

# Filling the gaps in geometry and relativity

Thomas F. C. Bastos

May 24, 2024

## Abstract

The interplay between geometry and physics has been known for many decades, but it's only in fundamental physics where we can fully appreciate this relationship. In these notes, we will uncover one of these theories, General Relativity, and see how differential geometry plays an essential role. In the first part, we will strive to understand what curvature is, using ideas from Gauss, Cartan and Riemann. The second part is concerned with gravity: how it is an effect of the curvature of spacetime, and how such ideas came into being. It is meant to be an introduction to the introduction of general relativity, and for those who are already acquainted with the theory, it should fill some gaps that may have been left when you studied GR for the first time.

## 0 The physicist approach

If I asked people what they think a physicist would do to study geometry, they would probably say: "By throwing particles at it". That's precisely what we are going to do! No, seriously. Consider a free particle in space  $\mathbb{R}^3$ . It's lagrangian is just kinetic energy  $L = m \frac{v^2}{2}$ . Using cartesian coordinates the Euler-Lagrange equation gives us

$$\frac{d}{dt} \left( \frac{dx^i}{dt} \right) = 0$$

which is, unsurprisingly, a straight line in euclidean space. Consider now the case of a free particle constrained to the surface of a sphere. Using spherical coordinates we have:

$$v^2 = R^2 \dot{\theta}^2 + R^2 \dot{\phi}^2 \sin^2 \theta$$

The Euler-Lagrange equations for  $\theta$  and  $\phi$  are:

$$\begin{aligned} \ddot{\theta} - \sin \theta \cos \theta \dot{\phi}^2 &= 0 \\ \ddot{\phi} + 2 \cot \theta \dot{\theta} \dot{\phi} &= 0 \end{aligned}$$

One solution to this problem is  $\theta(t) = \frac{\pi}{2}$  and  $\phi(t) = \omega t$ , which is just a particle in the equator with constant velocity. In fact, from the symmetry we can rotate this solution in such a way that for every two points, the path the particle will take lies

in a great circle<sup>1</sup> that passes through both points. Observe that in every case the free particle takes the path of minimal length, i.e takes the path of geodesics, and it's easy to see that this is a general property of particles constrained to surfaces since the action is just  $S = \frac{m}{2} \int v^2 dt \sim \int ds$ .

A particle on a general surface may be parametrized by two degrees of freedom  $u^1, u^2$

$$(x(u^1, u^2), y(u^1, u^2), z(u^1, u^2))$$

To measure the infinitesimal distance between neighboring points on surfaces, the distance must be something of the form:

$$ds^2 = dx^2 + dy^2 + dz^2 = \sum_{ij} g_{ij} du^i du^j$$

where we used the chain rule  $dx^j = \sum_{i=1}^2 \frac{\partial x^j}{\partial u^i} du^i$  and defined the metric coefficients  $g_{ij}(x)$  which are symmetric and depends on the coordinates we are using to describe the surface. For example, on 'flat' space the line element is

$$ds^2 = dx^2 + dy^2 + dz^2$$

and on the sphere

$$ds^2 = R^2 d\theta^2 + R^2 \sin^2 \theta d\phi^2$$

It's convenient to write  $ds^2 = g_{ij} dx^i dx^j$  suppressing the sum while keeping in mind that every pair of repeated indices must be summed. Note that  $ds^2$  is a distance so it does not depend on the choice of coordinates. In some sense, the information of the curvature on the surface must be encoded in the metric coefficients: the way we measure distances must give information about the geometry.

Now we seek for the equation of a geodesic in any surface, so we need to find the curve that connect two points<sup>2</sup> which makes the distance minimal. This is equivalent to finding the path a free particle will take on the surface, so the lagrangian of this problem is:

$$L = \frac{1}{2} g_{ij} \dot{x}^i \dot{x}^j = \frac{1}{2} v^2$$

where we consider a particle with  $m = 1$  without loss of generality. Therefore

$$\frac{\partial L}{\partial x^i} = \frac{1}{2} \frac{\partial g_{ij}}{\partial x^i} \dot{x}^i \dot{x}^j = \frac{1}{2} \frac{d}{dt} \left( g_{ij} \frac{\partial \dot{x}^i}{\partial \dot{x}^l} \dot{x}^j + g_{ij} \frac{\partial \dot{x}^j}{\partial \dot{x}^l} \dot{x}^i \right) = \frac{1}{2} \frac{d}{dt} (g_{ij} \delta_{il} \dot{x}^j + g_{ij} \delta_{jl} \dot{x}^i) = \frac{d}{dt} (g_{il} \dot{x}^i)$$

where we use the fact that  $g_{ij}$  is symmetric and does not depend on the derivative of the coordinates. Hence

$$\frac{d}{dt} (g_{il} \dot{x}^i) = \frac{\partial g_{il}}{\partial x^j} \dot{x}^j \dot{x}^i + g_{il} \ddot{x}^i = \frac{1}{2} \left( \frac{\partial g_{il}}{\partial x^j} \dot{x}^j \dot{x}^i + \frac{\partial g_{jl}}{\partial x^i} \dot{x}^i \dot{x}^j \right) + g_{il} \ddot{x}^i$$

---

<sup>1</sup>A great circle is the circle that lies in the intersection of a sphere with a plane that passes through the center of the sphere

<sup>2</sup>Actually there are some issues in considering *any* two points because the extremal of the functional may not be a minimum, so in general we take two points in a neighborhood close enough to give a unique curve with minimal length.

Provided that the metric is non-degenerate and its inverse is given by  $g^{ij}$ , we can set

$$\Gamma_{jl}^i = \frac{1}{2}g^{il} \left( \frac{\partial g_{il}}{\partial x^j} + \frac{\partial g_{jl}}{\partial x^i} - \frac{\partial g_{ij}}{\partial x^l} \right)$$

the geodesic equations are given by

$$\ddot{x}^i + \Gamma_{jl}^i \dot{x}^j \dot{x}^l = 0 \quad (1)$$

Such quantities  $\Gamma_{jl}^i$  are called the Levi Civita connection. For a physicist, equation 1 is just the acceleration of a particle in funny coordinates, or the acceleration of a particle on a surface.<sup>3</sup> We can make a slight change in equation 1

$$\dot{x}^l \left( \frac{\partial \dot{x}^i}{\partial x^l} + \Gamma_{jl}^i \dot{x}^j \right) = \dot{x}^l \nabla_l \dot{x}^i = 0$$

where  $\nabla_l = \frac{\partial}{\partial x^l} + \Gamma_{jl}^i$  is called the covariant derivative. Now this looks like the equation of a 'straight line', the only difference being that we changed our notion of derivative  $\partial_k \rightarrow \nabla_k$  when geometry comes at play. We still can't tell if this occurs because of the curvature of the surface or just because we choose a funny coordinate system. Partial derivatives always commute  $\partial_i \partial_j = \partial_j \partial_i$  while covariant derivatives may not. With a bit more work and some extra definitions you can show that

$$(\nabla_i \nabla_j - \nabla_j \nabla_i) x^k = R_{ijl}^k x^l$$

where the Riemann tensor is

$$R_{ijl}^k = \partial_j \Gamma_{il}^k + \Gamma_{jr}^k \Gamma_{il}^r - (j \leftrightarrow l)$$

The failure of covariant derivatives to commute is something you can only find in curved geometries: parallel transport a vector in a loop and see what happens both on a piece of paper and a sphere. It turns out that the Riemann tensor contains all the information about the curvature.

In section 2 we will see that gravity can be completely incorporated by the change  $\partial_\mu \rightarrow \nabla_\mu$  once we let spacetime be an arbitrary four dimensional manifold equipped with a lorentzian metric that satisfy Einstein's field equations. If you ever tried to do the quantum mechanics of a particle in a magnetic field, you may remember that you had to change the derivative  $\partial_k \rightarrow \partial_k - iqA_k$  in Schrodinger's equations to accommodate the magnetic field. Under a gauge transformation  $A \rightarrow A + \partial\Omega$  the wave function transforms like  $\psi \rightarrow e^{i\Omega(x)}\psi$  to keep gauge invariance. This looks like the action of the circle group  $U(1) \simeq S^1$  on the set of wave functions at each spacetime point. Now imagine that we attached a copy of the circle  $S^1$  at each point in spacetime, so we can keep track of the position and phase factor. Our space would look like figure 1

It turns out that electromagnetism is the theory of a principal  $U(1)$ -bundle over Minkowski spacetime, where the connection is precisely the gauge potential  $A_\mu$ . In fact, every interaction in the Standard Model is a principal  $G$ -bundle over spacetime where the gauge group is  $G = SU(n)$ . We also require that the gauge fields obey the Yang-Mills equations, which in some sense is just a choice of a connection that minimizes the curvature

A physicist would be proud.

---

<sup>3</sup>For a mathematician  $\Gamma_{jl}^i$  is the pull back of a Yang-Mills connection on the frame bundle. I guess physicists are winning this one.

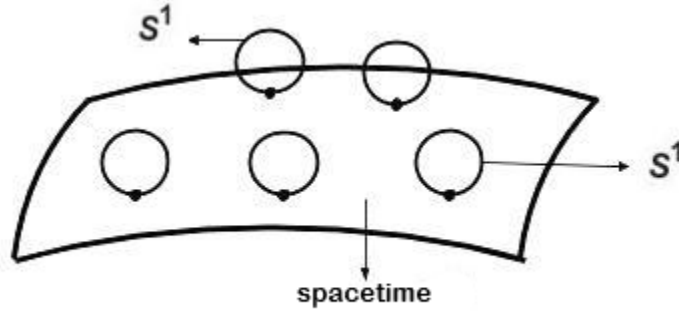


Figure 1:  $U(1)$ -bundle

## 1 What is curvature?

In Book XI of the Confessions (397) Saint Augustine was trying to understand time. There he said something that struck me for a while:

“What then is time? If no one asks me, I know; if I want to explain it to a questioner, I do not know.” (Augustine, p. 242)

I believe the same is true for many concepts, but it is specially true for curvature: you know exactly what it is, except when you have to compute it. Instead of looking at intuition as being a bug in our brains, I like to think it’s actually a feature. It’s time to roll up our sleeves and try to build curvature and geometry from our intuition. We will thank Augustine for this insight latter.

We start with one dimensional geometries, i.e. curves, and quickly exhaust everything there is to know about them using frame fields. It turns out that this method is extremely powerful and, provided some suitable modifications, it will carry over to any number of dimensions. In our path we will naturally encounter concepts such as manifolds, tangent spaces, metrics, connections and covariant derivatives. Hopefully this will make you see the whole picture and don’t be bothered asking yourself questions like: why in the world is curvature a  $(1,3)$  tensor?

### 1.1 Curves and the power of Frenet Frames

We start with the simplest possible geometry: a parametrized curve  $\gamma(t) : I \rightarrow \mathbb{R}^3$  in three dimensions. In figures 1 and 2, we can clearly see that 1 has no curvature while 2 has some curved sections. Curvature is a local property. Another basic fact is that a circle with a small radius is more curved than a circle with a large radius. This is easy to see: if you zoom in on any circle, it starts to look like a straight line. Therefore the curvature  $k$  of a circle should be inversely proportional to its radius  $R$ :

$$k \sim \frac{1}{R} \quad (2)$$

Now consider the helix in figure. It is intuitive that the helix not only have curvature but also twist as it moves, a feature that is not shared with the circle. Notice that all

of these features are defined by how the curve changes its directions through space. Let's formalize all of these statements.

Since we are only interested in the directions, we may only consider paths with unit speed  $\|\gamma'\| = 1$  without loss of generality.

**Definition 1.** A regular curve is a differentiable map  $\gamma : I \rightarrow \mathbb{R}^3$  such that  $\|\gamma'(s)\| = 1$

There is a standard procedure to take any path  $\gamma(t)$  without cusps and make it a regular curve. First calculate the integral

$$s = \int_0^t \|\gamma'\| dt$$

and then take the inverse of the function  $s = s(t)$  so that  $t = t(s)$  and  $\gamma(s) = \gamma(t(s))$  where  $\|\gamma'(s)\| = 1$ . This is sometimes called arc length parametrization. Making the rather confusing relabel  $T \doteq \gamma'$  we are ready to define curvature:

**Definition 2.** The curvature  $k(s)$  of a regular curve  $\gamma$  is given by:

$$T' = kN \tag{3}$$

where  $N(s)$  is a unit vector field  $\|N\| = 1$ .

At first this may look abstract but it's precisely what we wanted:  $k$  measures the amount of change of direction of the unit tangent  $T$ . We can use the inner product  $\langle T, T \rangle = 1$  to show that the vector fields are orthogonal:

$$\langle T, T \rangle' = 2k\langle T, N \rangle = 0$$

Now define another vector field  $B \doteq T \times N$  which is again unitary and orthogonal to both  $T$  and  $N$ . You should picture these vector fields  $T, N, B$  as moving frames attached to the curve.

**Definition 3.** Given a regular curve  $\gamma$  its Frenet frame is the set of vector fields  $\{T, N, B\}$ .

You can use the inner product to show (exercise) that  $B' = -\tau N$  where  $\tau(s) \in \mathbb{R}$  is called the torsion. The most important thing about the Frenet frames is that their derivatives  $T', N', B'$  are expressed in terms of themselves  $T, N, B$ . In this way we can keep track of all the changes in all directions of the curve.

**Theorem 1.1.** If  $\gamma$  is a regular curve and  $T, N, B$  are its Frenet frame fields, then

$$\begin{pmatrix} T' \\ N' \\ B' \end{pmatrix} = \begin{pmatrix} 0 & k & 0 \\ -k & 0 & \tau \\ 0 & -\tau & 0 \end{pmatrix} \begin{pmatrix} T \\ N \\ B \end{pmatrix} \tag{4}$$

*Proof.* The first and last equations of 4 are just definitions. Since the Frenet frame fields are orthonormal:

$$N' = \langle N', T \rangle T + \langle N', N \rangle N + \langle N', B \rangle B$$

From the identity  $\langle N, T \rangle = 0$  it follows that  $\langle N', T \rangle = -\langle N, T' \rangle = -k$  and by similar arguments  $\langle N', B \rangle = \tau$  □

At this point we should check if this definition agrees with our expectations. A straight line  $\gamma(s) = a + bs$  where  $\|b\| = 1$  has curvature  $k = \|b'\| = 0$ . A circle can be parametrized by  $\gamma(t) = (R \cos t, R \sin t, 0)$  but its arc length parametrization is different:

$$s = \int_0^t \|\gamma'(t)\| dt = Rt$$

so that  $\gamma(s) = (R \cos \frac{s}{R}, R \sin \frac{s}{R}, 0)$ . The curvature is then

$$k = \|T'\| = \|\gamma''(s)\| = \left\| -\frac{1}{R}(\cos \frac{s}{R}, \sin \frac{s}{R}, 0) \right\| = \frac{1}{R}$$

and the torsion is  $\tau = 0$  since  $B' = 0$ . This is exactly what we expected from equation 2 and also from the lack of twisting of a circle. An arc length parametrization of the helix is (exercise):

$$\gamma = (a \cos \frac{s}{c}, a \sin \frac{s}{c}, \frac{b}{c}s)$$

where  $c^2 = a^2 + b^2$ . Then  $k = \frac{a}{c^2}$  and  $\tau = \frac{b}{c^2}$ . Notice that when  $b = 0$  the torsion goes to zero and the helix collapse into a circle.

These are some nice results, but there is a strong statement about curvature and torsion when you integrate the Frenet formulas 4:

**Theorem 1.2.** *Given two scalar fields  $k, \tau : I \rightarrow \mathbb{R}$  there exists one regular curve (up to isometries)  $\gamma : I \rightarrow \mathbb{R}^3$  such that its curvature and torsion are given by  $k, \tau$ .*

In other words,  $k$  and  $\tau$  completely determines the geometry of a curve. Well, that was kinda easy. The take away should be that if we want to understand curvature, we better have a way to construct frame fields attached to every point of our geometric object and express their derivatives in terms of themselves like in 4.

The thing is that there is nothing special about curves. We could choose a frame field that is defined on every point of  $\mathbb{R}^3$ . We would expect them to describe the curvature and torsion of the whole space  $\mathbb{R}^3$ , so it's interesting to generalize the formalism. But before we jump into that, we need to define extra objects that will greatly help us. These are covariant derivatives and forms. Let's start with some terminology:

**Definition 4.** *For each point  $p \in \mathbb{R}^3$  the set of all vectors with base at  $p$  is denoted by  $T_p(\mathbb{R}^3)$  and is called the tangent space at  $p$ . The set of all tangent spaces is called the tangent bundle  $T\mathbb{R}^3$ .*

**Definition 5.** *A smooth vector field is a smooth map  $V : \mathbb{R}^3 \ni p \mapsto V(p) \in T_p(\mathbb{R}^3)$ .*

Since we want to see how vector fields change in any direction, it's useful to define derivatives of these objects.

**Definition 6.** *Let  $W$  and  $V$  be smooth vector fields. The covariant derivative of  $W$  in the direction of  $V$  is another vector field  $\nabla_V W$  such that:*

$$\nabla_V W(p) = W(p + V(p)t)'|_{t=0}$$

The definition of a covariant derivative looks like something you have never seen before, but you actually did. Take a point  $p$ , a vector  $V(p)$  pointing somewhere with base on  $p$  and compute the directional derivative  $W(p+V(p)t)'|_{t=0}$  of  $W$  in the direction  $p+V(p)t$ . For example, let  $V = -y\hat{x} + x\hat{y}$  and  $W = \cos(x)\hat{x} + \sin(y)\hat{y}$  be vector fields and  $p = x\hat{x} + y\hat{y} + z\hat{z}$  be any point in  $\mathbb{R}^3$ . Then

$$p + V(p)t = (x - yt)\hat{x} + (y + xt)\hat{y} + z\hat{z}$$

$$\nabla_V W(p) = W(p + V(p)t)'|_{t=0} = y \sin(x)\hat{x} + x \cos(y)\hat{y}$$

Remember when your calculus teacher said that the differential

$$df = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy + \frac{\partial f}{\partial z} dz$$

is not well defined? They lied to you (for a good reason). We can define them using forms, which turns out to be *extremely* useful in physics and mathematics.<sup>4</sup>

**Definition 7.** A one-form is a linear map  $\omega_p : T_p(\mathbb{R}^3) \rightarrow \mathbb{R}$ . A one-form field is a linear map  $\omega : T\mathbb{R}^3 \rightarrow \mathbb{R}$ .

**Definition 8.** Let  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  be a differentiable function. The differential  $df$  is a one-form such that:

$$df(v) = v(f) = f(p + vt)'|_{t=0}$$

where  $v \in T_p(\mathbb{R}^3)$ .

**Lemma 1.3.** The differentials  $dx, dy, dz$  form a basis<sup>5</sup> of the vector space  $T_p^*\mathbb{R}^3$  of all one forms at a point  $p \in \mathbb{R}^3$

*Proof.* Notice that  $dx(v) = v(x) = v_x$  and the same for the other differentials. For any one form  $\omega_p \in T_p^*\mathbb{R}^3$  and vector  $v = v_x\hat{x} + v_y\hat{y} + v_z\hat{z} \in T_p(\mathbb{R}^3)$  we have that:

$$\omega_p(v) = \omega_p\left(\sum_i v_i \hat{x}_i\right) = \sum_i v_i \omega_p(\hat{x}_i) = \left(\sum_i \omega_p(\hat{x}_i) dx_i\right)(v)$$

□

**Corollary 1.3.1.** The differential can be written as  $df = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy + \frac{\partial f}{\partial z} dz$

*Proof.* Trivial. □

Given two one-forms  $\omega, \phi$  we can make their product  $\omega \wedge \phi$  which is a two-form. The simplest way to think about them is to consider their action on the differentials and the distributive property

**Definition 9.** The wedge product  $\wedge$  is a bilinear map that takes two one-forms and gives a two-form such that

$$dx_i \wedge dx_j = -dx_j \wedge dx_i$$

---

<sup>4</sup>That is an understatement.

<sup>5</sup>Recall that the dual  $V^*$  of a finite dimensional vector space  $V$  is again another vector space of the same dimension.

Given any two one-forms  $\omega = \sum_i \omega_i dx_i$  and  $\phi = \sum_j \phi_j dx_j$  their product is

$$\omega \wedge \phi =$$

Needs a lot of attention. Mention the electromagnetic field as a one-form, and faraday tensor as a two-form  $F = dA$ .

**Definition 10.** A frame field is a set of vector fields  $E_1, E_2, E_3$  in  $\mathbb{R}^3$  such that:

$$\langle E_i, E_j \rangle = \delta_{ij}$$

where  $\delta_{ij}$  is the Kronecker delta.

We know that something funny happens when we express derivatives of these frames in terms of themselves. Let  $V$  be any vector field and

$$\begin{aligned}\nabla_V E_1(p) &= \omega_{11} E_1(p) + \omega_{12} E_2(p) + \omega_{13} E_3(p) \\ \nabla_V E_2(p) &= \omega_{21} E_1(p) + \omega_{22} E_2(p) + \omega_{23} E_3(p) \\ \nabla_V E_3(p) &= \omega_{31} E_1(p) + \omega_{32} E_2(p) + \omega_{33} E_3(p)\end{aligned}$$

In a shorter notation using Einstein's convention of repeated indices  $\nabla_V E_i(p) = \omega_{ij} E_j(p)$ , the coefficients  $\omega_{ij}$  clearly depends on the vector  $V$  and looks like a one-form field:

**Lemma 1.4.** Let  $E_i$  be frames fields in  $\mathbb{R}^3$ ,  $v \in T_p(\mathbb{R}^3)$  any vector and

$$\omega_{ij}(v) = \langle \nabla_v E_i(p), E_j(p) \rangle$$

Then  $\omega_{ij}$  is a matrix valued one-form and  $\omega_{ij} = -\omega_{ji}$ . They are called connection coefficients.

*Proof.* You can use definition 6 and the chain rule to prove that  $\nabla_{av+bu} E = a \nabla_v E + b \nabla_u E$  for any vectors  $u, v \in T_p(\mathbb{R}^3)$  and numbers  $a, b \in \mathbb{R}$ . Therefore

$$\omega_{ij}(av + bu) = \langle \nabla_{av+bu} E_i(p), E_j \rangle = a \omega_{ij}(v) + b \omega_{ij}(u)$$

To prove the antisymmetry of  $\omega$  notice that

$$\nabla_v \langle E_i, E_j \rangle = \langle \nabla_v E_i, E_j \rangle + \langle E_i, \nabla_v E_j \rangle = 0$$

□

Now we can use the dual 1-forms  $\theta_i$  of the frame fields defined as:

$$\theta_i(v) = \langle v, E_i \rangle$$

to show a very interesting result that will clarify the geometry of  $\mathbb{R}^3$ . Notice that  $\theta_i$  forms a basis of  $T_p^*(\mathbb{R}^3)$ .

**Theorem 1.5** (Cartan Structural Equations). Let  $E_i$  be frame fields,  $\theta_i$  its duals and  $\omega$  the connection coefficients. Then

$$d\theta_i = \omega_{ij} \wedge \theta_j \tag{5}$$

$$d\omega_{ij} = \omega_{ik} \wedge \omega_{kj} \tag{6}$$



*Proof.* Note that any frame can be written in terms of the canonical vector field like  $E_i = a_{ij}\hat{x}_j$ , and so too the dual forms  $\theta_i = a_{ij}dx_j$ . On a more compact notation  $\theta = Ad\xi$  where:

$$\begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \begin{pmatrix} dx_1 \\ dx_2 \\ dx_3 \end{pmatrix} \quad (7)$$

and  $A \in O(3)$  is an orthogonal matrix. We have that

$$\omega_{ij}(v) = \langle \nabla_v E_i(p), E_j(p) \rangle = da_{ik}a_{jk}$$

or in matrix notation  $\omega = dAA^T$ . The Cartan equations then follows

$$d\theta = dAd\xi + Ad^2\xi = dAd\xi = dAA^T Ad\xi = \omega\theta$$

$$d\omega = d(dAA^T) = -dAdA^T = -dAA^T AdA^T = -\omega\omega^T = \omega\omega$$

□

We will see later that the second Cartan equation 6 means that  $\mathbb{R}^3$  is flat. All the geometric properties are found in the connection coefficients! The next logical step in our investigation is to go up a dimension.

## 1.2 Surfaces, isometries and why the Earth isn't flat

What exactly is a 2 dimensional surface? This is an important question to will be addressed very soon. We will naturally run into the notion of a manifold. For now we just need a simplified version of the answer.

In figure we have an example of what we generally think is a surface  $M$ . A surface is something that we can describe (at least locally) by two coordinates. At a point  $p \in M$  in the surface, there exist a tangent plane which is the best first approximation to  $M$ .

**Definition 11.** *The set of all tangent vectors to the surface based at point  $p \in M$  is called the tangent space  $T_p(M)$*

It's clear that  $T_p(M)$  is really a 2 dimensional vector space. This gives us a very natural way to attach a frame field to the surface: choose  $E_1, E_2$  to lay on the tangent space and  $E_3$  to be the unit normal vector field at each point of  $M$ . Notice that we can only hope for this construction to be true locally <sup>6</sup>. To write the structural equations for the surface we only apply derivatives wrt. vectors in the tangent space  $T_p(M)$ . An immediate consequence is that  $\theta_3(v) = \langle v, E_3 \rangle = 0$ . In total, there are just two dual-forms  $\theta_1, \theta_2$  and applying equations 5, 6 we have all the geometrical information of the surface.

---

<sup>6</sup>Look up the hairy ball theorem

**Lemma 1.6.** *Let  $M$  be a surface,  $E_i \in TM$  the attached frame fields and  $\theta_i \in T^*M$  its dual-forms. The structural equations are*

$$d\theta_1 = \omega_{12} \wedge \theta_2 \quad (8)$$

$$d\theta_2 = \omega_{21} \wedge \theta_1 \quad (9)$$

$$\omega_{31} \wedge \theta_1 + \omega_{32} \wedge \theta_2 = 0 \quad (10)$$

$$d\omega_{12} = \omega_{13} \wedge \omega_{32} \quad (11)$$

$$d\omega_{13} = \omega_{12} \wedge \omega_{23} \quad (12)$$

$$d\omega_{23} = \omega_{21} \wedge \omega_{13} \quad (13)$$

We are done! However, it's not immediate how to get the curvature scalar out of those equations. The following lemma prescribes an algorithm for this

**Lemma 1.7.** *Let  $K \doteq \omega_{13}(E_1)\omega_{23}(E_2) - \omega_{13}(E_2)\omega_{23}(E_1)$  be the Gaussian curvature. Then*

$$d\omega_{12} = -K\theta_1 \wedge \theta_2 \quad (14)$$

*Proof.* Note that any two-form defined on the surface can be written as  $\phi = \phi(E_1, E_2)\theta_1 \wedge \theta_2$  since  $\theta_3 = 0$ . We have that

$$\omega_{13} \wedge \omega_{23}(E_1, E_2) = \omega_{13}(E_1)\omega_{23}(E_2) - \omega_{13}(E_2)\omega_{23}(E_1) = K$$

Using the second structural equation

$$d\omega_{12} = \omega_{13} \wedge \omega_{32} = -K\theta_1 \wedge \theta_2 \quad (15)$$

□

*Exercise:* Choose some frame field for a sphere of radius  $R$ , find the dual-forms and use the structural equations to show that  $K = \frac{1}{R^2}$ . Do the same for a plane and a cylinder and show that both have  $K = 0$ .

At this point it is worthwhile to stop and reflect about what we just did. There are several problems that we didn't consider. Do we get different curvatures by choosing different frames? We only wanted one curvature but it looks like we got three of them, one for each structural equation  $d\omega_{ij} = \omega_{ik} \wedge \omega_{kj}$ . Why did we get so much more than we put in? It is not clear how to generalize what we did for higher dimensions because we don't know what a surface is. In all cases we considered geometries embedded in a higher dimensional space  $\mathbb{R}^3$ . How to make a geometry that is *intrinsic* to the surface?

We will solve these problems all at once by adjusting the formalism to be intrinsic in nature. You can show that the gaussian curvature  $K$  is the only invariant under local isometries transformations. These are maps between surfaces  $F : M \rightarrow M'$  that preserves the inner product of tangent vectors. Think of them as bending without stretching. It's clear that a cylinder is isometric to a plane, so it shouldn't be a surprise that both have zero gaussian curvature. It also means that we cannot possibly obtain an isometry between a sphere a plane: they have different gaussian curvatures. That's one reason why we can't have a faithful map of the Earth. What we call intrinsic

geometry is really the set of objects that are invariant or keep it's structure under local isometries. Simply put, isometric geometries are the same. A two-dimensional inhabitant of a surface would only be able to measure the intrinsic distance, hence detecting only a Gaussian curvature and all his vectors and forms would transform under isometries. Just like they, we can't tell if we live embedded in a 11 dimensional geometry<sup>7</sup> and it should not matter for the questions we are trying to answer. In the next chapter we will construct intrinsic geometry for any dimensions, and finally get the general notion of curvature via the Riemann tensor.

### 1.3 Higher dimensions, Riemann curvature and the modern stuff

Before jumping into the hard math, let's make a wish list. First we need to define what truly is an abstract intrinsic surface. Our definition has to be sensible enough to let us have a tangent space at each point, since this enable us to do intrinsic calculus and build frames. Our next step is to smoothly assign an inner product  $\langle, \rangle$  on each tangent space, since we need them to define the connection coefficients. These ingredients are enough to define curvature and torsion, choose a sensible covariant derivative, define geodesics and all the good things. Let's go.

A 2-d surface is just a set  $M$  where you can ascribe 2 coordinates. In mathematical lingo we say that for each point  $p \in M$  there exists an open set  $V \subset M$  that contains  $p$  and a homeomorphism<sup>8</sup>  $x : V \rightarrow \mathbb{R}^2$ . We call  $x$  a coordinate system, or chart, of the patch  $V$ . Notice that  $\mathbb{R}^3$  didn't appear in the definition of a surface! We can choose  $M$  to be any set, not only a subset  $M \subset \mathbb{R}^3$ .

It's clear how we generalize the notion of a surface for any number of dimensions:

**Definition 12.** *A smooth  $n$ -dimensional manifold is a nice topological space  $M$  such that:*

(i) *For every point  $p \in M$  there exists an open set  $V \subset M$  and a homeomorphism*

$$x : V \rightarrow U \subset \mathbb{R}^n$$

(ii) *Any overlap  $U \cap V$  of different charts  $x : U \rightarrow \mathbb{R}^n$  and  $y : V \rightarrow \mathbb{R}^n$  are smoothly joined together. This means that the transition functions:*

$$x \circ y^{-1} : y(U \cap V) \subset \mathbb{R}^n \rightarrow x(U \cap V) \subset \mathbb{R}^n$$

*are  $C^\infty$  functions in the usual sense of  $\mathbb{R}^n$  analysis.*

The first condition is obvious. The second condition (ii) is there to guarantee we are going to have a sensible definition of calculus on the manifold. For instance, we define a function  $f : M \rightarrow \mathbb{R}$  to be smooth if and only the map

$$f \circ x^{-1} : \mathbb{R}^n \rightarrow \mathbb{R}$$

is smooth in the usual  $\mathbb{R}^n$  sense. This is well defined by property: any overlap of charts  $x, y$  in  $U \cap V$  implies

---

<sup>7</sup>Shout out to all of the string people

<sup>8</sup>See the appendix for topology.

$$f \circ x^{-1} = (f \circ y^{-1}) \circ (y \circ x^{-1})$$

which is still a smooth function since  $y \circ x^{-1}$  is smooth by property (ii). We are basically transferring notions of differentiability in conventional calculus to define it on abstract manifolds. The rest is procedural.

**Definition 13.** A curve  $\lambda : \mathbb{R} \rightarrow M$  is of class  $C^\infty$  if for every chart  $(U, x)$  the map  $x \circ \lambda : \mathbb{R} \rightarrow x(U) \subset \mathbb{R}^n$  is of class  $C^\infty$ .

**Definition 14.** Given a  $C^\infty$  curve  $\lambda : \mathbb{R} \rightarrow M$ , the tangent vector  $X_p$  to a point  $p = \lambda(t_0)$  is the operator which maps functions  $f : M \rightarrow \mathbb{R}$ <sup>9</sup> to the directional derivative of  $f$  by the curve  $\lambda$  calculated at the point  $p$ :

$$X_p : f \mapsto \frac{d(f \circ \lambda)}{dt}(t_0)$$

we write symbolically  $X_p = (\frac{d}{dt})_p$

It may seem strange to define tangent vectors as operators acting on smooth functions, but that's something you have to complain to your mathematician friend. To be honest this is not such a weird thing considering that linear operators live in a vector space, just like honest looking vectors.

**Proposition 1.8.** Let  $M$  be a smooth manifold and  $p \in M$ . The set of all tangent vectors to  $p$  forms a vector space  $T_p(M)$  called the tangent space, where the sum of operators is defined pointwise  $(X_p + Y_p)(f) = X_p(f) + Y_p(f)$ .

*Proof.* Let  $\mu, \nu$  be curves in  $M$  such that  $\mu(t_0) = \nu(t_0) = p$ , and  $X_p, Y_p$  be the induced tangent vectors. We have to show that there exists a curve  $\lambda$  passing through  $p = \lambda(t_0)$  such that the tangent vector induced by  $\lambda$  is  $Z_p = X_p + Y_p$ .

Let  $\bar{\lambda} : t \mapsto x \circ \mu(t) + x \circ \nu(t) - x(p)$ , then we see that  $\lambda = x^{-1} \circ \bar{\lambda}$  satisfies the properties that we want:  $\lambda(t_0) = p$  and for every function  $f$

$$\frac{d(f \circ \lambda)}{dt} = \frac{d(f \circ x^{-1} \circ \bar{\lambda})}{dt} = \frac{\partial f \circ x^{-1}}{\partial x_i} \frac{d(x \circ \mu_i)}{dt} + \frac{\partial f \circ x^{-1}}{\partial x_i} \frac{d(x \circ \nu_i)}{dt}$$

using the chain rule in the opposite direction we have

$$Z_p(f) = X_p(f) + Y_p(f)$$

The same argument is used for the closure by multiplication of scalars, where we use the curve  $\mu : t' \rightarrow \lambda(\alpha.t' + t_0)$  where  $\alpha$  is a constant.  $\square$

We can talk about a basis for  $T_p(M)$  and one natural choice is the chart-induced basis: for every vector  $X_p \in T_p(M)$

$$X_p(f) = \frac{df \circ \lambda}{dt} = \frac{\partial f \circ x^{-1}}{\partial x^k} \frac{d(x \circ \lambda)^k}{dt} = \frac{\partial f \circ x^{-1}}{\partial x^k} \frac{d(x^k \circ \lambda)}{dt}$$

Set a notation for these derivatives is going to save us some sanity:

---

<sup>9</sup>Every time we talk about functions it is implied that its domain is the manifold.

$$\frac{\partial f \circ x^{-1}}{\partial x^k} \doteq \frac{\partial f}{\partial x^k} = \frac{\partial}{\partial x^k} f$$

we can rewrite equations of this type to a simpler form:

$$X_p(f) = \frac{\partial x^k \circ \lambda}{\partial t} \frac{\partial}{\partial x^k} f$$

$$X_p = X_p(x^k) \left( \frac{\partial}{\partial x^k} \right) |_p = X^k \frac{\partial}{\partial x^k}$$

In fact, the operator  $\frac{\partial}{\partial x^k}$  is a tangent vector once we consider curves  $\lambda_k : M \rightarrow \mathbb{R}$  such that  $x \circ \lambda_k = (x^1(p), \dots, x^k(p) + t, \dots, x^n(p))$ . That they are linearly independent is trivial. This shows that the dimension of  $T_p(M)$  is  $n$ , so the dimension of the manifold is precisely the dimension of its tangent space. Now we talk about dual-vectors:

**Definition 15.** Let  $T_p(M)$  be the tangent space in  $p$ . The dual space  $T_p^*(M)$  is the space of all linear functionals in  $T_p(M)$ , i.e., the space of all linear map  $\omega : T_p(M) \rightarrow \mathbb{R}$

With the point sum defined as  $(\omega^1 + \omega^2)(X) = \omega^1(X) + \omega^2(X)$  where  $X \in T_p(M)$ , it's clear that  $T_p^*(M)$  is a vector space and from basic linear algebra  $\dim(T_p^*(M)) = n$ . The elements of  $T_p^*(M)$  are the beloved one-forms. The differential of  $f$  is again defined as a map  $df : T_p(M) \rightarrow \mathbb{R}$  such that:

$$df(X) = X(f)$$

**Definition 16.** A dual basis to  $\{\frac{\partial}{\partial x^k}\}$  are those forms  $\{dx^k\}$  which satisfy:

$$dx^k \left( \frac{\partial}{\partial x^i} \right) = \left( \frac{\partial}{\partial x^i} \right) x^k = \delta_i^k$$

They form a basis of  $T_p^*(M)$ .

Now something of utmost importance for our discussion: tensors. Take a point  $p$ , and define

$$\Pi_r^s \doteq T_p^* \times \dots \times T_p^* \times T_p \times \dots \times T_p = (T_p^*)^r \times (T_p)^s$$

**Definition 17.** A tensor of type  $(r, s)$  at  $p \in M$  is a map  $T : \Pi_r^s \rightarrow \mathbb{R}$  that is linear in each argument.

The space of all tensors  $T_s^r(p)$  defined in  $\Pi_r^s$  is a vector space once we define the sum of tensors and product by a scalar as:

$$(T + T')(\omega^1, \dots, \omega^r, X_1, \dots, X_s) = T(\omega^1, \dots, \omega^r, X_1, \dots, X_s) + T'(\omega^1, \dots, \omega^r, X_1, \dots, X_s)$$

$$(a.T)(\omega^1, \dots, \omega^r, X_1, \dots, X_s) = a.T(\omega^1, \dots, \omega^r, X_1, \dots, X_s)$$

We see that it has dimension  $n^{s+r}$ .

**Definition 18.** Let  $R \in T_s^r$  and  $S \in T_l^k$  be tensors, the tensor product  $R \otimes S \in T_{s+l}^{r+k}$  is again another tensor defined as:

$$(R \otimes S)(\omega^1, \dots, \omega^{r+k}, X_1, \dots, X_{s+l}) = R(\omega^1, \dots, \omega^r, X_1, \dots, X_s) \cdot S(\omega^{r+1}, \dots, \omega^{r+k}, X_{s+1}, \dots, X_{s+l})$$

As an example, take  $R = Y \in T_0^1 = T_p(M)$  and  $S = \eta \in T_1^0 = T_p^*(M)$  then  $(R \otimes S)(\omega, X) = R(\omega) \cdot S(X) = \omega(Y) \cdot \eta(X)$ . These properties are more easily seen in terms of it's components when we have a basis for  $T_s^r$ :

**Proposition 1.9.** If  $\{\frac{\partial}{\partial x^i}\}, \{dx^i\}$  are dual basis of  $T_p$  and  $T_p^*$ , then the set:

$$\{\frac{\partial}{\partial x^{a_1}} \otimes \dots \otimes \frac{\partial}{\partial x^{a_r}} \otimes dx^{b_1} \otimes \dots \otimes dx^{b_s}\}$$

is a basis of  $T_s^r(p)$ .

Writing  $T = T^{a_1 \dots a_r}_{b_1 \dots b_s} \frac{\partial}{\partial x^{a_1}} \otimes \dots \otimes \frac{\partial}{\partial x^{a_r}} \otimes dx^{b_1} \otimes \dots \otimes dx^{b_s}$ , it's an exercise for the reader to show that this is in fact a basis. In regions of overlapping where it may be possible to change the coordinates, the coordinate-induced components will change as follows:

**Proposition 1.10.** Let  $p \in M$  and  $(U, x), (V, y)$  be charts such that  $p \in U \cap V \neq \emptyset$ . Using the coordinates-induced basis, the components transforms like:

- (i)  $X_{(y)}^j = X_{(x)}^i \frac{\partial y^j}{\partial x^i}$  for a vector.
- (ii)  $\omega_k^{(y)} = \omega_i^{(x)} \frac{\partial x^i}{\partial y^k}$  for a one-form.
- (iii) products of (i) and (ii) for an  $(r,s)$ -tensor.

*Proof.* In the overlapping region  $U \cap V$ , the vector  $X \in T_p(M)$  can be expressed as

$$X = X_{(x)}^i \frac{\partial}{\partial x^i} = X_{(y)}^j \frac{\partial}{\partial y^j}$$

therefore a simple calculation shows that

$$\frac{\partial}{\partial x^i} f = \frac{\partial f \circ y^{-1}}{\partial y^j} \frac{\partial (y \circ x^{-1})^j}{\partial x^i} = \frac{\partial y^j}{\partial x^i} \frac{\partial}{\partial y^j} f$$

leading to

$$\frac{\partial}{\partial x^i} = \frac{\partial y^j}{\partial x^i} \frac{\partial}{\partial y^j}$$

In terms of the coefficients we see that

$$X_{(x)}^i \frac{\partial y^j}{\partial x^i} \frac{\partial}{\partial y^j} = X_{(y)}^j \frac{\partial}{\partial y^j} \implies X_{(y)}^j = X_{(x)}^i \frac{\partial y^j}{\partial x^i}$$

Using similar techniques you can also show (ii) and (iii). □

**Definition 19.** A  $p$ -form is a completely antisymmetric  $(0,p)$ -tensor.

Let's now pause and remember our recipe. We need to assign an inner product to the tangent spaces. Instead of calling it  $\langle, \rangle$  we will write it as  $g$ :

**Definition 20.** Let  $M$  be a smooth manifold. A metric  $g$  is a choice of a smooth symmetric non-degenerate tensor field of type  $(0,2)$ . In other words:

- (i)  $g(X, Y) = g(Y, X)$
- (ii)  $g(X, Y) = 0$  for all  $X \in T_p(M)$  implies  $Y = 0$

In a chart, using coordinate-basis we know that  $g$  takes the form:

$$g = g_{\mu\nu} dx^\mu \otimes dx^\nu$$

This is the rigorous definition of the metric  $ds^2 = g_{\mu\nu} dx^\mu dx^\nu$  that we gave in section 0. This correspondence is more explicitly seen by taking any two vectors  $X, Y \in T_p$  so that

$$g(X, Y) = g_{\mu\nu} dx^\mu(X) dx^\nu(Y) = ds^2$$

And we are finally able to define the curvature. Let's gather everything we learned in previous sections.

## 1.4 Frame bundle and principal fiber bundles

---

## 2 Gravity enters the scene...

Now that we've explored the intricacies of geometry and curvature, we are in a position to confront the second big question: How exactly does gravity manifest as the curvature of spacetime? How did Einstein come up with this idea in the first place?

[lots to be revised. change the order and add spacetime diagrams](#)

### 2.1 It's simple once you know geometry

The idea that gravity is not a force is needed for Newtonian physics to be consistent! Consider what Newton's first law have to say:

**(Law of Inertia): A body follows uniform straight motion unless acted on by a force.**

Now imagine a universe with only a single particle. How can this particle tell that it is moving at all? Well, besides being an extremely boring universe it can't tell that it's moving. We need at least two particles, one the observer and other the observed. The observer will check with its coordinates and clocks if Newton's law holds for the other particle. But wait a second: both particles have mass, so there is a gravitational force that deviates the uniform straight motion. And if one includes more particles it only gets worse! It is clear that the force of gravity and Newton's first law cannot both be true since this leads to a contradiction. How can we resolve this problem?

It turns out that gravity is not a force, so it does not deviate particles from straight motion in spacetime. We may think this is untrue because clearly a particle in a straight line is very different from an orbit. So for this to be true we must loose a bit our notion

of *straight*, which is equivalent to consider a curved geometry. That is no problem for us! The straightest possible paths are geodesics (locally) in a more general geometry. If we can show that the effect of gravity is the same as a geodesic path in some sort of curved space, there are no more contradictions with Newton's first law since gravity is no longer a force.

Let us look at the simplest physical system where gravity is present: a free falling body in a gravitational potential  $\Phi$ . The equation of motion is:

$$m\ddot{x} = -m\nabla\Phi$$

Notice that the mass appears on both sides. This is the difference of gravity to other forces like the electromagnetic one: all particles follows the same path regardless of the masses. A positively charged particle will follow a trajectory that is very different than a negatively charged one when put on the same electromagnetic field. Rewriting the equation of motion

$$\ddot{x}^i + (\nabla\Phi)^i = 0 \tag{16}$$

We want this to look like a geodesic equation, but first derivatives are lacking on the second term. We are stuck: there is only three dimensions of space and none of them can be put together to build a geodesic out of equation (16). Turns out it's 21 century and there is another dimension we can play with: time. If we include the extra coordinate  $x^0 = t$  of time on our recipe, it obviously obeys  $\ddot{x}^0 = 0$  and we get a system of equations:

$$\begin{cases} \ddot{x}^0 = 0 \\ \ddot{x}^i + (\nabla\Phi)^i \dot{x}^0 \dot{x}^0 = 0 \end{cases}$$

That is exactly a geodesic equation on *spacetime*, not just space! Notice that the law of inertia mentions space *and* time: a uniform straight motion. These are necessarily straight lines on spacetime diagrams, not just straight lines on space<sup>10</sup>. So it should not come as a surprise that we had to include time in the last step.

With the geodesic equation we can obtain the connection  $\Gamma_{00}^i = (\nabla\Phi)^i$ , and  $\Gamma_{jk}^i = 0$  otherwise. With a connection at hand we find the Riemann curvature tensor  $R_{0j0}^i = -\partial_j \partial^i \Phi$ . Therefore we conclude that gravity may be cast as an effect of the curvature of spacetime.

A refinement of Newton's first law is:

**(Enhanced Law of Inertia): In the absence of forces all particles follows geodesics in spacetime.**

Now a small digression for the more advanced reader. The idea above is somewhat related to the Newton-Cartan gravity and [can be constructed](#) more rigorously by defining a Newtonian spacetime  $(M, \tau, h, \nabla)$  where  $M$  is a manifold,  $\tau$  is a closed form called the clock form,  $h$  is a 3 dimensional metric and  $\nabla$  is a Newtonian connection. They satisfy  $\tau_\mu h^{\mu\nu} = 0$  and  $\nabla\tau = \nabla h = 0$ . The form  $\tau$  gives a notion of absolute time  $d\tau = 0$

---

<sup>10</sup>Think about some particle moving on a straight line with constant acceleration.



and causality. Test particles moves on geodesics of  $\nabla$  as well. There is even a field equation relating the mass density to the Riemann tensor:

$$R_{\mu\nu} = 4\pi G \rho \tau_\mu \tau_\nu$$

Although this may look artificial, it turns out that this theory is completely natural once we understand the structure of spacetime with relativity. The Newton-Cartan gravity is the formal non-relativistic limit  $c \rightarrow \infty$  of general relativity. You can sense this by looking at the formal Laurent series of a lorentz metric  $g = g_{\mu\nu} dx^\mu \otimes dx^\nu$  in terms of  $c^{-1}$

$$g = -c^2 \tau \otimes \tau + h - 2\phi \tau \otimes \tau + O(c^{-1})$$

where  $\phi$  is some function. We see the structures  $\tau, h$  pop out of the leading terms of the expansion. For more details please look at the reference.

## 2.2 Einstein's happy thought

If you are not convinced by the arguments above, let's try Einstein's own thought experiment. Imagine that you are on a completely closed elevator with a scale, like any normal person would do. If someone with bad intentions cut the cables, you will see and sense your weight disappearing along with everything in the elevator. For a brief moment it will be just like floating in vacuum away of gravity and danger. Einstein took this idea to a whole new level: he conjectures that there is no physical experiment you can possibly do to differentiate free fall from no gravity at all. Free fall is just another type of inertial motion.

## 2.3 Generalize Special relativity

In the section above we gave some physical arguments to conclude that gravity must be the curvature of some geometry. But it may seem unnatural to mix space and time together, specially because we were dealing with a galilean structure of flat space and absolute time. A more natural arena is, of course, Minkowski spacetime of special relativity where the geometry is explicitly given to us by a metric  $\eta_{\mu\nu}$ . Our goal is to introduce gravity in special relativity. This will inevitably lead us to general relativity.

We know that in special relativity it is best to view space and time together, i.e. spacetime, as a set of points  $X^\mu = (ct, x, y, z)$  that resembles the vector space  $\mathbb{R}^4$  but has a very different notion of distance between points:

$$ds^2 = -c^2 dt^2 + dx^2 + dy^2 + dz^2 = \eta_{\mu\nu} dX^\mu dX^\nu \quad (17)$$

where the metric is defined as  $\eta_{\mu\nu} = \text{diag}(-1, 1, 1, 1)$ .

Much of special relativity is encoded in the expression (17). For instance, the causal structure of spacetime is the statement that causally connected points satisfy  $ds^2 \leq 0$ . When we look for coordinate transformations  $X'^\mu = \Lambda^\mu_\nu X^\nu$  that preserves (17) we obtain the famous Lorentz transformations  $\eta_{\mu\nu} \Lambda^\mu_\alpha \Lambda^\nu_\beta = \eta_{\alpha\beta}$ . They form a group  $\Lambda^\mu_\nu \in SO(1, 3)$  where composition is given by the usual matrix multiplication.<sup>11</sup>

---

<sup>11</sup>The set of coordinate transformations  $X'^\mu = \Lambda^\mu_\nu X^\nu + a^\mu$  that preserves the metric (17) is called

Another consequence is that all equations in SR that does not depend on a particular frame of reference  $X^\mu$  are made with stuff that transforms nicely under Lorentz transformations. By 'stuff' I mean tensor fields, which are objects that we saw in the geometry part, and by transforms nicely I mean the coordinate transformations of 1.10. For example, the proper time:

$$d\tau = \frac{\sqrt{-ds^2}}{c}$$

is the time measured by a particle's clock in its rest frame  $X'^\mu = (ct', 0, 0, 0)$ . Different observers will experience different times  $t$  but they all agree on the value of  $\tau$ . With these ideas in mind we can easily reconstruct the dynamics of particles in a relativistic fashion. We define the 4-velocity and 4-momentum as:

$$U^\mu = \frac{dX^\mu}{d\tau}$$

$$P^\mu = mU^\mu$$

Let  $X^\mu, Y^\mu$  be 4-vectors, which are just  $(1,0)$ -tensor fields, in spacetime. We call  $X_\mu Y^\mu \doteq \eta_{\mu\nu} X^\mu Y^\nu$  the contraction of  $X$  with  $Y$ . It is easy to show that  $U^2 = U_\mu U^\mu = -c^2$  and  $P^2 = -(mc)^2$ . Newton's second law should be something like:

$$\frac{dP^\mu}{d\tau} = F^\mu \quad (18)$$

Although the above expression is neat, it is not clear how to include some force that maintain Lorentz invariance. The problem gets more tamed if we look at the Lagrangian formulation of SR: Since Lorentz transformations are a symmetry of the system, we could look for a relativistic invariant action which is far more approachable. The simplest thing to come up with is the following:

$$S = -m \int d\tau \quad (19)$$

sine we already know that  $\tau$  is an invariant scalar. From here on I'll use natural units where  $c = 1$  just to simplify matters. To get a better grip at equation (19) first notice that a particle describes a path in spacetime  $X^\mu(\sigma) : I \subset \mathbf{R} \rightarrow \mathbf{R}^{1,3}$  parametrized by some  $\sigma$ . The parameter does not have to be the proper time, we can always reparametrize  $\tau = \tau(\sigma)$  the path just like we did when analysing curves in the Geometry part of the notes. In practice this means that we can rewrite equation (19):

$$S = -m \int d\sigma \frac{d\tau}{d\sigma} = -m \int d\sigma \sqrt{-\eta_{\mu\nu} \frac{dX^\mu}{d\sigma} \frac{dX^\nu}{d\sigma}} \quad (20)$$

The action is fully relativistic and has a gauge redundancy given by the reparametrization invariance: changing the parameter should not affect the underlying physics, i.e. the worldline of the particle. You can check this explicitly by making an arbitrary

---

the Poincaré group  $\mathcal{P} \simeq O(1,3) \oplus \mathbf{R}^{1,3}$ . The reader should notice that this is the group of isometries of spacetime. The Lorentz transformations are the subgroup  $SO(1,3)$  of the Poincaré group connected to the identity.

change  $\sigma = \sigma(\lambda)$  in the action. At the end, we want this redundancy in order to calculate the euler-lagrange equations without too much heat. The Euler-Lagrange equations gives us:

$$\frac{d}{d\sigma} \left( m \frac{\partial \sqrt{-\eta_{\mu\nu} \dot{X}^\mu \dot{X}^\nu}}{\partial \dot{X}^\rho} \right) = m \frac{d^2 X_\rho}{d\tau^2} = 0$$

You should check this expression for yourself. In the middle of calculation we switched to proper time where  $-\eta_{\mu\nu} \dot{X}^\mu \dot{X}^\nu = -U^2 = 1$ . The equation above is precisely what we were expecting from our previous equation 18. With the Lagrangian formalism it is straightforward to insert a potential in the theory:

$$S = -m \int d\sigma \sqrt{-\eta_{\mu\nu} \frac{dX^\mu}{d\sigma} \frac{dX^\nu}{d\sigma}} - \int d\sigma \Phi$$

But this does not keep the reparametrization invariance! Any change in the parameter  $\sigma \rightarrow \sigma'$  will be felt by the jacobian  $\frac{\partial \sigma}{\partial \sigma'}$  in the second term. We can get around this by considering a four-potential  $A_\mu$  instead of the scalar potential  $\Phi$ , and contract it with the four-velocity:

$$S = -m \int d\sigma \sqrt{-\eta_{\mu\nu} \frac{dX^\mu}{d\sigma} \frac{dX^\nu}{d\sigma}} - \int d\sigma q A_\mu \dot{X}^\mu$$

where  $q$  is just a constant measuring the coupling with the potential. Notice that a reparametrization don't change the action because of the 4-velocity term! And everything keeps lorentz invariance as well. The suggestive notation is going to make more sense once we derive the Euler Lagrange equations for the system:

$$m \frac{d^2 X_\mu}{d\tau^2} = q \left( \frac{\partial A_\mu}{\partial X^\nu} - \frac{\partial A_\nu}{\partial X^\mu} \right) \dot{X}^\nu = q F_{\mu\nu} \dot{X}^\nu$$

This looks exactly like the (covariant) equation of a particle in an electromagnetic field, where  $F_{\mu\nu} = \partial_\nu A_\mu - \partial_\mu A_\nu$  is the electromagnetic tensor! You can easily show that it satisfies  $\partial_{[\rho} F_{\mu\nu]} = 0$  which is equivalent to the two homogeneous Maxwell's equations. Of course, it will only be Maxwell's EM when the other two inhomogeneous equations are provided  $\partial^\mu F_{\mu\nu} = J_\nu$  where  $J_\nu = (\rho, \mathbf{J})$  is the 4-current.

Does that mean that relativistic gravity is somehow a kind of gravitomagnetism? You could go on this route but eventually you would find serious inconsistencies in the solutions of the equations. Stuff like infinite energy, instability of simple closed orbits and worse.(reference here)

What now? Maybe we could go on with the previous idea and insert a (0,2)-tensor field  $h_{\mu\nu}$  instead of the 4-potential  $A_\mu$  and see where that leads us. We would actually get linearized general relativity! But the procedure is much more involved and subtle. Since this is a field theory we would like to have positive kinetic terms (free of ghosts, if these words even make sense) and kill some degrees of freedom. In the appendix we show how you can achieve this, but for know we go through a much more simple route.

We already had a clue of the geometric nature of gravity, so instead of subtracting a potential we change the metric:

$$S = -m \int d\sigma \sqrt{-g_{\mu\nu} \frac{dX^\mu}{d\sigma} \frac{dX^\nu}{d\sigma}} \quad (21)$$

where now  $g_{\mu\nu} = \eta_{\mu\nu}$  except at the 00 component  $g_{00} = \eta_{00} - 2\Phi$ . This choice of lagrangian is inspired by what we did in the section above since this leads to the same connection  $\Gamma_{00}^i = (\nabla\Phi)^i$ .

*Exercise: Show that in the non relativistic limit  $v \ll c$  the Lagrangian in (21) simplifies to  $L \approx mv^2/2 - m\Phi$*

The Lagrangian (21) not only reproduces all of the effects of gravity, but also introduces some new phenomena. Take a photon with a certain frequency and shoot up from the bottom of a high building. At the top, someone measures the same photon and realizes the frequency has changed! The gravitational redshift is predicted by the new metric  $g_{\mu\nu}$ . The photon 4-momentum is:

$$k^\mu = (E, 0, 0, E)$$

and the 4-velocities of the observers at the bottom and top of the building are  $U_{\text{bottom}}^\mu = U_{\text{top}}^\mu = (1, 0, 0, 0)$ . The ratio of the measured frequencies is:

$$\frac{\omega_2}{\omega_1} = \frac{k^\mu U_\mu(\text{top})}{k^\mu U_\mu(\text{bottom})} = \frac{g_{00} k^0 U^0(\text{top})}{g_{00} k^0 U^0(\text{bottom})} = \frac{1 - \frac{2GM}{r_2}}{1 - \frac{2GM}{r_1}}$$

Where it's clear that the frequency of the top  $\omega_2$  gets redshifted since  $r_2 > r_1$ . At this point there is nothing stopping us from thinking that the metric  $g_{\mu\nu}$  may be anything, or at least anything that is consistent with the matter content of the physical system. Each  $g_{\mu\nu}$  is interpreted as the geometry of spacetime or of just a portion of spacetime. The Euler-Lagrange equations of 21 gives us geodesics just like we did in the Geometry part. We thus arrive at the same conclusions drawn in the section above: gravity is the geometry of spacetime, and free particles move on geodesics. Furthermore, we have the necessary tools of differential geometry to make precise statements of these ideas in a much more general setting.

*Definition: Spacetime  $(M, g)$  is a four-dimensional smooth manifold equipped with a Lorentzian signature metric.*

The four dimensions should be obvious. A smooth manifold structure is the least we can impose for something to look like spacetime without having any kind of weird topological phenomena.<sup>12</sup> It also assures that everything is made without any reference to coordinate systems. A Lorentzian signature is necessary so that we maintain Lorentz invariance locally, just like we showed that a Riemannian metric is trivial  $g_{\mu\nu}(p) = \delta_{\mu\nu}$

---

<sup>12</sup>By topological I mean continuity, paracompactness etc. The paracompactness of a manifold is a sufficient condition for the existence of a Riemannian metric. For a Lorentzian metric we have to impose one further condition: it has to admit a non-vanishing vector field. For a non-compact manifold that is no problem, but for compact manifolds this is equivalent to have zero Euler characteristic. Compact spacetimes could be 4-torus!

at a point. It assures that when gravity (curvature) is not present, we go back to our beloved flat Minkowski spacetime full of stairs and paradoxes.

How can we possibly know the metric from the matter content in this general setting? Just like in electromagnetism where Maxwell's equations relate the electromagnetic field and the charge distribution, we need field equations relating  $g_{\mu\nu}$  to the stress-energy-momentum tensor  $T_{\mu\nu}$ . You see, the energy momentum tensor is just a way of encoding the matter content in a tensorial, relativistic fashion. I will just say what it is: it's a matrix where the  $T_{\mu\nu}$  component is the flux of  $P^\mu$  momentum across the hypersurface of constant  $x^\nu$ . For example, the 00 component is just the energy density where the 11 component is the pressure along the x direction

$$T_{00} = \frac{P^0}{\Delta X \Delta Y \Delta Z} = \rho$$

$$T_{11} = \frac{P^1}{\Delta T \Delta Y \Delta Z} = p_x$$

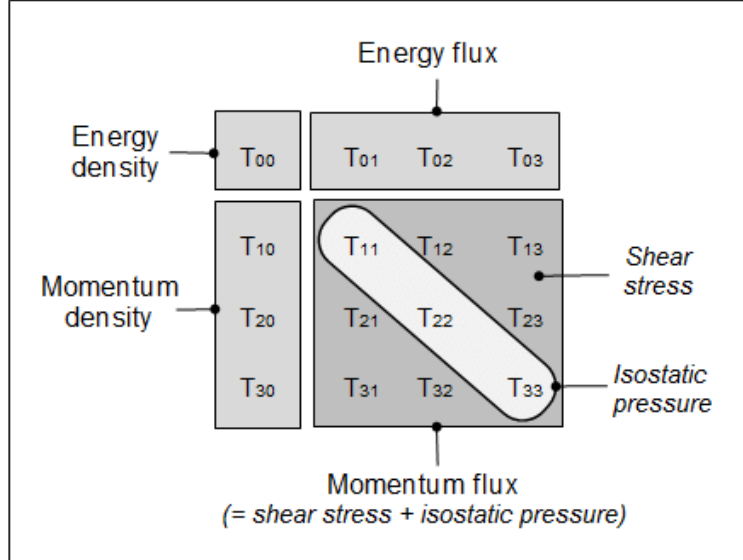


Figure 2: Energy momentum in all of its glory jumpscare.

If you know some field theory you may think that  $T_{\mu\nu}$  is the Noether current associated to spacetime translations. This is not exactly the same tensor we use in general relativity, since the former is defined only in the context of special relativity. Hence it is not covariant under general coordinate transformations.

Recall that in non-relativistic mechanics, Poisson's equation  $\nabla^2 \Phi = 4\pi G \rho$  determines the gravitational potential of a source. We know that  $\Phi$  is related to the 00 component of the metric tensor, so a natural generalization of Poisson's equation is

$$\nabla^2 g_{\mu\nu} = k T_{\mu\nu}$$

But the covariant derivative of the metric is zero by metric compatibility. Furthermore, we want the energy momentum tensor to be conserved  $\nabla_\mu T^{\mu\nu} = 0$ . On the left

hand side we should have a symmetric tensor such that  $\nabla_\mu G^{\mu\nu} = 0$  where  $G$  is made of second derivatives of the metric. We know just the guy to do it!

$$R_{\mu\nu} - \frac{R}{2}g_{\mu\nu} = kT_{\mu\nu}$$

Now we can make the non-relativistic, static weak field approximation to get the constant  $k$ . I'll let you finish the details!

**Definition 21.** *Our spacetime  $(M, g)$  must satisfy the Einstein's field equations:*

$$R_{\mu\nu} - \frac{R}{2}g_{\mu\nu} = \frac{8\pi G}{c^2}T_{\mu\nu}$$

Some comments must be draw. Notice that I switched back to normal units of  $c$  just so you remember where the right stuff goes to. We could also add a constant term of the form  $\Lambda g_{\mu\nu}$  on the left-hand side, where  $\Lambda$  is called the cosmological constant. As it turns out this terms fits the data of cosmological observations.

If you are stubborn (like you should) you want an action for spacetime instead of the the Einstein's equations. Here it is:

$$S = \frac{1}{16\pi G} \int d^4x \sqrt{-g} R \tag{22}$$

[finish section, black holes](#)

## References