

Contents

1	Background Material	3
1.1	Preliminaries	3
1.1.1	Functional Data Analysis	3
1.1.2	Penalised Regression	8
1.1.3	Finte Dimensionalisation: the General Case	11
1.2	The Two-Stage Parameter Cascade	15
1.3	The Data2LD Package	18
1.3.1	Reflux Data	18
1.3.2	Partial Differential Equation - the Transport Equation	22
1.3.3	Quasi-linear Differential Equations	23
1.4	Three Stage Parameter Cascade	27
2	Dissecting the Data2LD Package	29
2.1	Line Search Methods	29
2.1.1	Wolfe Conditions	30
2.1.2	Extrapolation Failure	31
2.2	How Data2LD Works in Practice	32
3	Derivative-Free Optimisation and the Parameter Cascade	35
3.1	Overview of Quadratic Optimisation Methods	35
3.1.1	Newton's Method	35
3.1.2	Secant Method	36
3.1.3	Successive Parabolic Interpolation	37
3.1.4	Discussion	38
3.2	Bisection Methods	39
3.3	Brent's Method	39
3.4	Estimation Of Parameters For A Standard Cauchy Distribution Using Brent's Method	40
3.5	Robust ODE Parameter Estimation	43
3.6	The Parameter Cascade and Brent's Method	46
3.7	Implicit Filtering	48
3.7.1	Description Of Implicit Filtering	48
3.7.2	Implicit Filtering and Brent's Method	52

3.7.3	Using Implicit Filtering to Fit an ODE to the Melanoma Data .	52
4	A Two Level L_1 Parameter Cascade Using the MM Alogrithm.	57
4.1	L_1 Estimation for the Inner Problem	57
4.1.1	An MM Algorithm For Computing the Median	58
4.1.2	Penalised L_1 Fitting	60
4.1.3	Discussion	61
4.1.4	Testing the Algorithm on the Melanoma Data	64
4.2	The Two Level Parameter Cascade with L_1 Norm	68
4.3	Acclerating The Rate of Convergence	72
4.3.1	Illustrative Example: Using Imputation to Fit an ANOVA Model With Missing Data	73
4.3.2	Generalisations	73
4.3.3	Accelerating the L_1 Fitting Alogrithm	78
5	Fitting Non Linear Regression Models With the Parameter Cascade	81
5.1	Fitting Linear Homogenous ODEs Using the Parameter Cascade	82

Chapter 1

Background Material

1.1 Preliminaries

1.1.1 Functional Data Analysis

Functional data analysis (FDA) is a field of statistics where it is assumed that the data observed at a given set of independent observation times (or coordinates etc.) represent noisy observations of some underlying function.

The approach taken here is to assume that an unknown differential equation can adequately - though not necessarily exactly - describe the process producing the data.

Specification of Function Spaces

The functions in question are generally all assumed to be members of some countably infinite dimensional vector space, such as the set of all functions $f(\cdot)$ such that $\int_0^T |f''(t)|^2 dt < \infty$ over some interval $[0, T]$.

This assumption implies that any given function can be represented as a countably infinite combination of basis elements, which are themselves functions. This means for a chosen set of basis elements $\{\phi_1(t), \phi_2(t), \dots\}$ and any given function $f(t)$, there is a set of coefficients $\{a_1, a_2, \dots\}$ such that:

$$f(t) = a_1\phi_1(t) + a_2\phi_2(t) + \dots$$

Functional Data Analysis can thus be regarded as a generalisation of multivariate statistics where the number of dimensions is potentially infinite.

Substantial complications are introduced into the statistical analysis because functions are generally much richer objects than real numbers or vectors. A function will generally have a different value for each input value, and the number of non-integer numbers on any interval - and hence potential inputs - is infinite. Functions cannot be trivially represented on paper or in computer memory in a similar fashion as real numbers or vectors.

For the purposes of this thesis, it is assumed that the functions in question are continuous mappings.

In practice one attempts to resolve this difficulty by finding or otherwise constructing a discrete problem that resembles the functional problem, and then solve this approximate problem.

It might be that case that the approximate problem can itself only be solved approximately using numerical methods.

Statistical problems that involve differential equations are particularly difficult. More naive approaches force the practitioner to solve the ordinary differential equation (ODE) numerically everytime it is desired to evaluate the goodness of fit. For these situations, it is necessary by definition to use numerical analytic techniques to construct a proxy problem that resembles the original problem sufficiently well and that is sufficiently easy to work with.

For example, consider the problem of parametric estimation for a stochastic differential equation (SDE) of the form

$$dX = f(X; \theta)dt + \sigma dW.$$

Here $X(t)$ is the stochastic process being modelled, $f(\cdot; \theta)$ is a known function with a parameter θ to be estimated, σ is a volatility parameter, and $W(t)$ is a standard Brownian motion.

This SDE is equivalent to asserting for any time t and increment h that

$$X(t+h) = X(t) + \int_t^{t+h} f(X(s); \theta)ds + \sigma[W(t+h) - W(t)].$$

Suppose there are observations X_1, X_2, \dots, X_N of $X(t)$ at evenly spaced times, and that h is the distance between the time points. The integral formulation of the SDE suggests that if h is small enough, then

$$X_{k+1} \approx X_k + f(X_k; \theta)h + \sigma h Z_{k+1}.$$

The Z_k here are i.i.d standard Normal random variables. This is known as the *Euler-Maruyama Approximation*.

Instead of attempting to estimate parameters for the SDE, we can fit parameters for a non-linear AR(1) process that acts as a proxy problem for the original SDE. This is a much more tractable problem than the original SDE.

In FDA, the assumption is usually made that all the functions can be represented as a linear combination from some chosen *finite* set of basis functions. Rather than discretise the differential operator as in the above example, the space of functions is discretised instead.

A differential equation (or a similar problem) over some finite dimensional space of functions with n dimensions can be represented as a problem over the Euclidean space R^n , this is a discrete problem.

The modelling process for functional data as described in Figure 1.1 can be more complex than standard statistical problems.

As is the case for typical statistical problems, the first step is to . Here one must only be certain that the model at used is sufficiently broad or well-specified to be able to actually capture the phenomena at hand.

Formulate a Model.

Construct a Discretised Model that Approximates the Original Model. Unless the statistical model is trivial, the next step is to construct a proxy model. This generally requires ideas from Numerical Analysis.

Conduct Statistical Analysis Using the Discretised Model. While the discretised model tends to be simpler than the original model, this task is not necessarily trivial as shall be seen.

Check the Approximation Error in Discretised Model. If the discretised model is too poor an approximation, then the results of any statistical analysis conducted could be substantially biased as a result of the approximation error introduced, even if the original model were perfectly sound. If the original model is biased, then the approximate one might be even more so.

Therefore, one should consider conducting post hoc checks. For example, running the analysis again with a more accurate approximate model and comparing the results with the original model. If both agree, it is evidence the approximate models are both reasonably accurate. Constructing a more accurate approximation is generally a straightforward and intuitive process, with the exact approach depending on the situation at hand.

In the context of FDA, this generally entails increasing the number of basis functions so that the associated approximation error is smaller. This is too complex a question to consider in sufficient detail at this point, so here is an illustrative example of how an analysis would be conducted, that only makes use of elementary ideas from Statistics and Numerical Analysis:

Suppose that one were attempting to estimate the parameters of an ODE by means of least squares, and one was using a finite-difference solver to compute the fitted values, and hence to determine the goodness-of-fit.

Once the fitting algorithm had converged, one might run the solver again with a smaller stepsize and the same parameters and check if this has made a substantial change in the the sum of squared residuals.

If there has been a substantial change as a result of the stepsize reduction, then one would have to consider running the entire fitting procedure again starting from the previously computed parameter estimate, except with the smaller stepsize, taking the

new parameter estimates, and then checking again if decreasing the stepsize yet again produces substantial change in the goodness-of-fit statistic.

This procedure can even be automated. The Implicit Filtering algorithm computes an approximate gradient using finite differences and uses this to perform optimisation. If the algorithm cannot produce a decrease in the objective function, or it cannot be certain that the true gradient isn't in fact zero, it reduces the stepsize. The algorithm terminates when the change in the objective function between changes in the stepsize has fallen below a chosen tolerance level.

If the fitting method used is slow however, then these such approaches can potentially be very slow due to the need to solve the same problem over and over again at increasing levels of precision.

Fortunately, Functional Data Analysis does not always require the recomputation of the curve in such a fashion whenever the parameters are changed. Instead of being implicitly represented as solutions of an ODE, functions are explicitly represented as elements in some finite dimensional vector space. As shall be seen, the objective function is generally a mapping from some vector space R^n to R that can often be evaluated reasonably easily, or at least more easily than having to run an ODE solver.

Check If Results of Statistical Analysis Are Consistent With Discretised Model. In the previous step, one checked that the approximate model was actually acting as a proxy for the original model. One must then check that the statistical analysis conducted using the approximate model is valid in its own right.

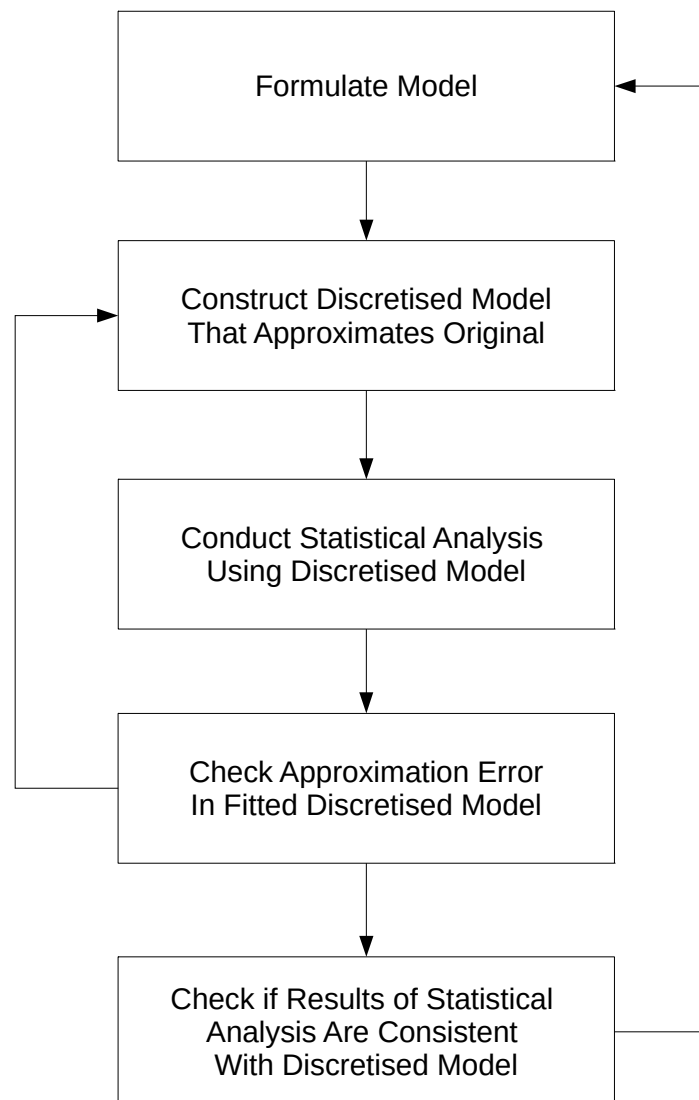


Figure 1.1: Statistical Modelling Process For Functions

1.1.2 Penalised Regression

Suppose we have N noisy observations y_i at times t_i from some function $f(t)$, and we wish to estimate $f(t)$, from the data. A naive approach would be to estimate $f(t)$ by minimising a least squares criterion:

$$SSE(f) = \sum_{i=1}^N (y_i - f(t_i))^2$$

Here, $SSE(\cdot)$ is a function that assigns a real number to every real-valued function that is defined for all the t_i .

There is an obvious problem with this criterion - it does not have a unique minimiser. Any function $g(t)$ such that $g(t_i) = y_i$ will minimise $SSE(\cdot)$. There are an infinite number of degrees of freedom, but only a finite number of observations.

To ensure uniqueness, it is necessary to impose a further condition to discriminate between different candidates, a way to choose between different functions that interpolate a given set of points.

Smoothing Splines

One potential criterion is to introduce a second order penalty. If two functions fit the observed data equally well, the more regular or less "wiggly" function is chosen. There are several ways of translating this intuition into a formal fitting procedure.

A common choice is to measure the degree of irregularity by using the integral of the second derivative over a chosen interval $[0, T]$. The upper limit T should be chosen to allow for all observation times to be included.

$$\int_0^T |f''(t)|^2 dt.$$

For a given set of points, the smooth interpolating curve that minimises the energy integral above is given by an interpolating cubic spline.

Choosing the most regular interpolating curve is not necessarily a very good estimation strategy however because it strongly prioritises goodness-of-fit above all other considerations. If the data is noisy, there is a risk of overfitting and poor predictive power. There is a trade-off between bias and variance.

In practice, a joint estimation strategy is pursued that attempts to find a good balance between fidelity to the observed data and reasonably regular behavior. This involves minimising the following penalised least squares criterion:

$$PENSSSE(f; \lambda) = \sum_{i=1}^N (y_i - f(t_i))^2 + \lambda \int_0^T |f''(t)|^2 dt$$

The λ term dictates the trade-off between fidelity to the data and regularity.

Suppose there were a candidate function $g(t)$, then by taking the cubic spline such that its value at t_i is equal to $g(t_i)$, we can produce a curve $s(t)$ that has the same least-squares error as $g(t)$, but with $\int [s''(t)]^2 dt \leq \int [g''(t)]^2 dt$. Thus, the curve that minimises *PENSSE* can be assumed to be a cubic spline.

To find the minimiser of $PENSSE(\cdot; \lambda)$ first, assume that $f(t)$ can be represented as a linear combination of K cubic spline functions $\phi_i(t)$ that can represent any cubic spline with knots at the t_i . This implies that

$$f(t) = \sum_{i=1}^K c_i \phi_i(t).$$

Note that it is only required that the set of basis splines only possess enough resolution to represent the function that minimises *PENSSE*, it is not required that this set of splines is minimal.

A perfect basis is generally not needed, the basis need only be a reliable workhorse. This obviates what otherwise could prove to be a distracting nuisance for a practitioner who wishes to conduct data analysis instead of getting entangled in minutiae.

Let the design matrix Φ be defined by $\Phi_{ij} = \phi_i(t_j)$, and let the weight matrix \mathbf{R} be defined by $\mathbf{R}_{ij} = \int_0^T \phi_i''(t) \phi_j''(t) dt$. Then *PENSSE* can be written in terms of the vector of coefficients \mathbf{c} and observations \mathbf{y} as:

$$PENSSE(\mathbf{c}; \lambda) = \|\mathbf{y} - \Phi \mathbf{c}\|^2 + \lambda \mathbf{c}' \mathbf{R} \mathbf{c}$$

The problem has been discretised into one on R^K .

The optimal value of \mathbf{c} is given by

$$\hat{\mathbf{c}} = (\Phi' \Phi + \lambda \mathbf{R})^{-1} \Phi' \mathbf{y}$$

This is an exact solution to the original problem because the span of the $\{\phi_i(t)\}$ contains the function that minimises *PENSSE*. The coefficient vector $\hat{\mathbf{c}}$ is the set of coordinates of the optimal function within this finite-dimensional vector space.

Piecewise Trigonometric Interpolation

Consider a more difficult penalised regression problem, that gives a sense of the limits of the approach employed so far.

$$PENSSE(f; \lambda) = \sum_{i=1}^N (y_i - f(t_i))^2 + \lambda \int_0^T |f''(t) - f(t)|^2 dt$$

PENSSE can be minimised in this case taking by a piecewise function consisting of linear combinations of $\sin(t)$ and $\cos(t)$ over each interval, and matching them together.

Note that a function of the form $a_0 + a_1 \cos(t) + b_1 \sin(t) + a_2 \cos(2t) + b_2 \sin(2t) + \dots$ can be written as a polynomial in e^{it} and e^{-it} . For this reason, such a piecewise trigonometric function can also be referred to as a piecewise trigonometric polynomial.

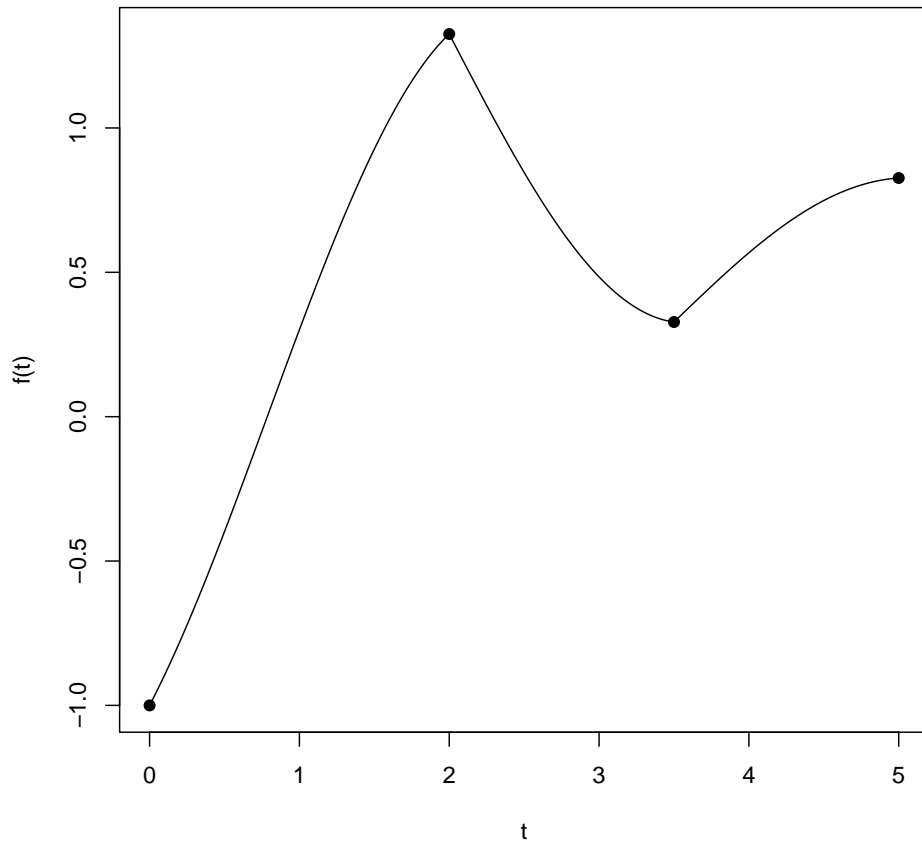


Figure 1.2: Plot of a Piecewise Trigonometric Curve. Note the kinks between segments.

As can be seen in Figure 1.2, a piecewise trigonometric polynomial of second degree generally fails to be smooth at the boundary points, and thus has a kinked appearance. For the purposes of statistical modelling, it is strongly desirable to impose the additional constraint that $f(t)$ must at least everywhere differentiable.

1.1.3 Finte Dimensionalisation: the General Case

To find an exact solution to the two problems in the previous section, it was necessary to construct a finite dimensional function space that contained the minimal function. However it is not guaranteed that this is always possible. In practice, one would hope that the optimal function can be approximated sufficiently well by taking a linear combination from some chosen set of functions. Spline bases tend to be a reliable workhorse that are effectively the default choice. They provide a good balance between being well behaved as objects for regression and having good approximating power.

For comparison, Chebyshev Polynomials can often provide better approximation power for a given number of basis functions and are used in Computational Fluid Dynamics for this purpose. Unfortunately, they can be poorly behaved statistically because they consist of high order polynomials that are difficult to fit to data.

Functional Data Analysis thus consists of the following steps:

1. Formulate a model for $f(t)$. Usually, this takes the form of a penalised regression model, where $f(t)$ is defined as the function that minimises some kind of penalised error
2. Assume that $f(t)$ can be written as a finite combination of chosen basis functions. In practice, this is only approximately true, so it is important to ensure that our basis can actually approximate the optimal $f(t)$ sufficiently well. The function $f(t)$ can thus be written:

$$\begin{aligned} f(t) &= \sum_{i=1}^K c_i \phi_i(t) \\ &= [c_1, \dots, c_K]' [\phi_1(t), \dots, \phi_K(t)] \\ &= \mathbf{c}' \boldsymbol{\phi}(t) \end{aligned}$$

Note that $f(t)$ is now defined by the coefficient vector \mathbf{c} .

3. Formulate the model in terms of the coefficient vector \mathbf{c} . A statistical problem over some given functional space has been transformed into a statistical problem over R^K .

As is done in some texts, we will immediately provide an example with a very small basis to illustrate these steps [1]. Consider the following non linear penalised regression problem:

$$PENSSE(f; \lambda) = \sum_{i=1}^N (y_i - f(t_i))^2 + \lambda \int_0^1 |t^2 f'' - f|^2 dt$$

The differential equation associated with the penalty term is known as Euler's Equation, and can be regarded here as a toy equation representative of equations such

as Bessel's equation. The solution is given by $f(t) = at^{r_1} + bt^{r_2}$, where r_1 and r_2 are the roots of the quadratic equation $2r^2 - 2r - 1 = 0$. Thus, $r_1 \approx 0.36$ and $r_2 \approx 1.36$.

Some texts in the field, for the sake of illustration it will be assumed that that $f(t)$ can be written as a quadratic - a linear combination of the basis functions $\{1, t, t^2\}$:

$$f(t) = at^2 + bt + c$$

Then:

$$\begin{aligned} \int_0^1 |f'f' - f|dt &= \int_0^1 |t^2(at + b) - (at^2 + bt + c)|^2 dt \\ &= \int_0^1 |at^3 + bt^2 - at^2 - bt - c|^2 dt \\ &= \int_0^1 |at^3 + (b - a)t^2 - bt - c|^2 dt \\ &= [a \ (b - a) \ (-b) \ (-c)]' \mathbf{H} [a \ (b - a) \ (-b) \ (-c)] \end{aligned}$$

The elements of the matrix \mathbf{H} are defined by $\mathbf{H}_{ij} = \int_0^1 t^i t^j dt = 1/(i + j + 1)$.

To conclude, the penalised error is given by:

$$\begin{aligned} PENSSE(a, b, c; \lambda) &= \sum_{i=1}^N (y_i - at_i^2 - bt_i - c)^2 \\ &\quad + \lambda [a \ (b - a) \ (-b) \ (-c)]' \mathbf{H} [a \ (b - a) \ (-b) \ (-c)] \end{aligned}$$

Thus, we now gone from a problem specified in terms of functions, to a nonlinear least squares problem in the three coefficients a, b and c . The quality of this approximate model as λ gets larger and larger depends on how well the functions $t^{0.36}$ and $t^{1.36}$ can be respectively approximated by quadratics over the interval $[0, 1]$.

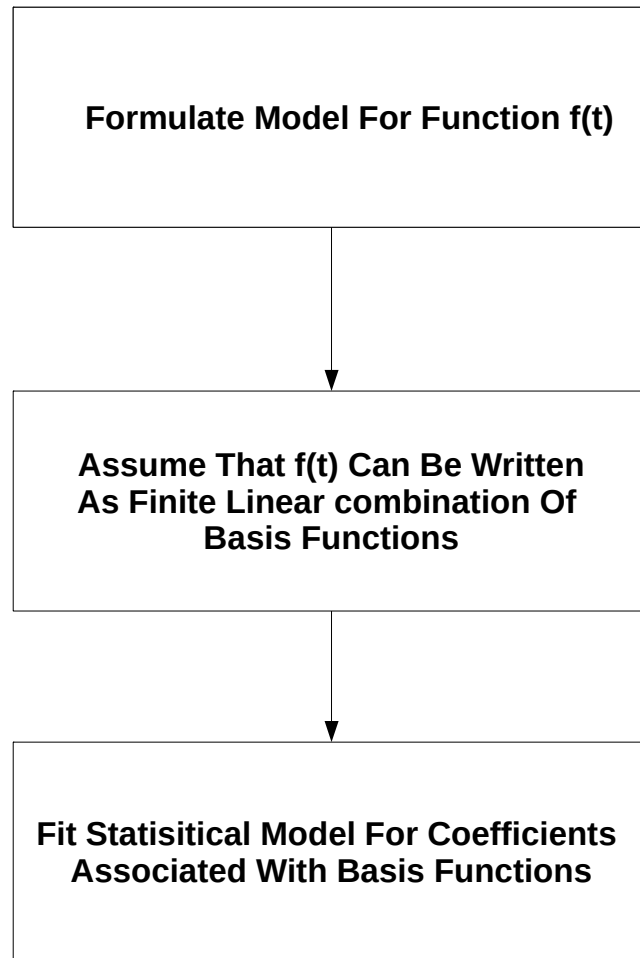


Figure 1.3: Statistical Modelling Process For Functional Data Analysis

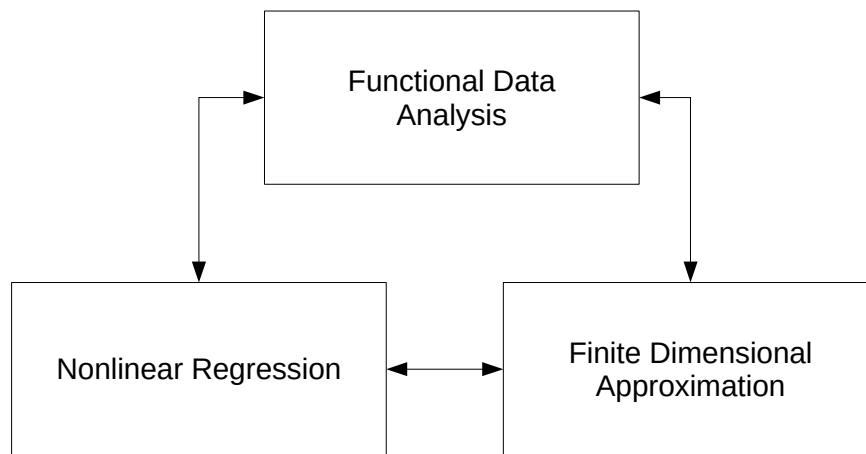


Figure 1.4: Elements of Functional Data Analysis

1.2 The Two-Stage Parameter Cascade

Consider the following penalised regression problem:

$$PENSSE(f, \theta) = \sum_{i=1}^N (y_i - f(t_i))^2 + \lambda \int_0^1 |T_\theta f|^2 dt.$$

Here T_θ is some differential operator, that is parameterised by an unknown θ that is to be estimated.

T_θ can be an ordinary differential operator or a partial differential operator; and whether it is linear, quasi-linear, or nonlinear.

There are two statistical objects to be estimated here: the parameter θ , and the function $f(t)$.

Ramsay and Cao propose the following hierarchical approach to estimation[4]:

Given a fixed value of θ , let $f(t|\theta)$ denote the function that minimises $PENSSE(f, \theta)$

For a given value of θ , it's associated mean square error is defined by:

$$MSE(\theta) = \sum_{i=1}^N [y_i - f(t_i|\theta)]^2$$

By making $f(t)$ dependent on θ , the fitting problem has been reduced to a problem in non-linear least squares.

This leaves the issue of estimating the optimal value of θ - Ramsay and Cao propose the use of gradient descent.

For a given value of θ , $f(t|\theta)$ is found. These two values together are then used to compute $MSE(\theta)$ and $\nabla MSE(\theta)$. Finally, a new value of θ is computed by perturbing θ in the direction of the gradient. This scheme is sketched out in Figure 1.5.

It is assumed that $f(t)$ can be represented by a finite vector \mathbf{c} associated with an appropriate basis. This leads to a pair of nested optimisation problems: the *Inner Optimisation* involves finding the value of \mathbf{c} that minimises the penalised least squares criterion given θ , and the *Middle Optimisation* entails finding the value of θ that minimises $MSE(\theta)$.

There is thus a "cascade" of estimation problems, where the results of the lower level estimation problem feeds back in to the higher level one.

Note that every time a new value of θ is introduced, the associated function $f(t|\theta)$ must be computed from scratch. The middle optimisation can thus generate many inner optimisation subproblems as the parameter space is explored, and these in turn could require multiple iterations to complete if no explicit formula for \mathbf{c} given θ is available.

Figure 1.5 is a high level overview of the Parameter Cascade, in particular, the step of computing $f(t|\theta)$ is presented as a single atomic and organic step, even though it could be a complex process in its own right. This risks masking some of the computational work that is happening. A more complete description is provided in Figure 1.6.

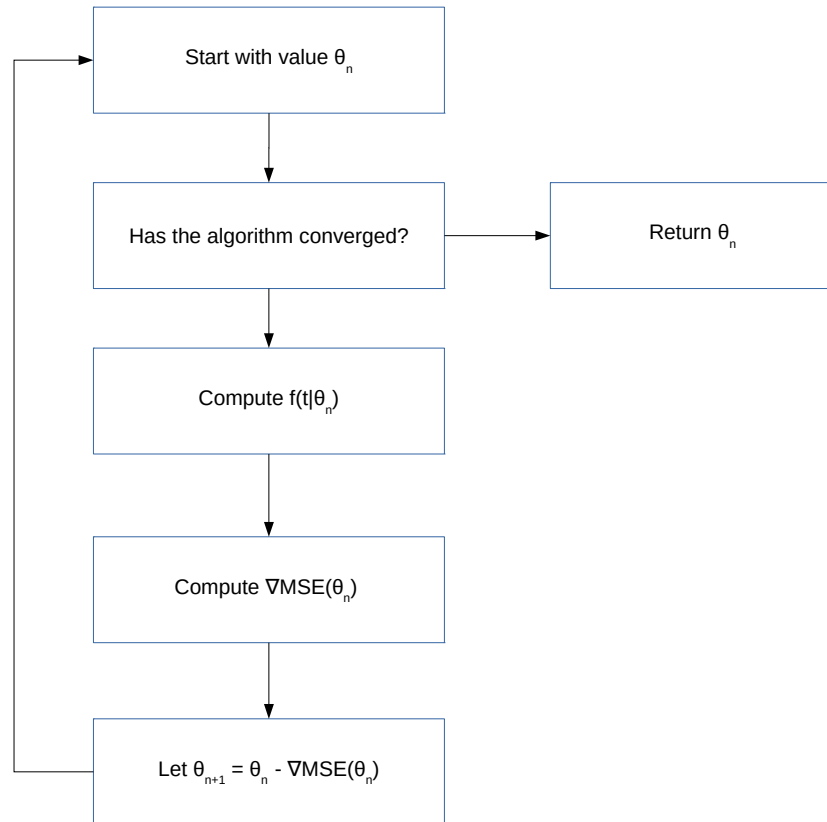


Figure 1.5: Two Stage Parameter Cascade (Simplified)

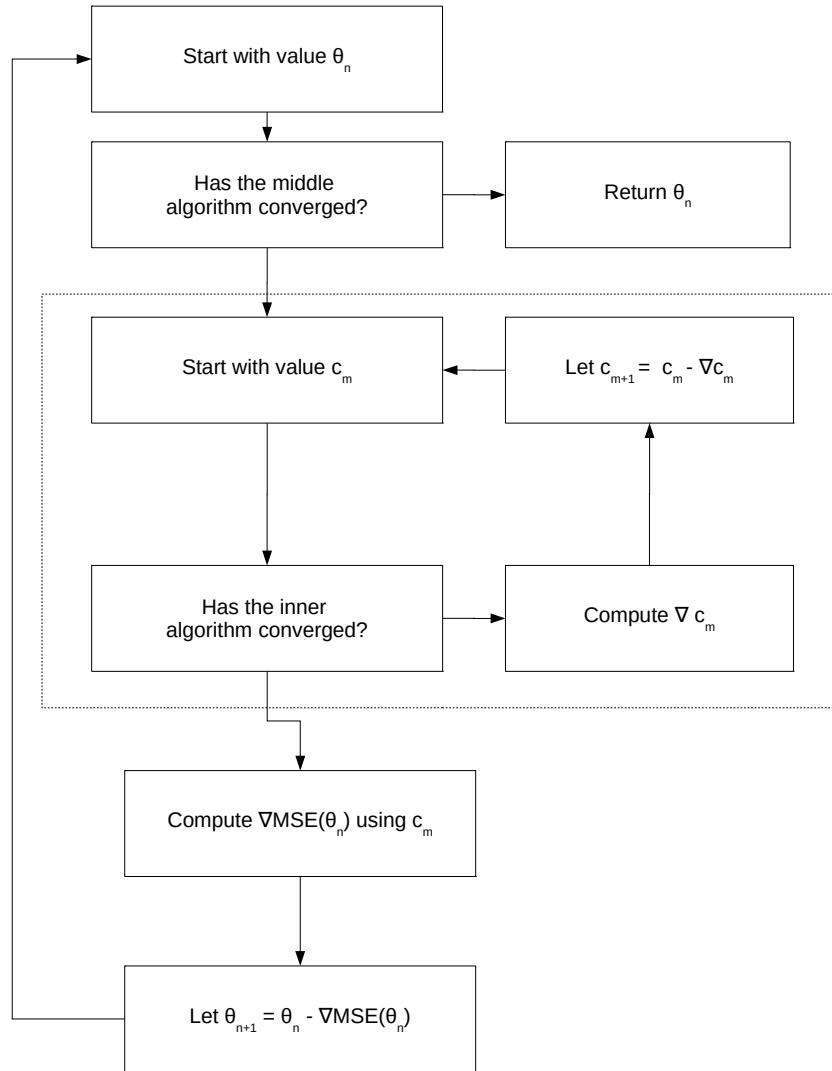


Figure 1.6: Schematic of the Two Stage Parameter Cascade With the Inner Optimisation Visible

1.3 The Data2LD Package

The Data2LD package is an R package intended to perform smoothing using the Paramter Cascade to fit linear differential operators with a forcing function, that is, ODEs of the form:

$$\sum \beta_i(t) D^i f(t) = u(t)$$

1.3.1 Reflux Data

The Reflux data, plotted in Figure 1.7, describes the output of an oil refining system. A given fraction of oil is being distilled into a specific tray, at which point it flows out through a valve. At a given time, the valve is switched off, and distillate starts to accumulate in the tray [26].

ODE Fit

$$\begin{cases} y'(t) = -\beta y(t) & t \leq t_0 \\ y'(t) = -\beta y(t) + u_0 & t \geq t_0 \\ y(0) = 0 \\ y'(0) = 0 \end{cases}$$

This ODE admits an exact solution. Letting $\gamma = u_0/\beta$ and C be an arbitray constant, then the solution is given by

$$y(t) = \begin{cases} 0 & t \leq t_0 \\ \gamma + Ce^{-\beta t} & t \geq t_0 \end{cases}$$

The exponential term is of the form $Ce^{-\beta t}$ rather than $Ce^{-\beta(t-t_0)}$ because the $Ce^{-\beta(t-t_0)} = Ce^{-\beta t}e^{-\beta t_0} = [Ce^{-\beta t_0}e^{-\beta t}]$, the $e^{-\beta t_0}$ term is thus absorbed into the unfixed constant term.

In order to ensure that $y(t)$ is continous at t_0 and monotone increasing, we require that $\gamma + C = 0$ and that $\beta > 0$

This model is perfectly adaaquate, but it turns out that the constraint $C = -\gamma$ is unsuitable from the point of view of numerical parameter estimation.

However, if we allow t_0 to vary, we can allow C to assume any negative value while preserving monotonicity and continuity.

Assume that $y(t)$ is instead given by:

$$\tilde{y}(t) = \max(0, \gamma + Ce^{-\beta(t-t_0)})$$

The function $\tilde{y}(t)$ satisfies the same ODE and initial conditions as $y(t)$ except that the change point t_0 is shifted to t'_0 defined by:

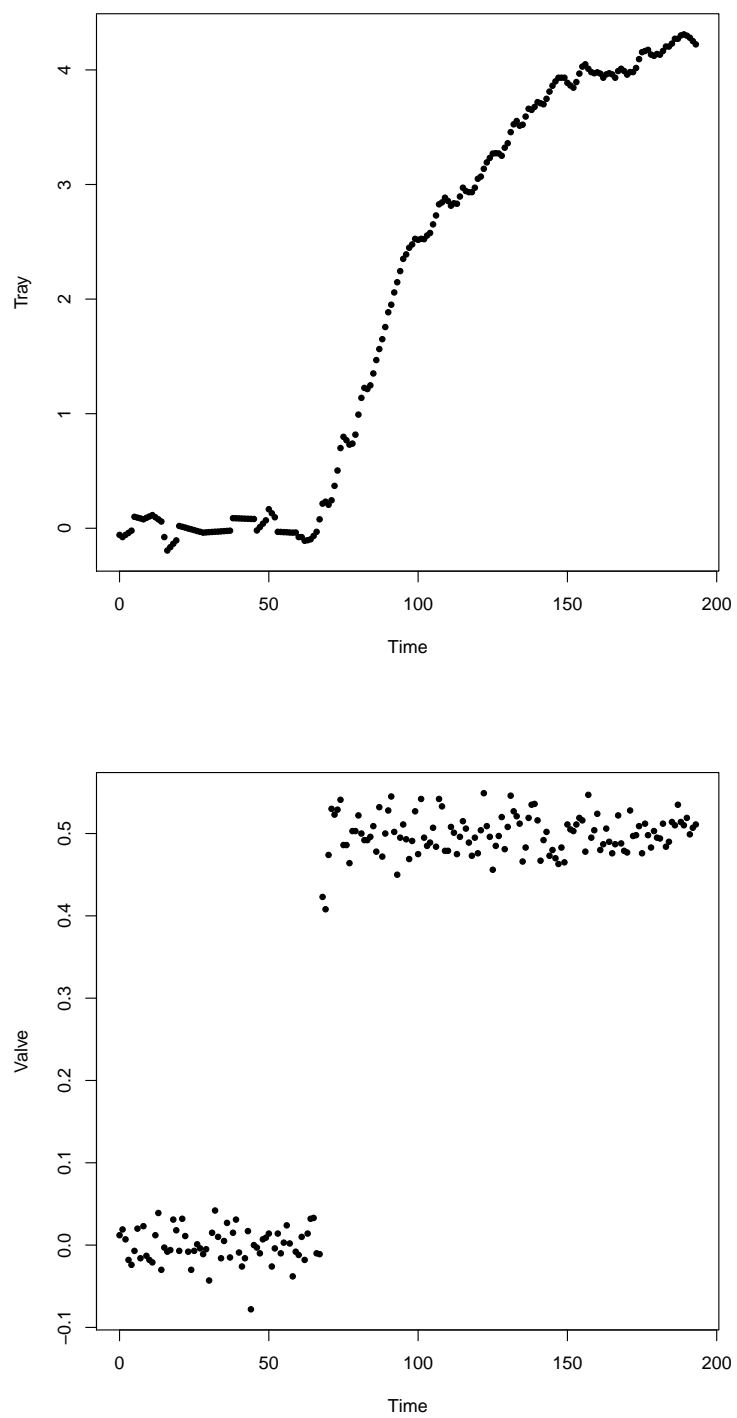


Figure 1.7: Reflux Data

$$t'_0 = \max \left(t_0, t_0 - \frac{1}{\beta} \ln \left(\frac{-\gamma}{C} \right) \right)$$

The function $\tilde{y}(t)$ is a combination of simpler functions, joined together using the maximum operator instead of the addition operator, see Figure 1.8.

Some might be sceptical at the ad hoc fashion in which we are proceeding, but this often proves necessary in the course of mathematical modelling with ODEs. Often, one can only make an educated or inspired guess about the coarse behaviour in advance, and then refine the model depending on how the proposed ODE behaves.

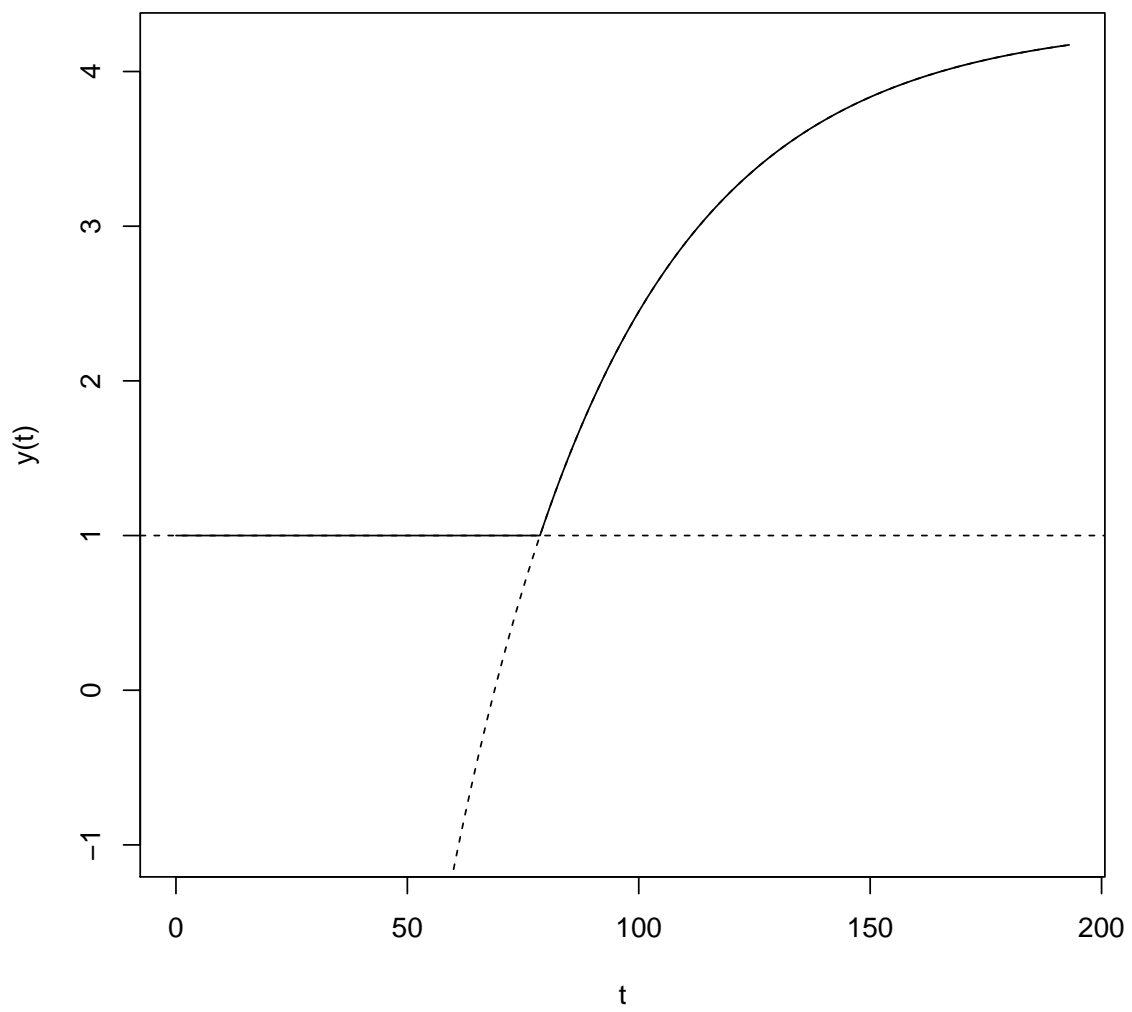


Figure 1.8: Plot of $\tilde{y}(t)$ and its constituent functions

Parameter Estimation

We assume that the breakpoint t_0 is known in advance. Then our model for $y(t)$

$$y(t) = \begin{cases} 0 & t \leq t_0 \\ \beta_0 + \beta_1 e^{\beta_2 t} & t \geq t_0 \end{cases}$$

Note that this function might not be well defined at t_0 , we will address the question of matching later on. We must estimate the three unknown coefficients $\beta_0, \beta_1, \beta_2$.

Estimating β_0 from the data The plot suggests that $\beta_2 < 0$, and $\beta_1 < 0$, under this assumption, we have that:

$$\lim_{t \rightarrow \infty} y(t) = \beta_0$$

Where the convergence happens monotonically from below

So an initial estimate for β_0 is given by $\hat{\beta}_0 = \max(y_i)$

Estimating β_1 and β_2 given β_0 The model can be rearranged so that

$$\log(\beta_0 - y(t)) = \log(\beta_1) - \beta_2 t$$

This equation is only valid so long as the left hand side is well defined however. For our initial estimate of β_0 , it is only necessary to exclude the largest observed value of y .

Simultaneous Estimation Once we have reasonable estimates for β_0, β_1 and β_2 , we can use non linear regression to estimate all three.

Matching

We have now produced separate estimates for $y(t)$ at $t \leq t_0$ and $t \geq t_0$, these distinct functions do not necessarily agree at $t = t_0$.

To stitch the two functions together, we let $\hat{y}(t) = \max(0, \beta_0 + \beta_1 e^{\beta_2 t})$. This is a continuous function that entirely satisfies the original ODE, except for the precise location of the transition point.

This seems very "ad hoc", but is consistent with how ODE's are stitched together in Applied Maths texts.

1.3.2 Partial Differential Equation - the Transport Equation

A linear PDE that would be analogous to the linear ODE above would be the Transport Equation:

$$\frac{\partial u(x, t)}{\partial t} + \beta \frac{\partial u(x, t)}{\partial x} = 0$$

The ODE $y'(t) + \beta y(t) = 0$ can be thought of as a simplification of the Transport Equation, where it is assumed that $u(x, t)$ only varies with time, and not with space.

A general solution to the Transport Equation is given by:

$$u(x, t) = f(x - \beta t)$$

The function $f(\cdot)$ is unspecified. The solution $u(x, t)$ is constant along the rays $x = \beta t + C$

The solution is in effect an animation of the shape $f(x)$ moving to the right at fixed speed β .

Statistically speaking, fitting the Transport Equation to observed data is a semi-parametric problem because one of the parameters to be estimated is a function. This is also a transformation model, since the plot of $u(x, t)$ with respect to x at a fixed time t is a transformed version of $f(x)$, the curve at $t = 0$.

If the parameter governing the transformation process - β - is known, $f(\cdot)$ is reasonably easy to estimate. If β is unknown, one might attempt to maximise a profiled objective function in β .

Suppose there are n observed values y_i at time t_i and location x_i .

The value observed at a point x at time t depends only on $x - \beta t$. The function $f(\cdot)$ could thus be estimated by non-parametrically regressing the observed values at y_i against $x_i - \beta t_i$

What if β were unknown? The above discussion suggests a hierarchical approach to estimation: for a given choice of β , to fit an associated function $f(\cdot|\beta)$ using an appropriate non-parametric estimation method, and compute the associated least squares error. Define the function that associates each β with its sum of squared error:

$$H(\beta) = \sum_{i=1}^n [y_i - f(x_i - \beta t_i | \beta)]^2$$

(In case the left hand side might be slightly unclear - for the i th observation, the associated function $f(\cdot|\beta)$ is evaluated at $x_i - \beta t_i$.)

This is a non-linear least squares problem in β . To estimate β , one would attempt to find the value of β that minimises $H(\beta)$.

The Parameter Cascade Algorithm entails a similar hierarchical approach.

1.3.3 Quasi-linear Differential Equations

Up until this point, it has been possible to use techniques from Applied Mathematics to construct solution strategies on a case by case basis. As differential equations become more complex, this approach begins to rapidly become non-viable. In this section, we will introduce quasi-linear variations of the previous two examples.

The difference between a quasi-linear and a linear differential equation is that the coefficients in a quasi-linear equation are allowed to depend on the unknown function. Instead of an ODE such as $y' = \beta(t)y$, one would have an ODE such as $y' = \beta(y, t)y$. As

we shall see, though quasi-linear problems tend to be reminiscent of linear ones, they are nonetheless substantially more complicated, and require more technical knowledge, and even ingenuity to tackle.

For a quasi-linear variation of a linear ODE, consider the Van Der Pol Equation:

$$y''(t) + \beta(1 - y(t))^2 y'(t) + y(t) = 0$$

This ODE has no obvious solution.

Even if a solution exists, an estimation strategy might be difficult to derive. Consider the inviscid Burger's Equation:

$$\frac{\partial u(x, t)}{\partial t} + \beta u(x, t) \frac{\partial u(x, t)}{\partial x} = 0$$

This equation is identical to the Transport Equation except that the rate term is equal to $\beta u(x, t)$. The solution is given by:

$$u(x, t) = f(s)$$

Here $f(\cdot)$ is some arbitrary function as before, and s is implicitly defined as the solution of the equation $x = \beta f(s)t + s$. Since $s = x - \beta u t$, this can be written as:

$$u(x, t) = f(x - \beta u t)$$

Fitting this model is substantially trickier than the Transport Equation. There is no clean separation between the problem of estimating $f(\cdot)$ and β since $u(x, t)$ appears on the righthand side and scales β .

A further complication is that $u(x, t)$ might only define a *relation*, instead of a function. There might be multiple values of u associated with a given (x, t) that satisfy the solution equation. Physically speaking, multiple values correspond to shock waves.

Discussion We see that the level of knowledge required to devise fitting strategies can increase substantially even with seemingly modest increases in the complexity of the differential equation.

Consider the following quasi-linear model of genetic drift in a population proposed by R.A. Fisher (in the *Annals of Eugenics*) [8]:

$$\frac{\partial u(x, t)}{\partial t} + \beta_1 \frac{\partial u(x, t)}{\partial x} = \beta_2 u(x, t)(1 - u(x, t))$$

This problem is similar to the previous two PDEs we discussed, it even admits travelling wave solutions of the form $f(x + Ct)$ as Fisher himself noted. Nonetheless, it is a much more difficult problem than the previous two despite the apparently modest increase in complexity. One would likely have to consult a textbook that covers non-linear PDEs that can generate waves in fair degree of detail to be able to devise a fitting strategy. This is quite a specialised subject!

As the complexity increases, a practioner will find themselves spending more and more time doing Applied Mathematics, and less and less time doing Statistics.

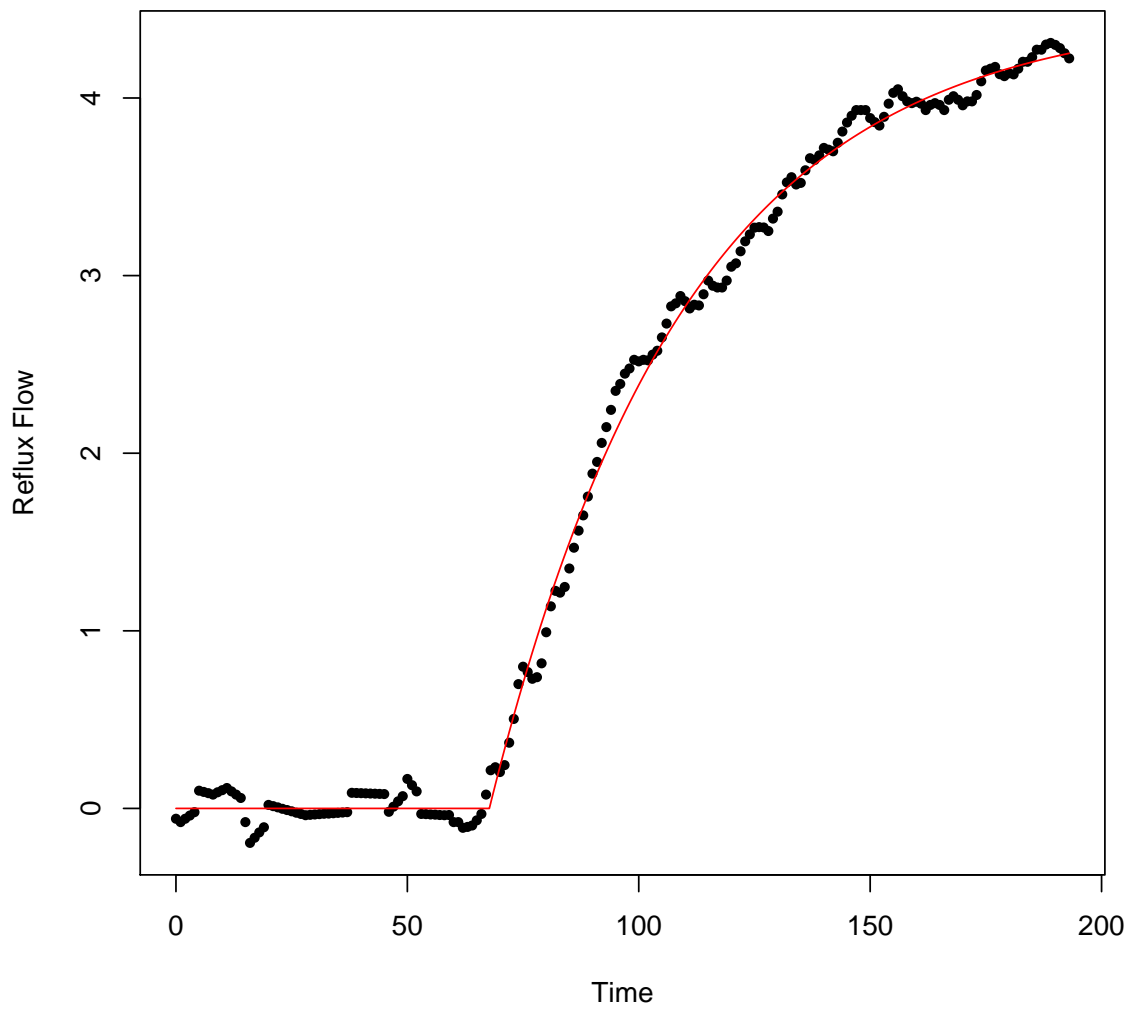


Figure 1.9: Fitted curve using NLS

Functional Model

Instead of the previous approach, which requires a considerable amount of domain-specific knowledge, the functional model is generally employed by Statisticians instead. It has the advantage of being more broadly applicable.

The functional model asserts that

$$y'(t) \approx -\beta y(t) + u(t)$$

Where $y(\cdot)$ and $u(\cdot)$ are functions to be estimated, and β is a single scalar parameter. It is assumed that $u(t)$ is a step function of the form

$$u(t) = a\mathbb{I}_{[0,t_0)}(t) + b\mathbb{I}_{[t_0,\infty)}(t)$$

And that $y(t)$ can be expanded as a linear combination of B-Splines. The knots are duplicated at t_0 so that the first derivative is discontinuous. This model was fitted using the Data2LD package.

1.4 Three Stage Parameter Cascade

Up to this point, the structural parameter λ has been treated as fixed. But it is possible to extend the Parameter Cascade to estimate λ .

It is necessary to an *Outer Criterion* $F(\lambda)$ that determines how good a given choice of λ is.

A common choice of outer criterion is the so-called Generalised Cross Validation.

Just as the problem of fitting a function $f(\cdot|\theta)$ can generate an optimisation subproblem, so that of fitting a third level in the cascade can generate a series of subproblems to find the best parameter choice associated with a given value of λ , which in turn generates a series of subproblems to find the fitted function as the parameter space is explored.

Many packages do not implement the three state parameter cascade. They instead expect practioners to find the best choice of λ by cycling through a set of predetermined values or even just employing manual adjustment.

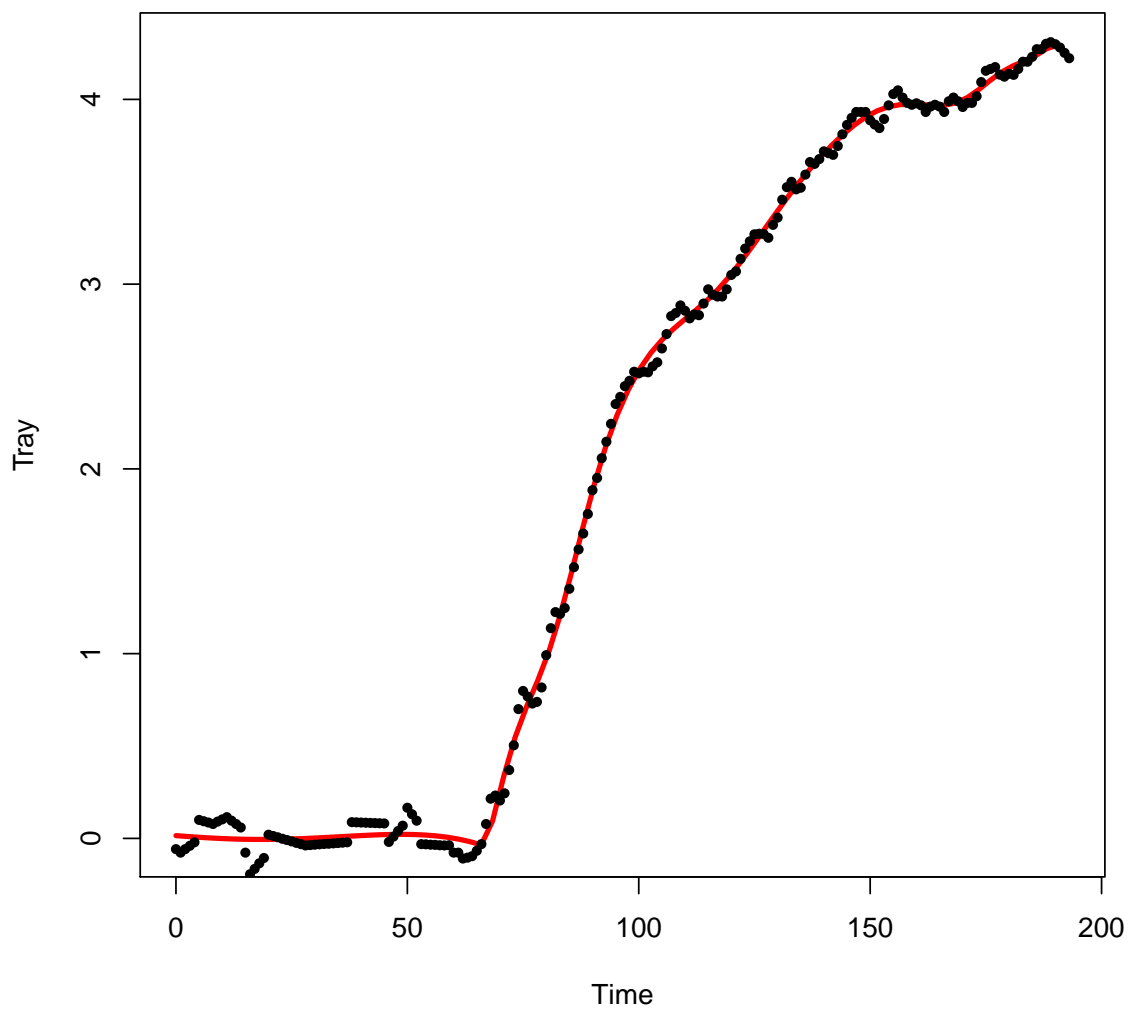


Figure 1.10: Fitted curve using FDA

Chapter 2

Dissecting the Data2LD Package

2.1 Line Search Methods

A common methodology to minimise $f(x)$ is as follows: given a point x_0 , an approximate function $\tilde{f}(x|x_0)$ is constructed, the point that minimises $\tilde{f}(x|x_0)$ becomes the new candidate point, the process is then repeated until convergence.

The Newton Raphson Method employed for Maximum Likelihood Estimation is a well-known example of this approach. The log-likelihood $\ell(\theta)$ is approximated by a quadratic in the neighbourhood of some point θ_0 :

$$\ell(\theta) \approx \ell(\theta_0) + S(\theta_0)(\theta - \theta_0) + \frac{1}{2}I(\theta_0)(\theta - \theta_0)^2$$

It can be shown easily that this approximation is minimised at $\theta_1 = \theta_0 - I^{-1}(\theta_0)S(\theta_0)$

There is a serious weakness with this approach, there is no guarantee that the model will be valid over a sufficiently large radius. It could be the case that the next point advocated by the model could be so far away that the model's opinion on what it regards as the best possible point is completely wrong.

For a one dimensional problem, one can often plot the function at hand if one suspects irregular behaviour. For a function with many input variables the problem is more difficult. Usually, one must conduct post hoc tests to check if the approximation is acting as a good guide as to how the objective function is actually behaving.

The situation needs to be approached more systemically.

The Data2LD package employs Line Search Methods, which work by picking a direction, and then deciding how far in that direction to go. For comparison, Trust-Region Methods set a maximum allowable distance, and then decide on the direction.

For a smooth function $f(x)$ and a search direction p , the following holds:

$$f(x_0 + p) \approx f(x_0) + \nabla f(x_0)'p$$

Two important results follow from this.

The first is that a point x_0 can only be a local minimum or maximum of $f(x)$ if the gradient is zero at that point, otherwise one could find a point nearby where $f(x)$ were higher or lower than $f(x_0)$ by taking a small step in the appropriate direction.

The second, is that if one is at some point x_0 and wishes to reduce $f(x)$ from its current value, one should ensure that the direction p satisfies $\nabla f(x_0)'p < 0$.

The naive choice would be to set $p = -\nabla f(x_0)$, this is known as the gradient descent method.

Gradient descent can be very slow. In particular, the Steepest Descent Method, which produces the point exactly minimise $f(x)$ along the line through x_0 in the direction of $\nabla f(x_0)$, produces steps that are perpendicular to each other, so that the path take by the algorithm forms a zig-zag pattern.

More generally, given a positive definite matrix Q , the vector $p = -Q\nabla f(x_0)$ will produce a descent direction. The obvious choice of Q would be to use the Hessian matrix of second partial derivatives, this is Newton's Method, and can be regarded as a slight generalisation of the Newton-Raphson Method.

There are other choices. For example, if one were trying to maximise a likelihood function, one could use the expected Fisher Information in place of the observed Fisher Information. This is known as Fisher's Method of Scoring.

Alternatively, one could use a matrix of the form $\hat{\Sigma}^{-1}$, where $\hat{\Sigma}$ were an estimate of the variance of the parameters to be estimated. The Cramer-Rao bound suggests that this approximation might produce good results in some cases - it is approximately equal to Fisher's Method of Scoring if the data comes from a Poisson distribution for example and an unbiased estimate of the variance is used.

R's `optim` command makes use of Quasi-Newton Methods that attempt to slowly build up an approximation to the true Hessian as the algorithm progresses. These methods produce a sequence of psuedo-Hessians B_n that satisfy the so-called Secant Condition:

$$B_n(x_n - x_{n-1}) = \nabla f(x_n) - \nabla f(x_{n-1})$$

2.1.1 Wolfe Conditions

In order to ensure the line search method converges, the steps are required to satisfy the Wolfe Conditions. There are other conditions, but these are standard. Data2LD actually makes use of the Strong Wolfe Conditions instead of the standard ones:

$$\begin{aligned} f(\mathbf{x}_k + t_k \mathbf{p}_k) &\leq f(\mathbf{x}_k) + c_1 t_k \mathbf{p}_k' \nabla f(\mathbf{x}_k) \\ |\mathbf{p}_k' \nabla f(\mathbf{x}_k + t_k \mathbf{p}_k)| &\leq c_2 |\mathbf{p}_k' \nabla f(\mathbf{x}_k)| \end{aligned}$$

The first condition ensures sufficient decrease in the objective function, the second ensures a sufficient decrease in the gradient between steps.

2.1.2 Extrapolation Failure

The line search method employs an approximation to the objective function at each iteration instead of working on the true function. As illustrated in Figure 2.1, the model can fail if one takes too big a step. For complex estimation problems, the objective function often has multiple peaks and troughs, so one must be careful that one has not wandered out of the range of validity of the locally constructed approximation.

For example, Data2LD checks whether the slope at a candidate point is positive - the line search method's approximation would regard this as an impossibility. Should this be the case, it attempts cubic interpolation at that point to find a minimum, and if that is impossible, it falls back on linear interpolation.

Data2LD is thus required to verify that the objective function is actually behaving as predicted, this is should not be necessary according to standard implementations of line searches [22]

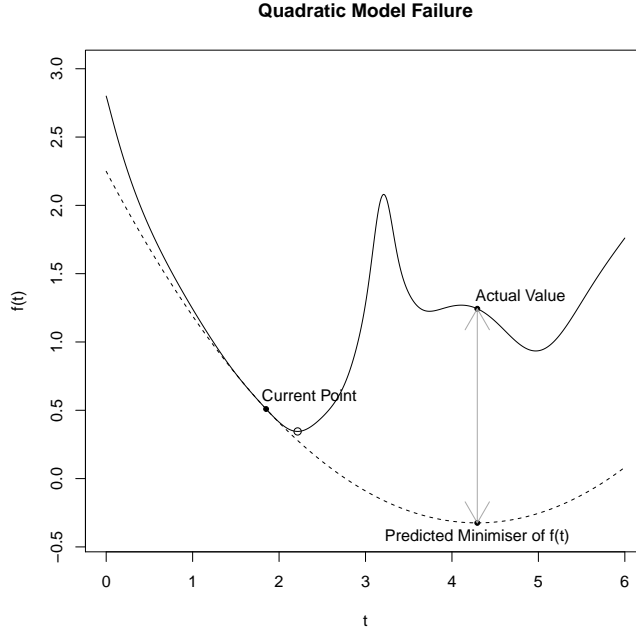


Figure 2.1: Extrapolating too far out can lead to disaster

The code in Data2LD hardcodes unnamed constants into the code. For example putting the number 3.14159 into code without context instead of π . Allowing such 'Magic Numbers' is strongly discouraged because it makes code more error prone and difficult to understand [20].

2.2 How Data2LD Works in Practice

The search direction used by Data2LD is given by:

$$\mathbf{p}_k = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|} \quad (2.1)$$

The objective function $f(\cdot)$ is then probed along the line segment $\{\mathbf{x}_k + \alpha_k \mathbf{p}_k | \alpha \geq 0\}$ to find the next point. Let $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{p}_k)$ denote the value of $f(\cdot)$ restricted to the search direction \mathbf{p}_k . Note that $\phi(0) = f(\mathbf{x}_k)$ and $\phi'(\alpha) = \mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \alpha \mathbf{p}_k)$. This in turn implies that:

$$\begin{aligned}
\phi'(0) &= \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\
&= \left[-\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|} \right]^\top [\nabla f(\mathbf{x}_k)] \\
&= -\frac{\nabla f(\mathbf{x}_k)^\top \nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|} \\
&= -\frac{\|\nabla f(\mathbf{x}_k)\|^2}{\|\nabla f(\mathbf{x}_k)\|} \\
&= -\|\nabla f(\mathbf{x}_k)\|
\end{aligned}$$

Since $\phi(\alpha) = \phi(0) + \phi'(0)\alpha + \mathcal{O}(\alpha^2)$ and $\phi'(0) < 0$, if α is small enough, then $\phi(\alpha) < \phi(0)$. In other words, it is always possible to reduce $\phi(\alpha)$ so long as the step taken is small enough.

Data2LD uses four tests to determine how good a step is:¹

- First Wolfe Condition - compares the decrease in the value of the objective function to an estimate of what it should be.

$$f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \quad (\mathbf{T1})$$

- Second Wolfe Condition (Strong Version) - tests whether the gradient has decreased sufficiently relative to the previous value

$$|\mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)| \leq c_2 |\mathbf{p}_k^\top \nabla f(\mathbf{x}_k)| \quad (\mathbf{T2})$$

- Has the function even decreased compared to the previous iteration?

$$f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}_k) \quad (\mathbf{T3})$$

- Has the slope along the search direction remained nonnegative?

$$\mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq 0 \quad (\mathbf{T4})$$

If **T1** and **T2** are satisfied, then the line search has converged completely. If **T3** has failed, this represents a total failure because it means the line search has failed to actually produce any improvement in the objective function. A failure in **T4** means the function has overshoot a critical point.²

Depending on the outcome of the tests, Data2LD chooses the stepsize as follows:

¹Data2LD actually tests for the negative of **T3** and **T4**, but they are presented so here so that passing a test consistently good and failing is consistently bad.

²If **T4** fails, this implies that $\mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)$ and $\mathbf{p}_k^\top \nabla f(\mathbf{x}_k)$ are of opposite sign since \mathbf{p}_k is chosen so that $\mathbf{p}_k^\top \nabla f(\mathbf{x}_k) < 0$. The Intermediate Value Theorem means there is an $\bar{\alpha}$ between 0 and α_k such that $\mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \bar{\alpha} \mathbf{p}_k) = 0$, so that there is a critical point on the line segment between \mathbf{x}_k and $\mathbf{x}_k + \alpha_k \mathbf{p}_k$.

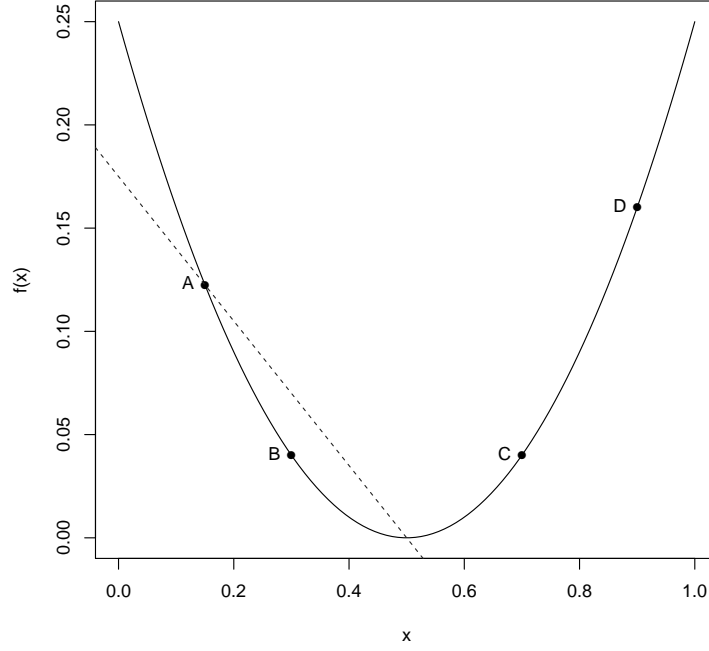


Figure 2.2: Point A is the initial point. Point B passes **T1** with $c_1 = 0.5$ and passes **T2** with $c_2 = 0.9$. Point C fails **T1** with $c_1 = 0.5$ and also fails **T3**, but passes **T4**, and passes **T2** with $c_2 = 0.9$. Point D fails all four tests.

- If **T1**, **T2**, and **T3** are passed, the algorithm terminates.
- If **T1** and **T1** are passed, or **T4** is passed; but **T3** is failed, it means that the slope is satisfactory, but the function has increased rather than decreased. Data2LD reduces the step size
- If all four tests are failed, then the newest point is unsuitable entirely. Data2LD falls back on interpolation to try to find a critical point of $\phi(\alpha)_k = f(\mathbf{x}_k + \alpha \mathbf{p}_k)$, falling back on the secant method if necessary.

Chapter 3

Derivative-Free Optimisation and the Parameter Cascade

The Parameter Cascade requires the computation of derivatives for the various levels of the problem to perform optimisation. However, computing these can be a time-consuming and complex task. In some cases the derivative might not even exist. In this chapter, the use of derivative-free methods for optimisation. Derivative-free methods are used to tackle a series of increasingly complex problems, culminating in fitting one level of the Parameter Cascade without derivatives.

3.1 Overview of Quadratic Optimisation Methods

A large class of numerical optimisation methods rely on constructing a quadratic approximation to objective function $f(\theta)$. Given an iterate θ_n and possibly some associated data, a quadratic approximation $m_n(\theta)$ to the objective function is constructed. The next iterate θ_{n+1} is then found by minimising $m_n(\theta)$. Constructing the approximate quadratic and then minimising it tends to be straightforward. If the next iterate $\theta_{n+1}(\theta)$ is unsatisfactory, a new quadratic model function $m_{n+1}(\theta)$ is minimised, producing a new iterate θ_{n+2} . Ideally, the θ_n will approach the optimal point and the sequence of quadratic models will become increasingly accurate approximations so that the process can be repeated until convergence. [22]

3.1.1 Newton's Method

The Newton-Raphson Method is a well-known member of this class. Newton's method constructs the approximation using a second-order Taylor expansion around θ_n :

$$f(\theta) \approx m_n(\theta) = f(\theta_n) + f'(\theta_n)(\theta - \theta_n) + \frac{1}{2}f''(\theta_n)(\theta - \theta_n)^2$$

It is not difficult to show that the critical point of $m_n(\theta)$ is given by $\theta_{n+1} = \theta_n - f'(\theta)/f''(\theta)$, which is the usual Newton formula [12, 16, 17, 22].

For a point close to θ_n , the difference between $f(\theta)$ and $m_n(\theta)$ is roughly equal to $[f'''(\theta_n)/3!](\theta - \theta_n)^3$ so long as $f(\theta)$ is sufficiently well behaved [12]. This formula suggests that if θ_n is close to the optimal point θ^* so that $|\theta^* - \theta_n|$ is sufficiently small, then $|\theta_n - \theta^*|^3$ will be very small indeed and so the quadratic model will be a very accurate approximation of $f(\theta)$ around θ^* . As a result, θ_{n+1} will be quite close to θ^* . The next model $m_{n+1}(\theta)$ will thus be substantially better than $m_n(\theta)$ at approximating $f(\theta)$ around θ^* , and so θ_{n+2} will be much closer to θ^* than θ_{n+1} . Newton's method converges very rapidly so long as one is sufficiently close to θ^* to start with. In fact, it can be shown that Newton's method exhibits *quadratic convergence* subject to technical conditions. This means that $|\theta_{n+1} - \theta^*| \approx C|\theta_n - \theta^*|^2$. For example, if the error in the first iterate is approximately 0.1, the next iterate will have error on the order of 10^{-2} , the next will have error on the order of 10^{-4} , and so on [17, 22].

Newton's method is a very effective estimation algorithm so long as the derivatives $f'(\theta)$ and $f''(\theta)$ can be computed, and so long as the initial starting value is not too far from the optimal value. Choosing a good initial value is thus very important. For maximum likelihood estimation for example, a method of moments estimator or the median could be used to provide an initial starting value.

3.1.2 Secant Method

If the second derivative is difficult to calculate, one can approximate it with a difference quotient instead [12, 22]¹:

$$f''(\theta) \approx \frac{f'(\theta_n) - f'(\theta_{n-1})}{\theta_n - \theta_{n-1}} \quad (3.1)$$

This leads to the quadratic approximation:

$$m_n(\theta) = f(\theta_n) + f'(\theta_n)(\theta - \theta_n) + \frac{1}{2} \left(\frac{f'(\theta_n) - f'(\theta_{n-1})}{\theta_n - \theta_{n-1}} \right) (\theta - \theta_n)^2$$

And the update formula:

$$\begin{aligned} \theta_{n+1} &= \theta_n - \left[\frac{f'(\theta_n) - f'(\theta_{n-1})}{\theta_n - \theta_{n-1}} \right]^{-1} f'(\theta_n) \\ &= \theta_n - \frac{f'(\theta_n)[\theta_n - \theta_{n-1}]}{f'(\theta_n) - f'(\theta_{n-1})} \\ &= \frac{\theta_{n-1}f'(\theta_n) - \theta_n f'(\theta_{n-1})}{f'(\theta_n) - f'(\theta_{n-1})} \end{aligned}$$

The Secant Method is straightforward to implement, and only requires first derivatives. Relying on $f(\theta_n)$, $f'(\theta_n)$ and $f'(\theta_{n-1})$ instead of $f(\theta)$, $f'(\theta_n)$ and $f''(\theta_n)$ has a drawback however. The Secant Method's model is less accurate because θ_{n-1} tends to

¹The Secant Method is denoted as the *Method of False Position* in [12]

be further from θ^* than θ_n . More formally, the error for the model is roughly equal to $[f'''(\theta_n)/3!](\theta_n - \theta)^2(\theta_{n-1} - \theta)$. If the sequence is converging to θ^* , substituting in the $(\theta - \theta_{n-1})$ term inflates the error relative to Newton's Method and acts as a drag on convergence. It can be shown that the Secant Method only has a convergence rate of 1.618, but avoiding the cost of computing a second derivative on each step means that more iterations can be completed in a given period of time. The Secant Method is comparable with Newton's Method, and can be faster if computing the second derivative is difficult.

The Secant Method is a widely used method that provides a good trade-off between convergence speed and ease of implementation. R's `optim` routine uses a multivariate generalisation of the Secant Method for gradient-based optimisation [25]. Multivariate versions of the Secant Method for optimisation are usually referred to as *Quasi-Newton Methods* [16, 17, 22].

For multivariate problems, the second derivative is in the form of a matrix, so there is not enough information to construct a full approximation afresh on each iteration. Rather the approximate Hessian is partially updated using one of several approaches. No matter the exact approach, the approximate Hessians \mathbf{B}_{n+1} are required to satisfy the multivariate secant condition [16, 17, 22]:

$$\nabla \mathbf{f}(\theta_{n+1}) - \nabla \mathbf{f}(\theta_n) = \mathbf{B}_{n+1}(\theta_{n+1} - \theta_n)$$

This is a generalisation of Equation 3.1. A further condition generally imposed in the multivariate case is that the approximate Hessians must be positive-definite. This ensures that the approximate quadratics don't have any saddle points or surfaces on which the second derivative is zero so that a well defined search direction is guaranteed [16, 17, 22].

R's `optim` routine uses the BFGS method to compute the next approximate Hessian [25]. BFGS uses the symmetric matrix \mathbf{B}_{n+1} satisfying the secant condition such that the inverse \mathbf{B}_{n+1}^{-1} minimises a weighted Frobenius distance between itself and the previous inverse \mathbf{B}_n^{-1} . A low memory variant of BFGS known as L-BFGS is also available [22, 25].

3.1.3 Successive Parabolic Interpolation

Parabolic interpolation goes one step further than the Secant Method and dispenses with derivatives entirely. Instead, a model function is constructed by interpolation through the points $(\theta_n, f(\theta_n))$, $(\theta_{n-1}, f(\theta_{n-1}))$, and $(\theta_{n-2}, f(\theta_{n-2}))$ [22, 30].

$$\begin{aligned}
m_n(\theta) = & f(\theta_n) \frac{(\theta - \theta_{n-1})(\theta - \theta_{n-2})}{(\theta_n - \theta_{n-1})(\theta_n - \theta_{n-2})} \\
& + f(\theta_{n-1}) \frac{(\theta - \theta_n)(\theta - \theta_{n-2})}{(\theta_n - \theta_n)(\theta_n - \theta_{n-2})} \\
& + f(\theta_{n-2}) \frac{(\theta - \theta_{n-1})(\theta - \theta_n)}{(\theta_n - \theta_{n-1})(\theta_n - \theta_n)}
\end{aligned}$$

This model has an approximate error of $[f'''(\theta_n)/3!](\theta - \theta_n)(\theta - \theta_{n-1})(\theta - \theta_{n-2})$. By relying on the past two iterates, the rate of convergence is slowed further. Parabolic interpolation has an order of convergence of 1.32.

An issue with parabolic interpolation is providing enough initial points to seed the method[22]. This is more acute for multivariate problems in particular. One approach is to provide enough points at the start and run the algorithm from there. Alternatively, one can start off with just enough points needed to estimate an ascent or descent direction and construct a linear approximation, and then run the optimisation routine using a sequence of linear approximations until there are enough points to construct a parabola. If one is using a linear approximation, one must impose a limit on the maximum distance that the routine can travel on each iteration since linear functions do not have a minimum or maximum and diverge off to infinity.

3.1.4 Discussion

All three approaches are governed by the same fundamental theory of approximating functions by polynomials. The only difference is the precise inputs used to construct an approximation. This means that if a problem is suitable for Newton's Method, the other two methods will very likely perform well. If one applies parabolic interpolation to a sufficiently smooth objective function, then one is in a sense automatically employing Newton's Method even if one made no effort to investigate the differentiability of the objective function.

On the other hand, the methods all share the same fundamental handicap as well, these methods are not guaranteed to converge unless the starting point is close to the optimal value. Local convergence doesn't necessarily imply global convergence. The error terms in the quadratic approximations are all something like $(\theta - \theta_n)^3$. If $|(\theta - \theta_n)|$ and any other error terms are small, the error in the approximation will be much smaller since it is proportional to the product of three such errors. If however the errors are large, their product might be so large that the method fails to converge. [12, 17, 22]

This is less academic than it might seem. Suppose one had a complicated likelihood function $L(\theta)$. Perhaps to evaluate the likelihood one must numerically integrate some kind of complex marginal distribution that depends on θ . Instead of attempting to find explicit formulae for the score and information functions, if one could produce a crude estimate $\hat{\theta}$ and crude estimate of the error $\hat{\sigma}_\theta$, then one could use successive parabolic

interpolation with $\{\hat{\theta}, \hat{\theta} - 2\hat{\sigma}_{\theta}, \hat{\theta} + 2\hat{\sigma}_{\theta}\}$ as a set of starting points. If $L(\theta)$ is in fact a well behaved smooth function, then parabolic interpolation will find the value of θ that maximises $L(\theta)$ fairly quickly. It is necessary to provide plausible starting values for θ because the quadratic model is only certain to be valid if one is already near the optimal value.

3.2 Bisection Methods

In contrast to the methods discussed above, bisection methods tend to be slow, but are guaranteed to ensure consistent and steady progress towards the optimal point provided the function is continuous and does not have more than one local minima or maxima and that the interval of interest is always shrunk by a given amount on each iteration. One starts with an interval $[a, b]$ and a third point c such that $f(c) < f(a)$ and $f(c) < f(b)$. A fourth point d within the interval $[a, b]$ is selected, and $f(d)$ is computed. If d is between a and c , and $f(d) < f(a)$ and $f(d) < f(c)$, then $[a, c]$ becomes the new interval and d becomes the new provisional minimum. If $f(c) < f(d)$, then the new interval becomes $[d, b]$, - c remains the provisional minimum, but the interval has been narrowed. A similar approach applies if d is between c and b . The whole process is plotted in 3.1

The most common bisection method is known as Golden-Section Search, where the point d is chosen so that the width of the new interval is equal to that of the old one divided by 1.618.[17]

3.3 Brent's Method

Brent's Method is a hybrid of successive parabolic interpolation and golden-section search [3]. If parabolic interpolation is failing to provide a sufficiently rapid decrease in the objective function, a bisection step is performed. While the bisection steps might not produce as much progress as the parabolic steps, they are certain to produce a consistent rate of improvement no matter how close the algorithm is to the optimal point, while parabolic interpolation is only certain to work if one is already within a neighbourhood of the optimal point as noted in Section 3.1.4. Brent's method will also perform a bisection step if the interpolating parabola is ill-conditioned, or if a bisection step has not been performed recently.

The hybrid method is robust as a result of the golden section steps, and the parabolic steps ensure it performs well when applied to smooth functions along with a decent starting value.

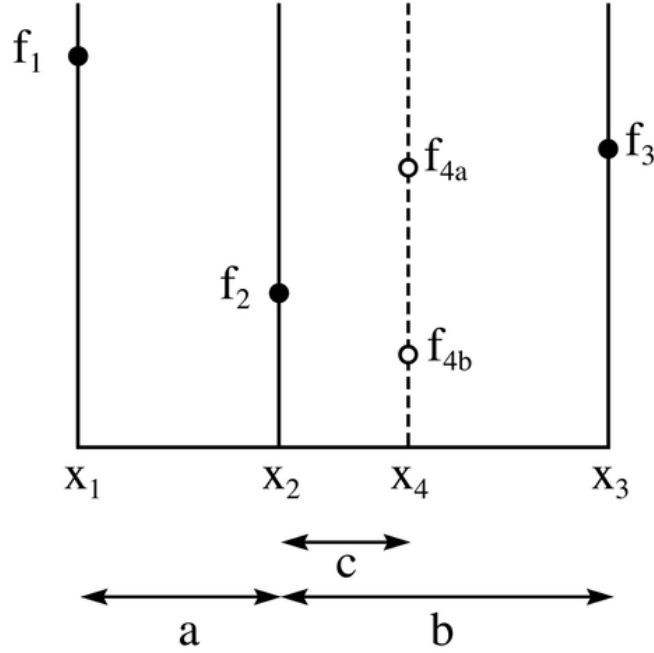


Figure 3.1: Diagram of golden-section search taken from Wikipedia. **TODO:** Find a better plot???

3.4 Estimation Of Parameters For A Standard Cauchy Distribution Using Brent's Method

To illustrate how Brent's method is employed in practice it will be used on a straightforward estimation problem first. Consider the question of fitting a Cauchy distribution to some data. Given n observations x_1, \dots, x_n from an unknown Cauchy distribution, the likelihood function is given by:

$$L(\mu, \sigma) = \prod_{i=1}^n \frac{1}{\pi \sigma \left[1 + \left(\frac{x_i - \mu}{\sigma} \right)^2 \right]}$$

Attempting to maximise this likelihood by the usual method entails solving a fairly complex system of rational equations in μ and σ . Our purpose is to demonstrate that Brent's Method can tackle this problem without much difficulty.

Brent's Method can only optimise a function in one dimension at a time, so it is necessary to attempt to optimise for μ and σ separately. The profile log-likelihood of σ is computed:

$$\ell(\sigma) = \sup_{\mu} \log(L(\mu, \sigma))$$

R can evaluate $\ell(\sigma)$ straightforwardly by using Brent's method to optimise $L(\mu, \sigma)$

with respect to μ and holding σ constant. The function $\ell(\sigma)$ can then be in turn optimised with respect to σ to find the optimal value of σ . This procedure is illustrated in Figure 3.2.

One subtlety with optimising a Cauchy likelihood is that the likelihood function can have multiple local maxima since the likelihood function is the ratio of two multivariate polynomials in μ and σ . To ensure that the algorithm was sufficiently close to the MLE, the median was used as an initial estimate of μ , and half the interquartile range was used as an initial estimate for σ . Given these somewhat crude estimates $\tilde{\mu}$ and $\tilde{\sigma}$, the standard error of the median $\sigma_{\tilde{\mu}}$ is approximately given by:

$$\hat{\sigma}_{\tilde{\mu}} \approx \frac{1}{2f(\tilde{\mu}; \tilde{\mu}, \tilde{\sigma})\sqrt{n}}$$

Where $f(x; \mu, \sigma)$ is the Cauchy density function with location parameter μ and scale parameter σ . The values $\tilde{\mu} \pm 2\hat{\sigma}_{\tilde{\mu}}$ are then used to provide the initial lower and upper bounds for the optimiser. The aim is to construct a confidence interval that is highly likely to contain the MLE for μ (rather than the actual true parameter), but isn't so wide that the interval is in danger of containing multiple local maxima for the likelihood.

Not only can the likelihood be maximised without derivatives, but asymptotic inference can be done without derivatives as well. Given the score function and the Fisher information at the maximum likelihood estimates, it is possible in principle to compute an approximate confidence interval for σ and μ [24]. Instead of analytic methods, one can use finite differences to approximately compute the necessary derivatives to the desired degree of accuracy [9, 19]. This was successful at producing a valid approximation for the profile likelihood, shown as a red dotted parabola in Figure 3.2.

It is thus possible to compute a confidence interval using the Score test. The test statistic $S(\sigma)^2/I(\sigma)$ could be accurately approximated using finite differences. One takes the value of σ for which the test statistic is less than or equal to the appropriate critical value from a chi-squared distribution. By inspecting the plot in Figure 3.3 and then solving for σ , an approximate confidence interval for σ can be computed. It is approximately the case that σ lies in $(0, 2.20)$ with 95 percent confidence.

An important assumption underpinning such asymptotic confidence intervals is that the two term quadratic Taylor expansion based on the score and information functions is valid over the range of interest. This is not the case here as can be seen in the spike in the score statistic on the left caused by the Fisher information changing sign at approximately $\sigma = 2.35$. This indicates that the confidence interval might be wider than the range of for which a quadratic approximation around the MLE is valid, and should perhaps be treated with some scepticism.

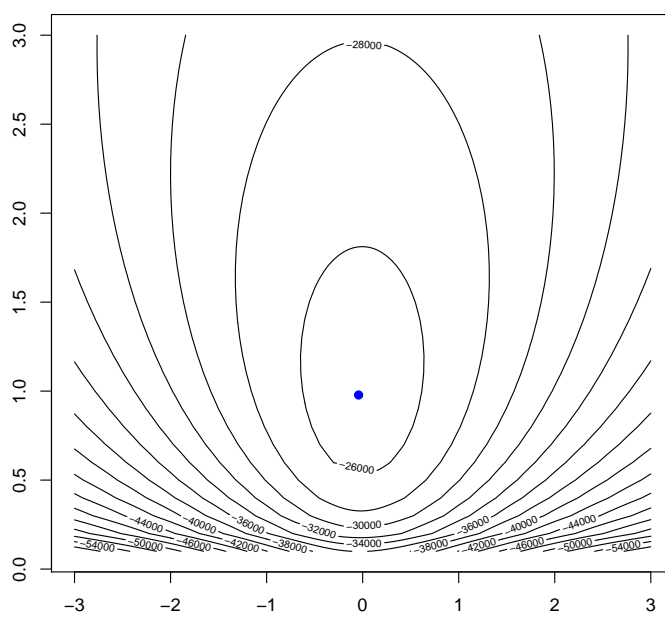
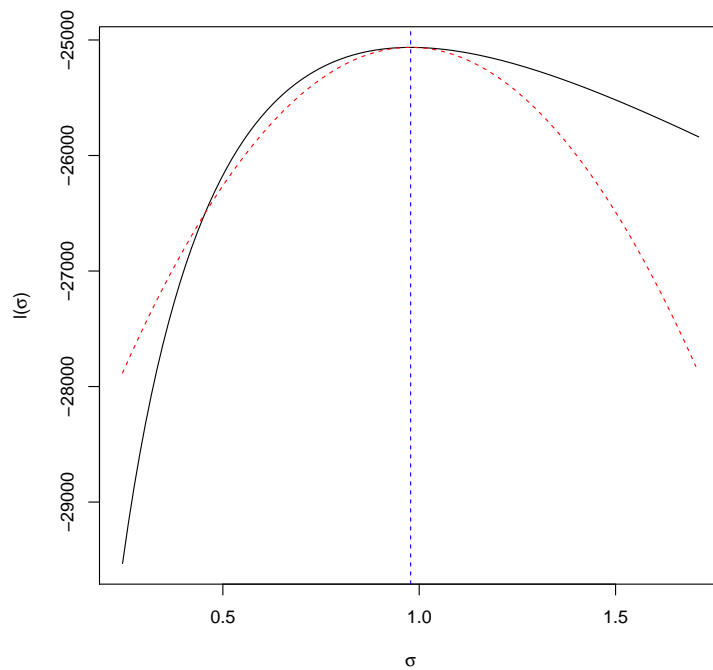


Figure 3.2: Profile log likelihood in σ , and contour plot of the joint log likelihood.

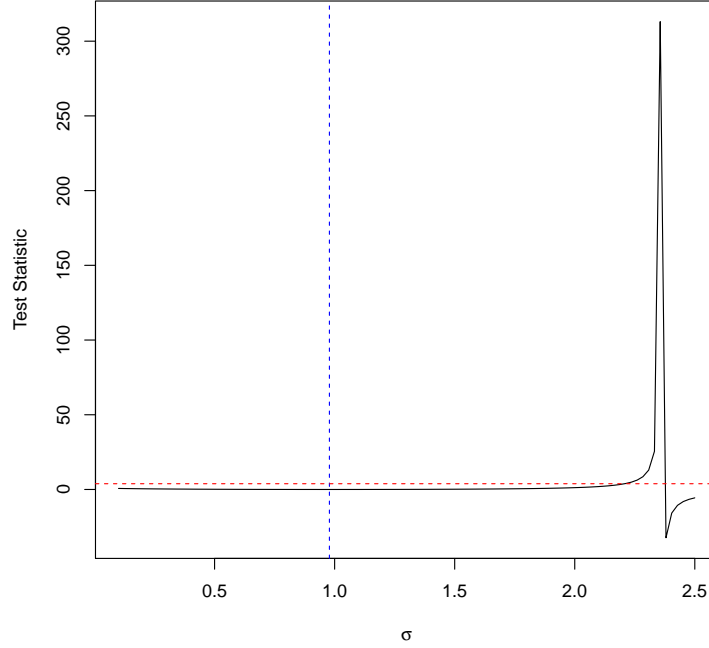


Figure 3.3: Plot of profile score statistic.

3.5 Robust ODE Parameter Estimation

If observations of values from an ODE are subject to heavy-tailed noise such as in the Cauchy case, least squares regression becomes unsuitable. An obvious candidate is L1 regression, which attempts to minimise the sum of the absolute values of the residuals instead of the sum of the squared residuals. An important property of L1 regression is that median is naturally associated with this approach; the sample median of a set of numbers is the constant value that minimises the L1 error just as the sample mean is the constant value that minimises the least squares error[28]². L1 regression can greatly complicate the process of estimation however, because the the function $|x|$ is not everywhere differentiable. This means that the usual gradient-based approaches to nonlinear regression such as gradient descent should not be applied. Even methods that attempt to numerically approximate the derivatives such as parabolic interpolation are either entirely unsuitable at worst, or not guaranteed to converge quickly at best.

Brent’s Method can tackle such problems however, being robust against non differentiability. For nonlinear L1 regression, the objective function tends to be piecewise smooth - between the “kinks”, the function is differentiable and amenable to parabolic interpolation. Once the bisection steps have reached a neighbourhood of the optimal value, parabolic interpolation will find it fairly quickly.

²This is discussed in more detail in the next chapter

Consider for example, the following ODE with $\beta = -0.5$:

$$\begin{cases} y'' - \beta(1 - y^2)y' + y = 0 \\ y(0) = 1 \\ y'(0) = 0 \end{cases} \quad (3.2)$$

This ODE describes a non-linear oscillator, and is representative of quasi-linear mathematical models that can't be tackled by the FDA package or Data2LD. Note that this ODE is of the form $y'' + \beta(y)y' + y = 0$ with $\beta(y) = -\beta(1 - y^2)$. By definition, the linear ODEs usually used in FDA cannot model systems where the $\beta(\cdot)$ terms have y as a dependent variable, they can only model situations where the parameters vary with time alone (and/or space in the case of a linear PDE).

We wish to investigate the problem of estimating β from noisy observations.

The **desolve** package [29] was used to numerically find the values of $y(t)$ at choosen time points $\{t_1, \dots, t_K\}$. The values of $y(t)$ at these points - corrupted by random Cauchy noise - were independently sampled N times. This produced a set of KN observations: $\{y_{11}, y_{12}, \dots, y_{1N}, \dots, y_{K1}, \dots, y_{KN}\}$. Because the data is heavy-tailed, least squares regression is inappropriate. Instead, the goodness of fit associated with a given choice of β was measured by the sum of absolute errors associated with a given choice of β :

$$SAE(\beta) = \sum_{i=1}^K \sum_{j=1}^N |y(t_i; \beta) - y_{ij}|$$

Here $y(t; \beta)$ denotes the solution of Equation 3.2 for a given choice of β . To evaluate $SAE(\beta)$ at a given value of β , it is necessary to use **desolve** to numerically find the values of $y(t_i|\beta)$. Brent's method was used to find the number $\hat{\beta}$ that minimised $SAE(\beta)$. Figure 3.4 shows the original curve, the generated points, the realisation of $SAE(\beta)$, and the fitted curve generated by $\hat{\beta}$

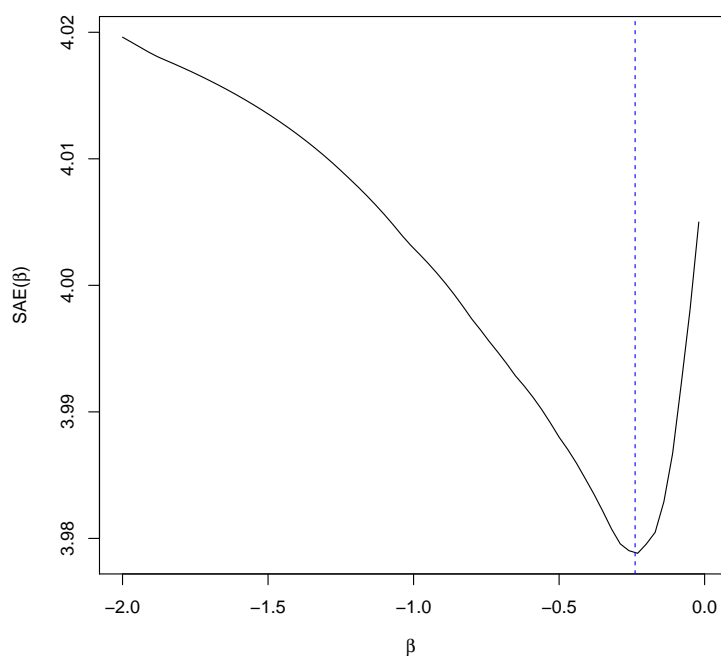
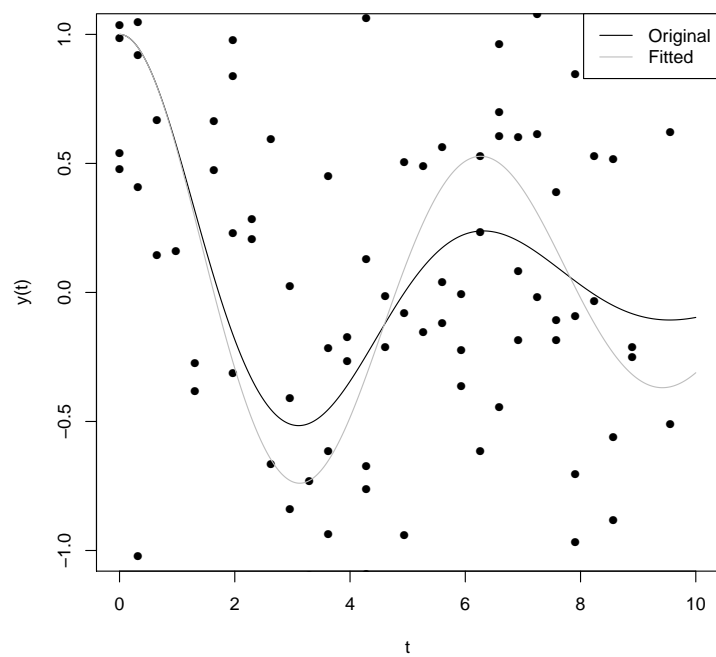


Figure 3.4: Original curve, fitted curve, and objective function.

3.6 The Parameter Cascade and Brent's Method

Recall that the Parameter Cascade has three levels.

First, the inner problem. There is a given functional $J(f; \theta, \lambda)$ that takes a function f and associated parameters θ and λ and returns a real number. Usually, the function f is represented by a vector of coefficients with a given associated basis. The function $\hat{f}(t|\omega, \lambda)$ that optimises $J(\cdot; \theta, \lambda)$ is then found. Outside of toy cases, this problem cannot be solved analytically. The problem is nearly always solved numerically by restricting the space of functions to the span of some set of chosen basis functions and optimising over that.

This in turn defines the middle problem, $H(\theta, \hat{f}(t|\omega, \lambda); \lambda) = H(\theta; \lambda)$, which is usually defined as the least squares error associated with the optimal f given θ and λ :

$$H(\theta; \lambda) = \sum [x_i - \hat{f}(t_i|\omega, \lambda)]^2$$

As suggested in the previous section on fitting an ODE with Cauchy noise, the middle error might be another loss function besides least squares error such as the sum of absolute errors. As before, value of θ that optimises $H(\cdot)$ holding λ constant, defined by $\hat{\theta}(\lambda)$, is computed.

And finally, the outer problem attempts to determine the value of λ that minimises the prediction error (generalisation error) by minimising another function $F(\lambda, \hat{\theta}(\lambda), \hat{f}(t|\omega, \lambda)) = F(\lambda)$. There are several plausible choices for $F(\cdot)$, one could use leave-one-out cross-validation, one could partition the data set into a training set and a validation one, and let $F(\lambda)$ be the associated error for the validation set, one could use Generalised Cross-Validation. This criterion is in turn optimised to find the optimal λ .

Note that the three levels are somewhat isolated from each other and only interact by exchanging parameters downwards and optimal values upwards. The middle function $H(\cdot)$ for example only requires the value of the optimal $f(\cdot)$ evaluated at the chosen points t_i , and does not care about how these values were found or how $f(\cdot)$ is represented.

The inner problem consists of finding a function that minimises a certain criterion for a given set of parameters. As previously discussed, the complexity of such problems can increase fairly rapidly and require a considerable degree of non-Statistical expert knowledge and often must be essentially developed from scratch if the differential penalty changes too much. It is thus desirable that the inner problem can be solved with already existing methods and tools such as the FDA package or Data2LD to avoid the effort of having to develop one's own. Ideally, it should be possible for one to plug in existing code that can compute $H(\cdot)$ and the optimal function as required.

There is thus a considerable degree of potential modularity present in the Parameter Cascade that is not fully investigated in Ramsay and Cao's paper [4], and research that inherits that framework. The Parameter Cascade can be adapted to heavy-tailed errors for example, by using appropriate loss functions for the various levels of the problem.

Not only is it good research practice to have mostly independent components that can be tackled and verified separately before being combined, it is also good practice

from a software engineering perspective because the potential for complex interactions between different parts of code is reduced. This tends to save on debugging and testing requirements, which can be quite high when implementing codes for FDA.

The Data2LD package is fairly tightly coupled. Rather than use R's built-in routines to implement the Quasi-Newton algorithm for example to optimise the associated middle problem, the authors wrote their own code. With Brent's method however, there is more separation, which makes it very easy to build optimisation routines on top of other code. This substantially elides the cost and effort of tackling the inner problem and allows one to concentrate on the statistical questions such as fitting the model to data.

Melanoma Data

This derivative free optimisation strategy was applied to fitting the melanoma dataset with a parameterised linear differential operator:

$$L_\omega = D^2 - \omega^2 D^4. \quad (3.3)$$

The inner problem consists of finding the function $f(t)$ that minimises a penalised regression problem of the form:

$$PENSSE(f; \omega, \lambda) = \sum (x_i - f(t_i))^2 + \lambda \int |L_\omega f(t)|^2 dt$$

The penalty term measures the extent to which a function lies outside of the span of the functions $\{1, t, \cos(\omega t), \sin(\omega t)\}$.

The FDA package has routines that can do the numerical work of fitting the data with differential penalty given in (3.3) for given choices of λ and ω , and then report the associated mean square error.

Using Brent's method, the function $H(\omega; \mathbf{x}, \lambda)$ can be optimised with respect to ω for a given fixed λ . In turn, the outer objective function can be parameterised in terms of λ and the associated optimal choice of ω . This defines an objective function that can be again optimised to find the optimal choice of λ .

For ω , Figure 3.5 shows that the error is not particularly sensitive to small deviations from the optimal value even for fairly high values of λ . This suggests that the fitted curve will be adjusted to ensure no substantial increase in the error so long as ω isn't altered too much from the optimal value.

Heuristically speaking, a flat objective function in the neighbourhood of the optimal point as can be seen in Figure 3.5 increases the uncertainty in estimation because it is more difficult to argue that the optimal value is definitively better than adjacent ones. The loss function associated with a given fitting problem only approximates the 'true' loss function as the sample size goes to infinity.

If λ is set too low, the optimal value of ω is numerically indistinguishable from zero. This is the case when ω is optimised for the value of λ that minimises the GCV, Brent's method reports zero as the optimal value to within its default tolerance.

For λ , the curve has two critical points, with an asymptote as λ tends to infinity.

A huge advantage of this approach compared to Data2LD's use of quasi-Newton methods is that it allows for the use of more robust loss functions since no use at all is made of derivatives.

Suppose one wanted to choose ω to minimise the Median Absolute Deviation - $\text{median}(|y_i - \hat{f}(t_i|\omega, \lambda)|)$ - instead of the least squares error. This loss function is chosen instead of the usual L1 error for the sake of demonstration because the L1 error might sometimes be tackled using a generalised version of gradient descent known as the subgradient method, while getting any kind of a derivative for MAD is difficult. It is quite simple, one just replaces the code that computes the least squares error with a few lines of R code that computes the MAD and run the optimisation routine again. It can be seen in Figures 3.6 and 3.7 that the MAD gives similar results to the usual least squares criterion, which suggests that both estimators are mutually consistent with each other.

3.7 Implicit Filtering

There are other methods for fitting without derivatives besides Brent's method and Parabolic Interpolation. One method that was investigated is known as the Implicit Filtering algorithm. Implicit Filtering will only be briefly covered here because it proved to be inferior to other optimisation methods.³

3.7.1 Description Of Implicit Filtering

The implicit filtering algorithm is designed for optimising problems where the exact value of the objective function $f(\cdot)$ is unavailable. Instead, $f(\cdot)$ can only be evaluated up to an arbitrary degree of accuracy. It is assumed that there is parameter h which controls the degree of the accuracy - the lower h , the lower the error. It is usually the case that getting a higher degree of accuracy means a higher run time.

For example, if the objective function is an expectation being approximated using a Monte Carlo method for example, it would be reasonable to set $h = 1/\sqrt{N}$ where N is the number of samples used so that the standard deviation is proportional to h . On the other hand if the expectation were being approximated by numerical integration, h would be set to the step size.

Let $f(\mathbf{x}; h)$ denote the result of approximately evaluating $f(\cdot)$ at the point \mathbf{x} with precision level h . To generate a search direction, Implicit Filtering uses an approximation $\nabla_h f(\mathbf{x})$ to the gradient $\nabla f(\mathbf{x})$ that depends on h . The simplest such approximation employs forward differencing to approximate the gradient:

$$[\nabla_h f(\mathbf{x})]_i = \frac{f(\mathbf{x} + h\mathbf{e}_i; h) - f(\mathbf{x}; h)}{h} \quad (3.4)$$

³Interested readers are pointed towards [13, 14, 22].

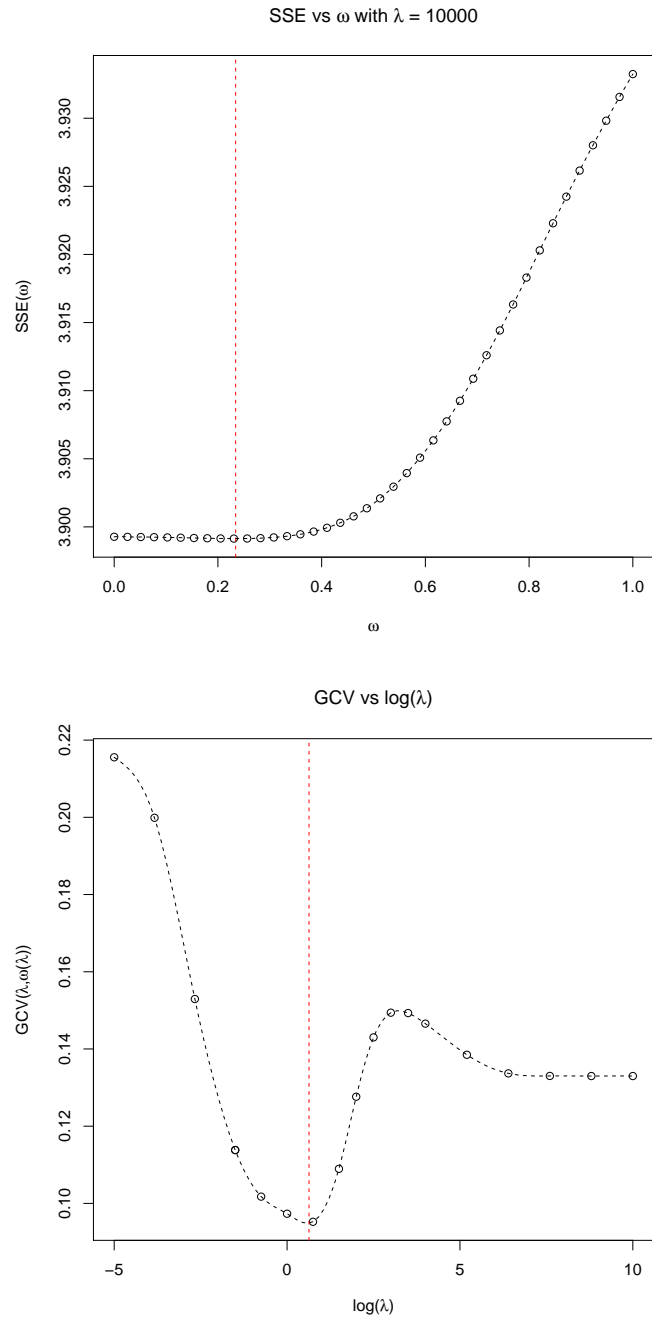


Figure 3.5: Plots of the middle and outer optimisation problems.

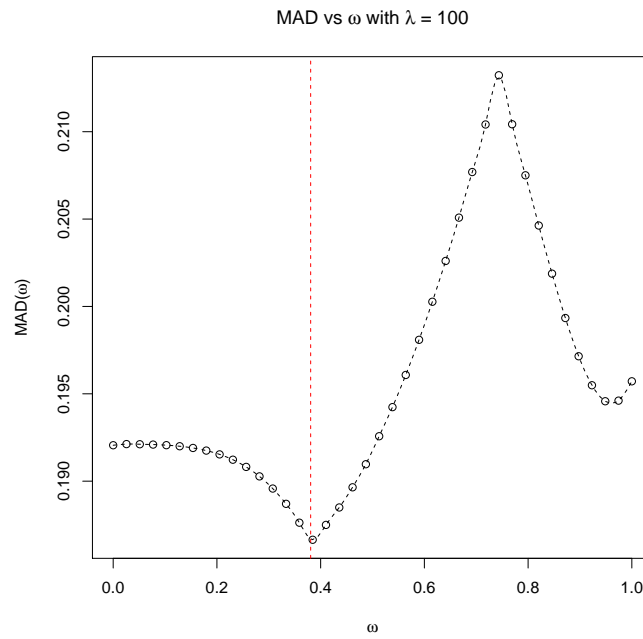


Figure 3.6: Plot of the middle optimisation problem with MAD used as a loss function

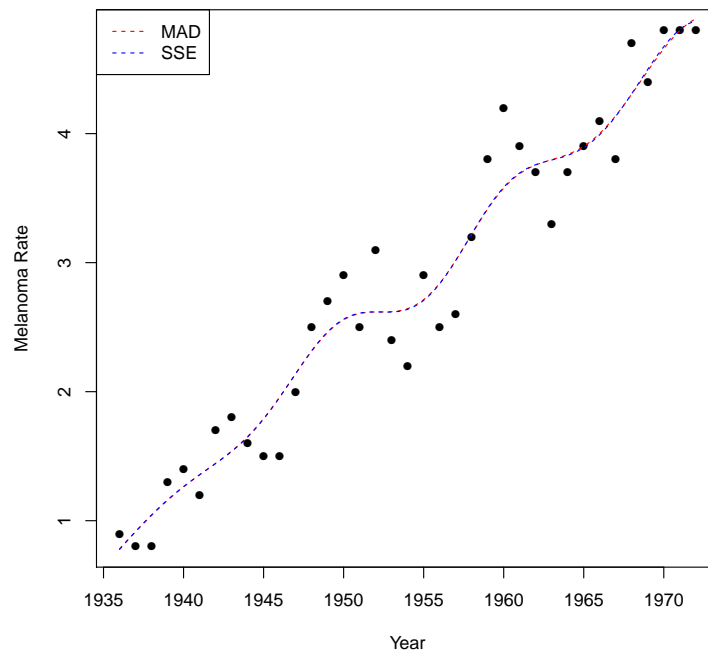


Figure 3.7: Comparison of fits for MAD and SSE criteria for middle problem

Here $[\nabla_h f(\mathbf{x})]_i$ denotes the i th component of $\nabla_h f(\mathbf{x})$ and \mathbf{e}_i is the i th basis vector. This approximate gradient is then used to define a search direction. The algorithm proceeds to conduct a line search along this direction until a point that achieves a sufficient reduction is found, which then becomes the latest iterate, and a new search direction is computed.⁴

In the event of any of the following occurring, it is deemed that more precision is needed:

- A point achieving sufficient reduction cannot be found after a maximum number of iterations.
- A clear descent direction cannot be identified.
- The approximate gradient is of a similar order of magnitude to h , so that one can't be confident that true gradient isn't in fact zero.

If any of these conditions hold, the value of h is shrunk so that $h \leftarrow \delta h$ with $0 < \delta < 1$. The algorithm then proceeds again with this higher level of precision.

The algorithm terminates when the change in the value of the objective function produced by reducing the value of h and running again is within a chosen tolerance.

The sequence of approximate values returned by Implicit Filtering is monotone decreasing so that if $m < n$ then $f(\mathbf{x}_m; h_1) \geq f(\mathbf{x}_n; h_2)$, where h_1 is the precision used with \mathbf{x}_m and h_2 is the precision used with \mathbf{x}_n . Bear in mind however that since Implicit Filtering only ever approximately evaluates the objective function, it is not necessarily the case that $m < n$ implies $f(\mathbf{x}_m) \geq f(\mathbf{x}_n)$.

⁴More precisely, the next point \mathbf{x}_{k+1} is required to satisfy a condition of the form $f(\mathbf{x}_{k+1}; h) \leq f(\mathbf{x}_k; h) - c\mathbf{d}_k^\top [\nabla_h f(\mathbf{x}_k)]$, where $\mathbf{d}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ and $0 < c < 1$. Note that there is no requirement to decrease the norm of the approximate gradient $\nabla_h f(x)$.

3.7.2 Implicit Filtering and Brent's Method

It was determined that there are many disadvantages with Implicit Filtering compared to Brent's Method:

- Implicit Filtering is much more complex to code, and is thus more difficult to maintain and debug. The R code used to fit the ODE (3.5) by Implicit Filtering came out at a little over 300 lines long. The code in the Data2LD package that performs optimisation is over 600 lines long. Code that uses Brent's Method tends to be much shorter.
- Implicit Filtering proved to be very slow when applied to the test problem in Section 3.7.3 below.
- The results of the fitting are sensitive to the value of the shrink factor δ chosen.
- It can be necessary to add a penalty term to the objective function to ensure convergence.

3.7.3 Using Implicit Filtering to Fit an ODE to the Melanoma Data

To test Implicit Filtering, the following quasi-linear fourth order ODE was fitted to the melanoma data:⁵

$$y^{(4)} = \mu^2[1 - \sin(\pi y'')]y''' - \omega^2 y'' \quad (3.5)$$

The objective function used was a penalised sum of squared errors of the form:

$$PENSSE(f(t), \omega, \mu) = \rho \sum [y_i - f(t_i)]^2 + (1 - \rho) \int |f''(t)|^2 dt$$

The value of $PENSSE$ is influenced by ω and μ because $f(t)$ is required to be a solution of (3.5) with given values of ω and μ . The Implicit Filtering algorithm will not converge correctly without the penalty term as illustrated in Figure 3.9.

To compute $PENSSE$, the package `deSolve` was used to numerically solve (3.5) with the appropriate values of ω and μ . The precision factor h determined the stepsize used. The $\int |f''(t)|^2 dt$ term was approximated by taking the vector of computed values of $f''(t)$ returned by `deSolve`, and then finding the sum of squared values. As $h \rightarrow 0$, this approximation becomes more and more accurate.

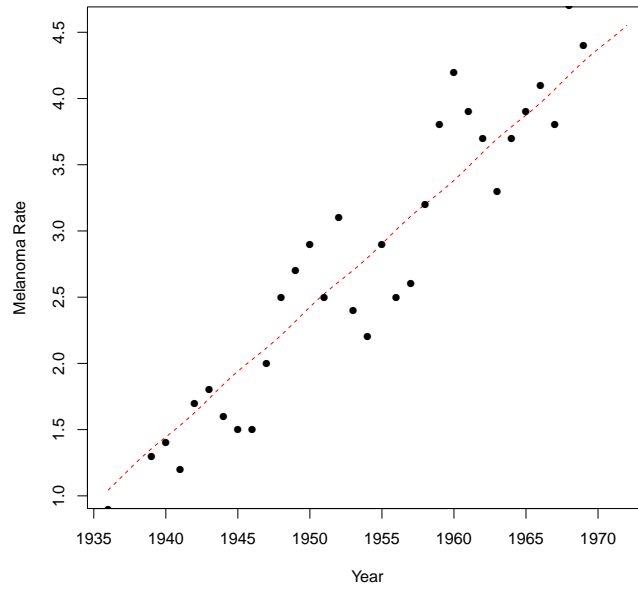
As can be seen in Table 3.1, the algorithm takes a long time to run. It can be seen in both the table and Figure 3.8 that changing the value of δ can introduce qualitative changes in behaviour. The algorithm is much quicker for $\delta = 0.9$, presumably because

⁵The version of Implicit Filtering used is actually a modified version of that described above. A Quasi-Newton algorithm was used instead of naive gradient descent to compute search directions, and central differences were used to estimate the gradient instead of forward differences as in (3.4).

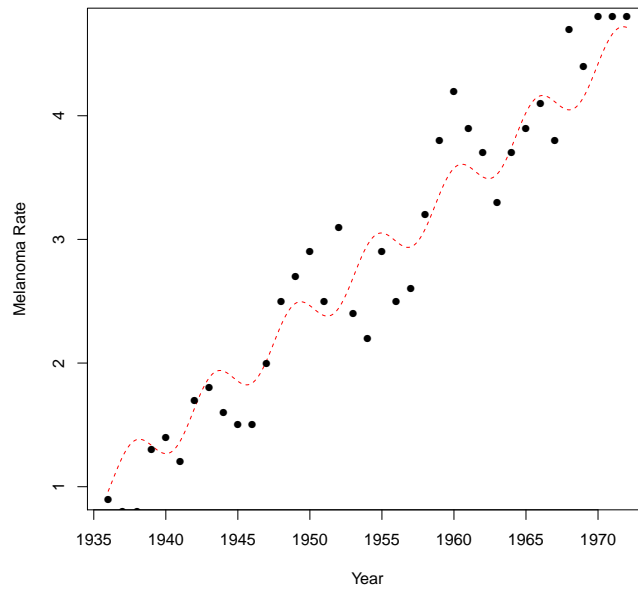
the algorithm is converging to a different fit than for the other cases. For the fastest case where $\delta = 0.9$, 200 values of the *PENSSE* sequence are generated before the sequence converges to within a tolerance of 10^{-4} . This statistic substantially underestimates the actual amount of work done since Implicit Filtering rejects many evaluations as being inadequate in the course of its execution and further evaluations are needed to compute the approximate gradients $\nabla_h f(\cdot)$. For case where $\delta = 0.9$, *PENSSE* was computed over 3700 times with various values of h used.

δ	Running Time (Seconds)	Running Time (Minutes)
0.7	1717.807	28.63
0.8	1611.459	26.85
0.9	1013.165	16.88

Table 3.1: Time taken for Implicit Filtering to fit (3.5) to the melanoma data for various values of δ .



(a)



(b)

Figure 3.8: Fitting the ODE (3.5) to the Melanoma data. The exact value of the shrink value δ effects the fit the Implicit Filtering algorithm converges to. For $\delta = 0.7$ the computed fit in (a) resembles a straight line, but $\delta = 0.9$ results in a sinusoidal plus linear trend as can be seen in (b).

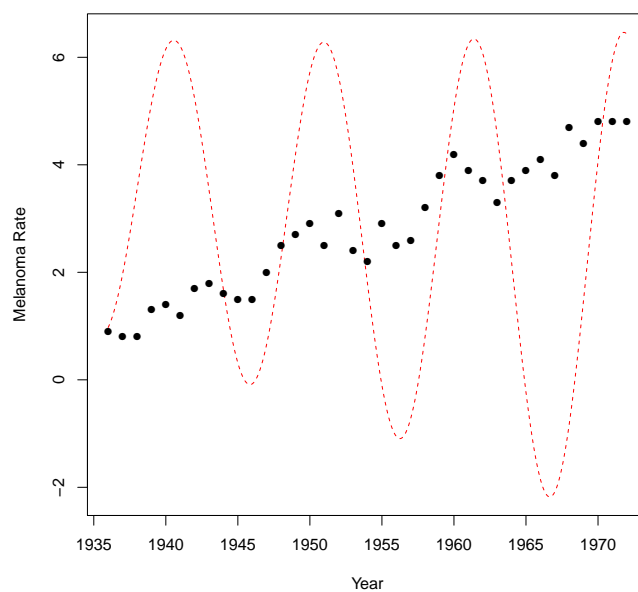


Figure 3.9: Without a penalty term, Implicit Filtering entirely fails to fit the ODE (3.5) to the melanoma data.

Chapter 4

A Two Level L_1 Parameter Cascade Using the MM Alogrithm.

In the previous chapter, Brent's method was introduced and was used to tackle Parameter Cascade problems where the middle level uses a loss function which is difficult or to differentiate or simply has no well-defined derivative everywhere. It was remarked that Brent's method ensures that the different elements of the Parameter Cascade tend to be loosely coupled from each other and this allows one to combine different fitting methodologies for different levels with straightforwardly.

Here these ideas are developed futher. First, the L_2 based penalised fitting method is extended to the L_1 case. This new method is then used alongside Brent's method to implement a two level Paremeter Cascade with L_1 loss functions at both levels.

4.1 L_1 Estimation for the Inner Problem

Brent's method is designed to optimise real-valued functions over a real interval. In the previous it was extended to functions that take more than one real argument by optimising over each coordinate seperately when the Cauchy likelihood was optimised. However, there is no guarantee that this approach will perform well, and it can even fail entirely for functions that have an exotic topography or multiple local optima arranged unusually¹. Even in the best case, optimising over each coordinate generates its own optimisation subproblem, which has the cumulative effect of increasing the running time of the alogrithm. Brent's method further requires the specification of a bounding box that contains the optimal point since it uses bisection, and that is harder and harder to do as the number of dimensions increases. All of these considerations mean that Brent's Method is highly unsuitable for peforming L_1 fitting over a space of functions which tend to have a large number of dimensions - by definition, there is one

¹Consider for example the problem of finding the minimum of the function $f(x) = x \sin(x)$ over the interval $[0, 13]$. It is easy to see that the minimum is not on the boundary points of the interval because $f(0) = 0$, $f(6) = -1.67$, and $f(13) = 5.45$. Brent's method fails to find the minimum. It claims the optimal value is given by $f(4.9) = -4.81$ though $f(11) = -10.99$.

dimension introduced for each basis function used. Likewise, the non-differentiability of the absolute value function means that other approaches that implicitly rely on differentiability such as parabolic interpolation are inadvisable.

Instead a different approach will be employed, a generalisation of the Iteratively Reweighted Least Squares algorithm for computing the L_1 median of a set of N items $\{x_1, \dots, x_N\}$ to which an L_1 norm can be associated. The L_1 median is defined as the object x that minimises $\sum_{i=1}^N |x - x_i|$. We will start by describing how IWLS can be used to compute the L_1 median of a set of real numbers. We will further show that this as an example of what is known as an MM algorithm, and then proceed to straightforwardly extend this MM algorithm to produce a modified Penalised Sum of Squares problem that can be iteratively solved and reweighted to find the function that minimises a penalised L_1 norm.

4.1.1 An MM Algorithm For Computing the Median

Suppose that given a set of numbers $\{x_1, \dots, x_N\}$, one wished to find the number x that minimised the L_1 distance between them:

$$SAE(x) = \sum_{i=1}^N |x_i - x|$$

It is well known that $SAE(x)$ is minimised by the sample median of the numbers [28]². The usual approach to computing the sample median - sorting the numbers and taking the one in the middle - can't be generalised to FDA problems, so we will use a different approach. The main difficulty is that the function $SAE(x)$ is not everywhere differentiable, which means that the usual derivative-based techniques such as gradient descent or Newton's method can't work. Instead an approach known as Majorise-Minimise or the MM Algorithm will be used [11, 16, 17]. For a given iterate x_n , a function $M(x|x_n)$ is required with the following properties:

$$\begin{aligned} M(x|x_n) &\geq SAE(x) \\ M(x_n|x_n) &= SAE(x_n) \end{aligned}$$

The function $M(x|x_n)$ is said to majorise $SAE(x)$. The next iterate x_{n+1} is then found as the value of x that minimises $M(x|x_n)$. Thus:

²This is asserted without proof in [28], probably because the proof tends to be simultaneously awkward but trivial to those familiar with it. The amount of work required to demonstrate that the sample median minimises $SAE(x)$ is greatly reduced if one notes that $SAE(x)$ is a convex function in x . Any local minimum of a convex function is also a global minimum [2], so one only needs to show that for any sufficiently small ϵ that $SAE(\bar{x}) \leq SAE(\bar{x} \pm \epsilon)$, where \bar{x} is the sample median.

$$\begin{aligned}
SAE(x_{n+1}) &\leq M(x_{n+1}|x_n) \\
&\leq M(x_n|x_{n+1}) \\
&= SAE(x_n)
\end{aligned}$$

If such a function $M(x|y)$ could be determined such that $M(x|y)$ would be straightforward to minimise, it is then possible to easily produce a sequence of iterates x_n such that $SAE(x_{n+1}) \leq SAE(x_n)$ for all n . This pattern of monotone improvement in the objective function is similar to the EM Algorithm. In fact, the EM algorithm is a special case of the MM algorithm[33]³.

The most important question associated with the MM algorithm is the construction of the majorising function because this tends to take up the bulk of the effort. Once the majoriser has been found, the algorithm is generally straightforward to implement, as will be seen shortly [11, 18]. Verifying a potential majoriser is usually straightforward, finding one in the first place is more difficult. The EM algorithm for example takes advantage of the probabilistic structure of the problem and Jensen's inequality⁴. For an L_1 problem, the usual approach is to employ the Arithmetic Mean-Geometric Mean inequality [18]. Only the AM-GM inequality in its simplest form is required here, that the geometric mean of two numbers is less than or equal to their arithmetic mean:

$$\sqrt{xy} \leq \frac{x+y}{2}$$

It's worth noting that the AM-GM inequality is in fact a special case of Jensen's Inequality since the log function is concave:

$$\begin{aligned}
\log\left(\frac{x+y}{2}\right) &\geq \frac{\log x}{2} + \frac{\log y}{2} \\
&= \log \sqrt{x} + \log \sqrt{y} \\
&= \log \sqrt{xy}
\end{aligned}$$

It is possible to exploit the AM-GM inequality to majorise an L_1 regression problem by a weighted L_2 problem. One can represent the L_1 norm as a geometric mean, which then allows for the L_1 norm to be majorised and separated by a weighted sum of squares. Given an iterate x_n , the AM-GM inequality implies that:

³When applied to maximisation problems, MM instead stands for Minorise-Maximise. This case is the same except the surrogate function is required to be less than or equal to the objective function and it is maximised on each iteration. Thus, each iteration drives the objective function upwards.

⁴The EM algorithm is intended to maximise the log-likelihood and drives it upwards on each iteration, so it's an example of a Minorise-Maximise algorithm.

$$\begin{aligned}
|y - x| &= \sqrt{(y - x)^2} \\
&= \sqrt{\frac{(y - x)^2}{|y - x_n|} |y - x_n|} \\
&\leq \frac{1}{2} \left(\frac{(y - x)^2}{|y - x_n|} + |y - x_n| \right)
\end{aligned}$$

This in turn implies that:

$$\begin{aligned}
\sum |x_i - x| &\leq \frac{1}{2} \sum \left(\frac{(x_i - x)^2}{|x_i - x_n|} + |x_i - x_n| \right) \\
&= \frac{1}{2} \sum \frac{(x_i - x)^2}{|x_i - x_n|} + \frac{1}{2} \sum (|x_i - x_n|)
\end{aligned}$$

The L_1 problem is thus majorised by a weighted least squares problem. The $\frac{1}{2} \sum |x_i - x_n|$ term is constant with respect to x , so neglecting it makes no difference to the choice of x that is optimal. Likewise, multiplying the weighted least squares problem by a positive constant doesn't change the optimal value either, so the $\frac{1}{2}$ term can be eliminated by multiplying by 2. The optimal value x_{n+1} can thus be found by minimising this weighted least squares score:

$$\sum \frac{(x_i - x)^2}{|x_i - x_n|}$$

The algorithm thus consists of finding the value of x that minimises the least squares error inversely weighted by the residuals from the previous iteration.

4.1.2 Penalised L_1 Fitting

For the case of penalised regression, the penalised sum of absolute errors is defined by:

$$PENSAE(f|\theta, \lambda) = \sum |x_i - f(t_i)| + \lambda \int |Tf|^2 dt$$

Here T is used instead of L to denote a differential operator that might not necessarily be linear ⁵. As before, this can be majorised by a weighted sum of a residual-weighted penalised sum of squared errors, and a vestigial $\sum |x_i - f_n(t_i)|$ term that is only included for completeness and can be safely ignored in the course of the actual optimisation.

⁵In some situations T could even be an integral operator. This could easily be the case for example if the observed values were measured velocities of a vehicle, and the penalty was intended to impose constraints on quantities such as the distance travelled or fuel consumed

$$PENSAE(f) \leq \frac{1}{2} WPENSSE(f|f_n, \theta, 2\lambda) + \frac{1}{2} \left(\sum |x_i - f_n(t_i)| \right) \quad (4.1)$$

$$= \frac{1}{2} \left(\sum \frac{[x_i - f(t_i)]^2}{|x_i - f_n(t_i)|} + 2\lambda \int |Tf|^2 dt \right) + \frac{1}{2} \left(\sum |x_i - f_n(t_i)| \right) \quad (4.2)$$

To find the function that minimises the penalised L_1 error, one repeatedly finds the function that minimises $WPENSSE$ with the previous set of residuals used as inverse weights. This produces a sequence of fitted functions for which the penalised sum of absolute errors is monotonically forced downwards.

4.1.3 Discussion

The sequence of penalised errors $PENSAE(f_n)$ is monotone decreasing but cannot be less than zero, so it is a bounded monotone sequence. The Monotone Convergence Theorem for sequences of real numbers[27] thus guarantees that a given generated sequence $PENSAE(f_n)$ will always converge to a limit. There are two caveats. First, the sequence might converge to a different point depending on the starting values - there is no guarantee that the sequence will converge to the lowest possible value of $PENSAE$. Second, there is no guarantee that the underlying sequence of functions will converge, and may just oscillate between several points. The sequence $-1, 1, -1, \dots$ does not converge but the associated sequence of absolute values $1, 1, 1, \dots$ does.

This approach of associating the objective function with surrogate problem that is a more standard is employed in the literature on the EM Algorithm. For example, in the introductory chapter of [21], the authors discuss how a multinomial estimation problem can be transformed into a binomial problem with missing data by artificially splitting one of the cells; they then construct a simple iterative EM scheme that can then be repeatedly iterated to estimate parameters for the original multinomial. They even remark that the surrogate problems associated with EM algorithms tend to be easy to solve using existing tools in the field. Likewise, the L_1 problem has been replaced here with a surrogate sequence of weighted L_2 problems that can easily be solved using the **FDA** package and so actual code for implementing the L_1 regression is brief.

The literature on the MM algorithm remarks that it is simple to implement and good at tackling high dimensional penalised regression, though convergence can be slow [33]. These claims are borne out when the convergence of the method is examined in Section 4.1.4 below.

Estimating the Worst Case Running Time

The arguments used to prove the Monotone Convergence Theorem can be extended a little to give a crude and weak estimate of how long the MM algorithm runs before the associated $PENSAE(f_n)$ sequence converges to within a given tolerance. Suppose there were a monotone decreasing sequence bounded between a and b so that $a \geq$

$x_0 \geq x_1 \geq \dots \geq x_n \geq \dots \geq b$. By the Monotone Convergence Theorem, this sequence converges to some value x^* such that $a \geq x^* \geq b$ and $x_n \geq x^*$ for all n .

Given $\epsilon > 0$, it is not hard to see that there exists an associated number $N(\epsilon)$ such that there exists an $n \leq N(\epsilon)$ with the property that $|x_{n+1} - x_n| \leq \epsilon$, where $N(\epsilon)$ is defined by $N(\epsilon) = \lceil (a - b)/\epsilon \rceil + 1$.⁶ This result gives a crude and weak bound on how long it takes an algorithm with a descent property achieves a tolerance threshold. If for example, the MM algorithm is programmed to terminate as soon as $|PENSAE(f_{n+1}) - PENSAE(f_n)| \leq 10^{-6}$, then one can expect the algorithm to terminate within around 10^6 iterations in the worst case.

Such woefully slow performance is unusual, but not entirely unprecedented. An example of a simple Poisson estimation problem where the EM Algorithm exhibits $\mathcal{O}(1/n)$ convergence⁷ is given in [17]. Gradient descent methods can likewise exhibit the same rate of convergence in some cases [32]. The Implicit Filtering fitting algorithm previously discussed is another example of an incredibly slow method, requiring around 200 iterates before the sequence of values objective function converges to within 10^{-4} of each other⁸.

Linear Convergence

A more common rate of convergence for iterative methods in Statistics is known as *linear convergence*, which is defined as follows [6]: a vector-valued sequence \mathbf{x}_n is said to converge linearly to \mathbf{x}^* if there is a constant $0 < \mu < 1$ such that:

$$\lim_{n \rightarrow \infty} \frac{\|\mathbf{x}_{n+1} - \mathbf{x}^*\|}{\|\mathbf{x}_n - \mathbf{x}^*\|} = \mu$$

If a sequence \mathbf{x}_n converges linearly with constant μ , then $\|\mathbf{x}_{n+1} - \mathbf{x}^*\| \approx \mu \|\mathbf{x}_n - \mathbf{x}^*\|$ for n sufficiently large. This rate of convergence is slow, but not as catastrophically slow as $\mathcal{O}(1/n)$ convergence.⁹ For example, if $\mu = 1/2$, then it would require on the order of 21 iterations for a sequence to converge to within a tolerance of 10^{-6} .

Suppose there is a sequence $\{x_0, x_1, \dots, x_n, \dots\}$ converging to x^* defined by $x_n = F(x_n)$ and $F(\cdot)$ is differentiable around x^* , then for large n :

$$x_{n+1} - x^* \approx F'(x^*)(x_n - x^*)$$

If $0 < F'(x^*) < 1$, the sequence will converge linearly. In the case of a multivariate sequence \mathbf{x}_n converging to \mathbf{x}^* where $\mathbf{x}_{n+1} = \mathbf{F}(\mathbf{x}_n)$, the ordinary derivative $F'(x^*)$ is

⁶If $|x_{n+1} - x_n| > \epsilon, \forall n \leq N(\epsilon)$, it would be the case that $x_{N(\epsilon)} < b$, which contradicts the assumption that every element x_n of the sequence is bounded below by b .

⁷A possibly vector-valued sequence $\{\mathbf{x}_n\}_{n=1}^{\infty}$ is said to exhibit $\mathcal{O}(1/n)$ convergence or converges at rate $\mathcal{O}(1/n)$ if there exists constants $c_1 > c_2 > 0$ and a limit value \mathbf{x}^* such that $c_1/n \geq \|\mathbf{x}_n - \mathbf{x}^*\| \geq c_2/n$.

⁸The exact order of convergence of Implicit Filtering is difficult to determine. Experimentation suggests that the rate of convergence is sublinear as defined in⁹.

⁹If a sequence $\{\mathbf{x}_n\}$ converges at rate $\mathcal{O}(1/n)$ to x^* , then $\|\mathbf{x}_{n+1} - \mathbf{x}^*\|/\|\mathbf{x}_n - \mathbf{x}^*\| \rightarrow 1$ as $n \rightarrow \infty$. Any sequence with such a property is said to exhibit *sublinear convergence*. [6]

replaced with the Jacobian matrix derivative $J_{\mathbf{F}}(\mathbf{x}^*)$, so that:

$$\mathbf{x}_{n+1} - \mathbf{x}^* \approx J_{\mathbf{F}}(\mathbf{x}^*)(\mathbf{x}_{n+1} - \mathbf{x}^*)$$

The question of the rate of convergence is more subtle than in the univariate case. So long as there is a matrix norm such that $\|J_{\mathbf{F}}(\mathbf{x}^*)\| < 1$ the sequence will converge linearly at worst, though faster than linear convergence is possible if 0 is an eigenvalue of $J_{\mathbf{F}}(\mathbf{x}^*)$.

As a heuristic rule, any iterative statistical method that relies on solving the same problem over and over again will converge linearly unless it has been specifically designed to ensure that $J_{\mathbf{F}}(\mathbf{x}^*) = 0$. [12, 16, 21, 22]. The EM Algorithm, MM Algorithm, and Block Relaxation methods generally only achieve linear convergence.

The MM algorithm for L_1 fitting does not quite converge linearly though as shall be seen in Section 4.1.4 below.

4.1.4 Testing the Algorithm on the Melanoma Data

Since minimising *PENSAE* is a problem over many dimensions, plotting the objective function to verify that the optimal function has been found isn't possible. Instead the MM algorithm described in Section 4.1.2 will be tested by applying it to the melanoma data perturbed by random noise. Further, the convergence of the algorithm for the original melanoma dataset will be examined.

Figure 4.1 presents the L_1 and L_2 inner fits to the melanoma data corrupted by Cauchy distributed noise. The value of ω is held fixed at the reasonable value of 0.3, which was chosen as being roughly the average of the two different estimates of ω from computed in the previous chapter. It is apparent from the Figure 4.1 that the MM fit is robust against outliers, tends to ignore more deviant points, and even manages to remain similar to the original fit. The least-squares fit tends to chase the heavy-tailed noise on the other hand. This is strong evidence that the curve that minimises *PENSAE* has been found and that the method has been implemented correctly.

Figures 4.2 and 4.4 plot the convergence of *SAE* and *PENSAE* over the course of the algorithm. Note that the *PENSAE* statistic doesn't quite actually converge monotonically as the theoretical analysis predicted. Instead, it fluctuates before settling down to the typical and expected pattern of languid monotone decline. Upon investigation, it was determined that over the first handful of iterations the range of the weights applied to the observations on each iteration, that were computed using the residuals from the previous iteration, grew very rapidly. By the fourth iteration, the lowest weight is equal to 1.48, and the highest was equal to 4.8×10^6 . It seems that this rapid and large change produces qualitative changes in behaviour before the algorithm manages to 'burn in'. It is likely that observations with low weights are being effectively censored after a few iterations due to roundoff error. It was found that imposing a minimum threshold for the weights by adding a constant to all the residuals before proceeding to computing the weights smooths out this behaviour, but doesn't eliminate it entirely.

Figure 4.3 plots the convergence of the coefficient vectors \mathbf{c}_n . This log-plot suggests that the sequence of fitted coefficient vectors \mathbf{c}_n converges linearly since $\|\mathbf{c}_{n+1} - \mathbf{c}_n\| \approx C\|\mathbf{c}_n - \mathbf{c}_{n-1}\|$ as $n \rightarrow \infty$.

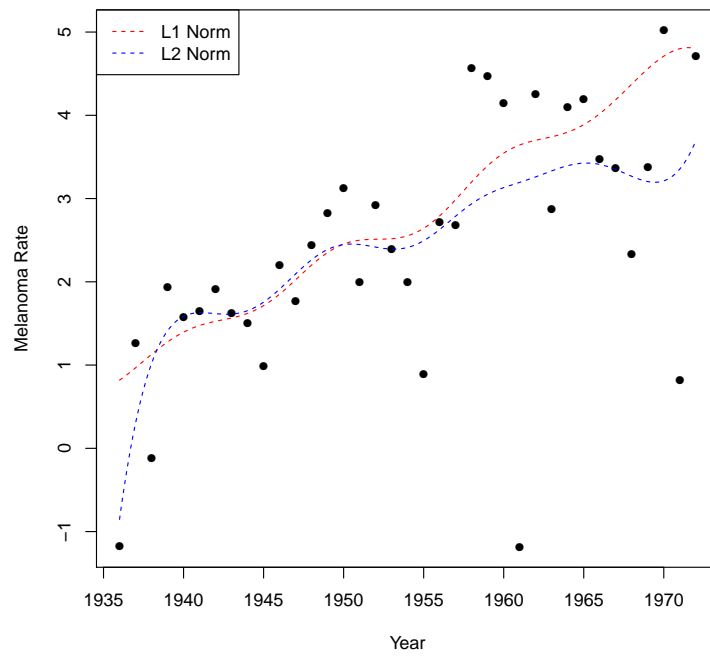


Figure 4.1: Comparison of L_1 and L_2 inner fits to Cauchy perturbed data with ω fixed at 0.3

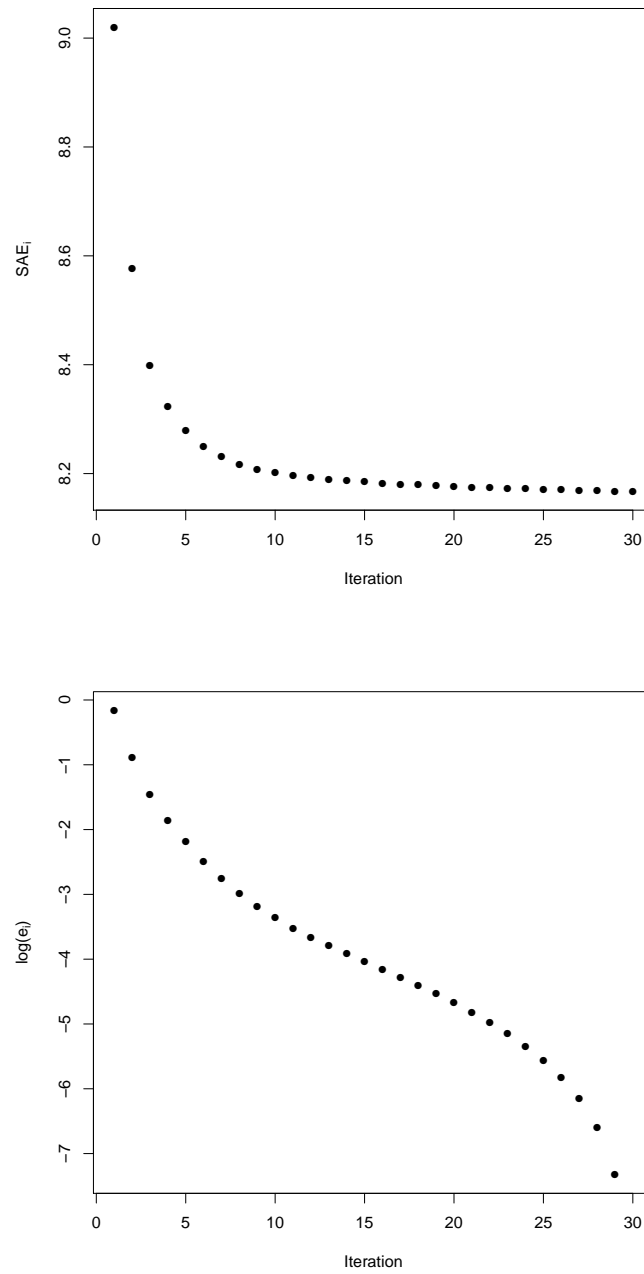


Figure 4.2: Plot of values and log differences for SAE Statistic

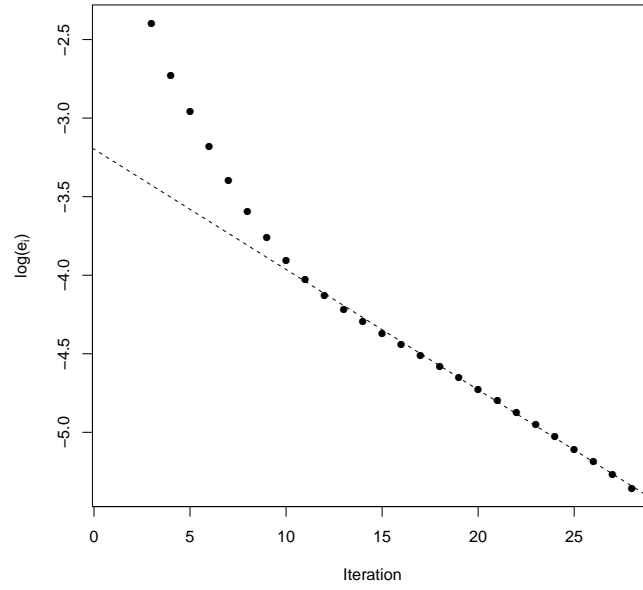


Figure 4.3: Plot of log norm differences for coefficients. Note that they tend to settle on a line.

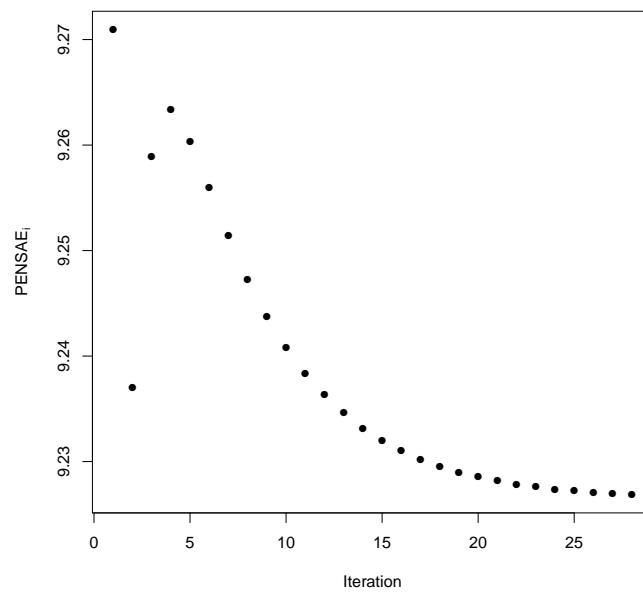


Figure 4.4: Plot of PENSAT statistic as the algorithm proceeds.

4.2 The Two Level Parameter Cascade with L_1 Norm

The inner problem of the parameter cascade is a semiparametric least squares regression model. The fitted function is modeled as a weighted sum of a solution to the differential equation (parametric), and a non-parametric residual. The lambda term governs how big the residual is allowed to be relative to the least squares error term

If the usual least-squares error function is used, the inner problem will probably struggle with outliers and heavy tailed errors as is the case for any form of least-squares regression.

For high order differential operators like that used to model the melanoma data, there are many degrees of freedom associated with the differential operator's solution set. The omega and lambda parameters don't strongly constrain the lower level of the cascade. There is thus little capacity for the higher levels of the cascade to restrain the lowest level through altering the lambda and omega parameters and the parameter cascade must use robust estimation at every level.

In the previous chapter, it was discussed how Brent's Method can be used to tackle the middle problem without derivatives and then used this approach to optimise a highly irregular loss function. In the previous section, the MM algorithm was used to optimise the inner problem with an L_1 norm.

Combining the two methods, it is very straightforward to implement a two-level parameter cascade with L_1 errors at both levels.

In Figure 4.5, the result of fitting a two level L_1 Parameter Cascade with L_1 errors is plotted. It can be seen that the $SAE(\omega)$ function is irregularly shaped. In Figure 4.6 both the L_1 and L_2 fits to Cauchy-perturbed melanoma data are shown. Figure 4.7 plots the results of applying the L_1 and L_2 Parameter Cascades to the original and perturbed Melanoma data, alongside mixed versions where the L_1 loss function is used for the inner fitting and L_2 loss function for the middle fitting and vice versa.

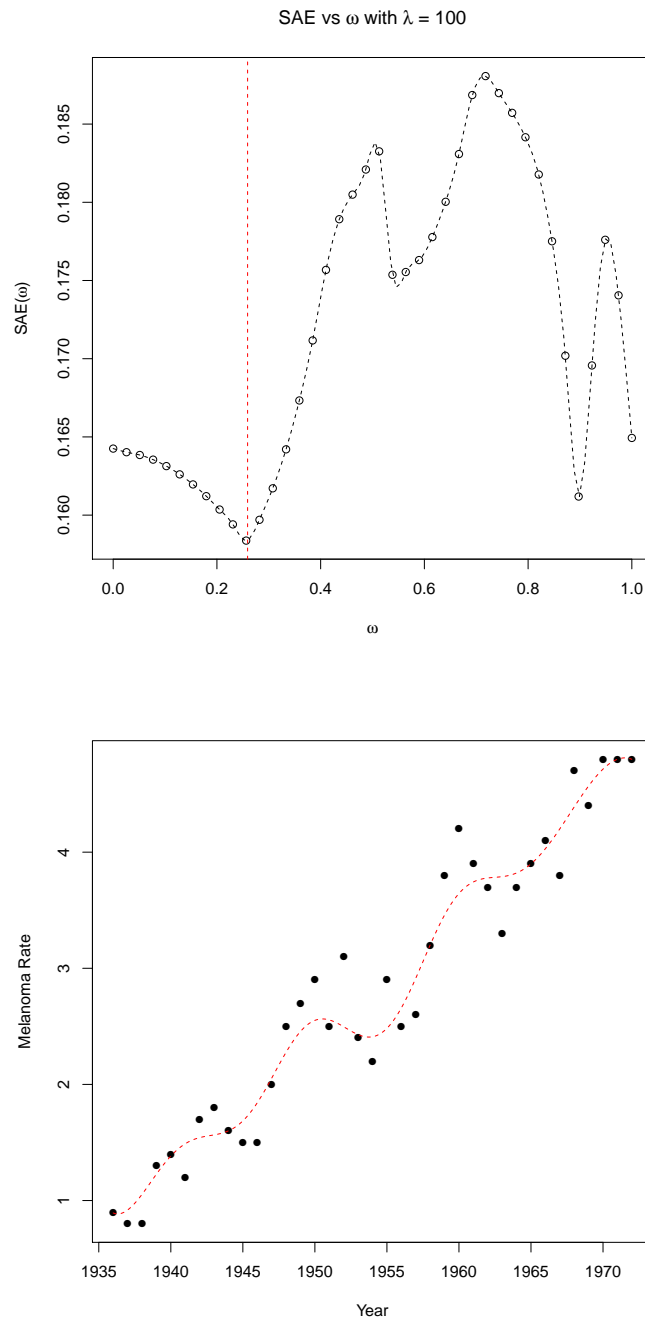


Figure 4.5: Fitting an L_1 Parameter Cascade to the Melanoma Data

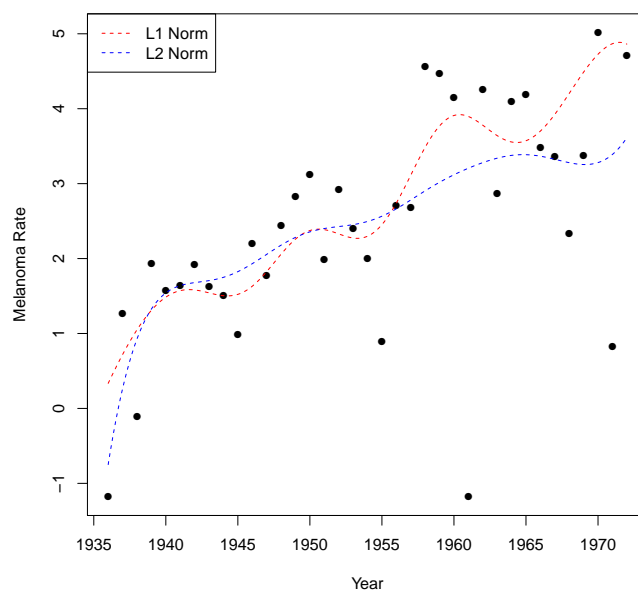


Figure 4.6: L_1 and L_2 Parameter Cascades with the same perturbed data as in Figure 4.1. Compare the L_1 curve in this plot with the one in Figure 4.5.

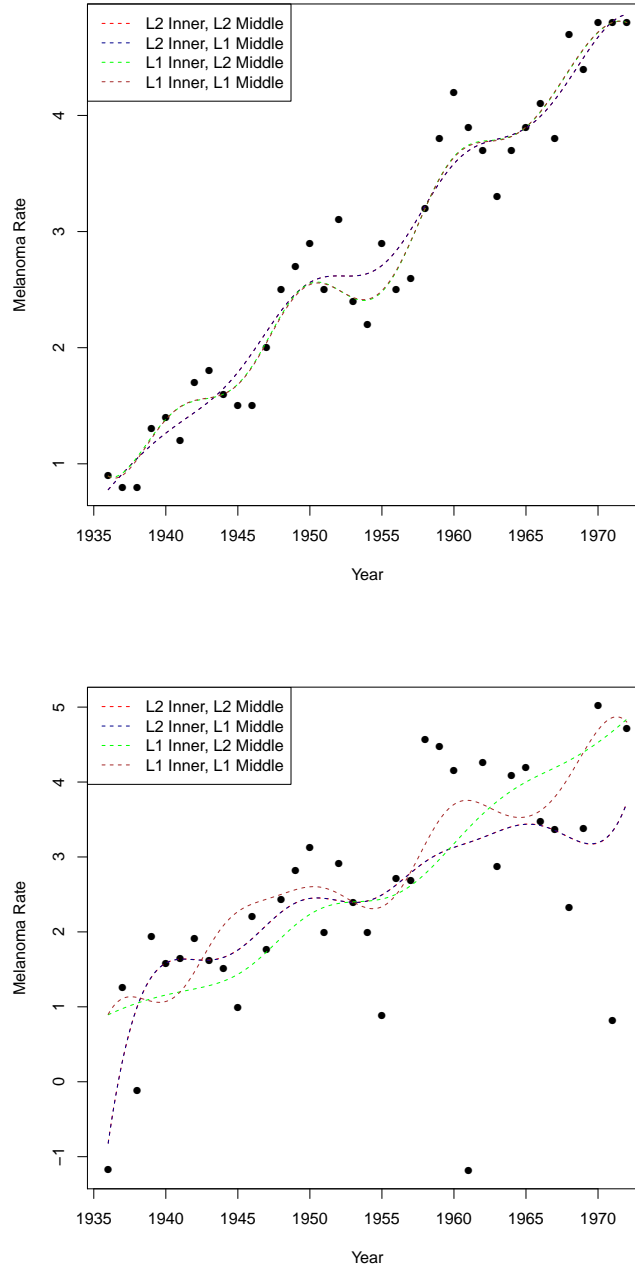


Figure 4.7: All possible combinations of L_1 and L_2 loss functions that can be used for the Parameter Cascade. The top plot applies them to the original melanoma data, the bottom to the same perturbed data as in Figures 4.1 and 4.6

4.3 Accelerating The Rate of Convergence

The MM Algorithm is very sluggish, and this is a well known weakness of both itself and the EM algorithm. The literature however suggest that this problem could be easily ameliorated in this particular case because of a special feature present. In practice one doesn't want to fit the full model, but only wants to compute an associated summary statistic that determines how good a given choice of parameters is. It will often be the case that only the value of *PENSSAE* or *GCV* or *SAE* associated with a given choice of parameters is required as inputs to an optimisation routine, and it is not desirable to iterate until the full model converges if this effort can be avoided.

MacLanan and Krishnan discuss the situation where one only wants to compare the likelihoods between a restricted model and a full model. They suggest the use of sequence acceleration methods to rapidly extract the likelihoods [21] instead of running the EM algorithm to completion since the full models aren't needed. The literature on the MM algorithm claims that acceleration methods for the EM algorithm translate quite easily to the MM case [33]. On this basis, we explored whether this approach might be applied here.

The textbook approach employed is known as Aitken Acceleration[7, 21]. Suppose that there is a sequence x_0, x_1, x_2, \dots converging to a limit x^* . Aitken's method makes the ansatz that $x_{n+1} - x^* \approx C(x_n - x^*)$ for some constant C . Many iterative algorithms in statistics exhibit this pattern as discussed in Section 4.1.3. This suggests the following equation:

$$\frac{x_{n+1} - x^*}{x_n - x^*} \approx \frac{x_n - x^*}{x_{n-1} - x^*}$$

Solving for x^* gives the accelerated sequence.

There is an equivalent definition that is easier to generalise [12]. Consider a sequence defined by functional iteration so that $x_{n+1} = F(x_n)$ for some function $F(\cdot)$. Define the error sequence by $e_n = x_{n+1} - x_n = F(x_n) - x_n$. The function $g(x) = F(x) - x$ returns the error associated with any value, and the limit of the sequence satisfies $g(x^*) = 0$. Suppose one knew the inverse of $g(x)$, which will be denoted by $h(e)$. Then x^* could be found by evaluating $f(0)$. The next best thing would be to use the values of the sequence to approximate $h(e)$, and then evaluate this approximate function at zero instead. The Aitken method approximates $h(e)$ by linear interpolation between (e_n, x_n) and (e_{n-1}, x_{n-1}) , and then evaluates this approximation at $e = 0$.

4.3.1 Illustrative Example: Using Imputation to Fit an ANOVA Model With Missing Data

For illustrative purposes, we will make use of an example from chapter 2 of [21]. The authors discuss fitting an ANOVA model to a factorial experiment where some of the values are missing. They proceed by using the fitted model to estimate the missing values; fitting the model again with the new imputed values; and using the new fitted values in turn to again update the estimates of missing values. The process is repeated until convergence. In the text, the authors do not work here with likelihood or any probabilistic models and treat the question as purely a regression problem. This is similar to our L_1 fitting problem.

The authors' example was implemented again in R.¹⁰ For each iteration, the SSE statistic was computed. This defines an associated sequence $\{SSE_1, SSE_2, \dots, SSE_n, \dots\}$. Applying Aitken's method to this sequence produces a new sequence $\{ASSE_n\}$. As can be seen in Figure 4.8 and Table 4.1, the accelerated sequence converges far more quickly to the limit of the $\{SSE_i\}$ sequence than the original sequence.

4.3.2 Generalisations

Outside of more specialised texts, Statistics is generally content with Aitken's method and multivariate generalisations. Exploring more powerful methods can be justified in two circumstances. The first is that if one is running the algorithm over and over again such that an increase in speed over many iterations means the effort invested is worth it. This might be the case for example if one wanted to use the bootstrap to model the distribution of an likelihood ratio statistic computed using the EM Algorithm as previously described. The second is if the sequence is difficult to accelerate. In this situation, it shall be seen that both conditions apply.

As a field of study, sequence acceleration is closely related to time series analysis. A generic first order autoregressive model is given by:

$$x_{n+1} = f(x_n, n) + \epsilon_n$$

Consider the case where there are both no random errors so that ϵ_n is always zero, and the sequence converges to a limit. Here, the problem of determining the long term value of the sequence from a set of observations is equivalent to that of accelerating the sequence. If the specific form of $f(x_n, n)$ is known, there can often be a specific acceleration method that can exactly extract the limit. For illustration, suppose there were a sequence of the following form, but the parameters β_0 and β_1 were unknown:

$$x_n = \beta_0 + \frac{\beta_1}{n} \tag{4.3}$$

¹⁰The SSE statistic here is different from the RSS statistic presented in the text. The new code converges to the same estimates as in the original, so the example has been re-implemented correctly. It was not possible to determine with what degrees of freedom RSS was associated with.

As n goes to infinity, x_n converges to β_0 . It is not difficult to show that the limit β_0 can be found by applying the following sequence transformation:

$$\begin{cases} \hat{\beta}_{1,n} = \frac{x_n - x_{n+1}}{\left(\frac{1}{n} - \frac{1}{n+1}\right)} \\ \tilde{x}_n = x_n - \frac{\hat{\beta}_{1,n}}{n} \end{cases} \quad (4.4)$$

If the transformation (4.4) is applied to a sequence of values $x_1, x_2, \dots, x_n, \dots$ that is of form (4.3), then the transformed sequence $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n, \dots$ will have the property that $\tilde{x}_n = \beta_0$ for all n . Likewise, the Aitken method is exact for sequences of the form $x_{n+1} = \beta_0 + \beta_1 x_n$, and so can be thought of as the deterministic analogue of an $AR(1)$ model.

The process of acceleration isn't quite so neat in practice because sequences don't adhere perfectly to these simple forms. Instead, the best that can be realistically hoped for is that the transformed sequence converges to the same limit as the original, but the rate of convergence is higher. For example, if transformation (4.4) is applied to a sequence of the form $y_n = \beta_0 + \frac{\beta_1}{n} + \frac{\beta_2}{n^2}$, then the transformed sequence is now of the form $\tilde{y}_n = \beta_0 + \mathcal{O}(\frac{1}{n^2})$, which converges to β_0 more quickly than the original sequence.¹¹

Suppose a convergent sequence is of the form $x_{n+1} = f(x_n)$ with $f(\cdot)$ differentiable and x^* is the limit. Using a first order Taylor expansion, it can be seen that for sufficiently large n , $x_{n+1} \approx x^* + f'(x^*)(x_n - x^*)$. In this case, Aitken acceleration has a decent chance of accelerating the sequence so long as it has 'burned in' sufficiently.

One generalisation, proposed in [12] is to use higher order polynomials to model the inverse error function $h(e)$. So $h(e)$ would be approximated by a quadratic through (e_n, x_n) , (e_{n-1}, x_{n-1}) and (e_{n-2}, x_{n-2}) . Making e the independent variable here instead of x means the estimated limit can simply be found by evaluating the approximating quadratic at $e = 0$ instead of having to find the correct root of a quadratic to compute each element of the accelerated sequence.

Another approach is to simply apply Aitken Acceleration to the sequence twice.

Both these approaches were attempted for the missing data model, and the results can be seen in Table 4.1 and Figure 4.9. It can be seen that both methods improve convergence, though double Aitken acceleration is more effective (and easier to implement).

One can take the process further. For the missing values linear model, these higher-order methods converge very rapidly and are prone numerically instability thereafter due to the error terms being so small, so plotting or tabulating them was not worth the additional clutter. If the Aitken method is applied three times to the original sequence, the first entry yields the limit immediately and there is no need to go any further. Applying the quadratic method twice in a row produces a new sequence for which the first entry is within 10^{-12} of the limit.

¹¹Doing the algebra, it can be seen that it is now the case that $\hat{\beta}_{1,n} = \beta_1 + \beta_2 \left[\frac{2n+1}{n(n+1)} \right]$, and so $\tilde{y}_n = \beta_0 + \beta_2 \left[-\frac{2n+1}{n^2(n+1)} + \frac{1}{n^2} \right] = \beta_0 + \beta_2 \left[\frac{-n^3 - n^2 + n}{n^4(n+1)} \right]$.

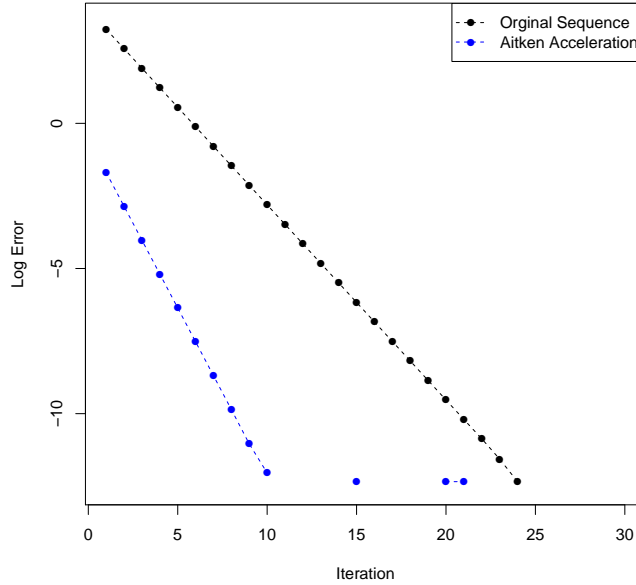


Figure 4.8: Log Errors for original sequence of SSE values and the accelerated one

Other Approaches: There are alternative approaches besides those described here. For example, the EM and MM algorithms generate a sequence of coefficient vectors $\{\mathbf{c}_0, \mathbf{c}_1, \dots, \mathbf{c}_n, \dots\}$ with $\mathbf{c}_{n+1} = \mathbf{F}(\mathbf{c}_n)$ for some function $\mathbf{F}(\cdot)$. In our particular situation, the function $\mathbf{F}(\cdot)$ would denote the operator that takes a coefficient vector and returns the coefficient vector that minimises the associated *WPENSSE* problem (4.1). The limit of this sequence - should it exist - is a solution to the equation $\mathbf{c} = \mathbf{F}(\mathbf{c})$. It is proposed in the literature to use Newton or Quasi-Newton methods such as those described in the chapter on Brent's method to numerically solve this fixed point equation [5, 7]. The idea is that such methods will find the fixed point more rapidly than simply iterating $\mathbf{F}(\cdot)$ until one gets sufficiently close to the limit. These methods have the disadvantage of being more complex and time consuming to implement than the univariate acceleration methods.

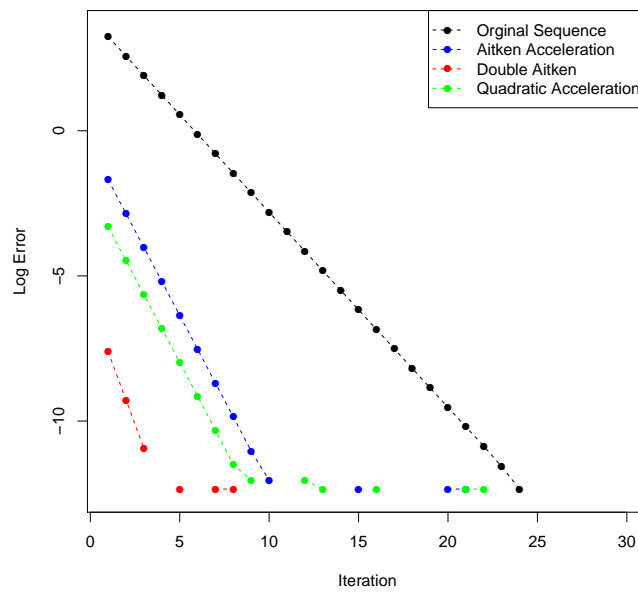


Figure 4.9: More sophisticated acceleration methods can provide a further boost to convergence. There are gaps in the plot because the more accelerated iterations have no valid log error since R cannot numerically distinguish them from the final limit.

Table 4.1: Iterations of the original sequence SSE_n , the accelerated sequence $ASSE_n$, the quadratically accelerated sequence $QASSE_n$, and the doubly accelerated sequence $DASSE_n$.

	SSE_n	$ASSE_n$	$QASSE_n$	$DASSE_n$
1	4949.6944444444	3203.8032711619	3203.7834325738	3203.7829457122
2	3575.1658950617	3203.7843303225	3203.7829788622	3203.7829457359
3	3282.7945625667	3203.7830400346	3203.7829479917	3203.7829457364
4	3220.5935028609	3203.7829521582	3203.7829458900	3203.7829457364
5	3207.3596279977	3203.7829461738	3203.7829457469	3203.7829457364
6	3204.5439391293	3203.7829457662	3203.7829457371	3203.7829457364
7	3203.9448588925	3203.7829457385	3203.7829457365	3203.7829457364
8	3203.8173952920	3203.7829457366	3203.7829457364	3203.7829457364
9	3203.7902754193	3203.7829457364	3203.7829457364	3203.7829457364
10	3203.7845052414	3203.7829457364	3203.7829457364	3203.7829457364
11	3203.7832775456	3203.7829457364	3203.7829457364	3203.7829457364
12	3203.7830163340	3203.7829457364	3203.7829457364	3203.7829457364
13	3203.7829607572	3203.7829457364	3203.7829457364	3203.7829457364
14	3203.7829489323	3203.7829457364	3203.7829457364	3203.7829457364
15	3203.7829464164	3203.7829457364	3203.7829457364	3203.7829457364
16	3203.7829458811	3203.7829457364	3203.7829457364	3203.7829457364
17	3203.7829457672	3203.7829457364	3203.7829457364	3203.7829457364
18	3203.7829457430	3203.7829457364	3203.7829457364	3203.7829457364
19	3203.7829457378	3203.7829457364	3203.7829457364	3203.7829457364
20	3203.7829457367	3203.7829457364	3203.7829457364	3203.7829457364
21	3203.7829457365	3203.7829457364	3203.7829457364	3203.7829457364
22	3203.7829457364	3203.7829457364	3203.7829457364	3203.7829457364
∞	3203.7829457364	3203.7829457364	3203.7829457364	3203.7829457364

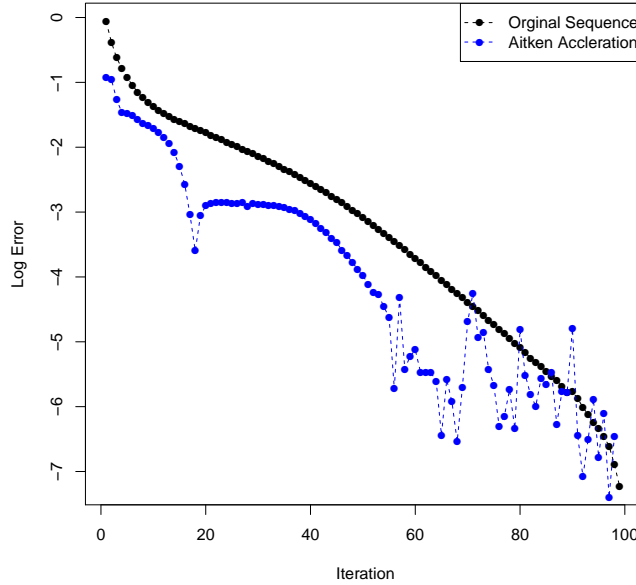


Figure 4.10: Accelerating the SAE sequence generating by the L_1 fitting algorithm using Aitken’s Method. The improvement in convergence is mediocre.

4.3.3 Accelerating the L_1 Fitting Algorithm

The L_1 fitting algorithm is much more difficult to accelerate as can be seen in Figures 4.10 and 4.11. Even advanced acceleration methods recommended for slowly converging sequences that the Aitken method cannot tackle such as the Epsilon Algorithm[10, 23] or Lubkin’s W transform [23, 31] yield much improvement¹². The *SAE* sequence is apparently either numerically ill-behaved or of a very unusual form.

To conclude, it is possible to save some time by acceleration, but the scope for doing so is limited and the process would have to be monitored carefully to ensure that numerical instability isn’t causing trouble.

¹²Several other methods proposed in [23] were such as Levin transforms were also attempted. They proved similarly unsatisfactory

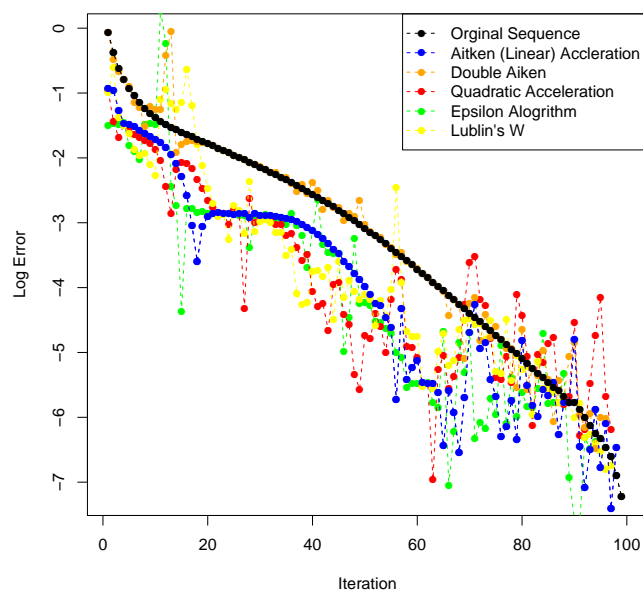


Figure 4.11: Accelerating the SAE sequence using multiple methods. Aitken's method performs the best, despite it's lack of sophistication.

Chapter 5

Fitting Non Linear Regression Models With the Parameter Cascade

Consider the following non-linear regression model where observed values y_i are values observed at times t_i :

$$y_i = \alpha + \beta e^{\gamma t_i} + \epsilon_i \quad (5.1)$$

All of the difficulty here comes from the $e^{\gamma t}$ term. If γ were known, α and β could be found through simple linear regression with the $e^{\gamma t_i}$ term acting as an independent variable predicting the y_i

This suggests the following regression strategy. Define a function $H(\gamma)$ to be the sum of squared errors from performing simple linear regression on the y_i against $e^{\gamma t_i}$. That is:

$$H(\gamma) = \min_{\alpha, \beta} \sum [y_i - \alpha - \beta e^{\gamma t_i}]^2$$

This defines a middle problem, with the inner problem being that of minimising the simple linear regression problem given γ . The non-linear model can be fitted by using Brent's Method to fit the middle problem.

This approach was applied to simulated data with $\alpha = 100$, $\beta = 4$, and $\gamma = 1$, and the results can be seen in Figure 5.1.

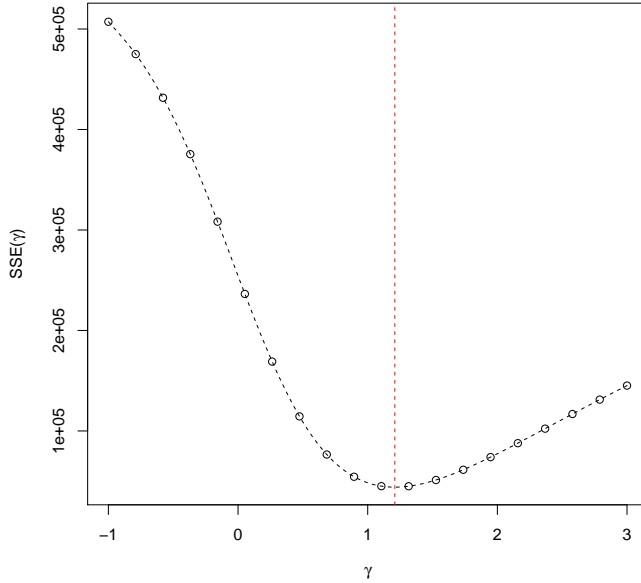


Figure 5.1: Profile Plot and Fitted Curve

5.1 Fitting Linear Homogenous ODEs Using the Parameter Cascade

Recall a linear homogenous ODE of order n is given by:

$$\frac{d^n y}{dt^n} = \sum_{k=0}^{n-1} a_k(t; \theta) \frac{d^k y}{dt^k}$$

Under some mild technical conditions, the set of solutions to such an ODE is an n dimensional vector space and has a unique solution for each set of initial conditions. It is often more convenient to work with ODEs in matrix form from now on. Any homogenous linear ODE can be represented in matrix form:

$$\frac{d\mathbf{y}}{dt} = \mathbf{A}(t; \theta)\mathbf{y}$$

For example, the ODE $y'' = -\omega^2 y$ with the initial conditions $y(0) = y_0$ and $y'(0) = v_0$ can be represented as:

$$\frac{d}{dt} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -\omega^2 & 0 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

with the initial condition $\mathbf{y}(0) = (y_0, v_0)'$.

A basis for the solution set of any linear ODE can be formed from the set of solutions associated with the initial conditions $\mathbf{y}(0) = \mathbf{e}_i$, where \mathbf{e}_i denotes the i th basis vector.

So this suggests the following cascade algorithm: given a set of parameters, find the set of solutions $\{y_1(t|\theta), \dots, y_n(t|\theta)\}$, where $y_k(t|\theta)$ denotes the solution with the k th basis vector as an initial condition. Then perform regression to fit the $\{y_i(t|\theta)\}$ to the observed data. The inner problem consists of fitting a weighted sum of the $\{y_i(t|\theta)\}$ to the observed data and reporting the associated error given a choice of parameters. The middle problem consists of finding the set of parameters that minimises this associated error.

For a problem where the ODE can be solved explicitly, things proceed as in the previous section. Consider again the ODE $y'' - \omega^2 y = 0$. The solutions generated by the initial conditions $(1, 0)$ and $(0, 1)$ is given by $\{A \cos \omega t + B \sin \omega t | A, B \in \mathbb{R}\}$. So the middle least squares criterion is given by:

$$H(\omega) = \min_{a,b} \sum [y_i - a \cos \omega t_i - b \sin \omega t_i]^2$$

Finding the optimal a and b given ω is an inner problem that can be solved using least squares regression as before. In fact, such a problem has already been encountered: the nonlinear model given in Equation 5.1 is associated with the ODE $y'' - \gamma y' = 0$.

For ODE problems that cannot be explicitly solved, the trajectories $y_n(t|\theta)$ must be instead found by a numerical solver for each choice of θ . The inner problem then consists of linear regressing the computed solutions against the observed data.

To illustrate the method, it was applied to the following ODE with $\alpha = -0.3$ and $\beta = -1.0$. To minimise the middle problem, the Nelder-Mead method was used - Brent's method was felt to be unsuitable because of the awkward topography. The results can be seen in Figure 5.2.

$$y''(t) = \alpha \sqrt{t} y(t) + \beta \sin(2t) y'(t) \quad (5.2)$$

The advantage of the parameter cascade here is that it is noticeably faster than trying to optimise everything in one go. The linear regression steps mean that the ODE needs to be solved numerically fewer times, so that the algorithm runs around 30% faster for the ODE in Equation 5.2. However, the `nls` command is faster than the Parameter Cascade for Equation 5.1 even when no derivatives are provided.

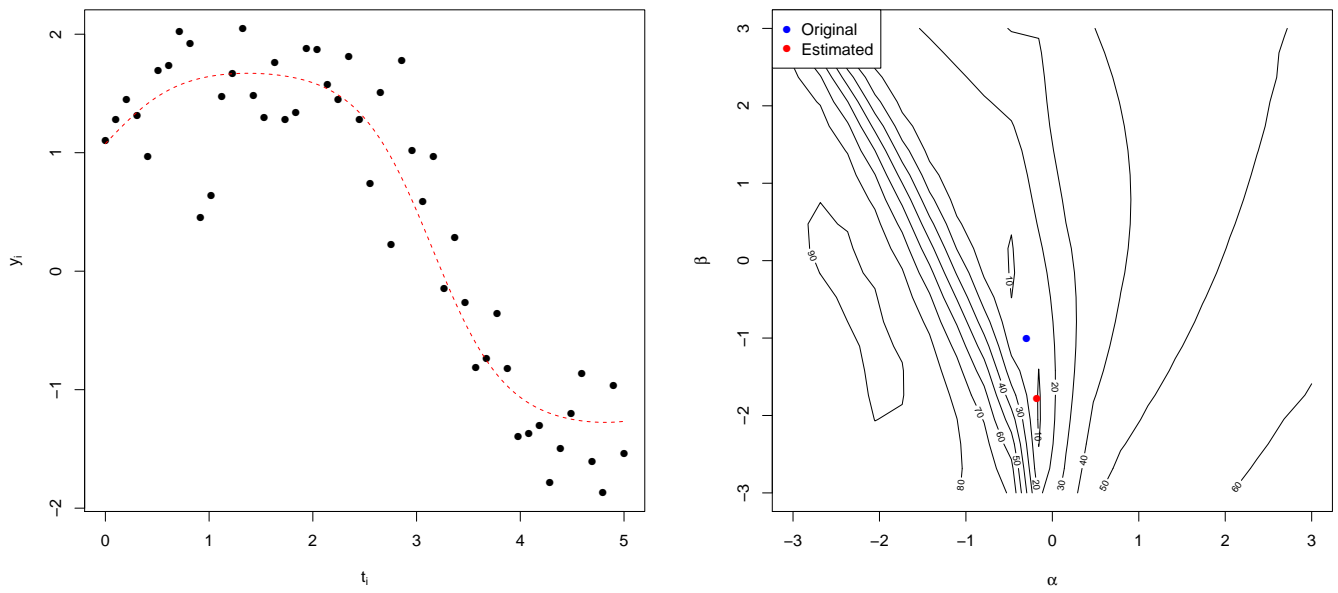


Figure 5.2: Plot of fit to simulated data, and contour plot of SSE against α and β .
Blue dot is true values, red is estimated values.

Bibliography

- [1] John P Boyd. *Chebyshev and Fourier spectral methods*. Courier Corporation, 2001.
- [2] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [3] Richard P Brent. *Algorithms for minimization without derivatives*. Courier Corporation, 2013.
- [4] Jiguo Cao and James O Ramsay. Parameter cascades and profiling in functional data analysis. *Computational Statistics*, 22(3):335–351, 2007.
- [5] Kwun Chuen Gary Chan. Acceleration of expectation-maximization algorithm for length-biased right-censored data. *Lifetime data analysis*, 23(1):102–112, 2017.
- [6] Edwin KP Chong and Stanislaw H Zak. *An introduction to optimization*, volume 76. John Wiley & Sons, 2013.
- [7] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- [8] Ronald Aylmer Fisher. The wave of advance of advantageous genes. *Annals of eugenics*, 7(4):355–369, 1937.
- [9] Bengt Fornberg. Generation of finite difference formulas on arbitrarily spaced grids. *Mathematics of computation*, 51(184):699–706, 1988.
- [10] PR Graves-Morris, DE Roberts, and A Salam. The epsilon algorithm and related topics. *Journal of Computational and Applied Mathematics*, 122(1-2):51–80, 2000.
- [11] David R Hunter and Kenneth Lange. A tutorial on MM algorithms. *The American Statistician*, 58(1):30–37, 2004.
- [12] Eugene Isaacson and Herbert Bishop Keller. *Analysis of numerical methods*. Courier Corporation, 2012.
- [13] Carl T Kelley. *Implicit filtering*, volume 23. SIAM, 2011.

- [14] C.T. Kelley. A brief introduction to implicit filtering. <https://projects.ncsu.edu/crsc/reports/ftp/pdf/crsc-tr02-28.pdf>, 2002. [Online; accessed 12-October-2019].
- [15] Jack Kiefer. Sequential minimax search for a maximum. *Proceedings of the American mathematical society*, 4(3):502–506, 1953.
- [16] Kenneth Lange. *Optimization*. Springer, 2004.
- [17] Kenneth Lange. *Numerical analysis for statisticians*. Springer Science & Business Media, 2010.
- [18] Kenneth Lange. The MM algorithm. <https://www.stat.berkeley.edu/~aldous/Colloq/lange-talk.pdf>, April 2007. [Online, accessed 18-September-2019].
- [19] Randall J LeVeque. *Finite difference methods for ordinary and partial differential equations: steady-state and time-dependent problems*, volume 98. Siam, 2007.
- [20] Steve McConnell. *Code complete*. Pearson Education, 2004.
- [21] Geoffrey McLachlan and Thiriyambakam Krishnan. *The EM algorithm and extensions*, volume 382. John Wiley & Sons, 2007.
- [22] J Nocedal and SJ Wright. *Numerical Optimisation*. Springer verlag, 1999.
- [23] Naoki Osada. *Acceleration methods for slowly convergent sequences and their applications*. PhD thesis, PhD thesis, Nagoya University, 1993.
- [24] Yudi Pawitan. *In all likelihood: statistical modelling and inference using likelihood*. Oxford University Press, 2001.
- [25] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2013. ISBN 3-900051-07-0.
- [26] James Ramsay. Functional data analysis. *Encyclopedia of Statistics in Behavioral Science*, 2005.
- [27] Walter Rudin et al. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1964.
- [28] Christopher G Small. A survey of multidimensional medians. *International Statistical Review/Revue Internationale de Statistique*, pages 263–277, 1990.
- [29] Karline Soetaert, Thomas Petzoldt, and R. Woodrow Setzer. Solving differential equations in R: Package deSolve. *Journal of Statistical Software*, 33(9):1–25, 2010.

- [30] Keller Vandebogart. Method of quadratic interpolation. http://people.math.sc.edu/kellerlv/Quadratic_Interpolation.pdf, September 2017. [Online; accessed 13-September-2019].
- [31] Jet Wimp. *Sequence transformations and their applications*. Elsevier, 1981.
- [32] Stephen Wright. Optimization for data analysis. In Michael W. Mahoney, John C. Duchi, and John C. Duchi, editors, *The Mathematics of Data*, chapter 2, pages 49–98. American Mathematical Society and IAS/Park City Mathematics Institute and Society for Industrial and Applied Mathematics, 2018.
- [33] Tong Tong Wu, Kenneth Lange, et al. The MM alternative to EM. *Statistical Science*, 25(4):492–505, 2010.