# Paramter Estimation In IGW

By Huasheng Liao

**The Goal:**

The goal is to find the parameters, $\vec{a}$ ($\vec{a}$ is a vector) in a given model or multiple models such that the given models' values f($\vec{x}$, $\vec{a}$) fit best some data points

**To achieve the goal:**

Given a model with unknown parameters, $\vec{a} =(a_1,a_2)$ to be estimsted:

Model: $\gamma(a_1,a_2) = \dfrac{a_1 x}{a_2 + x} = f(\vec{x}, \vec{a})$ (or complex numerical models)

And observations: $\gamma_i$ at $x_i$, i=1,2,…,m

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| x | 0.038 | 0.194 | 0.425 | 0.626 | 1.253 | 2.500 | 3.740 |
| $\gamma$ | 0.050 | 0.127 | 0.094 | 0.2122 | 0.2729 | 0.2665 | 0.3317 |

One can have m residual functions:

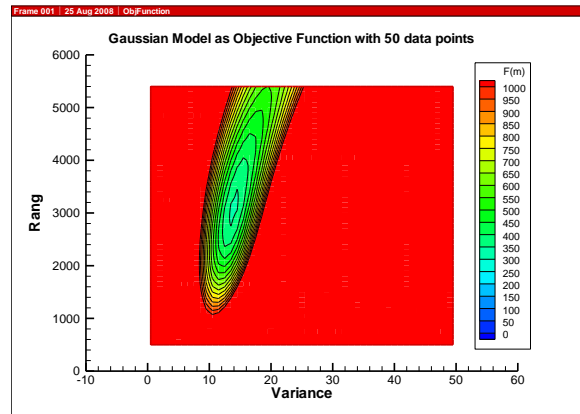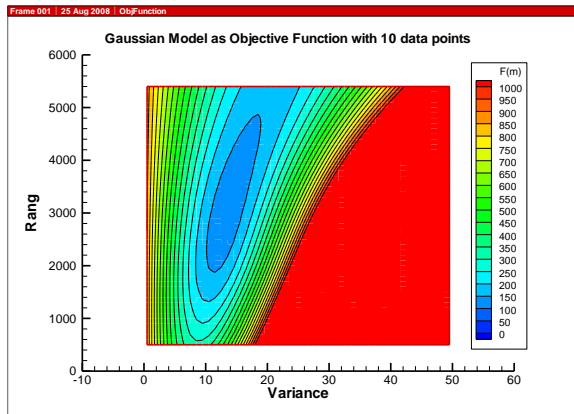$$r_i(a_1,a_2) = \gamma_i - \gamma(a_1,a_2,x_i), \quad i=1,2,…,m$$

Note that, $r_i(a_1,a_2)$ is function of unknown parameters, $\vec{a} =(a_1,a_2)$ to be estimsted.

Best fitting means that those residuals at data points tend to minimum or zero: **Least Square Minimization (LSM) or Weighted Least Squares Minimization (WLSM)** is one of the good method to be used in parameters estimation.
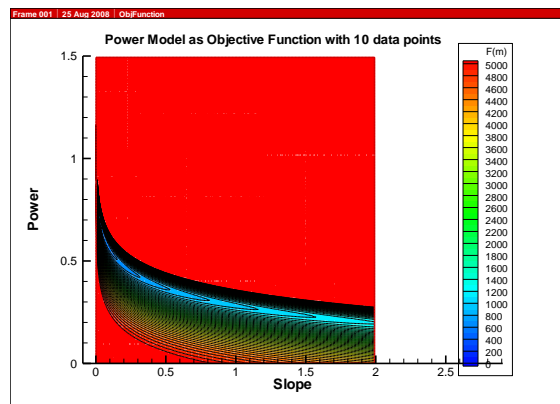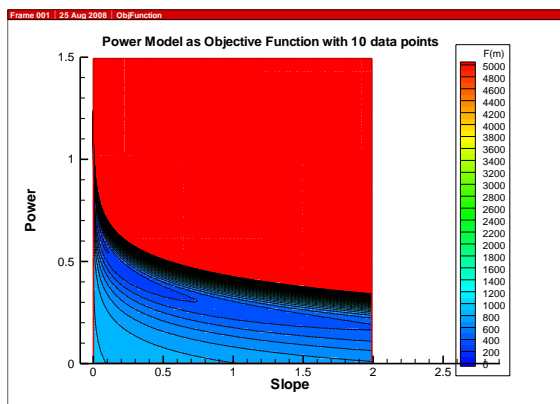
LSM : the unknown values of the parameters in the models are estimated by finding the numerical values for the parameter estimates that minimize **the sum of the squared deviations** between the observed responses and the functional portion of the model.

Sum of squares residuals: $S = \displaystyle\sum_{i=1}^{m} r_i^2(a_1,a_2) = S(a_1,a_2)$
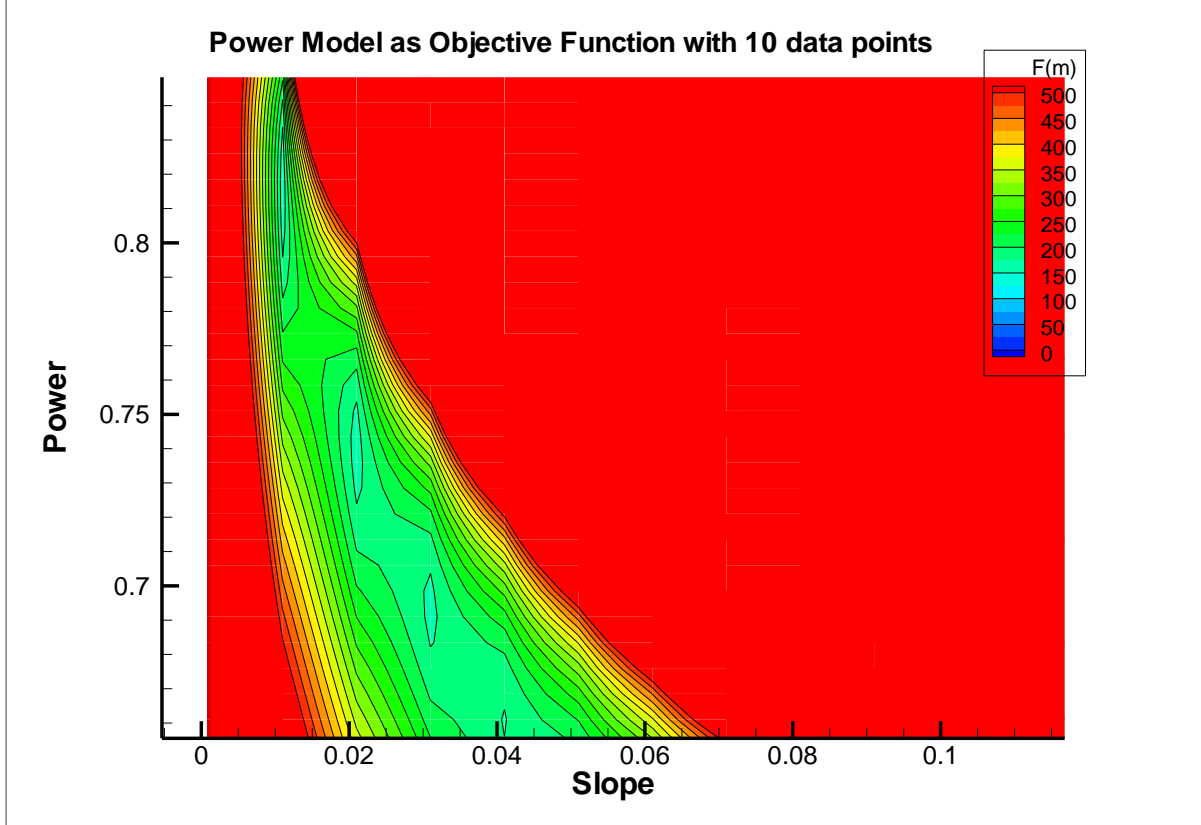
What does S look like?

Gaussian model



Power Model

**Power Model as Objective Function with 10 data points**

Local traps

# Multiple minima

Multiple minima can occur in a variety of circumstances some of which are:

- A parameter is raised to a power of two or more. For example, when fitting data to a Lorentzian curve

$$f(x_i, a) = \frac{\alpha}{1 + (\frac{\gamma - xi}{a})^2}$$

where α is the height, γ is the position and a is the half-width at half height, there are two solutions for the half-width, and which give the same optimal value for the objective function.

- Two parameters can be interchanged without changing the value of the model. A simple example is when the model contains the product of two parameters, since αβ will give the same value as βα.
- A parameter is in a trigonometric function, such as sina , which has identical values at a'+2nπ. See Levenberg-Marquardt algorithm for an example.

Not all multiple minima have equal values of the objective function. False minima, also known as local minima, occur when the objective function value is greater than its value at the so-called global minimum. For example, the model

$$f(x_i, a) = (1 - 3a + a^3)x_i$$

has a local minimum at a= 1 and a global minimum at a= -3

To be certain that the minimum found is the global minimum, the refinement should be started with widely differing initial values of the parameters. When the same minimum is found regardless of starting point, it is likely to be the global minimum.

When multiple minima exist there is an important consequence: the objective function will have a maximum value somewhere between two minima. The normal equations matrix is not positive definite at a maximum in the objective function, as the gradient is zero and no unique direction of descent exists. Refinement from a point (a set of parameter values) close to a maximum will be ill-conditioned and should be avoided as a starting point. For example, when fitting a Lorentzian the normal equations matrix is not positive definite when the half-width of the band is zero.[7]

The goal turn out that we have solve a problem of finding a minimum of a function S:

$$\begin{cases} \min \ S(a_1, a_2) = \sum_{i=1}^{m} r_i^2(a_1, a_2) \\ s.j.to \\ Constrains \ 1 \\ constrains \ 2 \\ \dots \\ constrain \ n_c \end{cases}$$

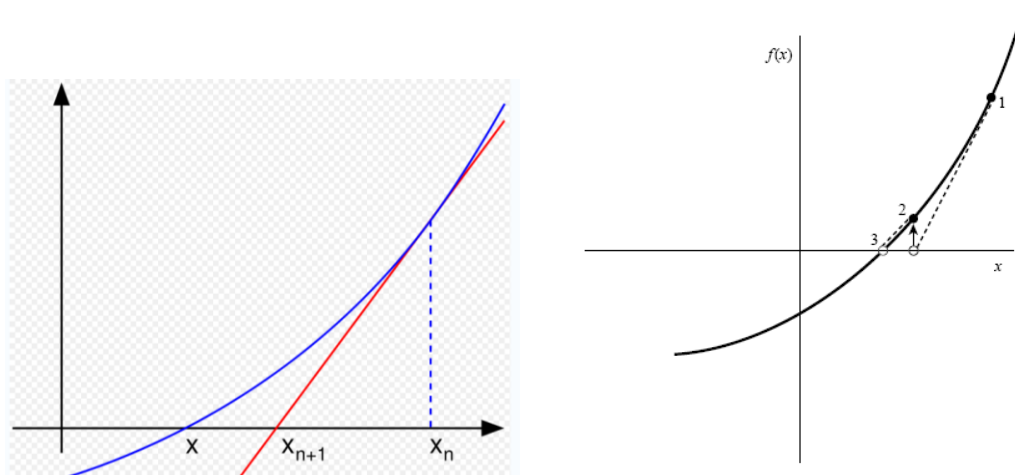this is a nonlinear optimization problem and data-dependent

**Method used in IGW**

The **Gauss–Newton algorithm** is a method used to solve [non-linear least squares](non-linear least squares) problems. Unlike Newton's method has the advantage that second derivatives, which can be challenging to compute, are not required.

=======================================

**General Newton Method:** Newton method , also know as the Newton-Raphson method, is perhaps the best know mehod for finding successively better approximations to the zeros (or roots) of a real-valued function. Newton's method can often converge remarkably quickly, especially if the iteration begins "sufficiently near " the desired root.

The idea of the method is as follows: one starts with an initail guess which is reasonably close to the true root, then the function is approximated by its tangent line, which can be computed using the tools of calculus, and one computes the x-intercept of this tangent line, which is easily done with elementary algebra. This x-intercept will typically be a better approximation to the function's root than the original guess, and the method can be iterated.



An illustration of one iteration of Newton's method (the function $f$ is shown in blue and the tangent line is in red). We see that $x_{n+1}$ is a better approximation than $x_n$ for the root $x$ of the function $f$.

Given a nonlinear function f(x) defined on the interval [a,b], is tangent line at point x0 can be written as

$$f_{linear}(x) = f(x_0) + \Delta x f'(x_0) \qquad\qquad (1)$$

Eq(1) can also be derived using Taylor seriers expansion by neglecting the second derivative terms and higher orders terms.

Finding the x-intercept of Eq(1) becomes to solve the following linear equation

$$f_{linear}(x) = f(x_0) + \Delta x f'(x_0) = 0 \tag{2}$$

This gives

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} \tag{3}$$

Then we can derive the formula for a better approximation, $x_2$, by referring to the diagram above:

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)} \tag{4}$$

The final iterative formula can be written as

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} = x_k - f'(x_k)^{-1} f(x_k) \tag{5}$$

---

**Newton Method for Optimization**

The necessary condition for the function $S = \sum_{i=1}^{m} r_i^2(a_1, a_2)$ to have an extremum is that the partial derivatives vanish i.e.

$$\sum_{i=1}^{m} r_i \frac{\partial r_i}{\partial a_j} = 0 \ (j = 1, 2, \cdots, n) \tag{6}$$

or, equivalently,

$$J^T(\vec{a}) \vec{r}(\vec{a}) = 0 \tag{1}$$

Where $J_{ij} = \dfrac{\partial r_i}{\partial a_j}$ is Jacobian matrix of residual vector $\vec{r} = (r_1, r_2, \cdots, r_m)$.

Equation (1) is usually a system of non-linear equations that numerically, can be solved by Netown-Rasphson's method. ➔ how ?

Denoting nonlinear system eq(6) as

$$\frac{\partial S}{\partial a_j} = F_j(\vec{a}) = 0 \ \ (j = 1, 2, \cdots, n)$$

Where $F_j$ is nonlinear funcion of vector $\vec{a}$. Referring to eq(5),

$$F_j(\vec{a} + \Delta \vec{a}) = F_j(\vec{a}) + \sum_{i=1}^{n} \frac{\partial F_j}{\partial a_i} \Delta \vec{a} + O(\Delta \vec{a}^2) \ \ \ (j = 1, 2, \cdots, n)$$

The matrix of partial derivatives appearing in above equation is the Jacobian matrix J:

$$J_{ij} = \frac{\partial F_i}{\partial a_j}$$

In matrix notation,

$$F(\vec{a}^{(k+1)}) = F(\vec{a}^{(k)}) + J_F((\vec{a}^{(k)})(\vec{a}^{(k+1)} - \vec{a}^{(k)}) + O(\Delta \vec{a}^2)$$

By neglecting terms of order $O(\Delta \vec{a}^2)$ and higher and by setting $F(\vec{a}^{(k+1)}) = 0$, one obtains the similar iterative formula (Rather than actually computing the inverse of this matrix, one can save time by solving the system of linear equations)

$$J_F((\vec{a}^{(k)})(\vec{a}^{(k+1)} - \vec{a}^{(k)}) = -F(\vec{a}^{(k)})$$

Recalling that $F_j(\vec{a}) = \dfrac{\partial S}{\partial a_j}$ and $J_{F,ij} = \dfrac{\partial F_i(\vec{a})}{\partial a_j} = \dfrac{\partial}{\partial a_j}(\dfrac{\partial S(\vec{a})}{\partial a_i}) = \dfrac{\partial^2 S(\vec{a})}{\partial a_i \partial a_j}$, from above equation, gives

$$H(\vec{a}^{(k+1)} - \vec{a}^{(k)}) = -\nabla S$$

Where $\vec{g} = \nabla S$ is gradient, $H$ is Hessian matrix of S, or
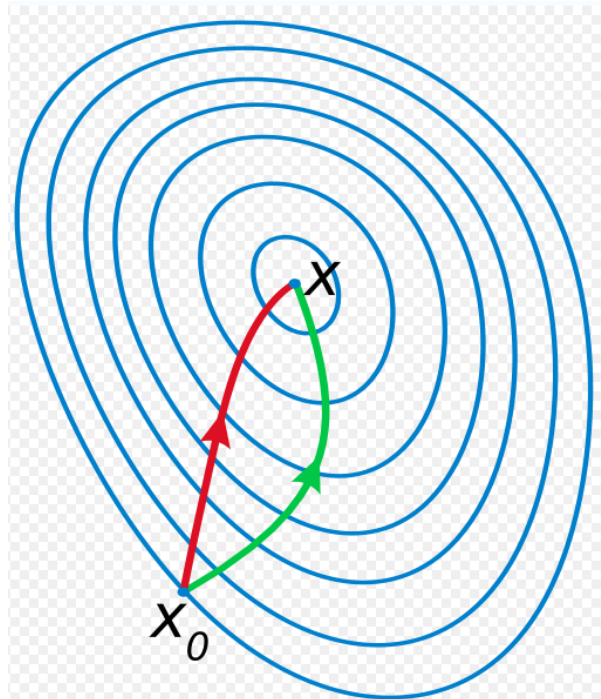
$$\vec{a}^{(k+1)} = \vec{a}^{(k)} - H^{-1}g$$

This is the recurrence relation for Newton's method for minimizing a function $S$ of parameters, where **g** denotes the gradient vector of $S$ and **H** denotes the Hessian matrix of $S$.

As can be seen , Newton's method needs to calculate the Hessian, which is always an expensive operation, so it is very hard to apply to those cases with difficulties in the second derivatives calculations.   For these cases an approximation for Hessian is used instead, which leads to various modified Newton methods such as Gauss-Newton, quasi-Newton and others.
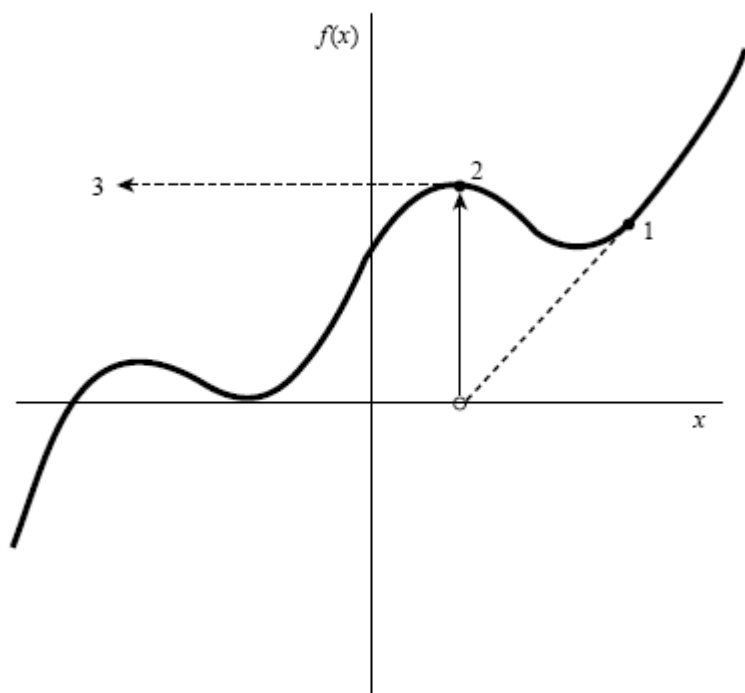
**Note 1:**

Newton's method is different from Gradient descent method: $\vec{a}^{(k+1)} = \vec{a}^{(k)} - \alpha^{(k)}g$, which is

based on the observation that if the real-valued function $S(\vec{a}^{(k)})$ is defined and differentiable in a neighborhood of point $\vec{a}^{(k)}$ , then S decreases fastest if one goes from $\vec{a}^{(k)}$ in the

direction of the negative gradient of S at $\vec{a}^{(k)}$ , that is, $S(\vec{a}^{(k+1)}) \le S(\vec{a}^{(k)}) \le S(\vec{a}^{(k-1)}) \le \cdots$,
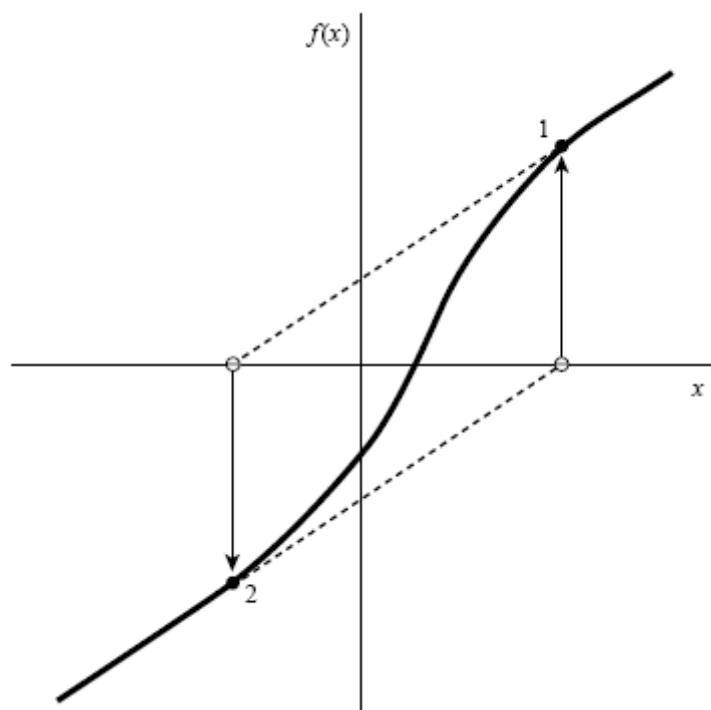
where $\alpha^{(k)}$ >0, a small enough number.



A comparison of gradient descent (green) and Newton's mehod(red) for minimizing a function (with small step sizes). Newton's method uses curvature information to take a more direct route.

**Note 2:**



Zero derivative



same derivative at different point

## Gauss–Newton algorithm : Derivation course I

Since $S = \sum_{i=1}^{m} r_i^2 (a_1, a_2)$, the gradient is given by

$$g_j = 2 \sum_{i=1}^{m} r_i \frac{\partial r_i}{\partial a_j}$$

Elements of the Hessian are calculated by differentiating the gradient elements, $g_j$, with respect to $\vec{a}$ :

$$H_{jk} = 2 \sum_{i=1}^{m} \left( \frac{\partial r_i}{\partial a_j} \frac{\partial r_i}{\partial a_k} + r_i \frac{\partial^2 r_i}{\partial a_j \partial a_k} \right)$$

The Gauss–Newton method is obtained by ignoring the second-order derivative terms (the second term in this expression). That is, the Hessian is approximated by

$$H_{jk} \approx 2 \sum_{i=1}^{m} \frac{\partial r_i}{\partial a_j} \frac{\partial r_i}{\partial a_k} = 2 \sum_{i=1}^{m} J_{ij} J_{ik}$$

The gradient and the approximate Hessian can be written in matrix notation as

$$\vec{g} = 2 J^T \vec{r}$$
$$H \approx 2 J^T J$$

⬅== at least $J_{ij}$ need to be evaluated ➔ sensitivities Mtx

These expressions are substituted into the recurrence relation above (Eq(2) to obtain the operational equations, the resulting linear system with unknown vector, $\Delta \vec{a}^{(k)}$

$$J^T (\vec{a}^{(k)}) J(\vec{a}^{(k)}) \Delta \vec{a}^{(k)} = -J(\vec{a}^{(k)}) \vec{r}(\vec{a}^{(k)})$$

Once unknown vector, $\Delta \vec{a}^{(k)}$ has been solved, the following
$$\vec{a}^{(k+1)} = \vec{a}^{(k)} - \alpha \Delta \vec{a}^{(k)}$$

is taken as the new approximation of the optimum. The relaxation factor $\alpha$ is a parameter of the method.

## Gauss–Newton algorithm : Derivation cousre II

Consider a set of m data points, (x1,y1),(x2,y2), …, (xm,ym), and a model function (a curve) y=f(x,a), that in addition to the variable x also dependes on n parameters, a=(a1,a2,…,an), with m>=n. It is desired to find the vector a of parameters such that the model (curve) fits best the given data in the least squares sense, that is, the sum of squares

$$S = \sum_{i=1}^{m} r_i^2 (a_1, a_2)$$

Is minimized, where the residuals (errors) ri are given by

$$r_i (a_1, a_2) = y_i - f(a, x_i), \quad i=1,2,…,m$$

The minimum valuesof S occurs when the gradient is zero. Since the model contains n parameters there are n gradient equations:

$$\frac{\partial S}{\partial a_j} = 2 \sum_{i=1}^{m} r_i \frac{\partial r_i}{\partial a_j} = 0 \ (j = 1,2,\cdots,n)$$

In a non-linear system, the derivatives $\dfrac{\partial r_i}{\partial a_j}$ are functions of both the independent variable

and the parameters, so these gradient equations do not have a closed solution. Instead, initial values must be chosen for the parameters. Then, the parameters are refined interatively, that is, the values are obtained by successive approximation,

$$a_j \approx a_j^{k+1} = a_j^k + \Delta a_j$$

Here, k is an iteration number and the vector of increments, $\Delta a$ is known as the shift vector. At each iteration **the model is linearized** by approximation to a first-order Taylor series expansion about $a_j^k$

$$f(x_i, a) = f(x_i, a^k) + \sum_{j}^{n} \frac{\partial f(x_i, a^k)}{\partial a_j} (a_j - a_j^k) + \frac{1}{2} \sum_{j}^{n} \sum_{l}^{n} \frac{\partial^2 f(x_i, a^k)}{\partial a_j \partial a_l} (a_j - a_j^k)(a_l - a_l^k) + \cdots$$

$$= f(x_i, a^k) + \sum_{j}^{n} J_{ij} \Delta a_j + \frac{1}{2} \sum_{j}^{n} \sum_{l}^{n} H_{jl_{(i)}} \Delta a_j \Delta a_l + \cdots$$

Or approximately

$$f(x_i, a) \approx f(x_i, a^k) + \sum_j^n J_{ij} \Delta a_j$$

The jacobian, J is a function of constants, the independent variables and the parameters, so it changes from iteration to the next. Thus, **in term of the linearized model**, $\dfrac{\partial r_i}{\partial a_j} = J_{ij}$ and the residuals are given by

$$r_i = y_i - f(a, x_i) = y_i - f(x_i, a^k) - \sum_j^n J_{ij} \Delta a$$

Substituting these expressions into the gradient equations, they become

$$-2\sum_{i=1}^m r_i \frac{\partial r_i}{\partial a_j} = -2\sum_{i=1}^m \frac{\partial r_i}{\partial a_j}[y_i - f(x_i, a^k) - \sum_l^n J_{il}\Delta a] = 0 \ (j = 1,2,\cdots,n)$$

Or

$$-2\sum_{i=1}^m J_{ij}[y_i - f(x_i, a^k) - \sum_l^n J_{il}\Delta a] = 0 \ (j = 1,2,\cdots,n)$$

Which, on rearrangement, become n simultaneous linear equations, the normal equations

$$\sum_{i=1}^m \sum_l^n J_{ij}J_{il}\Delta a_l = \sum_{i=1}^m J_{ij}[y_i - f(x_i, a^k) \ (j = 1,2,\cdots,n)$$

The normal equations are written in matrix notation as

$$\boldsymbol{J}^T(\vec{a}^{(k)})\boldsymbol{J}(\vec{a}^{(k)})\Delta\vec{a}^{(k)} = \boldsymbol{J}^T(\vec{a}^{(k)})\vec{r}(\vec{a}^{(k)})$$

When the observations are not equally reliable, a weighted sum of squares may be minimized,

$$S = \sum_{i=1}^m W_{ii} r_i^2(a_1, a_2)$$

Each element of the diagonal weight matrix, W should , ideally, be equal to the reciprocal of the variance of the measurement. This implies that the observations are uncorrelated. If the observations are correlated, the expression

$$S = \sum_{i=1}^m \sum_j^m r_i W_{ij} r_j$$

Applies. In this case the weight matrix should ideally be equal to the inverse if the variance-covariance matrix of the observations. The normal equations are then

$$J^T(\vec{a}^{(k)})WJ(\vec{a}^{(k)})\Delta\vec{a}^{(k)} = J^T(\vec{a}^{(k)})W\vec{r}(\vec{a}^{(k)})$$
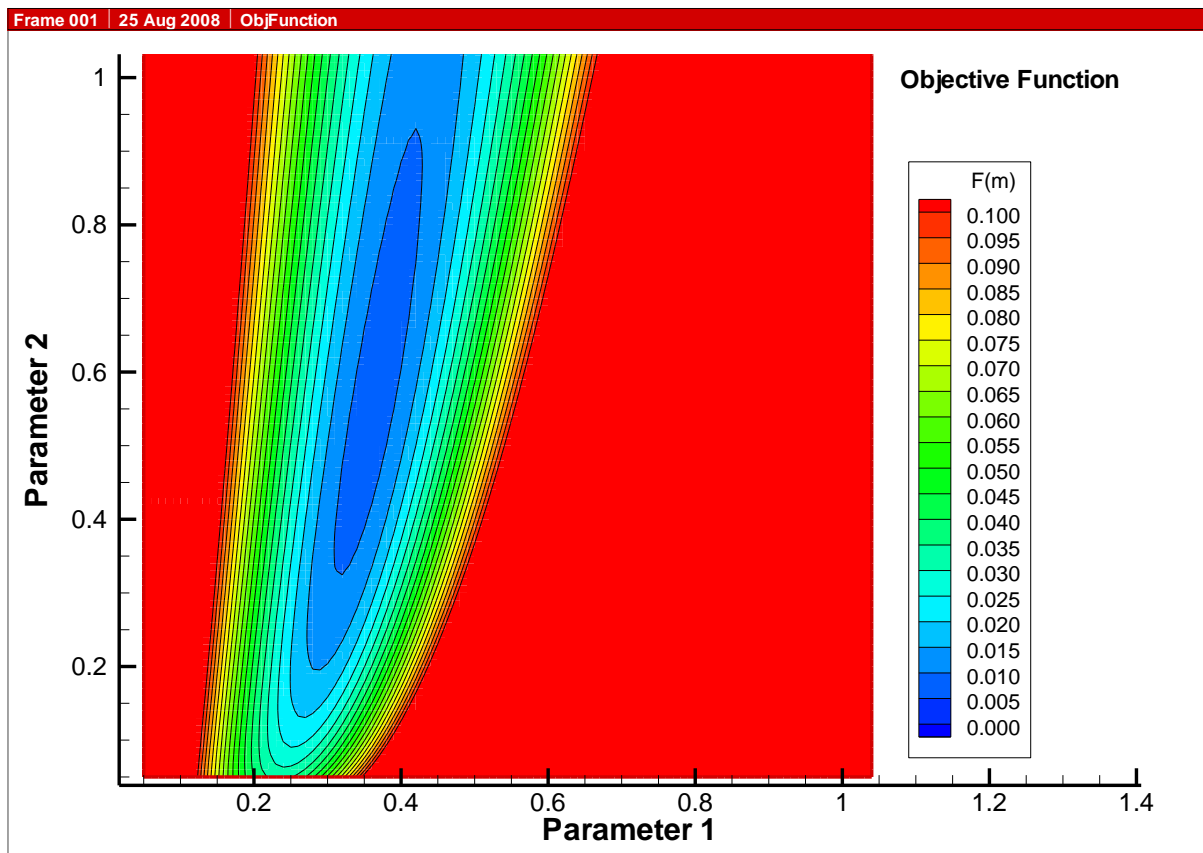
These equations form the basis for the Gauss-Newton algorithm for a non-linear least squares problem.

**Example**

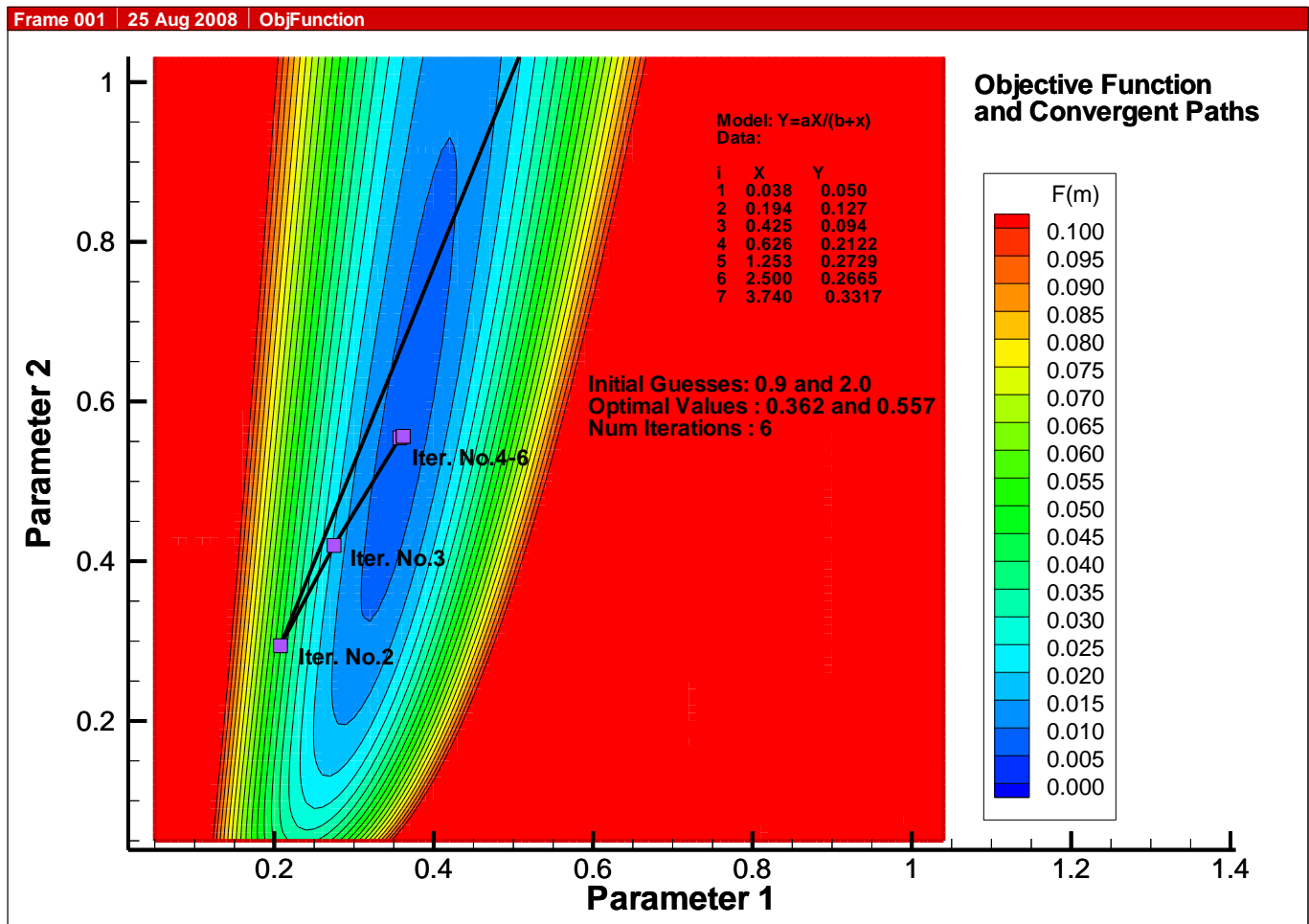Model: $\gamma(a_1, a_2) = \dfrac{a_1 x}{a_2 + x}$
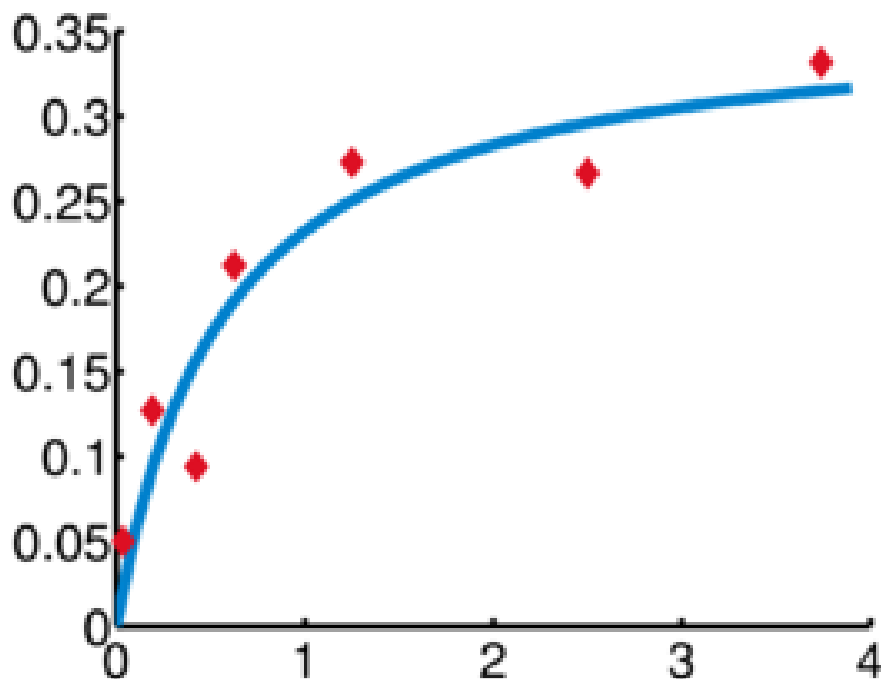
And observations:

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-----|-------|-------|-------|--------|--------|--------|--------|
| x | 0.038 | 0.194 | 0.425 | 0.626 | 1.253 | 2.500 | 3.740 |
| $\gamma$ | 0.050 | 0.127 | 0.094 | 0.2122 | 0.2729 | 0.2665 | 0.3317 |



Objective Function within $0<a_1<1.05$ and $0<a_2<1.05$

**Objective Function and Convergent Paths**

Model: Y=aX/(b+x)
Data:

| i | X | Y |
|---|---|---|
| 1 | 0.038 | 0.050 |
| 2 | 0.194 | 0.127 |
| 3 | 0.425 | 0.094 |
| 4 | 0.626 | 0.2122 |
| 5 | 1.253 | 0.2729 |
| 6 | 2.500 | 0.2665 |
| 7 | 3.740 | 0.3317 |

Initial Guesses: 0.9 and 2.0
Optimal Values : 0.362 and 0.557
Num Iterations : 6

Convergent Paths starting with initial guesses $a_1$=0.9 and $a_2$=2.0
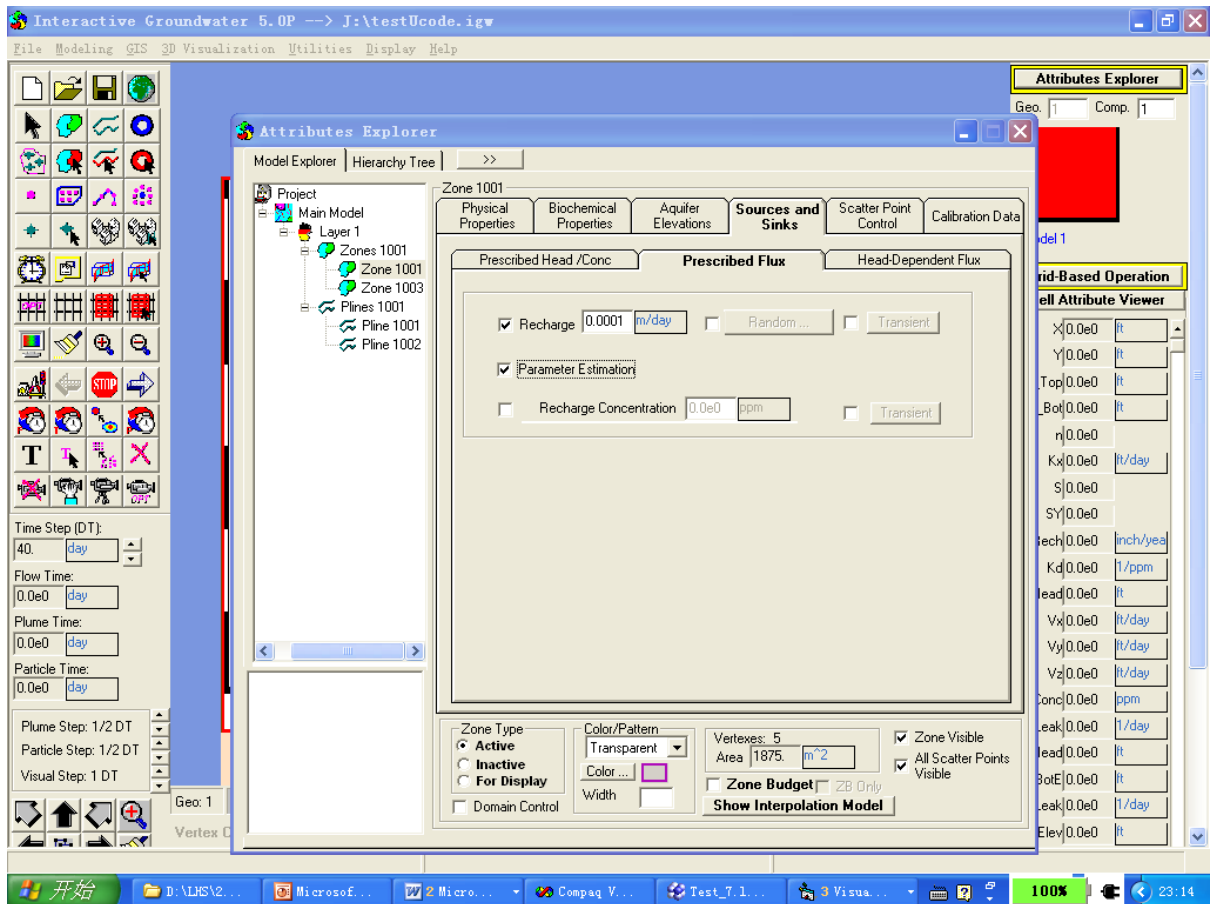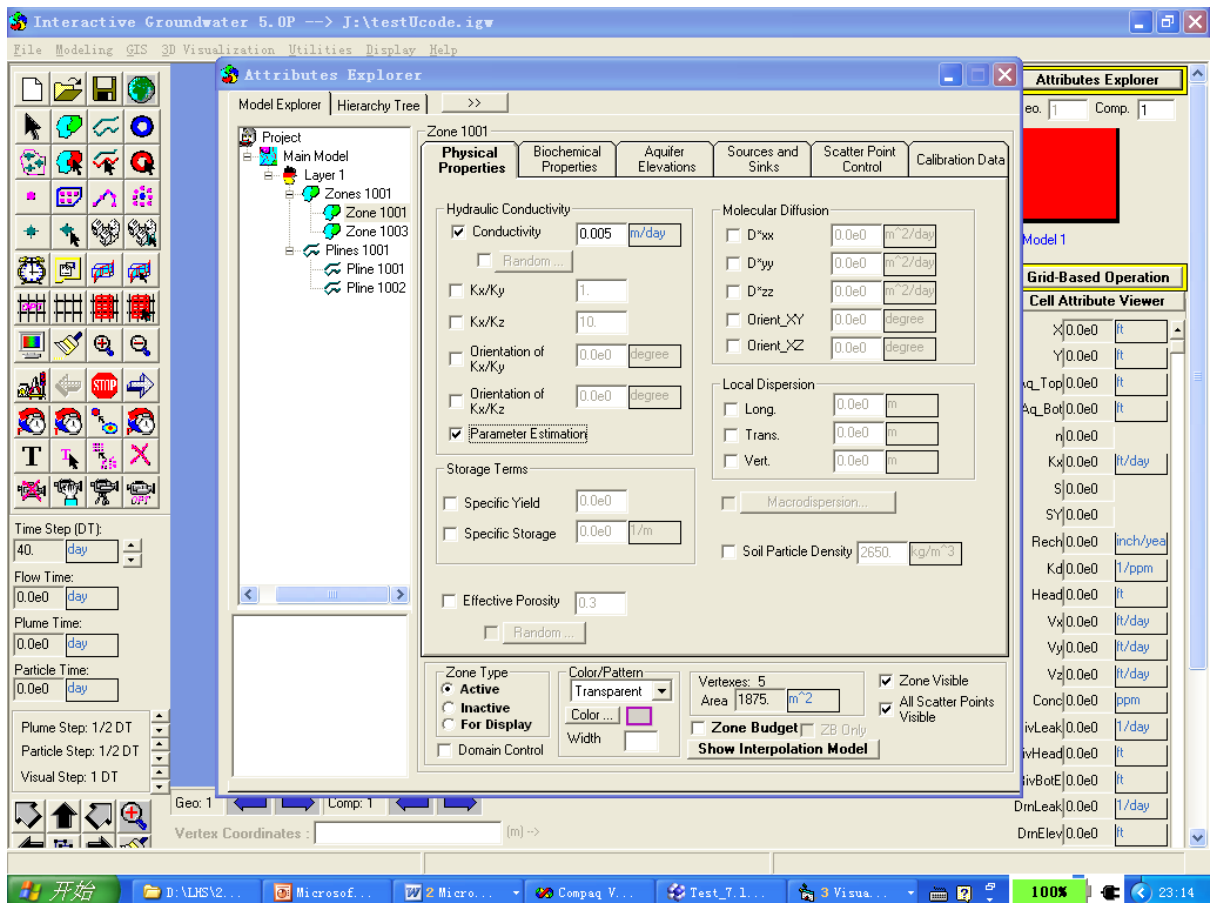Optimal paramters are $a_1$=0.362 and $a_2$=0.557

# 1 Automatic Variogram Fitting



Button "BestFit"

Parameters Estimation window

# 2 Parameter Estimation by using IGW

(1) Seleting parameter zones: Check those polygon-based parameters to be estimated when create your conceptual model (only conductivity and recharge are available right now)
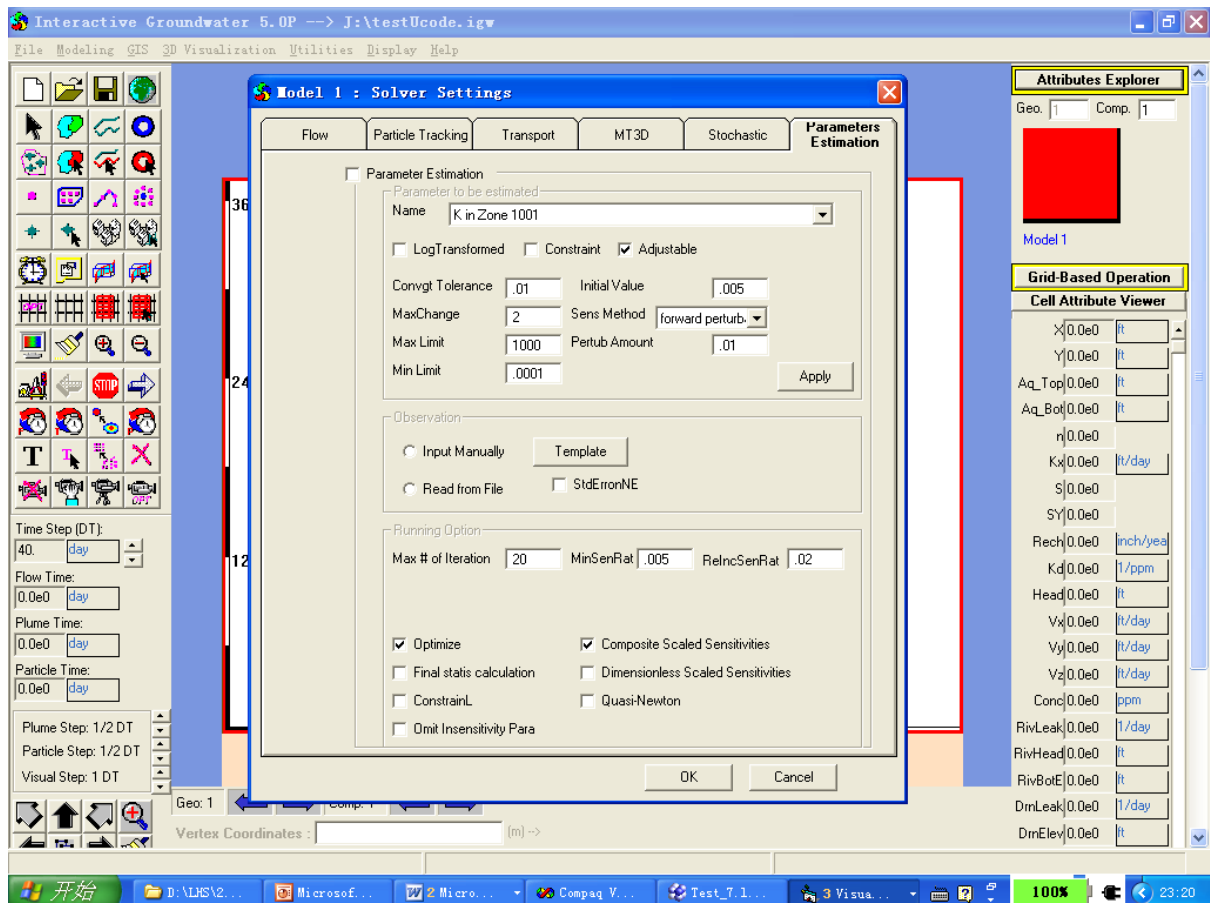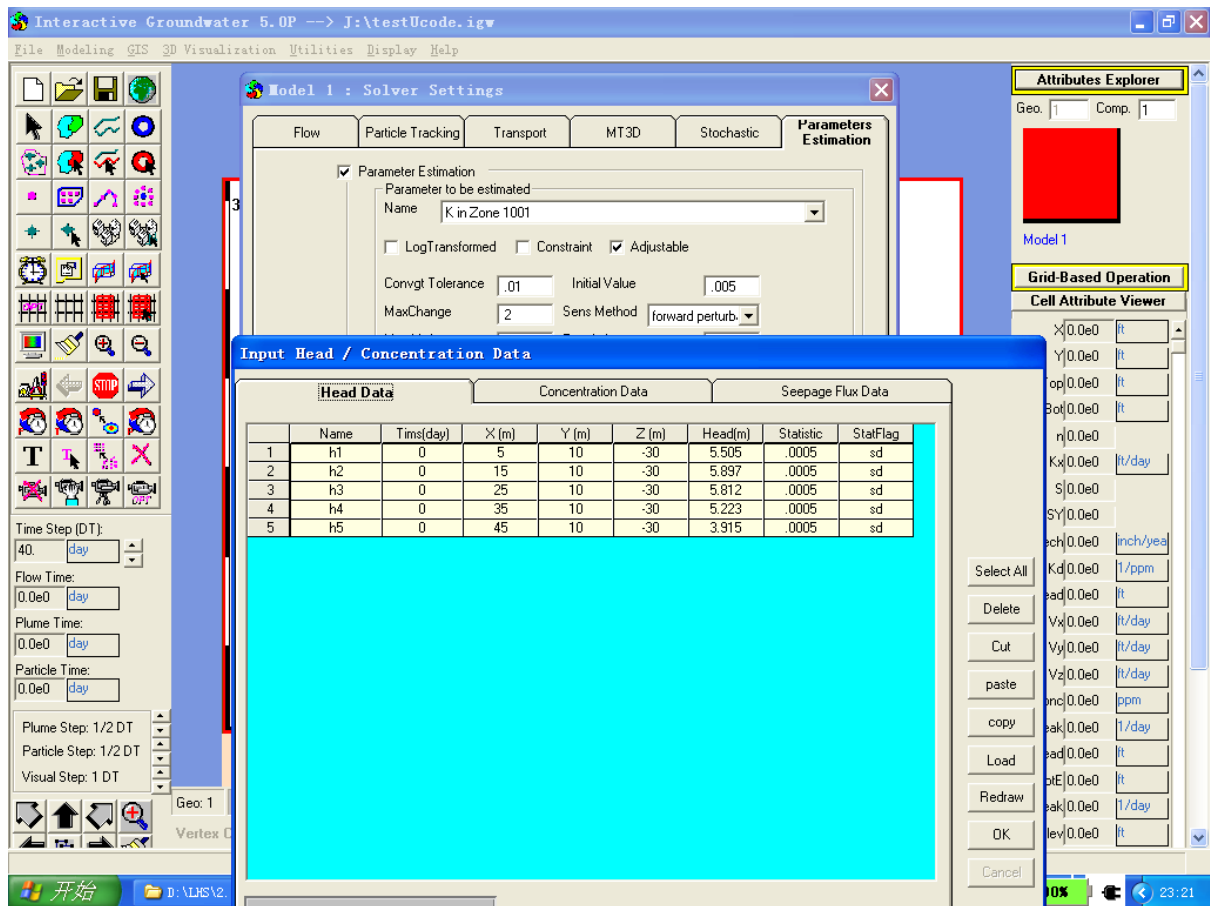
(2) Select Solver: hit  lead to



New tab in "Solver" toobar: Parameter Estimation
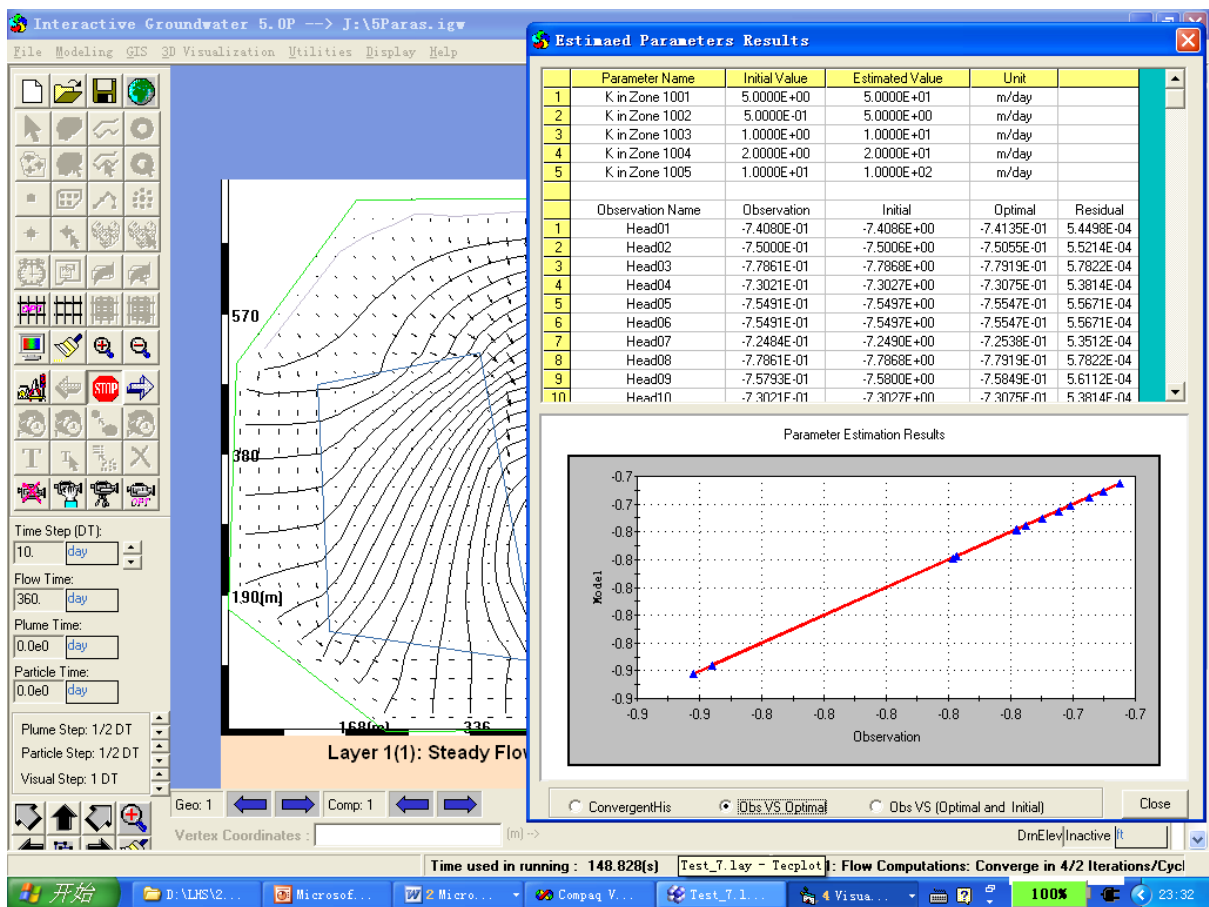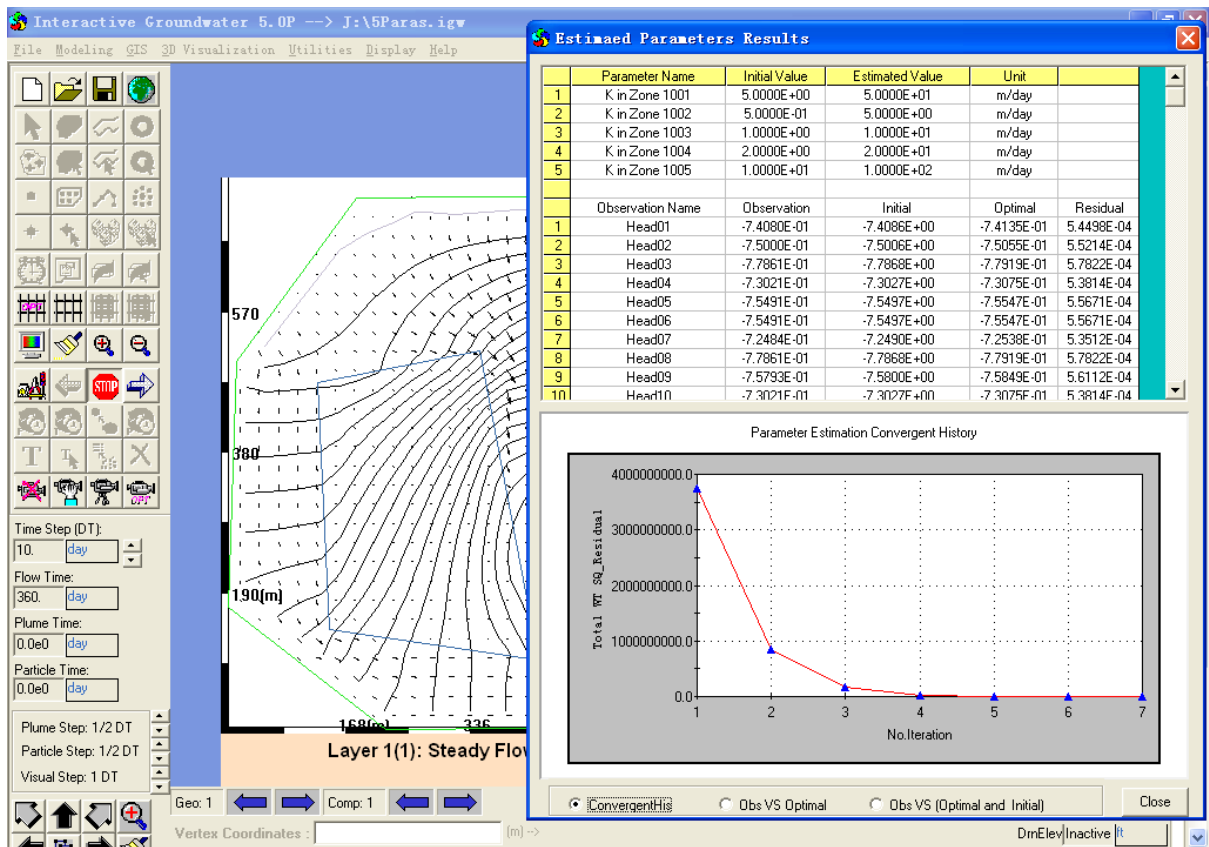
(3) Check "parameter estimation" check box

In the new tab: Newton method option, constrains, method to calculate senstivity matrix, …

(4) Observation input

Observation input: file/ manually input/copy-paste
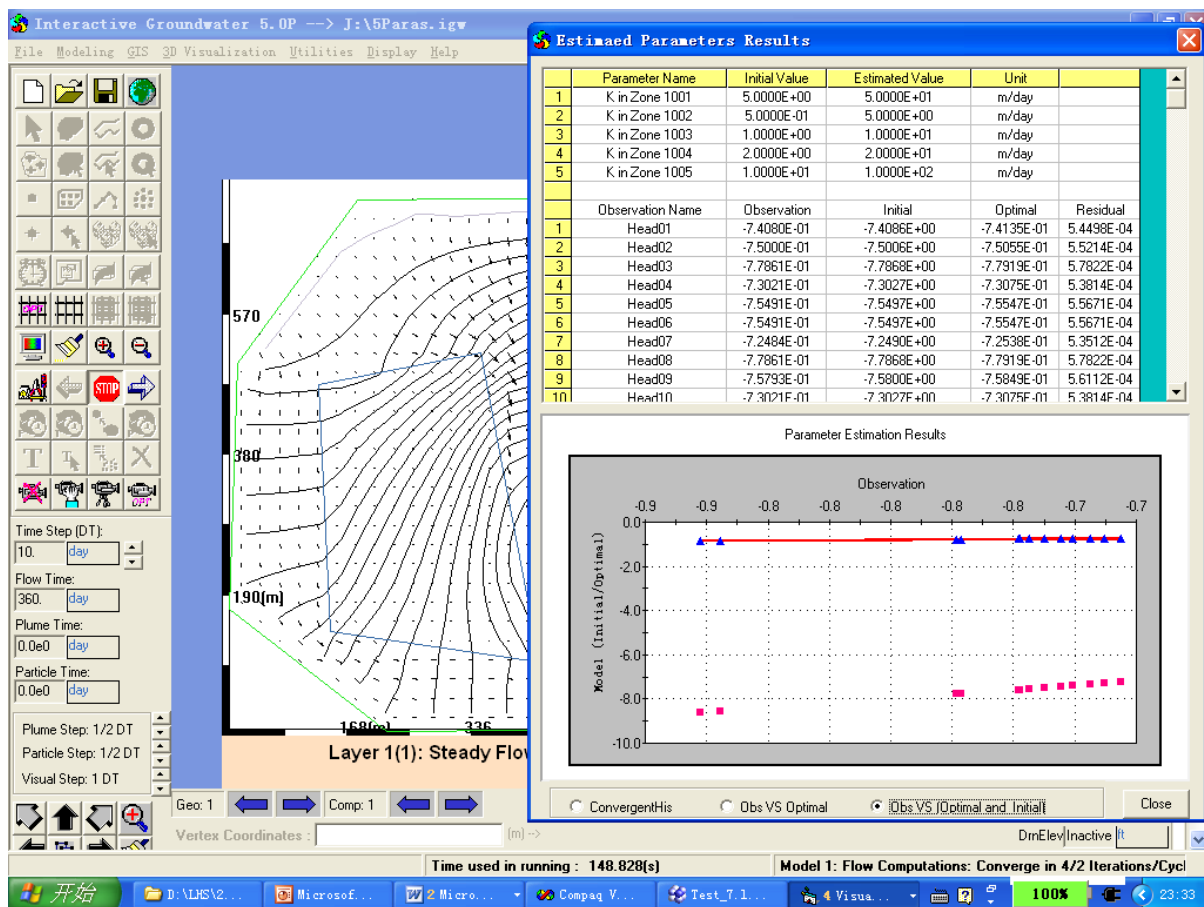
(5) Results --- parameter estimation results window

Table: Parameters estimated/residuals

Charts:

Iteration history

Observations VS model values

Initial guesses VS Optimal values

**Demonstration:**

1 GIS data variogram model

2 5 K parameters estimated by using IGW flow model

3 5 K + 1 recharge parameters estimated by using IGW flow model