# The Brachistochrone and the Calculus of Variations

Aditya Mittal

May 2022

## 1  Introduction

Often when we are dealing with finding the maximum or minimum of a function, it is not too difficult using its derivative (or its higher dimensional equivalents). If asked to

$$\max(f(x) : x = [-5, 5])$$

we can readily do so algebraically using *Fermat's theorem*: the local extrema of a function occur when its derivative is 0.
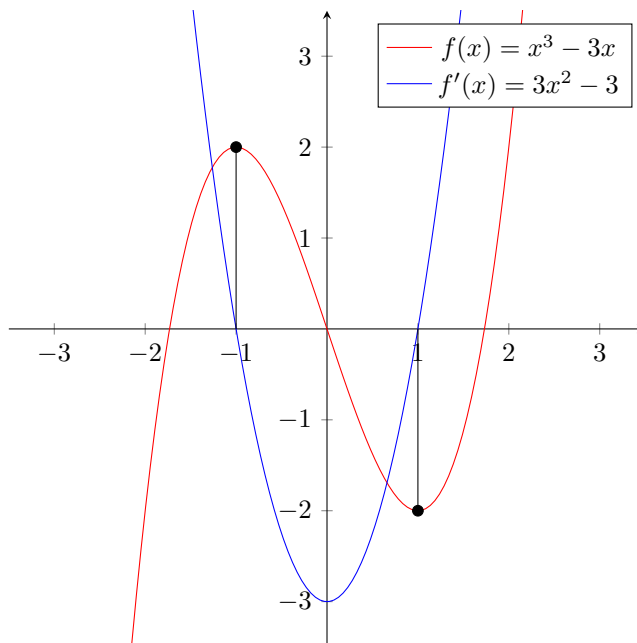


Figure 1: Graphical representation of using the derivative to find the local extrema of a function

This makes some intuitive sense, since the derivative of a function measures the slope, or rate at which a function changes locally. So if the derivative is positive or negative, that implies that the function is increasing or decreasing. If a function is still increasing, then it can't be at a local maximum, and if it is decreasing, then it can't be at a local minimum. Thus, our *critical points* occur when the derivative is 0 (where the graph's tangent line is flat). This strategy works since we essentially are doing an efficient guess-and-check method: we are letting our variables (in this case, only $x$, but could be more), slide over their domain until they satisfy a condition we observe to be true (derivative equalling 0).

In general, for any function $f(x, y, z, \dots)$, a small change in the function $df$ is

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

At an extreme point $(x_0, y_0, z_0, \dots)$, $df = 0$ since any infinitesimal step in any direction (i.e. $dx$, $dy$, $dz$, etc.) should not increase or decrease $f$. Thus, local maximums and minimums occur when

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} = \frac{\partial f}{\partial z} = \dots = 0$$

But sometimes, certain max-/minimization problems prove this to be an infeasible strategy.

## 2  Bernoulli's Challenge to the Mathematical World

Our story begins in 1696, in Switzerland. Johann Bernoulli, well-respected mathematician and feeding on his ego, issues the following to the rest of the mathematical world (really, just Europe):

> I, Johann Bernoulli, address the most brilliant mathematicians in the world. Nothing is more attractive to intelligent people than an honest, challenging problem, whose possible solution will bestow fame and remain as a lasting monument. Following the example set by Pascal, Fermat, etc., I hope to gain the gratitude of the whole scientific community by placing before the finest mathematicians of our time a problem which will test their methods and the strength of their intellect. If someone communicates to me the solution of the proposed problem, I shall publicly declare him worthy of praise.

Basically undermining every respected scholar at the time, Bernoulli provoked many—especially his brother Jakob Bernoulli—to prove his worth. Insecurities aside, what even was this problem that many bothered concerning themselves with?

**Problem.** *Given two points A and B in a vertical plane, what is the curve traced out by a point acted on only by gravity, which starts at A and reaches B in the shortest time.*
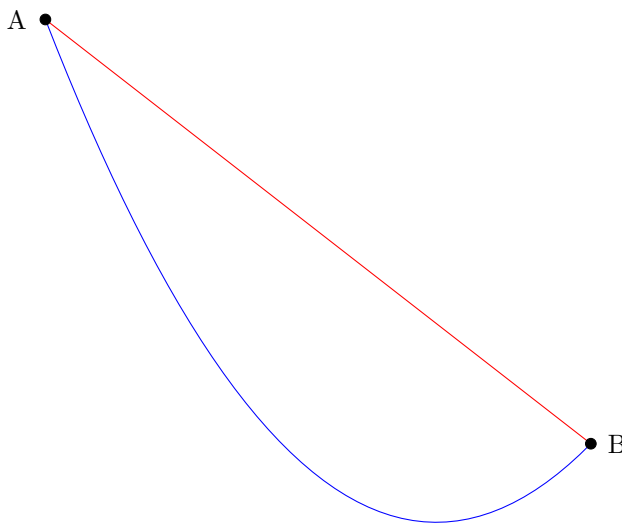


Figure 2: While a straight line (red) might be the *shortest distance*, the blue path sacrifices distance with a steeper path to build speed. Which will be faster?

This looks like a minimization problem! Let's try setting it up, and hopefully we can use our derivative method like before. What we are looking to find is known as the *brachistochrone*, deriving literally from the Greek word for "shortest time".

# 3   First Attempt

We want to find the curve that minimizes the time of descent for ball rolling between point A and point B. Let's call that curve $h(x)$ (with the positive $y$-axis as down, and A at the origin). So, as our ball rolls and is at a point $x$, it will be at a height $h(x)$ below A.

## 3.1   Physics Predicaments

Before we can tackle this problem, we need to cover some concepts from physics. Most importantly, the idea of *conservation of energy*. As our ball falls from a certain height, as we intuitively know, it will gain speed, but how much speed? Conservation of energy tells us that the ball's initial potential energy (in this case, height above a point) is equivalent to $U_g = mgh$, where $m$ is mass of the ball and $g$ is the acceleration due to gravity. This is then converted to kinetic energy, which known to be $K = \frac{1}{2}mv^2$, where $v$ is velocity. Since the total energy of a system is constant, the total amount of potential energy lost must be converted into kinetic energy, allowing us to equate $mgh = \frac{1}{2}mv^2$, and thus

$$\boxed{v = \sqrt{2gh}}$$

The fact mass divides out is a good sign, as then it does not matter the mass of our ball when rolling down our optimal curve.

Another key idea we'll utilize comes from kinematics, and that is

$$\text{velocity} = \frac{\text{distance}}{\text{time}} = \frac{ds}{dt}$$

This will be the key equation that guides our the final expression we will try to minimize.

---

Now that we have a relationship between time, velocity, and distance, we can construct an expression to minimize total time travelled. We already have an expression for velocity: $v = \sqrt{2gh(x)}$. For a small step along our curve $h(x)$, we can approximate it as a line, considering the small step in the $x$ direction, $dx$, and the small step in the $y$ direction, $dy$.

$$ds = \sqrt{(dx)^2 + (dy)^2} = \sqrt{1 + \left(\frac{dy}{dx}\right)^2}\ dx = \sqrt{1 + (h'(x))^2}\ dx$$

Putting this all together, we have that

$$v = \frac{ds}{dt} \Rightarrow \sqrt{2gh(x)} = \sqrt{1 + (h'(x))^2}\ \frac{dx}{dt}$$

This is a separable differential equation.

$$dt = \frac{\sqrt{1 + (h'(x))^2}}{\sqrt{2gh(x)}}\ dx$$

Integrating over the total horizontal distance travelled nets us that the total time of descent along our curve, $h(x)$ due to gravity is

$$T = \int_0^{x_B} \frac{\sqrt{1 + (h'(x))^2}}{\sqrt{2gh(x)}}\ dx$$

...so how do we take the derivative of this?

# 4 The Calculus of Variations

In our previous minimization problems, we have a function $f(x)$ and are solving for an $x$ value. Here, our value of time $T$ is not a function of $x$, but rather is a function *of another function, $h(x)$, which we are trying to find; $T$ inputs a function, and outputs a number. We therefore don't call $T$ a function, but rather a *functional* of $h(x)$.

Functionals are key to "the calculus of variations", as we'll see in a bit, instead of considering small nudges to inputs (i.e. numbers) of functions to find extrema, we instead consider, or vary, values in functions to find derivative-esque equivalents.

We can rewrite our integrand as a "multivariable" function, say $L$:

$$T = \int_0^{x_B} L(x, h, h') \ dx$$

To find our value-minimizing function, we can use a well-known identity in the calculus of variations: the *Euler-Lagrange Equation*

## 4.1 The Euler-Lagrange Equation

Let's suppose $h(x)$ is the function that minimizes (in general, can maximize as well). So we have the following integral (with slightly modified bounds for generality).

$$I = \int_a^b L(x, h, h') \ dx$$

We can now take an approach similar to our previous, simple functions (not functionals). Instead of considering small steps in our variables (i.e. $x$, $y$, $z$, etc.), let's consider a small change in our desired function. We'll define a new function $F(x)$ to do so. Let

$$F(x) = h(x) + \epsilon \eta(x)$$

The only important aspect of $F(x)$ is that it has the same boundaries as $f(x)$. In other words, $F(a) = f(a)$ and $F(b) = f(b)$, or that $\eta(a) = \eta(b) = 0$.
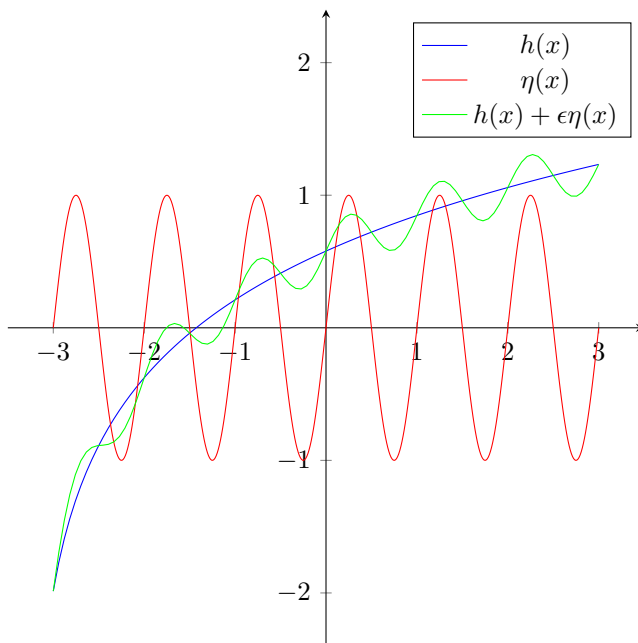
Figure 3: What it looks like to "vary" over a function $h(x)$. The green function is *almost* the blue function, but not quite. Remember, $\eta(x)$ can be whatever you want it to be, so long as the boundaries are the same.

The scalar $\epsilon$ is only there to remind us that $\eta(x)$ is truly a nudge to our $h(x)$; $\epsilon$ is really, really small, and will tend to 0 as we consider our candidate function $h(x)$. The key idea here, is that if $h(x)$ is a max-/minimizing function of our functional $I$, then $F(x)$, which is our slight nudge, should be a *worse* function, that either decreases/increases $I$ away from a maximum/minimum. So if we consider the same integral with $F(x)$,

$$I_\epsilon = \int_a^b L(x, F, F') \, dx = \int_a^b L_\epsilon \, dx$$

We can then take the derivative of $I$ with respect to $\epsilon$, but not just any derivative: the *total derivative*.

$$\frac{\mathrm{d}I_\epsilon}{\mathrm{d}\epsilon} = \frac{\mathrm{d}}{\mathrm{d}\epsilon} \int_a^b L_\epsilon \, dx = \int_a^b \frac{\mathrm{d}L_\epsilon}{\mathrm{d}\epsilon} \, dx$$

Since changing $\epsilon$ affects multiple variables of $L_\epsilon$ once, we can't make the assumption that our other variables remain constant in the way partial derivatives employ; we basically need to remember the chain rule for partial derivatives.

$$\begin{aligned}
\frac{\mathrm{d}L_\epsilon}{\mathrm{d}\epsilon} &= \frac{\partial L_\epsilon}{\partial \epsilon} + \frac{\partial L_\epsilon}{\partial x}\frac{dx}{d\epsilon} + \frac{\partial L_\epsilon}{\partial F}\frac{dF}{d\epsilon} + \frac{\partial L_\epsilon}{\partial F'}\frac{dF'}{d\epsilon} \\
&= \frac{\partial L_\epsilon}{\partial F}\frac{dF}{d\epsilon} + \frac{\partial L_\epsilon}{\partial F'}\frac{dF'}{d\epsilon} \\
&= \frac{\partial L_\epsilon}{\partial F}\eta(x) + \frac{\partial L_\epsilon}{\partial F'}\eta'(x)
\end{aligned}$$

Again, the above was all chain rule (note $\frac{dx}{d\epsilon} = 0$ since $x$ does not depend on $\epsilon$). At the end of all this, we have

$$\frac{\mathrm{d}I_\epsilon}{\mathrm{d}\epsilon} = \int_a^b \frac{\partial L_\epsilon}{\partial F}\eta(x) + \frac{\partial L_\epsilon}{\partial F'}\eta'(x) \, dx$$

As we let $\epsilon \to 0$, $F(x) = f(x)$, $L_\epsilon = L(x, h, h')$, and most importantly, our integral $I_\epsilon = I$ and has its extremum value (i.e. maximum or minimum). So, we can therefore conclude that at $\epsilon = 0$,

$$\lim_{\epsilon \to 0} \frac{\mathrm{d}I_\epsilon}{\mathrm{d}\epsilon} = \int_a^b \frac{\partial L}{\partial h}\eta(x) + \frac{\partial L}{\partial h'}\eta'(x) \, dx = 0$$

This should look familiar, as it is the equivalent of our very first max-/minimization problem: we found the derivative and set it equal to 0 (in this case, the 0 came from the fact that at $\epsilon = 0$, $I_\epsilon$ doesn't depend on $\epsilon$ anymore, therefore $\frac{\mathrm{d}I_\epsilon}{\mathrm{d}\epsilon} = 0$). We can do integration-by-parts on the right-hand side of the integrand.

$$\int_a^b \frac{\partial L}{\partial h'}\eta'(x) \, dx = \eta(x)\frac{\partial L}{\partial h'}\Big|_a^b - \int_a^b \eta(x)\frac{\mathrm{d}}{\mathrm{d}x}\frac{\partial L}{\partial h'} \, dx$$

Combining this with our original integral, we get that

$$\lim_{\epsilon \to 0} \frac{\mathrm{d}I_\epsilon}{\mathrm{d}\epsilon} = \int_a^b \left[\frac{\partial L}{\partial h} - \frac{\mathrm{d}}{\mathrm{d}x}\frac{\partial L}{\partial h'}\right]\eta(x) \, dx + \eta(x)\frac{\partial L}{\partial h'}\Big|_a^b = 0$$

Recall our boundary condition that specifies $\eta(a) = \eta(b) = 0$:

$$\lim_{\epsilon \to 0} \frac{\mathrm{d}I_\epsilon}{\mathrm{d}\epsilon} = \int_a^b \left[\frac{\partial L}{\partial h} - \frac{\mathrm{d}}{\mathrm{d}x}\frac{\partial L}{\partial h'}\right]\eta(x) \, dx = 0$$

The only way that last equality holds for *any* function $\eta(x)$ is if $\frac{\partial L}{\partial h} - \frac{\mathrm{d}}{\mathrm{d}x}\frac{\partial L}{\partial h'} = 0$. This is intuitive in the sense that if you have any strictly non-zero, continuous function $\eta(x)$, it will have a strictly non-zero integral

too, and the only way it can become 0 is if the function is scaled by 0. Here's a more in-depth sketch of the *fundamental lemma of the calculus of variations* (yes, it is that important).

Finally, for our integral $I$, the function $h(x)$ that max-/minimizes $I$ only does so if it satisfies

$$\frac{\partial L}{\partial h} - \frac{\mathrm{d}}{\mathrm{d}x}\frac{\partial L}{\partial h'} = 0$$

What we've just derived is the *Euler-Lagrange Equation*, named after—two of the most famous mathematicians who ever lived—who discovered it, Leonard Euler and Joseph-Louis Lagrange.

---

To go back to our original goal, we are looking to minimize

$$T = \int_0^{x_B} L(x, h, h') \ dx = \int_0^{x_B} \frac{\sqrt{1 + (h'(x))^2}}{\sqrt{2gh(x)}} \ dx$$

Notice how our $L(x, h, h')$ is only defined in terms of $h(x)$ and $h'(x)$; there is no explicit relation between $L$ and $x$. This allows us to use a fancy restatement of the Euler-Lagrange Equation to our advantage.

## 4.2   Beltrami's Identity

First, note that

$$\frac{\mathrm{d}L}{\mathrm{d}x} = \frac{\partial L}{\partial x} + \frac{\partial L}{\partial h}\frac{dh}{dx} + \frac{\partial L}{\partial h'}\frac{dh'}{dx}$$

Again, can't forget the chain rule. Since our functional $L$ does not depend on $x$, $\frac{\partial L}{\partial x} = 0$. This combined with the definition of the derivatives of $h(x)$ and $h'(x)$, we therefore obtain

$$\frac{\mathrm{d}L}{\mathrm{d}x} = \frac{\partial L}{\partial h}h' + \frac{\partial L}{\partial h'}h''$$

Combining this with the Euler-Lagrange Equation gives us

$$\frac{\mathrm{d}L}{\mathrm{d}x} = h'\frac{\mathrm{d}}{\mathrm{d}x}\frac{\partial L}{\partial h'} + \frac{\partial L}{\partial h'}h''$$

The right-hand side is the equivalent result of an application of the product rule of derivatives!

$$\frac{\mathrm{d}L}{\mathrm{d}x} = \frac{\mathrm{d}}{\mathrm{d}x}\left(h'\frac{\partial L}{\partial h'}\right)$$

Integrating both sides and rearranging. . .

$$L - h'\frac{\partial L}{\partial h'} = C$$

(Can't forget the $+C$) This special case of the Euler-Lagrange Equation is known as *Beltrami's identity*.

---

Plugging in the respective functions to Beltrami's identity, we get

$$\frac{\sqrt{1 + (h'(x))^2}}{\sqrt{2gh(x)}} - \frac{(h'(x))^2}{\sqrt{2gh(x)(1 + (h'(x))^2)}} = \frac{1}{\sqrt{2gh(x)(1 + (h'(x))^2)}} = C$$

Squaring both sides and rearranging,

$$\frac{1}{\sqrt{2gh(x)(1 + (h'(x))^2)}} = C \Rightarrow h(x)(1 + (h'(x))^2) = C$$

We then get the final differential equation

$$h'(x) = \frac{dh}{dx} = \sqrt{\frac{C}{h(x)} - 1} = \sqrt{\frac{C - h(x)}{h(x)}}$$

Solving for $x$ results in

$$x = \int \sqrt{\frac{h(x)}{C - h(x)}} \ dh$$

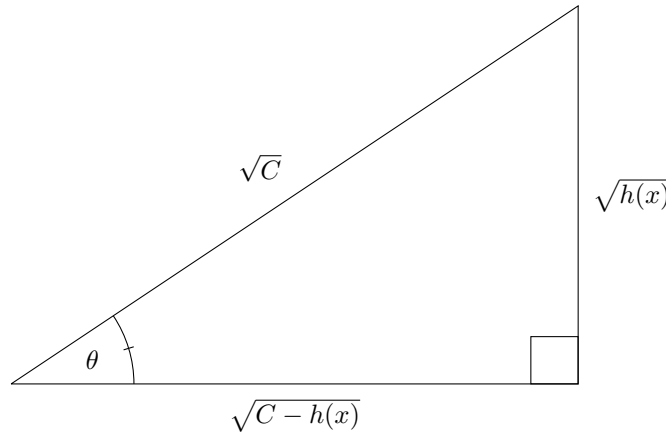This looks like something we could perform a trigonometric substitution on!



Figure 4: Visualization of our trigonometric substitutions.

Using the above triangle as a reference point, we can set

$$\tan(\theta) = \sqrt{\frac{h(x)}{C - h(x)}}$$

Moreover, we can also conclude that

$$\sin(\theta) = \sqrt{\frac{h(x)}{C}} \Rightarrow dh = 2C \sin(\theta) \cos(\theta) \ d\theta$$

Therefore, our new integral is

$$x = \int \sqrt{\frac{h(x)}{C - h(x)}} \ dh = \int \tan(\theta) \cdot 2C \sin(\theta) \cos(\theta) \ d\theta = 2C \int \sin^2(\theta) \ d\theta$$

Using the identities $\sin^2(\theta) + \cos^2(\theta) = 1$ and that $\cos(2\theta) = \cos^2(\theta) - \sin^2(\theta)$, we can rewrite this as

$$x = 2C \int \sin^2(\theta) \ d\theta = C \int 1 - \cos(2\theta) \ d\theta = C\theta - \frac{1}{2}\sin(2\theta) + D$$

Since we said our curve starts at the origin $(0, 0)$, that implies that $D = 0$. Also, from our trigonometric substitution, we also get that $y = h(x) = C \sin^2(\theta)$, netting us the final parametric equations of our curve:

$$x = C\theta - \frac{C}{2}\sin(2\theta)$$

$$y = \frac{C}{2} - \frac{C}{2}\cos(2\theta)$$

Let $\theta \to \frac{\theta}{2}$:

$$x = C(\frac{\theta}{2} - \frac{1}{2}\sin(\theta))$$

$$y = C(\frac{1}{2} - \frac{1}{2}\cos(\theta))$$

These equations famously describe the cycloid! The value $C$ describes the radius of the circle that generates the cycloid, and you set it as needed to have the curve pass through the final point $B$ at $(x_b, y_b)$.
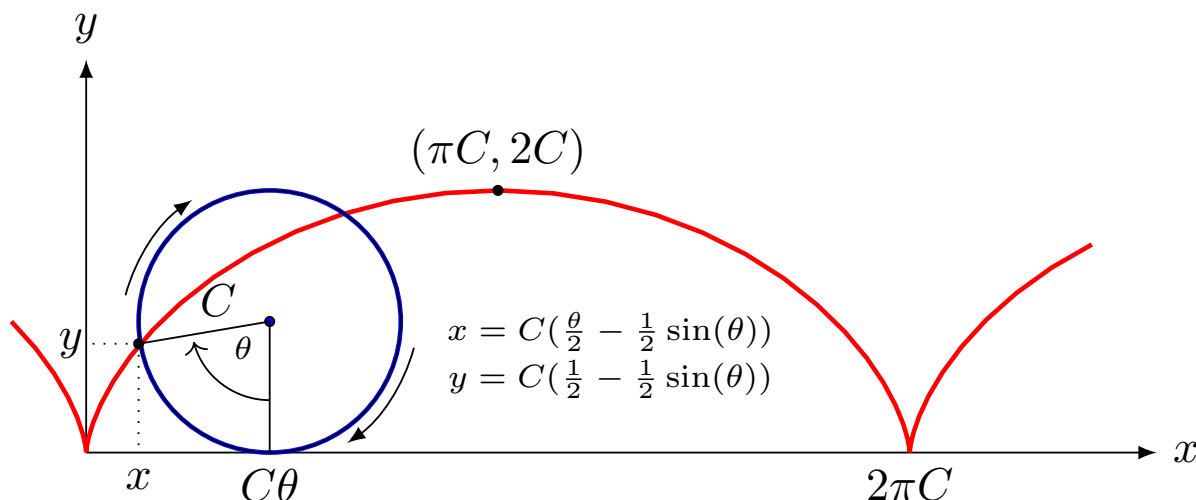


Figure 5: A cycloid is a type of roulette curve generated by following a point fixed a long a rolling wheel.

That concludes the ironically long path to rediscovering the path of shortest descent, the brachistochrone.

## 5  What Next?

While we found the cycloid that is generates the path of fastest descent, where else could we take this?

### 5.1  Calculating Time of Descent

While we found the curve, we didn't actually calculate the time itself it would take for a ball to roll down it.

Recall our definition of the total time travelled:

$$T = \int dt = \int \frac{ds}{v} = \int \frac{\sqrt{(dx)^2 + (dy)^2}}{\sqrt{2gy}}$$

If we factor a $dt$ out of the numerator, we can express our time of total descent in terms of our parametric equations as

8

$$T = \int \frac{\sqrt{(dx)^2 + (dy)^2}}{\sqrt{2gy}} = \int_0^{\theta_B} \frac{\sqrt{(\frac{dx}{d\theta})^2 + (\frac{dy}{d\theta})^2}}{\sqrt{2gy}} \, d\theta$$

where $\theta_B = \sin^{-1}(\sqrt{\frac{y_B}{C}})$. This gives the total (but the shortest!) amount of time to get between two points in space.

## 5.2   The Catenary

The calculus of variations may seem like a niche mathematical tool, but it is an extremely useful one when looking at these physical models like we had with the brachistochrone. Another famous problem is the one of the catenary:

**Problem.** *If a chain or cable hangs under its own weight between two points, what shape does it assume?*

The key idea from physics here is that a system always tends towards its most stable state. In other words, the chain will hang such that it minimizes its total potential energy. The potential energy of mass $m$ hanging a height $y$ above the ground is $U = mgy$ where $g$ is the acceleration due to gravity. Therefore,

$$dU = gy \, dm = gh(\rho \, ds)$$

$\rho$ is the linear mass density of our cable; by definition, $\rho = \frac{dm}{ds}$. We've already derived from before that $ds = \sqrt{1 + (\frac{dy}{dx})^2} \, dx$, so integrating $dU$ to get the total potential energy results in

$$U = \int_0^D gy\sqrt{1 + \left(\frac{dy}{dx}\right)^2} \, dx$$

$y$ is our function that describes the desired catenary, and $D$ is how far apart the two end points of our cable is hanging from. However, we're not quite done formulating this problem: we must subject this integral to the constraint that

$$\int_0^D \sqrt{1 + \left(\frac{dy}{dx}\right)^2} \, dx = L$$

since our cable is of a fixed length $L$. Just like with our brachistochrone, we are minimizing a functional $U$ that is dependent on $y$ and $y'$ (and implicitly $x$ since $y$ is a function of $x$). In a similar manner, the Beltrami identity of the Euler-Lagrange Equation can be used to solve this minimization problem, and the answer ends up being

$$y = a\cosh(\frac{x}{a})$$

The parameter $a$ is the measure of the lowest point of the curve from the ground. Also, note the locations of $a$ in the equation: there's one outside the main cosh function that scales the $y$-axis and one inside it dividing $x$ that scales the $x$-axis as well. This means that all catenary curves are similar to each other since all of them are just uniformly scaled to a different constant in all directions.