

How are theories and laws in physics created?

Most people who have not studied much of physics perhaps think that theories in physics are created *directly* from experiments, i.e., adjusting theories to fit the experimental data that had never been previously explained. I also used to hold a similar view before I started seriously studying physics. However, nothing can be farther from the truth than this view. In this essay, I explain how theories and laws in physics are *actually* created by providing some examples. If you read our earlier essay “The mathematical beauty of physics: simplicity, consistency, and unity,” you would remember how physicists tried to explain an anomaly in moon’s motion by modifying Newtonian theory of gravity only to find that the Newtonian theory could explain it. In other words, in that essay, I gave you an example of a beautiful, successful theory passing a further test, overcoming the danger to be modified into an ugly form. There, I also gave another example to show how complicated looking formulas that describe our nature actually turned out to be simple, beautiful looking formulas. In this essay, I will give you examples of how *new* theories and laws are created.

Let’s begin with the physics laws, which high school students are most familiar with. Newton’s laws. There are three Newton’s laws: Newton’s first law, Newton’s second law, and Newton’s third law. Among these, I will focus on the second, because it highlights my point best. Newton’s second law tells you how much an object is accelerated when a force is exerted on this object. It is given by $a = F/m$, where a is the acceleration of the object, F is the force exerted on the object, and m is the mass of the object. Naively, you may guess that Newton or somebody before him performed many careful experiments to find this relation by exerting objects of various masses with various forces and carefully measuring the accelerations of each object. Then, you may perhaps think, from such experiments, Newton found out that a is proportional to F instead of F^2 and inversely proportional to m instead of \sqrt{m} .

Nothing is farther from the truth. He never performed such experiments. I would say he simply postulated Newton’s second law from the insight he gained from the observations of the world by him and his predecessors, such as Galileo Galilei. Going further, I would say, Newton’s second law *defined* what force is. If we write Newton’s second law slightly differently, we can write $F = ma$. In other words, the magnitude of the force exerted on an object is defined by the product of its mass and its acceleration so obtained. Of course, the word “force” existed perhaps since the beginning of human language, long before the discovery of Newton’s laws, but there is a difference between how physicists use the word “force” and how others use the word “force” in daily lives. In physics, every term has a very narrow and specific meaning, unlike the everyday usage of language, which can be vague, metaphorical and poetic.

Then, how did Newton prove his second law? We will have occasions to talk about it in detail later. At this point, I would just note that he proved his laws mostly from the observed data of orbits of planets and the Moon. Nevertheless, I would say he never *deduced* his second law from them, which is impossible. He started with his postulates,

one of which was Newton's second law, and applied them to explain the observations and proved them by the successful explanation.

Our second example is Einstein's theory of relativity. Through the popular movie *Interstellar*, most of the people who may have never studied Einstein's theory of relativity have heard of one of its key predictions: that for some, time goes more slowly than the others. Actually, such an effect was very well tested and confirmed. For example, in 1971, Joseph Hafele and Richard Keating sent two atomic clocks around the world through commercial airplanes, one eastward and the other westward. Then, they were compared with two atomic clocks that remained on the ground. Remarkably, their time differences agreed with the ones predicted by Einstein's theory of relativity. Of course, such differences are noticeable only by very precise clocks such as atomic clocks. In the case of the experiment in 1971, they were as small as 0.0000003 seconds.

Anyhow, Einstein didn't come up with the theory of relativity to explain such an effect. When Einstein discovered the theory of relativity, such an effect had never been known, because atomic clocks were not available then. How Einstein discovered the theory of relativity is as follows. In the 19th century, the constancy of the speed of light was experimentally "proven" by Michelson and Morley. (I put the quotation mark around proven, because other interpretations for Michelson-Morley experiments were possible, which most of the physicists then believed.) Let me briefly explain what the constancy of the speed of light means. Naively, you might think that the speed of light would be observed as 10,000 km/s if you followed a 300,000 km/s beam of light at a speed of 290,000 km/s. However, it turns out that the speed of light still appears to be 300,000 km/s no matter how fast you follow the light. This is what the constancy of the speed of light means. In 1905, Einstein discovered special relativity using the constancy of the speed of light as one of its key basis. Among many other things, he also calculated how much time needs to slow down for moving objects to make the speed of light constant. In 1915, Einstein discovered general relativity from the Riemannian geometry, which was discovered by the German mathematician Riemann more than half a century ago. As the first test of general relativity, he explained the orbit of Mercury, which could not be explained from the existing theory, i.e., Newtonian gravity. Then, he predicted how much the trajectory of starlight would be bent by the gravity of the Sun, which was confirmed in 1919.

So, Einstein discovered both theories of relativity from principles such as the constancy of the speed of light and that our nature must be described by elegant mathematics such as the Riemannian geometry. Here, I want to emphasize that Einstein could not "tweak" his theory to explain the orbit of Mercury.¹ He could test that his theory was correct, only after he had done all the calculations for the orbit of Mercury. Let me explain why he could not tweak his theory. It is because it is so simple that there is no room for tweaking. (Or, more precisely, I should say the orbit of Mercury can be explained without unnecessary tweaking.) Einstein equation is given by

$$G_{ab} = \frac{8\pi G}{c^4} T_{ab} \quad (1)$$

Here, G_{ab} and T_{ab} are variables, while G and c are constants. G_{ab} is a kind of curvature of spacetime, which one can calculate from the Riemannian geometry, which has been known since the middle of the 19th century. T_{ab} is "energy-momentum tensor," a generalization of mass, energy, and momentum, which has been known long before

¹Those of you who are familiar with the history of astronomy may argue that Newton "tweaked" his theory to explain the precession of equinoxes. No, he didn't. He tweaked his data, not his theory.

Einstein. Could Einstein have changed the definition of G_{ab} or T_{ab} to explain the orbit of Mercury? No. They have already been known.² Could he change G and c to explain the orbit of Mercury? No. They have been measured very precisely. Could he change the coefficient 8π into 9π ? No, then Einstein equation would not be reduced to Newton's law of universal gravitation, failing in all other predictions made by Newton's law of universal gravitation.

Well, let me be honest about the history. According to his own account, Einstein first tried to force his equations to fit with the orbit of Mercury. However, all of his such efforts failed. Only when he tried to make his equation as simple and natural as possible could he come up with the correct equation (1).

Of course, I would like to remark that the actual calculations of the orbit of Mercury and the light trajectory near the Sun are much more complicated than (1). Nevertheless, one must remember that their origin was as simple as (1), and there is no leeway in the actual calculations once the form of the equation (1) is written this way.

Let me give you the third example. In the 1970s, it was discovered that the stars rotate around the center of galaxies at speeds much faster than the Newtonian gravity (or Einstein's general relativity) predicts. Thus, some propose that there is dark matter in the galaxy that, due to its mass, affects the speeds of stars, yet remains unobserved and therefore neglected during mass measurements and the speed calculations based on these measurements, which is now known as "missing mass problem." On the other hand, in the 1980s, the Israeli physicist, Mordehai Milgrom proposed the MOND (Modified Newtonian Dynamics) to explain this phenomenon.

First, he tried to explain the Tully-Fisher relation, which says that the outermost stars in a galaxy rotate with speed v proportional to the square root of the square root of the total mass of the galaxy M . Mathematically, this can be expressed as

$$v = b\sqrt{\sqrt{M}} \quad (2)$$

for some constant b that can be determined from observations. He noticed that this relation can be explained if he modified Newton's second law as follows: When a is as small as, or in the order of a_M , which is $1.2 \times 10^{-10} \text{m/s}^2$, instead of Newton's law $F = ma$, we have $F = ma^2/a_M$. Here, a_M , the Milgrom's constant, is calculated from the observed value of b . (This calculation is not difficult at all. All you need is high school physics. We will show the calculation in a later article.)

Of course, we cannot have $F = ma^2/a_M$ when a is much bigger than a_M as well, because $F = ma$ has been already well tested in such cases by Newtonian mechanics for more than 300 years. Therefore, when a is much bigger than a_M , we must have $F = ma$ as original Newton's law.

Summarizing, Milgrom's proposal is as follows:

$$F = \begin{cases} ma^2/a_M, & \text{if } a \ll a_M \\ ma, & \text{if } a \gg a_M. \end{cases} \quad (3)$$

So, the Tully-Fisher relation does not guide us what the modified Newton's law would look like when a is neither much larger than a_M nor much smaller than a_M . What should

²Strictly speaking, he could have put R_{ab} instead of G_{ab} on the left-hand side of his equation, and he did consider such a form of equation, but it is ruled out by the physical condition that energy is conserved, so he corrected it to G_{ab} . Therefore, I would say (1) has the second simplest form that can be constructed from the Riemannian geometry. Nevertheless, if we follow David Hilbert's formulation of general relativity, it has the single simplest form the Riemannian geometry would allow.

we do? It seems that we need to use a lot of data in this range, including the ones for rotation speeds of stars that are not outermost, but inside the galaxy. Maybe, by doing so, we can deduce what F would look like for such a region of a . However, the problem is that there are many formulas that satisfy the property (3). Here are four examples of such formulas.

$$F = \frac{ma}{1 + a_M/a} \quad (4)$$

$$F = \frac{ma}{\sqrt{1 + (a_M/a)^2}} \quad (5)$$

$$a = \frac{F/m}{1 - e^{-\sqrt{F/(ma_M)}}} \quad (6)$$

$$F = \begin{cases} ma^2/a_M, & \text{if } a \leq a_M \\ ma, & \text{if } a \geq a_M. \end{cases} \quad (7)$$

Let me explain. In the first example (4), when a is much smaller than a_M , a_M/a is a number much bigger than 1. Thus, the term 1 in the denominator is negligible compared to the term a_M/a in the denominator. Thus, in such a case, we have

$$F \approx \frac{ma}{a_M/a} = ma^2/a_M \quad (8)$$

as in (3). On the other hand, when a is much bigger than a_M , a_M/a is a number much smaller than 1. Thus, the term a_M/a in the denominator is negligible compared to the term 1 in the denominator. Thus, in such a case, we have simply $F = ma$ as in (3). Similarly, the other three examples can be easily shown to satisfy (3).

However, one formula will not particularly fit the observational data better than all the others; many other formulas that satisfy the property (3) will fit the data equally well or equally bad. Perhaps, one formula could fit the data *slightly* better or *slightly* worse than others, but no formula can ever perfectly fit the data, as the data themselves have quite big measurement errors, as big as around 30%. In other words, even if the correct theory existed in the MOND framework, it would *not* be able to *perfectly* fit the data. Then, which one should we choose? Data *alone* cannot be the judge to determine which of the MOND candidates is the correct one.

The problem with the MOND, I believe, is that Milgrom “tweaked” (i.e., modified) Newton’s law to fit the data. In my opinion, such an approach (i.e., tweaking the theory to explain data) is never promising. Only theories derived from principles are promising. For example, Verlinde gravity, partially based on the holographic principle, explains the galaxy rotation curve, as Prof. Ho Seong Hwang, Prof. Jong-Chul Park and I showed, even though it cannot fit the data perfectly well due to the big measurement errors as mentioned. Actually, despite being the “correct” theory, it fits the data no better than MOND does.

Anyhow, as was the case with the orbit of Mercury and Einstein equation, the actual galaxy rotation curve can be compared with the Verlinde gravity prediction only after all the relevant calculations are performed from the original equations that are based on Verlinde’s simple principle, i.e., that there is a volume law contribution to entropy in addition to the well-known area law contribution. The actual calculations are quite tedious, and the resulting formula is much more complicated than each of the four earlier examples of MOND and any of their other variants.

The correct explanation of galaxy rotation curve was never meant to be discovered by guessing and fitting method such as the MOND ones; in the earlier four examples, the acceleration a is determined effectively only by one variable, i.e., F/m , but in Verlinde gravity more variables, such as how much gravity changes as you change your position, how much this change changes as you change your position and the mass density of galaxy are needed as you can see here.³

$$g = \sqrt{g_B^2 + \frac{a_0}{6} \left(2g_B - \vec{n} \cdot \nabla \left(\frac{2g_B^2}{4\pi G\rho_B + \vec{n} \cdot \nabla g_B} \right) \right)} \quad (9)$$

When there are many variables and the formula looks complicated as this case, it is virtually impossible to guess this correct formula by controlling, analyzing and tuning them systematically to compare the calculation with the data.

In the frontier of fundamental physics, guessing formulas to fit the data doesn't help, as long as the correct formulas, even though they may seem more complicated, are only the ones that are derived from simple principles. Those of you who read my earlier essay "The mathematical beauty of physics: simplicity, consistency, and unity" will understand what I mean. The MOND equations seem look "simple" to untrained laymen, but God chose Verlinde gravity instead as the correct theory, because it is based on a simple principle.

Of course, this simple principle was not an old principle that had been applied so far in theory of gravity. Actually, Verlinde pointed out why he had come up with a theory of gravity with a whole new idea: Einstein didn't come up with general relativity by modifying Newton's law of universal gravitation. He had to come up with a whole new idea to develop general relativity. If you learn general relativity you will understand what Verlinde meant. As we similarly mentioned, Einstein could not guess the Einstein field equation (1) from the fact that it becomes Newton's law of universal gravitation, when gravity is weak and things are moving much slower than speed of light. He first guessed (1), and by a lengthy calculation showed that (1) becomes Newton's law of universal gravitation in the appropriate limit. Thus, Verlinde said that he was convinced that a whole new approach and idea was necessary to explain the missing mass problem from this lesson of history.

Of course, that only a whole new approach can be successful in solving the missing mass problem never means that Milgrom's MOND was in vain. Prof. Verlinde first showed that his Verlinde gravity can explain the Tully-Fisher relation, and I strongly suspect that he would never have been able to do so if he hadn't known Milgrom's proposal. Moreover, if Prof. Verlinde hadn't known that a_M is in the same order as the Hubble's constant, which tells the expansion speed of our Universe, he would have never been able to suspect the connection between the two, which resulted in Verlinde gravity. (I highly suspect that he knew it.)

Anyhow, I want to emphasize that there is nothing one can "tweak" in the Verlinde gravity calculation, unlike in the MOND. The statement that there is a volume law contribution to entropy is a statement without ambiguity. One cannot make even two different models in which the volume law contributes to entropy differently, once one understands the reasoning behind Verlinde's argument that there should be volume law contribution to entropy. Such a simple statement fixes even the details of the theory.

³Precisely speaking, I am using a different notation here than the four earlier MOND expressions. If I use this different notation they become more simple. For example, (6) becomes $g = \frac{g_B}{1 - e^{-\sqrt{g_B/a_M}}}$.

To argue against what I have just said, those of you who know what happened with my research on Verlinde gravity may point out that I used a different formula than Prof. Verlinde presented in his paper on Verlinde gravity. Namely, I used

$$g = \sqrt{g_B^2 + g_D^2} \quad (10)$$

while Prof. Verlinde used

$$g = g_B + g_D \quad (11)$$

That I used a different formula is true, but it's not because I used a different "model" than Prof. Verlinde's "model," both within the framework of Verlinde gravity. Logic and mathematics must determine which one of the two is the correct formula, and which one of the two is the wrong formula, if one is faithful to the principle Prof. Verlinde used to construct Verlinde gravity. Verlinde gravity, with the correct formula (10) replaced by the wrong formula (11), while not touching all the other formulas, will necessarily suffer from mathematical inconsistencies. (Please read our earlier essay "The mathematical beauty of physics: simplicity, consistency and unity" to learn about what inconsistency (or consistency) here means.)

Back to our original point, Einstein was serious when he said that there is no way to come up with a theory from experimental data. Instead, one has to rely on "principles" which are usually not *directly* deduced from experiments.⁴ For example, the holographic principle was discovered and developed by theoretical particle physicists, mainly, by string theorists. (Another crucial reason to refute that string theory is useless.) In other words, it is a principle discovered by theoretical calculations, which in turn are based on wisdom physicists learned from long theoretical research; it was never discovered by specific observations or particular experiments. Of course, before Verlinde's application of the holographic principle to his theory of gravity, the holographic principle didn't seem to be related to galaxy rotation curves, which are measured by astronomers, who do not understand string theory calculations any more than laymen.

There are countless examples like this in the history of physics. In our earlier essay "The mathematical beauty of physics: simplicity, consistency, and unity," we briefly mentioned that Yang-Mills theory had been invented without any experimental input. I would say, Yang-Mills theory is a mathematical generalization of Maxwell theory, which describes electromagnetism. However, when Yang and Mills developed their theory in 1954, they found out that their theory predicted yet unknown massless particle that has never been experimentally detected before. When Yang gave a talk on their new theory at Institute for Advanced Study (IAS), Pauli called it "not even wrong." Pauli had already been thinking on the same lines as Yang and Mills earlier, but abandoned it as he also found out that it predicted a new massless particle. When Pauli raised this problem to Yang, he could not give a satisfactory answer. When Pauli attacked him again few minutes later claiming that it was no excuse, Oppenheimer, the head of IAS had to tell him that he should let Yang proceed. In the 1960s and in the 1970s, it was found out that weak force and strong force, the two forces out of the four forces that describe our nature, were described by Yang-Mills theory. How about the massless particle? The mechanism that the massless particle that mediates the weak force gains

⁴An exception would be the constancy of the speed of light, which had been experimentally confirmed before the discovery of special relativity. But, still, I believe that many did suspect the constancy of the speed of light after Maxwell's correct theoretical calculation of the speed of light, before the Michelson-Morley experiment.

mass was proposed, and proven correct by experiments. Regarding the strong force, the massless particle was discovered in the late 1970s.

Now, let me respond to two criticisms some people without deep physics education may have on physics.

First, in my earlier article “The mathematical beauty of physics: simplicity, consistency, and unity,” I explained that the Coulomb’s law and Newton’s law of universal gravitation obey very simple inverse square law, which can be recasted into simple Gauss’s law, shows the beauty of physics. However, some people without physics education then may point out that we now know that actual gravitation doesn’t obey the simple inverse square law, as Newton’s law of universal gravitation cannot explain Mercury’s orbit or galaxy rotation curve; physicists had to invent new theories that seem to be more complicated than the simple inverse square law. To this criticism, I would point out that Einstein equation and Verlinde’s ideas are more beautiful and simpler than Newton’s law of universal gravitation. It is like our example of Maxwell’s equations in my earlier article on the mathematical beauty of physics; if we know more difficult math, we can express the same equation simpler (in case of Maxwell’s equations) or come up with a slightly different, but “more correct” equation that’s simpler (in case of Einstein equation). In *A Mathematical Journey*, mathematician Stanley Gudder wrote, “The essence of mathematics is not to make simple things complicated, but to make complicated things simple.” The more we study mathematics, we can understand our universe from the more beautiful, simpler mathematical equations.

Second, I will respond to a criticism relativists have about physics. Relativism is the belief that there is no absolute truth. (It has nothing to do with the word “relativity” in “Einstein’s theory of relativity.”) They may think that Newton was proven wrong by Einstein, Einstein was proven wrong by Verlinde, and so on, which shows that there is no absolute truth. I already explained what is wrong with such an argument in our earlier essays, “The mathematical beauty of physics: simplicity, consistency, and unity” and “Did Einstein really prove that Newton was wrong?” Still, some relativists could point out that Einstein’s theory of general relativity was based on the equivalence principle, which was proven wrong by Verlinde gravity. Likewise, they may argue, the theories we have now are not on good footing, because they are based on principles or assumptions that may be proven wrong in the future.

Let me put it this way. We now know that time flows at different rates from person to person due to the relativistic effect, but nobody knew this, before Einstein came up with the theory of relativity. As it is impossible to notice this from our everyday lives, everybody thought that time flows at the same rate, and that was what Newton naturally assumed when he proposed Newton’s laws and the law of universal gravitation. Therefore, in a sense, this Newtonian mechanics is “wrong,” because it is based on a wrong assumption.⁵

However, does this mean that Newton should not have started from this wrong assumption to construct Newtonian mechanics? No. The experiments that confirmed the constancy of the speed of light, without which Einstein’s theory of relativity would not have been discovered, could not be performed in Newton’s time due to technological difficulties. Therefore, he could not have started from the correct assumption, i.e., the constancy of the speed of light, and instead had to rely on a wrong assumption, the absoluteness of time, that time flows at the same rate for all observers, which is nevertheless “almost correct.” He had no choice, but to accept a wrong assumption.

⁵Of course, the word “wrong” here must be carefully interpreted. Please read our earlier essay “Did Einstein really prove that Newton was wrong?” to understand how you should interpret this word.

I believe that principles that we assume and rely on in physics are often this kind of nature. We can rely on them because we know that they are almost correct, and we must rely on them because we will get nowhere without relying on them. If Newton hadn't assumed that the time flows at the same rate for everyone, at how different rates should he have assumed that the time flows? Surely, he wouldn't be able to guess the correct differing rates of time flow, not knowing Einstein's special theory of relativity. Moreover, experiments and observations at his day can be analyzed well assuming the absoluteness of time, instead of assuming that time flows at different rates, as their experimental apparatuses were not sensitive enough. In other words, for Newton, considering the possibility that time flows at different rates would *not* have helped his research *at all*.

Without Newtonian mechanics, Einstein's theory of relativity, which contradicts Newtonian mechanics and the very assumption on which it was based, would never have been discovered. Thus, only by supposing a "wrong," but "almost correct" assumption can physics be developed further, and only by doing so can this original, wrong assumption be proven wrong. That's how physics develops.

The last comment. For those of you who are interested in Einstein's model for constructing a scientific theory, please see [1,2]. I first encountered this model from a Japanese book that was translated into Korean. (Japanese people write lots of good science books for the general public. No wonder many Japanese students grow up to be great scientists.) I was only ten years at the moment and knew nothing beyond very basic physics, so I could not understand his model. Twenty years later, I encountered this model again from [1] and found myself agreeing with this model.

References

- [1] p 51-57. DONGEN, JEROEN VAN. "Einstein's Unification.": Cambridge Univ Press, 2010.
- [2] HOLTON, GERALD. "Constructing a Theory: Einstein's Model." The American Scholar 48, no. 3 (1979): 309-40. <http://www.jstor.org/stable/41210527>.