

1:30

# Chapter 2

## Least Square Fitting

Let us assume that  $y = f_{c_1, \dots, c_k}(x_1, \dots, x_m)$  is a real-valued function of  $(x_1, \dots, x_m)$  which depends upon  $k$  parameters  $c_1, \dots, c_k$ . These parameters are unknown to us. However, suppose that we can perform repeated experiments that for given values of  $(x_1, \dots, x_m)$  allows us to measure output values for  $y$ . How can we estimate the parameter  $c_1, \dots, c_k$  that best correspond with this information?

Let us assume that the experiment measuring the value  $y = f_c(x)$  for specific input values  $x = (x_1, \dots, x_m)$  is repeated  $n$ -times. We will then obtain a system of equations

$$\begin{aligned} f_{c_1, \dots, c_k}(x_{11}, x_{12}, \dots, x_{1m}) &= y_1 \\ f_{c_1, \dots, c_k}(x_{21}, x_{22}, \dots, x_{2m}) &= y_2 \\ &\vdots \\ f_{c_1, \dots, c_k}(x_{n1}, x_{n2}, \dots, x_{nm}) &= y_n \end{aligned}$$

where  $x_{ij}$  and  $y_i$  are the measurements of  $x_j$  and  $y$  in the  $i$  experiment. These experimental results gives us information about the unknown coefficients  $c_i$ 's.

Since we may perform the experiments as many times as we wish, we may end up with more relations of the type above than unknowns (that is,  $n$  is larger than  $k$ ). The larger the  $n$ , the more information we have collected about the coefficients. However, even if the experiments are carried out with great care, they unavoidably will contain some error. The question remains: how may we estimate judiciously the coefficients  $c_1, \dots, c_k$  using the collected information about  $y = f_c(x)$ ? What is the best fit?

A very common method to respond to this question is known as the *method of least-squares*. The idea is simple: per realization of the experiment, we measure the fitting error by the distance from the real number  $f_c(x_1, x_2, \dots, x_m)$  and the observed value of  $y$ . The best fit for the distance though will also lead to a best fit for the square of the distance. To avoid absolute values, we change our viewpoint slightly and measure the fitting error by  $(f_{c_1, \dots, c_k}(x_1, x_2, \dots, x_m) - y)^2$ .

The error function that considers all the information obtained from the  $n$  experiments is then

$$E(c_1, \dots, c_k) = \sum_{i=1}^n ((f_{c_1, \dots, c_k}(x_{i1}, x_{i2}, \dots, x_{im}) - y_i)^2,$$

where  $(x_{i1}, x_{i2}, \dots, x_{im})$  and  $y_i$  are the  $i$ -th input for  $x$  and value of the measurement for  $y$ .

This turns out to be a function of  $c = (c_1, \dots, c_k)$ . Mathematically, our best fit problem is now reduced to finding the value of  $c$  which produces a minimum for this error function  $E$ . The details of how this can be done depends intrinsically upon the assumed form of the function  $f$ , and its relation to the parameters  $c_1, \dots, c_k$ .

## 1 Fitting a line to data

Suppose that we are in the situation described above, with  $m = 1$ , and that the function  $f$  is a polynomial of degree 1. Hence, the number of parameters is two, and for convenience we write

$$f(x) = mx + b,$$

in order to interpret the parameters  $m$  and  $b$  as the slope and  $y$ -intercept of the graph of  $f$ . We have data points  $(x_1, y_1), \dots, (x_n, y_n)$ . Our problem is to find the values of  $m$  and  $b$  which will best fit this data. Or to put it simply and geometrically, what is the straight line that best fits the information contained in the data  $(x_1, y_1), \dots, (x_n, y_n)$ ?

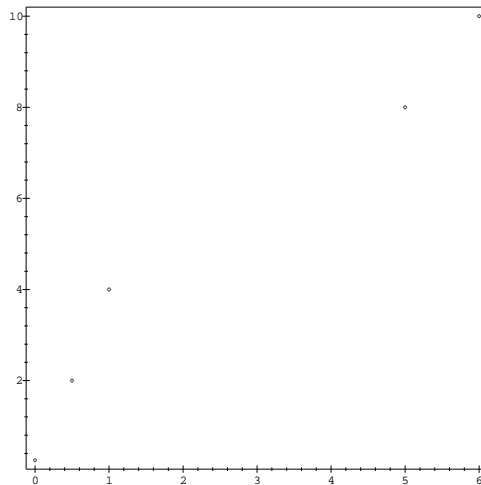
The contribution to the error function coming from the  $i$ -th piece of data,  $(x_i, y_i)$ , is given by  $(mx_i + b - y_i)^2$ , and therefore, the total error is

$$E(m, b) = \sum_{i=1}^n (mx_i + b - y_i)^2.$$

The problem then becomes that of searching for the values of  $m$  and  $b$  where the error  $E(m, b)$  achieves a minimum. Let us solve this problem using Maple.

First, let us type in some data to try to fit our line to:

```
> pts := [ [0,0.25], [1/2,2], [1,4], [5,8], [6,10] ]:
> plot(pts, style=point, axes=boxed);
```



In order to solve the problem at hand, we could use Maple's built-in regression package to find the answer, but that would not help much in explaining what is going on. However, let us do this anyway.

Since the least-squares package expects the  $x$  and  $y$  values in separate lists, we have to make a slight adjustment to our data:

```
> with(stats):
   fit[leastsquare][[x,y]]([[ pts[i][1] $i=1..5], [ pts[i][2] $i=1..5]]);
      y = 1.271370968 + 1.431451613x
```

This says that the best values of  $m$  and  $b$  are  $m = 1.431451613$  and  $b = 1.271370968$ , respectively. Let us do this *by hand*.

First, we define the error function according to the data points to fit that we have. Notice again that this function is the square of the distance from the line to the data:

```
> E := (m,b)-> sum( ((m*pts[i][1]+b)-pts[i][2] )^2, i=1..5);
```

$$E := (m, b) \rightarrow \sum_{i=1}^5 (m \text{ pts}_{i1} + b - \text{pts}_{i2})^2$$

For example, had we guessed that the line  $y = 2x + 1$  is the answer to our problem, we could compute the distance:

```
> E(2,1);
      19.5625
```

However, decreasing  $m$  and increasing  $b$  produces a smaller value of  $E$ , so that guess cannot be the best fit:

```
> E(1.5,1.2);
      3.125000000
```

We want to minimize the error function. Since minima of differentiable functions occur at critical points, we begin by solving for the values of  $(m, b)$  which annihilate the two partial derivatives. Notice that Maple happily computes partial derivatives:

```
> diff( E(m,b), m);
```

$$\frac{249}{2}m + 25b - 210$$

We may then ask Maple to solve the system of equations obtained by equating the two partial derivatives by

```
> solve( { diff( E(m,b), m)=0, diff( E(m,b), b)=0 }, {m,b});
      {m = 1.431451613, b = 1.271370968}
```

Maple thus finds only one solution, and not surprisingly, it coincides with the solution found using the statistics package.

If we now want to use this values of  $m$  and  $b$  without retyping it, one way to do that is to use the command `assign(%)` to let  $b$  and  $m$  be given these constant values.<sup>1</sup>

```
> assign(%);
> assign(%);
```

Maple executes this assignment silently, without reporting the assignment. However, now both  $m$  and  $b$  have the assigned values:

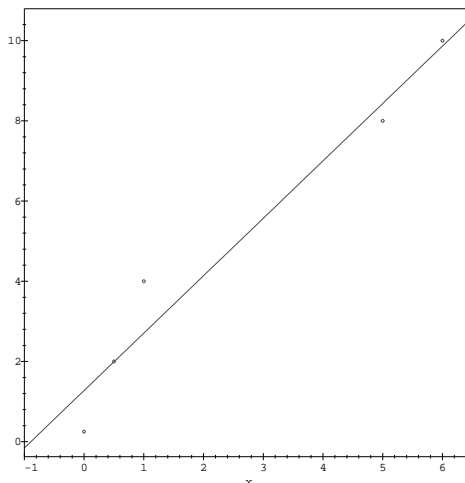
```
> m;
      1.431451613

> E(m,b);
      2.929334677
```

The value of  $E(m, b)$  calculated above, is the value of  $E$  when  $(m, b)$  is the critical point found earlier.

We may visualize how good our fit is, using `plot` to display the data points and the fit line on the same graph. We load the plots package, so that we can use `display`:

```
> with(plots):
lplot := plot(m*x+b, x=-1..6.5, axes=boxed):
pplot := plot(pts, style=point):
display({lplot, pplot});
```




---

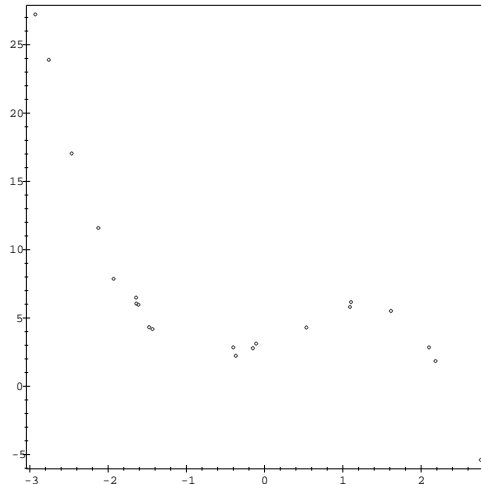
<sup>1</sup>We are using Maple's *ditto operator* `%` to refer to the result of the previous command. In versions of Maple prior to release 5 of MapleV, the double quote `"` was used for this operator instead.



Now we can put the data to be fitted into a list, and visualize the result with `plot`:

```
> data:=cubic_pts();
   data := [[2.100491206, 2.842523280], [-1.642321526, 6.490510559],
            [-2.930111988, 27.22116242], [1.613687708, 5.512132122],
            [-2.757747815, 23.89666247], [-.108617942, 3.124512438],
            [.533554657, 4.302312866], [-2.465340511, 17.04297836],
            [1.103974870, 6.158134050], [-1.613219905, 5.973490363],
            [-1.929997474, 7.866127538], [-.150544582, 2.786384816],
            [2.761660877, -5.391947131], [-1.433219263, 4.186307388],
            [-2.124410507, 11.59447629], [1.090163214, 5.804094012],
            [-.400537225, 2.845969065], [-.367977866, 2.228820732],
            [-1.639166049, 6.054218944], [-1.475967686, 4.324516109],
            [2.183284715, 1.845821478]]
```

```
> plot(data, style=point, axes=boxed);
```



Don't worry about the way we generated the data points. As far as the problem to solve is concerned, we have a set of 21 data points which are to be fitted to a cubic polynomial. A polynomial function of degree  $n$  is determined by  $n + 1$  coefficients. So we define

```
> cub := (x,a,b,c,d)->a*x^3+b*x^2+c*x+d;
      cub := (x, a, b, c, d) → ax3 + bx2 + cx + d
```

and then define the error function:

```
> E:=(a,b,c,d)->sum((cub(data[i][1],a,b,c,d)-data[i][2])^2,i=1..nops(data));
      E := (a, b, c, d) → ∑i=1nops(data) (cub(datai1, a, b, c, d) - datai2)2
```

We have four parameters to determine, the coefficients of the function `cub`. We find its values by

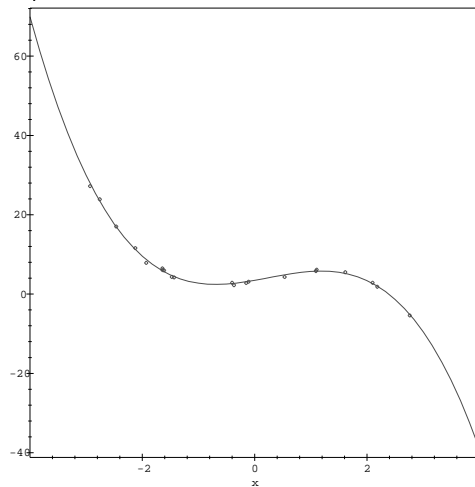
```
> assign(solve({diff(E(a,b,c,d),a)=0,
                diff(E(a,b,c,d),b)=0,
                diff(E(a,b,c,d),c)=0,
                diff(E(a,b,c,d),d)=0},
                {a,b,c,d}));
```

Of course, we do not see the answer, but the resulting values have already being assigned to the coefficients  $a$ ,  $b$ ,  $c$ ,  $d$ :

```
> cub(x,a,b,c,d);
-1.006141915x3 + .7505869176x2 + 2.487673553x + 3.430571969
```

Finally, let's make a picture of the result, showing both the data and the fitted curve:

```
> with(plots):
cplot := plot(cub(x,a,b,c,d), x=-4..4, axes=boxed):
pplot := plot(data, style=point):
display({cplot,pplot});
```



For comparison, let us find the answer using Maple's built-in statistical package:

```
> with(stats):
fit[leastsquare][[x,y],y=A*x^3+B*x^2+C*x+D]
([ [data[i][1]$i=1..nops(data)], [data[i][2]$i=1..nops(data)]]);
y = -1.006141912x3 + .7505869193x2 + 2.487673537x + 3.430571965
```

For this particular version of least square fitting, we may once again interpret the error function  $E$  geometrically as the square of the Euclidean distance between  $a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d}$  and the vector  $\vec{y}$ , where  $\vec{x}_j = (x_1^j, x_2^j, \dots, x_n^j)$ ,  $\vec{d} = d(1, 1, \dots, 1)$  and  $\vec{y} = (y_1, \dots, y_n)$ . The best fit is thus achieved when the coefficients  $a$ ,  $b$ ,  $c$  and  $d$  are chosen so that  $a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d}$  is the orthogonal projection of the vector  $\vec{y}$  onto the subspace of  $\mathbb{R}^n$  spanned by  $\vec{x}_3$ ,  $\vec{x}_2$ ,  $\vec{x}_1$  and  $(1, 1, \dots, 1)$ . That is to say, the best fit occurs when the vectors  $\vec{y} - (a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d})$  and  $a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d}$  are perpendicular. You should verify that this is the case for the solution given above.

In both cases, you should consider why there is only one critical point for the error function, and why it is of necessity a global minimum.

### 3 Fitting other types of functions

In the previous constructions, we dealt with functions that depended linearly upon the parameters to be fitted. We found them by solving a linear system of equations, a relatively easy task.

However, as explained at the beginning of this chapter, the same scheme works equally well for any function. In the general case, though, the resulting system that we must solve to find the best fit could depend non-linearly on the parameters, and the mere existence of solutions to such systems is not a trivial problem to settle. Take for instance, a function such as  $y = \sum_{i=1}^n \sin(m_i x)$ . Assuming measurements  $(x_1, y_1), \dots, (x_k, y_k)$ , we may take as our error the function

$$E(m_1, \dots, m_n) = \sum_{j=1}^k (y_j - \sum_{i=1}^n \sin(m_i x_j))^2.$$

Its partial derivatives are quite easy to calculate, but it is far from clear if there are values of  $m_1, \dots, m_n$  where they all vanish simultaneously. For a particular set of values, we may ask Maple to solve the resulting system and get absolutely nothing. This either indicates an inability of Maple to handle such a system, or worse yet, the fact that such a system has no solution at all. Even when the latter happens, Maple itself might not be able to tell us so.

Sometimes, even if the function does not depend linearly on the coefficients, it can be transformed to one which does. For example, if our data points  $\{(x_i, y_i)\}$  were believed to approximate an exponential function of the form  $y = ae^{kx}$ , then setting

$$E(a, k) = \sum_{i=1}^m (y_i - ae^{kx_i})^2.$$

would require us to solve the system

$$\sum_{i=1}^m (y_i - ae^{kx_i})e^{kx_i} = 0 \quad \sum_{i=1}^m (y_i - ae^{kx_i})ax_i e^{kx_i} = 0.$$

While this isn't impossible, it is much more straightforward to make a change of variables.

Assuming the  $y_i$  are all positive (which is reasonable since we believe the data is exponential), we can let  $z_i = \ln y_i$ . Then we are trying to fit a function  $z = \ln(ae^{kx})$ , which reduces to the line  $z = \ln a + kx$ . This is familiar territory, and we can just proceed as before.

## 4 Fitting a circle

In a slightly different vein, suppose that the data points  $\{(x_i, y_i)\}$  lie near a circle of unknown center and radius. Note that if we did know the center, this problem would be very simple: the desired radius would be the average distance from the points to the center. If you don't think very hard, you might suspect that the desired center would be very near the center of mass of the points  $\{(x_i, y_i)\}$ , but this will only be the case if the points are evenly distributed around the circle.

Instead, we can come up with a function that measures the error between an arbitrary circle and our data points.

Recall that the general equation of a circle centered at  $(a, b)$  with a radius of  $r$  is

$$(x - a)^2 + (y - b)^2 = r^2.$$

Unlike the previous cases, there is no independent variable. (Note that we *could* try to fit  $y_i \simeq b \pm \sqrt{r^2 - (x_i - a)^2}$ , but not only would the resulting equations be messy, this would bias things very badly; do you see why?). Nevertheless, we press on.

One reasonable measure of the distance between the points and a circle is the “area difference”, that is

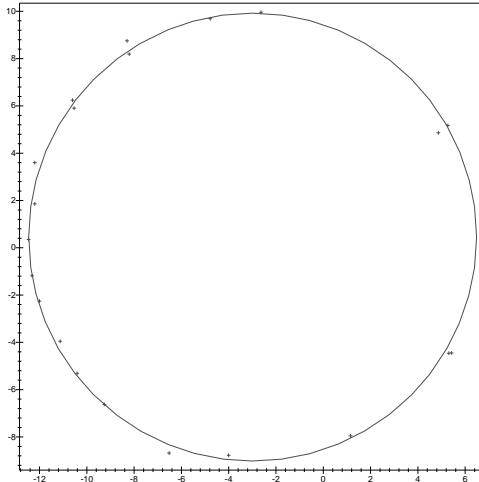
$$E(a, b, r) = \sum_i ((x_i - a)^2 + (y_i - b)^2 - r^2)^2.$$

One difficulty with this is that  $E$  is not quadratic in  $a$ ,  $b$ , and  $r$ . However, Maple is able to solve the resulting equations anyway:

```
> cpts:=circle_pts():
  epsilon:=(pt,a,b,r) -> (pt[1]-a)^2 + (pt[2]-b)^2 - r^2;
  ε := (pt, a, b, r) → (pt1 - a)2 + (pt2 - b)2 - r2
> E:= (a,b,r) -> sum( epsilon(cpts[i],a,b,r)^2, i=1..nops(cpts));
  E := (a, b, r) → ∑i=1nops(cpts) ε(cptsi, a, b, r)2
> sol:= solve({diff(E(a,b,r),a)=0,
  diff(E(a,b,r),b)=0,
  diff(E(a,b,r),r)=0},
  {a,b,r});
sol := {a = -7.279165053 - 12.95193570 I, b = -2.210665541 + 2.912161761 I, r = 0},
  {a = -7.279165053 + 12.95193570 I, b = -2.210665541 - 2.912161761 I, r = 0},
  {b = .4328495120, a = -4.481107455, r = 0},
  {b = -14.49583861 - 13.00051962 I, a = 4.676168966 - 18.08449394 I, r = 0},
  {a = 4.676168966 + 18.08449394 I, b = -14.49583861 + 13.00051962 I, r = 0},
  {r = 9.471322781, a = -2.987792224, b = .4467421770},
  {r = -9.471322781, a = -2.987792224, b = .4467421770}
```

Here we find a number of critical points for the function  $E$ , but from physical considerations, the only reasonable choice is the circle with center at  $(-2.987792224, .4467421770)$  and radius 9.471322781. In order to see the fit, we plot the points and the circle. Note that we represent the circle parametrically, and use `subs` to substitute the desired solution.

```
> pp:=plot(cpts,style=point,scaling=constrained,axes=boxed):
  cp:=plot(subs(sol[6],[a+r*cos(t),b+r*sin(t),t=0..2*Pi])):
  plots[display](pp,cp);
```



The fact that there is only one interesting critical point is no accident, however. If we let  $k = a^2 + b^2 - r^2$ , then the resulting error function  $H(a, b, k) = E(a, b, \sqrt{a^2 + b^2 - k})$  is quadratic in  $a$ ,  $b$ , and  $k$ , and so we really only need to solve a linear system. It is easy to check that this new functional  $H$  still satisfies all the criteria we wanted (that is  $H(a, b, k) = 0$  if and only if all the data points lie on the circle  $C(a, b, k)$ ,  $H$  is non-negative, and it is smooth), so fitting a circle becomes a linear problem after all. The reader should verify that, in fact, the unique minimum of  $H$  corresponds exactly to the critical points the original function  $E$  for which  $r \neq 0$ .

The interested reader might want to consider a similar approach to fitting other conic sections, such as an ellipse or a hyperbola, to given data. Do you expect to be able to make the problem linear, as in the case of the circle?

## 5 Robust fitting

Let us return momentarily to the problem discussed in §1: given some data  $(x_1, y_1), \dots, (x_n, y_n)$ , we found the best line that fits it. This was accomplished by measuring the square of the *vertical* distance from each data point and the line  $y = mx + b$ , which led us to consider the error function

$$E(m, b) = \sum_{i=1}^n (mx_i + b - y_i)^2.$$

The problem was solved by finding the values of  $m$  and  $b$  where  $E$  achieves its minimum.

However, the square of the vertical distance is only one of many reasonable ways of measuring the distance from the data points to the line  $mx + b$ . For example, it is perfectly reasonable to consider instead the expression  $1 + (mx_i + b - y_i)^2$ . If  $(x_i, y_i)$  is on the line, this value would be equal to 1, and its logarithm would then be zero. Thus, the function

$$E(m, b) = \sum_{i=1}^n \ln(1 + (mx_i + b - y_i)^2),$$

is also a reasonable way of measuring the distance from the data points to the line  $y = mx + b$ , and its minimum would produce a best fit line in this other sense. Similarly, We may argue that the function

$$E(m, b) = \sum_{i=1}^n |mx_i + b - y_i|.$$

is another reasonable alternative, though this time  $E$  is not even differentiable in the region of interest (if a data point  $(x_i, y_i)$  happens to be exactly on the line,  $mx_i + b - y_i = 0$  and the absolute value is not differentiable at 0).

There is a priori no particular reason to prefer one error function over another. And for each choice of such we could end up with a problem whose solution exists and produces a line that best fit the data. Then we could compute the value of the error function for the best fit line, value that depends upon the choice of function made. A *canonical* way of choosing the best error function could be that for which this number is the smallest. But finding such a function is quite a difficult, if not impossible, problem to solve. If we limit our attention to *linear estimators*, in a sense to be explained in the last section of this chapter, the solution obtained in §1 is optimal. In here, we pursue further the problem of finding the best fit according to the two error functions introduced above.

Notice that if  $\epsilon_i = mx_i + b - y_i$ , data with large errors  $\epsilon_i$  exert a huge influence on the original error function of §1 (because we use  $\sum \epsilon_i^2$ ). The two error functions mentioned above grow linearly (or sub-linearly) with  $|\epsilon_i|$ , causing the tails to count much less heavily. We will see this when comparing the different results.

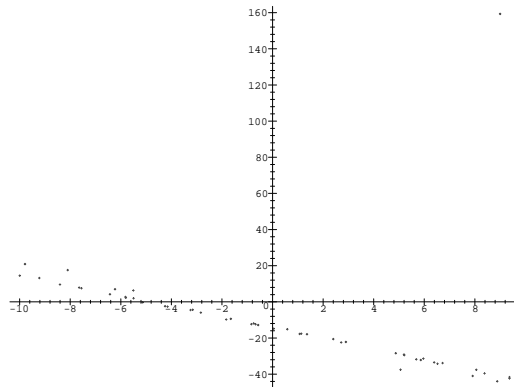
In order to solve the problem now at hand, we first generate some “noisy” data. We load the routines in `lsq_data.txt`:

```
> read('lsq_data.txt');

defined line\_pts(), bad\_line\_pts(), quadratic\_pts(), cubic\_pts(),
and circle\_pts()
```

The data and its graphical visualization can then be obtained by:

```
> pts:=bad_line_pts();
> plot(pts,style=point);
```



```
> epsilon:= (pt,m,b) -> (m*pt[1]+b-pt[2]);
```

$$\varepsilon := (pt, m, b) \rightarrow m pt_1 + b - pt_2$$

In order to compare results, let us first find the line given by least squares:

```
> H:=(m,b)->sum(epsilon(pts[i],m,b)^2,i=1..nops(pts));
```

$$H := (m, b) \rightarrow \sum_{i=1}^{\text{nops}(pts)} \varepsilon(pts_i, m, b)^2$$

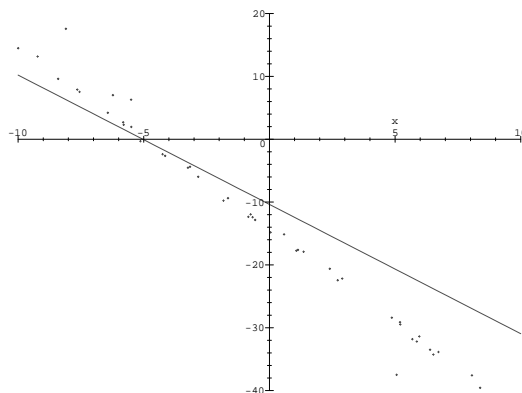
```
> sol:=solve(\{diff(H(m,b),m)=0, diff(H(m,b),b)=0\},\{m,b\});
```

$$sol := \{m = -2.058151695, b = -10.36674445\}$$

```
> p := plot(pts,style=point):
```

```
l := plot(subs(sol,m*x+b),x=-10..10):
```

```
plots[display] (p,l,view=-40..20);
```



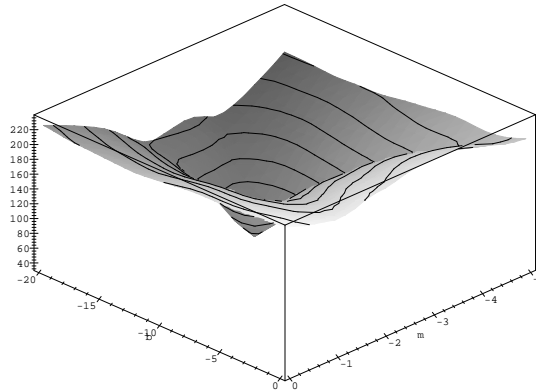
Now, let's try it minimizing  $\ln(1 + \varepsilon_i^2)/2$ , which behaves like  $\varepsilon^2$  for small errors, but grows very slowly for large ones.

```
> R:=(m,b)->sum(ln(1+epsilon(pts[i],m,b)^2/2),i=1..nops(pts));
```

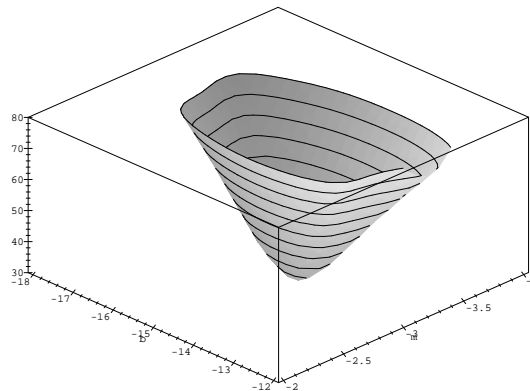
$$R := (m, b) \rightarrow \sum_{i=1}^{\text{nops}(pts)} \ln\left(1 + \frac{1}{2} \varepsilon(pts_i, m, b)^2\right)$$

Here's what the functional we want to minimize looks like:

```
> plot3d(R(m,b),m=-5..0,b=-20..0,style=patchcontour, axes=boxed);
```



```
> plot3d(R(m,b),m=-4..-2,b=-18..-12,view=30..80, style=patchcontour, axes=boxed);
```

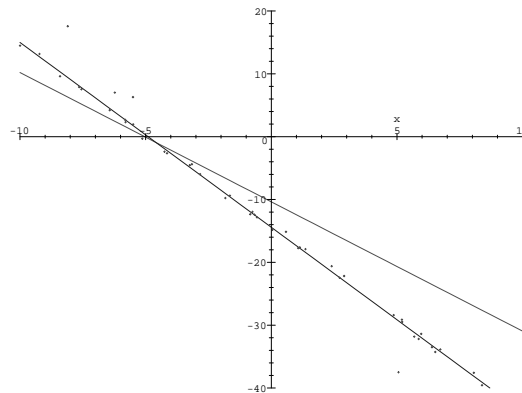


The equations we want to solve are not linear. In fact, they're messy enough that Maple needs some help to find the minimum. From the pictures above, we can choose an appropriate region.

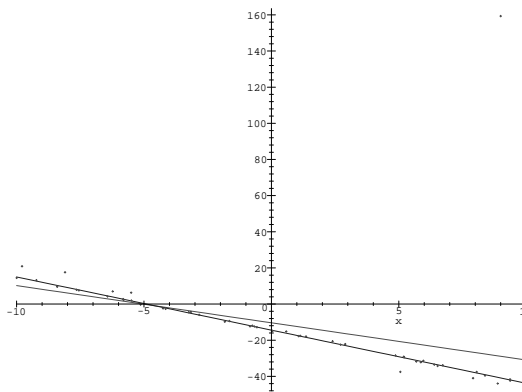
```
> rs:=fsolve( \{diff(R(m,b),m)=0, diff(R(m,b),b)=0\},
              \{m,b\},m=-3.1..-2.8, b=-15..-14);
              rs := {b = -14.41947456, m = -2.942250792}
```

And here we can compare the two lines found.

```
> p := plot(pts,style=point):
    l := plot(subs(sol,m*x+b),x=-10..10):
    s := plot(subs(rs,m*x+b),x=-10..10,color=blue):
    plots[display](p,l,s,view=-40..20);
```



```
> plots[display](p,l,s);
```



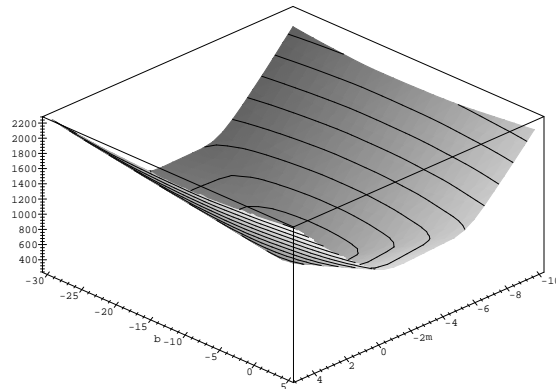
Let us now try to minimize  $\sum |\epsilon_i|$ . Since this function is not even differentiable near  $\epsilon_i = 0$ , we have to work harder to get less:

```
> A:=(m,b)->sum(abs(epsilon(pts[i],m,b)),i=1..nops(pts));
```

$$A := (m, b) \rightarrow \sum_{i=1}^{\text{nops}(pts)} |\epsilon(pts_i, m, b)|$$

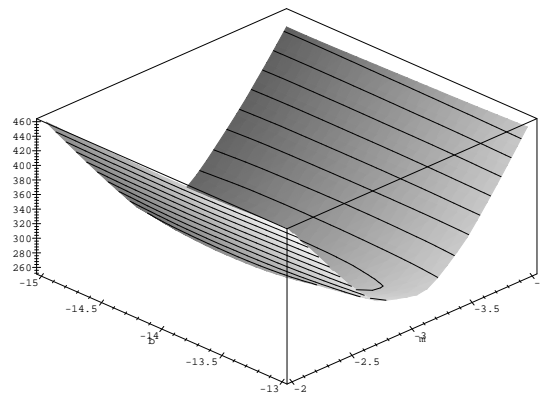
Displaying the graph of this function gives us an idea about where its minimum lies:

```
> plot3d(A(m,b),m=-10..5,b=-30..5,style=patchcontour, axes=boxed);
```



But coercing Maple into finding a numerical approximation to it is not so easy. There *are* reliable and efficient ways to determine the location of the minimum to any precision (at least given some hints), but they are inappropriate for the scope of this chapter, and will not be considered here. Instead we may zero-in on the minimum by constraining the domain over which we display the graph:

```
> plot3d(A(m,b),m=-4..-2,b=-15..-13,style=patchcontour, axes=boxed);
```



Repeatedly narrowing in on the minimum graphically will locate a good choice of  $m$  and  $b$ .

## 6 A nod toward statistics

In the problem we have treated so far, there have been a distinction made between the  $x$  and  $y$  coordinate of a data point  $(x_i, y_i)$ : the  $x$  coordinate is thought as the predictor variable, while the  $y$  coordinate is the predicted value. The distinction is significant. If, for example, we think of  $x$  as a function of  $y$ , the best line that fits the data using the method of least square is, in general, different than the corresponding line when thinking of  $y$  as a function of  $x$ . Indeed, if we use the same data as in §1, the best line  $x = ny + c$  that fits it is given by  $x = 0.6677953348y - 0.738807374$ , whose inverse is the line  $y = 1.497464789x + 1.106338028$ ; that inverse differs from the line we found in §1. In this case, we use as error the square of

the *horizontal* distance, and it is somewhat perturbing that this similarly reasonable approach leads to a best fit that differs from the one obtained when employing vertical distances.

The situation can be made even worse if neither  $x$  nor  $y$  is a predictor variable. In this case, we want to minimize the shortest distance to a line  $y = mx + b$  rather than the vertical (or horizontal) distance. The resulting equations will not be linear, nor can they be made linear. However, Maple will be able to find the critical points with no trouble. There will always be at least two (there may be a third with a huge slope and intercept)—one is the minimum, and the other is a saddle point. It is worth thinking for a few minutes about the geometric interpretation of the saddle point in terms of the problem at hand.

In practice, the perturbing fact that different but reasonable error functions lead to best fits that are different is resolved by knowing what part of the data is the input and what part is predicted from it. This happens often enough, though not always. We thus can make a good choice for  $E$  and solve a minimization problem. That settled, we are still left with the problem of demonstrating why our choice for  $E$  is a good one.

It is rather easy to write down the solution to the problem in §1: if

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i,$$

then

$$\hat{m} = \frac{n \sum_{i=1}^n x_i y_i - (\sum_{j=1}^n x_j)(\sum_{j=1}^n y_j)}{n \sum_{i=1}^n x_i x_i + (\sum_{j=1}^n x_j)^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{b} = \bar{y} - m\bar{x}.$$

If  $y_i$  are assumed to be the values of random variables  $Y_i$  which depend *linearly* upon the  $x_i$ ,

$$Y_i = mx_i + b + \varepsilon_i,$$

with errors  $\varepsilon_i$  that are independent from each other, are zero on average and have some fixed variance, then the value of  $m$  given above is the *best* estimator of the slope among all those linear estimators that are unbiased. Best here is measured by calculating the deviation from the mean, and this best estimator is the one that produces the smallest such deviation.

This conclusion follows by merely making assumptions about the inner products of the data points  $(x_1, \dots, x_n)$  and  $(y_1, \dots, y_n)$ . Statisticians often would like to answer questions such as the degree of accuracy of the estimated value of  $m$  and  $b$ . For that one would have to assume more about the probability distribution of the error variables  $Y_i$ . A typical situation is to assume that the  $\varepsilon_i$  above are normally distributed, with mean 0 and variance  $\sigma^2$ . Under these assumptions, the values of  $m$  and  $b$  given above are the so called *maximum likelihood estimators* for these two parameters, and there is yet one such estimator for the variance  $\sigma^2$ . But, since we assumed more, we can also say more. The estimators  $\hat{m}$  and  $\hat{b}$  are normally distributed and, for example, the mean of  $\hat{m}$  is  $m$  and its variance is  $\sigma^2 / \sum_{i=1}^n (x_i - \bar{x})^2$ . With this knowledge, one may embark into determining the confidence we could have on our estimated value for the parameters. We do not do so in here, but want just to plant the idea in the interested reader, whom we refer to books on the subject.