

Module 6

Doing least squares regression analysis in practice

In this session we going to make some final comments on the least squares regression fitting method of data

and we going to look at how we going really do this in real world.

Using Computational tools like MATLAB or Python or R

A few comments to make before we move on; in reality there are a huge number of solvers for nonlinear least squares problems

But we can observe, that if we do a Taylor series expansion for x^2 then, the 2nd term, the 2nd derivative is the Hessian.

which gives us information about the
 curvature, or gradient of gradient,
 gradient of the function

and therefore we can shoot directly from
 where the function is zero, just as in NR,
 using that 2nd derivative

Now using the Hessian will be faster, than simply
 taking steps along the steepest descent
 algorithm

Effectively, we would be using the Hessian
 to give us a guess as to the size of the
 step we should take in gradient descent

The problem is that often, the Hessian is
 not very stable, especially far from
 the minimum

The Levenberg-Marquardt method uses steepest
 descent far from minimum, and
 then switches to using the Hessian

3
As it gets close to the minimum, based on the
criterion of whether X^2 is getting better or
not.

• ~~Better of~~
If it's getting better, it uses the Hessian, if
it's in trouble, it uses steepest descent.

There are also the Gauss Newton method and
BFGS method amongst many others that
either use the Hessian directly or build
up information about the Hessian over
successive iterations

and depending on the Convergence, different
methods may be better than others

Robust fitting is another topic you should be
aware of in case you need to look it up later

If we come back to Anscombe's quartet,
we see

that we see the bottom left data set in
 our problem, has just that 1 flying data point.
 A truly Robust fitting method, will be unbothered
 by such a data point

One approach to robust fitting, is minimising,
 instead of the Cost Square, the absolute, or
 the Square deviation

So it does not weight the points that are far away ^{from line},
 as strongly

It means it fits a little bit more usually, look a
 bit better.

Now, let's turn to look at how you do this
 in real world.

In Matlab...

In Python is very nearly as simple:

See how in the Scientific Python, SciPy set of modules the optimizer module, includes a least squares fitting minimizer, curve_fit

"Scipy.optimize.curve_fit"...

See example provided in course.

3 lines to do: ^{first define}
 \Rightarrow 2 lines to do ^{the function}
 \Rightarrow 1 line to do the fit.

rest is about plotting, and importing the data

```
def func(x, a, b, c):  
    return a * np.exp(-b * x) + c
```

~~pop, popt~~

popt, pcov = curve_fit(func, xdata, ydata)

So only a few lines in python to bring it to life

6

Now, into a python Code block, to plot the
Gaussian distribution Shantler (pic)

You will need to give it a starting guess, and we
will give you the input data for the height distribution
in the population

It's important to note here, why we need to have
a starting guess

If we started with a guess here for mean
of 100cm, the model curve would not
overlap with the data at all.

So when we did a little more to 10, we
get no change.

and therefore the gradient of χ^2 w.r.t mean b would be zero

so the algorithm would not know what direction ~~to~~ to

go in

we would not get a sensible answer, for the

location or for grad

and therefore our algorithm would not

know where to go to find ~~the~~ the minimum.

So in doing any of the ~~the~~ data fitting, it's vital to come up with a good means for generating a starting guess

here it's easy, you pick the biggest value

so what we have done in the session, we have finished our little discussion in using vectors and multivariate calculus together

to help us do optimizations of functions and to fit data to functions

①
It seemed it turned out to be very easy
Computation.

Remember in Python, it's just a few lines
to achieve what we have covered in these
2 sessions.

But Now, you understand something of how
these algorithms work under the hood.

But means we will be much better off at
figuring out how to fix them when they
go wrong.