

Elements of General Relativity

Roberto Casadio

July 9, 2018

Contents

1	Special relativity	1
1.1	Newtonian relativity	1
1.1.1	Observers and frames	1
1.1.2	Galilean transformations	2
1.1.3	Conservative forces	4
1.1.4	Electromagnetism	5
1.1.5	Alternative explanations	6
1.2	Foundations of special relativity	11
1.2.1	Two new principles	11
1.2.2	Newtonian space and time	11
1.2.3	Relativity of simultaneity and space-time	12
1.3	Relativistic kinematics	14
1.3.1	Lorentz transformations	14
1.3.2	Space-time diagrams	16
1.3.3	Addition of velocities	21
1.3.4	Invariance of the phase of a wave	23
1.3.5	Twin paradox	26
1.4	Old fashioned covariant formalism	28
1.5	Relativistic dynamics	32
1.5.1	Relativistic momentum and mass	32
1.5.2	Elastic collisions	33
1.5.3	Inelastic collisions	37
1.5.4	Equivalence of mass and energy	38
1.5.5	Relativistic force law	39
1.6	Electromagnetism	39
1.6.1	Electric charge and current	40
1.6.2	Transformations for \vec{E} and \vec{B}	42
1.6.3	Maxwell equations redux	44
1.6.4	Nature and relativistic fields	47
2	Differentiable manifolds and tensors	49
2.1	Differentiable manifolds	50
2.1.1	Manifolds and coordinates	50

2.1.2	Curves	54
2.1.3	Functions	55
2.1.4	Vectors and vector fields	55
2.1.5	Vector fields and integral curves	60
2.1.6	One-forms	66
2.1.7	Tensors and tensor fields	68
2.2	Length and angles	71
2.2.1	Metric tensor	72
2.2.2	Metric tensor field	74
2.3	Lie derivative and symmetry	75
2.3.1	Passive and active transformations	75
2.3.2	Congruences and Lie dragging	77
2.3.3	Lie derivatives	80
2.3.4	Symmetry and vector fields	87
2.4	Differential forms	92
2.4.1	P -forms	92
2.4.2	Area and volume	93
2.5	Covariant derivatives	97
2.5.1	Parallelism and covariant derivative	97
2.5.2	Geodesics	101
2.5.3	Riemann tensor and curvature	104
2.5.4	Metric connection	106
3	General Relativity	111
3.1	Arbitrary observers and gravity	111
3.2	Gravitational equations	116
3.2.1	Gravity and test particles	116
3.2.2	Source of gravity and Einstein equations	118
3.2.3	Classical tests of General Relativity	123
3.3	Black holes	125
3.3.1	The Schwarzschild metric	125
3.3.2	Radial geodesics	128
3.3.3	General orbits	131
3.3.4	Light-like geodesics	132
3.3.5	Gravitational red-shift	132
3.3.6	Radially infalling probe	134
3.3.7	The (event) horizon and black holes	137
3.4	Cosmology	139
3.4.1	Friedman-Robertson-Walker metric	140
3.4.2	Cosmic fluids	142
3.4.3	Friedmann equations	144
3.4.4	Cosmic Microwave Background	147
3.4.5	Cosmological redshift	148

3.4.6	Luminosity-distance relation	149
3.4.7	Hubble law	150
3.4.8	The Universe today	151
A	Symmetries and group theory	153
A.1	Abstract groups	153
A.2	Matrix representations and Lie groups	154
A.3	Rotations in N dimensions	157
A.3.1	Rotations in 2 dimensions: $SO(2)$ and $U(1)$	158
A.3.2	Rotations in 3 dimensions: $SO(3)$ and $SU(2)$	160
A.4	Lorentz group: $SO(3, 1)$	165
A.4.1	Irreducible representations: bosons and fermions	166
A.4.2	Poincaré group: $SO(4, 1)$	167

Chapter 1

Special relativity

We start by briefly reviewing the main concepts at the heart of classical Newtonian mechanics, and the inconsistencies that arise when trying to incorporate Maxwell's electromagnetism. A successful description of the latter will lead us to accept Special Relativity as the new general framework for studying the motion of objects with velocities comparable to the speed of light. However, this comes at the expense of one of Newton's greatest achievements: his universal law of gravity (for one of the fundamental interactions) is incompatible with the global Lorentz transformations, so that we gain one and lose one.

1.1 Newtonian relativity

The laws of Newtonian mechanics can be derived in a very specific framework, which nowadays goes under the rather generic name of “classical physics”, underlining which there are very precise (albeit often understated) notions of observers and observations.

1.1.1 Observers and frames

One of the key concepts in this course is that of the *observer*: physicists, more or less implicitly, divide the universe into the specific *object* of study (for example, a moving ball inside the room or planets around the sun) and the *observer*, all the rest being included in the co-called *environment*, whose effects on both object and observer are neglected (as a simplifying assumption). Much of the progress achieved in physics, during the last century or so, can be measured by our increasing ability to describe the object, but its origin is arguably related to improved descriptions of the observer, and the way the latter interacts (or affects) with the former (the *measurement*). In fact, mathematically, we are taught to think of an observer as a geometrical *reference frame*, whereas its physical (experimental) meaning is that of an *apparatus* to locate objects in space and time. Confusing the two meanings can be hazardous: many are the situations in which the actual apparatus can just cover a small portion of the observed phenomena and not all mathematical reference frames can be physically realised. One must therefore beware of the physical relevance of

mathematical computations carried out in mathematically convenient frames ¹.

Classical mechanics deals mostly with point-like objects and their motion. In Newtonian physics, extended bodies (which represents more realistic objects) are then just collections of points with mass kept together by some internal force, and the observer's details (along with those of the measurement process) are assumed to remain irrelevant: the observer is a space-filling “mathematical frame”, with the notion of absolute time attached, which can measure everything of the objects without affecting their physical status. In Special Relativity the situation is more complicated, since one starts to consider the observer as an apparatus and the measurement's limitations which follow, and becomes even more so in General Relativity, in which the physical localisation of (extended) objects is highly non-trivial.

We will not review the details here, but recall that the three laws of Newtonian mechanics introduce a family of preferred frames: the *inertial observers*. Their very definition is logically a loophole: the first principle defines an inertial observer given the notion of (absence of) forces, whereas the second principle defines the (effect of a) force given the notion of inertial observer, namely

$$\vec{F} = m \vec{a} , \quad (1.1.1)$$

where m is the mass and \vec{a} the acceleration of a body as measured by an inertial observer, whereas \vec{F} stands for the expression describing a specific force (like Newton's law of gravity or the Lorentz electromagnetic force). However, this is typical of the physicists' pragmatism. In practical terms, one considers observers (frames) of suitable size (not too big, nor too small) for the problem at hand and views them as inertial. For example, the bench in a laboratory is good enough an inertial frame for studying collisions of ping-pong balls, whereas the solar system is good enough to study the motion of the Earth around the sun ².

1.1.2 Galilean transformations

Once the notion of inertial frames is accepted, the principle of Newtonian (or Galilean) Relativity can be phrased as follows:

Galilean Relativity: “The laws of mechanics are the same for all inertial observers.”

This idea can be made mathematically more precise by introducing sets of suitable coordinate transformations.

Given two frames $S = \{x, y, z\}$ and $S' = \{x', y', z'\}$, the latter moving with velocity $\vec{v} = (v, 0, 0)$ with respect to the former, the coordinates transform according to

$$\begin{cases} x' = x - vt \\ y' = y \\ z' = z \end{cases} \quad \Leftrightarrow \quad \begin{cases} x = x' + vt' \\ y = y' \\ z = z' \end{cases} , \quad (1.1.2)$$

¹Mathematical convenience is still a powerful guideline, particularly in a theorist's mind.

²Ideally, one would still like to be able to define the concepts at the heart of our physical theories unambiguously. We shall see how General Relativity helps in this respect, at the end of the course.

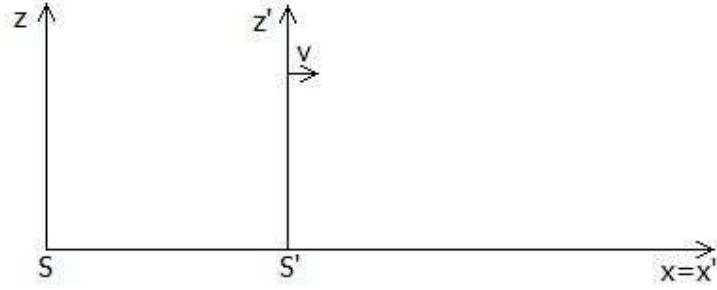


Figure 1.1: Parallel and transverse axes for the frames S and S' .

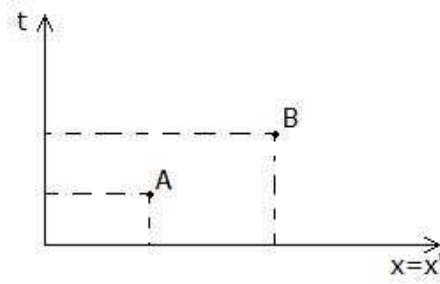


Figure 1.2: Two events in space-time.

in which we implicitly assumed *spatial homogeneity*, and also granted time is absolute: all observers can use the same synchronised clock(s), or, in mathematical terms

$$t' = t . \quad (1.1.3)$$

Let us first consider a (apparently very) simple problem: how to measure the length of, say, a rod. From the above transformations, it follows that, for two events A and B (points in space and time; see Fig. 1.2), one finds

$$t'_A - t'_B = t_A - t_B \quad \text{and} \quad x'_A - x'_B = x_A - x_B - v(t_A - t_B) . \quad (1.1.4)$$

If the events are the end-points of a rod, and coordinates have a meaning as lengths, the length of the rod is then the same in both frames provided *the measurements are taken simultaneously*, that is

$$t_A = t_B , \quad (1.1.5)$$

otherwise one could measure the position of the end-points at different times *only* in the frame in which the rod is at rest. This is a first elementary consideration that shows why the rest frames are physically *privileged*, in a practical sense.

Let us now consider a point-like object which can possibly change its position in space over time. From (1.1.2) and (1.1.3), assuming the (necessarily) finite differences in measurements

of the space and time coordinates can be approximated by their mathematical limit as derivatives, we immediately obtain the law of velocity composition

$$\left\{ \begin{array}{l} \frac{dx'}{dt'} = \frac{dx}{dt} - v \\ \frac{dy'}{dt'} = \frac{dy}{dt} \\ \frac{dz'}{dt'} = \frac{dz}{dt} \end{array} \right. \Leftrightarrow \quad \vec{u}' = \vec{u} - \vec{v} \quad (1.1.6)$$

and

$$\left\{ \begin{array}{l} \frac{du'_x}{dt'} = \frac{du_x}{dt} \\ \frac{du'_y}{dt'} = \frac{du_y}{dt} \\ \frac{du'_z}{dt'} = \frac{du_z}{dt} \end{array} \right. \Leftrightarrow \quad \vec{a}' = \vec{a} , \quad (1.1.7)$$

which implies

$$m \vec{a} = m \vec{a}' , \quad (1.1.8)$$

having also assumed all inertial observes measure the same value for the mass of the same object (*invariance of the mass*),

$$m' = m . \quad (1.1.9)$$

Note Eq. (1.1.8) is not yet enough to guarantee invariance of the second of Newton's laws. For that, we need to show the laws describing specific forces are also invariant.

1.1.3 Conservative forces

For conservative forces there exist, by definition, a potential energy $U = U(\vec{x})$ such that

$$\vec{F} = m \frac{d^2 \vec{x}}{dt^2} = -\vec{\nabla} U = -m \vec{\nabla} V , \quad (1.1.10)$$

where $\vec{\nabla} = (\partial_x, \partial_y, \partial_z)$ is the “gradient”. Consider, in particular, two particles respectively located at P_1 and P_2 [or $P_i = (x_i, y_i, z_i)$, with $i = 1, 2$], which interact through a potential

$$V = V(r) , \quad (1.1.11)$$

with

$$\begin{aligned} r &= \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2} \\ &= \sqrt{(x'_1 - x'_2)^2 + (y'_1 - y'_2)^2 + (z'_1 - z'_2)^2} = r' , \end{aligned} \quad (1.1.12)$$

in which we used the notion of absolute time to compute the “distance” r between the two particles *at the same time in both frames* S and S' ,

$$t_1 = t_2 = t'_1 = t'_2 . \quad (1.1.13)$$

Then, Eq. (1.1.4) implies $x_1 - x_2 = x'_1 - x'_2$, $y_1 - y_2 = y'_1 - y'_2$, $z_1 - z_2 = z'_1 - z'_2$ and the force acting on the particle at P_1 is obtained from

$$\vec{\nabla} V(r) \Big|_{r=r_1} = \frac{dV(r)}{dr} \vec{\nabla} r = \frac{dV(r')}{dr'} \vec{\nabla} r' = \frac{dV(r')}{dr'} \vec{\nabla}' r' = \vec{\nabla}' V(r') , \quad (1.1.14)$$

where $\vec{\nabla}' = (\partial_{x'_1}, \partial_{y'_1}, \partial_{z'_1}) = (\partial_{x_1}, \partial_{y_1}, \partial_{z_1}) = \vec{\nabla}$ when Eq. (1.1.13) holds. From Eq. (1.1.8), we can therefore conclude that the second law of classical mechanics is *Galilean invariant in form* for conservative forces, which means it takes the same mathematical form in all inertial frames,

$$\vec{F} = m \vec{a} = m \vec{a}' = \vec{F}' . \quad (1.1.15)$$

In particular, the above result implies that Newton’s law of gravity is Galilean invariant in form, since it can be derived from the potential

$$V_G = -\frac{GM}{r} . \quad (1.1.16)$$

Given the very good accuracy with which Eq. (1.1.16) describes the motion of planets and other objects in the solar system, and that gravity was the only known (fundamental) interaction at his times, Newton was practically right in assuming time is an absolute “concept”.

1.1.4 Electromagnetism

The electrostatic Coulomb force can be derived from a potential of the same form as the gravitational expression in Eq. (1.1.16). However, not all of the electromagnetic interactions admit a similar description. In fact, we shall here show the non-invariance of Maxwell’s equations (in particular, of the consequent wave equation for light propagation) under (Galilean) addition of velocities.

Let us first recall that the speed of light is related to the vacuum electromagnetic constants by

$$c = \frac{1}{\sqrt{\epsilon_0 \mu_0}} \simeq 3 \cdot 10^8 \text{ m/s} , \quad (1.1.17)$$

and does not apparently refer to any preferred frame or observer. It is therefore not clear whether the law of velocity composition (1.1.6) applies to light signals. As we shall see later on, Maxwell’s equations imply that such signals propagate according to the wave equation

$$\left(\frac{\partial^2}{\partial t^2} - c^2 \frac{\partial^2}{\partial x^2} \right) \Psi(t, x) = 0 , \quad (1.1.18)$$

where Ψ is any of the electromagnetic field components and, for simplicity, we assumed the wave is plane-symmetric (so that it carries no dependence on y and z). From Eqs. (1.1.2), using the chain rule, we have

$$\begin{aligned}\frac{\partial^2}{\partial t^2} &= \frac{\partial}{\partial t} \left(\frac{\partial t'}{\partial t} \frac{\partial}{\partial t'} + \frac{\partial x'}{\partial t} \frac{\partial}{\partial x'} \right) = \frac{\partial}{\partial t} \left(\frac{\partial}{\partial t'} - v \frac{\partial}{\partial x'} \right) = \left(\frac{\partial}{\partial t'} - v \frac{\partial}{\partial x'} \right) \left(\frac{\partial}{\partial t'} - v \frac{\partial}{\partial x'} \right) \\ &= \frac{\partial^2}{\partial t'^2} - 2v \frac{\partial}{\partial t'} \frac{\partial}{\partial x'} + v^2 \frac{\partial^2}{\partial x'^2},\end{aligned}\tag{1.1.19}$$

and

$$\begin{aligned}\frac{\partial^2}{\partial x^2} &= \frac{\partial}{\partial x} \left(\frac{\partial t'}{\partial x} \frac{\partial}{\partial t'} + \frac{\partial x'}{\partial x} \frac{\partial}{\partial x'} \right) = \frac{\partial}{\partial x} \left(\frac{\partial}{\partial x'} \right) = \left(\frac{\partial}{\partial x'} \right) \left(\frac{\partial}{\partial x'} \right) \\ &= \frac{\partial^2}{\partial x'^2}.\end{aligned}\tag{1.1.20}$$

Substituting into Eq. (1.1.18), we obtain (note the dimension of all operators is time^{-1})

$$\left[\frac{\partial^2}{\partial t'^2} - (c - v)^2 \frac{\partial^2}{\partial x'^2} \right] \Psi(t', x') = 2v \frac{\partial}{\partial x'} \left[(c - v) \frac{\partial}{\partial x'} + \frac{\partial}{\partial t'} \right] \Psi(t', x').\tag{1.1.21}$$

Because of the non-zero right hand side, the above form differs from (1.1.18) (for two observers with relative velocity $v \neq 0$). Note also that, by introducing a rescaled time variable $w = ct$ (with units of length), the above equation can be rewritten as

$$\left[\frac{\partial^2}{\partial w'^2} - \left(1 - \frac{v}{c}\right)^2 \frac{\partial^2}{\partial x'^2} \right] \Psi(w', x') = 2 \frac{v}{c} \frac{\partial}{\partial x'} \left[\left(1 - \frac{v}{c}\right) \frac{\partial}{\partial x'} + \frac{\partial}{\partial w'} \right] \Psi(w', x'),\tag{1.1.22}$$

so that, in the approximation in which $v/c \ll 1$ and negligible, the form (1.1.18) is recovered. This simple observation suggests that we are considering potential effects of order (at least) v/c , which is necessarily very small for experiments physically realisable on Earth.

1.1.5 Alternative explanations

As we shall see, Special Relativity modifies Galilean relativity in order to accommodate for electromagnetism. However, many experiments were conducted based on the assumption that Newton was right. For example, Maxwell's equations were modified in the so called “emission theory” and the existence of a preferred frame called the “aether” was inspected. We shall here limit ourselves to a short list of the main attempts:

1. Absolute frame (aether) and the Michelson-Morley experiment.
2. Aether-drag hypothesis and Fresnel formula.
3. Aether and Lorentz-Fitzgerald contraction.

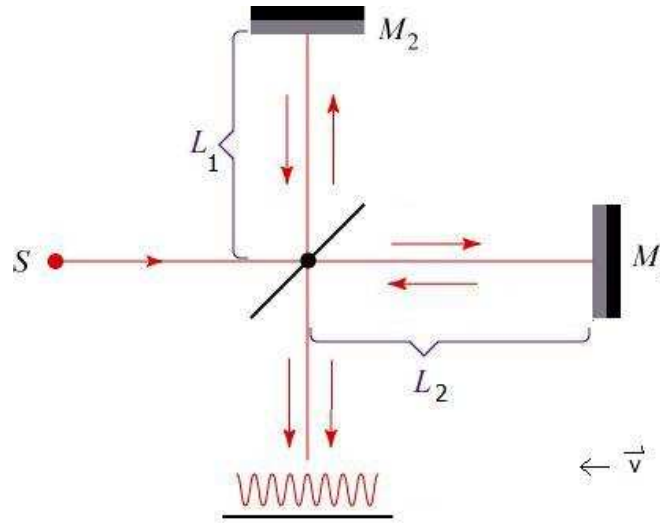


Figure 1.3: Michelson interferometer

Michelson-Morley experiment

One of the main experiments was conducted by Michelson and Morley in order to measure the velocity of moving bodies with respect to the aether. Earth moves around the sun at 30 km/s and the speed of light (1.1.17) is four orders of magnitude larger. The square of the ratio of these two speeds is therefore

$$\left(\frac{v}{c}\right)^2 \simeq 10^{-8}, \quad (1.1.23)$$

and very difficult to measure directly. Michelson thought that interference patterns might provide a convenient means, and realised what is now known as Michelson's interferometer. A schematic view is provided in Fig. 1.3, where S is a monochromatic light source, M_1 and M_2 are mirrors, and there are a semi-reflective mirror in the centre and a screen where interference patterns can be seen at the bottom.

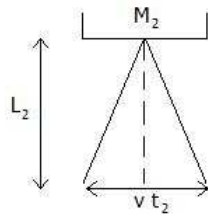


Figure 1.4: Light path for ray reflected by M_2 .

In a reference frame at rest with the apparatus, the aether would be moving with velocity \vec{v} , which we first assume is parallel to the arm of length L_1 . We can then determine the length of the return light paths along L_1 and L_2 , and corresponding travel times. Along L_1 ,

we have

$$t_1 = \frac{L_1}{c-v} + \frac{L_1}{c+v} = 2 \frac{L_1}{c} \frac{1}{1 - \frac{v^2}{c^2}} . \quad (1.1.24)$$

And, from Pythagoras' theorem (see Fig. 1.4), the traveling time along L_2 is given by

$$c t_2 = 2 \sqrt{\left(\frac{v t_2}{2}\right)^2 + L_2^2} \Rightarrow t_2 = \frac{2 L_2}{c} \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} . \quad (1.1.25)$$

From the above, we obtain the time difference

$$\Delta t = t_2 - t_1 = \frac{2}{c} \left(\frac{L_2}{\sqrt{1 - \frac{v^2}{c^2}}} - \frac{L_1}{1 - \frac{v^2}{c^2}} \right) . \quad (1.1.26)$$

Rotating the apparatus by 90 degrees, we likewise obtain

$$\Delta t' = \frac{2}{c} \left(\frac{L_2}{1 - \frac{v^2}{c^2}} - \frac{L_1}{\sqrt{1 - \frac{v^2}{c^2}}} \right) . \quad (1.1.27)$$

The two time differences differ by the amount

$$\Delta t' - \Delta t = \frac{2}{c} (L_1 + L_2) \left(\frac{1}{1 - \frac{v^2}{c^2}} - \frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} \right) \simeq \frac{L_1 + L_2}{c} \left(\frac{v}{c} \right)^2 , \quad (1.1.28)$$

which should result in a shift of the interference pattern on the screen. Since the wave period is $T = \lambda/c$, the shift is given by a (possibly fractional) number of wavelengths equal to

$$\Delta N = \frac{\Delta t' - \Delta t}{T} \simeq \frac{L_1 + L_2}{\lambda} \left(\frac{v}{c} \right)^2 . \quad (1.1.29)$$

In the original experiment, the two arms were 22 m long and light with a wavelength of $5.5 \cdot 10^{-7}$ m was used. One therefore expected $\Delta N = 0.4$, a fairly large quantity which was however not seen. The conclusion was then that $\vec{v} = 0$, which implied that either the aether did not exist or the earth moves along with it.

It is worth mentioning that the Michelson interferometer has survived as a useful apparatus to present. For example, the largest earth based gravitational detectors now active (LIGO and Virgo) are just an upscaled (albeit much refined) version of the Michelson's design.

Aether dragging hypothesis: Fizeau's experiment

Another hypothesis assumed that the Earth was at rest with the laboratory, the latter therefore being dragged along by moving bodies, so that rotating the Michelson interferometer would lead to no shift in the interference pattern.

If light moves in a medium of refractive index n , and the medium moves with respect to our reference frame, then Fresnel's empirical law³ was known to give the correct velocity of light in the laboratory frame,

$$v = \frac{c}{n} \pm v_w \left(1 - \frac{1}{n^2}\right), \quad (1.1.30)$$

in which v_w is the speed of the medium, for instance water flowing inside a pipe.

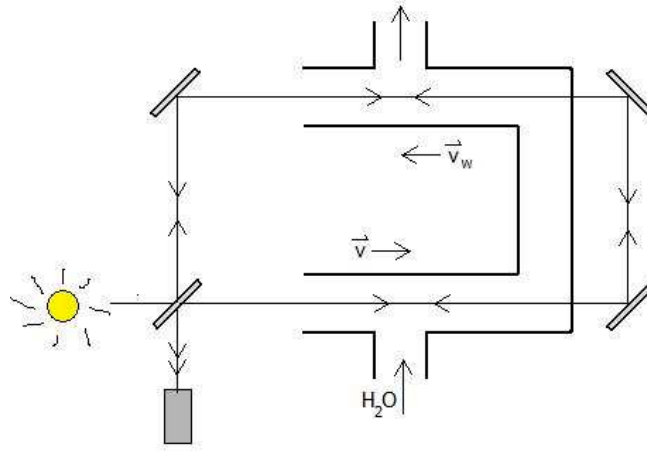


Figure 1.5: Fizeau experiment. Water enters the pipe from the bottom valve and exits from the top one, running counter-clockwise.

Fizeau tried to assume the aether is dragged by moving transparent objects with a size that fits in a laboratory. If this is the case, the second term inside the bracket in Eq. (1.1.30), that is $-1/n^2$, must be dropped. To prove his idea, he used an apparatus made by several mirrors and a pipe filled with running water (see Fig. 1.5). The experiment was intended to show that water and aether move with the same speed in the laboratory, but the effect was not observed (meaning Fresnel's law held) and Fizeau's hypothesis was discarded.

Aether dragging hypothesis: aberration of light

The results obtained with Fizeau's experiment could still be explained if the aether is dragged only by very large and massive celestial bodies, such as the Earth itself in its motion around the sun.

³A formula that describes specific effects but has not been derived from a theory or model of the system in consideration.

However, this explanation can be easily ruled out because of the so-called aberration of light, an effect astronomers had known for long: distant light sources, such as the stars in our galaxy, appear to move along ellipses during a solar year.

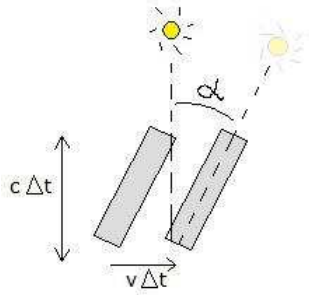


Figure 1.6: Aberration of light. The angle α depends on relative velocity of Earth and distant star.

The simple explanation for this effect is that, when aiming at a star, one should slightly incline the telescope so that the light traveling inside it is not absorbed by the sides before reaching the viewfinder (see Fig. 1.6). If the aether were dragged along with the Earth, so would be light rays and no such adjustment would be needed.

Lorentz-Fitzgerald contraction hypothesis

An alternative explanation of Michelson's results was that a body actually shrinks in the direction of motion with respect to the aether. Due to some complicated electron reaction, the actual length of a moving body would be related to its rest length by

$$L_1 = L_1^0 \sqrt{1 - \frac{v^2}{c^2}} , \quad (1.1.31)$$

One could then have

$$\Delta t = \frac{2}{c} (L_2^0 - L_1^0) \left(\frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} \right) = \Delta t' . \quad (1.1.32)$$

Michelson's experiment was therefore repeated using different angles. This should (supposedly) yield different speeds, v_1 and v_2 , with respect to the aether along the two arms, so that one expected a corrected shift given by

$$\Delta N \simeq 2 \frac{L_2^0 - L_1^0}{\lambda} \left(\frac{v_1^2}{c^2} - \frac{v_2^2}{c^2} \right) , \quad (1.1.33)$$

which again was not observed.

Before giving up the aether, it was hypothesised that the speed of light depends on the nature of the source and type of mirrors. Observations were conducted of distant binary stars

and corresponding variations of their orbits, which however were never confirmed. Finally, Michelson’s experiment was repeated using extra-terrestrial light sources, but no evidence of the existence of aether was ever found.

It is quite remarkable that the contraction (1.1.31) was meant to be mathematically described by the Lorentz’s transformations we shall see below, although in a conceptually very different context which denies the existence of a preferred medium. It is also curious, then, that the idea of a preferred physical frame has resurfaced much more recently in cosmology, as we shall see in Section 3.4.

1.2 Foundations of special relativity

1.2.1 Two new principles

In 1905, Einstein formulated two new principles:

The principle of relativity: “The laws of physics are the same for all inertial observers. No preferred inertial system exists.”

The principle of the constancy of the speed of light: “The speed of light in free space (vacuum) has the same value c in all inertial systems.”

These two postulates will lead to a re-thinking of our fundamental view of space and time. In particular, since space and time are related by the second postulate, the kinematics must be rebuilt from the onset, as we shall see in Section 1.3, where the arena of Minkowski space-time will be introduced. In this geometrical description, modified laws of dynamics will also arise and described in Section 1.5.

Before we proceed, we remark once more that giving up Galilean invariance means the Newtonian law of gravity cannot be correct (*i.e.*, we gain electromagnetism but lose gravity). Indeed, we shall see that from assuming c is the same in all inertial frames, one immediately finds that c is also the maximum attainable speed for all signals (in vacuum). Therefore, the idea of instantaneous interactions at a distance becomes fully questionable, since it would allow to move information faster than c , and conservative forces (described by a space-dependent potential function) are correspondingly banned.

1.2.2 Newtonian space and time

Let us first review the “space-time” diagrams for Newtonian mechanics. Instead of the “snapshot” type of diagram of Fig. 1.7, it is often more convenient to consider a diagram like the one in Fig. 1.8, in which space and time axes of the frame S are orthogonal and the axes of the moving system S' are represented as follows: the axis t' is given by the trajectory of the origin O' , that is by the condition $x' = 0$, from which one immediately finds

$$t = \pm \frac{x}{v} . \quad (1.2.1)$$

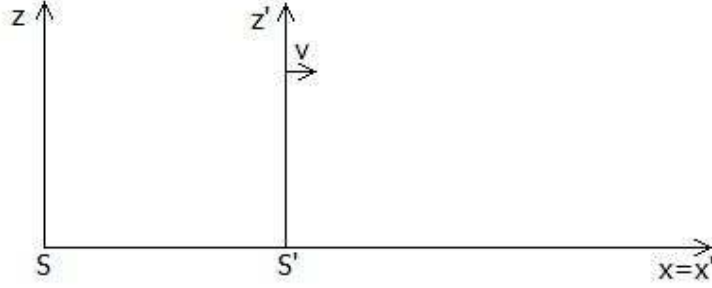


Figure 1.7: Parallel and transverse axes for the frames S and S' .

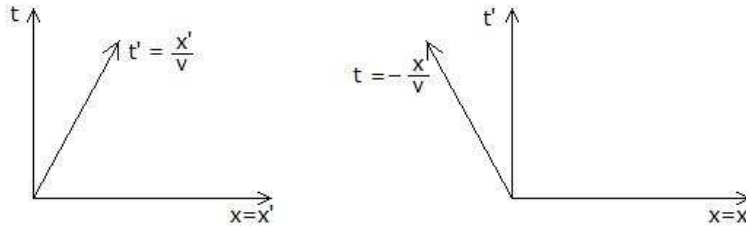


Figure 1.8: Space and time axes for the frames S and S' .

Likewise, the axis x' is represented by the condition $t' = \text{constant}$, or, usually, $t' = 0$. According to Galilean transformations, this means $t = 0$, so that the axis x' is parallel the axis x .

Now, consider that the axes x and x' are mathematical representations of a graduated rod with clocks attached. Therefore, as time evolves, both axes shift upward, with the origin O moving along the axis t , and the origin O' along t' . Note that the axes t and t' do not represent a physical apparatus that moves in the same sense, so that space and time remain *physically* distinct. In order to determine the coordinates of a given point (or event) A , one should move the axis x from its position at the time $t = 0$ until A lies on it. This will determine $x(A)$ and $t(A)$. In practice, it is more convenient to move backward (or forward) in time the point A as if it were at rest in S : the point A is projected onto the axis x at $t = 0$ parallelly to the axis t in order to determine $x(A)$, whereas $t(A)$ is determined by projecting A onto the axis t parallelly to the axis x . Likewise, in order to determine $x'(A)$ and $t'(A)$, one projects A parallelly to the axes t' and x' .

We shall next see how such diagrams change according to the principles of Special Relativity.

1.2.3 Relativity of simultaneity and space-time

In order to measure the length of a moving object, one must determine the positions of its two ends simultaneously. Let us consider two reference frames, say S and S' , moving with relative velocity \vec{v} directed along the $x = x'$ axes, and the five events A , B , C , D and E of Fig. 1.9: according to Galilean relativity, events represented by A , B and C occur at the

same position in S , thus they are displaced from each other in S' (see the projections parallel to the t' axis). The distance $D - E$ between two simultaneous events however appears the same in both frames, since time is absolute.

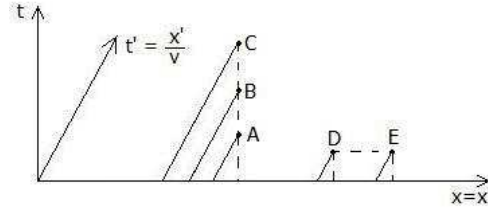


Figure 1.9: Events A , B , C , D and E as discussed in the main text.

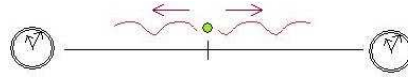


Figure 1.10: Synchronizing clocks in one frame.

Let us then analyze how two clocks placed at a fixed distance from each other can be synchronized, and keep in mind that no signal can travel faster than light (see Fig. 1.10). One can, for example, use two electromagnetic signals emitted *simultaneously* from a source placed at the midpoint between the clocks (or from a generic position, by taking in suitable account the signal time of travel). Note that the simultaneity of multiple emissions from *one* point can hardly be questioned by different observers, so that no ambiguity arises in this respect.

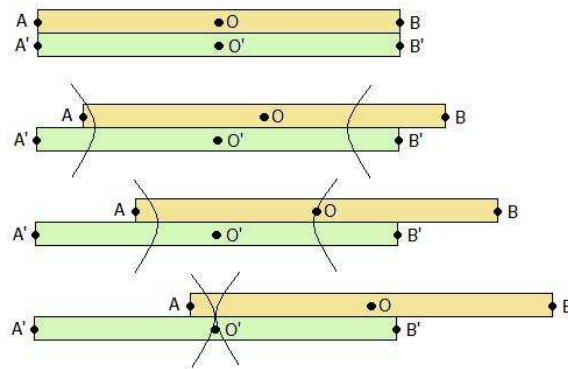


Figure 1.11: Synchronizing clocks in different frames.

The situation is however more tricky if we consider a frame in which the clocks are moving (see Fig. 1.11). For example, let us consider two light sources at the ends A and B of a bar which moves with constant velocity \vec{v} in S' and is at rest in S . Let the origins of the two systems, O and O' , coincide at the time $t = t' = 0$, with O at equal distance

from A and B . Precisely at that instant, a light signal is emitted from both sources, and will then be received from O and O' where there are two clocks. The two signals arrive in O' at the same time and an observer in S' will therefore conclude they have been emitted simultaneously. On the other hand, the clock in O will receive the two signals at different times and an observer in S will not say the emissions were simultaneous. One can however reverse the role of O and O' and the conclusion becomes then necessary that the very concept of simultaneity and clock synchronization is observer dependent.

1.3 Relativistic kinematics

The first step in the development of a relativistic kinematics is to determine the new coordinate transformations between two inertial observers. From these relations, several interesting consequences will follow.

1.3.1 Lorentz transformations

We shall here derive Lorentz transformations from Einstein's Principle of Relativity and the assumptions of *space-time homogeneity and isotropy* [2]. An important role will be also played by a *correspondence principle*: Newtonian mechanics must be recovered in the experimental contexts in which it is verified.

Let us consider two frames S and S' moving with constant relative velocity \vec{v} , and such that their origins coincide at $t = 0$. We further assume space is *homogeneous and isotropic* in both systems and all times, so that the change of coordinates is linear,

$$\begin{cases} x' = a_{11} x + a_{12} y + a_{13} z + a_{14} t \\ y' = a_{21} x + a_{22} y + a_{23} z + a_{24} t \\ z' = a_{31} x + a_{32} y + a_{33} z + a_{34} t \\ t' = a_{41} x + a_{42} y + a_{43} z + a_{44} t \end{cases} \quad (1.3.1)$$

where the coefficients a_{ij} may only depend on \vec{v} . In fact, suppose, for example, that

$$x' = \alpha x^2, \quad (1.3.2)$$

where the coefficient α has dimensions of length^{-1} . The observer S' would therefore see space as endowed with an intrinsic length α^{-1} , hence not homogeneous, as can be easily seen by considering displacements in the two frames, for example

$$x'_1 - x'_2 = \alpha (x_1^2 - x_2^2) \neq \alpha (x_1 - x_2)^2. \quad (1.3.3)$$

Note then that the above a_{ij} are all dimensionless, except for the a_{i4} , which have dimensions of a velocity. But we already know there is a fundamental velocity, the speed of light c , in our theory.

For $\vec{v} = 0$, we require that off-diagonal coefficients vanish and the $a_{ii} = 1$ (which can always be achieved by a suitable choice of time and length units in the two systems). Moreover, we expect to recover Galilean relativity for small relative speed $v \ll c$. From isotropy, we can always rotate S' so that the axes x and x' are parallel,

$$\begin{array}{ll} x\text{-axis : } & y = z = 0 \\ x'\text{-axis : } & y' = z' = 0 \end{array} \Rightarrow \begin{cases} y' = a_{22} y + a_{23} z \\ z' = a_{32} y + a_{33} z . \end{cases} \quad (1.3.4)$$

Moreover, the planes x - y and x' - y' are parallel as well, as are the planes x - z and x' - z' ,

$$\begin{array}{ll} xy\text{-plane : } & z = 0 \\ x'y'\text{-plane : } & z' = 0 \\ xz\text{-plane : } & y = 0 \\ x'z'\text{-plane : } & y' = 0 \end{array} \Rightarrow \begin{cases} y' = a_{22} y \\ z' = a_{33} z . \end{cases} \quad (1.3.5)$$

From these relations it follows that, if we place in S an object of length L with one end at the origin O and the other end at a point A on the y -axis, the coordinate of A in S' is $y' = a_{22} L$. If the same object were at rest in S' , the coordinate of the second end A would instead be $y = \frac{L}{a_{22}}$ in S (and analogously for objects placed on the z and z' -axis)⁴. One must therefore have

$$a_{22} = \frac{1}{a_{22}} = 1 = a_{33} = \frac{1}{a_{33}} . \quad (1.3.6)$$

From isotropy, we also expect $t' = a_{22}x + a_{44}t$ and, from the small velocity agreement with Galilean invariance, $x' = a_{11}(x - vt)$, so that

$$\begin{cases} x' = a_{11}(x - vt) \\ y' = y \\ z' = z \\ t' = a_{22}x + a_{44}t . \end{cases} \quad (1.3.7)$$

So far we have not yet considered the propagation of light and the principle of constancy of c . Suppose then that at $t = t' = 0$, a flash of light is emitted from the coinciding origins $O = O'$. The path of such a pulse is given in the two frames by

$$\begin{cases} x^2 + y^2 + z^2 = c^2 t^2 \\ x'^2 + y'^2 + z'^2 = c^2 t'^2 . \end{cases} \quad (1.3.8)$$

⁴Note we are implicitly assuming here that the length of an object measured at rest is an intrinsic property and does not depend on the observer.

Upon substituting for Eq. (1.3.7) in the second relation, we obtain

$$(a_{11}^2 - c^2 a_{41}^2) x^2 + y^2 + z^2 - 2(v a_{11}^2 + c^2 a_{41} a_{44}) x t = (c^2 a_{44}^2 - v^2 a_{11}^2) t^2 , \quad (1.3.9)$$

which must equal the first one, that is

$$\begin{cases} a_{11}^2 - c^2 a_{41}^2 = 1 \\ v a_{11}^2 + c^2 a_{41} a_{44} = 0 \\ c^2 a_{44}^2 - v^2 a_{11}^2 = c^2 . \end{cases} \quad (1.3.10)$$

We finally obtain the Lorentz transformations

$$\begin{cases} x' = \frac{x - v t}{\sqrt{1 - \frac{v^2}{c^2}}} \\ y' = y \\ z' = z \\ t' = \frac{t - \frac{v}{c^2} x}{\sqrt{1 - \frac{v^2}{c^2}}} . \end{cases} \quad (1.3.11)$$

Note that, as required, these transformation laws reduce to the Galilean ones for non-relativistic speed $v \ll c$ (correspondence principle). For later convenience, we also define a new time variable w with units of length⁵ and a parameter β ,

$$\begin{cases} \beta = \frac{v}{c} \\ w = c t \end{cases} \Rightarrow \begin{cases} x' = \frac{x - \beta w}{\sqrt{1 - \beta^2}} \\ y' = y \\ z' = z \\ w' = \frac{w - \beta x}{\sqrt{1 - \beta^2}} , \end{cases} \quad (1.3.12)$$

which makes the symmetry between space and time more apparent.

1.3.2 Space-time diagrams

We now introduce Minkowski space-time diagrams for two inertial frames, and use calibration hyperbolae to set units of length and time in both frames [3].

On the planes wx and $w'x'$, we can use the calibrating hyperbolae,

$$x^2 - w^2 = \pm 1 , \quad (1.3.13)$$

⁵Later in these notes, we shall set $c = 1$, corresponding to a more natural choice of units, and use the same symbol t for the time with dimensions of a length.

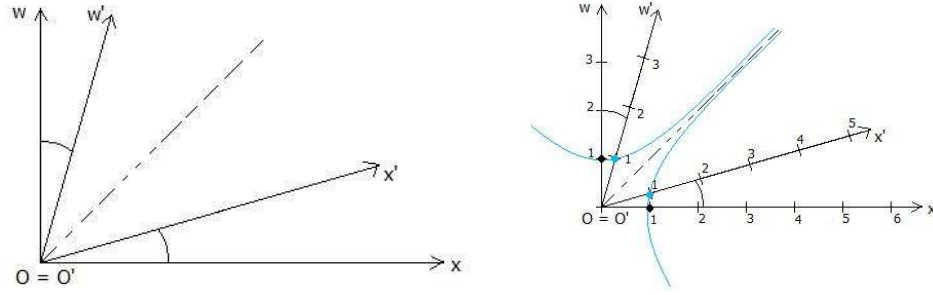


Figure 1.12: Space-time axes for S and S' (left panel) and calibration hyperbolae (right panel).

to relate length and time units in the two systems S and S' . For example, the point $P = (w = 0, x = 1)$, which represents the unit of length in S , is mapped by the hyperbola (with + sign) into a point $P' = (w' = 0, x' = 1)$ with the same meaning in S' . The value of x' can be easily determined by using the Lorentz transformations

$$0 \equiv w'(P') = \frac{w(P') - \beta x(P')}{\sqrt{1 - \beta^2}} \Rightarrow x(P') = \frac{w(P')}{\beta}, \quad (1.3.14)$$

and Eq. (1.3.13), which yield

$$\begin{cases} w(P') = \frac{\beta}{\sqrt{1 - \beta^2}} \\ x(P') = \frac{1}{\sqrt{1 - \beta^2}} \end{cases} \Rightarrow x'(P') = \frac{x(P') - \beta w(P')}{\sqrt{1 - \beta^2}} = 1, \quad (1.3.15)$$

as it should be. A similar argument for Eq. (1.3.13) with the $-$ sign leads to an analogous conclusion for the unit of time.

There are two famous effects predicted by Special Relativity which can now be easily derived from this graphical construction: length contraction and time dilation.

Length contraction

Let us consider a bar of length $L' = 1$ at rest in S' (see Fig. 1.13). After determining the units of length according to the previous section, one can easily show graphically that L' is always larger than the length L measured by the observer S in relative motion with respect to the object.

The length of an object, measured by an observer at rest with the object itself (S' in this case), is called “proper length” and defines an intrinsic property of the object which remains observer independent (all observers agree on the value of the proper length of an object!). However, the length of the same object as measured by an observer in relative motion with respect to it, S for example, will always be shorter and this effect is therefore known as “length contraction”.

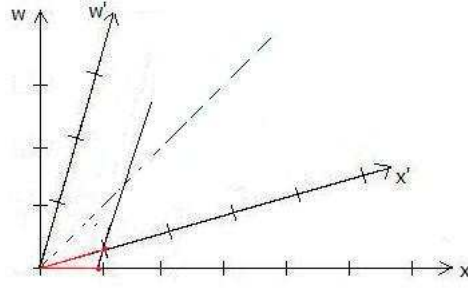


Figure 1.13: Length contraction.

It is easy to find the expression that describes analytically this effect. In S , we have $\Delta t = 0$, since the positions of the two ends of the bar are measured simultaneously therein, and Lorentz transformations (1.3.12) with $\Delta x' = L'$ then yield

$$\Delta x = L = L' \sqrt{1 - \beta^2}, \quad (1.3.16)$$

which, incidentally, represents a correction of order β^2 with respect to the Galilean result $L' = L$.

Time dilation

Let us now consider a time interval Δt in S : it is easy to show that Δt is always shorter than $\Delta t'$, the duration of the same time interval as measured by an observer in relative motion (see Fig. 1.14).

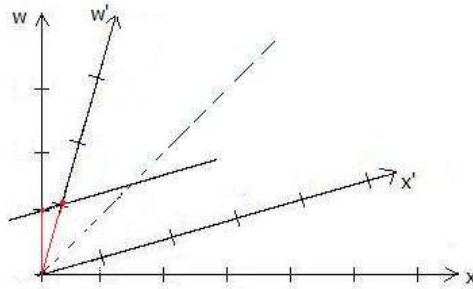


Figure 1.14: Time dilation.

In analogy with lengths, the time interval separating two events that occur at the same spatial location in a given reference frame is called “proper time”. The proper time therefore appears as the shortest possible measured time separating two given events, since any observer in relative motion will measure for the same separation a larger value. This effect is called “time dilation”.

It is also easy to find the expression that describes analytically this effect. In S , we have $\Delta x = 0$, since the two events occur at the same position, and Lorentz transformations

(1.3.12) with $\Delta t = T$ then yield

$$\Delta t' = T' = \frac{T}{\sqrt{1 - \beta^2}}, \quad (1.3.17)$$

which, analogously, represents a correction of order β^2 with respect to the Galilean result $T' = T$. Note that, unlike S , the two events marking the time interval do not occur at the same position in S' , that is $x'(0) \neq x'(T')$.

Minkowski diagrams

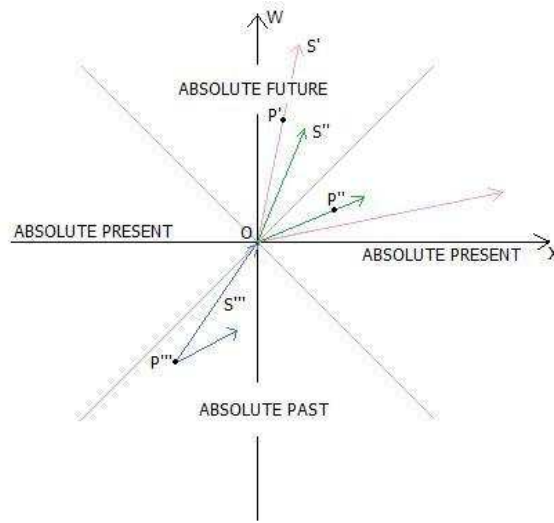


Figure 1.15: Minkowski diagram of two-dimensional space-time with the light cone of the origin.

The graphical representation of the two-dimensional space-time w - x is called Minkowski space. Light trajectories with $v = \pm c$ in this diagram are represented by straight lines at $\pm \frac{\pi}{4}$ rad, along which the units of time and length become infinite (in fact, the calibrating hyperbolae approach these lines asymptotically for $x \rightarrow \pm\infty$). Each pair of such lines starting from a given event say the origin O in S , form what is named the “light cone” of O and divide the space in three regions:

- Absolute future: for any given point P' inside this region, it is always possible to find a reference frame S' such that P' lies on the w' -axis. In this frame, P' occurs in the same spatial position as O but at a later time $t' > 0$. It is then easy to see that there is no inertial frame in which P' occurs before O .
- Absolute past: this is just the time-reverse of the absolute future: for any given point P''' inside this region, it is always possible to find a reference frame S''' such that P''' lies on the w''' -axis at time $t''' < 0$. It is likewise easy to see that there is no inertial frame in which P''' occurs after O .

- Absolute present: for any given point P'' inside this region, it is always possible to find a reference frame S'' , moving with speed β with respect to S , such that P'' lies on the x'' -axis. In this frame, P'' occurs at the same time as O , but at a different place $x'' \neq 0$. Moreover, if we chose a frame corresponding to a velocity smaller (larger) than β , we would obtain a frame in which P'' occurs after (before) O . In other words, there is no fixed temporal order between P'' and O .

Consider now a generic event P of coordinates (w, x) in a certain frame S . For the line segment \overline{OP} , one can have the three case:

- $w^2 - x^2 = \tau^2 > 0$

P is in the absolute future or past of O and \overline{OP} is said “time-like”. A physical signal can reach P starting from O (or conversely, depending on the time order between the two events).

- $w^2 - x^2 = \tau^2 < 0$

P is in the absolute present of O and \overline{OP} is said “space-like”. A physical signal cannot reach P starting from O (or conversely).

- $w^2 - x^2 = \tau^2 = 0$

P is on the light-cone of O and \overline{OP} is said “light-like” or “null”. Only an electromagnetic signal can travel from P to O , or conversely, depending on the time order between the two events.

It is also important to note that, as a consequence of Lorentz transformations, $\tau^2 = \tau'^2$ for any pair of systems S and S' , so that the quantity τ is invariant. If $\tau^2 > 0$, in a system in which $x = 0$, $\tau = t$ and is said the “proper time” between the origin and the given event.

Garage paradox

This is a neat example to explain the absolute present. Suppose we are at rest with a box of length L and a car of proper length ℓ_0 is moving towards it at a given speed β . With respect to the box, the car will have a length $\ell = \ell_0 \sqrt{1 - \beta^2}$. If $\ell = L$, the car will just fit in the box and the box door can be closed before the car hits the end of the box. However, in the frame of the car, the length of the box is

$$L' = L \sqrt{1 - \beta^2} = \ell_0 (1 - \beta^2) < \ell, \quad (1.3.18)$$

and the car will hit the end of the box “before” the door can be closed. Clearly, the front of the car hitting the box (event A) and the door being closed (event B) are in the present of each other (suggestion: check explicitly that no physical signal can connect A with B).

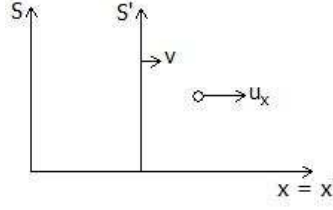


Figure 1.16: Addition of velocities.

1.3.3 Addition of velocities

Let us consider two frames S and S' moving with relative velocity \vec{v} along the x -axis. A particle moves with velocity \vec{u}'_x in S' , parallel to the axes x and x' (see Fig. 1.16). Galilean relativity predicts $u_x = u'_x + v$ in S . However, from the Lorentz transformations (1.3.12), one instead obtains

$$\Delta x = \frac{\Delta x' + v \Delta t'}{\sqrt{1 - \beta^2}} \quad (1.3.19)$$

$$\Delta t = \frac{\Delta t' + \frac{v}{c^2} \Delta x'}{\sqrt{1 - \beta^2}}, \quad (1.3.20)$$

so that ⁶

$$\begin{aligned} u_x &= \frac{\Delta x}{\Delta t} = \frac{\Delta x' + v \Delta t'}{\Delta t' + \frac{v}{c^2} \Delta x'} = \frac{\frac{\Delta x'}{\Delta t'} + v}{1 + \frac{v}{c^2} \frac{\Delta x'}{\Delta t'}} = \frac{u'_x + v}{1 + \frac{v}{c^2} u'_x} \\ &= \frac{u'_x + \beta}{1 + \beta u'_x}, \end{aligned} \quad (1.3.21)$$

It is important to note now that in the limit $\beta \rightarrow 1$ or $u'_x \rightarrow c$ (or both) one has $u_x \rightarrow c$, and there is therefore no way to go past the speed of light by changing (inertial) reference frame. On the other hand, for small velocities $|u'_x| \sim c \beta \ll 1$, we obtain

$$u_x \simeq u'_x + \beta, \quad (1.3.22)$$

in agreement with the Galilean expression.

Moreover, the orthogonal components u'_y e u'_z also change. Since $\Delta t \neq \Delta t'$, one obtains

$$u_y = u'_y \frac{\sqrt{1 - \frac{v^2}{c^2}}}{1 + \frac{v}{c^2} u'_x} = u'_y \frac{\sqrt{1 - \beta^2}}{1 + \beta u'_x} \quad (1.3.23)$$

$$u_z = u'_z \frac{\sqrt{1 - \frac{v^2}{c^2}}}{1 + \frac{v}{c^2} u'_x} = u'_z \frac{\sqrt{1 - \beta^2}}{1 + \beta u'_x}, \quad (1.3.24)$$

⁶Note that in the last expression we actually display $\Delta x / \Delta t$, which equals u_x if we set $c = 1$.

which, besides being an effect of order $(v/c)^2$, is still *qualitatively* different from the Galilean result $u_y = u'_y$ and $u_z = u'_z$.

Finally, the relativistic acceleration takes the following form

$$a_x = a'_x \left(\frac{\sqrt{1-\beta^2}}{1+\beta u'_x} \right)^3 \quad (1.3.25)$$

$$a_y = \frac{1-\beta^2}{(1+\beta u'_x)^2} \left(a'_y - a'_x \frac{\beta u'_y}{1+\beta u'_x} \right) \quad (1.3.26)$$

$$a_z = \frac{1-\beta^2}{(1+\beta u'_x)^2} \left(a'_z - a'_x \frac{\beta u'_z}{1+\beta u'_x} \right) . \quad (1.3.27)$$

Note that $\vec{a}' = 0$ only if $\vec{a} = 0$: although the value of the acceleration depends on the frame, it can only vanish in a frame if it is zero in all frames. The fact that an object is accelerated or not (that is, subject to a force or not) is still an absolute concept in Special Relativity as it was in the Galilean framework. This result is crucial, in that it allows for the very existence of *inertial observers*.

Fresnel formula

From Eq. (1.3.21), we can derive Fresnel's formula which was previously used to describe Fizeau's experiment. Since light and water move along the same direction, it is sufficient to replace $u' = c/n$ and $v = v_w$ and expand for v_w/c small. This yields

$$u \simeq \left(\frac{c}{n} + v_w \right) \left(1 - \frac{v_w}{c^2} \frac{c}{n} \right) \simeq \frac{c}{n} + v_w \left(1 - \frac{1}{n^2} \right) , \quad (1.3.28)$$

in which we neglected terms of order v_w/c or higher. Fresnel's formula is therefore explained as an approximation of the relativistic law of velocity addition.

Aberration of light

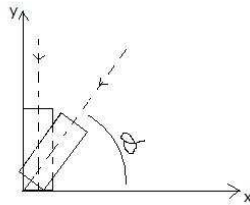


Figure 1.17: Aberration of light.

The relativistic addition law for velocities (1.3.21)-(1.3.24) also provides an easy description of the aberration of light. Since a star emits light in all directions, assuming it is at rest in

S' , the light rays have velocity

$$\begin{cases} u'_x = c \cos \theta' \\ u'_y = c \sin \theta' \end{cases} \Rightarrow \begin{cases} u_x = \frac{u'_x + v}{1 + \frac{v}{c^2} u'_x} \\ u_y = u'_y \frac{\sqrt{1 - \frac{v^2}{c^2}}}{1 + \frac{v}{c^2} u'_x} \end{cases}, \quad (1.3.29)$$

where S is the earth frame. It is now easy to determine the angle that a light ray forms with the x -axis in S ,

$$\tan \theta = \frac{u_y}{u_x} = \frac{\sin \theta' \sqrt{1 - \beta^2}}{\cos \theta' + \beta}. \quad (1.3.30)$$

In particular, for $\theta' = \pi/2$, we have

$$\tan \theta = \frac{\sqrt{1 - \beta^2}}{\beta} = \sqrt{\frac{c^2}{v^2} - 1} \simeq \frac{c}{v}. \quad (1.3.31)$$

Note that, for $v \rightarrow 0$, one then has $\theta \rightarrow \pi/2$ as expected.

1.3.4 Invariance of the phase of a wave

We shall now derive the relativistic Doppler effect from the invariance of the phase of a wave,

$$\Phi = \frac{2\pi}{\lambda} (x \cos \theta + y \sin \theta - \lambda \nu t) \rightarrow \Phi' = \frac{2\pi}{\lambda'} (x' \cos \theta' + y' \sin \theta' - \lambda' \nu' t) . \quad (1.3.32)$$

The invariance of the phase follows, for example, from the requirement that the number of cycles at the source point between two fixed times must be independent of the observer. Let us assume the source is located at $x = y = z = 0$ at all times, and starts to emit at $t = 0$. The number of oscillations of the source at a later time $t > 0$ is thus given by $N(t) = \Phi(\vec{0}, t)/2\pi$ for the observer S at rest with the source, and by $N'(t) = \Phi(\vec{x}'(t), t'(t))/2\pi$ for any other inertial observer S' , where $\vec{x}' = \vec{x}'(t)$ and $t' = t'(t)$ are the coordinates of the source in S' at the later time. If $N'(t) \neq N(t)$, there might exist an observer S' which does not see any cycle, and determinism would be totally lost, since S' would not see anything happen. A similar argument immediately leads to the conclusion that the number of oscillations between two events (space-time points) A and B along the path of a light signal must be independent of the observer (although the coordinates of the two space-time points of course do depend on the given observer). Since this number is given by $[\Phi(A) - \Phi(B)]/2\pi$, that phase difference must be independent of the observer.

The phase of a wave is an example of what is called a *scalar* quantity. By common definition, a scalar is a quantity which, under a change of coordinates $x \rightarrow x'$, changes according to

$$\Phi'(x') = \Phi(x) , \quad (1.3.33)$$

where Φ' denotes a possibly different functional form with respect to Φ (for the phase above, $\Phi' = \Phi_{\lambda'\nu'}$ and $\Phi = \Phi_{\lambda\nu}$). Consider, for example, a quantity represented by the real functions Φ and Φ' of the real axis in two different frames S and S' . The coordinates in the two frames are related by the transformation $x \rightarrow x'(x)$, meaning that the same point P has coordinate x in S and x' in S' (*passive transformation*). But we can also consider this transformation as one mapping the point P to a different point Q in one coordinate frame, say S (*active transformation*), where $x'(P) = x(Q)$. In the passive interpretation, we must then have

$$\Phi'(x'(P)) = \Phi(x(P)) , \quad (1.3.34)$$

meaning that the quantity at a given point P has the same value for both observers, S and S' . Likewise, in the active interpretation,

$$\Phi'(x'(P)) = \Phi'(x(Q)) = \Phi(x(P)) , \quad (1.3.35)$$

for exactly the same reason: the quantity we are measuring conserves its value even if the point is moved. Note that both Eqs. (1.3.34) and (1.3.35) can formally be written as the defining Eq. (1.3.33). This might look somewhat confusing, and we shall indeed spend a good deal of time later on in the course to clarify such transformation laws. For now, we just need to consider the passive interpretation and no confusion should arise.

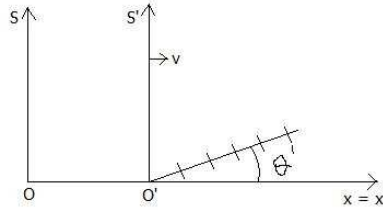


Figure 1.18: Emission of light and Doppler effect.

Let us now go back to the problem of wave transmission. Consider the two usual inertial frames S and S' moving with relative velocity \vec{v} along the x -axis. A light source is placed at O' which emits plane waves in the direction forming an angle θ' with respect to the axis x' (see Fig. 1.18). In S' , the signal is described by the wave-function

$$\Psi'(x', y', t') = A' \cos \left[\frac{2\pi}{\lambda'} (x' \cos \theta' + y' \sin \theta' - \lambda' \nu' t') \right] \equiv A' \cos(\Phi') . \quad (1.3.36)$$

The invariant number of cycles between the origin $t' = t = x' = x = y' = y = z' = z = 0$ (where we can assume the source is located at the time of emission) and a second arbitrary point of coordinates (t, x, y, z) in S [and equivalent to (t', x', y', z') in S'] on the wave path is represented by the difference between the arguments of the cosine evaluated at the two points. We thus must have

$$\Phi_{\lambda\nu}(x, y, z) = \Phi_{\lambda'\nu'}(x', y', z') , \quad (1.3.37)$$

or, neglecting the irrelevant direction z ,

$$\frac{2\pi}{\lambda'}(x' \cos \theta' + y' \sin \theta' - \lambda' \nu' t') = \frac{2\pi}{\lambda}(x \cos \theta + y \sin \theta - \lambda \nu t) , \quad (1.3.38)$$

where λ and $\nu = c/\lambda$ are to be determined. Clearly, this equality implies

$$\frac{x' \cos \theta' + y' \sin \theta'}{\lambda'} - \nu' t' = \frac{x \cos \theta + y \sin \theta}{\lambda} - \nu t . \quad (1.3.39)$$

From the Lorentz transformations from S' to S in Eq. (1.3.12), we then obtain

$$\frac{\cos \theta' + \beta}{\lambda' \sqrt{1 - \beta^2}} x + \frac{\sin \theta'}{\lambda'} y - \frac{\beta \cos \theta' + 1}{\sqrt{1 - \beta^2}} \nu' t = \frac{x \cos \theta + y \sin \theta}{\lambda} - \nu t . \quad (1.3.40)$$

Upon equating the coefficients of x , y and t , we finally obtain the laws of transformation for the frequency $\nu = c/\lambda$,

$$\nu = \nu' \frac{1 + \beta \cos \theta'}{\sqrt{1 - \beta^2}} , \quad (1.3.41)$$

and the angle θ ,

$$\begin{cases} \frac{\cos \theta}{\lambda} = \frac{\cos \theta' + \beta}{\lambda' \sqrt{1 - \beta^2}} \\ \frac{\sin \theta}{\lambda} = \frac{\sin \theta'}{\lambda'} . \end{cases} \quad (1.3.42)$$

The first of the above expressions gives the relativistic Doppler effect, whereas the other two describe the aberration of light.

Doppler effect

If we set $\theta' = 0$ (source approaching observer, $\cos \theta' = +1$) or π (source moving away from observer, $\cos \theta' = -1$) in Eq. (1.3.41), we obtain the *longitudinal Doppler effect*,

$$\nu_L = \nu' \frac{1 \pm \beta}{\sqrt{1 - \beta^2}} = \nu' \sqrt{\frac{1 \pm \beta}{1 \mp \beta}} , \quad (1.3.43)$$

whereas for $\theta' = \pm \frac{\pi}{2}$ (so that $\cos \theta' = 0$), we obtain the *transversal Doppler effect*,

$$\nu_T = \frac{\nu'}{\sqrt{1 - \beta^2}} . \quad (1.3.44)$$

Note that the longitudinal effect reproduces the Newtonian result at leading order (first order in $\beta = v/c$),

$$\nu_L = \nu' (1 \pm \beta) , \quad (1.3.45)$$

whereas the transversal effect is of order β^2 and, in fact, was not known in Newtonian dynamics.

Aberration of light

From Eqs. (1.3.42), we easily obtain

$$\tan \theta = \frac{\sin \theta' \sqrt{1 - \beta^2}}{\cos \theta' + \beta} . \quad (1.3.46)$$

Setting, for example, $\theta' = -\pi/2$ (light emitted by S' “straight down”), we obtain the angle

$$\tan \theta = -\sqrt{\frac{1}{\beta^2} - 1} \simeq -\frac{1}{\beta} \quad (1.3.47)$$

so that the telescope must be tilted in S by an angle equal to $\theta - \pi$ (upward).

1.3.5 Twin paradox

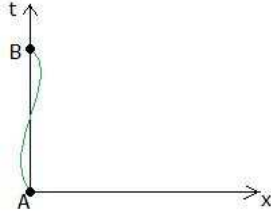


Figure 1.19: Twin paradox.

Let us consider twins (or two identically prepared clocks) initially placed in the origin O of our reference frame [4]. One of the twins leaves the other at $t = 0$ (the event A in Fig. 1.19), moves away (mostly) along an inertial trajectory, then reverts its direction of motion and comes back again (mostly) at constant velocity until it meets his twin (at point B). What is the difference in time measured between separation (A) and reunion (B) by the two twins?

For the twin who stays at O , the total proper time between the two events is simply $\Delta t = t_B - t_A$. For the twin who travelled, we could approximate its trajectory with a sequence of connected, sufficiently short, straight lines in the Minkowski diagram, and apply the law of time dilation (1.3.17) to each piece. This way, we would find that the proper time measured by the travelling twin, $\Delta\tau = \tau_b - \tau_A$ is necessarily shorter than Δt . The “paradox” arises when one tries to switch the points of view of the two twins, and neglects the fact that only one of them can be represented by an inertial frame. If both could be inertial, then we would also find $\Delta\tau > \Delta t$. In fact, deciding which one is an inertial observer would require a proper definition of being an inertial observer, which is still missing. However, once we have decided the twin at rest on earth is inertial, the travelling twin cannot be inertial, since it necessarily undergoes periods of acceleration (when it leaves earth, when it reverts its direction of motion and when he finally comes back to earth).

Approximating the trajectory of the travelling twin with inertial segments is rather a cumbersome way of doing the math, and we would instead like to have a general method

to compute the proper time along any trajectory. To this purpose, let us first review the two-dimensional Euclidean metric and rotations. The action of the latter on two-dimensional vectors (“spatial displacements”) $\vec{V} = (\Delta x, \Delta y)$ can be represented by a 2×2 matrix R , which defines the linear transformation

$$\begin{pmatrix} \Delta x' \\ \Delta y' \end{pmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix}. \quad (1.3.48)$$

From $(AB)^T = B^T A^T$,

$$(\Delta x', \Delta y') = (\Delta x, \Delta y) \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad (1.3.49)$$

so that

$$\begin{aligned} (\Delta x'_1, \Delta y'_1) \begin{pmatrix} \Delta x'_2 \\ \Delta y'_2 \end{pmatrix} &= (\Delta x_1, \Delta y_1) \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{pmatrix} \Delta x_2 \\ \Delta y_2 \end{pmatrix} \\ &= (\Delta x_1, \Delta y_1) \begin{pmatrix} \Delta x_2 \\ \Delta y_2 \end{pmatrix}, \end{aligned} \quad (1.3.50)$$

since

$$R^T R = R^T \mathbb{I} R = \mathbb{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (1.3.51)$$

The Cartesian scalar product is therefore invariant under rotations

$$\Delta x_1 \Delta x_2 + \Delta y_1 \Delta y_2 = \Delta x'_1 \Delta x'_2 + \Delta y'_1 \Delta y'_2. \quad (1.3.52)$$

Introduce next the Lorentz boost, for example along the direction x , which acts on space-time displacements $(\Delta x, \Delta w)$ as

$$\begin{pmatrix} \Delta x' \\ \Delta w' \end{pmatrix} = \begin{bmatrix} \frac{1}{\sqrt{1-\beta^2}} & \frac{-\beta}{\sqrt{1-\beta^2}} \\ \frac{-\beta}{\sqrt{1-\beta^2}} & \frac{1}{\sqrt{1-\beta^2}} \end{bmatrix} \begin{pmatrix} \Delta x \\ \Delta w \end{pmatrix} \quad (1.3.53)$$

and the Minkowski metric

$$\eta = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}. \quad (1.3.54)$$

It is easy to check that

$$M^T \eta M = \eta \quad (1.3.55)$$

and the invariant quantity is now precisely the proper time

$$\Delta \tau^2 = \Delta w_1 \Delta w_2 - \Delta x_1 \Delta x_2 = \Delta w'_1 \Delta w'_2 - \Delta x'_1 \Delta x'_2 = \Delta \tau'^2. \quad (1.3.56)$$

We remark once again that results of measurements are *scalars*: what is measured by a given observer cannot be argued by others. For example, if S measures a rod's length is L , no observers S' can claim S saw the rod is long $L' \neq L$, although they would measure the rod's length is L' . Proper time is a particular case: it is the time measured by a comoving observer. We will see that things become even more complicated in General Relativity where a measurement involves (at least) two scalars: the quantity measured by a given observer and a scalar that defines the space-time point where the measurement is taken.

Now that we have introduced the Minkowski metric, it is easy to see that, given the two events A and B , any physical trajectory connecting them must be represented by a time-like curve, along which the proper time is determined by

$$\tau_B - \tau_A = \int_A^B \sqrt{dw^2 - dx^2} = \int_A^B \sqrt{1 - \left(\frac{dx}{dw}\right)^2} dw \leq c(t_B - t_A) , \quad (1.3.57)$$

in which we replaced $\Delta x \rightarrow dx$ and $\Delta t \rightarrow dt$, $x = x(t)$ and w (or t) being the coordinates along the chosen trajectory in the inertial frame where A and B occur at the same position ($x_A = x_B = 0$). It follows that the straight line representing the twin who remained in the origin is the longest possible proper time $t_B - t_A$: the twin who travelled will necessarily be younger by the time they meet again.

1.4 Old fashioned covariant formalism

The twin paradox has led us to a new concept of “distance” given by the proper time of observers, and “vectors” in the Minkowski space-time. In fact, one can formulate a mathematical description entirely based on the Lorentz invariance of this distance and quantities which are “well-behaved” under Lorentz transformations (for more mathematical details about the Lorentz group, see Appendix A).

We first introduce a compact notation for (four-)vectors:

$$V^\mu = (V^0, V^1, V^2, V^3) = (V^t, V^x, V^y, V^z) = (V^0, V^i) = (V^0, \vec{V}) , \quad (1.4.58)$$

in a specific frame S with coordinates $\{t, x, y, z\}$. Consider then a *linear* coordinate transformation to another frame S' , that is

$$x^{\mu'} = \sum_{\nu=0}^3 M^{\mu'}_{\nu} x^{\nu} \equiv M^{\mu'}_{\nu} x^{\nu} = \left(\frac{\partial x^{\mu'}}{\partial x^{\nu}} \right) x^{\nu} \Leftrightarrow x' = M x , \quad (1.4.59)$$

where repeated indices always appear one up and one down and are implicitly summed over (Einstein's notation), and the matrix M is invertible, that is

$$\det(M) \neq 0 \quad \Rightarrow \quad \exists M^{-1} = \left(\frac{\partial x^{\mu}}{\partial x^{\nu'}} \right) \equiv M^{\mu}_{\nu'} , \quad (1.4.60)$$

such that

$$M^{\mu'}_{\mu} M^{\mu}_{\nu'} = \delta^{\mu'}_{\nu'} \quad \text{and} \quad M^{\mu}_{\nu'} M^{\nu'}_{\nu} = \delta^{\mu}_{\nu} , \quad (1.4.61)$$

where δ_ν^μ (respectively $\delta_{\nu'}^{\mu'}$) is the Kronecker delta in S (respectively S').

We can now define the following quantities:

Scalars: Quantities that “do not change”, like *numbers* and functions that satisfy

$$\Phi'(x') = \Phi(x) . \quad (1.4.62)$$

All scalars we have seen so far are numbers (the proper mass and charge of a particle), except for the phase of a wave (which is a scalar field).

Vectors: Quantities V^μ that transform like the coordinates:

$$V^{\mu'} = M_{\nu}^{\mu'} V^\nu = \left(\frac{\partial x^{\mu'}}{\partial x^\nu} \right) V^\nu . \quad (1.4.63)$$

The prototype vectors are in fact given by the displacements $\Delta x^\mu = x_B^\mu - x_A^\mu$ between two points A and B (note that the coordinates of each point are not a vector, as it will become clear later on). A vector field is a set of vectors defined in a region of space, and they transform like

$$V^{\mu'}(x') = M_{\nu}^{\mu'} V^\nu(x) , \quad (1.4.64)$$

where it is important to recall that the matrix $M_{\nu}^{\mu'}$ does not depend on the position.

Covectors: Quantities ω_μ that contracted with a vector yield a scalar:

$$\omega_\mu V^\mu = f . \quad (1.4.65)$$

It is easy to see that ω must then transform with M^{-1} ,

$$\omega_{\mu'} = M_{\mu'}^{\alpha} \omega_{\alpha} = \left(\frac{\partial x^{\alpha}}{\partial x^{\mu'}} \right) \omega_{\alpha} , \quad (1.4.66)$$

since then

$$\omega_{\mu'} V^{\mu'} = \left(\frac{\partial x^{\alpha}}{\partial x^{\mu'}} \right) \omega_{\alpha} \left(\frac{\partial x^{\mu'}}{\partial x^{\beta}} \right) V^{\beta} = \delta_{\beta}^{\alpha} \omega_{\alpha} V^{\beta} = \omega_{\alpha} V^{\alpha} . \quad (1.4.67)$$

Tensors: A general (n, m) tensor is a quantity that transforms like

$$T_{\nu'_1 \nu'_2 \dots \nu'_m}^{\mu'_1 \mu'_2 \dots \mu'_n} = \left(\frac{\partial x^{\mu'_1}}{\partial x^{\alpha_1}} \right) \left(\frac{\partial x^{\mu'_2}}{\partial x^{\alpha_2}} \right) \dots \left(\frac{\partial x^{\mu'_n}}{\partial x^{\alpha_n}} \right) \left(\frac{\partial x^{\beta_1}}{\partial x^{\nu'_1}} \right) \left(\frac{\partial x^{\beta_2}}{\partial x^{\nu'_2}} \right) \dots \left(\frac{\partial x^{\beta_m}}{\partial x^{\nu'_m}} \right) T_{\beta_1 \beta_2 \dots \beta_m}^{\alpha_1 \alpha_2 \dots \alpha_n} . \quad (1.4.68)$$

It is now easy to see that all operations defined in a general vector space can be applied to the present case. For example, multiplication of a (n, m) tensor by a scalar does not change its transformation properties, that is

$$\Phi T_{\beta_1 \beta_2 \dots \beta_m}^{\alpha_1 \alpha_2 \dots \alpha_n} = R_{\beta_1 \beta_2 \dots \beta_m}^{\alpha_1 \alpha_2 \dots \alpha_n} , \quad (1.4.69)$$

is still a (n, m) tensor, and tensors of same rank can be added and subtracted.

Further, by multiplying the components of a (n, m) tensor T by the components of a (p, q) tensor Q , one obtains a $(n + p, m + q)$ tensor R ,

$$T^{\alpha_1 \alpha_2 \dots \alpha_n}_{\beta_1 \beta_2 \dots \beta_m} Q^{\sigma_1 \sigma_2 \dots \sigma_p}_{\gamma_1 \gamma_2 \dots \gamma_q} = R^{\alpha_1 \alpha_2 \dots \alpha_n \gamma_1 \gamma_2 \dots \gamma_p}_{\beta_1 \beta_2 \dots \beta_m \sigma_1 \sigma_2 \dots \sigma_q} . \quad (1.4.70)$$

On the other hand, by *contracting* a rank (n, m) tensor T with a $(p < m, q < n)$ tensor Q , one obtains a $(n - q, m - p)$ tensor R . For example,

$$T^{\alpha_1 \alpha_2 \dots \alpha_n}_{\beta_1 \beta_2 \dots \beta_m} Q^{\beta_1 \beta_2 \dots \beta_p}_{\alpha_1 \alpha_2 \dots \alpha_q} = R^{\alpha_{q+1} \alpha_{q+2} \dots \alpha_n}_{\beta_{p+1} \beta_{p+2} \dots \beta_m} . \quad (1.4.71)$$

An example of covector is given by the (four-)gradient. From the chain rule, we in fact have

$$\frac{\partial}{\partial x^{\mu'}} = \left(\frac{\partial x^\alpha}{\partial x^{\mu'}} \right) \frac{\partial}{\partial x^\alpha} = M^\alpha_{\mu'} \frac{\partial}{\partial x^\alpha} . \quad (1.4.72)$$

For this reason, we shall often use the more compact notation

$$\frac{\partial}{\partial x^\mu} = \partial_\mu . \quad (1.4.73)$$

It then follows that, if $F^{\mu\nu}$ is a $(2, 0)$ tensor,

$$\partial_\alpha F^{\alpha\mu} = J^\mu \quad (1.4.74)$$

is a vector.

The above formalism holds for any linear transformation. To make contact with relativity, we require that $M = \Lambda$ be a Lorentz transformation. We shall see later on what this really means, but for now it is enough to know that, from

$$\eta_{\mu'\nu'} = \Lambda^\alpha_{\mu'} \Lambda^\beta_{\nu'} \eta_{\alpha\beta} \quad \Leftrightarrow \quad \Lambda^T \eta \Lambda = \eta \quad (1.4.75)$$

one finds that to each vector V^μ there can be associated a co-vector by contracting with the Minkowski $(0, 2)$ metric tensor,

$$V_\mu = \eta_{\mu\alpha} V^\alpha , \quad (1.4.76)$$

where

$$\eta_{\mu\nu} = \text{diag}(-1, 1, 1, 1) . \quad (1.4.77)$$

In fact,

$$V_{\mu'} = \eta_{\mu'\alpha'} V^{\alpha'} = \Lambda^\alpha_{\mu'} \Lambda^\beta_{\alpha'} \eta_{\alpha\beta} \Lambda^{\alpha'}_{\gamma} V^\gamma = \Lambda^\alpha_{\mu'} \eta_{\alpha\beta} \delta^\beta_{\gamma} V^\gamma = \Lambda^\alpha_{\mu'} (\eta_{\alpha\beta} V^\beta) = \Lambda^\alpha_{\mu'} V_\alpha . \quad (1.4.78)$$

Keep in mind that V^μ and V_μ are (mathematically and perhaps physically) *different* objects. The above metric tensor clearly equals its matrix inverse. We therefore define the inverse metric $\eta^{\mu\nu}$ by

$$\eta_{\mu\alpha} \eta^{\alpha\nu} = \delta^\nu_\mu , \quad (1.4.79)$$

and remark that we can likewise associate a vector to a covector according to

$$V^\mu = \eta^{\mu\alpha} V_\alpha . \quad (1.4.80)$$

As a simple mnemonic rule, to go from upper indices to lower indices (and vice versa) one just needs to change the sign of 0 components. For example,

$$V_\mu = (-V^0, \vec{V}) , \quad T_{\mu\nu} = \begin{bmatrix} T^{00} & -T^{0i} \\ -T^{i0} & T^{ij} \end{bmatrix} . \quad (1.4.81)$$

In general, contracting one index of a (n, m) tensor with the $(0, 2)$ metric produces a $(n-1, m+1)$ tensor as well as the (inverse) $(2, 0)$ metric will produce a $(n+1, m-1)$ tensor. Note the total number of indices does not change: $(n-1)+(m+1) = (n+1)+(m-1) = n+m$.

A simple example that shows how useful the covariant formalism can be is given by the following: consider a particle's four-velocity,

$$u^\mu = \frac{dx^\mu}{d\tau} . \quad (1.4.82)$$

Since $d\tau$ is a Lorentz scalar and dx^μ is a Lorentz vector, u^μ is also a vector, that is

$$u^{\mu'} = \Lambda^{\mu'}_\nu u^\nu . \quad (1.4.83)$$

As such, u^μ can be computed in any inertial reference frame and its components in any other inertial frame will be given by Lorentz transformations. In particular, in the (instantaneous) rest frame of the particle, since $dt = d\tau$,

$$u^\mu = (1, 0, 0, 0) , \quad (1.4.84)$$

which implies the scalar relation

$$u_\mu u^\mu = -1 , \quad (1.4.85)$$

the latter being just a re-statement of the mass-shell condition $p_\mu p^\mu = -m_0^2$ (having set $c = 1$). By differentiating the above with respect to $d\tau$, we obtain a vector relation,

$$0 = \frac{du_\mu}{d\tau} u^\mu + u_\mu \frac{du^\mu}{d\tau} = 2 u_\mu \frac{du^\mu}{d\tau} , \quad (1.4.86)$$

which implies that the four-acceleration is always orthogonal to the 4-velocity [unlike the usual acceleration defined in Eq. (1.5.40)]. Note though that, beside being a neat result, the above has not much physical sense since the four-velocity and four-acceleration are not the quantity we actually measure. Eq. (1.4.86) in the rest frame of the particle simply means the acceleration is purely spatial,

$$a^\mu = (0, a_x, a_y, a_z) . \quad (1.4.87)$$

Since the type of vector' is invariant, we deduce a^μ is space-like for all observers, like u^μ is time-like.

1.5 Relativistic dynamics

The very concept of Newtonian force between separate bodies implies an “action at a distance”, which is incompatible with the principles of Special Relativity. On the other hand, contact interactions, such as those involved in collisions, are perfectly acceptable, since they occur when two bodies touch at one point in space-time. In fact, the outcome of a collision is never determined by analysing the forces acting among the colliding objects: one instead neglects the size of these objects and the duration of the collision, for simplicity, and employs the conservation of linear momentum and energy. Which brings us to the question of how these very important quantities are modified in Special Relativity [5].

1.5.1 Relativistic momentum and mass

It should already be clear from the law of addition of velocities (1.3.21)-(1.3.24), that the Newtonian momentum $\vec{p} = m_0 \vec{u}$ of a point-like particle of (rest or proper) mass m_0 is not going to be a very useful quantity in Special Relativity, since it takes different forms in different inertial frames, that is $m_0 \vec{u} \not\rightarrow m_0 \vec{u}'$. In particular, we anticipate that, from the study of collision processes, the role played by \vec{p} will be taken over by the so-called *four-momentum*

$$P^\mu = \left(\frac{m_0 c}{\sqrt{1 - u^2/c^2}}, \frac{m_0 \vec{u}}{\sqrt{1 - u^2/c^2}} \right), \quad (1.5.1)$$

where the index $\mu = 0, 1, 2, 3$ for t, x, y and z -components, respectively.

Using the law of velocity composition (1.3.21)-(1.3.24), one can explicitly verify that this four-momentum indeed transforms like a space-time displacement Δx^μ under Lorentz transformations, that is

$$\left\{ \begin{array}{l} \vec{P} = \frac{m_0 \vec{u}}{\sqrt{1 - \frac{u^2}{c^2}}} \\ \frac{E}{c} = \frac{m_0 c}{\sqrt{1 - \frac{u^2}{c^2}}} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} P'_x = \frac{m_0 u'_x}{\sqrt{1 - \frac{u'^2}{c^2}}} = \frac{P_x - \left(\frac{v}{c}\right) \frac{E}{c}}{\sqrt{1 - \frac{v^2}{c^2}}} \\ P'_y = \frac{m_0 u'_y}{\sqrt{1 - \frac{u'^2}{c^2}}} = P_y \\ P'_z = \frac{m_0 u'_z}{\sqrt{1 - \frac{u'^2}{c^2}}} = P_z \\ \frac{E'}{c} = \frac{m_0 c}{\sqrt{1 - \frac{u'^2}{c^2}}} = \frac{\frac{E}{c} - \frac{v}{c} P_x}{\sqrt{1 - \frac{v^2}{c^2}}} \end{array} \right. \quad (1.5.2)$$

It is instead much easier to note that

$$U^\mu = \frac{\Delta x^\mu}{\Delta \tau} \rightarrow \frac{dx^\mu}{d\tau} \quad (1.5.3)$$

is a four-vector, since the proper time $\Delta\tau = \sqrt{1 - u^2/c^2} \Delta t$ is a scalar, as we have seen previously, and Δx^μ is the displacement four-vector.

Note however that the above four-velocity U^μ is *not* the quantity measured by any observer. The latter is given by the usual \vec{u} in the reference frame of the observer: this is the reason some texts prefer to introduce the “relativistic mass”

$$m(u) = \frac{m_0}{\sqrt{1 - u^2/c^2}} , \quad (1.5.4)$$

where m_0 is the proper mass of the particle, as measured by an observer at rest with the particle itself. This choice is obviously equivalent to introducing the “relativistic velocity” \vec{U} in the expression of the relativistic momentum,

$$\vec{P} = m \vec{u} = \frac{m_0 \vec{u}}{\sqrt{1 - \frac{u^2}{c^2}}} = m_0 \vec{U} . \quad (1.5.5)$$

In the particular case in which the particle is at rest in the frame S (that is, $\vec{u} = \vec{U} = 0$), we straightforwardly have that, in S ,

$$U^\mu = (1, 0, 0, 0) , \quad (1.5.6)$$

since $dt = d\tau$ and the observer does not move with respect to itself. Further, Eq. (1.5.1) clearly satisfies a correspondence principle, since its spatial components reduce to \vec{p} for $|u| \ll c$.

In the following, we shall show that it is P^μ which is conserved, besides being relativistically invariant (in form). We shall also see the meaning of P^0 , which has no counterpart in the Newtonian momentum \vec{p} .

1.5.2 Elastic collisions

Let us first consider an elastic collision, in which both energy and momentum are conserved (in a given inertial frame, at least), between two particles A and B with the same proper mass m_0 . We choose the frame S' so that initial total momentum is zero, that is $u'_{Bx} = -u'_{Ax}$ and $u'_{By} = -u'_{Ay}$. In S , on the other hand, we assume A does not move along the x -axis, so that $v = -u'_{Ax} = u'_{Bx}$, and $U_{Ax} = u_{Ax} = 0$ (see Fig. 1.20). In S' , momentum conservation yields

$$U'_{Ax} = u'_{Ax} = -u'_{Bx} = -U'_{Bx} \quad (1.5.7)$$

$$U'_{Ay} = -u'_{Ay} = u'_{By} = -U'_{By} , \quad (1.5.8)$$

so that both energy and momentum are indeed conserved (only the y -components flip sign and the speeds do not change after the collision).

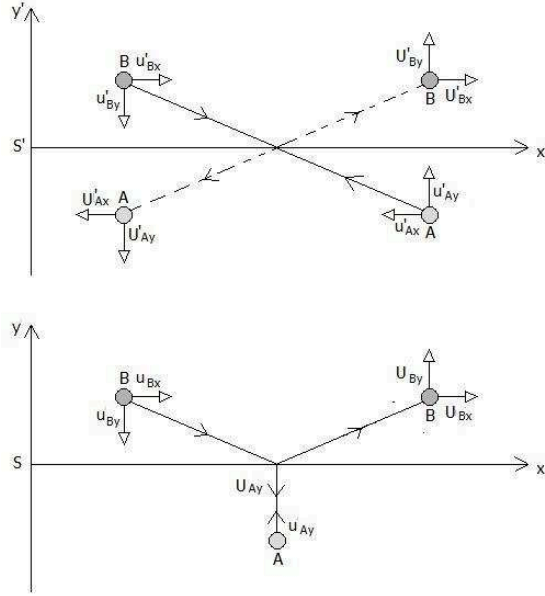


Figure 1.20: Example of elastic collision.

In Newtonian mechanics, velocities perpendicular to the direction of motion of S' with respect to S should not change. From the above Eq. (1.5.8), we would then readily find

$$U_{Ay} = -u_{Ay} = u_{By} = -U_{By} . \quad (1.5.9)$$

We next note that Lorentz transformations instead affect the components of the momentum along perpendicular directions, and yield

$$\begin{cases} u'_{By} = u_{By} \frac{\sqrt{1 - \beta^2}}{1 - \frac{v}{c^2} u_{Bx}} \\ u'_{Ay} = u_{Ay} \sqrt{1 - \beta^2} , \end{cases} \quad (1.5.10)$$

where we used $u_{Ax} = 0$. Eqs. (1.5.10) imply that Eq. (1.5.8) is not compatible with Eq. (1.5.9) for $v \neq 0$. This is evidence that the Newtonian definition of momentum, $\vec{p} = m_0 \vec{u}$, cannot represent a conserved quantity in all inertial frames, and must be modified. This can be achieved by modifying either its dependence on the speed or the mass, in such a way that it reproduces the usual expression for small velocities $|v| \ll c$. For example, one can introduce the relativistic mass of Eq. (1.5.4) or, simply, the relativistic momentum (1.5.1).

We conceived the above collision having in mind Newtonian concepts (for which it is clearly elastic). However, we should now show that the collision is indeed elastic by means of truly relativistic quantities. Let us first note that the new momentum satisfies the following relation

$$m_0^2 c^2 = m^2 c^2 - \vec{P} \cdot \vec{P} = m^2 c^2 - m^2 u^2 , \quad (1.5.11)$$

and, upon differentiating both sides [note that $d(m_0^2 c^2) = 0$], we obtain

$$c^2 dm = m u du + u^2 dm . \quad (1.5.12)$$

We can then derive an expression for the (change in) relativistic kinetic energy, by requiring it equals the work done on the particle, that is the rate of change in momentum integrated along the path of the particle. For example, given the trajectory $x = x(t)$ in one spatial dimension, one has

$$\begin{aligned} \Delta K &= \int_{x_A}^{x_B} \frac{d}{dt}(m u) dx \\ &= \int_{x_A}^{x_B} \left(m \frac{du}{dt} dx + \frac{dm}{dt} u dx \right) \\ &= \int_{t_A}^{t_B} \left(m \frac{du}{dt} u dt + \frac{dm}{dt} u^2 dt \right) \\ &= \int_{u_A}^{u_B} \left(m u du + u^2 \frac{dm}{du} du \right) \\ &= c^2 \int_{m_A}^{m_B} dm = c^2 (m_B - m_A) , \end{aligned} \quad (1.5.13)$$

where we repeatedly changed integration variable and finally used Eq. (1.5.12). For a particle initially at rest, we set $m_A = m_0$, and rename $m_B = m$. We then find

$$K = c^2 (m - m_0) = m_0 c^2 \left[\frac{1}{\sqrt{1 - \frac{u^2}{c^2}}} - 1 \right] = m_0 c^2 \left(1 - \frac{u^2}{2c^2} \dots - 1 \right) . \quad (1.5.14)$$

In the non-relativistic limit $u \ll c$,

$$K \simeq \frac{1}{2} m_0 u^2 \quad (1.5.15)$$

and, in general,

$$K = m_0 c^2 \left[\frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} - 1 \right] , \quad (1.5.16)$$

which is a first hint the total energy of the particle is given by

$$E = K + m_0 c^2 = m c^2 = c P^0 . \quad (1.5.17)$$

Note that for $u \rightarrow c$, the kinetic energy K diverges. In fact, it is sensible that one needs an infinite amount of energy to accelerate a particle to the speed of light.

The four-momentum can now be rewritten as

$$P^\mu = \left(\frac{E}{c}, \vec{P} \right) , \quad (1.5.18)$$

and, from Eq. (1.5.11), it is easy to see that its components satisfy

$$E^2 - c^2 \vec{P} \cdot \vec{P} = m_0^2 c^4 , \quad (1.5.19)$$

which is known as the *mass-shell relation*.

We can now consider the system in Fig. 1.20 again, and note that, in S' , the total Newtonian momentum \vec{p} vanishes, by construction (some components of the velocities change sign, but speeds are conserved), both before and after the collision. This immediately implies that the spatial components of the total four-momentum in S' also vanish and are conserved,

$$\vec{P}'_{\text{in}} = 0 = \vec{P}'_{\text{fin}} . \quad (1.5.20)$$

Moreover, it is easy to see that

$$c P'^0 = (m_0 c^2 + K'_A) + (m_0 c^2 + K'_B) = 2 m_0 c^2 + (K'_A + K'_B) = 2 (m_0 c^2 + K'_A) , \quad (1.5.21)$$

is also conserved, meaning its initial (before collision) and final (after collision) values are the same, since speeds are again conserved by construction in S' ,

$$P'^0_{\text{in}} = P'^0_{\text{fin}} , \quad (1.5.22)$$

from which

$$K'_{\text{in}} = 2 K'_A = K'_{\text{fin}} . \quad (1.5.23)$$

We can thus conclude the collision is indeed elastic in the frame S' , and we have full relativistic momentum conservation,

$$P'^\mu_{\text{in}} = P'^\mu_{\text{fin}} . \quad (1.5.24)$$

Since the relativistic four-momentum transforms according to Eq. (1.5.2), Eq. (1.5.24) immediately implies that in S we must have

$$P^\mu_{\text{in}} = P^\mu_{\text{fin}} , \quad (1.5.25)$$

with ⁷

$$c P^0 = 2 m_0 c^2 + K_A + K_B , \quad (1.5.26)$$

so that

$$K_{\text{in}} = K_A + K_B = K_{\text{fin}} , \quad (1.5.27)$$

which proves that the collision is elastic also in the frame S . We therefore conclude the collision is *elastic in all inertial frames*.

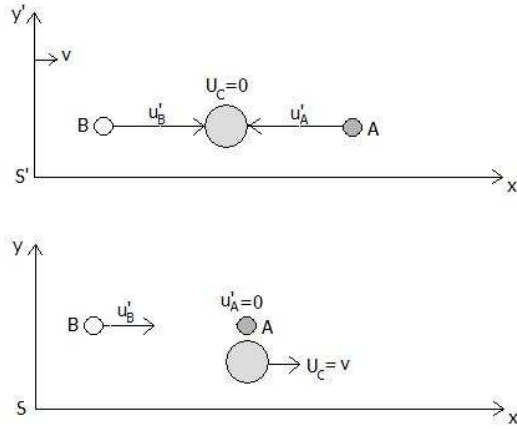


Figure 1.21: Example of inelastic collision.

1.5.3 Inelastic collisions

We have seen that the new definition of relativistic momentum hints at a new definition of total energy for a point-like particle, namely

$$E = c P^0 = m_0 c^2 + K . \quad (1.5.28)$$

This result would be purely academic if there were no ways to transform the new “proper energy” $m_0 c^2$ into a different form, for example kinetic energy (or the other way around).

For this purpose, we next consider a totally inelastic collision between two particles A and B with the same proper mass m_0 and velocities $u'_B = -u'_A = u'$ in S' . In the frame S , the particle A is at rest, so that $v = u'$. After the collision, the two particles remain attached and form a single particle C (see Fig. 1.21). From the Lorentz transformations (1.3.12), or velocity composition (1.3.21), we again obtain

$$u_B = \frac{u' + v}{1 + \frac{u'v}{c^2}} = \frac{2v}{1 + \frac{v^2}{c^2}} , \quad U_C = v , \quad (1.5.29)$$

and, consequently, the relativistic mass of B in S is given by

$$m_B = \frac{m_0}{\sqrt{1 - \frac{u_B^2}{c^2}}} = m_0 \frac{1 + \frac{v^2}{c^2}}{1 - \frac{v^2}{c^2}} . \quad (1.5.30)$$

Conservation of momentum along the x -axis in S implies

$$m_B u_B = \frac{m_0 u_B}{\sqrt{1 - \frac{u_B^2}{c^2}}} = \frac{M_0 U_C}{\sqrt{1 - \frac{U_C^2}{c^2}}} = M_C U_C , \quad (1.5.31)$$

⁷We do not need to compute K explicitly, as long as the proper mass m_0 does not depend on the reference frame.

where M_0 is the proper mass of C . Upon inserting the velocities from Eq. (1.5.29), we obtain

$$\frac{m_0 u_B}{\sqrt{1 - \frac{u_B^2}{c^2}}} = m_0 \frac{1 + \frac{v^2}{c^2}}{1 - \frac{v^2}{c^2}} \frac{2v}{1 + \frac{v^2}{c^2}} = \frac{2m_0 v}{1 - \frac{v^2}{c^2}} = \frac{M_0 v}{\sqrt{1 - \frac{v^2}{c^2}}} , \quad (1.5.32)$$

or

$$M_0 = \frac{2m_0}{\sqrt{1 - \frac{v^2}{c^2}}} , \quad (1.5.33)$$

which is larger than the sum of the proper masses of A and B . In fact, all kinetic energy converted into mass in S' ,

$$M_0 - 2m_0 = 2m_0 \left(\frac{1}{\sqrt{1 - \frac{v^2}{c^2}}} - 1 \right) = \frac{K'_B + K'_A}{c^2} = 2 \frac{K'_A}{c^2} , \quad (1.5.34)$$

and *energy is conserved* in both frames. In particular, in S' , initial and final energies are given by

$$\begin{cases} E'_{\text{in}} = 2(m_0 c^2 + K'_A) = \frac{2m_0 c^2}{\sqrt{1 - \frac{v^2}{c^2}}} \\ E'_{\text{fin}} = M_0 c^2 = \frac{2m_0 c^2}{\sqrt{1 - \frac{v^2}{c^2}}} , \end{cases} \quad (1.5.35)$$

and, from Eq. (1.5.29), in S we have

$$\begin{cases} E_{\text{in}} = m_0 c^2 + (m_0 c^2 + K_B) = 2m_0 c^2 + m_0 c^2 \left[\frac{1}{\sqrt{1 - \frac{u_B^2}{c^2}}} - 1 \right] = \frac{2m_0 c^2}{1 - \frac{v^2}{c^2}} \\ E_{\text{fin}} = M_0 c^2 + K_C = \frac{2m_0 c^2}{1 - \frac{v^2}{c^2}} . \end{cases} \quad (1.5.36)$$

Note also that $E > E'$: the energy of the system is the smallest in the frame at rest with the final particle C ⁸.

1.5.4 Equivalence of mass and energy

In the previous experiment we saw all kinetic energy turned into mass. This yields the physical meaning of the famous

$$E = m c^2 = m_0 c^2 , \quad (1.5.37)$$

⁸Note the duality with the proper time being the longest for an observer at rest.

which holds in the rest frame of a massive particle. Let us emphasise that energy is *always* conserved in Special Relativity, even in processes which are inelastic from the point of view of Newtonian mechanics, simply as a consequence of the linearity of Lorentz transformations in space-time. For example, it is clear that the total spatial momentum in the inelastic collision of Section 1.5.3 can be conserved in S only if the energy is conserved in the centre-of-mass frame S' (where the spatial momentum vanishes by definition).

Then, the actual possibility of converting mass into energy (or vice versa) is subjected to restrictions. For example, if the particle C cannot reduce its proper mass (like an elementary particle), then the inverse process of the inelastic collision previously discussed may not occur.

1.5.5 Relativistic force law

In the previous derivation (1.5.13) of the change in kinetic energy ΔK , we implicitly assumed the force \vec{F} acting on a particle is given by the time-derivative of the particle's relativistic momentum. If one insists on this interpretation, it immediately follows that \vec{F} is not parallel to the particle's acceleration,

$$\vec{F} = \frac{d}{dt}(m\vec{u}) = \left(\frac{d\vec{u}}{dt}\right)m + \vec{u}\left(\frac{dm}{dt}\right) = \left(\frac{dm}{dt}\right)\vec{u} + m\vec{a}, \quad (1.5.38)$$

since the first term above is parallel to the particle's velocity, and we recall that m is the (velocity dependent) relativistic mass. Moreover, from Eq. (1.5.17),

$$\frac{dm}{dt} = \frac{1}{c^2} \frac{dE}{dt} = \frac{1}{c^2} \frac{dK}{dt} = \frac{1}{c^2} \frac{d}{dt} (\vec{F} \cdot d\vec{l}) = \frac{1}{c^2} \vec{F} \cdot \vec{u}, \quad (1.5.39)$$

where we assumed \vec{F} is constant in the last step. Replacing the above into Eq. (1.5.38) yields

$$m\vec{a} = \vec{F} - \frac{\vec{F} \cdot \vec{u}}{c^2} \vec{u}, \quad (1.5.40)$$

which shows that the acceleration contains a component parallel to the force and a second component parallel to the particle's velocity. This second term acts against the force \vec{F} to restrain u from exceeding the speed of light.

From the transformation laws for the time $t \rightarrow t'$ and the product $m\vec{u} \rightarrow m'\vec{u}'$, one could easily derive the transformation laws for the force $\vec{F} \rightarrow \vec{F}'$. We shall however not need to display them here, since the force is no more a quantity that transforms nicely under the new coordinate transformations of Special Relativity.

1.6 Electromagnetism

We already know Special Relativity was designed to comply with Maxwell's laws of electromagnetism. We shall now see that this is indeed the case.

Maxwell's equations in the usual three-dimensional formalism ⁹

$$\begin{cases} \epsilon_0 \vec{\nabla} \cdot \vec{E} = \rho \\ \vec{\nabla} \times \vec{B} = \mu_0 \vec{J} + \mu_0 \epsilon_0 \frac{\partial \vec{E}}{\partial t} \end{cases} \quad (1.6.1)$$

$$\begin{cases} \vec{\nabla} \cdot \vec{B} = 0 \\ \vec{\nabla} \times \vec{E} = -\frac{\partial \vec{B}}{\partial t} , \end{cases} \quad (1.6.2)$$

do not make it particularly clear that they transform in a way that keeps them of the same form in all inertial reference frames. It is however easy to find how the source terms change under Lorentz transformations.

1.6.1 Electric charge and current

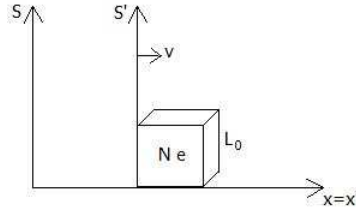


Figure 1.22: Charged cubic cell in motion.

Let us consider a cube of electrically charged matter at rest in S' , moving with velocity \vec{v} with respect to S (see Fig. 1.22). In particular, suppose the cube contains N electrons (with charge e each). In S' , we have the following (proper) densities of charge and current,

$$\begin{cases} \rho_0 = \frac{Ne}{L_0^3} \\ \vec{j}_0 = 0 . \end{cases} \quad (1.6.3)$$

According to Maxwell's equations (1.6.2), we therefore expect S' will see an electric field, but no magnetic field. In S , on the other hand, the cube is contracted along the x -axis and one has

$$\begin{cases} \rho = \frac{Ne}{L_0 \sqrt{1 - \beta^2} L_0^2} = \frac{\rho_0}{\sqrt{1 - \beta^2}} \\ \vec{j} = \frac{\rho_0 \vec{v}}{\sqrt{1 - \beta^2}} , \end{cases} \quad (1.6.4)$$

⁹We shall often use $c^{-2} = \mu_0 \epsilon_0 = 1$.

which implies that S should see both an electric and a magnetic field produced by the same charges.

The four-vector $J^\mu = (c\rho, \vec{j})$ is called four-current and is mathematically similar to the four-momentum P^μ (just replace m_0 with ρ_0 in the latter). In fact, it shares similar properties, such as the mass-shell relation which we recall here

$$c^2 t^2 - x^2 = c^2 \tau^2 \quad \Leftrightarrow \quad m^2 c^2 - p^2 = m_0^2 c^2 \quad \Leftrightarrow \quad c^2 \rho^2 - j^2 = c^2 \rho_0^2 . \quad (1.6.5)$$

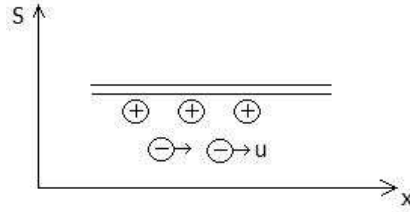


Figure 1.23: Current in a wire.

Let us next consider a current flowing through a thin wire at rest in S . The current is composed of electrons moving with velocity \vec{u} in S (see Fig. 1.23). Let n denote the number of electrons per unit volume (in S). Since the wire is electrically neutral, we must have

$$\begin{cases} \rho^- = -n e \\ \rho^+ = n e \end{cases} \quad \Rightarrow \quad \rho = \rho^+ + \rho^- = 0 , \quad (1.6.6)$$

and

$$\begin{cases} \vec{j}^- = -n e \vec{u} = \rho^- \vec{u} \\ \vec{j}^+ = 0 \end{cases} \quad \Rightarrow \quad \vec{j} = \vec{j}^+ + \vec{j}^- = -\rho^- \vec{u} , \quad (1.6.7)$$

where we can assume $\vec{u} = u_x$. In a frame S' moving with velocity \vec{v} parallel to \vec{u} with respect to S , we instead have

$$\begin{cases} \rho'^- = \frac{\rho^- - v \frac{j^-}{c^2}}{\sqrt{1 - \beta^2}} = \frac{\rho^- (1 - \frac{uv}{c^2})}{\sqrt{1 - \beta^2}} \\ \rho'^+ = \frac{\rho^+ - v \frac{j^+}{c^2}}{\sqrt{1 - \beta^2}} = \frac{\rho^+}{\sqrt{1 - \beta^2}} , \end{cases} \quad (1.6.8)$$

and

$$\begin{cases} j'^- = \frac{j^- - v \rho^-}{\sqrt{1 - \beta^2}} = \frac{\rho^- (\vec{u} - \vec{v})}{\sqrt{1 - \beta^2}} \\ j'^+ = \frac{j^+ - v \rho^+}{\sqrt{1 - \beta^2}} = -\frac{v \rho^+}{\sqrt{1 - \beta^2}} . \end{cases} \quad (1.6.9)$$

If, further, we have $u = v$, then $j'^- = 0$, as expected: in S the current is due to the motion of negative charges, whereas positive charges move in S' .

1.6.2 Transformations for \vec{E} and \vec{B}

We already mentioned that the spatial force, defined by

$$\vec{F} = \frac{d}{dt} \left(\frac{m_0 \vec{u}}{\sqrt{1 - \frac{u^2}{c^2}}} \right) , \quad (1.6.10)$$

does not transform nicely ¹⁰ into a new vector \vec{F}' because of the transformation law of $t \rightarrow t' \neq t$. There is one case which can be dealt with easily, namely the one in which the body subject to the force is (momentarily) at rest in S' , so that $dt' = d\tau$. The force acting upon it, as seen in a system S moving with velocity \vec{v} with respect to S' (and the body itself), is then given by

$$\begin{cases} F_x = F'_x \\ F_y = F'_y \sqrt{1 - \frac{v^2}{c^2}} \\ F_z = F'_z \sqrt{1 - \frac{v^2}{c^2}} , \end{cases} \quad (1.6.11)$$

which shows that there is no change in the longitudinal component (parallel to \vec{v}), but only in the orthogonal components.

We can now apply the above result to the case of an electromagnetic force acting on a test charge q , which is given by Lorentz law,

$$\vec{F} = q (\vec{E} + \vec{u} \times \vec{B}) , \quad (1.6.12)$$

and should transform according to the rules we mentioned previously (albeit never displayed in full). In particular, suppose the test charge q is at rest in S' but moves in S , with the velocity

$$u_x = v , \quad u_y = u_z = 0 . \quad (1.6.13)$$

The Lorentz force acting on it as measured in S' is simply

$$\vec{F}' = q \vec{E}' , \quad (1.6.14)$$

since $\vec{u}' = 0$. On the other hand, in S we have

$$\left(\vec{u} \times \vec{B} \right)_x = u_y B_z - u_z B_y = 0 \quad \Rightarrow \quad F_x = q E_x , \quad (1.6.15)$$

¹⁰Meaning that it is not part of a four-vector.

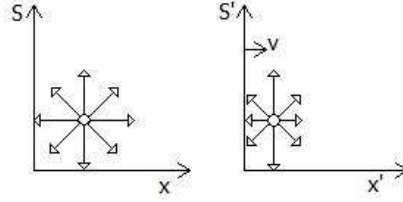


Figure 1.24: Electric field of a point-like charge.

and, from Eq. (1.6.11), one finds

$$F'_x = F_x \quad \Rightarrow \quad E'_x = E_x . \quad (1.6.16)$$

With similar arguments, for the y and z components, one can finally obtain the full set of transformation rules of the components of the electric field, namely

$$\begin{cases} E'_x = E_x \\ E'_y = \gamma(E_y - v B_z) \\ E'_z = \gamma(E_z + v B_y) , \end{cases} \quad (1.6.17)$$

where $\gamma = (1 - \beta^2)^{-1/2}$.

This result can be generalised for a charged body moving along a generic trajectory as

$$\begin{cases} \vec{E}'_{\parallel} = \vec{E}_{\parallel} \\ \vec{E}'_{\perp} = \gamma \left[\vec{E}_{\perp} + \left(\vec{v} \times \vec{B} \right)_{\perp} \right] \end{cases} \quad (1.6.18)$$

and, analogously, for the magnetic field

$$\begin{cases} \vec{B}'_{\parallel} = \vec{B}_{\parallel} \\ \vec{B}'_{\perp} = \gamma \left[\vec{B}_{\perp} + \left(\frac{\vec{v}}{c^2} \times \vec{E} \right)_{\perp} \right] , \end{cases} \quad (1.6.19)$$

where \parallel means the component parallel to the relative velocity \vec{v} and \perp those perpendicular to \vec{v} .

Instead of deriving the above transformations, we shall just show that the transformation of the electric field \vec{E} is indeed in agreement with Maxwell's equations. For this purpose, let us consider again a point-like particle with charge q . It will produce a spherically symmetric radial electric field in the frame in which it is at rest. However, in a frame moving with constant velocity, the field line will shrink along the direction of relative motion (see Fig. 1.24). This result can be easily derived from Gauss law, namely

$$\vec{\nabla} \cdot \vec{E} \propto \rho . \quad (1.6.20)$$

Upon integrating the above on a volume containing q , we obtain

$$q \propto \int_V \vec{\nabla} \cdot \vec{E} = \int_{\partial V} \vec{E} \cdot d\vec{s}, \quad (1.6.21)$$

which does not depend on the reference frame and the flux of the electric field is therefore invariant. Since we know the area parallel to the direction of motion shrinks, whereas that orthogonal is unaffected, we can immediately conclude that E_{\parallel} does not change, whereas E_{\perp} must increase to compensate for the reduced area.

It is clear from the expressions (1.6.18) and (1.6.19) that the fields \vec{E} and \vec{B} do not simply transform as Lorentz vectors, unlike the four-current J^{μ} they are sourced by. Their geometrical nature is in fact much more clear if we employ the tensor formalism, which will also allow us to write Maxwell's equation in a more compact (geometrical) form.

1.6.3 Maxwell equations redux

Let us now consider Maxwell's equations. The third and fourth equations (1.6.2) do not contain sources and allow us to introduce scalar and vector potentials for \vec{E} and \vec{B} respectively,

$$\begin{aligned} \vec{E} &= -\vec{\nabla}\phi - \frac{\partial \vec{A}}{\partial t} \\ \vec{B} &= \vec{\nabla} \times \vec{A}. \end{aligned} \quad (1.6.22)$$

If we define the four-vector potential

$$A^{\mu} = (\phi, \vec{A}), \quad (1.6.23)$$

Eq. (1.6.22) is then equivalent to

$$F_{\mu\nu} = \partial_{\mu} A_{\nu} - \partial_{\nu} A_{\mu}. \quad (1.6.24)$$

For example, if we write $\partial_0 \equiv \partial_t$ and so on, we have

$$\begin{aligned} F_{01} &= \partial_t A_x - \partial_x A_t = \partial_t A_x + \partial_x \phi = -E_x = -E^x \\ F_{12} &= \partial_x A_y - \partial_y A_x = B_z = B^z. \end{aligned} \quad (1.6.25)$$

We now see that the components of the electric and magnetic field form the Maxwell (or field-strength) tensor given by

$$F_{\mu\nu} = -F_{\nu\mu} = \begin{bmatrix} 0 & -E_x & -E_y & -E_z \\ E_x & 0 & B_z & -B_y \\ E_y & -B_z & 0 & B_x \\ E_z & B_y & -B_x & 0 \end{bmatrix}. \quad (1.6.26)$$

Finally, Eqs. (1.6.1), which contain the sources, takes the simple form

$$\partial_\alpha F^{\alpha\mu} = -J^\mu , \quad (1.6.27)$$

and encodes the true dynamics of the electromagnetic fields. Note that this differential equation is first order in the physical fields \vec{E} and \vec{B} , but second order in the potential A^μ . Also, since \vec{E} and \vec{B} are uniquely determined by the four components of A^μ , the number of degrees of freedom of the electromagnetic field is (at most) four. We will meet the same structure when we study the gravitational interaction later on in the course.

Wave equation and gauge freedom

Putting together Eq. (1.6.24) and (1.6.27), for $J^\mu = 0$, we obtain the wave equation for the propagation of light signals

$$\partial_\alpha \partial^\alpha A^\mu - \partial_\alpha \partial^\mu A^\alpha = \partial_\alpha \partial^\alpha A^\mu = -(\partial_t^2 - c^2 \nabla^2) A^\mu = 0 , \quad (1.6.28)$$

where we used the “gauge condition” $\partial_\alpha A^\alpha = 0$ (known as the *Lorenz gauge*).

In fact, we recall the Maxwell tensor does not determine the vector potential uniquely, namely for any *gauge transformation* of the form

$$A_\mu \rightarrow \bar{A}^\mu = A_\mu + \partial_\mu \Lambda \quad \Rightarrow \quad F_{\mu\nu} \rightarrow \bar{F}_{\mu\nu} = F_{\mu\nu} , \quad (1.6.29)$$

which follows from the skew-symmetry of $F_{\mu\nu}$ and the fact that partial derivatives commute when applied to a scalar Λ . This means that the “independent” components of A^μ (or the number of physical degrees of freedom of the electromagnetic field) are not four, but no more than three.

In particular, the Lorenz gauge is determined by a scalar Λ which satisfies ¹¹

$$\partial_\mu \bar{A}^\mu = 0 \quad \Rightarrow \quad \partial_\alpha \partial^\alpha \Lambda = -\partial_\alpha A^\alpha . \quad (1.6.30)$$

However, this gauge choice does not fix A^μ completely, since one can always further change the potential as

$$\bar{A}_\mu \rightarrow \tilde{A}^\mu = \bar{A}_\mu + \partial_\mu \Lambda \quad \text{with} \quad \partial^\mu \partial_\mu \Lambda = 0 , \quad (1.6.31)$$

and the new \tilde{A}^μ will still satisfy $\partial_\mu \tilde{A}^\mu = 0$.

It can in fact be shown that A^μ indeed contains only two independent components (degrees of freedom): the two independent polarisations of light. For this purpose, let us choose a plane-wave moving along the x^3 direction,

$$A^\mu = \epsilon^\mu e^{i k_\alpha x^\alpha} , \quad (1.6.32)$$

¹¹Second order partial differential equations always admit a local solution, which proves Λ exists for all A^μ and (at least piecewise) all space-time points.

where the wave-number $k^\mu = (k, 0, 0, k)$ is such that $k_\mu k^\mu = 0$, and ϵ^μ is the polarisation vector. The Lorenz condition then reads

$$0 = \partial_\mu A^\mu = i k (\epsilon^z - \epsilon^t) e^{i k_\mu x^\mu} , \quad (1.6.33)$$

which implies $\epsilon^z = \epsilon^t$. In order to completely fix the gauge, we can further choose the scalar

$$\Lambda = i \frac{e^z}{k} e^{i k_\mu x^\mu} , \quad (1.6.34)$$

which yields

$$\begin{aligned} A_\mu + \partial_\mu \Lambda &= [(-\epsilon^z, \epsilon^x, \epsilon^y, \epsilon^z) + (\epsilon^z, 0, 0, -\epsilon^z)] e^{i k_\mu x^\mu} \\ &= (0, \epsilon^x, \epsilon^y, 0) e^{i k_\mu x^\mu} . \end{aligned} \quad (1.6.35)$$

The components ϵ^x and ϵ^y therefore represent the two linear polarisations of the electric and magnetic fields.

Scalars and charge conservation

There are a few things we can learn from the covariant formalism: for example, there is a scalar quantity

$$F_{\mu\nu} F^{\mu\nu} = 2 (B^2 - E^2) , \quad (1.6.36)$$

and a *pseudo*-scalar¹² quantity

$$\epsilon_{\alpha\beta\gamma\delta} F^{\alpha\beta} F^{\gamma\delta} = 8 \vec{B} \cdot \vec{E} , \quad (1.6.37)$$

where $\epsilon_{\alpha\beta\gamma\delta}$ is the totally antisymmetric Levi-Civita (pseudo-)tensor, with $\epsilon_{0123} = 1$.

A more important result is obtained by taking the derivative of Eq. (1.6.27),

$$0 = \partial_\beta \partial_\alpha F^{\alpha\beta} = -\partial_\beta J^\beta , \quad (1.6.38)$$

which follows from $F_{\mu\nu} = -F_{\nu\mu}$ and implies *charge conservation*. In fact, in the rest frame of the charge, this becomes

$$J^\mu = (\rho, 0, 0, 0) \quad \Rightarrow \quad \partial_\mu J^\mu = \partial_t \rho = 0 . \quad (1.6.39)$$

Whenever we have a current J^μ with vanishing four-divergence $\partial_\mu J^\mu = 0$, we have a conservation law. For example, if $J^\mu = P^\mu$, four-momentum conservation for a point-like particle, $\partial_\mu P^\mu = 0$, in the rest frame of the particle implies

$$P^\mu = (m_0, 0, 0, 0) \quad \Rightarrow \quad \partial_\mu P^\mu = \partial_t m_0 = 0 , \quad (1.6.40)$$

which shows that four-momentum conservation is the relativistic consequence of proper mass conservation, exactly like electric current conservation follows from the invariance of the electric charge.

¹²Pseudo-scalars change sign under spatial reflections, $x \rightarrow -x$, unlike true scalars, which do not.

1.6.4 Nature and relativistic fields

Let us conclude this chapter about Special Relativity with a few important observations. One reason that led us to Special Relativity was precisely the fact that electromagnetism is not Galilean invariant, which consequently led us to discard interactions at a distance and conservative forces derived from a potential. Now, we have just seen that the electromagnetic field $F_{\mu\nu}$ can be derived from a four-potential, and the question comes immediately to mind how this can be consistent. The fact is that the electromagnetic field does not entail instantaneous interactions at a distance: changes in the state of sourcing charges generate perturbations in the field which travel at the speed of light in vacuum [see Eq. (1.6.28)] before affecting test charges.

In the modern view of the physical world, everything is indeed represented by fields, even matter. Questioning the nature of light (recall the aether theories, and the idea that light is made of waves in a medium) is therefore the same as wondering what matter is made of. In a sense, they are both “just real”.

A more concrete question however comes to mind. If real interactions may (and in fact can) be represented by mediating fields, what is the gravitational field? If it can be derived by a potential like electromagnetism, what is the gravitational potential? And what are the gravitational analogues of electromagnetic waves (light)? Before addressing these questions, we will need to take a rather long detour into mathematics and geometry.

Chapter 2

Differentiable manifolds and tensors

Preamble

The linear transformations we have seen so far are *global* in that they affect all space-time points at the same time. However, physical measurements are *local* and one may therefore want to be able to perform local coordinate transformations. Such changes would be represented by (in general) non-linear transformations

$$x^\mu \rightarrow y^\mu = y^\mu(x^\mu) ,$$

with the only restriction that the above functions must be invertible. Eventually, Einstein's General Relativity requires physical laws can be equivalently expressed in any general coordinate frame, so that we need a formalism to handle the case above.

One nice catch is that, for small (but in which sense?) variations $|\delta x^\mu| = |y^\mu - x^\mu|$, we can Taylor expand the coordinate transformation around a given point like

$$y^\mu(x^\mu) = x^\mu + \left. \frac{\partial y^\mu}{\partial x^\alpha} \right|_{y^\mu=x^\mu} \delta x^\alpha + \dots .$$

Since the δx^μ can be viewed as (local) vectors under the linear transformation defined by the matrix

$$M^\mu_\alpha = \left. \frac{\partial y^\mu}{\partial x^\alpha} \right|_{y^\mu=x^\mu}$$

which does not depend on y^μ , there is hope that we can partly recover our previous construction at least in a local sense.

This is what differential geometry is all about: apply all of the mathematical machinery of \mathbb{R}^N and $GL(N)$ to more general geometric spaces. In the process of introducing it, we will see many things change and one loses some and gains some [8].

2.1 Differentiable manifolds

In a nutshell, a (differential) manifold is a topological space which locally looks like (a portion of) the n -dimensional Euclidean space \mathbb{R}^n .

Brief review of \mathbb{R}^n

In the following, we shall assume most of the properties of the set \mathbb{R}^n of real n -tuples x^i , with $i = 1, \dots, n$. In particular, \mathbb{R}^n is a *vector space*, with vectors defined by the displacements $v^i = x^i - y^i$, upon which one can act with global linear transformations belonging to $GL(n)$. As a vector space, \mathbb{R}^n also admits the Cartesian scalar product

$$v^i w_i = \sum_{i=1}^n v^i w^i, \quad (2.1.1)$$

and is thus a (finite dimensional) *Hilbert space*. This scalar product induces a Euclidean norm

$$||x - y||^2 = \sum_{i=1}^n (x^i - y^i)^2. \quad (2.1.2)$$

which allows one to define a *ball* as the open set

$$||x - y|| < R. \quad (2.1.3)$$

These open sets yield \mathbb{R}^n the properties of a separable *topological space*. Finally, we shall assume knowledge of n -dimensional (real) calculus.

A general differential geometry is roughly defined by repeating the above steps backwards, that is we shall start from the topology and work all the way up to the differential calculus.

2.1.1 Manifolds and coordinates

We start by recalling that a topological space is a set of elements (points) in which the notion of “contiguity” is defined: two elements of the set are contiguous if they both belong to the same open subset (usually referred to as a “neighbourhood” of those elements). More precisely, given a set \mathcal{M} of “points”¹, the topological space $(\mathcal{M}, \{A_i\})$ is defined by a family of so-called open sets $\{A_i\}$, such that the empty set \emptyset and \mathcal{M} itself belong to $\{A_i\}$, as do an arbitrary union of (a finite or infinite number of) open sets $\cup_i A_i$, and the intersection of a finite number of open sets $\cap_i A_i$ ².

In particular, we shall be concerned with separable (or Hausdorff) topological spaces, in which, for any two arbitrary points P and Q , there always exist disjoint open sets $U \ni P$ and $V \ni Q$ (neighborhoods of P and Q , respectively), with $U \cap V = \emptyset$.

¹Sets and points are primitive concepts.

²Alternatively, one can define the topological space $(\mathcal{M}, \{B_i\})$ in terms of a family of closed sets, the latter being the complementary to open sets, $B_i = \mathcal{M} \setminus A_i$.

A *map* is in general an application from an open set $D \subseteq \mathcal{M}$ (the domain) to (a subset of) \mathbb{R}^n , that is $\phi : D \rightarrow \mathbb{R}^n$. Since \mathcal{M} is a topological space, the notion of continuity can also be defined and means that a map ϕ is continuous if it maps any open set $A \subseteq D$ in its domain into an open (sub)set of \mathbb{R}^n .

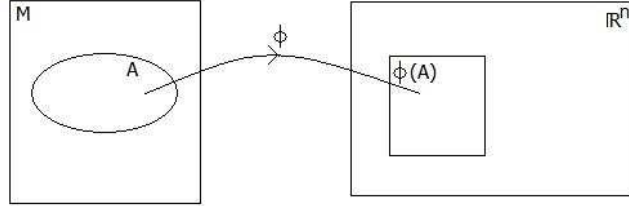


Figure 2.1: Coordinates of a point.

A *chart* is a pair (A, ϕ) , where $A \subseteq \mathcal{M}$ and ϕ is an invertible continuous map, $\phi : A \rightarrow \mathbb{R}^n$, which we often denote as $\phi(P) = x^i(P)$ (or, more concisely, $\phi = x^i$). In other words, the map of a chart is a set of n real coordinates for the open set $A \subseteq \mathcal{M}$.

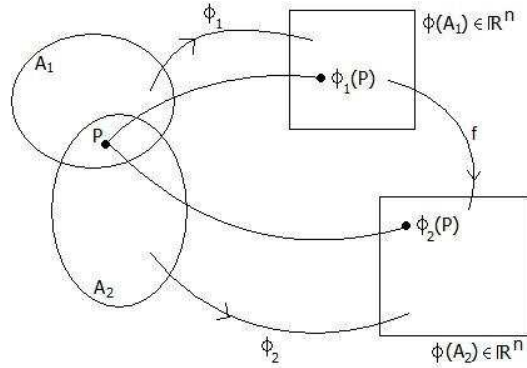


Figure 2.2: Consistency of maps.

A collection \mathcal{A} of charts is called an *atlas* if it covers the entire set \mathcal{M} , that is $\mathcal{A} = \{(A_i, \phi_i) : \cup_i A_i \supseteq \mathcal{M}\}$. Let then $P \in A_1 \cap A_2$, where A_1 and A_2 are the open sets of two charts ϕ_1 and ϕ_2 . There must then be a consistency relation between $\phi_1(P)$ and $\phi_2(P)$, in the form of an invertible application $f : \phi(A_1) \subseteq \mathbb{R}^n \rightarrow \phi(A_2) \subseteq \mathbb{R}^n$ such that

$$f(\phi_1(P)) = \phi_2(P) \quad \text{or} \quad (\phi_2^{-1} \circ f \circ \phi_1) = \mathbb{I} , \quad (2.1.4)$$

where \mathbb{I} is the identity in \mathcal{M} . Equivalently,

$$f^{-1}(\phi_2(P)) = \phi_1(P) \quad \text{or} \quad (\phi_1^{-1} \circ f^{-1} \circ \phi_2) = \mathbb{I} . \quad (2.1.5)$$

In layman terms, the application f is just a coordinate transformation in \mathbb{R}^n (see Fig. 2.2), and it immediately follows from the above conditions that the dimension n must be the same for all charts of a given \mathcal{M} . The integer n is therefore called the *dimension* of the manifold.

Moreover, if all the connecting functions $f \in C^p(\mathbb{R}^n)$, we can also say the manifold is p -times differentiable. We shall in general assume $f \in C^\infty(\mathbb{R}^n)$, meaning we can differentiate any functions as many times as we need.

Mathematically speaking, a *manifold* is an equivalence class of atlases: two atlases are equivalent if there exists a bijective correspondence between them. This puts on a firm mathematical basis the idea of a geometric space as a set in which coordinates may be introduced, but whose properties do not depend on the specific choice of coordinates we use to identify its points. This however does not mean that the concept of manifold and tools of differential geometry are restricted to sets with a natural geometrical interpretation. For example, the phase space and configuration space of classical mechanics are manifold; the three parameters (angles) of $O(3)$ form a three-dimensional manifold; vector spaces are manifolds of dimension equal to the number of basis vectors.

In order to prove that a given space is a manifold, it is sufficient to find one atlas which covers it. We shall now try to define atlases for the 2-dimensional sphere S^2 and the cone.

Example: the sphere

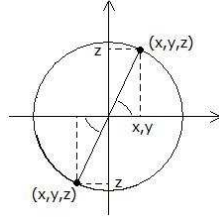


Figure 2.3: The sphere in cartesian coordinates.

We can embed the sphere S^2 in \mathbb{R}^3 with coordinates x , y and z by imposing the (smooth) condition

$$f_s = x^2 + y^2 + z^2 - R^2 = 0, \quad (2.1.6)$$

where R is the radius of S^2 . We can then cover the sphere by means of four charts (A_i, ϕ_i) . In particular, we cover the northern and southern (open) hemispheres by mapping them into the open disk $x^2 + y^2 < R$,

$$\left\{ \begin{array}{l} A_1 = \{x^2 + y^2 + z^2 = R^2 ; z > 0\} \\ \phi_1 = (x, y) \end{array} \right. \quad \left\{ \begin{array}{l} A_2 = \{x^2 + y^2 + z^2 = R^2 ; z < 0\} \\ \phi_2 = (x, y) \end{array} \right. \quad (2.1.7)$$

while the equator is covered by two open strips, namely

$$\left\{ \begin{array}{l} A_3 = \{x^2 + y^2 + z^2 = R^2 ; -z_0 < z < z_0 ; x > -x_0\} \\ \phi_3 = (x, \theta) \end{array} \right. \quad (2.1.8)$$

and

$$\begin{cases} A_4 = \{x^2 + y^2 + z^2 = R^2 ; -z_0 < z < z_0 ; x < x_0\} \\ \phi_4 = (x, \theta) \end{cases} \quad (2.1.9)$$

where $0 < z_0 < R$, $\tan(\theta) = y/x$ and $0 < x_0 < R$. This shows that S^2 is indeed a manifold.

Note that it is common practice to pretend A_3 and A_4 can be replaced by one strip with *periodic boundary condition*, namely

$$\begin{cases} A = \{x^2 + y^2 + z^2 = R^2 ; -z_0 < z < z_0\} \\ \phi = (z, \theta) , \quad \theta = \arctan(x/y) , \end{cases} \quad (2.1.10)$$

where $-\pi \leq \theta < \pi$. However, this is not an open (closed) subset of \mathbb{R}^2 . This periodic boundary condition is also used to define polar coordinates in \mathbb{R}^2 , or on the torus. However, for the sake of some more mathematical rigor, the plane \mathbb{R}^2 should be covered by two infinite (open) punctured disks $\|x - x_1\| > 0$ and $\|x - x_2\| > 0$ centered around $x_1 \neq x_2$. It is not uncommon that one meets with mathematical difficulties, for example, trying to solve differential equations, when such subtleties are overlooked.

Example: the cone

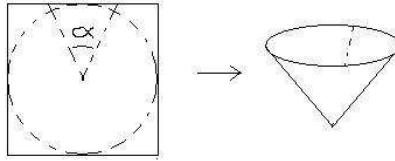


Figure 2.4: Mapping the cone into the plane (α is the deficit angle).

We first note the cone cannot be smoothly embedded in \mathbb{R}^3 . For example, one could describe the cone by means of the following condition

$$f_c = z - a \sqrt{x^2 + y^2} = 0 , \quad (2.1.11)$$

where $a > 0$ is a constant, and x and y may therefore look like valid coordinates on the cone. Although f_c is continuous in the three coordinates x , y and z , its first derivative with respect to x or y is not smooth at the tip ($x = y = 0$), since it shows a cusp there,

$$\frac{\partial f_c}{\partial x} = \frac{a^2 |x|}{f} . \quad (2.1.12)$$

We shall see later what this implies.

In fact, the cone cannot be smoothly mapped into one open subset A_1 of \mathbb{R}^2 . For example, suppose we cut the cone along a line starting from the apex and spread it flat onto \mathbb{R}^2 (see Fig. 2.4). The image of the cone on the plane would therefore be a disk minus the triangle within the so called deficit angle. However, such a set must be open to define a chart, which necessarily leaves out points along the cut. This can be seen by employing polar coordinates on the plane to cover A_1 , namely $r > 0$ and $0 < \theta < 2\pi - \alpha$. We can improve things a bit by employing two (or more) charts. For example, we can define another open subset A_2 covered by $r > 0$ and $-(\pi - \alpha/2) < \theta < \pi - \alpha/2$. However the apex ($r = 0$) must still be excluded, because it corresponds to all possible values of θ both in A_1 and A_2 . One must therefore conclude the cone is not a differentiable manifold.

It is again common (although improper) to refer to the above sets A_1 and A_2 as to one set A defined by $r > 0$ with any two points identified when their respective angles differ by $2\pi - \alpha$, or $0 \leq \theta < 2\pi - \alpha$ (which shows that this subset of \mathbb{R}^2 is neither open nor closed!).

2.1.2 Curves

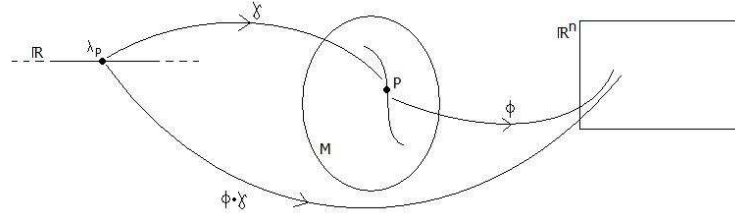


Figure 2.5: A curve

A curve is a (at least) continuous map from an interval $I \subseteq \mathbb{R}$ to the manifold \mathcal{M} (see Fig. 2.5),

$$\gamma : I \subseteq \mathbb{R} \rightarrow \mathcal{M} , \quad (2.1.13)$$

so that, given coordinates ϕ for the portion of \mathcal{M} including the curve,

$$\phi \circ \gamma : I \subseteq \mathbb{R} \rightarrow \mathbb{R}^n , \quad (2.1.14)$$

which is usually written as

$$x^i = x^i(\lambda) , \quad (2.1.15)$$

where λ is the real parameter which identifies points on the curve and x^i are its coordinates in the given chart. If the n functions $x^i = x^i(\lambda) \in C^p(\mathbb{R})$, then the curve γ is p -differentiable. Note that, according to our definition, a reparameterization of λ , that is

$$\lambda' = \lambda'(\lambda) , \quad (2.1.16)$$

defines a different curve γ' , although γ and γ' contain the same points ³.

³This will allow us to distinguish two particles following the same path at different speeds.

2.1.3 Functions

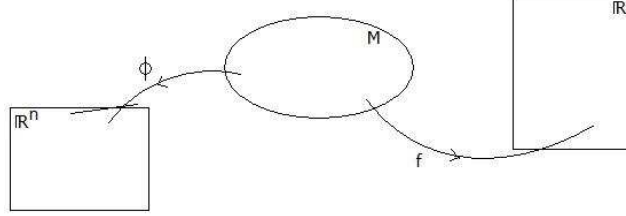


Figure 2.6: Function on a manifold.

A function on a manifold is simply an application

$$f : \mathcal{M} \rightarrow \mathbb{R} . \quad (2.1.17)$$

Given the existence of charts, it can be more easily studied by means of the composition $f \circ \phi^{-1} : \mathbb{R}^n \rightarrow \mathbb{R}$, where ϕ is a chart covering (part of) the domain of f (see Fig. 2.6). This definition gives a precise meaning to the commonly used notation

$$f = f(x^i) , \quad (2.1.18)$$

which, however, does not distinguish between f and its composition $f \circ \phi^{-1}$.

From now on, we shall assume all functions f we deal with are at least continuous (so as to preserve the topology) and differentiable as many times as necessary, which is denoted by the symbol $f \in C^\infty(\mathcal{M})$, and practically means

$$f \circ \phi^{-1} \in C^\infty(\phi(A_i) \subseteq \mathbb{R}^n) , \quad (2.1.19)$$

for all the charts (A_i, ϕ_i) of the given manifold \mathcal{M} .

2.1.4 Vectors and vector fields

We recall a vector in \mathbb{R}^n can be viewed as a displacement (an oriented straight path between two points), but also as the tangent to a curve. The first interpretation is difficult to make sense on a generic manifold, since displacements involve different points (arbitrarily separated) and the notion of a straight path between them is not necessarily given.

We shall instead generalise the concept of vector tangent to a curve (see Fig. 2.7) and define a vector \vec{v} at a point P of a manifold \mathcal{M} as an application which associates to any (differentiable) function f defined in a neighbourhood of P the derivative of that function along the curve, that is

$$\vec{v} : f \rightarrow \vec{v}_\gamma(f) = \left. \frac{df}{d\lambda} \right|_{\lambda=\lambda_P} \in \mathbb{R} , \quad (2.1.20)$$

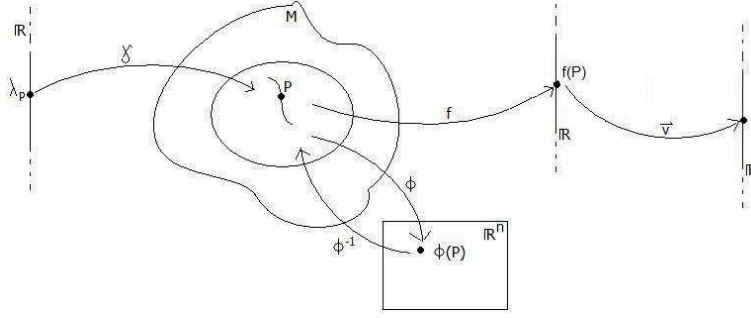


Figure 2.7: Vectors.

where γ is a given curve through P , such that $\gamma(\lambda_P) = P$ and f in the rightmost expression stands for $f \circ \gamma = f(\lambda)$. Since in any neighbourhood of P there is a chart ϕ into \mathbb{R}^n , we can also write

$$\begin{aligned}
 \vec{v}_\gamma(f) &= \left. \frac{d(f \circ \gamma)}{d\lambda} \right|_{\lambda=\lambda_P} = \left. \frac{d}{d\lambda} (f \circ \phi^{-1} \circ \phi \circ \gamma) \right|_{\lambda=\lambda_P} \\
 &= \left. \frac{d}{d\lambda} [f(x^i) \circ x^i(\lambda)] \right|_{\lambda=\lambda_P} \\
 &= \left. \frac{d}{d\lambda} f(x^i(\lambda)) \right|_{\lambda=\lambda_P} \\
 &= \sum_{i=1}^n \left. \frac{\partial f}{\partial x^i} \frac{dx^i}{d\lambda} \right|_{\lambda=\lambda_P} .
 \end{aligned} \tag{2.1.21}$$

If we now omit the generic function f (and the point P) and assume repeated indices are summed over, we can formally write the above expression in the familiar form

$$\vec{v}_\gamma = \frac{dx^i}{d\lambda} \frac{\partial}{\partial x^i} = \frac{d}{d\lambda} . \tag{2.1.22}$$

This gives a mathematically precise (and coordinate independent) meaning to the naive notion of a vector as the tangent to γ at P .

Note that the definition (2.1.20) immediately implies that a vector acts linearly on functions, since

$$\vec{v}_\gamma(a f + b g) = \frac{d}{d\lambda} (a f + b g) = a \frac{df}{d\lambda} + b \frac{dg}{d\lambda} , \quad \forall a, b \in \mathbb{R} , \tag{2.1.23}$$

and for all functions f and g defined in a neighbourhood of the point P . In fact, one could define a vector \vec{v} at a point P as a *linear* functional that acts on all the functions defined in a neighbourhood of the point P , and then, prove that there exists a curve $\gamma = \gamma(\lambda)$ such that Eq. (2.1.20) holds using Eq. (2.1.23). However, the latter (equivalent) definition is more formal and does not make clear from the very beginning that we are just generalising the notion of tangent to a curve.

We recall that, in the tensor formalism, we defined vectors as objects with special transformation properties under (certain global) coordinate transformations. It is now easy to see the true geometrical meaning of that definition by simply rewriting Eq. (2.1.22) as

$$\vec{v} = \frac{dx^i}{d\lambda} \frac{\partial}{\partial x^i} = v^i \vec{e}_i , \quad (2.1.24)$$

where

$$\vec{e}_i = \frac{\partial}{\partial x^i} , \quad (2.1.25)$$

is a coordinate (basis) vector, that is the vector tangent to the coordinate line defined by constant x^j for $j \neq i$ and passing through P . Under a general and local change of coordinates in a neighbourhood of P ,

$$y^i = y^i(x^j) , \quad (2.1.26)$$

we then have (take note of the position of the indices)

$$\begin{cases} \frac{dx^i}{d\lambda} = \frac{\partial x^i}{\partial y^j} \frac{dy^j}{d\lambda} \equiv \frac{\partial x^i}{\partial y^{j'}} v^{j'} \\ \frac{\partial}{\partial x^i} = \frac{\partial y^j}{\partial x^i} \frac{\partial}{\partial y^j} \equiv \frac{\partial y^{j'}}{\partial x^i} \vec{e}_{j'} , \end{cases} \quad (2.1.27)$$

where we temporarily returned to the old notation of primed indices. In other words, components transform like we discussed in the tensor formalism (and according to the kind of general linear transformation we studied in Chapter A), but the actual vector remains the same because the basis also changes (inversely).

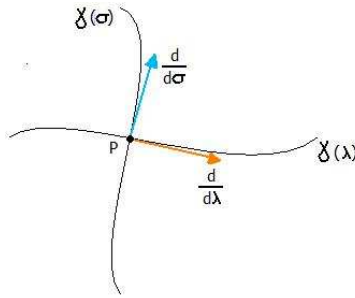


Figure 2.8: Two vectors at the same point.

We have not yet proved that vectors admit a basis. So let us consider two curves across a point P , respectively parameterized by λ and σ , which generate two different vectors $\vec{v} = \frac{d}{d\lambda}$ and $\vec{w} = \frac{d}{d\sigma}$ in P (see Fig. 2.8). It is easy to see that their linear combination can be expressed in terms of the same coordinate vectors \vec{e}_i ,

$$a \frac{d}{d\lambda} + b \frac{d}{d\sigma} = \left(a \frac{dx^i}{d\lambda} + b \frac{dx^i}{d\sigma} \right) \frac{\partial}{\partial x^i} = \left(a \frac{dx^i}{d\lambda} + b \frac{dx^i}{d\sigma} \right) \vec{e}_i . \quad (2.1.28)$$

Since there are n coordinates x^i , we can have n families of independent curves and parameters. This defines a vector space at each point P of a manifold \mathcal{M} (see Fig 2.9), called the *tangent space* T_P . Clearly, at each P , we have $T_P = \mathbb{R}^n$.

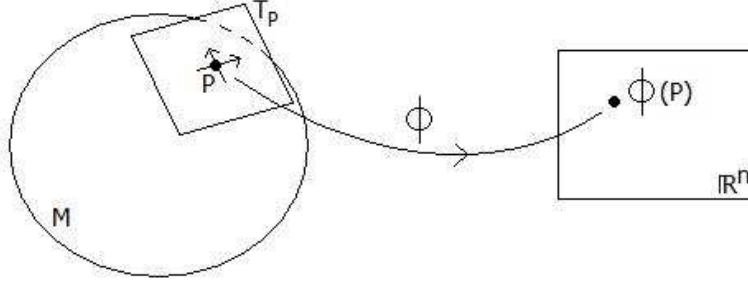


Figure 2.9: Tangent space.

Example: the cone redux

We can now see more clearly the problem with the cone. From Eq. (2.1.12), we infer that there is no unique definition of T_P at the tip ($x = y = z = 0$): the embedding condition $f_c(x, y, z) = 0$ implicitly defines a function $z = z(x, y)$ which must have support on all of the cone. However, depending on how we take the limit $x \rightarrow 0$, the coordinate vector $\frac{\partial}{\partial x}$ acting on z takes different values at the tip (namely ± 1), which is not allowed. And the same is of course true for $\frac{\partial}{\partial y}$.

Vector fields

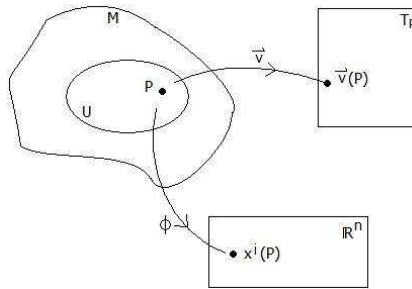


Figure 2.10: Vector field.

A vector field in an open set $U \subseteq \mathcal{M}$ is an application which maps each point $P \in U$ into a vector $\vec{v}(P) \in T_P$ (see Fig. 2.10). Given a chart in U , we can also write

$$\vec{v} \circ \phi^{-1} = \vec{v}(x^i) , \quad (2.1.29)$$

where, as usual, the notation does not distinguish between geometrical vectors and their coordinate representation.

Basis vectors and basis vector fields

We already anticipated that a basis of T_P is given by the coordinate vectors $\{\frac{\partial}{\partial x^i}\}$, since all vectors can be written as $\vec{v} = v^i \frac{\partial}{\partial x^i} = v^i \vec{e}_i$. This is a necessary condition, but, in order to prove that these \vec{e}_i form a basis, we need to show they are also linearly independent. The latter can be proven by recalling that the determinant of the Jacobian matrix for a change of coordinates $y^i = y^i(x^i)$ must not vanish, that is

$$J = \det \begin{bmatrix} \frac{\partial y^1}{\partial x^1} & \cdots & \frac{\partial y^1}{\partial x^n} \\ \frac{\partial y^2}{\partial x^1} & \cdots & \frac{\partial y^2}{\partial x^n} \\ \vdots & \cdots & \vdots \\ \frac{\partial y^n}{\partial x^1} & \cdots & \frac{\partial y^n}{\partial x^n} \end{bmatrix} \neq 0 . \quad (2.1.30)$$

It then follows that the n n -tuples of row (or column) entries are linearly independent, and so are the n vectors

$$\vec{e}_j = \frac{\partial}{\partial x^j} = \frac{\partial y^i}{\partial x^j} \frac{\partial}{\partial y^i} . \quad (2.1.31)$$

Since all coordinates $\phi = x^i$ are defined in open sets, the above definition of coordinate basis vectors can be naturally extended to define coordinate basis vector fields in the chart of ϕ . It is important however to remark that coordinate basis vectors at different points, say P and Q , belong to different tangent spaces, T_P and T_Q , and cannot be composed linearly, that is operations such as $a \vec{e}_i(P) + b \vec{e}_j(Q)$ are *not* allowed.

Fiber bundles

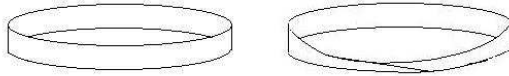


Figure 2.11: A simple band (left) and the Moebius strip (right).

The set of all tangent spaces to the points of a manifold together with the base manifold itself is called *tangent bundle* $T\mathcal{M}$. More precisely, $T\mathcal{M} = (\mathcal{M}, \{T_P : \forall P \in \mathcal{M}\})$, where the original manifold \mathcal{M} is now called the base manifold, and the tangent spaces T_P are the fibers. One can show that $T\mathcal{M}$ is also a manifold (thus continuity is a well defined property) and that vector fields can be viewed as sections of $T\mathcal{M}$.

An example of non-trivial tangent bundle is given by a closed band. Locally (at each point P of the band), the tangent bundle is simply given by $\mathbb{R}^2 \times \mathbb{R}^2$. However, the global tangent bundle is not necessarily the direct product of two manifolds: consider the Moebius strip obtained by cutting the band and twisting the edges before pasting them again (see Fig. 2.11). One therefore needs to travel twice along the strip in order to come back to the starting point. Spinors belong to such manifolds.

2.1.5 Vector fields and integral curves

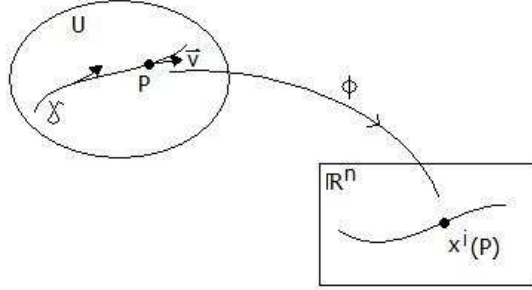


Figure 2.12: Integral curve of a vector field.

We call $\gamma = \gamma(\lambda)$ an *integral curve* of a given vector field \vec{Y} if its tangent vector $\frac{d}{d\lambda}$ is given by an element of \vec{Y} at all points $P \in \gamma$ (see Fig. 2.12), that is

$$\frac{d}{d\lambda} = \vec{Y} . \quad (2.1.32)$$

In particular, given a point P_0 , we can define the integral curve of \vec{Y} through P_0 by the system of equations

$$\begin{cases} \frac{d}{d\lambda} = \vec{Y} \\ \gamma(\lambda_0) = P_0 . \end{cases} \quad (2.1.33)$$

Upon introducing a chart $\phi = x^i$, in which the curve is represented by $x^i = x^i(\lambda)$, with $i = 1, \dots, n$, the above system becomes a system of n Cauchy problems for the coordinates of the curve, namely

$$\begin{cases} \frac{dx^i(\lambda)}{d\lambda} = Y^i(\lambda) \\ x^i(\lambda_0) = x^i(P_0) , \end{cases} \quad (2.1.34)$$

where the components of the vector field \vec{Y} are given at any point $P = \phi^{-1}(x^i(\lambda))$ by

$$\vec{Y}(P) = Y^i(x^j(\lambda)) \frac{\partial}{\partial x^i} = Y^i(\lambda) \frac{\partial}{\partial x^i} . \quad (2.1.35)$$

Theorems of calculus ensure that the problem (2.1.34) always admits one solution $x^i = x^i(\lambda)$ in a (sufficiently small) neighbourhood of the point P_0 , therefore integral curves of a vector field \vec{Y} always exist locally.

Exponential map

The formal solution to the n first order differential equations and initial conditions (2.1.34) can be written as

$$x^i = e^{(\lambda - \lambda_0) \vec{Y}} x^i \Big|_{\lambda = \lambda_0} , \quad (2.1.36)$$

which is called the *exponential map* and describes the flow of velocity \vec{Y} in a neighborhood of P_0 in the coordinate space \mathbb{R}^n .

Let us see in detail how the exponential of a vector field generates integral curves, and the (tangent) vectors \vec{Y} therefore act as generators of the displacements⁴. Given the vector field $\vec{Y} = \frac{d}{d\lambda}$, its integral curve $\gamma = \gamma(\lambda)$ across a point P_0 , and the chart ϕ which maps P_0 into $\phi(P_0) = x^i(\lambda_0)$, we can Taylor expand the n coordinates in a neighbourhood of P_0 along the integral curve γ as

$$\begin{aligned} x^i(\lambda_0 + \varepsilon) &= x^i(\lambda_0) + \varepsilon \frac{dx^i}{d\lambda} \Big|_{\lambda_0} + \frac{\varepsilon^2}{2} \frac{d^2 x^i}{d\lambda^2} \Big|_{\lambda_0} + \dots \\ &= \left(1 + \varepsilon \frac{d}{d\lambda} + \frac{\varepsilon^2}{2!} \frac{d^2}{d\lambda^2} + \dots \right) x^i \Big|_{\lambda_0} . \end{aligned} \quad (2.1.37)$$

Next, note that all terms in the above expansion are well-defined, because the coordinates of any point P , $x^i = x^i(P)$, are functions on the manifold. In the language of tensor calculus, *coordinates are scalars*, in agreement with the fact that the measurements by which an observer assigns coordinates to a point cannot be questioned by other observers⁵. The action of the vector \vec{Y} on the coordinates is then well-defined by the very definition of a vector and we can rewrite the above as

$$x^i(\lambda_0 + \varepsilon) = \exp \left\{ \varepsilon \frac{d}{d\lambda} \right\} x^i \Big|_{\lambda_0} = e^{\varepsilon \vec{Y}} x^i \Big|_{\lambda_0} . \quad (2.1.38)$$

A neat example is given by choosing the parameter along γ as one of the coordinates, say $\lambda = x^1$. We then have $\vec{Y} = \frac{\partial}{\partial x^1}$ and, setting $\epsilon = x^1 - x_0^1$, one easily obtains

$$\begin{aligned} x^i &= e^{\epsilon \frac{\partial}{\partial x^1}} x^i \Big|_{x^1 = x_0^1} \\ &= x_0^i + (x^1 - x_0^1) \frac{\partial x^i}{\partial x^1} \Big|_{x^1 = x_0^1} + \frac{(x^1 - x_0^1)^2}{2} \frac{\partial^2 x^i}{\partial (x^1)^2} \Big|_{x^1 = x_0^1} + \dots \\ &= \begin{cases} x_0^i + (x^1 - x_0^1) = x^1 , & i = 1 \\ x_0^i , & i \neq 1 . \end{cases} \end{aligned} \quad (2.1.39)$$

⁴The definition (A.2.7) of a Lie group is a special case of the exponential map, that generates all the elements of the group from the identity, and the generators of the Lie algebra play the role of the vectors \vec{Y} .

⁵Equivalently, since the coordinates identify a specific observer, the very definition of each observer is based on scalar quantities.

Using the chain rule for the derivation of composite functions, it is straightforward to generalise the above expression to any function (or field) defined in a neighbourhood of P . For example,

$$f(\lambda_0 + \varepsilon) = \exp \left\{ \varepsilon \frac{d}{d\lambda} \right\} f \Big|_{\lambda_0} = e^{\varepsilon \vec{Y}} f \Big|_{\lambda_0} , \quad (2.1.40)$$

where $f(\lambda) = f(x^i(\lambda)) = f \circ \phi^{-1} \circ \phi \circ \gamma$.

Lie brackets and non-coordinate basis

We now define the Lie brackets (commutator) for vector fields and the vanishing of the commutators as a necessary and sufficient condition for vector fields to generate a reference frame.

Let us consider two vector fields, $\vec{V} = \frac{d}{d\lambda}$ and $\vec{W} = \frac{d}{d\mu}$, and compute their commutator. For simplicity, let us just consider points inside one chart $(U, \phi = x^i)$, and use the corresponding coordinate basis in T_P for all $P \in U$, so that

$$\begin{aligned} [\vec{V}, \vec{W}] &= \frac{d}{d\lambda} \frac{d}{d\mu} - \frac{d}{d\mu} \frac{d}{d\lambda} \\ &= v^i \frac{\partial}{\partial x^i} \left(w^j \frac{\partial}{\partial x^j} \right) - w^i \frac{\partial}{\partial x^i} \left(v^j \frac{\partial}{\partial x^j} \right) \\ &= v^i w^j \left(\frac{\partial}{\partial x^i} \frac{\partial}{\partial x^j} - \frac{\partial}{\partial x^j} \frac{\partial}{\partial x^i} \right) + \left(v^i \frac{\partial w^j}{\partial x^i} - w^i \frac{\partial v^j}{\partial x^i} \right) \frac{\partial}{\partial x^j} \\ &= \left(v^i \frac{\partial w^j}{\partial x^i} - w^i \frac{\partial v^j}{\partial x^i} \right) \frac{\partial}{\partial x^j} , \end{aligned} \quad (2.1.41)$$

where all derivatives are computed at the same point P , and the results is therefore an element of T_P . This is a first remarkable result: the commutator of two vectors is still a vector ⁶. Moreover, the commutator vanishes if the two fields \vec{V} and \vec{W} are *coordinate vectors*, that is, if there exit coordinates, say x^1 and x^2 , such that $\vec{V} = \frac{\partial}{\partial x^1}$ and $\vec{W} = \frac{\partial}{\partial x^2}$. In fact, if this is the case, $v^j = \delta_1^j$ and $w^i = \delta_2^i$ are obviously constant and the bracket in the last line above vanishes.

Given two fields $\vec{X} = \frac{d}{d\lambda}$ and $\vec{Y} = \frac{d}{d\mu}$, the geometrical meaning of $[\vec{X}, \vec{Y}] \neq 0$ is explained in Fig. 2.13. By means of the exponential maps of the two fields, we easily obtain the coordinates of the points A and B reached by moving away from P along \vec{X} first and \vec{Y} next, or along \vec{Y} first and \vec{X} next,

$$\begin{aligned} x^i(A) &= \exp \left\{ \varepsilon \frac{d}{d\mu} \right\} \exp \left\{ \varepsilon \frac{d}{d\lambda} \right\} x^i \Big|_P \\ x^i(B) &= \exp \left\{ \varepsilon \frac{d}{d\lambda} \right\} \exp \left\{ \varepsilon \frac{d}{d\mu} \right\} x^i \Big|_P , \end{aligned} \quad (2.1.42)$$

⁶This property was assumed in the very definition of a Lie group G , and we can now see that the generators of a Lie algebra \mathcal{G} are “vectors” that live in the tangent space of the identity of the group manifold, $\mathcal{G} = T_1$.

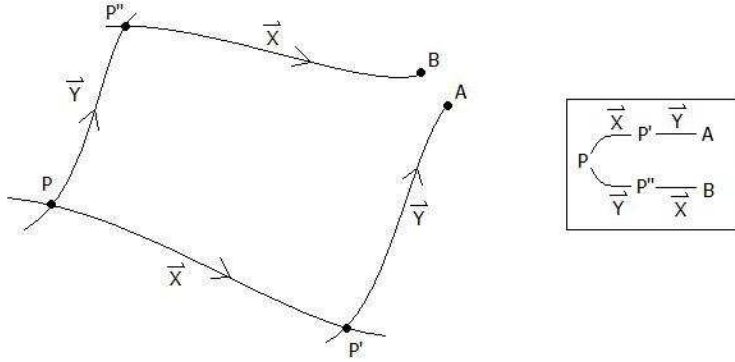


Figure 2.13: Commutator of two vector fields.

from which results the difference

$$x^i(B) - x^i(A) = \varepsilon^2 \left[\frac{d}{d\lambda}, \frac{d}{d\mu} \right] x^i \Big|_P + O(\varepsilon^3) . \quad (2.1.43)$$

If the commutator of the two fields does not vanish, $A \neq B$, and the path $PA \cup BP$ does not close.

Example: polar coordinates in \mathbb{R}^2

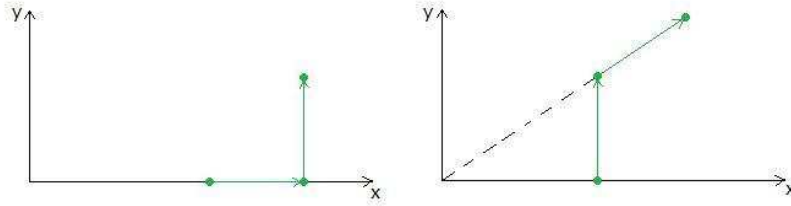


Figure 2.14: Polar coordinates on the plane.

Consider a given point $P_0 \in \mathbb{R}^2$ of cartesian coordinates (x_0, y_0) . In T_{P_0} , let us then define the four vectors

$$\begin{cases} \hat{x} = \frac{\partial}{\partial x} \\ \hat{y} = \frac{\partial}{\partial y} \end{cases} \quad \text{and} \quad \begin{cases} \hat{r} = \cos(\theta_0) \hat{x} + \sin(\theta_0) \hat{y} \\ \vec{\theta} = -\sin(\theta_0) \hat{x} + \cos(\theta_0) \hat{y} , \end{cases} \quad (2.1.44)$$

where $\theta_0 = \arctan(y_0/x_0)$. By extending the above definitions to all points P in a neighbourhood of P_0 , we can define four vector fields, and upon considering the integral curves of these four vector fields, one finds

$$[\hat{x}, \hat{y}] = 0 \quad [\hat{r}, \vec{\theta}] \neq 0 . \quad (2.1.45)$$

This is obvious if we choose particular curves. For example, let us start from P_0 on the x -axis and move first along \hat{r} and then in the direction of $\vec{\theta}$, and compare the result with the inverted steps (see the left and right panels in Fig. 2.14). This shows that $\{\hat{r}, \vec{\theta}\}$ do *not* form a coordinate basis (although they are a basis of T_P for all P in most of the plane).

Of course, the proper coordinate basis for polar coordinates is obtained by rescaling

$$\vec{\theta} \rightarrow \hat{\theta} = \vec{\theta}/r, \quad (2.1.46)$$

since motion along $\hat{\theta}$ corresponds to rotation of a given angle θ around the origin (whereas motion along $\vec{\theta}$ correspond to rotation of an arc of length $r\theta$).

Lie algebra of vector fields

We saw that, if \vec{A} and \vec{B} are coordinate vector fields, then $[\vec{A}, \vec{B}] = 0$. Let us now prove that the vanishing of the commutator is also a sufficient condition for the two fields to be coordinate, that is $[\vec{A}, \vec{B}] = 0$ implies that there exist two coordinates whose lines are tangent to \vec{A} and \vec{B} .

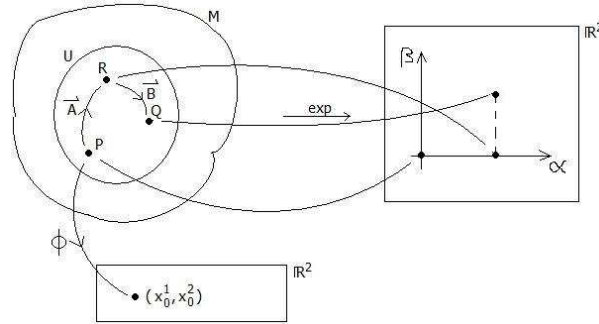


Figure 2.15: Exponential maps of \vec{A} and \vec{B} .

Let us consider a two-dimensional manifold for simplicity and assume the vector fields \vec{A} and \vec{B} are linearly independent in their domain of definition, with

$$\vec{A} = \frac{d}{d\lambda}, \quad \vec{B} = \frac{d}{d\mu}. \quad (2.1.47)$$

Let us start from a point P and first move along \vec{A} ,

$$P \rightarrow R : x^i(R) = \exp \left\{ \lambda_1 \frac{d}{d\lambda} \right\} x^i \Big|_P \quad (2.1.48)$$

and then along \vec{B} ,

$$R \rightarrow Q : x^i(Q) = \exp \left\{ \mu_1 \frac{d}{d\mu} \right\} \exp \left\{ \lambda_1 \frac{d}{d\lambda} \right\} x^i \Big|_P, \quad (2.1.49)$$

in which we assumed all relevant points are included in the same chart $(U, \phi = x^i)$ for simplicity. We should be able to look at λ_1 and μ_1 as coordinates of the final point Q ,

$$x^i(\alpha, \beta) = \exp \left\{ \beta \frac{d}{d\mu} \right\} \exp \left\{ \alpha \frac{d}{d\lambda} \right\} x^i \Big|_P . \quad (2.1.50)$$

Now, if α and β are to be coordinates, the corresponding basis vectors should be $\frac{\partial}{\partial\alpha}$ and $\frac{\partial}{\partial\beta}$, with

$$\begin{cases} \frac{\partial}{\partial\alpha} = \frac{\partial x^i}{\partial\alpha} \frac{\partial}{\partial x^i} \\ \frac{\partial}{\partial\beta} = \frac{\partial x^i}{\partial\beta} \frac{\partial}{\partial x^i} . \end{cases} \quad (2.1.51)$$

where the matrix of elements $[\frac{\partial x^i}{\partial\alpha} \frac{\partial x^i}{\partial\beta}]$ must be invertible, that is

$$J = \det \begin{bmatrix} \frac{\partial x^1}{\partial\alpha} & \frac{\partial x^2}{\partial\alpha} \\ \frac{\partial x^1}{\partial\beta} & \frac{\partial x^2}{\partial\beta} \end{bmatrix} \neq 0 . \quad (2.1.52)$$

In fact,

$$\begin{cases} \frac{\partial x^i}{\partial\alpha} = \frac{\partial}{\partial\alpha} \left(\exp \left\{ \beta \frac{d}{d\mu} \right\} \exp \left\{ \alpha \frac{d}{d\lambda} \right\} x^i \right)_P = \exp \left\{ \beta \frac{d}{d\mu} \right\} \exp \left\{ \alpha \frac{d}{d\lambda} \right\} \frac{dx^i}{d\lambda} \Big|_P \\ \frac{\partial x^i}{\partial\beta} = \frac{\partial}{\partial\beta} \left(\exp \left\{ \beta \frac{d}{d\mu} \right\} \exp \left\{ \alpha \frac{d}{d\lambda} \right\} x^i \right)_P = \exp \left\{ \beta \frac{d}{d\mu} \right\} \exp \left\{ \alpha \frac{d}{d\lambda} \right\} \frac{dx^i}{d\mu} \Big|_P , \end{cases} \quad (2.1.53)$$

which follows from

$$\frac{\partial}{\partial\alpha} \exp \left\{ \alpha \frac{d}{d\lambda} \right\} = \exp \left\{ \alpha \frac{d}{d\lambda} \right\} \frac{d}{d\lambda} , \quad (2.1.54)$$

and the hypothesis

$$\left[\frac{d}{d\lambda}, \frac{d}{d\mu} \right] = 0 . \quad (2.1.55)$$

Eq. (2.1.53) shows that $\frac{\partial}{\partial\alpha}$ is just the vector field $\frac{d}{d\lambda}$ evaluated at Q , as well as $\frac{\partial}{\partial\beta}$ is $\frac{d}{d\mu}$ evaluated at Q . Further, since $\frac{d}{d\lambda}$ and $\frac{d}{d\mu}$ were assumed linearly independent, the determinant (2.1.52) of the matrix built out from their components must be different from zero.

2.1.6 One-forms

One-forms are linear functionals on vectors and the geometrical counterparts of co-vectors.

Let us consider a point P on the manifold \mathcal{M} and the tangent space T_P . A 1-form at P is a linear functional \tilde{w} acting on vectors in T_P ,

$$\tilde{w} : T_P \rightarrow \mathbb{R} , \quad (2.1.56)$$

such that

$$\tilde{w}(\alpha \vec{v} + \beta \vec{u}) = \alpha \tilde{w}(\vec{v}) + \beta \tilde{w}(\vec{u})$$

$$(\alpha \tilde{w})(\vec{v}) = \alpha \tilde{w}(\vec{v}) \quad (2.1.57)$$

$$(\tilde{w} + \tilde{\sigma})(\vec{v}) = \tilde{w}(\vec{v}) + \tilde{\sigma}(\vec{v}) .$$

Note that linearity implies the action of a given 1-form on a generic vector is completely defined by its action on a basis of T_P . Several equivalent notations are in use, for example

$$\tilde{w}(\vec{v}) = \vec{v}(\tilde{w}) = \langle \tilde{w}, \vec{v} \rangle = \langle \tilde{w} \mid \vec{v} \rangle . \quad (2.1.58)$$

One-forms acting on the same T_P form a vector space T_P^* , dual to T_P , and the collection of all T_P^* forms the cotangent bundle $T^*\mathcal{M}$.

A 1-form field is an application which associates a 1-form from T_P^* for each point P of a manifold \mathcal{M} , and, as usual, we shall always assume such a map is sufficiently smooth.

The gradient of a function

The gradient of a given function is usually introduced as a vector, but its geometrical interpretation is in fact the prototype of a 1-form.

Let f be any function from a manifold \mathcal{M} to \mathbb{R} and \vec{V} a vector field. By definition, we have

$$\vec{V}(f) = \frac{df}{d\lambda} \in \mathbb{R} . \quad (2.1.59)$$

We can now use the same Eq. (2.1.59) to define the 1-form $\tilde{d}f\left(\frac{d}{d\lambda}\right)$ as the “reverse” operation, namely the 1-form which associates a real number to any \vec{V} , given a fixed function f ,

$$\vec{V}(f) = \frac{df}{d\lambda} = \tilde{d}f\left(\frac{d}{d\lambda}\right) . \quad (2.1.60)$$

The difference between the two interpretations is that, in Eq. (2.1.59), the vector \vec{V} is fixed and the function f is the generic argument, whereas in Eq. (2.1.60), f is fixed and \vec{V} can vary. The result (a real number) is however the same (for the same pair of vector and function). In

particular, by making use of a chart $\phi = x^i$, and still denoting with f the composed function $f \circ \phi^{-1}$, we obtain

$$\frac{df}{d\lambda} = \frac{\partial f(x)}{\partial x^i} \frac{dx^i}{d\lambda}, \quad \frac{\partial f(x)}{\partial x^i} = \nabla_i f = df_i, \quad (2.1.61)$$

where df_i are the components of the 1-form $\tilde{d}f$ we commonly call the gradient of the function f .

The geometrical interpretation of $\tilde{d}f$ is rather illuminating. In elementary calculus, one is taught that $\vec{\nabla}f$ is a vector which points along the direction of fastest increase of the function f , roughly $\vec{\nabla}f \simeq \Delta f / \Delta x^i$. This notion however requires the concept of distance (to define the length of Δx^i), since fastest means the rate of increase of the function for unit length is maximum. Without such a notion, we can indeed find a more general meaning of $\tilde{d}f$ from Eq. (2.1.60): suppose we draw a contour plot of f (where a line represents points along which the function f takes the same value, like in an elevation map) and then consider a generic vector \vec{V} at a point P in the domain of f . The application of $\tilde{d}f$ on \vec{V} equals the “number of contour lines” the vector \vec{V} crosses in an “infinitesimal neighbourhood” of P , as can be easily seen by choosing coordinates such that $\vec{V} = V \frac{\partial}{\partial x^1}$, and

$$\tilde{d}f(\vec{V}) = V \left. \frac{\partial f}{\partial x^1} \right|_P. \quad (2.1.62)$$

If the notion of an “infinitesimal neighbourhood” appears disturbing, one could actually consider the integral curve of \vec{V} through P with unit parametric length, that is the curve tangent to \vec{V} that starts at P and ends at Q , where

$$x^i(Q) = e^{\vec{V}} x^i(P). \quad (2.1.63)$$

This yields a more general interpretation of 1-forms in any dimension as a set of (level) surfaces (of a given function f) and its action on a vector as the number of surfaces the vector crosses.

Basis one-forms and one-form components

We shall now describe in more detail the space T_P^* and introduce dual bases. Analogously to T_P , we can also define the fiber bundle $T^*\mathcal{M}$.

Let us denote a basis in the tangent space T_P by $\{\vec{e}_i; i = 1, 2, 3, \dots, n\}$. A basis in T_P^* does not need to carry any relation with the \vec{e}_i . However, since the action of a given 1-form on a generic vector is completely defined by its action on a basis of T_P , we can conveniently introduce the so called *dual basis* of 1-forms in T_P^* by means of the conditions

$$\tilde{e}^i(\vec{v}) = \tilde{e}^i(v^j \vec{e}_j) = v^i, \quad (2.1.64)$$

that is, the i^{th} basis 1-form \tilde{e}^i associates to a vector its i^{th} component. Obviously, there are n such forms and the dimension of T_P^* equals the dimension of T_P and \mathcal{M} . Note that we can equivalently write Eq. (2.1.64) as

$$\tilde{e}^i(\vec{e}_j) = \delta_j^i. \quad (2.1.65)$$

It is easy to see that the above \tilde{e}^i are actually a basis, since, given any 1-form $\tilde{q} \in T_P^*$, we have

$$\tilde{q}(\vec{v}) = \tilde{q}(v^i \vec{e}_i) = v^i \tilde{q}(\vec{e}_i) \equiv v^i q_i . \quad (2.1.66)$$

On the other hand, if Eq. (2.1.64) holds, we have

$$v^i q_i = \tilde{e}^i(\vec{v}) q_i = q_i \tilde{e}^i(\vec{v}) , \quad (2.1.67)$$

or $\tilde{q} = q_i \tilde{e}^i$ for any 1-form.

The gradient, in particular, can be written as

$$\begin{aligned} \frac{df}{d\lambda} &= \frac{\partial f}{\partial x^i} \frac{dx^i}{d\lambda} = \frac{\partial f}{\partial x^i} \tilde{dx}^i \left(\frac{dx^j}{d\lambda} \frac{\partial}{\partial x^j} \right) \\ &= \frac{\partial f}{\partial x^i} \frac{dx^j}{d\lambda} \tilde{dx}^i \left(\frac{\partial}{\partial x^j} \right) = \frac{\partial f}{\partial x^i} \frac{dx^j}{d\lambda} \delta_j^i , \end{aligned} \quad (2.1.68)$$

which shows that \tilde{dx}^i is the dual basis to the coordinate basis vectors. We finally note that under a general change of coordinates, vectors and covectors do not change, only their components do and in a way that compensates so as to keep the above real number (a scalar) unchanged.

2.1.7 Tensors and tensor fields

The general definition of (n, m) tensors at P is that of linear functionals acting on n 1-forms and m vectors,

$$T : \underbrace{T_P^* \otimes \cdots \otimes T_P^*}_n \otimes \underbrace{T_P \otimes \cdots \otimes T_P}_m \rightarrow \mathbb{R} , \quad (2.1.69)$$

where \otimes is the usual cartesian product of vector spaces. It is however easier to build them from vectors and 1-forms (covectors) by means of the outer product, like we did when studying group theory.

Tensor components and outer product

We can now define a general tensor as a combination of vectors and covectors, where by combination we mean the *outer product*, likewise denoted by \otimes . For example, by multiplying two vectors and applying the result to *dual* basis covectors, we obtain

$$(\vec{V} \otimes \vec{W})(\tilde{e}^i, \tilde{e}^j) = \vec{V}(\tilde{e}^i) \vec{W}(\tilde{e}^j) = V^i W^j , \quad (2.1.70)$$

where we note that the second expression is simply the product of two numbers (for fixed i and j). We can therefore write

$$\vec{V} \otimes \vec{W} = V^i W^j \vec{e}_i \otimes \vec{e}_j , \quad (2.1.71)$$

where now V^i and W^j are just numbers. This means the outer product of two vectors is an application

$$\vec{V} \otimes \vec{W} : T_P^* \times T_P^* \rightarrow \mathbb{R} \quad (2.1.72)$$

and linear in both arguments.

We recall that a vector is a $(1, 0)$ tensor and a covector is a $(0, 1)$ tensor. A type (m, n) tensor is then given in terms of its components by

$$T = T_{j_1 j_2 \dots j_n}^{i_1 i_2 \dots i_m} \vec{e}_{i_1} \otimes \vec{e}_{i_2} \otimes \dots \otimes \vec{e}_{i_m} \otimes \tilde{e}^{j_1} \otimes \tilde{e}^{j_2} \otimes \dots \otimes \tilde{e}^{j_n} , \quad (2.1.73)$$

where the components are in turn defined by the action of the tensor on basis vectors and covectors,

$$T_{j_1 j_2 \dots j_n}^{i_1 i_2 \dots i_m} = T(\tilde{e}^{i_1}, \tilde{e}^{i_2}, \dots, \tilde{e}^{i_m}, \vec{e}_{j_1}, \vec{e}_{j_2}, \dots, \vec{e}_{j_n}) . \quad (2.1.74)$$

Basis transformations

We now study changes of basis in T_P and T_P^* .

Let us consider a point P on a manifold \mathcal{M} and the tangent space T_P , with $\{\vec{e}_i\}$ as a basis. A change of basis in T_P , namely $\{\vec{e}_i\} \rightarrow \{\vec{e}_{i'}\}$, is determined by a non-degenerate $n \times n$ matrix (of fixed real entries), that is an element of $GL(n)$. Such a matrix has in general no particular tensorial properties and just specifies a linear transformation ⁷

$$\vec{e}_{j'} = \Lambda_{j'}^i \vec{e}_i . \quad (2.1.75)$$

Let us further consider the dual space T_P^* and, given a basis $\{\tilde{e}^i\}$, determine how 1-forms change under the same transformation (2.1.75). In particular, this question (only) makes sense if we consider the dual basis

$$\tilde{e}^i(\vec{e}_k) = \delta_k^i , \quad (2.1.76)$$

so that we can write

$$\tilde{e}^i(\vec{e}_k) \Lambda_{j'}^k = \delta_k^i \Lambda_{j'}^k = \Lambda_{j'}^i . \quad (2.1.77)$$

Since 1-forms act linearly, the above expression defines the action of \tilde{e}^i on the transformed vector basis,

$$\tilde{e}^i(\vec{e}_k) \Lambda_{j'}^k = \tilde{e}^i(\vec{e}_k \Lambda_{j'}^k) = \tilde{e}^i(\vec{e}_{j'}) . \quad (2.1.78)$$

We then denote Λ^{-1} , the inverse of Λ , by $\Lambda^{i'}$, so that

$$\Lambda^{i'}_j \Lambda^j_{k'} = \delta_{k'}^{i'} , \quad \Lambda^i_{j'} \Lambda^{j'}_k = \delta_k^i . \quad (2.1.79)$$

⁷Note though that the following notation is in agreement with the old-fashioned tensorial calculus of Section 1.4.

Upon acting from the left with Λ^{-1} on Eq. (2.1.78), we obtain

$$\Lambda^{k'}_i \tilde{e}^i(\vec{e}_k) \Lambda^k_{j'} = \Lambda^{k'}_i \tilde{e}^i(\vec{e}_{j'}) , \quad (2.1.80)$$

and likewise from Eq. (2.1.77),

$$\Lambda^{k'}_i \tilde{e}^i(\vec{e}_k) \Lambda^k_{j'} = \Lambda^{k'}_i \Lambda^i_{j'} = \delta^{k'}_{j'} . \quad (2.1.81)$$

Equating the two results, we thus see that the transformed dual basis is precisely given by

$$\tilde{e}^{k'} = \Lambda^{k'}_i \tilde{e}^i , \quad (2.1.82)$$

that is, basis 1-forms transform according to the inverse matrix Λ^{-1} . Note that in the present notation, Λ^{-1} is also the matrix that transforms vector components, whereas 1-form components transform with Λ .

Tensor operations on components

Let us now summarize all operations that map tensors T of type (n, m) into tensors defined at the same point P , but of possibly different type:

$$\text{Scalar multiplication:} \quad T^{(n,m)} \rightarrow a T^{(n,m)} , \forall a \in \mathbb{R} \quad (2.1.83)$$

$$\text{Addition:} \quad T^{(n,m)} + Q^{(n,m)} = S^{(n,m)} \quad (2.1.84)$$

$$\text{Outer product:} \quad T^{(n,m)} \otimes Q^{(n',m')} = Z^{(n+n',m+m')} \quad (2.1.85)$$

$$\text{Saturation with 1-form:} \quad T^{(n,m)}(\dots, \tilde{\omega}, \dots) = T^{(n-1,m)} \quad (2.1.86)$$

$$\text{Saturation with vector:} \quad T^{(n,m)}(\dots, \vec{v}, \dots) = T^{(n,m-1)} . \quad (2.1.87)$$

The last two operations above can then be easily generalised to any saturation of (n, m) tensors with $(p \leq m, q \leq n)$ tensors.

Change of coordinates and coordinate basis

Let us now consider a point P on a manifold \mathcal{M} , the tangent space T_P , and two charts $\phi = x^i$ and $\psi = y^i$, connected by a bijective function f (see Fig. 2.16). We can then introduce two coordinate basis for the tangent space T_P , namely $\{\frac{\partial}{\partial x^i}\}$ and $\{\frac{\partial}{\partial y^i}\}$. As we have seen before, there must be a linear transformation between these basis, namely

$$\Lambda^i_{j'} = \frac{\partial x^i}{\partial y^{j'}} \quad (2.1.88)$$

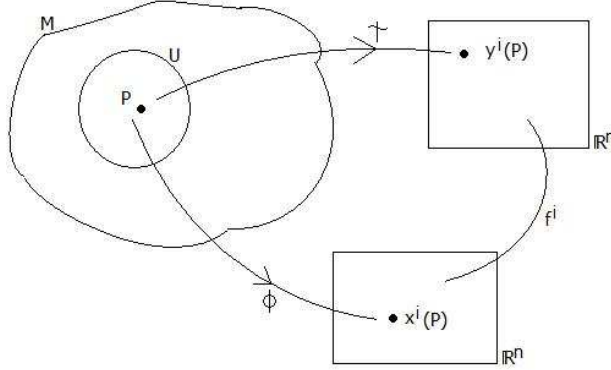


Figure 2.16: Tensor fields.

and its inverse

$$\Lambda_j^{k'} = \frac{\partial y^{k'}}{\partial x^j} . \quad (2.1.89)$$

Suppose we now (continuously ⁸) move from the point P to another point Q inside the open neighbourhood of P where both ϕ and ψ are defined. The coordinates x^i and y^i will therefore define local basis at all Q 's inside such a neighbourhood and the matrix transformation will become a (continuous) matrix field, $\Lambda = \Lambda(Q) = \Lambda(x^i(Q)) = \Lambda(y^i(Q))$. Since partial derivatives commute, we must necessarily have

$$\frac{\partial}{\partial x^i} \Lambda_j^{k'} = \frac{\partial}{\partial x^i} \frac{\partial y^{k'}}{\partial x^j} = \frac{\partial}{\partial x^j} \frac{\partial y^{k'}}{\partial x^i} = \frac{\partial}{\partial x^j} \Lambda_i^{k'} . \quad (2.1.90)$$

The conclusion is therefore that the condition

$$\frac{\partial \Lambda_j^{k'}}{\partial x^i} = \frac{\partial \Lambda_i^{k'}}{\partial x^j} \quad (2.1.91)$$

is necessary for a change of basis in the tangent space to correspond to a change of coordinates on the manifold. This is in fact a strong restriction, as we shall see better later.

2.2 Length and angles

We shall now introduce distance (length) and angles on a manifold. We first need define the length (or modulus) of a vector and the angle between two vectors belonging to the tangent space T_P of one point. Both quantities are obtained from a *scalar product* between vectors in T_P , which can in turn be introduced by means of a special tensor. Upon elevating this tensor to a field, we will finally be able to define the length of a path on \mathcal{M} .

⁸Recall the notion of continuity is well-defined in any topological set.

2.2.1 Metric tensor

A metric tensor is a type $(0, 2)$ tensor which maps any two vectors into a real number with the following properties:

1) it is symmetric

$$g(\vec{v}, \vec{w}) = g(\vec{w}, \vec{v}) = g_{ij} v^i w^j = \vec{v} \cdot \vec{w} , \quad (2.2.1)$$

where $g_{ij} = g(\vec{e}_i, \vec{e}_j)$;

2) it is non-degenerate

$$[g(\vec{v}, \vec{w}) = 0 , \quad \forall \vec{w} \in T_P \quad \Leftrightarrow \quad \vec{v} = 0] \quad \Leftrightarrow \quad \det(g_{ij}) \neq 0 . \quad (2.2.2)$$

Examples of metric tensors are the Euclidean metric $g_{ij} = \delta_{ij}$ and the Minkowski metric.

Any metric tensor automatically defines a scalar product with the expected properties. In particular, the squared modulus of a vector is given by

$$v^2 = g(\vec{v}, \vec{v}) = g_{ij} v^i v^j , \quad (2.2.3)$$

and the angle θ between two vectors by

$$g(\vec{v}, \vec{w}) = v w \cos \theta , \quad (2.2.4)$$

although the latter will only be properly defined for Euclidean metrics.

Canonical form and orthonormal bases

The components of any metric g at a point P , under a change of basis in the tangent space T_P , will change according to the matrix Λ we introduced before,

$$g' = \Lambda^T g \Lambda . \quad (2.2.5)$$

Since in given coordinates g_{ij} is a symmetric matrix, it can always be put in diagonal form. More precisely, we can always write Λ as the product of an orthogonal matrix $O^{-1} = O^T$ and a symmetric matrix $D = D^T$, such that

$$g' = D^T O^T g O D = D^T g_{(\text{diag})} D = D g_{(\text{diag})} D . \quad (2.2.6)$$

By a suitable choice of D we can finally set $|g'_{ij}| = 1$, that is, we obtain the *canonical form*

$$g'_{ij} = \pm \delta_{ij} . \quad (2.2.7)$$

The canonical form of the metric implicitly defines the orthonormal basis \vec{e}_i for vectors (and dual \tilde{e}^i for 1-forms) at the point P .

What we cannot change arbitrarily is the sign of each diagonal element, whose sum is called the *signature* of the metric. If all signs are positive (negative), the metric is positive (negative) definite and generically called *Riemannian*. The Euclidean metric is a special

case of Riemannian metric which can be put in canonical form *simultaneously* at all points of a manifold ⁹. If elements of both signs appear, the metric is said *pseudo-Riemannian*. In particular, if one element is negative (positive) and all the others are positive (negative), then is said to have *Lorentzian signature* (like the Minkowski metric).

Lowering and raising indices

Metric tensors allow us to define a map between the tangent space T_P and its dual T_P^* . Given a vector \vec{v} , we can in fact define the $(0, 1)$ tensor (or 1-form)

$$\tilde{v} = g(\vec{v}, \cdot) \in T_P^* , \quad (2.2.8)$$

whose components are given by

$$v_i = \tilde{v}(\vec{e}_i) = g(v^j \vec{e}_j, \vec{e}_i) = v^j g(\vec{e}_j, \vec{e}_i) = v^j g_{ji} = g_{ij} v^j . \quad (2.2.9)$$

In practical terms, the metric is used to lower the indices

$$v_i = g_{ij} v^j , \quad (2.2.10)$$

and this map is independent of any duality relation between \vec{e}_i and \tilde{e}^j .

Since g_{ij} is invertible, we denote its inverse with

$$g_{ij}^{-1} = g^{ij} \quad \Rightarrow \quad g_{ij} g^{jk} = \delta_i^k , \quad (2.2.11)$$

where

$$g^{-1}(\tilde{e}^i, \tilde{e}^j) = g^{ij} . \quad (2.2.12)$$

Eventually, this allows us to map a 1-form into a vector,

$$v^i = g^{ik} v_k . \quad (2.2.13)$$

Further, if the metric is in canonical form, the co-basis $\{\tilde{e}^j\}$ will also be orthonormal. We can then conclude that at a point P of a manifold where a metric tensor is given, vectors and 1-forms are indeed equivalent objects. For example, if the metric is Euclidean, we have $g_{ij} = \delta_{ij}$ and $v^i = v_i$.

The operation of raising or lowering indices is naturally generalised to tensors of any order. For example,

$$T^{ij} g_{jk} = T^i_k , \quad (2.2.14)$$

actually represents a map between (the components of) a $(2, 0)$ tensor to (the components of) a $(1, 1)$ tensor. Likewise,

$$T^{ij} g_{ij} = T^{ij} g_{jk} g^{kl} g_{li} = T^i_k \delta_i^k = T , \quad (2.2.15)$$

is a map between (the components of) a $(2, 0)$ tensor to scalars (the trace), and can be generalised to tensors of any order (n, m) to produce tensors of order $(n-2, m)$ or $(n, m-2)$ and also $(n-1, m-1)$.

⁹Note that, although it is always possible to diagonalize a symmetric matrix, it might not be possible to diagonalize a metric tensor field simultaneously at different points, since the required matrices Λ may not satisfy the condition (2.1.91). More on this in the following.

2.2.2 Metric tensor field

A metric tensor field is an application which maps each point of a manifold \mathcal{M} into a metric tensor $g = g(P)$. A manifold in which a metric tensor (field) is defined everywhere is called a *metric manifold*.

Locally flat metric

Assuming regularity of the metric tensor field, the components of g in a given frame can be expanded around a point P in a Taylor series of the coordinate displacements,

$$g_{ij}(x) = g_{ij}(x_P) + \left. \frac{\partial g_{ij}}{\partial x^k} \right|_{x=x_P} \delta x^k + \frac{1}{2} \left. \frac{\partial^2 g_{ij}}{\partial x^k \partial x^l} \right|_{x=x_P} \delta x^k \delta x^l + \dots, \quad (2.2.16)$$

where $x = x_P + \delta x$. We can then transform $g'_{ij} = \Lambda^T g_{ij} \Lambda$, so that the metric takes the canonical form in P , $g'_{ij}(x_P) = \delta_{ij}$. Moreover, by extending the same transformation Λ in a neighbourhood of P , that is with $\Lambda = \Lambda(x)$ which satisfies the condition (2.1.91), we can also obtain ¹⁰

$$\left. \frac{\partial g'_{ij}}{\partial x^k} \right|_{x=x_P} = 0. \quad (2.2.17)$$

The conclusion is thus that it is always possible, by a change of coordinates, to write a metric tensor field in the form

$$g'_{ij}(x) = \delta_{ij} + \frac{1}{2} \frac{\partial^2 g'_{ij}}{\partial x^k \partial x^l} \delta x^k \delta x^l + \dots, \quad (2.2.18)$$

around a given point P of coordinates $x(P)$. Equivalently, it is always possible to choose locally orthogonal coordinates at any given point P . In general, however, as we move away from P , the same coordinates will not be orthogonal, unless the manifold is \mathbb{R}^n or a subset of it: there exist no change of coordinates that can put a general metric tensor in canonical form everywhere on a manifold.

Length of a curve

We can finally define the concept of length of a path on a manifold by considering the integral curve of a vector field $\vec{v} = \frac{d}{d\lambda}$. We first define the (squared) length of an infinitesimal displacement along the vector field \vec{v} as

$$dl^2 = d\vec{x} \cdot d\vec{x} = (\vec{v} d\lambda) \cdot (\vec{v} d\lambda) = g(\vec{v} d\lambda, \vec{v} d\lambda) = g(\vec{v}, \vec{v}) d\lambda^2, \quad (2.2.19)$$

which is obviously a scalar quantity (since \vec{v} is a vector, $d\lambda$ a scalar and g a $(0, 2)$ tensor). Upon integrating along an integral curve γ of the vector field \vec{v} , we obtain the length of the

¹⁰Technically, Eq. (2.2.17) are first order partial differential equations for the components of $\Lambda(x)$, with Eq. (2.5.54) at P playing the role of (initial value) boundary conditions.

integral path between two points of parameters λ_1 and λ_2 ,

$$l(\lambda_1, \lambda_2) = \int_{\lambda_1}^{\lambda_2} \sqrt{g(\vec{v}, \vec{v})} d\lambda = \int_{\lambda_1}^{\lambda_2} \sqrt{g_{ij}(\lambda) v^i(\lambda) v^j(\lambda)} d\lambda . \quad (2.2.20)$$

The above expression can be made more explicit upon introducing coordinates $\phi = x^i$ that cover the region where the integral is performed, namely

$$l(\lambda_1, \lambda_2) = \int_{\lambda_1}^{\lambda_2} \sqrt{g_{ij}(\lambda) \frac{dx^i}{d\lambda} \frac{dx^j}{d\lambda}} d\lambda , \quad (2.2.21)$$

where we just used the definition of vector components $v^i = \frac{dx^i}{d\lambda}$.

2.3 Lie derivative and symmetry

Among the good tensor operation we have seen so far there is no derivative (except for the derivative of a function, which was used to define vectors and general tensors). An easy way to understand why, is to consider how the ordinary partial derivative of a vector field transforms under a change of coordinates $x^{a'} = x^{a'}(x^b)$, namely

$$\begin{aligned} \partial_{b'} T^{a'} &= \frac{\partial}{\partial x^{b'}} \left(\frac{\partial x^{a'}}{\partial x^b} T^b \right) = \frac{\partial x^c}{\partial x^{b'}} \frac{\partial}{\partial x^c} \left(\frac{\partial x^{a'}}{\partial x^b} T^b \right) \\ &= \frac{\partial x^{a'}}{\partial x^b} \frac{\partial x^c}{\partial x^{b'}} \partial_c T^b + \frac{\partial^2 x^{a'}}{\partial x^b \partial x^c} \frac{\partial x^c}{\partial x^{b'}} T^b , \end{aligned} \quad (2.3.1)$$

where we are implicitly decomposing tensor quantities in terms of coordinate basis vectors and the dual 1-form basis. Due to the presence of the second term in the last line, this quantity does not immediately look like a tensor (the way we defined tensors in Section 1.4).

The above argument, beside being inaccurate (as we shall see later on), does not clarify the real issue at stake here. Derivatives involve comparing quantities at different points, and also require a way to quantify the “difference” (the displacement) between those base points. For example, in the case of functions, we have seen from the onset that the derivative on a generic manifold requires a curve. Even if a curve is given, for other tensorial quantities, we then need a way to map such quantities between the tangent spaces at different points, an ingredient that is mathematically arbitrary. In particular, in this section we shall consider integral curves of a vector field as a flow for points on a manifold, thus implementing the active interpretation of (auto)diffeomorphisms of a manifold, which will turn out to represent Lie groups and allow us to introduce symmetry on differentiable manifolds [9].

2.3.1 Passive and active transformations

We have already seen that any function $f : \mathcal{M} \rightarrow \mathbb{R}$ is a scalar (see Fig. 2.17), including the coordinates in a chart $\phi = x^i$. The reason the definition (1.3.33) we gave in the covariant

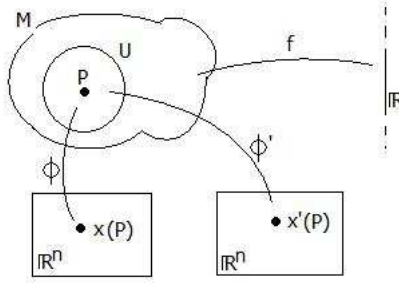


Figure 2.17: Scalars.

formalism of special relativity seems so much more involved can be seen by considering two charts $\phi = x$ and $\phi' = x'$ and their compositions Φ and Φ' with f , namely

$$(f \circ \phi^{-1}) = \Phi(x) = f(P) = \Phi'(x') = (f \circ \phi'^{-1}) \Rightarrow \Phi(x) = \Phi'(x') , \quad (2.3.2)$$

where x and x' represent the *same* point in *different* coordinate frames. We therefore see that it is not the function f that changes under coordinate transformations, but its composition with the chart ϕ . This describes the so-called *passive interpretation* of a *diffeomorphism* on a manifold \mathcal{M} : the points remain the same but their coordinates change. Technically, the diffeomorphism we are considering does not act on \mathcal{M} itself, but on the open subsets of \mathbb{R}^n that carry the coordinates for the manifold, that is,

$$\phi' = \phi'(\phi) , \quad (2.3.3)$$

whose existence is guaranteed by the very definition of differentiable manifold (see Section 2.1.1).

Like the composite functions Φ change (although one usually denotes them with the same symbol f), so do the components of any tensor field,

$$T'(x') = T(x) , \quad (2.3.4)$$

where T' is now short for the change of indices defined by $\Lambda = \frac{\partial x'}{\partial x}$ and its inverse Λ^{-1} . For example, vector field components change according to

$$V^{i'}(x') = \Lambda^{i'}_j(x(x')) V^j(x(x')) , \quad (2.3.5)$$

where the notation is meant to highlight that, in the new coordinate system $x' = x'(x)$, the new components $V^{i'}$ with respect to the new coordinate basis $\vec{e}_{i'} = \frac{\partial}{\partial x^{i'}}$ are linear combinations of the old components V^i in the old coordinate basis $\vec{e}_i = \frac{\partial}{\partial x^i}$, and the arguments of these functions must explicitly depend on the new coordinates x' .

Alternatively, one can consider the *active interpretation* of diffeomorphisms, which do act on the manifold so that points are actually moved along a flow $\psi : P \rightarrow P'$ (mathematically, this is an *automorphism* of \mathcal{M} into itself), and any function f is then dragged along. We can then define a new function f^* , called “pushed-forward” (or *Lie dragged*) of f , defined by

$$f^*(P') = f^*(\psi(P)) = f(P) . \quad (2.3.6)$$

Assuming points are not moved too far, one can use the same chart $\phi = x$ to cover both $U \subseteq \mathcal{M}$ and its image $U' = \psi(U)$ and, upon composing both f and f^* with the same chart ϕ , Eq. (2.3.6) implies

$$\Psi'(x) = (f^* \circ \psi \circ \phi^{-1})(x) = (f \circ \psi^{-1} \circ \phi^{-1})(x') = \Psi(x') \quad \Rightarrow \quad \Psi'(x) = \Psi(x'), \quad (2.3.7)$$

where Ψ and Ψ' clearly represent different composite functions with respect to the Φ and Φ' used in the previous paragraph, since x and x' here represent *different* points in the *same* coordinate frame ¹¹.

In both cases, active and passive mappings, a question remarkably relevant for physics remains (for now) unanswered: if coordinates are the only quantities that identify points, and we can change them freely, what is the physical meaning of points themselves? In other words, how can we say that two points P and Q on the same manifold are really different geometrical locations and distinguish active from passive transformations? A hint comes from considering coordinates as scalars: if we could give an operationally invariant meaning to the measurement of positions, points would be clearly identified by scalar quantities that do not change when we drag them around or change coordinates in the mathematical sense. In fact, only very selected coordinates can be given a physical meaning in this sense, whereas most charts will remain a mathematically useful, but otherwise formal, tool.

2.3.2 Congruences and Lie dragging

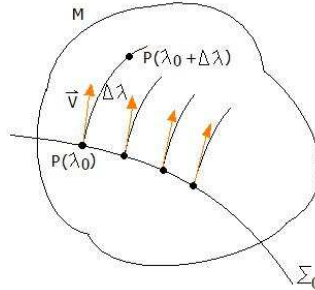


Figure 2.18: Congruence of a vector field.

Let us consider a two-dimensional manifold \mathcal{M} , a curve Σ_0 and a vector field $\vec{V} = \frac{d}{d\lambda} \in C^1(\mathcal{M})$. We call a *congruence* of the vector field \vec{V} the family of integral curves of \vec{V} which start from the curve Σ_0 (along which $\lambda = \lambda_0$) and cover (at least an open set U of) \mathcal{M} (see Fig. 2.18). By covering, we here mean that there is one (and only one) integral curve of \vec{V} across each point of U . Moving a point $P(\lambda_0)$ from Σ_0 along the corresponding congruence to the point $P(\lambda_0 + \Delta\lambda)$ is called “push forward” or “Lie dragging”. This operation can be

¹¹It is unfortunate that too many text-books do not distinguish these two compositions and the inherently different geometrical meanings.

straightforwardly generalized for any starting point $P(\lambda)$ and represents a continuous and invertible map of the manifold \mathcal{M} into itself,

$$\phi_{\Delta\lambda} : \mathcal{M} \rightarrow \mathcal{M} . \quad (2.3.8)$$

Clearly, this map also transforms Σ_0 in a new curve $\Sigma_{\Delta\lambda}$. If the vector field $\vec{V} \in C^\infty$, the maps $\phi_{\Delta\lambda}$, with $\Delta\lambda \in \mathbb{R}$, become diffeomorphisms (recall the active interpretation of changes of coordinates) and form a Lie group with respect to the usual composition law,

$$\phi_{\lambda_1} \circ \phi_{\lambda_2} = \phi_{\lambda_1 + \lambda_2} , \quad \phi_{\lambda}^{-1} = \phi_{-\lambda} , \quad \phi_{\lambda=0} = \mathbb{I} . \quad (2.3.9)$$

Of course, if \mathcal{M} is n -dimensional, we need an hypersurface Σ_0 of dimension $n - 1$ to define the $(n - 1)$ -dimensional congruence of integral curves of \vec{V} , which define the Lie group of the dragging ϕ_λ , where $\lambda \in \mathbb{R}$. We shall however leave this rather trivial generalisation for later developments.

Lie dragging functions

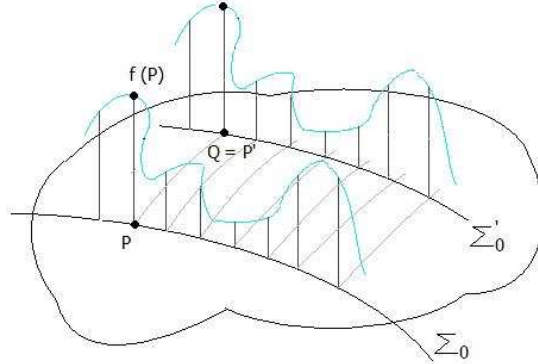


Figure 2.19: Lie-dragged function f^* .

Once we have defined how to drag a point, we can also define how to drag a function $f : \mathcal{M} \rightarrow \mathbb{R}$, by introducing the Lie dragged (or “pushed forward”) function f^* (see Fig. 2.19),

$$f_{\Delta\lambda}^*(Q) = f(P) , \quad \text{with} \quad \phi_{\Delta\lambda}(P) = Q . \quad (2.3.10)$$

In other words, f^* takes the same value at the Lie-dragged point Q the original function f takes at the point P . If Eq. (2.3.10) holds for all Q along the integral curve of $\vec{V} = \frac{d}{d\lambda}$ passing through P , it is clear that $f_{\Delta\lambda}^*$ as a function of $\Delta\lambda$ must be constant along such curve. Consequently, if $f_{\Delta\lambda}^*$ and f are the same for all values of $\Delta\lambda$, the function f must be constant along the lines of the congruence and $\frac{df}{d\lambda} = 0$.

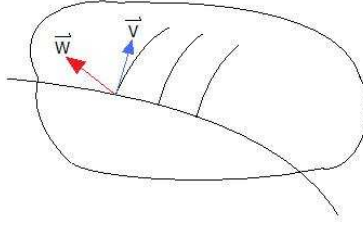


Figure 2.20: Congruence of \vec{V} and a second vector field \vec{W} .

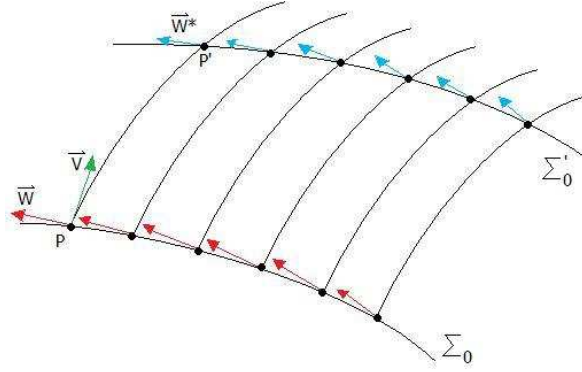


Figure 2.21: Lie-dragged vector field.

Lie dragging vector fields

We can next see how to Lie drag a vector field $\vec{W} = \frac{d}{d\mu}$ along a congruence of $\vec{V} = \frac{d}{d\lambda}$ (see Figs. 2.20-2.21). Quite naturally, we define the Lie dragged (or “pushed forward”) of \vec{W} by giving its action on an arbitrary function f , namely

$$\vec{W}_{\Delta\lambda}^*(f_{\Delta\lambda}^*)\Big|_Q = \vec{W}(f)\Big|_P, \quad \text{with } \phi_{\Delta\lambda}(P) = Q. \quad (2.3.11)$$

This is in analogy with the case of a function: the pushed forward vector \vec{W}^* applied to any Lie-dragged function f^* at the Lie-dragged point Q produces the same real number as the original vector \vec{W} applied to the original function f at the original point P . This can also be written as

$$\frac{df^*}{d\mu^*}\Big|_{\lambda_0+\Delta\lambda} = \frac{df}{d\mu}\Big|_{\lambda_0}, \quad (2.3.12)$$

where $\lambda(P) = \lambda_0$ and $\lambda(Q) = \lambda_0 + \Delta\lambda$.

As before, if one now views f as a fixed function and $\Delta\lambda$ as a variable in Eq. (2.3.12), it follows that $\frac{df^*}{d\mu^*}$ as a function of $\Delta\lambda$ is constant along congruences of $\vec{V} = \frac{d}{d\lambda}$. If \vec{W} and \vec{V} are linearly independent along an integral curve Σ_0 of $\vec{W} = \frac{d}{d\mu}$, we can use this Σ_0 as the initial curve (with $\lambda = \lambda_0$ and constant) for defining the congruence of \vec{V} . This initial curve

will then be mapped into a new curve $\Sigma = \Sigma_{\Delta\lambda}$ for each value of $\Delta\lambda$ ¹². Since $\lambda = \lambda_0$ along Σ_0 and $\lambda = \lambda_0 + \Delta\lambda$ along its image $\Sigma_{\Delta\lambda}$ are both constant values, λ can be naturally used as a coordinate. Further, the pushed forward vector field $\vec{W}^* = \frac{d}{d\mu^*}$ is, by definition, tangent to $\Sigma_{\Delta\lambda}$ and naturally defines a second parameter μ^* which (being the Lie-dragged function of μ on Σ_0) must be constant along the congruences of \vec{V} ¹³. This procedure therefore generates two coordinates (λ, μ^*) and the vector fields \vec{V} and \vec{W}^* , being coordinate vectors, must commute,

$$[\vec{V}, \vec{W}^*] = \left[\frac{d}{d\lambda}, \frac{d}{d\mu^*} \right] = 0, \quad (2.3.13)$$

for all points in a given open subset of \mathcal{M} (and, implicitly, for all functions in such a subset). This can indeed be taken as one of the defining equations for \vec{W}^* , as we shall see momentarily.

2.3.3 Lie derivatives

Lie derivatives are essentially derivatives along a congruence.

For functions

As before, let us consider the vector field $\vec{V} = \frac{d}{d\lambda} \in C^1$ and a function $f : \mathcal{M} \rightarrow \mathbb{R}$. We define the Lie derivative of the function f at the point $P = P(\lambda_0)$ as

$$\mathcal{L}_{\vec{V}} f|_{\lambda_0} = \lim_{\Delta\lambda \rightarrow 0} \frac{f_{-\Delta\lambda}^*(\lambda_0) - f(\lambda_0)}{\Delta\lambda} = \lim_{\Delta\lambda \rightarrow 0} \frac{f(\lambda_0 + \Delta\lambda) - f(\lambda_0)}{\Delta\lambda} = \left. \frac{df}{d\lambda} \right|_{\lambda_0} = \vec{V}(f), \quad (2.3.14)$$

in which we employed the “pull back” of f rather than the “push forward”. In other words, we used the flow generated by the congruence of \vec{V} to map $P(\lambda_0 + \Delta\lambda)$ to $P(\lambda_0)$,

$$\phi_{-\Delta\lambda}(P(\lambda_0 + \Delta\lambda)) = P(\lambda_0), \quad (2.3.15)$$

as well as the definition

$$f_{-\Delta\lambda}^*(\lambda_0) = f(\lambda_0 + \Delta\lambda). \quad (2.3.16)$$

Finally, note that this definition of Lie derivative naturally yields

$$\mathcal{L}_{\vec{V}} f = 0 \quad (2.3.17)$$

if f is constant along the congruence.

¹²Note that $\Sigma_{\Delta\lambda}$ is not in general an integral curve of \vec{W} for $\Delta\lambda \neq 0$. It will only be if $\vec{W} = \vec{W}_{\Delta\lambda}^*$ on Σ .

¹³One can easily see this by choosing the function $f = \mu$ along Σ_0 , so that $1 = \vec{W}(f)|_{\lambda_0} = \left. \frac{df^*}{d\mu^*} \right|_{\lambda_0 + \Delta\lambda}$.

For vector fields

Let us recall from Section 2.3.2 that the Lie dragged (push forward) of a vector field $\vec{W} = \frac{d}{d\mu}$ along congruences of $\vec{V} = \frac{d}{d\lambda}$ is characterized by the following property

$$\left. \frac{df^*}{d\mu^*} \right|_{\lambda_0 + \Delta\lambda} = \left. \frac{df}{d\mu} \right|_{\lambda_0}, \quad \forall \Delta\lambda \in \mathbb{R} \quad (2.3.18)$$

which implies

$$\left[\frac{d}{d\lambda}, \frac{d}{d\mu^*} \right] = 0, \quad (2.3.19)$$

Note that Eq. (2.3.18) applied to a given function f and for fixed $\Delta\lambda$ can be viewed as a proper set of n initial conditions for the unknown n -dimensional field $\vec{W}^* = \frac{d}{d\mu^*}$ satisfying the first order partial differential equation (2.3.19). In particular, for $\Delta\lambda = 0$ we have $f^* = f$ and Eq. (2.3.18) reads

$$\left. \frac{df}{d\mu^*} \right|_{\lambda_0} = \left. \frac{df}{d\mu} \right|_{\lambda_0}, \quad (2.3.20)$$

for a given f defined in a neighborhood of $P = P(\lambda_0)$. General theorems of calculus then guarantee that a \vec{W}^* solving Eq. (2.3.19) with initial condition (2.3.20) exists (at least) in a neighborhood of $P(\lambda_0)$.

Let us again consider the “pulled back” version of the above expressions by simply choosing as starting point $P = P(\lambda_0 + \Delta\lambda)$ instead of $P = P(\lambda_0)$, that is

$$\vec{W}_{-\Delta\lambda}^*(f_{-\Delta\lambda}^*) \Big|_{P(\lambda_0)} = \vec{W}(f) \Big|_{P(\lambda_0 + \Delta\lambda)}, \quad \text{with } \phi_{-\Delta\lambda}(P(\lambda_0 + \Delta\lambda)) = P(\lambda_0). \quad (2.3.21)$$

Upon acting on an arbitrary function f , the pulled back version of the initial condition (2.3.20) at the starting point $P(\lambda_0 + \Delta\lambda)$ reads

$$\left. \frac{df}{d\mu^*} \right|_{\lambda_0 + \Delta\lambda} = \left. \frac{df}{d\mu} \right|_{\lambda_0 + \Delta\lambda}, \quad (2.3.22)$$

and Eq. (2.3.19) can explicitly be rewritten as

$$\frac{d}{d\mu^*} \frac{d}{d\lambda} f = \frac{d}{d\lambda} \frac{d}{d\mu^*} f, \quad (2.3.23)$$

in the entire chosen neighborhood of $P(\lambda_0 + \Delta\lambda)$ [thus including $P(\lambda_0)$].

We can finally define the Lie derivative of a vector field \vec{W} as the limiting vector ¹⁴

$$\mathcal{L}_{\vec{V}} \vec{W}(f) \Big|_{\lambda_0} = \lim_{\Delta\lambda \rightarrow 0} \frac{\vec{W}_{-\Delta\lambda}^* - \vec{W}}{\Delta\lambda}(f) \Big|_{\lambda_0}. \quad (2.3.24)$$

¹⁴Note that both \vec{W} and \vec{W}^* act on the original function f (not on f^*).

In details, upon Taylor expanding around $\lambda_0 + \Delta\lambda$, we obtain

$$\vec{W}_{-\Delta\lambda}^*(f)\Big|_{\lambda_0} \equiv \frac{df}{d\mu^*}\Big|_{\lambda_0} = \frac{df}{d\mu^*}\Big|_{\lambda_0+\Delta\lambda} - \Delta\lambda \left(\frac{d}{d\lambda} \frac{d}{d\mu^*} f \right)_{\lambda_0+\Delta\lambda} + O(\Delta\lambda^2) . \quad (2.3.25)$$

On using the initial condition (2.3.22) and then expanding the first two terms in the right hand side around λ_0 , we next get

$$\begin{aligned} \frac{df}{d\mu^*}\Big|_{\lambda_0} &= \frac{df}{d\mu}\Big|_{\lambda_0+\Delta\lambda} - \Delta\lambda \left(\frac{d}{d\lambda} \frac{d}{d\mu^*} f \right)_{\lambda_0+\Delta\lambda} + O(\Delta\lambda^2) \\ &= \frac{df}{d\mu}\Big|_{\lambda_0} + \Delta\lambda \frac{d}{d\lambda} \frac{d}{d\mu} f \Big|_{\lambda_0} - \Delta\lambda \frac{d}{d\lambda} \frac{d}{d\mu^*} f \Big|_{\lambda_0} + O(\Delta\lambda^2) \\ &= \frac{df}{d\mu}\Big|_{\lambda_0} + \Delta\lambda \left(\frac{d}{d\lambda} \frac{d}{d\mu} f - \frac{d}{d\mu^*} \frac{d}{d\lambda} f \right)_{\lambda_0} + O(\Delta\lambda^2) , \end{aligned} \quad (2.3.26)$$

where we also took advantage of the commutator in Eq. (2.3.23). Since

$$\frac{dg}{d\mu^*} = \frac{dg}{d\mu} + O(\Delta\lambda) , \quad (2.3.27)$$

for any function g , we then have

$$\vec{W}_{-\Delta\lambda}^*(f)\Big|_{\lambda_0} = \vec{W}(f)\Big|_{\lambda_0} + \Delta\lambda \left[\frac{d}{d\lambda}, \frac{d}{d\mu} \right] f \Big|_{\lambda_0} + O(\Delta\lambda^2) , \quad (2.3.28)$$

from which

$$\begin{aligned} \mathcal{L}_{\vec{V}} \vec{W}(f)\Big|_{\lambda_0} &= \lim_{\Delta\lambda \rightarrow 0} \frac{\Delta\lambda \left[\frac{d}{d\lambda}, \frac{d}{d\mu} \right] f \Big|_{\lambda_0} + O(\Delta\lambda^2)}{\Delta\lambda} \\ &= \left[\frac{d}{d\lambda}, \frac{d}{d\mu} \right] f \Big|_{\lambda_0} , \end{aligned} \quad (2.3.29)$$

or, omitting the generic function f ,

$$\mathcal{L}_{\vec{V}} \vec{W} = \left[\vec{V}, \vec{W} \right] . \quad (2.3.30)$$

Note the above expression vanishes for a Lie dragged vector field $\vec{W} = \vec{W}^*$, as it reduces to the defining Eq. (2.3.19).

The Lie derivative has the following properties:

1. It vanishes if the components of \vec{W} are constant along the direction defined by \vec{V} . This can be easily seen by choosing a chart so that $\vec{V} = \frac{\partial}{\partial x^1}$, and then

$$\left(\mathcal{L}_{\vec{V}} \vec{W} \right)^i = V^j \frac{\partial}{\partial x^j} W^i - W^j \frac{\partial}{\partial x^j} V^i = \frac{\partial W^i}{\partial x^1} ; \quad (2.3.31)$$

2. It satisfies the Leibniz rule:

$$\mathcal{L}_{\vec{V}}(f \vec{W}) = (\mathcal{L}_{\vec{V}} f) \vec{W} + f \mathcal{L}_{\vec{V}} \vec{W} ; \quad (2.3.32)$$

3. It is linear:

$$\mathcal{L}_{\vec{V}} + \mathcal{L}_{\vec{W}} = \mathcal{L}_{\vec{V}+\vec{W}} ; \quad (2.3.33)$$

4. The commutator

$$[\mathcal{L}_{\vec{V}}, \mathcal{L}_{\vec{W}}] = \mathcal{L}_{[\vec{V}, \vec{W}]} , \quad (2.3.34)$$

so that:

5. It satisfies the Jacobi identity:

$$[[\mathcal{L}_{\vec{X}}, \mathcal{L}_{\vec{Y}}], \mathcal{L}_{\vec{Z}}] + [[\mathcal{L}_{\vec{Y}}, \mathcal{L}_{\vec{Z}}], \mathcal{L}_{\vec{X}}] + [[\mathcal{L}_{\vec{Z}}, \mathcal{L}_{\vec{X}}], \mathcal{L}_{\vec{Y}}] = 0 . \quad (2.3.35)$$

For one-forms and other tensors

Let us now consider a 1-form which maps any vectors at a point P into a real number. A 1-form field \tilde{w} applied to a vector field \vec{W} is then a map from (a subset of) the manifold \mathcal{M} into the real numbers, that is a scalar. Since we have already defined the Lie derivative of functions, we must consistently have

$$\mathcal{L}_{\vec{V}}(\tilde{w}(\vec{W})) = (\mathcal{L}_{\vec{V}} \tilde{w})(\vec{W}) + \tilde{w}(\mathcal{L}_{\vec{V}} \vec{W}) , \quad (2.3.36)$$

in which we assumed the Leibniz rule. From this, we can obtain the rather formal expression

$$(\mathcal{L}_{\vec{V}} \tilde{w})(\vec{W}) = \mathcal{L}_{\vec{V}}(\tilde{w}(\vec{W})) - \tilde{w}(\mathcal{L}_{\vec{V}} \vec{W}) . \quad (2.3.37)$$

If we now expand all vectors in a coordinate basis $\vec{e}_i = \partial/\partial x^i$ and the one-form in the dual basis, we can find the components of the Lie derivative one-form by applying it to basis vectors,

$$\begin{aligned} (\mathcal{L}_{\vec{V}} \tilde{w})_i &= (\mathcal{L}_{\vec{V}} \tilde{w})(\vec{e}_i) = \mathcal{L}_{\vec{V}}(\tilde{w}(\vec{e}_i)) - \tilde{w}(\mathcal{L}_{\vec{V}} \vec{e}_i) \\ &= \frac{d(\tilde{w}(\vec{e}_i))}{d\lambda} - \tilde{w}([\vec{V}, \vec{e}_i]) \\ &= V^k \frac{\partial w_i}{\partial x^k} - w_k \frac{\partial V^k}{\partial x^i} , \end{aligned} \quad (2.3.38)$$

and note that this expression reduces to

$$(\mathcal{L}_{\vec{V}} \tilde{w})_i = \partial_1 w_i \quad (2.3.39)$$

if we choose coordinates such that $\vec{V} = \partial_1$.

Upon noting that the full saturation of a type (n, m) tensor field is again a real function,

$$T_{(m)}^{(n)}(\tilde{w}_1, \tilde{w}_2 \dots \tilde{w}_n, \vec{W}^1, \vec{W}^2 \dots \vec{W}^m) : \mathcal{M} \rightarrow \mathbb{R} , \quad (2.3.40)$$

we can likewise obtain

$$\begin{aligned} \mathcal{L}_{\vec{V}} T(\tilde{w}_1, \tilde{w}_2 \dots \tilde{w}_n, \vec{W}^1, \vec{W}^2 \dots \vec{W}^m) = & (\mathcal{L}_{\vec{V}} T)(\tilde{w}_1, \tilde{w}_2 \dots \tilde{w}_n, \vec{W}^1, \vec{W}^2 \dots \vec{W}^m) \\ & + T(\mathcal{L}_{\vec{V}} \tilde{w}_1, \tilde{w}_2 \dots \tilde{w}_n, \vec{W}^1, \vec{W}^2 \dots \vec{W}^m) \\ & + T(\tilde{w}_1, \mathcal{L}_{\vec{V}} \tilde{w}_2 \dots \tilde{w}_n, \vec{W}^1, \vec{W}^2 \dots \vec{W}^m) + \dots \\ & + T(\tilde{w}_1, \tilde{w}_2 \dots \tilde{w}_n, \mathcal{L}_{\vec{V}} \vec{W}^1, \vec{W}^2 \dots \vec{W}^m) + \dots \\ & + T(\tilde{w}_1, \tilde{w}_2 \dots \tilde{w}_n, \vec{W}^1, \vec{W}^2 \dots \mathcal{L}_{\vec{V}} \vec{W}^m) . \end{aligned} \quad (2.3.41)$$

This is again a formal expression, which we can however simplify by a smart choice of coordinates.

Simple form of Lie derivatives

Let us choose $\lambda = x^1$ as one of the n coordinates, so that

$$\vec{V} = \frac{d}{d\lambda} = \frac{\partial}{\partial x^1} , \quad (2.3.42)$$

and review the different cases:

Scalars:

$$\mathcal{L}_{\vec{V}} f = \frac{df}{d\lambda} = \frac{\partial f}{\partial x^1} . \quad (2.3.43)$$

Vectors:

$$\left(\mathcal{L}_{\vec{V}} \vec{W} \right)^i = \left[\vec{V}, \vec{W} \right]^i = \frac{\partial W^i}{\partial x^1} . \quad (2.3.44)$$

Tensors:

$$\mathcal{L}_{\vec{V}} T = \lim_{\Delta\lambda \rightarrow 0} \frac{T_{-\Delta\lambda}^* - T}{\Delta\lambda} , \quad (2.3.45)$$

from which

$$(\mathcal{L}_{\vec{V}} T)^{i_1 i_2 \dots i_n}_{j_1 j_2 \dots j_m} = \frac{\partial T^{i_1 i_2 \dots i_n}_{j_1 j_2 \dots j_m}}{\partial x^1} . \quad (2.3.46)$$

To summarize, the Lie derivative is the *coordinate invariant* definition of partial derivatives. This result, incidentally, shows that the easy argument against partial derivatives being improper tensorial operations is inaccurate.

Example

Let us consider cartesian coordinates $\{x, y\}$ on the plane \mathbb{R}^2 and the vector field $\vec{V} \circ (x, y)$ given by

$$\vec{V}(x, y) = x^2 \frac{\partial}{\partial x} + \frac{\partial}{\partial y} , \quad (2.3.47)$$

which is well-defined in all of \mathbb{R}^2 . In order to compute any Lie derivative with respect to \vec{V} , we first define new coordinates $\{v = v(x, y), w = w(x, y)\}$, such that

$$\vec{V}(v, w) = \frac{\partial}{\partial v} . \quad (2.3.48)$$

The usual chain rule for partial derivatives then implies

$$\frac{\partial}{\partial v} = \left(\frac{\partial x}{\partial v} \right) \frac{\partial}{\partial x} + \left(\frac{\partial y}{\partial v} \right) \frac{\partial}{\partial y} \equiv x^2 \frac{\partial}{\partial x} + \frac{\partial}{\partial y} , \quad (2.3.49)$$

or

$$\begin{cases} \frac{1}{x^2} \frac{\partial x}{\partial v} = -\frac{\partial(x^{-1})}{\partial v} = 1 \\ \frac{\partial y}{\partial v} = 1 . \end{cases} \quad (2.3.50)$$

These equations are solved by

$$x(v, w) = \frac{1}{f(w) - v} \quad (2.3.51)$$

$$y(v, w) = v + g(w) , \quad (2.3.52)$$

where the functions $f = f(w)$ and $g = g(w)$ can be chosen freely, provided the transformation of coordinates is not singular [see Eq. (2.1.30)], that is

$$J \equiv \det \begin{bmatrix} \frac{\partial x}{\partial v} & \frac{\partial y}{\partial v} \\ \frac{\partial x}{\partial w} & \frac{\partial y}{\partial w} \end{bmatrix} = x^2 \left(\frac{dg}{dw} - \frac{df}{dw} \right) \neq 0 . \quad (2.3.53)$$

For example, we can set

$$f = 0 \quad \text{and} \quad g = w , \quad (2.3.54)$$

so that

$$\begin{cases} x = -\frac{1}{v} \\ y = v + w , \end{cases} , \quad \begin{cases} v = -\frac{1}{x} \\ w = y + \frac{1}{x} , \end{cases} \quad (2.3.55)$$

and the transformation matrices in the tangent space are given by

$$\Lambda^{-1} = \begin{bmatrix} \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial x} & \frac{\partial w}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{1}{x^2} & 0 \\ -\frac{1}{x^2} & 1 \end{bmatrix} \quad \text{and} \quad \Lambda = \begin{bmatrix} x^2 & 0 \\ 1 & 1 \end{bmatrix} . \quad (2.3.56)$$

In fact, upon comparing with Eq. (2.1.75), we have

$$\left(\frac{\partial}{\partial v}, \frac{\partial}{\partial w} \right) \begin{bmatrix} \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial x} & \frac{\partial w}{\partial y} \end{bmatrix} = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y} \right) , \quad (2.3.57)$$

which shows that Λ^{-1} maps the new coordinate basis vectors into the old ones, and will therefore map old components into the new components. Note that the new coordinates become singular for $x \rightarrow 0$ (where both v and w diverge, since the above mapping exchanges the origin in one frame with infinity in the other), whereas the Jacobian J vanishes for $x \rightarrow \infty$ (that is, the origin in $\{v, w\}$). As long as we avoid those two regions of \mathbb{R}^2 , the new coordinates v and w are fine.

Suppose we now want to compute the Lie derivative of the function $f \circ (x, y)$ given by

$$f(x, y) = x . \quad (2.3.58)$$

Direct application of the definition (2.3.14) for the vector (2.3.47) yields

$$\mathcal{L}_{\vec{V}} f = \vec{V}(f) = x^2 \frac{\partial f}{\partial x} + \frac{\partial f}{\partial y} = x^2 . \quad (2.3.59)$$

In the new coordinate system, we then have

$$f = -\frac{1}{v} , \quad (2.3.60)$$

and Eq. (2.3.43), with $x^1 = v$, gives

$$\mathcal{L}_{\vec{V}} f = \frac{\partial f}{\partial v} = \frac{1}{v^2} = x^2(v, w) . \quad (2.3.61)$$

Next we want to compute the Lie derivative of the vector field

$$\vec{W} = \frac{\partial}{\partial x} . \quad (2.3.62)$$

From the general formula (2.3.30), we immediately obtain

$$\mathcal{L}_{\vec{V}} \vec{W} = \left[x^2 \frac{\partial}{\partial x} + \frac{\partial}{\partial y}, \frac{\partial}{\partial x} \right] = -2x \frac{\partial}{\partial x} . \quad (2.3.63)$$

If we instead wish to apply the “simple” expression (2.3.44), we first need to express \vec{W} in the new coordinate system. Its components change with the matrix Λ^{-1} in Eq. (2.3.56) and the functional dependence on the coordinates according to Eq. (2.3.55),

$$\begin{bmatrix} \frac{1}{x^2} & 0 \\ -\frac{1}{x^2} & 1 \end{bmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{x^2} \\ -\frac{1}{x^2} \end{pmatrix}, \quad (2.3.64)$$

yielding

$$\vec{W} = v^2 \left(\frac{\partial}{\partial v} - \frac{\partial}{\partial w} \right). \quad (2.3.65)$$

Therefore,

$$(\mathcal{L}_{\vec{V}} W)^v = \frac{\partial(v^2)}{\partial v} = 2v = -\frac{2}{x}, \quad (2.3.66)$$

and

$$(\mathcal{L}_{\vec{V}} W)^w = -\frac{\partial(v^2)}{\partial v} = -2v = \frac{2}{x}, \quad (2.3.67)$$

or

$$\mathcal{L}_{\vec{V}} \vec{W} = 2v \left(\frac{\partial}{\partial v} - \frac{\partial}{\partial w} \right) = -2xv^2 \left(\frac{\partial}{\partial v} - \frac{\partial}{\partial w} \right) = -2x \frac{\partial}{\partial x}, \quad (2.3.68)$$

as it should.

The above example shows, among other things, that it is not always easier to use the simpler expressions of the Lie derivatives.

2.3.4 Symmetry and vector fields

We have seen how a vector field generates a “flow” on the manifold, and how the Lie derivative can be used to assess whether a given tensorial quantity remains unaffected by such a flow. It is therefore natural to associate the concept of “symmetry” to such flows generated by sets of vector fields. This is a deeply conceptual shift in perspective, in that the symmetry is no more a property of the manifold, but becomes a property of (tensorial) quantities defined on the manifold. Moreover, instead of relying on global coordinate transformations, the geometrical meaning of a symmetry will now reside on the local behaviour of relevant quantities under suitable displacements of the points. Such displacements can be further associated to “preferred foliations” of the manifold, and the latter then interpreted as “preferred observers”, thus coming to a closure with the old tensorial idea of symmetries as linear transformations.

Submanifolds and Lie algebras

We have just seen that one vector field on a manifold \mathcal{M} can generate congruences, that is a family of one-dimensional *submanifolds* of \mathcal{M} . Likewise, sets of vector fields can act as generators of submanifolds foliating a manifold.

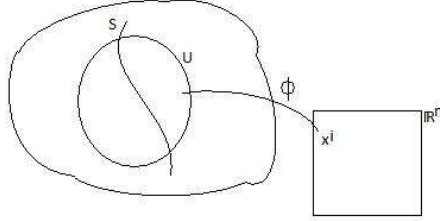


Figure 2.22: Submanifold.

Given a manifold \mathcal{M} of dimension n , one of its subsets \mathcal{S} is a submanifold of dimension $m \leq n$ if there exist charts U with coordinates $x \in \mathbb{R}^n$ such that $U \cap \mathcal{S} \subseteq \mathcal{M}$ and, for all points of \mathcal{S} (see Fig. 2.22),

$$x^1 = x^2 = \dots = x^{n-m} = 0, \quad m \leq n. \quad (2.3.69)$$

Given a point $P \in \mathcal{S}$, we can define the tangent space $T_P^{(\mathcal{S})}$ and, to each curve or vector on \mathcal{S} , we can associate a corresponding quantity in \mathcal{M} (see Fig. 2.23). We then have the

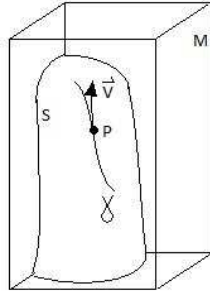


Figure 2.23: Embedding the tangent space.

following relations. First of all, we recall that

$$\dim(T_P^{(\mathcal{M})}) = n \geq \dim(T_P^{(\mathcal{S})}) = m. \quad (2.3.70)$$

A curve parameterized by $\lambda \in \mathbb{R}$ in \mathcal{S} naturally (and uniquely) maps to one in \mathcal{M} ,

$$\gamma_{\mathcal{S}} = (x^{n-m+1}(\lambda), \dots, x^n(\lambda)) \longleftrightarrow \gamma_{\mathcal{M}} = (0, 0, \dots, 0, x^{n-m+1}(\lambda), \dots, x^n(\lambda)), \quad (2.3.71)$$

and so do vectors,

$$\vec{V}_{\mathcal{S}} = (V^1, V^2, \dots, V^m) \longleftrightarrow \vec{V}_{\mathcal{M}} = (0, 0, \dots, 0, V^1, V^2, \dots, V^m). \quad (2.3.72)$$

However, the inverse maps are not uniquely defined: it is always possible to project a curve or vector from \mathcal{M} to \mathcal{S} , but the resulting curve or vector are the images of infinitely many curves and vectors. For example,

$$\vec{V}_{\mathcal{M}} = (0, \dots, 0, V^1, \dots, V^m) \longrightarrow \vec{V}_{\mathcal{S}} = (V^1, \dots, V^m) \longleftarrow \vec{V}'_{\mathcal{M}} = (1, \dots, 1, V^1, \dots, V^m) \quad (2.3.73)$$

We can likewise define the cotangent space $T_P^{*(\mathcal{S})}$, and a 1-form of \mathcal{S} will define a corresponding 1-form of \mathcal{M} , so that

$$\tilde{w}_{\mathcal{S}}(\vec{V}) = \tilde{w}_{\mathcal{M}}(0, \dots, 0, \vec{V}) \in \mathbb{R} , \quad (2.3.74)$$

where $\vec{V} \in T_P^{(\mathcal{S})}$. It is therefore clear that it is $\tilde{w}_{\mathcal{M}}$ which is not uniquely defined for a given $\tilde{w}_{\mathcal{S}}$, since, for example

$$\tilde{w}_{\mathcal{M}} = (0, \dots, 0, \vec{1}) \quad \text{and} \quad \tilde{w}_{\mathcal{M}} = (1, \dots, 1, \vec{1}) , \quad (2.3.75)$$

yield the same result for all $\vec{V} \in T_P^{(\mathcal{S})}$.

Since one vector field generates congruences (a foliation of \mathcal{M} in one-dimensional submanifolds), one could naively think $m \leq n$ vector fields define m -dimensional submanifolds. The general situation is rather different and stated by the very important

Frobenius theorem:

Given p linearly independent vector fields $\vec{V}^{(k)}$ ($k = 1, 2, \dots, p$) on the manifold \mathcal{M} , such that

$$[\vec{V}^{(i)}, \vec{V}^{(j)}] = c_{ij}^{ij} \vec{V}^{(k)} , \quad (2.3.76)$$

with c_{ij}^{ij} real constants, the integral curves of these fields form a family of submanifolds (or “foliation”) of \mathcal{M} , each of dimension $m \leq p$.

The meaning of the theorem is that a family of p vector fields could actually define a submanifold, but its dimension m is in general smaller than p . It is also important to stress that the linear independence of p vector fields means that there do not exist p constants a_i , $i = 1, \dots, p$, such that

$$\sum_{i=1}^p a_i \vec{V}^{(i)}(P) = 0 , \quad \forall P \in \mathcal{M} . \quad (2.3.77)$$

In particular, this does not mean that at a given point P the corresponding vectors are also linearly independent. In fact, if $p > n$, the manifold’s dimension, the above relation must hold at each P , but the coefficients then will depend on the point, that is $a_i = a_i(P)$.

Example (a): let us consider the manifold \mathbb{R}^3 and the vector fields $\vec{V}^{(1)} = \frac{\partial}{\partial x}$ and $\vec{V}^{(2)} = \frac{\partial}{\partial y}$. Integral curves of $\vec{V}^{(1)}$ are straight lines parallel to the x -axis, whereas integral curves of $\vec{V}^{(2)}$ are lines parallel to the y -axis. These two (obviously linearly independent) vector

fields together define a “foliation” of \mathbb{R}^3 by the planes of equation $z = \text{constant}$, which are 2-dimensional submanifolds \mathbb{R}^2 of \mathbb{R}^3 . Note that we have

$$[\vec{V}^{(1)}, \vec{V}^{(2)}] = 0 , \quad (2.3.78)$$

which implies that the two vectors $\vec{V}^{(i)}$ are actually a coordinate basis at all points in \mathbb{R}^3 (and \mathbb{R}^2).

Example (b): let us now consider the sphere $S \subseteq \mathbb{R}^3$. By introducing spherical coordinates with $\phi_z \in (0, 2\pi)$ the angle around the z -axis, one can immediately see that the vector field $\vec{\ell}_z = \frac{d}{d\phi_z} = \frac{\partial}{\partial \phi_z}$ generates “circles” of constant radius on the xy -planes. Likewise, by choosing the axis x and y we can also define the analogue vector fields $\vec{\ell}_x = \frac{d}{d\phi_x}$ and $\vec{\ell}_y = \frac{d}{d\phi_y}$. In any frame, the three vector fields $\vec{\ell}_z = \frac{d}{d\phi_z}$, $\vec{\ell}_x = \frac{d}{d\phi_x}$ and $\vec{\ell}_y = \frac{d}{d\phi_y}$ generate spheres of constant radius, which are 2-dimensional submanifolds of \mathbb{R}^3 . Note that this time

$$[\vec{\ell}_{(i)}, \vec{\ell}_{(j)}] \neq 0 , \quad (2.3.79)$$

and the three vector fields $\vec{\ell}_{(i)}$ do not define proper coordinates in \mathbb{R}^3 . Moreover, it is obvious that, at each P , there must exist real coefficients a_{ij} such that

$$\vec{\ell}_{(i)} = a_{i1} \frac{\partial}{\partial x} + a_{i2} \frac{\partial}{\partial y} + a_{i3} \frac{\partial}{\partial z} , \quad (2.3.80)$$

but these coefficients differ at different points, so that the $\vec{\ell}_{(i)}$ are independent vector fields.

Invariances and Lie algebras

Let us consider a tensor field T of type (p, q) on a manifold \mathcal{M} . A vector field \vec{V} is an “invariance”, or *symmetry*, of T if

$$\mathcal{L}_{\vec{V}} T = 0 , \quad (2.3.81)$$

that is, T is constant along congruences of \vec{V} .

The next important result is that vector fields leaving a set of tensors invariant generate a Lie algebra.

Theorem:

If we have a set of (linearly independent) tensors $T^{(k)}$, $k = 1, \dots, q$, and a set of (linearly independent) vectors $\vec{V}^{(i)}$, $i = 1, \dots, p$, such that

$$\mathcal{L}_{\vec{V}^{(i)}} T^{(k)} = 0 , \quad (2.3.82)$$

for all the combinations of such tensors and vectors, the vectors $\vec{V}^{(i)}$ form a Lie algebra.

Since the vectors $\vec{V}^{(i)}$ then satisfy the conditions of Frobenius theorem, they define a submanifold of dimension $m \leq p$.

This theorem is actually rather intuitive. Suppose we consider two of the symmetries $\vec{V}^{(i)}$, say $\vec{V}^{(1)}$ and $\vec{V}^{(2)}$, and compose the respective exponential maps to define two paths starting from one point P . As we saw in section 2.1.5, the end-point P' obtained by moving first along $\vec{V}^{(1)}$ and then along $\vec{V}^{(2)}$ will in general differ from the end-point P'' obtained by moving along $\vec{V}^{(2)}$ first and $\vec{V}^{(1)}$ after. However, P' and P'' will be connected by the exponential map of the vector $[\vec{V}^{(1)}, \vec{V}^{(2)}]$. Clearly, any of the tensors $T^{(k)}$ must be conserved along the closed path $P P' P'' P$, which implies that $[\vec{V}^{(1)}, \vec{V}^{(2)}]$ must also be a symmetry.

Isometries

An invariance of the metric g is called an *isometry*, and the vector associated with it is named a *Killing vector*,

$$\mathcal{L}_{\vec{V}}g = 0 . \quad (2.3.83)$$

It is a particularly important case for physics, since the above relation implies that Lie dragging points along congruences of \vec{V} preserves lengths and angles.

Example (a): the Euclidean metric in \mathbb{R}^3 ,

$$g = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} , \quad (2.3.84)$$

admits the Killing vectors $\frac{\partial}{\partial x}$, $\frac{\partial}{\partial y}$, $\frac{\partial}{\partial z}$, and $\frac{d}{d\phi_x}$, $\frac{d}{d\phi_y}$, $\frac{d}{d\phi_z}$ (which respectively generate translations and rotations around cartesian axes), and one can prove there is in fact no other (linearly independent vector field).

Example (b): the Minkowski metric in \mathbb{R}^3 ,

$$g = \begin{bmatrix} -1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} , \quad (2.3.85)$$

admits the Killing vectors $\frac{\partial}{\partial x}$ and $\frac{\partial}{\partial y}$. These vectors define planes (of constant t) corresponding to a given inertial observer evolving in time.

The above example allows us to draw a connection between the present geometrical formalism and the old tensor notation of Special Relativity. In Special Relativity, we started from the existence of *preferred reference frames* (the inertial observers) $S = \{x^\mu\}$ related by the Lorentz transformations

$$x^{\mu'} = \Lambda^{\mu'}_{\nu} x^{\nu} , \quad (2.3.86)$$

which represent the $SO(3,1)$ symmetry (*isometry*) of the Minkowski metric η ,

$$\Lambda^T \eta \Lambda = \eta , \quad (2.3.87)$$

and introduced tensors as mathematical objects (whose components) transform properly under the action of $SO(3,1)$. In differential geometry, tensors come along with the very definition of a manifold, and do not depend on any choice of coordinates (observers). This however does not prevent us from describing a family of preferred observers, which can be identified with a space-time foliation generated by Killing vectors \vec{V} . In this perspective, Eq. (2.3.87) is therefore replaced by the Killing condition

$$\mathcal{L}_{\vec{V}} g = 0 , \quad (2.3.88)$$

which, as we have just seen, mathematically implies the vectors \vec{V} form a Lie algebra [like Lorentz transformations belong to $SO(3,1)$] and generate a space-time foliation. If such vectors and foliation are space-like, they naturally identify an observer, or physical reference frame, at fixed time (a set of rulers to measure positions). This argument shows that the symmetry structure of Special Relativity can (at least in principle) be incorporated in the framework of differential geometry by simply assuming $g = \eta$ everywhere.

2.4 Differential forms

So far, we have defined tensors and, by means of a metric, also lengths and angles. We are still missing two important geometrical quantities, namely volume and area [10].

2.4.1 P -forms

A p -form is simply a type $(0,p)$ antisymmetric tensor. Since the condition of being skew-symmetric is preserved under linear combinations, it is clear that p -forms of fixed p , form a vector space.

Of course, we have already seen 1-forms and their definition does not involve skew symmetry. Given a type $(0,2)$ tensor \tilde{w} , we can build a 2-form by the following linear combination

$$\tilde{w}_A(\vec{V}, \vec{W}) = \frac{1}{2!} \left[\tilde{w}(\vec{V}, \vec{W}) - \tilde{w}(\vec{W}, \vec{V}) \right] . \quad (2.4.1)$$

Likewise, starting from a type $(0,3)$ tensor \tilde{w} , we can define

$$(\tilde{w}_A)_{ijk} = \frac{1}{3!} (\tilde{w}_{ijk} + \tilde{w}_{jki} + \tilde{w}_{kij} - \tilde{w}_{ikj} - \tilde{w}_{kji} - \tilde{w}_{jik}) \equiv \tilde{w}_{[ijk]} , \quad (2.4.2)$$

which defines the meaning of bracketed indices.

The above notation allows us to introduce the general expression for a p -form, that is

$$(\tilde{w}_A)_{i_1 i_2 \dots i_p} = \frac{1}{p!} (w_{i_1 i_2 \dots i_p} + \text{permutations}) \equiv w_{[i_1 i_2 \dots i_p]} . \quad (2.4.3)$$

The p -forms are vectors, in the sense that a p -form \tilde{w}_A on a n -dimensional manifold \mathcal{M} can be decomposed on C_p^n independent components, where

$$C_p^n = \frac{n!}{p!(p-n)!} , \quad p \leq n , \quad \text{with} \quad \sum_p C_p^n = n^2 . \quad (2.4.4)$$

Clearly, for $p > n$, this construction fails and there are no p -forms of that type. Moreover, the vector space of n -forms is 1-dimensional. It is also easy to see that p -forms at a point $P \in \mathcal{M}$ form a vector subspace of the space of $(0, p)$ tensors $(T_P^*)^p$, and that a C_p^n -dimensional basis is given by

$$\tilde{w}_A = \frac{1}{p!} w_{i_1 i_2 \dots i_p} \tilde{e}^{i_1} \wedge \tilde{e}^{i_2} \wedge \dots \wedge \tilde{e}^{i_p} , \quad (2.4.5)$$

where the wedge \wedge stands for the skew symmetric outer product. For example, a general 2-form can be written as

$$\tilde{w}_A = \frac{1}{2!} w_{ij} \tilde{e}^i \wedge \tilde{e}^j = \frac{1}{2} w_{ij} (\tilde{e}^i \otimes \tilde{e}^j - \tilde{e}^j \otimes \tilde{e}^i) . \quad (2.4.6)$$

The wedge product can be used to compose a p -form with a q -form, and obtain

$$(p\text{-form}) \wedge (q\text{-form}) = (p+q)\text{-form} , \quad (2.4.7)$$

provided $p+q \leq n$. Moreover, upon applying a p -form to a vector $\vec{V} \in T_P$, we obtain

$$\begin{aligned} \tilde{p}(\vec{V}, \cdot, \dots, \cdot) &= \left(\frac{1}{p!} w_{i_1 i_2 \dots i_p} \tilde{e}^{i_1} \wedge \tilde{e}^{i_2} \wedge \dots \wedge \tilde{e}^{i_p} \right) (V^k \vec{e}_k) \\ &= \frac{1}{p!} (w_{i_1 i_2 \dots i_p} V^k \tilde{e}^{i_1}(\vec{e}_k) \otimes \tilde{e}^{i_2} \otimes \dots \otimes \tilde{e}^{i_p} + \text{permutations}) \\ &= \frac{1}{(p-1)!} V^k w_{k i_2 \dots i_p} \tilde{e}^{i_2} \wedge \tilde{e}^{i_3} \wedge \dots \wedge \tilde{e}^{i_p} , \end{aligned} \quad (2.4.8)$$

which is a $(p-1)$ -form, and we used the dual basis of 1-forms to obtain the final expression.

2.4.2 Area and volume

There is a reason we used the index A to denote a p -form above: they can be used to define the area of a (sub)manifold. For example, given two vectors \vec{v} and \vec{w} , we can naturally define the area of the parallelogram they identify as

$$A = \vec{v} \wedge \vec{w} = |\vec{v}| |\vec{w}| \sin \theta , \quad (2.4.9)$$

where θ is the angle between \vec{v} and \vec{w} , whose definition requires a metric or is “implicitly” given by the choice of coefficients in the 2-form in Eq. (2.4.9) itself. Note that the above A can be either positive or negative, which means it represents an *oriented* area, and satisfies

$$A(\vec{v}, \vec{w}) + A(\vec{v}, \vec{b}) = A(\vec{v}, \vec{w} + \vec{b}) , \quad (2.4.10)$$

and

$$A(\vec{v}, \vec{w}) = -A(\vec{w}, \vec{v}) . \quad (2.4.11)$$

Given the interpretation of A as the area of a parallelogram, the antisymmetric property appears now necessary to ensure that $A(\vec{v}, a\vec{v}) = aA(\vec{v}, \vec{v}) = 0$, for all \vec{v} and $a \in \mathbb{R}$, without imposing further restrictions on the vectors the area 2-form acts upon.

Given a manifold \mathcal{M} of dimension n , a polyhedron is defined by n linearly independent vectors (which, for “infinitesimal” polyhedra, we can view as belonging to the tangent space of the same point) and its volume is simply a real number. We could therefore associate the volume to a type $(0, n)$ tensor. However, if we do not wish to restrict the n vectors and still assure that the volume vanishes if (at least) two of them are linearly dependent (roughly speaking, “parallel”¹⁵), we can instead define the volume as a n -form. Let us denote these n vectors as $\Delta\vec{x}_{(k)}$, with $k = 1, 2, \dots, n$. Since they all belong to the same T_P , we can expand them on the same coordinate basis,

$$\Delta\vec{x}_k = dx_{(k)}^i \frac{\partial}{\partial x^i} , \quad (2.4.12)$$

where $dx_{(k)}^i$ are just real numbers. All possible n -forms are proportional to each other, and will be given by

$$\tilde{\omega} = f \tilde{e}^1 \wedge \tilde{e}^2 \dots \wedge \tilde{e}^n , \quad (2.4.13)$$

where $f \in \mathbb{R}$. We then define the volume of the “infinitesimal polyhedron” (or cell) as

$$\tilde{\omega}(\Delta\vec{x}_{(1)}, \Delta\vec{x}_{(2)}, \dots, \Delta\vec{x}_{(n)}) = f \tilde{e}^1(\Delta\vec{x}_{(1)}) \tilde{e}^2(\Delta\vec{x}_{(2)}) \dots \tilde{e}^n(\Delta\vec{x}_{(n)}) + \text{permutations} . \quad (2.4.14)$$

If we, in particular, choose the n sides of the polyhedron along coordinate vectors, $\Delta\vec{x}_{(k)} = dx_{(k)}^i \frac{\partial}{\partial x^i} \equiv dx^k \frac{\partial}{\partial x^k}$ (no sum over k in the last expression), and the dual 1-form basis $\tilde{e}^i = \tilde{dx}^i$, we finally obtain the standard result

$$\begin{aligned} \tilde{\omega}(\Delta\vec{x}_{(1)}, \Delta\vec{x}_{(2)}, \dots, \Delta\vec{x}_{(n)}) &= f dx_{(1)}^1 dx_{(2)}^2 \dots dx_{(n)}^n + 0 + 0 + \dots + 0 \\ &= f dx^1 dx^2 \dots dx^n \equiv dV . \end{aligned} \quad (2.4.15)$$

If the n -form is a field in the chart $(U \subseteq \mathcal{M}, \phi = x^i)$, we can define the volume of U as simply

$$V = \int_U \tilde{\omega} = \int_{\phi(U)} f dx^1 dx^2 \dots dx^n , \quad (2.4.16)$$

where now $f \equiv f \circ \phi^{-1} = f(x^i)$ for $P(x^i) \in U$. It is important to check this expression is actually a scalar. So let us consider a change of coordinates $\phi = x^i \rightarrow y^i = y^i(x^i) = \phi'$ in U and, for simplicity, assume $n = 2$. We then have

$$\int_U \tilde{\omega} = \int_{\phi(U)} f(x^1, x^2) dx^1 dx^2 = \int_{\phi'(U)} \left(\frac{\partial x^1}{\partial y^1} \frac{\partial x^2}{\partial y^2} - \frac{\partial x^1}{\partial y^2} \frac{\partial x^2}{\partial y^1} \right) f(y^1, y^2) dy^1 dy^2 , \quad (2.4.17)$$

¹⁵We need a metric to define parallelism.

or, for general dimension $n \geq 2$,

$$V = \int_U \tilde{\omega} = \int_{\phi'(U)} f(y) J(y) d^n y , \quad (2.4.18)$$

where J is the determinant of the Jacobian matrix $\frac{\partial x}{\partial y}$, which shows that V is indeed coordinate independent.

Consider now a submanifold \mathcal{S} of dimension $n - 1$, and define the “infinitesimal” *area* of the $(n - 1)$ -dimensional hypersurface around the point P by means of a $(n - 1)$ -form. One could in principle use any $(n - 1)$ -form, but we also wish to maintain compatibility with the volume previously defined. We therefore take the volume form $\tilde{\omega}$ and apply it to a vector $\vec{v} \in T_P^{(\mathcal{M})} \notin T_P^{(\mathcal{S})}$, which means that \vec{v} is not a linear combination of vectors of $T_P^{(\mathcal{S})}$ ¹⁶. According to the expression (2.4.8), this defines the $(n - 1)$ -form $\tilde{A} = \tilde{\omega}(\vec{v}, \cdot, \dots, \cdot)$, which we can now apply to $n - 1$ vectors $\vec{w}_{(k)} \in T_P^{(\mathcal{S})}$, and obtain the area of the “infinitesimal cell”

$$\begin{aligned} \tilde{\omega}(\vec{v}, \vec{w}_{(1)}, \vec{w}_{(2)}, \dots, \vec{w}_{(n-1)}) &= \tilde{A}(\vec{w}_{(1)}, \dots, \vec{w}_{(n-1)}) \\ &= \frac{1}{(n-1)!} v f \tilde{e}^1(\vec{w}_{(1)}) \wedge \tilde{e}^2(\vec{w}_{(1)}) \wedge \dots \wedge \tilde{e}^{n-1}(\vec{w}_{(n-1)}) \\ &= \frac{1}{(n-1)!} v f dx^1 dx^2 \dots dx^{n-1} \equiv dA , \end{aligned} \quad (2.4.19)$$

in which we assumed $\vec{v} = v \vec{e}_n$ and $w_{(k)}^i = \delta_k^i dx^k$. The area of a portion $\Sigma \subseteq \mathcal{S}$ is then given by the integral

$$A = \int_{\Sigma} \tilde{A} = \int_{\phi(\Sigma)} f v dx^1 dx^2 \dots dx^{n-1} . \quad (2.4.20)$$

Note that under a change of coordinates, the above quantity does not change, since

$$\tilde{A} \rightarrow J^{(n-1)} \tilde{A}' , \quad (2.4.21)$$

with $J^{(n-1)}$ the Jacobian determinant of the transformation restricted on the hypersurface \mathcal{S} . Of course this construction can be further extended to lower and lower dimensional submanifolds.

Area and volume from the metric

As we expect, volume and area elements can be made compatible with the metric.

Let us assume there is a metric tensor field g on the manifold \mathcal{M} of dimension n , and that g is given the canonical form at the point P ,

$$g_{ij}(P) = \pm \delta_{ij} . \quad (2.4.22)$$

¹⁶One can naively think of \vec{v} as orthogonal to \mathcal{S} , although the notion of orthogonality again requires a metric, which we do not have in general at our disposal.

We recall that this implies both $\vec{e}_i = \frac{\partial}{\partial x^i}$ and the dual basis \tilde{e}^i , with $\tilde{e}^i(\vec{e}_j) = \delta_j^i$, are orthonormal with respect to g . The “natural volume” n -form is the one with $f = 1$ in this particular coordinate basis,

$$\tilde{\omega}_g = \tilde{e}^1 \wedge \dots \wedge \tilde{e}^n . \quad (2.4.23)$$

By a local change of coordinates around P ,

$$x^i \rightarrow y^i = y^i(x^i) , \quad (2.4.24)$$

we obtain that $\tilde{\omega}_g \circ \phi^{-1}$ transforms according to Eq. (2.4.18), that is ¹⁷

$$\tilde{\omega}_g = J \tilde{\omega}'_g = J \tilde{\sigma}^1 \wedge \dots \wedge \tilde{\sigma}^n , \quad (2.4.25)$$

where $\tilde{\sigma}_i = \frac{\partial}{\partial y^i}$, and $\tilde{\sigma}^i(\tilde{\sigma}_j) = \delta_j^i$. Now, observe that the determinant of a canonical metric is ± 1 , and from the transformation law

$$g' = \Lambda^T g \Lambda , \quad \text{with} \quad \Lambda^i_k = \frac{\partial y^i}{\partial x^k} , \quad (2.4.26)$$

we obtain

$$\det(g') = \det(\Lambda^T g \Lambda) = \det(g) \det(\Lambda^T \Lambda) = \det(g) \det^2(\Lambda) = \det(g) J^2 = \pm J^2 , \quad (2.4.27)$$

from which

$$J = \sqrt{|\det(g')|} . \quad (2.4.28)$$

We can finally write the volume of any subset $U \subseteq \mathcal{M}$ as

$$V = \int_U \tilde{\omega}_g = \int_{\phi(U)} \sqrt{|\det(g')|} dy^1 dy^2 \dots dy^n , \quad (2.4.29)$$

where $\phi = y^i$ is now a generic chart for U .

By employing the same n -form and repeating the argument which led us to define the area of an hypersurface Σ not containing \vec{e}_n , we get

$$\tilde{A}_g = \tilde{e}^1 \wedge \tilde{e}^2 \wedge \dots \wedge \tilde{e}^{n-1} , \quad (2.4.30)$$

where we can now say that \vec{e}_n is orthogonal to Σ , and

$$A = \int_{\Sigma} \tilde{A}_g = \int_{\phi(\Sigma)} \sqrt{|\det(g^{(n-1)})|} dx^1 dx^2 \dots dx^{n-1} , \quad (2.4.31)$$

where, for simplicity, we assumed the metric locally takes the form

$$g_{ij} = \begin{bmatrix} g_{ij}^{(n-1)} & 0 \\ 0 & \pm 1 \end{bmatrix} . \quad (2.4.32)$$

The overall conclusion is that we can use the metric to measure the length of a curve as well as the volume of any open sets of (sub)manifolds.

¹⁷Of course, the tensor $\tilde{\omega}_g$ does not change under a change of basis in T_P .

2.5 Covariant derivatives

On a manifold without the notion of angles (that is, without a metric), the only definition of parallelism can be given at a point P : two vectors of T_P are parallel if they are linearly dependent. But one then needs a way to confront vectors belonging to the tangent spaces at different points. One is in fact free to define this concept irrespectively of the metric. In particular, one can define how to transport a vector “parallelly” along a given path.

Let us consider again the example of the sphere S embedded in \mathbb{R}^3 . Being the latter a Euclidean space, there is a “natural” notion of parallel transport: a vector is parallelly transported if its angles with cartesian coordinate vectors remain constant. Consequently, a vector transported along a closed path returns into itself. From this notion of parallelism in \mathbb{R}^3 , we can induce a parallel transport on vectors on S . However, by transporting vectors along loops, we now find they do not return into themselves, in general [11].

2.5.1 Parallelism and covariant derivative

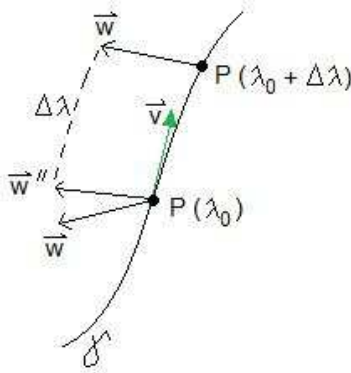


Figure 2.24: Parallel transport.

Let us first assume we have been given a rule to parallelly transport a vector \vec{W} along any curve γ tangent to a vector field \vec{V} on a manifold \mathcal{M} . This operation will associate to $\vec{W} \in T_{P(\lambda)}$ a second vector $\vec{W}''_{-\Delta\lambda} \in T_{P(\lambda_0)}$, where $\lambda = \lambda_0 + \Delta\lambda$ identifies a displaced point on the curve (see Fig. 2.24).

We then define the covariant derivative of the vector field \vec{W} with respect to \vec{V} at the point $P(\lambda_0)$ as the vector given by the limiting process

$$\nabla_{\vec{V}} \vec{W} \Big|_{\lambda_0} = \lim_{\Delta\lambda \rightarrow 0} \frac{\vec{W}''_{-\Delta\lambda}(\lambda_0) - \vec{W}(\lambda_0)}{\Delta\lambda}, \quad (2.5.1)$$

whose result is a vector, by definition, and vanishes if the parallelly transported vector coincides with the original vector in P . Note that, like for the Lie derivative, we are here transporting back the vector \vec{W} to the point $P(\lambda_0)$ from $P(\lambda_0 + \Delta\lambda)$. However, unlike the Lie derivative, we do not here need a whole congruence but just one curve.

Since functions do not identify a direction, it is natural to define the covariant derivative of a scalar to coincide with the Lie derivative, and thus with the total derivative

$$\nabla_{\vec{V}} f = \frac{df}{d\lambda} . \quad (2.5.2)$$

For general vectors and tensors, without specifying the actual transportation rule, we can still require the covariant derivative satisfies some formal properties. First of all, we want the following “Leibnitz rules” hold:

$$\nabla_{\vec{V}}(f \vec{W}) = \frac{df}{d\lambda} \vec{W} + f \nabla_{\vec{V}} \vec{W} \quad (2.5.3)$$

$$\nabla_{\vec{V}}(\vec{A} \otimes \vec{B}) = \vec{A} \otimes (\nabla_{\vec{V}} \vec{B}) + (\nabla_{\vec{V}} \vec{A}) \otimes \vec{B} \quad (2.5.4)$$

$$\nabla_{\vec{V}} [\tilde{\omega}(\vec{A})] = (\nabla_{\vec{V}} \tilde{\omega}) \vec{A} + \tilde{\omega} (\nabla_{\vec{V}} \vec{A}) . \quad (2.5.5)$$

We also assume that a change of parameterisation of the curve $\gamma \rightarrow \gamma'$, that is $\lambda \rightarrow \mu = \mu(\lambda)$, does not affect the notion of parallelism. Let $\vec{V} = \frac{d}{d\lambda}$ and $\vec{V}' = \frac{d}{d\mu}$ be the tangent vectors to γ and γ' respectively,

$$\frac{d}{d\lambda} = \frac{d\mu}{d\lambda} \frac{d}{d\mu} \equiv h \frac{d}{d\mu} . \quad (2.5.6)$$

We then impose that

$$\nabla_{h\vec{V}} \vec{W} = h \nabla_{\vec{V}} \vec{W} , \quad (2.5.7)$$

for all smooth functions h , so that $\nabla_{\vec{V}} \vec{W} = 0$ implies $\nabla_{h\vec{V}} \vec{W} = 0$. Finally, we want that, at a given point P ,

$$(\nabla_{\vec{V}} \vec{A})_P + (\nabla_{\vec{W}} \vec{A})_P = (\nabla_{\vec{V}+\vec{W}} \vec{A})_P , \quad (2.5.8)$$

so that

$$\nabla_{f\vec{V}+g\vec{W}} = f \nabla_{\vec{V}} + g \nabla_{\vec{W}} . \quad (2.5.9)$$

It is customary to name “covariant derivative of the vector \vec{W} ” the formal operator associated with the above derivative acting on a given \vec{W} , but with no specific curve (and thus for all vectors \vec{V}). This object at a point P can be viewed as a type $(1, 1)$ tensor,

$$\nabla \vec{W} : \vec{V} \rightarrow \nabla_{\vec{V}} \vec{W} , \quad (2.5.10)$$

which associates to any \vec{V} the corresponding covariant derivative of \vec{W} .

Affine connection

The formal properties introduced above allows one to obtain the components of the covariant derivative of a vector field in terms of the so-called Christoffel symbols (or affine connection).

We start by expanding both $\vec{V}(\lambda_0)$ and the difference between $\vec{W}(\lambda_0)$ and $\vec{W}''(\lambda_0)$ on a basis of $T_{P(\lambda_0)}$,

$$\begin{aligned}\nabla_{\vec{V}}\vec{W} &= \nabla_{V^i\vec{e}_i}(W^j\vec{e}_j) \\ &= V^i\nabla_{\vec{e}_i}(W^j\vec{e}_j) \\ &= V^i[(\nabla_{\vec{e}_i}W^j)\vec{e}_j + W^j(\nabla_{\vec{e}_i}\vec{e}_j)] .\end{aligned}\tag{2.5.11}$$

The second term in brackets above is called the *affine connection* (or Christoffel symbols),

$$\nabla_{\vec{e}_i}\vec{e}_j = \Gamma_{ji}^k\vec{e}_k ,\tag{2.5.12}$$

and, for fixed i and j , is obviously a vector in $T_{P(\lambda_0)}$. However, unlike its index notation might lead to think, Γ is not a type $(1,2)$ tensor. In fact, consider a reference frame with coordinates $\{x^i\}$ and coordinate basis $\{\vec{e}^i = \frac{\partial}{\partial x^i}\}$. We then see that, under a change of coordinates $x^{i'} = \Lambda^{i'}_j x^j$, the affine connection transforms according to

$$\Gamma_{j'i'}^{k'} = \Lambda_k^{k'}\Lambda_{i'}^i\Lambda_j^j\Gamma_{ji}^k + \Lambda_k^{k'}\Lambda_{i'}^i(\partial_i\Lambda_{j'}^k) .\tag{2.5.13}$$

Once a Γ is given, the covariant derivative becomes

$$\begin{aligned}\nabla_{\vec{V}}\vec{W} &= V^i[(\nabla_{\vec{e}_i}W^j)\vec{e}_j + W^j\nabla_{\vec{e}_i}\vec{e}_j] \\ &= V^i\left[\left(\frac{\partial W^j}{\partial x^i}\right)\vec{e}_j + W^j\Gamma_{ji}^k\vec{e}_k\right] \\ &= V^i\left[\frac{\partial W^k}{\partial x^i} + W^j\Gamma_{ji}^k\right]\vec{e}_k ,\end{aligned}\tag{2.5.14}$$

from which we can read out the components

$$\left(\nabla_{\vec{V}}\vec{W}\right)^k = V^i\frac{\partial W^k}{\partial x^i} + \Gamma_{ji}^k V^i W^j ,\tag{2.5.15}$$

and we recognise the first term on the right and side is the usual derivative along \vec{V} ,

$$V^i\frac{\partial W^k}{\partial x^i} = \frac{dW^k}{d\lambda} .\tag{2.5.16}$$

Since the vector \vec{V} enters only multiplicatively (that is, by contraction), it is customary (although quite improperly) to also call covariant derivative of a vector \vec{W} the type $(1,1)$ tensor (2.5.10), whose components are now given by

$$(\nabla W^k)_i = \frac{\partial W^k}{\partial x^i} + \Gamma_{ji}^k W^j .\tag{2.5.17}$$

Several different notations are in use for these components, for example

$$\nabla_i W^k = W^k_{;i} = W^k_{,i} + \Gamma_{ji}^k W^j .\tag{2.5.18}$$

Higher order tensors

The covariant derivative of tensors of any type (p, q) can be obtained in analogy with the procedure we employed for the Lie derivative, by simply starting from Eq. (2.5.2) and the formal properties of the covariant derivative given in Eqs. (2.5.3)-(2.5.8).

One can then obtain, for example, the components of the covariant derivative of a 1-form starting from the covariant derivative of the contraction $\tilde{W}(\vec{V}) = W_i V^i$, which is a function,

$$\partial_i (W_k V^k) = \nabla_i (W_k V^k) = (\nabla_i W_k) V^k + W_k (\nabla_i V^k) . \quad (2.5.19)$$

From the properties of partial derivatives and Eq. (2.5.17), we obtain

$$(\partial_i W_k) V^k + W_k (\partial_i V^k) = (\nabla_i W_k) V^k + W_k (\partial_i V^k + \Gamma_{ij}^k V^j) , \quad (2.5.20)$$

and, finally,

$$\nabla_i W_k = \frac{\partial W_k}{\partial x^i} - \Gamma_{ki}^j W_j . \quad (2.5.21)$$

By the same procedure, covariant derivatives of higher rank tensors are obtained.

Symmetric connection

An affine connection is *symmetric* if

$$\Gamma_{ij}^k = \Gamma_{ji}^k , \quad (2.5.22)$$

which implies the remarkable relation with the Lie derivative

$$\nabla_{\vec{V}} \vec{W} - \nabla_{\vec{W}} \vec{V} = [\vec{V}, \vec{W}] = \mathcal{L}_{\vec{V}} \vec{W} . \quad (2.5.23)$$

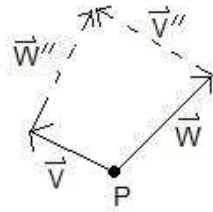


Figure 2.25: Absence of torsion.

The property of symmetry then straightforwardly implies that two linearly independent vectors \vec{V} and \vec{W} , defined at the same point $P \in \mathcal{M}$, and their parallelly transported versions (\vec{W}'' along \vec{V} and \vec{V}'' along \vec{W}) form a loop (see Fig. 2.25). In fact, since $\nabla_{\vec{V}} \vec{W}'' = \nabla_{\vec{W}} \vec{V}'' = 0$, it follows that $\mathcal{L}_{\vec{V}} \vec{W}'' = 0$, for “sufficiently small” vectors, and the parallelly transported

vectors define a reference frame. Therefore, moving P along \vec{V} first and then \vec{W}'' produces the same point as moving along \vec{W} first and then \vec{V}'' . If the connection is not symmetric,

$$T_{ji}^k = \Gamma_{ij}^k - \Gamma_{ji}^k \neq 0 , \quad (2.5.24)$$

the two paths yield in general two different images of P , and one says that parallelly transported vectors are subject to *torsion*.

2.5.2 Geodesics

A geodesic is a preferred curve along which the tangent vector to the curve itself is transported parallelly. This notion allows us to extend to a general manifold the concept of “straight line” and, eventually, of extremal curve (on metric manifolds).

Let $\vec{V} = \frac{d}{d\lambda}$ be the tangent vector to a curve γ parameterized by $\lambda \in \mathbb{R}$. Then, γ is a geodesic if \vec{V} satisfies

$$\nabla_{\vec{V}} \vec{V} \Big|_P = 0 , \quad \forall P \in \gamma , \quad (2.5.25)$$

and λ is then called an *affine parameter*¹⁸. From Eq. (2.5.7), it immediately follows that this definition is invariant under a change of parameterisation of γ (modulo singular points where the remapping fails), which implies that the same geodesic can be described by different affine parameters. Eq. (2.5.25) can be written in a local coordinate frame, in which $\gamma \in \mathcal{M}$ is mapped into $x^k = x^k(\lambda) \in \mathbb{R}^n$, as

$$\begin{aligned} \left(\nabla_{\vec{V}} \vec{V} \right)^k &= V^j \left(\frac{\partial V^k}{\partial x^j} + \Gamma_{ij}^k V^i \right) \\ &= \frac{dV^k}{d\lambda} + \Gamma_{ij}^k V^i V^j \\ &= \frac{d^2 x^k}{d\lambda^2} + \Gamma_{ij}^k \frac{dx^i}{d\lambda} \frac{dx^j}{d\lambda} = 0 , \end{aligned} \quad (2.5.26)$$

which is a set of n second-order differential equations for the variables $x^k = x^k(\lambda)$.

Note that, once the geodesic γ has been determined, we can use the affine parameter as a coordinate, say $x^1 = \lambda$, along that curve. The corresponding basis vector $\vec{e}_1 = \frac{\partial}{\partial x_1}$ will therefore be parallel transported along the geodesic by definition,

$$\nabla_{\vec{e}_1} \vec{e}_1 = 0 . \quad (2.5.27)$$

Moreover, we can also choose the remaining basis vectors \vec{e}_i , $i = 2, \dots, n$, so that they are also parallel transported along the geodesic, which eventually means we can introduce “adapted coordinates” around the geodesic, such that

$$0 = \nabla_{\vec{e}_i} \vec{e}_i = \Gamma_{1i}^k(P) , \quad \forall i, k = 1, \dots, n , \quad (2.5.28)$$

¹⁸The definition (2.5.25) is somewhat restrictive, as one could simply demand that the derivative of the tangent vector be parallel to the curve, or $\nabla_{\vec{V}} \vec{V} = \alpha \vec{V}$, with α a real function along the curve. This alternative definition is indeed more general, since there may be curves for which no parameterisation allows for $\alpha = 0$.

at any point $P \in \gamma$.

Normal frames

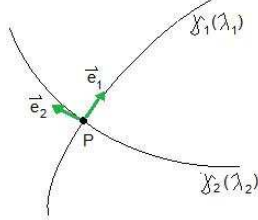


Figure 2.26: Geodesics and normal frame.

Given a point $P \in \mathcal{M}$ and *any* basis $\{\vec{e}_i^{(0)}\}$ of T_P , the n geodesic equations

$$\nabla_{\vec{e}_i} \vec{e}_i = 0 , \quad (2.5.29)$$

will admit a unique solution with $\vec{e}_i(P) = \vec{e}_i^{(0)}$ (see Fig. 2.26). We can then use any of these n geodesics to define n corresponding coordinates $\lambda_{(i)} = x^i$ having coordinate basis $\vec{e}_i = \frac{\partial}{\partial x^i}$. A very important result is that the affine connection totally vanishes at the origin P of any of these reference frames associated to geodesics through P . In fact, we obviously have that Eq. (2.5.28) must hold for all of the n directions,

$$\Gamma_{ij}^k|_P = 0 , \quad (2.5.30)$$

and the system is then called (*Gaussian*) *normal* around P .

It is now easier to see that the Γ 's indeed define how the coordinate basis vectors are parallel transported along the coordinate directions. In fact, from Eq. (2.5.30), we obtain that the corresponding coordinate basis at P satisfies

$$\nabla_{\vec{e}_i} \vec{e}_j|_P = 0 , \quad (2.5.31)$$

and Eq. (2.5.11) in this reference frame then becomes

$$\left(\nabla_{\vec{V}} \vec{W} \right)^j \Big|_P = V^i \frac{\partial W^j}{\partial x^i} \Big|_P , \quad (2.5.32)$$

so that the vector \vec{W} is parallelly transported along \vec{V} if its components do not change (along the direction of \vec{V}), which is the naive concept of parallel transport in \mathbb{R}^n . Note in fact that in such a reference frame

$$\nabla_{\vec{e}_i}|_P = \frac{\partial}{\partial x^i} \Big|_P = \mathcal{L}_{\vec{e}_i}|_P , \quad (2.5.33)$$

and the covariant derivatives coincide with the Lie derivatives along the coordinate vector fields (at P).

Of course, given any point $Q \neq P$, the above condition (2.5.30) will in general not hold, so that it is in general impossible to define a reference frame in which $\Gamma = 0$ in an arbitrary open set, or, equivalently, one in general has

$$\left. \frac{\partial \Gamma_{ij}^k}{\partial x^l} \right|_P \equiv \Gamma_{ij,l}^k|_P \neq 0 , \quad (2.5.34)$$

as well as higher order derivatives. It will however be possible to define a Gaussian normal frame (at least) around a given geodesic, so that $\Gamma(P) = 0$ for all P on the geodesic. This can be easily seen by simply repeating the above construction for all P of the geodesic.

Geodesic map

Another useful formula is the one which gives the parallelly transported vector starting from an initial vector \vec{A}_P at $P = \gamma(\lambda_0)$ along the curve of direction $\vec{V} = \frac{d}{d\lambda}$. Let $Q = \gamma(\lambda = \lambda_0 + \Delta\lambda)$ be a second point on the curve, then

$$\begin{aligned} \vec{A}(Q) &= \vec{A}_P + \Delta\lambda \nabla_{\vec{V}} \vec{A}_P + \frac{1}{2} \Delta\lambda^2 \nabla_{\vec{V}} \nabla_{\vec{V}} \vec{A}_P + \dots \\ &= e^{\Delta\lambda \nabla_{\vec{V}}} \vec{A}_P . \end{aligned} \quad (2.5.35)$$

It is immediate to understand the above expression if we introduce basis vectors along γ which are parallelly transported along \vec{V} , that is

$$0 = \nabla_{\vec{V}} \vec{e}_i|_{\gamma} = V^j \Gamma_{ji}^k \vec{e}_k|_{\gamma} . \quad (2.5.36)$$

This implies that the geodesic map becomes the exponential map along γ ,

$$\begin{aligned} A^i(Q) &= e^{\Delta\lambda V^j \partial_j} A_P^i \\ &= A_P^i + \Delta\lambda V^j \partial_j A_P^i + \mathcal{O}(\Delta\lambda^2) \\ &= A_P^i , \end{aligned} \quad (2.5.37)$$

where we used the compact notation $\vec{e}_i = \partial_i$ for the coordinate basis and the fact that partial derivatives of a vector defined at a point obviously vanish. The exponential map then defines a vector field $\vec{A}(Q) = \vec{A}(\lambda)$ by simply mapping the vector \vec{A}_P to all points Q , from λ to $\lambda_0 + \Delta\lambda$. Since this map does not affect the components, all the vectors $\vec{A}(\lambda)$ are indeed parallel to \vec{A}_P in the “naive” sense. In a generic reference frame, we must replace the partial derivative with the covariant derivative, with a connection Γ which will mix the coordinate basis vectors and the components of the vector according to our rule of parallel transport, namely

$$\vec{A}(Q) = A_P^i + \Delta\lambda V^j \Gamma_{jk}^i A_P^k + \mathcal{O}(\Delta\lambda^2) , \quad (2.5.38)$$

where we again used the fact that partial derivatives of a vector defined at a point vanish.

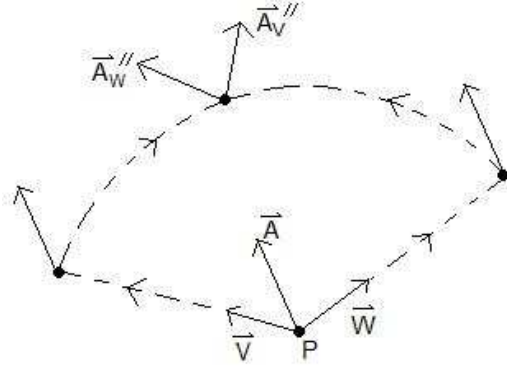


Figure 2.27: Commutator of two parallel transports.

2.5.3 Riemann tensor and curvature

The Riemann tensor is the mathematical quantity which allows one to define the *curvature* of a manifold as the effect of parallel transport of vector fields along loops.

Let us first consider two vector fields $\vec{V} = \frac{d}{d\lambda}$ and $\vec{W} = \frac{d}{d\mu}$, such that $[\vec{V}, \vec{W}] = 0$. They can therefore be used to introduce local coordinates, or, equivalently, to define a closed loop (see Fig. 2.27). Starting from a point $P \in \mathcal{M}$, we first parallelly move a third vector \vec{A} along \vec{V} , and then along \vec{W} , which yields a vector

$$\vec{A}''_{WV} = e^{\delta\mu \nabla_{\vec{W}}} e^{\delta\lambda \nabla_{\vec{V}}} \vec{A}. \quad (2.5.39)$$

We then repeat the same process reversing the order of \vec{V} and \vec{W} along which we transport \vec{A} and obtain

$$\vec{A}''_{VW} = e^{\delta\lambda \nabla_{\vec{V}}} e^{\delta\mu \nabla_{\vec{W}}} \vec{A}. \quad (2.5.40)$$

For infinitesimally small displacements $|\delta\lambda| \ll 1$ and $|\delta\mu| \ll 1$, respectively, along \vec{V} and \vec{W} , the difference between these two resulting vectors will be the vector

$$\delta\vec{A} = \vec{A}''_{WV} - \vec{A}''_{VW} = \delta\lambda \delta\mu [\nabla_{\vec{V}}, \nabla_{\vec{W}}] \vec{A} + O(3), \quad (2.5.41)$$

which further clarifies the meaning of the covariant derivative.

Given the above result, we now define the *Riemann tensor* as the type $(1, 3)$ tensor $R(\cdot, \cdot)$ which, given two directions \vec{V} and \vec{W} , produces a type $(1, 1)$ tensor $R(\vec{V}, \vec{W})$, whose *contraction*¹⁹ with a vector \vec{A} finally gives the vector

$$R(\vec{V}, \vec{W}) \vec{A} = [\nabla_{\vec{V}}, \nabla_{\vec{W}}] \vec{A} - \nabla_{[\vec{V}, \vec{W}]} \vec{A}, \quad (2.5.42)$$

where we specify that

$$R(\vec{V}, \vec{W}) \vec{A} \equiv \left[R(\vec{V}, \vec{W})^i{}_j A^j \right] \vec{e}_i, \quad (2.5.43)$$

¹⁹By contraction we here mean precisely the operation of contracting indices, so that the action of R on \vec{A} is a mere multiplication.

which allows us to write

$$\delta \vec{A} = \delta \lambda \delta \mu R(\vec{V}, \vec{W}) \vec{A} + O(3) , \quad (2.5.44)$$

or

$$\delta A^i = \delta \lambda \delta \mu R_{jkl}^i V^j W^k A^l + O(3) . \quad (2.5.45)$$

This yields the precise mathematical meaning of the concept of *intrinsic curvature* of a manifold: whenever the Riemann tensor does not vanish, parallelly transporting a vector along a closed path does not return the vector to its initial value. Conversely, if there exist loops such that vectors parallelly transported along them do not return into themselves, the manifold is curved. Note that this definition of curvature is “intrinsic” since it does not require embedding (viewing) the manifold \mathcal{M} into (from) a larger space. An equivalent definition of intrinsic curvature involves measuring the sum of the internal angles of a triangle, and was proposed by Gauss long ago as an experiment to measure the Earth’s curvature.

A simple example is again given by the sphere in \mathbb{R}^3 : one can of course define the “extrinsic curvature” radius R (and *extrinsic curvature* $1/R$) from the defining condition $x^2 + y^2 + z^2 = R^2$. However, the same conclusion can be drawn without referring to \mathbb{R}^3 at all, by simply noting that a vector transported along a loop starting from (say) the North pole, reaching the equator on a meridian, moving along the equator a distance $R\theta$, and coming back to the North pole along a meridian, will appear rotated of the angle θ . In this case, the intrinsic and extrinsic curvature radii coincide. However, in general, the two quantities may be different.

From the definition (2.5.43), it is easy to see that the Riemann tensor R has the following properties

$$R(\vec{V}, \vec{W})(f\vec{A}) = f R(\vec{V}, \vec{W}) \vec{A} \quad (2.5.46)$$

$$R(f\vec{V}, \vec{W}) \vec{A} = R(\vec{V}, f\vec{W}) \vec{A} = f R(\vec{V}, \vec{W}) \vec{A} , \quad (2.5.47)$$

and, at any point P , the Riemann tensor R can be written as

$$R_{ljk}^i \vec{e}^l \otimes \vec{e}_i = R(\vec{e}_j, \vec{e}_k)_l^i \vec{e}^l \otimes \vec{e}_i , \quad (2.5.48)$$

in the corresponding local basis.

Example: Parallel transport in \mathbb{R}^2

A neat example that can explain all the above formalism is given by the naive parallel transport of vectors in the plane \mathbb{R}^2 with global Cartesian coordinates $\{x, y\}$: a vector $\vec{A}(Q)$ is the parallel transported of \vec{A}_P if its angles with the coordinate axis are the same. Clearly, this occurs if

$$\vec{x}(A^i) = \vec{y}(A^i) = 0 , \quad (2.5.49)$$

where $\vec{x} = \partial_x$ and $\vec{y} = \partial_y$ are coordinate basis vectors in T_Q , for all $Q \in \mathbb{R}^2$. It then immediately follows that $\Gamma = 0$ in all of the plane, which is explicitly a flat manifold with this choice of parallel transport. It also follows that geodesics are curves at constant angle in this reference frame, that is straight lines, of which the coordinate axis are particular examples.

Imagine instead to chart \mathbb{R}^2 with polar coordinates $\{r, \theta\}$, with the same law of parallel transport. It is obvious that the new coordinate basis vectors $\vec{r} = \partial_r$ and $r\vec{\theta} = \partial_\theta$ are in general not parallel to \vec{x} and \vec{y} , and $\Gamma \neq 0$ in this frame.

It is finally interesting to note that the geodesics so defined in \mathbb{R}^2 are also the curves of minimum length between two points, say P and Q , and that the Euclidean metric is in the canonical form in the frame $\{x, y\}$, by definition. The metric is however not in canonical form in the coordinates $\{r, \theta\}$, and integral curves of $\vec{\theta}$ are in fact not geodesics. This immediately brings us to look deeper into the possible connection between parallel transport and the metric.

2.5.4 Metric connection

So far we have not specified any affine connections. But we are really interest in the case in which parallel transport preserves lengths and angles, which requires the manifold \mathcal{M} is endowed with a metric tensor g .

Let us then consider two vectors \vec{A} and \vec{B} , and assume they are transported parallelly along a curve of tangent \vec{V} , that is $\nabla_{\vec{V}}\vec{A} = \nabla_{\vec{V}}\vec{B} = 0$. It is natural to demand that the scalar product between these two vectors does not change along the curve,

$$\nabla_{\vec{V}} \left[g(\vec{A}, \vec{B}) \right] = 0, \quad \forall \vec{A}, \vec{B}, \vec{V} \quad \text{such that} \quad \nabla_{\vec{V}}\vec{A} = \nabla_{\vec{V}}\vec{B} = 0, \quad (2.5.50)$$

which, from the Leibniz rule, implies

$$\nabla_{\vec{V}} g = 0, \quad \forall \vec{V}, \quad (2.5.51)$$

or, more formally,

$$\nabla g = 0. \quad (2.5.52)$$

Upon expressing this equation in a specified coordinate frame, one finds that it is tantamount to an equation for the affine connection, namely

$$\Gamma_{ij}^k = \frac{1}{2} g^{kl} (g_{il,j} + g_{jl,i} - g_{ij,l}). \quad (2.5.53)$$

Since g is symmetric, one can immediately see that a metric connection is necessarily symmetric.

All expressions can be simplified by assuming the metric is in canonical form at a point P , so that it can be expanded as

$$g_{ij} = \pm \delta_{ij} + \frac{1}{2} \frac{\partial^2 g_{ij}}{\partial x^k \partial x^l} \bigg|_P \delta x^k \delta x^l + \dots \quad (2.5.54)$$

Eq. (2.5.53) above then implies that

$$g_{ij,k}|_P = 0 \quad \Rightarrow \quad \Gamma_{ij}^k|_P = 0 . \quad (2.5.55)$$

Starting from P , we can then consider n linearly independent directions and the corresponding geodesics will form a Gaussian normal reference frame around P (at least in a sufficiently small neighbourhood of P). As we showed in Section 2.5.2, in this particular reference frame, covariant derivatives along the coordinate directions are also Lie derivatives at P . The consistency condition (2.5.53) then implies that the coordinate basis vectors are also “Killing vectors at P ”, and the metric correspondingly admits (at least) n “point isometries at P ”. It is important to remark that, strictly speaking, Killing vectors are only defined as fields and the condition (2.5.55) should therefore hold in an open set of the manifold. Since we have seen that it is in general impossible to put the metric in canonical form in an open set, we are thus specifying “at P ” in order to stress this fact.

Once we have connected the parallel transport to the metric, we can also see that geodesics are indeed curves of *local extremal length*: suppose we take a specific geodesic γ of parameter λ and construct a Gaussian normal frame around it, like we mentioned before. In a (sufficiently small) neighbourhood of γ , the metric will take the form (2.5.54), so that moving off the geodesic from the point $P = P(\lambda_P)$ along each Gaussian direction η^i , with $i = 1, \dots, n-1$, one has

$$ds^2 \simeq \frac{1}{2} \frac{\partial^2 g_{ii}}{\partial (\eta^i)^2} \Big|_{\lambda=\lambda_P, \eta^1=\dots=\eta^{n-1}=0} (d\eta^i)^2 + \dots , \quad (2.5.56)$$

where there is no sum over the index i . Depending on the sign of the second derivative of the metric on the geodesic, this quantity will always be either positive or negative, in a sufficiently small portion of γ around P . One can therefore conclude that each portion of a geodesic is a local extremum for the length of a curve.

The above argument about geodesics can in fact be used in order to derive Eq. (2.5.53). Let us therefore consider the “length” of a curve between two fixed points A and B , namely²⁰

$$s = \int_A^B ds = \int_A^B \sqrt{g_{ij} \dot{x}^i \dot{x}^j} d\lambda \equiv \int_{\lambda_A}^{\lambda_B} \sqrt{2L(x^k, \dot{x}^l)} d\lambda , \quad (2.5.57)$$

where a dot denotes derivative with respect to the affine parameter λ and we note that $g_{ij} = g_{ij}(x^k)$. If we identify $\lambda = s$, we obviously have $2L = 1$ and varying the above action is equivalent to varying the action without the square root, that is

$$\delta s = \delta \int_{s_A}^{s_B} \sqrt{2L} ds = \int_{s_A}^{s_B} \frac{\delta L}{\sqrt{2L}} ds = \delta \int_{s_A}^{s_B} L(x^k, \dot{x}^l) ds . \quad (2.5.58)$$

By requiring $\delta s = 0$, we then find the Euler-Lagrange equations of motion

$$\frac{d}{ds} \left(\frac{\partial L}{\partial \dot{x}^m} \right) - \frac{\partial L}{\partial x^m} = 0 . \quad (2.5.59)$$

²⁰We assume for simplicity that the argument of the square root is positive. If it were not, one just needs to change its sign.

In particular, one finds

$$\frac{\partial L}{\partial \dot{x}^m} = \frac{1}{2} g_{jk,m} \dot{x}^j \dot{x}^k , \quad (2.5.60)$$

and

$$\frac{\partial L}{\partial \dot{x}^m} = g_{mj} \dot{x}^j , \quad (2.5.61)$$

from which

$$\begin{aligned} \frac{d}{ds} \left(\frac{\partial L}{\partial \dot{x}^m} \right) &= g_{mj} \ddot{x}^j + g_{mj,k} \dot{x}^j \dot{x}^k \\ &= g_{mj} \ddot{x}^j + \frac{1}{2} (g_{mj,k} + g_{jm,k}) \dot{x}^j \dot{x}^k . \end{aligned} \quad (2.5.62)$$

Putting now the two parts together and multiplying by g^{im} we obtain

$$\ddot{x}^i + \frac{1}{2} g^{il} (g_{lk,j} + g_{lj,k} - g_{jk,l}) \dot{x}^j \dot{x}^k = 0 , \quad (2.5.63)$$

which equals the geodesic equation

$$\ddot{x}^i + \Gamma_{jk}^i \dot{x}^j \dot{x}^k = 0 , \quad (2.5.64)$$

provided the Christoffel symbols are given by Eq. (2.5.53).

In a metric manifold, the Riemann tensor also describes what is usually referred to as the *geodesic deviation*. Suppose one considers two geodesics γ_1 and γ_2 , starting from P_1 and P_2 , and having parallel initial tangent vectors \vec{V} . This mean the tangent to γ_1 at the initial point P_1 is the parallely transported tangent vector to γ_2 at P_2 (and vice versa). We can then say both tangents are parametrised by the same λ . One can then consider the length of the curve between P_1 and P_2 and, in general, between the two points $\gamma_1(\lambda)$ and $\gamma_2(\lambda)$. If P_1 and P_2 are “infinitesimally” close to each other, this distance will be approximately equal to the length of a vector connecting $\gamma_1(\lambda)$ to $\gamma_2(\lambda)$. If the manifold is flat, this distance will remain constant and so will the length of the connecting vector \vec{Y} . Otherwise, the rate of change of this vector is given precisely by the Riemann tensor, and its modulus determines the rate of change of the distance between the two geodesics,

$$\ddot{Y}^i \simeq \frac{\delta V^i}{\delta \lambda^2} \simeq R_{jkl}^i V^j Y^k V^l , \quad (2.5.65)$$

where we used Eq. (2.5.41) with $\vec{A} = \vec{V}$, $\vec{Y} = \vec{W}$ and $\delta\mu = \delta\lambda$.

Since in a normal frame, covariant derivatives (at a point) can be replaced by partial derivatives (at the same point), one finds the Riemann tensor has components simply given by second derivatives of the metric (at a given point),

$$R_{ijkl} = \frac{1}{2} (g_{il,jk} - g_{ik,jl} + g_{jk,il} - g_{jl,ik}) , \quad (2.5.66)$$

and appears “block symmetric”,

$$R_{ijkl} = R_{klij} .$$

This property allows us to define the following tensors:

$$\text{Ricci (Curbastro) tensor: } R^k_{ikj} = R_{ij} , \quad \text{with } R_{ij} = R_{ji} \quad (2.5.67)$$

$$\text{Curvature scalar: } R = R^k_k \quad (2.5.68)$$

$$\text{Einstein tensor: } G_{ij} = R_{ij} - \frac{1}{2} R g_{ij} . \quad (2.5.69)$$

All expressions can then be generalised to any frames by simply replacing partial derivatives with covariant derivatives. For example, the general form of the Riemann tensor is given by ²¹

$$R_{ijkl} = \frac{1}{2} (g_{il;jk} - g_{ik;jl} + g_{jk;il} - g_{jl;ik}) , \quad (2.5.70)$$

which, once expanded explicitly, becomes rather involved.

In a normal frame, it is easy to check that the Einstein tensor satisfies the identity

$$\left(R^{ij} - \frac{1}{2} R g^{ij} \right)_{,i} = 0 . \quad (2.5.71)$$

Upon generalizing to any frames, we then obtain the important *Bianchi identity*

$$\nabla_i G^{ij} = 0 , \quad (2.5.72)$$

which resembles (and actually is) a conservation law, as we will elucidate further.

Another important expression we will make use of, is given by the Killing equation (2.3.83). We can easily obtain a more explicit expression by considering that $g(\vec{A}, \vec{B})$ is a scalar and therefore

$$\mathcal{L}_{\vec{V}} [g(\vec{A}, \vec{B})] = \nabla_{\vec{V}} [g(\vec{A}, \vec{B})] . \quad (2.5.73)$$

By applying Leibniz’s rule on both sides we obtain

$$[\mathcal{L}_{\vec{V}} g](\vec{A}, \vec{B}) + g(\mathcal{L}_{\vec{V}} \vec{A}, \vec{B}) + g(\vec{A}, \mathcal{L}_{\vec{V}} \vec{B}) = g(\nabla_{\vec{V}} \vec{A}, \vec{B}) + g(\vec{A}, \nabla_{\vec{V}} \vec{B}) , \quad (2.5.74)$$

which holds for a metric covariant derivative. In particular, we have

$$(\mathcal{L}_{\vec{V}} g)_{ij} \equiv [\mathcal{L}_{\vec{V}} g](\vec{e}_i, \vec{e}_j) = g(\nabla_{\vec{V}} \vec{e}_i, \vec{e}_j) + g(\vec{e}_i, \nabla_{\vec{V}} \vec{e}_j) - g(\mathcal{L}_{\vec{V}} \vec{e}_i, \vec{e}_j) - g(\vec{e}_i, \mathcal{L}_{\vec{V}} \vec{e}_j) . \quad (2.5.75)$$

On expanding $\vec{V} = V^k \vec{e}_k$, where $\vec{e}_i = \partial_i$, we see that

$$\mathcal{L}_{\vec{V}} \vec{e}_i = - [\vec{e}_i, \vec{V}] = -\partial_i V^k \vec{e}_k , \quad (2.5.76)$$

²¹We shall often use the notation introduced in Eq. (2.5.18).

and

$$\nabla_{\vec{V}} \vec{e}_i = V^k \nabla_{\vec{e}_k} \vec{e}_i = V^k \Gamma_{ik}^l \vec{e}_l , \quad (2.5.77)$$

from which

$$\begin{aligned} (\mathcal{L}_{\vec{V}} g)_{ij} &= (\partial_i V^k g_{kj} + V^k \Gamma_{ik}^l g_{lj}) + (\partial_j V^k g_{ki} + V^k \Gamma_{jk}^l g_{li}) \\ &= V_{j;i} + V_{i;j} . \end{aligned} \quad (2.5.78)$$

The Killing equation (2.3.83) finally reads

$$0 = (\mathcal{L}_{\vec{V}} g)_{ij} = V_{i;j} + V_{j;i} \equiv V_{(i;j)} . \quad (2.5.79)$$

Note that in a normal frame around the point P , the above becomes

$$V_{(i,j)}|_P = 0 \quad (2.5.80)$$

which clearly holds for the coordinate vectors $\vec{V} = \partial_i$, in agreement with what we discussed before. However, as we shall see, the true Killing vector fields of main interest are those defined on (an open subset of) the entire manifold, and those may usually not be used to define a global reference frame. For example, the Killing vectors corresponding to rotations $\vec{\ell}_{(i)}$ around the three orthogonal axes of \mathbb{R}^3 do not commute and cannot be used to generate a proper reference frame in all of \mathbb{R}^3 .

Chapter 3

General Relativity

We now have developed all the necessary mathematical tools to introduce a theory that does not rely on specific observers, and, at the same time, can reproduce the well-known and well tested results of Special Relativity in suitable “preferred frames” [12, 13, 14].

3.1 Arbitrary observers and gravity

We started from Newtonian mechanics and its Galilean invariance, that is a Principle of Relativity for the laws of mechanics with absolute time, which is compatible with Newton’s law of gravity. The non-invariance of Maxwell’s equations led us to replace the Galilean Principle of Relativity with an enlarged version, the Principle of Special Relativity, that covers electromagnetism and further requires invariance of the speed of light:

Galilean Relativity: “The laws of (Newtonian) mechanics are the same for all inertial observers (and time is absolute).”

- 1) Newton’s law of gravity (conservative forces): consistent and yields very a accurate description of astronomical observations.
- 2) Maxwell’s electromagnetism: incompatible.

↓

Special Relativity: “The laws of physics are the same for all inertial observers and the speed of light in vacuum is invariant.”

- 1) Newton’s law of gravity (action at a distance): incompatible.
- 2) Maxwell’s electromagnetism (field-mediated interactions): fully endorsed.

To summarize, Special Relativity has (at least) two drawbacks:

- 1) it still makes use of the ambiguous concept of inertial coordinate systems;
- 2) it claims to cover all of physics, but (the very accurately verified Newtonian theory of) gravitation is excluded.

From the mathematical point of view, Special Relativity is realized by assuming the existence of *global* (inertial) reference frames connected by Lorentz (Poincaré) transformations. The requirement that the laws of physics are the same is therefore given the mathematically precise meaning that physical laws may only involve quantities represented by tensors under the Lorentz (Poincaré) group and legit tensorial operations among them.

Ideally, a mathematical reference frame should be associated with a measuring apparatus. However, all physical measurements are carried out using detectors with finite spatial and temporal extension (with no *a priori* guarantee of being inertial), and should therefore be better described by generic *local* reference frames. For this reason we endeavoured the study of differential geometry, which provided us with mathematical tools (local charts, tensors and new tensorial operations) to write equations in any coordinate system, inertial or not. These tools turned out to be so powerful that we may now *write equations in the same form in any arbitrary reference frames*. It is thus tempting to speculate physics can be formulated in a way that is totally independent of the reference frame or, more physically, in a way that can be adapted to any measuring apparatus, regardless of its inertial nature. This is in essence the:

Principle of General Relativity: “The laws of physics are the same in all reference frames (for all observers).”

Assuming to each physical observer there can be associated a reference frame (and, quite ideally, also the other way around), the principle of General Relativity can be translated into the mathematical requirement that all physical laws must involve only tensors and tensorial operations in the sense of differential geometry (with no *a priori* connection with Lorentz transformations). We could actually go as far as saying that without the mathematical machinery of differential geometry, the principle of General Relativity would have remained an empty statement, as Einstein himself basically admitted when recognising the works of Ricci Curbastro and Levi-Civita ¹.

Of course, the principle of General Relativity does not tell us *what* the laws of physics are, but experiments show that Special Relativity works very well in describing phenomena in our laboratories. The question then naturally arises as to how General Relativity may be compatible with Special Relativity and solve its problems. Let us first go back to the original issue of consistently defining an inertial observer, and remove the assumption that reference frames and observers are equivalent. In fact, it is more realistic to think of observers as (possibly extended) physical apparati that move along trajectories (curves) in space-time, starting from which one can then define mathematical reference frames that cover larger portions of the space-time manifold ². In order to qualify any such apparatus as defining an inertial frame we would then need an independent way to determine whether an object is subject to a force. If we believe in our present knowledge of fundamental forces within Special Relativity, this is actually possible for electromagnetism and nuclear forces, because

¹Two of the founding fathers of the then-called “absolute tensorial calculus”, which, in modern terms, amounts to the introduction of the covariant derivative.

²Any reference to the exponential map and alike is clearly implied.

Standard model of elementary particles: “The strength of gauge (vector) field-mediated interactions ³ is governed by charges of both signs.”

Put another way, gauge interactions are both attractive and repulsive. By preparing an object with zero charge(s), we are therefore guaranteed that the only force acting on it could be gravity.

It is a fact that the gravitational attraction between two bodies cannot be made to vanish, however gravitational effects can be eliminated from the picture by considering a *freely falling* observer, which will not measure any gravitational acceleration in whatever experiment he carries on. The latter two observations are encoded in the

Equivalence Principle: “For all physical objects, the gravitational charge (mass) m_g equals the inertial mass m_i ⁴.”

This was first hypothesised by Galileo, who (presumably) verified it by letting objects fall from the Pisa tower and observing they reached the ground at the same time, independently of their mass, shape or chemical composition. Of course, since the Newtonian description of this experiment is sufficiently accurate, we can say that this result occurs because, from Newton’s second law for a massive particle in a homogeneous and constant gravitational acceleration field \vec{g} , one as

$$m_g \vec{g} = m_i \vec{a} \quad \Rightarrow \quad \vec{a} = \vec{g} , \quad (3.1.1)$$

if $m_g = m_i$ for all bodies. In particular, both the observer (a physical apparatus) and the test bodies will sustain the same acceleration and one cannot devise any local observation that can tell whether one is not subject to any gravitational attraction at all, or if one is inside an elevator falling freely towards the ground, which is Einstein’s version of Galileo’s experiment. Such an example makes it plainly clear that a freely-falling reference frame cannot be global but must “follow the line of force of gravity”, and will therefore be *local in space and time* in general.

Keep also in mind that non-gravitational forces can be strictly made to vanish only for *point-like* objects. In fact, consider for example a ruler we wish to use in order to define a spatial axis of our reference frame. Internal electromagnetic and nuclear forces will keep this ruler of a fixed length, so that, if its centre of mass is in free fall, the end-points will not, and the corresponding reference frame will be strictly inertial only along the trajectory of the centre of mass (a point in space). Of course, whereas the notion of zero in mathematics is precise, physically a quantity is zero if we cannot tell its measured value apart from zero within our experimental errors. One can therefore assume that freely falling, inertial frames can be defined in a sufficiently small neighbourhood U_P of each space-time point P , and the laws of Special Relativity, which may strictly hold only at each point P , will also be sufficiently good approximations of the true laws inside U_P for all inertial observers

³The gauge vector fields of electroweak and strong interactions.

⁴When the latter is not zero. This excludes photons and other massless particles, which cannot be stopped or accelerated.

defined therein. In fact, since the metric must (locally in space and time) reduce to the Minkowski form for freely falling observers, there must also exist corresponding “Killing vectors at a point P ” tangent to the trajectories of each freely falling observer (meaning the Lie derivatives of the metric tensor are given by partial derivatives *and* vanish along the *orthonormal* directions at each point where the inertial frame is defined). One can make use of such Killing vectors to build local reference frames starting from each point P of the trajectory of the freely falling observer.

These are the kind of frames we qualified as *Gaussian normal*, and the symmetries of Special Relativity will then hold exactly at P and (approximately) in a (sufficiently small neighbourhood) U_P . A particularly neat example is given by a space station orbiting the earth. Its trajectory can be (approximately) described by an ellipse in space-time from the point of view of an observer on the earth, however the station is in free fall and one could place rulers on the inner walls of a living area inside the station to define a triad of space-like vectors and from these generate a local inertial frame. From the point of view of the earth observer, these three vectors and the one time-like vector tangent to the station’s trajectory “rotate” along the ellipse, although they truly define a parallel transport along the station’s trajectory ⁵. If we next consider a second space station orbiting the earth not far from the previous one, we can repeat the same construction and build a second locally inertial reference frame. We can then paste together these two frames “smoothly”. However, since the two stations orbit with different angular velocities (from the point of view of the earth observer), it is clear that their tetrads will “rotate” with different speeds as well, and a vector parallelly transported along a closed path in this reference frame will consequently not coincide with itself ⁶. In other words, we know that the Newtonian theory predicts the appearance of gravitational tidal forces between the two stations (which could be measured, for example, by means of a spring connecting them). These forces reflect in the non-vanishing of the Riemann tensor in General Relativity, thus space-time curvature, as a consequence of parallel transport being tied to local inertial observers.

To summarise, from the mathematical point of view, the Equivalence Principle means that freely falling observers are the true inertial observers, for which one finds

$$\mathcal{L}_{\bar{e}_\mu} \simeq \partial_\mu \simeq \nabla_\mu \equiv \nabla_{\partial_\mu} , \quad (3.1.2)$$

along normal directions in suitably small neighbourhoods U_P , and physics in these frames must be locally (and, in the worst case, only at a space-time point P) described, according to Special Relativity, by tensorial equations, in the sense of the local Lorentz group. (General translations are of course lost since they connect observers at different locations.) The Lorentz group $SO(3,1)$ hence remains a symmetry of physics (strictly speaking) in the tangent space T_P at all points P of the space-time manifold \mathcal{M} and approximately in the sufficiently small neighbourhoods U_P . According to the principle of General Relativity, different (non-inertial) observers will then see the laws of physics of Special Relativity as

⁵These four basis vectors are called *tetrads* or *vierbien* and must be explicitly introduced to describe the spin in General Relativity.

⁶Note that this operation cannot be realised physically, since nothing can travel along a closed path in space-time without violating causality (perhaps)!

tensorial equations in the sense of differential geometry, with the partial derivatives replaced by covariant derivatives. This is precisely encoded in yet another principle:

Principle of General Covariance: “The laws of physics in a general reference frame are obtained from the laws of Special Relativity by replacing tensor quantities of the Lorentz group with tensor quantities of the space-time manifold.”

In practical terms, this means that one takes a law of physics in the locally inertial frame at a point P as given by Special Relativity and:

- a) re-interpret tensorial indices of the Lorentz group as representing the components of tensors under general coordinate transformations;
- b) further, the Minkowski metric (used to raise, lower and contract indices) must be replaced by a general metric tensor with the same signature

$$\eta_{\mu\nu} \rightarrow g_{\mu\nu} , \quad (3.1.3)$$

where, from now on, we shall only consider four-dimensional space-time manifolds \mathcal{M} with coordinates $x^\mu = (x^0, x^1, x^2, x^3) = (x^0, x^i)$ and metric signature $(-, +, +, +)$, unless differently specified;

- c) General Covariance and Eq. (3.1.3) then imply that partial derivatives must also be replaced by the metric covariant derivative

$$\partial_\mu \equiv \frac{\partial}{\partial x^\mu} \rightarrow \nabla_\mu . \quad (3.1.4)$$

For example, Maxwell’s equations in General Relativity will simply read

$$\partial_\mu F^{\mu\nu} = -J^\nu \rightarrow \nabla_\mu F^{\mu\nu} = -J^\nu . \quad (3.1.5)$$

From the physical point of view, this seemingly simple mathematical replacement should not undermine the fact that there are now two terms in the left hand side: one representing the usual flat space gradient, and the second entailing the effect of curved space-time (gravity) on the propagation of electromagnetic degrees of freedom (the photons).

It really cannot be emphasised enough that, for this construction to work, it is crucial that one can always put the metric tensor in canonical form locally, with its first derivatives (the metric connection) locally vanishing at the same time, so that Eq. (3.1.2) holds. Without this general property of metric manifolds, one could not embed Special Relativity inside General Relativity and, given the experimental success of the former, the latter could not be made into a physical theory at all. It is not hard to see that this mathematical property must have been a true source of inspiration for Einstein’s ideas.

3.2 Gravitational equations

The above construction covers all the interactions of the Standard Model of particle physics, but does not yet explicitly include a description of gravity at all. This means two questions are still open:

Q1) how do we describe the action of gravity on a test particle?

Q2) what sources gravity and how do we determine gravity from its sources?

There is a natural answer for Q1) that follows from the stated principles (like Newton's force law naturally follows from the principle of Galilean invariance), whereas Q2) must be addressed as an independent issue (much like Maxwell's equations and the equations governing any fundamental interaction do not follow from relativity principles).

3.2.1 Gravity and test particles

First of all, the Newtonian idea that gravity is represented by an acceleration field \vec{g} cannot work, since \vec{g} is a three-vector for which one can hardly conceive a “temporal component” to build up a four-vector g^μ . However, Newtonian gravity describes the motion of celestial bodies with very high accuracy and this implies that General Relativity must reduce to the Newtonian theory in some suitable limit.

Freely falling observers and test particles

Since a freely falling observer is “inertial”, the local metric in its own reference frame is the canonical Minkowski metric all along (of course, only in a sufficiently small neighbourhood of each point P of the observer's trajectory). Let $u^\mu = dx^\mu/d\tau$ denote the four-velocity of a test particle subject to no other force (but gravity). In the freely falling frame, it must then move along a straight line,

$$0 = \gamma^2 \frac{d^2 x^\alpha}{dt^2} = \frac{d^2 x^\alpha}{d\tau^2} = \frac{d^2 x^\alpha}{d\tau^2} + \Gamma_{\mu\nu}^\alpha \frac{dx^\mu}{d\tau} \frac{dx^\nu}{d\tau} = u^\mu \nabla_\mu u^\alpha, \quad (3.2.1)$$

since $\gamma\tau = x^0 \equiv t$, where $\gamma = (1 - u^2/c^2)^{-1/2}$ is the usual special relativistic factor for a particle moving with (constant) velocity $\vec{u} = \frac{d\vec{x}}{dt}$, and $\Gamma_{\mu\nu}^\alpha \sim g_{\mu\nu,\beta} = 0$ in the coordinate system of the freely falling observer. Note that this argument, strictly speaking, only holds at a space-time point, say P , where the trajectory of the freely-falling observer and the trajectory of the test particle happen to cross. However, the final result (3.2.1) is frame-independent and we can simply say that test particles follow geodesics of the given space-time metric. These trajectories are usually referred to as *world-lines*.

This argument further implies that the inertial observer itself moves along a geodesic of the space-time metric, as can be simply deduced by considering the observer as a test particle at rest with respect to itself. In a different frame (equivalently, for a different observer), the Christoffel symbols will not be zero at P , and this suggests that *the metric $g_{\mu\nu}$ can be viewed as a potential for the gravitational interaction*,

$$\Gamma_{\mu\nu}^\alpha \sim g_{\mu\nu,\beta}, \quad (3.2.2)$$

like the four-vector A^μ is a potential for the electromagnetic field, $F^\mu{}_\nu \sim A^\mu{}_{,\nu}$.

The Newtonian limit

The previous conclusion can be further supported by considering the *weak field limit* and *non-relativistic limit* of the geodesic equation, in which we expect Newton's law of gravity is recovered. Non-relativistic means we expect the test particle moves much slower than $c = 1$ in the relevant reference frame, and the weak field limit means that the metric $g_{\mu\nu}$ is static and very close to $\eta_{\mu\nu}$, in the same frame.

Giving these limits a precise meaning is however far from straightforward, and we shall employ a mathematical trick to formalise our procedure. Namely, let us introduce a parameter $0 < \epsilon \leq 1$ such that the particle's spatial velocity can be written as $\vec{u} = \epsilon \vec{v}$, where $v < c = 1$, and the metric is given by

$$g_{\mu\nu} = \eta_{\mu\nu} + \epsilon h_{\mu\nu} . \quad (3.2.3)$$

Both the non-relativistic limit and the weak field limit can now be implemented by Taylor expanding in ϵ all of our expressions and keep only the first order. (We can then formally set $\epsilon = 1$ at the end of the computation in the truncated expressions.) In particular, the four-velocity becomes

$$\begin{aligned} u^\mu &= (1 + \mathcal{O}(\epsilon^2), \epsilon \vec{v} + \mathcal{O}(\epsilon^2)) \\ &= (1, \vec{0}) + \epsilon (0, \vec{v}) + \mathcal{O}(\epsilon^2) , \end{aligned} \quad (3.2.4)$$

so that

$$\frac{d^2 x^\alpha}{d\tau^2} = \epsilon \left(0, \frac{d\vec{v}}{dt} \right) + \mathcal{O}(\epsilon^2) , \quad (3.2.5)$$

and the Christoffel symbols read

$$\begin{aligned} \Gamma_{\mu\nu}^\alpha &= \frac{1}{2} g^{\alpha\beta} (g_{\mu\beta,\nu} + g_{\nu\beta,\mu} - g_{\mu\nu,\alpha}) \\ &= \frac{\epsilon}{2} g^{\alpha\beta} (h_{\mu\beta,\nu} + h_{\nu\beta,\mu} - h_{\mu\nu,\alpha}) \\ &= \frac{\epsilon}{2} \eta^{\alpha\beta} (h_{\mu\beta,\nu} + h_{\nu\beta,\mu} - h_{\mu\nu,\alpha}) + \mathcal{O}(\epsilon^2) , \end{aligned} \quad (3.2.6)$$

where the derivatives of the metric are different from zero only if they are not taken with respect to time, and we must recall that $\eta^{\alpha\beta}$ is diagonal. This implies that the only non-trivial components of the geodesic equation at order ϵ are given by (no sum over i)

$$\begin{aligned} \frac{d^2 x^\alpha}{d\tau^2} + \Gamma_{\mu\nu}^\alpha \frac{dx^\mu}{d\tau} \frac{dx^\nu}{d\tau} &\simeq \epsilon \frac{d^2 x^i}{dt^2} + \Gamma_{\mu\nu}^i \delta_0^\mu \delta_0^\nu \\ &\simeq \epsilon \left(\frac{d^2 x^i}{dt^2} - \frac{1}{2} \eta^{ii} h_{00,i} \right) , \end{aligned} \quad (3.2.7)$$

in which we used

$$\Gamma_{00}^i \simeq \frac{\epsilon}{2} \eta^{ii} (h_{0i,0} + h_{0i,0} - h_{00,i}) = -\frac{\epsilon}{2} \eta^{ii} h_{00,i} . \quad (3.2.8)$$

Now, if we call “gravitational potential” the function

$$V = -\frac{1}{2} \eta^{ii} h_{00} \simeq -\frac{1}{2} h_{00} , \quad (3.2.9)$$

and set $\epsilon \rightarrow 1$, we have recovered the Newton-like equation

$$\frac{d^2 x^i}{dt^2} = -\frac{\partial V}{\partial x^i} . \quad (3.2.10)$$

We shall see later that V in Eq. (3.2.9) exactly reproduces Newton’s potential once a solution for $g_{\mu\nu}$ is obtained outside a spherically symmetric body.

An important remark is now in order. We have not yet explicitly considered *light*, that is a signal which propagates at a speed equal to c and can therefore be associated with a massless particle. In Special Relativity, light propagates along the null cone, which is a geodesic of the Minkowski metric, and one can easily show that this result generalises to any space-time metric. In fact, the modulus $u^\mu u_\mu = C$ of the (parallel transported) tangent vector u^μ to a geodesic is conserved along the geodesics itself, since

$$u^\nu \partial_\nu C = 2 u^\nu u^\mu \nabla_\nu u_\mu = 2 u^\mu (u^\nu \nabla_\nu u_\mu) = 0 . \quad (3.2.11)$$

Given a point P along a physical geodesic, its four-velocity must satisfy $u^\mu u_\mu = C$, where $C = -1$ for massive particles and $C = 0$ for light, in a locally inertial reference frame at P . The principle of General Covariance then implies that $u^\mu u_\mu = C$ in any reference frame and Eq. (3.2.11) ensures the modulus is conserved along the trajectory. We can therefore conclude that light also propagates along geodesics, although there is no (affine) parameter along such geodesic that can be identified with a “proper time”⁷. It also follows from Eq. (3.2.11) that the metric $g_{\mu\nu}$ encodes much more information than a usual “potential” field, like, for example, the four-vector A^μ of electromagnetism: it determines the causal structure of space-time by governing the propagation of light and of any other signal. This is the very essence of Einstein’s “geometric view” of gravity.

3.2.2 Source of gravity and Einstein equations

Now, answering Q2) is a lot more of a guesswork. First of all, $g_{\mu\nu}$ is symmetric and contains (at most) 10 independent components. We therefore need ten equations for such components and we would like they be at most second order partial differential equations, like is the case for Maxwell’s equations. We therefore need a tensor constructed solely from $g_{\mu\nu}$ and up to

⁷Keep in mind that geodesics are defined modulo an arbitrary reparametrization along the world-line. Only one of such affine parameters will coincide with the proper time (for a massive particle).

its second partial derivatives and there are not many choices: the Riemann tensor and its contractions. One possibility is the $(0, 2)$ Einstein tensor,

$$G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2} R g_{\mu\nu} , \quad (3.2.12)$$

which is obviously symmetric and therefore contains 10 independent components. We actually recall that $G_{\mu\nu}$ is covariantly conserved (contracted Bianchi identities),

$$\nabla_\mu G^\mu{}_\nu = \nabla_\mu \left(R^\mu{}_\nu - \frac{1}{2} R g^\mu{}_\nu \right) , \quad (3.2.13)$$

and the above property reduces the number of components of the Einstein tensor from $(4 \times 5)/2 = 10$ to $10 - 4 = 6$ (which, incidentally, is the number of components of the spatial metric).

Energy-momentum tensor

If the tensor (3.2.12) is to be the left hand side of the equation which determines the metric, the source on the right hand side must have the same mathematical properties: it must be a symmetric and covariantly conserved $(0, 2)$ tensor built out of the matter content of the system. One such tensor is the *energy-momentum tensor*, which for a perfect fluid with four-velocity u^μ , is given by

$$T^{\mu\nu} = \rho u^\mu u^\nu + p (g^{\mu\nu} + u^\mu u^\nu) = (\rho + p) u^\mu u^\nu + p g^{\mu\nu} , \quad (3.2.14)$$

where ρ is the (proper) density and p the (proper) pressure, both measured by an observer comoving with the fluid ⁸, which makes such quantities true scalars. Note the tensor multiplying the pressure is orthogonal to u^μ (recall that $u^\mu u_\mu = -1$),

$$(g^\mu{}_\nu + u^\mu u_\nu) u^\nu = u^\mu - u^\mu = 0 . \quad (3.2.15)$$

In fact, in the frame comoving with the fluid, the four-velocity $u^\mu = (1, 0, 0, 0)$ and $d\tau = dt$, which means $g_{00} = g^{00} = -1$ and $g_{0i} = g^{0i} = 0$. We then have

$$T^{\mu\nu} = \begin{bmatrix} \rho & 0 \\ 0 & p g^{ij} \end{bmatrix} , \quad (3.2.16)$$

which implies

$$T^\mu{}_\nu = \text{diag} (-\rho, p, p, p) . \quad (3.2.17)$$

In order to better understand the meaning of the energy-momentum tensor, let us first consider a particle of four-momentum p^μ and an observer with four-velocity U^μ . It is easy to see that the energy of the particle as measured by the observer is given by

$$E^{(0)} = -p^\mu U_\mu . \quad (3.2.18)$$

⁸A necessary requirement is therefore that the fluid particles be massive. We shall see later on how to deal with radiation.

In fact, let us chose a freely falling reference frame in which

$$U^\mu = (1, 0, 0, 0) \equiv V_{(0)}^\mu , \quad (3.2.19)$$

at the space-time point P where the measurement takes place ⁹. Then

$$p^\mu V_{(0)\mu} = -m \frac{dt}{d\tau} = -\frac{m}{\sqrt{1-u^2}} , \quad (3.2.20)$$

where $\vec{u} = \frac{dx^i}{d\tau}$ is the particle's three-velocity in the locally inertial frame of choice and we used $U_0 = -U^0 = -1$. We can of course complete the *vierbein* of basis vectors of T_P with three spatial vectors orthogonal to U^μ , say $V_{(i)}^\mu$. It is then not difficult to see that we can choose local coordinates such that $V_{(i)}^\mu = \delta_i^\mu$ and

$$p_\mu V_{(i)}^\mu = \frac{p^i}{\sqrt{1-u^2}} = \frac{m u^i}{\sqrt{1-u^2}} \quad (3.2.21)$$

yield the spatial components of the four-momentum of the particle in the rest frame of the observer. Note that quantities measured by a given observer are correctly represented by scalars, namely

$$E_{(a)} = p_\mu V_{(a)}^\mu , \quad a = 0, 1, 2, 3 , \quad (3.2.22)$$

which, furthermore, form the components of a Lorentz vector (under local Lorentz transformations of the tetrad at P ¹⁰).

Let us now consider the density term in the energy-momentum tensor, and contract it with $U^\mu = (1, 0, 0, 0)$ twice,

$$(\rho u_\alpha u_\mu U^\mu) U^\alpha = \left(\frac{\rho}{\sqrt{1-u^2}} u_\alpha \right) U^\alpha = \frac{\rho}{1-u^2} . \quad (3.2.23)$$

The two factors of $\sqrt{1-u^2}$ in the denominator are easily explained by first recalling that the proper density (as measured by an observer comoving with the fluid) is defined by

$$\rho = \frac{m}{V_0} , \quad (3.2.24)$$

where m is the proper mass (equal to the energy) of fluid particles contained in a proper (or comoving) volume V_0 , both of which scale with factors of $\sqrt{1-u^2}$ when measured by the observer moving with relative speed u ,

$$m \rightarrow \frac{m}{\sqrt{1-u^2}} \quad \text{and} \quad V_0 \rightarrow V_0 \sqrt{1-u^2} , \quad (3.2.25)$$

⁹Such a locally inertial observer will coincide with our chosen observer only at the space-time point P . The metric in P will have the canonical Minkowski form.

¹⁰Incidentally, this further shows how the Lorentz group of Special Relativity is embedded as a local symmetry into General Relativity.

since only lengths parallel to \vec{u} are contracted. A similar analysis holds for the (spatial) pressure term.

An important property of fluids is the *continuity equation*, which in a locally inertial frame reads

$$\partial_t \rho + \vec{\nabla} \cdot \vec{p} = 0 , \quad (3.2.26)$$

and just tells us that energy is conserved: the loss of energy per unit time $-\partial_t \rho$ of a portion of fluid inside a proper volume V_0 equals the work per unit time done by that fluid to expand V_0 (which, for example, is given by $-\partial_x p^x$ in the x -direction). Another way to look at Eq. (3.2.26) is by integrating it over a cell of volume V_0 . Let us for example assume the cell is a cubic box with $a_- \leq x, y, z \leq a_+$, and $\vec{p} = (p^x(x), 0, 0)$, so that

$$\int_{V_0} \partial_t \rho \, dx \, dy \, dz = \partial_t \int_{V_0} \rho \, dx \, dy \, dz = \partial_t E , \quad (3.2.27)$$

where E is the energy inside the volume $V_0 = (a_+ - a_-)^3$, and

$$\int_{V_0} (\partial_x p^x \, dx) \, dy \, dz = A_0 [p^x(a_+) - p^x(a_-)] = F^x(a_+) - F^x(a_-) , \quad (3.2.28)$$

where $A_0 = (a_+ - a_-)^2$ is the area of the square surfaces at $x = a_{\pm}$ and $F^x = A_0 p^x$ is the force acting on such surfaces. If matter is not created or destroyed inside the cell, we can apply the usual theorem relating kinetic energy to the work of the force, $F^x = \partial E / \partial x$, and finally obtain

$$\partial_t E = \left. \frac{\partial E}{\partial x} \right|_{x_-} - \left. \frac{\partial E}{\partial x} \right|_{x_+} . \quad (3.2.29)$$

The above means that the rate of energy increase inside the cell equals the amount of energy which enters from the left (through the surface at $x = a_-$) minus the energy which exits (through the surface at $x = a_+$). In a general reference frame, (3.2.26) means that the energy-momentum tensor is covariantly conserved,

$$\nabla_\mu T^{\mu\nu} = 0 , \quad (3.2.30)$$

where we recall the covariant derivative now contains extra terms (with respect to the Minkowski case) which represent the effect of gravity on the fluid. This makes $T^{\mu\nu}$ the natural candidate as the source of gravity.

Einstein equations and Newtonian approximation

In order to recover the Newtonian approximation in this framework, we must assume the local curvature is small, so that the metric can be written as in Eq. (3.2.3). The Ricci scalar then takes the simple form

$$R = \epsilon (\Box h - \partial^\mu \partial^\nu h_{\mu\nu}) + \mathcal{O}(\epsilon^2) , \quad (3.2.31)$$

where

$$\square = -\partial_t^2 + \Delta \quad (3.2.32)$$

is the d'Alembertian in flat space, the trace $h = \eta_{\mu\nu} h^{\mu\nu}$, and the linearised Einstein field equation is given by

$$\epsilon \left(-\square h_{\mu\nu} + \eta_{\mu\nu} \square h + \partial_\mu \partial^\lambda h_{\lambda\nu} + \partial_\nu \partial^\lambda h_{\lambda\mu} - \eta_{\mu\nu} \partial^\lambda \partial^\rho h_{\lambda\rho} - \partial_\mu \partial_\nu h \right) \simeq 16 \pi G_N T_{\mu\nu} , \quad (3.2.33)$$

which shows that gravity at order ϵ couples to matter at order ϵ^0 or, equivalently, that one must view Newton's constant G_N as a quantity of order ϵ . In the de Donder gauge,

$$2 \partial^\mu h_{\mu\nu} = \partial_\nu h , \quad (3.2.34)$$

the trace of the field equation yields

$$\epsilon \square h = 16 \pi G_N T , \quad (3.2.35)$$

where $T = \eta^{\mu\nu} T_{\mu\nu}$, and Eq. (3.2.33) reduces to

$$-\epsilon \square h_{\mu\nu} = 16 \pi G_N \left(T_{\mu\nu} - \frac{1}{2} \eta_{\mu\nu} T \right) . \quad (3.2.36)$$

In addition to the weak field limit, we assume that all matter in the system moves with a characteristic velocity much slower than the speed of light in the (implicitly) chosen reference frame $x^\mu = (t, \mathbf{x})$. The only relevant component of the metric is therefore $h_{00}(\mathbf{x})$, and its time derivatives are also neglected¹¹. The stress-energy tensor is accordingly determined by the energy density,

$$T^{\mu\nu}(\mathbf{x}) = \delta_0^\mu \delta_0^\nu \rho(\mathbf{x}) \simeq T_{00} \simeq -T , \quad (3.2.37)$$

and the Ricci scalar reduces to

$$R \simeq -\Delta h_{00}(\mathbf{x}) . \quad (3.2.38)$$

In this approximation, Eq. (3.2.36) takes the very simple form

$$\Delta h_{00}(\mathbf{x}) = -8 \pi G_N T_{00}(\mathbf{x}) = -8 \pi G_N \rho(\mathbf{x}) . \quad (3.2.39)$$

Since the Newtonian potential V_N is generated by the mass density ρ according to the Poisson Equation

$$\Delta V_N = 4 \pi G_N \rho , \quad (3.2.40)$$

we can finally identify $h_{00} = -2 V_N$.

¹¹For static configurations, the gauge condition (3.2.34) is always satisfied.

Gravity and geometry

Let us summarise the principles of General Relativity and their connection with the mathematical background we have developed:

$$\begin{array}{ccccc} \text{G. R.} & \Leftrightarrow & \text{Differential Geometry} & \Leftrightarrow & \text{Tensor calculus} \\ & & & & \\ \left. \begin{array}{l} \text{E. P.} \\ \text{G. C.} \end{array} \right\} & \Leftrightarrow & g_{\mu\nu} = \eta_{\mu\nu} + O(2) & \Leftrightarrow & \text{Local } SO(3, 1) \end{array}$$

There is therefore an explicit connection between this description of gravity and the geometry of space-time, which deserves some clarifications.

By looking at the equation of motion of massive test particles, it is clear that the concept of “straight line” is replaced by that of geodesic line. However, the geodesic equation is a second order differential equation like Newton’s law of mechanics, and one may just look at the connection term as a “force” acting on the particle. We have in fact seen that such term reduces to that of a conservative force in the weak field limit. Since no particle (massive or massless) sees a (globally) flat space-time ¹²,

$$\nabla_{\vec{v}} \vec{v} = \frac{d^2 x^\mu}{d\lambda^2} + \Gamma_{\nu\alpha}^\mu \frac{dx^\nu}{d\lambda} \frac{dx^\alpha}{d\lambda} = 0, \quad \text{with } v^\mu v_\mu = 0, \text{ or } -1, \quad (3.2.41)$$

we can start to think of gravity in terms of pure geometry.

The above conclusion is further supported by the Einstein field equations, according to which matter sources determine the space-time curvature, which in turn affects the matter’s motion. This is all encoded in ten (with only six independent) partial differential equations of a highly non-linear kind: unlike Newton’s law, the effect of two gravitational sources is not just their sum. A graphical picture of this geometrical view of gravity is given by the so-called Einstein’s billiard: the space-time is represented by a sheet of elastic material upon which rest the sources, and motion of test particles therefore follow curvy lines.

3.2.3 Classical tests of General Relativity

There are three historical tests conducted within the Solar system which strongly support General Relativity.

Perihelion precession of Mercury

By solving the equation of motion for a test particle in the gravitational field of a much more massive body, one finds almost elliptic orbits, similar to those predicted by Newtonian mechanics. The difference can be modeled by a rotation of the axes of the ellipse or, equivalently, of the point of minimum distance from the source, around the source itself. This point, for planets in the solar system moving around the sun, is called perihelion, and

¹²Note λ can be the proper time for massive particles and is a generic affine parameter for light signals.

its motion was already predicted in studies of Newtonian celestial mechanics, due to the presence of the other planets in the solar system.

A small fraction of the observed precession for the perihelion of the orbit of Mercury, however, remained unexplained until Einstein employed his new theory to find an astonishingly good agreement: General Relativity was able to explain a difference with respect to the Newtonian theory of just 43'' per century.

Deflection of light (gravitational lensing)

By again studying the motion of a massless particle (photon) around the sun, one can see that it will be affected by the gravitational field and move along a trajectory very close to an hyperbola, like any massive particles would.

This effect was first seen by Arthur Eddington and collaborators who organised expeditions to Brazil and Africa during a total solar eclipse in 1919. The eclipse allowed to see stars whose image was close enough to the sun's surface to amplify the deflection to measurable values.

Gravitational redshift

Photons loose energy when they climb up into a gravitational field. This effect was first observed by Pound and Rebka in 1959 using gamma-rays traveling along the 72 m tall tower of the Jefferson Physical Laboratory in Harvard. When an atom makes a transition from an excited state to the ground state, it emits a photon with a specific frequency and energy. Conversely, when the same atom in its ground state hits a photon with that same frequency and energy, it will absorb the photon and jump to the same excited state. If the photon's frequency is even slightly different, the atom will not absorb it. When the photon travels through a gravitational field, its frequency and therefore its energy will change due to the gravitational redshift and, as a result, the receiving atom can no longer absorb it. But if the emitting atom moves with just the right speed relative to the receiving atom, the resulting Doppler shift will cancel out the gravitational shift and the receiving atom will be able to absorb the photon. The relative speed of the atoms is therefore a measure of the gravitational shift. Pound and Rebka measured the gravitational blueshift by moving the emitter at the top of the tower away from the receiver placed at its bottom. Their experiment also involved the Mössbauer effect to detect the recoil of the atoms which actually absorbed the photons.

At least ideally, it is possible that the gravitational redshift of a photon moving away from a matter source cancels out the initial photon energy completely. In this case, we have a *black hole*. We shall see that the Newtonian counterpart of this situation is

$$\text{escape velocity} > \text{light speed} , \tag{3.2.42}$$

and was hypothesized well before General Relativity.

3.3 Black holes

Soon after Einstein proposed the equations (??) for the gravitational interaction, Karl Schwarzschild found a spherically symmetric solution which carries his name, and is the prototype of a *black hole* space-time. We shall here sketch the derivation and review some of its main features by studying geodesics.

3.3.1 The Schwarzschild metric

Let us consider a spherically symmetric source, such as would be the earth or the sun to first approximation. Of course, modeling the interior of an astrophysical object is anything but easy. However, if we are just interested in the region outside the source, everything simplifies significantly. In fact, outside the source, the space-time is empty and $T_{\mu\nu} = 0$. Upon taking the trace of the Einstein tensor, we obtain for the curvature scalar

$$R - \frac{1}{2} R g^\mu{}_\mu = -2 R = 0 , \quad (3.3.1)$$

and Eq. (??) simplifies to

$$R_{\mu\nu} = 0 . \quad (3.3.2)$$

We can now try and solve Eq. (3.3.2), and, in doing so, we expect the general solution will depend on free parameters that we could later fix by means of information coming from the region where $T_{\mu\nu} \neq 0$. In the specific case at hand here, we will in fact see that the assumed symmetry of the space-time is strong enough to reduce the arbitrariness to just one parameter, whose physical meaning can be obtained from the weak-field limit, regardless of the details of the source. We shall call this one-parameter family of spherically symmetric solutions to Eq. (3.3.2) the *Schwarzschild (metric) manifold*, or Schwarzschild space-time.

Finding solutions to Eq. (3.3.2) can be greatly eased by making use of isometries, that is, by assuming the existence of Killing vectors. First of all, we shall require the metric is static, so that there exists a time-like Killing vector \vec{K}_t , and a suitable coordinate t , so that we can write

$$\vec{K}_t = \frac{\partial}{\partial t} . \quad (3.3.3)$$

Moreover, since the source is spherically symmetric, we also assume the existence of three space-like Killing vectors corresponding to rotations around axes with origin at the centre of the source,

$$\vec{K}_i = \frac{d}{d\theta_i} , \quad i = 1, 2, 3 . \quad (3.3.4)$$

The above three vectors must be conserved in time, which means they must commute with \vec{K}_t ,

$$\left[\frac{\partial}{\partial t}, \frac{d}{d\theta_i} \right] = 0 . \quad (3.3.5)$$

We may therefore assume the metric is such that rotations are orthogonal to \bar{K}_t , like they are in the Minkowski space-time \mathbb{R}^{1+3} , and that we can use the analogue of polar coordinates on surfaces of constant t . This allows us to write the metric in diagonal form

$$\begin{aligned} ds^2 &= -A(r) dt^2 + B(r) dr^2 + C(r) (d\theta^2 + \sin^2 \theta d\phi^2) \\ &= -A(r) dt^2 + B(r) dr^2 + r^2 d\Omega^2 , \end{aligned} \quad (3.3.6)$$

where we have also used the freedom to rescale the radial coordinate r so that $C = r^2$. Since the metric tensor element only depend on r (and trivially on θ), one can always redefine the coordinate r so that the above holds *locally*. However, one should keep in mind that the rescaling *may* change the domain of definition of r , if the transformation becomes singular at some point. With this choice, the area of a sphere of coordinate radius r is given by

$$A(r) = \int d\Omega^2 = \int \sqrt{\det g_{(2)}} d\theta d\phi = r^2 \int \sin \theta d\theta d\phi = 4\pi r^2 . \quad (3.3.7)$$

For this reason r is commonly called the *areal radius*. Note now that the proper length of the radius of such a sphere is not r ,

$$R(r) = \int_0^r \sqrt{g_{rr}} dx = \int_0^r \sqrt{B(x)} dx , \quad (3.3.8)$$

unless $B = 1$, and the space-time outside the source is in general curved. Finally, let us note that, in order for $\{t, r, \theta, \phi\}$ to be a proper chart, we should also define their domain of definition (technically, the image of the open subset of the Schwarzschild manifold covered by these coordinates). The assumed isometries again help here, since time-translation invariance means we can suppose

$$-\infty < t < +\infty , \quad (3.3.9)$$

and the spatial volumes (at fix t) can be foliated by two-dimensional spheres (the submanifolds generated by the rotations) on which the angular coordinates take their usual values ¹³

$$0 \leq \theta \leq \pi , \quad 0 \leq \phi < 2\pi . \quad (3.3.10)$$

Not much can yet be said for the radial coordinate r , for which we can only expect it goes up to infinity, where the metric presumably approaches the Minkowski form. Of course, this does not assure us that the chosen set of coordinates covers the whole Schwarzschild manifold. In fact, this point represents a very significant difference with respect to other field equations of physics: unlike Maxwell equations, for example, which are *a priori* defined everywhere on a given manifold (such as \mathbb{R}^{1+3} in Special relativity), Einstein equations (??) [or the vacuum version (3.3.2)] implicitly define the manifold on which the metric lives. In other words, we can say that *the unknown determined by Eq. (??) is not just the metric tensor but also the manifold on which it exists*.

¹³We do not need to be particularly picky here about the “singularity” of such coordinates at the poles.

By inserting the above expression for the metric into Eq. (3.3.2), one obtains the following four equations

$$R_{00} = -\frac{A''}{2B} + \frac{A'}{4B} \left(\frac{A'}{A} + \frac{B'}{B} \right) - \frac{A'}{rB} = 0 \quad (3.3.11)$$

$$R_{11} = \frac{A''}{2A} - \frac{A'}{4A} \left(\frac{A'}{A} + \frac{B'}{B} \right) - \frac{B'}{rB} = 0 \quad (3.3.12)$$

$$R_{22} = \frac{1}{B} - 1 + \frac{r}{2B} \left(\frac{A'}{A} - \frac{B'}{B} \right) = 0 \quad (3.3.13)$$

$$R_{33} = R_{22} \sin^2 \theta, \quad (3.3.14)$$

where a prime denotes the derivative with respect to r . One can then notice that

$$0 = \frac{B}{A} R_{00} + R_{11} = A' B + B' A = (AB)' , \quad (3.3.15)$$

which implies the product AB is a constant or

$$A = B^{-1} , \quad (3.3.16)$$

where we set the dimensionless constant to one (this can always be achieved by rescaling the time variable). Finally, Eq. (3.3.13) reads $A + r A' = 1$ or

$$(rA)' = 1 , \quad (3.3.17)$$

which yields

$$A = 1 - \frac{2K}{r} , \quad (3.3.18)$$

where K has dimensions of a length and emerges as an integration constant. Note that the metric reduces to Minkowski for $K = 0$ but, at this point, does not contain any information about the source, nor the gravitational constant G_N [which in fact does not appear in the vacuum equation (3.3.2)].

In order to understand the physical meaning of K , let us look at the weak field limit, Eq. (3.2.7)¹⁴,

$$\frac{d^2 x^i}{dt^2} = \frac{1}{2} \eta^{ii} h_{00,i} = -V_{,i} . \quad (3.3.19)$$

Clearly,

$$g_{00} = \eta_{00} + \frac{2K}{r} \quad \Rightarrow \quad V = -\frac{K}{r} , \quad (3.3.20)$$

¹⁴No summation over the index i .

which implies

$$\frac{d^2 r}{dt^2} = -\frac{K}{r^2}, \quad (3.3.21)$$

and Newton's law of gravity is properly recovered if $K = G_N M$ and for

$$r \gg 2 G_N M \equiv R_H, \quad (3.3.22)$$

where we can now interpret M as the total mass of the source as measured by a distant observer.

Finally,

$$ds^2 = -\left(1 - \frac{2 G_N M}{r}\right) dt^2 + \left(1 - \frac{2 G_N M}{r}\right)^{-1} dr^2 + r^2 (d\theta^2 + \sin^2 \theta d\phi^2), \quad (3.3.23)$$

is the famous Schwarzschild metric. The length R_H plays a crucial role in the Schwarzschild space-time and is called the *Schwarzschild radius*. Note that for $r \gg R_H$, the Schwarzschild metric approaches the Minkowski flat metric, a property referred to as *asymptotic flatness*. Static observers associated with the Schwarzschild coordinates $\{t, r, \theta, \phi\}$, and placed at fixed r , are therefore asymptotically inertial (at $r \gg R_H$), but they depart from being inertial the more they approach R_H . Note also that the proper radius of a sphere of area $4\pi r^2$ is given by

$$\begin{aligned} R(r) &= \int_0^{r_H} \frac{dx}{\sqrt{\frac{r_H}{x} - 1}} + \int_{r_H}^r \frac{dx}{\sqrt{1 - \frac{r_H}{x}}} \\ &= r \sqrt{1 - \frac{r_H}{r}} + \frac{r_H}{2} \log \left[2 \frac{r}{r_H} \left(1 + \sqrt{1 - \frac{r_H}{r}} \right) - 1 \right], \end{aligned} \quad (3.3.24)$$

for $r > r_H$, so that $R(r_H) = 0$ and $R(r) > r$ for $r \gtrsim 3 r_H/2$ (see Fig. 3.1). However, for $r < r_H$, we have that $g_{rr} < 0$ and the coordinate r becomes time-like. The above result (3.3.24) is therefore of little physical meaning, and the geometry of the Schwarzschild space-time will be better understood by studying its geodesics.

3.3.2 Radial geodesics

The Schwarzschild metric (3.3.23) does not carry any dependence on the size of the spherically symmetric source that generates it. We can therefore suppose that the source has a very small areal radius $r_s \ll R_H$ (ideally reducing to a point), in which case the metric (3.3.23) shows a clearly disturbing feature:

$$r \rightarrow R_H^\pm \quad \Rightarrow \quad g_{tt} = 0 \quad \text{and} \quad g_{rr} \rightarrow \pm\infty. \quad (3.3.25)$$

Further, the signs of g_{tt} and g_{rr} change across $r = R_H$, so that inside the Schwarzschild radius t becomes a spatial coordinate and r is a time (but note that the Killing vector \bar{K}_t is associated with t everywhere, whereas there is no Killing vector $\bar{K}_r = \partial_r$ anywhere).

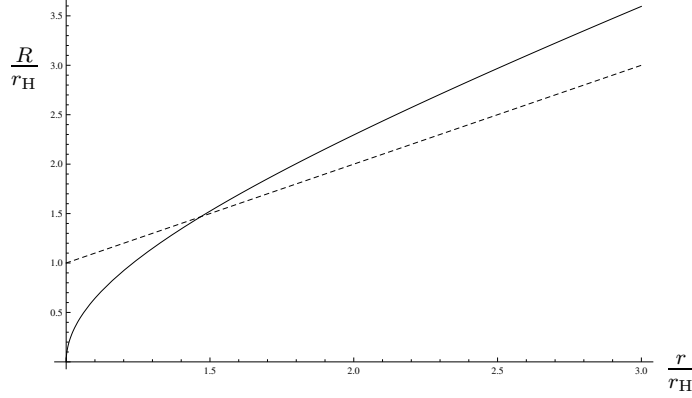


Figure 3.1: Proper radius $R = R(r)$ in units of r_H .

The above observations lead us to consider two charts for the Schwarzschild manifold, both defined by the same Schwarzschild coordinates $\{t, r, \theta, \phi\}$, but with domains:

- 1) \mathcal{M}^- , for $0 < r < R_H$, which describes the interior, and
- 2) \mathcal{M}^+ , for $r > R_H$, which describes the exterior geometry.

Both charts are properly defined on open sets and, in order to cover the entire manifold, we should then define a chart that includes the Schwarzschild sphere with $r = R_H$. From the physical point of view, this construction is consistent if a test particle moving in the Schwarzschild space-time can travel from \mathcal{M}^+ to \mathcal{M}^- without meeting any obstructions. We therefore proceed by directly studying (radial) geodesics which start from $r > R_H$.

Instead of considering the geodesic equations directly, it is more convenient to employ the property of geodesics being extremal of the length, namely vary the action for a massive test particle

$$S[x^\mu(\tau)] = m \int_0^\tau \sqrt{-g_{\mu\nu} \dot{x}^\mu \dot{x}^\nu} d\tau' \equiv m \int_0^\tau \sqrt{2T} d\tau' , \quad (3.3.26)$$

and recover the geodesic equations as the Euler-Lagrange equations of motion¹⁵. In Schwarzschild space-time and using the mass-shell condition for massive particles, namely

$$2T = \left(1 - \frac{2G_N M}{r}\right) \dot{t}^2 - \left(1 - \frac{2G_N M}{r}\right)^{-1} \dot{r}^2 - r^2 \left(\dot{\theta}^2 + \sin^2 \theta \dot{\phi}^2\right) = 1 , \quad (3.3.27)$$

where all coordinates are now functions of the particle proper time τ , with $\dot{f} = df/d\tau$, one finds

$$\delta S = m \delta \int_0^\tau \sqrt{2T} d\tau' = m \int_0^\tau \frac{\delta T}{\sqrt{2T}} d\tau' = m \delta \int_0^\tau T d\tau' . \quad (3.3.28)$$

The equations of motion can therefore be written as

$$\frac{d}{d\tau} \left(\frac{\partial T}{\partial \dot{x}^\mu} \right) = \frac{\partial T}{\partial x^\mu} , \quad (3.3.29)$$

¹⁵The advantage of an action principle is that it makes the effect of symmetries more apparent via Noether theorem.

thus avoiding to deal with the square root. Since T does not depend on t and ϕ (because of the existence of the Killing vectors \bar{K}_t for time translations and \bar{K}_ϕ for one of the three rotations), from Noether's theorem, we expect two integrals of motion, namely

$$E = \frac{m}{2} \frac{\partial T}{\partial \dot{t}} = m \left(1 - \frac{2 G_N M}{r} \right) \dot{t} , \quad (3.3.30)$$

which reduces to the proper mass m in the weak field limit $r \gg R_H$ and $\dot{t} = 1$, and

$$J = -\frac{m}{2} \frac{\partial T}{\partial \dot{\phi}} = m r^2 \sin^2 \theta \dot{\phi} , \quad (3.3.31)$$

which gives the angular momentum around the axis that defines the angle ϕ if we choose $\theta = \pi/2$. The latter choice is always possible, since the directions of the axes are arbitrary, due to the existence of the three Killing vectors \bar{K}_i , and $(\theta = \pi/2, \dot{\theta} = 0)$ clearly solves the corresponding equation of motion

$$\frac{d}{d\tau} \left(\frac{\partial T}{\partial \dot{\theta}} \right) = \frac{\partial T}{\partial \theta} \quad \Rightarrow \quad \ddot{\theta} = \sin \theta \cos \theta \dot{\phi}^2 . \quad (3.3.32)$$

Note that this implies that the particle motion occurs on a plane passing from the origin of the reference frame.

We are now left with just the equations of motion for $\phi = \phi(\tau)$ and $r = r(\tau)$, which we consider for the simple case of radial in-fall, namely $\phi = 0$ and constant. It is then easier to use Eq. (3.3.30) and write

$$2T = 1 = \left(1 - \frac{2 G_N M}{r} \right)^{-1} \left(\frac{E^2}{m^2} - \dot{r}^2 \right) . \quad (3.3.33)$$

This yields

$$\dot{r}^2 = \frac{2 G_N M}{r} - \left(1 - \frac{E^2}{m^2} \right) , \quad (3.3.34)$$

which we recall must be solved along with

$$\dot{t} = \left(1 - \frac{2 G_N M}{r} \right)^{-1} \frac{E}{m} . \quad (3.3.35)$$

The latter means that

$$\frac{dt}{d\tau} \sim \frac{1}{r - R_H} , \quad (3.3.36)$$

and diverges for a trajectory that approaches R_H . In particular, an asymptotically inertial observer placed at $r \gg R_H$ which measures the time t would see a particle falling (radially) take an infinite amount of (coordinate) time t to reach the Schwarzschild radius.

Upon deriving the above equation (3.3.34) with respect to τ , we obtain

$$\ddot{r} = -\frac{G_{\text{N}} M}{r^2} , \quad (3.3.37)$$

which is again Newton's law but is now valid for all values of $r > 0$. Note that this result does not imply there is no General Relativity correction to Newtonian trajectories: the radial coordinate r is *not* the proper distance from the source's centre and τ is not the any observer's time t . In fact, a better understanding of Eq. (3.3.37) can be obtained by considering two radially falling geodesics starting from the same sphere at $r = r(0) \equiv r_0$ and with the same initial velocity, so that they will measure the same proper time τ . The length of the arc connecting these two geodesics on the common sphere at the proper time τ is precisely proportional to $r = r(\tau)$. We can therefore conclude that Eq. (3.3.37) is proportional to the acceleration between such two geodesics, and measure the tidal forces “perpendicular” to the line of free fall.

3.3.3 General orbits

When the conserved quantity $J \neq 0$, Eq. (3.3.34) is replaced by

$$\dot{r}^2 = \frac{2 G_{\text{N}} M}{r} - \left(1 - \frac{E^2}{m^2}\right) - \left(1 - \frac{2 G_{\text{N}} M}{r}\right) \frac{4 J^2}{m^2 r^2} , \quad (3.3.38)$$

which will allow for bound orbits (with $\dot{r} = 0$ for $r = r_0 < \infty$) only for

$$1 - \frac{E^2}{m^2} = \frac{2 G_{\text{N}} M}{r_0} - \left(1 - \frac{2 G_{\text{N}} M}{r_0}\right) \frac{4 J^2}{m^2 r_0^2} . \quad (3.3.39)$$

Since we expect a planet orbits at a relatively large $r_0 > R_{\text{H}}$, and with a speed $v \ll 1$, the second term in the right hand side, $J \sim v^2$, is negligible, and the first term is smaller than one. One thus concludes that bound orbits for planets only exist if for

$$E < m . \quad (3.3.40)$$

Otherwise, one will have scattering trajectories for $E \geq m$. The quantity E can thus be viewed as the sum of the proper particle mass m and its (negative) gravitational potential energy.

General time-like geodesics with $J > 0$ can now be studied. One then usually expresses the radial coordinate $r = r(\tau)$ in terms of the angle $\phi = \phi(\tau)$ by means of the chain rule, and obtains

$$\left(\frac{dr}{d\phi}\right)^2 = \frac{G_{\text{N}} M m^2}{2 J^2} r^3 - \left(1 - \frac{E^2}{m^2}\right) \frac{m^2}{4 J^2} r^4 - \left(1 - \frac{2 G_{\text{N}} M}{r}\right) r^2 , \quad (3.3.41)$$

but we shall not go any further here. We just wish to mention that such an analysis explains the “*anomalous*” *precession of Mercury's perihelion*: in Newtonian gravity, even including

the effects of other planet, one determines a precession which is about 43'' (arc seconds) per century off (total precession is about 5600'' per century). Einstein then showed that this small discrepancy can be precisely accounted for in General Relativity by making use of the Schwarzschild metric. Historically, this is recognised as one of the (three) classical tests of General Relativity.

3.3.4 Light-like geodesics

Null trajectories can likewise be studied by means of the same action principle (3.3.26), but with $T = 0$, from which one can derive an analogue to Eq. (3.3.41). Again, we shall not go into details and just mention that such an equation could be used to describe the second of the three classical tests of General Relativity, namely the *deflection of light rays* around stellar sources. This effect was first “verified” by Eddington during a famous expedition to observe a solar eclipse in 1919.

The only result we will need in the following is the equation for radial geodesics, which simply follows from setting $ds^2 = d\theta^2 = d\phi^2 = 0$ in Eq. (3.3.23),

$$\frac{dr}{dt} = \pm \left(1 - \frac{2 G_N M}{r} \right) , \quad (3.3.42)$$

where the sign of course depends on whether the light ray is falling toward or climbing away from the central gravitational source.

3.3.5 Gravitational red-shift

Let us again consider the particular case of a radial geodesic, but this time for photons, which means $T = 0$ in Eq. (3.3.26). Without using the explicit form of the geodesic equation, one can already derive a simple and very general expression for the gravitational red-shift experienced by a photon which travels in a static space-time, such as it would be a light ray moving radially in the Schwarzschild metric. This result is particularly important because it will allow us to model a more realistic kind of observation. In fact, given the fact that the Schwarzschild metric departs from the flat Minkowski metric significantly for $r \simeq R_H$, it is very unlikely that we can place a static measuring apparatus there (see Footnote 18). What we could instead do more easily is to look at a in-falling particle from far away. This means receiving light signals from such a freely-falling probe which, as we shall see momentarily, are increasingly weakened.

First of all we note that, if there exists a time-like Killing vector $\bar{K}_{(t)}$, the “Killing energy”¹⁶

$$E = -K_{(t)}^\mu u_\mu , \quad (3.3.43)$$

¹⁶Properly speaking, this E is not the energy measured by any observer, unless the space-time is everywhere flat.

is conserved along geodesics. In fact, let u^μ be the 4-velocity of a particle which moves along a geodesic, then

$$-\frac{dE}{d\lambda} = u^\nu \nabla_\nu \left(K_{(t)}^\mu u_\mu \right) = K_{(t)}^\mu (u^\nu \nabla_\nu u_\mu) + u^\nu u^\mu (\nabla_\nu K_{(t)\mu}) = 0 , \quad (3.3.44)$$

where we used the geodesic equation and the definition of Killing vectors (2.5.79),

$$K_{\mu;\nu} + K_{\nu;\mu} = 0 , \quad (3.3.45)$$

and λ parameterizes the geodesic (for time-like geodesics, λ can be taken the proper time, for null geodesics it is a generic affine parameter). We then observe that the 4-velocity of a static observer in a static space-time must be proportional to the Killing vector $\bar{K}_{(t)}$, which implies that

$$U^\mu = \frac{K_{(t)}^\mu}{\sqrt{-K_{(t)}^\nu K_{(t)\nu}}} \equiv \frac{K_{(t)}^\mu}{|K_{(t)}|} , \quad (3.3.46)$$

since $U^\mu U_\mu = -1$, whereas $\bar{K}_{(t)}$ is not normalised. The energy measured by the static observer is therefore given by

$$\omega = -U^\mu u_\mu = -\frac{K_{(t)}^\mu u_\mu}{\sqrt{-K_{(t)}^\nu K_{(t)\nu}}} = \frac{E}{|K_{(t)}|} . \quad (3.3.47)$$

Let us then consider a photon which crosses two different static observers placed at $r = r_1$ and $r = r_2$ respectively. They will measure the photon's energies

$$\frac{\omega_1}{\omega_2} = \frac{|K_{(t)}(r_2)|}{|K_{(t)}(r_1)|} . \quad (3.3.48)$$

In particular, observers placed at constant r in a Schwarzschild space-time have 4-velocity

$$U^\mu = (\dot{t}(r), 0, 0, 0) = \dot{t} K_{(t)}^\mu , \quad (3.3.49)$$

and Eq. (3.3.48) yields

$$\frac{\omega_1}{\omega_2} = \frac{\dot{t}(r_1)}{\dot{t}(r_2)} = \sqrt{\frac{1 - R_H/r_2}{1 - R_H/r_1}} < 1 , \quad (3.3.50)$$

where we assumed the photon was emitted at $r = r_2$ and subsequently observed at $r = r_1 > r_2$. Static observers therefore see the photon lose energy as it “climbs up” the gravitational potential of the Schwarzschild space-time. This effect is the subject of the third classical test of General Relativity, the *Pound-Rebka experiment*, performed in 1959 inside the Jefferson tower at Harvard University.

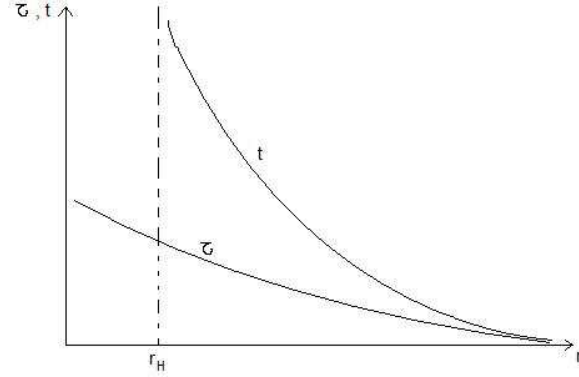


Figure 3.2: Radial geodesic in Schwarzschild t and proper time τ .

Let us conclude by mentioning that the above effect is not the entire story regarding the redshift of a signal emitted by a probe falling toward the Schwarzschild radius. In fact, if ω_s is the frequency emitted by the probe, as measured by an observer located on the probe itself, to the effect (3.3.50) one must add (meaning, multiply) the Doppler effect (1.3.43) due to the falling observer moving with a given velocity, say $v_2 = v(r_2)$, with respect to the static observer at the point of emission $r = r_2$,

$$\omega_1 \simeq \omega_s \sqrt{\frac{1-v_2}{1+v_2}} \sqrt{\frac{1-R_H/r_2}{1-R_H/r_1}}. \quad (3.3.51)$$

The Doppler contribution is arbitrary, in a sense, since the velocity v depends on the specific initial conditions of the probe's trajectory, whereas the gravitational contribution (3.3.62) is uniquely determined by the geometry. For this reason, one usually omits the former in discussing the gravitational redshift.

3.3.6 Radially infalling probe

We first remark that Eq. (3.3.34) allows for radial trajectories that fall from $r \gg R_H$ to $r \leq R_H$ in a *finite* amount of proper time τ . However, as we mentioned above, due to Eq. (3.3.35), the static asymptotic observer at $r_o \gg R_H$ would appear to see such a trajectory approach R_H and never cross it (see Fig. 3.2). One could associate with the reference frame $\{t, r, \theta, \phi\}$ a “static observer” defined by a set of (distinguished) clocks placed at fixed values of the areal radius r and synchronised in such a way that the angles between fixed vectors and their modulus remains constant in time (which explains the meaning of a time-independent metric¹⁷). Such an observer would experimentally determine the trajectory of the in-falling particle by recording the particle's subsequent positions r and related times t . The data so collected could then be plotted as the line $t = t(r)$ of Fig. 3.2. Note however that t is *not* the proper time T of an observer placed at constant r (as well as θ and ϕ), since the metric

¹⁷Of course, for a general synchronisation, the Schwarzschild metric would not appear time-independent.

element $g_{tt} \neq 1$ implies that

$$\frac{dT}{dt} = \sqrt{1 - \frac{2G_N M}{r}} \sim \left(1 - \frac{2G_N M}{r}\right)^{-1/2} \frac{d\tau}{dt}. \quad (3.3.52)$$

We can also consider an observer comoving with the particle, namely a clock sitting on the particle itself. The line $\tau = \tau(r)$ would then be built from the data collected by this observer, which could record the “positions” r by looking at the set of distinguished clocks of the static observer and annotating its own (proper) times. This makes it clear that in both cases we should build a (extremely unrealistic) measuring apparatus (the static observer) which extends over a large volume of space. Note also that the two observers would move with increasing relative velocity $v = v(r) = \frac{dr}{dt}$, and one therefore expects a relativistic time dilation with increasing

$$\gamma = \frac{dt}{d\tau} = \frac{1}{\sqrt{1 - v^2}}. \quad (3.3.53)$$

If we push this argument towards the Schwarzschild radius R_H , it seems that $\gamma \rightarrow \infty$ there, and $v \rightarrow 1$ (the speed of light). However this conclusion relies on the use a specific observer (the static one), which might not be physically realisable¹⁸. We shall therefore need a better (and more realistic) way to describe the physics we can see for $r \gtrsim R_H$.

Instead of insisting in building up reference frames connected with unrealistic observers, let us just describe explicitly what is going to happen to a probe sent radially toward the centre of our system, as it would be perceived by a static observer placed precisely at the distance $r = r_o \gg R_H$ from which the probe is dropped with zero initial velocity. This probe will emit a signal of frequency ω_s at fixed intervals $\Delta\tau_s$. Since the probe is a locally inertial observer, we can assume both ω_s and $\Delta\tau_s$ are constant, and the points of emission are therefore given by $S_n = (\tau_n = n \Delta\tau, r_n)$, with $n = 0, 1, 2, \dots$. We want to determine the points of detection of the signals, that is $O_n = (t_n, r_o)$, and the corresponding frequencies ω_n as measured by the asymptotic observer at $r = r_o$.

The relevant equations for the in-falling probe’s trajectory are given by

$$\frac{dr_s}{d\tau} = -\sqrt{\frac{R_H}{r_s} - 1 + \frac{E^2}{m^2}} \quad \Leftrightarrow \quad \frac{d^2 r_s}{d\tau^2} = -\frac{R_H}{2r_s^2}, \quad (3.3.54)$$

where $E^2/m^2 = 1 - R_H/r_0 \simeq 1$, and

$$\frac{dt}{d\tau} = \left(1 - \frac{R_H}{r_s}\right)^{-1} \frac{E}{m} \simeq \left(1 - \frac{R_H}{r_s}\right)^{-1}. \quad (3.3.55)$$

The out-going trajectory of the light signal is instead governed by

$$\frac{dr_\gamma}{dt} = 1 - \frac{2G_N M}{r_\gamma}, \quad (3.3.56)$$

¹⁸Due to the gravitational pull of the central source, it is likely such an observer would need a very powerful rocket to stay at constant r . One should then compute how powerful such a rocket should be for $r \rightarrow R_H$, and compare the required energy with the mass of the source.

and its frequency ω_n received at $r = r_o$ changes according to

$$\omega_n = \omega'_n \sqrt{\frac{1 - R_H/r_n}{1 - R_H/r_o}} \simeq \omega'_n \sqrt{1 - \frac{R_H}{r_n}} , \quad (3.3.57)$$

where ω'_n is the frequency measured by a (fictitious) static observer placed at the point of emission, $r = r_n$.

Let us first determine the times of detection. First of all, Eq. (3.3.55) implies that the interval $\Delta\tau$ translates into a difference

$$\Delta t_n^{(1)} = \frac{\Delta\tau}{1 - R_H/r_n} , \quad (3.3.58)$$

which depends on the emission point $r_s = r_n$. To the above difference, we need to add the difference in travel time for the light between two successive emissions. In fact, during the proper time $\Delta\tau$, the probe falls a radial difference

$$\begin{aligned} \Delta r_n \equiv r_n - r_{n-1} &= \int_{\tau_{n-1}}^{\tau_n} \frac{dr_s}{d\tau} d\tau \\ &= - \int_{\tau_{n-1}}^{\tau_n} \sqrt{\frac{R_H}{r_s(\tau)} - 1 + \frac{E^2}{m^2}} d\tau \\ &\simeq - \int_{\tau_{n-1}}^{\tau_n} \sqrt{\frac{R_H}{r_s(\tau)}} d\tau \simeq \sqrt{\frac{R_H}{r_n}} \Delta\tau , \end{aligned} \quad (3.3.59)$$

where we used the geodesic equation (3.3.54), and assumed the interval $\Delta\tau$ is short compared to the typical rate of change of the radial coordinate. The difference in the time of travel between the two points of emission and the asymptotic observer is thus given by

$$\begin{aligned} \Delta t_n^{(2)} \equiv t_n - t_{n-1} &= \int_{r_n}^{r_o} \frac{dt}{dr_\gamma} dr_\gamma - \int_{r_{n-1}}^{r_o} \frac{dt}{dr_\gamma} dr_\gamma = - \int_{r_{n-1}}^{r_n} \frac{dt}{dr_\gamma} dr_\gamma \\ &= - \int_{r_{n-1}}^{r_n} \left(1 - \frac{R_H}{r_\gamma}\right)^{-1} dr_\gamma \\ &\simeq - \left(1 - \frac{R_H}{r_n}\right)^{-1} \Delta r_n \simeq \sqrt{\frac{R_H}{r_n}} \frac{\Delta\tau}{1 - R_H/r_n} , \end{aligned} \quad (3.3.60)$$

where we used the null geodesic equation (3.3.56) and Eq. (3.3.59). Adding up the two results, we obtain

$$\Delta t_n = \Delta t_n^{(1)} + \Delta t_n^{(2)} \simeq \left(1 + \sqrt{\frac{R_H}{r_n}}\right) \frac{\Delta\tau}{1 - R_H/r_n} , \quad (3.3.61)$$

which implies that the observer will receive less and less pulses per unit time as the probe falls down. Eventually, as the probe gets close to the Schwarzschild radius, the observer will have to wait an asymptotically infinite amount of time in between pulses.

Secondly, having chosen a static observer at $r = r_o$, we can use Eq. (3.3.57) to describe the total redshift for the emitted signals,

$$\omega_n \simeq \omega_s \sqrt{\frac{1-v_n}{1+v_n}} \sqrt{1 - \frac{R_H}{r_n}} , \quad (3.3.62)$$

where we can estimate the velocity with respect to the locally static observer in the Doppler effect as

$$v_n = \left. \frac{dr_s}{dT} \right|_{t_n} = \left. \frac{dr_s}{d\tau} \frac{d\tau}{dt} \frac{dt}{dT} \right|_{t_n} \simeq \sqrt{\frac{R_H}{r_n} \left(1 - \frac{R_H}{r_n} \right)} . \quad (3.3.63)$$

In particular, for $r_s \rightarrow R_H$ one finds $v_n \rightarrow 0$, and the observed frequency vanishes precisely according to the gravitational redshift formula (3.3.50). Since any physically realisable receiver has a lower minimum sensitivity ω_c , there is always going to be a radius $r_s > R_H$ such that $\omega_n < \omega_c$, thus lower than the minimum sensitivity. At that point, the probe would simply “black out” and become invisible to any static observer.

3.3.7 The (event) horizon and black holes

The gravitational redshift in Eq. (3.3.62) implies that a photon emitted near the Schwarzschild radius would spend all of its energy to escape. This leads us to interpret the Schwarzschild sphere as an *event horizon*. The precise physical nature of the Schwarzschild sphere could in fact be fully understood by studying light cones starting at different areal radii. One would then discover that, whereas at $r > R_H$ there exists both in-going (contracting) and out-going (expanding) light cones, for $r = R_H$ the out-going light cone is stuck at $r = R_H$ (which is therefore a null surface) and for $r < R_H$ it also contracts. This is the defining property of an *apparent horizon* or *trapping horizon*: no signal, including light, can escape from within it. This concept is the General Relativity realization of an older conjecture made in Newtonian gravity by Michell and Laplace in the 18th century: by simply equating the Newtonian escape velocity to the speed of light,

$$\frac{1}{2} m v_\infty^2 = m \frac{G_N M}{r} , \quad \text{with } v_\infty = 1 , \quad (3.3.64)$$

one finds a limiting mass M (independent of m) above which even a signal travelling at the speed of light cannot escape from a star with given radius r_s . Of course, it is questionable that the above derivation makes sense for $m = 0$ ¹⁹, but it is quite interesting that this conjecture exactly leads to $r_s = R_H$, the Schwarzschild radius of the star of mass M . Our previous analysis of geodesic motion and Eq. (3.3.37) in particular, already clarifies the coincidence, as it should be clear that the coordinate r used to describe the proper geodesic motion is *not* the same as the radius r in the Newtonian argument above. It is in fact the analysis of the probe in Section 3.3.6 that gives Michell and Laplace’s old idea a proper status in General Relativity.

¹⁹If one takes the limit $m \rightarrow 0$ before solving Eq. (3.3.64), one does not obtain any escape velocity.

If the horizon remains static, it is then called the *event horizon* because the region inside it will never be able to communicate with the outer region. The inner region was named a *black hole* by J. A. Wheeler in 1968²⁰, and we now know only a very limited variety may (mathematically) exist. They include the spherically-symmetric but electrically charged Reissner-Nordstroem metric (found in 1916-18), the axially-symmetric and rotating Kerr metric (found in 1963), and the electrically charged and rotating Kerr-Newman metric (finally discovered in 1965). A common feature of all these solutions of the Einstein equations is the existence of one (or more) horizons, approaching which photon frequencies are red-shifted to zero. It is important to remark this is a purely kinematical effect. In fact, tidal forces (that is, geodesic deviation, in the General Relativity jargon) described by the Riemann tensor remain small for $r \sim R_H$ and only diverge for $r \rightarrow 0$ (thus known as the *real singularity*, whereas the horizon is also named a *coordinate singularity* since it can be removed by a suitable change of coordinates). Rotating black hole metrics also display a frame-dragging effect, so that the space-time appears to be dragged by the angular momentum of the source. In particular, if one carries a vector (for example, a spinning top) along a spatially closed geodesic (in the sense that certain coordinate positions are the same at the beginning and end of the trajectory), the vector will appear rotated with respect to its initial direction. This effect is extremely small (of the order of 10^{-13}), but has been recently tested around our planet by the Probe-B satellite.

In a dynamical situation, there in general exist *dynamical horizons* which evolve in time and may (or not) give rise to an event horizon. We remark the former are defined by the local causal structure (that is, the light cones around a point), whereas the latter is global in nature: whether a gravitational system possesses an event horizon or not requires the knowledge of the whole space-time. In particular, for a collapsing body (like a supernova) one, in principle, needs to know the entire future of the remnant body after the initial explosion.

Black hole space-times would however be just a mathematical curiosity if they were not realized in nature. In the 40's, Oppenheimer and collaborators described very simple models of collapsing spheres of dust which ended into forming black holes, thus showing that similar objects might be the final outcome of supernovae explosions. Further, the same authors and Chandrasekhar later provided limiting masses (of order $3 M_\odot$) above which neutron stars would not be stable but collapse to a point-like singularity. In particular, Chandrasekhar found in the early '30 that a collapsing star will not produce a stable *white dwarf* if its mass exceeds the "Chandrasekhar limit" of about $1.5 M_\odot$. The gravitational attraction will in fact be strong enough to force electrons to merge protons and give rise to a *neutron star*, kept together by the quantum mechanical pressure of the degenerate Fermi gas of neutrons. Neutron stars are the best candidate for *pulsars*. Later, in 1939, Oppenheimer and Volkoff employed the results of Tolman and further found that neutron stars will not be stable if their mass exceeds $3 M_\odot$. At that point, there is no (known) force which could prevent the collapse of matter to a point. Currents estimates of such limits vary, but the general picture

²⁰The term "frozen star" was largely used previously, with a clear reference to the fact that a distant observer would never see the surface of a collapsing star cross the Schwarzschild radius. Also, the term "black hole" had been previously used by a journalist in the early 60's.

stands and the existence of black holes in our universe is widely accepted. For example, astronomers have found evidence of very large black holes (with $M \sim 10^6 M_\odot$) at the centre of galaxies, including our Milky Way.

Finally, let us mention that, in the early 70's, Hawking discovered that black holes actually emit particles like black bodies, as a quantum effect, albeit at a very small effective temperature (smaller than the typical temperature of the CMB radiation). This result makes black holes one of the most interesting (theoretical) arenas for General Relativity and quantum physics, with a possible conceptual clash between the two, which produced many interesting speculations about the possible quantum theory of gravity.

3.4 Cosmology

The other case of interest is cosmology which, quite remarkably, like the study of black holes and General Relativity overall, was for a long time after the initial excitement looked upon as a mathematical (if not merely metaphysical) subject [13].

Modern cosmology is mainly based on the following

Copernican Principle: “We are not a preferred observer in the Universe.”

In other words, it is reasonable to assume that the Universe would look to any other observers like it looks to us. From the practical point of view, this principle is of limited use. However, and although it goes a long way ahead to infer from the above, one eventually relies on the

Cosmological Principle (CP): “The Universe is homogeneous and isotropic.”

Isotropy is here taken as an observational statement, whereas homogeneity then follows from assuming isotropy is independent of the observation point according to the Copernican principle.

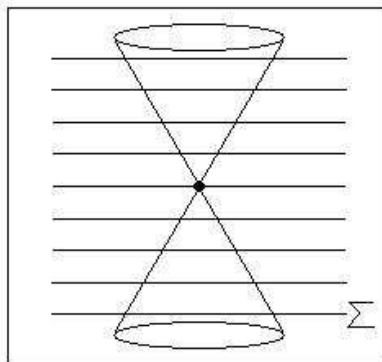


Figure 3.3: Past and future light-cones originating from us now. Σ_t are hypersurfaces of constant time t (to be defined).

Let us first review these (apparently) simple symmetries. Looking at the sky at night, it is pretty obvious that it does not appear particularly isotropic: our solar system is pretty

much empty space with a few planets and asteroids, and one can clearly see scattered stars forming Constellations and clusters, and even a strip we call the Milky Way. Nonetheless, one would like to think that, could we detect all the matter in the Universe, the average distribution (on a suitably sampled area of the night sky) is the same in all directions. In brief, isotropy is therefore “assumed” as much as it is observed. It is further worth noting that what we see in the sky is not the Universe at a given instant of time, but the image of it produced by light-cones that reach us at the time of observation (see Fig. 3.3). Saying that the Universe is homogenous and isotropic therefore means that there exists a time t such that the Universe is homogenous and isotropic on each time slice Σ_t . The matter content on each Σ_t then affects the light propagation, which shows that the entire construction needs to be experimentally self-supported. Moreover, a signal traveling along light cones may have generated at different times Δt in the past from different distances Δs , and we therefore need a separate way of determining either Δt or Δs in order to verify specific models.

The CP may therefore be taken as a working assumption upon which explicit models of the Universe are built and (hopefully) verified *a posteriori*. In particular, we expect there is a minimum scale above which the Universe appears homogeneous, but in the following we shall consider an idealized view in which *galaxies form a homogeneous fluid filling up the entire space*.

3.4.1 Friedman-Robertson-Walker metric

Like with the case of the Schwarzschild metric, the form of the cosmological metric can be partly fixed by assuming the existence of Killing vectors. In particular, we will now have three space-like Killing vectors generating spatial translations (which mathematically defines *homogeneity*), and three space-like Killing vectors generating rotations (which mathematically defines *isotropy*). It can also be proven that isotropy with respect to an arbitrary point is equivalent to homogeneity, and that there may not be any further isometry on a three-dimensional (space-like) foliation of the space-time ²¹, but we shall not go into these details. It is just worth pointing out that we have no time-like Killing vector since we want to describe an evolving Universe, since in 1929 E. Hubble and M. Humason observed the farer galaxies recedes from us faster and faster.

Homogeneity and isotropy uniquely identifies the FRW (Friedmann, Robertson and Walker) metric

$$ds^2 = -dt^2 + a^2(t) \left[\frac{dr^2}{1 - k r^2} + r^2 (d\theta^2 + \sin^2 \theta d\phi^2) \right], \quad (3.4.1)$$

where the origin $r = 0$ is totally arbitrary, t is the proper time of an observer moving along with the homogenous and isotropic cosmic fluid (the idealized representation of galaxies) at

²¹This property is called *maximal symmetry*.

r , θ and ϕ constant. We shall also call

$$\begin{array}{ll} \{t, r, \theta, \phi\} & \text{comoving coordinates} \\ a(t) & \text{(cosmic) scale factor} \\ k = 0, \pm 1 & \text{curvature constant.} \end{array} \quad (3.4.2)$$

Note that k could in fact take any real value, however, if $k \neq 0$, a suitable rescaling of r and a ²²,

$$\begin{aligned} dr &\rightarrow \sqrt{\pm k} dr \\ a^2 &\rightarrow k^2 a^2 , \end{aligned} \quad (3.4.3)$$

will always allow to have $k = \pm 1$ [but no rescaling allows to change between the three integer values in (3.4.2)]. Depending on the value of the curvature scalar, we can introduce new coordinates such that the topology of the hypersurface Σ_t is apparent from the line element $d\sigma^2$:

- **Flat Universe:** for $k = 0$, the coordinate r is very similar to the usual radial coordinate in \mathbb{R}^3 ,

$$d\sigma^2 = dr^2 + r^2 d\Omega^2 = dx^2 + dy^2 + dz^2 , \quad (3.4.4)$$

and Σ_t is flat (zero spatial curvature).

- **Closed Universe:** for $k = +1$, the proper radius $R_{(3)}$ is bounded from above and

$$r = \sin(X) \quad \Rightarrow \quad d\sigma^2 = dX^2 + \sin^2(X) d\Omega^2 , \quad (3.4.5)$$

and Σ_t is a three-dimensional sphere.

- **Open Universe:** for $k = -1$, one can write

$$r = \sinh(\Psi) \quad \rightarrow \quad d\sigma^2 = d\Psi^2 + \sinh^2(\Psi) d\Omega^2 , \quad (3.4.6)$$

and Σ_t is a three-dimensional hyperboloid.

Finally, the meaning of the coordinate r is very different from that we used in the Schwarzschild space-time. If we write

$$ds^2 = -dt^2 + a^2(t) d\sigma^2 , \quad (3.4.7)$$

we see that the areal radius in FRW is given by

$$r_A = a(t) r , \quad (3.4.8)$$

²²In this Section we assume r is dimensionless, with t and a being lengths.

and the area of surfaces of constant r therefore depends on time. Likewise, the proper distance between two points is given by

$$dR = a(t) \frac{dr}{\sqrt{1 - k r^2}} = a(t) dR_{(3)} , \quad (3.4.9)$$

where $R_{(3)}$, the rescaled proper distance on Σ_t , can be bounded.

Observations suggest that the distance between galaxies increases in time, whereas their typical size remains the same. We can therefore claim that the Universe is expanding, with the farer galaxies moving faster away from us, like dots on an inflating balloon ²³. This picture can be mathematically modeled by a modified FRW metric which locally (around matter sources such as a galaxy) looks like the Schwarzschild metric: local lengths are mostly affected by the localized sources and do not appreciably change in time, whereas the distance between sources increases because of the increasing scale factor ²⁴.

3.4.2 Cosmic fluids

As we wrote above, we assume the Universe is filled with a (perfect) fluid of matter and energy. Its energy-momentum tensor then takes the form

$$T^\mu_\nu = \text{diag}(-\rho, p, p, p) , \quad (3.4.10)$$

where

$$\rho = \rho(t) \quad \text{and} \quad p = p(t) , \quad (3.4.11)$$

and satisfies the continuity equation $\nabla_\mu T^\mu_\nu = 0$. The 00-component of this equation yields energy conservation ²⁵,

$$-\nabla_\mu T^\mu_0 = \dot{\rho} + 3H(\rho + p) = 0 , \quad (3.4.12)$$

where

$$H = \frac{\dot{a}}{a} \quad (3.4.13)$$

is the so-called *Hubble constant* (rather improperly, since it is not at all constant in general) and $\dot{a} = da/dt$. In fact, one can rewrite Eq. (3.4.12) as

$$\frac{d}{dt}(a^3 \rho) = -p \frac{d}{dt}(a^3) , \quad (3.4.14)$$

²³Alternative scenarios were proposed, in which the Universe is stationary and matter is produced continuously during the expansion so as to keep the average density constant.

²⁴This picture is still being debated, and is the topic of the so-called Einstein-Straus problem in General Relativity.

²⁵Note the Christoffel symbols do not vanish for the FRW metric.

or, on noting that the spatial volume $V \sim a^3$, one can integrate over a cubic cell and obtain

$$\frac{dE}{dt} = - \sum_i F_i \frac{dx^i}{dt} , \quad (3.4.15)$$

where dx^i is the displacement of the i^{th} face of the cube.

Now, let us assume an equation of state for the fluid,

$$p = \omega \rho , \quad (3.4.16)$$

where ω is a constant. Energy conservation then reads

$$\frac{\dot{\rho}}{\rho} = -3(1+\omega) \frac{\dot{a}}{a} \quad \Rightarrow \quad \rho \propto a^{-3(1+\omega)} . \quad (3.4.17)$$

The simplest components of cosmic fluids are given by *dust* (pressureless matter, or non-relativistic matter almost exactly at rest with the cosmic frame) and *radiation* (massless matter, or highly-relativistic matter).

- **Dust:** in this case no force is present, beside gravity, and $\omega = 0$ (so that $p = 0$). Eq. (3.4.17) therefore yields

$$\rho_{\text{dust}} = \frac{E}{V} \propto a^{-3} , \quad (3.4.18)$$

which agrees with the fact that the proper mass of dust particles, $E \sim m_0$, is an invariant and the volume element scales like

$$V \propto a \times a \times a . \quad (3.4.19)$$

One can consider dust particles (stars, galaxies, etc.) as being located at fixed r , θ and ϕ , or the chosen reference frame as *comoving* with the cosmic fluid. Moreover, since dust particles are only subject to gravity, lines of constant r , θ and ϕ are also geodesics.

- **Radiation:** since mass is totally negligible, so is the trace of the energy-momentum tensor ²⁶.

$$T = (-\rho + 3p) = 0 . \quad (3.4.20)$$

We then find

$$p = \frac{1}{3} \rho \quad \Rightarrow \quad \omega = \frac{1}{3} , \quad (3.4.21)$$

²⁶Recall that the trace of a $(0,2)$ tensor is invariant under rotations and this result is lifted to any reference frame in General Relativity. For a plane wave moving along the x axis, the energy-momentum tensor is $T^\mu_\nu = \text{diag}[-R, p, 0, 0]$, where E is the energy and p the momentum, and the “mass-shell” condition $E = p$ implies $T = 0$.

and

$$\rho_{\text{radiation}} = \frac{E}{V} \propto a^{-4} . \quad (3.4.22)$$

This result can be understood by noting that the volume scales again like in Eq. (3.4.19), and photon energy redshifts according to

$$E \propto a^{-1} . \quad (3.4.23)$$

Of course, it makes no sense to consider the chosen reference frame as comoving with the photons. The only possible definition of the reference frame in use is then provided by the CP, or that the coordinates are such that $\rho_{\text{radiation}} = \rho_{\text{radiation}}(t)$ and $p_{\text{radiation}} = p_{\text{radiation}}(t)$.

For a long time it was thought that we now live in a matter (dust)-dominated Universe, whereas in the early stages, the Universe dynamics was controlled by radiation, since the density of the latter increases faster (going backward in time). We now know that the Universe expansion is presently accelerating ($\ddot{a} > 0$), which is not compatible with the effect of dust.

- **Vacuum or dark energy:** Among possible sources, we may also include a fluid with equation of state

$$\rho = -p = \frac{\Lambda}{8\pi G_{\text{N}}} , \quad \omega = -1 , \quad \rho_{\Lambda} \propto 1 , \quad (3.4.24)$$

where Λ is the famous *cosmological constant* first introduced by Einstein, who later defined it his biggest mistake (but is now necessary to explain the current accelerated expansion of the Universe).

3.4.3 Friedmann equations

The specific form of the FRW metric (3.4.1) reduces the Einstein equations to just two Friedmann equations,

$$G_{00} = 8\pi G_{\text{N}} T_{00} \quad \Rightarrow \quad 3 \left[\left(\frac{\dot{a}}{a} \right)^2 + \frac{k}{a^2} \right] = 8\pi G_{\text{N}} \rho \quad (3.4.25)$$

\Downarrow

$$G_{ii} = 8\pi G_{\text{N}} T_{ii} \quad \Rightarrow \quad 3 \frac{\ddot{a}}{a} = -4\pi G_{\text{N}} (\rho + 3p) . \quad (3.4.26)$$

The first one, Eq. (3.4.25), is technically a constraint, which selects the possible combinations of initial conditions $a(t_0) = a_0$ and $\dot{a}(t_0) = \dot{a}_0$ for the truly dynamical (second order) Eq. (3.4.26) for the scale factor $a = a(t)$, given a specific matter content. However, the constraint is preserved at all times, as can be seen by deriving Eq. (3.4.25) with respect

to time and using the continuity Eq. (3.4.17) to obtain Eq. (3.4.26). In details, the time derivative of Eq. (3.4.25) yields

$$6 \left(\frac{\ddot{a}}{a} - \frac{\dot{a}^2}{a^2} - \frac{k}{a} \right) = 8 \pi G_N \dot{\rho} . \quad (3.4.27)$$

We can then replace the second and third term inside the brackets using Eq. (3.4.25) and express the time derivative of the density using Eq. (3.4.17), that is

$$3 \frac{\ddot{a}}{a} - 8 \pi G_N \rho = -12 \pi G_N (\rho + p) , \quad (3.4.28)$$

which is Eq. (3.4.26). To summarise, for a fluid satisfying the continuity Eq. (3.4.17), it is easier to just solve for the constraint (3.4.25) at all times $t \geq t_0$.

We further define

$$\begin{aligned} q &= -\frac{a \ddot{a}}{\dot{a}^2} && \text{deceleration parameter} \\ \Omega &= \frac{8 \pi G_N}{3 H^2} \rho = \frac{\rho}{\rho_{\text{critical}}} && \text{density parameter} \end{aligned} \quad (3.4.29)$$

where

$$\rho_{\text{critical}} = \frac{3 H^2}{8 \pi G_N} . \quad (3.4.30)$$

The Friedmann equation (3.4.25) can then be written as

$$\Omega - 1 = \frac{k}{H^2 a^2} , \quad (3.4.31)$$

and the following conclusions can be drawn:

- $\rho < \rho_{\text{critical}} \quad \Leftrightarrow \quad \Omega < 1 \quad \Leftrightarrow \quad k = -1 \quad \Leftrightarrow \quad \text{Open Universe}$
- $\rho = \rho_{\text{critical}} \quad \Leftrightarrow \quad \Omega = 1 \quad \Leftrightarrow \quad k = 0 \quad \Leftrightarrow \quad \text{Flat Universe}$
- $\rho > \rho_{\text{critical}} \quad \Leftrightarrow \quad \Omega > 1 \quad \Leftrightarrow \quad k = +1 \quad \Leftrightarrow \quad \text{Closed Universe}$

The spatial curvature k then determines the evolution of the scale factor (see Fig. 3.4, which displays the cosmic evolution for dust, and similar behaviours would also occur for radiation).

Observations suggest that our Universe is very close to $k = 0$. For a flat, matter dominated Universe, one has

$$\left\{ \begin{array}{l} \rho_{\text{dust}} \sim \frac{1}{a^3} \\ a \dot{a}^2 \sim 1 \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \frac{\dot{a}^2}{a^2} \sim \frac{1}{a^3} \\ \sqrt{a} da \sim dt \end{array} \right. \Rightarrow a^{\frac{3}{2}} \sim t . \quad (3.4.32)$$

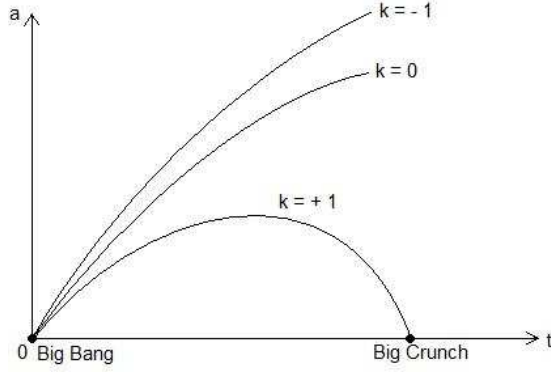


Figure 3.4: Evolution of $a = a(t)$ for $k = 0, \pm 1$.

For a flat, but radiation dominated Universe,

$$\begin{cases} \rho_{\text{rad}} \sim \frac{1}{a^4} & \Rightarrow & \frac{\dot{a}^2}{a^2} \sim \frac{1}{a^4} & \Rightarrow & a^2 \sim t . \\ a^2 \dot{a}^2 \sim 1 & \Rightarrow & a \, da \sim dt \end{cases} \quad (3.4.33)$$

Finally, for a flat and empty Universe, with only a positive vacuum energy present, one obtains the exact solution

$$\begin{cases} \rho_{\text{vacuum}} \sim \Lambda & \Rightarrow & \frac{\dot{a}^2}{a^2} \sim \frac{\Lambda}{3} \\ \sqrt{\frac{\Lambda}{3}} \sim \frac{\dot{a}}{a} = H_0 \end{cases} \quad \Rightarrow \quad a \sim e^{H_0 t} \quad (\text{de Sitter Universe}) , \quad (3.4.34)$$

where H_0 is now a true cosmological *constant*.

It is finally interesting to compare the Friedmann equation (3.4.26) with what would be predicted by the Newtonian theory of gravity. Let us consider the motion of a point-like particle of mass m located on the surface of a sphere of homogeneous density ρ and radius $R = r a$. The radial Newtonian acceleration this particle is subjected to would be

$$\ddot{R} = r \ddot{a} = -\frac{G_N (4 \pi \rho / 3) R^3}{R^2} = -\frac{4}{3} \pi G_N \rho r a , \quad (3.4.35)$$

or

$$3 \frac{\ddot{a}}{a} = -4 \pi G_N \rho , \quad (3.4.36)$$

which coincides with Eq. (3.4.26) only for $p = 0$. In other words, the pressure does not gravitate according to the Newtonian theory, but does so according to the Einstein theory.

3.4.4 Cosmic Microwave Background

Many observations have confirmed that the Universe is filled with an almost homogeneous relic radiation (CMB). We believe this radiation was generated in the very early times, when matter and radiation decoupled (and photons could therefore start to travel freely) on the *surface of last scattering*. Looking back at the energy densities (3.4.18) and (3.4.22), we see that the Universe must have been much denser and hot during its early stages. At those high energies, the mean free path of photons was very short since they have enough energy to produce pairs of (oppositely charged) particles, and photons were in (approximate) thermal equilibrium with electrons and positrons (among others). As the photon energy decreased below the threshold for electron pair productions (corresponding to an average energy of 1 MeV or temperature of $3 \cdot 10^3$ K), photons became essentially free and those are the oldest light signals we can detect now, at a temperature of about 3 K. Only gravitational waves, if they exist, could have decoupled far earlier and arrive to us from earlier times.

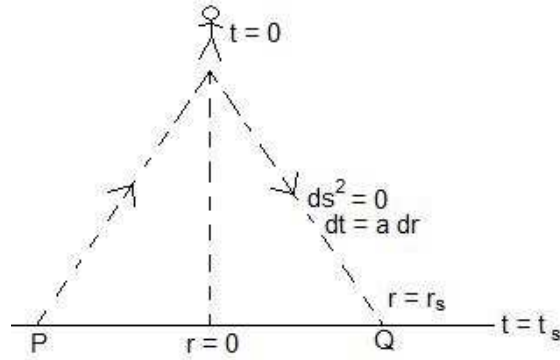


Figure 3.5: Signals coming from opposite directions in the sky.

In light of the above description of the early Universe, the homogeneity (actually, the isotropy from our point of observation) of the CMB is surprising. Suppose we look along two opposite directions in the sky. The light we receive now from those directions will have originated from very distant places, and one then wonders how such points could have been at the same temperature. In fact, the two points could have never been in causal contact before (no signal may have yet travelled between them; see Fig. 3.5). From the metric (3.4.1), we find that light-cones are defined by the equation

$$ds^2 = 0 \quad \Rightarrow \quad dt = a \, dr \quad \Rightarrow \quad dr = \frac{dt}{a} . \quad (3.4.37)$$

Suppose we place ourselves at $r = 0$ and integrate the above expression (along the light-cone) from $t = -t_s$ to now ($t = 0$). We thus find the comoving radial coordinate of the point of origin,

$$r_s = \int_0^{r_s} dr \sim \int_{-t_s}^0 \frac{dt}{a(t)} . \quad (3.4.38)$$

If the Universe is either matter or radiation dominated, we have

$$a(t) \sim t^\alpha, \quad 0 < \alpha < 1 \quad \Rightarrow \quad r_s \sim t_s^{1-\alpha}. \quad (3.4.39)$$

The proper distance travelled by that photon is thus

$$R \sim a(-t_s) r_s = t_s^\alpha t_s^{1-\alpha} = t_s, \quad (3.4.40)$$

which, quite remarkably, coincides with the Minkowskian result in flat space. One also has

$$\dot{a} \sim t^{\alpha-1} \quad \Rightarrow \quad H = \frac{\dot{a}}{a} \sim t^{-1}, \quad (3.4.41)$$

and the so-called *particle horizon*,

$$R_H \sim \frac{1}{H(-t_s)} \sim t_s, \quad (3.4.42)$$

grows with time (more of the Universe comes into causal contact with a given observer).

Let us now note that in a vacuum dominated Universe, we instead have

$$a(t) \sim e^{H_0 t} \quad \Rightarrow \quad r_s \sim \frac{e^{-H_0 t}}{H_0}, \quad (3.4.43)$$

and the particle horizon therefore appears to remain at the same distance from the central observer,

$$R \sim e^{H_0 t_s} r_s \sim \frac{1}{H_0}. \quad (3.4.44)$$

The latter result can explain the CMB homogeneity: the Universe started out very small, enough so that all of its parts had time to come into causal contact. It then underwent an early phase of rapid (almost exponential) expansion, called *inflation*, during which the initial state of matter was almost frozen. The CMB originated after the end of inflation, which explains why we do not yet see the entire Universe, but the CMB is homogeneous.

3.4.5 Cosmological redshift

We would like to assess which type of Universe we live in and the value of H_0 by direct observations. For this purpose, we can derive the cosmological redshift as a time-dilation between two peaks of the same wave, or the emission of two successive signals moving at the speed of light.

Let us denote with t_1 the time of the first emission and with $t_1 + \delta t_1$ the time of the second emission. The two signals will be received, respectively, at the time t_2 and $t_2 + \delta t_2$, after having travelled the same coordinate distance, that is

$$\int_{t_1}^{t_2} \frac{dt}{a(t)} = \int_{r_1}^{r_2} dr = \int_{t_1 + \delta t_1}^{t_2 + \delta t_2} \frac{dt}{a(t)}, \quad (3.4.45)$$

which implies

$$\int_{t_1}^{t_1+\delta t_1} \frac{dt}{a(t)} = \int_{t_2+\delta t_2}^{t_2+\delta t_2} \frac{dt}{a(t)} . \quad (3.4.46)$$

We can now assume, like with the probe falling in Schwarzschild, that the times δt_1 and δt_2 are short enough that the cosmic factor does not change appreciably. Eq. (3.4.46) then yields

$$\frac{\delta t_1}{a(t_1)} = \frac{\delta t_2}{a(t_2)} , \quad (3.4.47)$$

which immediately gives

$$\frac{\omega(t_1)}{\omega(t_2)} = \frac{a(t_2)}{a(t_1)} \equiv \frac{a_2}{a_1} . \quad (3.4.48)$$

It is customary to express this wavelength change in terms of the redshift

$$z = \frac{\lambda_o - \lambda_s}{\lambda_s} = \frac{a_o}{a_s} - 1 , \quad (3.4.49)$$

where the subscript o means the quantity is taken at the time of observation and s at the (earlier) time of emission. Unlike the Doppler effect, this redshift is not caused by the relative motion of emitter and observer, but by the space-time expansion, and can be directly measured.

It is important to remark now that the value of $a = a(t)$ at a given instant of time is not physically meaningful, since a can always be rescaled by an arbitrary constant. However, the ratio $a(t_1)/a(t_2)$, for any two times $t_1 \neq t_2$, is instead measurable, in principle, and by means of Eq. (3.4.49), also in practice. It also gives us a way of measuring distances (indirectly).

3.4.6 Luminosity-distance relation

In astronomy, measuring distances is of course not trivial, but one can measure the apparent luminosity of an object. One method to estimate distances is then to use the *luminosity-distance relation*, denoted by d_L , for specific light sources (stars, galaxies, cluster of galaxies, etc).

To explain this better, let us denote by F the flux of energy (energy E per unit time T and area A) measured by an observer and first consider Minkowski space-time. Since the energy E is conserved during light propagation (there is no gravitational redshift in this case), the total energy that crosses any concentric sphere per unit time $L = \frac{E}{T}$ does not depend on the sphere radius and equals the intrinsic luminosity of the source, $L_0 = E_0/T_0$. The flux observed on a portion of unit area of this sphere will then be

$$F = \frac{L}{A} = \frac{L_0}{4\pi R^2} \quad \Rightarrow \quad d_L^2 = R^2 = \frac{L_0}{4\pi F} , \quad (3.4.50)$$

which is the trivial luminosity-distance relation for flat space-time, with d_L simply equal to the sphere's radius.

In a FRW space-time, photons are redshifted by a factor of $(1+z)$ while they propagate, according to Eqs. (3.4.48) and (3.4.49). Moreover, if we assume the the cosmic evolution does not affect the microscopic mechanisms by which light is emitted, the frequency at which the observer registers the arrival of photons is likewise reduced with respect to the (previous) frequency at which they were emitted. In fact, let δt be the time between two “discrete” emissions from the source. In an expanding Universe, the observer will (later) detect these two subsequent signals a time $(1+z)\delta t$ apart. We therefore have that $L \simeq L_0/(1+z)^2$ and the measured flux

$$F = \frac{L}{A} \simeq \frac{L_0}{4\pi (a_0 r)^2 (1+z)^2} \equiv \frac{L_0}{4\pi d_L^2}, \quad (3.4.51)$$

where a_0 is the scale factor at the time of observation, and $a_0 r$ the proper radius of the sphere centred on the source and upon which the observer is placed. From the above, one immediately obtains

$$d_L = a_0 r (1+z), \quad (3.4.52)$$

which can therefore be used to determine the cosmic scale factor from the measurement of z and $d_L \simeq \sqrt{L_0/F}$.

In order to measure the redshift, we need to know the original frequency of the light emitted from the source. Luckily, most astrophysical sources show clear spectral bands of emission or absorption. We also need to know the intrinsic luminosity L_0 of the source. For this purpose, one can use variable stars which show a specific relation between the period of oscillation of their (apparent) luminosity and the absolute luminosity (defined as the intrinsic luminosity measured from a standard distance). For larger distances, one can instead use galaxies with similar properties. Altogether, these preferred sources are thereby called *standard candles*, and form the so-called *cosmic distance ladder*. Estimating their proper intrinsic luminosity is therefore very important, since any error would introduce a systematic bias in the measurement of distances across the universe.

3.4.7 Hubble law

Having built the reference model for the evolution of the Universe, we can now derive the famous Hubble law.

We have seen that a photon is gravitationally redshifted according to

$$z = \frac{a_0}{a_s} - 1, \quad (3.4.53)$$

where a_0 and a_s are the scale factors at the time of detection ($t = t_0$, today) and emission $t = t_s$, respectively. Then, along a null ray, and for $k r^2 \ll 1$, we easily find

$$ds^2 = -dt^2 + \frac{a^2 dr^2}{1 - k r^2} = 0 \quad \Rightarrow \quad \int_t^{t_0} \frac{dt'}{a(t')} = \int_0^r \frac{dr'}{(1 - k r'^2)^{1/2}} \simeq \int_0^r dr'. \quad (3.4.54)$$

Upon expanding the cosmic factor for the emission time $t_s = t$ around t_0 ,

$$a(t) = a_0 - \dot{a}_0 (t - t_0) + \frac{1}{2} \ddot{a}_0 (t - t_0)^2 + \dots , \quad (3.4.55)$$

we then obtain

$$r = \frac{1}{a_0} \left[(t_0 - t) + \frac{1}{2} H_0 (t_0 - t)^2 + \dots \right] , \quad (3.4.56)$$

and the redshift is

$$1 + z = \frac{a_0}{a_s} = 1 + H_0 (t - t_0) - \frac{1}{2} q_0 H_0^2 (t - t_0)^2 + \dots , \quad (3.4.57)$$

where the deceleration parameter today is given by

$$q_0 = -\frac{a_0 \ddot{a}_0}{\dot{a}_0^2} = \frac{1 + 3\omega}{2} \Omega . \quad (3.4.58)$$

For small values of $H_0 (t - t_0)$, we can write

$$t_0 - t = \frac{1}{H_0} \left[z - \left(1 + \frac{q_0}{2} \right) z^2 + \dots \right] . \quad (3.4.59)$$

Replacing the above into the expression for r , we finally obtain the *Hubble law*

$$d_L = \frac{1}{H_0} \left[z + \frac{1}{2} (1 - q_0) z^2 + \dots \right] \simeq \frac{z}{H_0} , \quad (3.4.60)$$

that is, (not too far) galaxies recede from us with a velocity $v \sim z$ (directly) proportional to the distance d_L . The constant of proportionality is the Hubble constant H_0 , whose inverse is therefore representative of the age of the Universe.

3.4.8 The Universe today

We need to remind us that the distance that appears in the Hubble relation (3.4.60) is in fact the luminosity-distance (3.4.52) obtained from the measured flux F and estimated intrinsic luminosity L_0 by means of the Eq. (3.4.51). And the “velocity” z is not quite the same quantity that determines the Doppler effect in flat space, but the cosmological redshift (3.4.49). This said, from separate measurements of d_L and z one can deduce the values of H_0 and q_0 in Eq. (3.4.60).

These observations, along with others that we shall not discuss, have led us to picture the current status of the Universe as spatially flat, with $\Omega \simeq 1$, corresponding to an average density

$$\rho_0 = \rho_{\text{critical}} \simeq 10^{-29} \text{ g/cm}^3 , \quad (3.4.61)$$

equivalent to about 6 protons per square cubic meter. In particular, three different sources have been identified to contribute to ρ_0 :

- **Regular massive matter**, well approximated by a dust fluid, and estimated through the luminosity of galaxies in the cosmos,

$$\frac{\rho_{\text{matter}}}{\rho_0} \simeq 5\% , \quad (3.4.62)$$

corresponding to about 1 proton per 4 square cubic meters.

- **Dark matter**, which again behaves like dust but is not directly detected,

$$\frac{\rho_{\text{DM}}}{\rho_0} \simeq 25\% . \quad (3.4.63)$$

The existence of dark matter is required, for example, in order to explain how stars rotate inside galaxies, but its nature is not clear yet.

- **Dark energy**, with the equation of state of the vacuum,

$$\frac{\rho_{\text{DE}}}{\rho_0} \simeq 70\% , \quad (3.4.64)$$

and which is required by the present negative value of the deceleration parameter $q_0 < 0$ (thus $\ddot{a}_0 > 0$).

Explaining dark matter and dark energy are in fact the two biggest puzzles in present cosmology.

Appendix A

Symmetries and group theory

Lorentz transformations are a special case of linear coordinate transformations. Such transformations have special properties, which may be useful both for computing and conceptual purposes [6].

A.1 Abstract groups

Mathematical group: let G be a set of objects for which a binary operation is defined (we shall mostly use “multiplicative notation”)

$$\forall g_1, g_2 \in G, \quad \exists g_3 \in G : \quad g_1 \cdot g_2 = g_3 . \quad (\text{A.1.1})$$

The couple (G, \cdot) forms a group if the following properties hold

1) \cdot is associative

$$(g_1 \cdot g_2) \cdot g_3 = g_1 \cdot (g_2 \cdot g_3) = g_1 \cdot g_2 \cdot g_3 ; \quad (\text{A.1.2})$$

2) there exists a neutral element \mathbb{I} (identity) such that

$$g \cdot \mathbb{I} = \mathbb{I} \cdot g = g , \quad \forall g \in G ; \quad (\text{A.1.3})$$

3) all elements have an inverse

$$\forall g \in G, \quad \exists g^{-1} \in G : \quad g \cdot g^{-1} = g^{-1} \cdot g = \mathbb{I} . \quad (\text{A.1.4})$$

4) the group is *Abelian* if

$$g_1 \cdot g_2 = g_2 \cdot g_1 , \quad \forall g_1, g_2 \in G . \quad (\text{A.1.5})$$

The “additive notation” goes as follows: the operation is denoted by

$$\forall g_1, g_2 \in G, \quad \exists g_3 \in G : \quad g_1 + g_2 = g_3 . \quad (\text{A.1.6})$$

and the defining properties become

1') $+$ is associative

$$(g_1 + g_2) + g_3 = g_1 + (g_2 + g_3) = g_1 + g_2 + g_3 ; \quad (\text{A.1.7})$$

2') there exists a neutral element 0 (zero) such that

$$g + 0 = 0 + g = g , \quad \forall g \in G ; \quad (\text{A.1.8})$$

3') all elements have an opposite

$$\forall g \in G , \quad \exists (-g) \in G : \quad g + (-g) = (-g) + g = 0 . \quad (\text{A.1.9})$$

The above definition has *a priori* nothing to do with transformations and is therefore more general. For example, the prototype of multiplicative groups is the set of rational numbers \mathbb{Q} , whereas the prototype of additive groups is the set of integer numbers \mathbb{Z} . What characterizes a (formal) group is the formal set of elements and the operation between them. If the elements of two groups can be put in correspondence in such a way that the corresponding operations also yield corresponding results, then the two groups are *formally the same*. On the other hand, the same formal group may be realized in different ways. For example, the formal group \mathbb{Z} can be realized by any set of elements that can be added (and subtracted) like apples and money (provided we define the negative of an apple by the need of one and the negative of money as a debt).

A.2 Matrix representations and Lie groups

Linear changes of coordinates in a N -dimensional space can be represented by matrices with real (or complex) entries. A particular example of multiplicative groups is thus given by square matrices of constant numbers with non-vanishing determinant. Such group is called $GL(N)$ for *General Linear* in N dimensions. Note that each matrix is defined by $N \times N - 1$ elements (the -1 coming from the determinant condition). The operation defined for matrices is the usual matrix multiplication,

$$\sum_j A_{ij} B_{jk} = C_{ik} . \quad (\text{A.2.1})$$

Elements of this group naturally act on N -dimensional (real or complex) vectors V^i , $i = 1, \dots, N$ (the *fundamental representation*). One can however consider the same group $GL(N)$ acting on objects other than vectors. For example, the action of $GL(N)$ on scalars f is simply represented by the multiplication by 1: $f \rightarrow 1 \cdot f$. For $(2, 0)$ tensors T_{ij} , we note that they have N^2 components, and the action of $GL(N)$ on them, rearranged so as to form a N^2 -dimensional vector, must be realized by $N^2 \times N^2$ matrices (and so on and so forth). These explain the need to distinguish between the formal group $GL(N)$ and its (many) realizations. And also that different formal groups may share some realizations (the action of all multiplicative groups on scalars is the multiplication by 1).

Given a class of transformations, one therefore has a formal group (G, \cdot) , which entails all of their properties, and a map from it to its *realization* $(D(G), \times)$,

$$\forall g_i \in G, \quad \exists D(g_i) \in D(G), \quad (\text{A.2.2})$$

such that

$$D(g_1) \times D(g_2) = D(g_1 \cdot g_2), \quad (\text{A.2.3})$$

where the elements $D(g)$ are usually matrices and act upon vectors V , but one can have different structures (see below for examples). Notice also that the order of factors is crucial and we used different symbols for the formal multiplication and multiplication between elements of the realizations, since they are *different* concepts (in the following we shall instead use the same symbol for notational simplicity).

Let us also recall that a vector space \mathcal{V} is a set of objects we can add together and multiply by real (or complex) numbers (scalars):

$$\forall V_1, V_2 \in \mathcal{V} \text{ and } a, b \in \mathbb{R} (\mathbb{C}) \Rightarrow a V_1 + b V_2 = V_3 \in \mathcal{V}. \quad (\text{A.2.4})$$

A fundamental property of \mathcal{V} is that there exists a (finite or infinite) *basis* of elements V_i which linearly generate all of \mathcal{V} . The number of basis elements is the *dimension* of the vector space.

Note that by *representation* of a group one actually means the set \mathcal{V} of objects $D(g)$ acts upon, where $D(G)$ is a realization of the group G . For example, for the Lorentz transformations, the matrices (1.4.59) are the realization of the Lorentz transformation on the vector representation $\mathcal{V} = \{V^\mu\}$. It is then natural to assume the element obtained by means of the operation (A.2.1) is still a coordinate transformation, since the above mathematical formula simply means we apply the transformation B to a given vector followed by the transformation A on the resulting vector.

Groups can have finite or infinite number of elements. An example of a finite group is given by *parity* transformation,

$$P : x \rightarrow -x, \quad (\text{A.2.5})$$

with the property

$$P^2 = \mathbb{I}, \quad (\text{A.2.6})$$

so that the inverse exists and $P = P^{-1}$. The elements of the parity group are therefore just $\{P, \mathbb{I}\}$.

Of particular interest are the infinite groups with a finite number of *generators*, also known as *Lie groups*. Roughly speaking, a Lie group is a group whose elements can be continuously parametrized by a finite set of real variables $\theta_i \in \mathbb{R}$, with $i = 1, 2, \dots, d$, in such a way that one can write any element of G (or, equivalently, of a realisation of G) in the exponentiated form

$$D(g) \leftrightarrow g = e^{\sum \theta_i J_i} = \sum_{n=0}^{\infty} \frac{(\sum_{i=1}^d \theta_i J_i)^n}{n!}, \quad (\text{A.2.7})$$

where J_i are the generators of G . It follows that the set $\mathcal{G} = \{j\}$ of all linear combinations of the J_i 's must be endowed with three operations in order to recover the group multiplication and make sense of Eq. (A.2.7):

- i) an associative operation $+$ between elements $j \in \mathcal{G}$;
- ii) a multiplicative operation \times (usually omitted) by (real or complex) scalars and
- iii) a multiplicative operation \cdot between elements $j \in \mathcal{G}$.

This means that $(\mathcal{G}, +, \times, \cdot)$ form an *algebra* (the *Lie algebra* of the group G): $(\mathcal{G}, +, \times)$ is a (real or complex) vector space and (\mathcal{G}, \cdot) is a group. Further, multiplication \cdot and addition $+$ are mutually compatible, meaning they satisfy the distribution property

$$j_1 \cdot (j_2 + j_3) = j_1 \cdot j_2 + j_1 \cdot j_3 , \quad \forall j_i \in \mathcal{G} . \quad (\text{A.2.8})$$

It is also conventional to set $\theta_i = 0$ for the identity of G : $g(0) = \mathbb{I}$.

For general matrices in $GL(N)$, the following important (Baker-Campbell-Hausdorf) formula holds

$$e^A e^B = e^{A+B+[A,B]/2+\dots} \neq e^{A+B+[B,A]/2+\dots} = e^B e^A \quad \text{iff} \quad [A, B] \neq 0 , \quad (\text{A.2.9})$$

where

$$[A, B] = A B - B A , \quad (\text{A.2.10})$$

is the *commutator*. The above matrix property implies

$$D(g_1) \cdot D(g_2) = e^{\sum \theta_i^{(1)} J_i} e^{\sum \theta_j^{(2)} J_j} \simeq e^{\sum_i \theta_i^{(1)} J_i + \sum_j \theta_j^{(2)} J_j + \sum_{ij} \theta_i^{(1)} \theta_j^{(2)} [J_i, J_j]} . \quad (\text{A.2.11})$$

This means that, if the Lie group is not Abelian, $J_i J_j \neq J_j J_i$ (otherwise the algebra is also called *Abelian*), one must have

$$J_i J_j - J_j J_i = [J_i, J_j] = \sum_k c_{ij}{}^k J_k , \quad (\text{A.2.12})$$

where $c_{ij}{}^k$ are the *structure constants* of the algebra, otherwise the product of two (or more) J_i 's would not belong to the algebra, nor would the product on the left hand side of Eq. (A.2.11) belong to the (realisation of the) group G . In fact, the above commutation relations imply that products of elements of \mathcal{G} are linear combinations of the generators, as required by the fact \mathcal{G} is a vector space generated by the J_i . The number d of (linearly independent) J_i is called the *dimension* of the Lie algebra and group. Eq. (A.2.12) uniquely specifies a Lie algebra, meaning that two algebras with the same commutator structure are *the same* mathematical object.

Two groups G_1 and G_2 whose algebras are equal are the same group (at least) *near the identity* (or *locally*). However, the two groups can be *globally* different. For example, the parameters θ_i may have different ranges, like $\theta_i \in \mathbb{R}$ for G_1 and $|\theta_i| < 1$ for G_2 (which is then a compact group – this concept requires the notion of manifold and will be clarified later in the course).

Another important concept is that of *irreducible representations*. A representation of a group G is said irreducible if the corresponding realization $D(G)$ cannot be put in block diagonal form,

$$D(g) = \begin{bmatrix} D_1(g) & 0 & 0 & \dots \\ 0 & D_2(g) & 0 & \dots \\ \dots & & & \end{bmatrix} . \quad (\text{A.2.13})$$

Note that if such a block diagonal form exists, since each $D_i(G)$ is a realization of G , the corresponding representation must be given by a vector space that is the cartesian product of separate vector spaces

$$\mathcal{V} = \mathcal{V}_1 \times \mathcal{V}_2 \times \dots , \quad (\text{A.2.14})$$

so that the D_i act on elements of \mathcal{V}_i . An important result is that *all the representations of a group can be build out of irreducible ones*.

Before considering the Lorentz group, let us see the simpler case of rotations.

A.3 Rotations in N dimensions

We already saw that the defining equations for rotations is

$$R^T R = R^T \mathbb{I} R = \mathbb{I} = \begin{bmatrix} 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ \dots & & & \end{bmatrix} , \quad (\text{A.3.1})$$

so that the Cartesian scalar product is invariant

$$\sum_{i=1}^N x_i y_i = \sum_{i=1}^N x'_i y'_i . \quad (\text{A.3.2})$$

It follows from the above definition that the inverse of any element R exists and coincides with the transposed matrix, $R^{-1} = R^T$. It is also easy to verify that the property (A.3.1) is preserved by the matrix product, so that these matrices indeed form a (in general non-Abelian) group denoted by $O(N)$ for *Orthogonal* in N dimensions.

From (A.3.1) it also follows that

$$1 = \det(R^T R) = (\det R)^2 \quad \Rightarrow \quad \det R = \pm 1 . \quad (\text{A.3.3})$$

A particular case is thus given by orthogonal matrices with positive unit determinant,

$$\det(R) = 1 , \quad (\text{A.3.4})$$

which is denoted by $SO(N)$ (*Special Orthogonal*). From (A.3.1) one obtains (note that $[J_i, J_i] = 0$ and $\theta_i^T = \theta_i$)

$$\mathbb{I} = (e^{\sum \theta_i J_i})^T e^{\sum \theta_k J_k} = e^{\sum \theta_i (J_i^T + J_i)} \quad \Rightarrow \quad J_i^T = -J_i . \quad (\text{A.3.5})$$

which tells us the generators J_i are realized by *skew-symmetric* matrices. From (A.3.4), one likewise obtains

$$1 = \det(e^{\sum \theta_i J_i}) = e^{\theta_i \text{tr}(J_i)} \quad \Rightarrow \quad \text{tr}(J_i) = 0 , \quad (\text{A.3.6})$$

that is, the generators are *traceless*. Of course, this second condition is not new and in fact follows from the previous (A.3.5) (but not the other way around), which is reminiscent of the argument in Eq. (A.3.3).

A.3.1 Rotations in 2 dimensions: $SO(2)$ and $U(1)$

Let us start with rotations in 2 dimensions, which are defined by

$$R^T R = R^T \mathbb{I} R = \mathbb{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} . \quad (\text{A.3.7})$$

The fundamental *realization* of this group is given by the 2×2 matrices

$$R(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} , \quad (\text{A.3.8})$$

whose determinant is also equal to one,

$$\det \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} = \cos^2 \theta + \sin^2 \theta = 1 , \quad (\text{A.3.9})$$

and we identify

$$R(2n\pi + \theta) = R(\theta) , \quad n \in \mathbb{N} . \quad (\text{A.3.10})$$

This means the group $SO(2)$ is *compact*, since its Lie parameter $\theta \in (0, 2\pi)$.

One can easily check that $(R(\theta), \cdot)$ form a group:

$$\begin{aligned} R(\theta_1) R(\theta_2) &= \begin{bmatrix} \cos \theta_1 & \sin \theta_1 \\ -\sin \theta_1 & \cos \theta_1 \end{bmatrix} \begin{bmatrix} \cos \theta_2 & \sin \theta_2 \\ -\sin \theta_2 & \cos \theta_2 \end{bmatrix} \\ &= \begin{bmatrix} \cos \theta_1 \cos \theta_2 - \sin \theta_1 \sin \theta_2 & \cos \theta_1 \sin \theta_2 + \sin \theta_1 \cos \theta_2 \\ -\sin \theta_1 \cos \theta_2 + \cos \theta_1 \sin \theta_2 & -\sin \theta_1 \sin \theta_2 + \cos \theta_1 \cos \theta_2 \end{bmatrix} \\ &= \begin{bmatrix} \cos(\theta_1 + \theta_2) & \sin(\theta_1 + \theta_2) \\ -\sin(\theta_1 + \theta_2) & \cos(\theta_1 + \theta_2) \end{bmatrix} = R(\theta_1 + \theta_2) . \end{aligned} \quad (\text{A.3.11})$$

Further,

$$\mathbb{I} = \begin{bmatrix} \cos 0 & \sin 0 \\ -\sin 0 & \cos 0 \end{bmatrix} = R(0) , \quad (\text{A.3.12})$$

and

$$R^{-1}(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}^{-1} = \begin{bmatrix} \cos(-\theta) & \sin(-\theta) \\ -\sin(-\theta) & \cos(-\theta) \end{bmatrix} = R(-\theta) = R^T(\theta) . \quad (\text{A.3.13})$$

The matrices R cannot be put in block diagonal forms: this means the 2-vectors $V^i \in \mathbb{R}^2$ are an irreducible representation of $SO(2)$ ¹. In fact this is the fundamental *representation* of the group $SO(2)$, and all the other representations can be built from it. For example,

$$\begin{aligned} \underbrace{T^{ij}}_{2 \times 2 = 4} = V^i W^j &= \underbrace{\frac{1}{2} V_k W^k \delta^{ij}}_1 + \underbrace{\frac{1}{2} (V^i W^j + V^j W^i - V_k W^k \delta^{ij})}_2 + \underbrace{\frac{1}{2} (V^i W^j - V^j W^i)}_1 \\ &= S + T^{(ij)} + T^{[ij]} , \end{aligned} \quad (\text{A.3.14})$$

is the $(2, 0)$ -tensor representation of $SO(2)$, which one can show reduces to a combination of one scalar (the trace), the traceless symmetric and skew-symmetric parts. It is indeed easy to see that the trace of T^{ij} is invariant under rotation, since

$$\text{Tr} (R^T T R) = \text{Tr} (R R^T T) = \text{Tr} (T) . \quad (\text{A.3.15})$$

Moreover, by contracting $T^{(ij)}$ [or $T^{[ij]}$] with twice the same rotation, one equally obtains a symmetric (antisymmetric) matrix. Such properties are not peculiar to $N = 2$ but extends to all dimensions.

Note that T^{ij} contains 4 free entries. The matrix S has 1 entry, $T^{(ij)}$ has 2 entries and can therefore be mapped into a 2-vector, $T^{[ij]}$ has 1 entry and is a (pseudo)-scalar, for a total of 4 independent entries, as it should. Suppose we arrange these 4 elements of T^{ij} into a vector V^A with $A = 1, \dots, 4$. The action of $SO(2)$ on such vectors should be given by a 4×4 matrix with block diagonal form

$$M^A_B V^B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \bar{R} \end{bmatrix} \begin{pmatrix} S \\ T^{[ij]} \\ \bar{T}^{(ij)} \end{pmatrix} , \quad (\text{A.3.16})$$

where the matrices M^A_B therefore realize $SO(2)$ and the vectors V^A represent it. In particular, if we write the symmetric and traceless part as

$$T^{(ij)} = \frac{1}{2} \begin{bmatrix} V^1 W^1 - V^2 W^2 & V^1 W^2 + V^2 W^1 \\ V^1 W^2 + V^2 W^1 & V^2 W^2 - V^1 W^1 \end{bmatrix} \equiv \begin{bmatrix} b & a \\ a & -b \end{bmatrix} , \quad (\text{A.3.17})$$

we can then map it into the 2-vector

$$T^{(ij)} \rightarrow \bar{T}^{(ij)} \equiv \begin{bmatrix} b \\ a \end{bmatrix} , \quad (\text{A.3.18})$$

and the matrix $\bar{R} = R(2\theta)$. This result is somewhat expected: if we rotate vectors by an angle θ , the product of two vectors will rotate twice the same angle (and so on).

¹One might notice however that rotations do not mix vectors of different norm, and one should therefore expect that “normalised” vectors (or *rays*) can be treated like forming their own vector space provided the linear combinations are suitably modified.

Clearly, there is only one parameter in $SO(2)$, the angle of rotation θ , and the group is therefore *one-dimensional* and Abelian. One can see this by determining the one generator J of the algebra $so(2)$ by means of a Taylor expansion of (A.3.8) about $\theta = 0$,

$$\begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \theta \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} + \frac{\theta^2}{2} \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} + \dots = \sum_n \frac{\theta^n}{n!} J^n \quad (\text{A.3.19})$$

where

$$J = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad \text{with} \quad J^2 = -\mathbb{I}, \quad (\text{A.3.20})$$

which is skew-symmetric, obviously commutes with itself and is traceless as required.

There is another one-dimensional Lie group whose algebra has the same property, namely $U(1)$, the group of complex numbers with module equal to one (the *Unitary* group)

$$z = e^{i\theta} = e^{\theta i}, \quad \theta \in \mathbb{R}, \quad (\text{A.3.21})$$

with the usual multiplication

$$z_1 z_2 = e^{\theta_1 i} e^{\theta_2 i} = e^{(\theta_1 + \theta_2) i}, \quad (\text{A.3.22})$$

and the property that

$$z^* z = 1. \quad (\text{A.3.23})$$

To this group one can associate a formal generator $J = i$ with the property that

$$z^{-1} z = z^* z = e^{-\theta i + \theta i} = 1, \quad i^* = -i \quad \text{and} \quad i^2 = -1. \quad (\text{A.3.24})$$

There is therefore a mathematical equivalence between the algebras $so(2)$ and $u(1)$ given by interpreting θ as an angle. Note that this realization of the group $U(1)$ does not naturally involve operators acting on any vector space ² and does not have a naive representation (in classical physics!), unlike the group $SO(2)$ which is realized by matrix transformations and is represented by 2-vectors.

A.3.2 Rotations in 3 dimensions: $SO(3)$ and $SU(2)$

This is the group generated by rotations around each of the three cartesian axes, namely

$$\begin{aligned} R_1 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \theta_1 & -\sin \theta_1 \\ 0 & \sin \theta_1 & \cos \theta_1 \end{bmatrix}, & R_2 &= \begin{bmatrix} \cos \theta_2 & 0 & \sin \theta_2 \\ 0 & 1 & 0 \\ -\sin \theta_2 & 0 & \cos \theta_2 \end{bmatrix}, \\ R_3 &= \begin{bmatrix} \cos \theta_3 & -\sin \theta_3 & 0 \\ \sin \theta_3 & \cos \theta_3 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \end{aligned} \quad (\text{A.3.25})$$

²Of course, one can consider the product $z v \in \mathbb{C}$ as the action of $z \in U(1)$ on the vector space \mathbb{C} .

all of which have $\det R_i = 1$ and, on a “trivial” side, note that

$$R_i(2n\pi + \theta) = R_i(\theta) , \quad n \in \mathbb{N} , \quad (\text{A.3.26})$$

which shows that $SO(3)$ is also compact. The corresponding fundamental representation is given by 3-vectors $V^i \in \mathbb{R}^3$. Since the above rotation matrices cannot be simultaneously put in block diagonal form, \mathbb{R}^3 is the fundamental irreducible representation of $SO(3)$.

All other representations can be built out of vectors, like for $SO(2)$, and they are in general reducible. For example, the product of two vectors,

$$\underbrace{T^{ij}}_{3 \times 3 = 9} = V^i W^j = \underbrace{S}_1 + \underbrace{T^{(ij)}}_5 + \underbrace{T^{[ij]}}_3 \quad (\text{A.3.27})$$

“breaks” into a scalar S (the trace), a pseudo-vector V (the skew-symmetric part $T^{[ij]}$) and an irreducible $(2,0)$ tensor Q (the traceless symmetric part $T^{(ij)}$). Note that $\dim T = 9$, $\dim S = 1$, $\dim V = 3$ and $\dim Q = 5$, with $1 + 3 + 5 = 9$. A representation of such tensors as 9-dimensional vectors V then requires a realisation of $SO(3)$ by means of 9×9 matrices M of the following block diagonal form

$$M^A_B V^B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & R_{(3)} & 0 \\ 0 & 0 & R_{(5)} \end{bmatrix} \begin{pmatrix} S \\ V \\ Q \end{pmatrix} , \quad (\text{A.3.28})$$

where $R_{(3)}$ is a usual 3×3 rotation matrix in 3 dimensions and $R_{(5)}$ a suitable 5×5 matrix. Again, such examples of different realisations of $SO(3)$ prove the necessity of distinguishing between a formal group and its realisations.

Irreducible representations of $SO(3)$ are identified by an integer number $s = 0, 1, \dots$, with $s = 0$ for the trivial scalar representation ($\mathbb{I} = 1$ acting on elements of \mathbb{R}) and $s = 1$ for the (fundamental) vector representation. The integer s can then be related with the angular momentum of a spinning body, the corresponding realisation of $SO(3)$ being given by the operators that generate rotations. In particular, an object with $s = 0$ will always look the same regardless of the amount and direction of rotation; an object with $s = 1$ will appear the same after a rotation of $\theta_i = 2\pi$ around the i^{th} -axis; an object with $s = 2$ needs a rotation of $\theta_i = \pi$. In general, an object with a given $s \in \mathbb{N}$ requires a rotation of $\theta_i = 2\pi/s$ to return to the initial configuration.

From the above matrices (A.3.25) one obtains the (skew-symmetric and traceless) generators of $so(3)$, namely

$$J_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} , \quad J_2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix} , \quad J_3 = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} , \quad (\text{A.3.29})$$

with commutators (from now on, we shall employ Einstein summation convention on repeated indices)

$$[J_i, J_j] = \epsilon_{ij}^{k} J_k , \quad (\text{A.3.30})$$

where ϵ_{ijk} is the Levi-Civita symbol. For example,

$$J_1 J_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (\text{A.3.31})$$

and

$$J_2 J_1 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} , \quad (\text{A.3.32})$$

so that

$$[J_1, J_2] = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = J_3 . \quad (\text{A.3.33})$$

The same algebra (A.3.30) holds for the generators of $SU(2)$, the group of 2-dimensional unitary matrices

$$U^\dagger U = \mathbb{I} , \quad (\text{A.3.34})$$

with positive determinant. If we write (with $i = 1, 2, 3$ and $\theta_i \in \mathbb{R}$)

$$U = e^{-i\theta^i \sigma_i / 2} , \quad (\text{A.3.35})$$

we find (since $[\sigma_i, \sigma_i] = 0$ and $\theta^* = \theta$)

$$\mathbb{I} = U^\dagger U = e^{i\theta^i (\sigma_i^\dagger - \sigma_i) / 2} \Rightarrow (\sigma_i^T)^* = \sigma_i^\dagger = \sigma_i , \quad (\text{A.3.36})$$

and

$$1 = \det(U) = e^{-i\theta^i \text{tr}(\sigma_i)} \Rightarrow \text{tr}(\sigma_i) = 0 , \quad (\text{A.3.37})$$

which defines the Pauli matrices as the traceless Hermitian generators of $su(2)$,

$$\sigma_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} , \quad \sigma_2 = \begin{bmatrix} 0 & -i \\ i & 0 \end{bmatrix} , \quad \sigma_3 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} . \quad (\text{A.3.38})$$

Note that the matrices $J_i = -i\sigma_i/2$ satisfy Eq. (A.3.30), namely

$$[\sigma_i, \sigma_j] = 2i\epsilon_{ij}{}^k \sigma_k \Rightarrow \left[-i\frac{\sigma_i}{2}, -i\frac{\sigma_j}{2} \right] = \epsilon_{ij}{}^k \left(-i\frac{\sigma_k}{2} \right) , \quad (\text{A.3.39})$$

which manifests the correspondence between the algebras $su(2)$ and $so(3)$, meaning one can find common representations (that is, an equivalence between corresponding representations).

An explicit construction is the following: let us map 3-vectors $\vec{x} = (x, y, z)$ to 2×2 complex matrices (the indices $a, b = 1, 2$ in the following)

$$\vec{x} \leftrightarrow h_{ab}(\vec{x}) = \vec{x} \cdot \vec{\sigma} = x \sigma_1 + y \sigma_2 + z \sigma_3 = \begin{bmatrix} z & x - i y \\ x + i y & -z \end{bmatrix} . \quad (\text{A.3.40})$$

For transformations belonging to $SU(2)$, this is a particular $(0, 2)$ tensor, which must transform according to

$$h_{a'b'} = U_{a'}^c(\theta_i) U_{b'}^d(\theta_i) h_{cd} \quad \Leftrightarrow \quad h' = U^T h U , \quad (\text{A.3.41})$$

where U is given in Eq. (A.3.35) and one can then check that h' is equivalent to

$$\vec{x}' = R_i(\theta_i) \vec{x} , \quad (\text{A.3.42})$$

since $\vec{x}' \leftrightarrow h'$.

The fundamental representation of $SU(2)$ is given by 2-dimensional complex vectors (called *spinors*),

$$\psi = (z_1, z_2) , \quad z_a \in \mathbb{C} . \quad (\text{A.3.43})$$

We will see the role such vectors play in modern physics later. For now, let us just note that

$$\sigma_i^2 = 1 , \quad \forall i = 1, 2, 3 , \quad (\text{A.3.44})$$

which allows us to easily compute, for example,

$$\begin{aligned} U_3(\theta) &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + (-i\theta/2) \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} + \frac{(-i\theta/2)^2}{2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{(-i\theta/2)^3}{3!} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} + \dots \\ &= \begin{bmatrix} 1 + (-i\theta/2) + \frac{(-i\theta/2)^2}{2!} + \frac{(-i\theta/2)^3}{3!} + \dots & 0 \\ 0 & 1 + (i\theta/2) + \frac{(i\theta/2)^2}{2!} + \frac{(i\theta/2)^3}{3!} + \dots \end{bmatrix} \\ &= \begin{bmatrix} e^{-i\theta/2} & 0 \\ 0 & e^{+i\theta/2} \end{bmatrix} . \end{aligned} \quad (\text{A.3.45})$$

Note then that

$$U_3(2\pi) = \begin{bmatrix} e^{-i\pi} & 0 \\ 0 & e^{+i\pi} \end{bmatrix} = -1 . \quad (\text{A.3.46})$$

In fact, in order to rotate a *spinor* (A.3.43) back to its initial value, we need to go around the axis *twice*, $\theta = 4\pi$ (unlike a vector!).

Irreducible representations of $SU(2)$ are identified by one parameter s called the *spin*, which can only take (half)-integer values. The physical meaning of the (non-zero) spin is that $1/s$ equals the number of complete (2π) rotations needed to map the object into itself. Spin $1/2$ objects require 2 full rotations, spin 1 objects are vectors and require 1 full rotation, $(2, 0)$ tensors are spin 2 objects and require $1/2$ a full rotation (rotation of an angle $\theta = \pi$).

Before moving on to the Lorentz group, let us summarize our findings for $SO(3)$ and $SU(2)$:

- $SO(3)$ irreducible representations

$$s = 0, 1, 2, 3, \dots ; \quad (\text{A.3.47})$$

- $SU(2)$ irreducible representations

$$s = 0, \frac{1}{2}, 1, \frac{3}{2}, 2, \dots . \quad (\text{A.3.48})$$

We have explicitly seen that irreducible representations of the two groups with $s = 1$ can be put in correspondence by constructing the (traceless) 2-tensor (A.3.40) of $SU(2)$, which is equivalent to an $SO(3)$ vector. Since the latter defines the fundamental representation of $SO(3)$ out of which all representations of $SO(3)$ are built, it follows that each irreducible representation of $SO(3)$ is equivalent to an irreducible representation of $SU(2)$ with the same integer s . Of course, the other way around does not hold, since there is no equivalent of (half-integer) spinors in $SO(3)$.

In general, both for $SO(3)$ and $SU(2)$, the dimension of the representation is given by

$$d = 2s + 1 , \quad (\text{A.3.49})$$

as can be easily checked for scalars ($s = 0, d = 1$), spinors ($s = 1/2, d = 2$), vectors ($s = 1, d = 3$), etc. One also has the following composition rule for the tensor product of two irreducible representations (which generalizes the cases of the product of two vectors and two spinors we have seen before):

$$(s) \times (s') = (s - s') + (s - s' + 1) + \dots + (s + s') , \quad (\text{A.3.50})$$

where (s) denotes the irreducible representation of spin s and we assumed $s \geq s'$. We can easily check the above formula for the cases we saw before:

- 1) by composing two vectors like in Eq. (A.3.27), we find

$$(1) \times (1) = (0) + (1) + (2) , \quad (\text{A.3.51})$$

that is, a scalar, a (pseudo-)vector and an irreducible 2-tensor;

- 2) by composing two spinors like in Eq. (A.3.40), we find

$$(1/2) \times (1/2) = (0) + (1) , \quad (\text{A.3.52})$$

that is, a scalar and a vector.

As a final remark, let us mention a different realisation/representation of $SO(3) \sim SU(2)$: the angular momentum operators \hat{J}_i in quantum mechanics satisfy the same algebra as the J_i of $so(3)$. In fact, they are the generators of rotations in the Hilbert space of state vectors (wave-functions). However, the \hat{J}_i are *not* matrices and act on wave-functions, not on vectors of \mathbb{R}^3 .

A.4 Lorentz group: $SO(3, 1)$

We are now ready to study the Lorentz group of matrices which satisfy

$$\Lambda^T \eta \Lambda = \eta , \quad (A.4.1)$$

where the Minkowski metric tensor η has replaced \mathbb{I} for the rotation group and is symbolised by the notation $SO(4) \rightarrow SO(3, 1)$. We shall restrict ourselves to the so called *proper orthochronous* subgroup, which consists of matrices Λ that also satisfy

$$\det \Lambda = 1 , \quad \Lambda^0_0 \geq 1 , \quad (A.4.2)$$

the latter being an invariant condition (only affected by time reversal).

We first note that by introducing the following notation

$$\gamma = \frac{1}{\sqrt{1 - \beta^2}} = \cosh \phi , \quad \beta \gamma = \sinh \phi , \quad (A.4.3)$$

the boosts along the direction x^i can be realized by the 4×4 matrices

$$\begin{aligned} B_1 &= \begin{bmatrix} \cosh \phi_1 & \sinh \phi_1 & 0 & 0 \\ \sinh \phi_1 & \cosh \phi_1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} , & B_2 &= \begin{bmatrix} \cosh \phi_2 & 0 & \sinh \phi_2 & 0 \\ 0 & 1 & 0 & 0 \\ \sinh \phi_2 & 0 & \cosh \phi_2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} , \\ B_3 &= \begin{bmatrix} \cosh \phi_3 & 0 & 0 & \sinh \phi_3 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ \sinh \phi_3 & 0 & 0 & \cosh \phi_3 \end{bmatrix} , \end{aligned} \quad (A.4.4)$$

acting on 4-vectors $V^\mu \in \mathbb{R}^4$. The above matrices show a remarkable similarity with rotation matrices in 3 dimensions, except that the “angles” of rotation are “imaginary” ($\sin \rightarrow \sinh$ and $\cos \rightarrow \cosh$). In fact, the above matrices satisfy the defining equations

$$\cosh^2 \phi - \sinh^2 \phi = 1 \quad \Rightarrow \quad B \eta B^T = \eta , \quad \det B = 1 , \quad (A.4.5)$$

as well as the rotation matrices R_i do, where now

$$R_i = \begin{bmatrix} 1 & 0 \\ 0 & R_i^{(3)} \end{bmatrix} , \quad R_i^{(3)} \in SO(3) . \quad (A.4.6)$$

The set $\{B_i, R_i\}$ therefore represent a realisation of the Lorentz group $SO(3, 1)$. Further, since not all of them can be put in block diagonal form simultaneously, the 4-vectors are an irreducible representation of $SO(3, 1)$, in fact the fundamental one.

One can immediately obtain the corresponding generators, that is for the rotations

$$J_i = \begin{bmatrix} 0 & 0 \\ 0 & J_i^{(3)} \end{bmatrix} , \quad J_i^{(3)} \in so(3) , \quad (A.4.7)$$

and for the boosts

$$K_1 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad K_2 = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad K_3 = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}. \quad (\text{A.4.8})$$

Note that $[J_i, K_j] \neq 0$, and

$$[K_i, K_j] = -\epsilon_{ij}^{\quad k} J_k \quad (\text{A.4.9})$$

$$[J_i, J_j] = \epsilon_{ij}^{\quad k} J_k \quad (\text{A.4.10})$$

$$[J_i, K_j] = \epsilon_{ij}^{\quad k} K_k. \quad (\text{A.4.11})$$

In particular, from Eq. (A.4.10), we see that one can still consider rotations by themselves, since applying rotations in different orders leads to just more rotations, so to say. However, as soon as we wish to deal with the boosts, we cannot just consider their generators K_i , since rotations inevitably come into play according to Eq. (A.4.9).

In order to proceed with our analysis, we define

$$A_i = \frac{1}{2} (J_i + i K_i), \quad B_i = \frac{1}{2} (J_i - i K_i), \quad (\text{A.4.12})$$

for which one finds

$$\begin{aligned} [A_i, B_j] &= \frac{1}{4} ([J_i, J_j] - i [J_i, K_j] + i [K_i, J_j] + [K_i, K_j]) \\ &= \frac{1}{4} (\epsilon_{ij}^{\quad k} J_k - i \epsilon_{ij}^{\quad k} K_k - i \epsilon_{ji}^{\quad k} K_k - \epsilon_{ij}^{\quad k} J_k) = 0, \end{aligned} \quad (\text{A.4.13})$$

as well as

$$[A_i, A_j] = \epsilon_{ij}^{\quad k} A_k \quad (\text{A.4.14})$$

$$[B_i, B_j] = \epsilon_{ij}^{\quad k} B_k.$$

The two sets of generators therefore belong to two copies of $su(2)$ and we conclude that the Lie algebra of the Lorentz group $so(3, 1) \sim su(2) \times su(2)$.

A.4.1 Irreducible representations: bosons and fermions

The representations of $SO(3, 1)$ can thus be obtained by composing the fundamental representations of $SU(2)$. Let us denote with $(1/2, 0)$ the fundamental representation of the $SU(2)$ generated by the A_i 's and with $(0, 1/2)$ the fundamental representation of the $SU(2)$ generated by the B_i 's. The reason for this notation is that the dimension of both the representations $(s, 0)$ and $(0, s)$ is of course $d = 2s + 1$. For a generic representation (s, s') the dimension is given by the product

$$d = (2s + 1)(2s' + 1). \quad (\text{A.4.15})$$

One then finds from the composition rule (A.3.50) that, for example,

$$(1/2, 0) \times (0, 1/2) = (1/2, 1/2) , \quad (\text{A.4.16})$$

is a 4-vector ($d = 4$). Further

$$(1/2, 0) \times (1/2, 0) = (0, 0) + (1, 0) , \quad (\text{A.4.17})$$

where $(0, 0)$ is a scalar ($d = 1$) and $(1, 0)$ a 3-vector ($d = 3$), that is the skew-symmetric part of a $(2, 0)$ tensor.

We shall not go into further details, however a contact with physics is in order. From the physical point of view, one wants *fundamental particles* to appear of the same species for all inertial observers. This can be accomplished if such objects mathematically correspond to irreducible representations of $SO(3, 1)$. This far we skipped a detail, which is worth brining up in light of this observation. Irreducible representations of $SO(3, 1)$ [or rather $SU(2) \times SU(2)$] are uniquely identified by *two* parameters:

- 1) the mass $m \geq 0$ ³, and
- 2) the spin $s = 0, 1/2, 1, 3/2, \dots$

We already saw s , so the question is where does m come from. One can justify this second parameter formally by introducing the notion of Casimir operators for the algebra $su(2) \times su(2)$. However, we can just give a simpler physical answer to this question: consider for simplicity the vector representation of $SO(3, 1)$. We already know that the Minkowski modulus of a 4-vector is a scalar, so that, for the 4-momentum we have

$$P^\mu P_\mu = -m^2 , \quad (\text{A.4.18})$$

where m here denotes the proper mass. Since 4-momenta with different m are not transformed into each other by Lorentz transformations, it appears natural to consider that m contributes to distinguish different particles as well as s does. And that the corresponding vector spaces $\mathcal{V}(m, s)$ and $\mathcal{V}(m', s')$ are physically distinct (at least before we allow for interactions).

If we trust the mathematical structure that arises from the principle of relativity we therefore expect that there *may* exist two kinds of particles:

- 1) the *bosons* with integer spin and
- 2) the *fermions* with half-integer spin.

As a matter of fact, both such kinds do exist. The historical and physical reasons for their names however go beyond the scope of this course.

A.4.2 Poincaré group: $SO(4, 1)$

The group of Lorentz transformations and space-time translations can be represented by 5×5 matrices of the form

$$P = \begin{bmatrix} \Lambda & a \\ 0 & 1 \end{bmatrix} , \quad (\text{A.4.19})$$

³We have not introduced this parameter before, but one can argue about its existence since both rotations and unitary transformations do not mix vectors of different norm (as we noted in footnote 1 of this chapter).

where $\Lambda \in SO(3,1)$, $a \in \mathbb{R}^4$ and $1 \in \mathbb{R}$, and the fundamental representation is given by 5-vectors of the form

$$V^i = (V^\mu, 1) . \quad (\text{A.4.20})$$

In fact,

$$P^i_j V^j = \begin{bmatrix} \Lambda^\mu_\nu & a^\mu \\ 0 & 1 \end{bmatrix} \begin{pmatrix} V^\mu \\ 1 \end{pmatrix} = \begin{pmatrix} \Lambda^\mu_\nu V^\nu + a^\mu \\ 1 \end{pmatrix} , \quad (\text{A.4.21})$$

One can check that the matrix multiplication actually reproduces the expected action on 4-vectors.

Bibliography

- [1] R. Resnick, *Introduction to special relativity*, J. Wiley and Sons (1968).
- [2] R. Resnick, *Introduction to special relativity*, J. Wiley and Sons (1968), Section 2.2.
- [3] R. Resnick, *Introduction to special relativity*, J. Wiley and Sons (1968), Appendix A.
- [4] R. Resnick, *Introduction to special relativity*, J. Wiley and Sons (1968), Appendix B.
- [5] R. Resnick, *Introduction to special relativity*, J. Wiley and Sons (1968), Section 3.2.
- [6] M. Kaku, *Quantum field theory: a modern introduction*, Oxford Univ. Press (1993), Chapter 2.
- [7] B. Schutz, *Geometrical methods of mathematical physics*, Cambridge Univ. Press (1980).
- [8] B. Schutz, *Geometrical methods of mathematical physics*, Cambridge Univ. Press (1980), Chapters: 1.1-1.6, 2.1-2.9, 2.12-2.17, 2.19-2.30.
- [9] B. Schutz, *Geometrical methods of mathematical physics*, Cambridge Univ. Press (1980), Chapters: 3.1-3.7, 3.10-3.11.
- [10] B. Schutz, *Geometrical methods of mathematical physics*, Cambridge Univ. Press (1980), Chapters: 4.1-4.3.
- [11] B. Schutz, *Geometrical methods of mathematical physics*, Cambridge Univ. Press (1980), Chapters: 6.1-6.12.
- [12] B. Schutz, *A first course in general relativity*, Cambridge Univ. Press (2009).
- [13] S. Carroll, *Spacetime and geometry*, Addison-Wesley (2004). See also arXiv:gr-qc/9712019.
- [14] L. Landau e E. Lifshits, *Teoria dei campi*, Editori Riuniti (1976).