# 4

# GRAVITATIONAL  FIELD  THEORY

**Introduction**. It is by "instantaneous action at a distance"—more specifically: by exerting attractive central forces of strength

$$F = Gm_1m_2/r^2 \qquad (1)$$

upon one another—that, according to Newton (*Principia Mathematica* 1686), bodies interact gravitationally. Building upon this notion (to which many of his continental contemporaries objected on *philosophical* grounds), he was able to account theoretically for Kepler's emperical "laws of planetary motion" (1610) and to lay the foundation for a famously successful celestial mechanics.[1]

Special Relativity (1905) declared Newtonian dynamics to be defective because not Lorentz covariant, and Newton's Law of Universal Gravitation to be untenable because it drew upon a concept—"distant simultaneity"— which relativity had rendered obsolete. It placed upon physicists the burden of devising a "relativistic theory of gravitation"...not to account for some disagreement between Newtonian theory and observation (of which, if we set aside a little problem concerning the precession of Mercury's orbit, there were none), but to achieve logical consistency.

---

[1] "Celestial mechanics" is an archaic term for what we would today call "planetary mechanics," and contains an echo of Newton's question (fairly radical for the time): Does gravity—the gravity which causes apples to fall— extend all the way to the moon? Is the moon "falling" around the earth? For a good account of this ancient history (the details of which are much more convoluted/interesting than one might imagine) see G. E. Christenson, *This Wild Abyss: The Story of the Men Who Made Modern Astronomy* (1978).

   At p. 1265 in the index of Misner, Thorne & Wheeler's *Gravitation*[2], under the heading "Gravitation, theories of," one encounters mention of (among others)

- Bergmann's theory
- Cartan's theory
- Coleman's theory
- Kustaanheimo's theory
- Ni's theory
- Nordstrøm's theory
- Papapetrou's theory
- Whitehead's theory
- the flat spacetime theories of Gupta, Kraichnan, Thirring, Feynman, Weinberg, Deser and others

and at p. 1049 a description of a "parameterized post-Newtonian formalism" which was, as an aid to experimentalists, developed in the 1970's to provide *simultaneous* expression of most of the theories listed above...each of which was originally put forward in response to the logical consistency problem just mentioned. Consulting the list of Einstein's publications[3] we find that he himself first addressed the problem in 1907; a second paper appeared in 1911, five more in 1912, and those (though he remained intensely involved in a variety of other problems) were followed by a flood of gravitational papers up until the publication, in 1915, of an account of "general relativity" in its finished form.

   Einstein's own point of departure was provided by what he called "the happiest thought of my life."[4] It was November of 1907 and Einstein was, as he later wrote,

> "...*sitting in a chair in the patent office at Bern when all of a sudden a thought occurred to me: 'If a person falls freely he will not feel his own wwight.' I was startled. This simple thought made a deep impression upon me. It impelled me toward a theory of gravitation.*"

Thus did the Principle of Equivalence spring into being. It holds "gravitational force" and "the fictitious force which arises from acceleration relative to the local inertial frame(s)"—in short: it holds gravitation and non-inertiality—to be physically indistinguishable; i.e., to be different names for the same thing. The Principle of Equivalence contributed little or not at all to most of the early efforts to construct a "relativistic theory of gravitation," but exerted a

---

   [2] I will have frequent occasion to refer to this "Black Bible" of gravitation theory, which appeared in 1973 but remains in many respects definitive. I will, on such occasions, use the abbreviation *MTW*.

   [3] Such a list—not quite complete—can be found on pp. 689–760 in Paul Schlipp's *Albert Einstein: Philosopher-Scientist* (1951). See also C. Lanczos, *The Einstein Decade (1905–1915)* (1974) and A. Pais, *'Subtle is the Lord...': The Science and the Life of Albert Einstein* (1982) for annotated bibliographic information.

   [4] "...die Glücklichste Gedanke meines Lebens;" see Chapter 9 in Pais.[3]

powerful guiding influence upon Einstein's own thought.[5] It was, however, not immediately evident just *how* his happy thought was to be folded into a fully developed theory of gravitation; the journey from special to general relativity took nearly a decade to complete, and was marked by many hesitations, retreats, amendations. The voyage reached its end on 25 November 1915, when Einstein submitted to the Prussian Academy a paper which presented the gravitational field equations in their final form. Five days previously, David Hilbert had submitted to Gesellschaft der Wissenschaft in Göttingen a manuscript containing identical equations.[6]

**Field-theoretic aspects of Newton's theory of gravitation**. Newton's theory was presented as a theory of 2-body interaction. But it can, by importation of concepts and methods borrowed from electrostatics, readily be portrayed as a rudimentary field theory, and it is in that form that it is most conveniently compared to the full-blown field theories which would supplant it, and from which it must (for well-established observational reasons) be recoverable as the leading approximation.

Let the density function $\rho(\boldsymbol{\xi})$ describe (relative to an inertial frame) the instantaneous distribution of "gravitating matter," and let a test particle of mass (which is to say: of "gravitational charge") $m$ reside momentarily at $\boldsymbol{x}$. To describe the force experienced by the test particle write

$$\boldsymbol{F} = m\boldsymbol{g}(\boldsymbol{x}) \tag{2}$$

---

[5] It accounts, in particular, for the circumstance that his second gravitational paper (1911) bore the title "Bemerkung zu dem Gesetz von Eötvös." The "Eötvös experiments" (1889 and 1922) looked to the relationship of "inertial mass" to "gravitational mass (or charge)" and established that the ratio

$$\frac{\text{gravitational mass}}{\text{inertial mass}}$$

is "universal" in the sense that it does not vary from material to material by more than 5 parts in $10^9$. In the 1960's Robert Dicke used more modern techniques to establish that departures from the Principle of Equivalence cannot exceed one part in $10^{11}$. See *MTW* §38.3 for more detailed discussion.

[6] See Pais,[3] §14d. Hilbert, who considered physics to be "too difficult for physicists," imagined himself to be constructing an axiomatic theory of the world (an ambition which Einstein considered to be "too great an audacity... since there are still so many things which we cannot yet remotely anticipate"), and in his grandly titled "Die Grundlagen der Physik" imagined that he had achieved a unified theory of gravitation and electromagnetism. Hilbert's theory is distinguished most importantly from Einstein's by the more prominent role which he assigned to variational principles; we recall that he had retained Emmy Noether to assist him in this work, and it is Noether whom in 1924 he credited for some of his paper's most distinctive details.

Evidently $\boldsymbol{g}(\boldsymbol{x})$ is the gravitational analog of an electrostatic $\boldsymbol{E}$-field.  The force-law proposed by Newton is conservative, so $\boldsymbol{\nabla}\times\boldsymbol{g} = \boldsymbol{0}$, from which follows the possibility of writing

$$\boldsymbol{g} = -\boldsymbol{\nabla}\varphi \tag{3}$$
$$[\varphi] = (\text{velocity})^2$$

In mimicry of the electrostatic equation

$$\boldsymbol{\nabla}{\cdot}\boldsymbol{E} = \rho \quad : \quad \begin{cases} \text{charge density regulates the divergence} \\ \text{of the electrostatic field} \end{cases}$$

we write

$$\boldsymbol{\nabla}{\cdot}\boldsymbol{g} = -4\pi G\rho \quad : \quad \begin{cases} \text{mass density regulates the convergence} \\ \text{of the gravitostatic field} \end{cases} \tag{4}$$

where the minus sign reflects the fact that the gravitational interaction is *attractive*, and the $4\pi$ was inflicted upon us when Newton neglected to install a $\frac{1}{4\pi}$ in (1). Introducing (3) into (4) we have the gravitational Poisson equation

$$\nabla^2\varphi = 4\pi G\rho \tag{5}$$

which in integral formulation reads

$$4\pi G\iiint_{\mathcal{R}} \rho\, dx^1 dx^2 dx^3 = \text{total mass interior to } \mathcal{R}$$
$$= \iiint_{\mathcal{R}} \nabla^2\varphi\, dx^1 dx^2 dx^3$$
$$= -\iiint_{\mathcal{R}} \boldsymbol{\nabla}{\cdot}\boldsymbol{g}\, dx^1 dx^2 dx^3$$
$$= -\iint_{\partial\mathcal{R}} \boldsymbol{g}{\cdot}\boldsymbol{dS}$$
$$= \text{gravitational influx through } \partial\mathcal{R} \tag{5}$$

Take $\mathcal{R}$ to be, in particular, a sphere centered on a point mass $M$; then (5) gives $4\pi r^2 g(r) = 4\pi GM$ whence $g(r) = GM/r^2$ and we have

$$\left. \begin{aligned} \boldsymbol{g}(r) = -GM\,\hat{\boldsymbol{r}}/r^2 = -\boldsymbol{\nabla}\varphi(r) \\ \varphi(r) = -GM/r \end{aligned} \right\} \tag{6}$$

which describe the gravitational field of an isolated point mass.  For a distributed source we have

$$\varphi(\boldsymbol{x}) = -G\iiint \left\{ \rho(\boldsymbol{\xi})/|\boldsymbol{x} - \boldsymbol{\xi}| \right\} d\xi^1 d\xi^2 d\xi^3 \tag{7.1}$$

which gives back (6) in the case $\rho(\boldsymbol{\xi}) = M\delta(\boldsymbol{\xi})$.

The distinction between "gravitostatics" and "gravitodynamics" did not exist for Newton, since he considered gravitational effects to propagate instantaneously. To describe the gravitational potential engendered by a *moving* mass distribution he would have written

$$\varphi(\boldsymbol{x}, t) = -G \iiint \left\{ \rho(\boldsymbol{\xi}, t) / |\boldsymbol{x} - \boldsymbol{\xi}| \right\} d\xi^1 d\xi^2 d\xi^3 \qquad (7.2)$$

which is to say: he would simply have repeated (7.1) at each incremented value of $t$. Such a program makes no provision for the "retardation" effects which distinguish electrodynamics from electrostatics.

To describe the motion of a mass point $m$ in the presence of such an imposed field, Newton writes

$$m\ddot{\boldsymbol{x}} = -m\boldsymbol{\nabla}\varphi(\boldsymbol{x}, t) \qquad (8)$$

from which the $m$-factors (an inertial mass on the left, a gravitational charge on the right) drop away.

To complete the theory Newton would be obligated by his own $3^{\text{rd}}$ Law to describe the action of $m$ back upon $\rho(\boldsymbol{x}, t)$—else to argue that it can, in the specific instance, be neglected—and, more generally, to construct (borrow from fluid dynamics?) a field equation descriptive of $\rho(\boldsymbol{x}, t)$; the latter assignment presents one with a continuous analog of the gravitational $n$-body problem... where the "problem" is not to write but to *solve* the equations of motion.

**Special relativistic generalizations of Newtonian gravitation**. If we look upon the potential $\varphi$ as the object most characteristic of Newtonian gravity—i.e., if we imagine ourselves to be looking for a *relativistic scalar field theory* which gives back Newton's theory in the non-relativistic limit, then it becomes natural in place of (5) to write

$$\left\{ \left( \tfrac{1}{c} \tfrac{\partial}{\partial t} \right)^2 - \nabla^2 \right\} \varphi(x) = -4\pi G \rho(x) \quad ; i.e. \quad \Box \varphi = -4\pi G \rho$$

since this familiar equation is manifestly covariant, and gives back (5) in the limit $c \uparrow \infty$. In place of (7.2) one would then obtain

$$\varphi(x) = -G \iiiint D_R(x - \xi) \rho(\xi) d\xi^0 d\xi^1 d\xi^2 d\xi^3$$

where the retarded Green's function $D_R(x - \xi)$ vanishes except on the lightcone which extends backward from $x$.[7]

Relativistic generalization of (8) is more interesting because a bit less straightforward. Notice first that we can *not* simply write $ma^\mu = -m\partial^\mu \varphi$

---

[7] See ELECTRODYNAMICS (1980), pp. 379–389 for details.

because the Minkowski force on the right is not velocity-dependent, therefore cannot satisfy $K \perp u$, as required.[8] We are led thus to write

$$ma^\mu = K^\mu \quad \text{with (tentatively)} \quad K_\mu = m(\partial_\alpha \varphi)[\delta^\alpha{}_\mu - u^\alpha u_\mu/c^2]$$

because ($i$) the proposed $K_\mu$ depends linearly upon the derivatives of $\varphi$ and ($ii$) clearly does yield $K_\mu u^\mu = 0$. Noticing that we have

$$\begin{aligned} K_\mu &= m\big\{\varphi_{,\mu} - (1/c^2)u_\mu \tfrac{d}{d\tau}\varphi\big\} \\ &= m\big\{\varphi_{,\mu} - (1/c^2)\tfrac{d}{d\tau}(u_\mu\varphi) + (1/c^2)\varphi a_\mu\big\} \end{aligned}$$

and that the final term on the right (since it is itself normal to $u$) could be abandoned without compromising the normality of what remains . . . we do so, obtaining a refined equation of motion

$$\tfrac{d}{d\tau}\big[m(1 + \varphi/c^2)u^\mu\big] = m\partial^\mu\varphi$$

which can be notated

$$\gamma\tfrac{d}{dt}\big[m(1 + \varphi/c^2)\gamma\begin{pmatrix} c \\ \boldsymbol{v} \end{pmatrix}\big] = m\begin{pmatrix} \tfrac{1}{c}\partial_t\varphi \\ -\boldsymbol{\nabla}\varphi \end{pmatrix}$$

The spatial part of the preceding equation gives back (8) in the non-relativistic limit $c \uparrow \infty$, which was the point of the exercise. The form of the equation makes it natural to introduce

$$m^* \equiv m(1 + \varphi/c^2) \equiv \text{effective inertial mass}$$

It is interesting that in this theory "gravitational charge" $m$ is invariable, while the "effective inertial mass" is environmentally contingent, and that the two do *not* cancel each other out.[9]

If, on the other hand, we look upon the gravitational 3-vector $\boldsymbol{g}$ as the object most characteristic of Newtonian gravity then it becomes natural to suppose that it possesses a heretofore unnoticed companion $\boldsymbol{h}$—a gravitational analog of magnetism—and to proceed in direct imitation of Maxwellian electrodynamics to an "antisymmetric tensor theory of relativistic gravitation."

These and a number of other purported "special relativistic generalizations of Newtonian gravitation" are discussed in Chapter 7 of *MTW*. All require formal flights of fancy which take one away from the secure observational base of the theory, all at one point or another are contradicted by the observational facts, and some have been found to be internally inconsistent. All, that is, except some of the most recent, which were found—somewhat surprisingly—to be round-the-bush *reconstructions of Einstein's general relativity*, and not the intended alternatives to it. I suspect that Einstein himself considered all such efforts misguided for the simple reason that, in sanctifying special relativity, they created no place at the table for the Principle of Equivalence—no place for what he sometimes called the "relativity of non-uniform motion."

---

[8]  See again the discussion preliminary to (3–51).

[9]  For further discussion see RELATIVISTIC DYNAMICS (1967), pp. 17–20.

**Theoretical program evidently implicit in the Principle of Equivalence.** Here in this laboratory we erect a Cartesian $\boldsymbol{y}$-frame, which we naively take to be an inertial frame, with 3-axis pointing up. We have turned off all force fields except gravity

$$\boldsymbol{g} = \begin{pmatrix} 0 \\ 0 \\ -g \end{pmatrix}$$

(for which we could find no switch, no shielding...else we would have turned if off too). At $t = 0$ we launch pellets in all directions, with all velocities. Each
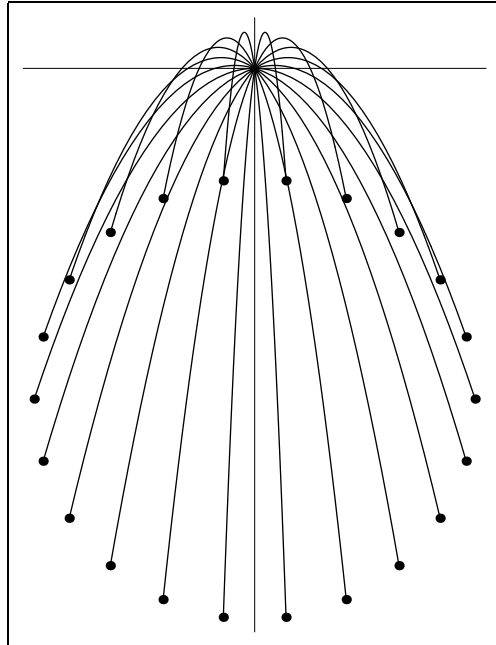


FIGURE 1: *Pellets are launched (all with the same initial speed) in various directions at $t = 0$, and trace parabolic arcs as they fall in the uniform gravitational field. In the figure they have been arrested at the same instant, and are seen to lie on a circle. The envelope of the family of trajectories appears to be parabolic, with the origin at the focus.*

pellet traces a parabolic arc, and the whole display looks like the 4$^\text{th}$ of July. To describe any particular pellet we write $m\ddot{\boldsymbol{y}} = m\boldsymbol{g}$ and obtain $\boldsymbol{y}(t) = \boldsymbol{v}t + \frac{1}{2}\boldsymbol{g}t^2$, of which

$$\begin{pmatrix} y^1(t) \\ y^2(t) \\ y^3(t) \end{pmatrix} = \begin{pmatrix} v^1 t \\ v^2 t \\ v^3 t - \frac{1}{2}gt^2 \end{pmatrix}$$

provides a more explicit rendition. Now introduce the Cartesian $\boldsymbol{x}$-frame of an observer who (irrotationally) *drops* from the origin at the moment of the

explosion. To describe the time-dependent relationship between the $\boldsymbol{y}$-frame and the $\boldsymbol{x}$-frame we write $\boldsymbol{y} = \boldsymbol{x} + \frac{1}{2}\boldsymbol{g}t^2$ and from $m\ddot{\boldsymbol{y}} = m(\ddot{\boldsymbol{x}} + \boldsymbol{g}) = m\boldsymbol{g}$ conclude that the falling observer writes $m\ddot{\boldsymbol{x}} = \boldsymbol{0}$ to describe the motion of the pellets, which he sees to be moving uniformly and rectilinearly (which is to say: in accordance with Newton's $1^{\text{st}}$ Law): $\boldsymbol{x}(t) = \boldsymbol{v}t$. The relationship between our view of the display and the view presented to the falling observer is illustrated in Figure 2.

Einstein argues that it is the observer in (irrotational) free fall who is the *inertial* observer in this discussion; that we in the laboratory have had to
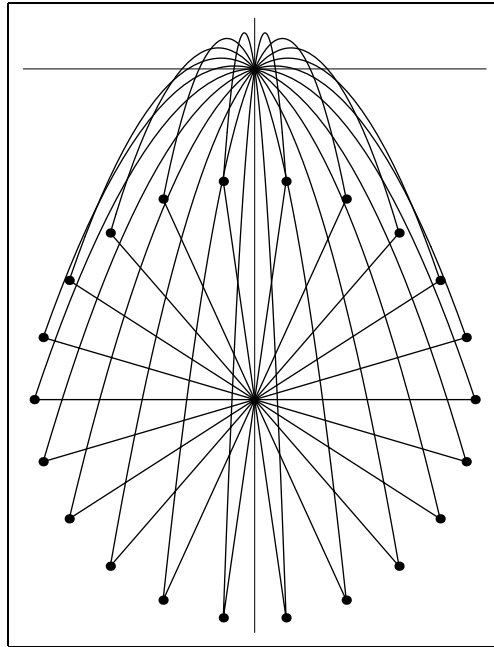


FIGURE 2: *A falling observer experiences no gravitational field, and sees the pellets to be in radial recession. Knowing them to have been launched with the same speed, he is not surprised to observe that they lie on a uniformly expanding circle (sphere).*

"invent gravity" in order to compensate for the circumstance that, relative to the inertial observer, we are accelerating upwards with acceleration $g$, and that therefore we should not be surprised when it is "discovered" that what we misguidedly call "gravitational charge" is proportional through a universal factor to inertial mass.[10]

---

[10] Einstein's viewpoint is nicely developed in §2 of *Spacetime Physics*, by Edwin Taylor & John Wheeler (1963). It is advanced on somewhat different grounds in Chapter 1 ("Ground to Stand on: Inertiality & Newton's First Law") of CLASSICAL MECHANICS (1983).

The perception of uniform/rectilinear pellet motion—i.e., of the absence of a gravitational field—would be shared also by other observers $\mathcal{O}', \mathcal{O}'', \ldots$ who are in states of (irrotational) unaccelerated motion with respect to our falling observer $\mathcal{O}$. It is this *population* of observers which in Newtonian physics is interlinked by Galilean transformations, and in relativistic physics by Lorentz transformations.

We nod indulgently at $\mathcal{O}$'s account of events, then observe that "Of course, you will see pellets (and the observers who ride them) to move uniformly/rectilinearly only for awhile—only until they have ventured far enough away to sense non-uniformity of the gravitational field, by which time tidal effects will have begun to distort their spherical pattern."[11] We might write

$$\boldsymbol{g}(\boldsymbol{x}) = \boldsymbol{g}(\boldsymbol{0}) + \boldsymbol{g}_i(\boldsymbol{0})x^i + \tfrac{1}{2}\boldsymbol{g}_{ij}(\boldsymbol{0})x^i x^j + \cdots$$
$$= \boldsymbol{g} + \text{tidal terms}$$

to underscore the force of our remark.

$\mathcal{O}$'s—Einstein's—response to our remark must necessarily be radical, for he cannot reasonably enter into discussion of the higher order properties of something he has already declared to be a delusion. Einstein's embrace of the Principle of Equivalence leads him...

  • to borrow from Newton the idea that inertial motion is the motion that results when all forces/interactions that can in principle be turned off/shielded have been;[12]

---

[11] What we have initially in mind is simply that the "down vector" here is not parallel to the "down vector" over there, but that remark is subject to interesting refinement. Assume the earth to be a (non-rotating) homogeneous sphere of radius $R$ and mass $M$. The pellets then move in an attractive central force of strength

$$g(r) = \begin{cases} GM\,r^{-2} & \text{if } r \geqslant R \\ (GM/R)\,r^{+1} & \text{if } r \leqslant R \end{cases}$$

The "parabolic arcs" mentioned previously are really sections of Keplerean ellipses (or hyperbolæ, if $v$ is sufficiently great), with the earth's center at one focus. While passing through the earth (which pellets do because we have "turned off all force fields," and "observers" find easy to do because they are mythical) the pellets move as though attached to an isotropic spring, and trace a section of an ellipse with *center* at the center of the earth.

[12] This simple-seeming thought rests upon some heavy idealization. It is one thing to *imagine* turning off the phenomenological forces which would impede a pellet's passage through the earth, but how in practical fact would one turn off the fundamental interactions which underlie those phenomenological forces? We are in something like the predicament of the classical physicist who finds it convenient to "turn off quantum mechanics," and is forced to pretend that he has not thereby precluded the existence of meter sticks.

- to recognize that "gravitational forces" *cannot* be turned off/shielded (can only be *transformed* away); at this point "inertial motion" has become synomymous with what Newton called "free fall;"

- to recognize that inertial observers can engage in special relativistic dialog only so long as—fleetingly—they share the same spatio-temporal "neighborhood," where the bounds of neighborhood are breached when the relative motion of $O$ and $O'$ is no longer uniform/rectilinear; in the absence of gravitation all neighborhoods would be co-extensive and infinite (i.e., there would be only one neighborhood, and we could dispense with the concept), but in the presence of gravitation they become local... like little platelets tangent to a curved surface.

Newton's "rectilinear" can be phrased "geodesic" in Euclidean 3-space. His "uniform rectilinear" could be similarly phrased if 4-dimensional spacetime were suitably metrized, and (as Minkowski was the first to emphasize) it was such a "suitable metrization of spacetime" which lay at the heart of Einstein's accomplishment when he invented special relativity. It became therefore natural for Einstein

- to associate the inertial motion of free-falling test particles with *geodesics in a spacetime of suitably altered geometry*.

Newton's assertion that "masses cause one another to depart from inertiality by exerting gravitational forces upon one another" becomes, from this point of view, an assertion that masses cause *no* "departure from inertiality," but instead *alter the geometry* of the spacetime upon which their respective inertial geodesics are inscribed. But Einstein had already established the equivalence of mass and energy. It became therefore natural for him

- to anticipate that the geometry of spacetime is conditioned by the distribution of mass/energy (which is itself in free-fall controlled by the geometry: the world has thus become "self-interactive geometry").[13]

Prior to (and well into) 1912 Einstein had concentrated on *scalar* theories of gravitation. He had achieved what he considered to be good results in the static theory, but was finding the dynamical theory to be "devilishly difficult." On August 10 Einstein registered as a resident of Zürich, to which he had, upon the invitation of Marcel Grossmann, moved from Prague in order to accept an

---

[13] I must emphasize that it is *as a rhetorical device*—the better to clarify my expository intent—that I have allowed myself to impute motivations to Einstein for which I can, in some instances, provide no specific documentation. My remarks are not (!) intended to be read as "encapsulated history of science." Made-up history is the worst kind of history, and always does violence to the entangled facts. In the present instance it would be a mistake to lose sight of the fact that it took Einstein *several years* to create general relativity, that during those years his motivation was marked by frequent shifts and turns, and that it was not entirely clear where he was headed until he got there.

appointment to the faculty of the ETH. Grossman (1878–1936), it is invariably remarked, had loaned class notes to Einstein when both were students at the ETH. His father had helped Einstein gain employment at the patent office in Bern, while he himself had gone on to become a professor of geometry and (recently) dean of the mathematics–physics section at the ETH. It was, according to Pais,[14] sometime between August 10 and August 16 that Einstein pleaded "Grossmann, Du musst mir helfen, sonst werd' ich verrückt!"[15] and was made aware for the first time of Riemannian geometry, and of the tensor analysis of Ricci and Levi-Civita. Whereupon Einstein recalled that he had, in fact, already been exposed to the Gaussian theory of surfaces in the classroom on one Geisler (whose successor at the ETH was Weyl).[16]

Thus it came about that the scalar gravitation of a Saturday had become a tensor theory by the next Friday. In October of 1912 Einstein wrote to Sommerfeld that

> "At present I occupy myself exclusively with the problem of gravitation and now believe that I shall master all difficulties with the help of a friendly mathematician here [Grassmann]. But one thing is certain: in all my life I have labored not nearly as hard, and I have become imbued with great respect for mathematics, the subtler part of which I had in my simple-mindedness regarded as pure luxury until now. Compared with this problem, the original relativity is child's play."

Einstein had entered upon what were to be three years of the most intense work of his life. We have now to consider what he was up to.

---

[14] See *Subtle is the Lord*,[3] p. 210.

[15] Grossmann, you must help me or else I'll go crazy!

[16] Grossmann had not previously published in the areas in question, but had a good academic's familiarity with developments in his field. Einstein, on the other hand, did not possess a deep command of the literature, and was unaware that aspects of his train of thought had been anticipated decades before. In 1854 Riemann (1826–1866), in a *Habilitation* lecture entitled "Über die Hypothesen welche der Geometrie zugrunde liegen," had suggested that matter might be the cause of geometrical structure, and had in support of that conception described the outlines of "Riemannian geometry." That work was not published until 1867—the year following Riemann's death. In 1873 Clifford (1845–1879) arranged for an English translation of Riemann's essay to be published in *Nature*, and in his own "On the space-theory of matter" (1876) carried the idea even further: by the time he wrote Chapter 4 of *The Common Sense of the Exact Sciences* (1885) he was prepared to argue that not only mechanics but also electrodynamics—the whole of classical physics—are manifestations of the curvature of space. Similar ideas (if somewhat differently motivated) were advanced by Hertz (1857–1894) in his *The Principles of Mechanics* (1894). But these prescient thinkers worked in ignorance of special relativity, so contemplated the physical geometry not of spacetime but of space. Nor were they in position to draw guidance from the Principle of Equivalence.

General relativity lies on the other side of the mathematical thicket before which we now stand, and through which we—like Einstein in 1912, and like every student of gravitation since 1915—are obliged now to thread our way. The path was not yet clearly marked in Einstein's day (though the thicket had been in place—neglected by the generality of mathematicians—since before he was born) but has by now been very clearly mapped by any number of authors. I hesitate to add to that vast literature. Were it my option I would say "Find a book, as highbrow or lowbrow, as abstractly elegant or specifically concrete as seems most comfortable to you, and come back when you have mastered it." But that might take a while, and when you did come back we would almost certainly find that we had developed a language/notation problem. So...I attempt now to visit the principal landmarks, and to describe them in terms calculated to serve my immediate practical needs but which make no claim to mathematical modernity.

<div style="border:1px solid black; text-align:center;">

MATHEMATICAL DIGRESSION
**Tensor analysis on Riemannian manifolds**

</div>

When Riemann devised "Riemannian geometry" (1854)—which he did at the instigation of Gauss, who had selected the least favored of the three *Habilitation* topics proposed by Riemann—he built upon earlier work done by Gauss himself, and was influenced also by the then fairly recent non-Euclidean geometries of Lobachevsky and Bolyai ($\sim$1830)...but managed to say what he had to say entirely without reference to "tensor analysis" (which hadn't been invented yet). The first steps toward the creation of the latter subject were taken by Elwin Christoffel (1829–1900), who was a disciple of Riemann, and in 1882 was motivated to invent what we now call the "covariant derivative." Gregorio Ricci-Curbastro (1853–1925) was for the last forty-five years of his life a professor of mathematical physics at the University of Padua,[17] and it was there that, during the years 1884–1894 and drawing inspiration from Riemann and Christoffel, he single-handedly invented what he called the "absolute differential calculus." During the latter phases of that work he was joined by his student, Tullio Levi-Civita (1873–1941), and together they wrote the monograph "Méthodes de calcul différentiel absolute et leurs applications" (Mathematische Annalen 1900) which brought tensor analysis to a recognizably modern form (though it was not until 1917 that Levi-Civita invented the important concept of "parallel transport"). This accomplishment was not widely applauded, and in some quarters inspired hostility; it was in reaction to Ricci's work that in 1899 Élie Cartan published the paper which laid the foundation for what was to become the "exterior calculus." For an account of

---

[17] Early in his career he had, at the instigation of Betti, published a memoir in *Nuovo Cimento* which introduced Italian physicists to the electrodynamics of Maxwell, and during a post-doctoral year in Munich (1877/8) he had come under the influence of Felix Klein.

the relationship between tensor analysis and the exterior calculus—an account which contains still some echo of that ancient tension—see §1.2 in H. Flanders' *Differential Forms* (1963).

I offer the preceding thumbnail history in order to underscore this fact: in opting to construct a *simultaneous* account of the relevant essentials of Riemannian geometry and tensor analysis I am melding two semi-independent subjects, one of which is fully thirty years older than the other (but both of which had been in place for nearly twenty years by the time Einstein was motivated to draw upon them).

**Metrically connected manifolds**. To start with the concrete: let $x^1, x^2, x^3$ refer to a Cartesian fraame in Euclidean 3-space. To describe the distance between two neighboring points write

$$(ds)^2 = (dx^1)^2 + (dx^2)^2 + (dx^3)^2$$

$$= \delta_{ij} dx^i dx^j \quad \text{with} \quad \|\delta_{ij}\| \equiv \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \tag{9}$$

and agree to call $\delta_{ij}$ the "metric connection." Let equations $x^i = x^i(y^1, y^2, y^3)$ serve to describe the introduction $x \to y$ of an arbitrary (and in the general case curvilinear) recoordinatization of 3-space.[18] Evidently

$$x \to y \quad \text{induces} \quad dx^i \to dy^i = M^i{}_j dx^j \tag{10.1}$$
$$M^i{}_j \equiv \partial y^i / \partial x^j$$

which is to say:

> *coordinate differentials transform as components of a*
> *<u>contravariant tensor</u> of first rank (contravariant vector)*

The inverse transformation $x \leftarrow y$ induces

$$dx^i = W^i{}_j dy^j \leftarrow dy^i \tag{10.2}$$
$$W^i{}_j \equiv \partial x^i / \partial y^j$$

Consistently with the elementary observation that $x \to y \to x$ must be the identity transformation, we (by the chain rule) have

$$W^i{}_a M^a{}_j = \sum_a \frac{\partial x^i}{\partial y^a} \frac{\partial y^a}{\partial x^j} = \partial x^i / \partial x^j = \delta^i{}_j$$

which is to say: $\mathbb{W}\mathbb{M} = \mathbb{I}$. Our assumption that $x \to y$ is invertible can be expressed $\det \mathbb{M} \neq 0$, which assures the existence of $\mathbb{W} = \mathbb{M}^{-1}$.

---

[18] We require $x \to y$ to be invertible (maybe not globally, but at least) in a neighborhood containing the point $P$ where, for the purposes of this discussion, we have elected to live.

The geometrical *meaning* of $(ds)^2$ is clearly independent of coordinatized language we elect to speak when *describing* it. At (9) we spoke in $x$-language. To say the same thing in $y$-language we write

$$(ds)^2 = g_{ij}(y)dy^i dy^j \qquad (11)$$
$$g_{ij}(y) \equiv \delta_{ab}\frac{\partial x^a}{\partial y^i}\frac{\partial x^b}{\partial y^j}$$

which is to say:

$$x \to y \quad \text{induces} \quad \delta_{ij} \to g_{ij} = W^a{}_i W^b{}_j \delta_{ab}$$

In words,

<div align="center"><em>the metric connection transforms as<br>a <u>covariant tensor</u> of second rank</em></div>

In matrix notation we have $\|\delta_{ij}\| \to \|g_{ij}\| = \mathbb{W}^{\mathsf{T}}\|\delta_{ab}\|\mathbb{W}$ from which it follows that

$$x \to y \quad \text{induces} \quad \det\|\delta_{ij}\| \to \det\|g_{ij}\| = W^2 \cdot \det\|\delta_{ab}\|$$
$$W \equiv \det\mathbb{W}$$

In words,

<div align="center"><em>the determinant $g \equiv \|g_{ij}\|$ of the metric connection<br>transforms as a scalar <u>density</u> of weight $w = 2$</em></div>

More generally, we would say of the multiply-indexed objects $X^{ijk}{}_{mn}$ that they transform "as components of a mixed tensor of
- contravariant rank 3 (number of superscripts)
- covariant rank 2 (number of subscripts) and
- weight $w$"

if and only they respond to $x \to y$ by the rule

$$X^{ijk}{}_{mn} \to Y^{ijk}{}_{mn} = W^w \cdot M^i{}_a M^j{}_b M^k{}_c W^d{}_m W^e{}_n X^{abc}{}_{de} \qquad (12)$$

which generalizes straightforwardly to arbitrary covariant/contravariant rank. It remains to be established that geometry/physics present a vast number of objects which *do* transform by this rule (as well as a population of multiply-indexed objects made all the more interesting by the fact that they don't!). And it is important to appreciate the significance of the "as components of" in the sentence which led to (12); it is important, that is to say, not to confuse the geometrical/physical object $\mathfrak{X}$ with the set $X^{ijk}{}_{mn}$ of numbers which, *relative to a coordinate system*, serve to describe it (its "coordinates").

From (12) follow a number of propositions—collectively, the subject matter of "tensor algebra"—which are in each instance either self-evident or so easy to prove that I state them without proof:

• Tensors[19] can be added/subtracted (to yield again a tensor) if and only if they possess the same covariant/contravariant ranks and weight (for otherwise they would come unstuck when transformed).

• $A^{\cdots}{}_{\cdots} = B^{\cdots}{}_{\cdots}$ means $A^{\cdots}{}_{\cdots} - B^{\cdots}{}_{\cdots} = 0$, which if valid in one coordinate system is, by the design of (12), clearly valid in all. Such tensor equations require that $A^{\cdots}{}_{\cdots}$ and $B^{\cdots}{}_{\cdots}$ have the same ranks and weight, and provide coordinatized expression of statements of the form $\mathfrak{A} = \mathfrak{B}$.

• Tensors can be multiplied (to yield again a tensor) even if they possess distinct ranks and weights; the resulting tensor will have
  ∗ contravariant/covariant rank equal to the sum of the respective ranks of the factors
  ∗ weight equal to the sum of the weights of the factors.

• If $A^{\cdots i \cdots}{}_{\cdots j \cdots}$ is a tensor density of contravariant rank $r$ and covariant rank $s$ then $A^{\cdots a \cdots}{}_{\cdots a \cdots}$ transforms as a tensor density of the same weight, and of the respective ranks $r - 1$ and $s - 1$. We say the $^i$ has been "contracted" into the $_j$. Attempted contraction of a superscript into a superscript (or of a subscript into subscript) would, on the other hand, yield an object which fails to transform tensorially.

• It makes transform-theoretic good sense to say of a tensor that it is symmetric/antisymmetric with respect to some specified pair of superscripts/subscripts

$$A^{\cdots i \cdots j \cdots}{}_{\cdots} = \pm A^{\cdots j \cdots i \cdots}{}_{\cdots} \quad \text{or} \quad A^{\cdots}{}_{\cdots i \cdots j \cdots} = \pm A^{\cdots}{}_{\cdots j \cdots i \cdots}$$

but statements of the form $A^{\cdots i \cdots}{}_{\cdots j \cdots} = \pm A^{\cdots j \cdots}{}_{\cdots i \cdots}$ come unstuck when transformed.

• It is in this light remarkable that Kronecker's mixed tensor—defined

$$\delta^i{}_j \equiv \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

in some coordinate system—retains that description in all coordinate systems.

Returning again to (11), we are in position now to understand the coordinate-independence of $(ds)^2$ to be a result of our having contracted the second rank product $dy^i dy^j$ of a pair of contravariant vectors into a second rank covariant tensor $g_{ij}$. And to observe that (as a result of the "Pythagorean symmetry" originally attributed to $\delta_{ij}$) the metric tensor $g_{ij}(y)$ will in all coordinate systems be symmetric:

$$g_{ij} = g_{ji}$$

---

[19] . . . of, it need hardly be added, the *same dimension*. On says of a tensor that it is "$N$-dimensional" if the indices range on $\{1, 2, \ldots, N\}$.

To $g_{ij}(y)$ we assign a contravariant companion $g^{ij}(y)$, defined by the square array of linear equations

$$g^{ia}g_{aj} = \delta^i{}_j$$

The numbers $g^{ij}$ are in effect just the elements of the matrix $\|g_{ij}\|^{-1}$. We observe that necessarily

$$g^{ij} = g^{ji}$$

and that $\det\|g^{ij}\| = 1/g$ transforms as a scalar density of weight $w = -2$.

Recalling that at the beginning of this conversation we placed ourselves in Euclidean 3-space, let the $y$-coordinates be, for purposes of illustration, spherical: then

$$\left.\begin{aligned} x^1 &= r\sin\theta^1\cos\theta^2 \\ x^2 &= r\sin\theta^1\sin\theta^2 \\ x^3 &= r\cos\theta^1 \end{aligned}\right\} \tag{13}$$

give

$$dx^1 = \sin\theta^1\cos\theta^2 \cdot dr + r\cos\theta^1\cos\theta^2 \cdot d\theta^1 - r\sin\theta^1\sin\theta^2 \cdot d\theta^2$$

$$dx^2 = \sin\theta^1\sin\theta^2 \cdot dr + r\cos\theta^1\sin\theta^2 \cdot d\theta^1 + r\sin\theta^1\cos\theta^2 \cdot d\theta^2$$

$$dx^3 = \quad\;\; \cos\theta^1 \cdot dr - \quad\quad\; r\sin\theta^1 \cdot d\theta^1$$

whence (squaring and simplifying; it would in the present context be cheating to simply read from a figure)

$$(ds)^2 = (dr)^2 + r^2(d\theta^1)^2 + (r\sin\theta^1)^2(d\theta^2)^2$$

$$\big|$$

Let us now take up residence *on the 2-dimensional surface* of the sphere of radius $R$ (see the following figure); we then have

$$\downarrow$$

$$(ds)^2 = R^2(d\theta^1)^2 + (R\sin\theta^1)^2(d\theta^2)^2 \tag{14}$$

$$= g_{ij}(\theta)d\theta^i d\theta^j \quad\text{with}\quad \|g_{ij}\| = \begin{pmatrix} R^2 & 0 \\ 0 & (R\sin\theta^1)^2 \end{pmatrix}$$

which describes in $\theta$-coordinates the non-Euclidean metric structure that the spherical surface has inherited from the enveloping Euclidean 3-space. Notice that

$$g \equiv \det\|g_{ij}\| = R^4\sin^2\theta^1 \text{ vanishes at the poles}$$

where $g^{ij}$ therefore fails to exist, essentially because (13) becomes non-invertible on the polar axis.

We have come upon a particular instance of a class of objects—surfaces $\Sigma^2$ embedded in Euclidean 3-space $E^3$—which (especially with regard to their curvature properties) were the subject of Gauß' pioneering *Disquisitiones generales circa superficies curvas* (1827). It is to expose the magnitude of Riemann's accomplishment that I look now very briefly to Gauß' "theory of surfaces."
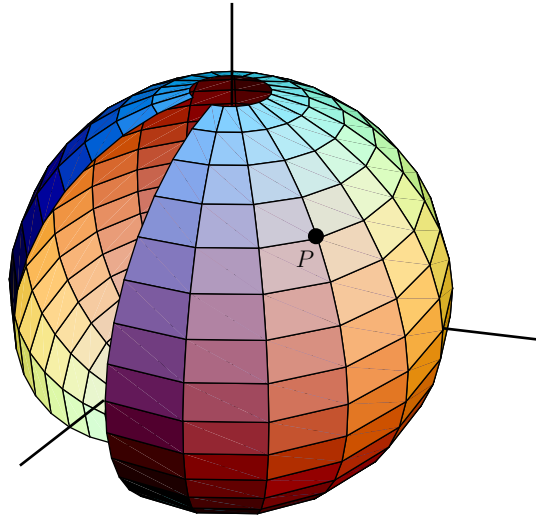
FIGURE 3: *Spherical surface from which polar caps and a sector have been excised to avoid complications which result from the circumstance that the $\{\theta^1, \theta^2\}$ coordinate system (13) becomes singular at the poles (where $\cos\theta^1 = \pm 1$) and is periodic in $\theta^2$. The surface inherits its metric properties from the enveloping Euclidean space.*

Write $\boldsymbol{x}(u, v)$ to present a parametric description of such a surface $\Sigma^2$ (which $u$ and $v$ serve to coordinatize). Let $P$ and $Q$ mark a pair of neighboring points on $\Sigma^2$. To describe the squared Euclidean length of the interval separating $Q$ from $P$, Gauss writes

$$(ds)^2 = d\boldsymbol{x}\cdot d\boldsymbol{x} \quad \text{with} \quad d\boldsymbol{x} = \frac{\partial\boldsymbol{x}}{\partial u}\, du + \frac{\partial\boldsymbol{x}}{\partial v}\, dv + \cdots \equiv \boldsymbol{x}_u\, du + \boldsymbol{x}_v\, dv$$

and obtains

$$
\begin{aligned}
(ds)^2 &= \boldsymbol{x}_u\cdot\boldsymbol{x}_u\, dudu + 2\boldsymbol{x}_u\cdot\boldsymbol{x}_v\, dudv + \boldsymbol{x}_v\cdot\boldsymbol{x}_v\, dvdv \\
&\equiv E\, dudu + 2F\, dudv + G\, dvdv \\
&\equiv \text{FIRST FUNDAMENTAL FORM}
\end{aligned}
\tag{15.1}
$$

of which (14) presents a particular instance. Carrying the expansion of $d\boldsymbol{x}$ to higher order

$$d\boldsymbol{x} = \big\{\boldsymbol{x}_u\, du + \boldsymbol{x}_v\, dv\big\} + \tfrac{1}{2!}\big\{\boldsymbol{x}_{uu}\, dudu + 2\boldsymbol{x}_{uv}\, dudv + \boldsymbol{x}_{vv}\, dvdv\big\} + \cdots$$

Gauss observes that the leading term on the right lies in the plane tangent to $\Sigma^2$ at $P$, and that the second term provides lowest-order indication of how $\Sigma^2$

*curves away from* that plane; taking $\hat{\boldsymbol{n}}$ to be the unit normal to the tangent plane, Gauss constructs

$$
\begin{aligned}
\hat{\boldsymbol{n}} \cdot d\boldsymbol{x} &= \tfrac{1}{2}\big\{\hat{\boldsymbol{n}} \cdot \boldsymbol{x}_{uu}\, dudu + 2\hat{\boldsymbol{n}} \cdot \boldsymbol{x}_{uv}\, dudv + \hat{\boldsymbol{n}} \cdot \boldsymbol{x}_{vv}\, dvdv\big\} \\
&\equiv \tfrac{1}{2}\big\{e\, dudu + 2f\, dudv + g\, dvdv\big\} \\
&\equiv \tfrac{1}{2} \cdot \text{SECOND FUNDAMENTAL FORM}
\end{aligned}
\tag{15.2}
$$

where $\hat{\boldsymbol{n}}$ can be described

$$
\begin{aligned}
\hat{\boldsymbol{n}} &= \frac{\boldsymbol{x}_u \times \boldsymbol{x}_v}{\sqrt{(\boldsymbol{x}_u \times \boldsymbol{x}_v) \cdot (\boldsymbol{x}_u \times \boldsymbol{x}_v)}} = \boldsymbol{x}_u \times \boldsymbol{x}_v \Bigg/ \sqrt{\det\begin{pmatrix} \boldsymbol{x}_u \cdot \boldsymbol{x}_u & \boldsymbol{x}_u \cdot \boldsymbol{x}_v \\ \boldsymbol{x}_v \cdot \boldsymbol{x}_u & \boldsymbol{x}_v \cdot \boldsymbol{x}_v \end{pmatrix}} \\
&= \boldsymbol{x}_u \times \boldsymbol{x}_v \Big/ \sqrt{EG - F^2}
\end{aligned}
$$

Possession of the first and second fundamental forms placed Gauss in position to construct an account of the *curvature of* $\Sigma^2$ *at* $P$ which builds upon the elementary theory of the curvature of plane curves. Consider the planes which stand normal to $\Sigma^2$ at $P$ (in the sense that each contains $\hat{\boldsymbol{n}}$). The intersection of $\Sigma^2$ with such a plane inscribes a curve upon that plane, which at $P$ has curvature $\kappa$. Gauss shows[20] that the least/greatest values of $\kappa$ to be found within that set of planes can be obtained as the respective roots of a quadratic polynomial

$$
(EG - F^2)\kappa^2 - (Eg - 2Ff + Ge)\kappa + (eg - f^2) = 0
$$

which combines information written into the fundamental forms. Let those roots be denoted $\kappa_1$ and $\kappa_2$; one has then the definitions

$$
\begin{aligned}
\text{GAUSSIAN CURVATURE} &\quad : \quad K \equiv \kappa_1 \kappa_2 = \frac{eg - f^2}{EG - F^2} \\
\text{MEAN CURVATURE} &\quad : \quad M \equiv \tfrac{1}{2}(\kappa_1 + \kappa_2) = \frac{Eg - 2Ff + ge}{2(EG - F^2)}
\end{aligned}
\tag{16}
$$

which serve in their respective ways to describe the curvature of $\Sigma^2$ at $P$. Looking to the dimensional implications of the definitions (15) we find that $[\kappa] = (\text{length})^{-1}$, as required by the familiar relation

$$
\text{curvature} = \frac{1}{\text{radius of curvature}}
$$

---

[20] See, for example, §§32–38 in H. Lass, *Vector & Tensor Analysis* (1950) or §§9.14–9.16 in K. Rektorys, *Survey of Applicable Mathematics* (1969) for the elementary details, and accessible surveys of related matters. More modern discussions—such, for example, as can be found in John McCleary, *Geometry from a Differentiable Viewpoint* (1994), Chapter 9 or R. Darling, *Differential Forms & Connections* (1994) §§4.7–4.11—tend generally to be vastly more elegant, but (in proportion to their self-conscious modernism) to impose such heavy formal demands upon the reader as to be almost useless except to people who have already gained a geometrical sense of the subject from other sources.

Look back again, by way of illustration, to the spherical surface shown in Figure 3. We have

$$\boldsymbol{x}(u,v) = R\,\hat{\boldsymbol{n}}(u,v) \quad \text{with} \quad \hat{\boldsymbol{n}} = \begin{pmatrix} \sin u \cos v \\ \sin u \sin v \\ \cos u \end{pmatrix}$$

giving

$$\boldsymbol{x}_u = R \begin{pmatrix} +\cos u \cos v \\ +\cos u \sin v \\ -\sin u \end{pmatrix}, \qquad \boldsymbol{x}_v = R \begin{pmatrix} -\sin u \sin v \\ +\sin u \cos v \\ 0 \end{pmatrix}$$

$$\boldsymbol{x}_{uu} = -\boldsymbol{x}, \qquad \boldsymbol{x}_{uv} = R \begin{pmatrix} -\cos u \sin v \\ +\cos u \cos v \\ 0 \end{pmatrix}, \qquad \boldsymbol{x}_{vv} = R \begin{pmatrix} -\sin u \cos v \\ -\sin u \sin v \\ 0 \end{pmatrix}$$

from which we compute

$$\begin{aligned} E &= R^2 & e &= -R \\ F &= 0 & f &= 0 \\ G &= R^2 \sin^2 u & g &= -R \sin^2 u \end{aligned}$$

from which it follows by (16) that

$$K = \frac{1}{R^2} \quad : \quad \text{all values of } u \text{ and } v$$

The specialness of the result reflects, in an obvious sense, the specialness of spheres.

The functions

$$\begin{aligned} g_{uu}(u,v) &= E(u,v) \\ g_{uv}(u,v) = g_{vu}(u,v) &= F(u,v) \\ g_{vv}(u,v) &= G(u,v) \end{aligned}$$

are "intrinsic to $\Sigma^2$" in the sense that they are accessible to determination by a flat mathematician who has a good understanding of the Euclidean geometry of $E^2$ but no perception of the enveloping $E^3$. The $\hat{\boldsymbol{n}}$ which enters into the definitions of $e(u,v)$, $f(u,v)$ and $g(u,v)$ presents, on the other hand, an explicit reference to the enveloping space. It is, in this light, remarkable—in Latin: *egregium*—that

Gaussian curvature is an *intrinsic* property of $\Sigma^2$

which is the upshot of Gauss' THEOREMA EGREGIUM, the capstone of his *Disquisitiones*. The technical point is that $e$, $f$ and $g$ can be expressed in terms of $E$, $F$ and $G$; when this is done, (16) becomes

$$K = -\frac{1}{4D^4}\begin{vmatrix} E & E_u & E_v \\ F & F_u & F_v \\ G & G_u & G_v \end{vmatrix} - \frac{1}{2D}\left\{\frac{\partial}{\partial v}\frac{E_v - F_u}{D} - \frac{\partial}{\partial u}\frac{F_v - G_u}{D}\right\} \qquad (17)$$

with $D \equiv \sqrt{EG - F^2}$.[21] Gauss stressed also the fact that (17) is structurally invariant with respect to recoordinatizations

$$\left.\begin{matrix} u \\ v \end{matrix}\right\} \longrightarrow \left\{\begin{matrix} u' = u'(u, v) \\ v' = v'(u, v) \end{matrix}\right.$$

of the surface $\Sigma^2$.

   To consult the literature is to encounter the names of people like G. Mainardi (1800–1879), D. Codazzi (1824–1875) and O. Bonnet (1819–1892) who were productive participants in the flurry of geometrical activity which followed publication of *Disquisitiones*. But it is clear in retrospect—and was clear already to Gauss[22]—that the most profoundly creative of that second generation of geometers was Riemann. Gauss had been led to geometry from physical geodesy, which may account for why he concentrated on the geometry of surfaces in 3-space (as did those who followed in his steps). Riemann, on the other hand, imagined himself to be exploring the "foundations of geometry," and in doing so to be pursuing a complex physico-philosophical agenda;[23] he found it natural[24] to look upon Gaussian surfaces as special cases of much more general ($N$-dimensional) structures. "Manifolds" he called them, the properties of which were internal to themselves—developed without reference to any enveloping Euclidean space.

---

[21] In the spherical case we would be led on this basis to write

$$K = -0 - \frac{1}{2R^2\sin u}\left\{\frac{\partial}{\partial v}0 - \frac{\partial}{\partial u}(-2\cos u)\right\} = \frac{1}{R^2}$$

which checks out.

[22] Gauss died in 1855, only one year after Riemann—at age 28—had (after delays caused by Gauss' declining health) presented the geometrical lecture in which Gauss found so much to praise. By 1860 Riemann himself was dead.

[23] See E. T. Bell captures the flavor of that agenda in Chapter 26 of his *Men of Mathematics* (1937). Or see Michael Spivak's translation of "Riemann's *Habilatationsvortrag*: On the hypotheses which lie at the foundations of geometry," which is reprinted in McCleary,[20] which is pretty heavy sledding, but ends with the remark that "This leads us away into the domain of another science, the realm of physics, into which the nature of the present occasion does not allow us to enter."

[24] There is no accounting for genius, but could he have been influeced by "Riemannian surfaces of $N$-sheets" which had played a role in his dissertation, "Grundlagen für eine allegemeine Theorie der Functionen einer veränderlichen complexen Grösse" (1851)?

With a vagueness worthy of Riemann himself, I will understand an $N$-dimensional manifold to be an "$(x^1, x^2, \ldots, x^N)$-coordinatized continuum" which is sufficiently structured to permit us to do the things we want to do. The manifold becomes a *Riemannian manifold* when endowed with functions

$$g_{ij}(x^1, x^2, \ldots, x^N) \quad : \quad i, j \in \{1, 2, \ldots, N\}, \ g_{ij} = g_{ji}$$

which permit one to write

$$(ds)^2 = g_{ij}(x) \, dx^i dx^j \tag{18}$$

to lend postulated *metric structure* to the manifold. In (18) we maintain the form of (12) but abandon the notion that (18) is the curvilinear expression of some more primitive metric axiom (Euclidean metric, either of the manifold itself or—as it was at (14), and always was for Gauss—of a Euclidean space within which the manifold is imagined to be embedded). Riemann retains the service of the "first fundamental form," but his program obligates him to abandon the "second fundamental form," and therefore to find some way to circumvent the mathematics which for Gauss culminated in the Theorema Egregium.

Somewhat idiosyncratically, I use the word "connection" to refer to *any* ancillary device we "smear on a manifold" in order to permit us to *do* things there, and say of a manifold $\mathcal{M}$ that has been endowed with a $g_{ij}(x)$ that it has been "metrically connected."[25]

The coordinate-independence of (18) requires that $g_{ij}$ transform as a tensor. At each of the points $P$ of $\mathcal{M}$ we erect, as our mathematical/physical interest may dictate, also populations of other tensors of various ranks and weights. Those live in multivector spaces which are "tangent" to $\mathcal{M}$ at $P$. If $X^i$ refers (in $x$-coordinates) to a contravariant vector defined at $P$, then $g_{ia}X^a$ and $g_{ai}X^a$ refer similarly to covariant vectors defined at $P$. And these are, by the symmetry of $g_{ij}$, the *same* covariant vector, which we may agree to denote $X_i$. Moreover, $g^{ia}X_a = g^{ia}g_{ab}X^b = \delta^i{}_b X^b = X^i$ gives back the contravariant vector from which we started. The idea extends naturally to tensors of arbitrary rank and weight

$$X^{\cdots}{}_i{}^{\cdots}{}_{\ldots} = g_{ia}X^{\cdots a \cdots}{}_{\ldots} \quad ; \quad X^{\cdots i \cdots}{}_{\ldots} = g^{ia}X^{\cdots}{}_a{}^{\cdots}{}_{\ldots}$$

On metrically connected manifolds we can agree to press the metric connection $g_{ij}$ into secondary service as the *universal index manipulator*. The fundamental Reimannian axiom (18) can by this convention be written in a way

$$(ds)^2 = dx_i \, dx^i$$

which renders $g_{ij}$ itself covert.

---

[25] My usage is arguably consistent with that employed by Schrödinger (See Chapter 9 in his elegant little book, *Spacetime Structure*(1954)), but departs from that favored by most differential geometers. As used by the latter, the concept originates in work (1916) of Gerhard Hassenberg, a German set theorist.

**Geodesics**.　　Let $x^i(t)$ describe a $t$-parameterized curve $\mathcal{C}$ which has been inscribed on $\mathcal{M}$. Assume

$$x^i(0) = \text{coordinates of a point } P$$
$$x^i(1) = \text{coordinates of a point } Q$$

i.e., that the curve $\mathcal{C}_{P \to Q}$ links $P$ to $Q$. Riemann's axiom (18) places us in position to write

$$\text{length of } \mathcal{C}_{P \to Q} = \int_0^1 \sqrt{g_{ab} v^a v^b}\, dt$$
$$v^a \equiv \tfrac{d}{dt} x^a(t)$$

"Geodesics" are curves of extremal length, and by straightforward appeal to the calculus of variations are found to satisfy

$$\left\{ \frac{d}{dt} \frac{\partial}{\partial v^i} - \frac{\partial}{\partial x^i} \right\} \sqrt{g_{ab} v^a v^b} = 0 \tag{19.1}$$

Working out the implications of (19.1) is a task made complicated by the presence of the radical. Those complications[26] can, however, be circumvented; one can show without difficulty[27] that if $\dot{s} = 0$ —i.e., if

$$t = (\text{constant}) \cdot (\text{arc length}) + (\text{constant})$$

—then the $\sqrt{\phantom{xxx}}$ can be discarded. So we adopt *arc-length parameterization* and obtain

$$\left\{ \frac{d}{ds} \frac{\partial}{\partial u^i} - \frac{\partial}{\partial x^i} \right\} g_{ab} u^a u^b = 0 \tag{19.2}$$
$$u^a \equiv \tfrac{d}{ds} x^a(s)$$

Quick calculation gives $g_{ia} \frac{d}{ds} u^a + \tfrac{1}{2} \left\{ \partial_a g_{ib} + \partial_b g_{ia} - \partial_i g_{ab} \right\} u^a u^b = 0$ which can be written

$$\tfrac{d}{ds} u^i + \left\{ {}^{\,i}_{ab} \right\} u^a u^b = 0 \tag{20.2}$$
$$\left\{ {}^{\,i}_{ab} \right\} \equiv g^{ij} \cdot \tfrac{1}{2} \left\{ \partial_a g_{jb} + \partial_b g_{ja} - \partial_j g_{ab} \right\} \tag{21}$$

To recover the result to which we would have been led had we proceeded from (19.1)—i.e., had we elected to use *arbitrary* parameterization—we use $\frac{d}{ds} = (\dot{s})^{-1} \frac{d}{dt}$ (whence $u^i = (\dot{s})^{-1} v^i$) and obtain

$$\tfrac{d}{dt} v^i + \left\{ {}^{\,i}_{ab} \right\} v^a v^b = v^i \tfrac{d}{dt} \log \tfrac{ds}{dt} \tag{20.1}$$
$$\tfrac{ds}{dt} = \sqrt{g_{ab} v^a v^b}$$

---

[26]　See CLASSICAL DYNAMICS (1964), Chapter 2, p. 105.

[27]　The simple argument can be found in "Geometrical mechanics: Remarks commemorative of Heinrich Hertz" (1994) at p. 14.

Comparison of (19.2) with 19.1 , (20.2) with (20.1) underscores the marked *simplification which typically results from arc-length parameterization.*

**Covariant differentiation on affinely connected manifolds**. If $X(x)$ responds to $x \to y$ as a scalar density of zero weight

$$X(x) \to Y(y) = X(x(y))$$

then

$$\frac{\partial Y}{\partial y^i} = \frac{\partial x^a}{\partial y^i}\frac{\partial X}{\partial x^a} \quad \text{which we abbreviate} \quad Y_{,i} = W^a{}_i X_{,a}$$

In short: the gradient of a scalar transforms tensorially (i.e., as a covariant vector, and is standardly held up as the exemplar of such an object). But if $X_i$ transforms as a covariant vector $(Y_i = W^a{}_i X_a)$ then

$$Y_{i,j} = W^a{}_i W^b{}_j X_{a,b} + X_a \cdot \frac{\partial^2 x^a}{\partial y^i \partial y^j}$$

shows that $X_{i,j}$ *fails to transform tensorially* (unless we impose upon $x \to y$ the restrictive requirement that $\partial^2 x^a / \partial y^i \partial y^j = 0$). A similar remark pertains generally: if $Y^{i_1 \cdots i_r}{}_{j_1 \cdots j_s} = W^w \cdot M^{i_1}{}_{a_1} \cdots M^{i_r}{}_{a_r} W^{b_1}{}_{j_1} \cdots W^{b_s}{}_{j_s} X^{a_1 \cdots a_r}{}_{b_1 \cdots b_s}$ then

$$Y^{i_1 \cdots i_r}{}_{j_1 \cdots j_s, k} = \left[ W^w \cdot M^{i_1}{}_{a_1} \cdots M^{i_r}{}_{a_r} W^{b_1}{}_{j_1} \cdots W^{b_s}{}_{j_s} \right] W^c{}_k X^{a_1 \cdots a_r}{}_{b_1 \cdots b_s, c}$$
$$+ X^{a_1 \cdots a_r}{}_{b_1 \cdots b_s} \frac{\partial}{\partial y^k} \left[ \text{etc.} \right]$$

This circumstance severely constrains our ability to do ordinary differential calculus on manifolds; it limits us on transformation-theoretic grounds to such "accidentally tensorial" constructs as $X_{i,j} - X_{j,i}$ in which the unwanted terms (by contrivance) cancel.[28]

There exists, however, an elegant work-around, devised early in the present century by Ricci and Levi-Civita (who harvested the fruit of a seed planted by Christoffel in 1869) and which has much in common with the spirit of gauge field theory. In place of $\partial_j$ we study a *modified* operator $\mathcal{D}_j$, the action of which can, in the simplest instance, be described

$$\mathcal{D}_j X_i \equiv X_{i,j} - X_a \Gamma^a{}_{ij} \quad ; \quad \text{denoted } X_{i;j} \tag{22}$$
$$\underset{\text{semi-colon instead of comma}}{\uparrow\!\!\!\_\!\!\!\_}$$

The idea is to impose upon the "affine connection" $\Gamma^a{}_{ij}$ such transformation properties as are sufficient to *insure* that $X_{i;j}$ transform tensorially

$$\bar{X}_{i;j} = W^a{}_i W^b{}_j X_{a;b}$$

---

[28] For lists of such "accidentally tensorial constructs"—of which the exterior calculus provides a systematic account (and which are in themselves rich enough to support much of physics)—see §4 of "Electrodynamical Applications of the Exterior Calculus" (1996) and pp. 22-24 of Schrödinger.[25]

(It has at this point become more natural to write $x \to \bar{x}$ where formerly we wrote $x \to y$.) This by

$$\bar{X}_{i,j} - \bar{X}_a \bar{\Gamma}^a{}_{ij} = \left\{ W^a{}_i W^b{}_j X_{a,b} + X_a \cdot \frac{\partial^2 x^a}{\partial y^i \partial y^j} \right\} - W^b{}_a X_b \bar{\Gamma}^a{}_{ij}$$

$$= W^a{}_i W^b{}_j X_{a,b} - X_c \left\{ W^c{}_a \bar{\Gamma}^a{}_{ij} - \frac{\partial^2 x^c}{\partial y^i \partial y^j} \right\}$$

$$= W^a{}_i W^b{}_j \left( X_{a,b} - X_c \Gamma^c{}_{ab} \right)$$

requires that $W^c{}_a \bar{\Gamma}^a{}_{ij} - \partial^2 x^c / \partial y^i \partial y^j = W^a{}_i W^b{}_j \Gamma^c{}_{ab}$. Multiplication by $M^k{}_c$ leads to the conclusion that $\Gamma^k{}_{ij}$ must transform

$$\Gamma^k{}_{ij} \to \bar{\Gamma}^k{}_{ij} = M^k{}_c W^a{}_i W^b{}_j \Gamma^c{}_{ab} + \frac{\partial y^k}{\partial x^c} \frac{\partial^2 x^c}{\partial y^i \partial y^j} \tag{23}$$

if it is to do the job we require of it. To assign natural meaning to $X^i{}_{;j}$ we shall REQUIRE that *for weightless scalars "covariant differentiation" reduces to ordinary differentiation*

$$X_{;i} = X_{,i} \tag{24.1}$$

and that for *for tensor products covariant differentiation satisfies the product rule*

$$(X^{\cdots}...Y^{\cdots}...)_{;i} = X^{\cdots}...{}_{;i} Y^{\cdots}... + X^{\cdots}...Y^{\cdots}...{}_{;i} \tag{24.2}$$

Look in this light to the contracted product $Y_a X^a$; we have

$$(Y_{a,j} - Y_i \Gamma^i{}_{aj}) X^a + Y_i X^i{}_{;j} = (Y_i X^i)_{,j} = Y_{a,j} X^a + Y_i X^i{}_{,j}$$

for all $Y_i$, from which we obtain the second of the following equations (the first being simply a repeat of (22)):

$$\left. \begin{array}{l} X_{i;j} = X_{i,j} - X_a \Gamma^a{}_{ij} \\ X^i{}_{;j} = X^i{}_{,j} + X^a \Gamma^i{}_{aj} \end{array} \right\} \tag{25.1}$$

A somewhat more intricate argument—which I omit[29]—leads to the conclusion that for scaral *densities* the natural thing to write is

$$X_{;i} = X_{,i} - w X \Gamma^a{}_{ai} \tag{25.2}$$

which gives back (24.1) in the case $w = 0$. Using (24) and (25) in combination one can describe the covariant derivatives of tensors of all ranks and weights; for example, we find

$$X^{ij}{}_{k;l} = X^{ij}{}_{k,l} + X^{aj}{}_k \Gamma^i{}_{al} + X^{ia}{}_k \Gamma^j{}_{al} - X^{ij}{}_a \Gamma^a{}_{kl} - w X^{ij}{}_k \Gamma^a{}_{al}$$

---

[29] See Chapter 2, p. 59 in the notes recently cited,[26] or Schrödinger,[25] p. 32.

which—by construction—transforms as a mixed tensor density of
- unchanged contravariant rank,
- augmented covariant rank, and
- unchanged weight.

All of which becomes available as a mechanism for obtaining tensors by the (covariant) differentiation of tensors, and for constructing transformationally well-behaved tensorial differential equations...if and only if the underlying manifold $\mathcal{M}$ has been endowed with an affine connection; i.e., if and only if (in some coordinate system) functions $\Gamma^k{}_{ij}(x)$ have been prescribed at every point, in which case we say that $\mathcal{M}$ is "affinely connected." Differentiation of objects attached to such a manifold is *relative to the prescribed affine connection*: replace one affine connection with another, and the meaning of all derivatives changes. Description of $\Gamma^k{}_{ij}(x)$ requires the specification of $N^3$ functions; its description in other coordinate systems (reached by $x \to \bar{x}$) is then accomplished by appeal to (23), in connection with which we observe that

- $\Gamma^k{}_{ij}(x) \to \bar{\Gamma}^k{}_{ij}(\bar{x})$ is tensorial only if $x \to \bar{x}$ is linear (as, we note in passing, are the inhomogeneous Lorentz transformations); in more general cases $\Gamma^k{}_{ij}(x)$ transforms distinctively—"like an affine connection."

- The rule (23) is "transitive" in the sense that if $\Gamma^k{}_{ij}(x) \to \bar{\Gamma}^k{}_{ij}(\bar{x})$ and $\bar{\Gamma}^k{}_{ij}(\bar{x}) \to \bar{\bar{\Gamma}}^k{}_{ij}(\bar{\bar{x}})$ conform to it, then so does $\Gamma^k{}_{ij}(x) \to \bar{\bar{\Gamma}}^k{}_{ij}(\bar{\bar{x}})$

- If $\Gamma^k{}_{ij}(x)$ vanishes in some coordinate system (call it the $x$-system) then it is given in other coordinate systems by

$$\bar{\Gamma}^k{}_{ij}(\bar{x}) = \frac{\partial \bar{x}^k}{\partial x^c} \frac{\partial^2 x^c}{\partial \bar{x}^i \partial \bar{x}^j}$$

which is $ij$-symmetric.

- Let $\Gamma^k{}_{ij}$ be resolved $\Gamma^k{}_{ij} = \frac{1}{2}(\Gamma^k{}_{ij} + \Gamma^k{}_{ji}) + \frac{1}{2}(\Gamma^k{}_{ij} - \Gamma^k{}_{ji})$ into its $ij$-symmetric and $ij$-antisymmetric (or "torsional") parts. It follows from (23) that, while the symmetric part transforms "like an affine connection," the torsional part transforms *tensorially*.[30]

---

[30] In Riemannian geometry, and in gravitational theories based upon it, one (tacitly) assumes the affine connection to be "torsion free" (i.e., that $\Gamma^k{}_{ij}$ is symmetric), and *MTW* remark (p. 250) that to assume otherwise would be inconsistent with the Principle of Equivalence. But, beginning in the 1920's, Einstein and others studied various generalizations of general relativity which involved relaxation of that assumption. For a useful summary of that work see §§17d&e in Pais.[3] According to MTW (§39.2), all such generalizations can now be dismissed on observational grounds except (possibly) for the torsional theory described by E. Cartan in 1922/23.

On metrically connected manifolds there exists a "natural" (symmetric) affine connection. It arises when one imposes the requirement that
- covariant differentiation and
- index manipulation

be *compatible* operations, performable in either order:

$$(g_{ia}X^{\cdots a \cdots}{}_{\cdots})_{;k} = g_{ia}(X^{\cdots a \cdots}{}_{\cdots ;k})$$

This, by (24.2), amounts to requiring that $g_{ij;k} = 0$; i.e., that

$$g_{ij,k} - g_{aj}\Gamma^a{}_{ik} - g_{ia}\Gamma^a{}_{jk} = 0$$

Let this equation be notated $\Gamma_{jik} + \Gamma_{ijk} = g_{ij,k}$ and draw upon the symmetry assumption to write $\Gamma_{jki}$ in place of $\Gamma_{jik}$. Cyclic permutation on $ijk$ leads then to a trio of equations which can be displayed

$$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} \Gamma_{ijk} \\ \Gamma_{jki} \\ \Gamma_{kij} \end{pmatrix} = \begin{pmatrix} g_{jk,i} \\ g_{ki,j} \\ g_{ij,k} \end{pmatrix}$$

from which it follows by matrix inversion that

$$\begin{pmatrix} \Gamma_{ijk} \\ \Gamma_{jki} \\ \Gamma_{kij} \end{pmatrix} = \tfrac{1}{2} \begin{pmatrix} -1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & -1 \end{pmatrix} \begin{pmatrix} g_{jk,i} \\ g_{ki,j} \\ g_{ij,k} \end{pmatrix}$$

Making free use of metric symmetry ($g_{ij} = g_{ji}$) we are led thus—by an argument which is seen to hinge on the "permutation trick" employed already once before (see again the derivation of (2–83))—to cyclic permutations of the following basic statement:

$$\Gamma_{kij} = g_{ka}\Gamma^a{}_{ij} = \tfrac{1}{2}\big\{g_{ki,j} + g_{kj,i} - g_{ij,k}\big\} \tag{26.1}$$

$$\Gamma^k{}_{ij} = g^{ka} \cdot \tfrac{1}{2}\big\{g_{ai,j} + g_{aj,i} - g_{ij,a}\big\} \tag{26.2}$$

To compute the covariant derivatives of densities it helps to know also that

$$\Gamma^a{}_{ai} = \frac{\partial}{\partial x^i} \log \sqrt{g} \tag{26.3}$$

but the demonstration is somewhat intricate and will be postponed.[31]

The expressions which stand on the right side of (25) were introduced into the literature by Christoffel (1869); they are called "Christoffel symbols," and are (or used to be) standardly notated

$$[ij, k] \equiv \text{right side of (25.1)} \quad : \quad \text{Christoffel symbol of 1}^{\text{st}} \text{ kind}$$

$$\left\{{}^{k}_{ij}\right\} \equiv \text{right side of (25.2)} \quad : \quad \text{Christoffel symbol of 2}^{\text{nd}} \text{ kind}$$

---

[31] See p. 73 in some notes previously cited,[28] or (58) below.

The Christoffel symbol $\left\{ {k \atop ij} \right\}$ is the symmetric affine connection—the "Christoffel connection"—most natural to metrically connected manifolds (Riemannian manifolds), and its valuation is, according to (26), latent in specification of the metric $g_{ij}(x)$. It sprang spontaneously to our attention already at (20), when we were looking to the description of Riemannian geodesics. We confront therefore the question: What has covariant differentiation to do with geodesic design?

**Parallel transport**. When asked to "differentiate a tensor"—let it, to render the discussion concrete, be a covariant vector—one's first instinct might be to write

$$\frac{X_i(x + \delta x) - X_i(x)}{\delta x}$$

Such a program is, however, foredoomed. For, while it makes transformational good sense to write $Y_i(x) - X_i(x)$, it is not permissible to add/subtract vectors which are *attached to distinct points on the manifold* (vectors which therefore live in distinct vector spaces, and transform a bit differently from one another).
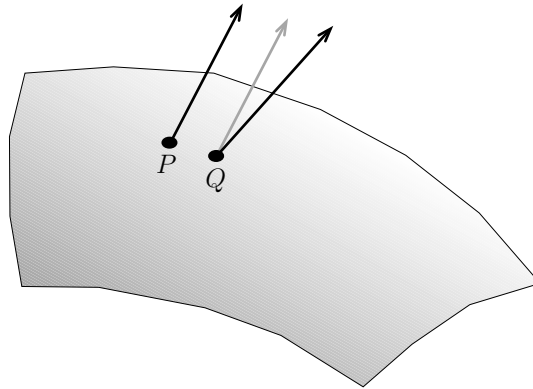


FIGURE 4: *The vector $X_i(x)$ is attached to the manifold at $P$, and $X_i(x + \delta x)$ at the neighboring point $Q$. The affine connection $\Gamma^k{}_{ij}(x)$ defines the sense in which the grey vector $\mathfrak{X}_i(x + \delta x)$ is "parallel" to $X_i(x)$. The difference $X_i(x + \delta x) - \mathfrak{X}_i(x + \delta x)$ is tensorially meaningful, and gives rise to the covariant derivative in the limit.*

We need to identify at $x + \delta x$ a "stand-in" for $X_i(x)$—a vector $\mathfrak{X}_i(x + \delta x)$ obtained by (in Weyl's phrase) the "parallel transport" of $X_i(x)$ from $x$ to $x + \delta x$. With the aid of such an object we would construct

$$\frac{X_i(x + \delta x) - \mathfrak{X}_i(x + \delta x)}{\delta x}$$

which provides escape from our former transformational problem. Adopt this *definition* of the infinitesimial parallel transport process:

$$\mathfrak{X}_i(x + \delta x) \equiv X_i(x) + X_a(x)\Gamma^a{}_{ij}(x)\delta x^j \qquad (27)$$

We then have

$$\begin{aligned}
X_i(x + \delta x) - \mathfrak{X}_i(x + \delta x) &= \{X_i(x + \delta x) - X_i(x)\} - X_a(x)\Gamma^a{}_{ij}(x)\delta x^j \\
&= \{X_{i,j}(x) - X_a(x)\Gamma^a{}_{ij}(x)\}\delta x^j \\
&\equiv X_{i;j}(x)\delta x^j
\end{aligned}$$

which gives back (24.1) and permits reconstruction of all that has gone before.

The parallel transport concept permits formulation of a valuable *metric-independent theory of geodesics*, which I now sketch. Let $\mathcal{C}$ refer as before to an arbitrarily parameterized[32] curve $x^i(t)$ which has been inscribed on $\mathcal{M}$. Indexed objects $v^i(t) \equiv \frac{d}{dt}x^i(t)$ are associated with the points of $\mathcal{C}$, and are readily seen to transform as contravariant vectors (essentially because the differentials $dx^i$ do). It is natural to

call $v^i(t)$ the vector "tangent" to $\mathcal{C}$ at $x^i(t)$

Introduce

$$\mathcal{V}^i(t + \delta t) \equiv v^i(t) - v^a(t)\Gamma^i{}_{ab}(x(t))v^b(t)\delta t \qquad (28)$$

to describe the result of parallel transporting $v^i(t)$ from $x^i(t)$ to the neighboring point $x^i(t) + v^i(t)\delta t$ and ask: How does $\mathcal{V}^i(t + \delta t)$ compare to $v^i(t + \delta t)$? If $\mathcal{C}$ is geodesic we expect those two to be parallel; i.e., we *expect geodesics to be generated by parallel transportation of a tangent*. One is tempted to look to the implications of $\mathcal{V}^i(t + \delta t) = v^i(t + \delta t)$ —i.e., of

$$\frac{d}{dt}v^i + \Gamma^i{}_{ab}\,v^a v^b = 0 \qquad (29)$$

—but that equation is *not stable with respect to parametric regraduation*: an adjustment $t \to \tau = \tau(t)$ would cause (29) to become (see again (20), and agree in the present instance to write $w^i \equiv dx^i/d\tau = (d\tau/dt)^{-1}\,v^i = v^i/\dot{\tau}$)

$$\frac{d}{d\tau}w^i + \Gamma^i{}_{ab}\,w^a w^b = w^i \frac{d}{d\tau}\log\frac{dt}{d\tau} = -w^i\,\ddot{\tau}/\dot{\tau}^2$$

which is structurally distinct from (29). Evidently the most we can require is that $\mathcal{V}^i(t + \delta t) \sim v^i(t + \delta t)$; i.e., that there exists a $\varphi(t)$ such that

$$v^i(t) - v^a(t)\Gamma^i{}_{ab}(x(t))v^b(t)\delta t = [1 - \varphi(t)\delta t] \cdot v^i(t + \delta t)$$

Then in place of (29) we have

$$\frac{d}{dt}v^i + \Gamma^i{}_{ab}\,v^a v^b = v^i \cdot \varphi \qquad (30.1)$$

which upon reparameterization $t \to \sigma = \sigma(t)$ becomes (write $u^i \equiv \frac{d}{d\sigma}x^i = v^i/\dot{\sigma}$)

$$\frac{d}{d\sigma}u^i + \Gamma^i{}_{ab}\,u^a u^b = u^i \cdot (\dot{\sigma}\phi - \ddot{\sigma})/\dot{\sigma}^2 \qquad (30.2)$$

---

[32] "Arc-length parameterization" is, in the absence of a metric, not an option.

where $\phi(\sigma(t)) \equiv \varphi(t)$. Notice that if we are given any instance of (30.1) then we have only to take $\sigma(t)$ to be any solution of $\dot{\sigma}\varphi - \ddot{\sigma} = 0$[33] to bring (30.2) to the form

$$\tfrac{d}{d\sigma}u^i + \Gamma^i{}_{ab}\, u^a u^b = 0 \qquad (30.3)$$

so even in the absence of $g_{ij}(x)$ a "natural parameterization" is available, and when it is employed the geodesic condition (30.1) reads (30.3), which *does* possess the design which at (29) was originally conjectured.[34]

On *metrically* connected manifolds the specialization $\Gamma^i{}_{ab} \mapsto \left\{{}^{i}_{ab}\right\}$ becomes natural. It becomes natural, moreover, to adopt $s$-parameterization, and from $(ds)^2 = g_{ab}\, dx^a dx^b$ we conclude that

> the tangent vector $u^i(s) \equiv \tfrac{d}{ds}x^i(s)$ is a *unit* vector: $g_{ab}u^a u^b = 1$

Equation (30.3) gives back (20.2), but with this added information: *The unit tangents to a Riemannian geodesic are parallel transports of one another*. The question which motivated this discussion—What has covariant differentiation to do with geodesic design?—has now an answer. But the concept of parallel tranport (introduced by Levi-Civita in 1917—too late to do Einstein any immediate good, but immediately put to general relativistic work by Weyl) has yet other things to teach us.

---

[33] The solution resulting from initial conditions $\sigma(0) = 0$, $\dot{\sigma}(0) = 1$ can be described

$$\sigma(t) = \int_0^t \exp\left\{\int_0^{t'} \varphi(t'')\, dt''\right\} dt' = t + \tfrac{1}{2}\varphi(0)t^2 + \cdots$$

[34] Given a tensor field $X^i(x)$, and a curve $\mathcal{C}$ inscribed by $x^i(t)$ on an affinely connected manifold $\mathcal{M}$, the "absolute derivative"

$$\tfrac{\delta}{\delta t}X^i \equiv \tfrac{d}{dt}X^i + \Gamma^i{}_{ab}\, X^a v^b$$

(here $v^i \equiv dx^i/dt$ and $\tfrac{d}{dt}X^i = X^i{}_{,a}v^a$) describes the rate at which $X^i$ is seen to change as one progresses $t \to t + \delta t$ along the curve. It is easy to show that $\delta X^i/\delta t$ responds
  • tensorially to recoordinatization $x \to y$, and
  • by the chain rule $\tfrac{\delta}{\delta\tau} = \tfrac{dt}{d\tau}\tfrac{\delta}{\delta t}$ to reparameterization $t \to \tau$.
so statements of the form $\tfrac{\delta}{\delta t}X^i = 0$ are transformationally stable. Why, therefore, was (29), *not* stable? Because $X^i \mapsto v^i$ introduces an object which carries its own parameter-dependence. See "Non-Riemannian Spaces," the final chapter in J. L. Synge & A. Schild's *Tensor Calculus* (1952), for a detailed account of the absolute derivative and its applications. It was, by the way, from this text that I myself learned tensor calculus while an undergraduate, and in my view it remains the text best calculated to serve the needs of young physicists.

**Curvature**. Inscribe a closed curve $\mathcal{C}$ on the Euclidean plane. Parallel transport of a vector $\boldsymbol{X}$ around $\mathcal{C}$ (taken in Figure 5 to be a triangle) returns $\boldsymbol{X}$ to its
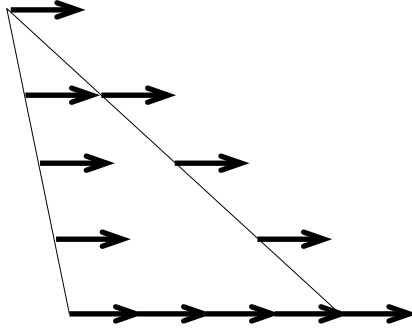


FIGURE 5: *Parallel transport of a vector along a closed path (here taken to be triangular) inscribed on the Euclidean plane. The vector returns home unchanged by its adventure.*

initial value. That this is, in general, *not* the case if $\mathcal{C}$ has been inscribed on a *curved* surface $\Sigma^2 \in E^3$ is illustrated in Figure 6. In that figure the apex of the spherical triangle sits at the pole ($\theta = 0$) and the equatorial base points differ in longitude by $\varphi = \phi_2 - \phi_1$; $\boldsymbol{X}$, upon return to its point of departure (where $\theta = \frac{\pi}{2}, \phi = \phi_1$), has been rotated $\circlearrowleft$ through angle $\varphi$. Spherical trigonometry supplies the information that

$$\text{(sum of interior angles)} - \pi \equiv \text{``spherical excess''} = \text{area}/R^2 \qquad (31)$$

which acquires interest from the observation that

$$\text{``spherical excess''} = \varphi = \text{angle through which } \boldsymbol{X} \text{ is rotated}$$

This result—though special to the case illustrated—has a "generalizable look about it" ... and indeed: the "Gauss-Bonnet theorem" (which I will discuss in a moment: see Figure 6) asserts that

$$2\pi - \text{(sum of exterior angles)} - \int_{\mathcal{C}} \kappa_g \, ds = \iint_{\mathcal{D}} K \, dS \qquad (32)$$

which in the case illustrated—where
  - $\kappa_g = 0$ because $\mathcal{D}$ is bounded by geodesics;
  - $K = 1/R^2$ everywhere because the surface is spherical
—gives back (31).

The meanings of the terms which enter into statement of the Gauss-Bonnet theorem are evident with one exception; I refer to the "geodesic curvature," the meaning of which is explained in the caption to Figure 8.
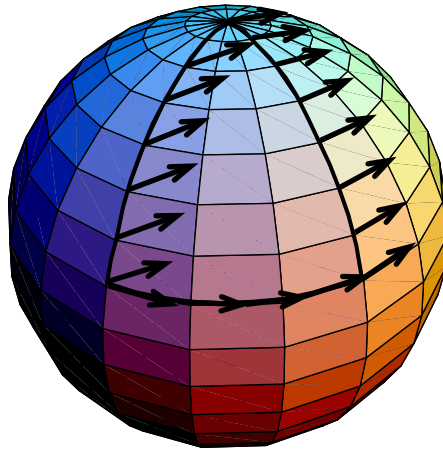
FIGURE 6: *Parallel transport of a vector along a triangular path inscribed on a sphere of radius R. The base points sit on the equator, the apex is at the pole, the sides are geodesic, and the circulation sense ↺ has been counterclockwise.*
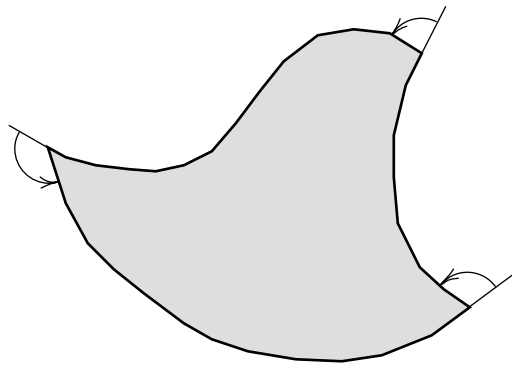


FIGURE 7: *The shaded region $\mathcal{D}$ is bounded by a closed contour $\mathcal{C}$ which has been inscribed on a Gaussian surface $\Sigma^2 \in E^3$. The Gauss-Bonnet theorem (32) refers to situations in which the number of vertices is arbitrary, the bounding arcs need not be geodesic, and the Gaussian curvature $K$ may vary from point to point within $\mathcal{D}$.*
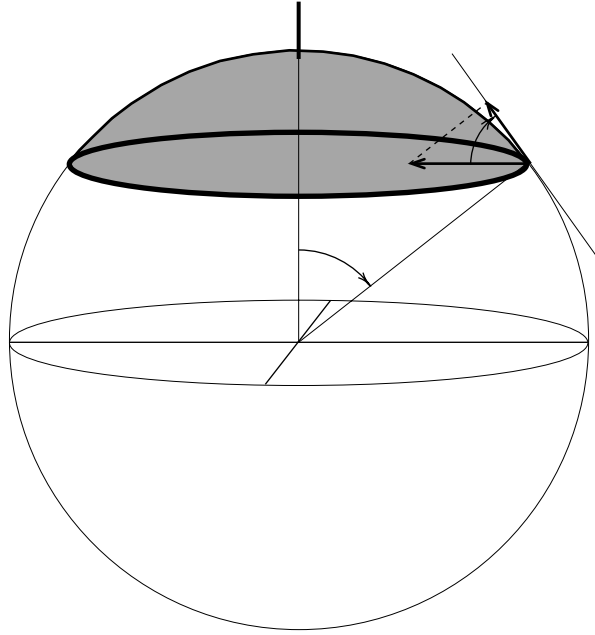
FIGURE 8: *The shaded cap $\mathcal{D}$, bounded by a circle of constant $\theta$, has area $2\pi R^2(1 - \cos\theta)$, and its boundary $\mathcal{C}$ presents no vertices. The Gauss-Bonnet theorem (32) therefore asserts that*

$$2\pi - \int_{\mathcal{C}} \kappa_g \, ds = 2\pi(1 - \cos\theta) \qquad\qquad (i)$$

*The curvature $\kappa$ of the boundary $\mathcal{C}$—thought of as a space curve—arises from*

$$\tfrac{d}{ds}\boldsymbol{x}(s) = \boldsymbol{T}(s) \quad \text{and} \quad \tfrac{d}{ds}\boldsymbol{T}(s) = \kappa(s)\,\boldsymbol{n}(s)$$

*(here $\boldsymbol{T}$ is the unit tangent to $\mathcal{C}$ at $\boldsymbol{x}(s)$, and $\boldsymbol{n}$ is the unit normal in the plane of $\mathcal{C}$), and in the present instance has constant value*

$$\kappa = 1/(\textit{radius of curvature}) = 1/(R\sin\theta)$$

*The "geodesic curvature" $\kappa_g$ refers to the length of the projection of $\kappa\,\boldsymbol{n}$ onto the tangent plane, and is in the present instance given by*

$$\kappa_g = \cos\theta/(R\sin\theta)$$

*Therefore*

$$\int_{\mathcal{C}} \kappa_g \, ds = \kappa_g \cdot 2\pi R\sin\theta = 2\pi\cos\theta$$

*—in precise agreement with $(i)$.*

Generally, $\mathcal{C}$ is locally geodesic if $\boldsymbol{n}$ stands normal to the local tangent plane; it follows that

$$\kappa_g = 0 \;\; \text{everywhere on a geodesic}$$

(whence the name), and that the geodesics inscribed on a Gaussian surface $\Sigma^2$ are "as uncurved as possible." The Gauss-Bonnet theorem (32), insofar as it entails

- exploration of the boundary $\mathcal{C} = \partial\mathcal{D}$ on the left
- exploration of the interior of $\mathcal{D}$ on the right,

bears a family resemblance to Stokes' theorem.[35]   One expects it to be the case—as, indeed, it is—that $\kappa_g$ can be described by operations intrinsic to the Gaussian surface.[36]

The circumstance discussed above—and illustrated in Figure 6 as it pertains to one particular Gaussian surface—was recognized by Riemann to pertain in metrically connected spaces of any dimension, and in fact it presumes only the affine connection needed to lend meaning to "parallel transport." The
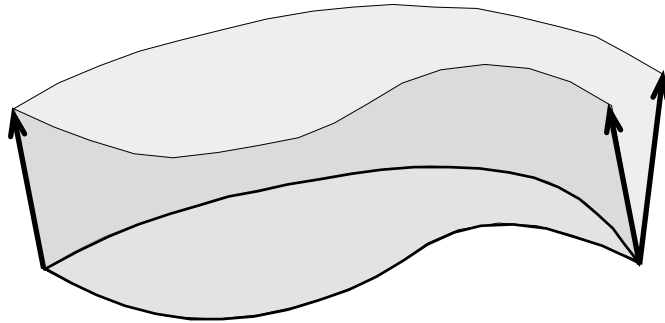


FIGURE 9:  *The result of parallel transport from $P$ to $Q$ in on affinely connected manifold is typically path-dependent, and that fact can be taken to be the defining symptom of "curvature."*

basic phenomenon is illustrated in the preceding figure, and can be approached analytically in several ways.   One might, for example, develop a formal description of the $X_Q^i$ which results from parallel transport of $X_P^i$ along $\mathcal{C}$, and then examine the $\delta X_Q^i$ which results from variation of the path.[37]   It is far easier and more efficient, however, to look to the $\delta X^i$ which results from comparison of a pair of *differential paths* (Figure 10); one's interest is then

---

[35] The point is developed on pp. 45–48 of "Ellipsometry" (1999).   The motivation there comes not from general relativity but from optics.

[36] For indication of how this can be accomplished, see the *Encyclopedic Dictionary of Mathematics* (1993), p. 1731.

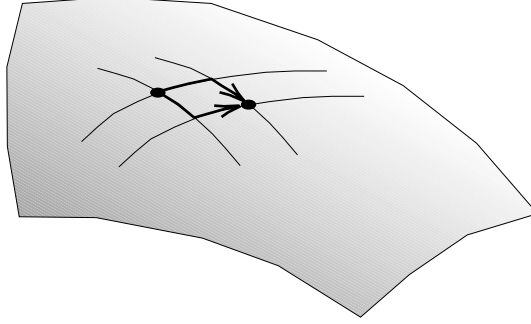[37] This is the program pursued on pp. 134–137 in notes previously cited.[28]

FIGURE 10:  *Alternative paths* $\{\delta u \text{ then } \delta v\}$ *and* $\{\delta v \text{ then } \delta u\}$ *linking point* $P$ *to a neighboring point in* $\mathcal{M}$.

directed to the expression

$$\delta X^i = \left\{ \frac{\delta}{\delta v}\frac{\delta}{\delta u} - \frac{\delta}{\delta u}\frac{\delta}{\delta v} \right\} X^i \delta u \delta v$$

$$= \left\{ X^i{}_{;jk} - X^i{}_{;kj} \right\} \frac{\partial x^j}{\partial u} \frac{\partial x^k}{\partial v} \delta u \delta v$$

and thus by simple calculations to the statements

$$\left\{ X^i{}_{;jk} - X^i{}_{;kj} \right\} = (X^i{}_{,j} + X^a \Gamma^i{}_{aj})_{,k} + (X^b{}_{,j} + X^a \Gamma^b{}_{aj}) \Gamma^i{}_{bk}$$
$$- \text{ ditto with } j \text{ and } k \text{ interchanged}$$
$$= X^a{}_{,k}\Gamma^i{}_{aj} + X^b{}_{,j}\Gamma^i{}_{bk} + X^a(\Gamma^i{}_{aj,k} + \Gamma^b{}_{aj}\Gamma^i{}_{bk})$$
$$- \text{ ditto with } j \text{ and } k \text{ interchanged}$$
$$= -X^a R^i{}_{ajk} \tag{33.1}$$
$$R^i{}_{ajk} \equiv \Gamma^i{}_{ak,j} - \Gamma^i{}_{aj,k} + \Gamma^b{}_{ak}\Gamma^i{}_{bj} - \Gamma^b{}_{aj}\Gamma^i{}_{bk} \tag{34}$$
$$\left\{ X_{i;jk} - X_{i;kj} \right\} = +X_a R^a{}_{ijk} \tag{33.2}$$

In the matrix notation $I\!\Gamma_i \equiv \|\Gamma^a{}_{bi}\|$ the definition (33) becomes

$$\mathbb{R}_{jk} = \partial_j I\!\Gamma_k - \partial_k I\!\Gamma_j + I\!\Gamma_j I\!\Gamma_k - I\!\Gamma_k I\!\Gamma_j \tag{35}$$

It is an implication of the design of (33) that $R^i{}_{ajk}$ transforms tensorially (as a mixed fourth-rank tensor of zero weight), even though it has been assembled from objects which do *not* transform tensorially. On *metrically* connected manifolds it becomes natural in place of (34) to write

$$R^i{}_{ajk} = \frac{\partial}{\partial x^j}\left\{ {i \atop ak} \right\} - \frac{\partial}{\partial x^k}\left\{ {i \atop aj} \right\} + \left\{ {i \atop bj} \right\}\left\{ {b \atop ak} \right\} - \left\{ {i \atop bk} \right\}\left\{ {b \atop aj} \right\} \tag{36}$$

This is the *Riemann curvature tensor*, which Riemann allegedly obtained

without benefit either of a theory of connections or of a tensor calculus.[38] The right side of (34) describes its (metric-independent) affine generalization.

From (34) it follows readily that

$$R^i{}_{ajk} = -R^i{}_{akj} \quad : \quad \text{antisymmetry in last subscripts} \quad (37.1)$$

$$R^i{}_{ajk} + R^i{}_{jka} + R^i{}_{kaj} = 0 \quad : \quad \text{cyclic symmetry} \quad (37.2)$$

On metrically connected manifolds we can use $g_{ij}$ to construct

$$R_{iajk} \equiv g_{ib} R^b{}_{ajk} \quad (38)$$

which is found to possess, in addition to those symmetry properties, also two others:

$$R_{iajk} = -R_{aijk} \quad : \quad \text{antisymmetry in first subscripts} \quad (37.3)$$

$$R_{iajk} = +R_{jkia} \quad : \quad \text{symmetry in first/last pair-interchange} \quad (37.4)$$

Careful counting,[39] based upon those overlapping statements, shows that in the $N$-dimensional case $R_{iajk}$ possesses a total of $\# \equiv \frac{1}{12} N^2(N^2 - 1)$ independent components; as $N$ ranges on $\{2, 3, 4, 5, \dots\}$ $\#$ ranges on $\{1, 6, 20, 50, \dots\}$.

Direct computation establishes also that the Riemann tensor satisfies also a population of first derivative conditions which can be written

$$\mathbb{R}_{ij;k} + \mathbb{R}_{jk;i} + \mathbb{R}_{ki;j} = \mathbb{O} \quad (39)$$

and are called *Bianchi identities*.

It was remarked in connection with (22) that the replacement

$$\partial_j X_i \rightarrow \mathcal{D}_j X_i \equiv \partial_j X_i - X_a \Gamma^a{}_{ij}$$

"has much in common with the spirit of gauge field theory." I had then in mind the minimal coupling adjustment $\partial_\mu \psi \rightarrow \mathcal{D}_\mu \psi \equiv \partial_\mu \psi - ig\psi A_\mu$ which is basic to the latter theory. The parallel became more striking when it developed (compare (23) with (3–8)) that both the connection $\Gamma^a{}_{ij}$ and the gauge field $A_\mu$ had to *transform by acquisition of an additive term* if they were to perform their respective "compensator" roles. I draw attention now to the striking fact that (35) had a precise precursor in the equation

$$\boldsymbol{F}_{\mu\nu} \equiv (\partial_\mu \boldsymbol{A}_\nu - \partial_\nu \boldsymbol{A}_\mu) - ig(\boldsymbol{A}_\mu \boldsymbol{A}_\nu - \boldsymbol{A}_\nu \boldsymbol{A}_\mu)$$

---

[38] I say "allegedly" because I can find no hint of any such thing in Riemann's collected works! Eddington—I suspect with good historical cause—refers in his *Mathematical Theory of Relativity* (1923) always to the "Riemann-Christoffel tensor."

[39] See p. 86 in Synge & Schild.[34]

which at (3–101) served to define the non-Abelian analog of the electromagnetic field tensor. And the Bianchi identities (39) were anticipated at (1–126). Authors—writing in the presumption that their readers possess some familiarity with general relativity—most typically point to those formal parallels in an effort to make gauge field theory seem less alien.[40] My own intent has been the reverse. But pretty clearly, there must exist some sufficiently elevated viewpoint from which, in sufficiently fuzzy focus, gauge field theory and Riemannian geometry (if not general relativity itself) appear to be "the same thing."

Continuing in the presumption that $\mathcal{M}$ is metrically connected (i.e., that we are doing Riemannian geometry, and excluding from consideration the more relaxed structures contemplated by Weyl and others[34]), we observe that

$$g^{ab}R_{abij} = g^{ab}R_{ijab} = 0$$

and

$$g^{ab}R_{aijb} = -g^{ab}R_{aibj} = -g^{ab}R_{iajb} = g^{ab}R_{iabj}$$

are consequences of (37); there is, in other words, only one way (up to a sign) to contract the metric into the curvature tensor. The *Ricci tensor* is defined[41]

$$R_{ij} \equiv g^{ab}R_{aijb} = R^a{}_{ija} \tag{40}$$

and is, by (37), symmetric. More obviously, there is only one way to contract the metric into the Ricci tensor; it yields the *curvature invariant*:

$$R \equiv g^{ij}R_{ij} = R^i{}_i \tag{41}$$

The Bianchi identities (39), with input from (37), can be written

$$R_{abij;k} - R_{bajk;i} - R_{abik;j} = 0 \tag{42.1}$$

and, when contracted into $g^{ai}g^{bj}$, give $R_{;k} - 2R^i{}_{k;i} = 0$ which can be expressed

$$(R^i{}_k - \tfrac{1}{2}R\delta^i{}_k)_{;i} = 0 \tag{42.2}$$

The preceding equations, known as the "contracted Bianchi identities," can be read as an assertion that the (covariant) divergence of the *Einstein tensor*

$$G_{ij} \equiv R_{ij} - \tfrac{1}{2}Rg_{ij} \quad : \quad G_{ij} = G_{ji} \tag{43}$$

---

[40] See, for example, p. 638 in M. Kaku, *Quantum Field Theory* (1993), or §15.1 "The geometry of gauge invariance" in M. E. Peskin & D. V. Schroeder, *An Introduction to Quantum Field Theory* (1995).

[41] Careful! I have followed the conventions of Synge & Schild, but many modern authors write $R^i{}_{ajk}$ where I have written $R^i{}_{akj}$; i.e., they work with the *negative* of my Riemann tensor. But where I write $R_{ij} \equiv R^a{}_{ija}$ they write $R_{ij} \equiv R^a{}_{iaj}$, so we are in agreement on the definition of the Ricci tensor, and of $R$.

vanishes: $G^{ij}{}_{;i} = 0$. The contracted identities (42.2) were first noted by Aurel Voss in 1880, rediscovered by Ricci in 1889 and rediscovered again by Luigi Bianchi in 1902, the argument in each instance being heavily computational; the observation that (42.2) follows quickly from "the" Bianchi identities (42.1) was not made until 1922. The geometrical statements (42.2) strike the eye of a physicist as a quartet of *conservation laws*, and it is as such that they have come to play an important role in gravitational theory. But they remained unknown to Einstein and Hilbert (and to most other experts, except for Weyl) during the years when they were most needed.[42]

But what have the Riemann curvature tensor and its relatives got to do with "curvature" in any familiar geometrical sense? To start with the almost obvious: *If* there exists *some* coordinate system $x$ with respect to which the components of $g_{ij}$ *all become constant*, then (see again the definitions (26) of the Christoffel symbols) the $\left\{ {}^{i}_{jk} \right\}$ all *vanish* in those coordinates, and so also therefore (see again the definition (36)) do the components $R^{a}{}_{ijk}$ of the curvature tensor. Coordinate adjustment $x \rightarrow y$ will, in general, cause the metric to no longer be constant, and the Christoffel symbols to no longer vanish, but because $R^{a}{}_{ijk}$ transforms *as a tensor* it will *vanish in all coordinate systems*. It turns out that the converse is also true:

$$R^{a}{}_{ijk} = 0 \text{ if and only if there exists a coordinate}$$
$$\text{system in which the metric } g_{ij} \text{ becomes constant}$$

In such a coordinate system $g_{ij}$ can, by rotation, be diagonalized, and by dilation one can arrange to have only $\pm 1$'s appear on the principal diagonal. If all signs are positive, then one has, in effect, "Cartesian coordinatized Euclidean $N$-space," and in all cases it becomes sensible to interpret $R^{a}{}_{ijk} = 0$ as a "flatness condition." Which is to interpret $R^{a}{}_{ijk} \neq 0$ as the condition that the Riemannian manifold $\mathcal{M}$ be "not flat, or curved." To test the plausibility of the interpretation we must look to circumstances in which we have some *prior* conception of curvature. Though we might look to "spheres in $N$-space,"[43] concerning which we have some sense of what the constant curvature should be, it makes more sense (and is easier) to look to the case $N = 2$, where our intuitions are vivid, and where additionally we have Gauss' analytical accomplishments to guide us.

In the case $N = 2$ the Riemann tensor has only one independent element;

$$R_{abij} = \begin{cases} \pm R_{1212} \\ \text{else } 0 \end{cases} \quad \text{on 2-dimensional manifolds}$$

Look to the case $(ds)^2 = r^2(d\theta)^2 + (r\sin\theta)^2(d\phi)^2$, which we saw at (14) refers to the spherical coordinatization of a globe of radius $r$ (Figure 3). In *MTW* (at p. 340 in Chapter 14: "Calculation of Curvature") we are walked through

---

[42] See §15c in Pais for further historical commentary.

[43] See §4 in "Algebraic theory of spherical harmonics" (1996) for indication of how this might be done.

the demonstration that
- $\Gamma^{\theta}{}_{\phi\phi} = -\sin\theta\cos\theta$; $\Gamma^{\phi}{}_{\phi\theta} = \Gamma^{\phi}{}_{\theta\phi} = \cot\theta$; other $\Gamma$'s vanish;
- $R_{\theta\phi\theta\phi} = r^2\sin^2\theta$, with other elements determined by symmetries;
- $R^{\theta}{}_{\theta} = R^{\phi}{}_{\phi} = 1/r^2$; $R^{\theta}{}_{\phi} = R^{\phi}{}_{\theta} = 0$;
- $R = 2/r^2$

So at least in this hallowed case the curvature tensor speaks (through the curvature scalar $R$) directly to what we are prepared to call the "curvature of the sphere." More generally,[44] the Gaussian curvature can be described

$$K = \frac{R_{1212}}{g}$$
$$= \frac{r^2\sin^2}{r^4\sin^2\theta} = 1/r^2 \quad \text{in the spherical case just considered}$$

We conclude on this evidence that $R^{a}{}_{ijk}$ reproduces what we already knew about curvature, and puts us in position to say things we didn't already know. More particularly, it achieves a vast generalization of ideas pioneered by Gauss —ideas to which the young Einstein had been exposed, but had paid no special attention.

**Variational approach to the gravitational field equations**.   Once Einstein had acquired the conviction that
- gravitation is an artifact of the (Riemannian) geometry of spacetime, and
- general covariance is an essential formal property of any theory of gravity

he set out discover the field equations to which $g_{\mu\nu}$ must be subject.[45]   His principal guidance was supplied by Newton; i.e., by Poisson's equation

$$\nabla^2\varphi = 4\pi G\rho \tag{5}$$

from which he inferred that $g_{\mu\nu}(x)$ should satisfy a generally covariant coupled system of second-order partial differential equations—equations which are, in particular, linear and homogeneous in the second partials.[46] The tensor analytic theory of Riemannian manifolds supplies, as we have seen, only limited material, and led him to contemplate field equations of the narrowly constrained design

$$R_{\mu\nu} + \alpha R g_{\mu\nu} + \Lambda g_{\mu\nu} = \kappa T_{\mu\nu} \tag{44}$$

where $\alpha$ and $\Lambda$ are adjustable constants, $T_{\mu\nu}$ is the (necessarily symmetric) stress-energy tensor of such non-gravitational (material, electromagnetic) fields as may inhabit spacetime, and $\kappa$ is to be fixed by dimensional considerations (units).   It was to achieve the conservation law (42.2) that he—somewhat

---

[44]  See McCleary,[20] p. 155 for the demonstration.

[45]  I revert to Greek indices to emphasize that we now inhabit not some arbitrary Riemannian manifold, but the *physical spacetime manifold*, which is understood to be 4-dimensional.

[46]  See his *The Meaning of Relativity* (3ʳᵈ revised edition 1950), p. 84.

tentatively—set $\alpha = -\frac{1}{2}$ and $\Lambda$ (the so-called "cosmological constant") equal to zero.

Hilbert, on the other hand, chose to look upon the problem posed by Einstein as a straightforward problem in Lagrangian field theory; i.e., to work from the assumption that the field equations can be derived from a variational principle.[47] His initial effort, therefore, was to build Einstein's requirements into the design of a Lagrange density. Before we attempt to follow in Hilbert's steps, let us. . .

Recall again (from elementary calculus) that integrals respond to a change of variables $x \to y$ by "picking up a Jacobian"

$$\iiiint_{\text{bubble}} f(x)\,dx^0 dx^1 dx^2 dx^3$$
$$= \iiiint_{\text{image bubble}} \underbrace{f(x(y)) \left| \frac{\partial(x^0, x^1, x^2, x^3)}{\partial(y^0, y^1, y^2, y^3)} \right|}_{F(y)} dy^0 dy^1 dy^2 dy^3$$

The general covariance of an action functional

$$S \equiv \iiiint_{\mathcal{R}} \mathcal{L}\,dx^0 dx^1 dx^2 dx^3 = \iiiint_{\mathcal{R}'} \mathcal{L}'\,dx'^0 dx'^1 dx'^2 dx'^3$$

therefore requires that $\mathcal{L}$ respond to $x \to x'$ by the rule

$$\mathcal{L} \to \mathcal{L}' = \mathcal{L} \cdot \left| \tfrac{\partial x}{\partial x'} \right|$$

i.e., as a *scalar density of unit weight*. In special relativity the Jacobian has value $\pm 1$, so (as was remarked already at p. 9 in Chapter 2) to draw attention to the point just made is to underscore a "distinction without a difference." But in general relativity the point assumes non-trivial importance. It was remarked previously that in metric theories $\sqrt{g}$ supplies the "prototypical instance" of a scalar density of unit weight. We expect therefore to be able to write

$$\mathcal{L} = \sqrt{g} \cdot (\text{generally covariant scalar of zero weight})$$

Another preparatory detail. Assume for the moment that $[x^\mu] = length$ for all $\mu$. Then the $g_{\mu\nu}$ are all dimensionless, giving[48]

$$[R_{aijk}] = [R_{ij}] = [R] = 1/(\text{length})^2$$

---

[47] Pauli, writing in *1921*, considered reliance upon a variational principle to be a defect, "unacceptable to physicists". . . as at the time it assuredly was; today most physicists would consider this to be the principal virtue of Hilbert's approach.

[48] The following statements remain in force even if some of the generalized coordinates are (for example) angles, since compensating dimensionality— sufficient to preserve $[ds] = length$ —will attach then to associated components of the metric. For that same reason, $[\sqrt{g}\,dx^0 dx^1 dx^2 dx^3] = (length)^4$ in all cases.

We want to achieve $[\mathcal{L}] = \text{energy density} = ML^{-1}T^{-2}$, and have only $R$, $c$ (a velocity) and $G$ (of dimension $M^{-1}L^3T^{-2}$) to work with. Write

$$\mathcal{L} \sim G^\alpha c^\beta R^\sigma \quad : \quad ML^{-1}T^{-2} = (M^{-1}L^3T^{-2})^\alpha(LT^{-1})^\beta L^{-2\sigma}$$
$$= M^{-\alpha}L^{3\alpha+\beta-2\sigma}T^{-2\alpha-\beta}$$

and find that necessarily $\alpha = -1$, $\beta = 4$ and $\sigma = 1$. We expect therefore to have

$$\mathcal{L} = (\text{dimensionless numeric})(c^4/G) \cdot \sqrt{g}\left\{R - 2\Lambda\right\} \tag{45}$$

where $\Lambda$ is a dimensioned constant ($[\Lambda] = [R]$) and the anticipates recovery of (44). If we possessed a "natural length" $\lambda$ then we could install also a factor of the form

$$\left\{1 + a(\lambda^2 R) + b(\lambda^2 R)^2 + \cdots\right\}$$

but $(G, c)$-theory presents no such object; $(G, c, \hbar)$-theory, on the other hand, does: we have the

$$\text{Planck length } \lambda \equiv \hbar G/c^3 = 1.616 \times 10^{-33}\,\text{cm}$$

which is too small to be of any direct relevance to macroscopic theory.[49] The prefactor in (45) will come into play when we undertake to establish detailed contact with Poisson's equation (5); until then we will ignore it, writing[50]

$$\mathcal{L} \sim \sqrt{g}\left\{R - 2\Lambda\right\} \tag{46}$$

In the same sense that (say) a sphere is "locally Euclidean" in the neighborhood of every point, we expect the spacetime manifold $\mathcal{M}$ to be "locally Lorentzian." More specifically, we expect to be able to write

$$\|g_{\mu\nu}(x)\| = \mathbb{A}^\mathsf{T}(x)\,\mathbb{G}\,\mathbb{A}(x) \quad \text{with} \quad \mathbb{G} \equiv \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix} \tag{47.1}$$

with $\mathbb{A}(x)$ determined only to within a local Lorentz transformation:

$$\mathbb{A}(x) \to \mathbb{A}'(x) \equiv \mathbb{L}\mathbb{A}(x) \quad : \quad \mathbb{L}^\mathsf{T}\mathbb{G}\mathbb{L} = \mathbb{G} \tag{47.2}$$

---

[49] In this connection see S. Weinberg, *Gravitation and Cosmology: Principles and Applications of the General Theory of Relativity* (1972), p. 365.

[50] In four dimensions—exceptionally—it is possible to contemplate a scalar density of the design $\epsilon^{\alpha\mu\nu\sigma}R_{\alpha\mu\nu\sigma}$. But antisymmetry properties of the Levi-Civita tensor can be used to write

$$\epsilon^{\alpha\mu\nu\sigma}R_{\alpha\mu\nu\sigma} = \tfrac{1}{3}\epsilon^{\alpha\mu\nu\sigma}(R_{\alpha\mu\nu\sigma} + R_{\alpha\nu\sigma\mu} + R_{\alpha\sigma\mu\nu})$$
$$= 0 \quad \text{by (37.2)}$$

so such a scalar can play no role in the theory.

The implication is that we expect to have

$$g(x) \equiv \det \|g_{\mu\nu}(x)\| = -(\det \mathbb{A})^2 < 0 \tag{48}$$

It is in light of this circumstance, and to preserve manifest reality, that the general relativistic literature bristles with $\sqrt{-g}$ factors. I find the minus signs a distraction, so will (in the good company of Weinberg) drop them (i.e., I will write $\sqrt{-g} = i\sqrt{g}$ and absorb the $i$ into the meaning of "$\sim$"), and pick them up again only when their presence makes a difference.

If we interpret the field functions to be $g_{\mu\nu}(x)$, and recall from (26) the definitions of the Christoffel symbols, then (46) assumes the design

$$\mathcal{L}(g, \partial g, \partial\partial g)$$

which, owing to the presence of the second derivatives, requires—or appears to require—an extension of Lagrangian field theory.[51] Several familiar tricks are available: we might, for example, expand the number fields, writing $\mathcal{L}(g, h, \partial h)$ ... though this, so far as I am aware, is never done. Alternatively, one can borrow a trick from Procca field theory (a trick which is often useful also in electrodynamics and in many other applications): we found at (2–31) that it is formally advantageous to consider $U^\mu$ and $G^{\mu\nu} \equiv \partial^\mu U^\nu - \partial^\nu U^\mu$ to be *independent* fields, even though it is clearly impossible to vary $G^{\mu\nu}$ while holding $U^\mu$ constant. "Palatini's method"[52] proceeds similarly: one opts to look upon $g_{\mu\nu} = g_{\nu\mu}$ and $\Gamma^\alpha{}_{\mu\nu} = \Gamma^\alpha{}_{\nu\mu}$ as formally independent fields (10 of the former, 40 of the latter), and hopes to recover the definitions (26.2) as forced implications of an expanded set of field equations

$$\left\{ \partial_\sigma \frac{\partial}{\partial g_{\mu\nu,\sigma}} - \frac{\partial}{\partial g_{\mu\nu}} \right\} \mathcal{L} = 0 \tag{49.1}$$

$$\left\{ \partial_\sigma \frac{\partial}{\partial \Gamma^\rho{}_{\mu\nu,\sigma}} - \frac{\partial}{\partial \Gamma^\rho{}_{\mu\nu}} \right\} \mathcal{L} = 0 \tag{49.2}$$

where the Lagrangian has now the design

$$\mathcal{L}(g_{..}, \Gamma^{\cdot}{}_{..}, \partial\Gamma^{\cdot}{}_{..}) \sim \sqrt{g} \left\{ g^{\alpha\beta} R_{\alpha\beta}(\Gamma^{\cdot}{}_{..}, \partial\Gamma^{\cdot}{}_{..}) - 2\Lambda \right\} \tag{50}$$

We look first to (49.1), where the absence from $\mathcal{L}$ of any $\partial g$-dependence results in some welcome simplification. We have

$$\sqrt{g}\, R_{\alpha\beta} \frac{\partial}{\partial g_{\mu\nu}} g^{\alpha\beta} + \left\{ R - 2\Lambda \right\} \frac{\partial}{\partial g_{\mu\nu}} \sqrt{g} = 0$$

and—drawing upon the soon-to-be-established information that

$$\frac{\partial}{\partial g_{\mu\nu}} g^{\alpha\beta} = -g^{\mu\alpha} g^{\beta\nu} \quad \text{and} \quad \frac{\partial}{\partial g_{\mu\nu}} \sqrt{g} = \tfrac{1}{2}\sqrt{g}\, g^{\mu\nu} \tag{51}$$

---

[51] A clever way to circumvent this problem—due to Dirac—will be described later.

[52] A. Palatini, "Deduzione invariantiva delle equazioni gravitazioni dal principio di Hamilton," Rend. Circ. Mat. Palermo **43**, 203 (1919); A. Einstein, "Einheitliche Feldtheorie von Gravitation und Elektrizität," Preussische Akademie der Wissenschaften (1925), p. 414.

—obtain $\sqrt{g}\left\{-R^{\mu\nu}+\frac{1}{2}(R-2\Lambda)g^{\mu\nu}\right\}=0$, or again

$$R^{\mu\nu}-\tfrac{1}{2}Rg^{\mu\nu}+\Lambda g^{\mu\nu}=0 \tag{52}$$

This is an "empty universe" instance ($T^{\mu\nu}=0$) of Einstein's gravitational field equations (44). Einstein's tentative $\alpha=-\frac{1}{2}$ is seen to have been *forced by the variational principle*. Notice that the $\frac{1}{2}$ has its origin in the $\sqrt{\phantom{.}}$; i.e., in the circumstance that $\mathcal{L}$ transforms as a scalar *density*. It is interesting in the light of more recent developments that Einstein's $\Lambda=0$ is *not* forced.[53]

I digress to establish equations (51). Observe first that differentiation of $x^{-1}x=1$ gives $x\frac{d}{dx}x^{-1}+x^{-1}=0$ and provides a demonstration (as if one were needed) that $\frac{d}{dx}x^{-1}=-x^{-2}$. Similarly... let $\mathbb{A}\equiv\|a_{ij}\|$ be an invertible square matrix, and let $\mathbb{A}^{-1}$ be notated $\|a^{ij}\|$. Differentiation of $a^{ik}a_{kj}=\delta^{i}{}_{j}$ gives

$$\frac{\partial a^{ik}}{\partial a_{pq}}a_{kj}+a^{ik}\frac{\partial a_{kj}}{\partial a_{pq}}=\frac{\partial a^{ik}}{\partial a_{pq}}a_{kj}+a^{ik}\delta^{p}{}_{k}\delta^{q}{}_{j}=0$$

$$\frac{\partial a^{ik}}{\partial a_{pq}}a_{kh}=-a^{ip}\delta^{q}{}_{h}$$

Multiplication by $a^{hj}$ gives

$$\frac{\partial a^{ij}}{\partial a_{pq}}=-a^{ip}a^{qj}$$

which establishes the first part of (51). To gain insight into the second part, suppose that $\mathbb{A}$ were $2\times2$; we would then have

$$a\equiv\det\mathbb{A}=a_{11}a_{22}-a_{12}a_{21}$$

and

$$\mathbb{A}^{-1}=a^{-1}\begin{pmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{pmatrix}=\begin{pmatrix} \frac{1}{a}\frac{\partial a}{\partial a_{11}} & \frac{1}{a}\frac{\partial a}{\partial a_{21}} \\ \frac{1}{a}\frac{\partial a}{\partial a_{12}} & \frac{1}{a}\frac{\partial a}{\partial a_{22}} \end{pmatrix}$$

giving

$$a^{ij}=\frac{1}{a}\frac{\partial a}{\partial a_{ji}}=\frac{\partial}{\partial a_{ji}}\log\det\mathbb{A} \tag{53}$$

----

[53] The "cosmological term" $\Lambda g^{\mu\nu}$ was introduced into (52) by Einstein in 1917, in an effort to make general relativity conform to what he imagined to be the steady state of the universe; the event took place at equation (13a) in a paper "Kosmologische Betrachtungen zur allgemeinen Relativitätstheorie" of which an English translation can be found in the Dover edition of *The Principle of Relativity: A Collection of Original Memoirs on the Special and General Theory of Relativity*. He had become disenchanted with the term already by 1923, and officially abandoned it in 1931 (see Pais' §15.e). For indication of why there is renewed interest in the term, see Alan Guth, *The Inflationary Universe* (1997), p. 283.

whence

$$\frac{\partial}{\partial g_{\mu\nu}}\sqrt{g} = \frac{1}{2\sqrt{g}}\frac{\partial}{\partial g_{\mu\nu}}g = \frac{1}{2\sqrt{g}}\,g\,g^{\nu\mu}$$

from which (in the 2-dimensional case) the desired result follows by $g^{\nu\mu} = g^{\mu\nu}$. To establish the general validity of (53) we use the Laplace expansion

$$a = a_{j1}A_{j1} + a_{j2}A_{j2} + \cdots + a_{jN}A_{jN}$$

to obtain

$$\frac{\partial a}{\partial a_{ji}} = A_{ji} \equiv \text{cofactor of } a_{ji}$$

But by Cramer's Rule

$$a^{ij} = a^{-1}A_{ji}$$

which completes the argument. We will have need in a moment also of the second of these corollary of (51):

$$(\sqrt{g})_{,\sigma} = \tfrac{1}{2}\sqrt{g}\,g^{\mu\nu}g_{\mu\nu,\sigma} = -\tfrac{1}{2}\sqrt{g}\,g_{\mu\nu}g^{\mu\nu}{}_{,\sigma} \tag{54}$$

Look now to (49.2), the second set of "field equations," which can be written

$$\left\{\partial_\sigma\frac{\partial}{\partial\Gamma^\rho{}_{\mu\nu,\sigma}} - \frac{\partial}{\partial\Gamma^\rho{}_{\mu\nu}}\right\}\sqrt{g}\,g^{\alpha\beta}R_{\alpha\beta} = 0$$

with[54]

$$R_{\alpha\beta} = \Gamma^i{}_{\alpha k,\beta}\,\delta^k{}_i - \Gamma^i{}_{\alpha\beta,k}\,\delta^k{}_i + \Gamma^i{}_{\alpha j}\Gamma^j{}_{i\beta} - \Gamma^i{}_{\alpha\beta}\Gamma^j{}_{ik}\,\delta^k{}_j$$

Carefully performing the indicated differentiations,[55] we obtain

$$\partial_\sigma\left\{\sqrt{g}\left[\tfrac{1}{2}g^{\mu\sigma}\delta^\nu{}_\rho + \tfrac{1}{2}g^{\nu\sigma}\delta^\mu{}_\rho - g^{\mu\nu}\delta^\sigma{}_\rho\right]\right\}$$
$$= \sqrt{g}\left[g^{\mu\alpha}\Gamma^\nu{}_{\rho\alpha} + g^{\nu\alpha}\Gamma^\mu{}_{\rho\alpha} - g^{\mu\nu}\Gamma^\alpha{}_{\rho\alpha} - \tfrac{1}{2}g^{\alpha\beta}\Gamma^\mu{}_{\alpha\beta}\delta^\nu{}_\rho - \tfrac{1}{2}g^{\alpha\beta}\Gamma^\nu{}_{\alpha\beta}\delta^\mu{}_\rho\right]$$

which by reorganization (I make free use of the $\mu\nu$-symmetry of $g_{\mu\nu}$ and $\Gamma^\alpha{}_{\mu\nu}$, and have underscored a couple of places where I have introduced a term promptly to subtract it again) becomes

$$\left[\sqrt{g}\,g^{\mu\nu}\right]_{,\rho} + \sqrt{g}\left[g^{\alpha\nu}\Gamma^\mu{}_{\alpha\rho} + g^{\mu\alpha}\Gamma^\nu{}_{\alpha\rho} - g^{\mu\nu}\Gamma^\alpha{}_{\alpha\rho}\right]$$
$$= \tfrac{1}{2}\left\{\left[\sqrt{g}\,g^{\mu\sigma}\right]_{,\sigma} + \sqrt{g}\left[g^{\alpha\sigma}\Gamma^\mu{}_{\alpha\sigma} + \underline{g^{\mu\alpha}\Gamma^\sigma{}_{\alpha\sigma}} - g^{\mu\sigma}\Gamma^\alpha{}_{\alpha\sigma}\right]\right\}\delta^\nu{}_\rho$$
$$+ \tfrac{1}{2}\left\{\left[\sqrt{g}\,g^{\nu\sigma}\right]_{,\sigma} + \sqrt{g}\left[g^{\alpha\sigma}\Gamma^\nu{}_{\alpha\sigma} + \underline{g^{\nu\alpha}\Gamma^\sigma{}_{\alpha\sigma}} - g^{\nu\sigma}\Gamma^\alpha{}_{\alpha\sigma}\right]\right\}\delta^\mu{}_\rho$$

[54] Depleted ranks here force me to press some Roman soldiers into Greek service. The $\delta$'s have been introduced to prevent repeated indices from appearing on any $\Gamma$, which simplifies calculation in the present context.

[55] The assumed symmetry of $\Gamma^\alpha{}_{\mu\nu}$ presents a formal problem similar to that confronted/resolved in connection with (1–44). The equations which follow have been written in such a way as to display *manifest $\mu\nu$-symmetry*.

Recalling from (25) some defining properties of the covariant derivative, we find that the preceding equation can be expressed

$$\left[\sqrt{g}\,g^{\mu\nu}\right]_{;\rho} = \tfrac{1}{2}\left[\sqrt{g}\,g^{\mu\sigma}\right]_{;\sigma}\delta^\nu{}_\rho + \tfrac{1}{2}\left[\sqrt{g}\,g^{\nu\sigma}\right]_{;\sigma}\delta^\mu{}_\rho$$

or again (here I adopt the notation of $CTW$, p. 502)

$$\mathfrak{g}^{\mu\nu}{}_{;\rho} - \tfrac{1}{2}\delta^\mu{}_\rho\,\mathfrak{g}^{\nu\sigma}{}_{;\sigma} - \tfrac{1}{2}\delta^\nu{}_\rho\,\mathfrak{g}^{\mu\sigma}{}_{;\sigma} = 0 \tag{55}$$

where $\mathfrak{g}^{\mu\nu} \equiv \sqrt{g}\,g^{\mu\nu}$ defines a "metric density" or unit weight. Look upon (55) as a homogeneous linear system of 40 equations in 40 unknowns:

$$\begin{pmatrix} \bullet & \bullet & \cdots & \bullet \\ \bullet & \bullet & \cdots & \bullet \\ \vdots & \vdots & & \vdots \\ \bullet & \bullet & \cdots & \bullet \end{pmatrix} \begin{pmatrix} \mathfrak{g}^{00}{}_{;0} \\ \mathfrak{g}^{01}{}_{;0} \\ \vdots \\ \mathfrak{g}^{44}{}_{;4} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Writing the $40{\times}40$ matrix into *Mathematica* we discover it to be non-singular, so (55) will be satisfied if and only if

$$\mathfrak{g}^{\mu\nu}{}_{;\sigma} = 0 \tag{56}$$

It remains ($i$) to show that this result implies (and is implied by) $g_{\mu\nu;\sigma} = 0$, and ($ii$) to discuss the remarkable significance of that fact. When written out in detail, (56) reads

$$(\sqrt{g}g^{\mu\nu})_{,\sigma} + \sqrt{g}g^{\alpha\nu}\Gamma^\mu{}_{\alpha\sigma} + \sqrt{g}g^{\mu\alpha}\Gamma^\nu{}_{\alpha\sigma} - \sqrt{g}g^{\mu\nu}\Gamma^\alpha{}_{\alpha\sigma} = 0$$

which when contracted into $g_{\mu\nu}$ gives

$$(\sqrt{g})_{,\sigma}N + \sqrt{g}g_{\mu\nu}g^{\mu\nu}{}_{,\sigma} + \sqrt{g}\Gamma^\alpha{}_{\alpha\sigma} + \sqrt{g}\Gamma^\alpha{}_{\alpha\sigma} - N\sqrt{g}\Gamma^\alpha{}_{\alpha\sigma} = 0$$

where $N = \delta^\alpha{}_\alpha =$ dimension of the spacetime manifold. Drawing now upon (54), we have

$$(N-2)\Big\{(\sqrt{g})_{,\sigma} - \sqrt{g}\Gamma^\alpha{}_{\alpha\sigma}\Big\} = 0$$

which by (25.2) becomes $(N-2)(\sqrt{g})_{;\sigma} = 0$ and (if $N \neq 2$) entails

$$(\sqrt{g})_{;\sigma} = 0 \tag{57}$$

This result, we not in passing, supplies the following often-useful information:

$$\Gamma^\alpha{}_{\alpha\sigma} = \frac{1}{\sqrt{g}}\frac{\partial}{\partial x^\sigma}\sqrt{g} = \frac{\partial}{\partial x^\sigma}\log\sqrt{g} \tag{58}$$

Now construct $\mathfrak{g}^{\alpha\beta}g_{\beta\nu} = \sqrt{g}\,\delta^\alpha{}_\nu$ and covariantly differentiate:

$$\mathfrak{g}^{\alpha\beta}{}_{;\sigma}g_{\beta\nu} + \mathfrak{g}^{\alpha\beta}g_{\beta\nu;\sigma} = (\sqrt{g})_{;\sigma}\delta^\alpha{}_\nu + \sqrt{g}\,\delta^\alpha{}_{\nu;\sigma}$$
$$\Downarrow$$
$$\mathfrak{g}^{\alpha\beta}g_{\beta\nu;\sigma} = 0$$

Contract into $g_{\mu\alpha}$ and obtain

$$g_{\mu\nu;\sigma} = 0 \tag{59}$$

Clearly, the order of the argument could be reversed: (56) $\Leftrightarrow$ (59).

At (59) we have recovered—now as "field equations"—the conditions which were previously seen to be necessary and sufficient for $\Gamma^{\cdot}{}_{\cdot\cdot}$-mediated covariant differentiation and $g_{\cdot\cdot}$-mediated index manipulation to performable in either order, conditions which at (26.2) were seen to entail

$$\Gamma^{\sigma}{}_{\mu\nu} = \tfrac{1}{2}g^{\sigma\alpha}\big\{g_{\alpha\mu,\nu} + g_{\alpha\nu,\mu} - g_{\mu\nu,\alpha}\big\}$$

That system of equations serves to locate "Riemannian geometry" within the broader class of (torsion-free) "affine geometries," and makes precise the sense in which (in Riemannian geometry) *metric structure dictates affine structure*. Einstein, having been led to embrace the Principle of General Covariance, found (actually, Hilbert found) the design (46) of the Lagrangian $\mathcal{L}$ to be essentially forced, but while he

- needed $\Gamma^{\sigma}{}_{\mu\nu}$ to construct $R^{\sigma}{}_{\mu\nu\rho}$
- needed $g_{\mu\nu}$ to construct $R \equiv g^{\sigma\rho}g^{\mu\nu}R^{\sigma}{}_{\mu\nu\rho}$ and to supply a density

he did not need to *assume the compatability* of those connections; compatability was (at least within the Palatini formalism) *automatically enforced by the variational principle*.

To describe the geometry of spacetime in a world *not* devoid of matter Einstein (Hilbert) makes the *ad hoc* adjustment

$$\mathcal{L}_{\text{gravitational field}} \quad \longmapsto \quad \mathcal{L}_{\text{gravitational field}} + \mathcal{L}_{\text{matter}} \tag{60}$$

Note the absence of an explicit "interaction term." The coupling of matter to gravitation is accomplished implicitly, through in the requirement that $\mathcal{L}_{\text{matter}}$ be a *generally covariant density* (i.e., by introducing occasional $\sqrt{g}$-factors, and replacing some commas with semi-colons!). In the simplest case one might write (compare (2–11))

$$\mathcal{L}_{\text{matter}} \sim \sqrt{g}\big\{g^{\mu\nu}\varphi_{;\mu}\varphi_{;\nu} - \varkappa^2\varphi^2\big\}$$

on the assumption that $\varphi$ is a weightless scalar field (in which case $\varphi_{;\mu} = \varphi_{,\mu}$), or still more simply

$$\mathcal{L}_{\text{matter}} \sim \big\{g^{\mu\nu}\varphi_{;\mu}\varphi_{;\nu} - \varkappa^2\varphi^2\big\}$$

on the assumption that $\varphi$ transforms as a scalar density of weight $w = \tfrac{1}{2}$ (in which case $\varphi_{;\mu} = \varphi_{,\mu} - \tfrac{1}{2}\varphi\Gamma^{\alpha}{}_{\alpha\mu}$). I will return later to discussion of some of the general relativistic ramifications of (60).

I have organized the preceding discussion in a way intended to emphasize that gravitational field theory is (in at least its variational aspects) *classical field theory like any other*. All followed from (49). I cannot account for the fact that the authors of the standard monographs prefer *not* to work from those elegant equations, but to "reinvent" variational methodology as they go along.[56] This

---

[56] See, for example, *MTW* §21.2; Weinberg,[49] Chapter 12; J. L. Anderson, *Principles of Relativity Physics* (1967), §10-4.

is in marked contrast to the tradition established by Einstein himself, whose variational remarks are stylistically more similar to my own.[57]

But Lagrangian field theory is an elastic vessel. I sketch now an alternative approach due to Dirac.[58] Dirac proceeds from $\mathcal{L} \sim \sqrt{g}\,R$ but (unlike Palatini) is prepared to assume at the outset that $\Gamma^{\sigma}{}_{\mu\nu} = \frac{1}{2} g^{\sigma\alpha} \{ g_{\alpha\mu,\nu} + g_{\alpha\nu,\mu} - g_{\mu\nu,\alpha} \}$; i.e., that general relativity is an exercise in *Riemannian* geometry, nothing more abstruse. So he confronts a Lagrangian of the design $\mathcal{L}(g, \partial g, \partial\partial g)$, but writing

$$R = \underbrace{g^{\mu\nu}(\Gamma^{\sigma}{}_{\mu\sigma,\nu} - \Gamma^{\sigma}{}_{\mu\nu,\sigma})}_{Q} - \underbrace{g^{\mu\nu}(\Gamma^{\sigma}{}_{\mu\nu}\Gamma^{\rho}{}_{\rho\sigma} - \Gamma^{\sigma}{}_{\mu\rho}\Gamma^{\rho}{}_{\nu\sigma})}_{R^*}$$

he observes that the offending $\partial\partial g$ terms are present only in the $Q$ term (into which they enter linearly, and) from which they can be gauged away. More particularly, Dirac shows that

$$\sqrt{g}\,Q = 2\sqrt{g}\,R^* + \underbrace{(\sqrt{g}\,[g^{\rho\mu}\Gamma^{\sigma}{}_{\rho\sigma} - g^{\rho\sigma}\Gamma^{\mu}{}_{\rho\sigma}])_{,\mu}}_{\text{gauge term}}$$

so

$$\mathcal{L} \sim \sqrt{g}\,R = \sqrt{g}\,R^* + \text{gauge term}$$

Dirac *abandons the gauge term* (where all $\partial\partial g$ terms reside), electing to work from

$$\mathcal{L}^* \sim \sqrt{g}\,R^* = \tfrac{1}{4}\sqrt{g}\big\{ (g^{\mu\alpha}g^{\nu\beta} - g^{\mu\nu}g^{\alpha\beta})g^{\rho\sigma} \\ - 2(g^{\mu\rho}g^{\alpha\beta} - g^{\mu\alpha}g^{\beta\rho})g^{\nu\sigma} \big\} g_{\mu\nu,\rho}g_{\alpha\beta,\sigma} \tag{61}$$

which we see to be homogeneous of degree two in $\partial g$. This (because the rewards are so great) he is content to do even though $\mathcal{L}^*$ *does not transform as a scalar density*; indeed, he considers the latter circumstance to be not a defect of the formalism but evidence that "four-dimensional symmetry is not a fundamental property of the physical world." Working from (61) he computes

$$\left\{ \partial_\sigma \frac{\partial}{\partial g_{\mu\nu,\sigma}} - \frac{\partial}{\partial g_{\mu\nu}} \right\} \sqrt{g}\,R^* = -\sqrt{g}\,(R^{\mu\nu} - \tfrac{1}{2} R g^{\mu\nu}) \tag{62.1}$$

---

[57] "Hamiltonsches Princip und allgemeine Relativitätstheorie" (1916), of which an English translation can be found in the Dover collection cited earlier.[52] Einstein cites, in addition to Hilbert, four papers by Lorentz (1915 & 1916).

[58] See Chapter 26—two and one half pages long—in his elegantly slim *General Theory of Relativity* (1996), to which I refer my reader for all the omitted details. Dirac presented his argument ("Theory of gravitation in Hamiltonian form," Proc. Roy. Soc. **A246**, 333 (1958)) for its methodological interest, as an illustrative application of ideas developed in "Generalized Hamiltonian dynamics," Proc. Roy. Soc. **A246**, 326 (1958).

I draw belated attention to the fact that in $N$-dimensional spacetime ($N \neq 2$)

$$R^{\mu\nu} - \tfrac{1}{2}Rg^{\mu\nu} = 0 \quad \Rightarrow \quad R = 0 \quad \Rightarrow \quad R^{\mu\nu} = 0 \tag{63}$$

With this fact in mind, Dirac observes that

$$\left\{ \partial_\sigma \frac{\partial}{\partial \mathfrak{g}^{\mu\nu}{}_{,\sigma}} - \frac{\partial}{\partial \mathfrak{g}^{\mu\nu}} \right\} \sqrt{g}\, R^* = R_{\mu\nu} \tag{62.2}$$

**Harmonic coordinates: Gravitational analog of the Lorentz gauge condition**. By way of preparation, look back again to equations (2–24), where it is observed that the electromagnetic field equations can be expressed

$$\Box A^\nu - \partial^\nu(\partial_\mu A^\mu) = J^\nu \tag{64.1}$$

This looks like a system of four equations in four unknown fields $A^\mu$. But the expressions on the left are (trivially and automatically) subject to a differential identity

$$\left[ \Box A^\nu - \partial^\nu(\partial_\mu A^\mu) \right]_{,\nu} = 0 \tag{64.2}$$

—the solitary electromagnetic analog the contracted Bianchi identities (42.2). So the $A^\nu$ which satisfy (64.1) still *retain one degree of freedom*, familiar to us as gauge freedom $A^\nu \mapsto A'^\nu = A^\nu + \partial^\nu\chi$. That freedom can be exploited in various ways to achieve simplifications.[59] For example, we can install the Lorentz gauge condition $\partial_\nu A^\nu = 0$, replacing (64.1) by a *quintet* of equations

$$\Box A^\nu = J^\nu \quad \text{and} \quad \partial_\mu A^\mu = 0 \tag{64.3}$$

The identity (614.2) still pertains, but the $A^\nu$ which satisfy the expanded set of field equations are unique (retain no degrees of freedom).[60] Similarly...

The Einstein equations

$$R^{\mu\nu} - \tfrac{1}{2}Rg^{\mu\nu} = \kappa T^{\mu\nu} \tag{65.1}$$

---

[59] It might be interesting, on another occasion, to consider whether the non-Abelian gauge transformation

$$\boldsymbol{A}_\mu \longrightarrow \boldsymbol{A}'_\mu = \boldsymbol{S}\boldsymbol{A}_\mu\boldsymbol{S}^{-1} + i\,\tfrac{1}{g}\,\boldsymbol{S}_{,\mu}\boldsymbol{S}^{-1} \tag{3–93}$$

can be used to achieve useful simplification of

$$\left. \begin{array}{l} \partial_\mu \boldsymbol{F}^{\mu\nu} = \tfrac{1}{c}\boldsymbol{s}^\nu \\ \qquad \boldsymbol{s}^\nu \equiv igc[\boldsymbol{F}^{\nu\alpha}, \boldsymbol{A}_\alpha] \end{array} \right\} \tag{3–106}$$

where $\boldsymbol{F}_{\mu\nu} \equiv (\partial_\mu \boldsymbol{A}_\nu - \partial_\nu \boldsymbol{A}_\mu) - ig\,(\boldsymbol{A}_\mu\boldsymbol{A}_\nu - \boldsymbol{A}_\nu\boldsymbol{A}_\mu)$.

[60] But see the cautionary statement two pages farther along!

appear on their face to be a system of ten equations in ten unknown fields $g^{\mu\nu}$. But the expressions on the left are (not quite trivially, but necessarily) subject to the quartet (42.2) of contracted Bianchi identities

$$\left[R^{\mu\nu} - \tfrac{1}{2}Rg^{\mu\nu}\right]_{;\nu} = 0 \tag{65.2}$$

So the $g^{\mu\nu}$ which satisfy (65.2) still *retain* $10 - 4 = 6$ *degrees of freedom*. This is associated in generally covariant theory with our freedom to recoordinatize, which is accomplished by the presentation of four (more or less) arbitrary functions: $x^{\mu} \mapsto \bar{x}^{\mu} = f^{\mu}(x)$.

To illustrate how "recoordinatization freedom" can be exploited to achieve simplifications, let $\varphi$ be a scalar field and look to the generally covariant construction

$$\begin{aligned}
\Box\varphi \equiv (g^{\mu\nu}\varphi_{;\nu})_{;\mu} &= (g^{\mu\nu}\varphi_{,\nu})_{;\mu} \\
&= (g^{\mu\nu}_{\ \ ;\mu})\varphi_{,\nu} + g^{\mu\nu}(\varphi_{,\mu\nu} - \varphi_{,\alpha}\Gamma^{\alpha}_{\ \mu\nu})
\end{aligned}$$

The argument which gave (59) can be tweaked to give $g^{\mu\nu}_{\ \ ;\sigma} = 0$, so we have

$$= g^{\mu\nu}\varphi_{,\mu\nu} - \Gamma^{\alpha}\varphi_{\alpha} \tag{66.1}$$

$$\Gamma^{\alpha} \equiv g^{\mu\nu}\Gamma^{\alpha}_{\ \mu\nu} \tag{66.2}$$

It follows easily from (23) that $\Gamma^{\alpha}$ transforms by the rule

$$\bar{\Gamma}^{\mu} = \frac{\partial\bar{x}^{\mu}}{\partial x^{\alpha}}\Gamma^{\alpha} - g^{\alpha\beta}\frac{\partial^{2}\bar{x}^{\mu}}{\partial x^{\alpha}\partial x^{\beta}} \tag{66.3}$$

This last equation shows what (in principle) one must do to arrive in a coordinate system $\bar{x}$ in which $\bar{\Gamma}^{\mu} = 0$, and where, according to (66.1),

<p align="center">generally covariant d'Alembertian = ordinary d'Alembertian</p>

One clearly *gives up general covariance to achieve such simplification*, but some freedom does survive: it follows from (66.3) that if $\Gamma^{\alpha} = 0$ then so does $\bar{\Gamma}^{\mu} = 0$ provided $x \mapsto \bar{x}$ is "harmonic" in the sense that

$$g^{\alpha\beta}\frac{\partial^{2}\bar{x}^{\mu}}{\partial x^{\alpha}\partial x^{\beta}} = 0$$

. . . All of which is precisely mimiced in electrodynamics, where in place of (66.3) one has

$$A'^{\mu} = A^{\mu} + \partial^{\mu}\chi$$

and to achieve the simplicity of Lorentz gauge $\partial_{\mu}A'^{\mu} = 0$ requires that $\chi$ be a solution of

$$\Box\chi = -\partial_{\mu}A^{\mu}$$

Gauge freedom has been sacrificed, but *some freedom does survive*: if $\partial_\mu A^\mu = 0$ then so does $\partial_\mu A'^\mu = 0$ provided $A^\mu \mapsto A'^\mu$ is "harmonic" in the sense that

$$\square \chi = 0$$

How does this development square with the "uniqueness" claim made on the basis of (64.3)? It doesn't! …the lesson being that

> *"degree of freedom counting" is delicate business, and generally unreliable unless side conditions (initial and boundary data) have been specified.*

Another way to characterize the simplification achieved by the introduction of harmonic coordinates: we have

$$\begin{aligned}
\Gamma^\sigma \equiv g^{\mu\nu}\Gamma^\sigma{}_{\mu\nu} &= \tfrac{1}{2}g^{\mu\nu}g^{\sigma\alpha}\{g_{\alpha\mu,\nu} + g_{\alpha\nu,\mu} - g_{\mu\nu,\alpha}\} \\
&= -\tfrac{1}{2}g^{\mu\nu}g_{\alpha\mu}g^{\sigma\alpha}{}_{,\nu} - \tfrac{1}{2}g^{\mu\nu}g_{\alpha\nu}g^{\sigma\alpha}{}_{,\mu} - \tfrac{1}{2}g^{\sigma\alpha}g^{\mu\nu}g_{\mu\nu,\alpha} \\
&= -g^{\sigma\alpha}{}_{,\alpha} - g^{\sigma\alpha}\cdot\tfrac{1}{\sqrt{g}}(\sqrt{g})_{,\alpha} \quad \text{by (54)} \\
&= -\tfrac{1}{\sqrt{g}}(\sqrt{g}g^{\sigma\alpha})_{,\alpha}
\end{aligned}$$

From (56) follow the generally covariant statements

$$\mathfrak{g}^{\mu\alpha}{}_{;\alpha} = 0 \tag{67.1}$$

but the preceding manipulations show that in harmonic coordinate systems we have

$$\mathfrak{g}^{\mu\alpha}{}_{,\alpha} = 0 \tag{67.2}$$

These last conditions are trivially satisfied if $g^{\mu\nu}$ (whence also $\sqrt{g}$) are constant, as would be the case if, in the absence of gravitation, we installed Cartesian coordinates in flat space (or any coordinates harmonically related to them). We are not surprised, therefore, to find that harmonic coordinate systems lend themselves especially well to the description of *weak* gravitational fields, to discussion of the curved spacetime physics *in relation to flat spacetime physics*.[61]

It will be appreciated that, while the Principle of General Covariance served to guide the creation of general relativity, the abandonment of general covariance in favor of some specialized class of coordinate systems (harmonic coordinates, for example) *does no violence to the physics*: it has not the nature

---

[61] See, for example, *MTW*'s Chapter 18. It is curious that the index of the Black Bible contains no entry at "coordinates, harmonic;" what most authors refer to as "installation of harmonic coordinates" Misner, Thorne & Wheeler prefer to call "imposition of the Lorentz gauge condition"—ill-advisedly, in my view, for I think it deceptive to conflate *gauge freedom present in the physics written upon spacetime* with *freedom to recoordinatize spacetime itself*.

of an approximation...though it may facilitate the formulation of useful approximations. The situation here is similar to one encountered in classical mechanics, where imposition of a general covariance requirement leads from Newton's equations to Lagrange's, but the power of the Lagrangian formalism resides in the circumstance that it permits one to work in the special coordinates best suited to the specific problem at hand; it is in the spending that money reveals its value, in its abandonment that general covariance sometimes declares its utility.