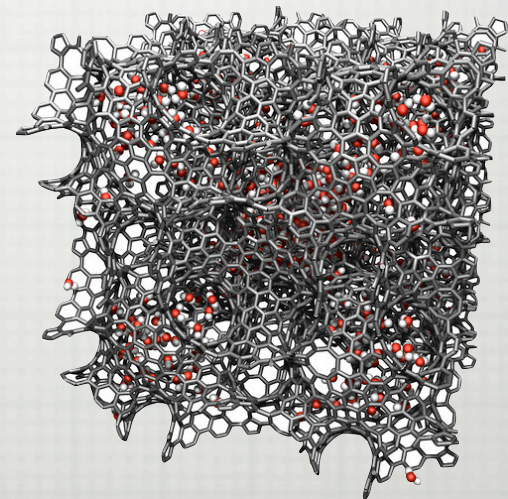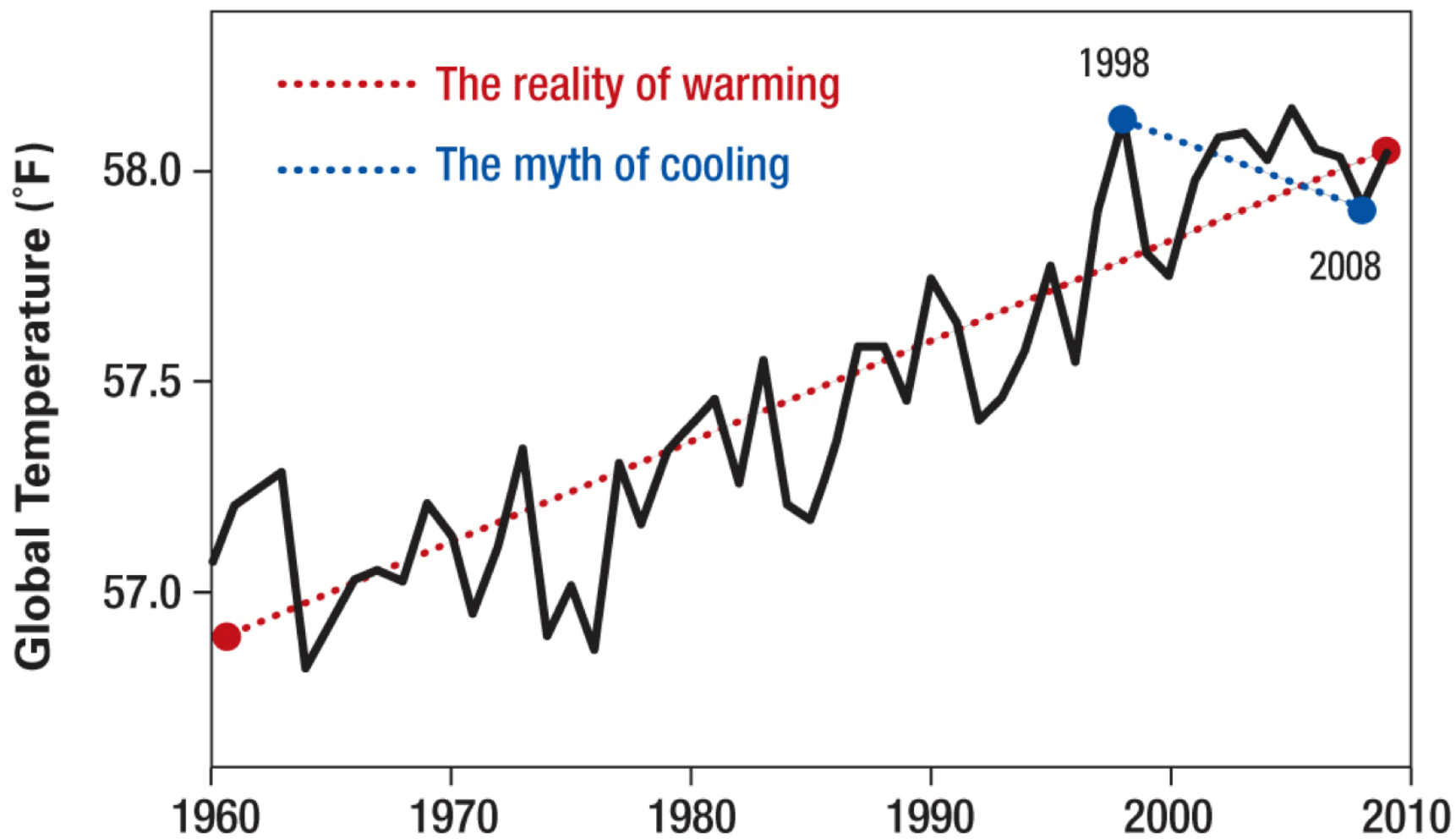# PHY-4810
# COMPUTATIONAL PHYSICS
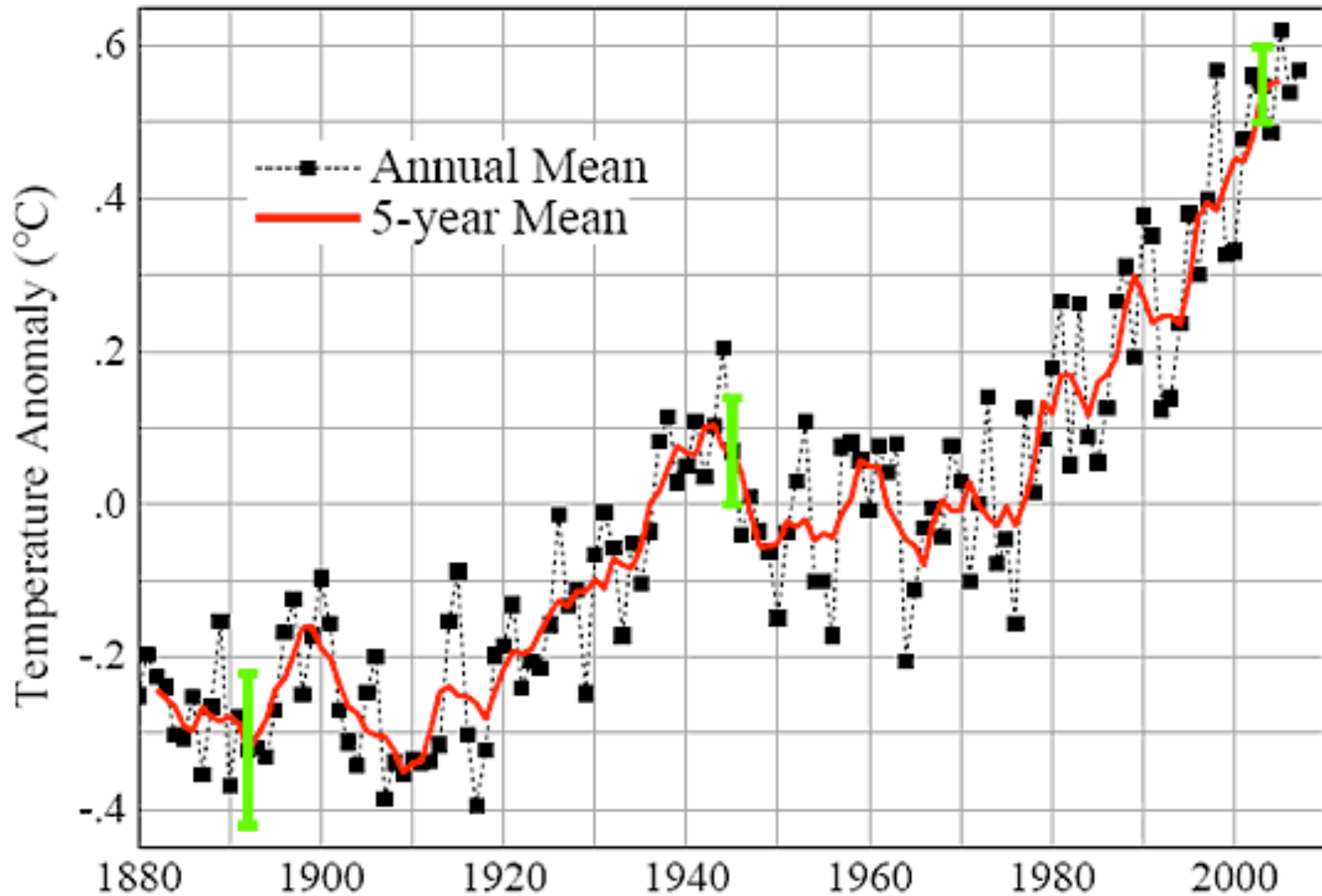
## LECTURE 7: DATA FITTING

**Global Surface Temperature Over the Last Fifty Years**

HTTP://DATA.GISS.NASA.GOV/GISTEMP/TABLEDATA/GLB.TS+DSST.TXT

# GLOBAL LAND-OCEAN TEMPERATURE INDEX IN 0.01 DEGREES CELSIUS

Global Temperature Land-Ocean Index

# ON THE MENU

- ☐ **Lagrange multipliers**

- ☐ **Spline**

- ☐ **Fit to theory**

- ☐ **How to use Numerical Recipes**

# PART 1: LAGRANGE INTERPOLATION

# METHOD AND ALGORITHM

☐ **Any function g(x) can be approximated as a polynomial of degree (n − 1) in each interval i for a sufficiently high *n*:**

$$g_i(x) \simeq a_0 + a_1 x + a_2 x^2 + \cdots + a_{n-1} x^{n-1}$$

☐ **Because our fit is <u>local</u>, we do not assume that one g(x) can fit all the data in a list, but instead will use a different polynomial, that is, a different set of <u>a</u>i values, for each region of the list.**

# LAGRANGE INTERPOLATION FORMULA

☐ **Lagrange figured out a closed-form approach that directly fits the *(n – 1)*-order polynomial to *n* values of the function *g(x)* evaluated at the points $x_i$. The formula is written as the sum of polynomials:**

$$g(x) \simeq g_1\lambda_1(x) + g_2\lambda_2(x) + \cdots + g_n\lambda_n(x)$$

$$\lambda_i(x) = \prod_{j(\neq i)=1}^{n} \frac{x - x_j}{x_i - x_j} = \frac{x - x_1}{x_i - x_1} \frac{x - x_2}{x_i - x_2} \cdots \frac{x - x_n}{x_i - x_n}$$

☐ **For <u>three</u> points, this formula provides a <u>second-degree polynomial</u>, while for <u>eight</u> points it gives a <u>seventh-degree</u> polynomial.**

# EXAMPLE: POLYNOMIAL OF ORDER 3



| $x_i$ | 0 | 2 | 5 | 7 |
|---|---|---|---|---|
| $g(x_i)$ | 9 | 12 | 14 | 8 |

$$g(x) \simeq g_1\lambda_1(x) + g_2\lambda_2(x) + \cdots + g_n\lambda_n(x)$$

$$\lambda_i(x) = \prod_{j(\neq i)=1}^{n} \frac{x - x_j}{x_i - x_j} = \frac{x - x_1}{x_i - x_1} \frac{x - x_2}{x_i - x_2} \cdots \frac{x - x_n}{x_i - x_n}$$



| $x_i$ | 0 | 2 | 5 | 7 |
|-------|---|----|----|---|
| $g(x_i)$ | 9 | 12 | 14 | 8 |

$$g(x) = g_1 \lambda_1(x) + g_1 \lambda_2(x) + g_1 \lambda_3(x)$$

$$\lambda_1(x) = \frac{(x-2)(x-5)(x-7)}{-2 \times -5 \times -7}$$

$$\lambda_2(x) = \frac{x(x-5)(x-7)}{2 \times -3 \times -5}$$

$$\lambda_3(x) = \frac{x(x-2)(x-7)}{5 \times 3 \times -2}$$

$$\lambda_4(x) = \frac{x(x-2)(x-5)}{7 \times 5 \times 2}$$



$$g(x) = \frac{-9}{70}(x-2)(x-5)(x-7) + \frac{2}{5}x(x-5)(x-7) + \frac{-7}{15}x(x-2)(x-7) + \frac{4}{35}x(x-2)(x-5)$$

# …HOWEVER…RESIST THE TEMPTATION… (INTERPOLATION VS. EXTRAPOLATION)

# 10 DATA POINTS: POLYNOMIAL OF ORDER 9

# POLYNOMIAL FIT

# LET'S HAVE A LOOK AT THE C++ CODE

```cpp
#INCLUDE <IOSTREAM>
USING NAMESPACE STD;

INT MAIN () {
    // INSERT CODE HERE...
    INT I, J, K;
    INT KMAX=393;
    DOUBLE XX=0, GXX=0;
    DOUBLE XXMIN, XXMAX;
    //
    INT NPOINTS;
    CIN >> NPOINTS;
    //
    DOUBLE X[NPOINTS], G[NPOINTS];
    DOUBLE LAMBDA[NPOINTS];

    //
    FOR (I=0;I<NPOINTS;I++){
        CIN >> X[I] >> G[I];
    }
    //
    XXMIN=X[0];
    XXMAX=X[NPOINTS-1];
    //
    FOR(K=0;K<KMAX;K++){
        XX=XXMIN+(DOUBLE)K*(XXMAX-XXMIN)/((DOUBLE)KMAX-1);
        GXX=0.0;
        FOR (I=0;I<NPOINTS;I++){
            LAMBDA[I]=1.;
            FOR (J=0;J<NPOINTS;J++){
                IF(J==I) CONTINUE;
                LAMBDA[I]*=(XX-X[J])/(X[I]-X[J]);
            }
            GXX+=LAMBDA[I]*G[I];
        }
        COUT << XX <<  " " << GXX << ENDL;
    }
    RETURN 0;
}
```

# EXAMPLE OF POLYNOMIAL OBTAINED USING NEVILLE'S ALGORITHM

# NEVILLE'S ALGORITHM

- A much better algorithm to construct the (unique) polynomial is called Neville's algorithm

```cpp
//NEVILLE's ALGORITHM
#include <iostream>
using namespace std;

int main () {
    int i,k,m;
    int kmax=50;
    double xx=0, gxx=0;
    double xxmin, xxmax;
    double sum;
    //
    int npoints;
    cin >> npoints;
    //
    double x[npoints], g[npoints];
    double P0[npoints], P[npoints];


    //
    for (i=0;i<npoints;i++){
            cin >> x[i] >> g[i];
    }
    //
    xxmin=x[0];
    xxmax=x[npoints-1];
    //
    for(k=0;k<kmax;k++){
            xx=xxmin+(double)k*(xxmax-xxmin)/((double)kmax-1);
            gxx=0.0;
            sum=0;

            for(m=0;m<npoints;m++){

                if(m==0){
                        for (i=0;i<npoints;i++){
                                P[i]=g[i];
                        }
                }
                else {
                        for (i=0;i<npoints-m;i++){
                                P[i]=((xx-x[i+m])*P0[i]+(x[i]-xx)*P0[i+1])/(x[i]-x[i+m]);
                        }
                }
                for (i=0;i<npoints-m;i++){
                        P0[i]=P[i];
                }

            }
            cout << xx <<  " " << P[0] << endl;
    }

    return 0;
}
```

# PRACTICAL CONSIDERATIONS

☐ **If the data contain little noise, this polynomial can be used with some confidence within the range of data, but with risk beyond the range of data.**

☐ **Notice that Lagrange interpolation makes no restriction that the points in the table be evenly spaced.**

☐ **As a check, it is also worth noting that the sum of the Lagrange multipliers equals one, $\sum \lambda_i = 1$.**

☐ **Usually the Lagrange fit is made to only a small region of the table with a small value of n, even though the formula works perfectly well for fitting a high-degree polynomial to the entire table.**

# LOCAL INTERPOLATION

☐ **Here, we interpolate the function, small interval by small interval, with a low order polynomial (here: 3-points or parabola)**

# 4-POINT SEGMENTATION
# (THIRD ORDER POLYNOMIAL)

# PROBLEMS: DISCONTINUITY IN DERIVATIVE!

# PART2:
# SPLINE

# SPLINE: INTRODUCTION

- [ ] **Fitting parabolas (three-point interpolation) within a table may avoid the erroneous and possibly catastrophic deviations of a high-order formula.**

- [ ] **(Two-point interpolation, which connects the points with straight lines, may not lead you far astray, but it is rarely pleasing to the eye or precise.)**

- [ ] **A sophisticated variation of n = 4 interpolation, known as cubic splines, often leads to surprisingly eye-pleasing fits.**

# SPLINE: INTRODUCTION (II)

☐ **In this approach, cubic polynomials are fit to the function in each interval, with the additional constraint that the first and second derivatives of the polynomials must be continuous from one interval to the next. This continuity of slope and curvature is what makes the spline fit particularly eye-pleasing.**

# SPLINE: C'TD

- The series of cubic polynomials obtained by spline-fitting a table can be integrated and differentiated, and is guaranteed to have well-behaved derivatives.

- The complexity of simultaneously matching polynomials and their derivatives over all the interpolation points leads to many simultaneous, linear equations to be solved.

- This makes splines unattractive for hand calculations, yet easy for computers.

- The basic approximation of splines is the representation of the function g(x) in the subinterval [$x_i$, $x_{i+1}$] with a cubic polynomial:

$$g(x) \simeq g_i(x) \qquad \text{for } x_i \leq x \leq x_{i+1}$$

$$g_i(x) = g_i + g_i'(x - x_i) + \tfrac{1}{2}g_i''(x - x_i)^2 + \tfrac{1}{6}g_i'''(x - x_i)^3$$

- **This representation makes it clear that the coefficients in the polynomial equal the values of g(x) and to its first, second, and third derivatives at the tabulated points $x_i$. Derivatives beyond the third vanish.**

  - The computational chore is to determine these derivatives in terms of the N tabulated values $g_i$.

  - The matching of $g_i$ from one interval to the next (at the nodes) provides the equations

$$g_i(x_{i+1}) = g_{i+1}(x_{i+1}) \qquad i = 1, N - 1$$

  - The matching of the first and second derivatives at each subinterval's boundary provides the equations

$$g_{i-1}'(x_i) = g_i'(x_i) \qquad g_{i-1}''(x_i) = g_i''(x_i)$$

To provide the additional equations needed to determine all constants, the third derivatives at adjacent nodes are matched. Values for the third derivatives are found by approximating them in terms of the second derivatives (our old FD approximation):

$$g_i''' \simeq \frac{g_{i+1}'' - g_i''}{x_{i+1} - x_i}$$

It is straightforward though complicated to solve for all the parameters.

# BOUNDARIES

**Matching at the boundaries of the intervals results in only N − 2 linear equations for N unknowns.** **Further input is required.**

It usually is taken to be the boundary conditions at the endpoints $a = x_1$ and $b = x_N$, specifically, the second derivatives $g''(a)$, and $g''(b)$.

● <u>Natural spline:</u>   Set $g''(a) = g''(b) = 0$, (This is "natural" because the derivative vanishes for the flexible spline drafting tool (its ends being free).)

● <u>Input values for g' at boundaries</u>:   The computer uses $g'(a)$ to approximate $g''(a)$. If you do not know the first derivatives, you can calculate them numerically from the table of gi values.

● <u>Input values for g'' at boundaries</u>:   Knowing values is of course better than assuming values, but it requires more input. If the values of $g''$ are not known, they can be approximated by applying a forward-difference approximation to the tabulated values

**HARD WORK**

NUMERICAL RECIPES
The Art of Scientific Computing
THIRD EDITION

William H. Press
Saul A. Teukolsky
William T. Vetterling
Brian P. Flannery

Includes CD-ROM
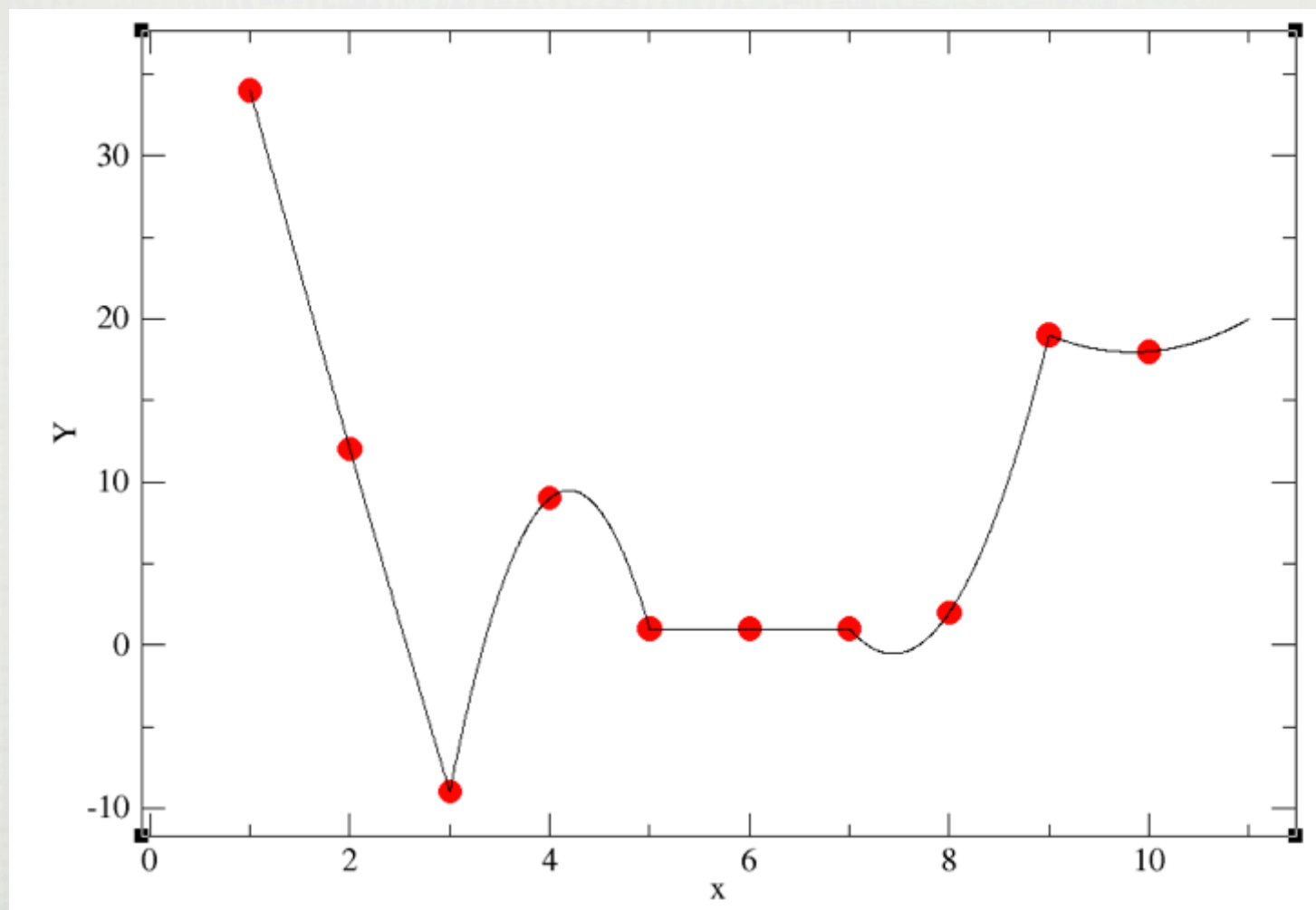
```cpp
#include <iostream>
#include "nr3.h"
#include "interp_1d.h"
using namespace std;
int main () {
    int i,k;
    int kmax=50;
    double xx=0, gxx=0;
    double xxmin, xxmax;
    //
    int npoints;

//GET INPUT
    cin >> npoints;
//allocate memory
    VecDoub x(npoints), g(npoints);
//read in data points
    for (i=0;i<npoints;i++){
        cin >> x[i] >> g[i];
    }
//CREATE SPLINE OBJECT
    Spline_interp myspline (x,g);

//CREATE INTERPOLATION
    xxmin=x[0];
    xxmax=x[npoints-1];
      //
    for(k=0;k<kmax;k++){
        xx=xxmin+(double)k*(xxmax-xxmin)/((double)kmax-1);
        gxx= myspline.interp(xx);
        cout << xx << " " << gxx <<endl;
    }
}
```
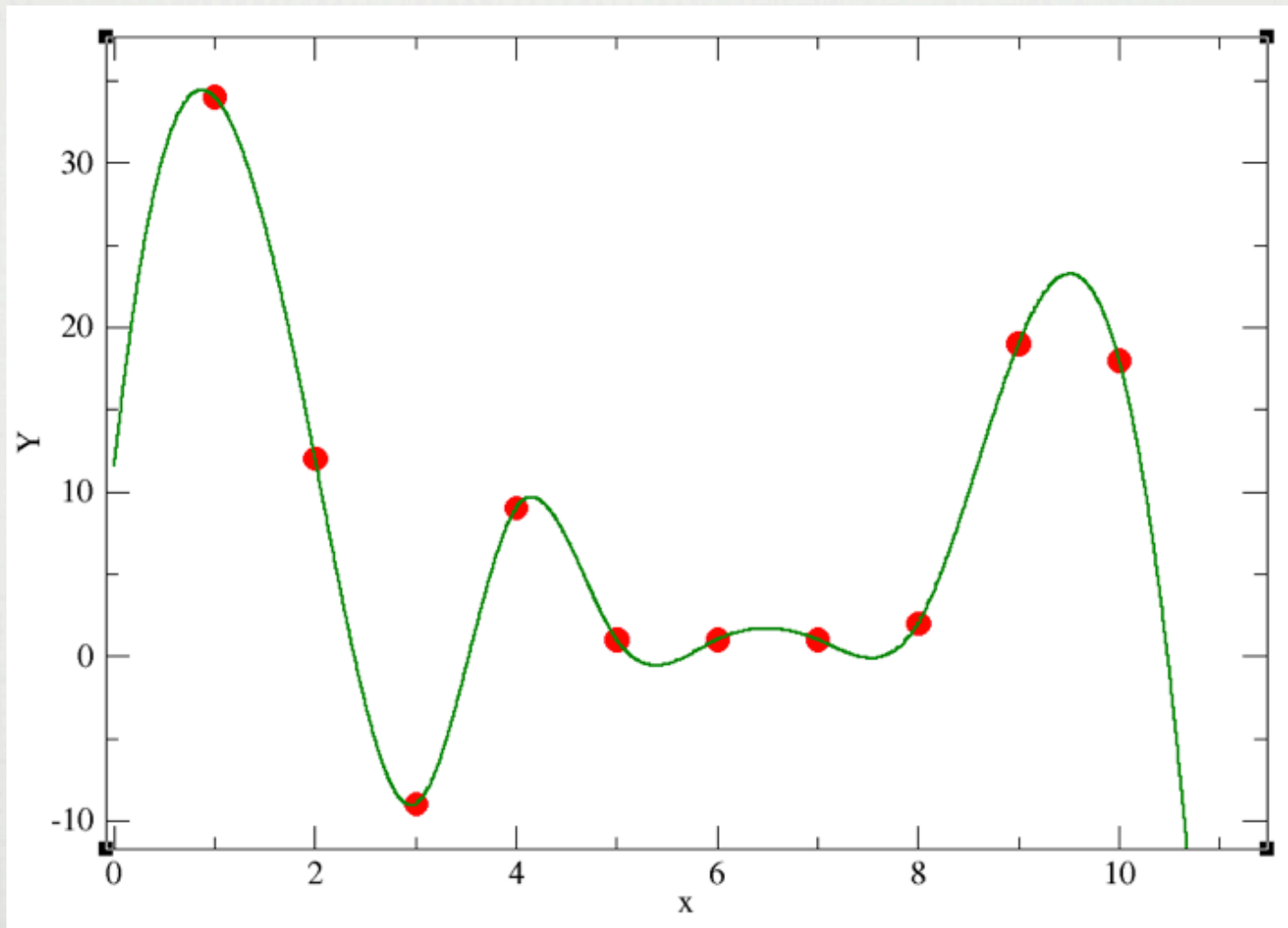
**SPLINE WITH NUMREC**

SPLINE!

# USING SPLINES FOR INTEGRATION

☐ A powerful integration scheme is to fit an integrand with splines, and then integrate the cubic polynomials analytically.

☐ If you have the ability to actually calculate the function for arbitrary x, Gaussian quadrature may be preferable. We know that the spline fit to g in each interval is the cubic (9.8),

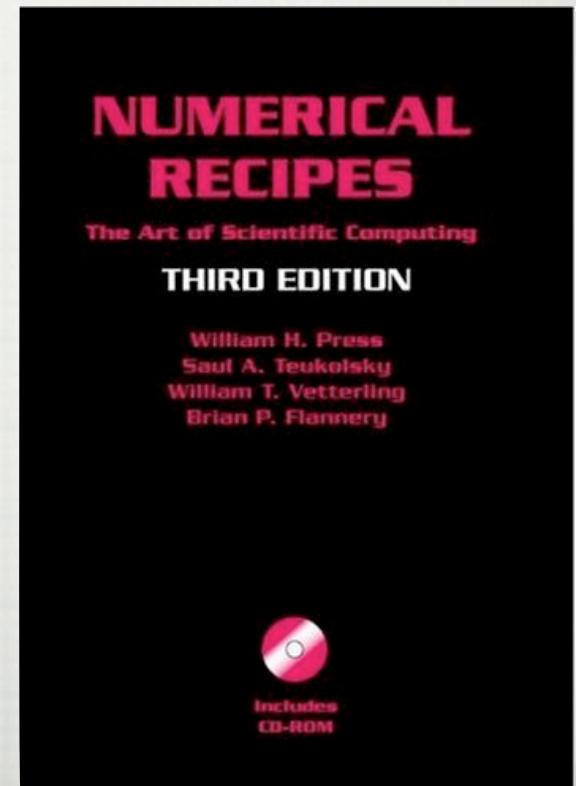$$g(x) \simeq g_i + g_i'(x - x_i) + \tfrac{1}{2}g_i''(x - x_i)^2 + \tfrac{1}{6}g_i'''(x - x_i)^3$$

☐ It is easy to integrate this to obtain the integral of g for this interval, and then to sum over all intervals:

$$\int_{x_i}^{x_{i+1}} g(x)dx \simeq \left( g_i x + \tfrac{1}{2}g_i' x_i^2 + \tfrac{1}{6}g_i'' x^3 + \tfrac{1}{24}g_i''' x^4 \right) \Big|_{x_i}^{x_{i+1}}$$

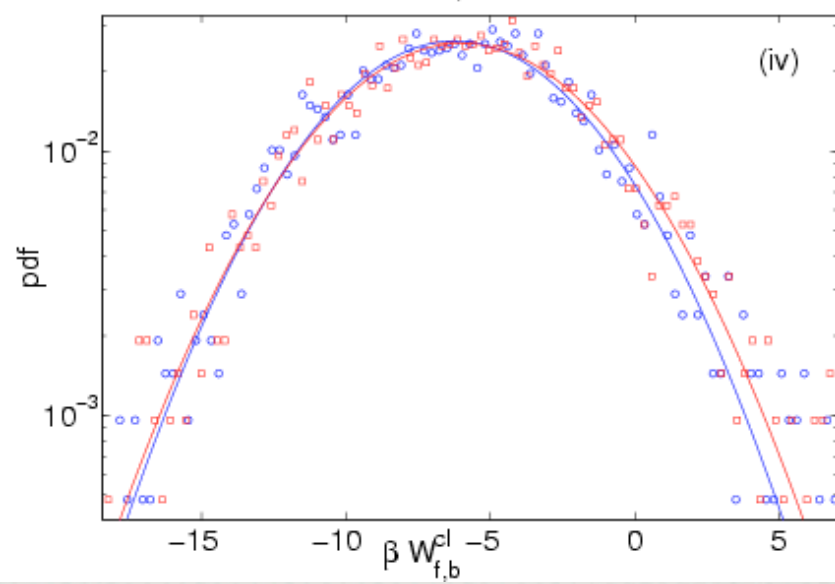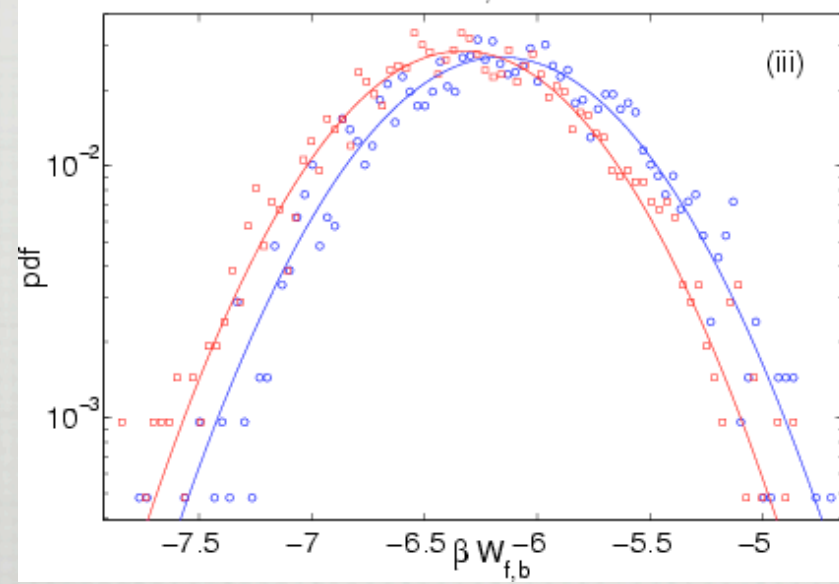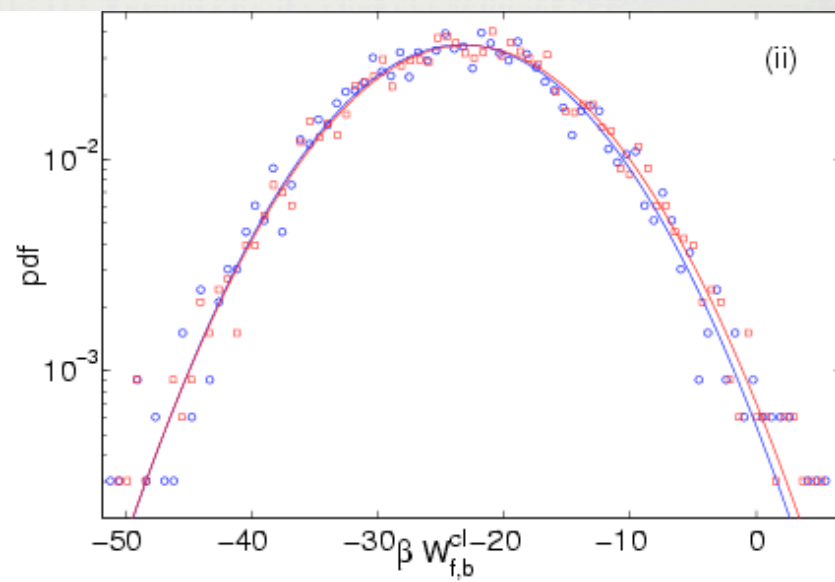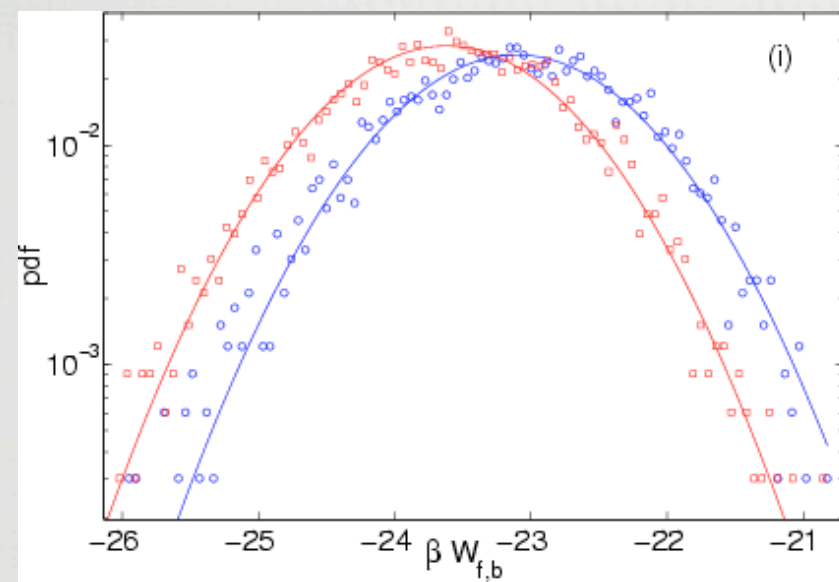$$\int_{x_j}^{x_k} g(x)dx = \sum_{i=j}^{k} \int_{x_i}^{x_{i+1}} g(x)dx$$

☐ Making the intervals smaller does not necessarily increase precision as subtractive cancellations may get large.

# HOW DO WE USE THESE METHODS?

☐ **Preferred method: using numerical libraries**

☐ **Go online to nr.com or look in a book...or download the routines...**

☐ **Example: spline**

☐ **http://www.nr.com/dependencies/index.php**

**NUMERICAL RECIPES**

The Art of Scientific Computing

**THIRD EDITION**

William H. Press
Saul A. Teukolsky
William T. Vetterling
Brian P. Flannery

Includes CD-ROM

# PART 3:
# FIT TO A FUNCTION ("FIT TO THEORY")

# LEAST SQUARE FITTING: "GOODNESS OF A FIT"

☐  Imagine that you have measured $N_D$ data values of the independent variable *y* as a function of the dependent variable *x*:

$$(x_i, y_i \pm \sigma_i) \qquad i = 1, N_D$$

☐  where **±$\sigma_i$ is the uncertainty** in the ith value of *y*.

☐  Our goal is to determine how well a mathematical function *y = g(x)* (also called theory or model) can describe these data.

$$g(x) = g(x; \{a_1, a_2, \ldots, a_{M_P}\}) = g(x; \{a_m\})$$

* We assume that the model function g(x) contains, in addition to the functional dependence on x, an additional dependence upon MP parameters $\{a_1, a_2, \ldots, a_{MP}\}$.

* The parameters $\{a_m\}$ are not variables, in the sense of numbers read from a meter, but rather are parts of the theoretical model such as the size of a box, the mass of a particle, or the depth of a potential well.

# MINIMIZATION "CHI-SQUARE"

$$\chi^2 \stackrel{\text{def}}{=} \sum_{i=1}^{N_D} \left( \frac{y_i - g(x_i; \{a_m\})}{\sigma_i} \right)^2$$

* **We use the chi-squared (χ²) measure as a gauge of how well a theoretical function g reproduces data**

$$\chi^2 \;\overset{\text{def}}{=}\; \sum_{i=1}^{N_D} \left( \frac{y_i - g(x_i; \{a_m\})}{\sigma_i} \right)^2$$

* **The sum is over the ND experimental points $(x_i, y_i \pm \sigma_i)$.**

* **The definition is such that smaller values of χ² are better fits, with χ² = 0 occurring if the theoretical curve went through the center of every data point.**

* **Notice also that the $1/\sigma_i^2$ weighting means that measurements with larger errors contribute less to χ².**

- **Least-squares fitting refers to adjusting the theory until a minimum in $\chi^2$ is found**

- **The goal is to find a curve that produces the least value for the summed squares of the deviations of the data from the function *g(x)*.**

- **In general, this is the best fit possible or the best way to determine the parameters in a theory.**

- **The MP parameters {$a_m$, m = 1, MP} that make $\chi^2$ an extremum are found by solving the MP equations:**

$$\frac{\partial \chi^2}{\partial a_m} = 0 \quad \Rightarrow \quad \sum_{i=1}^{N_D} \frac{[y_i - g(x_i)]}{\sigma_i^2} \frac{\partial g(x_i)}{\partial a_m} = 0 \quad (m = 1, M_P)$$

# GOODNESS OF A FIT

☐ When the deviations from theory are due to random errors and when these errors are described by a Gaussian distribution, there are some useful rules of thumb to remember.

☐ You know that your fit is good if the value of $\chi^2$ calculated via the definition is approximately equal to the number of degrees of freedom $\chi^2 \approx ND - MP$, where ND is the number of data points and MP the number of parameters in the theoretical function.

☐ If your $\chi^2$ is much less than $ND - MP$, it does not mean that you have a "great" theory or a really precise measurement; instead, you probably have too many parameters or have assigned errors ($\sigma_i$ values) that are too large. In fact, too small a $\chi^2$ may indicate that you are fitting the random scatter in the data rather than missing $\sim 1$ of the error bars, as expected for Gaussian statistics.

# FITTING TO A STRAIGHT LINE: LINEAR REGRESSSION

☐ The MP simultaneous equations simplify considerably if the functions g(x; {am}) depend linearly on the parameter values $a_i$:

$$g\left(x; \{a_1, a_2\}\right) = a_1 + a_2 x$$

☐ In this case there are MP = 2 parameters, **the slope $a_2$ and the y intercept $a_1$.**

☐ A unique solution is not possible unless the number of data points is equal to or greater than the number of parameters.

☐ For this linear case, there are just two derivatives,

$$\frac{\partial g(x_i)}{\partial a_1} = 1 \qquad \frac{\partial g(x_i)}{\partial a_2} = x_i$$

☐ and after substitution, the $\chi^2$ minimization equations can be solved

# LINEAR REGRESSION (2)

$$a_1 = \frac{S_{xx}S_y - S_x S_{xy}}{\Delta} \qquad a_2 = \frac{SS_{xy} - S_x S_y}{\Delta}$$

$$S = \sum_{i=1}^{N_D} \frac{1}{\sigma_i^2} \qquad S_x = \sum_{i=1}^{N_D} \frac{x_i}{\sigma_i^2} \qquad S_y = \sum_{i=1}^{N_D} \frac{y_i}{\sigma_i^2}$$

$$S_{xx} = \sum_{i=1}^{N_D} \frac{x_i^2}{\sigma_i^2} \qquad S_{xy} = \sum_{i=1}^{N_D} \frac{x_i y_i}{\sigma_i^2} \qquad \Delta = SS_{xx} - S_x^2$$

- **Statistics also gives you an expression for the variance or uncertainty in the deduced parameters:**

$$\sigma_{a_1}^2 = \frac{S_{xx}}{\Delta} \qquad \sigma_{a_2}^2 = \frac{S}{\Delta}$$

```cpp
#include <iostream>
#include "nr3.h"
#include "gamma.h"
#include "incgammabeta.h"
#include "fitab.h"

using namespace std;

int main () {
    int i,k;
    int kmax=500;
    double xx=0, gxx=0;
    double aa, bb, xxmin, xxmax;
    //
    int npoints;
     cin >> npoints;
    //
    VecDoub x(npoints), g(npoints);
    //
    for (i=0;i<npoints;i++){
        cin >> x[i] >> g[i];
    }
    Fitab myreg(x,g);
//
    xxmin=x[0];
    xxmax=x[npoints-1];
     //
    aa=myreg.a;
    bb=myreg.b;
    cout << aa << " " << bb << " " << myreg.chi2<< endl;

    for(k=0;k<kmax;k++){
        xx=xxmin+(double)k*(xxmax-xxmin)/((double)kmax-1);
        gxx=bb*xx+aa;
        cout << xx << " " << gxx <<endl;
    }
}
```

LINEAR REGRESSION

# GLOBAL LAND-OCEAN TEMPERATURE INDEX IN 0.01 DEGREES CELSIUS

GLOBAL Land-Ocean Temperature Index in 0.01 degrees Celsius    base period: 1951-1980

sources: GHCN 1880-01/2011 + SST:   1880-11/1981 HadISST1
12/1981-01/2011 Reynolds v2
using elimination of outliers and homogeneity adjustment
Notes: 1950 DJF = Dec 1949 - Feb 1950 ;  ***** = missing

Divide by 100 to get changes in degrees Celsius (deg-C).
Multiply that result by 1.8(=9/5) to get changes in degrees Fahrenheit (deg-F).

Best estimate for absolute global mean for 1951-1980 is  14.0 deg-C or 57.2 deg-F,
so add that to the temperature change if you want to use an absolute scale
(this note applies to global annual means only, J-D and D-N !)

Example        --      Table Value :    40
change :   0.40 deg-C  or  0.72 deg-F
abs. scale if global annual mean :   14.40 deg-C  or 57.92 deg-F

```cpp
#include <iostream>
#include "nr3.h"
#include "gamma.h"
#include "incgammabeta.h"
#include "fitab.h"
```

**WITH EXPERIMENTAL ERROR!**

```cpp
using namespace std;
double randDouble(double low, double high)
{
    double temp;
    temp = ((double) rand() / (static_cast<double>(RAND_MAX) + 1.0))* (high - low) + low;
    return temp;
}
int main () {
    int i,k;
    int kmax=500;
    double xx=0, gxx=0;
    double aa, bb, xxmin, xxmax;
    //
    int npoints=20;
    srand((unsigned) time(0));
    VecDoub x(npoints), g(npoints), sigma(npoints);

    for (i=0;i<npoints;i++){
        x[i]=double(i);
        g[i]=randDouble(-1.0,1.0)+x[i];
        sigma[i]=randDouble(-1.,1.0);
        cout << x[i] << " " << g[i] << " " << sigma[i]<< endl;
    }
    Fitab myreg(x,g,sigma);

     //
    xxmin=x[0];
    xxmax=x[npoints-1];
     //
    aa=myreg.a;
    bb=myreg.b;
    cout << aa << " " << bb << " " << myreg.chi2<< endl;
    for(k=0;k<kmax;k++){
        xx=xxmin+(double)k*(xxmax-xxmin)/((double)kmax-1);
        gxx=bb*xx+aa;
        cout << xx << " " << gxx <<endl;
    }
}
```
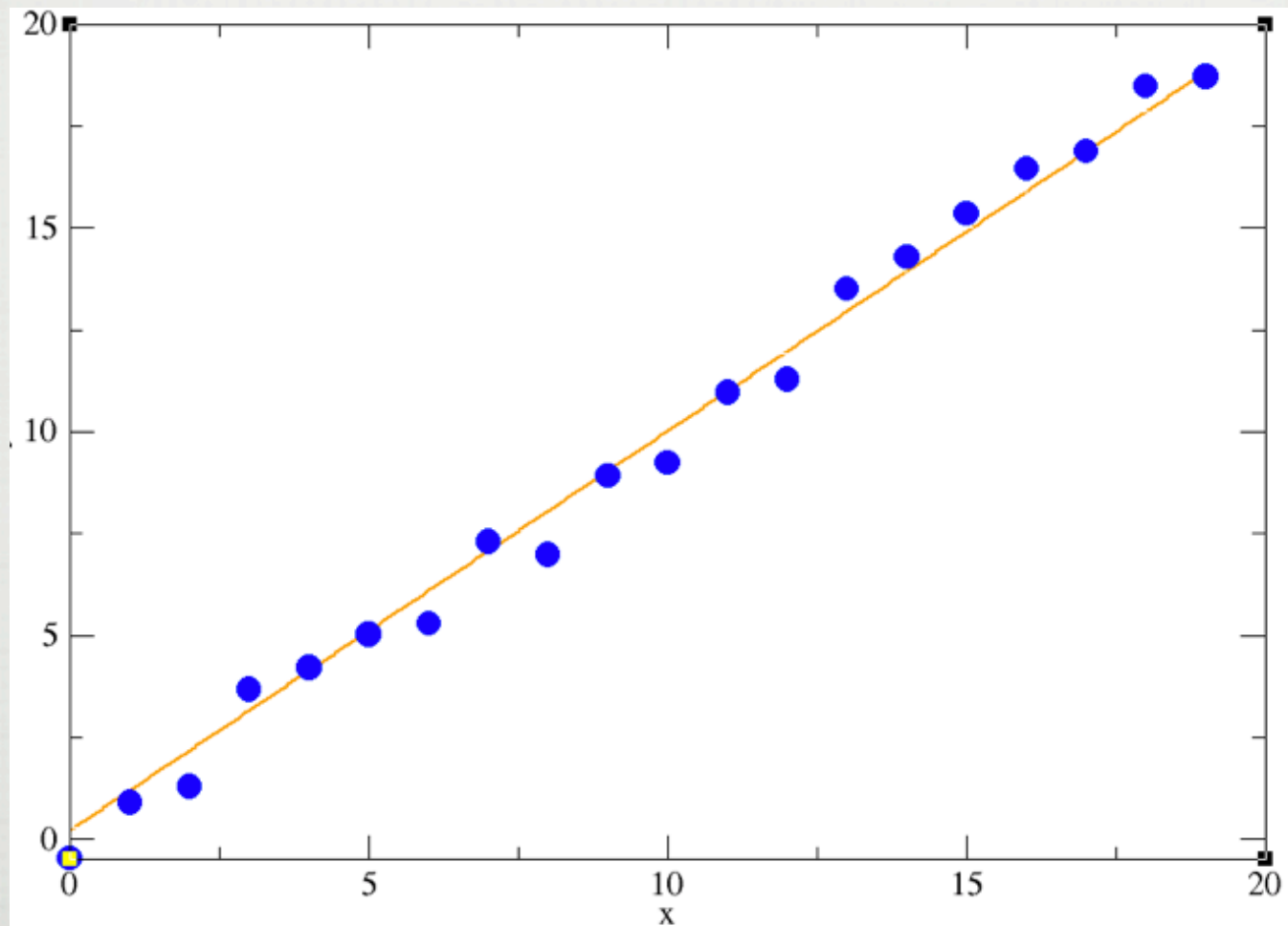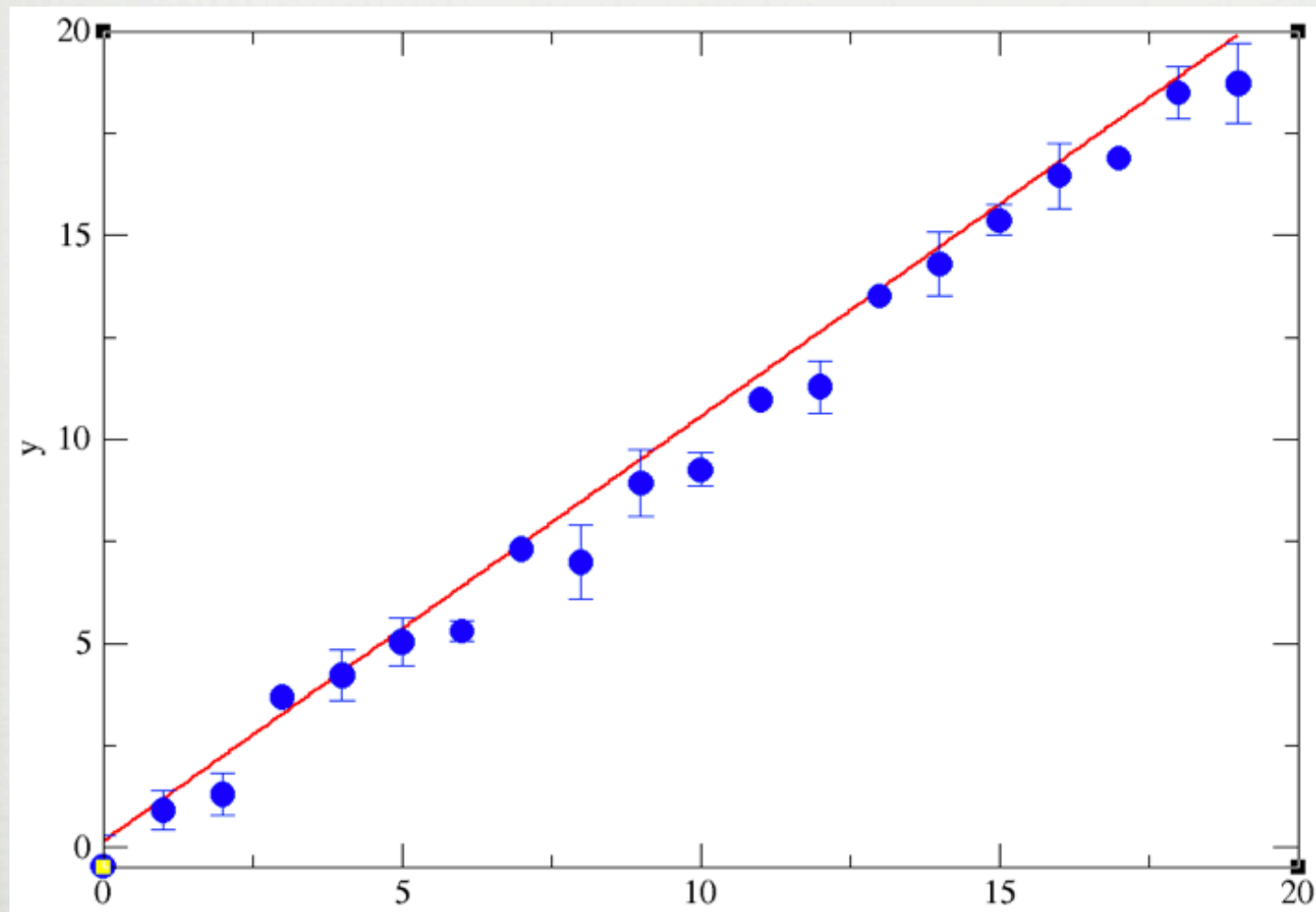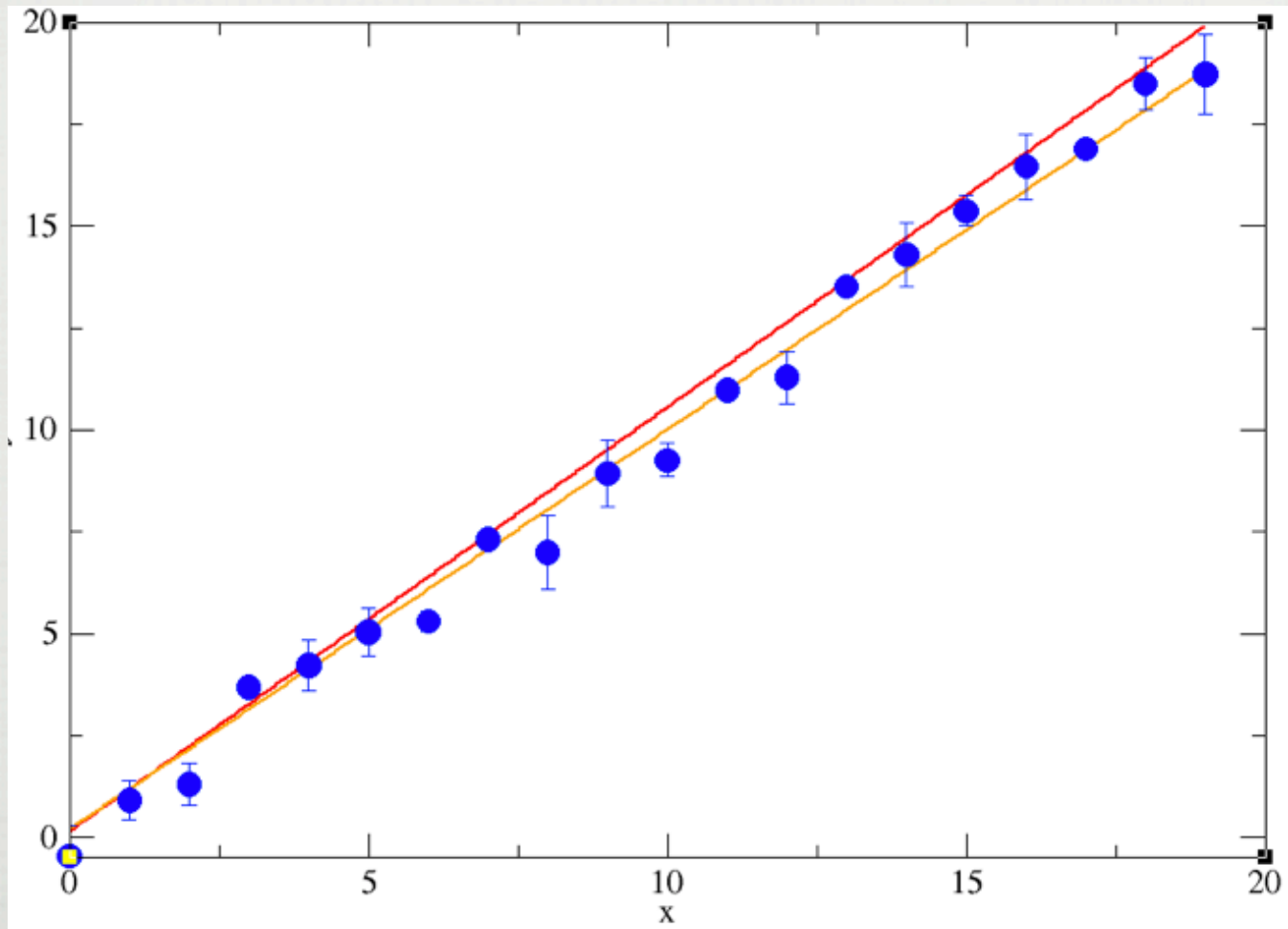
DATA SET

# SUMMARY

- **Polynomial fits**

    - **Lagrange formula and Neville's algorithm**

    - **Can fit polynomial of order N-1 passing exactly through N points**

- **Spline fits**

    - **More pleasing for the eye**
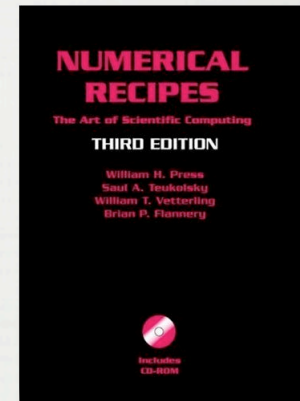
    - **3rd order polynomial with continuity of slope and curvature**

- **Least $Xi^2$ fits**

    - **best fit does not necessarily goes though all the points**

    - **very powerful when we know the analytical form of the solution**

# IN PRACTICE



- [ ] **We use "numerical recipes"**

- [ ] **Today: we will work on computing splines using numrec.**

- [ ] **See LMS**

```cpp
#include <iostream>
#include "nr3.h"
#include "interp_1d.h"
using namespace std;
int main () {
    int i,k;
    int kmax=50;
    double xx=0, gxx=0;
    double xxmin, xxmax;
    //
    int npoints;

//GET INPUT
    cin >> npoints;
//allocate memory
    VecDoub x(npoints), g(npoints);
//read in data points
    for (i=0;i<npoints;i++){
        cin >> x[i] >> g[i];
    }
//CREATE SPLINE OBJECT
    Spline_interp myspline (x,g);

//CREATE INTERPOLATION
    xxmin=x[0];
    xxmax=x[npoints-1];
      //
    for(k=0;k<kmax;k++){
        xx=xxmin+(double)k*(xxmax-xxmin)/((double)kmax-1);
        gxx= myspline.interp(xx);
        cout << xx << " " << gxx <<endl;
    }
}
```

**SPLINE WITH NUMREC**