
FBMath

Math library for FreeBasic

Jean Debord

May 5, 2007

Contents

1	Installation and compilation	7
1.1	Introduction	7
1.2	Unpacking the archive	7
1.3	Installation under Windows	7
1.3.1	Automatic installation	7
1.3.2	Manual installation	8
1.3.3	Recompiling the library	8
1.4	Installation under Linux	8
1.4.1	Compilation of the library	8
1.4.2	Installation	8
1.4.3	Demo programs	9
1.5	Compilation of a program	9
1.5.1	Linking to the whole libraries	9
1.5.2	Using specific modules	10
2	Numeric precision	11
2.1	Numeric precision	11
2.2	Demo program	12
3	Elementary functions	13
3.1	Constants	13
3.2	Error handling	14
3.3	Minimum and maximum	14
3.4	Rounding functions	14
3.5	Logarithms and exponentials	15
3.6	Trigonometric functions	15
3.7	Hyperbolic functions	15
3.8	Demo programs	16
3.8.1	Function accuracy	16
3.8.2	Computation speed	16
3.8.3	Function plotter	16

3.8.4	Contour plot	17
4	Special functions	18
4.1	Factorial	18
4.2	Gamma function	18
4.3	Polygamma functions	19
4.4	Beta function	20
4.5	Error function	20
4.6	Lambert's function	20
4.7	Demo programs	21
5	Probability distributions	22
5.1	Binomial distribution	22
5.2	Poisson distribution	23
5.3	Standard normal distribution	23
5.4	Student's distribution	24
5.5	Khi-2 distribution	25
5.6	Snedecor's distribution	25
5.7	Exponential distribution	26
5.8	Beta distribution	26
5.9	Gamma distribution	27
5.10	Demo programs	27
6	Matrices and linear equations	28
6.1	Programming conventions	28
6.2	Error codes	29
6.3	Gauss-Jordan elimination	29
6.4	LU decomposition	30
6.5	QR decomposition	31
6.6	Singular value decomposition	32
6.7	Cholesky decomposition	33
6.8	Eigenvalues and eigenvectors	33
6.8.1	Definitions	33
6.8.2	Symmetric matrices	34
6.8.3	General square matrices	34
6.9	Demo programs	35
6.9.1	Determinant and inverse of a square matrix	35
6.9.2	Hilbert matrices	36
6.9.3	Gauss-Jordan method: single constant vector	36
6.9.4	Gauss-Jordan method: multiple constant vectors	37
6.9.5	LU, QR and SV decompositions	38

6.9.6	Cholesky decomposition	38
6.9.7	Eigenvalues of a symmetric matrix	38
6.9.8	Eigenvalues of a general square matrix	38
6.9.9	Eigenvalues and eigenvectors of a general square matrix	39
7	Function minimization	41
7.1	Functions of one variable	41
7.2	Functions of several variables	42
7.2.1	Minimization along a line	42
7.2.2	Newton-Raphson method	43
7.2.3	Approximate gradient and hessian	45
7.2.4	Marquardt method	45
7.2.5	BFGS method	46
7.2.6	Approximate gradient	47
7.2.7	Simplex method	47
7.3	Demo programs	47
7.3.1	Function of one variable	47
7.3.2	Minimization along a line	48
7.3.3	Newton-Raphson method	48
7.3.4	Approximate gradient and hessian	49
7.3.5	Other programs	49
8	Nonlinear equations	50
8.1	Equations in one variable	50
8.1.1	Bisection method	50
8.1.2	Secant method	51
8.1.3	Newton-Raphson method	51
8.2	Equations in several variables	52
8.2.1	Newton-Raphson method	52
8.2.2	Approximate jacobian	54
8.2.3	Broyden's method	54
8.3	Demo programs	55
8.3.1	Equations in one variable	55
8.3.2	Equations in several variables	55
9	Polynomials	56
9.1	Polynomials	56
9.2	Rational fractions	56
9.3	Roots of polynomials	56
9.3.1	Analytical methods	56
9.3.2	Iterative method	57

9.4	Ancillary functions	57
9.5	Demo programs	58
9.5.1	Evaluation of a polynomial	58
9.5.2	Evaluation of a rational fraction	58
9.5.3	Roots of a polynomial	58
10	Numerical integration and differential equations	59
10.1	Integration	59
10.1.1	Trapezoidal rule	59
10.1.2	Gauss-Legendre integration	59
10.2	Convolution	60
10.3	Differential equations	60
10.4	Demo programs	63
11	Fast Fourier Transform	65
11.1	Introduction	65
11.2	Programming	66
11.2.1	Array dimensioning	66
11.2.2	FFT procedures	66
11.3	Demo program	67
12	Random numbers	69
12.1	Random numbers	69
12.1.1	Introduction	69
12.1.2	Generic functions	70
12.1.3	Specific functions	70
12.1.4	Gaussian random numbers	71
12.2	Markov Chain Monte Carlo	72
12.3	Simulated Annealing	75
12.4	Genetic Algorithm	78
12.5	Demo programs	80
12.5.1	Test of MWC generator	80
12.5.2	Test of MT generator	80
12.5.3	Test of UVAG generator	80
12.5.4	File of random numbers	80
12.5.5	Gaussian random numbers	80
12.5.6	Multinormal distribution	81
12.5.7	Markov Chain Monte-Carlo	81
12.5.8	Simulated Annealing	82
12.5.9	Genetic Algorithm	82

13 Statistics	83
13.1 Descriptive statistics	83
13.2 Comparison of means	85
13.2.1 Student's test for independent samples	85
13.2.2 Student's test for paired samples	86
13.2.3 One-way analysis of variance (ANOVA)	86
13.2.4 Two-way analysis of variance	88
13.3 Comparison of variances	89
13.3.1 Comparison of two variances	89
13.3.2 Comparison of several variances	90
13.4 Non-parametric tests	90
13.4.1 Mann-Whitney test	91
13.4.2 Wilcoxon test	91
13.4.3 Kruskal-Wallis test	92
13.5 Statistical distribution	93
13.6 Comparison of distributions	93
13.6.1 Observed and theoretical distributions	93
13.6.2 Several observed distributions	94
13.7 Demo programs	95
13.7.1 Descriptive statistics, comparison of means and variances	95
13.7.2 Student's test for paired samples	95
13.7.3 One-way analysis of variance	96
13.7.4 Two-way analysis of variance	96
13.7.5 Statistical distribution	96
13.7.6 Comparison of distributions	97
14 Linear regression	98
14.1 Straight line fit	98
14.2 Analysis of variance	100
14.3 Precision of parameters	101
14.4 Probabilistic interpretation	101
14.5 Weighted regression	102
14.6 Programming	103
14.6.1 Regression procedures	103
14.6.2 Quality of fit	104
14.7 Demo programs	104
14.7.1 Unweighted linear regression	105
14.7.2 Weighted linear regression	105

15 Multilinear regression and principal component analysis	106
15.1 Multilinear regression	106
15.1.1 Normal equations	106
15.1.2 Analysis of variance	107
15.1.3 Precision of parameters	108
15.1.4 Probabilistic interpretation	108
15.1.5 Weighted regression	108
15.1.6 Programming	109
15.2 Principal component analysis	109
15.2.1 Theory	109
15.2.2 Programming	110
15.3 Demo programs	111
15.3.1 Multilinear regression	111
15.3.2 Polynomial regression	112
15.3.3 Principal component analysis	113
16 Nonlinear regression	114
16.1 Theory	114
16.2 Monte-Carlo simulation	116
16.3 Demo programs	117
16.3.1 Nonlinear regression	117
16.3.2 Monte-Carlo simulation	118
17 String functions	119
17.1 Fill functions	119
17.2 Character replacement	119
17.3 Parsing	119
17.4 Formatting functions	120

Chapter 1

Installation and compilation

1.1 Introduction

Welcome to **FBMath**, a mathematical package for the FreeBasic compiler. **FBMath** is entirely written in FreeBasic and does not depend on external libraries. **FBMath** is comprised with two independent libraries:

- a general library (**libmath**) for numerical analysis, including mathematical functions, probabilities, random numbers, matrices, linear and nonlinear equations, optimization, statistics and graphics.
- a ‘Large Integer’ library (**liblargeint**), contributed by Sjoerd J. Schaper.

This chapter explains how to install the **FBMath** package and how to compile a program which uses it.

1.2 Unpacking the archive

Extract the archive **fbmat[...].zip** (where [...] stands for version number) in a given directory.

Be sure to preserve the directory structure. For instance, if you use **pkunzip**, add the option **-d** (i. e. **pkunzip -d fbmat[...].zip**).

1.3 Installation under Windows

1.3.1 Automatic installation

Run the **compil.bat** script which is located in the **demo** subdirectory of the installation directory.

This script will compile all demo programs and install the library files in the relevant directories.

Note: The script assumes that the FreeBasic compiler has been installed in `C:\Program Files\FreeBasic`. If this is not the case, you will have to edit the batch file.

1.3.2 Manual installation

- Copy the library files `libmath.a` and `liblargeint.a` to the `lib` subdirectory of the FreeBasic directory (usually `\Program Files\FreeBasic\lib\win32`)
- Copy the include files `math.bi` and `largeint.bi` to the `inc` subdirectory of the FreeBasic directory (usually `\Program Files\FreeBasic\inc`)

1.3.3 Recompiling the library

If you have to recompile the library (e. g. after modifying the source code), run the `compil.bat` script which is located in the `modules` subdirectory of the installation directory.

You may have to edit the batch file, for instance to modify the paths for the compiler and the library manager (`ar.exe`), or if you add a directory to the source code.

1.4 Installation under Linux

1.4.1 Compilation of the library

Run the shell script `compil.sh` located in the `modules` subdirectory. This will create the two library files, `libmath.a` and `liblargeint.a`.

1.4.2 Installation

- Copy the library files into the appropriate FreeBasic subdirectory (usually `/usr/share/freebasic/lib/linux`).
- Copy the include files `math.bi` and `largeint.bi` into the appropriate FreeBasic subdirectory (usually `/usr/share/freebasic/inc`).

You will probably have to become root to perform these operations.

1.4.3 Demo programs

The script `compil.sh` located in the `demo` subdirectory will compile all demo programs.

1.5 Compilation of a program

1.5.1 Linking to the whole libraries

Add the following line at the beginning of the program:

```
#INCLUDE "math.bi"
```

or

```
#INCLUDE "largeint.bi"
```

depending on which library is used (it is possible to use both in the same program).

Note: If you want to use the graphic functions, you must also add:

```
#INCLUDE "fbgfx.bi"
```

Then compile the program in the usual way, e. g. use the following command in a DOS box or Linux terminal:

```
fbc prog.bas
```

Note: Under Windows, the complete path to the compiler must be included in the environment variable `PATH`, which is defined in the `AUTOEXEC.BAT` file located in the root directory of the computer. This file should therefore contain a line like:

```
PATH= ... C:\Progra~1\FreeBasic; ...
```

(or equivalent)

1.5.2 Using specific modules

You can also select from the source files the modules containing the functions of interest and add them to your project. Note that some modules require functions from other modules (these functions are listed in the ‘External functions’ section of each module). Be sure to add these additional modules too, if necessary.

Example: We have a program `prog.bas` which computes some factorials by calling function `Fact` which is defined in module `fact.bas`. However, looking at the ‘External function’ section in this module indicates that it requires function `Gamma` which is defined in module `gamma.bas`. So, we have to add both `fact.bas` and `gamma.bas` to our project (and only these two, because looking at `gamma.bas` shows that no external function is needed).

Then, we place the declaration of function `Fact` in our main module, `prog.bas`:

```
DECLARE FUNCTION Fact(N AS INTEGER) AS DOUBLE
```

We no longer need to add `#INCLUDE "math.bi"`

We can then compile our program, assuming that all modules have been placed in the same working directory:

```
fbx prog.bas fact.bas gamma.bas
```

(Note that the main module is the first).

Chapter 2

Numeric precision

2.1 Numeric precision

All computations are performed in double precision (type `DOUBLE`, i.e. 8-byte real).

FBMath defines the following constants:

Constant	Meaning
<code>MachEp</code>	The smallest real number such that $(1.0 + \text{MachEp})$ has a different representation (in the computer memory) than 1.0; it may be viewed as a measure of the numeric precision which can be reached within the given floating point type.
<code>MaxNum</code>	The highest real number which can be represented.
<code>MinNum</code>	The lowest positive real number which can be represented.
<code>MaxLog</code>	The highest real number X for which $\text{Exp}(X)$ can be computed without overflow.
<code>MinLog</code>	The lowest (negative) real number X for which $\text{Exp}(X)$ can be computed without underflow.
<code>MaxFac</code>	The highest integer for which the factorial can be computed.
<code>MaxGam</code>	The highest real number for which the Gamma function can be computed.
<code>MaxLgm</code>	The highest real number for which the logarithm of the Gamma function can be computed.

2.2 Demo program

The program `testmach.bas` located in the `demo\fmath` subdirectory checks that the machine-dependent constants are correctly handled by the computer.

This program lists the values of the machine-dependent constants and computes the following quantities:

<code>Exp(MinLog)</code>	Should be approximately equal to <code>MinNum</code>
<code>Ln(MinNum)</code>	Should be approximately equal to <code>MinLog</code>
<code>Exp(MaxLog)</code>	Should be approximately equal to <code>MaxNum</code>
<code>Ln(MaxNum)</code>	Should be approximately equal to <code>MaxLog</code>
<code>Fact(MaxFac)</code>	
<code>Gamma(MaxGam)</code>	Should be computed without overflow.
<code>LnGamma(MaxLgm)</code>	

The following results were obtained with FreeBasic 0.16b:

```
MachEp          = 2.220446049250313e-016

MinNum          = 2.225073858507202e-308
Exp(MinLog)     = 2.225073858507263e-308

MinLog          = -708.3964185322641
Ln(MinNum)      = -708.3964185322641

MaxNum          = 1.797693134862315e+308
Exp(MaxLog)     = 1.797693134862273e+308

MaxLog          = 709.782712893384
Ln(MaxNum)      = 709.782712893384

MaxFac          = 170
Fact(MaxFac)    = 7.257415615308285e+306

MaxGam          = 171.624376956302
Gamma(MaxGam)   = 1.797693134855531e+308

MaxLgm          = 2.556348e+305
LnGamma(MaxLgm) = 1.795136671459441e+308
```

Chapter 3

Elementary functions

This chapter describes the mathematical constants and elementary mathematical functions available in `FBMath`.

3.1 Constants

The following mathematical constants are defined in `math.bi`:

Constant	Value	Meaning
Pi	3.14159...	π
Ln2	0.69314...	$\ln 2$
Ln10	2.30258...	$\ln 10$
LnPi	1.14472...	$\ln \pi$
InvLn2	1.44269...	$1/\ln 2$
InvLn10	0.43429...	$1/\ln 10$
TwoPi	6.28318...	2π
PiDiv2	1.57079...	$\pi/2$
SqrtPi	1.77245...	$\sqrt{\pi}$
Sqrt2Pi	2.50662...	$\sqrt{2\pi}$
InvSqrt2Pi	0.39894...	$1/\sqrt{2\pi}$
LnSqrt2Pi	0.91893...	$\ln \sqrt{2\pi}$
Ln2PiDiv2	0.91893...	$(\ln 2\pi)/2$
Sqrt2	1.41421...	$\sqrt{2}$
Sqrt2Div2	0.70710...	$\sqrt{2}/2$
Gold	1.61803...	Golden Ratio = $(1 + \sqrt{5})/2$
CGold	0.38196...	

In addition, the logical constants `True` (-1) and `False` (0) are defined.

3.2 Error handling

The function `MathErr()` returns the error code from the last function evaluation. It must be checked immediately after a function call:

```
Y = f(X)    ' f is one of the functions of the library
if MathErr = FOk then ...
```

If an error occurs, a default value is attributed to the function. The possible error codes are the following:

Error code	Value	Meaning
FOk	0	No error
FDomain	-1	Argument domain error
FSing	-2	Function singularity
FOverflow	-3	Overflow range error
FUnderflow	-4	Underflow range error

3.3 Minimum and maximum

- Function `Min(X, Y)` returns the minimum of two real numbers X, Y .
- Function `Max(X, Y)` returns the maximum of two real numbers X, Y .
- Function `IMin(X, Y)` returns the minimum of two integer numbers X, Y .
- Function `IMax(X, Y)` returns the maximum of two integer numbers X, Y .

3.4 Rounding functions

- Function `Round(X, Digit_Count)` will round X to `Digit_Count` decimal places. `Digit_Count` must be between 0 and 16. The default value is 0, so that `Round(X)` is equivalent to `Round(X, 0)`.
- Function `Floor(X)` returns the lowest integer $\geq X$
- Function `Ceil(X)` returns the highest integer $\leq X$

3.5 Logarithms and exponentials

The functions `Expo` and `Ln` may be used instead of the standard functions `Exp` and `Log`, when it is necessary to check the range of the argument. The new function performs the required tests and calls the standard function if the argument is within the acceptable limits (for instance, $X > 0$ for `Ln(X)`); otherwise, the function returns a default value and `MathErr()` will return the appropriate error code.

Other logarithmic and exponential functions are:

Function	Definition
<code>Exp2(X)</code>	2^X
<code>Exp10(X)</code>	10^X
<code>Log2(X)</code>	$\log_2 X$
<code>Log10(X)</code>	$\log_{10} X$
<code>LogA(X, A)</code>	$\log_A X$

3.6 Trigonometric functions

FBMath provides two additional functions:

Function	Definition
<code>Pythag(X, Y)</code>	$\sqrt{X^2 + Y^2}$
<code>FixAngle(Theta)</code>	Returns the angle <code>Theta</code> in the range $[-\pi, \pi]$

3.7 Hyperbolic functions

The following functions are available:

Function	Definition
Sinh(X)	$\frac{1}{2}(e^X - e^{-X})$
Cosh(X)	$\frac{1}{2}(e^X + e^{-X})$
Tanh(X)	$\frac{\sinh X}{\cosh X}$
ASinh(X)	$\ln(X + \sqrt{X^2 + 1})$
ACosh(X)	$\ln(X + \sqrt{X^2 - 1}) \quad X > 1$
ATanh(X)	$\frac{1}{2} \ln \frac{X+1}{X-1} \quad -1 < X < 1$

In addition, the subroutine **SinhCosh(X, SinhX, CoshX)** computes the hyperbolic sine and cosine simultaneously, saving the computation of one exponential.

3.8 Demo programs

These programs are located in the `demo\fmath` subdirectory.

3.8.1 Function accuracy

Program `testfunc.bas` checks the accuracy of the elementary functions. For each function, 20 random arguments are picked, then the function is computed, the reciprocal function is applied to the result, and the relative error between this last result and the original argument is computed. This error should be around 10^{-15} in double precision.

3.8.2 Computation speed

Program `speed.bas` measures the execution time of the built-in mathematical functions, as well as the additional functions provided in **FBMath**. The results are printed on the screen and saved in a text file named `speed.out`.

3.8.3 Function plotter

Program `plotfunc.bas` allows to plot a function in linear or logarithmic coordinates. The square root function is taken as example, since its graph becomes a straight line in log-log coordinates.

3.8.4 Contour plot

Program `contour.bas` draws a contour plot of a function of two variables. It uses the `ConRec` algorithm developed by Paul Bourke (<http://astronomy.swin.edu.au/~pbourke/projection/conrec/>).

The example function is:

$$f(x, y) = \sin \sqrt{x^2 + y^2} + \frac{1}{2\sqrt{(x + c)^2 + y^2}}$$

where c is a constant which may be varied to modify the aspect of the graph.

Chapter 4

Special functions

This chapter describes the special functions available in **FBMath**. Most of them have been adapted from C codes in the Cephes library by S. Moshier (<http://www.moshier.net>).

4.1 Factorial

Function **Fact(N)** returns the factorial of the non-negative integer N , also noted $N!$:

$$N! = 1 \times 2 \times \cdots \times N \quad 0! = 1$$

The constant **MaxFac** defines the highest integer for which the factorial can be computed (See chapter 2, p. 11).

4.2 Gamma function

- Function **Gamma(X)** returns the Gamma function, defined by:

$$\Gamma(X) = \int_0^{\infty} t^{X-1} e^{-t} dt$$

This function is related to the factorial by:

$$N! = \Gamma(N + 1)$$

The Gamma function is indefinite for $X = 0$ and for negative integer values of X . It is positive for $X > 0$. For $X < 0$ the Gamma function changes its sign whenever X crosses an integer value. More precisely, if X is an even negative integer, $\Gamma(X)$ is positive on the interval $]X, X+1[$, otherwise it is negative.

- Function **SgnGamma(X)** returns the sign of the Gamma function for a given value of X .
- Function **LnGamma(X)** returns the natural logarithm of the Gamma function.
- Function **Stirling(X)** approximates **Gamma(X)** with Stirling's formula, for $X \geq 30$.
- Function **StirLog(X)** approximates **LnGamma(X)** with Stirling's formula, for $X \geq 13$.

The constants **MaxGam** and **MaxLgm** define the highest values for which the Gamma function and its logarithm, respectively, can be computed (See chapter 2, p. 11).

- Function **IGamma(A, X)** returns the incomplete Gamma function, defined by:

$$\frac{1}{\Gamma(A)} \int_0^X t^{A-1} e^{-t} dt \quad A > 0, X > 0$$

- Function **JGamma(A, X)** returns the complement of the incomplete Gamma function, defined by:

$$\frac{1}{\Gamma(A)} \int_X^\infty t^{A-1} e^{-t} dt$$

Although formally equivalent to $1.0 - \text{IGamma}(A, X)$, this function uses specific algorithms to minimize roundoff errors.

- Function **InvGamma(A, Y)** returns X such that $\text{IGamma}(A, X) = Y$

4.3 Polygamma functions

The polygamma function of order n , denoted $\psi_n(x)$, is the n -th derivative of the logarithm of the gamma function:

$$\psi_n(x) = \frac{d^n}{dx^n} \ln \Gamma(x)$$

The cases $n = 1$ and $n = 2$ are implemented in **FBMath** as **DiGamma(X)** and **TriGamma(X)**

4.4 Beta function

- Function `Beta(X, Y)` returns the Beta function, defined by:

$$\mathcal{B}(X, Y) = \int_0^1 t^{X-1} (1-t)^{Y-1} dt = \frac{\Gamma(X)\Gamma(Y)}{\Gamma(X+Y)}$$

(Here \mathcal{B} denotes the uppercase greek letter ‘Beta’ !)

- Function `IBeta(A, B, X)` returns the incomplete Beta function, defined by:

$$\frac{1}{\mathcal{B}(A, B)} \int_0^X t^{A-1} (1-t)^{B-1} dt \quad A > 0, B > 0, 0 \leq X \leq 1$$

- Function `InvBeta(A, B, Y)` returns X such that `IBeta(A, B, X) = Y`

4.5 Error function

- Function `Erf(X)` returns the error function, defined by:

$$\text{erf}(X) = \frac{2}{\sqrt{\pi}} \int_0^X \exp(-t^2) dt$$

- Function `Erfc(X)` returns the complement of the error function, defined by:

$$\text{erfc}(X) = \frac{2}{\sqrt{\pi}} \int_X^\infty \exp(-t^2) dt$$

4.6 Lambert’s function

Lambert’s W function is the reciprocal of the function xe^x . That is, if $y = W(x)$, then $x = ye^y$. Lambert’s function is defined for $x \geq -1/e$, with $W(-1/e) = -1$. When $-1/e < x < 0$, the function has two values; the value $W(x) > -1$ defines the *upper branch*, the value $W(x) < -1$ defines the *lower branch*.

The function `LambertW(X, UBranch, Offset)` computes Lambert’s function.

- X is the argument of the function (must be $\geq -1/e$)

- **UBranch** is a boolean parameter which must be set to **True** for computing the upper branch of the function and to **False** for computing the lower branch.
- **Offset** is a boolean parameter indicating if **X** is an offset from $-1/e$. In this case, $W(X - 1/e)$ will be computed (with $X > 0$). Using offsets improves the accuracy of the computation if the argument is near $-1/e$.

The default values of **UBranch** and **Offset** are respectively **True** and **False**, so that, for instance, **LambertW(X)** is equivalent to **LambertW(X, True, False)**

The code for Lambert's function has been translated from a Fortran program written by Barry *et al* (<http://www.netlib.org/toms/743>).

4.7 Demo programs

- Programs **testfact.bas**, **testgam.bas**, **testigam.bas**, **testerf.bas**, **testbeta.bas**, **testibet.bas**, located in the **demo\fmath** subdirectory, check the accuracy of the functions **Fact**, **Gamma**, **IGamma**, **Erf**, **Beta** and **IBeta**, respectively.

These programs use reference data from *Numerical Recipes* (<http://www.nr.com>), but the reference values have been re-computed to 20 significant digits with the **Maple** software (<http://www.maplesoft.com>) and the **Gamma** values for negative arguments have been corrected.

Each program computes the values of a given function for a set of predefined arguments and compares the results to the reference values. Then it displays the number of correct digits found. This number should be between 14 and 16 in double precision.

- Program **testw.bas** checks the accuracy of the Lambert function.

The program computes Lambert's function for a set of pre-defined arguments and compares the results with reference values. It displays the number of exact digits found. This number should correspond with the numeric precision used (14-16 digits in double precision).

This program has been translated from a Fortran program written by Barry *et al* (<http://www.netlib.org/toms/743>).

Chapter 5

Probability distributions

This chapter describes the functions available in `FBMath` to compute probability distributions. Most of them are applications of the special functions studied in chapter 4.

5.1 Binomial distribution

Binomial distribution arises when a trial has two possible outcomes: ‘failure’ or ‘success’. If the trial is repeated N times, the random variable X is the number of successes.

- Function `Binomial(N, K)` returns the binomial coefficient $\binom{N}{K}$, which is defined by:

$$\binom{N}{K} = \frac{N!}{K!(N-K)!} \quad 0 \leq K \leq N$$

- Function `PBinom(N, P, K)` returns the probability of obtaining K successes among N repetitions, if the probability of success is P .

$$\text{Prob}(X = K) = \binom{N}{K} P^K Q^{N-K} \quad \text{with } Q = 1 - P$$

- Function `FBinom(N, P, K)` returns the probability of obtaining at most K successes among N repetitions, i. e. $\text{Prob}(X \leq K)$. This is called the *cumulative probability function* and is defined by:

$$\text{Prob}(X \leq K) = \sum_{k=0}^K \binom{N}{k} P^k Q^{N-k} = 1 - I_{\mathcal{B}}(K+1, N-K, P)$$

where $I_{\mathcal{B}}$ denotes the incomplete Beta function.

The mean of the binomial distribution is $\mu = NP$, its variance is $\sigma^2 = NPQ$. The standard deviation is therefore $\sigma = \sqrt{NPQ}$.

5.2 Poisson distribution

The Poisson distribution can be considered as the limit of the binomial distribution when $N \rightarrow \infty$ and $P \rightarrow 0$ while the mean $\mu = NP$ remains small (say $N \geq 30$, $P \leq 0.1$, $NP \leq 10$)

- Function `PPoisson(Mu, K)` returns the probability of observing the value K if the mean is μ . It is defined by:

$$\text{Prob}(X = K) = e^{-\mu} \frac{\mu^K}{K!}$$

- Function `FPoisson(Mu, K)` gives the cumulative probability function, defined by:

$$\text{Prob}(X \leq K) = \sum_{k=0}^K e^{-\mu} \frac{\mu^k}{k!} = 1 - I_{\Gamma}(K+1, \mu)$$

where I_{Γ} denotes the incomplete Gamma function.

5.3 Standard normal distribution

The normal distribution (a. k. a. Gauss distribution or Laplace-Gauss distribution) corresponds to the classical bell-shaped curve. It may also be considered as a limit of the binomial distribution when N is sufficiently ‘large’ while P and Q are sufficiently different from 0 or 1. (say $N \geq 30$, $NP \geq 5$, $NQ \geq 5$).

The normal distribution with mean μ and standard deviation σ is denoted $\mathcal{N}(\mu, \sigma)$ with $\mu = NP$ and $\sigma = \sqrt{NPQ}$. The special case $\mathcal{N}(0, 1)$ is called the standard normal distribution.

- Function `DNorm(X)` returns the probability density of the standard normal distribution, defined by:

$$f(X) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{X^2}{2}\right)$$

The graph of this function is the bell-shaped curve.

- Function **FNorm(X)** returns the cumulative probability function:

$$\Phi(X) = \text{Prob}(U \leq X) = \int_{-\infty}^X f(x)dx = \frac{1}{2} \left[1 + \text{erf} \left(X \frac{\sqrt{2}}{2} \right) \right]$$

where U denotes the standard normal variable and erf the error function.

- Function **PNorm(X)** returns the probability that the standard normal variable exceeds X in absolute value, i. e. $\text{Prob}(|U| > X)$.
- Function **InvNorm(P)** returns the value X such that $\Phi(X) = P$.

5.4 Student's distribution

Student's distribution is widely used in Statistics, for instance to estimate the mean of a population from a sample taken from this population. The distribution depends on an integer parameter ν called the *number of degrees of freedom* (in the mean estimation problem, $\nu = n - 1$ where n is the number of individuals in the sample). When ν is large (say > 30) the Student distribution is approximately equal to the standard normal distribution.

- Function **DStudent(Nu, X)** returns the probability density of the Student distribution with **Nu** degrees of freedom, defined by:

$$f_{\nu}(X) = \frac{1}{\nu^{1/2} \mathcal{B}\left(\frac{\nu}{2}, \frac{1}{2}\right)} \cdot \left(1 + \frac{X^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

where \mathcal{B} denotes the Beta function.

- Function **FStudent(Nu, X)** returns the cumulative probability function:

$$\Phi_{\nu}(X) = \text{Prob}(t \leq X) = \int_{-\infty}^X f_{\nu}(x)dx = \begin{cases} I/2 & \text{if } X \leq 0 \\ 1 - I/2 & \text{if } X > 0 \end{cases}$$

where t denotes the Student variable and $I = I_{\mathcal{B}}\left(\frac{\nu}{2}, \frac{1}{2}, \frac{\nu}{\nu+X^2}\right)$

- Function **PStudent(Nu, X)** returns the probability that the Student variable t exceeds X in absolute value, i. e. $\text{Prob}(|t| > X)$.
- Function **InvStudent(Nu, P)** returns the value X such that $\Phi_{\nu}(X) = P$.

5.5 Khi-2 distribution

The χ^2 distribution is a special case of the Gamma distribution (see below). It depends on an integer parameter ν which is the number of degrees of freedom.

- Function `DKhi2(Nu, X)` returns the probability density of the χ^2 distribution with `Nu` degrees of freedom, defined by:

$$f_\nu(X) = \frac{1}{2^{\frac{\nu}{2}} \Gamma\left(\frac{\nu}{2}\right)} \cdot X^{\frac{\nu}{2}-1} \cdot \exp\left(-\frac{X}{2}\right) \quad (X > 0)$$

- Function `FKhi2(Nu, X)` returns the cumulative probability function:

$$\Phi_\nu(X) = \text{Prob}(\chi^2 \leq X) = \int_0^X f_\nu(x) dx = I_\Gamma\left(\frac{\nu}{2}, \frac{X}{2}\right)$$

where I_Γ denotes the incomplete Gamma function.

- Function `PKhi2(Nu, X)` returns the probability that the χ^2 variable exceeds X , i. e. $\text{Prob}(\chi^2 > X)$.
- Function `InvKhi2(Nu, P)` returns the value X such that $\Phi_\nu(X) = P$.

5.6 Snedecor's distribution

The Snedecor (or Fisher-Snedecor) distribution is used mainly to compare two variances. It depends on two integer parameters ν_1 and ν_2 which are the degrees of freedom associated with the variances.

- Function `DSnedecor(Nu1, Nu2, X)` returns the probability density of the Snedecor distribution with `Nu1` and `Nu2` degrees of freedom, defined by:

$$f_{\nu_1, \nu_2}(X) = \frac{1}{\mathcal{B}\left(\frac{\nu_1}{2}, \frac{\nu_2}{2}\right)} \cdot \left(\frac{\nu_1}{\nu_2}\right)^{\frac{\nu_1}{2}} \cdot X^{\frac{\nu_1}{2}-1} \cdot \left(1 + \frac{\nu_1}{\nu_2} X\right)^{-\frac{\nu_1+\nu_2}{2}} \quad (X > 0)$$

- Function `FSnedecor(Nu1, Nu2, X)` returns the cumulative probability function:

$$\Phi_{\nu_1, \nu_2}(X) = \text{Prob}(F \leq X) = \int_0^X f_{\nu_1, \nu_2}(x) dx = 1 - I_B\left(\frac{\nu_2}{2}, \frac{\nu_1}{2}, \frac{\nu_2}{\nu_2 + \nu_1 X}\right)$$

where F denotes the Snedecor variable.

- Function `PSnedecor(Nu1, Nu2, X)` returns the probability that the Snedecor variable F exceeds X , i. e. $\text{Prob}(F > X)$.
- Function `InvSnedecor(Nu1, Nu2, P)` returns the value X such that $\Phi_{\nu_1, \nu_2}(X) = P$.

5.7 Exponential distribution

The exponential distribution is used in many applications (radioactivity, chemical kinetics...). It depends on a positive real parameter A .

- Function `DExpo(A, X)` returns the probability density of the exponential distribution with parameter A , defined by:

$$f_A(X) = A \exp(-AX) \quad (X > 0)$$

- Function `FExpo(A, X)` returns the cumulative probability function:

$$\Phi_A(X) = \int_0^X f_A(x) dx = 1 - \exp(-AX)$$

5.8 Beta distribution

The Beta distribution is often used to describe the distribution of a random variable defined on the unit interval $[0, 1]$. It depends on two positive real parameters A and B .

- Function `DBeta(A, B, X)` returns the probability density of the Beta distribution with parameters A and B , defined by:

$$f_{A,B}(X) = \frac{1}{\mathcal{B}(A, B)} \cdot X^{A-1} \cdot (1 - X)^{B-1} \quad (0 \leq X \leq 1)$$

- Function `FBeta(A, B, X)` returns the cumulative probability function:

$$\Phi_{A,B}(X) = \int_0^X f_{A,B}(x) dx = I_B(A, B, X)$$

- Function `InvBeta(A, B, P)` returns the value X such that $\Phi_{A,B}(X) = P$.

5.9 Gamma distribution

The Gamma distribution is often used to describe the distribution of a random variable defined on the positive real axis. It depends on two positive real parameters A and B .

- Function `DGamma(A, B, X)` returns the probability density of the Gamma distribution with parameters A and B , defined by:

$$f_{A,B}(X) = \frac{B^A}{\Gamma(A)} \cdot X^{A-1} \cdot \exp(-BX) \quad (X > 0)$$

- Function `FGamma(A, B, X)` returns the cumulative probability function:

$$\Phi_{A,B}(X) = \int_0^X f_{A,B}(x) dx = I_\Gamma(A, BX)$$

The χ^2 distribution is a special case of the Gamma distribution, with $A = \nu/2$ and $B = 1/2$.

5.10 Demo programs

- Program `binom1.bas` computes the binomial coefficient `Binomial(N, K)` for several values of N and K and compares the results with reference values.
- Program `binom2.bas` compares the cumulative probabilities of the binomial distribution, estimated by function `FBinom`, with the values obtained by summing up the individual probabilities.

Chapter 6

Matrices and linear equations

This chapter describes the procedures and functions available in `FBMath` to perform vector and matrix operations, and to solve systems of linear equations.

6.1 Programming conventions

All subroutines dealing with arrays in `FBMath` apply the following conventions:

- Unless otherwise noted, the `FBMath` subroutines use arrays of type `DOUBLE`. For instance:

```
DIM AS DOUBLE V(1 TO 10)           ' For a vector
DIM AS DOUBLE A(1 TO 10, 1 TO 10)  ' For a matrix
```

- Arrays may be static or dynamic.
- The subroutines do not allocate the arrays present in their parameter lists. These allocations must therefore be performed by the main program, by means of the appropriate `DIM` statements, before calling the subroutines.
- The array dimensions are not passed to the subroutines; instead they are obtained inside the subroutines by means of the `LBOUND` and `UBOUND` functions. So, any array which appear in the parameter list must have its exact dimensions before calling the subroutine.

6.2 Error codes

The following error codes are defined in the include file `math.bi`:

Error code	Value	Meaning
<code>MatOk</code>	0	No error
<code>MatNonConv</code>	-1	Non-convergence
<code>MatSing</code>	-2	Quasi-singular matrix
<code>MatErrDim</code>	-3	Non-compatible dimensions
<code>MatNotPD</code>	-4	Matrix not positive-definite

6.3 Gauss-Jordan elimination

If $\mathbf{C}(n \times n)$ and $\mathbf{B}(n \times m)$ are two real matrices, the Gauss-Jordan elimination can compute the inverse matrix \mathbf{C}^{-1} , the solution \mathbf{X} to the system of linear equations $\mathbf{CX} = \mathbf{B}$, and the determinant of \mathbf{C} .

This procedure is implemented as subroutine `GaussJordan(A(), Det)` where:

- On input, **A** is the global matrix $[\mathbf{C}|\mathbf{B}]$, which means that:
 - the first n columns of **A** contain the matrix \mathbf{C}
 - the other columns of **A** contain the matrix \mathbf{B}
- On output, **A** is transformed into the global matrix $[\mathbf{C}^{-1}|\mathbf{X}]$, which means that:
 - the first n columns of **A** contain the inverse matrix \mathbf{C}^{-1}
 - the other columns of **A** contain the solution matrix \mathbf{X}
- **Det** is the determinant of \mathbf{C}

Notes:

- **B** may be a vector, in this case $m = 1$ and \mathbf{X} is also a vector.
- The original matrix **A** is overwritten by the subroutine. If necessary, the calling program must save a copy of it.

Subroutine `LinEq(A(), B(), Det)` is a simplified version of the Gauss-Jordan procedure where **A** is a square matrix and **B** is a vector. This subroutine solves the system $\mathbf{AX} = \mathbf{B}$. On output, **A** contains the inverse matrix and **B** contains the solution vector.

After a call to `GaussJordan` or `LinEq`, function `MathErr` will return the error code:

- `MatOk` if there is no error.
- `MatSing` if **A** is singular (or quasi-singular)
- `MatErrDim` if the matrices have incompatible dimensions

6.4 LU decomposition

The LU decomposition algorithm factors the square matrix **A** as a product **LU**, where **L** is a lower triangular matrix (with unit diagonal terms) and **U** is an upper triangular matrix.

The linear system $\mathbf{AX} = \mathbf{B}$ is then solved by:

$$\mathbf{LY} = \mathbf{B} \quad (6.1)$$

$$\mathbf{UX} = \mathbf{Y} \quad (6.2)$$

System 6.1 is solved for vector **Y**, then system 6.2 is solved for vector **X**. The solutions are simplified by the triangular nature of the matrices.

`FBMath` provides the following subroutines:

- subroutine `LU_Decom(A())` performs the LU decomposition of matrix **A**.

The matrices **L** and **U** are stored in **A**, which is therefore destroyed.

After a call to `LU_Decom`, the function `MathErr` will return one of the following error codes:

- `MatOk` if no error
- `MatErrDim` if **A** is not square or if its lower bounds in the two dimensions are different (i. e. `LBOUND(A, 1) <> LBOUND(A, 2)`)
- `MatSing` if **A** is quasi-singular
- subroutine `LU_Solve(A(), B(), X())` solves the system $\mathbf{AX} = \mathbf{B}$, where **X** and **B** are real vectors, once the matrix **A** has been transformed by `LU_Decom`.

6.5 QR decomposition

This method factors a matrix \mathbf{A} as a product of an orthogonal matrix \mathbf{Q} by an upper triangular matrix \mathbf{R} :

$$\mathbf{A} = \mathbf{QR}$$

The linear system $\mathbf{AX} = \mathbf{B}$ then becomes:

$$\mathbf{QRX} = \mathbf{B}$$

Denoting the transpose of \mathbf{Q} by \mathbf{Q}' and left-multiplying by this transpose, one obtains:

$$\mathbf{Q}'\mathbf{QRX} = \mathbf{Q}'\mathbf{B}$$

or:

$$\mathbf{RX} = \mathbf{Q}'\mathbf{B}$$

since the transpose of an orthogonal matrix is equal to its inverse.

The last system is solved by making advantage of the triangular nature of matrix \mathbf{R} .

Note : The QR decomposition may be applied to a rectangular matrix $n \times m$ (with $n > m$). In this case, \mathbf{Q} has dimensions $n \times m$ and \mathbf{R} has dimensions $m \times m$. For a linear system $\mathbf{AX} = \mathbf{B}$, the solution minimizes the norm of the vector $\mathbf{AX} - \mathbf{B}$. It is called the *least squares* solution.

FBMath provides the following subroutines:

- subroutine `QR_Decomp(A(), R())` performs the QR decomposition on the input matrix \mathbf{A} .

The matrix \mathbf{Q} is stored in \mathbf{A} , which is therefore destroyed.

After a call to `QR_Decomp`, the function `MathErr` will return one of the following error codes:

- `MatOk` if no error
 - `MatErrDim` if $n > m$ or if \mathbf{A} has different lower bounds in the two dimensions (i. e. `LBOUND(A, 1) <> LBOUND(A, 2)`)
 - `MatSing` if \mathbf{A} is quasi-singular
- subroutine `QR_Solve(Q(), R(), B(), X())` solves the system $\mathbf{QRX} = \mathbf{B}$.

6.6 Singular value decomposition

Singular value decomposition (SVD) factors a matrix \mathbf{A} as a product:

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}'$$

where \mathbf{U} et \mathbf{V} are orthogonal matrices. \mathbf{S} is a diagonal matrix. Its diagonal terms S_{ii} are all ≥ 0 and are called the *singular values* of \mathbf{A} . The *rank* of \mathbf{A} is equal to the number of non-null singular values.

- If \mathbf{A} is a regular matrix, all S_{ii} are > 0 . The inverse matrix is given by:

$$\mathbf{A}^{-1} = (\mathbf{U}\mathbf{S}\mathbf{V}')^{-1} = (\mathbf{V}')^{-1}\mathbf{S}^{-1}\mathbf{U}^{-1} = \mathbf{V} \times \text{diag}(1/S_{ii}) \times \mathbf{U}'$$

since the inverse of an orthogonal matrix is equal to its transpose.

So the solution of the system $\mathbf{A}\mathbf{X} = \mathbf{B}$ is given by $\mathbf{X} = \mathbf{A}^{-1}\mathbf{B}$

- If \mathbf{A} is a singular matrix, some S_{ii} are null. However, the previous expressions remain valid provided that, for each null singular value, the term $1/S_{ii}$ is replaced by zero.

It may be shown that the solution so calculated corresponds:

- in the case of an under-determined system, to the vector \mathbf{X} having the least norm.
- in the case of an impossible system, to the least-squares solution.

Note : Just like the QR decomposition, the SVD may be applied to a rectangular matrix $n \times m$ (with $n > m$). In this case, \mathbf{U} has dimensions $n \times m$, \mathbf{S} and \mathbf{V} have dimensions $m \times m$. For a linear system $\mathbf{A}\mathbf{X} = \mathbf{B}$, the SVD method gives the least squares solution.

FBMath provides the following subroutines:

- subroutine `SV_Decomp(A(), S(), V())` performs the singular value decomposition on the input matrix \mathbf{A} .

The matrix \mathbf{U} (such that $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}'$) is stored in \mathbf{A} , which is therefore destroyed.

After a call to `SV_Decomp`, the function `MathErr` will return one of the following error codes:

- **MatOk** if no error
 - **MatErrDim** if $n > m$ or if **A** has different lower bounds in the two dimensions (i. e. $\text{LBOUND}(\mathbf{A}, 1) \neq \text{LBOUND}(\mathbf{A}, 2)$)
 - **MatNonConv** if the iterative process does not converge
- procedure **SV_SetZero**(**S**(), **Tol**) sets to zero the singular values S_i which are lower than a fraction **Tol** of the highest singular value. This procedure may be used when solving a system with a near-singular matrix.
 - procedure **SV_Solve**(**U**(), **S**(), **V**(), **B**(), **X**()) solves the system $\mathbf{USV}'\mathbf{X} = \mathbf{B}$.
 - procedure **SV_Approx**(**U**(), **S**(), **V**(), **A**()) approximates a matrix **A** by the product \mathbf{USV}' , after the lowest singular values have been set to zero by **SV_SetZero**.

6.7 Cholesky decomposition

The symmetric matrix **A** is said to be *positive definite* if, for any vector **x**, the product $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is positive.

For such matrices, it is possible to find a lower triangular matrix **L** such that:

$$\mathbf{A} = \mathbf{L}\mathbf{L}^\top$$

L can be viewed as a kind of ‘square root’ of **A**.

Subroutine **Cholesky**(**A**(), **L**()) performs the Cholesky decomposition on **A**. After a call to the subroutine, function **MathErr** returns the error code:

- **MatOk** if there is no error.
- **MatNotPD** if **A** is not positive definite.

6.8 Eigenvalues and eigenvectors

6.8.1 Definitions

A square matrix **A** is said to have an eigenvalue λ , associated to an eigenvector **V**, if and only if:

$$\mathbf{A} \cdot \mathbf{V} = \lambda \cdot \mathbf{V}$$

A symmetric matrix of size n has n distinct real eigenvalues and n orthogonal eigenvectors.

A non-symmetric matrix of size n has also n eigenvalues but some of them may be complex, and some may be equal (they are said to be degenerate).

6.8.2 Symmetric matrices

Subroutine `Jacobi(A(), MaxIter, Tol, V(), Lambda())` computes the eigenvalues and eigenvectors of the real symmetric matrix `A`, using the iterative method of Jacobi.

`MaxIter` is the maximum number of iterations, `Tol` is the required precision on the eigenvalues.

The eigenvectors are returned in matrix `V`; the eigenvalues are returned in vector `Lambda`.

The eigenvectors are stored along the columns of `V()`. They are normalized, with their first component always positive.

After a call to `Jacobi`, function `MathErr` returns one of two error codes:

- `MatOk` if all goes well.
- `MatNonConv` if the iterative process does not converge.

This procedure destroys the original matrix `A`.

6.8.3 General square matrices

- Subroutine `EigenVals(A(), Lambda())` computes the eigenvalues of the real square matrix `A`.

Eigenvalues are stored in the complex vector `Lambda`. The real and imaginary parts of the i^{th} eigenvalue are stored in `Lambda(i).X` and `Lambda(i).Y`, respectively. The eigenvalues are unordered, except that complex conjugate pairs appear consecutively with the value having the positive imaginary part first.

Function `MathErr` returns the following error codes:

- 0 if no error
- (-i) if an error occurred during the determination of the i^{th} eigenvalue. The eigenvalues should be correct for the indices $> i$.

This procedure destroys the original matrix **A**.

- Subroutine **EigenVect(A(), Lambda(), V())** computes the eigenvalues and eigenvectors of the real square matrix **A**.

Eigenvalues are stored in the *complex* vector **Lambda**, just like with **EigenVals**.

Eigenvectors are stored along the columns of the *real* matrix **V**.

If the i^{th} eigenvalue is real, the i^{th} column of **V** contains its eigenvector. If the i^{th} eigenvalue is complex with positive imaginary part, the i^{th} and $(i+1)^{th}$ columns of **V** contain the real and imaginary parts of its eigenvector. The eigenvectors are unnormalized.

Function **MathErr** returns the same error codes than **EigenVals**. If the error code is not null, none of the eigenvectors has been found.

This procedure destroys the original matrix **A**.

6.9 Demo programs

These programs are located in the **demo\matrices** subdirectory.

6.9.1 Determinant and inverse of a square matrix

Program **detinv.bas** computes the determinant and inverse of a square matrix. The inverse matrix is re-inverted and the result (which should be equal to the original matrix) is printed.

The example matrix is:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 0 & -1 \\ -1 & 4 & 3 & -0.5 \\ 2 & 2 & 1 & -3 \\ 0 & 0 & 3 & -4 \end{bmatrix}$$

The inverse is:

$$\mathbf{A}^{-1} = \begin{bmatrix} -\frac{41}{21} & \frac{4}{21} & \frac{11}{7} & -\frac{5}{7} \\ \frac{16}{21} & \frac{1}{21} & -\frac{5}{14} & \frac{1}{14} \\ -\frac{40}{21} & \frac{8}{21} & \frac{8}{7} & -\frac{3}{7} \\ -\frac{10}{7} & \frac{2}{7} & \frac{6}{7} & -\frac{4}{7} \end{bmatrix}$$

or, in approximate form:

$$\mathbf{A}^{-1} \approx \begin{bmatrix} -1.9523 & 0.1905 & 1.5714 & -0.7143 \\ 0.7619 & 0.0476 & -0.3571 & 0.0714 \\ -1.9048 & 0.3810 & 1.1429 & -0.4286 \\ -1.4286 & 0.2857 & 0.8571 & -0.5714 \end{bmatrix}$$

The determinant is -21.

6.9.2 Hilbert matrices

Program `hilbert.bas` tests the Gauss-Jordan method by solving a series of Hilbert systems of increasing order. Such systems have matrices of the form:

$$\mathbf{A} = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \cdots & \frac{1}{N} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \cdots & \frac{1}{N+1} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \cdots & \frac{1}{N+2} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \cdots & \frac{1}{N+3} \\ \vdots & & & & \ddots & \\ \frac{1}{N} & \frac{1}{N+1} & \frac{1}{N+2} & \frac{1}{N+3} & \cdots & \frac{1}{2N-1} \end{bmatrix}$$

Each element of the constant vector is equal to the sum of the terms in the corresponding line of the matrix :

$$B_i = \sum_{j=1}^N A_{ij}$$

The solution of such a system is $[1, 1, 1, \dots, 1]$

The determinant of the Hilbert matrix tends towards zero when the order increases. The program stops when the determinant becomes too low with respect to the numerical precision of the floating point numbers. This occurs at order 13 in double precision.

6.9.3 Gauss-Jordan method: single constant vector

Program `lineq1.bas` solves the linear system $\mathbf{AX} = \mathbf{B}$. After a call to `LinEq`, \mathbf{A} contains the inverse matrix and \mathbf{B} contains the solution vector.

The example system matrix is:

$$\mathbf{A} = \begin{bmatrix} 2 & 1 & 5 & -8 \\ 7 & 6 & 2 & 2 \\ -1 & -3 & -10 & 4 \\ 2 & 2 & 2 & 1 \end{bmatrix}$$

The constant vector is:

$$\mathbf{B} = \begin{bmatrix} 0 \\ 17 \\ -10 \\ 7 \end{bmatrix}$$

The solution vector is:

$$\mathbf{X} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

The determinant is -135

6.9.4 Gauss-Jordan method: multiple constant vectors

Program `lineqm.bas` solves a series of linear systems with the same system matrix and several constant vectors. The system matrix is stored in the first n columns of matrix **A**; the constant vectors are stored in the following columns. After a call to `GaussJordan`, the first n columns of **A** contain the inverse matrix and the following columns contain the solution vectors.

The example system matrix from the previous program is used. The matrix of constant vectors is:

$$\begin{bmatrix} 0 & -15 & 14 & -13 & 5 \\ 17 & 50 & 1 & 84 & 30 \\ -10 & -5 & -12 & -51 & -15 \\ 7 & 17 & 1 & 37 & 10 \end{bmatrix}$$

The solution matrix is:

$$\begin{bmatrix} 1 & 2 & 1 & 4 & 0 \\ 1 & 5 & -1 & 5 & 5 \\ 1 & 0 & 1 & 6 & 0 \\ 1 & 3 & -1 & 7 & 0 \end{bmatrix}$$

6.9.5 LU, QR and SV decompositions

The demo programs `test_lu.bas`, `test_qr.bas` and `test_svd.bas` solve the linear system used by `lineq1.bas` (paragraph 6.9.3) with the LU, QR, and singular value decompositions, respectively.

6.9.6 Cholesky decomposition

Program `cholesk.bas` performs the Cholesky decomposition of a positive definite symmetric matrix. The matrix is decomposed then the program computes the product $\mathbf{L}\mathbf{L}^\top$ which must give the original matrix.

The example matrix is:

$$\mathbf{A} = \begin{bmatrix} 60 & 30 & 20 \\ 30 & 20 & 15 \\ 20 & 15 & 12 \end{bmatrix}$$

Its Cholesky factor is:

$$\mathbf{L} = \begin{bmatrix} 2\sqrt{15} & 0 & 0 \\ \sqrt{15} & \sqrt{5} & 0 \\ \frac{2}{3}\sqrt{15} & \sqrt{5} & \frac{1}{3}\sqrt{3} \end{bmatrix}$$

or, in approximate form:

$$\mathbf{L} \approx \begin{bmatrix} 7.745967 & 0 & 0 \\ 3.872983 & 2.236068 & 0 \\ 2.581989 & 2.236068 & 0.577350 \end{bmatrix}$$

6.9.7 Eigenvalues of a symmetric matrix

Program `eigensym.bas` computes the eigenvalues and eigenvectors of Hilbert matrices (see program `hilbert.bas`) by the method of Jacobi. Such matrices are very ill-conditioned, which can be seen from the high ratio between the highest and lowest eigenvalues (the *condition number*).

6.9.8 Eigenvalues of a general square matrix

Program `eigenval.bas` computes the eigenvalues of a general square matrix.

The example matrix from the `detinv.bas` program is used. It has two real and two complex (conjugate) eigenvalues:

```

-1.075319 +      1.709050 * i
-1.075319 -      1.709050 * i
-1.000000
5.150639

```

6.9.9 Eigenvalues and eigenvectors of a general square matrix

Program `eigenvec.bas` computes both the eigenvalues and eigenvectors of a general square matrix. The same example matrix is used.

The eigenvectors are stored columnwise in a matrix `V(1..N, 1..N)`. In order to retrieve the eigenvectors associated with complex eigenvalues, the program takes into account the following properties:

- Complex conjugate pairs of eigenvalues are stored consecutively in vector `Lambda`, with the value having the positive imaginary part first.
- If the i^{th} eigenvalue is complex with positive imaginary part, the i^{th} and $(i+1)^{th}$ columns of matrix `V` contain the real and imaginary parts of its eigenvector.
- Eigenvectors associated with complex conjugate eigenvalues are themselves complex conjugate.

Hence the algorithm:

```

IF Lambda(I).Y = 0 THEN
  ' Eigenvector is in column I of V
ELSEIF Lambda(I).Y > 0 THEN
  ' Real and imag. parts of eigenvector are in columns I and (I+1)
  ' For component K: real part = V(K,I), imag. part = V(K,I+1)
ELSE
  ' Real and imag. parts of eigenvector are in columns (I-1) and I
  ' For component K: real part = V(K,I-1), imag. part = - V(K,I)
END IF

```

The results obtained with the example matrix are the following:

Eigenvalue:

```

-1.075319 +      1.709050 * i

```


Eigenvector:

$$\begin{array}{rcl} -0.220224 & + & 0.394848 * i \\ 0.078289 & - & 0.303345 * i \\ 0.029348 & + & 0.787594 * i \\ 0.374358 & + & 0.589119 * i \end{array}$$

Eigenvalue:

$$-1.075319 - 1.709050 * i$$

Eigenvector:

$$\begin{array}{rcl} -0.220224 & - & 0.394848 * i \\ 0.078289 & + & 0.303345 * i \\ 0.029348 & - & 0.787594 * i \\ 0.374358 & - & 0.589119 * i \end{array}$$

Eigenvalue:

$$-1.000000$$

Eigenvector:

$$\begin{array}{r} 2.605054 \\ -1.042021 \\ 3.126065 \\ 3.126065 \end{array}$$

Eigenvalue:

$$5.150638$$

Eigenvector:

$$\begin{array}{r} 0.345194 \\ 0.788801 \\ 0.441744 \\ 0.144823 \end{array}$$

Chapter 7

Function minimization

This chapter describes the procedures and functions available in **FBMath** to minimize functions of one or several variables. Only deterministic optimizers are considered here. Stochastic optimization will be studied in a next chapter.

7.1 Functions of one variable

Let **Func** be a function of a real variable **X**. In **FBMath** such a function is declared as:

```
FUNCTION Func(X AS DOUBLE) AS DOUBLE
```

The problem is to find the real **Xmin** for which **Func(X)** is minimal.

Subroutine **GoldSearch(Func, A, B, MaxIter, Tol, Xmin, Ymin)** performs the minimization by the ‘golden search’ method. This means that, at each iteration, the number **Xmin** is ‘bracketed’ by a triplet (**A**, **B**, **C**) such that:

- $A < B < C$
- A, B, C are within the golden mean ϕ , i.e.

$$\frac{B - A}{C - B} = \frac{C - A}{B - A} = \phi = \frac{1 + \sqrt{5}}{2} \approx 1.618$$

- $\text{Func}(B) < \text{Func}(A)$ and $\text{Func}(B) < \text{Func}(C)$.

The user must provide two numbers **A** and **B** which define the ‘unit vector’ on the **X** axis. The number **C** is found by the program itself. It is not necessary that the interval [**A**, **B**] contains the minimum.

The user must also provide:

- the maximum number of iterations **MaxIter**
- the tolerance **Tol** with which the minimum must be located. This value should not be higher than the square root of the machine precision ($\text{MachEp}^{1/2} \approx 1.5 \times 10^{-8}$ in double precision)

The subroutine returns the coordinates (**Xmin**, **Ymin**) of the minimum.

After a call to **GoldSearch**, function **MathErr()** will return one of two error codes:

- **OptOk** if no error occurred
- **OptNonConv** (non-convergence) if the number of iterations exceeds the maximum value **MaxIter**

The determination of the bracketing triplet **A**, **B**, **C** is performed within **GoldSearch** by a call to a subroutine **MinBrack**. This subroutine may be called independently. Its syntax is:

```
MinBrack(Func, A, B, C, Fa, Fb, Fc)
```

The user must provide the first two numbers **A** and **B**. The number **C** is found by the subroutine. The corresponding values of the function are returned in **Fa**, **Fb**, **Fc**.

7.2 Functions of several variables

Let f be a function of a real vector \mathbf{x} such that $\mathbf{x} = [x_1, x_2, \dots]$. In **FBMath** such a function is declared as:

```
FUNCTION Func(X() AS DOUBLE) AS DOUBLE
```

The problem is to find the vector **X()** for which **Func(X())** is minimal.

7.2.1 Minimization along a line

If \mathbf{x}^0 is a starting point and $\delta\mathbf{x}$ is a constant vector, minimizing f from \mathbf{x}^0 along the direction specified by $\delta\mathbf{x}$ is equivalent to finding the number r such that $g(r) = f(\mathbf{x}^0 + r \cdot \delta\mathbf{x})$ is minimal.

Subroutine **LinMin(Func, X(), DeltaX(), R, MaxIter, Tol, Fmin)** will minimize function **Func** from **X()** in the direction specified by **DeltaX()**. **R**

is the initial step in that direction, expressed as a fraction of the norm of `DeltaX()`. If `R` is set to 0 or a negative value, the subroutine will use the default value `R = 1`. The user must also provide the maximum number of iterations `MaxIter` and the tolerance `Tol`, as for subroutine `GoldSearch`.

On output, `LinMin` returns:

- the coordinates of the minimum in `X()`
- the step corresponding to the minimum in `R`
- the function value at the minimum in `Fmin`

After a call to `LinMin`, function `MathErr()` will return one of the error codes `OptOk` or `OptNonConv`, as with `GoldSearch`.

7.2.2 Newton-Raphson method

The Newton-Raphson method starts with an approximation \mathbf{x}^0 for the coordinates of the minimum and generates a new approximation \mathbf{x} by using the second-order Taylor series expansion of function f around \mathbf{x}^0 :

$$f(\mathbf{x}) = f(\mathbf{x}^0) + (\mathbf{x} - \mathbf{x}^0)^\top \cdot \mathbf{g}(\mathbf{x}^0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^0)^\top \cdot \mathbf{H}(\mathbf{x}^0) \cdot (\mathbf{x} - \mathbf{x}^0) \quad (7.1)$$

\mathbf{g} denotes the gradient vector (vector of first partial derivatives) and \mathbf{H} denotes the hessian matrix (matrix of second partial derivatives). For instance, for a function of two variables $f(x_1, x_2)$:

$$\mathbf{g}(\mathbf{x}^0) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x_1^0, x_2^0) \\ \frac{\partial f}{\partial x_2}(x_1^0, x_2^0) \end{bmatrix}$$

$$\mathbf{H}(\mathbf{x}^0) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x_1^0, x_2^0) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x_1^0, x_2^0) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x_1^0, x_2^0) & \frac{\partial^2 f}{\partial x_2^2}(x_1^0, x_2^0) \end{bmatrix}$$

By differentiating eq. (1) we obtain the gradient of f at point \mathbf{x} :

$$\mathbf{g}(\mathbf{x}) = \mathbf{g}(\mathbf{x}^0) + \mathbf{H}(\mathbf{x}^0) \cdot (\mathbf{x} - \mathbf{x}^0) \quad (7.2)$$

If \mathbf{x} is sufficiently close to the minimum, $\mathbf{g}(\mathbf{x}) \approx 0$ so:

$$\mathbf{x} = \mathbf{x}^0 - \mathbf{H}^{-1}(\mathbf{x}^0) \cdot \mathbf{g}(\mathbf{x}^0)$$

In practice, it is better to determine the step k which minimizes the function in the direction specified by $-\mathbf{H}^{-1}(\mathbf{x}^0) \cdot \mathbf{g}(\mathbf{x}^0)$:

$$\mathbf{x} = \mathbf{x}^0 - k \cdot \mathbf{H}^{-1}(\mathbf{x}^0) \cdot \mathbf{g}(\mathbf{x}^0)$$

The determination of k is performed by line minimization.

Subroutine `Newton(Func, HessGrad, X(), MaxIter, Tol, Fmin, G(), Hinv(), Det)` minimizes function `Func` by the Newton-Raphson method.

The user must provide a subroutine `HessGrad` to compute the gradient `G()` and the hessian `H()` of the function at point `X()`. This subroutine is declared as:

```
SUB HessGrad(X() AS DOUBLE, G() AS DOUBLE, H() AS DOUBLE)
```

`MaxIter` and `Tol` have their usual meaning.

On output, `Newton` returns:

- the coordinates of the minimum in `X()`
- the function value at the minimum in `Fmin`
- the gradient at the minimum in `G()` (should be near 0)
- the inverse hessian matrix at the minimum in `H()`
- the determinant of the hessian matrix at the minimum in `Det`

After a call to `Newton`, function `MathErr()` will return one of three error codes:

- `OptOk` if no error occurred
- `OptNonConv` (non-convergence) if the number of iterations exceeds the maximum value `MaxIter`
- `OptSing` if the hessian matrix is quasi-singular

7.2.3 Approximate gradient and hessian

Although it is recommended to compute the gradient and hessian from analytical derivatives, approximate values may be found using finite difference approximations:

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) \approx \frac{f(x_i + h_i) - f(x_i - h_i)}{2h_i}$$

$$\frac{\partial^2 f}{\partial x_i^2}(\mathbf{x}) \approx \frac{f(x_i + h_i) + f(x_i - h_i) - 2f(x_i)}{h_i^2}$$

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) \approx \frac{f(x_i + h_i, x_j + h_j) - f(x_i + h_i, x_j) - f(x_i, x_j + h_j) + f(x_i, x_j)}{h_i h_j}$$

The increment h_i is such that $h_i = \eta |x_i|$ where η is a constant which should not be less than the cube root of the machine epsilon ($\text{MachEp}^{1/3} \approx 6.06 \times 10^{-6}$ in double precision).

Subroutine `NumHessGrad(Func, X(), Eta, G(), H())` performs these computations for function `Func` at point `X()` using the relative increment `Eta`. The approximate gradient and hessian are returned in `G()` and `H()`.

To use `Newton` with numerical derivatives, you must imbed `NumHessGrad` into another subroutine (e. g. `HessGrad`) which will be passed to `Newton`, for instance:

```
SUB HessGrad (X() AS DOUBLE, G() AS DOUBLE, H() AS DOUBLE)
  CONST Eta = 1D-6
  NumHessGrad @Func, X(), Eta, G(), H()
END SUB
```

where function `Func` is defined in the main program (see demo program `testnewt.bas` for an example).

It is not possible to pass `NumHessGrad` directly to `Newton` because the parameter list of `NumHessGrad` do not match the parameter list of `HessGrad` as defined in `Newton`'s declaration.

7.2.4 Marquardt method

This method is a variant of the Newton-Raphson method, in which each diagonal term of the hessian matrix is multiplied by a scalar equal to $(1 + \lambda)$, where λ is the Marquardt parameter. This parameter is initialized at some

small value (e.g. 10^{-2}) at the beginning of the iterations, then it is decreased by a factor 10 if the iteration leads to a decrease of the function, otherwise it is increased by a factor 10. This procedure usually improves the convergence of the Newton-Raphson method.

If the method converges, λ should reach a very small value, so that the Marquardt and Newton-Raphson algorithms should produce identical results for the inverse hessian matrix. However, this is not guaranteed, so that, if a precise inverse hessian is required, it may be useful to perform a single iteration of the Newton-Raphson method once Marquardt's algorithm has successfully terminated (see demo program `testmarq.bas`).

This procedure is implemented as: `Marquardt(Func, HessGrad X(), MaxIter, Tol, Fmin, G(), Hinv(), Det)`

It is used like `Newton`, except that an additional error code, `OptBigLambda`, may be returned by `MathErr` if the Marquard parameter increases beyond a predefined value (10^3 in this implementation).

7.2.5 BFGS method

The BFGS (Broyden-Fletcher-Goldfarb-Shanno) method is another variant of the Newton method in which the hessian matrix does not need to be computed explicitly. It is said a *quasi-Newton* method.

The BFGS algorithm uses the following formula to construct the inverse hessian matrix iteratively:

$$\mathbf{H}_{i+1}^{-1} = \mathbf{H}_i^{-1} + \frac{\delta \mathbf{x} \cdot \delta \mathbf{x}^\top}{\delta \mathbf{x}^\top \cdot \delta \mathbf{g}} - \frac{(\mathbf{H}_i^{-1} \cdot \delta \mathbf{g}) \cdot (\mathbf{H}_i^{-1} \cdot \delta \mathbf{g})^\top}{\delta \mathbf{g}^\top \cdot \mathbf{H}_i^{-1} \cdot \delta \mathbf{g}} + (\delta \mathbf{g}^\top \cdot \mathbf{H}_i^{-1} \cdot \delta \mathbf{g}) \cdot \mathbf{u} \cdot \mathbf{u}^\top$$

with:

$$\delta \mathbf{x} = \mathbf{x}_{i+1} - \mathbf{x}_i \quad \delta \mathbf{g} = \mathbf{g}(\mathbf{x}_{i+1}) - \mathbf{g}(\mathbf{x}_i) \quad \mathbf{u} = \frac{\delta \mathbf{x}}{\delta \mathbf{x}^\top \cdot \delta \mathbf{g}} - \frac{\mathbf{H}_i^{-1} \cdot \delta \mathbf{g}}{\delta \mathbf{g}^\top \cdot \mathbf{H}_i^{-1} \cdot \delta \mathbf{g}}$$

The algorithm is usually started with the identity matrix ($\mathbf{H}_0^{-1} = \mathbf{I}$).

This procedure is implemented as: `BFGS(Func, Gradient, X(), MaxIter, Tol, Fmin, G(), Hinv())`

The user must provide a subroutine `Gradient` to compute the gradient `G()` of the function at point `X()`. This subroutine is declared as:

```
SUB Gradient (X() AS DOUBLE, G() AS DOUBLE)
```

The other parameters have the same meaning than in `Newton`.

7.2.6 Approximate gradient

Subroutine `NumGrad(Func, X(), Eta, G())` computes the gradient of function `Func` using finite difference approximations. `Eta` is the relative increment used to compute derivatives. It should not be less than the square root of the machine epsilon (about 1.5×10^{-8}).

You can use `NumGrad` with BFGS by defining a subroutine such that:

```
SUB Gradient (X() AS DOUBLE, G() AS DOUBLE)
  CONST Eta = 1D-6
  NumGrad @Func, X(), Eta, G()
END SUB
```

(see demo program `testbfgs.bas` for an example).

As usual, it is recommended to use analytical derivatives whenever possible.

7.2.7 Simplex method

Unlike previous methods, the simplex method of Nelder and Mead does not use derivatives to locate the minimum. Instead it constructs a geometrical figure (the ‘simplex’) having $(n + 1)$ vertices, if n is the number of variables. For instance, in the two-dimensional space ($n = 2$), the simplex would be a triangle. Depending on the function values at the vertices, the simplex is reduced or expanded until it comes close to the minimum.

This method is implemented as: `Simplex(Func, X(), MaxIter, Tol, Fmin)`, where the parameters have their usual meaning.

7.3 Demo programs

These programs are located in the `demo\optim` subdirectory.

7.3.1 Function of one variable

Program `minfunc.bas` performs the golden search minimization on the function:

$$f(x) = e^{-2x} - e^{-x}$$

The minimum is at $(\ln 2, -1/4)$.

The minimum found by `GoldSearch` is compared with the true minimum.

7.3.2 Minimization along a line

Program `minline.bas` applies line minimization to the function of 3 variables (taken from the *Numerical Recipes* example book) :

$$f(x_1, x_2, x_3) = (x_1 - 1)^2 + (x_2 - 1)^2 + (x_3 - 1)^2$$

The minimum is $f(1, 1, 1) = 0$, i. e. for a step $r = 1$ from $\mathbf{x} = [0, 0, 0]$ in the direction $\delta\mathbf{x} = [1, 1, 1]$.

The program tries a series of directions:

$$\delta\mathbf{x} = \left[\sqrt{2} \cos\left(i \frac{\pi}{20}\right), \sqrt{2} \sin\left(i \frac{\pi}{20}\right), 1 \right] \quad i = 1..10$$

For each pass, the location of the minimum, and the value of the function at the minimum, are printed. The true minimum is found at $i = 5$.

7.3.3 Newton-Raphson method

Program `testnewt.bas` uses the Newton-Raphson method to minimize Rosenbrock's function (H. Rosenbrock, *Comput. J.*, 1960, **3**, 175):

$$f(x, y) = 100(y - x^2)^2 + (1 - x)^2$$

for which the gradient and hessian are:

$$\mathbf{g}(x, y) = \begin{bmatrix} -400(y - x^2)x - 2 + 2x \\ 200y - 200x^2 \end{bmatrix}$$

$$\mathbf{H}(x, y) = \begin{bmatrix} 1200x^2 - 400y + 2 & -400x \\ -400x & 200 \end{bmatrix}$$

and the determinant of the hessian is:

$$\det \mathbf{H}(x, y) = 80000(x^2 - y) + 400$$

The minimum is $f(1, 1) = 0$, where:

$$\mathbf{g}(1, 1) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\mathbf{H}^{-1}(1, 1) = \begin{bmatrix} \frac{1}{2} & 1 \\ 1 & \frac{401}{200} \end{bmatrix}$$

$$\det \mathbf{H}(1, 1) = 400$$

In the demo program, the gradient and hessian are computed analytically. You can compare with the numerical computations using `NumHessGrad` by commenting off the relevant subroutine in the program.

7.3.4 Approximate gradient and hessian

Program `testnum.bas` computes the gradient and hessian of Rosenbrock's function numerically, at a given point, using subroutine `NumHessGrad`.

7.3.5 Other programs

Programs `testmarq.bas`, `testbfgs.bas` and `testsimp.bas` minimize Rosenbrock's function with the Marquardt, BFGS and Simplex methods, respectively.

Chapter 8

Nonlinear equations

This chapter describes the procedures available in **FBMath** to solve nonlinear equations in one or several variables. Only general methods are considered here. Polynomial equations will be studied in the next chapter.

8.1 Equations in one variable

The goal is to solve the nonlinear equation $f(x) = 0$, or, in other terms, find a root of function f .

8.1.1 Bisection method

Subroutine **Bisect**(**Func**, **X**, **Y**, **MaxIter**, **Tol**, **F**) finds a root of function **Func** by the bisection method. At each iteration, the root is bounded by two numbers (**X**, **Y**) such that the function has opposite signs. Then, a new approximation to the root is generated by taking the mean of these numbers.

The function **Func** must be declared as:

```
FUNCTION Func(X AS DOUBLE) AS DOUBLE
```

The user must provide initial values for **X** and **Y**. It is not necessary that the interval [**X**, **Y**] contains the root.

The user must also provide:

- the maximum number of iterations **MaxIter**
- the tolerance **Tol** with which the root must be located.

The subroutine returns the refined values of **X** and **Y** and the function value **Func(X)** in **F**.

After a call to **Bisect**, function **MathErr()** will return one of two error codes:

- **OptOk** if no error occurred
- **OptNonConv** (non-convergence) if the number of iterations exceeds the maximum value **MaxIter**

If the starting interval **[X, Y]** does not contain the root, **Bisect** will expand it by calling a subroutine **RootBrack**. This subroutine may be called independently. Its syntax is:

RootBrack(Func, X, Y, FX, FY)

The user must provide initial values for the two numbers **X** and **Y**, which will be refined by the subroutine. The corresponding function values are returned in **FX** and **FY**.

8.1.2 Secant method

The secant method also starts with two approximations x and y and generates a new approximation z from the formula:

$$z = \frac{xf(y) - yf(x)}{f(y) - f(x)}$$

z is the intersection of the Ox axis with the line connecting the points $(x, f(x))$ and $(y, f(y))$, i. e. the secant.

This method is implemented as:

Secant(Func, X, Y, MaxIter, Tol, F)

The parameters and error codes are the same than in **Bisect**. Here too, it is not necessary that the interval **[X, Y]** contains the root.

8.1.3 Newton-Raphson method

The Newton-Raphson method starts with an approximate root x^0 and generates a new approximation x by using the first-order Taylor series expansion of function f around x^0 :

$$f(x) \approx f(x^0) + f'(x^0) \cdot (x - x^0)$$

If x is sufficiently close to the root, $f(x) \approx 0$ so:

$$x = x^0 - \frac{f(x^0)}{f'(x^0)}$$

This method is implemented as:

`NewtEq(Func, Deriv, X, MaxIter, Tol, F)`

where `Func` and `Deriv` are the procedures which compute the function and its derivative, respectively (they have the same syntax). The user must provide the initial approximation `X`.

After a call to `NewtEq`, function `MathErr()` will return one of three error codes:

- `OptOk` if no error occurred
- `OptNonConv` (non-convergence) if the number of iterations exceeds the maximum value `MaxIter`
- `OptSing` if the derivative becomes zero

8.2 Equations in several variables

The goal is to solve a system of n nonlinear equations in n unknowns x_1, x_2, \dots, x_n :

$$\begin{aligned} f_1(x_1, x_2, \dots, x_n) &= 0 \\ f_2(x_1, x_2, \dots, x_n) &= 0 \\ &\dots\dots\dots \\ f_n(x_1, x_2, \dots, x_n) &= 0 \end{aligned}$$

or, in matrix notation:

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}$$

where \mathbf{f} is a function vector.

8.2.1 Newton-Raphson method

The Newton-Raphson method starts with an approximate root \mathbf{x}^0 and generates a new approximation \mathbf{x} by using the first-order Taylor series expansion of function \mathbf{f} around \mathbf{x}^0 :

$$\mathbf{f}(\mathbf{x}) \approx \mathbf{f}(\mathbf{x}^0) + \mathbf{D}(\mathbf{x}^0) \cdot (\mathbf{x} - \mathbf{x}^0)$$

\mathbf{D} denotes the jacobian matrix (matrix of first partial derivatives). For instance, for a system of 2 equations in two variables:

$$f_1(x_1, x_2) = 0$$

$$f_2(x_1, x_2) = 0$$

the jacobian matrix is:

$$\mathbf{D}(\mathbf{x}^0) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(x_1^0, x_2^0) & \frac{\partial f_1}{\partial x_2}(x_1^0, x_2^0) \\ \frac{\partial f_2}{\partial x_1}(x_1^0, x_2^0) & \frac{\partial f_2}{\partial x_2}(x_1^0, x_2^0) \end{bmatrix}$$

If \mathbf{x} is sufficiently close to the root, $\mathbf{f}(\mathbf{x}) \approx 0$ so:

$$\mathbf{x} = \mathbf{x}^0 - \mathbf{D}^{-1}(\mathbf{x}^0) \cdot \mathbf{f}(\mathbf{x}^0)$$

In practice, it is better to determine a step k in the direction specified by $\mathbf{D}^{-1}(\mathbf{x}^0) \cdot \mathbf{f}(\mathbf{x}^0)$:

$$\mathbf{x} = \mathbf{x}^0 - k \cdot \mathbf{D}^{-1}(\mathbf{x}^0) \cdot \mathbf{f}(\mathbf{x}^0)$$

The determination of k is performed by line minimization applied to the sum of squared functions:

$$S(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})^2$$

This method is implemented as:

```
NewtEqs(Equations, Jacobian, X(), MaxIter, Tol, F())
```

where `Equations` and `Jacobian` are the subroutines which compute the function vector and the jacobian matrix, respectively. Their syntaxes are:

```
SUB Equations(X() AS DOUBLE, F() AS DOUBLE)
```

```
SUB Jacobian(X() AS DOUBLE, D() AS DOUBLE)
```

The user must provide the initial approximations to the roots in vector `X()`. After refinement by the subroutine, the corresponding function values are returned in `F()`.

The possible error codes returned by `MathErr` are:

- `OptOk` if no error occurred
- `OptNonConv` (non-convergence) if the number of iterations exceeds the maximum value `MaxIter`
- `OptSing` if the jacobian matrix is quasi-singular

8.2.2 Approximate jacobian

Approximate values of the jacobian matrix may be computed using finite difference approximations:

$$\frac{\partial f_i}{\partial x_j}(\mathbf{x}) \approx \frac{f_i(x_j + h_j) - f_i(x_j - h_j)}{2h_j}$$

The increment h_j is such that $h_j = \eta |x_j|$ where η is a constant which should not be less than the square root of the machine epsilon ($\text{MachEp}^{1/2}$).

This method is implemented as:

```
NumJacobian(Equations, X(), Eta, D())
```

You can use this procedure with `NewtEqs` by defining a subroutine such that:

```
SUB Jacobian (X() AS DOUBLE, D() AS DOUBLE)
  CONST Eta = 1D-6
  NumJacobian @Equations, X(), Eta, D()
END SUB
```

(see demo program `testnr.bas` for an example).

As usual, it is recommended to use analytical expressions for the derivatives whenever possible.

8.2.3 Broyden's method

This method is similar to the BFGS method of function minimization. It can also be viewed as a multidimensional version of the secant method.

Broyden's algorithm uses the following formula to construct the inverse jacobian matrix iteratively:

$$\mathbf{D}_{i+1}^{-1} = \mathbf{D}_i^{-1} + \frac{[(\delta\mathbf{x} - \mathbf{D}_i^{-1} \cdot \delta\mathbf{f}) \cdot \delta\mathbf{x}^\top] \cdot \mathbf{D}_i^{-1}}{\delta\mathbf{x}^\top \cdot \mathbf{D}_i^{-1} \cdot \delta\mathbf{f}}$$

with:

$$\delta\mathbf{x} = \mathbf{x}_{i+1} - \mathbf{x}_i \quad \delta\mathbf{f} = \mathbf{f}(\mathbf{x}_{i+1}) - \mathbf{f}(\mathbf{x}_i)$$

The algorithm is usually started with the identity matrix ($\mathbf{D}_0^{-1} = \mathbf{I}$).

This method is implemented as: `Broyden(Equations, X(), MaxIter, Tol, F())`, where the parameters have the same significance than in `NewtEqs`.

The possible error codes returned by `MathErr` are `OptOk` and `OptNonConv`.

8.3 Demo programs

8.3.1 Equations in one variable

The demo programs `testbis.bas`, `testsec.bas` and `testnr1.bas` demonstrate the bisection, secant and Newton-Raphson methods, respectively, on the equation:

$$f(x) = x \ln x - 1 = 0$$

for which the derivative is:

$$f'(x) = \ln x + 1$$

The true solution is $x = 1.763222834\dots$

8.3.2 Equations in several variables

The demo programs `testnr.bas` and `testbrdn.bas` demonstrate the Newton-Raphson and Broyden methods, respectively, on the following system (taken from the *Numerical Recipes* example book) :

$$f(x, y) = x^2 + y^2 - 2 = 0$$

$$g(x, y) = \exp(x - 1) + y^3 - 2 = 0$$

for which the jacobian is:

$$\mathbf{D}(x, y) = \begin{bmatrix} 2x & 2y \\ \exp(x - 1) & 3y^2 \end{bmatrix}$$

The true solution is $(x, y) = (1, 1)$.

Chapter 9

Polynomials

This chapter describes the procedures and functions related to polynomials and rational fractions.

9.1 Polynomials

Function `Poly(X, Coef(), Deg)` evaluates the polynomial:

$$P(X) = \text{Coef}(0) + \text{Coef}(1) \cdot X + \text{Coef}(2) \cdot X^2 + \cdots + \text{Coef}(\text{Deg}) \cdot X^{\text{Deg}}$$

9.2 Rational fractions

Function `RFrac(X, Coef(), Deg1, Deg2)` evaluates the rational fraction:

$$F(X) = \frac{\text{Coef}(0) + \text{Coef}(1) \cdot X + \cdots + \text{Coef}(\text{Deg1}) \cdot X^{\text{Deg1}}}{1 + \text{Coef}(\text{Deg1} + 1) + \cdots + \text{Coef}(\text{Deg1} + \text{Deg2}) \cdot X^{\text{Deg2}}}$$

9.3 Roots of polynomials

Analytical methods can be used to compute the roots of polynomials up to degree 4. For higher degrees, iterative methods must be used.

9.3.1 Analytical methods

- Function `RootPol1(A, B, X)` solves the linear equation $A + BX = 0$. The function returns 1 if no error occurs ($B \neq 0$), -1 if X is undetermined ($A = B = 0$), -2 if there is no solution ($A \neq 0, B = 0$).

- Functions `RootPolN(Coef(), Z())`, with $N = 2, 3, 4$, solve the equation:

$$\text{Coef}(0) + \text{Coef}(1) \cdot X + \text{Coef}(2) \cdot X^2 + \cdots + \text{Coef}(N) \cdot X^N = 0$$

The roots are stored in the complex vector `Z`. The real part of the i^{th} root is in `Z(i).X`, the imaginary part in `Z(i).Y`.

If no error occurs, the function returns the number of real roots, otherwise it returns (-1) or (-2) just like `RootPol1`.

9.3.2 Iterative method

Function `RootPol(Coef(), Deg, Z())` solves the polynomial equation:

$$a_0 + a_1x + a_2x^2 + \cdots + a_nx^n = 0$$

by the method of the *companion matrix*.

The companion matrix **A** is defined by:

$$\mathbf{A} = \begin{bmatrix} -\frac{a_{n-1}}{a_n} & -\frac{a_{n-2}}{a_n} & \cdots & -\frac{a_1}{a_n} & -\frac{a_0}{a_n} \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}$$

It may be shown that the eigenvalues of this matrix are equal to the roots of the polynomial (Eigenvalues will be treated in the next chapter).

The coefficients of the polynomial are passed in vector `Coef`, such that `Coef(0) = a0`, `Coef(1) = a1` etc. The degree of the polynomial is passed in `Deg`. The roots are returned in the complex vector `Z` as described before.

If no error occurred, the function returns the number of real roots.

If an error occurred during the search for the i^{th} root, the function returns (-i). The roots should be correct for indices `(i+1)..Deg`. The roots are unordered.

9.4 Ancillary functions

Two subroutines have been added to facilitate the handling of polynomials roots:

- Function `SetRealRoots(Deg, Z(), Tol)` allows to set the imaginary part of a root to zero if it is less than a fraction `Tol` of the real part. The function returns the total number of real roots.

Due to roundoff errors, some real roots may be computed with a very small imaginary part, e.g. $1 + 10^{-8}i$. The function `SetRealRoots` tries to correct this problem.

- Subroutine `SortRoots(Deg, Z())` sort the roots such that:
 1. The `N` real roots are stored in elements `1..N` of vector `Z`, in increasing order.
 2. The complex roots are stored in elements `(N + 1) .. Deg` of vector `Z` and are unordered.

9.5 Demo programs

9.5.1 Evaluation of a polynomial

Program `evalpoly.bas` evaluates a polynomial for a series of user-specified values. Entering 0 stops the program.

9.5.2 Evaluation of a rational fraction

Program `evalfrac.bas` performs the same task as the previous program, but with a rational fraction.

9.5.3 Roots of a polynomial

Program `polyroot.bas` computes the roots of a polynomial with real coefficients. Analytical methods are used up to degree 4, otherwise the method of the companion matrix is used.

The example polynomial is:

$$x^6 - 21x^5 + 175x^4 - 735x^3 + 1624x^2 - 1764x + 720$$

for which the roots are 1, 2 ... 6

Chapter 10

Numerical integration and differential equations

This chapter describes the procedures available in **FBMath** to integrate a function of one variable, and to solve systems of differential equations.

10.1 Integration

10.1.1 Trapezoidal rule

The trapezoidal rule approximates the integral I of a tabulated function by the formula:

$$I \approx \frac{1}{2} \sum_{i=1}^{N-1} (x_{i+1} - x_i)(y_{i+1} + y_i)$$

where N is the number of points and (x_i, y_i) the coordinates of the i^{th} point.

This procedure is implemented as function **TrapInt(X(), Y())**. Note that the lower bound of the arrays does not need to be 1.

10.1.2 Gauss-Legendre integration

This method approximates the integral of a function f in an interval $[a, b]$ by:

$$\int_a^b f(x)dx \approx \frac{b-a}{2} \sum_{i=1}^N w_i f(y_i)$$
$$y_i = \frac{b-a}{2}x_i + \frac{b+a}{2}$$

The abscissae x_i and weights w_i are predefined values for a given number of points N .

This method is implemented as function `GausLeg(Func, A, B)` for $N = 16$. Function `Func` must be declared as:

```
FUNCTION Func(X AS DOUBLE) AS DOUBLE
```

For the special case $A = 0$ there is a variant `GausLeg0(Func, B)`.

10.2 Convolution

The convolution product of two functions f and g is defined by:

$$(f * g)(t) = \int_0^t f(u)g(t-u)du$$

This product is often used to describe the output of a linear system when $f(t)$ is the input signal (function of time) and $g(t)$ is the impulse response of the system.

Function `Convol(Func1, Func2, T)` approximates the convolution product of the two functions `Func1` and `Func2` at time `T` by the Gauss-Legendre method. The functions must be declared as above.

10.3 Differential equations

The Runge-Kutta-Fehlberg (RKF) method allows to compute numerical solutions to systems of first-order differential equations of the form:

$$\begin{aligned} y_1'(t) &= f_1[t, y_1(t), y_2(t), \dots] \\ y_2'(t) &= f_2[t, y_1(t), y_2(t), \dots] \\ &\dots\dots\dots \end{aligned}$$

where the f_i are known functions and the y_i are to be determined.

The RKF procedure is an extension of the classical Runge-Kutta method. For instance, in the case of a single differential equation

$$y'(t) = f[t, y(t)]$$

this method generates a sequence $\{t_n, y_n\}$ which approximates the function $y(t)$.

The order of the method corresponds to the number of points used in the interval $[t_n, t_{n+1}]$. For instance, the sequence generated by the 4-th order Runge-Kutta method is defined by:

$$y_{n+1} = y_n + \frac{k_1}{6} + \frac{k_2}{3} + \frac{k_3}{3} + \frac{k_4}{6}$$

with:

$$\begin{aligned} k_1 &= h \cdot f(t_n, y_n) \\ k_2 &= h \cdot f\left(t_n + \frac{h}{2}, y_n + \frac{k_1}{2}\right) \\ k_3 &= h \cdot f\left(t_n + \frac{h}{2}, y_n + \frac{k_2}{2}\right) \\ k_4 &= h \cdot f(t_n + h, y_n + k_3) \end{aligned}$$

with $h = t_{n+1} - t_n$

In the RKF method, the step size h is automatically varied so as to maintain a given level of precision on the estimated y values.

The implementation used in `FBMath` is a translation of a Fortran program by H. A. Watts and L. F. Shampine (http://www.csit.fsu.edu/~burkardt/f_src/rkf45/rkf45.f90). It is intermediate between the 4-th and 5-th order Runge-Kutta methods, hence the name RKF45.

In order to use RKF45 you must:

1. Define the following constants and variables (the names are optional and the values are given as examples, except for `Flag` which must be initialized to 1):

```

CONST Neqn = 2                                ' Number of equations

DIM AS DOUBLE  Y(1 TO Neqn) = {1, 0}         ' Functions and initial cond.
DIM AS DOUBLE  Yp(1 TO Neqn)                 ' Derivatives

DIM AS DOUBLE  Tstart = 0, Tstop = 10         ' Integration interval
DIM AS INTEGER Nstep = 5                      ' Number of steps

DIM AS DOUBLE  AbsErr = 1D-6, _
              RelErr = 1D-6                   ' Absolute and relative errors

DIM AS INTEGER Flag = 1                       ' Error flag

```

2. Define a subroutine for computing the system of differential equations:

```
SUB DiffEq(T AS DOUBLE, Y() AS DOUBLE, Yp() AS DOUBLE)
  Yp(1) = Y(2)
  Yp(2) = - Y(1)
END SUB
```

3. Compute the step size and call RKF45 for each integration step:

```
DIM AS DOUBLE StepSize = (Tstop - Tstart) / Nstep
DIM AS DOUBLE T        = Tstart
DIM AS DOUBLE Tout
DIM AS INTEGER I

FOR I = 1 TO Nstep
  Tout = T + StepSize
  RKF45 @DiffEq, Y(), Yp(), T, Tout, RelErr, AbsErr, Flag
  T = Tout
NEXT I
```

Upon return from the RKF45 subroutine:

- `Y()`, `Yp()` contain the values of the functions and their first derivatives at `Tout`
- `Flag` contains an error code:
 - * 2 : no error
 - * 3 : too small `RelErr` value
 - * 4 : too much function evaluations needed
 - * 5 : too small `AbsErr` value
 - * 6 : the requested accuracy could not be achieved
 - * 7 : the method was unable to solve the problem
 - * 8 : invalid input parameters

If an error occurs, it should be possible in most cases to restart the computation, using the values returned by the subroutine in `RelErr` and `AbsErr`.

Note : RKF45 may be used to compute a definite integral:

$$\int_a^b f(t)dt = F(b) - F(a)$$

since this is equivalent to integrate the differential equation:

$$F'(t) = f(t)$$

between a and b , with the initial condition specified by $f(a)$.

10.4 Demo programs

- Program `trap.bas` applies the trapezoidal rule to a tabulated function.

The example function $f(x) = e^{-x}$ is tabulated for $x = 0$ to 1 by steps of 0.1. The integral is:

$$\int_0^1 e^{-x} dx = 1 - e^{-1} \approx 0.6321$$

- Program `gauss.bas` demonstrates the Gauss-Legendre integration method.

The example function is $f(x) = xe^{-x}$. The integral is:

$$\int_0^x f(t)dt = 1 - (x + 1)e^{-x}$$

- Program `conv.bas` computes the convolution of two functions by the Gauss-Legendre method.

The example functions are $f(x) = xe^{-x}$ and $g(x) = e^{-2x}$. The convolution product is:

$$(f * g)(x) = \int_0^x f(u)g(x - u)du = e^{-2x} \int_0^x ue^u du = (x - 1)e^{-x} - e^{-2x}$$

- Program `test_rkf.bas` solves 3 systems of differential equations by the RKF method:

1. A single nonlinear equation:

$$y'(t) = 0.25 \cdot y(t) \cdot [1 - 0.05 \cdot y(t)]$$

with the initial condition $y(0) = 1$.

The analytic solution is:

$$y(t) = \frac{20}{1 + 19 \exp(-0.25t)}$$

2. A system of two linear equations:

$$y_1'(t) = y_2(t)$$

$$y_2'(t) = -y_1(t)$$

with the initial conditions $y_1(0) = 1, y_2(0) = 0$.

The analytic solution is:

$$y_1(t) = \cos t \quad y_2(t) = -\sin t$$

3. A system of 5 equations with one nonlinear:

$$y_1'(t) = y_2(t)$$

$$y_2'(t) = y_3(t)$$

$$y_3'(t) = y_4(t)$$

$$y_4'(t) = y_5(t)$$

$$y_5'(t) = \frac{45 \cdot y_3(t) \cdot y_4(t) \cdot y_5(t) - 40[y_4(t)]^3}{9[y_3(t)]^2}$$

with initial conditions $y_i(0) = 1 \quad \forall i$

The program prints the numeric solution, and, if possible, the analytic one.

Chapter 11

Fast Fourier Transform

11.1 Introduction

Fourier transform is a mathematical method which allows to determine the frequency spectrum of a given signal (for instance a sound). The mathematical definition is the following :

$$y(f) = \int_{-\infty}^{\infty} x(t) \exp(2\pi i f t) dt = \int_{-\infty}^{\infty} x(t) (\cos 2\pi f t + i \sin 2\pi f t) \quad (11.1)$$

where $x(t)$ is the input signal (function of time), f the frequency, and i the complex number such that $i^2 = -1$. y is the Fourier transform of x .

The input signal may have real or complex values. However, the Fourier transform is always a complex number. For each frequency f , the modulus of $y(f)$ represents the energy associated with this frequency. A plot of this modulus as a function of f gives the frequency spectrum of the input signal.

If the input signal is sampled as a sequence of n values x_0, x_1, \dots, x_{n-1} , taken at constant time intervals, the Fourier transform is a sequence of complex number y_0, y_1, \dots, y_{n-1} , such that:

$$y_p = \sum_{k=0}^{n-1} x_k \left[\cos \left(2\pi \frac{kp}{n} \right) + i \sin \left(2\pi \frac{kp}{n} \right) \right] \quad (11.2)$$

This formula allows, in principle, to compute the transform y_p at any point. In practice, a faster algorithm called the Fast Fourier Transform (FFT) is used.

11.2 Programming

11.2.1 Array dimensioning

The FFT algorithm requires that the number of points n is a power of 2. Moreover, the arrays must be dimensioned from 0 to n . For instance:

```
CONST NumSamples = 512           ' Buffer size must be power of 2
CONST MaxIndex   = NumSamples - 1 ' Max. array index

DIM AS Complex InArray(0 TO MaxIndex) ' FFT input
DIM AS Complex OutArray(0 TO MaxIndex) ' FFT output
```

where type `Complex` is defined in the include file `math.bi` as follows:

```
TYPE Complex
  X AS DOUBLE
  Y AS DOUBLE
END TYPE
```

11.2.2 FFT procedures

- Subroutine `FFT(InArray(), OutArray())` calculates the Fast Fourier Transform of the array of complex numbers `InArray` to produce the output complex numbers in `OutArray`
- Subroutine `IFFT(InArray(), OutArray())` calculates the Inverse Fast Fourier Transform of the array of complex numbers represented by `InArray` to produce the output complex numbers in `OutArray`

In other words, this subroutine reconstitutes the input signal from its FFT.

- Function `CalcFrequency(FrequencyIndex, InArray())` calculates the complex frequency sample at a given index directly, by means of eq. 11.2. Use this instead of `FFT` when you only need one or two frequency samples, not the whole spectrum. It is also useful for calculating the Fourier Transform of a number of data which is not an integer power of 2. For example, you could calculate the transform of 100 points instead of rounding up to 128 and padding the extra 28 array slots with zeroes.

11.3 Demo program

Program `test_fft.bas`, located in the `demo\fourier` subdirectory, shows how the Fourier transform may be used to filter a signal. The program plots several graphics and writes its results to the output file `fftout.txt`

The example is a 200 Hz sine wave contaminated by a 2000 Hz parasitic signal. The sampling frequency `SamplingRate` is 22050 Hz, the number of points `NumSamples` is 512 ($= 2^9$). These two numbers determine the time and frequency units:

```
CONST DT = 1 / SamplingRate      ' Time unit
CONST DF = SamplingRate / NumSamples ' Frequency unit
```

so that the entry `InArray(I)` in the input array of subroutine `FFT` corresponds to the signal value at time $I * DT$, and that the entry `OutArray(I)` in the output array corresponds to the Fourier transform at frequency $I * DF$.

The highest frequency which may be detected is equal to `SamplingRate/2` and is called *Nyquist's frequency*. Hence, only the first half of array `OutArray` needs to be plotted (the second half contains symmetric values).

The program generates the input signal, plots it, then performs the FFT and plots the real and imaginary parts as a function of frequency. The plot shows two peaks, corresponding to the 5-th and 46-th entries in `OutArray` (as seen from the file `fftout.txt`). The corresponding frequencies are:

$$5 \times \frac{22050}{512} \approx 215\text{Hz}$$

$$46 \times \frac{22050}{512} \approx 1981\text{Hz}$$

The high peak corresponds to the main signal and the small peak to the parasite. To filter the last one, the program sets to zero all the FFT values corresponding to the frequencies higher than 1000 Hz, according to the following code:

```
CONST MidIndex = NumSamples \ 2

FreqIndex = 1000 / DF
SymIndex = NumSamples - FreqIndex

FOR I = 0 TO MaxIndex
```

```

      IF (I > FreqIndex AND I < MidIndex) OR _
        (I >= MidIndex AND I < SymIndex) THEN
          OutArray(I).X = 0
          OutArray(I).Y = 0
      END IF
NEXT I

```

(note that the two halves of the output array, on either side of Nyquist's frequency, must be treated).

The program then calls subroutine `IFFT` to compute the inverse Fourier transform of the modified data and plots the result, showing that the parasite has been removed, at the expense of a slight distortion of the main signal.

In addition, the program performs a direct computation of the Fourier transform of a set of random complex values, using function `CalcFrequency`, and stores the results in the output file, for comparison with the FFT computed on the same data.

Chapter 12

Random numbers

This chapter describes the procedures and functions available to generate random numbers and perform stochastic simulation and optimization.

12.1 Random numbers

12.1.1 Introduction

FBMath provides three random number generators (RNG) :

- the ‘Multiply With Carry’ (MWC) generator of George Marsaglia.
- the ‘Mersenne Twister’ (MT) generator of Takuji Nishimura and Makoto Matsumoto (<http://www.math.sci.hiroshima-u.ac.jp/~m-mat/MT/emt.html>). The FreeBasic version was contributed by Greg Lyon.
- the ‘Universal Virtual Array Generator’ (UVAG) contributed by Alex Hay.

The first method produces a sequence $\{I_n\}$ of integer numbers by means of the following recurrence relationships:

$$I_{n+1} = (aI_n + c_n) \bmod b$$

$$c_{n+1} = (aI_n + c_n) \operatorname{div} b$$

where a is the *multiplier*, b the *base* (here $b = 2^{32}$), c_n the *carry*. One may start with $c_0 = 0$.

If a is properly chosen, the period of the generator is $a \times 2^{31} - 1$. In our implementation, $a = 2051013963$ so that the period is approximately 4.4×10^{18} .

The second method is more complex and slightly slower but it may be safer for intensive simulations since it has a much longer period (about 10^{6000}) and produces uncorrelated numbers in 623 dimensions.

We have verified that the three generators pass Marsaglia's DIEHARD battery of tests (<http://stat.fsu.edu/pub/diehard/>).

12.1.2 Generic functions

These functions may be used with any of the three generators.

Choice of generator

The default generator is MWC. This selection can be modified by using the statements `SetRNG RNG_MT`, `SetRNG RNG_UVAG` or `SetRNG RNG_MWC`.

Initialization

The selected generator can be initialized with the statement `InitGen Seed` where `Seed` is a 32-bit integer. The default seeds are 123456789 for MWC and 5489 for MT.

Note : The initialization is optional for MT and MWC (the default initialization will be used) but mandatory for UVAG.

Uniform random numbers

The following functions are available:

Function	Type	Bits	Interval
<code>IRanGen</code>	<code>UInteger</code>	32	[0, 4294967295]
<code>IRanGen31</code>	<code>Long</code>	31	[0, 2147483647]
<code>RanGen1</code>	<code>Double</code>	32	[0, 1]
<code>RanGen2</code>	<code>Double</code>	32	[0, 1)
<code>RanGen3</code>	<code>Double</code>	32	(0, 1)
<code>RanGen53</code>	<code>Double</code>	53	[0, 1)

12.1.3 Specific functions

The following functions are specific of a given generator:

- Subroutine `InitMWC(Seed)` initializes the MWC generator with a 32-bit integer

- Function `IRanMWC` returns a 32-bit random integer from the MWC generator
- Subroutine `InitMT(Seed)` initializes the MT generator with a 32-bit integer
- Subroutine `InitMTbyArray(InitKey())` initializes the MT generator with an array `InitKey(0 to N)` with $N < 624$
- Function `IRanMT` returns a 32-bit random integer from the MT generator
- Subroutine `InitUVAGbyString(KeyPhrase)` initializes the 32-bit UVAG generator with a string
- Subroutine `InitUVAG64byString(KeyPhrase)` initializes the 64-bit UVAG generator with a string
- Subroutine `InitUVAG(Seed)` initializes the 32-bit UVAG generator with a 32-bit integer
- Subroutine `InitUVAG64(Seed)` initializes the 64-bit UVAG generator with a 64-bit integer
- Function `IRanUVAG` returns a 32-bit random integer from the UVAG generator
- Function `IRanUVAG64` returns a 64-bit random integer from the UVAG generator

12.1.4 Gaussian random numbers

These functions use the selected generator (MWC by default).

Normal distribution

- Function `RanGaussStd` generates a random number from the standard normal distribution.

The Box-Muller algorithm is used: if x_1 and x_2 are two uniform random numbers $\in (0, 1)$, the two numbers y_1 and y_2 defined by:

$$y_1 = \sqrt{-2 \ln x_1} \cos 2\pi x_2 \quad y_2 = \sqrt{-2 \ln x_1} \sin 2\pi x_2$$

follow the standard normal distribution.

- Function `RanGauss(Mu, Sigma)` generates a random number from the normal distribution with mean `Mu` and standard deviation `Sigma`.

Multinormal distribution

- Subroutine `RanMult(M(), L(), X())` generates a random vector `X` from a multidimensional normal distribution. `M` is the mean vector, `L` is the Cholesky factor of the variance-covariance matrix.

To simulate the n -dimensional multinormal distribution $\mathcal{N}(\mathbf{m}, \mathbf{V})$, where \mathbf{m} is the mean vector and \mathbf{V} the variance-covariance matrix, the following algorithm is used:

1. Let \mathbf{u} be a vector of n independent random numbers following the standard normal distribution,
 2. Let \mathbf{L} be the lower triangular matrix resulting from the Cholesky factorization of matrix \mathbf{V} ,
 3. Vector $\mathbf{x} = \mathbf{m} + \mathbf{L}\mathbf{u}$ follows the multinormal distribution $\mathcal{N}(\mathbf{m}, \mathbf{V})$.
- Subroutine `RanMultIndep(M(), S(), X())` generates a random number from an uncorrelated multidimensional distribution. Here `S` is simply the vector of standard deviations.

12.2 Markov Chain Monte Carlo

It is not always possible to simulate the distribution of a random variable with a direct algorithm such as the ones used for normal or multinormal distributions.

However, there exist iterative algorithms which generate a sequence of random variables for which the distribution tend towards the desired distribution, after starting from a standard distribution (e. g. uniform).

These random sequences are known as *Markov chains* and the iterative simulation method is therefore known as *Markov chain Monte-Carlo* (MCMC).

There are several MCMC variants. Here we will present the Metropolis-Hastings method.

Let \mathbf{X} a vector of random variables and $P(\mathbf{X})$ its probability density function (p.d.f.), which is to be simulated. The classical formulation of the Metropolis-Hastings algorithm is the following:

1. Choose an initial parameter vector \mathbf{X}_0
2. At iteration n :
 - (a) Draw a vector \mathbf{u} from the multinormal distribution $\mathcal{N}(\mathbf{X}_{n-1}, \mathbf{V})$ where \mathbf{V} is the variance-covariance matrix
 - (b) If $r = P(\mathbf{u})/P(\mathbf{X}_{n-1}) > 1$, set $\mathbf{X}_n = \mathbf{u}$
 otherwise if $\text{Random}(0, 1) < r$, set $\mathbf{X}_n = \mathbf{u}$
 where $\text{Random}(0, 1)$ denotes a uniform random number in the interval $[0, 1]$
3. Set $n = n + 1$; goto 2

It is convenient to introduce a function $F(\mathbf{X})$ such that:

$$P(\mathbf{X}) = C \exp \left[-\frac{F(\mathbf{X})}{T} \right] \quad \Longleftrightarrow \quad F(\mathbf{X}) = -T \ln \frac{P(\mathbf{X})}{C} \quad (12.1)$$

where C and T are positive constants. By analogy with statistical thermodynamics, T is known as the *temperature*.

From this equation, it may be seen that:

$$r = \frac{P(\mathbf{u})}{P(\mathbf{X}_{n-1})} = \exp \left(-\frac{\Delta F}{T} \right)$$

where

$$\Delta F = F(\mathbf{u}) - F(\mathbf{X}_{n-1})$$

so, the Metropolis-Hastings algorithm may be rewritten as:

1. Choose an initial parameter vector \mathbf{X}_0
2. At iteration n :
 - (a) Draw a vector \mathbf{u} from the multinormal distribution $\mathcal{N}(\mathbf{X}_{n-1}, \mathbf{V})$
 Set $\Delta F = F(\mathbf{u}) - F(\mathbf{X}_{n-1})$
 - (b) if $\Delta F < 0$, set $\mathbf{X}_n = \mathbf{u}$
 otherwise if $\text{Random}(0, 1) < \exp(-\Delta F/2)$, set $\mathbf{X}_n = \mathbf{u}$
3. Set $n = n + 1$; goto 2

The initial variance-covariance matrix \mathbf{V} may be diagonal and its elements may be given large values, so that the initial distribution spans a relatively large space. When the iterations progress, the matrix converges to the variance-covariance matrix of the simulated distribution. It is often useful to perform several cycles of the algorithm, with the variance-covariance matrix being re-evaluated at the end of each cycle.

The vector \mathbf{X} corresponding to the lowest value of F is recorded; hence, the algorithm may be used as a stochastic optimization algorithm for minimizing the function F . The advantage of such an algorithm is that it can ‘escape’ from a local minimum (with a probability equal to $e^{-\Delta F/T}$) and has therefore more chances to reach the global minimum, unlike the deterministic optimizers studied in chapter 7, for which only decreases of the function are acceptable. This application is however restricted by the fact that the function F must be linked to a p.d.f. by means of eq. (12.1).

This method is implemented in **FBMath** as:

```
Hastings(Func, T, X(), V(), Xmat(), Xmin(), Fmin)
```

The user must provide :

- the function **Func** to be minimized (defined as in paragraph 7.2, p. 42)
- the temperature **T**
- a starting vector **X()**
- a starting variance-covariance matrix **V()**.

On output, **Hastings** returns:

- the mean of the simulated distribution in **X()**
- its variance-covariance matrix in **V()**
- a matrix of simulated vectors in **Xmat()** (one vector by line)
- the vector which minimizes the function in **Xmin()**
- the value of the function at the minimum in **Fmin** (corresponds to the mode of the simulated distribution).

The behavior of the algorithm can be controlled with the following procedure:

`InitMHPParams(NCycles, MaxSim, SavedSim)`

where:

- `NCycles` is the number of cycles (default = 10)
- `MaxSim` is the maximum number of simulations at each cycle (default = 1000)
- `SavedSim` is the number of simulated vectors which are saved in matrix `Xmat()`. Only the last `SavedSim` vectors from the last cycle are saved. (default = 1000)

After a call to `Hastings`, function `MathErr` will return one of the following codes:

- `OptOk` if no error occurred
- `MatNotPD` if the variance-covariance matrix is not positive definite

The random number generator is re-initialized at the start of the algorithm, so that a different result will be obtained for each call of the subroutine.

12.3 Simulated Annealing

Simulated annealing (SA) is an extension of the Metropolis-Hastings algorithm which tries to find the global minimum of any function (not necessarily a p.d.f.). Here the temperature starts from a high value and is progressively decreased as the algorithm progresses towards the minimum. The optimized parameters may then be refined with a local optimizer (chapter 7).

There are several implementations of this algorithm. The one used in `FBMath` is a modification of a Fortran program written by B. Goffe (<http://www.netlib.org/simann>).

With the notations:

$F(\mathbf{X})$: function to be minimized
$\delta\mathbf{X}$: range of \mathbf{X}
F_{min}	: minimum of $F(\mathbf{X})$
T	: temperature
N_T	: number of loops at constant T
N_S	: number of loops before ajustement of $\delta\mathbf{X}$
R_T	: temperature reduction factor
N_{acc}	: number of accepted function increases

the algorithm may be described as follows:

-
- initialize $T, \mathbf{X}, \delta\mathbf{X}$
 - repeat
 - repeat N_T times
 - ★ repeat N_S times
 - for each parameter X_i :
 - ◊ pick a random value X'_i in the interval $X_i \pm \delta X_i$
 - ◊ compute $F(X'_i)$
 - ◊ accept the new value X'_i according to the Metropolis criterion
 - ◊ update N_{acc}
 - ◊ update F_{min} if necessary
 - ★ adjust step length δX_i so as to maintain an acceptance ratio of about 50%
 - $T \leftarrow T \cdot R_T$
 - until $N_{acc} = 0$ or $T < T_{min}$ or $|F_{min}| < \epsilon$
-

The threshold values T_{min} and ϵ are fixed at 10^{-300} in our implementation.

At the beginning of the iterations, while we are away from the minimum, it makes sense to choose a high probability of acceptance, for instance $p = \frac{1}{2}$. It is then possible to perform a given number of random drawings and to compute the median M of the increases of function F , from which the initial temperature T_0 is deduced by:

$$p = \exp\left(-\frac{M}{T_0}\right) = \frac{1}{2} \quad \Rightarrow \quad T_0 = \frac{M}{\ln 2}$$

This procedure is implemented in the following subroutine:

`SimAnn(Func, X(), Xmin(), Xmax(), Fmin)`

where:

- `Func` is the function to be minimized (defined as in paragraph 7.2, p. 42)

- `X()` is the parameter vector
- `Xmin()`, `Xmax()` are the bound values of `X()`

The optimized parameters are returned in `X()` and the corresponding function value in `Fmin`

The user must provide reasonable values of `Xmin()` and `Xmax()` as well as a starting value for `X()`. It is convenient to pick a random value in the range specified by `Xmin()` and `Xmax()`.

Upon return from the subroutine, function `MathErr()` will return one of the following error codes:

- `MatOk` if no error
- `MathErrDim` if `X()`, `Xmin()`, `Xmax()` don't have the same bounds

The behavior of the algorithm can be controlled with the following procedure:

```
InitSAParams(NT, NS, RT, NCycles)
```

where:

- `NT`, `NS`, `RT` correspond to the variables N_T , N_S and R_T in the algorithm. Default values are 5, 15 and 0.9 respectively.
- `NCycles` is the number of cycles (default = 1).

In some difficult situations, it may be useful to perform several cycles of the algorithm. Each cycle will start with the optimized parameters `X()` from the previous cycle and the temperature will be re-initialized (the bound values `Xmin()`, `Xmax()` remaining the same).

It is possible to record the progress of the iterations in a log file. This file is created with:

```
SA_CreateLogFile(LogFileName)
```

The default name is `simann.txt`

If the file is created, the following information will be stored:

- iteration number (each iteration corresponds to a single temperature)

- temperature value
- lowest function value obtained at this temperature
- number of function increases
- number of accepted increases

The file will be automatically closed upon return from **SimAnn**.

12.4 Genetic Algorithm

Genetic Algorithms (GA) are another class of stochastic optimization methods which try to mimick the law of natural selection in order to optimize a function $F(\mathbf{X})$.

There are several implementations of these algorithms. We use a method described by E. Perrin *et al.* (*Recherche operationnelle / Operations Research*, 1997, **31**, 161-201). In this version, the vector \mathbf{X} is considered as the ‘phenotype’ of an ‘individual’ belonging to a ‘population’. This phenotype is determined by two ‘chromosomes’ \mathbf{C}_1 and \mathbf{C}_2 and a vector of ‘dominances’ \mathbf{D} such that:

$$X_i = D_i C_{1i} + (1 - D_i) C_{2i} \quad (0 < D_i < 1) \quad (12.2)$$

A population is defined by a matrix \mathbf{P} , such that each row of the matrix corresponds to a vector \mathbf{X} .

The population is initialized by taking vectors \mathbf{C}_1 and \mathbf{C}_2 at random in a given interval, vector \mathbf{D} at random in (0,1) then applying eq. (12.2) to obtain the corresponding \mathbf{X} vectors.

At each step (‘generation’) of the algorithm:

1. The function values $F(\mathbf{X})$ are computed for each vector \mathbf{X} and the N_S individuals having the lowest function values (the ‘survivors’) are selected.
2. The remaining individuals are discarded and replaced by new ones, generated as follows:
 - (a) Two ‘parents’ are chosen at random in the selected sub-population and a ‘child’ is generated by:
 - taking the vectors \mathbf{C}_1 and \mathbf{C}_2 at random from the parents

- generating a new vector \mathbf{D}
- computing the new \mathbf{X} according to eq. (12.2)

This process is repeated until the function value for the child is lower than the lowest function value of the two parents.

- (b) The child is ‘mutated’ (i. e. its vectors are reinitialized at random) with a probability M_R
- (c) The child is made ‘homozygous’ (i. e. its vectors \mathbf{C}_1 and \mathbf{C}_2 are made identical to its vector \mathbf{X}) with a probability H_R

This procedure is implemented in the following subroutine:

```
GenAlg(Func, X(), Xmin(), Xmax(), Fmin)
```

where the parameters have the same meaning as in **SimAnn**. The error codes returned by function **MathErr** are also the same.

The behavior of the algorithm can be controlled with the following procedure:

```
InitGAParams(NP, NG, SR, MR, HR)
```

where:

- NP is the population size (default = 200)
- NG is the number of generations (default = 40)
- SR is the survival rate (default = 0.5)
- MR is the mutation rate (default = 0.1)
- HR is the probability of homozygosis (default = 0.5)

A log file may also be created with:

```
GA_CreateLogFile(LogFileName)
```

The default name is **genalg.txt**. The file will contain the iteration (generation) number and the optimized function value for this generation.

12.5 Demo programs

12.5.1 Test of MWC generator

Program `testmwc.bas` picks 2000 random numbers and displays the next 6 together with the correct values obtained with the default initialization, i.e. `InitMWC(123456789)`. These values have been checked with the `Maple` software.

12.5.2 Test of MT generator

Program `testmt.bas` writes 1000 integer numbers and 1000 real numbers from functions `IRanGen` and `RanGen2`, after initialization with a vector of 4 integers:

```
DIM Init(0 TO 3) AS UINTEGER = {&H123, &H234, &H345, &H456}
SetRNG RNG_MT
InitMTbyArray Init()
```

The output of this program should be similar to the contents of file `mt_out.txt`

12.5.3 Test of UVAG generator

Program `testuvag.bas` writes 1000 integer numbers (32-bit or 64-bit) from functions `IRanUVAG` or `IRanUVAG64`, after initialization with the string "abcd". The output should be similar to the contents of file `uvag32.txt` or `uvag64.txt`

12.5.4 File of random numbers

Program `randfile.bas` generates a binary file of 32-bit random integers to be used as input for the DIEHARD program. The user must specify the number of random integers to be generated (default is 3,000,000).

12.5.5 Gaussian random numbers

Program `testnorm.bas` picks a random sample of size N from a gaussian distribution with known mean and standard deviation (SD), estimates mean (m) and SD (s) from the sample, and computes a 95% confidence interval

for the mean (i.e. an interval which has a probability of 0.95 to include the true mean), using the formula:

$$\left[m - 1.96 \frac{s}{\sqrt{N}}, m + 1.96 \frac{s}{\sqrt{N}} \right]$$

This formula is valid for $N > 30$.

12.5.6 Multinormal distribution

Program `ranmul.bas` simulates a multi-normal distribution. The example is a 3-dimensional distribution with the following means, standard deviations, and correlation matrix:

$$\mathbf{m} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad \mathbf{s} = \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} 1 & 0.25 & 0.5 \\ 0.25 & 1 & -0.25 \\ 0.5 & -0.25 & 1 \end{bmatrix}$$

After 1000 simulations with the default random number initialization, the following estimations were obtained:

$$\hat{\mathbf{m}} = \begin{bmatrix} 0.99 \\ 2.00 \\ 2.98 \end{bmatrix} \quad \hat{\mathbf{s}} = \begin{bmatrix} 0.095 \\ 0.201 \\ 0.302 \end{bmatrix} \quad \hat{\mathbf{R}} = \begin{bmatrix} 1 & 0.229 & 0.462 \\ 0.229 & 1 & -0.281 \\ 0.462 & -0.281 & 1 \end{bmatrix}$$

12.5.7 Markov Chain Monte-Carlo

Although MCMC methods are best suited when there is no direct simulation algorithm available, we will use the Metropolis-Hastings method to simulate the previous multinormal distribution (program `testmcmc.bas`).

First, we have to define the function to be optimized. The probability density for a n -dimensional normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{V})$ is:

$$P(\mathbf{X}) = \frac{1}{\sqrt{(2\pi)^n |\mathbf{V}|}} \exp \left[-\frac{1}{2} (\mathbf{X} - \mathbf{m})^\top \mathbf{V}^{-1} (\mathbf{X} - \mathbf{m}) \right]$$

So, according to eq. 12.1, $T = 2$ and:

$$F(\mathbf{X}) = (\mathbf{X} - \mathbf{m})^\top \mathbf{V}^{-1} (\mathbf{X} - \mathbf{m})$$

Then, we have to define a starting vector \mathbf{X}_{sim} and variance-covariance matrix \mathbf{V}_{sim} . In order to show that the algorithm can converge from a point

chosen relatively far away from the optimum, we have chosen $\mathbf{X}_{sim} = 3\mathbf{m}$ and $\mathbf{V}_{sim} = \text{diag}(10V_{ii})$.

With the default initializations (10 cycles of 1000 simulations each), the results of a typical run were:

$$\hat{\mathbf{m}} = \begin{bmatrix} 1.01 \\ 2.02 \\ 3.01 \end{bmatrix} \quad \hat{\mathbf{s}} = \begin{bmatrix} 0.099 \\ 0.210 \\ 0.320 \end{bmatrix} \quad \hat{\mathbf{R}} = \begin{bmatrix} 1 & 0.286 & 0.467 \\ 0.286 & 1 & -0.299 \\ 0.467 & -0.299 & 1 \end{bmatrix}$$

12.5.8 Simulated Annealing

Program `simann.bas` uses simulated annealing to minimize a set of 10 notoriously difficult functions (most of them presenting multiple minima). Several successive runs of the program may be necessary to have all functions minimized (the random number generator being reinitialized at each call of the `SimAnn` subroutine).

12.5.9 Genetic Algorithm

Program `genalg.bas` optimizes the same functions than the previous program but with genetic algorithm. Here, too, it may be necessary to run the program several times.

Chapter 13

Statistics

This chapter describes some of the statistical functions available in **FBMath**. The specific problem of curve fitting will be considered in subsequent chapters.

13.1 Descriptive statistics

The following functions are available:

- Function **Mean**(**X**()) returns the mean of sample **X**, defined by:

$$m = \frac{1}{n} \sum_{i=1}^n x_i$$

where n is the size of the sample.

- Function **Median**(**X**(), **Sorted**) returns the median of **X**, defined as the number x_{med} which has equal numbers of values above it and below it. If the array **X** has been sorted, the median is:

$$x_{med} = x_{\frac{n+1}{2}} \quad (n \text{ odd})$$

$$x_{med} = \frac{1}{2} \left(x_{\frac{n}{2}} + x_{\frac{n}{2}+1} \right) \quad (n \text{ even})$$

The parameter **Sorted** indicates if array **X** has been sorted before calling function **Median**. If not, it will be sorted within the function (the array **X** will therefore be modified).

Sorting is performed by calling a subroutine `QSort(X())` which implements the ‘Quick Sort’ algorithm. Of course, this subroutine may be called outside function `Median`.

The default value of parameter `Sorted` is `True` (assuming that the array has already been sorted). So, `Median(X())` is equivalent to `Median(X(), True)`.

- Function `StDev(X(), M)` returns the estimated standard deviation of the population from which sample `X` is extracted, `M` being the mean of the sample. This standard deviation is defined by:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2}$$

These estimated standard deviations are used in statistical tests.

- Function `StDevP(X(), M)` returns the standard deviation of `X`, *considered as a whole population*. This standard deviation is defined by:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - m)^2}$$

- Function `Skewness(X(), M, Sigma)` returns the skewness of `X`, with mean `M` and standard deviation `Sigma`. This parameter is defined by:

$$\gamma_1 = \frac{1}{n\sigma^3} \sum_{i=1}^n (x_i - m)^3$$

Skewness is an indicator of the symmetric nature of the distribution. It is zero for a symmetric distribution (e. g. Gaussian), and positive (resp. negative) for an assymmetric distribution with a tail extending towards positive (resp. negative) x values.

- Function `Kurtosis(X(), M, Sigma)` returns the kurtosis of `X`, with mean `M` and standard deviation `Sigma`. This parameter is defined by:

$$\gamma_2 = \frac{1}{n\sigma^4} \sum_{i=1}^n (x_i - m)^4 - 3$$

Kurtosis is an indicator of the flatness of the distribution. It is zero for a Gaussian distribution, and positive (resp. negative) if the distribution is more (resp. less) ‘sharp’ than the Gaussian.

13.2 Comparison of means

13.2.1 Student's test for independent samples

We have 2 independent samples with sizes n_1, n_2 , means m_1, m_2 , standard deviations s_1, s_2 . It is assumed that the samples are taken from gaussian populations with means μ_1, μ_2 and equal variances. The sample means are compared by computing the t -statistic:

$$t = \frac{m_1 - m_2}{s \sqrt{1/n_1 + 1/n_2}}$$

where s^2 is the estimation of the common variance:

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

If $n_1 \geq 30$ and $n_2 \geq 30$, the conditions of normality and equal variances are no longer required and the formula begins:

$$t = \frac{m_1 - m_2}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

The null hypothesis is $(H_0) : \mu_1 = \mu_2$

The alternative hypothesis (H_1) depends on the test:

One-tailed test	$(H_1) : \mu_1 > \mu_2 \Rightarrow \text{reject } (H_0) \text{ if } t > t_{1-\alpha}$
	$(H_1) : \mu_1 < \mu_2 \Rightarrow \text{reject } (H_0) \text{ if } t < t_{1-\alpha}$
Two-tailed test	$(H_1) : \mu_1 \neq \mu_2 \Rightarrow \text{reject } (H_0) \text{ if } t > t_{1-\alpha/2}$

where $t_{1-\alpha}$ is the value of the Student variable such that the cumulative probability function $\Phi_\nu(t) = 1 - \alpha$ at $\nu = n_1 + n_2 - 2$ d.o.f. (cf. chap. 5).

If H_0 is rejected, the difference of the means is considered significant at risk α

This procedure is implemented in the following subroutine :

`StudIndep(N1, N2, M1, M2, S1, S2, T, DoF)`

where (N1, N2) are the sizes of the samples, (M1, M2) their means and (S1, S2) the estimated standard deviations (computed with `StDev`). The procedure returns Student's t in T and the number of degrees of freedom in DoF.

13.2.2 Student's test for paired samples

If the samples are paired (e. g. the same patients before and after a treatment), the t -statistic becomes:

$$t = \frac{m_d}{s_d} \sqrt{n}$$

where m_d and s_d are, respectively, the mean and standard deviations of the differences $(x_{1i} - x_{2i})$ between the paired values in the two samples, and n is the common size of the samples.

Apart from this, the test is carried out as with the independent case, with $(n - 1)$ d. o. f.

This procedure is implemented in the following subroutine :

`StudPaired(X(), Y(), T, DoF)`

where `X()`, `Y()` are the two samples. The procedure returns Student's t in `T` and the number of degrees of freedom in `DoF`.

After a call to this procedure, function `MathErr` returns one of the following error codes:

- `F0k` (0) if no error occurred
- `FSing` (-2) if $s_d = 0$
- `MatErrDim` (-3) if `X()` and `Y()` have different sizes.

13.2.3 One-way analysis of variance (ANOVA)

We have k independent samples with sizes n_i , means m_i , standard deviations s_i . It is assumed that the samples are taken from gaussian populations with means μ_i and equal variances. The goal is to compare the k means.

The following equation holds:

$$SS_t = SS_f + SS_r \quad (13.1)$$

with:

$$SS_t = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 \quad SS_f = \sum_{i=1}^k n_i (m_i - \bar{x})^2 \quad SS_r = \sum_{i=1}^k (n_i - 1) s_i^2$$

- \bar{x} is the global mean:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k n_i m_i \quad n = \sum_{i=1}^k n_i$$

- SS_t is the *total sum of squares*; it has $(n - 1)$ degrees of freedom
- SS_f is the *factorial sum of squares*; it has $(k - 1)$ degrees of freedom.
- SS_r is the *residual sum of squares*; it has $(n - k)$ degrees of freedom

Note that the degrees of freedom (d.o.f.) are additive, just like the sums of squares:

$$(n - 1) = (k - 1) + (n - k)$$

The *variances* are defined by dividing each sum of squares by the corresponding number of d.o.f.

$$V_t = \frac{SS_t}{n - 1} \quad V_f = \frac{SS_f}{k - 1} \quad V_r = \frac{SS_r}{n - k}$$

These are the total, factorial, and residual variances, respectively. Note that the variances, unlike the sum of squares, are *not* additive!

The comparison of means is performed by computing the F -statistic:

$$F = \frac{V_f}{V_r}$$

The null hypothesis is $(H_0) : \mu_1 = \mu_2 = \dots = \mu_k$

(H_0) is rejected if $F > F_{1-\alpha}$ where $F_{1-\alpha}$ is the value of the Fisher-Snedecor variable such that the cumulative probability function $\Phi_{\nu_1, \nu_2}(F) = 1 - \alpha$ at $\nu_1 = k - 1$ and $\nu_2 = n - k$ d.o.f. (cf. chap. 5).

This procedure is implemented in the following subroutine :

`AnOVa1(N(), M(), S(), V_f, V_r, F, DoF_f, DoF_r)`

where `N()` are the sizes of the samples, `M()` their means and `S()` the estimated standard deviations (computed with `StDev`). The procedure returns the factorial and residual variances in `V_f` and `V_r`, their ratio in `F` and their numbers of d. o. f. in `DoF_f` and `DoF_r`.

After a call to this procedure, function `MathErr` returns one of the following error codes:

- `F0k` (0) if no error occurred
- `FSing` (-2) if $n - k \leq 0$
- `MatErrDim` (-3) if the arrays have non-compatible dimensions

13.2.4 Two-way analysis of variance

We assume here that the means of the samples depend on two factors A and B, such that the sample corresponding to the i -th level of A and the j -th level of B has mean m_{ij} and standard deviation s_{ij} .

It is also assumed that all samples are taken from gaussian populations with equal variances, and that they have the same size n .

The previous equations become:

$$\bar{x} = \frac{1}{npq} \sum_{i=1}^p \sum_{j=1}^q nm_{ij}$$

$$SS_t = \sum_{i=1}^p \sum_{j=1}^q (x_{ij} - \bar{x})^2 \quad SS_f = \sum_{i=1}^p \sum_{j=1}^q n(m_{ij} - \bar{x})^2 \quad SS_r = \sum_{i=1}^p \sum_{j=1}^q (n-1)s_{ij}^2$$

with $npq - 1$, $pq - 1$, and $(n - 1)pq$ d.o.f., respectively.

In addition, the factorial sum of squares can be splitted into three terms:

$$SS_A = qn \sum_{i=1}^p (m_{i.} - \bar{x})^2 \quad ; \quad (p - 1) \text{ d.o.f.}$$

$$SS_B = pn \sum_{j=1}^q (m_{.j} - \bar{x})^2 \quad ; \quad (q - 1) \text{ d.o.f.}$$

$$SS_{AB} = n \sum_{i=1}^p \sum_{j=1}^q (m_{ij} - m_{i.} - m_{.j} + \bar{x})^2 \quad ; \quad (p - 1)(q - 1) \text{ d.o.f.}$$

where $m_{i.}$ and $m_{.j}$ are the conditional means:

$$m_{i.} = \frac{1}{q} \sum_{j=1}^q m_{ij} \quad m_{.j} = \frac{1}{p} \sum_{i=1}^p m_{ij}$$

that is, the means of the lines and columns of matrix $[m_{ij}]$

These sums of squares represent, respectively, the influence of factor A, the influence of factor B, and the interaction of the two factors (that is, the fact that the influence of one factor depends on the level of the other factor).

The variances are computed as before:

$$V_A = \frac{SS_A}{p - 1} \quad V_B = \frac{SS_B}{q - 1} \quad V_{AB} = \frac{SS_{AB}}{(p - 1)(q - 1)} \quad V_r = \frac{SS_r}{(n - 1)pq}$$

There are three null hypotheses:

- $(H_0)_A$: The populations means do not depend on factor A
- $(H_0)_B$: The populations means do not depend on factor B
- $(H_0)_{AB}$: There is no interaction between the two factors

Each hypothesis is tested by computing the corresponding F -statistic (for instance, $F_A = V_A/V_r$ for testing $(H_0)_A$) and comparing with the critical value $F_{1-\alpha}$

Special case: $n = 1$. If there is only one observation per sample, the residual variance is zero. The null hypotheses $(H_0)_A$ and $(H_0)_B$ are tested with $F_A = V_A/V_{AB}$ and $F_B = V_B/V_{AB}$. The interaction of the factors cannot be tested.

This procedure is implemented in the following subroutine :

AnOVa2(N, M(), S(), V(), F(), DoF())

where **N** is the common size of the samples, **M()** the matrix of means and **S()** the matrix of standard deviations, such that the lines correspond to factor A and the columns to factor B.

The procedure returns the variances in vector **V(1..4)** = $[V_A, V_B, V_{AB}, V_r]$, the variance ratios in **F(1..3)** = $[F_A, F_B, F_{AB}]$, and the degrees of freedom in **DoF(1..4)**. If $N = 1$, the last element of each vector disappears.

After a call to this procedure, function **MathErr** returns one of the following error codes:

- **F0k** (0) if no error occurred
- **MatErrDim** (-3) if the arrays have non-compatible dimensions.

13.3 Comparison of variances

13.3.1 Comparison of two variances

We have 2 independent samples with sizes n_1, n_2 , standard deviations s_1, s_2 . It is assumed that the samples are taken from gaussian populations with variances σ_1^2, σ_2^2 .

Snedecor's test uses the following statistic:

$$F = \frac{\max(s_1^2, s_2^2)}{\min(s_1^2, s_2^2)}$$

which is compared with the critical value $F_{1-\alpha/2}$ (two-tailed test).

This procedure is implemented in the following subroutine :

`Snedecor(N1, N2, S1, S2, F, DoF1, DoF2)`

where (N1, N2) are the sizes of the samples and (S1, S2) the estimated standard deviations. The procedure returns the variance ratio in F and the numbers of d. o. f. in DoF1 and DoF2.

13.3.2 Comparison of several variances

We have k independent samples with sizes n_i , standard deviations s_i . It is assumed that the samples are taken from gaussian populations with variances σ_i^2 . The goal is to compare the k variances.

Bartlett's test uses the following statistic:

$$B = \frac{1}{\lambda} \left[(n - k) \ln V_r - \sum_{i=1}^k (n_i - 1) \ln s_i^2 \right]$$

$$\lambda = 1 + \frac{1}{3(k - 1)} \left[\sum_{i=1}^k \frac{1}{n_i - 1} - \frac{1}{n - k} \right]$$

where $n = \sum n_i$ and V_r is the residual variance, as defined previously (§ 13.2.3).

The null hypothesis is:

$$(H_0) : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$$

Under (H_0) , B follows approximately the χ^2 distribution with $(k - 1)$ d. o. f. The hypothesis is tested by comparing B with the value $\chi_{1-\alpha}^2$ such that the cumulative probability function $\Phi_\nu(\chi^2) = 1 - \alpha$ at $\nu = k - 1$ d.o.f. (cf. chap. 5).

This procedure is implemented in the following subroutine :

`Bartlett(N(), S(), Khi2, DoF)`

where `N()` are the sizes of the samples and `S()` the estimated standard deviations. The procedure returns Bartlett's statistic in `Khi2` and the number of d. o. f. in `DoF`. The error codes are the same than for `AnOVA1`

13.4 Non-parametric tests

Non-parametric tests are used when the assumptions needed by the classical tests (gaussian populations with equal variances) are not fulfilled. They are also called *rank tests* because they work with the ranks of the values, rather than the values themselves.

13.4.1 Mann-Whitney test

This test compares the means of two independent samples. It is the non-parametric analog of Student's test for independent samples.

The test uses the following statistic:

$$U = \min(u_1, u_2)$$

with:

$$u_1 = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - r_1 \quad ; \quad u_2 = n_1 n_2 + \frac{n_2(n_2 + 1)}{2} - r_2$$

where (n_1, n_2) are the sample sizes, (r_1, r_2) the sums of the ranks of the two samples.

If $n_1 \geq 20$ and $n_2 \geq 20$, the variable:

$$\epsilon = \frac{U - \mu}{\sigma}$$

with:

$$\mu = \frac{n_1 n_2}{2} \quad ; \quad \sigma = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$$

follows the standard normal distribution under (H_0) .

This procedure is implemented in the following subroutine :

`Mann_Whitney(X1(), X2(), U, Eps)`

where `X1()`, `X2()` are the two samples. The procedure returns Mann-Whitney's statistic in `U` and the associated normal variable in `Eps`.

13.4.2 Wilcoxon test

This test compares the means of two paired samples. It is the non-parametric analog of Student's test for paired samples.

The test uses the following statistic:

$$T = \min(T_+, T_-)$$

where T_+ and T_- are the sums of the ranks of the positive and negative differences between the values of the two samples.

If the sample size is $N > 25$, the variable:

$$\epsilon = \frac{T - \mu}{\sigma}$$

with:

$$\mu = \frac{N(N+1)}{4} \quad ; \quad \sigma = \sqrt{\frac{N(N+1)(2N+1)}{24}}$$

follows the standard normal distribution under (H_0) .

This procedure is implemented in the following subroutine :

`Wilcoxon(X(), Y(), Ndiff, T, Eps)`

where `X()`, `Y()` are the two samples. The procedure returns the number of non-zero differences in `Ndiff`, Wilcoxon's statistic in `T` and the associated normal variable in `Eps`.

13.4.3 Kruskal-Wallis test

This test compares the means of several independent samples. It is the non-parametric analog of one-way ANOVA.

The test uses the following statistic:

$$H = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{r_i^2}{n_i} - 3(n+1)$$

where k is the number of samples, n_i the size of sample i , r_i the sum of the ranks for sample i and n the total size.

If $n_i > 5 \forall i$, H follows the χ^2 distribution with $k - 1$ d.o.f.

This procedure is implemented in the following subroutine :

`Kruskal_Wallis(N(), X(), H, DoF)`

where `N()` is the vector of sizes and `X()` the sample matrix (with the samples as columns). The procedure returns the Kruskal-Wallis statistic in `H` and the number of d. o. f. in `DoF`.

13.5 Statistical distribution

A statistical distribution is generated by binning data into a set of statistical classes $]x_i, x_{i+1}]$. Each class is characterized by the following parameters:

- its bounds x_i, x_{i+1}
- the number of values n_i contained in the class
- the frequency $f_i = n_i/N$ where N is the total number of values
- the density $d_i = f_i/(x_{i+1} - x_i)$

This structure is implemented in FBMath as:

```
TYPE StatClass
  Inf AS DOUBLE    ' Lower bound
  Sup AS DOUBLE    ' Upper bound
  N   AS INTEGER   ' Number of values
  F   AS DOUBLE    ' Frequency
  D   AS DOUBLE    ' Density
END TYPE
```

A distribution is generated with the following subroutine:

```
Distrib(X(), A, H, Ntot, C())
```

where $X()$ is the original set of values, A the lower bound of the first class and H the common width of the classes. The total number of values in the classes is returned in $Ntot$ and the distribution is returned in $C()$ which is an array of type `StatClass`.

13.6 Comparison of distributions

13.6.1 Observed and theoretical distributions

An observed distribution may be compared to a theoretical one by using the following statistics:

- Pearson's χ^2 :

$$\chi^2 = \sum_{i=1}^p \frac{(O_i - C_i)^2}{C_i}$$

- Woolf's G :

$$G = 2 \sum_{i=1}^p O_i \ln \frac{O_i}{C_i}$$

where O_i and C_i denote the observed and theoretical numbers of values in class i , and p the number of classes.

The null hypothesis is (H_0): the observed distribution conforms to the theoretical one (it is a test for conformity)

Under (H_0), both statistics follow the χ^2 distribution with $(p - 1 - N_e)$ d. o. f., where N_e is the number of parameters which have been estimated to compute the C_i values (e. g. $N_e = 2$ if the mean and standard deviation of the distribution have been estimated).

(H_0) is rejected if the chosen statistic is higher than the critical value $\chi^2_{1-\alpha}$ for the chosen risk α .

Pearson's statistic is an approximation of Woolf's statistic. It is usually recommended to use it only if $C_i \geq 5 \forall i$.

These procedures are implemented as:

```
Khi2_Conform(Obs(), Calc(), N_estim, Khi2, DoF)
```

```
Woolf_Conform(Obs(), Calc(), N_estim, G, DoF)
```

where `Obs()` and `Calc()` denote the observed and theoretical distributions, and `N_estim` the number of estimated parameters. The statistic is returned in `Khi2` or `G` and the number of d. o. f. in `DoF`.

13.6.2 Several observed distributions

To compare several observed distributions, we can group them into a *contingency table* **O** such that O_{ij} denotes the number of values for class i in the j -th distribution.

The Pearson and Woolf statistics may then be computed as:

$$\chi^2 = \sum_{i=1}^p \sum_{j=1}^q \frac{(O_{ij} - C_{ij})^2}{C_{ij}}$$

$$G = 2 \sum_{i=1}^p \sum_{j=1}^q O_{ij} \ln \frac{O_{ij}}{C_{ij}}$$

where p the number of classes, q the number of distributions, and C_{ij} the theoretical value of O_{ij} , computed as:

$$C_{ij} = \frac{N_{i.}N_{.j}}{N}$$

where $N_{i.}$ is the sum of terms in line i , $N_{.j}$ is the sum of terms in column j , and N the global sum of all terms in the matrix ($N = \sum_i N_{i.} = \sum_j N_{.j}$).

The null hypothesis is (H_0): the observed distributions come from the same population (it is a test for homogeneity or independence).

Under (H_0), both statistics follow the χ^2 distribution with $(p-1)(q-1)$ d. o. f.

These procedures are implemented as:

```
Khi2_Indep(Obs(), Khi2, DoF)
```

```
Woolf_Indep(Obs(), G, DoF)
```

where `Obs()` is the matrix of observed distributions. The statistic is returned in `Khi2` or `G` and the number of d. o. f. in `DoF`.

13.7 Demo programs

These programs are located in the `demo\stat` subdirectory of the `FBMath` directory.

13.7.1 Descriptive statistics, comparison of means and variances

Program `stat.bas` performs a statistical analysis of hemoglobin concentrations in two samples of 30 men and 30 women. The computed parameters are the mean, standard deviation, skewness and kurtosis. The means are compared by Student's test (two-tailed) and Mann-Whitney's test, and the variances are compared by Snedecor's test.

13.7.2 Student's test for paired samples

Program `student.bas` compares the means of two paired samples, using Student's and Wilcoxon's two-tailed tests.

13.7.3 One-way analysis of variance

Program `anova1.bas` compares the means of 5 independent samples, each with 12 observations, using one-way ANOVA and the Kruskal-Wallis test. In addition, the variances of the samples are compared with Bartlett's test.

13.7.4 Two-way analysis of variance

- Program `anova2.bas` compares the means of 4 samples, depending on two factors, using two-way ANOVA. Each sample contains 12 observations.
- Program `anova2a.bas` performs two-way ANOVA with one observation per sample.

13.7.5 Statistical distribution

Program `distrib.bas` uses the hemoglobin data from program `stat.bas` (men data) to generate a statistical distribution.

The first step determines a suitable range for the data. This is done by calling procedure `Interval` :

```
Interval X(1), X(N), 5, 10, XMin, XMax, XStep
```

The arguments 5 and 10 represent the minimal and maximal number of classes which is desired.

The second step generates the distribution, using the ranges determined in the previous step:

```
Nc = (Xmax - Xmin) / XStep  
DIM AS StatClass C(1 TO Nc)  
Distrib X(), Xmin, XStep, Ntot, C()
```

This distribution is then compared with the normal distribution, using both χ^2 and Woolf's tests. The theoretical C_i values are computed from the cumulative probability function for the normal distribution having the same mean and standard deviation than the observed distribution.

The program plots an histogram of the observed distribution, together with the curve corresponding to the normal distribution. This curve is generated from the probability density function:

```

FUNCTION PltFunc(X AS DOUBLE) AS DOUBLE
  PltFunc = DNorm((X - M) / S) / S
END FUNCTION

```

where M, S are the mean and standard deviation of the observed distribution, and DNorm is the probability density of the standard normal distribution (see chapter 5). Note that the histogram is constructed with the class densities as ordinates, so that a comparison with the plotted curve can be made.

13.7.6 Comparison of distributions

Program `khi2.bas` performs both χ^2 and Woolf's tests, first to compare an observed distribution with a theoretical one, and then to analyse a contingency table.

Chapter 14

Linear regression

This chapter describes the routines available in **FBMath** for fitting a straight line by linear regression. Other types of curve fitting will be described in subsequent chapters.

14.1 Straight line fit

The problem is to determine the equation of the line which comes closest to a set of points.

The model is defined by the equation:

$$y = a + bx$$

- x is the independent (or ‘explicative’) variable
- y is the dependent (or ‘explained’) variable
- a and b are the model parameters

Assume that the n points $(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)$ are perfectly lined, so that each of them verifies the equation of the straight line:

$$y_1 = a + bx_1$$

$$y_2 = a + bx_2$$

.....

$$y_n = a + bx_n$$

Or in matrix form:

$$\mathbf{y} = \mathbf{X}\beta \quad \Longleftrightarrow \quad \mathbf{y} - \mathbf{X}\beta = \mathbf{0}$$

where:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \quad \beta = \begin{bmatrix} a \\ b \end{bmatrix}$$

In the general case, the points are not exactly lined, so that:

$$\mathbf{y} - \mathbf{X}\beta = \mathbf{r}$$

where \mathbf{r} is the vector of residuals:

$$\mathbf{r} = [r_1, r_2 \dots r_n]^\top = \mathbf{y} - \hat{\mathbf{y}}$$

where $\hat{\mathbf{y}} = \mathbf{X}\beta$

It is possible to compute β so that $\|\mathbf{r}\|$ is minimal (*least squares criterion*).

$$\|\mathbf{r}\|^2 = \mathbf{r}^\top \mathbf{r} = r_1^2 + r_2^2 + \dots + r_n^2 = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = SS_r$$

where $\hat{y}_i = a + bx_i$ and SS_r is the *sum of squared residuals*.

Several methods allow the determination of β under the least squares criterion. The QR and SVD algorithms have been described previously. Here we will study the *method of normal equations*.

It may be shown that β is the solution of the system:

$$\mathbf{A}\beta = \mathbf{c}$$

with:

$$\mathbf{A} = \mathbf{X}^\top \mathbf{X} \quad \mathbf{c} = \mathbf{X}^\top \mathbf{y}$$

so:

$$\beta = (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y})$$

14.2 Analysis of variance

The following equation holds:

$$SS_t = SS_e + SS_r \quad (14.1)$$

with:

$$SS_t = \sum_{i=1}^n (y_i - \bar{y})^2 \quad SS_e = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad SS_r = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- \bar{y} is the mean of the y values:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

- SS_t is the *total sum of squares*; it has $(n - 1)$ degrees of freedom
- SS_e is the *explained sum of squares*; it has 1 degree of freedom.
- SS_r is the *residual sum of squares*; it has $(n - 2)$ degrees of freedom

Note that the degrees of freedom (d.o.f.) are additive, just like the sums of squares:

$$(n - 1) = 1 + (n - 2)$$

The *variances* are defined by dividing each sum of squares by the corresponding number of d.o.f.

$$V_t = \frac{SS_t}{n - 1} \quad V_e = \frac{SS_e}{1} \quad V_r = \frac{SS_r}{n - 2}$$

These are the total, explained, and residual variances, respectively. Note that the variances, unlike the sum of squares, are *not* additive!

The following quantities are derived from the above equations:

- the **coefficient of determination** r^2

$$r^2 = \frac{SS_e}{SS_t}$$

r^2 represents the percentage of the variations of y which are ‘explained’ by the independent variable. It is always comprised between 0 and 1. A value of 1 indicates a perfect fit.

- the **correlation coefficient** r

It is the square root of the coefficient of determination, with the sign of the slope b . It is therefore comprised between -1 and 1.

- the **residual standard deviation** s_r

It is the square root of the residual variance ($s_r = \sqrt{V_r}$). It is an estimate of the error made on the measurement of the dependent variable y . It should be 0 for a perfect fit.

- the **variance ratio** F

It is the ratio of the explained variance to the residual variance ($F = V_e/V_r$). It should be infinite for a perfect fit.

14.3 Precision of parameters

The matrix:

$$\mathbf{V} = V_r \cdot \mathbf{A}^{-1} = V_r \cdot (\mathbf{X}^\top \mathbf{X})^{-1}$$

is called the **variance-covariance matrix** of the parameters. It is a symmetric matrix with the following structure:

$$\mathbf{V} = \begin{bmatrix} \text{Var}(a) & \text{Cov}(a, b) \\ \text{Cov}(a, b) & \text{Var}(b) \end{bmatrix}$$

The diagonal terms are the variances of the parameters, from which the standard deviations are computed by:

$$s_a = \sqrt{\text{Var}(a)} \quad s_b = \sqrt{\text{Var}(b)}$$

The off-diagonal term is the covariance of the two parameters, from which the correlation coefficient r_{ab} is computed by:

$$r_{ab} = \frac{\text{Cov}(a, b)}{s_a s_b}$$

14.4 Probabilistic interpretation

It is assumed that the residuals ($y_i - \hat{y}_i$) are identically and independently distributed according to a normal distribution with mean 0 and standard deviation σ (estimated by s_r).

It may be shown that the regression parameters (a, b) are distributed according to a Student distribution with $(n - 2)$ d.o.f.

It is therefore possible to compute a confidence interval for each parameter, for instance:

$$\left[a - t_{1-\alpha/2} \cdot s_a \quad , \quad a + t_{1-\alpha/2} \cdot s_a \right]$$

where $t_{1-\alpha/2}$ is the value of the Student variable corresponding to the chosen probability α (usually $\alpha = 0.05$). This interval has a probability $(1 - \alpha)$ to contain the ‘true’ value of the parameter.

It is also possible to compute a ‘critical’ value $F_{1-\alpha}$ from the Fisher-Snedecor distribution with 1 and $(n - 2)$ d.o.f. The fit is considered satisfactory if the variance ratio F exceeds 4 times the critical value.

Note : for the straight line fit, $F_{1-\alpha} = \left(t_{1-\alpha/2} \right)^2$

14.5 Weighted regression

It is assumed here that the variance $v_i = \sigma_i^2$ of the measured value y_i is not constant.

The sums of squares become:

$$SS_t = \sum_{i=1}^n w_i (y_i - \bar{y})^2 \quad SS_e = \sum_{i=1}^n w_i (\hat{y}_i - \bar{y})^2 \quad SS_r = \sum_{i=1}^n w_i (y_i - \hat{y}_i)^2$$

where w_i denotes the ‘weight’, equal to $1/v_i$, and \bar{y} denotes the weighted mean:

$$\bar{y} = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i}$$

The regression parameters \mathbf{b} are estimated by:

$$\mathbf{b} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{W} \mathbf{y})$$

where \mathbf{W} is the diagonal matrix of weights:

$$\mathbf{W} = \text{diag}(w_1, w_2, \dots, w_n) = \begin{bmatrix} w_1 & 0 & \cdots & 0 \\ 0 & w_2 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & w_n \end{bmatrix}$$

The values of r_2 , s_r and F , as well as the variance-covariance matrix, are computed as above (§ 14.2). The normalized residual for the i -th observation is:

$$\frac{y_i - \hat{y}_i}{\sigma_i} = (y_i - \hat{y}_i)\sqrt{w_i}$$

These normalized residuals should follow the standard normal distribution.

14.6 Programming

14.6.1 Regression procedures

The following subroutines are available:

- `LinFit(X(), Y(), B(), V())` for unweighted linear regression
- `WLinFit(X(), Y(), S(), B(), V())` for weighted linear regression

The input parameters are:

- `X()`, `Y()` : coordinates of points
- `S()` : standard deviations of `Y` values (noted σ_i in paragraph 14.5)

The output parameters are:

- `B()` : regression parameters
- `V()` : inverse of the matrix of normal equations (noted \mathbf{A}^{-1} in paragraph 15.1.3). This is **not** the variance-covariance matrix. This one will be computed by the routines described in the next paragraph.

After a call to one of these procedures, function `MathErr` returns one of the following error codes:

- `MatOk` if no error occurred
- `MatSing` if the matrix of normal equations is quasi-singular

14.6.2 Quality of fit

The parameters used to test the quality of the fit are grouped in a user-defined type:

```
TYPE TRegTest
  Vr    AS DOUBLE    ' Residual variance
  R2     AS DOUBLE    ' Coefficient of determination
  R2a    AS DOUBLE    ' Adjusted coeff. of determination
  F      AS DOUBLE    ' Variance ratio (explained/residual)
  Nu1    AS INTEGER   ' D.o.f. of explained variance
  Nu2    AS INTEGER   ' D.o.f. of residual variance
END TYPE
```

They are computed by the following subroutines:

- `RegTest(Y(), Ycalc(), V(), Test)` for unweighted regression
- `WRegTest(Y(), Ycalc(), S(), V(), Test)` for weighted regression

The input parameters are:

- `Y()` : ordinates of points
- `Ycalc()` : Y values computed from the regression equation, using the fitted parameters `B()`. This computation must be done before calling `RegTest` or `WRegTest`.
- `V()` : the inverse matrix of the normal equations, as returned by the regression procedures.

The output parameters are:

- `V()` : the variance-covariance matrix of the fitted parameters
- `Test` : variable of type `TRegTest`, as defined above.

14.7 Demo programs

These programs are located in the `demo\curfit` subdirectory of the `FBMath` directory.

14.7.1 Unweighted linear regression

Program `reglin.bas` performs the least squares fit of a straight line, according to the following equation:

$$Y = B(0) + B(1) * X$$

The parameter vector and variance-covariance matrix are therefore declared as:

```
DIM AS DOUBLE B(0 TO 1), V(0 TO 1, 0 TO 1)
```

The program calls procedure `LinFit`, then computes the theoretical Y values:

```
FOR I = 1 TO N
    Ycalc(I) = B(0) + B(1) * X(I)
NEXT I
```

Note that this computation must be done before calling procedure `RegTest`

The critical values of Student's t and Snedecor's F are computed for the chosen probability `Alpha` by using the functions from chapter 5.

```
T = InvStudent(N - 2, 1 - 0.5 * Alpha)
F = InvSnedecor(1, N - 2, 1 - Alpha)
```

The output shows the standardized residuals, equal to $(y_i - \hat{y}_i)/\sigma$, where σ is estimated by s_r . They should be distributed according to the standard normal distribution.

14.7.2 Weighted linear regression

Program `wreglin.bas` performs the weighted least squares fit of a straight line. Here the standard deviations σ_i of the observed y values are stored in a vector `S()` defined by the user.

The computations involve the same steps as with the previous program, except that procedures `WLinFit` and `WRegTest` are used, and that the standardized residuals are computed as $(y_i - \hat{y}_i)/\sigma_i$

The plot shows the error bars, corresponding to $y_i \pm \sigma_i$ for each point.

Chapter 15

Multilinear regression and principal component analysis

This chapter describes the routines available in `FBMath` for multilinear regression, polynomial regression and principal component analysis.

15.1 Multilinear regression

15.1.1 Normal equations

The regression model is:

$$y = a + bx_1 + cx_2 + \cdots$$

where the x_i are m independent variables.

The method of normal equations, studied in chapter 14, is still applicable with:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1m} \\ 1 & x_{21} & x_{22} & \cdots & x_{2m} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix}$$

There are $p = m + 1$ parameters. The number of observations n must be such that $n > p$.

Special case: The x_i may be functions of another variable x , as long as these functions do not contain parameters.

Examples:

- Polynomial: $y = a + bx + cx^2 + \cdots$

- Fourier series: $y = a + b \sin x + c \sin 2x + \dots$

In such cases, the matrix \mathbf{X} , the matrix of normal equations $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$ and the constant vector $\mathbf{c} = \mathbf{X}^\top \mathbf{y}$ will have special forms. For instance with polynomial regression, if d is the degree of the polynomial:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ 1 & x_2 & x_2^2 & \dots & x_2^d \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} n & \Sigma x_i & \Sigma x_i^2 & \dots & \Sigma x_i^d \\ \Sigma x_i & \Sigma x_i^2 & \Sigma x_i^3 & \dots & \Sigma x_i^{d+1} \\ \dots & \dots & \dots & \dots & \dots \\ \Sigma x_i^d & \Sigma x_i^{d+1} & \Sigma x_i^{d+2} & \dots & \Sigma x_i^{2d} \end{bmatrix}$$

$$\mathbf{c} = \begin{bmatrix} \Sigma y_i \\ \Sigma x_i y_i \\ \dots \\ \Sigma x_i^d y_i \end{bmatrix}$$

It is possible to use these special forms to simplify the computations. For instance, only the first line and the last column of the above matrix \mathbf{A} need to be computed; the others terms are deduced by shifting.

15.1.2 Analysis of variance

Equation 14.1 still holds with the following modifications:

- the explained sum of squares SS_e has $(p - 1)$ degrees of freedom.
- the residual sum of squares SS_r has $(n - p)$ degrees of freedom

Note that the degrees of freedom are still additive:

$$(n - 1) = (p - 1) + (n - p)$$

The explained and residual variances become:

$$V_e = \frac{SS_e}{p - 1} \quad V_r = \frac{SS_r}{n - p}$$

The quantities r^2, s_r, F are derived as in § 14.1, but here the correlation coefficient r is always positive.

In multilinear regression, the use of r^2 may be misleading because it is always possible to artificially increase its value by adding more independent variables or using a higher degree polynomial. To overcome this drawback, the **adjusted coefficient of determination** may be used instead:

$$r_a^2 = 1 - (1 - r^2) \frac{n - 1}{n - p}$$

15.1.3 Precision of parameters

The variance-covariance matrix \mathbf{V} is computed as in chapter 14. It is a $p \times p$ symmetric matrix such that:

- the diagonal term V_{ii} is the variance of the i -th parameter
- the off-diagonal term V_{ij} is the covariance of the i -th and j -th parameters

The correlation coefficient r_{ij} is computed by:

$$r_{ij} = \frac{V_{ij}}{\sqrt{V_{ii}V_{jj}}}$$

15.1.4 Probabilistic interpretation

Assuming that the residuals are identically and independently distributed according to a normal distribution, the regression parameters are distributed according to a Student distribution with $(n - p)$ d.o.f. Confidence intervals may be computed as in chapter 14.

The ‘critical’ value $F_{1-\alpha}$ is computed from the Fisher-Snedecor distribution with $(p - 1)$ and $(n - p)$ d.o.f. However, the relationship $F_{1-\alpha} = \left(t_{1-\alpha/2}\right)^2$ does not hold if $p > 2$.

15.1.5 Weighted regression

Weighted multilinear regression may be performed as for the simple linear case (chap. 14).

15.1.6 Programming

The following subroutines are available:

- `MulFit(X(), Y(), B(), V())` for unweighted multilinear regression
- `WMulFit(X(), Y(), S(), B(), V())` for weighted multilinear regression
- `PolFit(X(), Y(), B(), V())` for unweighted polynomial regression
- `WPolFit(X(), Y(), S(), B(), V())` for weighted polynomial regression

where the parameters have the same meaning than in the simple linear case (chap. 14), except that `X()` is a matrix in multilinear regression.

After a call to one of these procedures, function `MathErr` returns one of the following error codes:

- `MatOk` if no error occurred
- `MatSing` if the matrix of normal equations is quasi-singular
- `MatErrDim` if the array dimensions do not match

15.2 Principal component analysis

15.2.1 Theory

The goal of Principal Component Analysis (PCA) is to replace a set of m variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$, which may be correlated, by another set $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_m$, called the *principal components* or *principal factors*. These factors are independent (uncorrelated) variables.

Usually, the algorithm starts with the *correlation matrix* \mathbf{R} which is a $m \times m$ symmetric matrix such that R_{ij} is the correlation coefficient between variable \mathbf{x}_i and variable \mathbf{x}_j .

The eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_m$ (in decreasing order) of matrix \mathbf{R} are the variances of the principal factors. Their sum $\sum_{i=1}^p \lambda_i$ is equal to m . So, the percentage of variance associated with the i -th factor is equal to λ_i/m .

If \mathbf{C} is the matrix of eigenvectors of \mathbf{R} , the correlation coefficient between variable \mathbf{x}_i and factor \mathbf{f}_j (sometimes called *loading*) is:

$$R_{C_{ij}} = C_{ij} \sqrt{\lambda_j}$$

The coordinates of the principal factors (sometimes called *scores*) are such that:

$$\mathbf{F} = \mathbf{ZC}$$

where \mathbf{Z} denotes the matrix of scaled original variables:

$$Z_{ij} = \frac{X_{ij} - m_j}{s_j}$$

where m_j and s_j are the mean and standard deviation of the j -th variable.

Note that the reduced variables have means 0 and variances 1, while the principal factors have means 0 and variances λ_i .

In most cases, a limited number of principal factors represent the most part of the total variance. It is therefore possible to neglect the other factors and to replace the m original (partially correlated) variables by a smaller set of independent variables. These variables can then be used in a regression analysis instead of the original ones (*orthogonalized regression*).

15.2.2 Programming

The following subroutines are available:

- `VecMean(X(), M())` computes the mean vector $\mathbf{M}()$ from matrix $\mathbf{X}()$.
 $\mathbf{M}(J)$ is the mean of the J -th column of the matrix.
- `VecSD(X(), M(), S())` computes the standard deviations $\mathbf{S}()$ from matrix $\mathbf{X}()$ and mean vector $\mathbf{M}()$.
 $\mathbf{S}(J)$ is the standard deviation of the J -th column of the matrix.
- `ScaleVar(X(), M(), S(), Z())` computes the scaled variables $\mathbf{Z}()$ from the original variables $\mathbf{X}()$, the means $\mathbf{M}()$ and the standard deviations $\mathbf{S}()$.
- `MatVarCov(X(), M(), V())` computes the variance-covariance matrix $\mathbf{V}()$ from matrix $\mathbf{X}()$ and mean vector $\mathbf{M}()$.
 $\mathbf{V}(I, J)$ is the covariance of the I -th and J -th column of $\mathbf{X}()$.
- `MatCorrel(V(), R())` computes the correlation matrix $\mathbf{R}()$ from the variance-covariance matrix $\mathbf{V}()$.

- `Pca(R(), Lambda(), C(), Rc())` performs the principal component analysis of the correlation matrix `R()`, which is destroyed. The eigenvalues are returned in vector `Lambda()`, the eigenvectors in the columns of matrix `C()`. The matrix `Rc()` contains the correlation coefficients (loadings) between the original variables (rows) and the principal factors (columns).
- `PrinFac(Z(), C(), F())` computes the principal factors (scores) `F()` from the scaled variables `Z()` and the matrix of eigenvectors `C()`.

After a call to these procedures, function `MathErr` returns one of the following error codes:

- `MatOk` if no error occurred
- `MatErrDim` if the array dimensions do not match
- `MatNonConv` if the iterative procedure (Jacobi method) did not converge in subroutine `PCA`

15.3 Demo programs

These programs are located in the `demo\curfit` subdirectory of the `FBMath` directory.

15.3.1 Multilinear regression

Program `regmult.bas` performs a multilinear least squares fit with `Nvar = 4` independent variables, according to the following equation:

$$Y = B(0) + B(1) * X1 + B(2) * X2 + B(3) * X3 + B(4) * X4$$

The data are stored in a matrix `X()` and a vector `Y()`. In agreement with the presence of a constant term `B(0)`, the first column of matrix `X()` is made of 1's.

The parameter vector and variance-covariance matrix are declared as:

```
DIM AS DOUBLE B(0 TO Nvar), V(0 TO Nvar, 0 TO Nvar)
```

The program calls procedure `MulFit`, then computes the theoretical `Y` values:


```

FOR I = 1 TO N
  Ycalc(I) = 0
  FOR J = 0 TO Nvar
    Ycalc(I) = Ycalc(I) + B(J) * X(I, J)
  NEXT J
NEXT I

```

Note that this computation must be done before calling procedure **RegTest**

The critical values of Student's t and Snedecor's F are computed for the chosen probability **Alpha** by using the functions from chapter 5.

```

T = InvStudent(Test.Nu2, 1 - 0.5 * Alpha)
F = InvSnedecor(Test.Nu1, Test.Nu2, 1 - Alpha)

```

where **Test.Nu1** and **Test.Nu2** are the numbers of d.o.f., returned by procedure **RegTest**.

The output shows the standardized residuals, equal to $(y_i - \hat{y}_i)/\sigma$, where σ is estimated by s_r . They should be distributed according to the standard normal distribution.

Due to the multi-dimensional nature of the relationship, a plot of y as a function of the x 's is not possible. Rather, the program plots a diagram of the observed and computed values of y , together with the theoretical line $\hat{y} = y$.

15.3.2 Polynomial regression

Program **regpoly.bas** performs a polynomial least squares fit. The structure of the program is very similar to the previous one, with the degree of the polynomial (**Deg**) playing the role of the number of variables (**Nvar**).

Here, only a vector **X()** is needed to store the values of the independent variable, since the powers of x are computed by the polynomial regression routine **PolFit**.

The theoretical Y values are computed by means of function **Poly**, studied in chapter 9.

The program plots the fitted curve by calling the plotting subroutine **PlotFunc**. The function which is passed to this subroutine is defined as:

```

FUNCTION PltFunc(X AS DOUBLE) AS DOUBLE
  PltFunc = Poly(X, B(), Deg)
END FUNCTION

```

The definition of procedure `PlotFunc` does not allow additional parameters for `PltFunc`. This is the only reason why the parameter vector `B()` is `DIM SHARED`.

15.3.3 Principal component analysis

There are three programs which use the same data set.

- Program `pca.bas` performs the principal component analysis of the data.
- Program `pcaplot1.bas` plots the loadings in the plane defined by the first two principal factors. The original variables appear as points in the plane (each variable is identified by its number). For each point, the coordinates (r_1, r_2) are the correlation coefficients of the corresponding original variable with the two principal factors. The program plots the *correlation circle*, defined by $r_1^2 + r_2^2 = 1$ (correlation with the first two factors only).
- Program `pcaplot2.bas` plots the scores in the plane defined by the first two principal factors. Each observation appears as a point in the plane.

Chapter 16

Nonlinear regression

This chapter describes the routines available in **FBMath** for fitting models which are nonlinear with respect to their parameters. For instance, the exponential model $y = ae^{-bx}$ is nonlinear with respect to the parameter b .

16.1 Theory

The regression model is:

$$y = f(x; a, b \dots)$$

where f is a nonlinear function of the parameters $a, b \dots$

Assume that we have a first estimate $(a^0, b^0 \dots)$ of the parameters. Let us write the Taylor series expansion of y in the vicinity of this estimate:

$$y = y^0 + y'_a \cdot (a - a^0) + y'_b \cdot (b - b^0) + \dots$$

where:

$$\begin{aligned} y^0 &= f(x; a^0, b^0 \dots) \\ y'_a &= \frac{\partial f}{\partial a}(x; a^0, b^0 \dots) \\ y'_b &= \frac{\partial f}{\partial b}(x; a^0, b^0 \dots) \\ &\dots\dots\dots \end{aligned}$$

The equation may be rewritten as:

$$y - y^0 = y'_a \cdot (a - a^0) + y'_b \cdot (b - b^0) + \dots$$

which corresponds to the linear regression problem:

$$\mathbf{z} = \mathbf{J} \cdot \delta$$

with:

$$\mathbf{z} = \begin{bmatrix} y_1 - y_1^0 \\ y_2 - y_2^0 \\ \dots \\ y_n - y_n^0 \end{bmatrix} \quad \mathbf{J} = \begin{bmatrix} y'_{a1} & y'_{b1} & \dots \\ y'_{a2} & y'_{b2} & \dots \\ \dots & \dots & \dots \\ y'_{an} & y'_{bn} & \dots \end{bmatrix} \quad \delta = \begin{bmatrix} a - a^0 \\ b - b^0 \\ \dots \end{bmatrix}$$

where \mathbf{J} is the *Jacobian matrix*, such that $y'_{ai} = \partial f(x_i; a^0, b^0 \dots) / \partial a$ etc.

Application of the linear regression relationships leads to:

$$\delta = (\mathbf{J}^\top \mathbf{J})^{-1} (\mathbf{J}^\top \mathbf{z}) \quad (16.1)$$

Knowing the correction vector δ , it is possible to compute better estimates a and b of the parameters. The process is repeated until convergence of the parameter estimates.

The method so described is known as the *Gauss-Newton* method. It is usually combined with nonlinear optimization, usually Marquardt's method, in order to minimize the sum of squared residuals:

$$SS_r = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \Phi(a, b \dots)$$

In this case, the gradient vector \mathbf{g} and hessian matrix \mathbf{H} of function Φ are computed by the following relationships:

$$\mathbf{g} = -\mathbf{J}^\top \mathbf{z} \quad \mathbf{H} = \mathbf{J}^\top \mathbf{J} \quad (16.2)$$

so that the Gauss-Newton formula (16.1) becomes equivalent to the Newton-Raphson formula for nonlinear optimization (p. 43).

Note that, in the previous formula:

1. \mathbf{g} and \mathbf{H} are scaled by a factor 1/2 since this factor cancels during the computations.
2. The expression of \mathbf{H} is only approximate, since a factor containing the term $(y_i - \hat{y}_i)$ is neglected during the computation of the second partial derivatives:

$$\frac{\partial^2 \Phi}{\partial a \partial b} = \sum_{i=1}^n \left[\frac{\partial \hat{y}_i}{\partial a} \frac{\partial \hat{y}_i}{\partial b} - (y_i - \hat{y}_i) \frac{\partial^2 \hat{y}_i}{\partial a \partial b} \right]$$

The residual variance is:

$$V_r = \frac{SS_r}{n - p}$$

where p is the number of parameters in the model.

It is still possible to compute r^2 and F , as well as confidence intervals, but their interpretation is less straightforward since the ANOVA relationship (§ 14.1) does not hold for nonlinear models. In this case, r^2 may be > 1 ! Moreover, the distribution of the parameters is only approximately described by the Student distribution.

16.2 Monte-Carlo simulation

The distribution of the regression parameters may be simulated by the MCMC method discussed in § 12.2 p. 72).

Let β denote the vector of model parameters. According to Bayes' theorem, the posterior probability density $\pi(\beta)$ of these parameters is given by:

$$\pi(\beta) = \frac{L(\beta)P(\beta)}{\int L(\beta)P(\beta)d\beta} = \frac{L(\beta)P(\beta)}{N}$$

where $P(\beta)$ denotes the prior probability density of the parameters and $L(\beta)$ denotes the likelihood, i.e. the probability of observing the experimental results (x_i, y_i) given the parameters.

The integral which appears in the denominator is usually too complex to be calculated and is therefore treated as a normalizing constant N .

Assuming that, for a given β , the residuals $(y_i - \hat{y}_i)$ are identically and independently distributed according to a normal distribution with variance σ^2 , the likelihood is given by:

$$L(\beta) = \prod_{i=1}^n \left(\frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2} \right] \right)$$

If we choose a uniform prior probability $P(\beta)$ over an interval \mathcal{B} , the posterior probability becomes:

$$\pi(\beta) = C \prod_{i=1}^n \exp \left[-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2} \right]$$

where C is a constant.

In order to use the Metropolis-Hastings algorithm, as described in chapter 12.2, we define the function:

$$F(\beta) = \begin{cases} -2 \ln \frac{\pi(\beta)}{C} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 & \text{if } \beta \in \mathcal{B} \\ \infty & \text{otherwise} \end{cases} \quad (16.3)$$

It is the same objective function than for the nonlinear regression algorithm, except that it is bounded on the interval \mathcal{B} .

16.3 Demo programs

These programs are located in the `demo\curfit` subdirectory.

16.3.1 Nonlinear regression

Program `regnlin.bas` performs a nonlinear least squares fit of the exponential model:

$$y = ae^{-bx}$$

which is coded as:

```
Y = B(1) * EXP(- B(2) * X)
```

The partial derivatives used to compute the Jacobian are:

$$\frac{\partial y}{\partial a} = e^{-bx} \quad \frac{\partial y}{\partial b} = -axe^{-bx}$$

Initial estimates of the parameters `B()` are obtained by linearization:

$$\ln y = \ln a - bx$$

However, this transformation modifies the standard deviations of the independent variables:

$$\sigma(\ln y) \approx d \ln y = \frac{dy}{y} \approx \frac{\sigma(y)}{y}$$

It is therefore recommended to use weighted linear regression for this step.

Subroutine **ApproxFit** selects the data points for which the transformation is appropriate (i. e. $y > 0$) and stores the transformed coordinates and standard deviations in 3 vectors `X1()`, `Y1()`, `S1()` which are passed to the weighted linear regression subroutine **WLinFit**. The results are returned in the global arrays `B()` and `V()`. The parameters are transformed back to the original form of the model:

```

B(1) = EXP(B(1))
B(2) = - B(2)

```

Marquardt's method is then used to perform nonlinear minimization of the residual sum of squares. Subroutine **Marquardt** needs two other procedures:

- a function **ObjFunc(B())** which computes the objective function to be minimized
- a subroutine **HessGrad(B(), G(), H())** which computes the Gradient and Hessian of the objective function

Since the parameter lists of these procedures cannot be modified, the other variables which they must access (such as the point coordinates **X()**, **Y()**) are made global by means of **DIM SHARED** statements.

The results of the minimization are printed as with the linear regression programs, except that the correlation coefficients are shown only if $r \leq 1$.

The program may be adapted to another regression model by changing the following parts:

- the function name (constant **FuncName**)
- the constants **FirstParam** and **LastParam** which define the bounds of the parameter array **B()**
- the subroutine **ApproxFit** which computes the initial estimates of the parameters
- the definition of the regression model in functions **ObjFunc** and **PltFunc**
- the definition of derivatives in subroutine **HessGrad**

16.3.2 Monte-Carlo simulation

Program **mcsim.bas** simulates the posterior distribution of the regression parameters for the previous exponential model. The interval \mathcal{B} must be defined by the user. The objective function is defined as with the previous program, except that it is set to a high value if one of the parameters goes outside the bounds. When the simulation is done, the program draws a plot showing the distribution of the parameters.

Chapter 17

String functions

Some string functions have been added to **FBMath**, mainly to help printing results.

17.1 Fill functions

- function **RFill(S, L)** returns string **S** completed with trailing blanks for a total length **L**
- function **LFill(S, L)** returns string **S** completed with leading blanks for a total length **L**
- function **CFill(S, L)** returns string **S** completed with leading blanks so as to center the string on a total length **L** (80 characters by default)

17.2 Character replacement

Subroutine **Replace(S, C1, C2)** replaces in string **S** all the occurrences of character **C1** by character **C2**

17.3 Parsing

Subroutine **Parse(S, Delim, Elem(), N)** parses string **S** into its constitutive fields (separated by **Delim**). The fields are returned in array **Elem()**. The number of fields is returned in **N**.

The array **Elem()** must be dimensioned by the calling program. If it is not large enough, the number of extracted fields will be limited to the size of the array.

17.4 Formatting functions

These functions are based on the built-in function `Format`, which is declared in the include file `string.bi`. This file is automatically added each time you include `math.bi`.

- function `FloatToStr(X, Ntot, Ndec, E)` converts the double precision number `X` to a string of `Ntot` characters, with `Ndec` decimals. The boolean parameter `E` must be set to `True` (-1) to get exponential notation (e. g. `6.626E-034`).

The default values are: `Ntot = 10`, `Ndec = 4` and `E = False`, so that, for instance, `FloatToStr(Pi)` will return the string `3.1416` with 4 leading blanks.

- function `RemZero(S)` will remove the non-significant zeroes from a string `S` which represents a number (usually the output of `FloatToStr`).
e.g. `RemZero("1.2300")` will return `1.23` and `RemZero("1.2300E+00")` will return `1.23E+00`
- function `IntToStr(N, Ntot)` converts the integer `N` to a string of `Ntot` characters. Default is `Ntot = 10`.