the inverse to be $\omega^{\mu\nu}$ so that $\omega^{\mu\nu} \omega_{\nu\rho} = \delta_\rho^\mu$, then

$$X_H^\mu = \omega^{\mu\nu} \partial_\nu H. \quad (3.113)$$

The integral curves generated by X_H obey

$$\frac{dx^\mu(t)}{dt} = X_H^\mu = \omega^{\mu\nu} \partial_\nu H. \quad (3.114)$$

These are the general form of Hamilton's equations, just written without reference to canonical coordinates. If we let $x^\mu = (q^i, p_j)$ and choose the symplectic form to have block diagonal form

$$\omega^{\mu\nu} = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \quad \Leftrightarrow \quad \omega = dp^i \wedge dq_i \quad (3.115)$$

then the integral curves reduce precisely to Hamilton's equations (3.108).

To define the Poisson structure, we first note that we can repeat the map for obtaining a vector from a scalar for any function f , to obtain a vector field X_f . Then

$$\{f, g\} = \omega(X_f, X_g) = -\omega(X_g, X_f). \quad (3.116)$$

This may be written in a multitude of different ways, we have

$$\{f, g\} = -i_{X_f} \omega(X_g) = df(X_g) = X_g(f). \quad (3.117)$$

It follows that the equation of motion in Poisson bracket structure is then

$$\dot{f} = \{f, H\} = X_H(f) = \mathcal{L}_{X_H} f. \quad (3.118)$$

We see that the Lie derivative along X_H generates time evolution!

So far we have not explained why the symplectic two-form was taken to be closed. This is required in order for the Poisson bracket to obey the Jacobi identity. It is also a necessary (and sufficient) condition for the symplectic form to be invariant under Hamiltonian flow.

3.4.4 Integration

We have learnt how to differentiate on a manifold using a vector field X , what about integration? What can we integrate on a manifold and how? It turns out that it is differential forms that we can integrate.

To begin we need to define an *orientation* on a manifold. Let M be a connected m -dimensional differentiable manifold. At a point $p \in M$ the tangent space $T_p(M)$ is spanned by the basis $\{e_\mu\} = \{\frac{\partial}{\partial x^\mu}\}$ where x^μ is the local coordinate on the chart U_i which contains p . Take U_j to be another chart such that $U_i \cup U_j \neq \emptyset$ and such that $p \in U_i \cup U_j$. Then the tangent space $T_p(M)$ is spanned by both $\{e_\mu\}$ or $\{\tilde{e}_\nu\} = \{\frac{\partial}{\partial y^\nu}\}$. The change of basis is

$$\tilde{e}_\nu = \frac{\partial x^\mu}{\partial y^\nu} e_\mu \equiv J^\mu{}_\nu e_\mu. \quad (3.119)$$

If $\det(J) > 0$ on $U_i \cup U_j$, the two bases $\{e_\mu\}$ and $\{\tilde{e}_\nu\}$ are said to defined the *same orientation* on $U_i \cup U_j$. If on the other hand $\det(J) < 0$ then they define the *opposite orientation*.

Definition Ley M be a connected manifold covered by $\{U_i\}$. The manifold M is *orientable* if for any overlapping charts U_i, U_j there exist local coordinates $\{x^\mu\}$ for U_i and $\{y^\nu\}$ for U_j such that $\det(\frac{\partial x^\mu}{\partial y^\nu}) > 0$.

If M is non-orientable, J cannot be made positive in all intersections of charts. An example of a non-orientable manifold is the Möbius strip, see figure 10. To construct a Möbius strip take two rectangles and glue them together with a twist of π on one of the edges to glue.

If an m -dimensional manifold M is orientable there exists an m -form ω which is nowhere vanishing, called the *volume form* or *volume element*. It plays the role of the measure when we integrate a function $f \in \mathcal{F}(M)$ over M . Two volume elements are said to be *equivalent* if there exists a strictly positive function $h \in \mathcal{F}(M)$ such that $\omega = h\omega'$. A negative-definite

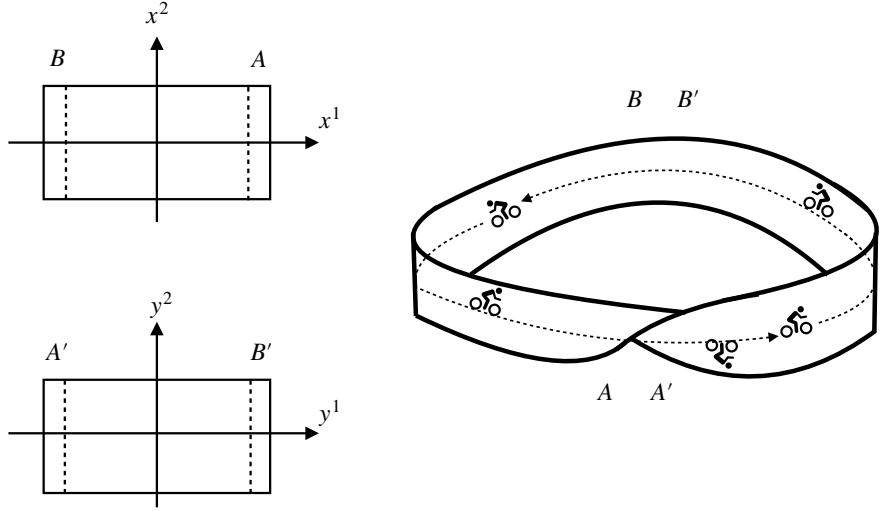


Figure 10: To construct the Möbius strip we glue two rectangles together: A with A' and B with B' . When joining A with A' we twist by π . The coordinate transformation on the A, A' intersection $y^1 = x^1$ and $y^2 = -x^2$, which has Jacobian -1 and is thus not orientable. We see that the cyclist going around the Möbius strip end up “up-side down” as they travel around the strip.

function $h' \in \mathcal{F}(M)$ gives and inequivalent orientation to M . Therefore for any orientable manifold there are two inequivalent orientations, we may refer to one of them as right-handed and the other as left-handed.

Since the volume form is a top form, and thus may be written locally as

$$\omega = h(x)dx^1 \wedge \dots \wedge dx^m, \quad (3.120)$$

with the requirement that $h(x) \neq 0$. We must then be able to patch this over the whole manifold without the handedness changing. Suppose that we have two sets of coordinates x^μ and y^ν on the charts U_i and U_j respectively, then in the new coordinates we have

$$\omega = h(x) \frac{\partial x^1}{\partial y^{\nu_1}} dy^{\nu_1} \wedge \dots \wedge \frac{\partial x^n}{\partial y^{\nu_m}} dy^{\nu_m} = h(x) \det \left(\frac{\partial x^\mu}{\partial y^\nu} \right) dy^1 \wedge \dots \wedge dy^m, \quad (3.121)$$

which makes clear that we may only define a volume form when the manifold is orientable. For the Möbius strip we see that we begin with volume form $\omega = dx \wedge dy$ but as we change charts this becomes $\omega = -dx \wedge dy$ and so ω is not defined uniquely on the Möbius strip.

Now we can define an integration of a function $f : M \rightarrow \mathbb{R}$ over an orientable manifold M . Let us take the volume form to be ω . Then in a coordinate neighbourhood U_i with

coordinates x^μ we define the integration of an m -form $f\omega$ to be

$$\int_{U_i} f\omega \equiv \int_{\varphi(U_i)} f(\varphi_i^{-1}(x)) h(\varphi_i^{-1}(x)) dx^1 \dots dx^m. \quad (3.122)$$

Notice that the right-hand side is an ordinary integration, albeit in m variables. Once the integral of f over U_i is defined it can be extended to an integration over all of M by making use of a *partition of unity*.

Definition Take an open covering $\{U_i\}$ on M such that each point of M is covered with a finite number of U_i . If this is always possible we call M *paracompact*.¹⁵ If a family of differentiable functions $\epsilon_i(p)$ satisfies

1. $0 \leq \epsilon_i(p) \leq 1$,
2. $\epsilon_i(p) = 0$ if $p \notin U_i$,
3. $\epsilon_1(p) + \epsilon_2(p) + \dots = 1$ for every point $p \in M$.

The family $\{\epsilon_i(p)\}$ is called a *partition of unity* for the covering $\{U_i\}$.

From condition (3) it follows that

$$f(p) = \sum_i f(p)\epsilon_i(p) = \sum_i f_i(p), \quad f_i(p) \equiv \epsilon_i(p)f(p). \quad (3.123)$$

Hence given a point $p \in M$ assumed paracompactness ensures that there are only a finite number of terms in the summation over i . For each of the $f_i(p)$ we may define the integral over U_i via (3.122), and therefore we have

$$\int_M f\omega \equiv \sum_i \int_{U_i} f_i\omega. \quad (3.124)$$

Though a different choice of atlas gives a different set of coordinates and a different partition of unity the integral as defined above stays the same.

Example Let us consider integrating a function on the circle. Let us take the atlas as given in (3.6) and (3.7). Let $U_1 = S^1 - \{(1, 0)\}$ and $U_2 = S^1 - \{(-1, 0)\}$. Then we may give a partition of unity by fixing $\epsilon_1(\theta) = \sin^2 \frac{\theta}{2}$ and $\epsilon_2(\theta) = \cos^2 \frac{\theta}{2}$. Note that $\epsilon_1(0) = 0$ and $\epsilon_2(\pi) = 0$ and therefore they vanish at the removed points as required. Moreover $\epsilon_1(\theta) + \epsilon_2(\theta) = 1$ as required. Thus $\{\epsilon_i(\theta)\}$ furnishes us with a partition of unity subordinate to $\{U_i\}$. We can, for an example, integrate the function $f = \cos^2 \theta$. Of course we know

$$\int_0^{2\pi} d\theta \cos^2 \theta = \pi, \quad (3.125)$$

¹⁵We will assume this is the case whenever we integrate something in this course.

but we should check with our partition of unity that we obtain the same result. We find

$$\int_{S^1} d\theta \cos^2 \theta = \int_0^{2\pi} d\theta \sin^2 \frac{\theta}{2} \cos^2 \theta + \int_{-\pi}^{\pi} d\theta \cos^2 \frac{\theta}{2} \cos^2 \theta = \frac{1}{2}\pi + \frac{1}{2}\pi = \pi. \quad (3.126)$$

So far we have left the function $h(x)$ appearing in the volume-form arbitrary. Since this gets multiplied by the Jacobian it changes between different coordinate patches and therefore there is no canonical way to pick this. Once we endow the manifold with a metric, as we require to GR, there is a canonical choice that we can make.

We can also integrate forms over sub-manifolds of M , rather than the full manifold. A manifold Σ with dimension $k < n$ is a *sub-manifold* of M if we can find a map $\sigma : \Sigma \rightarrow M$ which is one-to-one and $\sigma_* : T_p(\Sigma) \rightarrow T_{\sigma(p)}(M)$ is also one-to-one. We can then integrate a k -form ω on M over a k -dimensional sub-manifold Σ by pulling the form back to Σ :

$$\int_{\sigma(\Sigma)} \omega = \int_{\Sigma} \sigma^* \omega. \quad (3.127)$$

For example consider a one-form A living on M and take C to be a one-dimensional manifold in M . We can introduce a map $\sigma : C \rightarrow M$ which defines a non-intersecting curve $\sigma(C)$ which is a sub-manifold of M . We can then pull-back A onto the curve and integrate to obtain,

$$\int_{\sigma(C)} A = \int_C \sigma^* A. \quad (3.128)$$

Let the curve trace out a path $x^\mu(\tau)$ in M then, in coordinates this reads

$$\int_C \sigma^* A = \int d\tau A_\mu(x) \frac{dx^\mu}{d\tau}, \quad (3.129)$$

which is precisely the way in which a worldline of a particle couples to the electromagnetic field.

Stokes Theorem Until now our focus has been on smooth manifolds without boundary. We saw that this can be extended to manifolds with a boundary in section 3.1. There we have charts $\varphi : M \rightarrow U_i$ where U_i is an open subset of $\mathbb{R}^m = \{(x^1, \dots, x^m) | x^m \geq 0\}$. The boundary is denoted by ∂M , and is the sub-manifold fixed by $x^m = 0$. Then for a manifold M with a boundary, for any $(m-1)$ -form ω we have

$$\int_M d\omega = \int_{\partial M} \omega. \quad (3.130)$$

Stoke's theorem is the mother of all integral theorems. You may be familiar with the divergence theorem, Green's theorem, etc. for example, this is the generalisation of those. *Exercise: show that this reduces to Stoke's theorem on \mathbb{R}^3 .*

4 Riemannian geometry

We now have all the necessary pre-requisites to introduce the M.V.P. of general relativity: the metric. The introduction of a metric brings a whole slew of new objects that we can define.

4.1 The metric

Definition: Let M be a differentiable manifold. A *Riemannian metric* g on M is a type $(0, 2)$ tensor field on M which at each point $p \in M$ satisfies

- Symmetric: $g_p(X, Y) = g_p(Y, X)$,
- $g_p(X, X) \geq 0$ with equality iff $X = 0$

with $X, Y \in T_p(M)$. A tensor field g of type $(0, 2)$ is a *pseudo-Riemannian metric* if it satisfies the first condition and

- Non-degenerate. If for any $p \in M$ $g_p(X, Y) = 0$ for all $Y \in T_p(M)$ then $X_p = 0$,

We may extend the tensor g_p over the full manifold. With a choice of coordinates we can write the metric as

$$g = g_{\mu\nu}(x)dx^\mu \otimes dx^\nu. \quad (4.1)$$

We will often write this as the line elements ds^2 ,

$$ds^2 = g_{\mu\nu}(x)dx^\mu dx^\nu. \quad (4.2)$$

Strictly this is not a metric, since the metric is a tensor, however we will often use this abuse of terminology as is common in the field.

One can extract out the components by evaluating the metric on a pair of basis elements

$$g_{\mu\nu}(x) = g\left(\frac{\partial}{\partial x^\mu}, \frac{\partial}{\partial x^\nu}\right). \quad (4.3)$$

We may view $g_{\mu\nu}$ as a matrix, which by the symmetry property above is symmetric. This implies that the matrix is diagonalisable, with real eigenvalues. If there are i positive eigenvalues and j negative eigenvalues the pair (i, j) is called the *index* of the metric. If $j = 1$ the metric is called a *Lorentz metric*, for $j = 0$ we have a *Euclidean metric*. The number of negative entries is called the *signature* and by Sylvester's law of inertia¹⁶, this is independent of the choice of basis.

¹⁶This has nothing to do with inertia, Sylvester just wanted a law of inertia like Newton.

For most applications of differential geometry, we are interested in manifolds with signature 0, i.e. a Riemannian manifold. The simplest example which you are probably familiar with, though maybe not in this language, is the metric on Euclidean space \mathbb{R}^m , which in Cartesian coordinates has the metric

$$g = dx^1 \otimes dx^1 + \dots + dx^m \otimes dx^m, \quad (4.4)$$

which in components reads $g_{\mu\nu} = \delta_{\mu\nu}$.

4.1.1 Riemannian metric

A general Riemannian metric is a useful object to have in ones tool belt. It gives us a way of measuring the length of a vector X at each point

$$|X| = \sqrt{g(X, X)}. \quad (4.5)$$

Moreover we may measure the angle between two vectors

$$g(X, Y) = |X||Y|\cos\theta. \quad (4.6)$$

Furthermore it can be used to measure the distance between two points p and q along a curve in M . For the curve $\sigma : [a, b] \rightarrow M$ with $\sigma(a) = p$ and $\sigma(b) = q$ the distance between the two points along the curve is

$$d(p, q) = \int_a^b dt \sqrt{g(X, X)|_{\sigma(t)}}, \quad (4.7)$$

where X is the tangent vector field of the curve. If the curve has coordinates $x^\mu(t)$ then $X^\mu = \frac{dx^\mu}{dt}$ and the distance is

$$d(p, q) = \int_a^b dt \sqrt{g_{\mu\nu} \frac{dx^\mu(t)}{dt} \frac{dx^\nu(t)}{dt}}. \quad (4.8)$$

This distance does not depend on the parametrisation of the curve.

4.1.2 Lorentzian manifolds

For our purposes Riemannian manifolds are not what we want to consider, instead we want to consider Lorentzian manifolds. The simplest example is Minkowski space. This is $\mathbb{R}^{1,m-1}$ equipped with the metric

$$\eta = -dx^0 \otimes dx^0 + dx^1 \otimes dx^1 + \dots + dx^{m-1} \otimes dx^{m-1}, \quad (4.9)$$

which has components $\eta_{\mu\nu} = \text{diag}(-1, 1, \dots, 1)$. Note that on a Lorentzian manifold we take the index to run over $0, 1, \dots, m-1$.

At any point p on a general Lorentzian manifold it is always possible to find an orthonormal basis $\{e_\mu\}$ of $T_p(M)$ such that locally the metric looks like the Minkowski metric

$$g_{\mu\nu}|_p = \eta_{\mu\nu}. \quad (4.10)$$

This is closely related to the equivalence principle we discussed previously.

The fact that locally the metric looks locally like Minkowski space allows us to import some of the ideas of special relativity, namely we can classify the elements of $T_p(M)$ into three classes

- $g(X, X) > 0 \longrightarrow X$ is *spacelike* ,
- $g(X, X) = 0 \longrightarrow X$ is *lightlike* or *null* ,
- $g(X, X) < 0 \longrightarrow X$ is *timelike* .

At each point on M we can then draw light cones which are the null tangent vectors at that point. The novelty is that the directions of these light cones can vary smoothly as we move around the manifold. This specifies the causal structure of spacetime which determines which regions of spacetime can interact together.

As in the Riemannian case we can use the metric to determine the length of curves. The nature of a curve is inherited from the nature of its tangent vector. A curve is called *timelike* if its tangent vector is everywhere timelike. We then measure the proper time

$$\tau = \int_a^b dt \sqrt{-g_{\mu\nu} \frac{dx^\mu}{dt} \frac{dx^\nu}{dt}}. \quad (4.11)$$

4.1.3 Why is the metric useful

The existence of a metric comes with a large number of benefits.

The metric as an isomorphism The metric gives a natural isomorphism between vectors and covectors, $g : T_p(M) \rightarrow T_p^*(M)$ for each p . In a coordinate basis we can write $X = X^\mu \partial_\mu$, and map it to a one-form $X = X_\mu dx^\mu$, as

$$X_\mu = g_{\mu\nu} X^\nu. \quad (4.12)$$

We will usually say that we use the metric to lower (or raise) an index. What we really mean is that the metric provides an isomorphism between a vector space and its dual. Since g is non-degenerate and is thus invertible we also have the inverse map. We take the inverse

of $g_{\mu\nu}$ to be $g^{\mu\nu}$ so that $g^{\mu\nu}g_{\nu\rho} = \delta_\rho^\mu$. This can then be thought of as the components of a symmetric $(2,0)$ tensor

$$\hat{g} = g^{\mu\nu}\partial_\mu \otimes \partial_\nu. \quad (4.13)$$

Then

$$X^\mu = g^{\mu\nu}X_\nu. \quad (4.14)$$

In Euclidean space since $g_{\mu\nu} = \delta_{\mu\nu}$ one does not immediately notice the distinction between vectors and one-forms.

The Volume form The metric also gives a natural volume form on the manifold M . On a Riemannian manifold we take the volume form to be

$$\text{vol}(M) = \sqrt{\det(g_{\mu\nu})}dx^1 \wedge \dots \wedge dx^m, \quad (4.15)$$

and we use the shorthand $\sqrt{\det(g_{\mu\nu})} = \sqrt{g}$. On a Lorentzian manifold the determinant is negative and therefore we take the volume form to be

$$\text{vol}(M) = \sqrt{-g}dx^0 \wedge dx^1 \wedge \dots \wedge dx^{n-1}. \quad (4.16)$$

As it is written it looks coordinate dependent however it is not. To see this recall that if we change coordinates $y = y(x)$ we have (see (3.29))

$$dx^\mu = \Lambda^\mu{}_\nu dy^\nu, \quad \Lambda^\mu{}_\nu = \frac{\partial x^\mu}{\partial y^\nu}. \quad (4.17)$$

Then

$$\begin{aligned} dx^1 \wedge \dots \wedge dx^m &= \Lambda^1{}_{\nu_1} \dots \Lambda^m{}_{\nu_m} dy^{\mu_1} \wedge \dots \wedge dy^{\mu_m} \\ &= \sum_{P \in S_m} \text{sgn}(P) \Lambda^1{}_{P(1)} \dots \Lambda^m{}_{P(m)} dy^1 \wedge \dots \wedge dy^m \\ &= \det(\Lambda) dy^1 \wedge \dots \wedge dy^m, \end{aligned} \quad (4.18)$$

where in the penultimate line we have used the properties of the wedge product and in the last line used the definition of the determinant. The metric components transform as

$$g_{\mu\nu} = \frac{\partial y^\rho}{\partial x^\mu} \frac{\partial y^\sigma}{\partial x^\nu} \tilde{g}_{\rho\sigma}, \quad (4.19)$$

and therefore

$$\det(g_{\mu\nu}) = \det(\tilde{g}_{\mu\nu}) (\det(\Lambda))^{-2}, \quad (4.20)$$

and therefore this cancels with the transformation of the wedge product leaving

$$\text{vol}(M) = \sqrt{|g|}dy^1 \wedge \dots \wedge dy^m. \quad (4.21)$$

We may rewrite the volume form as

$$\text{vol}(M) = \frac{1}{m!} v_{\mu_1 \dots \mu_m} dx^{\mu_1} \wedge \dots \wedge dx^{\mu_m}, \quad \text{where} \quad v_{\mu_1 \dots \mu_m} = \sqrt{|g|} \epsilon_{\mu_1 \dots \mu_m}. \quad (4.22)$$

It follows that $v_{\mu_1 \dots \mu_m}$ is a tensor, while $\epsilon_{\mu_1 \dots \mu_m}$ is not, instead it is a tensor density (one needs to multiply by the square root of the determinant to obtain a tensor). Note that we define $\epsilon^{\mu_1 \dots \mu_m}$ to again be the totally anti-symmetric tensor with $\epsilon^{1 \dots m} = 1$, i.e. we do not raise the indices on ϵ with the metric. Instead we have

$$v^{\mu_1 \dots \mu_m} = g^{\mu_1 \nu_1} \dots g^{\mu_m \nu_m} v_{\nu_1 \dots \nu_m} = \pm \frac{1}{\sqrt{|g|}} \epsilon^{\mu_1 \dots \mu_m}. \quad (4.23)$$

Hodge dual On an oriented manifold M we can use the totally anti-symmetric tensor density to define a map which takes a p -form $\omega \in \Omega^p(M)$ to a $(m-p)$ -form $\star \omega \in \Omega^{m-p}(M)$. We define this map to be

$$(\star \omega)_{\mu_1 \dots \mu_{m-p}} = \frac{1}{p!} \sqrt{|g|} \epsilon_{\mu_1 \dots \mu_{m-p} \nu_1 \dots \nu_p} \omega^{\nu_1 \dots \nu_p}. \quad (4.24)$$

This is called the *Hodge dual* and is independent of coordinates. One can see that it satisfies

$$\star(\star \omega) = \pm (-1)^{p(m-p)} \omega, \quad (4.25)$$

with + for a Riemannian metric and - for a Lorentzian.¹⁷

With the Hodge dual in tow we can define an inner product on each vector space $\Omega^r(M)$. If $\omega, \eta \in \Omega^r(M)$ then

$$\langle \eta, \omega \rangle \equiv \int_M \eta \wedge \star \omega. \quad (4.26)$$

With such an inner product one can look at operators on $\Omega^r(M)$ and their adjoints. The differential operator we have introduced on r -forms is the exterior derivative. For $\omega \in \Omega^r(M)$ and $\alpha \in \Omega^{r-1}(M)$ the adjoint is defined via

$$\langle d\alpha, \omega \rangle = \langle \alpha, d^\dagger \omega \rangle, \quad (4.27)$$

where the adjoint operator $d^\dagger : \Omega^r(M) \rightarrow \Omega^{r-1}(M)$ is given by

$$d^\dagger = \pm (-1)^{m(r+1)-1} \star d \star. \quad (4.28)$$

¹⁷One has actually seen the Hodge dual before, it was just hidden from view. Consider two vectors \vec{a} and \vec{b} in \mathbb{R}^3 , We can take the cross product to obtain a third vector \vec{c} as $\vec{a} \times \vec{b} = \vec{c}$. This however mixes a lot of different objects. This is equivalent in our new language to first use the metric to relate the vectors to one-forms. The cross product is really the wedge product of the two one-forms to give a two-form. We then take the Hodge dual of this two-form to obtain a one-form and then use the metric once again to extract out a vector. This more complicated route is hidden since the metric is just the Kronecker delta and so we can raise and lower indices with impunity. Going to curved space and a non-trivial metric these subtleties become relevant.

One can then define a Laplacian $\square : \Omega^r(M) \rightarrow \Omega^r(M)$ defined as¹⁸

$$\square = (\mathrm{d} + \mathrm{d}^\dagger)^2 = \mathrm{d}\mathrm{d}^\dagger + \mathrm{d}^\dagger\mathrm{d}. \quad (4.29)$$

It can be defined on both Riemannian manifolds and Lorentzian, however it is only positive definite on Riemannian manifolds. On a function f the Laplacian acts as

$$\square f = -\frac{1}{\sqrt{|g|}}\partial_\nu\left(\sqrt{|g|}g^{\mu\nu}\partial_\mu f\right). \quad (4.30)$$

Aside: There is a beautiful interplay between the Eigenforms and Eigenvalues of the Laplacian and the topology of the space that we will not cover. If one defines a harmonic form to be one which is annihilated by the Laplacian $\square\omega = 0$, then there is an isomorphism between the set of all harmonic forms and the cohomology group:

$$\mathrm{Harm}^r(M) \cong H^r(M). \quad (4.31)$$

The Betti numbers which were the dimensions of the cohomology groups are then just the dimension of the group of harmonic r -forms on the manifold.

4.2 Connections and curvature

A vector field X is a directional derivative acting on a function $f \in \mathcal{F}(M)$. However so far we have not introduced such a derivative for tensors of type (q, r) . The Lie derivative is not quite what we want since it also involves derivatives of the vector defining the direction. This other derivative is more useful than the Lie derivative, but requires the introduction of a *connection* to map the vector spaces at one point to vector spaces at another. The resultant object is known as the *covariant derivative* and is distinct from the Lie derivative that we introduced previously.

An *affine connections* ∇ is a map $\nabla : \mathcal{X}(M) \times \mathcal{X}(M) \rightarrow \mathcal{X}(M)$, $(X, Y) \mapsto \nabla_X Y$ which satisfies

$$\nabla_X(Y + Z) = \nabla_X Y + \nabla_X Z, \quad (4.32)$$

$$\nabla_{(fX+gY)}Z = f\nabla_X Z + g\nabla_Y Z, \quad (4.33)$$

$$\nabla_X(fY) = X[f]Y + f\nabla_X Y, \quad (4.34)$$

for vector fields $X, Y, Z \in \mathcal{X}(M)$ and functions $f, g \in \mathcal{F}(M)$.

¹⁸You may also see the Laplacian denoted by \triangle rather than \square .

Let us take a chart (U, φ) with coordinate $x = \varphi(p)$ and define m^3 functions $\Gamma^\mu_{\nu\rho}$ called the *connection coefficients* by

$$\nabla_\nu e_\mu \equiv \nabla_{e_\nu} e_\mu = e_\lambda \Gamma^\lambda_{\nu\mu}, \quad (4.35)$$

where $\{e_\mu\} = \{\frac{\partial}{\partial x^\mu}\}$ is the coordinate basis in $T_p(M)$. The connection coefficients specify how the basis vectors change from point to point, i.e. how to map the tangent space $T_p(M)$ to $T_q(M)$. Using the properties of the connection we can work out the general covariant derivative of a vector field

$$\begin{aligned} \nabla_X Y &= \nabla_X(Y^\mu e_\mu) \\ &= X[Y^\mu]e_\mu + Y^\mu \nabla_X e_\mu \\ &= X^\nu \partial_\nu(Y^\mu)e_\mu + X^\nu Y^\mu \nabla_\nu e_\mu \\ &= X^\nu \left(\partial_\nu Y^\mu + \Gamma^\mu_{\nu\rho} Y^\rho \right) e_\mu. \end{aligned} \quad (4.36)$$

We can strip off the overall X^ν to write

$$(\nabla_\nu Y)^\mu = \frac{\partial Y^\mu}{\partial x^\nu} + \Gamma^\mu_{\nu\rho} Y^\rho, \quad (4.37)$$

so that

$$(\nabla_X Y)^\mu = X^\nu \nabla_\nu Y^\mu \quad (4.38)$$

In a function the covariant derivative coincides with both the Lie derivative and the regular partial derivative, however its action on vectors differs. While the Lie derivative $\mathcal{L}_X Y$ depends on both X and its first derivative, the covariant derivative depends only on X . This is the natural generalisation of the partial derivative on curved space.

We will often be sloppy and write

$$(\nabla_X Y)^\mu = \nabla_\nu Y^\mu. \quad (4.39)$$

Typically in older books, though some still like to use this stupid convention, one may see the semi-colon notation

$$\nabla_\nu Y^\mu = Y^\mu_{;\nu}. \quad (4.40)$$

We will refrain from using this convention to preserve our sanity.

At the moment the connection $\Gamma^\mu_{\nu\rho}$ is somewhat abstract. One may guess that it is a tensor however this is not correct. To see this let us consider how it transforms under a change of coordinates. Recall that the basis elements transform as

$$\tilde{e}_\mu = \Lambda^\mu_{\nu} e_\mu, \quad \text{with} \quad \Lambda^\mu_{\nu} = \frac{\partial y^\mu}{\partial x^\nu}. \quad (4.41)$$

Recall that a $(1, 2)$ tensor $T^\mu_{\nu\rho}$ transforms as

$$\tilde{T}^{\mu_1}_{\nu_1\rho_1} = (\Lambda^{-1})^{\mu_1}_{\mu_2} \Lambda^{\nu_2}_{\nu_1} \Lambda^{\rho_2}_{\rho_1} T^{\mu_2}_{\nu_2\rho_2}. \quad (4.42)$$

We can compute the transformation of the connection. In the basis $\{\tilde{e}_\mu\}$ we have

$$\begin{aligned} \nabla_{\tilde{e}_\rho} \tilde{e}_\nu &= \tilde{\Gamma}^\mu_{\nu\rho} \tilde{e}_\mu \\ &= \nabla_{\Lambda^\sigma_\rho} e_\sigma \left(\Lambda^\tau_\nu e_\tau \right) \\ &= \Lambda^\sigma_\rho \left(\nabla_\sigma (\Lambda^\tau_\nu) e_\tau + \Lambda^\tau_\nu \nabla_\sigma e_\tau \right) \\ &= \Lambda^\sigma_\rho \left(\Lambda^\tau_\nu \Gamma^\kappa_{\sigma\tau} + \partial_\sigma \Lambda^\kappa_\nu \right) e_\kappa \\ &= \Lambda^\sigma_\rho \left(\Lambda^\tau_\nu \Gamma^\kappa_{\sigma\tau} + \partial_\sigma \Lambda^\kappa_\nu \right) (\Lambda^{-1})^\mu_\kappa \tilde{e}_\mu. \end{aligned} \quad (4.43)$$

From this we obtain

$$\tilde{\Gamma}^\mu_{\nu\rho} = (\Lambda^{-1})^\mu_\kappa \Lambda^\sigma_\rho \Lambda^\tau_\nu \Gamma^\kappa_{\sigma\tau} + (\Lambda^{-1})^\mu_\kappa \Lambda^\sigma_\rho \partial_\sigma \Lambda^\kappa_\nu. \quad (4.44)$$

The first term is the expected transformation term of a $(1, 2)$ tensor, however there is an additional piece. This additional piece is independent of Γ and depends only on the $\partial\Lambda$. This is the characteristic transformation of a connection.

Differentiating other tensors We can use the properties of the covariant derivative to extend its action to any tensor field. Consider a one-form ω . If we differentiate ω we will get another one-form $\nabla_X \omega$, so we should check its action on a vector field $Y \in \mathcal{X}(M)$. We impose that the connection obeys the Leibniz identity

$$\nabla_X(\omega(Y)) = (\nabla_X \omega)(Y) + \omega(\nabla_X Y). \quad (4.45)$$

Since $\omega(Y)$ is a function we know that

$$\nabla_X(\omega(Y)) = X[\omega(Y)]. \quad (4.46)$$

Using the Leibniz condition we have

$$(\nabla_X \omega)(Y) = X(\omega(Y)) - \omega(\nabla_X Y), \quad (4.47)$$

and reducing to coordinates we find

$$\begin{aligned} X^\mu (\nabla_\mu \omega)_\nu Y^\nu &= X^\mu \partial_\mu (\omega_\nu Y^\nu) - \omega_\nu X^\mu (\partial_\mu Y^\nu + \Gamma^\nu_{\mu\rho} Y^\rho) \\ &= X^\mu \left(\partial_\mu \omega_\nu - \Gamma^\nu_{\mu\rho} \omega_\nu \right) Y^\rho. \end{aligned} \quad (4.48)$$

We may then write

$$(\nabla_\mu \omega)_\rho \equiv \nabla_\mu \omega_\rho = \frac{\partial}{\partial x^\mu} \omega_\rho - \Gamma^\nu{}_{\mu\rho} \omega_\nu. \quad (4.49)$$

We can now extend this argument to an arbitrary tensor of rank (q, r) and we find

$$\begin{aligned} \nabla_\mu T^{\nu_1 \dots \nu_q}_{\rho_1 \dots \rho_r} &= \frac{\partial}{\partial x^\mu} T^{\nu_1 \dots \nu_q}_{\rho_1 \dots \rho_r} + \Gamma^{\nu_1}{}_{\mu\sigma} T^{\sigma \dots \nu_q}_{\rho_1 \dots \rho_r} + \dots + \Gamma^{\nu_q}{}_{\mu\sigma} T^{\nu_1 \dots \nu_{q-1}\sigma}_{\rho_1 \dots \rho_r} \\ &\quad - \Gamma^\sigma{}_{\mu\rho_1} T^{\nu_1 \dots \nu_q}_{\sigma \dots \rho_r} - \dots - \Gamma^\sigma{}_{\mu\rho_r} T^{\nu_1 \dots \nu_q}_{\rho_1 \dots \rho_{r-1}\sigma}. \end{aligned} \quad (4.50)$$

In words, you first differentiate the tensor and then for each upper index you add in a $+\Gamma T$ and for every down index a $-\Gamma T$.

4.3 Torsion and curvature

Even though the connection is not a tensor we can use it to construct two tensors. The first is a rank $(1, 2)$ tensor T known as *Torsion*, the second is a rank $(1, 3)$ tensor known as *curvature* or the *Riemann tensor*. The torsion tensor acts on $X, Y \in \mathcal{X}(M)$ and $\omega \in \Omega^1(M)$ by

$$T(\omega : X, Y) = \omega(\nabla_X Y - \nabla_Y X - [X, Y]). \quad (4.51)$$

We may equivalently think of this as a map $T : \mathcal{X}(M) \times \mathcal{X}(M) \rightarrow \mathcal{X}(M)$ defined by

$$T(X, Y) = \nabla_X Y - \nabla_Y X - [X, Y]. \quad (4.52)$$

The curvature acts on $X, Y, Z \in \mathcal{X}(M)$ and $\omega \in \Omega^1(M)$ as

$$R(\omega : X, Y, Z) = \omega(\nabla_X \nabla_Y Z - \nabla_Y \nabla_X Z - \nabla_{[X, Y]} Z) \quad (4.53)$$

As for the torsion we may think of this as a map $\mathcal{X}(M) \times \mathcal{X}(M)$ to a differential operator acting on $\mathcal{X}(M)$ as

$$R(X, Y) = \nabla_X \nabla_Y - \nabla_Y \nabla_X - \nabla_{[X, Y]}. \quad (4.54)$$

Exercise: Check that both the Torsion and Curvature tensors are actually tensors. You should show that they are linear in all arguments, for example show $T(\omega : fX, Y) = fT(\omega : X, Y)$ for all $f \in \mathcal{F}(M)$.

Component form We can evaluate the tensors in a basis to obtain the component form. Let $\{f^\rho\} = \{dx^\rho\}$ then the components of the torsion tensor are

$$\begin{aligned} T^\rho_{\mu\nu} &= T(f^\rho : e_\mu, e_\nu) \\ &= f^\rho (\nabla_\mu e_\nu - \nabla_\nu e_\mu - [e_\mu, e_\nu]) \\ &= f^\rho (\Gamma^\sigma{}_{\mu\nu} - \Gamma^\sigma{}_{\nu\mu}) e_\sigma \\ &= \Gamma^\rho{}_{\mu\nu} - \Gamma^\rho{}_{\nu\mu}. \end{aligned} \quad (4.55)$$

So despite $\Gamma^\sigma_{\mu\nu}$ not being a tensor, the anti-symmetrised part is! The torsion tensor is clearly anti-symmetric in the two lowered indices. We see that connections $\Gamma^\sigma_{\mu\nu}$ which are symmetric in the lowered indices have $T^\rho_{\mu\nu} = 0$ and are called *torsion-free*.

A similar computation for the Riemann tensor gives

$$R^\sigma_{\rho\mu\nu} = \partial_\mu \Gamma^\sigma_{\nu\rho} - \partial_\nu \Gamma^\sigma_{\mu\rho} + \Gamma^\lambda_{\nu\rho} \Gamma^\sigma_{\mu\lambda} - \Gamma^\lambda_{\mu\rho} \Gamma^\sigma_{\nu\lambda}. \quad (4.56)$$

Consider the commutator of covariant derivatives acting on a vector field

$$\begin{aligned} \nabla_{[\mu} \nabla_{\nu]} X^\sigma &= \partial_{[\mu} (\nabla_{\nu]} X^\sigma) + \Gamma^\sigma_{[\mu|\lambda]} \nabla_{\nu]} Z^\lambda - \Gamma^\rho_{[\mu\nu]} \nabla_\rho X^\sigma \\ &= \partial_{[\mu} \partial_{\nu]} X^\sigma + (\partial_{[\mu} \Gamma^\sigma_{\nu]\rho}) X^\sigma + (\partial_{[\mu} X^\rho) \Gamma^\sigma_{\nu]\rho} + \Gamma^\sigma_{[\mu|\lambda]} \partial_{\nu]} X^\lambda \\ &\quad + \Gamma^\sigma_{[\mu|\lambda]} \Gamma^\lambda_{\nu]\rho} X^\rho - \Gamma^\rho_{[\mu\nu]} \nabla_\rho X^\sigma. \end{aligned} \quad (4.57)$$

The first term vanishes, while the third and fourth cancel. The Second and fifth combine to give the Riemann tensor while the last gives the torsion, we have

$$2\nabla_{[\mu} \nabla_{\nu]} X^\sigma = R^\sigma_{\rho\mu\nu} X^\rho - T^\rho_{\mu\nu} \nabla_\rho X^\sigma. \quad (4.58)$$

This is the *Ricci identity*. Similar identities hold when acting on other tensors.

4.3.1 Levi–Civita connection

So far the discussion has not required a metric. When a metric exists we have

Theorem There exists a unique, torsion free, connection that is compatible with the metric g :

$$\nabla_X g = 0, \quad (4.59)$$

for all vector fields X .

To prove this we first show uniqueness before constructing the connection. Suppose that such a connection exists, then we have

$$X(g(Y, Z)) = \nabla_X(g(Y, Z)) = (\nabla_X g)(Y, Z) + g(\nabla_X Y, Z) + g(Y, \nabla_X Z). \quad (4.60)$$

Since $\nabla_X g = 0$ we have

$$X(g(Y, Z)) = g(Y, \nabla_X Z) + g(\nabla_X Y, Z). \quad (4.61)$$

We may use our favourite trick and cyclically permute X, Y, Z to find

$$\begin{aligned} Y(g(Z, X)) &= g(Z, \nabla_Y X) + g(\nabla_Y Z, X), \\ Z(g(X, Y)) &= g(X, \nabla_Z Y) + g(\nabla_Z X, Y). \end{aligned} \quad (4.62)$$

By the no torsion condition we have

$$\nabla_X Y - \nabla_Y X = [X, Y], \quad (4.63)$$

and therefore

$$\begin{aligned} X(g(Y, Z)) &= g(\nabla_Y Z, X) + g(\nabla_Y X, Z) + g([X, Y], Z), \\ Y(g(Z, X)) &= g(\nabla_Z Y, X) + g(\nabla_Z X, Y) + g([Y, Z], X), \\ Z(g(X, Y)) &= g(\nabla_X Z, Y) + g(\nabla_X Y, Z) + g([Z, X], Y), \end{aligned} \quad (4.64)$$

Adding the first and second and subtracting the third we find

$$g(\nabla_Y X, Z) = \frac{1}{2} \left[X(g(Y, Z)) + Y(g(Z, X)) - Z(g(X, Y)) \right. \\ \left. - g([X, Y], Z) - g([Y, Z], X) + g([Z, X], Y) \right] \quad (4.65)$$

With a non-degenerate metric this specifies the connection uniquely.

It remains to be seen that the connection as defined does satisfy the properties of a connection. We will present one of the terms to check. The most finicky one is $\nabla_{fX} Y = f\nabla_X Y$, so let us present that one

$$\begin{aligned} g(\nabla_{fY} X, Z) &= \frac{1}{2} \left[X(g(fY, Z)) + fY(g(Z, X)) - Z(g(X, fY)) \right. \\ &\quad \left. - g([X, fY], Z) - g([fY, Z], X) + g([Z, X], fY) \right] \\ &= \frac{1}{2} \left[fX(g(Y, Z)) \color{red}{+ X(f)g(Y, Z)} + fY(g(Z, X)) - fZ(g(X, Y)) \color{blue}{- Z(f)g(X, Y)} \right. \\ &\quad \left. - fg([X, Y], Z) \color{red}{- X(f)g(Y, Z)} - fg([Y, Z], X) \color{blue}{+ Z(f)g(Y, X)} + fg([Z, X], Y) \right] \\ &= g(f\nabla_Y X, Z). \end{aligned} \quad (4.66)$$

The coloured terms in the penultimate line cancel amongst themselves, leaving just the black terms as required. The other properties follow similarly. This then proves the uniqueness and has explicitly constructed such a connection.

In components we can evaluate

$$g(\nabla_\nu e_\mu, e_\rho) = \Gamma^\lambda_{\nu\mu} g_{\lambda\rho} = \frac{1}{2} (\partial_\mu g_{\nu\rho} + \partial_\nu g_{\mu\rho} - \partial_\rho g_{\mu\nu} - \partial_\mu g_{\nu\rho}). \quad (4.67)$$

Multiplying by the inverse metric we have

$$\Gamma^\lambda_{\mu\nu} = \frac{1}{2} g^{\lambda\rho} (\partial_\mu g_{\nu\rho} + \partial_\nu g_{\mu\rho} - \partial_\rho g_{\mu\nu} - \partial_\mu g_{\nu\rho}). \quad (4.68)$$

The connection compatible with the metric is called the *Levi–Civita connection* while the components of the Levi–Civita connection are called the *Christoffel symbols*.

There is a nice expression if you contract two indices of the Christoffel symbols, we have

$$\Gamma^\mu_{\mu\nu} = \frac{1}{\sqrt{|g|}} \partial_\nu \sqrt{|g|} \quad (4.69)$$

To see this note

$$\Gamma^\mu_{\mu\nu} = \frac{1}{2} g^{\mu\rho} \partial_\nu g_{\mu\rho} = \frac{1}{2} \text{tr}(g^{-1} \partial_\nu g) = \frac{1}{2} \text{tr}(\partial_\nu \log g), \quad (4.70)$$

for diagonalisable matrices we have $\text{tr} \log A = \log \det(A)$ and therefore we find

$$\Gamma^\mu_{\mu\nu} = \frac{1}{2} \partial_\nu \log \det(g) = \frac{1}{\sqrt{\det(g)}} \partial_\nu \sqrt{\det(g)}. \quad (4.71)$$

This implies that

$$\sqrt{|g|} \nabla_\mu X^\mu = \sqrt{|g|} (\partial_\mu X^\mu + \Gamma^\mu_{\mu\nu} X^\nu) + \sqrt{|g|} \left(\partial_\mu X^\mu + X^\nu \frac{1}{\sqrt{|g|}} \partial_\nu \sqrt{|g|} \right) = \partial_\mu (\sqrt{|g|} X^\mu). \quad (4.72)$$

Using this result we can prove the divergence theorem:

$$\int_M d^m x \sqrt{|g|} \nabla_\mu X^\mu = \int_{\partial M} d^{n-1} x \sqrt{\gamma} n_\mu X^\mu, \quad (4.73)$$

where γ_{ij} is the pull-back of the metric to ∂M , $\gamma = \det(\gamma_{ij})$ and n_μ is an outward pointing unit vector orthogonal to ∂M . On a Lorentzian manifold this holds provided that ∂M is either purely spacelike or purely timelike, which guarantees that $\gamma \neq 0$.

4.4 Parallel transport and geodesics

We have introduced the connection but we are yet to explain what it connects. It connects tangent spaces, or more generally any vector space at different points of the manifold. This map is called *parallel transport*. Take a vector field X with some associated curve γ with coordinates $x^\mu(\lambda)$ such that

$$X^\mu|_\gamma = \frac{dx^\mu(\lambda)}{d\lambda}. \quad (4.74)$$

We say that a tensor field T is parallel transported along γ if

$$\nabla_X T = 0. \quad (4.75)$$

Let γ connect two points $p, q \in M$. The condition (4.75) provides a map from the vector space defined at p to the vector space defined at q . Consider a second vector field Y . In components (4.75) reads

$$X^\nu (\partial_\nu Y^\mu + \Gamma^\mu_{\nu\rho} Y^\rho) = 0. \quad (4.76)$$

If we evaluate it on the curve γ , we can write $Y^\mu = Y^\mu(x(\lambda))$ and therefore the condition is

$$\frac{dY^\mu}{d\lambda} + X^\nu \Gamma^\mu_{\nu\rho} Y^\rho = 0. \quad (4.77)$$

This defines a set of coupled ordinary differential equations, given an initial condition at $p = \gamma(\lambda = 0)$ for example these can be solved to find a unique vector field at each point along the curve. This is path dependent and depends on the connection and the underlying path which was characterised by X here.

There is a subtle difference between what we are doing here and what we did with the push-forward and pull-back, which we used to define the Lie derivative. Here X only appears to define the map, there are no derivatives applied to X^μ as was for those maps. The connection does the work of relating the vector spaces along the curve and not the vector X .

4.4.1 Geodesics

A *geodesic* is a curve tangent to a vector field X that obeys

$$\nabla_X X = 0. \quad (4.78)$$

Along the curve γ with coordinates x^μ and tangent vector X this implies

$$\frac{d^2x^\mu}{d\lambda^2} + \Gamma^\mu_{\nu\rho} \frac{dx^\nu}{d\lambda} \frac{dx^\rho}{d\lambda} = 0. \quad (4.79)$$

This is the same geodesic equation one obtains by varying the action

$$S = \int d\lambda \sqrt{-g_{\mu\nu}(x) \frac{dx^\mu}{d\lambda} \frac{dx^\nu}{d\lambda}}, \quad (4.80)$$

and picking an affine parameter if ∇ is the Levi–Civita connection.

If we choose the Levi–Civita connection, since $\nabla_X g = 0$ it follows that for any vector field Y which is parallel transported along a geodesic defined by X we have

$$\frac{d}{d\lambda} g(X, Y) = 0. \quad (4.81)$$

The vector field Y makes the same angle with the tangent vector at each point along the geodesic. Further, this holds true if we replace Y by X in the expression above. Since the norm of the vector field X tangent to the geodesic classifies the character of the geodesic, (timelike/null/spacelike), if we define a geodesic using a metric compatible connection, then the nature of the geodesic does not change. This statement relies on us using a metric compatible connection though, in this course we will always take such a connection and therefore the nature of a geodesic is preserved throughout all spacetime.

Let us consider a timelike geodesic. When we vary the action what are we extremising and is it a maximum or minimum? From our definition of the proper time we see that we are extremising the proper time, it turns out that geodesics maximise the proper time. Why is this true? Well given any time-like curve we can approximate it to arbitrary accuracy by a null curve. We should consider jagged null curves that follow the time-like one, see figure 11. As we increase the number of null curves the approximation gets better and better, while still having zero length. Timelike curves cannot therefore be curves with minimal proper time since they are infinitesimally close to curves of zero length (and therefore zero proper time). They must therefore maximise the proper time. This is why the twin who remains home in the twin paradox ages more, they are on a geodesic (pretty much). We should really say that this maximises the proper time locally. If we took a sphere, then there is more than one geodesic between two points, we can either go the short way around or the long way around. One is longer than the other (assuming the points are not opposite each other, i.e. picking the poles), but both maximise locally the length functional.

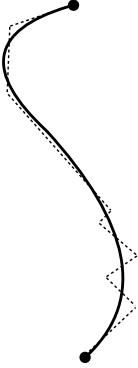


Figure 11: We approximate the time-like curve with null curves. As we increase the number of null curves the approximation gets better and better.

4.4.2 Normal coordinates

Geodesics allow for the construction of a particularly useful coordinate system. This holds independently of whether the Levi–Civita connection is employed or not, however it takes a particularly simple form when it is used. On a Riemannian manifold, in the neighbourhood of a point $p \in M$ we can always find coordinates such that

$$g_{\mu\nu}(p) = \delta_{\mu\nu}, \quad \text{and} \quad \partial_\rho g_{\mu\nu}(p) = 0. \quad (4.82)$$

The same is true for Lorentzian manifolds with $\delta \rightarrow \eta$. These coordinates are known as *normal coordinates*. Since the first derivative of the metric vanishes at p it implies that the Christoffel symbols vanish there: $\Gamma^\mu_{\nu\rho}(p) = 0$. As we move away from p this does not need to continue to hold. It should be noted that one cannot generically make the second derivative of the metric vanish at p , it is only the first derivative. This means that it is not possible to pick the Riemann tensor to vanish at a given point.

We can brute force this. Start with a metric $\tilde{g}_{\mu\nu}$ in coordinates \tilde{x}^μ and try to find a new set of coordinates $x^\mu(\tilde{x})$ which satisfy the required conditions. In the new coordinates we have

$$\frac{\partial \tilde{x}^\rho}{\partial x^\mu} \frac{\partial \tilde{x}^\sigma}{\partial x^\nu} \tilde{g}_{\rho\sigma} = g_{\mu\nu} = \delta_{\mu\nu}. \quad (4.83)$$

We can take the point p to be the origin of both coordinate systems. Then we can Taylor expand around the point

$$\tilde{x}^\rho = 0 + \frac{\partial \tilde{x}}{\partial x^\mu} \Big|_{x=0} x^\mu + \frac{1}{2} \frac{\partial^2 \tilde{x}^\rho}{\partial x^\mu \partial x^\nu} \Big|_{x=0} x^\mu x^\nu + \dots \quad (4.84)$$

Inserting the expansion into (4.83) together with the Taylor expansion of $\tilde{g}_{\mu\nu}$ and then we can try to solve the resulting PDEs. The first order variation implies

$$\frac{\partial \tilde{x}^\rho}{\partial x^\mu} \Big|_{x=0} \frac{\partial \tilde{x}^\rho}{\partial x^\mu} \Big|_{x=0} \tilde{g}_{\rho\sigma}(p) = \delta_{\mu\nu}. \quad (4.85)$$

We can always find $\partial \tilde{x}/\partial x$ such that this is true, there are many choices. For $\dim M = m$ there are m^2 independent coefficients of $\partial \tilde{x}/\partial x$. The equation above contains $\frac{1}{2}m(m+1)$ conditions on these, leaving us with $\frac{1}{2}m(m-1)$ parameters still to play with. Notice that this remainder is precisely the same number of components of the rotational group of $\text{SO}(m)$ or $\text{SO}(1, m-1)$ that leaves the flat metric unchanged and so it is to be expected. Next consider the second order. There are $\frac{1}{2}m^2(m+1)$ independent components of $\partial^2 \tilde{x}^\rho/\partial x^\mu \partial x^\nu$ which is the same number of components of $\partial_\rho g_{\mu\nu}$ and so we can always choose the first derivative of the metric at p to vanish. Consider now the second derivative term, requiring $\partial_\rho \partial_\sigma g_{\mu\nu} = 0$ imposes $\frac{1}{4}m^2(m+1)^2$ constraints. However the next term in the Taylor expansion is $\partial^3 \tilde{x}^\rho/\partial x^\mu \partial x^\nu \partial x^\sigma$ which has only $\frac{1}{6}m^2(m+1)(m+2)$ independent coefficients: there are not enough independent coefficients to cancel all of the terms of the second derivative. The difference is the number of ways of characterising the second derivative of the metric that cannot be undone by coordinate transformations. This is precisely the number of independent components of the Riemann tensor, this is

$$\frac{1}{4}m^2(m+1)^2 - \frac{1}{6}m^2(m+1)(m+2) = \frac{1}{12}m^2(m+1)(m-1). \quad (4.86)$$

One can explicitly construct the normal coordinates using the exponential map and geodesics flowing through the point p . One can consider all affinely parametrised geodesics through p and label the point q at a small fixed distance of the affine parameter by the coordinates of the geodesic flowing through q . One then essentially uses geodesics to construct your basis vectors. We will not consider this construction here.

The Equivalence principle Normal coordinates play an important role in GR. Any observer at a point p who parametrises their immediate surroundings using normal coordinates will experience a locally flat metric.

This is the mathematics underling the Einstein equivalence principle. Any freely falling observer, performing local experiments will not experience a gravitational field. Here free falling means following a geodesic and therefore they will use normal coordinates. The lack of gravitational field is the statement that $g_{\mu\nu}(p) = \eta_{\mu\nu}$.

There are limitations to the equivalence principle and the important word is **local**. There is a way to distinguish whether there is a gravitational field or at p . We simply compute the Riemann tensor. This depends on the second derivative of the metric and will in general be non-vanishing. However to measure the effects of the Riemann tensor one typically has to compare the result of an experiment at p with the result at a nearby point q , this is then a “non-local” observable, according to the equivalence principle.

4.4.3 Path dependence: Curvature and Torsion

Let us take a vector $Z_p \in T_p(M)$ and parallel transport it along a curve C to some point $r \in M$. In addition condition another curve C' along which we can parallel transport Z_p to q . It is then natural to ask how do the resulting vectors differ?

Let us construct our curves from two segments, generated by linearly independent vector fields X, Y and let us take $[X, Y] = 0$. (Recall that this implies that the parallelogram constructed from the vectors closes, see section 3.3.2). We take the points to be close and pick normal coordinates $x^\mu = (\tau, \sigma, 0, \dots, 0)$ so that the starting point is at $x^\mu(p) = 0$, and the tangent vectors are aligned along the coordinates $X = \frac{\partial}{\partial \tau}$ and $Y = \frac{\partial}{\partial \sigma}$. The other corner points are $x^\mu(r) = (\delta\tau, 0, 0, \dots)$, $x^\mu(s) = (0, \delta\sigma, 0, \dots)$ and $x^\mu(r) = (\delta\tau, \delta\sigma, 0, \dots)$, with $\delta\tau$ and $\delta\sigma$ small, see figure 12.

First parallel transport Z_p along X to obtain Z_q . Along the curve, therefore Z^μ satisfies

$$\frac{dZ^\mu}{d\tau} + X^\nu \Gamma^\mu_{\rho\nu} Z^\rho = 0. \quad (4.87)$$

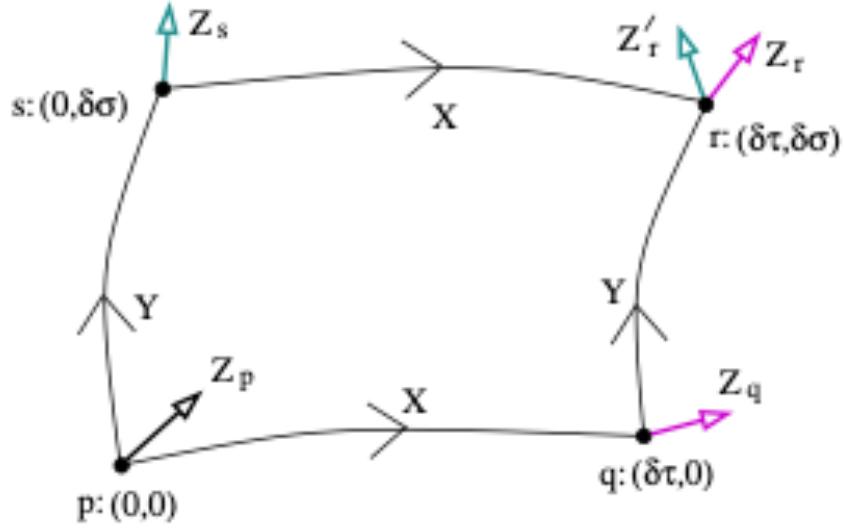


Figure 12: Parallel transporting a vector Z_p along two different paths does not give the same answer.

We can Taylor expand the solution as

$$Z_q^\mu = Z_p^\mu + \frac{dZ^\mu}{d\tau} \Big|_p \delta\tau_{\frac{1}{2}} \frac{d^2Z^\mu}{d\tau^2} \Big|_p \delta\tau^2 + \mathcal{O}(\delta\tau^3). \quad (4.88)$$

Using normal coordinates we have $\Gamma^\mu_{\rho\nu}(p) = 0$ and therefore $\frac{dZ^\mu}{d\tau} \Big|_p = 0$. To calculate the second derivative we differentiate (4.87), to obtain

$$\begin{aligned} \frac{dZ^\mu}{d\tau^2} \Big|_{\tau=0} &= - \left(X^\nu Z^\rho \frac{d\Gamma^\mu_{\rho\nu}}{d\tau} + \frac{dX^\nu}{d\tau} Z^\rho \Gamma^\mu_{\rho\nu} + X^\nu \frac{dZ^\rho}{d\tau} \Gamma^\mu_{\rho\nu} \right) \Big|_p \\ &= -X^\nu Z^\rho \frac{d\Gamma^\mu_{\rho\nu}}{d\tau} \Big|_p \\ &= -X^\nu X^\sigma Z^\rho \partial_\sigma \Gamma^\mu_{\rho\nu} \Big|_p \end{aligned} \quad (4.89)$$

To get to the second line we have used that we are working in normal coordinates at p and the final line is because τ parametrises the integral curve of X . We find

$$Z_q^\mu = Z_p^\mu - \frac{1}{2} X^\nu X^\sigma Z^\rho \partial_\sigma \Gamma^\mu_{\rho\nu} \Big|_p \delta\tau^2 + \dots \quad (4.90)$$

Now we parallel transport again, this time along Y to Z_r^μ . The Taylor expansion is

$$Z_r^\mu = Z_q^\mu + \frac{dZ^\mu}{d\sigma} \Big|_q \delta\sigma + \frac{1}{2} \frac{d^2Z^\mu}{d\sigma^2} \Big|_q \delta\sigma^2 + \mathcal{O}(\delta\sigma^3). \quad (4.91)$$

We can evaluate the first derivative $\frac{dZ^\mu}{d\sigma}|_q$ using the analogue of the parallel transport equation (4.87),

$$\frac{dZ^\mu}{d\sigma}|_q = -Y^\nu Z^\rho \Gamma^\mu_{\rho\nu}|_q, \quad (4.92)$$

however since our normal coordinates are at p and not q we cannot argue that this term immediately vanish, instead we can Taylor expand about p to get

$$Y^\nu Z^\rho \Gamma^\mu_{\rho\nu}|_q = Y^\nu Z^\rho X^\sigma \partial_\sigma \Gamma^\mu_{\rho\nu}|_q \delta\tau + \dots \quad (4.93)$$

One should also expand Y^ν and Z^ν however to leading order they multiply $\Gamma^\mu_{\rho\nu}(p) = 0$ ergo, only contribute at the next order. For the second order term in the Taylor expansion (4.91) there is a similar expression to before, we find

$$\begin{aligned} \frac{d^2 Z^\mu}{d\sigma^2}|_q &= -Y^\nu Y^\sigma Z^\rho \partial_\sigma \Gamma^\mu_{\rho\nu}|_q + \dots \\ &= -Y^\nu Y^\sigma Z^\rho \partial_\sigma \Gamma^\mu_{\rho\nu}|_p + \dots \end{aligned} \quad (4.94)$$

After the dust settles we have

$$Z_r^\mu = Z_q^\mu - Y^\nu Z^\rho X^\sigma \partial_\sigma \Gamma^\mu_{\rho\nu}|_p \delta\tau \delta\sigma - \frac{1}{2} Y^\nu Y^\sigma Z^\rho \partial_\sigma \Gamma^\mu_{\rho\nu}|_p \delta\sigma^2 + \dots \quad (4.95)$$

and therefore

$$Z_r^\mu = Z_p^\mu - \frac{1}{2} \partial_\sigma \Gamma^\mu_{\rho\nu}|_p \left[X^\nu X^\sigma Z^\rho \delta\tau^2 + 2Y^\nu Z^\rho X^\sigma \delta\sigma \delta\tau + Y^\nu Y^\sigma Z^\rho \delta\sigma^2 \right]|_p + \dots \quad (4.96)$$

with ... cubic and higher terms. We can now consider the same computation for the path C' . We merely need to swap the role of $\tau \leftrightarrow \sigma$ and $X \leftrightarrow Y$, so that

$$Z_r'^\mu = Z_p^\mu - \frac{1}{2} \partial_\sigma \Gamma^\mu_{\rho\nu}|_p \left[X^\nu X^\sigma Z^\rho \delta\tau^2 + 2X^\nu Z^\rho Y^\sigma \delta\sigma \delta\tau + Y^\nu Y^\sigma Z^\rho \delta\sigma^2 \right]|_p + \dots \quad (4.97)$$

and therefore

$$\begin{aligned} \Delta Z_r^\mu &= Z_r^\mu - Z_r'^\mu = - \left(\partial_\sigma \Gamma^\mu_{\rho\nu} - \partial_\nu \Gamma^\mu_{\rho\sigma} \right)|_p (Y^\nu Z^\rho X^\sigma)|_p \delta\sigma \delta\tau + \dots \\ &= R^\mu_{\rho\sigma\nu} Y^\nu Z^\rho X^\sigma|_p \delta\sigma \delta\tau. \end{aligned} \quad (4.98)$$

The final expression follows from the Riemann tensor expression in normal coordinates. Although our calculation was performed in a certain choice of coordinates since the end result is an equality between tensors it must hold in any coordinate system. This is a common trick, normal coordinates generally simplify expressions.

The Riemann tensor tells us the path dependence of parallel transport. This is related to the concept of *holonomy*. If we transport a vector around a closed loop we can ask how

it compares to the original vector. This is captured by the Riemann tensor. A particularly easy example is to consider a two-sphere. We can draw a loop by considering the intersection of three great circles. First go along the equator by $1/4$ of the circumference. Then make a $\pi/2$ turn and head to the north pole. At the north pole go south on another $\pi/2$ angle. You will end up with a triangle with angle $3\pi/2$. Now consider parallel transporting a vector along this loop. You will see that it changes direction when you get back to the start. Of course one could take any path and the direction you end up facing depends on the path. The set of all possible transformations of the vector at p along loops form a group known as the *holonomy group*. For a Riemannian manifold with a metric this is a subgroup of $\text{SO}(m)$ while for a Lorentzian manifold it is a subgroup of $\text{SO}(1, m - 1)$.

The meaning of Torsion Torsion will not play a role in GR for us but for completeness let us understand what is its geometric meaning.

Take two vectors $X, Y \in T_p(M)$ and let us use coordinates x^μ . Starting at $p \in M$ we can use these vectors to construct two points infinitesimally close to p , let them be r and s respectively:

$$r : x^\mu + \epsilon X^\mu \quad \text{and} \quad s : x^\mu + \epsilon Y^\mu. \quad (4.99)$$

We can now parallel transport X along Y to give a new vector $X' \in T_s(M)$ and similarly parallel transport Y along X to get a new vector $Y' \in T_r(M)$. The new vectors have components

$$X' = (X^\mu - \epsilon \Gamma^\mu_{\nu\rho} Y^\nu X^\rho) \partial_\mu, \quad Y' = (Y^\mu - \epsilon \Gamma^\mu_{\nu\rho} X^\nu Y^\rho) \partial_\mu. \quad (4.100)$$

Each now defines a new point. Starting from s and moving in the direction X' we get a new point

$$q : x^\mu + (X^\mu + Y^\mu)\epsilon - \epsilon^2 \Gamma^\mu_{\nu\rho} Y^\nu X^\rho. \quad (4.101)$$

Similarly if we sit at r and move in the direction of Y' we get to a typically different point t with coordinates

$$t : x^\mu + (X^\mu + Y^\mu)\epsilon - \epsilon^2 \Gamma^\mu_{\nu\rho} X^\nu Y^\rho. \quad (4.102)$$

The two points are not the same when $\Gamma^\mu_{\nu\rho} \neq \Gamma^\mu_{\rho\nu}$, i.e. when the connection has torsion. Torsion measures the failure of the parallelogram in figure 13 to close.

4.4.4 Geodesic deviation

Consider a one-parameter family of geodesics with coordinates $x^\mu(\tau : s)$. τ is the affine parameter along the geodesics, all of which are tangent to the vector field X . Thus, along

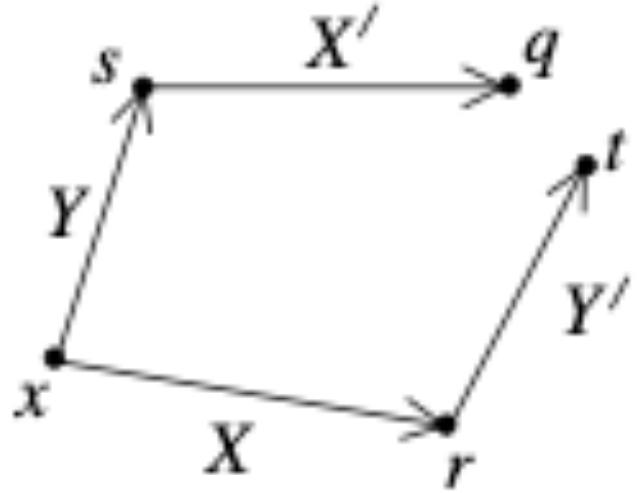


Figure 13: We transport the two vectors X and Y along each other. The failure for the parallelogram to close is measured by the torsion of the connection.

the surface spanned by $x^\mu(\tau : s)$ we have

$$\frac{\partial x^\mu}{\partial \tau} \Big|_s = X^\mu. \quad (4.103)$$

The parameter s labels the different geodesics, see figure 14. We can compute the tangent vector in the s direction to be generated by a second vector field S so that

$$S^\mu = \frac{\partial x^\mu}{\partial s} \Big|_\tau. \quad (4.104)$$

This tangent vector is known as the *deviation vector*, its job is to take us from one geodesic to a nearby geodesic with the same affine parameter τ .

The family of geodesics sweep out a 2d surface embedded in the manifold. We have freedom to choose coordinates so that on the surface $S = \frac{\partial}{\partial s}$ and $X = \frac{\partial}{\partial \tau}$ and $[X, S] = 0$.

We can ask how neighbouring geodesics behave, do they converge, diverge, remain the same distance apart? Consider a torsion free connection so that

$$\nabla_X S - \nabla_S X = [X, S]. \quad (4.105)$$

Since $[X, S] = 0$, we have

$$\nabla_X \nabla_X S = \nabla_S \nabla_S X = \nabla_S \nabla_X X + R(X, S)X, \quad (4.106)$$

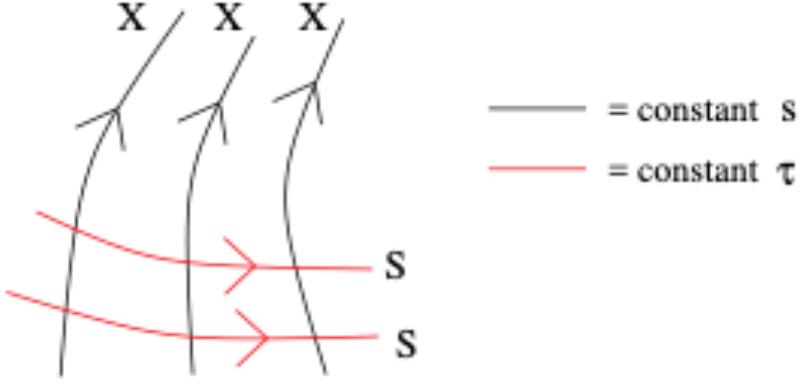


Figure 14: The black lines are geodesics generated by X while the red lines label constant τ and are generated by S with $[X, S] = 0$.

where we have used the expression for the Riemann tensor in (4.54). Since X is tangent to geodesics we have $\nabla_X X = 0$ and therefore

$$\nabla_X \nabla_X S = R(X, S)X. \quad (4.107)$$

In index notation we have

$$X^\nu \nabla_\nu (X^\rho \nabla_\rho S^\mu) = R^\mu_{\nu\rho\sigma} X^\nu X^\rho S^\sigma. \quad (4.108)$$

If we take an integral curve γ associated to X as before we have

$$\frac{D^2 S^\mu}{D\tau^2} = R^\mu_{\nu\rho\sigma} X^\nu X^\rho S^\sigma, \quad (4.109)$$

with $D/D\tau$ the covariant derivative along the curve γ , $D/D\tau = \frac{\partial x^\mu}{\partial \tau} \nabla_\mu$. The left hand side tells us how the deviation vector S changes as we move along the geodesic. It measures the relative acceleration of neighbouring geodesics. Relative acceleration is controlled by the Riemann tensor. Experimentally one observes this through *tidal forces*.

4.5 Riemann tensor and its symmetries

The components of the Riemann tensor are given in (4.56). It is not hard to see that it is anti-symmetric in the final two indices:

$$R^\sigma_{\rho\mu\nu} = -R^\sigma_{\rho\nu\mu}. \quad (4.110)$$

This does not exhaust the symmetries however. If we lower an index then we have

$$R_{\mu\nu\rho\sigma} = R_{\sigma\rho\mu\nu}, \quad (4.111)$$

$$R_{\mu[\nu\rho\sigma]} = 0, \quad (4.112)$$

$$\nabla_{[\mu} R_{\sigma\rho]\tau\nu} = 0. \quad (4.113)$$

These expressions can be proven by using normal coordinates.

4.5.1 Ricci and Einstein tensors

Given a rank (1, 3) tensor we can construct a rank (0, 2) tensor by contraction, for the Riemann tensor the resultant (0, 2)-rank tensor is called the *Ricci* tensor and is defined by

$$R_{\mu\nu} = R^{\rho}_{\mu\rho\nu}. \quad (4.114)$$

It inherits symmetry in its indices from the properties of the Riemann tensor

$$R_{\mu\nu} = R_{\nu\mu}. \quad (4.115)$$

We can create a scalar by contracting the indices again

$$R = g^{\mu\nu} R_{\mu\nu}. \quad (4.116)$$

The Bianchi identity implies that

$$\nabla^\mu \left(R_{\mu\nu} - \frac{1}{2} R g_{\mu\nu} \right) = 0, \quad (4.117)$$

which motivates us to define the covariantly constant tensor

$$G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2} R g_{\mu\nu}, \quad (4.118)$$

called the *Einstein tensor*. This will appear when we consider GR in the next section.

5 Einstein's equations

After defining all this mathematics we can now use it to introduce general relativity. Like the other forces, gravity is also mediated by some field, in this case it is the metric $g_{\mu\nu}$. It is a dynamical object, not something fixed and therefore there must be some rules as to how it can evolve. These are provided by the equations of motion following from the *Einstein–Hilbert* action.

5.1 The Einstein–Hilbert action

We want to write down an action for the gravity. Differential geometry places some rigid constraints on what this can be. We want the action to be diffeomorphism invariant, the physics should not depend on the choice of coordinates. This then implies that this is something intrinsic about the metric.

Spacetime is a manifold M equipped with a metric of Lorentzian signature. The action is an integral over M and so we require a volume-form. Thankfully the metric provides us with a canonical volume-form with which we can integrate any scalar. Given that we only have a metric there is not really much that we can do. The simplest non-trivial scalar we can construct is the Ricci scalar, and therefore we can guess the action

$$S_{\text{EH}} = \int d^4x \sqrt{-g} R. \quad (5.1)$$

As a quick check since the Ricci scalar takes the form $R \sim \partial\Gamma + \Gamma\Gamma$ and the Levi–Civita connection is $\Gamma \sim \partial g$ it follows that the action is second derivative in the metric. This is like all other actions that we have considered previously.

The equations of motion will follow from varying the action. We start with a fixed metric and see how the action varies as we shift

$$g_{\mu\nu}(x) \rightarrow g_{\mu\nu}(x) + \delta g_{\mu\nu}(x). \quad (5.2)$$

Writing the Ricci scalar as $R = g^{\mu\nu} R_{\mu\nu}$ the Einstein–Hilbert action changes as

$$\delta S = \int d^4x \left((\delta\sqrt{-g}) g^{\mu\nu} R^{\mu\nu} + \sqrt{-g} (\delta g^{\mu\nu}) R_{\mu\nu} + \sqrt{-g} g^{\mu\nu} \delta R_{\mu\nu} \right). \quad (5.3)$$

It turns out that it is simpler to consider the variation with respect to the inverse metric, this is of course equivalent to considering the variation with the metric since

$$g_{\mu\nu} g^{\nu\rho} = \delta_\mu^\rho, \quad \Rightarrow \quad (\delta g_{\mu\nu}) g^{\nu\rho} + g_{\mu\nu} \delta g^{\nu\rho} = 0, \quad \Rightarrow \quad \delta g^{\nu\rho} = -g^{\nu\sigma} g^{\rho\mu} \delta g_{\sigma\mu}. \quad (5.4)$$

The second term in the variation of the Einstein–Hilbert action is already proportional to $\delta g^{\mu\nu}$, we now want to consider the first and third terms. Let us first consider the variation of the determinant term. We want to show that

$$\delta\sqrt{-g} = -\frac{1}{2}\sqrt{-g} g_{\mu\nu} \delta g^{\mu\nu}. \quad (5.5)$$

To do this we must remember a few properties of a diagonalisable matrix A , namely

$$\log \det A = \text{tr} \log A. \quad (5.6)$$

(To prove this use that this is clearly true for a diagonal matrix since the determinant is the product of the eigenvalues while the trace is the sum of the eigenvalues. Since both the determinant and trace are invariant under conjugation it follows for any diagonalisable matrix.) Thus we have

$$\frac{1}{\det A} \delta \det A = \text{tr}(A^{-1} \delta A). \quad (5.7)$$

Applying this to the metric we have

$$\delta \sqrt{-g} = \frac{1}{2\sqrt{-g}} (-g) g^{\mu\nu} \delta g_{\mu\nu} = \frac{1}{2} \sqrt{-g} g^{\mu\nu} \delta g_{\mu\nu}. \quad (5.8)$$

Using the identity (5.4) we have that

$$\delta \sqrt{-g} = -\frac{1}{2} \sqrt{-g} g_{\mu\nu} \delta g^{\mu\nu}, \quad (5.9)$$

as claimed.

With this the variation of the Einstein–Hilbert action takes the form

$$\delta S = \int d^4x \sqrt{-g} \left(R_{\mu\nu} - \frac{1}{2} R g_{\mu\nu} \right) \delta g^{\mu\nu} + \sqrt{-g} g^{\mu\nu} \delta R_{\mu\nu}. \quad (5.10)$$

It remains to consider the final term. We claim that this term is a total derivative and can therefore be neglected by using Stoke's theorem under suitable assumptions of spacetime (no boundary). We wish to prove

$$\delta R_{\mu\nu} = \nabla_\rho \delta \Gamma^\rho_{\mu\nu} - \nabla_\nu \delta \Gamma^\rho_{\mu\rho}, \quad (5.11)$$

where

$$\delta \Gamma^\rho_{\mu\nu} = \frac{1}{2} g^{\rho\sigma} (\nabla_\mu \delta g_{\sigma\nu} + \nabla_\nu \delta g_{\mu\sigma} - \nabla_\sigma \delta g_{\mu\nu}). \quad (5.12)$$

We start by looking at the variation of the Christoffel symbols $\Gamma^\rho_{\mu\nu}$. Though the Christoffel symbol is not a tensor the variation $\delta \Gamma^\rho_{\mu\nu}$ is a tensor. This is because it is the difference of two Christoffel symbols, one computed using the metric $g_{\mu\nu}$ and one with $g_{\mu\nu} + \delta g_{\mu\nu}$ and the term in the transformation of the Christoffel which shows that it is not a tensor is independent of the metric and therefore cancels in the difference. This observation makes our lives a lot simpler. It implies that at any point $p \in M$ we can work in normal coordinates such that $\partial_\rho g_{\mu\nu} = 0$ and therefore $\Gamma^\rho_{\mu\nu} = 0$. To linear order in the variation the change in the Christoffel symbol evaluated at p is

$$\begin{aligned} \delta \Gamma^\rho_{\mu\nu} &= \frac{1}{2} g^{\rho\sigma} (\partial_\mu \delta g_{\sigma\nu} + \partial_\nu \delta g_{\sigma\mu} - \partial_\sigma \delta g_{\mu\nu}) \\ &= \frac{1}{2} g^{\rho\sigma} (\nabla_\mu \delta g_{\sigma\nu} + \nabla_\nu \delta g_{\sigma\mu} - \nabla_\sigma \delta g_{\mu\nu}) \end{aligned} \quad (5.13)$$

where we have used that in normal coordinates we can replace partial derivatives with covariant derivatives. Both the left and right hand side are tensors and therefore this holds in any coordinate system, moreover the point p was arbitrary and therefore this holds in all coordinate systems at all points $p \in M$.

Next consider the variation of the Riemann tensor. In normal coordinates we have

$$R^\sigma_{\rho\mu\nu} = \partial_\mu \Gamma^\sigma_{\nu\rho} - \partial_\nu \Gamma^\sigma_{\mu\rho}, \quad (5.14)$$

and the variation is

$$\delta R^\sigma_{\rho\mu\nu} = \partial_\mu \delta \Gamma^\sigma_{\nu\rho} - \partial_\nu \delta \Gamma^\sigma_{\mu\rho} = \nabla_\mu \delta \Gamma^\sigma_{\nu\rho} - \nabla_\nu \delta \Gamma^\sigma_{\mu\rho}, \quad (5.15)$$

where we have once again used that in normal coordinates we can replace partial derivatives with covariant derivatives. As before we have a tensorial equation and therefore this must hold in any coordinate system not just normal coordinates. We have

$$\delta R_{\rho\nu} = \nabla_\mu \delta \Gamma^\mu_{\nu\rho} - \nabla_\nu \delta \Gamma^\mu_{\mu\rho}. \quad (5.16)$$

It follows that

$$g^{\mu\nu} \delta R_{\mu\nu} = \nabla_\mu \left(g^{\rho\nu} \delta \Gamma^\mu_{\rho\nu} - g^{\mu\nu} \delta \Gamma^\rho_{\nu\rho} \right) = \nabla_\mu X^\mu \quad (5.17)$$

The variation of the Einstein–Hilbert action can then be written as

$$\delta S = \int d^4x \sqrt{-g} \left[\left(R_{\mu\nu} - \frac{1}{2} R g_{\mu\nu} \right) \delta g^{\mu\nu} + \nabla_\mu X^\mu \right]. \quad (5.18)$$

The final term is a total derivative after using the identity (4.72) and with suitable assumptions on spacetime can be neglected. Requiring that the action is extremised, so that $\delta S = 0$ we have the equations of motion

$$G_{\mu\nu} := R_{\mu\nu} - \frac{1}{2} R g_{\mu\nu} = 0. \quad (5.19)$$

These are the *Einstein field equations* in the absence of matter. We may further simplify them by first contracting with the inverse metric to find $R = 0$ and therefore in the absence of matter Einstein’s equations are simply

$$R_{\mu\nu} = 0. \quad (5.20)$$

Though this looks deceptively simple this holds a very rich set of solutions, in fact not all solutions to this equation are known.

We threw away the boundary term in the usual cavalier way one does with such variational principles. One can introduce the Gibbons–Hawking boundary term to allow for M to admit a boundary.

5.1.1 Newton's constant

As it stands the action we have given does not have the correct dimension. At this stage where we do not couple to matter this is not a problem however we wish to be able to couple to matter soon and therefore we must fix this. We take the coordinates to have dimension of length and therefore the metric is dimensionless. The Ricci scalar involves two derivatives and therefore it has dimension $[R] = L^{-2}$. Including the dimension of the integration measure the current action in (5.1) has dimension $[S] = L^2$. An action should have dimension of Energy \times time and therefore we should multiply the action by an appropriate dimensionful constant.

We take

$$S_{\text{EH}} = \frac{c^3}{16\pi G_N} \int d^4x \sqrt{-g} R, \quad (5.21)$$

where c is of course the speed of light, and G is Newton's constant

$$G_N c \sim 6.67 \times 10^{-11} m^3 kg^{-1} s^{-2}. \quad (5.22)$$

This will not change the equation of motion in the vacuum but once we couple matter will determine the strength of the coupling of the gravitational field to matter.

If we are just interested in phenomena related to gravity it is sensible to set $G_N = 1$. Instead if we want to consider other phenomena other than gravity this is not so sensible since it defines the coupling of the forces. Instead the more useful convention is to pick $\hbar = 1$, which equates energy with time. With this convention Newton's constant has dimension $[G] = m^{-2}$.

The corresponding energy scale is called the *Planck mass* and is given by

$$M_{pl}^2 = \frac{\hbar c}{8\pi G_N}. \quad (5.23)$$

It is around 10^{18} GeV which is a very high energy scale and far beyond anything we can probe experimentally. This is why the gravitational force is so weak.

5.1.2 Cosmological constant

We motivated the Einstein–Hilbert action as the simplest action one can write down. There is in fact a simpler term we may write down other than the Einstein–Hilbert term we considered previously. We may simply add a constant to the volume form. The resulting action is

$$S = \frac{1}{16\pi G_N} \int d^4x \sqrt{-g} (R - 2\Lambda). \quad (5.24)$$

The constant Λ is known as the *cosmological constant* and has dimension $[\Lambda] = L^{-2}$. The minus sign in the action comes from thinking of the Lagrangian as $T - V$ with the cosmological constant playing the role of the potential energy V .

Varying the action as before yields the Einstein equations

$$R_{\mu\nu} - \frac{1}{2}Rg_{\mu\nu} = -\Lambda g_{\mu\nu}. \quad (5.25)$$

This time if we contract with the inverse metric we get $R = 4\Lambda$. Substituting this back in the vacuum Einstein equations in the presence of a cosmological constant becomes

$$R_{\mu\nu} = \Lambda g_{\mu\nu}. \quad (5.26)$$

5.1.3 Higher derivative terms

The Einstein–Hilbert action with cosmological constant is the simplest thing we can write down. We may however construct other scalars from the Riemann tensor, they will however have higher derivative terms. For example there are three terms that we can add at four-derivative in the metric

$$S_{4-\text{deriv}} = \int d^4x \sqrt{-g} \left(c_1 R^2 + c_2 R_{\mu\nu} R^{\mu\nu} + c_3 R_{\mu\nu\rho\sigma} R^{\mu\nu\rho\sigma} \right), \quad (5.27)$$

with the c_i dimensionless constant. Generic choices of the constants will not give rise to higher derivative equations of motion with a well-defined initial value problem. Nonetheless there are certain combinations which conspire to keep the equations second order in derivatives. This goes by the name of *Lovelock’s theorem* and says that in four-dimensions the combination

$$\frac{1}{8\pi^2} \int_M d^4x \sqrt{g} (R^2 - 4R_{\mu\nu} R^{\mu\nu} + R_{\mu\nu\rho\sigma} R^{\mu\nu\rho\sigma}) = \chi(M), \quad (5.28)$$

where $\chi(M) \in \mathbb{Z}$, is the Euler character of M . In Lorentzian signature this is also a total derivative and therefore does not affect the classical equations of motion. Higher derivative terms only become relevant for fast varying fields. For us these will not be important and therefore we stick to the 2-derivative action.

5.1.4 Diffeomorphisms

A natural question to ask is how many degrees of freedom are there in the metric? Since it is a 4×4 symmetric matrix the naive guess is $\frac{1}{2} \times 4 \times 5 = 10$ however this is not quite correct. Not all of these 10 components are physical. Two metrics which are related by a change of coordinates $x^\mu \rightarrow \tilde{x}^\mu(x)$ describe the same physical spacetime. This means that there is a redundancy in any given representation of the metric which removes precisely 4 of the 10 degrees of freedom, leaving just 6 actual degrees of freedom.

This redundancy is implemented by diffeomorphisms. Recall that a diffeomorphism is a map $\phi : M \rightarrow M$. We may use it to map all fields, including the metric on M to a new set

of fields on M . The end result is physically indistinguishable from the original, it describes the same system just in a different set of coordinates. Such diffeomorphisms are analogous to the gauge transformations of a gauge theory, think Maxwell theory.

Let us look at how diffeomorphisms modify the action. Consider a diffeomorphism which takes a point with coordinate x^μ to a nearby point with coordinates

$$x^\mu \rightarrow \tilde{x}^\mu = x^\mu + \delta x^\mu. \quad (5.29)$$

We can view this either as an active change in which one point with coordinates x^μ is mapped to another point with coordinates $x^\mu + \delta x^\mu$ or as a passive transformation in which we use two different coordinate patches to label the same point. Either viewpoint leads to the same conclusion, here we will take the passive viewpoint.

We can think of the change of coordinates as being generated by an infinitesimal vector field X ,

$$\delta x^\mu = -X^\mu(x). \quad (5.30)$$

The metric transforms as

$$g_{\mu\nu}(x) \rightarrow \tilde{g}_{\mu\nu}(\tilde{x}) = \frac{\partial x^\rho}{\partial \tilde{x}^\mu} \frac{\partial x^\sigma}{\partial \tilde{x}^\nu} g_{\rho\sigma}(x). \quad (5.31)$$

We can invert the Jacobian matrix to find

$$\frac{\partial \tilde{x}^\mu}{\partial x^\rho} = \delta_\rho^\mu - \partial_\rho X^\mu \quad \Rightarrow \quad \frac{\partial x^\rho}{\partial \tilde{x}^\mu} = \delta_\mu^\rho + \partial_\mu X^\rho, \quad (5.32)$$

where the inverse holds to leading order in the variation X . Continuing to work infinitesimally we have

$$\begin{aligned} \tilde{g}_{\mu\nu}(\tilde{x}) &= (\delta_\mu^\rho + \partial_\mu X^\rho)(\delta_\nu^\sigma + \partial_\nu X^\sigma)g_{\rho\sigma}(x) \\ &= g_{\mu\nu}(x) + g_{\mu\rho}(x)\partial_\nu X^\rho + g_{\nu\rho}(x)\partial_\mu X^\rho. \end{aligned} \quad (5.33)$$

We can also Taylor expand the left-hand side to find

$$\tilde{g}_{\mu\nu}(\tilde{x}) = \tilde{g}_{\mu\nu}(x + \delta d) = \tilde{g}_{\mu\nu}(x) - X^\lambda \partial_\lambda \tilde{g}_{\mu\nu}(x). \quad (5.34)$$

Comparing the different metrics at the same point x we find that the metric undergoes the infinitesimal change

$$\delta g_{\mu\nu}(x) = \tilde{g}_{\mu\nu}(x) - g_{\mu\nu}(x) = X^\lambda \partial_\lambda g_{\mu\nu} + g_{\mu\rho}\partial_\nu X^\rho + g_{\nu\rho}\partial_\mu X^\rho. \quad (5.35)$$

This is precisely the Lie derivative of the metric. If we act with an infinitesimal diffeomorphism along X then the metric changes as

$$\delta g_{\mu\nu} = (\mathcal{L}_X g)_{\mu\nu}. \quad (5.36)$$

We may also rewrite this by lowering the index on X^ρ to find

$$\delta g_{\mu\nu} = \partial_\mu X_\nu + \partial_\nu X_\mu + X^\rho (\partial_\rho g_{\mu\nu} - \partial_\mu g_{\rho\nu} - \partial_\nu g_{\mu\rho}), \quad (5.37)$$

the last term is just the Christoffel symbols and therefore we have

$$\delta g_{\mu\nu} = \nabla_\mu X_\nu + \nabla_\nu X_\mu. \quad (5.38)$$

We may put this together to see how the action changes. Under a general change of the metric the Einstein–Hilbert action changes as

$$\delta S = \int d^4x \sqrt{-g} G^{\mu\nu} \delta g_{\mu\nu}, \quad (5.39)$$

where we have discarded the boundary term. Insisting that $\delta S = 0$ for any variation $\delta g_{\mu\nu}$ gives the equations of motion $G^{\mu\nu} = 0$. In contrast, symmetries of the action are those variations $\delta g_{\mu\nu}$ for which $\delta S = 0$ for any choice of metric. Since diffeomorphisms are symmetries we know that the action is invariant under changes of the form (5.38). Using the fact that $G_{\mu\nu}$ is symmetric we must have

$$\delta S = 2 \int d^4x \sqrt{-g} G^{\mu\nu} \nabla_\mu X_\nu = 0, \quad \text{for all } X_\mu(x). \quad (5.40)$$

After integrating by parts we find that this results in the Bianchi identity

$$\nabla_\mu G^{\mu\nu} = 0. \quad (5.41)$$

We learn that from the path integral perspective the Bianchi identity is a result of diffeomorphism invariance.

5.1.5 Coupling to matter

Until now the action has only involved gravity, and at most we can allow for test particles moving on geodesics. However matter is not just an actor doing what gravity says in spacetime, it also backreacts and affects the dynamics of spacetime. The first question to ask is how does matter couple to the metric? Let us take matter which is described by a Lagrangian.

Scalar Field Consider first a scalar field $\phi(x)$. In flat spacetime the action takes the form

$$S_{\text{scalar}} = \int d^4x \left(-\frac{1}{2} \eta^{\mu\nu} \partial_\mu \phi \partial_\nu \phi - V(\phi) \right), \quad (5.42)$$

with $\eta^{\mu\nu}$ the inverse Minkowski metric.¹⁹

¹⁹Note that the minus sign is due to our mostly plus signature convention, you may be more used to the opposite convention when considering a field theory. The Lagrangian will take the form of kinetic energy minus potential energy.

It is straightforward to generalise this to describe a field moving in curved spacetime, we simply need to replace the Minkowski metric with the curved metric, replace partial derivatives with covariant derivatives and introduce the volume form when we integrate in the action. This means that we take

$$S_{\text{scalar}} = \int d^4x \sqrt{-g} \left(-\frac{1}{2} g^{\mu\nu} \nabla_\mu \phi \nabla_\nu \phi - V(\phi) \right). \quad (5.43)$$

Despite upgrading the partial derivatives to covariant ones this is somewhat redundant here as they act the same on a scalar field: we keep it for later though.

Curved spacetime also introduces new possibilities for us to add to the action, for example we could add a term such as $\xi R\phi^2$ to the action which gives rise to extra couplings. We will not interest ourselves in such terms here however.

Maxwell Theory The action of Maxwell theory from special relativity is

$$S_{\text{Maxwell}} = -\frac{1}{4} \int d^4x \eta^{\mu\rho} \eta^{\nu\sigma} F_{\mu\nu} F_{\rho\sigma}, \quad (5.44)$$

with $F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu$. The electric and magnetic fields are encoded in F via

$$F_{\mu\nu} = \begin{pmatrix} 0 & -E_1 & -E_2 & -E_3 \\ E_1 & 0 & B_3 & -B_2 \\ E_2 & -B_3 & 0 & B_1 \\ E_3 & B_2 & -B_1 & 0 \end{pmatrix}, \quad (5.45)$$

and the Bianchi identity $dF = d^2A = 0$ yields two of the four Maxwell equations

$$\nabla \cdot \vec{B} = 0, \quad \nabla \times \vec{B} + \frac{\partial \vec{B}}{\partial t} = 0. \quad (5.46)$$

We may couple to curved space time through the minimal coupling outlined for the scalar theory. The action is

$$S_{\text{Maxwell}} = -\frac{1}{4} \int d^4x \sqrt{-g} g^{\mu\rho} g^{\nu\sigma} F_{\mu\nu} F_{\rho\sigma} = -\frac{1}{2} \int F \wedge \star F. \quad (5.47)$$

We again take $F = dA$, which in components reads $F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu = \nabla_\mu A_\nu - \nabla_\nu A_\mu$. Antisymmetry implies that we may replace the covariant derivatives with normal derivatives. The equations of motion are

$$\nabla^\mu F_{\mu\nu} = 0, \quad \Leftrightarrow \quad d \star F = 0. \quad (5.48)$$

We have now seen how to couple matter to gravity but how does the change the Einstein equations of the previous section. We need to consider the combined action

$$S = \frac{1}{16\pi G_N} \int d^4x \sqrt{-g} (R - 2\Lambda) + S_{\text{Matter}}, \quad (5.49)$$

where S_{Matter} is the action for any matter fields in the theory minimally coupled to gravity. When we vary the Einstein–Hilbert term we know that we will obtain the Einstein tensor, what about S_{Matter} ? We define the *Energy-Momentum tensor* to be

$$T_{\mu\nu} = -\frac{2}{\sqrt{-g}} \frac{\delta \mathcal{L}_{\text{Matter}}}{\delta g^{\mu\nu}}. \quad (5.50)$$

By construction $T_{\mu\nu}$ is symmetric. Varying the full action with respect to the metric we have

$$\delta S = \frac{1}{16\pi G_N} \int d^4x \sqrt{-g} (G_{\mu\nu} + \Lambda g_{\mu\nu}) \delta g^{\mu\nu} - \frac{1}{2} \int d^4x \sqrt{-g} T_{\mu\nu} \delta g^{\mu\nu}, \quad (5.51)$$

from which we may read the following equation of motion

$$G_{\mu\nu} + \Lambda g_{\mu\nu} = 8\pi G_N T_{\mu\nu}. \quad (5.52)$$

This constitutes the full Einstein equations describing gravity coupled to matter. Note that the presence of the energy-momentum tensor says that the matter distribution sources the curvature of the spacetime.

For the scalar theory above the energy-momentum tensor is

$$T_{\mu\nu} = \nabla_\mu \phi \nabla_\nu \phi - g_{\mu\nu} \left(\frac{1}{2} \nabla^\rho \phi \nabla_\rho \phi + V(\phi) \right). \quad (5.53)$$

If we restrict to flat space then

$$T_{00} = \frac{1}{2} \dot{\phi}^2 + \frac{1}{2} (\nabla \phi)^2 + V(\phi), \quad (5.54)$$

with ∇ the usual 3d spatial derivative. This is the energy density of a scalar field.

For the Maxwell action we have

$$T_{\mu\nu} = g^{\rho\sigma} F_{\mu\rho} F_{\nu\sigma} - \frac{1}{4} g_{\mu\nu} F^{\rho\sigma} F_{\rho\sigma}. \quad (5.55)$$

In flat space we have

$$T_{00} = \frac{1}{2} \left[\vec{E}^2 + \vec{B}^2 \right]. \quad (5.56)$$

This is the energy density of the magnetic and electric fields.

5.2 Newtonian gravity as a limit

We now want to see that this reduces correctly to Newtonian gravity in some limit. We will linearise Einstein's equations and work in an approximate regime where Newtonian gravity should hold. We consider a situation where the metric is approximately flat and set the

cosmological constant to vanish $\Lambda = 0$. The weakness of the gravitational field is expressed by decomposing the metric as

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}, \quad (5.57)$$

with $h_{\mu\nu} \ll 1$ small, that is each of the components of the metric is small. This assumption allows us to ignore anything that is higher than first order in this term. This allows us to immediately write the inverse metric (to first order) as

$$g^{\mu\nu} = \eta^{\mu\nu} - h^{\mu\nu}, \quad (5.58)$$

where $h^{\mu\nu} = \eta^{\mu\rho}\eta^{\nu\sigma}h_{\rho\sigma}$. We can now raise and lower indices with η since the corrections would be of higher order in the perturbation. We can think of this linearised theory as describing a theory of a symmetric tensor field $h_{\mu\nu}$ propagating on a flat background spacetime. We could equally think of a perturbation around some other background metric then the theory is that of the symmetric tensor field propagating on a curved background.

We want to consider equations of motion for the perturbations, which come from examining Einstein's equations to linear order. To begin we should work out the Christoffel symbols which take the form

$$\begin{aligned} \Gamma^\rho_{\mu\nu} &= \frac{1}{2}g^{\rho\sigma}(\partial_\mu g_{\sigma\nu} + \partial_\nu g_{\sigma\mu} - \partial_\sigma g_{\mu\nu}) \\ &= \frac{1}{2}\eta^{\rho\sigma}(\partial_\mu h_{\sigma\nu} + \partial_\nu h_{\sigma\mu} - \partial_\sigma h_{\mu\nu}) + \mathcal{O}(h^2). \end{aligned} \quad (5.59)$$

Since the Riemann tensor is of the form $R \sim \partial\Gamma + \Gamma\Gamma$ the first order contributions will come from the derivative terms and not the ‘squared’ terms. We have

$$\begin{aligned} R^\sigma_{\rho\mu\nu} &= \partial_\mu\Gamma^\sigma_{\nu\rho} - \partial_\nu\Gamma^\sigma_{\mu\rho} + \mathcal{O}(h^2) \\ &= \frac{1}{2}\eta^{\sigma\lambda}(\partial_\mu\partial_\rho h_{\nu\lambda} - \partial_\mu\partial_\lambda h_{\nu\rho} - \partial_\nu\partial_\rho h_{\mu\lambda} + \partial_\nu\partial_\lambda h_{\mu\rho}) + \mathcal{O}(h^2). \end{aligned} \quad (5.60)$$

It follows that the Ricci tensor is

$$R_{\mu\nu} = \frac{1}{2}(\partial^\sigma\partial_\nu h_{\sigma\mu} + \partial^\sigma\partial_\mu h_{\sigma\nu} - \square h_{\mu\nu} - \partial_\mu\partial_\nu h) + \mathcal{O}(h^2), \quad (5.61)$$

where $h = h^\mu_\mu$ is the trace and $\square = \partial^\mu\partial_\mu$. Moreover the Ricci scalar is

$$R = \partial^\mu\partial^\nu h_{\mu\nu} - \square h + \mathcal{O}(h^2). \quad (5.62)$$

Putting all of this together into the Einstein tensor we end up with

$$G_{\mu\nu} = \frac{1}{2}\left[\partial^\sigma\partial_\nu h_{\mu\sigma} + \partial^\sigma\partial_\mu h_{\nu\sigma} - \square h_{\mu\nu} - \partial_\mu\partial_\nu h - \eta_{\mu\nu}(\partial^\rho\partial^\sigma h_{\rho\sigma} - \square h)\right] + \mathcal{O}(h^2). \quad (5.63)$$

This can be obtained by varying the following Lagrangian with respect to $h_{\mu\nu}$,

$$\mathcal{L} = \frac{1}{2} \left[(\partial_\mu h^{\mu\nu}) \partial_\nu h + \frac{1}{2} \partial^\mu h^{\rho\sigma} \partial_\mu h_{\rho\sigma} - \partial^\mu h^{\rho\sigma} \partial_\rho h_{\mu\sigma} + \partial^\mu h \partial_\mu h \right]. \quad (5.64)$$

The full linearised equations of motion are then

$$\frac{1}{2} \left[\partial^\sigma \partial_\nu h_{\mu\sigma} + \partial^\sigma \partial_\mu h_{\nu\sigma} - \square h_{\mu\nu} - \partial_\mu \partial_\nu h - \eta_{\mu\nu} (\partial^\rho \partial^\sigma h_{\rho\sigma} - \square h) \right] = 8\pi G_N T_{\mu\nu}, \quad (5.65)$$

where $T_{\mu\nu}$ is assumed to be small.

Before we can proceed we must deal with gauge invariance. The demand that $g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu}$ does not completely fix the coordinate system on spacetime. Let us consider an infinitesimal change of coordinates

$$x^\mu \rightarrow x^\mu - \xi^\mu \quad (5.66)$$

with ξ assumed to be small. The metric changes by

$$\delta g_{\mu\nu} = (\mathcal{L}_\xi g)_{\mu\nu} = \nabla_\mu \xi_\nu + \nabla_\nu \xi_\mu. \quad (5.67)$$

When the metric takes the linearised form this should be understood as a transformation of $h_{\mu\nu}$. Since we assume that both h and ξ are small²⁰ it follows that we may replace covariant derivatives of g with covariant derivatives of η where the Christoffel symbols vanish. We then have

$$h_{\mu\nu} \rightarrow h_{\mu\nu} + (\mathcal{L}_\xi \eta)_{\mu\nu} = h_{\mu\nu} + \partial_\mu \xi_\nu + \partial_\nu \xi_\mu. \quad (5.68)$$

For those who have seen gauge theories this is precisely the form of a gauge transformation of Maxwell theory. There we shift the one-form A as $A \rightarrow A + d\Lambda$ which leaves the field strength (or curvature of the gauge bundle) $F = dA$ invariant. Similarly the above transformation leaves the linearised Riemann tensor invariant.

When we do computations in gauge theories we typically pick a gauge to work in. The most common gauge to take is the Lorentz gauge

$$\partial^\mu A_\mu = 0, \quad (5.69)$$

which reduces the Maxwell equation $d \star F = \star J$ with source to the wave equation

$$\square A_\nu = J_\nu. \quad (5.70)$$

²⁰If we did not restrict to small ξ then we could go to a region where $h_{\mu\nu}$ is not small by a coordinate transformation, clearly we do not want this.

There is a similar kind of gauge here called *de Donder gauge*. We take

$$\partial^\mu h_{\mu\nu} - \frac{1}{2}\partial_\nu h = 0. \quad (5.71)$$

To see that this is always possible suppose that you are given a metric where

$$\partial^\mu h_{\mu\nu} - \frac{1}{2}\partial_\nu h = f_\nu, \quad (5.72)$$

then after a gauge transformation we have

$$\partial^\mu h_{\mu\nu} - \frac{1}{2}\partial_\nu h + \square\xi_\nu = f_\nu, \quad (5.73)$$

and it amounts to finding ξ such that $\square\xi_\nu = f_\nu$.

De Donder gauge greatly simplifies our linearised equations of motion

$$\square h_{\mu\nu} - \frac{1}{2}\square h\eta_{\mu\nu} = -16\pi G_N T_{\mu\nu}. \quad (5.74)$$

It is useful to define

$$\bar{h}_{\mu\nu} = h_{\mu\nu} - \frac{1}{2}h\eta_{\mu\nu}, \quad (5.75)$$

so that the linearised Einstein equation becomes

$$\square\bar{h}_{\mu\nu} = -16\pi G_N T_{\mu\nu}. \quad (5.76)$$

To see that this is a sensible definition we see that from $\bar{h}_{\mu\nu}$ we can recover $h_{\mu\nu}$ since by taking the trace on both sides we have

$$\bar{h} = \eta^{\mu\nu}\bar{h}_{\mu\nu} = -h, \quad (5.77)$$

so we can reconstruct $h_{\mu\nu}$ as

$$h_{\mu\nu} = \bar{h}_{\mu\nu} - \frac{1}{2}\bar{h}\eta_{\mu\nu}. \quad (5.78)$$

Newtonian Limit We now are in a position to take the Newtonian limit. We require a low-density slowly moving distribution of matter. We will take a stationary matter configuration so that the Energy-momentum tensor is

$$T_{00} = \rho(\vec{x}), \quad (5.79)$$

with all other components vanishing. Via the stationary assumption we may replace the wave operator \square with the 3d Euclidean Laplacian $\square = -\partial_t^2 + \partial_i^2 = \nabla^2$. Einstein's equations then become

$$\nabla^2\bar{h}_{00} = -16\pi G_N\rho(\vec{x}), \quad \nabla^2\bar{h}_{0i} = 0, \quad \nabla^2\bar{h}_{ij} = 0. \quad (5.80)$$

With suitable boundary conditions the solutions are

$$\bar{h}_{00} = -4\Phi(\vec{x}) \quad \bar{h}_{0i} = \bar{h}_{ij} = 0, \quad (5.81)$$

where Φ is identified with the Newtonian potential obeying

$$\nabla^2 \Phi(\vec{x}) = 4\pi G_N \rho(\vec{x}). \quad (5.82)$$

Translating back to $h_{\mu\nu}$ we find

$$h_{00} = -2\Phi(\vec{x}), \quad h_{ij} = -2\Phi(\vec{x})\delta_{ij}, \quad h_{0i} = 0. \quad (5.83)$$

The final metric is then

$$ds^2 = -(1 + 2\Phi(\vec{x}))dt^2 + (1 - 2\Phi(\vec{x}))d\vec{x} \cdot d\vec{x}. \quad (5.84)$$

We conclude that we can recover Newtonian gravity from general relativity and therefore this is not complete craziness. We will see soon that if we replace $\Phi(\vec{x}) = -\frac{G_N M}{r}$, as would be expected for a point mass, then this is the leading expansion of the Schwarzschild solution.

One can also study gravitational waves using the linearised equations of motion for example. This has had recent experimental interest due to the observations of gravitational waves by LIGO. Theorists have also taken an interest in these experimental results with the hope that extra precision tests of GR and its quantum gravity extension can be performed using this data.

6 Schwarzschild solution

Black holes are one of the most enigmatic objects and probably the reason why most of you are here. We will take our first steps to understanding black holes here.

6.1 The Schwarzschild black hole

In 1915 Einstein had published his work on General relativity and made a comment saying that he was not optimistic that the equations he had found could be solved other than Minkowski space. Also in 1915 with the first world war raging in Europe, Karl Schwarzschild was in the German army on the Russian front performing ballistic calculations, and suffering from pemphigus a rare and painful autoimmune disease. Despite this he worked on finding solutions to general relativity and found the first exact (non-trivial) solution to Einstein's field equations.²¹ Schwarzschild's breakthrough was to use a convenient system of coordinates,

²¹Schwarzschild died in 1916 having left military service due to his illness.

taking a polar like coordinate system as opposed to Einstein's rectangular coordinate system. The metric that bears his name is

$$ds^2 = -\left(1 - \frac{2G_NM}{r}\right)dt^2 + \left(1 - \frac{2G_NM}{r}\right)^{-1}dr^2 + r^2(d\theta^2 + \sin^2\theta d\phi^2). \quad (6.1)$$

This solves Einstein's equations in a vacuum, $R_{\mu\nu} = 0$. The coordinate ranges are²²

$$t \in \mathbb{R}, \quad 0 < \theta < \pi, \quad 0 < \phi < 2\pi. \quad (6.2)$$

The range of r is slightly more subtle. At $r = 2G_NM$ something funky is happening since the prefactor of dt^2 and dr^2 vanish or diverge respectively. For the moment we will keep $2G_NM < r < \infty$ and we are then safe. This value of the radial coordinate is called the *Schwarzschild radius* and will play a prominent role later.

It depends on a single parameter M which is interpreted as the mass of the object. Indeed using our results above on the Linearised equations and their Newtonian limit we have

$$g_{00} = -(1 + 2\Phi), \quad (6.3)$$

with Φ the Newtonian potential. For the Schwarzschild metric we have

$$\Phi = -\frac{G_NM}{r}, \quad (6.4)$$

which is the Newtonian potential for a point mass M at the origin.

We can compute the mass of the black hole by using Komar integrals. The Schwarzschild solution admits a time-like Killing vector $K = \partial_t$: a Killing vector satisfies $\mathcal{L}_K g = 0$ which is equivalent to $\nabla_{(\mu} K_{\nu)} = 0$. Then to compute the Komar integral we must construct the dual one-form

$$K = g_{00}dt = -\left(1 - \frac{2GM}{r}\right)dt. \quad (6.5)$$

The Komar integral is given by

$$M_{\text{Komar}} = -\frac{1}{8\pi G_N} \int_{S^2} \star dK, \quad (6.6)$$

where the S^2 is any sphere with a radius larger than the horizon at $r = 2G_NM$ where the Killing vector has vanishing norm. Then

$$dK = -\frac{2G_NM}{r^2}dr \wedge dt \Rightarrow \star dK = -2G_NM \sin\theta d\theta d\phi. \quad (6.7)$$

²²The are singularities at $\theta = 0, \pi$ and $\phi = 0, 2\pi$ however these are just the expected singularities from considering a two-sphere and attempting to use just one coordinate patch. We should be careful about this but it is not a problem.

and therefore

$$M_{\text{Komar}} = M. \quad (6.8)$$

Note that $d \star dK = 0$ and therefore it obeys an equation similar to Maxwell's equations $d \star F = 0$. These are Maxwell's equations in the absence of any current and therefore one would expect the electric charge to vanish. Yet this electric charge is precisely the mass and this is non-zero. For the solution the mass is localised at the origin $r = 0$ where the field strength diverges. This allows for a non-trivial value.

We may thus expect that this describes something physical only when $M > 0$. For $M = 0$ we find Minkowski space while for $M < 0$ the metric becomes unphysical.

6.1.1 Birkhoff's theorem

The Schwarzschild solution turns out to be the unique spherically symmetric asymptotically flat solution to the vacuum Einstein equations, this fact is known as *Birkhoff's theorem*. This means that the Schwarzschild solution does not just describe the spacetime outside of a black hole but outside any non-rotating, spherically symmetric object such as a star or planet. We will sketch the proof of this fact since it allows us to get a feel for solving the Einstein equations.

The spherical symmetry of the metric means that it has an $\text{SO}(3)$ isometry. If you hold up a round sphere and rotate it it looks the same no matter which way you rotate it. If instead you did the same with a golf ball, which has dimples then this rotational symmetry is broken. The distinction between these two situations should be captured by the metric. The metric on a round two-sphere will look the same wherever you sit on the sphere whereas the metric on the golf ball will depend on where you are.

To define this mathematically we need to use the concept of a flow that we introduce a number of lectures ago. A flow on a manifold M is a one-parameter family of diffeomorphisms $\sigma_t : M \rightarrow M$, and may be associated to a vector field $K \in \mathcal{X}(M)$ at each point along the flow which is tangent to the flow

$$K^\mu = \frac{dx^\mu(\lambda)}{d\lambda}. \quad (6.9)$$

The flow is said to be an isometry if the metric looks the same at each point along a given flow line, mathematically this means that an isometry satisfies

$$\mathcal{L}_K g = 0, \quad \Leftrightarrow \quad \nabla_\mu K_\nu + \nabla_\nu K_\mu = 0. \quad (6.10)$$

A vector satisfying this equation is known as a *Killing vector field*. Sometimes it is simply to see that a vector is an isometry, particularly when it is an ignorable coordinate, i.e. the

metric does not depend on said coordinate. However sometimes the Killing vectors are not so obvious.

There is a group structure underlying the symmetries, well technically a Lie algebra structure. This follows since the Lie derivative satisfies

$$\mathcal{L}_X \mathcal{L}_Y - \mathcal{L}_Y \mathcal{L}_X = \mathcal{L}_{[X,Y]}. \quad (6.11)$$

Killing vectors form a Lie algebra of the isometry group of the manifold. (See problem sheet 3 where we consider the Killing vectors on the round three-sphere).

One can then prove that the $\text{SO}(3)$ isometry implies that the metric must take the form

$$ds^2 = g_{\tau\tau}(\tau, \rho)d\tau^2 + 2g_{\tau\rho}(\tau, \rho)d\tau d\rho + g_{\rho\rho}(\tau, \rho)d\rho^2 + r^2(\tau, \rho)ds^2(S^2), \quad (6.12)$$

where

$$ds^2(S^2) = d\theta^2 + \sin^2\theta d\phi^2, \quad (6.13)$$

is the metric on a round two-sphere. The $\text{SO}(3)$ isometry then acts on the two-sphere and leaves τ and ρ untouched. This is called a *foliation* of the space by S^2 leaves.

The size of the sphere is determined by $r(\tau, \rho)$ and it is convenient to redefine the coordinates such that r is a coordinate, we can then eliminate the ρ coordinate in favour of r , the metric becomes

$$ds^2 = g_{\tau\tau}(\tau, r)d\tau^2 + 2g_{\tau r}(\tau, r)d\tau dr + g_{rr}(\tau, r)dr^2 + r^2ds^2(S^2). \quad (6.14)$$

The only subtlety we could encounter in doing this change of coordinates is if it is not possible to exchange ρ with r , for example r could have been independent of ρ . We can rule out these cases by imposing that asymptotically the spacetime looks like Minkowski space.

We now want to get rid of a new coordinate which removes the cross term $d\tau dr$. If we pick $\tilde{t}(\tau, r)$ then we have

$$d\tilde{t} = \frac{\partial \tilde{t}}{\partial \tau}d\tau + \frac{\partial \tilde{t}}{\partial r}dr \quad (6.15)$$

and therefore we can pick a choice such that we can remove the cross term. The resultant metric is then

$$ds^2 = -e^{2\alpha(\tilde{t}, r)}d\tilde{t}^2 + e^{2\beta(\tilde{t}, r)}dr^2 + r^2ds^2(S^2). \quad (6.16)$$

We have included a minus sign since we are looking for a Lorentzian metric. This is the simplest form of the metric that we can achieve just through coordinate transformations and

we now need to plug this into Einstein's equations. We can compute the Christoffel symbols for the metric, the non-trivial ones are

$$\begin{aligned}\tilde{\Gamma}_{\tilde{t}\tilde{t}}^t &= \partial_{\tilde{t}}\alpha, & \tilde{\Gamma}_{\tilde{t}r}^t &= \partial_r\alpha, & \tilde{\Gamma}_{rr}^t &= e^{2\beta-2\alpha}\partial_{\tilde{t}}\beta, \\ \tilde{\Gamma}_{\tilde{t}\tilde{t}}^r &= e^{2\alpha-2\beta}\partial_r\alpha, & \tilde{\Gamma}_{\tilde{t}r}^r &= \partial_{\tilde{t}}\beta, & \tilde{\Gamma}_{rr}^r &= \partial_r\beta, \\ \tilde{\Gamma}_{r\theta}^\theta &= \frac{1}{r}, & \tilde{\Gamma}_{\theta\theta}^r &= -re^{-2\beta}, & \tilde{\Gamma}_{r\phi}^\phi &= \frac{1}{r}, \\ \tilde{\Gamma}_{\phi\phi}^r &= -re^{-2\beta}, & \tilde{\Gamma}_{\phi\phi}^\theta &= -\sin\theta\cos\theta, & \tilde{\Gamma}_{\theta\phi}^\phi &= \frac{\cos\theta}{\sin\theta}. \end{aligned}\quad (6.17)$$

It follows that the non-vanishing components of the Riemann tensor are

$$\begin{aligned}R_{r\tilde{t}r}^{\tilde{t}} &= e^{2\beta-2\alpha}\left(\partial_{\tilde{t}}^2\beta + (\partial_{\tilde{t}}\beta)^2 - \partial_{\tilde{t}}\alpha\partial_{\tilde{t}}\beta\right) + \left(\partial_r\alpha\partial_r\beta - \partial_r^2\alpha - (\partial_r\alpha)^2\right), \\ R_{\theta\tilde{t}\theta}^{\tilde{t}} &= -re^{-2\beta}\partial_r\alpha, \\ R_{\phi\tilde{t}\phi}^{\tilde{t}} &= -re^{-2\beta}\sin^2\theta\partial_r\alpha, \\ R_{\theta r\theta}^{\tilde{t}} &= -re^{-2\alpha}\partial_{\tilde{t}}\beta, \\ R_{\phi r\phi}^{\tilde{t}} &= -re^{-2\alpha}\sin^2\theta\partial_{\tilde{t}}\beta, \\ R_{\theta r\theta}^r &= re^{-2\beta}\partial_r\beta, \\ R_{\phi r\phi}^r &= re^{-2\beta}\sin^2\theta\partial_r\beta, \\ R_{\phi\theta\phi}^\theta &= (1-e^{-2\beta})\sin^2\theta. \end{aligned}\quad (6.18)$$

From the Riemann tensor we can construct the Ricci tensor finding the non-trivial components

$$\begin{aligned}R_{\tilde{t}\tilde{t}} &= \left(\partial_{\tilde{t}}^2\beta + (\partial_{\tilde{t}}\beta)^2 - \partial_{\tilde{t}}\alpha\partial_{\tilde{t}}\beta\right) + e^{2\alpha-2\beta}\left(\partial_r^2\alpha + (\partial_r\alpha)^2 - \partial_r\alpha\partial_r\beta + \frac{2}{r}\partial_r\alpha\right), \\ R_{rr} &= -\left(\partial_r^2\alpha + (\partial_r\alpha)^2 - \partial_r\alpha\partial_r\beta - \frac{2}{r}\partial_r\beta\right) + e^{2\beta-2\alpha}\left(\partial_{\tilde{t}}^2\beta + (\partial_{\tilde{t}}\beta)^2 - \partial_{\tilde{t}}\alpha\partial_{\tilde{t}}\beta\right) \\ R_{\tilde{t}r} &= \frac{2}{r}\partial_{\tilde{t}}\beta, \\ R_{\theta\theta} &= e^{-2\beta}\left(r(\partial_r\beta - \partial_r\alpha) - 1\right) + 1, \\ R_{\phi\phi} &= R_{\theta\theta}\sin^2\theta. \end{aligned}\quad (6.19)$$

Our job is to now solve Einstein's equations in the vacuum, $R_{\mu\nu} = 0$. There is an obvious component to consider first $R_{\tilde{t}r}$ which implies

$$\partial_{\tilde{t}}\beta = 0. \quad (6.20)$$

If we now take the \tilde{t} derivative of $R_{\theta\theta}$ and use the above condition we find

$$\partial_{\tilde{t}}\partial_r\alpha = 0, \quad (6.21)$$

and therefore we have

$$\beta = \beta(r), \quad \alpha = f(r) + g(\tilde{t}). \quad (6.22)$$

The first term in the metric is then

$$-e^{2f(r)+2g(\tilde{t})}d\tilde{t}^2, \quad (6.23)$$

and by a redefinition of \tilde{t} we can se

$$e^{g(\tilde{t})}d\tilde{t} = dt, \quad (6.24)$$

and we end up with the metric

$$ds^2 = -e^{2f(r)}dt^2 + e^{2\beta(r)}dr^2 + r^2ds^2(S^2), \quad (6.25)$$

and it remains to solve the remaining Einstein equations. Note that the metric is now independent of t , this naturally comes out of the Einstein equations, we did not impose this! This implies that *any spherically symmetric vacuum metric possesses a timelike Killing vector*. A metric with this property is called *stationary*, in fact the Schwarzschild metric is also *static* we will come back to this shortly.

We can now remove all \tilde{t} derivatives and exchange $\alpha \rightarrow f$ in the Ricci tensor components and where we se \tilde{t} replace with just t . We are free to add components and so we take the combination

$$0 = e^{2\beta-2f(r)}R_{tt} + R_{rr} = \frac{2}{r}(\partial_r f(r) + \partial_r \beta). \quad (6.26)$$

We then have

$$f(r) = -\beta(r) + \text{const}, \quad (6.27)$$

but we may rescale the time coordinate to set the constant to 0. Plugging this into $R_{\theta\theta}$ we find

$$e^{2f(r)}(2r\partial_r f(r) + 1) = 1 \Leftrightarrow \partial_r(re^{2f(r)}) = 1, \quad (6.28)$$

which has solution

$$e^{2f(r)} = 1 - \frac{R_S}{r}, \quad (6.29)$$

with R_S an undetermined constant which we will set to be $R_S = 2G_N M$. There is no remaining freedom except to set R_S to a certain value so the remaining components must vanish, and it turns out that they do, so we have solved Einstein's equations and derived the Schwarzschild solution.

Stationary vs Static There are two different meanings to time independence that we can use.

A spacetime is *stationary* if it admits an everywhere timelike Killing vector field K . We typically normalise it so that asymptotically $K^2 \rightarrow -1$.

A spacetime is *static* if, in addition to being stationary, it is invariant under $t \rightarrow -t$, where t is the coordinate along the integral curves of K . This rules out $dtdx$ cross terms in the metric with x any other coordinate except t .

Birkhoff's theorem tells us that spherical symmetry implies that the spacetime is necessarily static.

6.1.2 Geodesics

We now want to consider the geodesics of the Schwarzschild metric. We have computed the Christoffel symbols above and could just substitute this into the geodesic equation (4.79) however if one did not already have the Christoffel symbols this is not necessarily the quickest method. Instead one should use the Euler–Lagrange equations for the Lagrangian

$$\mathcal{L} = \sqrt{-g_{\mu\nu}\dot{x}^\mu\dot{x}^\nu}, \quad (6.30)$$

and use an affine parameter. With the choice of an affine parameter we can then compute the Euler–Lagrange equations of \mathcal{L}^2 instead and obtain the same equations of motion. We take

$$\begin{aligned} \mathcal{L} &= g_{\mu\nu} \frac{\partial x^\mu}{\partial \lambda} \frac{\partial x^\nu}{\partial \lambda} \\ &= - \left(1 - \frac{2G_NM}{r}\right) \dot{t}^2 + \left(1 - \frac{2G_NM}{r}\right)^{-1} \dot{r}^2 + r^2 \dot{\theta}^2 + r^2 \sin^2 \theta \dot{\phi}^2, \end{aligned} \quad (6.31)$$

with $\dot{\bullet} \equiv \frac{d\bullet}{d\lambda}$. Since we are using an affine parameter this is equal to a constant ϵ which we may take to be -1 for time-like geodesics, 0 for null and 1 for space-like geodesics.

Before we start with a brute force computation we should consider the conserved quantities. Ignorable coordinates, ones which do not appear explicitly, give rise to conserved quantities since from the Euler–Lagrange equations we find

$$\frac{d\mathcal{L}}{d\lambda} = 0 \quad \Rightarrow \quad \frac{d}{d\lambda} \frac{d\mathcal{L}}{d\lambda} = 0. \quad (6.32)$$

The action has two such ignorable coordinates t and ϕ : giving

$$\begin{aligned} 2l &= \frac{d\mathcal{L}}{d\dot{\phi}} = 2r^2 \sin^2 \theta \dot{\phi}, \\ -2E &= \frac{d\mathcal{L}}{dt} = -2 \left(1 - \frac{2G_N M}{r}\right) \dot{t}. \end{aligned} \quad (6.33)$$

Of course these should be identified with the angular momentum and energy respectively. Next consider the equation for θ , we find

$$\frac{d}{d\lambda}(r^2 \dot{\theta}) = r^2 \sin \theta \cos \theta \dot{\phi}^2. \quad (6.34)$$

Recall that in computing the motion in Newtonian gravity we noted that if we started the particle at $\theta = \frac{\pi}{2}$ with $\dot{\theta} = 0$ then it remained in the plane, the same is true here and so we can without loss of generality set $\theta = \frac{\pi}{2}$.

We can now plug this into (6.31) and equate with our constant parameter ϵ giving

$$\epsilon = - \left(1 - \frac{2G_N M}{r}\right)^{-1} E^2 + \left(1 - \frac{2G_N M}{r}\right)^{-1} \dot{r}^2 + r^{-2} l^2. \quad (6.35)$$

Rearranging we have

$$\frac{1}{2} \dot{r}^2 + V_{\text{eff}}(r) = \frac{E^2}{2}, \quad (6.36)$$

with

$$V_{\text{eff}}(r) = -\frac{\epsilon}{2} + \frac{\epsilon G_N M}{r} + \frac{l^2}{2r^2} - \frac{l^2 G_N M}{r^3}, \quad (6.37)$$

we should contrast this with the equivalent Newtonian expression in (2.80) for a massive particle which was

$$V_N(r) = -\frac{G_N M}{r} + \frac{l^2}{2r^2}. \quad (6.38)$$

We see that General relativity leads to additional corrections to the potential. The first term is simply a constant shift and so does not play much of a role since we can absorb it into a redefinition of the energy, the r^{-3} term is completely new however and changes the Newtonian potential at small distances. Note that the effective potential vanishes at $r = 2G_N M$ which is the Schwarzschild radius.

Let us reinstate the speed of light in the potential, we have

$$V_{\text{eff}}(r) = -\frac{\epsilon c^2}{2} + \frac{\epsilon G_N M}{r} + \frac{l^2}{2r^2} - \frac{l^2 G_N M}{r^3 c^2}, \quad (6.39)$$

then the equation for \dot{r} is

$$\frac{1}{2} \dot{r}^2 + V_{\text{eff}}(r) = \frac{1}{2} \frac{E^2}{c^2}. \quad (6.40)$$

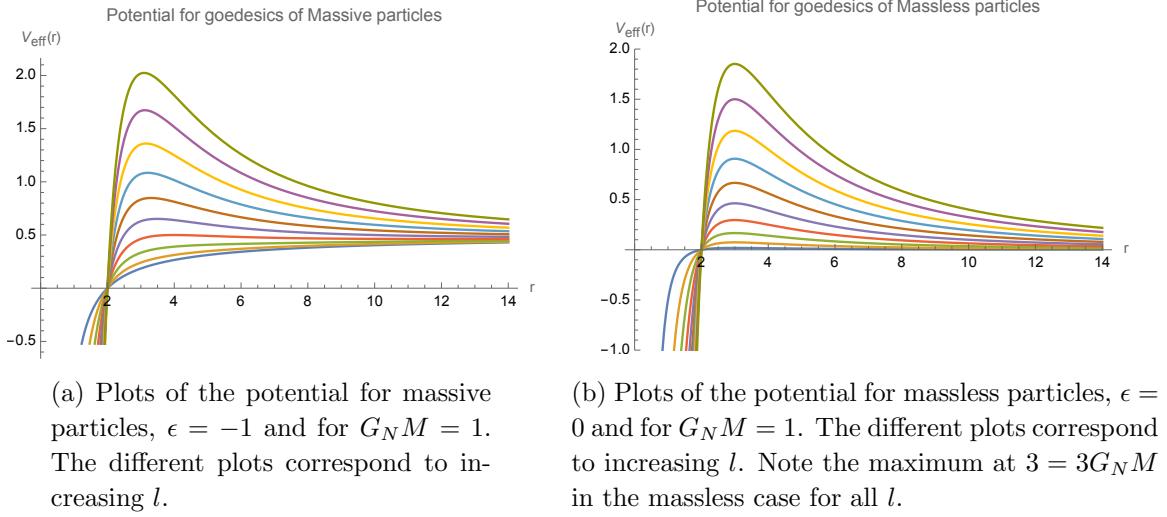


Figure 15: Plots of the potential for massive and massless particles. Note that the plots tends to $-\frac{\epsilon}{2}$ as $r \rightarrow \infty$. Moreover the potentials both vanish at $2 = 2G_N M$ which is the Schwarzschild radius.

We now want to analyse the different forms of trajectories that are possible. In figure 25 we have plotted the potential for various values of l , with fixed mass M .

Circular orbits will be at points where the potential has a turning point. Then we are stuck in a circular orbit, which is stable if it corresponds to a minimum of the potential and unstable if it corresponds to a maximum. Differentiating the potential we have

$$V'_{\text{eff}}(r) = \frac{1}{r^4} \left(3G_N l^2 M - l^2 r - G_N M \epsilon r^2 \right) \quad (6.41)$$

which potentially has two zeroes at

$$r_c = -\frac{l^2 \pm \sqrt{l^4 + 12G_N l^2 M \epsilon}}{2G_N M \epsilon}, \quad (6.42)$$

for $\epsilon \neq 0$ and

$$r_c = 3G_N M, \quad (6.43)$$

for $\epsilon = 0$.

For the massless photon the orbit is at a maximum and is therefore unstable. However a photon can orbit in a circular orbit forever around the black hole, but any perturbation will send it flying off to either $r = 0$ or $r = \infty$. It is known as the *photon sphere*. The focussing effects mean that much of the light emitted from an accretion disc around a non-rotating black hole emerges from the photon sphere. In practice, it seems likely that the photographs by the Event Horizon Telescope does not have the required resolution to see this.

For massive particles there are different regimes depending on the angular momentum. For large l there will be two circular orbits, one stable and one unstable. In the $l \rightarrow \infty$ regime they are at

$$r_c = \left(\frac{l^2}{G_N M}, 3G_N M \right). \quad (6.44)$$

The stable circular orbit gets further away while the unstable orbit approaches $3G_N M$. As we decrease l the two orbits come together and coincide when the discriminant of the quadratic in (6.41) vanishes. This is at

$$l = \sqrt{12}G_N M, \quad (6.45)$$

which gives

$$r_c = 6G_N M. \quad (6.46)$$

For smaller l there are no circular orbits and so $6G_N M$ is the smallest possible radius of a stable circular orbit of the Schwarzschild metric.

We have found that the Schwarzschild solution possesses stable circular orbits for $r > 6G_N M$ and unstable circular orbits for $3G_N M < r < 6G_N M$. We should comment that these are the motions of geodesics. For an accelerating observer such as a rocket ship, there is nothing stopping them from dipping below $r = 3G_N M$ and then reemerging, so long as they stay away from $r = 2G_N M$.

Most experimental test of general relativity involve the motion of test particles in the solar system. More recently, with the advancements in technology, using gravitational waves to test general relativity has also become possible. We will concentrate on three particular tests: the precession of perihelia, the bending of light and gravitational red-shift.

Perihelion precession We saw when we consider the orbits in Newtonian gravity that the non-circular orbits were closed ellipses. Observation of the orbit of Mercury showed that the closed elliptic orbits of Newtonian gravity were not realised, instead the orbit precessed. A non-trivial check of General Relativity is then to show that the orbits of the planets precess. We can approximate the metric of the sun to be Schwarzschild and take the planet to follow a geodesic of a massive particle.

The strategy is to describe the evolution of the radial coordinate r as a function of ϕ . If the orbit is a perfect ellipse $r(\phi)$ should be periodic with period 2π , for which the perihelion occurs at the same point every orbit. Instead for a non-closed ellipse the perihelion is shifted after every orbit. We will see that General Relativity gives a slight modification of the Newtonian result such that the orbit precesses. First consider the radial equation of motion

for a massive particle, (6.36), setting $\epsilon = -1$. To get an equation for $\frac{dr}{d\phi}$ we can use the chain rule and multiply the equation by

$$\left(\frac{d\phi}{d\lambda}\right)^{-2} = \frac{r^4}{l^2}, \quad (6.47)$$

yielding

$$\left(\frac{dr}{d\phi}\right)^2 + \frac{r^4}{l^2} - \frac{2G_N M}{l^2} r^3 + r^2 - 2G_N M r = \frac{E^2 r^4}{l^2} \quad (6.48)$$

We first define a new variable

$$x = \frac{l^2}{G_N M r}, \quad (6.49)$$

which for $x = 1$ gives rise to the Newtonian circular orbit. The equation of motion becomes

$$\left(\frac{dx}{d\phi}\right)^2 + \frac{l^2}{G_N^2 M^2} - 2x + x^2 - \frac{2G_N^2 M^2 x^3}{l^2} = \frac{E^2 l^2}{G_N^2 M^2}. \quad (6.50)$$

Next differentiate with respect to ϕ to obtain

$$\frac{d^2 x}{d\phi^2} - 1 + x = \frac{3G_N^2 M^2 x^2}{l^2}. \quad (6.51)$$

In the Newtonian calculation the last term would be absent and we could solve for x exactly. Here we will treat this as a perturbation around the Newtonian result.

We expand x into a Newtonian solution plus a small deviation

$$x = x_0 + x_1, \quad (6.52)$$

where the zeroth order part satisfies

$$\frac{d^2 x_0}{d\phi^2} - 1 + x_0 = 0, \quad (6.53)$$

leading to the equation for the first order part

$$\frac{d^2 x_1}{d\phi^2} + x_1 = \frac{3G_N^2 M^2}{l^2} x_0^2. \quad (6.54)$$

A solution to the zeroth order equation is (see (2.86))

$$x_0 = 1 + e \cos \phi, \quad (6.55)$$

which recall describes a perfect ellipse with eccentricity e , $e = 1 - \frac{b^2}{a^2}$ with a the semi-major axis, the distance from the centre to the farthest point on the ellipse and the semi-minor axis b the distance from the centre to the closest point. Plugging in the Newtonian solution into the first order equation of motion we find

$$\frac{d^2 x_1}{d\phi^2} + x_1 = \frac{3G_N^2 M^2}{l^2} (1 + e \cos \phi)^2. \quad (6.56)$$

A solution is given by

$$x_1 = \frac{3G_N^2 M^2}{l^2} \left[\left(1 + \frac{e^2}{2} \right) + e\phi \sin \phi - \frac{1}{6} e^2 \cos 2\phi \right]. \quad (6.57)$$

The first term is just a constant displacement while the third oscillates around 0. The important effect is contained within the second term which accumulates over successive orbits. Combining this term only with the zeroth-order solution we have

$$x = 1 + e \cos \phi + \frac{3G_N^2 M^2 e}{l^2} \phi \sin \phi. \quad (6.58)$$

We should emphasise that this is not a full solution, it is an approximation but it encapsulates the part we are interested in. We may write

$$x = 1 + e \cos ((1 - \alpha)\phi), \quad (6.59)$$

where

$$\alpha = \frac{3G_N^2 M^2}{l^2}. \quad (6.60)$$

where one should view this as a series expansion around $\alpha = 0$. It follows that during each orbit the perihelion advances by an angle

$$\Delta\phi = 2\pi\alpha = \frac{6\pi G_N^2 M^2}{l^2}. \quad (6.61)$$

We may replace the angular momentum in favour of the eccentricity by looking at the Newtonian solution. An ordinary ellipse satisfies

$$r = \frac{(1 - e^2)a}{1 + e \cos \phi}, \quad (6.62)$$

with a the semi-major axis. This leads us to identify

$$l^2 \sim G_N M (1 - e^2) a, \quad (6.63)$$

for the Newtonian orbit. Plugging this in and restoring the speed of light we find

$$\Delta\phi = \frac{6\pi G_N M}{c^2 (1 - e^2) a}. \quad (6.64)$$

historically the precession of mercury was the first test of GR. The apparent discrepancy between observation and Newtonian gravity was known long before the advent of GR, and

a number of solutions had been proposed including additional planets. For the motion of Mercury around the sun we have

$$\begin{aligned}\frac{G_N M_{\odot}}{c^2} &= 1.48 \times 10^3 m, \\ a &= 5.79 \times 10^{10} m, \\ e &= 0, 2056.\end{aligned}\tag{6.65}$$

This gives

$$\Delta\phi_{\text{Mercury}} = 5.01 \times 10^{-7} \text{ radians/orbit} = 0.103''/\text{orbit}\tag{6.66}$$

with " denoting arcseconds. Mercury orbits once every 88 days and therefore

$$\Delta\phi_{\text{Mercury}} = 43.0''/\text{century}.\tag{6.67}$$

The major axis of Mercury's orbit precesses at a rate of 43.0 arcseconds every 100 years. The observed value is 5601 arcseconds/100 years. Much of that is due to the precession of equinoxes in our geocentric coordinate system: 5025 arcseconds/100 years. The gravitational perturbations of the other planets contributes an additional 532 arcseconds/100 years leaving a 43 arcseconds/100 years to be explained by GR which is does quite well.

Bending of light We can now extend these results for null geodesics. We have seen that there is an unstable circular orbit for light. What about other orbits? The fate of other light rays depends on the relative value of their energy E to their angular momentum l . The maximum value of the potential is

$$V_{\text{null}}(r_*) = \frac{l^2}{54 G_N^2 M^2},\tag{6.68}$$

and therefore the physics depends on how this compares with the right-hand side of (6.36). There are two possibilities we need to consider

- $E < \frac{l}{\sqrt{27} G_N M}$. The energy of the light is lower than the angular momentum barrier. This means that light emitted from $r < r_*$ cannot escape to infinity; it will orbit the star before falling back towards the origin. For light coming from infinity it will not fall into the star but will instead bounce off the angular momentum barrier and return to infinity: the light will be scattered.
- $E > \frac{l}{\sqrt{27} G_N M}$. The energy of light is greater than the angular momentum barrier. Light can be emitted from $r < r_*$ and escape to infinity (this is only true for $R_s < r$). Meanwhile light coming from infinity is captured by the star/black hole.

Let us once again use the inverse parameter $u = \frac{1}{r}$. The equation of motion becomes

$$\left(\frac{du}{d\phi}\right)^2 + u^2(1 - 2G_NMu) = \frac{E^2}{l^2}. \quad (6.69)$$

Differentiating again we find

$$\frac{d^2u}{d\phi^2} + u = 3G_NMu^2. \quad (6.70)$$

We may once again work perturbatively. At zeroth order we can ignore the $G_N M$ term on the right-hand-side. Then to leading order we have

$$\frac{d^2u}{d\phi^2} + u = 0, \quad \Rightarrow \quad u = \frac{1}{b} \sin \phi, \quad (6.71)$$

for b a constant. Reinstating r we have $r \sin \phi = b$: which is the equation of a horizontal straight line, a distance b above the origin, see 16. The distance b is known as the *impact parameter*.

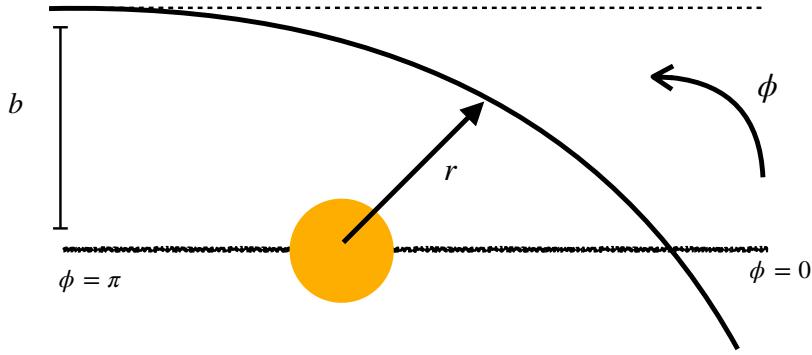


Figure 16: Light bending in the Schwarzschild metric. The dashed line at the top is the constant line $r \sin \phi = b$. The curved line is the geodesic.

With the zeroth order solution we can now solve (6.70) in an expansion around $\frac{G_N M}{b} = \beta$. We have

$$u = u_0 + \beta u_1 + \dots \quad (6.72)$$

At first order we need to solve

$$\frac{d^2u_1}{d\phi^2} + u_1 = \frac{3 \sin^2 \phi}{b} = \frac{3(1 - \cos 2\phi)}{2b}. \quad (6.73)$$

The general solution is

$$u_1 = A \cos \phi + B \sin \phi + \frac{1}{2b} (3 + \cos 2\phi), \quad (6.74)$$

where the first two parts are the solutions of the homogeneous part and A, B two integration constants. We should choose them so that the initial trajectory at $\phi = \pi$ agrees with the straight line u_0 . For this to hold we must take $A = \frac{2}{b}$ and $B = 0$ so that $u_1 \rightarrow 0$ as $\phi \rightarrow \pi$. To leading order in β the solution is

$$u = \frac{1}{b} \sin \phi + \frac{G_N M}{2b^2} (3 + 4 \cos \phi + \cos 2\phi). \quad (6.75)$$

What angle does the particle escape to $r = \infty \Leftrightarrow u = 0$. Before the correction this was at $\phi = 0$, within our perturbative approach we can approximate $\sin \phi \sim \phi$ and $\cos \phi \sim 1$ to find that the particle escapes at

$$\phi \sim -\frac{4G_N M}{b}. \quad (6.76)$$

This bending of light is known as *gravitational lensing*.

For the sun, $\frac{G_N M_\odot}{c^2} \sim 1.48$ km. If the light rays just graze the surface of the sun, then the impact parameter is the radius of the sun $R_\odot \sim 7 \times 10^5$ km. This gives a scattering angle of $\phi \sim 8.6 \times 10^{-5}$ radians or $\phi \sim 1.8''$. The Newtonian prediction gives only half of this contribution.

There is a difficulty in testing this prediction since things behind the sun are rarely visible. By a sheer coincidence, the size of the moon in the sky is about the same size of the sun. This means that during a solar eclipse the light from the sun is blocked allowing for the measurement of stars whose light passes nearby the Sun. This can then be compared with the usual positions of these stars.

The first measurement was carried out in 1919 by two expeditions lead by Arthur Eddington (we will see this name again shortly). Since then our evidence of the bending of light is more impressive. Clusters of galaxies have been seen to distort the light from a background source often revealing a distinct ring-like pattern of multiple copies of the light source. See figure 17.

Gravitational red shift Let us consider an observer with four velocity U^μ who is stationary in Schwarzschild coordinates, i.e. $U^i = 0$.²³ The four-velocity is normalised so that $U_\mu U^\mu = 1$,

²³We could allow for the observer to be moving, however the difference is just to superimpose the usual Doppler shift on top of the gravitational effect and therefore we consider the simpler example.



Figure 17: A diagram of light lensing picked up by the Hubble telescope. Notice that there are four copies of the distant quasar in the picture obtained by Hubble. Image credited to NASA, ESA and STScI.

which for our stationary observer in a Schwarzschild background implies

$$U^0 = \left(1 - \frac{2G_N M}{r}\right)^{-1/2}. \quad (6.77)$$

Such an observer measures the frequency of a photon following a null geodesic $x^\mu(\lambda)$ to be

$$\omega = -g_{\mu\nu} U^\mu \frac{dx^\nu}{d\lambda}. \quad (6.78)$$

We have

$$\begin{aligned} \omega &= \left(1 - \frac{2G_N M}{r}\right)^{1/2} \frac{dt}{d\lambda} \\ &= \left(1 - \frac{2G_N M}{r}\right)^{-1/2} E, \end{aligned} \quad (6.79)$$

where E was defined to be the conserved quantity associated to time translations when we worked out the geodesics. Since E is conserved it follows that ω will have different values when measured at different radial distances. For a photon emitted at r_1 and an observer at r_2 , the observed frequencies will be related by

$$\frac{\omega_2}{\omega_1} = \sqrt{\frac{1 - 2G_N M/r_1}{1 - 2G_N M/r_2}}. \quad (6.80)$$

This is the exact result for the frequency shift, in the limit $r \gg 2G_N M$ we have

$$\begin{aligned} \frac{\omega_2}{\omega_1} &= 1 - \frac{G_N M}{r_1} + \frac{G_N M}{r_2} \\ &= 1 + \Phi(r_1) - \Phi(r_2), \end{aligned} \quad (6.81)$$

with $\Phi = -G_N M/r$ the Newtonian potential.

We see that the frequency goes down as Φ increases which happens as we climb out of a gravitational field, leading to a red-shift. On the other hand photons which fall towards the gravitating body are blue shifted. Gravitational red-shift was first detected in 1960 by Pound and Rebka using gamma rays travelling a distance of 72-feet (about 22m) which was the height of the physics building at Harvard. Increasingly precise tests have found excellent agreement with GR.

There is a cosmological counterpart to this, where light is red-shifted in an expanding universe.

Time delay Since the temporal component of the metric is

$$g_{00}(x) = 1 + 2\Phi(x), \quad (6.82)$$

we see that there is a connection between time and gravity. Let us once again use the Schwarzschild solution. An observer sitting at a fixed distance r from the origin will measure a time interval

$$d\tau^2 = -g_{00}dt^2 = \left(1 - \frac{2G_N M}{r}\right)dt^2. \quad (6.83)$$

For an asymptotic observer at $r \rightarrow \infty$ who measures a time t , an observer at r will measure the time T

$$T(r) = t\sqrt{1 - \frac{2G_N M}{r}}. \quad (6.84)$$

It follows that time goes slower in the presence of a massive gravitating object. Notice that at $r = r_S$ that time seems to stop for the observer at r_s . We will come back to this later.

We can make this more quantitative by considering two observers: Alice and Bob. Bob has gone up in a hot air balloon while Alice is on the surface of the earth at r_A . Bob is at a distance $r_B = r_A + \Delta r$. The time measured by Bob is

$$\begin{aligned} T_B &= t\sqrt{1 - \frac{2G_N M}{(r_A + \Delta r)}} \sim t\sqrt{1 - \frac{2G_N M}{r_A} + \frac{2G_N M \Delta r}{r_A^2}} \\ &\sim t\sqrt{1 - \frac{2G_N M}{r_A} \left(1 + \frac{G_N M \Delta r}{r_A^2}\right)} = T_A \left(1 + \frac{G_N M \Delta r}{r_A^2}\right). \end{aligned} \quad (6.85)$$

A double expansion has been utilised where we assume $\Delta r \ll r_A$ and $\frac{2G_N M}{r_A} \ll 1$. If the hot air balloon flies a distance $\Delta r = 1000m$ above Alice then taking the radius of the Earth to be $r_A \approx 6000\text{km}$ the difference in times is about 10^{-12} and therefore over the whole day Bob ages by an extra 10^{-18} seconds or so. Clearly this is a small amount, in the vicinity of a black

hole this can be more pronounced. Recall that the smallest stable orbit was at $r = 3G_NM$ and such a person experiences time at a rate of $T = 3^{-1/2}t \approx 0.6t$ compared to an asymptotic observer at $r \rightarrow \infty$. For more dramatic results one would need to fly closer to the horizon and then return to asymptotic infinity.

This also gives a different perspective on the gravitational redshift. Bob doesn't like Alice and wants to ruin her day so he hovers above Alice and chuck's peanuts at her. He throws peanuts at time intervals ΔT_B . Alice, wise to Bob's antics, opens up an umbrella. The peanuts hit the umbrella at time intervals ΔT_A where as above

$$\Delta T_A = \Delta T_B \sqrt{\frac{1 + 2\Phi(r_A)}{1 + 2\Phi(r_B)}} \approx (1 + \Phi(r_A) - \Phi(r_B)) \Delta T_B. \quad (6.86)$$

We have that $r_A < r_B$ and therefore $\Phi(r_A) < \Phi(r_B) < 0$ and hence $\Delta T_A < \Delta T_B$. Alice receives the peanuts at a higher frequency than Bob threw them.

Having seen the peanuts hitting the umbrella Bob decides to instead shine a light down at Alice with a frequency $\omega_B \sim \Delta T_B^{-1}$. Alice will then receive the light at a frequency ω_A where

$$\omega_A \approx (1 + \Phi(r_A) - \Phi(r_B))^{-1} \omega_B. \quad (6.87)$$

This is a higher frequency $\omega_A > \omega_B$ and therefore a shorter wave-length. The light is therefore blue-shifted. In contrast if Alice retaliates and shines a light up to Bob then the frequency decreases and the light is redshifted.

6.2 Schwarzschild solution as a black hole

We have now studied some geodesics for the Schwarzschild solution and some phenomena. Each time we have carefully avoided the Schwarzschild radius $r_c = 2G_NM$ and also $r = 0$. At both of these points something funky happens with the metric, at least one of the components of the metric diverges or vanishes. The interpretation of the singularity is different for the two cases. The divergence at $r = 0$ is a *singularity*. General relativity breaks down here and we need a theory of quantum gravity. GR predicts its own death!

In contrast the divergence at $r = 2G_NM$ is a result of our choice of coordinates. This surface is referred to as the *event horizon* or simply the *horizon*. Many of the surprising properties of a black hole happen here.

There is a simple way to check whether a divergence is due to a singularity or a poor choice of coordinates. We can build scalar quantities, these are then independent of coordinates, if they diverge in one coordinate system they diverge in all and the spacetime is sick at

this point. One the other hand if it does not diverge we cannot say much, one would have to consider all possible scalar quantities to concretely say it is just a coordinate singularity. Since the Einstein equations in a vacuum set $R_{\mu\nu} = 0$ it follows that the simplest scalar quantities one can construct R and $R_{\mu\nu}R^{\mu\nu}$ both vanish. The next simplest is the *Kretschmann scalar* $R^{\mu\nu\rho\sigma}R_{\mu\nu\rho\sigma}$. For the Schwarzschild metric we find

$$R^{\mu\nu\rho\sigma}R_{\mu\nu\rho\sigma} = \frac{48G_N^2 M^2}{r^6}. \quad (6.88)$$

There is no pathology at $r = 2G_N M$ while there is at $r = 0$ where it diverges.

One way to understand the geometry of spacetime is to explore its causal structure as defined by light cones. We therefore consider radial null curves, i.e. those with constant θ, ϕ and $ds^2 = 0$, such that they satisfy

$$ds^2 = 0 = -\left(1 - \frac{2G_N M}{r}\right)dt^2 + \left(1 - \frac{2G_N M}{r}\right)^{-1}dr^2, \quad (6.89)$$

which gives

$$\frac{dt}{dr} = \pm \left(1 - \frac{2G_N M}{r}\right)^{-1}. \quad (6.90)$$

This measures the slope of the light cones on a spacetime diagram of the t - r plane. For large r the slope is ± 1 as it would be for flat spacetime. On the other hand as we approach $r = 2G_N M$ we get $\frac{dt}{dr} \rightarrow \pm\infty$ and the light cones close up, see figure 18. Thus a light ray which approaches $r = 2G_N M$ never seems to get there, at least in this coordinate system. This apparent inability to get to $r = 2G_N M$ is actually an illusion and an artefact of a bad choice of coordinates. An in-falling light ray or massive particle has no trouble reaching this radius. On the other hand an observer far away would never be able to tell. If we all hovered outside a black hole and one of your class mates jumped in the black hole sending back signals the whole way down we would simply see the signals reach us less frequently, see figure 19.

The fact that we never see them reach $r = 2G_N M$ is a meaningful statement but the fact that their trajectory in the t - r plane never reaches there is not: it is highly dependent on our coordinate system. We want to change coordinates to some that are better behaved at $r = r_S$. Note that we can solve (6.90) by introducing the *tortoise coordinate* r_*

$$r_* = r + 2G_N M \log\left(\frac{r - 2G_N M}{2G_N M}\right), \quad (6.91)$$

then

$$t = \pm r_* + \text{constant}, \quad (6.92)$$

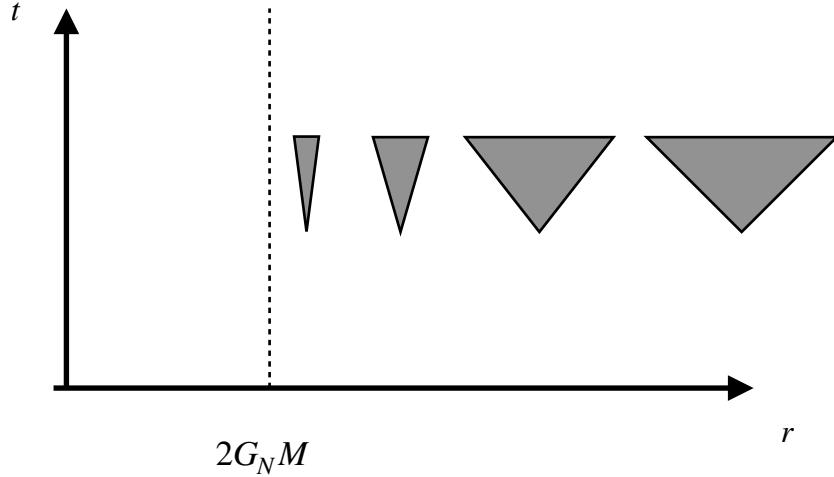


Figure 18: In Schwarzschild coordinates the light cones appear to close up as we approach the horizon. We will see that this is not quite correct.

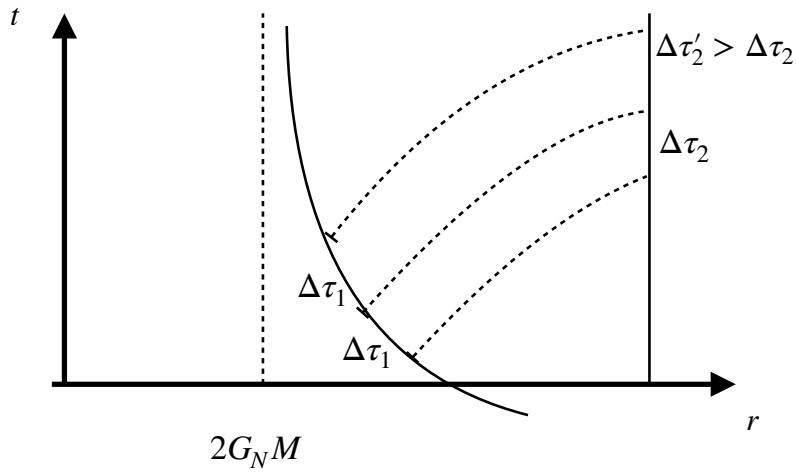


Figure 19: A beacon freely falling into a black hole emits signals at intervals of proper time $\Delta\tau_1$. An observer at fixed r receives these signals at a successively longer time intervals $\Delta\tau_2$.

and we see that this is well adapted to null radial geodesics. The plus sign corresponds to out-going geodesics and the negative to in-going geodesics²⁴. The metric with this new

²⁴The quick way to see this is to note that as $r \rightarrow \infty$ we have $r_* \rightarrow \infty$ and therefore we need the plus sign for out-going geodesics so that the radial direction increases with time.

coordinate becomes

$$ds^2 = \left(1 - \frac{2G_NM}{r}\right)(-dt^2 + dr_*^2) + r^2 ds^2(S^2). \quad (6.93)$$

Next we introduce a pair of null coordinates further adapted to the null geodesics:

$$v = t + r_*, \quad u = t - r_*. \quad (6.94)$$

We first consider the metric in (v, r) coordinates and then in (u, r) coordinates before biting the bullet and using (v, u) coordinates.

Ingoing Eddington–Finkelstein coordinates Eliminating t via $t = v - r_*(r)$ we find

$$ds^2 = -\left(1 - \frac{2G_NM}{r}\right)dv^2 + 2dvdr + r^2 ds^2(S^2). \quad (6.95)$$

This is the Schwarzschild solution in *ingoing Eddington–Finkelstein coordinates*. Even though the metric coefficient g_{vv} vanishes at $r = 2G_NM$ there is no real degeneracy. The determinant of the metric is

$$\det g = \det \begin{pmatrix} -\left(1 - \frac{2G_NM}{r}\right) & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & r^2 & 0 \\ 0 & 0 & 0 & r^2 \sin^2 \theta \end{pmatrix} = -r^4 \sin^2 \theta. \quad (6.96)$$

The cross terms stops the metric from being degenerate at the horizon. The metric is still degenerate at $r = 0$ and $\theta = 0, \pi$ however the latter are just the usual pole problems of the S^2 and nothing to worry about. This is the benefit of the Eddington–Finkelstein coordinates, the radial coordinate can be extended beyond the horizon.

To build further intuition we can look at the behaviour of light rays. We saw that the null radial geodesics were given by (6.92). The outgoing geodesics are

$$u = t - r_* = \text{const}. \quad (6.97)$$

Eliminating t in favour of v we have that the outgoing geodesics satisfy $v = 2r_* + \text{const}$. The solutions of this equation have a different behaviour depending on whether they are inside the horizon or outside. For $r > 2G_NM$ we can use the original definition of r_* in (6.91) to get

$$v = 2r + 4G_NM \log \left(\frac{r - 2G_NM}{2G_NM} \right) + \text{const}. \quad (6.98)$$

The Log term goes bad when $r < 2G_N M$, however we can simply modify the coordinate to take the norm of the argument of the log, so that

$$r_* = r + 2G_N M \log \left| \frac{r - 2G_N M}{2G_N M} \right|. \quad (6.99)$$

This means that r_* is multi-valued. Outside the horizon it takes values $r_* \in (-\infty, \infty)$ while inside the horizon it takes values $r_* \in (-\infty, 0)$. The singularity sits at $r_* = 0$. Outgoing geodesics inside the horizon obey

$$v = 2r + 4G_N M \log \left(\frac{2G_N M - r}{2G_N M} \right) + \text{const.} \quad (6.100)$$

Finally note that $r = 2G_N M$ is itself a null geodesic. This information can be captured in a *Finkelstein diagram*. It is designed so that ingoing null rays travel at 45° . This is simple to do if we label the coordinates of the diagram by t and r_* , however since r_* is not single valued we use r instead. We define a new temporal coordinate t_* by the requirement

$$v = t + r_* = t + * + r. \quad (6.101)$$

Thus ingoing null rays travel at 45° in the (t_*, r) -plane. See figure 20

The outgoing null geodesics that sit outside the horizon tend to infinity, whereas those inside the horizon don't actually go out, but rather go towards the singularity at $r = 0$. Each hits the singularity at some finite t_* . We can draw lightcones on the Finkelstein diagram. These are regions which are bounded by the in-going and out-going future pointing null geodesics. Any massive particle must follow a timelike path and this must then sit within these lightcones. We see that the light cones get tipped as we get closer to the horizon, and then once inside the horizon there is no way of getting back out. The causal structure of spacetime prevents this. The term black hole really refers to this area inside the horizon $r < 2G_N M$, any observer outside the horizon can never know what is happening inside the black hole.

We can also see what happens if we watch someone fall into a black hole. The person falls through the horizon without realising anything is wrong. However as they fall the light signals that come back to us take longer and longer to reach us. The actions of the in-falling person become increasingly slowed as they approach the horizon. In this way we continue to see the person forever, but we know nothing about their fate past the horizon. Since the light returns to us from a deeper and deeper gravitational well it appears increasingly red-shifted to us.

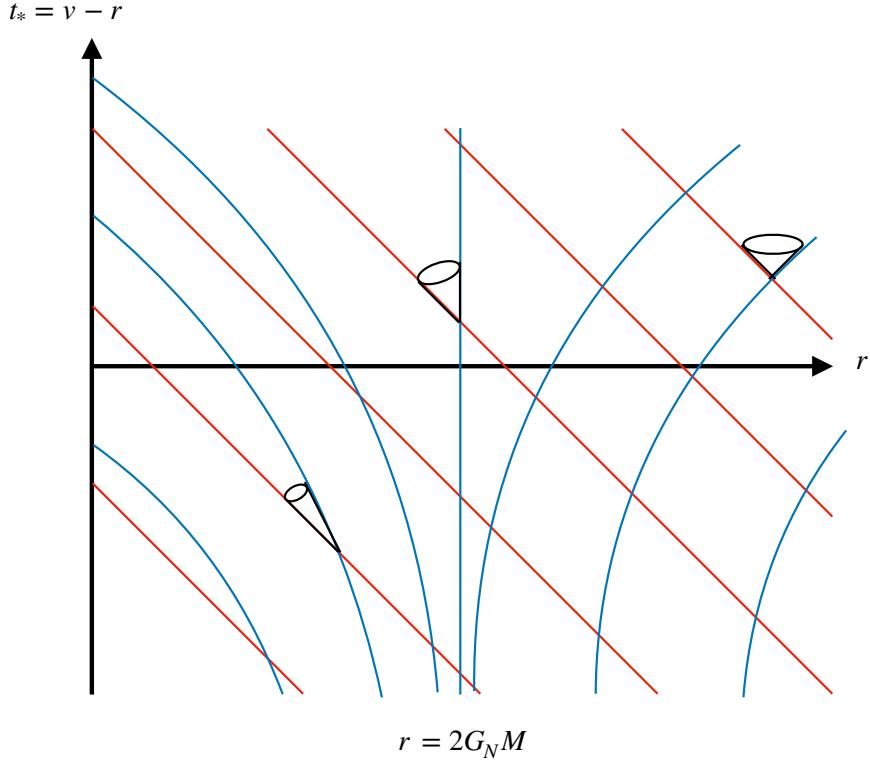


Figure 20: The Finkelstein diagram in in-going coordinates. The ingoing null geodesics are in red while the outgoing are in blue. Inside the horizon the outgoing geodesics never go past the horizon.

Out-going Eddington–Finkelstein coordinates We can also extend the exterior of the Schwarzschild black hole by replacing the time coordinate with the null coordinate

$$u = t - r_* . \quad (6.102)$$

Surfaces of constant u correspond to outgoing radial null geodesics. After the change of coordinates we have

$$ds^2 = -\left(1 - \frac{2G_N M}{r}\right)du^2 - 2dudr + r^2 ds^2(S^2) . \quad (6.103)$$

This is the Schwarzschild solution in *out-going Eddington–Finkelstein coordinates*. The only difference is in the sign of the cross term. This seemingly trivial modification changes the interpretation drastically.

As before the metric is smooth at the horizon and we can continue the metric down to the singularity at $r = 0$. However the region $r < 2G_N M$ now describes a different part of spacetime from the analogous region in ingoing Eddington–Finkelstein coordinates.

We again look at the ingoing and outgoing null radial geodesics. This time we pick coordinates so that the outgoing geodesics travel at 45° . This means that we take r and $t_* = u + r$ to be the axes.

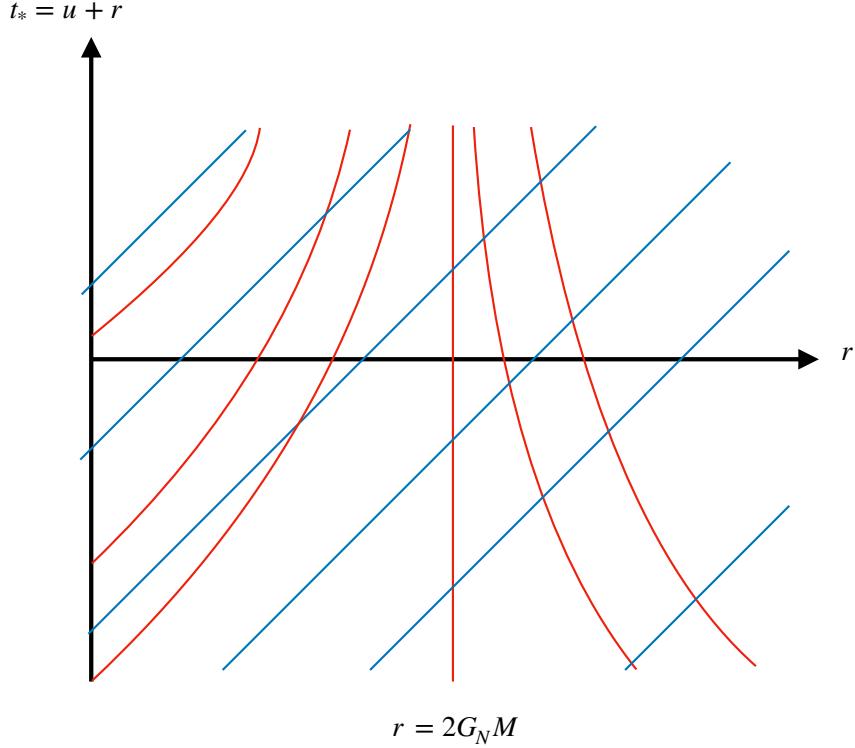


Figure 21: The Finkelstein diagram in out-going coordinates. The ingoing null geodesics are in red while the outgoing are in blue. Inside the horizon the ingoing geodesics never go past the horizon.

This time the ingoing null geodesics have the interesting property. Those which start outside are unable to reach the singularity, instead they pile up at the horizon. Those that start behind the horizon move towards the horizon, once again piling up there. What happens to massive particles that sit inside the horizon? Their trajectories must lie inside the future pointing light-cones. They cannot stay inside the horizon and the causal structure of spacetime requires them to be ejected outside of the horizon. This is a *white hole*, an object which expels matter. This is the time reversal of a black hole indeed the difference is purely a minus sign. Moreover if we flip white-hole upside down we get the black hole.

White holes are perfectly acceptable solutions of general relativity. Indeed they are implied by the time reversal invariance of Einstein's equations. However white holes are not

physically relevant since in contrast to a black hole they cannot be formed by collapsing matter.

6.2.1 Kruskal spacetime

We have seen that we can extend the $r \in (2G_N M, \infty)$ coordinate in two ways so that we gain the region $r \in (0, 2G_N M]$ which corresponds to two different parts of spacetime. We can write the Schwarzschild metric using both null (u, v) -coordinates, the metric is

$$ds^2 = -\left(1 - \frac{2G_N M}{r}\right)du dv + r^2 ds^2(S^2), \quad (6.104)$$

where r is a function of $u - v$. In these coordinates the metric is again degenerate at $r = 2G_N M$ so we need to perform another change of coordinates. We can introduce the *Kruskal-Szekeres coordinates*,

$$U = -\exp\left(-\frac{u}{4G_N M}\right), \quad V = \exp\left(\frac{v}{4G_N M}\right). \quad (6.105)$$

Both are null coordinates. The Schwarzschild black hole is parametrised by $U < 0$ and $V > 0$. Outside the horizon they satisfy

$$UV = -\exp\left(\frac{r_*}{2G_N M}\right) = \frac{2G_N M - r}{2G_N M} \exp\left(\frac{r}{2G_N M}\right), \quad (6.106)$$

and similarly

$$\frac{U}{V} = -\exp\left(-\frac{t}{2G_N M}\right). \quad (6.107)$$

The metric is then

$$ds^2 = -\frac{32(G_N M)^3}{r} e^{-\frac{r}{2G_N M}} dU dV + r^2 ds^2(S^2), \quad (6.108)$$

with $r(U, V)$ defined by inverting (6.106). The original Schwarzschild metric covers just $U < 0$ and $V > 0$ however there is no obstruction to extending $U, V \in \mathbb{R}$. Nothing bad happens at $r = 2G_N M$, the metric is smooth and non-degenerate. The Kruskal spacetime is the maximal extension of the Schwarzschild solution.

The Kruskal Diagram To find the location of the horizon in the new coordinates we can use equation (6.106). We see that this is at

$$r = 2G_N M \Rightarrow U = 0 \text{ or } V = 0. \quad (6.109)$$

The horizon is not just one null surface but 2 which intersect at $U = V = 0$. On the other hand the singularity is at

$$r = 0 \Rightarrow UV = 1. \quad (6.110)$$

This hyperbola has two disconnected, one with $U, V > 0$ and the other with $U, V < 0$. The former corresponds to the singularity of the black hole and the latter the singularity of the white hole, see figure 22 We can define $T = \frac{1}{2}(U + V)$ and $X = \frac{1}{2}(V - U)$ as the vertical and

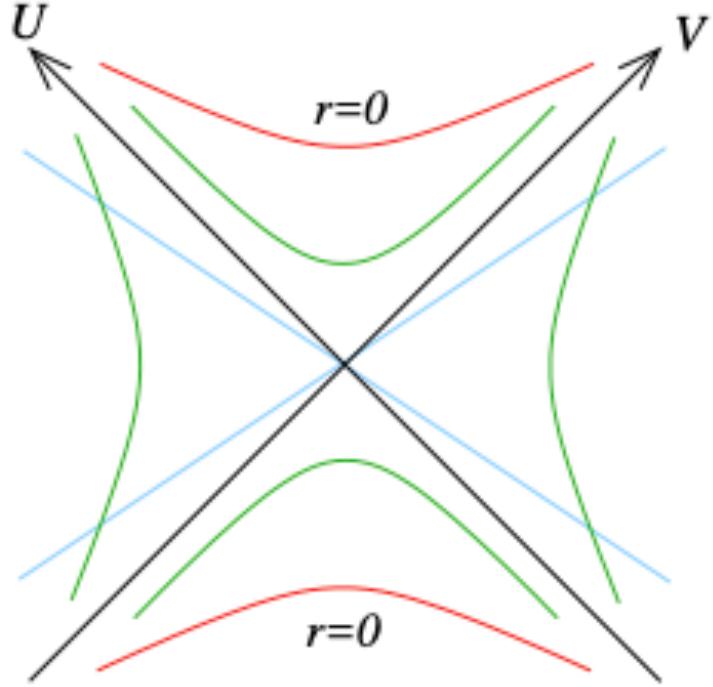


Figure 22: The Kruskal diagram. The U and V axes have been rotated 45° . They are the locations of the horizons at $r = 2G_N M$ and the red lines are the singularities at $r = 0$. Lines of constant r are in green and lines of constant t are in blue.

horizontal lines respectively. Lines of constant r are given by $UV = \text{constant}$ while lines of constant t are $U/V = \text{constant}$.

We see that the singularity is spacelike. Once you pass through the horizon the singularity lies in your future. You cannot avoid the singularity once you cross the horizon. Similarly the singularity of the white hole lies in the past, one could think of this as the singularity of the Big Bang.

We can understand three quadrants of the four. The right quadrant is the exterior of the black hole, the top quadrant is the black hole interior and the bottom quadrant is the interior of the white hole. The left hand quadrant is in fact another copy of the black hole exterior,

it is just covered by $U > 0$ and $V < 0$. To see this write

$$U = + \exp\left(-\frac{u}{4G_N M}\right), \quad V = - \exp\left(\frac{v}{4G_N M}\right). \quad (6.111)$$

Undoing all the coordinate transformations we see that this is precisely the metric of the Schwarzschild solution again.

Our spacetime contains two asymptotically flat regions joined by a black hole. Note that it is not possible for an observer to cross from one to the other, nor to send a signal from one region to the other. The causal structure of spacetime forbids this.

One could ask what the spatial geometry that connects the two regions is. Fix the $t = 0$ slice of Kruskal spacetime ($U = V = 0$). In our original Schwarzschild solution the spatial geometry is

$$ds^2 = \left(1 - \frac{2G_N M}{r}\right)^{-1} dr^2 + r^2 ds^2(S^2), \quad (6.112)$$

which is valid for $r > 2G_N M$. There is another copy of this that describes the geometry of the left-hand side and we can glue these two together at $r = 2G_N M$, giving a worm-hole like geometry. This is known as the *Einstein–Rosen bridge*. Before getting excited about travelling through the black hole you cannot travel through the worm-hole as the paths are space-like not time-like.

7 Cosmology

We have only considered one solution of Einstein’s equations so far in these lectures, we will consider another which describes the evolution of the universe. The basic idea behind this model is that the universe is pretty much the same everywhere. Since we inhabit an orbit close, in cosmological terms, to a star we do not see the similarity between our situation and the desolate cold of deep space and this assumption may seem somewhat crazy. This assumption is applied to the very largest scales, where local variations in density are averaged over. There are a number of observations which support this assumption. The most clear way of seeing this is by looking at the Cosmic Background Radiation (CMB), see figure 23. The microwave background radiation is not perfectly smooth but the deviations from regularity are of the order 10^{-15} or less. The radiation is consistent with that of a blackbody spectrum radiated in all directions. The spectrum has been redshifted due to the expansion of the universe and today the average temperature is $2.725K$.

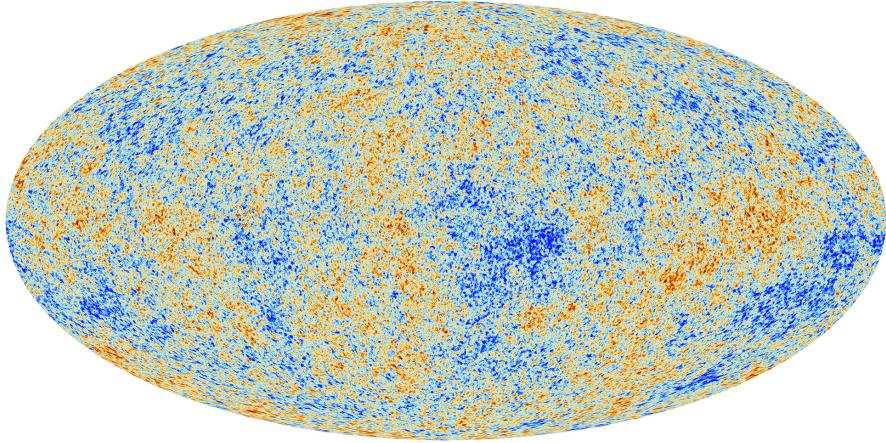


Figure 23: The anisotropies of the CMB as observed by Planck. It is a snapshot of the oldest light in the universe, coming from when the universe was just 380000 years old. It shows tiny temperature fluctuations that correspond to regions of slightly different densities and it is these regions which were the seeds for the stars and galaxies we see today. (Credit ESA for the picture).

7.1 FRW metric

We want to formalise this notion of the same everywhere in a more mathematical way. A manifold may have the properties of being *isotropic* and/or *homogeneous*: these are the necessary mathematical concepts which formalise our “same in every direction” comment.

Homogeneous A spacetime is spatially homogeneous if there exist a one-parameter family of space-like hypersurfaces Σ_t foliating spacetime, such that for each t and for any points $p, q \in \Sigma_t$ there exists an isometry of the spacetime metric $g_{\mu\nu}$ which takes p into q .

Isotropic A spacetime is isotropic at the point p if, for each pair of unit tangent vectors $X, Y \in T_p(M)$ there is an isometry which maps X to Y .

A spacetime can be isotropic around a point without being homogeneous. Conversely a spacetime can be homogenous without being isotropic ($\mathbb{R} \times S^2$ for example). If, however, a spacetime is isotropic around every point then it is homogeneous. Likewise if it is isotropic around any point, and homogeneous then it is isotropic everywhere.

Since there is ample observational data for isotropy (recall this is data about a point) and we are not so self-centred to think we are the centre of the universe we should assume that it is also homogeneous. The utility of these assumptions relies on the fact that a space which is both isotropic and homogeneous is maximally symmetric. (Think of isotropy as generalised rotations and homogeneity as generalised translations). This implies that the space has the

maximal number of Killing vectors. Now spacetime itself should not be maximally symmetric, we want it to evolve, instead we want spatial slices to be maximally symmetric.

For a maximally symmetric space with metric $g_{\mu\nu}$ the Riemann tensor takes the form

$$R_{\mu\nu\rho\sigma} = \kappa(g_{\mu\rho}g_{\nu\sigma} - g_{\mu\sigma}g_{\nu\rho}), \quad (7.1)$$

where κ is a normalised measure of the Ricci scalar

$$\kappa = \frac{R}{n(n-1)}, \quad (7.2)$$

which must be constant. These spaces are classified and for us the difference will arise in the sign of κ , either positive, negative or 0. We will consider our spacetime to be of the form $\mathbb{R} \times \Sigma$ with metric

$$ds^2 = -dt^2 + a^2(t)ds^2(\Sigma), \quad (7.3)$$

with t a time-like coordinate and $a(t)$ a function known as the *scale factor*. The metric used here which is free of cross terms with dt is known as *co-moving* coordinates. An observer who stays at fixed coordinate in Σ is said to be a *comoving observer*. Only a comoving observer sees the universe as isotropic. On Earth we are not quite comoving due to our motion around the sun.

We want a maximally symmetric 3d space, we can write the metric in the form

$$ds^2(\Sigma) = \frac{dr^2}{1-kr^2} + r^2(d\theta^2 + \sin^2\theta d\phi^2), \quad (7.4)$$

with $k = \{-1, 0, +1\}$.²⁵ The case $k = -1$ gives a constant negative curvature metric and is sometimes called *open*. The $k = 0$ case corresponds to no curvature on Σ and is sometimes called *flat*, while the case $k = +1$ corresponds to positive curvature and is sometimes called *closed*. Note that the $k = 1$ case is the only one which is compact (unless one makes certain identifications of the coordinates). We then have the metric

$$ds^2 = -dt^2 + a^2(t) \left[\frac{dr^2}{1-kr^2} + r^2(d\theta^2 + \sin^2\theta d\phi^2) \right]. \quad (7.5)$$

This is the *Friedmann–Robertson–Walker metric* (FRW).

To understand why $a(t)$ is called a scale factor consider the distance between two observers, one at $r = 0$ and another at $r = r_0$. Then the spatial distance between them is

$$d_{\text{prop}} = \int_0^{r_0} \sqrt{g_{rr}} dr = a(t) \int_0^{r_0} \frac{dr}{\sqrt{1-kr^2}} \equiv a(t)f(r_0). \quad (7.6)$$

²⁵Different values can be reduced to one of these three cases by redefining the radial coordinate r .

We see that the distance depends on the scale factor. We can look at the relative speed at which distance is changing with respect to time, we have

$$\dot{d}_{\text{prop}}(t) = \dot{a}(t)f(r_0) = \frac{\dot{a}}{a}d_{\text{prop}}(t) \equiv H(t)d_{\text{prop}}(t), \quad (7.7)$$

with

$$H(t) = \frac{\dot{a}(t)}{a(t)}, \quad (7.8)$$

the *Hubble parameter*. The value of the Hubble parameter at present is the Hubble constant H_0 . Current measurements give $H_0 = 70 \pm 10$ km/sec/Mpc. (Mpc is a megaparsec, $\sim 3.09 \times 10^{22} m$). Cosmology took off as a subject when the relative motions of the galaxies was first measured. We cannot actually determine the relative velocities of the galaxies now, i.e. at the same cosmological time, since we only have information about them at the time that the light left them. We are therefore not deducing $a(t)$ as it is now but rather as it was in the past. By looking at galaxies further away we can deduce the past history of $a(t)$.

7.1.1 Cosmological red-shift

Cosmological red-shift has a different origin to the gravitational red-shift we saw previously, however we can work it out in a similar manner. Assume that the light reaching us is on purely radial geodesics. Then we have

$$0 = -dt^2 + \frac{a(t)^2}{1-kt^2}dr^2, \quad (7.9)$$

and therefore

$$\frac{dt}{a(t)} = -\frac{dr}{\sqrt{1-kr^2}}, \quad (7.10)$$

where we picked the $-$ sign for the incoming radial geodesic (paths of decreasing r). The time emission t_1 and reception t_0 of the photon are given by

$$\int_{t_0}^{t_1} \frac{dt}{a(t)} = - \int_{r_0}^0 \frac{dr}{\sqrt{1-kr^2}} \equiv f(r_0). \quad (7.11)$$

Suppose that the next wave crest is emitted at time $t_1 + \delta t_1$ and received at $t_0 + \delta t_0$. Then since t is the proper time of stationary observers $\delta t_1 = \omega_1^{-1}$ and $\delta t_0 = \omega_0^{-1}$, with ω_i the frequency. Since the second photon leaves from r_0 and arrives at $r = 0$ it must also satisfy

$$\int_{t_1+\delta t_1}^{t_0+\delta t_0} \frac{dt}{a(t)} = f(r_0). \quad (7.12)$$

If δt_i are small then

$$f(r_0) = \int_{t_1+\delta t_1}^{t_0+\delta t_0} \frac{dt}{a(t)} = \left(\int_{t_1}^{t_0} + \int_{t_0}^{t_0+\delta t_0} - \int_{t_1}^{t_1+\delta t_1} \right) \frac{dt}{a(t)} \sim \int_{t_1}^{t_0} \frac{dt}{a(t)} + \frac{\delta t_0}{a(t_0)} - \frac{\delta t_1}{a(t_1)}$$

$$= f(r_0) + \frac{\delta t_0}{a(t_0)} - \frac{\delta t_1}{a(t_1)}. \quad (7.13)$$

Therefore we have

$$\frac{\delta t_0}{a(t_0)} \sim \frac{\delta t_1}{a(t_1)} \Rightarrow \omega_0 \sim \frac{a(t_1)}{a(t_0)} \omega_1. \quad (7.14)$$

The change in frequency is directly given by the ratio of the scale factors from when the light was emitted and when the light was received. The standard cosmologists definition of red-shift is through

$$z = \frac{\omega_1}{\omega_0} - 1 = \frac{a(t_0)}{a(t_1)} - 1. \quad (7.15)$$

Red shift is a direct measure of the change in separation of galaxies during the time the photon has taken to reach us. If a galaxy is at redshift 5 for example then it is 6 times further away than when the photon was emitted. Red shift does not give any direct information about the distance of the source, nor does it need to be faithful indicator of distance. Sources at different distances can have the same or similar red-shifts. If there was a period of time where the scale factor was essentially constant then any photons emitted during this period would appear to have the same red-shift. Similarly if there was a period of the scale factor decreasing then increasing again then sources at very different distances could give the same red-shift factor.

7.2 The Friedmann equations

Note that the Christoffel symbol $\Gamma^i_{tt} = 0$ and therefore the paths $\vec{x} = \text{const}$ are geodesics. The role of $a(t)$ is to change distances over time. There is a redundancy in the metric. If we rescale the coordinates as $a \rightarrow \lambda a$, $r \rightarrow \lambda$ and $k \rightarrow \lambda^{-2} k$ we leave the metric invariant. Of course now we are no longer fixed to take $k \in \{-1, 0, 1\}$. The non-zero components of the Ricci tensor are

$$\begin{aligned} R_{tt} &= -3 \frac{\ddot{a}}{a}, \\ R_{rr} &= \frac{a\ddot{a} + 2\dot{a}^2 + 2\kappa}{1 - kr^2}, \\ R_{\theta\theta} &= r^2(a\ddot{a} + 2\dot{a}^2 + 2\kappa), \\ R_{\phi\phi} &= r^2 \sin^2 \theta(a\ddot{a} + 2\dot{a}^2 + 2\kappa). \end{aligned} \quad (7.16)$$

It follows that the Ricci scalar is then

$$R = 6 \left[\frac{\ddot{a}}{a} + \left(\frac{\dot{a}}{a} \right)^2 + \frac{\kappa}{a^2} \right]. \quad (7.17)$$

The FRW metric is determined by the behaviour of $a(t)$. We want to plug this into Einstein's equations to derive the so called *Friedmann equations* which relates the scale factor to the energy-momentum of the universe. We choose to model the matter as a perfect fluid. If a fluid is isotropic in one frame and leads to an isotropic metric then it must be that the fluid is at rest in co-moving coordinates. The four-velocity is then

$$U^\mu = (1, 0, 0, 0), \quad (7.18)$$

and the energy momentum tensor is

$$T_{\mu\nu} = (\rho + p)U_\mu U_\nu + pg_{\mu\nu}. \quad (7.19)$$

With one index raised this becomes

$$T^\mu{}_\nu = \text{diag}(-\rho, p, p, p), \quad (7.20)$$

and the trace is

$$T = T^\mu{}_\mu = -\rho + 3p. \quad (7.21)$$

Before plugging into Einstein's equations it is useful to consider the conservation of the energy momentum tensor, in particular for the first component. We have

$$\begin{aligned} 0 &= \nabla_\mu T^\mu{}_0 \\ &= -\dot{\rho} - 3\frac{\dot{a}}{a}(\rho + p). \end{aligned} \quad (7.22)$$

7.2.1 Equation of state

To make progress we choose an equation of state, that is a relationship between p and ρ . The perfect fluids relevant to cosmology satisfy

$$p = w\rho, \quad (7.23)$$

with w a constant independent of time. The conservation of energy becomes

$$\frac{\dot{\rho}}{\rho} = -3(1+w)\frac{\dot{a}}{a}. \quad (7.24)$$

When w is constant this can be integrated to give

$$\rho \propto a^{-3(1+w)}. \quad (7.25)$$

For the vacuum to be stable²⁶ we need to pick $|w| \leq 1$. The two most popular cosmological fluids are known as *matter* and *radiation*.

²⁶This is beyond the scope of the course but one can read about this in Carroll chapter 4.

Matter is any set of collision-less non-relativistic particles which have zero pressure $p_M = 0$, i.e. $w = 0$. Examples include stars and galaxies for which the pressure is negligible. Matter also goes by the name of *dust* and universe whose energy density is mostly due to matter are known as *matter-dominated* universes. The energy density of matter falls off as

$$\rho_M \propto a^{-3}, \quad (7.26)$$

which is just interpreted as the decrease in number density of particles as the universe expands. For matter the energy density is dominated by the rest-energy which is proportional to the number density.

Radiation may be used to describe actual electromagnetic radiation or massive particles moving at relativistic velocities, close to the speed of light. The trace of the energy-momentum tensor of the electromagnetic field vanishes and therefore this fixes

$$p_R = \frac{1}{3}\rho_R \quad \Rightarrow \quad w = \frac{1}{3}. \quad (7.27)$$

In a *radiation dominated* universe the energy density falls off as

$$\rho_R \propto a^{-4}. \quad (7.28)$$

Thus the energy density of photons falls off slightly faster than that of matter. To understand why observe that the number density of photons decreases in the same way as for the slow moving massive particles, but in addition they lose energy due to cosmological red-shift of the previous section. When a is small radiation will dominate, while as a increases dust will dominate.

Vacuum energy also takes the form of a perfect fluid, that is a cosmological constant. In this case $p_\Lambda = -\rho_\Lambda$ and the energy density is constant,

$$\rho_\Lambda \propto a^0. \quad (7.29)$$

Since the energy density of both matter and radiation decreases as the universe expands if there is a non-zero vacuum energy it tends to dominate over the long term so long as the universe does not start contracting. If the vacuum energy begins to dominate then we say that the universe becomes *vacuum-dominated*. Examples of this are the maximal symmetric spaces de Sitter and anti-de Sitter.

7.2.2 Deriving the Friedmann equations

We can now substitute this into the Einstein equations

$$R_{\mu\nu} - \frac{1}{2}Rg_{\mu\nu} - \Lambda g_{\mu\nu} = 8\pi G_N T_{\mu\nu}. \quad (7.30)$$

The $\mu\nu = 00$ components give

$$\frac{3\dot{a}^2}{a^2} + \frac{3k}{a^2} - \Lambda = 8\pi G_N \rho, , \quad (7.31)$$

while the $\mu\nu = ij$ components give

$$\frac{2\ddot{a}}{a} + \left(\frac{\dot{a}}{a}\right)^2 + \frac{k}{a^2} - \Lambda = -8\pi G_N p. \quad (7.32)$$

There is only one distinct condition from the spatial part because of our isotropic assumption. From a linear combination of the two equations we find

$$\frac{\ddot{a}}{a} = \frac{\Lambda}{3} - \frac{4\pi G_N}{3}(\rho + 3p). \quad (7.33)$$

Note that the conservation of the energy momentum tensor,

$$\dot{\rho} + 3\frac{\dot{a}}{a}(\rho + p) = 0, \quad (7.34)$$

can be obtained from these two equations.

Equations (7.31) and (7.33) are known as the Friedmann equations and metrics of the form (7.5) satisfying these equations are FRW universes. If we know the dependence of ρ on a then the first can be solved.

7.3 Cosmological solutions

Let us consider some solutions. Before trying to solve anything let us analyse the behaviour of the function. With our equation of state the Friedmann equation becomes

$$\begin{aligned} \dot{a}^2 &= \frac{\Lambda a^2}{3} - k + \frac{8\pi G_N}{3} \rho a^2 \\ &= \frac{\Lambda a^2}{3} - k + \frac{C}{a^{1+3w}}, \end{aligned} \quad (7.35)$$

where C is a constant such that $8\pi G_N \rho = C a^{-3(1+w)}$.

We now want to analyse the form of $a(t)$. Note that qualitatively there is very little difference between dust and radiation, radiation is a little more dominant for small a but otherwise the overall structure is the same.

- For small a \dot{a}^2 is dominated by the term $C a^{-3(1+w)}$ and therefore $|\dot{a}| \rightarrow \infty$ as $a \rightarrow 0$.

This is then a period of rapid expansion or contraction. We have

$$\dot{a}^2 \sim a^{-3(1+w)}, \quad \Rightarrow \quad \dot{a} \sim \pm \sqrt{C} a^{-\frac{1+3w}{2}}, \quad (7.36)$$

which can be solved to give

$$a(t) \sim \text{constant} |t|^{\frac{2}{3(1+w)}}. \quad (7.37)$$

In both cases $a(t)$ will expand from zero to finite size, or collapse from finite size to 0 in finite time.

- For large a the behaviour depends on the sign of Λ and if this vanishes then on k .

We can now consider in more detail various cases.

7.3.1 Solutions with $k = 0$

Let us set $k = 0$. This is the most likely value for the current universe.

We can now distinguish the different behaviours depending on the sign of Λ , see figure 24.

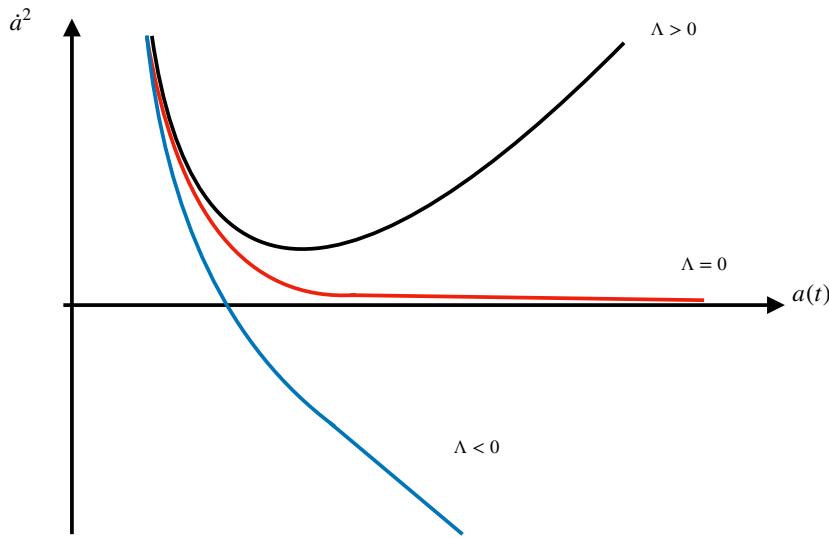


Figure 24: A plot of \dot{a}^2 for $k = 0$. Note the unphysical region for $\Lambda < 0$.

$\Lambda > 0$ For $\Lambda > 0$ \dot{a}^2 is never negative and therefore \dot{a} must always be positive or negative. For $\dot{a} > 0$ a starts off small with a rapid expansion which slows down to a minimum rate of expansion and then the rate of expansion increases again. See figure 25a.

For $\dot{a} < 0$ then the evolution is the opposite. a starts off large, collapsing quickly before the rate of collapse slows to a minimum before speeding up once again until the universe collapses again. See figure 25b.

$\Lambda = 0$ As before \dot{a}^2 is always positive so \dot{a} cannot change sign. For $\dot{a} > 0$ the universe starts off at zero size expands rapidly before the rate of expansion decreases, tending to zero but never reaching it. The opposite sign for \dot{a} is the time reversal of this. See figure 25c.

$\Lambda < 0$ In this case there is a critical value $a = a_c$ wt which $\dot{a} = 0$. One can show that at this point $\ddot{a} < 0$ and therefore if a is initially increasing it slows until it reaches a_c and then starts to decrease. The universe begins expanding before reaching a critical size before contracting again, all in finite time. See figure 25d.

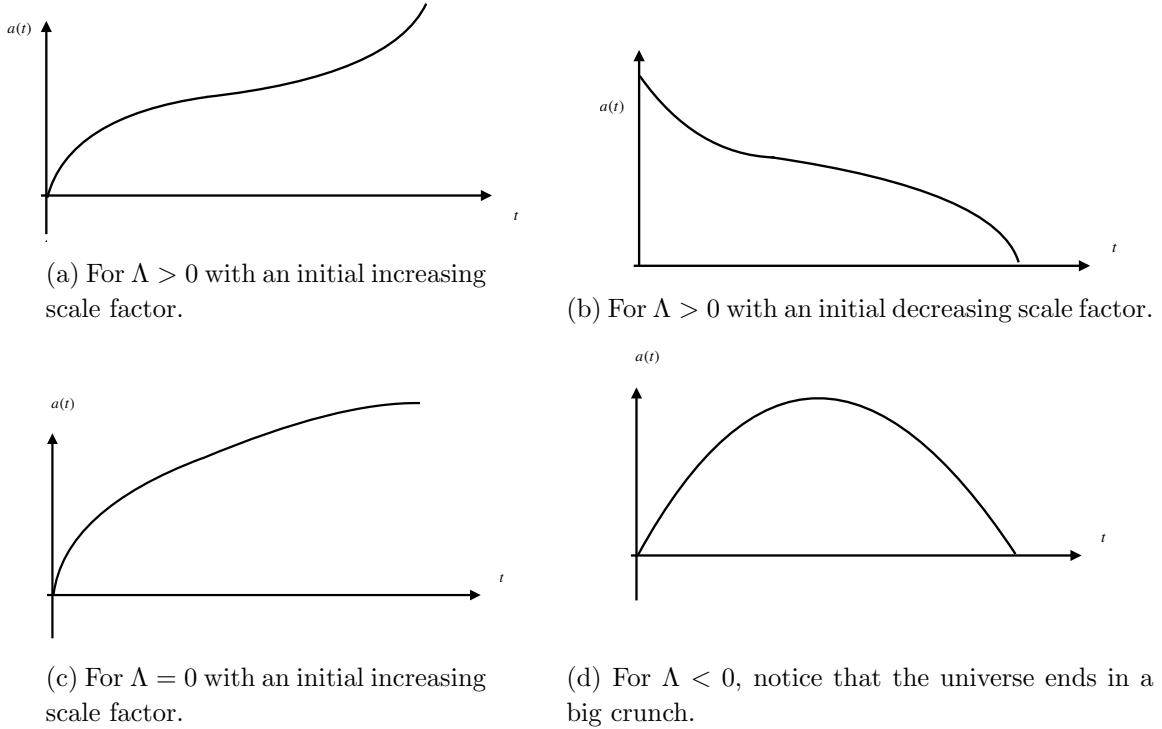


Figure 25: various plots of the scale factor for $k = 0$ and different choices of the cosmological constant.

We can in fact explicitly solve for $a(t)$. For dust, $w = 0$ we find

$$a(t) = \begin{cases} \left(\frac{3C}{\Lambda}\right)^{1/3} \sinh^{2/3}\left(\frac{\sqrt{3\Lambda}}{2}t\right) & \Lambda > 0, \\ \left(\frac{3\sqrt{C}}{2}\right)^{2/3} t^{2/3} & \Lambda = 0, \\ \left(-\frac{3C}{\Lambda}\right)^{1/3} \sin^{2/3}\left(\frac{\sqrt{-3\Lambda}}{2}t\right) & \Lambda < 0. \end{cases} \quad (7.38)$$

For radiation, $w = \frac{1}{3}$ we have

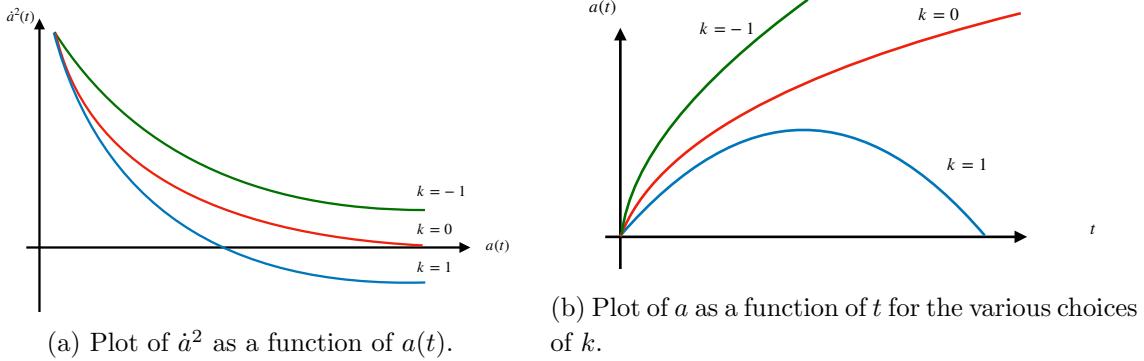
$$a(t) = \begin{cases} \left(\frac{3C}{\Lambda}\right)^{1/4} \sinh^{1/2}\left(\frac{\sqrt{3\Lambda}}{2}t\right) & \Lambda > 0, \\ \sqrt{2\sqrt{C}t^{1/2}} & \Lambda = 0, \\ \left(-\frac{3C}{\Lambda}\right)^{1/4} \sin^{1/2}\left(\frac{\sqrt{-3\Lambda}}{2}t\right) & \Lambda < 0. \end{cases} \quad (7.39)$$

7.3.2 Solutions with $\Lambda = 0$

We can now consider keeping k free, (well we can arrange for $k \in \{-1, 0, 1\}$ without loss of generality) and set the cosmological constant to vanish. We can again plot the qualitative features of $a(t)$.

- We have that for $k = 1$ there is a maximum value of a for which \dot{a}^2 is positive or zero and so we end up with an initial phase of expansion before reaching the critical value and then a subsequent contraction.
- If $k = 0$ or $k = -1$ then the universe continues to expand, but at different rates. For $a \rightarrow \infty$ we have that when $k = -1$ we have $\dot{a}^2 \rightarrow 1$ while for $k = 0$ we have $\dot{a} \rightarrow 0$.

We have plotted \dot{a}^2 in figure 26a while a is plotted in figure 26b.



One can again find full solutions to these equations however they are somewhat tedious to work out and best expressed in terms as parametric functions, for this reason we omit this.

7.3.3 The Big Bang

All of the solutions we have constructed have a region where $a = 0$. One can show that this is a generic feature of the Friedmann equations. From (7.33) we see that if the matter obeys the *strong energy condition*

$$\rho + 3p \geq 0, \quad (7.40)$$

then there is a singularity at a finite time t_{BB} where $a(t_{BB}) = 0$. This follows since the acceleration is necessarily negative. The universe is therefore decelerating, meaning it must have been accelerating faster at some point. If $\ddot{a} = 0$ then $a(t) = H_0 t + \text{const.}$

Suppose that $\ddot{a} = 0$, then $a(t) = H_0 t + \text{const.}$ This is the dotted line shown in figure 27. If this is the case then the Big bang occurs at $t_0 - t_{BB} = H_0^{-1}$. The strong energy condition ensures that $\ddot{a} \leq 0$ and so the dashed line provides an upper bound on the scale factor. In such a universe the Big Bang must occur at $t_0 - t_{BB} \leq H_0^{-1}$.

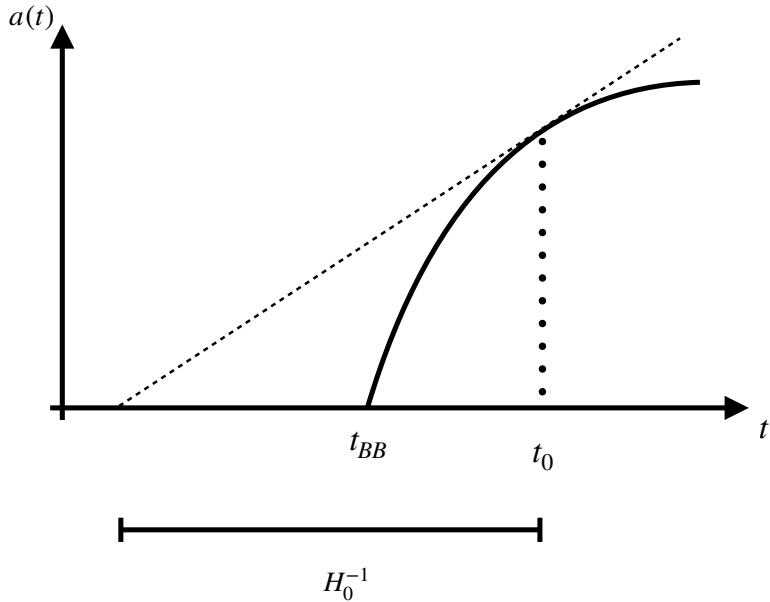


Figure 27: A plot of the scale factor showing the inevitability of the Big Bang.

The Big Bang refers to the creation of the universe from a singular state, not an explosion of matter into a pre-existing spacetime. One may wonder whether this singularity is an artefact of our choice of initial assumptions however it has been shown (by Hawking in his PhD thesis) that a singularity is a necessity even in the absence of such assumptions, given the strong energy condition.

The strong energy condition is obeyed by all conventional matter, including dust and radiation. However there are substances which violate it, leading to an accelerating universe. The single component pieces above still have a big bang however the above argument cannot rule out the possibility of more complicated solutions which avoid the Big Bang. In fact the leading theory at the moment is that in the very first moments after the Big Bang there was a period of exponential expansion.

All of the cosmological models we use predict a time in the past where the scale factor vanishes. The Big bang is a point in time not in space, it happens everywhere in space. We can get an estimate for the age of the universe by Taylor expanding $a(t)$ and truncating to linear order. Recall that we fixed $a(t_0) = 1$ then

$$a(t) \sim 1 + H_0(t - t_0). \quad (7.41)$$

This gives the estimate

$$t_0 - t_{BB} \sim H_0^{-1} \sim 4.4 \times 10^{17} s \sim 1.4 \times 10^{10} \text{ years}. \quad (7.42)$$

This is close to the 13.8 billion years which is widely accepted to be the age of the universe. Strictly speaking we should not trust the solution at $a(t_{BB}) = 0$ since the metric is singular there. Any matter in the universe will be squeezed into an infinite density object. In such a regime our classical equations are no longer any good and we need a quantum theory of gravity. Despite much effort such a theory of quantum gravity is lacking and so we are unable to answer many questions. Did time begin at t_{BB} ? Was there a previous phase of a contracting universe and we are another bounce?

7.3.4 Cosmological horizon

The existence of a special time t_{BB} means that there is a limit as to how far back we can look into the past. Let us set $t_{BB} = 0$ in the following.

The speed of light sets an upper bound on the local propagation velocity of any signal so at a given time t an observer at $r = 0$ can receive signals emitted at time t_1 only from radial coordinates $r < r_1$ where r_1 is the radial coordinate from which light signals emitted at time t_1 would just reach $r = 0$ at time t . We can determine r_1 as

$$\int_0^{r_1} \frac{dr}{\sqrt{1 - kr^2}} = \int_{t_1}^t \frac{dt'}{a(t')}. \quad (7.43)$$

If the t' integral diverges as $t_1 \rightarrow 0$ then it is in principle possible to receive signals emitted at sufficiently early times from any comoving particle in the universe. On the other hand if the t' -integral converges at $t_1 \rightarrow 0$ then our vision is limited by a so-called *particle horizon*: it is possible to receive signals from a comoving particles that lie within the radial coordinate $r_H(t)$ defined by

$$\int_0^{r_H(t)} \frac{dr}{\sqrt{1 - kr^2}} = \int_{t_1}^t \frac{dt'}{a(t')}. \quad (7.44)$$

The proper distance is

$$d_H(t) = a(t) \int_0^{r_H(t)} \frac{dr}{\sqrt{1 - kr^2}} = a(t) \int_{t_1}^t \frac{dt'}{a(t')} . \quad (7.45)$$

From (7.31) if ρ grows faster than $a^{-2-\epsilon}$ as $a \rightarrow 0$ then there will be a particle horizon.

We can play a similar game and ask if there are regions we will never see even if we wait long enough. If the t' integral diverges as $t \rightarrow \infty$ then in principle it is possible to receive signals from any event in the universe if we wait long enough. On the other hand if this is finite then it is only possible to receive signals for which

$$\int_0^{r_1} \frac{dr}{\sqrt{1 - kr^2}} \leq \int_{t_1}^{t_{\max}} \frac{dt'}{a(t')} . \quad (7.46)$$

Here t_{\max} can either be ∞ or the value of the next contraction to $a(t_{\max}) = 0$. This is known as an event horizon. It behaves in a similar way to falling inside the event horizon of a black hole, we will never be able to communicate with someone beyond the even-horizon.

This leads to some problems. We have assumed an isotropic universe, this is despite widely separated points being completely outside the event horizon of other points. Distinct patches of the CMB sky were causally disconnected. How then did they know ahead of time to coordinate their evolution (so that the CMB background looks isotropic) in the right way even though they were never in causal contact? One way of fixing this is by considering a period of inflation: an era of acceleration $\ddot{a} > 0$ in the very early universe, which is driven by some component other than matter or radiation.