# Chapter 1. First Order Equations

## 1.0 Introduction.

An ordinary differential equation (ODE) in the two variables $t, y$ is any equation that includes derivatives of the function $y$ with respect to the independent variable $t$. The ODE is solved for the function $y(t)$.

**Example 0.1.** $\frac{dy}{dt} = 2y$ can be solved by *separating the variables and integrating both sides.*

Algebraically, separate the variables:

Multiply both sides by the differential $dt$ giving $\frac{dy}{dt} dt = 2y dt$.

Cancelling $dt$ gives $dy = 2y dt$ and divide by $y$ to get $\frac{dy}{y} = 2 dt$ , which has the variables separated on the two sides of the equation. Use Calculus to integrate the differentials:
$$\int \frac{dy}{y} = 2 \int dt \;\Rightarrow\; \ln|y| = 2t + c \text{ (assuming } y \neq 0),$$
where the arbitrary constants from the integrations are combined into c.

The solution $\ln|y| = 2t + c$ is called an *implicit solution* because it shows a relation between the solution $y$ and the independent variable $t$ that has no derivatives remaining but is not yet solved for $y(t)$. The rest is Algebra.

Solve for $y$ by taking the exponential on both sides: $e^{\ln|y|} = e^{2t+c} \Rightarrow |y| = e^{2t+c}$

This does **not** mean $y = e^{2t+c}$. We need to remove the absolute value by introducing "$\pm$" so
$$y = \pm e^{2t+c}$$
is not a solution because we don't know which sign to take. Notice that at $t = 0$ we get
$$y(0) = \pm e^{c}$$
and because $e^{c} > 0$ for every $c$, $y(0)$ has two values for every possible constant $c$. No function $y(t)$ can do that. To fix this, we cleverly absorb the sign ambiguity into a new constant, say
$$k = \pm e^{c}$$
and the solution is the function $y = k e^{2t}$. We already observed that $y(0) = \pm e^{c}$ as well, so it is just as good to write $y = y_0 e^{2t}$. The subscripted $y_0$ is just the right constant so that $y(0) = y_0$. For the general solution, we just write:
$$y = y_0 e^{2t}, \; y(0) = y_0 . \qquad\qquad \square$$

**Definition 0.1.** A solution to the ODE: $\frac{dy}{dt} = f(y, t)$ is a differentiable function $y(t)$ that satisfies the ODE on an open interval $a < t < b$. A family of solutions $y(t, c)$ that depends on an arbitrary constant $c$ and where every solution to the ODE is equal to $y(t, c)$ for some fixed value of $c$, is called a *general solution to the ODE*.

**Theorem 0.1.** $\frac{dy}{dt} = ay$ has general solution, $y = y_0 e^{at}, \; y(0) = y_0$ .

Proof: $\frac{dy}{dt} = \frac{d}{dt}(y_0 e^{at}) = y_0 \frac{d}{dt} e^{at} = y_0(a e^{at}) = a(y_0 e^{at}) = ay$ so $y = y_0 e^{at}$ is a function that satisfies the ODE for any value of $y_0$ and for all values of $t$. Based on our first example, every solution satisfies the implicit relation: $\ln|y| = at + c$ (unless y=0) so every solution must be $y = y_0 e^{at}$ with some value of $y_0 \neq 0$ or else $y = 0$, in which case $y_0 = 0$ is a choice of $y_0$ that makes the formula $y = y_0 e^{at} = 0$ true identically. By Def 0.1, $y = y_0 e^{at}$ is a general solution to $\frac{dy}{dt} = ay$. □

As illustrated by our first example, the arbitrary constant $k = y_0$, it may be possible to write the general solution in such a way that the arbitrary constant specifies the value of the solution when $t = 0$. More generally, it may be convenient to write the solution with arbitrary constant $k = y(t_0)$, for some other "special" time $t_0$. In any case, it is customary to write the specification as an

$$Initial\ Condition\ (IC)\quad y(t_0) = y_0$$

even though the solution doesn't really "start" at time $t_0$. The *IC* specifies a condition the solution must satisfy. An ODE together with an *IC* is called an *initial value problem* (IVP).

Example 0.2. $\frac{dy}{dt} = -y$, $y(0) = 3$ is an IVP. Its solution is $y = 3e^{-t}$.

Example 0.3. $\frac{dy}{dt} = -y$, $y(2) = 3$ is an IVP. Without the *IC*, the general solution is $y = ke^{-t}$. Use
the *IC* to solve algebraically for $k$: $3 = ke^{-(2)} \Rightarrow 3e^2 = k$ so the solution to the IVP is
$$y = (3e^2)e^{-t} = 3e^{2-t}.$$
□

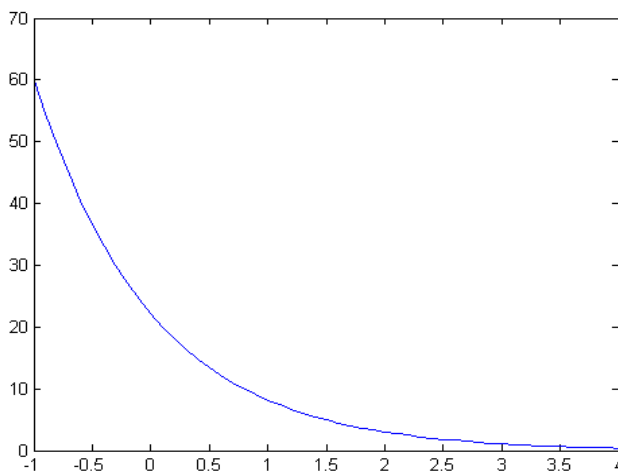You can make a graph of the solution using a calculator or software.



*Figure 0-1. Matlab plot of y = 3exp(2−t).*

MatLab code that produced the plot in Fig 1:

```
>> t=-1:0.1:4; y=3*exp(2-t);
>> plot(t,y)
```

The first line sets $t$ to be a list of values starting at -1 increasing by increments of 0.1 and ending at 4; y is set to be a list of the values of $3e^{2-t}$ at these values of $t$. Second line creates the plot.

Example 0.4. $\frac{dy}{dt} = 1 - 4t$, $y(0) = 1$ is an IVP. We can solve it by simple integration because the
right hand side is depends only on $t$. Multiplying both sides by $dt$, it says $\frac{dy}{dt} dt = (1 - 4t)dt$ and
we can integrate: $\int \frac{dy}{dt} dt = \int (1 - 4t)dt \Rightarrow y = t - 2t^2 + c$. The *IC*: $y(0) = 1$, requires $c = 1$,
so the solution to the IVP is: $y = t - 2t^2 + 1$.

For different *IC*s the solution will be the same in form with a different constant,

$$y(0) = 0 \implies y = t - 2t^2$$

$$y(0) = -1 \implies y = t - 2t^2 - 1$$

$$y(0) = -2 \implies y = t - 2t^2 - 2 \text{ and so on.}$$

We can plot these together using MatLab (or other software).

MatLab commands that produced Figure 0-2:

```
>> t=-2:0.1:2;
>> y1=t-2*t.^2+1*ones(size(t));
>> y0=t-2*t.^2;
>> y_1=t-2*t.^2-1* ones(size(t));
>> y_2=t-2*t.^2-2*ones(size(t));
>> h=plot(t,y1,'r-',t,y0,'g--',t,y_1,'b-.',t,y_2,'c:');
>> set(h,'LineWidth',2)
>> legend(h,'c = 1','c = 0','c = -1','c = -2')
```
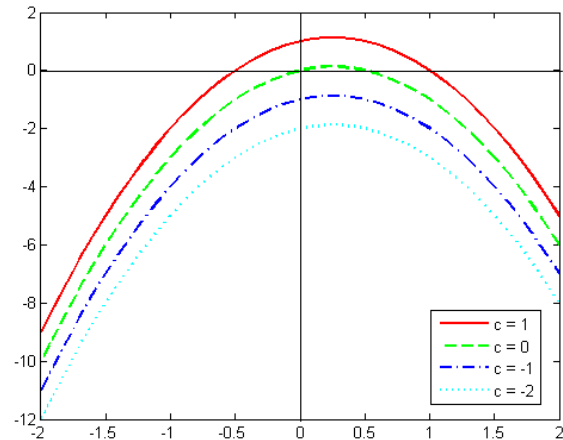
*Figure 0-2. MatLab plot of solutions to* $y' = 1 - 4t$ *for the four different IVPs of Example 0.4. The parabola* $y = t - 2t^2$ *is shifted vertically for the different initial values.*

The first line sets *t* to start at -2 and increase by increments of 0.1 ending at 2.The next four lines create the lists of *y* values for those values of *t*. The MatLab syntax for entrywise[1] exponentiation is ".^" not the usual "^". In lines two, four and five, when adding a constant to a matrix, we multiply the constant by a matrix full of ones that is the same size as the matrix we are adding it to, here that size is encoded as size(t). The line colors are 'r'=red, 'g'=green, 'b'=blue and 'c'=cyan. Line styles are '-'=solid, '--'=dash dash, '-.'=dash dot, ':'=dots. The axis lines were inserted in the plot using MatLab's Insert line tool. □

The entire *ty*-plane is covered by the graphs of the parabolic solutions to the ODE of Example 0.4. This means that the solutions for all *IC*s, $y(t_0) = y_0$ , are found by just specifying the correct value of the arbitrary constant of integration $c$ that picks out the parabola passing through the point $(t_0, y_0)$. This situation happens in general for first order linear equations: ODEs where the right hand side depends only linearly on $y$.

**Definition 0.2**. An ODE of the form $\frac{dy}{dt} = a(t)y$ is a *homogeneous first order linear equation*.

We will look at non-homogeneous 1st and 2nd order ODEs later. All linear ODEs have solutions that can be found more or less explicitly, unlike non-linear ODEs, which typically cannot. The homogeneous linear equation is the simplest of all. We have already solved the homogeneous linear equation $\frac{dy}{dt} = a(t)y$ when the function $a(t)$ is constant (the "constant coefficient" homogeneous linear equation). We solved it by the technique called "separation of variables". This technique is important enough to repeat.

---

[1] MatLab® is a matrix laboratory, so its operations are designed for matrices. Multiplication of matrices is denoted "*" but here, we need to multiply the *entries* of the matrix $t$ by themselves. To multiply matrices entry-wise is denoted ".*" so to square the entries in matrix $t$, we encode $\gg t.* t$; rather than $t * t$. The same goes for exponents; to square entries in $t$, we encode $\gg t.^2$; rather than $t^2$ which would fail because our $t$ is rectangular and such matrices cannot be raised to powers.

We algebraically transformed the constant coefficient equation into $\frac{dy}{y} = a\,dt$ so we could integrate the differentials on both sides getting the implicit solution $\ln|y| = at + c$. Taking the exponential function on both sides we found $|y| = e^{at+c}$. Removing the absolute value function introduced a sign ambiguity, which we absorbed into the arbitrary constant leaving $y = ke^{at}$. We can now do the same thing with the non-constant coefficient equation

$$\frac{dy}{dt} = a(t)y\,.$$

Separation gives $\frac{dy}{y} = a(t)\,dt$ and integration gives $\ln|y| = \int a(t)\,dt + c$. The same steps used in Example 0.1 lead to $y = ke^{\int a(t)\,dt}$ where $k$ is an arbitrary constant that depends on the *IC*.

**Theorem 0.2.** $\frac{dy}{dt} = a(t)y$ has general solution, $y = ke^{\int a(t)\,dt}$ where $k$ is an arbitrary constant.

Example 0.5. Let's find the general solution to the homogeneous linear equation $\frac{dy}{dt} = \cos(t)\,y$.

$$y = ke^{\int a(t)\,dt} = ke^{\int \cos(t)\,dt} = ke^{\sin(t)}$$

Remember that the arbitrary constant from the cosine integral is absorbed by the constant $k$. We can check the solution by taking the derivative $\frac{d}{dt}ke^{\sin(t)} = ke^{\sin(t)}\cos(t) = y\cos(t)$, so this solution is correct.

We used Matlab to plot the solutions for $k = 2, 1, \& 0.5$ in Figure 0-3.                □

Observe in Figure 3 that the local extrema for the solutions occur at the odd half-multiples of π: $-\frac{\pi}{2}, \frac{\pi}{2}, \frac{3\pi}{2}, \frac{5\pi}{2}$. This happens because their derivative $\frac{dy}{dt} = \cos(t)\,y$ vanishes at these *t*-values. In fact, we could have predicted this *without solving the equation*. The ODE says that at any point $(t, y)$, the solution's derivative $\frac{dy}{dt}$ equals $\cos(t)\,y$. We know $\cos\left(\frac{n\pi}{2}\right) = 0$ for all odd integers $n$, so $\frac{dy}{dt} = 0$ whenever $t$ is an odd half-multiple of π and the solution curve has a local extreme there. Furthermore,
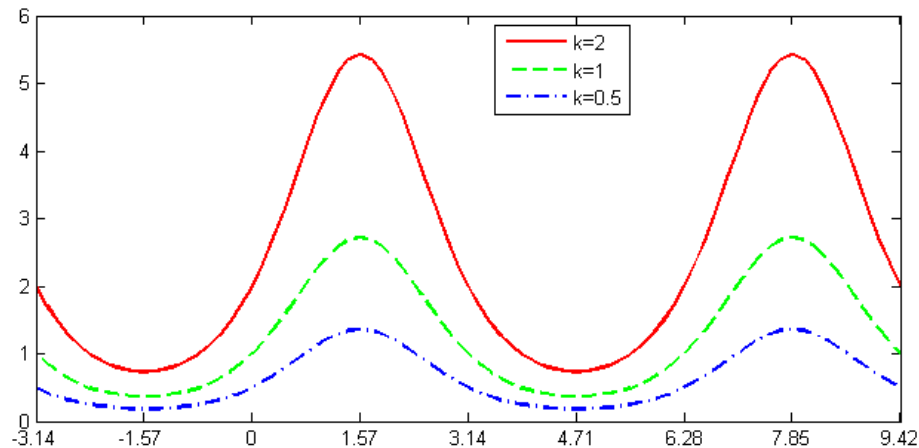


*Figure 0-3.Three solutions to y'= cos(t)y. t-axis scale is in half-multiples of pi.*

because the equation shows the formula for the derivative of every solution, it tells us the slope of every solution curve at every point. By plotting tiny tangent lines with these slopes, we can see the approximate shape of every solution curve *without solving the equation*. Such a plot is called a *slope field*.

4

Example 0.6. The slope field for $\frac{dy}{dt} = -2y$ is shown in Figure 0-4. Solutions are tangent to the tiny arrows. You can imagine solution curves passing through the slope field so they are tangent to the arrows. It may be easier to imagine the path of a particle moving through the slope field along a curve that just grazes arrows at it passes them. Several such paths are shown in Figure 0-5.  □



*Figure 0-4. Slope field for y' = -2y.*

The MatLab code that produced Figure 0-4 is:

```
>> [T,Y]=meshgrid(0:0.2:3,-2:0.2:2);
dydt=-2*Y;  I=ones(size(T));  quiver(T,Y,I,dydt);
```

The first line creates a pair of matrices holding the $t$-values (in T) and $y$-values (in Y) at which quiver will produce arrows. Two additional commands made Figure 0-4 into 0-5.

```
>> [iT,iY]=meshgrid(0:1:2, -2:2:3);
streamline(T,Y,I,dydt,iT(:),iY(:));
```
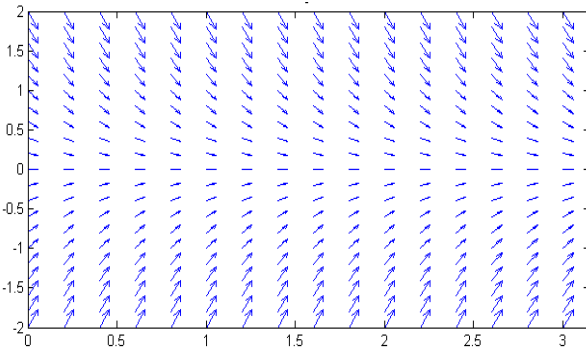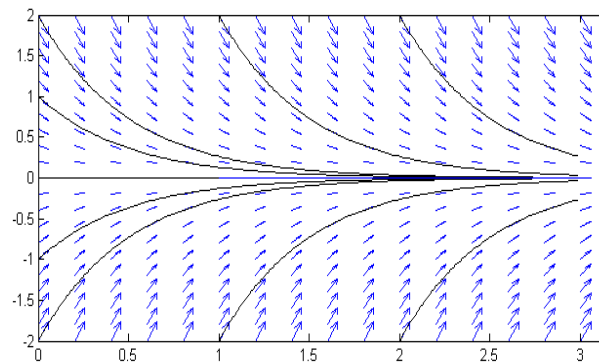


*Figure 0-5. Slope field for y' = -2y with a few solution*

The matrices iT and iY created in line 1 hold $t_0$ and $y_0$ values for solution curves that the streamline command will produce.[2] Theorem 0.1 tells us that the solution curves shown in Figure 0-5 are exponential functions because the solutions to $\frac{dy}{dt} = -2y$ are of the form $y = ke^{-2t}$. Because the slope field arrows and the solution curves tangent to them converge on the line $y = 0$, we see that the future of every initial condition is to approach $y = 0$ as $t \rightarrow \infty$. We think of the equation as describing a flow of initial conditions. The arrows point in the direction of the flow and the solution curves are *streamlines* of the flow. The streamlines give the paths of initial conditions under the influence of the flow determined by the differential equation.

One streamline is a straight line, it is the streamline $y = 0$ for the *IC* $y(0) = y_0 = 0$. Such a solution is called an *equilibrium*.[3]

---

[2] Four of the streamlines, ICs (1,-1), (1,1), (2,-1), (2,1), were removed using MatLab's plot Edit tools.

[3] The term comes from the origin of differential equations, which first appeared in Isaac Newton's monumental masterpiece *Philosophia Naturalis Principia Mathematica* of 1687 ("Mathematical Principles of Natural Philosophy", famously known as "Newton's Principia"). In his Principia, he gave the first comprehensive mathematical model for the motions of all massive objects in the universe based on his fundamental "Laws", the second of which was:

$$F = ma \quad \left(\text{Newton's 2}^{\text{nd}} \text{ Law of motion}\right).$$

Since $a = \frac{dv}{dt}$, this "Law" is a first order ODE $F = m\frac{dv}{dt}$ for the velocity $v$. If $y$ is position, then $v = \frac{dy}{dt}$, so $F = m\frac{d^2y}{dt^2}$ is a second order ODE in $y$. An object is in *equilibrium* when it is at rest or in a state of uniform motion (straight-line motion at constant velocity), in either of these cases the acceleration $a(t) = 0$.

**Definition 0.3.** A solution $y(t)$ to an ODE: $\frac{dy}{dt} = f(t, y)$ , is an *equilibrium* if $f(t, y(t)) = 0$, in which
case $y(t) = y_0$ for all time and the solution's streamline in the slope field is a horizontal line.

Throughout this course, we will consider physical systems that change in time. The evolution of such a
system is modeled by one or more differential equations. The evolution of the system in time is described
by a solution to the model equations. The *IC* for such a solution describes the state $y_0$ of the system at a
particular moment in time, denoted $t_0$. Each first order ODE describes the behavior of one of the numeric
properties of the physical system by equating its derivative $\frac{dy}{dt}$ to some function of the set of numeric
properties on which the derivative depends. Some systems are described by ODEs of higher order, e.g. a
second order equation will contain $\frac{d^2y}{dt^2}$. Basic mechanical examples of evolving properties are: *position*
*$y(t)$, velocity $v(t)$, acceleration $a(t)$*. Major goals of this text are to provide experience with a number of
very standard systems that occur in sciences and engineering, to develop the skill and background to be
able to solve and interpret their model equations, and the insight to be able to construct model equations
for novel systems.

## Exercises 1.0

1. Solve the IVP $\frac{dy}{dt} = ay$, $\quad y(t_0) = y_0$ and make a plot of the solution.

    *i*.) Use $a = -1$ and $y(0) = 10$.   *ii*.) Use $a = 3$ and $y(\ln 2) = 1$.   *iii*.) Using $a = -3$ and $y(1) = 4$.

2. Solve the IVP $\frac{dy}{dt} = te^{-at}$, $\quad y(t_0) = y_0$ and make a plot of the solution.

    *i*.) Use $a = 1$ and $y(0) = 2$.   *ii*.) Use $a = 3$ and $y(\ln 2) = 1$.   *iii*.) Using $a = -1$ and $y(1) = 5$.

# 1.1 First Order linear ODEs: The Easy Cases

Equations like those in Example 0.1, Theorem 0.1 and Definition 0.2 are examples of 1st order linear
ODEs. We give a general definition here:

**Definition 1.0**. A first order linear ODE can be written in the form: $\frac{dy}{dt} + p(t)y = f(t)$.

    If the equation written in this form has $f(t) = 0$, the equation is 1st order linear *homogeneous*,
    otherwise it in non-homogeneous.

Example 1.1. The equation $\frac{dy}{dt} = a(t)y$ from Definition 0.2 is first order linear *homogeneous* because it
    can be written as

$$\frac{dy}{dt} - a(t)y = 0.$$

    in the form of Def 1.0, with $p(t) = -a(t)$ and $f(t) = 0$.  This equation is easy to solve by
    separating the variables as done in the Introduction (see Theorem 0.2 and Example 0.5)         □

The equation $\frac{dy}{dt} = f(t)$ is the general form of Example 0.4 from the introduction. This equation is 1st
order linear (with $p(t) = 0$) and *non-homogeneous* but it is *separable*:

$$dy = f(t)dt \;\Rightarrow\; \int dy = \int f(t)dt \;\Rightarrow\; y = F(t) + c$$

where $F(t)$ is any anti-derivative of $f(t)$. Take a look at Example 0.4 where $f(t) = 1 - 4t$ so $F(t) = t - 2t^2$ or any of the other antiderivatives plotted in Figure 2. In some cases, you may not be able to find a simple formula for the antiderivative $F(t)$ but you can always write $F(t) = \int_{t_0}^{t} f(s)ds$ when you have an IVP.

**Theorem 1.0.** $f(t)$ continuous on an interval $(a, b)$ with $a < t_0 < b \Rightarrow \frac{dy}{dt} = f(t),\ y(t_0) = y_0$ has the unique solution

(1.0)                                              $y(t) = y_0 + \int_{t_0}^{t} f(s)\, ds\,.$

Proof. Definite integration in time implies (using the dummy variable $s$ in the definite integrals):

$$\int_{t_0}^{t} \frac{dy(s)}{ds}\, ds = \int_{t_0}^{t} f(s)\, ds\,.$$

The substitution, $y = y(s),\ dy = \frac{dy(s)}{ds} ds$, in the left side gives new endpoints $y(t_0)$ and $y(t)$, so

$$\int_{y(t_0)}^{y(t)} dy = \int_{t_0}^{t} f(s)\, ds\,,$$

Which upon integration becomes

$$y(t) - y(t_0) = \int_{t_0}^{t} f(s)\, ds\,.$$

Using $y(t_0) = y_0$ and moving the term to the right side gives equation 1.0, as desired.          □

Example 1.2. Let's solve the IVP $\frac{dy}{dt} = e^{-t^2},\ y(0) = 3$ .

Separate and integrate: $dy = e^{-t^2} dt \Rightarrow y = \int e^{-t^2} dt$. But there is no formula for the anti-derivative of $f(t) = e^{-t^2}$. To be honest, there is no formula using *elementary* functions. The integrand is continuous everywhere, so the integral exists on any bounded interval around $t_0 = 0$. So, it is correct to write:

$$y = 3 + \int_{0}^{t} e^{-s^2} ds\,.$$

where the dummy variable $s$ is introduced because $t$ appears as an endpoint in the integral.     □

The integral appearing in Example 1.2 cannot be simplified by the methods you learned in Calculus, however it has a value that can be calculated to great accuracy for any value of $t$. It is such a routinely needed integral that it has a special name and function symbol.
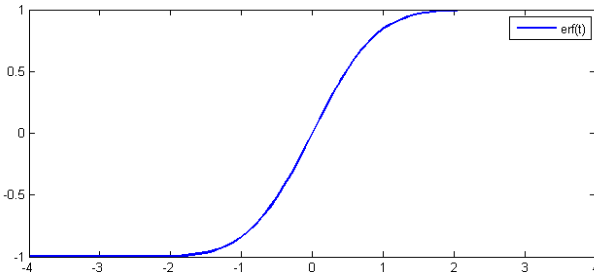
**Definition 1.1.** The *error function is*

$$erf(t) = \frac{2}{\sqrt{\pi}} \int_0^t e^{-s^2} ds.$$



*Figure 1-1. MatLab plot of the error function erf (t).*

```
>> t=-4:.1:4;
>> y=erf(t);
>> h=plot(t,y);set(h,'LineWidth',2);legend(h,'erf(t)')
```

The reason for the factor of $\frac{2}{\sqrt{\pi}}$ in front is to keep the values between -1 and 1. The error function is associated with the standard normal distribution in Probability.[4]

The solution to the IVP in Example 1.2 could also be written: $y = \frac{\sqrt{\pi}}{2} erf(t) + 3$.

Example 1.1 & 1.2 show the two special cases of $1^{st}$ order linear equations:

$$\frac{dy}{dt} + p(t)y = 0 \quad \text{and} \quad \frac{dy}{dt} = f(t).$$

In each case, the general solution is relatively simple to obtain by separation of variables. Aside from the case $\frac{dy}{dt} = f(t)$ of Theorem 1.0, non-homogeneous equations are not separable because both functions $p(t)$ and $f(t)$ appear. We need a new method of solution for such equations.

## Exercises 1.1

1. Give the general solution to $\frac{dy}{dt} + p(t)y = 0$
    i. $p(t) = t$     ii. $p(t) = -3t$     iii. $p(t) = -\sin^2 t$

2. Give the general solution to $\frac{dy}{dt} = f(t)$
    i. $f(t) = t \sin t$     ii. $f(t) = t^2 e^{-t}$     iii. $f(t) = e^{-t^2}$

3. Solve the IVP $\frac{dy}{dt} = f(t)$, $y(0) = y_0$ and plot your solution for $0 \le t \le t_{fin}$
    i. $f(t) = 2 - t^2$, $y_0 = 4$, $t_{fin} = 3$
    ii. $f(t) = \cos 2t$, $y_0 = 1$, $t_{fin} = \pi$
    iii. $f(t) = \cos 2t$, $y_0 = 0$, $t_{fin} = \pi$
    iv. $f(t) = t \sin t$, $y_0 = 0$, $t_{fin} = 3\pi$

4. Solve the IVP $\frac{dy}{dt} = ty$, $y(0) = 2$ and plot your solution for $-2 \le t \le 2$

5. Solve the IVP $\frac{dy}{dt} = -3ty$, $y(0) = 10$ and plot your solution for $0 \le t \le 5$

6. Solve the IVP $\frac{dy}{dt} = -ty$, $y(3) = 0$ and plot your solution for $0 \le t \le 6$

7. Solve the IVP $\frac{dy}{dt} = y \sin^2 t$, $y(0) = 2$ and plot your solution for $0 \le t \le 5$

---

[4] Symmetry of $e^{-s^2}$ implies $erf(t) = \frac{1}{\sqrt{\pi}} \int_{-t}^t e^{-s^2} ds$, which is the central integral of the Normal density $norm\left(0, \frac{1}{\sqrt{2}}\right)$ of variance ½.

## 1.2 Integrating Factor Method, Existence and Uniqueness of Solutions

The non-separable, non-homogeneous first order linear equation: $\frac{dy}{dt} + p(t)y = f(t)$ can always be
solved by multiplying by the *integrating factor*

$$\mu = e^{\int p(t)dt}$$

which according to Theorem 0.2, satisfies $\frac{d\mu}{dt} = p(t)\mu$. The equation multiplied on both sides by $\mu$
becomes

$$\mu\left(\frac{dy}{dt} + p(t)y\right) = \mu f(t)$$

The left side becomes:

$$\mu\frac{dy}{dt} + p(t)\mu\, y = \mu\frac{dy}{dt} + \frac{d\mu}{dt}\, y = \frac{d}{dt}(\mu y)$$

wherein the last equality comes from the product rule in Calculus. The resulting differential equation:

$$\frac{d}{dt}(\mu y) = \mu f(t)$$

Separates into

$$d(\mu y) = \mu f(t)dt$$

where integrating the left side gives $\mu y$ and the right side is a differential in $t$, so it can be integrated
too. Integration on both sides gives:

$$\mu y = \int \mu f(t)dt + c\,.$$

Dividing by $\mu$ gives the solution:

$$y = \frac{1}{\mu}\left(\int \mu f(t)dt + c\right).$$

The $\mu$ inside the integral **does not cancel the** $\mu$ outside the integral. There is also an arbitrary constant
that needs to multiply the $\frac{1}{\mu}$.

Example 2.1. Let's solve $\frac{dy}{dt} = 3t - 2y$ . Notice it is 1$^{st}$ order linear, general non-separable case,
because in the form $\frac{dy}{dt} + p(t)y = f(t)$ it becomes $\frac{dy}{dt} + 2y = 3t$. We need an integrating
factor

$$\mu = e^{\int p(t)dt} = e^{\int 2dt} = e^{2t}$$

$$\mu y = \int \mu(3t)dt + c = \int e^{2t}(3t)dt + c = 3\int te^{2t}dt + c \Rightarrow$$

$$e^{2t}y = 3\int te^{2t}dt + c \Rightarrow$$

$$y = 3e^{-2t} \int te^{2t} dt + ce^{-2t}.$$

We just need to do the integral $\int te^{2t} dt$, done with *integration by parts*, which we review here.

Recall the formula $\int u\, dv = uv - \int v\, du$ so for the integral $\int te^{2t} dt$, we set

$$u = t \quad \& \quad dv = e^{2t} dt$$

so

$$du = dt \quad \& \quad v = \int e^{2t} dt = \frac{e^{2t}}{2} \qquad \text{(taking the constant to be 0)}$$

Now

$$\int te^{2t} dt = \frac{te^{2t}}{2} - \frac{1}{2} \int e^{2t} dt = \frac{te^{2t}}{2} - \frac{e^{2t}}{4} \qquad \text{(the constant is already in } c \text{ )}$$

The solution is:

$$y = 3e^{-2t} \left( \frac{te^{2t}}{2} - \frac{e^{2t}}{4} \right) + ce^{-2t},$$

which simplifies to:   $y = 3\left(\frac{t}{2} - \frac{1}{4}\right) + ce^{-2t}.$                                      □


**Example 2.2.** Let's solve $\frac{dy}{dt} = 20 - \frac{4y}{50-2t}$. Rewrite it in the form $\frac{dy}{dt} + \frac{4y}{50-2t} = 20$ and recognize $p(t) = \frac{4}{50-2t}$ so the integrating factor is

$$\mu = e^{\int p(t)dt} = e^{\int \frac{4}{50-2t}dt}$$

Make a substitution, $u = 50 - 2t$, so $du = -2dt$, the integral becomes

$$4\int \frac{du}{-2u} = -2 \ln|u| \qquad \text{( again taking the constant to be 0)}$$

and

$$\mu = e^{\int p(t)dt} = e^{-2\ln|u|} = e^{\ln|u^{-2}|} = |u^{-2}| = \frac{1}{(50-2t)^2}.$$

The formula for the solution is

$$\mu y = \int \mu(20)dt + c = \int \frac{20}{(50-2t)^2} dt + c$$

Using the same $u = 50 - 2t$, so $du = -2dt$, this integral becomes: $20 \int u^{-2} \frac{du}{-2} = -\frac{10}{u}$
(constant already in $c$) so

$$\mu y = -\frac{10}{u} + c = \frac{-10}{50-2t} + c \implies$$

$$y = \frac{1}{\mu}\left(\frac{-10}{50-2t} + c\right) = (50-2t)^2\left(\frac{-10}{50-2t} + c\right) = 20t - 500 + c(50-2t)^2$$


Where either of the last two forms of the solution would be equally acceptable.                  □

Our discussion preceding these two examples all but proves the following theorem.

**Theorem 2.1.** The first order linear equation $\frac{dy}{dt} + p(t)y = f(t)$ with the functions $p$ & $f$ continuous on
an open interval $(a, b)$ such that $a < t_0 < b$ has a unique solution for every IC $y(t_0) = y_0$.

Proof. Because $p$ & $f$ continuous on $(a, b)$ and $a < t_0 < b$, we know the integral

$$\int_{t_0}^{t} p(s)ds$$

needed for the integrating factor $\mu(t) = e^{\int_{t_0}^{t} p(s)ds}$ in the interval $(a, b)$ exists. From the
Fundamental Theorem of Calculus we know the integral's dependence on its upper endpoint is
differentiable, so the integrating factor $\mu(t)$, as a composition of differentiable functions, is
then also differentiable as a function of $t$ in $(a, b)$.

The integral $\int \mu(t)f(t)dt$ in terms of the IC becomes

$$\int_{t_0}^{t} \mu(s)f(s)ds = \int_{t_0}^{t} e^{\int_{t_0}^{s} p(r)dr} f(s)ds$$

an integral of a continuous function, for $f$ is continuous on $(a, b)$ as well as $\mu$. Evidently, both
integrals exist for all $t$ in $(a, b)$ and both are differentiable functions of their upper endpoint.

Additionally, the function $\mu(t) = e^{\int_{t_0}^{t} p(s)ds}$ is an exponential, so it can never be zero; so its
reciprocal $\frac{1}{\mu(t)} = e^{-\int_{t_0}^{t} p(s)ds}$ is also differentiable on $(a, b)$.

The solution

$$y = \frac{1}{\mu(t)}\left(\int \mu(t)f(t)dt + c\right)$$

written in terms of the IC becomes:

(2.1)                     $y = e^{-\int_{t_0}^{t} p(s)ds}\left(\int_{t_0}^{t} e^{\int_{t_0}^{s} p(r)dr} f(s)ds + y_0\right)$

a product of functions differentiable for all $t$ in $(a, b)$ and so differentiable there itself. This,
together with our discussion prior to Example 2.1 shows that this is a solution on $(a, b)$ and it is
left to the reader to verify this solution satisfies the IC and so solves the IVP.

To prove this is the only solution to the IVP, we suppose there are two solutions and use the IVP
to see that they must be equal. Let $y$ be our solution (2.1) and $z$ be another solution to the IVP.

$$\frac{d(y-z)}{dt} + p(t)(y - z) = f(t) - f(t) = 0$$

because $\frac{d(y-z)}{dt} = \frac{dy}{dt} - \frac{dz}{dt}$ & $p(t)(y - z) = p(t)y - p(t)z$ , rearranging terms gives

$$\frac{dy}{dt} + p(t)y - \left(\frac{dz}{dt} + p(t)z\right) = f(t) - f(t) = 0$$

since both $y$ & $z$ are solutions. Now we have $(y - z)$ satisfying

$$\frac{d(y - z)}{dt} + p(t)(y - z) = 0$$

which by Theorem 0.2 implies $(y - z) = (y_0 - z_0)e^{-\int_{t_0}^{t} p(s)ds}$. But $y_0 = z_0$, so $(y - z) = 0$ or rather $y = z$. So, there is only one solution to the IVP. □

Theorem 2.1 is our first example of an *Existence and Uniqueness* (EU) theorem. Each time we consider a new form of equation, we will need such a theorem. The key points to notice are that it says:

**There exists a solution & there is only one solution to the IVP when $p$ & $f$ are continuous at the IC.**

### Exercises 1.2

1. Give the general solution to $\frac{dy}{dt} - ty = 4t$.

2. Give the general solution to $\frac{dy}{dt} = 10 - \frac{y}{5}$. Find the constant $c$ so that $y(0) = 30$.

3. Solve the IVP $\frac{dy}{dt} = -\frac{2y}{t} + e^{-2t}$, $y(2) = 4$ and plot the solution for $1 \le t \le 5$.

4. Following Example 2.2 as a guide, solve the IVP $\frac{dy}{dt} = 40 - \frac{5y}{100 - 4t}$ $y(0) = 10$.

## 1.3 Modeling Gravity and Drag

Example 3.1. According to Newton's second law of motion $F = ma$, where $F$ is the total force impressed on a moving object of mass $m$ and acceleration $a = \frac{dv}{dt}$. Newton's 2nd Law becomes a first order ODE for the object's velocity $v$.

Suppose the object is initially suspended at the height $y_0 = 1\ km$ above the surface of the Earth and at time $t_0 = 0$ it is dropped $(v(0) = 0)$. Can we predict the velocity of the object when it impacts the Earth? We need to solve the first order ODE:

$$m\frac{dv}{dt} = F$$

where $F$ is the force of gravity, $-mg$, plus the buoyancy of the object due to air resistance encountered as the object is falling. For this example, we suppose the object has no air resistance, so the equation becomes

$$m\frac{dv}{dt} = -mg \quad \Rightarrow \quad \frac{dv}{dt} = -g\ .$$

The equation separates and integrates to give:

$$\int dv = -\int g\ dt \quad \Rightarrow \quad v = v_0 - \int_0^t g(s)ds\ .$$

We just need a formula for $g(t)$. As it turns out, $g(t)$ depends on the height $y(t)$ but since $y(0) = 1\ km$ is pretty small compared to the size of the Earth, we can ignore the variation in the gravitational field of Earth[5] and assume $g(t) = 9.8055\ m/s^2$ which is nearly its average value between 0 & 1 $km$ altitudes.

---

[5] The acceleration due to Earth's gravity decreases in proportion to the square of the object's distance from the center of the Earth. Because the Earth has mean radius 6371 km, the formula for gravitational acceleration at height $y$ is $g(y) = 9.807\left(\frac{6371}{6371 + y}\right)^2 \approx (3404.3 - y)/347.32$. The change in acceleration is negligible over 1 km. Note $g(0) = 9.807 \Rightarrow g(0.5) = 9.8055$ & $g(1) = 9.8039$ within 0.0031 m/$s^2$ of the surface acceleration.

The equation is then: $\frac{dv}{dt} = -9.8055$, and since $v_0 = 0$ the solution is

$$v(t) = 0 - \int_0^t 9.8055 \, ds = -9.8055t \ (m/s)$$

Now, the height $y(t)$ in meters satisfies the ODE $\frac{dy}{dt} = v$ that has solution

$$y(t) = y_0 + \int_0^t v(t)dt = 1000 - 9.8055 \int_0^t s \, ds = 1000 - 9.8055\frac{t^2}{2} \ .$$

At impact, $y(t) = 0$, so we solve for the value of $t$ that makes

$$y(t) = 1000 - 9.8055\frac{t^2}{2} = 0 \ .$$

This easy quadratic has solutions: $t = \pm\sqrt{\frac{2000}{9.8055}} \approx \pm 14.282$, so the time until impact is about

14.28 sec. and $v(14.282) = -9.8055(14.282) \approx 140.4 \ m/s \approx 504 \ kph \ (\approx Mach \ 0.5)$.          □

In Example 3.1 we neglected air resistance. The air resistance of an object depends on its size, shape, and material density, whether and how it tumbles as it falls, how its temperature varies relative to the medium in which it falls, and so on. Given set values for all of these parameters, we still need to account for how the air resistance depends on the speed of the object. Air resistance is due to friction between the object and the air and resistance appears to increase at increasing speed. The simplest model supposes air resistance to increase *linearly* in proportion to the object's speed. This adds a force term to our ODE. The air resistance is a force that pushes in the opposite direction to the velocity, so we use a term $F_r = -kv$ where $k > 0$ is the "drag coefficient" of the object, based on its shape and density. This leads to a new equation for the falling object:

$$m\frac{dv}{dt} = -mg - kv \ \Rightarrow \ \frac{dv}{dt} = -g - \frac{k}{m}v \ .$$

This is again a first order equation in $v$. It can be written in the form of Theorem 2.1 as

$$\frac{dv}{dt} + \frac{k}{m}v = -g$$

which can be solved for any *IC* because the functions $p(t) = \frac{k}{m}$ & $f(t) = -g$ are continuous everywhere[6].

Using Theorem 2.1 the solution with $v(t_0) = v_0$ is:

$$v = e^{-\frac{k}{m}\int_{t_0}^t ds}\left(\int_{t_0}^t (-g) \exp\left(\frac{k}{m}\int_{t_0}^s dr\right) ds + v_0\right).$$

Example 3.2. As in Example 3.1, the object is initially suspended at the height $y_0 = 1 \ km$ above the surface of the Earth and at time $t_0 = 0$ it is dropped, so $t_0 = v_0 = 0$. The ODE for its velocity now depends on $\frac{k}{m}$. The solution from Theorem 2.1 is

$$v(t) = e^{-\frac{k}{m}t}\left((-g)\int_0^t \exp\left(\frac{k}{m}s\right) ds\right) = \frac{-gm}{k}\left(1 - e^{-\frac{k}{m}t}\right).$$

_____

[6] As you know from Calculus, all constant functions are continuous everywhere.

The solution $v(t) = -49.028(1 - e^{-0.2t})$ for $\frac{k}{m} = 0.2$ is shown in Figure 1-2. As $t$ increases, $v(t) \to \frac{-gm}{k} \approx -49.028$ asymptotically.

By $t = 25$ sec, $v(t)$ is nearly $\frac{-gm}{k} = -49.028$ m/$s$. This is the, so-called, *terminal velocity* of the falling object.

For the same mass $m$ but larger $k$, there would be more air resistance (drag), so the terminal velocity would be less, e.g. for $\frac{k}{m} = 0.4$, the terminal velocity would be half as fast, only about $-24.5$ m/$s$.

For the position function, solve $\frac{dy}{dt} = v(t),\ y(0) = y_0$ by separating and integrating to get:

$$y(t) = y_0 + \int_0^t v(s)ds$$

$$= 1000 - \frac{gm}{k}\int_0^t 1 - e^{-\frac{k}{m}s}\,ds$$

$$= 1000 - \frac{gm}{k}\left(t - \frac{m}{k}\left(1 - e^{-\frac{k}{m}t}\right)\right)$$

$$= 1000 - \frac{gm}{k}t + g\left(\frac{m}{k}\right)^2\left(1 - e^{-\frac{k}{m}t}\right).$$



*Figure 1-2. Velocity in meters/sec of an object dropped in air with k/m=0.2 for the first 25 sec of fall.*

For $\frac{k}{m} = 0.2$ and $g = 9.8055$, the height in meters is

$$y(t) \approx 1000 - 49.028\,t + 245.14\left(1 - e^{-0.2\,t}\right)$$

plotted in Figure 1-3 (solid line) along with (dashed line)

$$y(t) = 1000 - 9.8055\frac{t^2}{2}$$

from Example 3.1, neglecting drag.    □

In Example 3.2, we found formulas for both $v(t)$ and $y(t)$ but were not able to exactly solve the equation

$$y(t) = 0$$

to find the time of impact. This equation has both linear and exponential terms. The only way to solve it is to do so *numerically*. As is clear from the plot in Figure 1-3, impact will come shortly after $t = 25$ sec. The most obvious method is to calculate values of $y(t)$ just after $t = 25$ and find a first $t$-value where $y < 0$.



*Figure 1-3. Height y2(t) for object from Example 3.2, dropped from 1000 m in air (solid) compared to object from Example 3.1 dropped without drag y0(t) (dashed).*

We can easily find $y(25 + h)$ for $h = .01, .02, ...,$ until we find $y(25.36) = 0.2631$ and $y(25.37) = -0.2241$, so impact is at a time between $t = 25.36$ & $25.37$. Because the graph of $y(t)$ is approximately a straight line between these $t$ values, we can give an estimated impact time $t = 25.365$. For greater accuracy, we can calculate $y(t)$ values on a finer mesh between $t = 25.36$ & $25.37$.
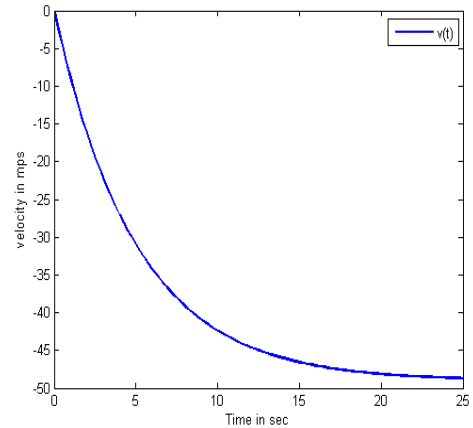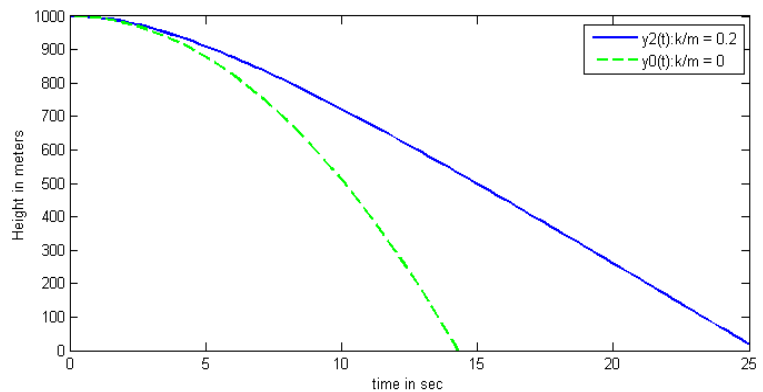
## Exercises 1.3

1. An object is initially suspended at the height $y_0 = 200\ m$ above the surface of the Earth and at time $t_0 = 0$ it is dropped. Neglecting air resistance, write an IVP for the object's velocity and solve it. Use $g = 9.8$.

2. An object is initially suspended at the height $y_0 = 200\ m$ above the surface of the Earth and at time $t_0 = 0$ it is hurled downwards, with initial velocity 10 $m/s$. Neglecting air resistance, write an IVP for the object's velocity and solve it. Use $g = 9.8$.

3. Proceed as in Example 3.1 to write and solve the IVP for the position $y(t)$ of the object in Problem 1. Find the time of impact and use it to find the object's velocity at impact.

4. Proceed as in Example 3.1 to write and solve the IVP for the position $y(t)$ of the object in Problem 2. Find the time of impact and use it to find the object's velocity at impact.

5. For the object in Problem 1, assume it has the coefficient $k/m$ given below. Include linear air resistance as in Example 3.2. Write and solve the IVP for $v(t)$. Use $g = 9.805$. Find $y(t)$.

   i. $k/m = 0.2$     ii. $k/m = 0.1$    iii. $k/m = 0.05$   iv. $k/m = 0.01$

6. For the object in Problem 2, assume it has the coefficient $k/m$ given below. Include linear air resistance as in Example 3.2. Write and solve the IVP for $v(t)$. Use $g = 9.805$. Find $y(t)$.

   i. $k/m = 0.2$     ii. $k/m = 0.1$    iii. $k/m = 0.05$   iv. $k/m = 0.01$

7. Make a plot comparing $y(t)$ from your solution to Problem 1 with $y(t)$ in each part of Problem 5.

8. Make a plot comparing $y(t)$ from your solution to Problem 2 with $y(t)$ in each part of Problem 6.

# 1.4 Mixing Problems

A common application of first order linear ODEs is to determine the evolution of a mixture of fluids or gasses in a container that receives input of one or more types of material as the mixture is also draining out. Many examples occur in manufacturing of fluid products and others in air and water quality estimation in environmental science.

Example 4.1. A tank of capacity 1000 L (litres) initially contains 300 L water with 2000 g salt dissolved in it.  Water with 2 g/L salt is mixed into the tank at the rate of 50 L/min.  At the same time, the mixed solution is draining from the tank at the rate of 50 L/min.

   a) Let's write the differential equation for $Q(t)$ = grams of salt in the tank at time = $t$ mins and give initial conditions.
   Fluids flow in and out of the tank each carrying salt at different concentrations, the amount of salt in the tank is changing as the tank fills.
   $$\frac{dQ}{dt} = salt\ inrate - salt\ outrate\ ,\quad Q(0) = 2000\ g$$
   Notice units matter, since $Q$ is in grams, $\frac{dQ}{dt}$ is in grams/min and so the same units should be used in
   $$salt\ inrate = (2\ g/L)(50\ L/min) = 100\ g/min$$
   The quantity expressed in $g/L$ is called the *concentration* of salt, or $[Salt]$. With this notation
   $$salt\ outrate = ([Salt(t)]\ g/L)(50\ L/min) = 50\ [Salt(t)]\ g/min.$$

   We need to find the salt concentration $[Salt(t)]\ g/L$ for the fluid draining out. The draining fluid has the same concentration of salt as the fluid in the tank:

We need to find the salt concentration $[Salt(t)]$ $g/L$ for the fluid draining out. The draining fluid has the same concentration of salt as the fluid in the tank:

$$[Salt(t)]\ g/L = \frac{salt\ in\ tank}{volume\ of\ fluid\ in\ tank} = \frac{Q(t)}{V(t)}\ g/L\,.$$

Because inflow and outflow rates are the same, volume of fluid in the tank, $V(t)$ stays constant at the initial volume of $300\ L$.

$$salt\ outrate = 50\,[Salt(t)]\ g/min = \frac{50\ Q(t)}{300}\ g/min = \frac{Q(t)}{6}\ g/min$$

Putting *inrate* and *outrate* into the equation gives the IVP for $Q = $ grams of salt in the tank:

$$\frac{dQ}{dt} = 100 - \frac{Q}{6}\ ,\quad Q(0) = 2000\ g$$

b) Now solve it. As a $1^{st}$ order linear ODE, put it in the form $\frac{dQ}{dt} + p(t)Q = f(t)$ of Theorem 2.1:

$$\frac{dQ}{dt} + \frac{Q}{6} = 100.$$

Using the integrating factor

$$r(t) = e^{\int \frac{1}{6}dt} = e^{\frac{t}{6}} \Rightarrow \frac{d}{dt}(rQ) = 100r = 100e^{\frac{t}{6}}$$

$$rQ = 100\int e^{\frac{t}{6}}\,dt = 600e^{\frac{t}{6}} + c$$

$$Q = \frac{600e^{\frac{t}{6}} + c}{r} = 600 + ce^{\frac{-t}{6}}$$

Using the *IC*:

$$Q(0) = 2000 = 600 + ce^{\frac{-0}{6}} \Rightarrow c = 1400$$

$$Q = 600 + 1400e^{\frac{-t}{6}}.$$

As $t \to \infty$, $Q(t) \to 600$ decreasing at all times.

c) Suppose the desired concentration of salt in the tank is $[Salt(t)] = 2.5 g/L$. We can find the time when this concentration is reached.

$$[Salt(t)] = \frac{Q(t)}{V} = \frac{600 + 1400e^{\frac{-t}{6}}}{300} = 2.5 \Rightarrow e^{\frac{-t}{6}} = \frac{150}{1400} \Rightarrow$$

$$t = -6\ln\frac{3}{28} \approx 13.4016\ min\,. \qquad\qquad \square$$

For a more interesting example, the volume of fluid in the tank can be increasing or decreasing in time because the inflow and outflow rates are different.

Example 4.2. A mixing tank of capacity 100 L  initially contains 64 L water with 10 g salt dissolved in it. A saltwater solution of 8 g/L salt flows into the vat at the rate of 5 L/min.  At the same time, the mixed solution is draining from the tank at the rate of 2 L/min.

a) Let's write the differential equation for $Q(t)$ = grams of salt in the tank at time = $t$ mins and give initial conditions.

$$\frac{dQ}{dt} = salt\ inrate - salt\ outrate\ ,\quad Q(0) = 10\ g$$

$$salt\ inrate = (8\ g/L)(5\ L/min) = 40\ g/min$$

and letting $[Salt(t)]$ represent concentration of salt in grams per litre,

$$salt\ outrate = [Salt(t)]\ g/L * 2\ L/min = 2\ [Salt(t)]\ g/min$$

The draining fluid has the same concentration $[Salt(t)]$ as the fluid in the tank:

$$[Salt(t)]\ g/L = \frac{salt\ in\ tank}{volume\ of\ fluid\ in\ tank} = \frac{Q(t)}{V(t)}g/L\ .$$

However, volume of fluid in the tank, $V(t)$ is changing in time:

$$V(t) = V(0) + net\ fluid\ inrate * time = 64 + (5-2)t = 64 + 3\ t$$

so,

$$salt\ outrate = 2\ [Salt(t)]\ g/min = \frac{2\ Q(t)}{64+3t}\ g/min$$

Putting *inrate* and *outrate* into the equation gives the IVP for $Q$ = grams of salt in the tank:

$$\frac{dQ}{dt} = 40 - \frac{2\ Q}{64+3t}\ ,\quad Q(0) = 10\ g$$

b) Now solve it. Put it in the form $\frac{dQ}{dt} + p(t)Q = f(t)$ of Theorem 2.1:

$$\frac{dQ}{dt} + \frac{2}{64+3t}Q = 40.$$

Using the integrating factor

$$r(t) = e^{\int 2(64+3t)^{-1}dt} = e^{\frac{2}{3}ln(64+3t)} = (64+3t)^{\frac{2}{3}}$$

$$\frac{d}{dt}(rQ) = 40r$$

$$rQ = 40 \int r\ dt = 40 \int (64+3t)^{\frac{2}{3}}\ dt = 40 * \frac{3}{5}\ (64+3t)^{\frac{5}{3}}\left(\frac{1}{3}\right) + c$$

$$= 8\ (64+3t)^{\frac{5}{3}} + c$$

$$Q = \frac{8\ (64+3t)^{\frac{5}{3}} + c}{r} = \frac{8\ (64+3t)^{\frac{5}{3}} + c}{(64+3t)^{\frac{2}{3}}} = 8(64+3t) + \frac{c}{(64+3t)^{\frac{2}{3}}}$$

$$\&\ Q(0) = 10 = 512 + \frac{c}{(64)^{\frac{2}{3}}}\ \Rightarrow\ c = -502(16) = -8032$$

$$Q = 8(64+3t) - \frac{8032}{(64+3t)^{\frac{2}{3}}}$$

c) Give the concentration of salt in the tank at the time of overflow.

$$V(t) = 64 + 3t = 100 \Rightarrow t = 12$$
$$Q(12) = 800 - \frac{8032}{(100)^{\frac{2}{3}}} = 427.19g$$

*Overflow concentration* $= 4.2719 \ g/L$   because the tank holds 1000 L.          □

## Exercises 1.4

1. A tank of capacity 500 L (litres) initially contains $V_0$ L water with 10 grams salt dissolved in it. Water with $[inflow \ salt]$ in g/L is mixed into the tank at the rate of 5 L/min. At the same time, the mixed solution is draining from the tank at the rate of 5 L/min. Write and solve the IVP for $Q(t)$, the quantity of salt in the tank at time $t$ for $V_0, [inflow \ salt]$ as given below.

   i. $V_0 = 100$, $[inflow \ salt] = 0$ g/L          ii. $V_0 = 200$, $[inflow \ salt] = 0$ g/L
   iii. $V_0 = 100$, $[inflow \ salt] = 50$ g/L       iv. $V_0 = 400$, $[inflow \ salt] = 25$ g/L

2. A tank of capacity 2000 L (litres) initially contains $V_0$ L water with 300 grams salt dissolved in it. Water with $[inflow \ salt]$ in g/L is mixed into the tank at the rate of 25 L/min. At the same time, the mixed solution is draining from the tank at the rate of 5 L/min. Write and solve the IVP for $Q(t)$, the quantity of salt in the tank at time $t$ for $V_0, [inflow \ salt]$ as given below.

   i. $V_0 = 1000$, $[inflow \ salt] = 0$ g/L          ii. $V_0 = 100$, $[inflow \ salt] = 10$ g/L
   iii. $V_0 = 1000$, $[inflow \ salt] = 50$ g/L       iv. $V_0 = 500$, $[inflow \ salt] = 20$ g/L

3. Make a plot of your solutions to Problem 1. Decide on a time interval that shows their asymptotic behavior. Include a legend.

4. Find the time of overflow for your solutions to Problem 2. Find the overflow concentrations. Make a plot of your solutions. Include a legend.

# 1.5 Non-linear First Order Autonomous Equations

The general form for a first order equation is:

(5.0)                                    $\frac{dy}{dt} = f(t, y)$

When $f$ is linear, we already know how to solve it. This section will cover the non-linear cases, where there is no general method to solve the equations. We first look at the *autonomous* case.

**Definition** 5.1. The equation $\frac{dy}{dt} = f(y)$ is *autonomous* (because $f(y)$ is independent of $t$).

Example 5.1. The equation $\frac{dy}{dt} = y^2$ is autonomous and non-linear. As done previously, we can separate variables, $\frac{dy}{y^2} = dt$ and integrate:

$$\frac{-1}{y} = \int \frac{dy}{y^2} = \int dt = t + c.$$

Solving algebraically for $y$ gives

$$y = \frac{-1}{t + c}.$$

Each value of $c$ gives a solution to a different IVP. Setting $t = 0$ in this seemingly general solution gives

$$y(0) = \frac{-1}{c}.$$

Evidently, every initial condition for which $y_0 \neq 0$ has a solution given by

$$y = \frac{-1}{t - \frac{1}{y_0}}.$$

But this formula gives no solution for which $y(0) = y_0 = 0$. There is such a solution though.
Clearly, $y(t) = 0$ works in the equation as an equilibrium solution.                                    □

Example 5.1 illustrates an important principle:

Non-linear equations may have solutions that _are not found_ by separation and integration. Furthermore:

**Principle 5.0.** Non-linear equations do not have general solutions.

In separating variables, we divide both sides of the equation $\frac{dy}{dt} = f(y)$ by $f(y)$. If there is some value $y_0$
for which $f(y_0) = 0$ then we divided by 0. That's why we only get solutions for $f(y_0) \neq 0$ when we solve
by separation. To clarify this, we can look at the slope-fields for a few different autonomous equations.
Remember the right hand side of the ODE gives the slope of the solution, so $f(y_0)$ is the slope of the
solution to the IVP

$$\frac{dy}{dt} = f(y), \qquad y(t_0) = y_0$$

at time $t_0$.

Example 5.2. Figure 1-4 shows the slope-field[7] for the non-linear autonomous equation

$$\frac{dy}{dt} = 1 - y^2.$$



The right hand side, $f(y) = 1 - y^2$ has two values that force $f(y_0) = 0$, namely,
$y_0 = \pm 1$. Since the right side is independent of time, solutions with ICs $y_0 = \pm 1$
are horizontal lines, these are the *equilibrium solutions* $y(t) = \pm 1$. A few other
streamlines are included to show the general shape of the non-equilibrium
solutions obtained by separation.

*Figure 1-4 Slope-field for y'=1-y² including solution curves as streamlines.*

The equilibrium solutions $y(t) = \pm 1$ are horizontal lines, along them, $y' = 0$ for
every value of $t$. The various streamlines in Figure 1 show the solutions $y(t)$ as
graphs of $y$-values evolving through time from a collection of different ICs.
Because the arrows in the slope-field are tangent to solutions, it is clear that for
any IC with $-1 < y_0 < 2$, the solution approaches the equilibrium $y(t) = 1$ as $t \to \infty$, so we
say the future limit of the solution is $\omega(y_0) = 1$. Similarly, the future limit $\omega(y_0) = -\infty$ for
solutions with $y_0 < -1$.

The past limit in Figure 1-4 for ICs with $-1 < y_0 < 1$ is $\alpha(y_0) = -1$, as is the past limit of
solutions with ICs $y_0 < -1$. The past limit $\alpha(y_0) = \infty$ for ICs $1 < y_0 < 2$.                    □

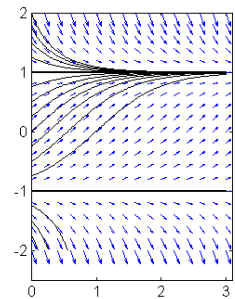The past and future limits of ICs describe the long-term behavior of solutions to an ODE.

---

[7] Look at Example 0.6 for MatLab code that makes a slope-field and inserts streamlines like those in Figure 1-4.

**Definition 5.2.** For a solution $y(t)$ to an IVP: $\frac{dy}{dt} = f(y)$, $y(t_0) = y_0$, the set of values $p$ such that there exists a decreasing sequence of times $t_k \to -\infty$ for which $\lim_{k\to\infty} y(t_k) = p$ is called a *past limit set* or *alpha limit set* $\alpha(y_0)$. The set of values $p$ for which there exists an increasing sequence of times $t_k \to \infty$ for which $\lim_{k\to\infty} y(t_k) = p$ is called a *future limit set* or *omega limit set* $\omega(y_0)$.

**Remark 5.0.** We allow $p$ to potentially take one or both of the values $\pm\infty$. Regardless, under this definition, some initial conditions will have empty alpha or omega limit sets, specifically when the solution $y(t) \to \pm\infty$ in finite time. We will consider such cases a little later.

In the case of autonomous ODEs, the alpha and omega limit sets are determined by the signs of the right hand side $f(y)$ on the $y$-intervals surrounding the equilibria.

You learned in Algebra how to graph a polynomial function $f$ by finding its zeros, $z_1, \ldots z_n$ and choosing sample points $c_0, \ldots c_n$ in the intervals surrounding them to check whether $f(c_k)$ would be positive or negative. Because $f$ is continuous, it can only change sign at a point $z_1, \ldots z_n$ where $f(z_k) = 0$, so knowing the sign of $f(c_k)$ determines the sign of $f$ in the interval between consecutive zeros $z_k, z_{k+1}$. We apply this to non-linear autonomous equations with polynomial right hand sides. It is a simple algebraic technique that yields a complete picture of the *behaviour of all solutions*.

Example 5.3. Consider the autonomous ODE,

$$\frac{dy}{dt} = y^2 - 5y + 6.$$

Factoring the right side gives

$$\frac{dy}{dt} = (y-2)(y-3).$$

The equilibria are $y(t) = 2$ and $y(t) = 3$. The graph of $y^2 - 5y + 6$ is a parabola opening in the positive direction (the sign on $y^2$ is positive). It has the zeros at $y = 2 \,\&\, 3$ so it is negative between them and positive "outside" them. The parabola is graphed in Figure 1-5.
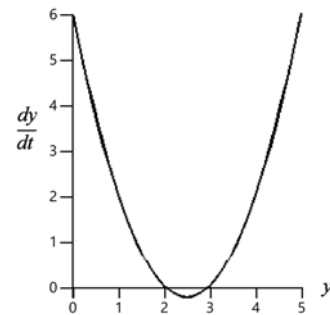
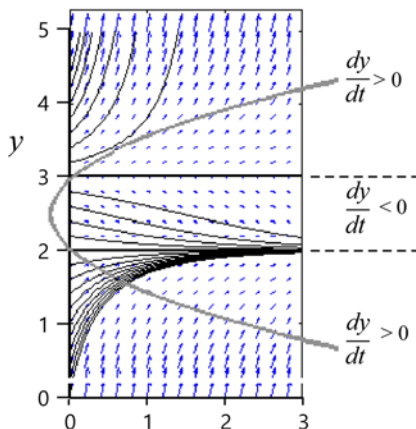*Figure 1-5. Parabola $y' = y^2 - 5y + 6$ shows $y' < 0$ for $2 < y < 3$ and $y' > 0$ "outside" the interval [2,3].*

Figure 1-6 shows the same parabola with axes flipped & the $y$-axis vertical, superimposed on the slope-field. Solutions with ICs between the equilibria, where $\frac{dy}{dt} < 0$, are decreasing; solutions with ICs $y_0 > 3$ or $y_0 < 2$ are increasing because $\frac{dy}{dt} > 0$.  □

*Figure 1-6. Parabola $y' = y^2 - 5y + 6$ superimposed on the slope-field.*

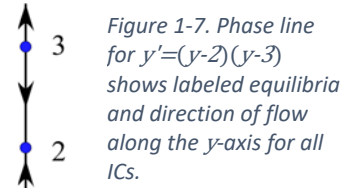### The Phase line and Stability of Equilibria

The situation depicted in Figure 1-6 can be summarized in a much simpler diagram, called a "phase diagram". The phase diagram for the equation $\frac{dy}{dt} = f(y)$ is just a vertical line representing the initial conditions with the equilibria indicated by labeled points. Along the phase line the direction of flow, increasing or decreasing, is indicated by arrows.

Example 5.4. The phase line for

$$\frac{dy}{dt} = (y-2)(y-3)$$

is pictured in Figure 1-7.                    □

*Figure 1-7. Phase line for y'=(y-2)(y-3) shows labeled equilibria and direction of flow along the y-axis for all ICs.*
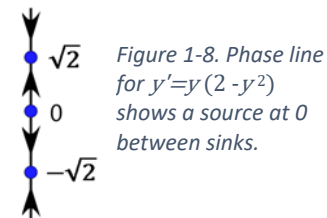
The *flow* on the phase line is induced by the motions of initial points $y_0$ carried along the line by the solutions to the IVPs with ICs $y(0) = y_0$. At time $t = 0$ the points $y_0$ in the interval between 2 & 3 have velocities given by the right side $(y_0 - 2)(y_0 - 3) < 0$ and as time goes forward, these points move downward along the phase line. Points near to the equilibria move more slowly than points further away. No point ever reaches the equilibrium at 2, but every point below 3 approaches 2 as a limit. For $y_0 < 3$, the equilibrium $y = 2$ is the omega limit set $\omega(y_0)$. We think of the equilibrium $y = 2$ as an *attractor* commonly called a *sink*. Similarly, for $y_0 > 2$, because $y = 3$ is the alpha limit set $\alpha(y_0)$, we can think of $y = 3$ as a *repellor* or *source*. Figuratively, we might imagine there is a compressible fluid (like air) flowing in a tube through which molecules must pass in single file. The molecule at $y = 3$ is stopped and the one at $y_0 = 2$ is also stopped. The other molecules have $y_0 = 3$ as their alpha limit, in their infinite past. Similar comments apply to $y = 2$, which is in their infinite future, so they never get there; it's their omega limit. The equilibria are blockages in the flow.

Example 5.5. The phase line for

$$\frac{dy}{dt} = y(2 - y^2)$$

is pictured in Figure 1-8.                    □

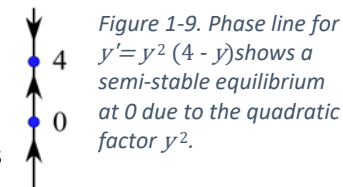*Figure 1-8. Phase line for y'=y(2 -y²) shows a source at 0 between sinks.*

We call sinks, *stable* equilibria because all solutions that start sufficiently nearby, stay nearby and can't wander off. We call sources *unstable* equilibria because, they are *not* stable: no matter how near to them an IC lies, it must eventually wander away. Not every equilibrium is either a source or a sink, as the next example shows.

Example 5.6. The cubic right side of

$$\frac{dy}{dt} = y^2(4 - y)$$

has the squared factor $y^2$ and so the equilibrium at $y = 0$ is unstable on one side and stable on the other as illustrated in Figure 1-9.                    □

*Figure 1-9. Phase line for y'=y² (4 - y) shows a semi-stable equilibrium at 0 due to the quadratic factor y².*

## Exercises 1.5 (See MatLab code in Example 0.6 for Problems 1-5. Problems 6&7 can be drawn by hand.)

1. Plot a slope-field for $\frac{dy}{dt} = y(5 - y)$. Show some streamlines.

2. Plot a slope-field for $\frac{dy}{dt} = 25 - y^2$. Show some streamlines

3. Plot a slope-field for $\frac{dy}{dt} = y^2 - 9$. Show some streamlines

4. Plot a slope-field for $\frac{dy}{dt} = y(25 - y^2)$. Show some streamlines

5. Plot a slope-field for $\frac{dy}{dt} = y(y^2 - 9)$. Show some streamlines

6. Make a phase-line diagram for each of these autonomous equations. Indicate equilibrium solutions and use arrows to indicate the direction of flow on each segment of the line.

$i. \dfrac{dy}{dt} = 25 - y^2$   $ii. \dfrac{dy}{dt} = y(25 - y^2)$   $iii. \dfrac{dy}{dt} = y^2(25 - y^2)$   $iv. \dfrac{dy}{dt} = y^3 - 2y^2 + y$

7. Make a phase-line diagram for each of these autonomous equations. Indicate equilibrium solutions and use arrows to indicate the direction of flow on each segment of the line.

$i. \dfrac{dy}{dt} = 9 - y^2$   $ii. \dfrac{dy}{dt} = y(y^2 - 9)$   $iii. \dfrac{dy}{dt} = y^2(y^2 - 2)$   $iv. \dfrac{dy}{dt} = -y^3 + 6y^2 - 9y$

8. Find the limit sets $\alpha(y_0)$ & $\omega(y_0)$ for initial conditions $y_0$ in each non-equilibrium interval of the phase lines in problem 6.

9. Find the limit sets $\alpha(y_0)$ & $\omega(y_0)$ for initial conditions $y_0$ in each non-equilibrium interval of the phase lines in problem 7.

## 1.6 Existence and Uniqueness

In this section, we look at theorems guaranteeing the existence and uniqueness of solutions to the IVP:

(6.0)                                   $\dfrac{dy}{dt} = f(t, y); \ \ y(t_0) = y_0.$

Linear IVPs with continuous coefficient functions have unique solutions for all ICs as we proved in Theorem 2.1. The non-linear IVP (6.0) will also have a solution when $f$ is continuous on an open domain containing the IC, however, it can happen in the non-linear case that the IVP has *more than one* solution. If we imagine the IVP is describing the motion of a point in a flow or the progress of some physical process, then non-uniqueness of the solution implies that we can't predict which of the various solutions the physical system will follow. This is problematic for an engineer or scientist trying to predict the evolution of a process from a known initial state — which solution does it follow? Before we give an example of non-uniqueness, let's settle the existence question.

**Theorem 6.1.** (Peano[8] ) If $f$ is continuous on a rectangle $R = \{(t, y) | \ |t - t_0| < a, |y - y_0| < b\}$ containing the initial point $(t_0, y_0)$, then there exists a solution $y(t)$ to the IVP (6.0) on an interval $|t - t_0| < \delta \le a$.

The proof of Peano's Existence Theorem is beyond the scope of our course but its consequences are not.

Example 6.1. The ODEs of examples 5.1-6 have solutions for every IC $y(t_0) = y_0$ because their right sides are polynomials continuous at every $y_0$.                                                    □

Example 6.2. The ODE $\dfrac{dy}{dt} = \sqrt[3]{y}$ has solutions for every IC $y(t_0) = y_0$ because the right side is continuous at every $y_0$. When we separate variables, we get $\int \dfrac{dy}{\sqrt[3]{y}} = \int dt$ which leads to

$$\frac{3}{2}y^{\frac{2}{3}} = t + c \ \Rightarrow \ y = \pm\left(\frac{2t + 2c}{3}\right)^{\frac{3}{2}} = \pm\left(\sqrt{\frac{2t + \xi}{3}}\right)^3$$

for which $y(0) = \left(\sqrt{\dfrac{\xi}{3}}\right)^3 \ \Rightarrow \ \xi \ge 0$ because of the square root. The solutions for $\xi = 0$ are

$$y(t) = \pm\left(\sqrt{\frac{2t}{3}}\right)^3$$

---

[8] First proved by Giuseppe Peano 1890, often called the Peano or the Cauchy-Peano Existence Theorem.

satisfying the IC $y(0) = 0$ but here the square root requires $t \geq 0$. These solutions are only defined to the right of $t_0 = 0$. Peano's theorem says there is a solution on both sides of the IC. There's another solution of course, the equilibrium solution $y(t) = 0$ passing through the IC. □

Most interesting in Example 6.2 is that there are three solutions exiting the IC on the right. This means that the equilibrium solution $y(t) = 0$ for $t \leq 0$ branches at $t_0 = 0$ and has three distinct futures on $t \geq 0$. This an example of *non-uniqueness* of solution to an IVP.[9]

Because of cases like Example 6.2, continuity of $f(t,y)$ in a neighborhood of the IC is not enough to guarantee uniqueness of solution to the IVP 6.0. Various conditions on $f(t,y)$ will guarantee uniqueness. The simplest, due to Augustin-Louis Cauchy (1789-1857)[10], is continuity of $\partial_y f(t,y)$ in a neighborhood of the IC. A standard weaker condition is *Lipschitz continuity of $f(t,y)$ as a function of $y$.*[11]

**Theorem 6.2.** Existence and Uniqueness of Solutions ("Cauchy's EU"). Suppose $f(t,y)$ and $\partial_y f(t,y)$ are continuous on a rectangle $R$ in the $ty$-plane and the point $(t_0, y_0)$ is in the interior of $R$, then the initial value problem (1.3) has unique solution for $t$ in an open interval $|t - t_0| < \delta$.

An equation like $\frac{dy}{dt} = \sqrt[3]{y}$ in Example 6.2, fails the condition of continuity for $\partial_y f(t,y)$ at $y_0 = 0$ for any $t_0$. This alone doesn't guarantee non-uniqueness. However, because the partial derivative $\partial_y \sqrt[3]{y} = \frac{1}{3} y^{-\frac{2}{3}}$ fails to exist at $y_0 = 0$, we know not to trust the equation to have unique solutions at such ICs. The Cauchy (and Cauchy-Lipschitz or Picard-Lindelöf) Existence and Uniqueness theorems give sufficient conditions for uniqueness of solutions, if these conditions are not met, it raises suspicion that there *may* be non-uniqueness. For our work and in most practical cases, that is enough of a warning to check carefully before proceeding with solutions that are likely not unique.

*Augustin-Louis Cauchy 1840 Lithograph (Wikimedia)*

## Exercises 1.6

1. Consider the IVP $\frac{dy}{dt} = f(t,y);\ y(t_0) = y_0$. In each case state whether or not Peano's theorem guarantees a solution exists and if not, explain why not.

   i. $f(t,y) = y\cos(t^2)$, $y(0) = 2$   ii. $f(t,y) = -\frac{t}{y}$, $y(1) = 0$   iii. $f(t,y) = -\frac{y^2}{2t}$, $y(1) = 0$

   iv. $f(t,y) = \sqrt[3]{y}$, $y(0) = 2$     v. $f(t,y) = \frac{1}{\sqrt{1+t^2}}$, $y(0) = 2$   vi. $f(t,y) = \frac{1}{\sqrt{1-t^2}}$, $y(1) = 2$

2. For each of the IVPs in Problem 1 for which Peano's theorem guarantees existence of a solution, state whether Cauchy's EU theorem guarantees uniqueness of the solution.

---

[9] At each IC, $y(t_0) = 0$, the equilibrium solution branches. The branching solution is $y(t) = \sqrt{\frac{2(t-t_0)}{3}^3}$

existing for $t \geq t_0$. It is also not hard to see that $y(t) = -\sqrt{\frac{2(t-t_0)}{3}^3}$ satisfies the IVP only for $t \geq t_0$, giving a second solution branching off the equilibrium solution at the IC, $y(t_0) = 0$.

[10] French mathematician Cauchy is the second most productive in history; his published works fill 27 large volumes.

[11] The Lipschitz condition is standard in upper level ODE courses. With it, our "Cauchy EU" Theorem becomes the Cauchy-Lipschitz Theorem or alternatively the Picard-Lindelöf because both Émile Picard (1856-1941 France) and Ernst Lindelöf (in 1890 at age 20 y. Finland) proved it using only Lipschitz continuity of $f$ by similar arguments.

3. The equation $\frac{dy}{dt} = y^{2/3}$ is a classic example discussed in many textbooks. Separate the variables and solve this equation. Discuss what happens at the IC $y(1) = 0$ and draw a picture of solutions at this IC. Consider both Peano's and Cauchy's theorems. Discuss briefly what they say about this equation in terms of existence of solutions at ICs $y(0) = y_0 \neq 0$ versus $y(0) = 0$.

## 1.7 Separable first order equations

**Definition** 7.1. The equation $\frac{dy}{dt} = f(t,y)$ is separable when $f(t,y) = a(t)b(y)$.

Example 7.1. The equation $\frac{dy}{dt} = ty^2$ is separable. As in Chapter 1, we separate variables, $\frac{dy}{y^2} = t\,dt$ and integrate:

$$\frac{-1}{y} = \int \frac{dy}{y^2} = \int t\,dt = \frac{t^2}{2} + c\,.$$

Solving algebraically for $y$ gives

$$y = \frac{-2}{t^2 + 2c}\,.$$

There is an additional solution though. Clearly, $y(t) = 0$ works in the equation as an equilibrium solution. Solutions exist and are unique for every IC by Cauchy's EU Theorem 6.2.                    □

In contrast to linear equations, non-linear separable equations may have equilibrium solutions not found by separation of variables. You just have to remember to check for them. In fact, it is a good policy to always find equilibrium solutions *first*. They are typically easy to find because you just need to find constant values of $y$ that make the right-hand side zero.

Example 7.2. Find all solutions to $\frac{dy}{dt} = (y^2 - 4)\cos t$.

It is separable and has a pair of equilibrium solutions $y(t) = \pm 2$. Separation gives

$$\int \frac{dy}{y^2 - 4} = \int \cos t\,dt = \sin t + c\,.$$

The left side integrates by partial fractions. The result is:

$$\int \frac{dy}{y^2 - 4} = \frac{1}{4}\ln\left|\frac{y - 2}{y + 2}\right| = \sin t + c$$

Upon exponentiating both sides and taking 4$^{\text{th}}$ powers,

$$\left|\frac{y - 2}{y + 2}\right| = e^{4(\sin t + c)} = e^{4c}e^{4\sin t}.$$

Consequently, $\frac{y-2}{y+2} = \pm e^{4c}e^{4\sin t}$, in which case we let a new arbitrary constant $k = \pm e^{4c}$ as explained in Example 0.1 of the Introduction in Chapter 1. Here, $k = \frac{y_0-2}{y_0+2}$. The implicit solution:

$$\frac{y - 2}{y + 2} = ke^{4\sin t}$$

can be solved algebraically for

$$y = 2\left(\frac{1 + ke^{4\sin t}}{1 - ke^{4\sin t}}\right).$$

Notice as $k \to 0$, $y(t) \to 2$, one equilibrium solution and the other, $y(t) = -2$, results by letting $k \to \infty$.

□

## Exercises 1.7

1. Solve $\frac{dy}{dt} = -\frac{t}{y}$. Plot a slope-field with a few streamlines for $-2 \le t \le 2$.

2. Solve $\frac{dy}{dt} = -\frac{t^2}{y}$. Plot a slope-field with a few streamlines for $-2 \le t \le 2$.

3. Solve the Logistic equation $\frac{dy}{dt} = 4y - y^2$. Plot a slope-field with a few streamlines for $0 \le t \le 4$. Be sure equilibrium solutions are included.

4. Solve the Logistic equation $\frac{dy}{dt} = y - 4y^2$. Plot a slope-field with a few streamlines for $0 \le t \le 4$. Be sure equilibrium solutions are included.

5. Solve $\frac{dy}{dt} = t(1 - y^2)$. Plot a slope-field with a few streamlines for $0 \le t \le 2$. Be sure equilibrium solutions are included.

6. Solve $\frac{dy}{dt} = (4 - y^2)\sin \pi t$. Plot a slope-field with a few streamlines for $0 \le t \le 10$. Be sure equilibrium solutions are included.

# 1.8 Numerical Integration

As we have seen, sometimes there is no elementary formula for the solution to an ODE (recall Example 1.2). Whenever this is the case, and only when we are *assured that a unique solution exists* for the IVP in question, we can use *a numerical method* to calculate the values of $y(t)$ at specific values of $t$.

In Example 1.2 we solved the IVP $\frac{dy}{dt} = e^{-t^2}$, $y(0) = 3$ by separating variables and integrating to find $y = \int_0^t e^{-s^2} ds + 3$. The integral has no elementary form. In Calculus class, you would say: "we can't do the integral of $e^{-t^2}$", because it has no antiderivative. In fact, there is no function made up of any combination of your familiar polynomial, rational, root, exponential, or trig functions that is an antiderivative for $e^{-t^2}$. We saw in Example 1.2 that the "*special*" function,

$$\text{erf}(t) = \frac{2}{\sqrt{\pi}} \int_0^t e^{-s^2} ds,$$

whose derivative equals $\frac{2}{\sqrt{\pi}} e^{-t^2}$, can be used to get the antiderivative as $\frac{\sqrt{\pi}}{2}\text{erf}(t)$. There are many such *special functions* that arise from integrals Calculus students "can't do". Rather than listing some of them, just observe that any function of the form

$$F(t) = \int_{t_0}^t f(s) ds$$

is differentiable[12] on the largest interval of continuity of $f$ that contains $t_0$. This means that for *any continuous* function $f$, the integral $\int_{t_0}^t f(s) ds$ depends *differentiably* on the value of $t$. If we calculate the integral using a Riemann sum or another method that partitions the interval into sub-intervals, we can

---

[12] By the Fundamental Theorem of Calculus, $f$ continuous on $(a, b) \Rightarrow F'(t) = f(t)$.

make estimates of the difference between our approximate calculation and the actual value of the integral by using a finer mesh partition of the interval $(t_0, t)$ to compare a closer approximation to $F(t)$.

Example 8.1. Let's calculate the value of $y(t) = \int_0^t e^{-s^2} ds$ at the point $t = 2$ using the trapezoid rule from Calculus class. To begin, we partition the interval $(0, t)$ using mesh size $\Delta t = 0.1$, which gives partition points $t_0, t_1, \ldots t_{20} = 0, 0.1, \ldots, 2.0$, as in Figure 1-10. The trapezoid rule sums the area of a the trapezoid whose slanted tops are secant lines of the curve between the points $(t_k, y(t_k))$ and $(t_{k+1}, y(t_{k+1}))$. The sum
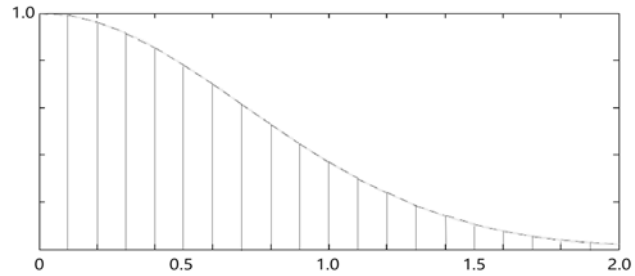


Figure 1-10. Integral of $y=exp(-t^2)$ partitoned into 20 trapezoids whose top edges are secants (dashed) joining consecutive points along the curve.

of these 20 trapezoid areas is 0.882020 … a close approximation to the area under $y = e^{-t^2}$ for $0 \leq t \leq 2$. The correct value of $y(2) = 0.882081$ …. good to six digits. The accuracy of the approximation would be improved using a finer mesh partition. □

In this way, we can calculate values of

$$y(t) = \int_0^t e^{-s^2} ds$$

to almost arbitrarily high precision. This calculation of values of $y(t)$ is called *numerical integration*. When $y(t) = \int_{t_0}^t f(s)ds + y_0$ is a solution to an IVP but the integral must be done numerically, we plot the results as a *numerical solution*. We know the plot is not exact but by specifying a finer mesh partition and comparing the two numerical solutions for different mesh sizes, we can estimate how much error there is in the approximation. For example using mesh size $\Delta t = 0.05$ for Example 8.1, we obtain $y(2) = 0.882066$ … so we can estimate the error in the latter to be less than their difference, 0.000046.

In the real world, numerical integration is carried out by adaptive methods, such as those implemented in MatLab®. Adaptive methods adjust the mesh of the partition to control error and improve efficiency. The general name for integration is *quadrature*[13]. To use MatLab's *quad* numerical integrator, the integrand needs to be defined in a separate file, a *function* "M-file" that is written by the user and stored in the location (called the "Path") shown in the box "Current Directory"[14] at the top of the MatLab command window. It is usually helpful to run >>doc *function* to read documentation before writing your first M-file function. Next, select the *Window* tab and choose *Editor* from the dropdown menu. Choose File, choose New, choose M-File. This is the template into which you will write your function.

---

[13]  The ancient Greeks recognized the problem of finding areas of curved regions as early as 450 BCE. Having already developed methods to find areas of polygons, they sought methods that would construct a quadrilateral polygon of the same area as a given circle or part thereof and hence determining its area..

[14] If you just open MatLab, it will probably be using the default Path, C:\Users\\*yourusername*\Documents\MATLAB. If you are not able to write to that directory, you will need to browse for another (click the ellipsis button) where you do have write permission.

Example 8.2. Let's use MatLab's numerical integrator *quad* to integrate $e^{-t^2}$ from $t = 0..2$. Open the MatLab Editor and type in the word *function* . It turns blue. Choose a descriptive name for the function like "BasicGaussian" and type in "y = BasicGaussian(t)" and then save the file as "BasicGaussian".     MatLab will save it as "BasicGaussian.m" in the current directory. Start a new line and type % followed by your name. The "%" tells MatLab the following characters are a "Comment" so it won't read them. Complete the function exactly as shown here and save it. Note we used the entry-wise form ".^" (also called "vector-form") for the exponentiation. Formulas in function files to be used by *quad* must have operations indicated entry-wise, so *quad* can evaluate the function on a vector $t = [t_0, t_1, \dots t_{fin}]$ of partition points it generates automatically. Read >>*doc quad* for explanation and more details.

```
function y = BasicGaussian(t)
% Your Name or Comments
y = exp(-t.^2);
end
```

Go back to the command line and call your function using >> *quad(@BasicGaussian,0,2)* . Notice your "ans" has only four places. To get more places use >>*format long* and run it again. □

The default error tolerance in *quad* is 1.0e-06 so your *ans* should have six good decimal places. Run >> *quad(@BasicGaussian,0,2, 1.0e-10)* and compare it with your previous result.

As we know, a first order linear ODE can always be put in the form:
$$\frac{dy}{dt} + p(t)y = f(t).$$

By Theorem 2.1 and its proof, such an equation, with *IC* $y(t_0) = y_0$, has solution:

$$y = e^{-\int_{t_0}^t p(s)ds} \left( \int_{t_0}^t e^{\int_{t_0}^s p(r)dr} f(s)ds + y_0 \right)$$

The integrals appearing in the solution may or may not be resolvable in terms of elementary functions. All of them can be found by numerical integration at a sequence of points $t_0, t_1, \dots t_{fin}$ leading up to any desired final time within the domain of continuity of $p$ and $f$. So all $1^{st}$ order linear equations with continuous $p$ and $f$ can be solved either exactly in terms of elementary functions or numerically by approximating the integrals at a sequence of points $t_0, t_1, \dots t_{fin}$. We should always remember that the exact solution is only a convenient formula with which to compute the value of the solution $y(t)$ at a sequence of time points. The formula used to find the values $y(t_0), y(t_1), \dots y(t_{fin})$ is really not the goal. We want to get good approximations to the solution at a specific sequence of points to determine the behavior of a physical system evolving forward in time from an initial condition.

## Exercises 1.8

1. Use a spreadsheet to calculate an approximation to the integral
   $y(2) = \int_0^2 e^{-s^2} ds$ using the trapezoid rule with 40 points.

2. Make a plot of $f(t) = \int_0^t e^{-s^2} ds$ for $0 \le t \le 2$. Use MatLab's quadrature function *quad* to get 10 good decimal places for $\int_0^2 e^{-s^2} dt$ .

3. Make a plot of $y = \sin(t^2)$ for $0 \le t \le \sqrt{\pi}$. Use MatLab's quadrature function *quad* to get 10 good decimal places for $\int_0^{\sqrt{\pi}} \sin(t^2) dt$ .

4. Solve the IVP $\frac{dy}{dt} + 2ty = \sin(t)$, $y(0) = 10$ and plot your solution for $0 \le t \le 4$.

5. Solve the IVP $\frac{dy}{dt} + 3t^2 y = \cos(t^2)$, $y(0) = 10$ and plot your solution for $0 \leq t \leq 5$.

# 1.9 Euler's Numerical Solution Method

By Theorem 2.1 and its proof, the solution to any 1$^{st}$ order linear IVP with continuous $p$ & $f$ can be written in terms of integrals and so solution values can be found either by simplifying the integrals to elementary functions or by numerical integration illustrated in 1.8. Separable equations can be integrated similarly. The general first order IVP

$$\frac{dy}{dt} = f(t,y), \qquad y(t_0) = y_0$$

may be of neither of these types.



Swiss postage stamp honoring Euler.

In such cases, we need a *numerical solution method*, a *numerical solver*. The most important conceptual tool for understanding numerical solvers is the slope-field discussed in the Introduction and illustrated in Example 0.6 and Exercises 1.5. The most basic numerical solver is due to Leonard Euler (1707-1783 Swiss)[15] and bears his name.

Euler's method works for any IVP, $\frac{dy}{dt} = f(t,y)$, $y(t_0) = y_0$. It starts from the IC $(t_0, y_0)$ and predicts a next point $(t_1, y_1)$ by following the line through $(t_0, y_0)$ at slope $f(t_0, y_0)$ out to a time $t_1$. Accuracy is maintained by keeping the time-step from $t_0$ to $t_1$ small. After the first step out to $(t_1, y_1)$ it uses that as a new IC and follows the line through it at slope $f(t_1, y_1)$ out to another time $t_2$. Figure 1.11 illustrates the idea.
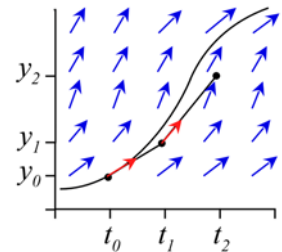


By repeating the process, an sequence of values, $y_0, y_1, \ldots y_{fin}$ are found for a sequence of times $t_0, t_1, \ldots t_{fin}$. The time steps are taken to be of equal length, $t_{i+1} - t_i = h$ for a small positive $h$, called the *step-size*. Taking smaller step-size $h$ gives a better approximation $y_{fin}$ to the exact $y(t_{fin})$, although more steps are needed to get there.

Figure 1-11. Slope-field arrows show slopes $y' = f(t,y)$. Red arrow at IC $(t_0, y_0)$ points to first Euler point $(t_1, y_1)$. Arrow at $(t_1, y_1)$ points to next Euler point $(t_2, y_2)$. Curve shows exact solution $y(t)$.

## Error estimation

Because Euler's method follows a straight line from $(t_0, y_0)$ to $(t_1, y_1)$, unless the exact solution is a straight line, the value Euler gives for $y_1$ will miss the target $y(t_1)$ on the solution curve, as depicted in Figure 1-11. If the solution has two continuous derivatives for $|t - t_0| < \delta$, then by Taylor's Theorem $y(t) \approx y(t_0) + y'(t_0)(t - t_0) + M(t - t_0)^2$ for some constant $M$ when $|t - t_0| < \delta$. Substituting $t = t_1$ and then $t_1 - t_0 = h$ into Taylor's estimate gives

$$y(t_1) \approx y(t_0) + y'(t_0)h + Mh^2 = y_1 + Mh^2$$

for $h < \delta$. This means, for small enough $h$, the error in one Euler step, $|y(t_1) - y_1|$, is proportional to $h^2$. However, when we take the second step, we are following a solution through $(t_1, y_1)$, a starting

---

[15] Euler, pronounced "Oiler", was arguably the most productive mathematician of all time. His published writings would fill over 60 large volumes of original Mathematics. (During his lifetime, Switzerland was known by its Latin name: *Helvetia*, appearing on the stamp.) Euler went blind in one eye but continued, even after losing sight in the other, to produce mathematical results until the day before he died.

point not on the solution with IC $(t_0, y_0)$ that we were trying to follow. We would expect the error to accumulate to a total of $n * Mh^2$, after $n$ Euler steps. Since the number of steps needed to reach $t_{fin}$ depends on the step size, as $n * h = t_{fin} - t_0$, so the error in going from $t_0$ to $t_{fin}$ accumulates to the *Global error*:

$$\left| y(t_{fin}) - y_{fin} \right| \approx n * Mh^2 = \frac{t_{fin} - t_0}{h}(Mh^2) = M(t_{fin} - t_0)h \,.$$

**Remark 9.0.** The *Global Error* $\left| y(t_{fin}) - y_{fin} \right|$ in Euler's method is proportional to $h$.

**Definition 9.0**. When a function $Q(h)$ is proportional to $h^p$ and $0 < \lim\limits_{h \to 0} \left| \frac{Q(h)}{h^p} \right| < \infty$, we say $Q(h) = O(h^p)$, read "$Q$ is big O of $h^p$" or "$Q$ is of the order of $h^p$ as $h$ goes to zero".

Using "Big O" notation, we have just proved:

**Theorem 9.0.** Euler's method global error $\left| y(t_{fin}) - y_{fin} \right| = O(h)$ as $h \to 0$.

By Theorem 9.0, cutting Euler's method step-size in half will cut the global error $\left| y(t_{fin}) - y_{fin} \right|$ in half (approximately).

Example 9.1. We illustrate Euler's method, as Euler would have done it by hand, for the IVP:

$$\frac{dy}{dt} = t^2 - y^3 \,, \qquad y(0) = 1.$$

Let's approximate $y(2)$ by taking four steps, size $h = 0.5$. It's best to be organized, so we make a table in a spreadsheet.

| step= | 0.5 | Euler | ydot=t^2-y^3 |
|---|---|---|---|
| k | t_k | y_k | ydot_k |
| 0 | 0 | 1 | -1 |
| 1 | 0.5 | 0.5 | 0.125 |
| 2 | 1 | 0.5625 | 0.822021484 |
| 3 | 1.5 | 0.9735107 | 1.327381318 |
| 4 | 2 | 1.6372014 | -0.388401179 |

The formulas are as follows:

$t_0 = 0, \quad y_0 = 1, \qquad y_0' = t_0^2 - y_0^3 = 0^2 - 1^3 = -1\,;$
$t_1 = 0.5, \quad y_1 = y_0 + y_0'h = 1 + (-1)(0.5)\,;$
$t_2 = 1, \quad y_2 = y_1 + y_1'h = 0.5 + (0.125)(0.5)\,;$
$t_3 = 1.5, \quad y_3 = y_2 + y_2'h = 0.5625 + (0.82202)(0.5)\,;$
$t_4 = 2, \quad y_4 = y_3 + y_3'h = 0.9735 + (1.3274)(0.5)$      □

In Example 9.1, $h = 0.5 \Rightarrow y_4 = 1.6372$. If we reduce the step-size to $h = 0.25$, we would expect to make half the global error.

Example 9.2. Euler's method using $h = 0.25$ for $\frac{dy}{dt} = t^2 - y^3$, $y(0) = 1$ gives $y(2) \approx \tilde{y}_8 = 2.695$.

| step= | 0.25 | Euler | ydot=t^2-y^3 |
|---|---|---|---|
| k | t_k | y_k | ydot_k |
| 0 | 0 | 1 | -1 |
| 1 | 0.25 | 0.75 | -0.359375 |
| 2 | 0.5 | 0.6601563 | -0.037700236 |
| 3 | 0.75 | 0.6507312 | 0.286947172 |
| 4 | 1 | 0.722468 | 0.62290062 |
| 5 | 1.25 | 0.8781931 | 0.885217086 |
| 6 | 1.5 | 1.0994974 | 0.920823566 |
| 7 | 1.75 | 1.3297033 | 0.711437135 |
| 8 | 2 | 1.5075626 | 0.573694745 |

Leaving the formulas to an exercise, we see that there is a dramatic difference between the predictions for $y(2)$. Following remark 9.0, we suspect that the global error for $h = 0.25$ is half of that in Example 9.1. Let's find a crude global error estimate.

$y_4 = 1.6372$ from Example 9.1 and the result in this example be $\tilde{y}_8 = 1.5076$. Let $\left| y(2) - y_4 \right| < \varepsilon$ and suppose $\left| y(2) - \tilde{y}_8 \right| < \varepsilon/2$ , then the worst case would be that $\left| \tilde{y}_8 - y_4 \right| \approx 3\varepsilon/2$, when $y(2)$ lies between $\tilde{y}_8$ and $y_4$.

Putting in the value $|\tilde{y}_8 - y_4| = |1.5076 - 1.6372| \approx 0.13$ gives a worst-case global error estimate: $\varepsilon \approx \frac{2}{3}(0.13) \approx 0.09$ , for $h = 0.5$ and $0.045$ when $h = 0.25$. Likely, both Euler approximations are poor because the step-sizes are too big. However, we wouldn't be sure how poor, until we tried an even smaller step-size.                                                   □

### Exercises 1.9

1. Construct a spreadsheet that finds an Euler approximation $\tilde{y}_{16}$ to $y(2)$ using step-size $h = 0.125$ for the IVP $\frac{dy}{dt} = t^2 - y^3$ , $y(0) = 1$. Plot the result in the spreadsheet.

2. Construct a spreadsheet that finds an Euler approximation $\tilde{y}_{50}$ to $y(5)$ using step-size $h = 0.1$ for the IVP $\frac{dy}{dt} = y \sin^2 t$ , $y(0) = 2$. Plot the result in the spreadsheet.

3. Write a MatLab function M-file that finds the Euler approximation to y(2) using $h = 0.0625$ for the IVP $\frac{dy}{dt} = t^2 - y^3$ , $y(0) = 1$. (Hint: Mildly challenging. You may nest a function inside that calculates dydt and call it as >>Euler(t0,tfin,y0,h). Also see doc for, doc if.)

# 1.10 Runge-Kutta Methods

Reducing step size in Euler's method will cut the global error, however, as Theorem 9.0 shows, making the error very small will require very small steps. Runge-Kutta methods greatly improve this situation; the classic Runge-Kutta method, RK4, offers a clever, if not elegant, solution to reducing the global error estimate to $O(h^4)$. Developed in the 1890s by German mathematicians Carl Runge and Martin Kutta, *Runge-Kutta methods* construct an average slope $\mu_i$ at each step that better approximates the correct slope needed to make an accurate Euler step to $(t_{i+1}, y_{i+1})$. These methods are fairly complicated to construct, so it is important to start small. We begin with a two-stage method, RK2, using two function evaluations per step that has global error $O(h^2)$ to help fix ideas.

## RK2 –" Midpoint" method

We will take one step forward from $(t_i, y_i)$ to $(t_{i+1}, y_{i+1})$ for $\frac{dy}{dt} = f(t, y)$ using stepsize $h$.

Before the step, we use the slope $k_1 = f(t_i, y_i)$ at the start point to take an Euler half-step to $t_i + \frac{h}{2}$ .

We evaluate $f$ there, getting the slope:  $k_2 = f\left(t_i + \frac{h}{2}, y_i + \frac{h}{2}k_1\right)$.

Next, take a full Euler step out to $t_{i+1}$ using slope $\mu_i = (k_1 + k_2)/2$ to get the RK2 approximation

$$y_{i+1} = y_i + \mu_i h = y_i + \left(\frac{k_1 + k_2}{2}\right)h .$$

Slope $k_1$ is at the at the start point and $k_2$ is at the Euler half-step. The idea is that the average of these should be better than just using $k_1 = f(t_i, y_i)$ as an approximation to the slope that points to the exact solution at $t_{i+1}$. Figure 1-12 illustrates the situation. In RK2 (midpoint), global error is $O(h^2)$, so it is called a *second order* method.



Figure 1-12. RK2 midpoint method. Green arrow slope used in the full Euler shot to $(t_{i+1}, y_{i+1})$ is the average of slopes $k_1$ and $k_2$ .

This type of Runge-Kutta method is called *explicit*. Explicit RK methods take an Euler step of size $h$, using an average slope, $\mu_i = b_1 k_1 + \cdots + b_s k_s$ to find
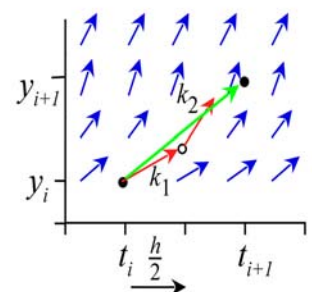
$y_{i+1} = y_i + \mu_i h$. The different slopes $k_1 \dots k_s$ are found successively. First $k_1 = f(t_i, y_i)$. Then a short Euler step, of size $c_2 h$ with $c_2 < 1$ using slope $k_1$ gives a point $(t_i + c_2 h, y_k + k_1 c_2 h)$ partway out to $t_{i+1}$ at which the second slope

$$k_2 = f(t_i + c_2 h, y_i + k_1 c_2 h).$$

Then another short Euler step size $c_3 h$ with $c_3 < 1$ and slope $a_{31} k_1 + a_{32} k_2$ gives another point:
$$(t_i + c_3 h, y_i + (a_{31} k_1 + a_{32} k_2) c_3 h)$$

At which the third slope is $k_3 = f(t_i + c_3 h, y_i + (a_{31} k_1 + a_{32} k_2) c_3 h)$. The number of slopes $k_1 \dots k_s$ found in this way and used to form $\mu_i$ is called the number of *stages* s. When all the slopes are found, a full Euler step finds $y_{i+1} = y_i + \mu_i h = y_i + (b_1 k_1 + \cdots + b_s k_s) h$ .

A Runge-Kutta method is characterized by: the number of stages $s$, the vector of fractions $c_2 \dots c_s$ for the short step sizes, the matrix of fractions $a_{nj}$ used to combine the previous $k_1, \dots k_{n-1}$ to get the slope for the (next) $n^{th}$ short step, and the vector of coefficients $b_1 \dots b_s$ used to construct $\mu_i$.

**Definition 10.1.** A Runge-Kutta method is a $p^{th}$ order method when its global error is $O(h^p)$.

**Property 10.0.** The number of stages $s$ and order $p$ of a Runge-Kutta method satisfy[16]

$$p < 5 \Rightarrow s \geq p \quad and \; 5 \leq p < 7 \Rightarrow s \geq p + 1 \; and \; p = 7 \Rightarrow s \geq p + 2 \; and \; 8 \leq p \Rightarrow s \geq p + 3.$$

We have already seen a version of RK2 (called Midpoint) with $s = p = 2$ and mentioned the classic RK4 with $s = p = 4$ , so for $p < 5$ there are Runge-Kutta methods known with $s = p$.

Example 10.1 RK3 – A third order method with three stages that has $s = p = 3$.

Let's take one step forward from $(t_i, y_i)$ to $(t_{i+1}, y_{i+1})$ for $\frac{dy}{dt} = f(t, y)$. We need to specify:

$(c_2, c_3) = \left(\frac{1}{2}, 1\right)$ and $a_{31} = -1$, $a_{32} = 2$ , which are good enough choices to give a 3$^{rd}$ order method.[17] Last, we need to pick the coefficients $b_1, b_2, b_3$. It seems reasonable to put more weight on the middle slope because it has solution curve on both sides in the interval $[t_i, t_{i+1}]$, while because the slopes $k_1$ and $k_3$ are found at the endpoints, they only represent the solution curve on one side in the interval $[t_i, t_{i+1}]$. It turns out $(b_1, b_2, b_3) = \left(\frac{1}{6}, \frac{2}{3}, \frac{1}{6}\right)$ will ensure $p = 3$. We summarize these choices:

$$\vec{c} = \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/2 \\ 1 \end{bmatrix}, \qquad \vec{k} = \begin{bmatrix} f(t_i, y_i) \\ f\left(t_i + \left(\frac{1}{2}h\right), y_i + \left(\frac{1}{2}h\right)k_1\right) \\ f(t_i + h, y_i + h(-k_1 + 2k_2)) \end{bmatrix}, \vec{b} = \begin{bmatrix} 1/6 \\ 2/3 \\ 1/6 \end{bmatrix}$$

Our RK3 has:

$$y_{i+1} = y_i + (b_1 k_1 + b_2 k_2 + b_3 k_3) h = y_i + \left(\frac{k_1 + 4k_2 + k_3}{6}\right) h . \qquad \square$$

## Implementation of error estimation and control

When using Runge-Kutta methods in practice, implementations in software like MatLab® use error estimation to control the growth of global error. The most common method compares paired RK methods

---

[16] Butcher, JC. *Numerical Methods for Ordinary Differential Equations*. New York. Wiley 2016.( p. 202)
[17] When combined with our choices for coefficients $b$. Generally, $a_{31} + a_{32} = c_3$ will do. Butcher, p.99.

that share function evaluations, but use different $a_{nj}$. When the two methods running together give values $y_{i+1}$ and $\tilde{y}_{i+1}$ that differ by more than a prescribed tolerance, the solver discards the step, returns to $t_i$ and tries to make a step using a smaller step size. MatLab® implements error estimation and control in all of its RK methods. The standard method, usually the default choice, is ode45, which uses a combination of an RK4 and an RK5 that reuses the function evaluations. The RK5 is the more accurate and when the RK4 solution is outside of the error tolerance, the step-size is reduced. If the retry with the smaller step fails, the step size is reduced again. Once the solver can proceed, after several steps result in sufficiently small errors, the step-size is allowed to increase. This type of solver adaptively implements *variable step size* to control error. Another of MatLab's solvers is ode23, which combines RK2 & RK3 to implement variable step-size error control and run faster than ode45. MatLab's documentation offers advice to the user as to which of its ode solvers might be best in different situations.

## Solving IVPs with MatLab® ode solvers

The ODE solvers take the function dydt as an external function, usually written in an M-file, called by its function handle as an input argument to the function.

Example 10.2. We solve the IVP $\frac{dy}{dt} = t^2 - y^3$, $y(0) = 1$ using ode45 by first writing a function M-file here,"myf.m" that ode45 can accept through its handle "@myf". Then we call ode45 on the command line with >> ode45(@myf,[0 2],1) and a figure opens[18] (Fig. 1-13) showing the solution.

```
function dydt=myf(t,y)
dydt=t.^2-y.^3;
end
```

To see the approximation to $y(2)$, call ode45 again with the syntax (using long format output)

```
>> [t,y45]=ode45(@myf,[0 2],1);
>> y45(end)
   ans =
     1.498062975349147
```



ode45 sol'n to dy/dt=t$^2$-y$^3$, y(0)=1

Only the first 5 or 6 digits are reliable.

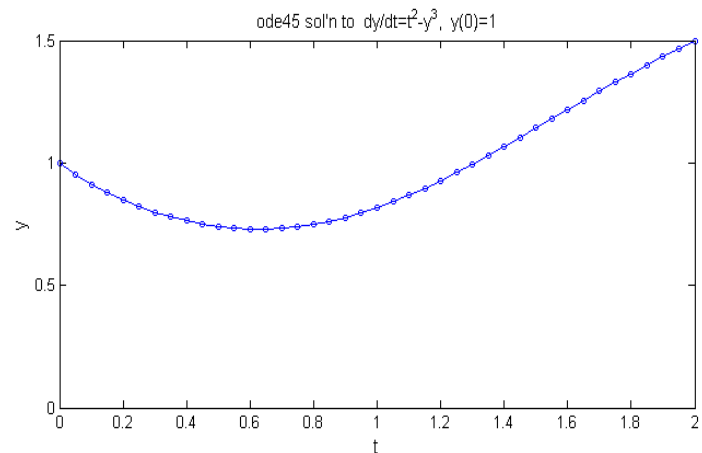To increase accuracy we need to modify the default error tolerances.                □

*Figure 1-13. MatLab ode45(@myf,[0 2],1) generated plot of solution to the IVP of Example 10.2 with $dy/dt = t^2 - y^3$ defined in the function myf.m.*

## Error Control

In adaptive solvers step-size is automatically set for the $n^{th}$ step by comparing results of two approximations $\tilde{y}_{n+1}$ and $y_{n+1}$ to estimate the local (per step) error $\varepsilon = \left| \tilde{y}_{n+1} - y_{n+1} \right|$, typically called "*abserr*". The tolerance *AbsTol* for *abserr* is the maximum value this $\varepsilon$ is allowed at any step, so when a step returns with $\varepsilon > AbsTol$, the step-size is reduced and the step retried until step-size is small enough to make $\varepsilon < AbsTol$.

A second control on error is the relative error tolerance *RelTol*. The relative error is

---

[18] Line and marker modifications, axis labels and title were added in Figure 1-13 using the MatLab figure editor.

$$relerr = \frac{\varepsilon}{|y_{n+1}|}.$$

By using both controls, *AbsTol* and *RelTol*, the solver adjusts step-size to make sure each step respects these tolerances. By default, in Matlab's explicit solvers, *AbsTol* = 1.0e-6 and *RelTol* = 1.0e-3. Step-size is reduced when either of these tolerances is violated. Roughly, the default local accuracy is *RelTol* <u>digits</u>, but as the solution approaches zero, when it becomes less than *AbsTol*, the step size will increase as long as the solution meets tolerances.

For small solutions, *AbsTol* may need to be reduced to maintain accuracy. In Example 10.2, the solution appears to be moderate in size, so we can improve accuracy by adjusting *RelTol*.

## Using *odeset* to modify tolerances

The MatLab command *odeset* is used to build an *options structure* that can be passed to any of the ode solvers.[19]

Example 10.3. We return to the IVP of Example 10.2 but run it with *RelTol*=1.0e-4.

```
>> optRel_4=odeset('RelTol',1.0e-4);
>> [t,y45_rel_4]=ode45(@myf,[0 2],1,optRel_4);
yRel_4=y45_rel_4(end)
yRel_4 =
   1.498074712054762
```

Comparing the result of Example 10.2:

```
yRel_3 =
    1.498062975349147
```

They agree in four places, so we suspect the tolerance was more than met with *RelTol* =1.0e-3 giving four good digits.                    □

## Exercises 1.10

1. Repeat the calculation of the solution to the IVP of Example 10.2 for [t0 tfin] = [0 10] using default *RelTol*=1.0e-3 with MatLab's solver ode23 called with no output arguments to get a plot. (Hint: Make the function file then use the same command syntax but use "23" instead of "45".)

2. Call ode23 with output arguments [t23, y23_3] and ode45 with output arguments [t45, y45_3]. Compare y23_3(end) with y45_3(end). How many more correct digits are in y45 than y23?

3. Repeat the calculation in `>>format long` using MatLab's solver ode45 for the solution to the IVP of Example 10.2 with *RelTol*=1.0e-*d* for *d* = 5,6,7,8 to see what happens. Show your code and briefly explain what you see in terms of correct decimal digits as *RelTol* is reduced.

4. Use ode23 and ode45 to solve the IVP $\frac{dy}{dt} = y \sin^2 t$ , $y(0) = 2$ for [t0 tfin] = [0 5] with output arguments [t23, y23] & [t45, y45] to compare y23(end) with y45 (end). Use `>>format long`. How many more correct digits in y45 than in y23? . Plot the solutions together.

5. Compare your results in Problem 4 with the results in Exercises 1.1 Problem 7. You will need to calculate the exact solution value $y(5)$ found in that earlier exercise. What are the actual errors?

6. Use ode23 and ode45 to solve the IVP $\frac{dy}{dt} = y \sin(t^2)$ , $y(0) = 2$ for [t0 tfin] = [0 20] with output arguments [t23, y23] & [t45, y45] to compare y23(end) with y45 (end). Use `>>format long`. How many correct digits are in y45? Plot the solutions together. Which solution is more accurate and why do you think so?

---

[19] Run doc odeset and read through for details. It is perhaps the most important page for users of the ode solvers.

7. Do Problem 6 using *RelTol* = 1.0e-4 and again plot the solutions together. Which solution is more accurate and why do you think so? Type up a brief original description of the problem of numerical instability as seen in these two problems and propose an approach that you would use to confirm the accuracy of a numerical solution in a situation where lives are at stake. Use < 250 words.

# 1.11 Exact Equations and Contour Plots

In this section, we follow the convention that uses $x$ instead of $t$ for the independent variable. This should cause no problem as, in practice, it could have almost any name and the literature on exact equations is consistent in the use of $x$.

In some equations, solutions can only be found in *implicit form*. This means the solutions lie along the level curves $F(x, y) = c = F(x_0, y_0)$ for some function $F$.

Example 11.1.  The separable equation $\frac{dy}{dx} = -\frac{x}{y}$ gives rise to $\int y dy = -\int x dx \Rightarrow \frac{y^2}{2} = -\frac{x^2}{2} + c$,

implicitly defining a solution $y(x)$. Rewrite the implicit solution as $\frac{x^2+y^2}{2} = c$, which says $F(x, y) = x^2 + y^2$ is constant along solutions, so solutions lie along level curves of $F(x, y)$. Such a function is called an *integral*[20] of the differential $dF = x\, dx + y\, dy$, in this case            □

Once we know an integral $F$, we can plot its level curve $F(x, y) = F(x_0, y_0)$, called an *integral curve*, for an IC $y_0 = y(x_0)$.  Example 1.11 was a particularly simple integration to find the integral $F(x, y)$. Before going to more difficult ODEs, Example 11.2 shows how to use MatLab to plot integral curves using the MatLab function *contour*.

Example 11.2. The MatLab commands shown will produce Figure 1-14.

```
>> [x,y]=meshgrid(-2:0.1:2,-2:0.1:2);
     F = x.^2+y.^2;
>> [C,h]=contour(x,y,F,[0:4]);
>> set(h,'LineWidth',2);
>> set(h,'LineColor','k');
```



A *meshgrid* using steps of 0.1 makes an acceptably smooth plot for *contour*. The integral F takes the third coordinate position in *contour*. The contours for levels 0,1,2,3,4 are requested, any list of values for F levels can be used.            □
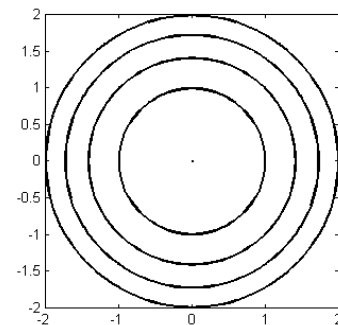
*Figure 1-14. Contour plot for Examples 11.1 & 11.2.*

The approach to exact equations begins by separating a differential equation $\frac{dy}{dx} = -\frac{M(x,y)}{N(x,y)}$ into its *differential form*:

$$M(x, y)\, dx + N(x, y)\, dy = 0$$

which is an *equation of differentials*.

---

[20] In this example, the level curves are circles centered at the origin. Note any function $F(x, y)$ with the same family of level curves would be an integral, only the values of $c = F(x_0, y_0)$ would be different.

Any separable equation can be implicitly solved by writing it in differential form:

$$\frac{dy}{dx} = a(x)b(y) \;\rightarrow\; a(x)dx - \frac{dy}{b(y)} = 0 \,.$$

Integrating the differentials, $\int a(x)dx + \int \frac{dy}{b(y)} = A(x) + B(y) = c$ (where $B(y) = -\int \frac{dy}{b(y)}$), implies $F(x,y) = A(x) + B(y)$ is an integral along whose level curves the solutions lie. It is more interesting that this approach generalizes to equations that are not separable. It won't implicitly solve all first order equations, only those that are *exact*.

**Definition 11.0.** A differential form $M(x,y)dx + N(x,y)dy$ is *exact* if there exists a function $F(x,y)$ such that: $\partial_x F(x,y) = M(x,y)$ and $\partial_y F(x,y) = N(x,y)$.

We call the form $Mdx + Ndy$ "exact" because it is the exact total differential of $F$:

$$dF = \frac{\partial F}{\partial x}dx + \frac{\partial F}{\partial y}dy = Mdx + Ndy \,.$$

The ODE $\frac{dy}{dx} = f(x,y)$ is called "exact" if $f(x,y) = -\frac{M(x,y)}{N(x,y)}$ and the differential form $Mdx + Ndy$ is exact.                                    □

Oddly, an equation $\frac{dy}{dx} = f(x,y)$ may not be exact, although multiplication of $f(x,y)$ by a fortuitously chosen form of the number $1 = \frac{\mu(x,y)}{\mu(x,y)}$ can sometimes result in an exact equation whose solutions are also solutions to $\frac{dy}{dx} = f(x,y)$. In such cases, the function $\mu(x,y)$ is called an integrating factor for the differential form of the equation. For this reason, the approach to exact equations is framed in terms of the differential form $M(x,y)dx + N(x,y)dy = 0$ of an ODE $\frac{dy}{dx} = -\frac{M(x,y)}{N(x,y)}$. Whereas, starting from an ODE, $\frac{dy}{dx} = f(x,y)$, it is not clear what function should play the role of $N(x,y)$ so as to determine $M(x,y) = -f(x,y)N(x,y)$. We leave this complication for later and revisit Example 11.1 to solidify ideas.

The ODE, $\frac{dy}{dx} = -\frac{x}{y}$ of Examples 11.1&2 has a differential form $xdx + ydy = 0$ having the integral:

$$F(x,y) = x^2 + y^2$$

whose level contours were plotted in Figure 1-14. Since the solutions lie in the level curves of $F$, the solution to the IVP with $y(x_0) = y_0$ lies in the level curve of $F$ that passes through $(x_0, y_0)$. So, it lies along a circular arc, centered at the origin, passing through $(x_0, y_0)$. The solution is only a semi-circular arc because it cannot be extended through the $x$-axis, where $y = 0$ and the ODE becomes undefined. In this example, we can solve $F(x,y) = k$ for $y$, giving implicit solutions:

$$y = \pm\sqrt{k - x^2}$$

for each $k > x^2$ and the $\pm$ is determined by the sign of $y(x_0) = y_0 \neq 0$. The explicit solution to the IVP is:

$$y = \text{sgn}(y_0)\sqrt{k - x^2} \,.$$

The situation is depicted in Figure 1-14 with open circles showing limit sets on the $x$-axis,

$$\alpha(x_0, y_0) = \left(-\sqrt{k}, 0\right), \qquad \omega(x_0, y_0) = \left(\sqrt{k}, 0\right),$$

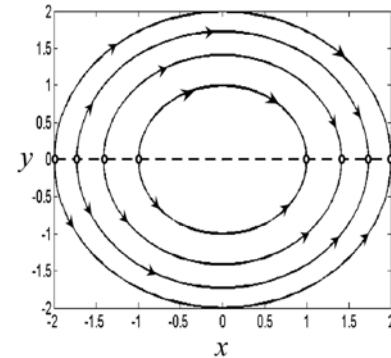where $k = r_0^2 = x_0^2 + y_0^2$ is the square of the radius of the circle.

*Figure 1-15. Circular contours containing semi-circular solutions to $y' = -x/y$. No solutions exist along the x-axis (dashed).*

The ODE, $\frac{dy}{dx} = -\frac{x}{y}$, whose solutions are shown in Figure 1-15, was separable. The next example shows a case, where it is not.

Example 11.3. The ODE $\frac{dy}{dx} = -\frac{x+y}{x+4y}$ has the differential form $(x + y)dx + (x + 4y)dy = 0$. This
   differential form is exact because there is an integral $F(x, y)$ with differential
$$dF = (x + y)dx + (x + 4y)dy.$$

   In other words,
$$\frac{\partial F}{\partial x} = x + y \quad \& \quad \frac{\partial F}{\partial y} = x + 4y.$$

   We don't know a formula for $F$ but let's suppose that there is such an $F$. Because its first partial
   derivatives are again differentiable (they are just polynomials here), we know from Calculus that
   mixed second partials taken in either order are equal:
$$\frac{\partial^2 F}{\partial y \partial x} = \frac{\partial^2 F}{\partial x \partial y}.$$
   So, by taking mixed partials, we know that if the mixed second partials are *not equal*, there is *no
   hope* to find a twice differentiable integral $F$. Let's check this to be sure:
$$\frac{\partial^2 F}{\partial y \partial x} = \frac{\partial}{\partial y}\frac{\partial F}{\partial x} = \frac{\partial}{\partial y}(x + y) = 1 \quad \& \quad \frac{\partial^2 F}{\partial x \partial y} = \frac{\partial}{\partial x}\frac{\partial F}{\partial y} = \frac{\partial}{\partial x}(x + 4y) = 1.$$

   They both equal 1, so there's hope. In fact, there's more than hope, there's a way to find an $F$. □


In Example 11.3, we can't "do" the integral $F = \int dF$ as in the previous example, because to do the
integral $\int_{x_0}^{x}(x + y)dx$ would require doing $\int_{x_0}^{x} y\, dx$ but we don't know $y$ as a function of $x$ until *after*
solving the ODE. Instead, because $x + y = \partial_x F$ is a *partial derivative*, we can use a *partial* integral in
which the value of $y$ is held constant, to find $F$.

**Definition 11.1.** The *partial* integral[21] $\int f(x,y)\partial x = \int f|_x(x,y)dx + a(y)$, where $f|_x(x,y)$ is $f(x,y)$
   with $y$ held constant and $a(y)$ is an arbitrary function.

Because $f|_x(x,y)$ is $f(x,y)$ with $y$ held constant, $f|_x(x,y)$ can be integrated as if it is a function of $x$
alone. The partial integral $\int f(x,y)\partial x$ gives every function $F(x,y) = \int f|_x(x,y)dx + a(y)$ for which
$\frac{\partial F}{\partial x} = f(x,y)$. So, $\int f(x,y)\partial x$ generalizes the idea of anti-differentiation to the partial $\frac{\partial}{\partial x}$.

Example 11.4. Let $1(x, y) = 1$ denote the constant function with value 1 at every $(x, y)$.

$$\int 1(x,y)\partial x = \int 1\, \partial x = x + a(y).$$

   Compare $\int 1|_x(x,y)dx = \int 1dx = x + c$, the ordinary anti-derivative of the constant function 1.
   The function $a(y)$ plays the role of the arbitrary constant in the partial integral with respect to $x$
   because $\frac{\partial a(y)}{\partial x} = 0$ so $a(y)$ behaves as if it were constant to $\frac{\partial}{\partial x}$. Different functions $a(y)$ give

---

[21] Really just finding partial anti-derivatives.

different antiderivatives of $1(x, y)$, for example $x + y$ gives another. We write $\int \partial x = x + a(y)$ when it is understood that we are working with functions of $(x, y)$.[22]

Similarly and perhaps more instructively, letting $f(x, y) = x + y$ in Def. 11.1 gives the partial integral

$$\int x + y\, \partial x = \frac{x^2}{2} + xy + a(y),$$

in which $y$ is a variable held constant during the integration. Specifying $a(y) = 2y^2$ gives one partial anti-derivative $F(x, y) = \frac{x^2}{2} + xy + 2y^2$ for the function $f(x, y) = x + y$.

Compare the ordinary integral $\int (x + y)|_x\, dx = \frac{x^2}{2} + xy + c$ where $y$ is an unknown parameter and $c$ is an arbitrary constant. Specifying $c = 2$ gives $\frac{x^2}{2} + xy + 2$ as one of the infinitely many ordinary anti-derivatives for $(x + y)|_x$.                    □

Returning to the equation $\frac{dy}{dx} = -\frac{x+y}{x+4y}$ of Example 11.3, recalling $\frac{\partial F}{\partial x} = x + y$ & $\frac{\partial F}{\partial y} = x + 4y$, we use partial integration to find the integral $F$ for the differential

$$dF = (x + y)dx + (x + 4y)dy.$$

From Example 11.4, $\frac{x^2}{2} + xy + a(y)$ is the general partial anti-derivative of $x + y$.

$$F = \int \frac{\partial F}{\partial x} \partial x = \int x + y\, \partial x = \frac{x^2}{2} + xy + a(y).$$

Similarly,

$$F = \int \frac{\partial F}{\partial y} \partial y = \int x + 4y\, \partial y = xy + 2y^2 + b(x)$$

where $a(y)$ & $b(x)$ are arbitrary functions of the variables held constant in the integrals. But these two partial anti-derivatives both equal $F$, so they are equal to each other.

$$F = \frac{x^2}{2} + xy + a(y)$$

$$= xy + 2y^2 + b(x).$$

So,

$$F - xy = \frac{x^2}{2} + a(y) = 2y^2 + b(x)$$

The easiest solution would be $a(y) = 2y^2$ & $b(x) = \frac{x^2}{2}$, giving $F - xy = \frac{x^2}{2} + 2y^2$. Finally,

$$F = \frac{x^2}{2} + xy + 2y^2$$

and we found the total integral of the total differential by matching the arbitrary functions $a(y)$ and $b(x)$ to corresponding uni-variate functions in the different partial integrals.

The differential form we got from the ODE was $dF = (x + y)dx + (x + 4y)dy = 0$ so when $y(x)$ is a solution to the ODE, $dF(x, y(x)) = 0$. The integral of $dF(x, y(x))$ is just the line integral of the *zero*

---

[22] For functions of more variables, $\int 1(x, y, z)\partial x = x + a(y, z)$ or $\int 1(x, y, z)\partial y = x + a(x, z)$ etc.

function $0(x, y) = 0$ along the solution curve $(x, y(x))$ from $x_0$ to $x$, which is zero, of course. We may write

$$\int_{x_0}^{x} dF(x, y(x)) = \int_{x_0}^{x} 0 \, dx = 0$$

and because $\int_{x_0}^{x} dF(x, y(x)) = F(x, y(x)) - F(x_0, y_0)$ for $y(x_0) = y_0$, we can substitute the integral $F(x, y)$ found above into $F(x, y(x)) - F(x_0, y_0) = 0$ to write (suppressing the dependence of $y$ on $x$):

$$F(x, y) - F(x_0, y_0) = \frac{x^2}{2} + xy + 2y^2 - F(x_0, y_0) = 0.$$

So, the final result is an implicit solution:

$$\frac{x^2}{2} + xy + 2y^2 = F(x_0, y_0)$$

and $F(x, y) = \frac{x^2}{2} + xy + 2y^2$ is constant along the solution $y(x)$ to the IVP $\frac{dy}{dx} = -\frac{x+y}{x+4y}$, $y(x_0) = y_0$.

The plot shown in Figure 1-16 shows level curves for the integral $F$. Solutions to the ODE

$$\frac{dy}{dx} = -\frac{x + y}{x + 4y}$$

lie along the level curves of $F = \frac{x^2}{2} + xy + 2y^2$. The expression $z_0 = \frac{x^2}{2} + xy + 2y^2$ gives a general implicit solution to the equation which, written for an IVP with IC $y(x_0) = y_0$, would have

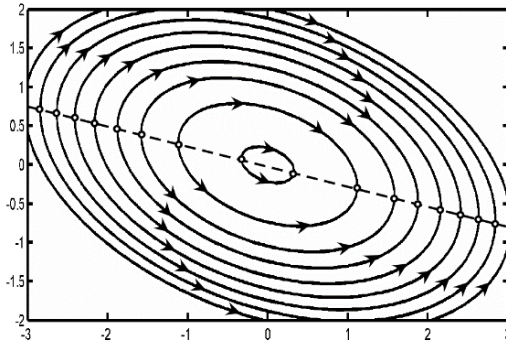$$z_0 = \frac{x_0^2}{2} + x_0 y_0 + 2y_0^2 = F(x_0, y_0).$$



Because $z_0 = \frac{x^2}{2} + xy + 2y^2$ is only quadratic in $y$, it can be solved for an explicit solution $y(x)$ using the quadratic formula $y = \frac{-b \pm \sqrt{b^2 - 4ac}}{2}$. Write the implicit solution as $0 = 2y^2 + xy + \frac{x^2}{2} - z_0$, so the standard form $ay^2 + by + c = 0 \Rightarrow a = 2, \; b = x, \; c = \frac{x^2}{2} - z_0$ are the coefficients $a, b, c$ for the quadratic formula, giving

$$y = \frac{-x \pm \sqrt{x^2 - 4(2)\left(\frac{x^2}{2} - z_0\right)}}{4} = \frac{-x \pm \sqrt{8z_0 - 3x^2}}{4}.$$

Substituting $\frac{x_0^2}{2} + x_0 y_0 + 2y_0^2$ for $z_0$ leads, after careful simplification to:

$$y = \frac{-x \pm \sqrt{4(x_0^2 - x^2 + 2x_0 y_0 + 4y_0^2) + x^2}}{4}$$

The $\pm$ sign is determined by the IC. To see how it is determined, solve for the indeterminate root:

$$\pm \sqrt{4(x_0^2 - x^2 + 2x_0 y_0 + 4y_0^2) + x^2} = 4y + x.$$

The positive root must be used for solutions where $4y + x > 0$, when the IC lies above the line $y = \frac{-x}{4}$, and the negative root for solutions below it. This line is dashed in Figure 1-16. Of course, solutions can only exist while $y \neq \frac{-x}{4}$ because $4y = -x$ makes the denominator in the ODE vanish. The direction of flow along level curves is left-to-right with increasing $x$ as shown by the arrows in Figure 1-16.

Generally, most forms $M(x,y)dx + N(x,y)dy$ are not exact, so we emphasize an easy way to check on exactness before we start integrating $M$ & $N$.

**Theorem:** $Mdx + Ndy$ with $M, N$ differentiable in $x$ and $y$, is exact $\Rightarrow \frac{\partial M}{\partial y} = \frac{\partial N}{\partial x}$.

Proof: Suppose $Mdx + Ndy$ is exact. By Definition 11.0, there is a differentiable function $F$ with $dF = Mdx + Ndy$, and $M = \frac{\partial F}{\partial x} = F_x \Rightarrow M_y = F_{xy}$ and similarly, $N = F_y \Rightarrow N_x = F_{yx}$. Since $F$ is twice differentiable, its mixed partials are equal, so: $M_y = F_{xy} = F_{yx} = N_x$.                    □

The way to remember the criterion is: Don't use $M$ & $N$, just always call them $F_x$ & $F_y$ . You know the mixed partials have to be equal, so when they aren't, it must not be an exact form. The following theorem is proved in virtually all Calculus textbooks where exact differentials are treated:

**Theorem:** $Mdx + Ndy$ is exact $\Leftrightarrow \frac{\partial M}{\partial y} = \frac{\partial N}{\partial x}$.

Now that we have the tools, let's work a few examples using them efficiently.

Example 11.5. Is the differential equation $\frac{dy}{dx} = -\frac{px+ry}{sx+qy}$ exact? If so, can we solve it?

The differential form is $(px + ry)dx + (sx + qy)dy = 0$, get the mixed partials:
$$\frac{\partial}{\partial y}(px + ry) = r, \quad \frac{\partial}{\partial x}(sx + qy) = s$$

These are only equal when $r = s$ so, although there are other methods to solve it[23], we can only use the exact approach when $r = s$.                    □

Example 11.6. Is the differential equation $\frac{dy}{dx} = -\frac{px+ry}{rx+qy}$ exact? If so, can we solve it?

The differential form is $(px + ry)dx + (rx + qy)dy = 0$ get the mixed partials:
$$\frac{\partial}{\partial y}(px + ry) = r, \quad \frac{\partial}{\partial x}(rx + qy) = r.$$

These are equal, so the form and the equation are exact. We can solve it implicitly by partial integrations.
$$F = \int px + ry \, \partial x = p\frac{x^2}{2} + rxy + a(y).$$

Similarly,
$$F = \int rx + qy \, \partial y = rxy + 2qy^2 + b(x).$$

They're equal, so matching univariate functions, $a(y) = 2qy^2$ and an integral is:
$$F = p\frac{x^2}{2} + rxy + 2qy^2.$$

A more convenient form without the fraction is
$$2F = z = px^2 + 2rxy + 4qy^2.$$

---

[23] Divide numerator and denominator by $x$ so $\frac{px+ry}{sx+qy} = \frac{p+ru}{s+qu}$ and make the substitution $ux = y$ differentiating by $x$ using the product rule leading to $\frac{du}{dx}x + u = -\frac{p+ru}{s+qu} \Rightarrow \frac{du}{dx}x = -\frac{qu^2+(s+r)u+p}{s+qu}$, which is ugly but separable.

Recall from Analytic Geometry that an equation of the form $z_0 = ax^2 + bxy + cy^2$ defines a quadratic curve of type determined by:

$b^2 - 4ac < 0 \Rightarrow$ elliptic; $b^2 - 4ac = 0 \Rightarrow$ parabolic; $b^2 - 4ac > 0 \Rightarrow$ hyperbolic. □

Example 11.7. The differential form $\frac{2y}{x} dx + \left(1 + \frac{y^2}{x^2}\right) dy = 0$ is not exact since

$$\frac{\partial}{\partial y}\left(\frac{2y}{x}\right) = \frac{2}{x} \neq \frac{\partial}{\partial x}\left(1 + \frac{y^2}{x^2}\right) = -2y^2 x^{-3}.$$

However, multiplying through by $\mu(x,y) = x^2$ gives a new form: $2xy\, dx + (x^2 + y^2)dy = 0$. This new and simpler form is exact because

$$\frac{\partial}{\partial y}(2xy) = 2x = \frac{\partial}{\partial x}(x^2 + y^2).$$

The factor $\mu(x,y) = x^2$ is an integrating factor for the original differential. Both differentials represent the same ODE:

$$\frac{dy}{dx} = \frac{-2xy}{x^2 + y^2}.$$

The first differential was not a natural choice, the second is better because it allows solution of the equation. Let's solve it using the exact differential form $2xy\, dx + (x^2 + y^2)dy = 0$.

$$F = \int 2xy\, \partial x = x^2 y + a(y)$$

$$= \int x^2 + y^2\, \partial y = x^2 y + \frac{y^3}{3} + b(x).$$

Evidently, $a(y) = \frac{y^3}{3}$ so $F = x^2 y + \frac{y^3}{3}$ is an integral. Solutions lie on the level contours of $F$ but this time we can't solve for $y$ explicitly because $F$ is cubic in $y$. We can still make the contour plot to see what the solutions look like.

Flow on contours in Figure 1-17 is left to right. □



The contours shown in Figure 1-17 provide $(x, y(x))$ data that numerically matches the solutions very closely. We can compare contours with numerical approximations found using ode45 or plots of exact solutions to validate results.
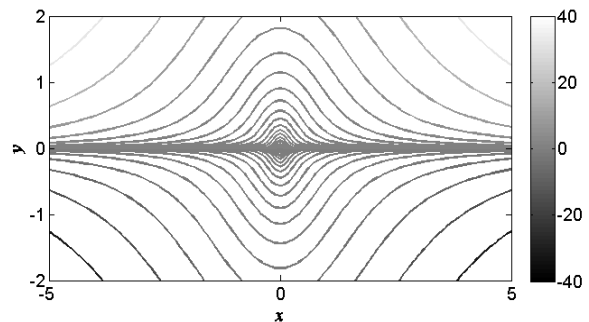
Figure1-17. Contours for $F = x^2y + y^3/3$ of Example 11.7. Levels of $z_0 = F(x_0,y_0)$ are indicated in the color bar on the right, keyed to the gray-tones of the contours.

## Exercises 1.11

1. For each differential equation, determine whether it is exact or not. Show why you reached your conclusion.
   i. $(2x + 3y)dx + (3x + 4y)dy = 0$   ii. $(x + 2y)dx + 2x\, dy = 0$   iii. $6x^2y^2 dx + 4x^3y\, dy = 0$
   iv. $xy^2 dx - x^2 y\, dy = 0$   v. $\frac{x\, dx}{x^2+y^2} + \frac{y\, dy}{x^2+y^2} = 0$   vi. $2x^2y^3 dx + 3x^2y^2 dy = 0$

2. If the differential is exact, find an integral. If not exact, show it is not exact.

   i. $(2x + y)dx + (x + 3y)dy = 0$   ii. $(x + y)^2 dx + (x + y)^2 dy = 0$   iii. $6x^2y^2 dx + 4x^3y\, dy = 0$

3. Write the ODE as an equation of differentials and if exact, find an integral $F$. Plot at least 5 contours of $F$ and draw arrows on the contours indicating the direction of flow.

   i. $\frac{dy}{dx} = \frac{x-2y}{2x+y}$   ii. $\frac{dy}{dx} = -\frac{y}{x}$   iii. $\frac{dy}{dx} = \frac{x+2y}{2x+y}$