

LECTURE NOTES
on
ELEMENTARY NUMERICAL METHODS

Eusebius Doedel

TABLE OF CONTENTS

Vector and Matrix Norms	1
Banach Lemma	20
The Numerical Solution of Linear Systems	25
Gauss Elimination	25
Operation Count	29
Using the LU-decomposition for multiple right hand sides	34
Tridiagonal Systems	37
Inverses	40
Practical Considerations	47
Gauss Elimination with Pivoting	53
Error Analysis	56
The Numerical Solution of Nonlinear Equations	73
Some Methods for Scalar Nonlinear Equations	77
Bisection	78
Regula Falsi	80
Newton's Method	83
The Chord Method	87
Newton's Method for Systems of Nonlinear Equations	92
Residual Correction	99
Convergence Analysis for Scalar Equations	102
Convergence Analysis for Systems	145

The Approximation of Functions	158
Function Norms	158
Lagrange Interpolation Polynomial	166
Lagrange Interpolation Theorem	176
Chebyshev Polynomials	185
Chebyshev Theorem	191
Taylor Polynomial	207
Taylor Theorem	211
Local Polynomial Interpolation	216
Numerical Differentiation	231
Best Approximation	240
Best Approximation in \mathbb{R}^3	240
Best Approximation in General	247
Gram-Schmidt Orthogonalization	256
Best Approximation in Function Space	259
Numerical Integration	268
Trapezoidal Rule	270
Simpson's Rule	273
Gauss Quadrature	287
Discrete Least Squares Approximation	296
Linear Least Squares	298
General Least Squares	306

Smooth Interpolation by Piecewise Polynomials	326
Cubic Spline Interpolation	330
Numerical Methods for Initial Value Problems	341
Numerical Methods	347
Stability of Numerical Approximations	355
Stiff Differential Equations	365
Boundary Value Problems in ODE	387
A Nonlinear Boundary Value Problem	403
Diffusion Problems	407
Nonlinear Diffusion Equations	420

VECTOR AND MATRIX NORMS

In later analysis we shall need a quantity (called *vector norm*) that measures the magnitude of a vector.

Let $\mathbf{x} \equiv (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$.

Then some examples of norms are :

$$\|\mathbf{x}\|_1 \equiv \sum_{k=1}^n |x_k|, \quad (\text{the “one-norm”})$$

$$\|\mathbf{x}\|_2 \equiv \left(\sum_{k=1}^n x_k^2 \right)^{\frac{1}{2}}, \quad (\text{the “two-norm”, or Euclidean length})$$

$$\|\mathbf{x}\|_\infty \equiv \max_{1 \leq k \leq n} |x_k|, \quad (\text{the “infinity-norm”, or “max-norm”})$$

$\| \mathbf{x} \|_1$ and $\| \mathbf{x} \|_2$ are special cases of

$$\| \mathbf{x} \|_p \equiv \left(\sum_{k=1}^n |x_k|^p \right)^{\frac{1}{p}},$$

while for any fixed vector \mathbf{x} we have

$\| \mathbf{x} \|_\infty$ is the limit of $\| \mathbf{x} \|_p$ as $p \rightarrow \infty$. (Prove this!)

For example if $\mathbf{x} = (1, -2, 4)^T$ then

$$\| \mathbf{x} \|_1 = 7, \quad \| \mathbf{x} \|_2 = \sqrt{21}, \quad \| \mathbf{x} \|_\infty = 4.$$

Vector norms are required to satisfy

$$(i) \quad \| \mathbf{x} \| \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n \quad \text{and} \quad \| \mathbf{x} \| = 0 \quad \text{only if } \mathbf{x} = \mathbf{0},$$

$$(ii) \quad \| \alpha \mathbf{x} \| = | \alpha | \| \mathbf{x} \|, \quad \forall \mathbf{x} \in \mathbb{R}^n, \quad \forall \alpha \in \mathbb{R},$$

$$(iii) \quad \| \mathbf{x} + \mathbf{y} \| \leq \| \mathbf{x} \| + \| \mathbf{y} \|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad (\textit{Triangle inequality}).$$

All of the examples of norms given above satisfy (i) and (ii). (Check!)

To check condition (iii) let

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T, \quad \mathbf{y} = (y_1, y_2, \dots, y_n)^T .$$

Then

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|_1 &= \sum_{k=1}^n |x_k + y_k| \leq \sum_{k=1}^n (|x_k| + |y_k|) \\ &= \sum_{k=1}^n |x_k| + \sum_{k=1}^n |y_k| = \|\mathbf{x}\|_1 + \|\mathbf{y}\|_1 . \end{aligned}$$

$$\|\mathbf{x} + \mathbf{y}\|_2 \leq \|\mathbf{x}\|_2 + \|\mathbf{y}\|_2 \quad \text{“by geometry” (Proof given later.)}$$

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|_\infty &= \max_k |x_k + y_k| \leq \max_k |x_k| + \max_k |y_k| \\ &= \|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty . \end{aligned}$$

EXERCISES:

- Let $\mathbf{x} = (1, -2, 3)^T$. Compute $\|\mathbf{x}\|_1$, $\|\mathbf{x}\|_2$, and $\|\mathbf{x}\|_\infty$.
- Graphically indicate all points $\mathbf{x} = (x_1, x_2)^T$ in \mathbb{R}^2 for which $\|\mathbf{x}\|_2 = 1$. Do the same for $\|\mathbf{x}\|_1$ and $\|\mathbf{x}\|_\infty$.
- Graphically indicate all points $\mathbf{x} = (x_1, x_2)^T$ in \mathbb{R}^2 for which $\|\mathbf{x}\|_2 \leq 1$. Do the same for $\|\mathbf{x}\|_1$ and $\|\mathbf{x}\|_\infty$.
- Graphically indicate all points $\mathbf{x} = (x_1, x_2, x_3)^T \in \mathbb{R}^3$ with $\|\mathbf{x}\|_2 = 1$. Do the same for $\|\mathbf{x}\|_1$ and $\|\mathbf{x}\|_\infty$.
- Prove that $\|\mathbf{x}\|_1 \leq n \|\mathbf{x}\|_\infty$.
- Prove that $\|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty$.
- Prove that $\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1$.
- Prove that $\lim_{p \rightarrow \infty} \|\mathbf{x}\|_p = \|\mathbf{x}\|_\infty$.

- Prove that $\| \mathbf{x} \|_1 \leq n \| \mathbf{x} \|_\infty$.

SOLUTION:

$$\| \mathbf{x} \|_1 = \sum_{i=1}^n |x_i| \leq n \max_i |x_i| = n \| \mathbf{x} \|_\infty .$$

- Prove that $\| \mathbf{x} \|_2 \leq \sqrt{n} \| \mathbf{x} \|_\infty$.

SOLUTION: From

$$\| \mathbf{x} \|_2^2 = \sum_{i=1}^n x_i^2 \leq n \max_i |x_i|^2 = n \| \mathbf{x} \|_\infty^2 ,$$

it follows that $\| \mathbf{x} \|_2 \leq \sqrt{n} \| \mathbf{x} \|_\infty$.

- Prove that $\| \mathbf{x} \|_2 \leq \| \mathbf{x} \|_1$.

SOLUTION: From

$$\| \mathbf{x} \|_1^2 = \left(\sum_{i=1}^n |x_i| \right)^2 \geq \sum_{i=1}^n x_i^2 = \| \mathbf{x} \|_2^2 ,$$

it follows that $\| \mathbf{x} \|_2 \leq \| \mathbf{x} \|_1$.

We also need a measure of the magnitude of a square matrix (*matrix norm*).

This is defined in terms of a given vector norm, namely,

$$\| \mathbf{A} \| \equiv \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A}\mathbf{x} \|}{\| \mathbf{x} \|} .$$

Thus $\| \mathbf{A} \|$ measures the maximum relative stretching in a given vector norm that occurs when multiplying all non-zero vectors $\mathbf{x} \in \mathbb{R}^n$ by \mathbf{A} .

From this definition it follows that for arbitrary $\mathbf{y} \in \mathbb{R}^n$ we have

$$\frac{\| \mathbf{A}\mathbf{y} \|}{\| \mathbf{y} \|} \leq \| \mathbf{A} \| ,$$

i.e.,

$$\| \mathbf{A}\mathbf{y} \| \leq \| \mathbf{A} \| \| \mathbf{y} \| .$$

For specific choices of vector norm it is convenient to express the corresponding induced matrix norm directly in terms of the elements of the matrix.

For the case of the $\|\cdot\|_\infty$ let

$$\mathbf{A} \equiv \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}, \quad \text{and let} \quad R \equiv \max_i \sum_{j=1}^n |a_{ij}| .$$

Thus R is the “*maximum absolute row sum*”. For $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x} \neq \mathbf{0}$, we have

$$\begin{aligned} \frac{\|\mathbf{Ax}\|_\infty}{\|\mathbf{x}\|_\infty} &= \frac{\max_i \left| \sum_{j=1}^n a_{ij} x_j \right|}{\|\mathbf{x}\|_\infty} \\ &\leq \frac{\max_i \sum_{j=1}^n |a_{ij}| \|x_j\|}{\|\mathbf{x}\|_\infty} \\ &\leq \frac{\max_i \left\{ \sum_{j=1}^n |a_{ij}| \|\mathbf{x}\|_\infty \right\}}{\|\mathbf{x}\|_\infty} = R . \end{aligned}$$

Next we show that for any matrix \mathbf{A} there always is a vector \mathbf{y} for which

$$\frac{\|\mathbf{A}\mathbf{y}\|_\infty}{\|\mathbf{y}\|_\infty} = R .$$

Let k be the row of \mathbf{A} for which $\sum_{j=1}^n |a_{kj}|$ is a maximum, *i.e.*,

$$\sum_{j=1}^n |a_{kj}| = R .$$

Take $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ such that

$$y_j = \begin{cases} -1 & \text{if } a_{kj} \geq 0 , \\ 1 & \text{if } a_{kj} < 0 . \end{cases}$$

Then

$$\frac{\|\mathbf{A}\mathbf{y}\|_\infty}{\|\mathbf{y}\|_\infty} = \|\mathbf{A}\mathbf{y}\|_\infty = \max_i \left| \sum_{j=1}^n a_{ij}y_j \right| = \sum_{j=1}^n |a_{kj}| = R .$$

Thus we have shown that

$\| \mathbf{A} \|_{\infty}$ is equal to *the maximum absolute row sum*.

For example, if

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & -3 \\ 1 & 0 & 4 \\ -1 & 5 & 1 \end{pmatrix},$$

then

$$\| \mathbf{A} \|_{\infty} = \max\{6, 5, 7\} = 7.$$

REMARK : In this example the vector \mathbf{y} is given by $\mathbf{y} = (-1, 1, 1)^T$.

For this vector we have

$$\frac{\| \mathbf{A}\mathbf{y} \|_{\infty}}{\| \mathbf{y} \|_{\infty}} = 7 = \text{maximum absolute row sum.}$$

Similarly one can show that

$$\begin{aligned}\| \mathbf{A} \|_1 &\equiv \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A} \mathbf{x} \|_1}{\| \mathbf{x} \|_1} = \max_j \sum_{i=1}^n | a_{ij} | \\ &= \textit{the maximum absolute column sum.}\end{aligned}$$

(Prove this!)

For example, for the matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & -3 \\ 1 & 0 & 4 \\ -1 & 5 & 1 \end{pmatrix},$$

we have

$$\| \mathbf{A} \|_1 = \max\{3, 7, 8\} = 8.$$

One can also show that

$$\| \mathbf{A} \|_2 \equiv \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A} \mathbf{x} \|_2}{\| \mathbf{x} \|_2} = \max_i \kappa_i(\mathbf{A}) ,$$

where the $\kappa_i(\mathbf{A})$ are the square roots of the eigenvalues of the matrix $\mathbf{A}^T \mathbf{A}$.

(These eigenvalues are indeed nonnegative).

The quantities $\{\kappa_i(\mathbf{A})\}_{i=1}^n$ are called the *singular values* of the matrix \mathbf{A} .

For example if

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix},$$

then

$$\mathbf{A}^T \mathbf{A} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix}.$$

The eigenvalues μ of $\mathbf{A}^T \mathbf{A}$ are obtained from

$$\det(\mathbf{A}^T \mathbf{A} - \lambda I) = \det \begin{pmatrix} 1 - \lambda & 1 \\ 1 & 2 - \lambda \end{pmatrix} = (1 - \lambda)(2 - \lambda) - 1 = \lambda^2 - 3\lambda + 1 = 0,$$

from which

$$\mu_1 = \frac{3 + \sqrt{5}}{2} \quad \text{and} \quad \mu_2 = \frac{3 - \sqrt{5}}{2}.$$

Thus we have

$$\|\mathbf{A}\|_2 = \sqrt{(3 + \sqrt{5})/2} \approx 1.618.$$

If \mathbf{A} is invertible then we also have

$$\| \mathbf{A}^{-1} \|_2 = \frac{1}{\min_i \kappa_i(\mathbf{A})} .$$

Thus if we order the square roots of the eigenvalues of $\mathbf{A}^T \mathbf{A}$ as

$$\kappa_1 \geq \kappa_2 \geq \cdots \kappa_n \geq 0 ,$$

then

$$\| \mathbf{A} \|_2 = \kappa_1, \quad \text{and} \quad \| \mathbf{A}^{-1} \|_2 = \frac{1}{\kappa_n} .$$

Thus in the previous example we have

$$\| \mathbf{A}^{-1} \|_2 = \frac{1}{\sqrt{(3 - \sqrt{5})/2}} \approx 1.618 (!)$$

EXERCISES:

○ Let $\mathbf{A} = \begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix}$. Compute $\|\mathbf{A}\|_2$.

○ Let $\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & 0 \end{pmatrix}$. Compute $\|\mathbf{A}\|_2$.

For a general n by n matrix \mathbf{A} :

○ Prove that $\|\mathbf{A}\|_2 \leq \sqrt{n} \|\mathbf{A}\|_\infty$.

○ Prove that $\|\mathbf{A}\|_1$ is equal to the maximum absolute column sum.

SOLUTIONS:

$$\circ \mathbf{A}^T \mathbf{A} = \begin{pmatrix} 0 & 0 \\ 2 & 0 \end{pmatrix} \begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 4 \end{pmatrix},$$

with eigenvalues $\lambda_1 = 0$ and $\lambda_2 = 4$, so that $\|\mathbf{A}\|_2 = \sqrt{4} = 2$.

$$\circ \mathbf{A}^T A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix},$$

which has three eigenvalues equal to 1. Thus $\|\mathbf{A}\|_2 = \sqrt{1} = 1$.

$$\circ \|\mathbf{A}\|_2 \leq \sqrt{n} \|\mathbf{A}\|_\infty :$$

By a previous exercise $\|\mathbf{y}\|_2 \leq \sqrt{n} \|\mathbf{y}\|_\infty$ for any vector \mathbf{y} .

Also it is easy to see that $\|\mathbf{x}\|_2 \geq \|\mathbf{x}\|_\infty$. Thus

$$\begin{aligned} \|\mathbf{A}\|_2 &\equiv \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2} \leq \sqrt{n} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_\infty}{\|\mathbf{x}\|_2} \\ &\leq \sqrt{n} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_\infty}{\|\mathbf{x}\|_\infty} = \sqrt{n} \|\mathbf{A}\|_\infty . \end{aligned}$$

○ $\| \mathbf{A} \|_1$ is equal to the maximum absolute column sum.

PROOF: Let $\mathbf{A} \equiv \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}$, $C \equiv \max_j \sum_{i=1}^n | a_{ij} |$.

Thus C is the maximum absolute column sum. First, for $\mathbf{x} \neq \mathbf{0}$, we have

$$\begin{aligned} \frac{\| \mathbf{Ax} \|_1}{\| \mathbf{x} \|_1} &= \frac{\sum_{i=1}^n | \sum_{j=1}^n a_{ij} x_j |}{\| \mathbf{x} \|_1} \leq \frac{\sum_{i=1}^n \sum_{j=1}^n | a_{ij} | | x_j |}{\| \mathbf{x} \|_1} \\ &= \frac{\sum_{j=1}^n \sum_{i=1}^n | a_{ij} | | x_j |}{\| \mathbf{x} \|_1} = \frac{\sum_{j=1}^n | x_j | \sum_{i=1}^n | a_{ij} |}{\| \mathbf{x} \|_1} \\ &\leq \frac{\sum_{j=1}^n | x_j | C}{\| \mathbf{x} \|_1} = \frac{C \| \mathbf{x} \|_1}{\| \mathbf{x} \|_1} = C. \end{aligned}$$

Secondly, we must show there is a vector \mathbf{x} for which $\| \mathbf{Ax} \|_1 / \| \mathbf{x} \|_1 = 1$. This vector is constructed as follows: If k is the index of the column having the maximum absolute column sum then let \mathbf{x} be the vector containing only zeroes, except for the k th element which is set to 1. Then $\| \mathbf{x} \|_1 = 1$ and $\| \mathbf{Ax} \|_1 = C$, so that indeed $\| \mathbf{Ax} \|_1 / \| \mathbf{x} \|_1 = C$.

EXERCISES:

- Let \mathbf{A} be any n by n matrix. For each of the following state whether it is true or false. If false then give a counter example.

$$\| \mathbf{A} \|_1 \leq \| \mathbf{A} \|_\infty \quad , \quad \| \mathbf{A} \|_\infty \leq \| \mathbf{A} \|_1 \quad .$$

- Prove that for any square matrix \mathbf{A} there always is a vector \mathbf{x} such that

$$\| \mathbf{Ax} \|_\infty = \| \mathbf{A} \|_\infty \| \mathbf{x} \|_\infty \quad .$$

- Prove that for any square matrix \mathbf{A} there always is a vector \mathbf{x} such that

$$\| \mathbf{Ax} \|_1 = \| \mathbf{A} \|_1 \| \mathbf{x} \|_1 \quad .$$

- Is there a vector \mathbf{x} such that

$$\| \mathbf{Ax} \|_1 > \| \mathbf{A} \|_1 \| \mathbf{x} \|_1 \quad ?$$

SOLUTIONS :

- $\| \mathbf{A} \|_1 \leq \| \mathbf{A} \|_\infty$ and $\| \mathbf{A} \|_\infty \leq \| \mathbf{A} \|_1$ are both False:

For $\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix}$ we have

$$\| \mathbf{A} \|_\infty = 2 \text{ and } \| \mathbf{A} \|_1 = 1, \text{ so } \| \mathbf{A} \|_1 < \| \mathbf{A} \|_\infty,$$

while for $\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}$ we have

$$\| \mathbf{A} \|_\infty = 1 \text{ and } \| \mathbf{A} \|_1 = 2, \text{ so } \| \mathbf{A} \|_\infty < \| \mathbf{A} \|_1.$$

All matrix norms defined in terms of (induced by) a given vector norm as

$$\| \mathbf{A} \| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A} \mathbf{x} \|}{\| \mathbf{x} \|}$$

automatically satisfy

$$(i) \quad \| \mathbf{A} \| \geq 0, \text{ and } \| \mathbf{A} \| = 0 \text{ only if } \mathbf{A} = \mathbf{O} \quad (\text{zero matrix}),$$

$$(ii) \quad \| \alpha \mathbf{A} \| = |\alpha| \| \mathbf{A} \|, \quad \forall \alpha \in \mathbb{R},$$

$$(iii) \quad \| \mathbf{A} + \mathbf{B} \| \leq \| \mathbf{A} \| + \| \mathbf{B} \|.$$

(Check Properties *(i)* and *(ii)* !)

PROOF of (iii):

$$\begin{aligned}\| \mathbf{A} + \mathbf{B} \| &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| (\mathbf{A} + \mathbf{B})\mathbf{x} \|}{\| \mathbf{x} \|} \\ &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{Ax} + \mathbf{Bx} \|}{\| \mathbf{x} \|} \\ &\leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{Ax} \| + \| \mathbf{Bx} \|}{\| \mathbf{x} \|} \\ &\leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{Ax} \|}{\| \mathbf{x} \|} + \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{Bx} \|}{\| \mathbf{x} \|} \\ &\equiv \| \mathbf{A} \| + \| \mathbf{B} \| . \quad \text{QED!}\end{aligned}$$

In addition we have

$$(iv) \quad \| \mathbf{AB} \| \leq \| \mathbf{A} \| \| \mathbf{B} \| .$$

PROOF of (iv):

$$\begin{aligned} \| \mathbf{AB} \| &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| (\mathbf{AB})\mathbf{x} \|}{\| \mathbf{x} \|} \\ &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A}(\mathbf{B}\mathbf{x}) \|}{\| \mathbf{x} \|} \\ &\leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A} \| \| \mathbf{B}\mathbf{x} \|}{\| \mathbf{x} \|} \\ &\leq \max_{\mathbf{x} \neq \mathbf{0}} \frac{\| \mathbf{A} \| \| \mathbf{B} \| \| \mathbf{x} \|}{\| \mathbf{x} \|} = \| \mathbf{A} \| \| \mathbf{B} \| . \quad \text{QED!} \end{aligned}$$

EXERCISES:

Let \mathbf{A} and \mathbf{B} be arbitrary n by n matrices.

- Is it true that

$$\|\mathbf{AB}\|_2 = \|\mathbf{A}\|_2 \|\mathbf{B}\|_2 \quad ?$$

If false then give a counterexample.

- Is it true that

$$\|\mathbf{A}\|_1 \|\mathbf{A}^{-1}\|_1 = 1 \quad ?$$

If false then give a counterexample.

- Let $\mathbf{A} = \begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix}$. Compute $\text{spr}(\mathbf{A})$.

Here the “*spectral radius*” $\text{spr}(\mathbf{A})$ is the absolute value of the largest eigenvalue of \mathbf{A} . Explain why $\text{spr}(\mathbf{A})$ is not a matrix norm.

SOLUTIONS:

- $\| \mathbf{AB} \|_2 = \| \mathbf{A} \|_2 \| \mathbf{B} \|_2$ is false:

Take $\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$, and $\mathbf{B} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$. Then \mathbf{AB} is the zero matrix.

So $\| \mathbf{AB} \|_2 = 0$, while $\| \mathbf{A} \|_2 = 1$ and $\| \mathbf{B} \|_2 = 1$.

- $\| \mathbf{A} \|_1 \| \mathbf{A}^{-1} \|_1 = 1$ is also false:

If $\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ then $\mathbf{A}^{-1} = \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix}$ and $\| \mathbf{A} \|_1 \| \mathbf{A}^{-1} \|_1 = 4$.

- Let $\mathbf{A} = \begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix}$. Compute $\text{spr}(\mathbf{A})$.

Here $\text{spr}(\mathbf{A}) = 0$ because both eigenvalues of \mathbf{A} are zero.

Since \mathbf{A} is not the zero matrix, it follows that the spectral radius of a matrix cannot serve as a matrix norm.

The Banach Lemma.

Let \mathbf{B} be an n by n matrix .

If in some induced matrix norm

$$\| \mathbf{B} \| < 1 ,$$

then

$\mathbf{I} + \mathbf{B}$ is nonsingular

and

$$\| (\mathbf{I} + \mathbf{B})^{-1} \| \leq \frac{1}{1 - \| \mathbf{B} \|} .$$

PROOF:

Suppose on the contrary that $\mathbf{I} + \mathbf{B}$ is singular.

Then

$$(\mathbf{I} + \mathbf{B})\mathbf{y} = \mathbf{0} ,$$

for some nonzero vector \mathbf{y} .

Hence

$$\mathbf{B}\mathbf{y} = -\mathbf{y} ,$$

and

$$1 = \frac{\|\mathbf{B}\mathbf{y}\|}{\|\mathbf{y}\|} \leq \frac{\|\mathbf{B}\| \|\mathbf{y}\|}{\|\mathbf{y}\|} = \|\mathbf{B}\| ,$$

which contradicts the assumption of the Lemma. (Why?)

Hence $\mathbf{I} + \mathbf{B}$ is nonsingular.

We now have

$$(\mathbf{I} + \mathbf{B})(\mathbf{I} + \mathbf{B})^{-1} = \mathbf{I} ,$$

from which

$$(\mathbf{I} + \mathbf{B})^{-1} = \mathbf{I} - \mathbf{B}(\mathbf{I} + \mathbf{B})^{-1} .$$

Hence

$$\| (\mathbf{I} + \mathbf{B})^{-1} \| \leq \| \mathbf{I} \| + \| \mathbf{B} \| \| (\mathbf{I} + \mathbf{B})^{-1} \| .$$

Since $\| \mathbf{I} \|$ is always 1 in any induced matrix norm, we get

$$(1 - \| \mathbf{B} \|) \| (\mathbf{I} + \mathbf{B})^{-1} \| \leq 1 ,$$

from which

$$\| (\mathbf{I} + \mathbf{B})^{-1} \| \leq \frac{1}{1 - \| \mathbf{B} \|} . \quad \text{QED!}$$

EXERCISES:

- Consider the n by n tridiagonal matrix $\mathbf{T}_n = \text{diag}[1, 3, 1]$.
For example,

$$\mathbf{T}_4 = \begin{pmatrix} 3 & 1 & 0 & 0 \\ 1 & 3 & 1 & 0 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & 1 & 3 \end{pmatrix} .$$

Use the Banach Lemma to show that \mathbf{T}_n is invertible for all positive integers n . Also compute an upper bound on $\|\mathbf{T}_n^{-1}\|_\infty$.

- Let \mathbf{A}_n be the n by n symmetric matrix

$$\mathbf{A}_n = \begin{pmatrix} 1 & \frac{1}{n} & \frac{1}{n} & \cdots & \frac{1}{n} & \frac{1}{n} \\ \frac{1}{n} & 1 & \frac{1}{n} & \cdots & \frac{1}{n} & \frac{1}{n} \\ \frac{1}{n} & \frac{1}{n} & 1 & \cdots & \frac{1}{n} & \frac{1}{n} \\ \cdot & \cdot & \cdot & \cdots & \cdot & \cdot \\ \frac{1}{n} & \frac{1}{n} & \frac{1}{n} & \cdots & \frac{1}{n} & 1 \end{pmatrix} .$$

Show \mathbf{A}_n for the case $n = 3$. Prove that \mathbf{A}_n is invertible for any dimension n , and determine an upper bound on $\|\mathbf{A}_n^{-1}\|_\infty$.

EXERCISES:

- Let \mathbf{A}_n be the n by n symmetric matrix

$$\mathbf{A}_n = \begin{pmatrix} 2n & 1 & 1 & \cdots & 1 & 1 \\ 1 & 2n & 1 & \cdots & 1 & 1 \\ 1 & 1 & 2n & \cdots & 1 & 1 \\ \cdot & \cdot & \cdot & \cdots & \cdot & \cdot \\ 1 & 1 & 1 & \cdots & 1 & 2n \end{pmatrix} .$$

Show \mathbf{A}_n for the cases $n = 1, 2, 3$. Prove that \mathbf{A}_n is invertible for any dimension n , and determine an upper bound on $\|\mathbf{A}_n^{-1}\|_\infty$.

- A square matrix is called *diagonally dominant* if in each row the absolute value of the diagonal element is greater than the sum of the absolute values of the off-diagonal elements. Use the Banach Lemma to prove that a diagonally dominant matrix is invertible.
- Derive an upper bound on $\|\mathbf{T}_n^{-1}\|_\infty$ for the n by n tridiagonal matrix $\mathbf{T}_n = \text{diag}[1, 2 + 1/n, 1]$.

SOLUTION :

- A diagonally dominant matrix \mathbf{A} is invertible:

Let

$$\mathbf{A} = \mathbf{D} + \mathbf{E} ,$$

where \mathbf{D} contains the diagonal entries of \mathbf{A} and \mathbf{E} the off-diagonal entries, that is,

$$\mathbf{D} = \begin{pmatrix} a_{11} & 0 & 0 & 0 & \cdot & 0 \\ 0 & a_{22} & 0 & 0 & \cdot & 0 \\ 0 & 0 & a_{33} & 0 & \cdot & 0 \\ 0 & 0 & 0 & a_{44} & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & \cdot & a_{nn} \end{pmatrix} \quad \mathbf{E} = \begin{pmatrix} 0 & a_{12} & a_{13} & a_{14} & \cdot & a_{1n} \\ a_{21} & 0 & a_{23} & a_{24} & \cdot & a_{2n} \\ a_{31} & a_{32} & 0 & a_{34} & \cdot & a_{3n} \\ a_{41} & a_{42} & a_{43} & 0 & \cdot & a_{4n} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ a_{n1} & a_{n2} & a_{n3} & a_{n4} & \cdot & 0 \end{pmatrix}$$

By assumption all diagonal entries a_{ii} must be greater than zero. Thus \mathbf{D} is invertible, and we can write

$$\mathbf{A} = \mathbf{D}(\mathbf{I} + \mathbf{B}) , \quad \text{where } \mathbf{B} = \mathbf{D}^{-1}\mathbf{E} .$$

SOLUTION: continued ...

The matrix $\mathbf{B} = \mathbf{D}^{-1}\mathbf{E}$ is given by

$$\mathbf{B} = \begin{pmatrix} 0 & \frac{a_{12}}{a_{11}} & \frac{a_{13}}{a_{11}} & \frac{a_{14}}{a_{11}} & \cdot & \frac{a_{1n}}{a_{11}} \\ \frac{a_{21}}{a_{22}} & 0 & \frac{a_{23}}{a_{22}} & \frac{a_{24}}{a_{22}} & \cdot & \frac{a_{2n}}{a_{22}} \\ \frac{a_{31}}{a_{33}} & \frac{a_{32}}{a_{33}} & 0 & \frac{a_{34}}{a_{33}} & \cdot & \frac{a_{3n}}{a_{33}} \\ \frac{a_{41}}{a_{44}} & \frac{a_{42}}{a_{44}} & \frac{a_{43}}{a_{44}} & 0 & \cdot & \frac{a_{4n}}{a_{44}} \\ a_{44} & a_{44} & a_{44} & \cdot & \cdot & a_{44} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \frac{a_{n1}}{a_{nn}} & \frac{a_{n2}}{a_{nn}} & \frac{a_{n3}}{a_{nn}} & \frac{a_{n4}}{a_{nn}} & \cdot & 0 \\ a_{nn} & a_{nn} & a_{nn} & a_{nn} & \cdot & 0 \end{pmatrix} .$$

Now

$$\|\mathbf{B}\|_{\infty} < 1 ,$$

because by the *diagonal dominance* assumption we have for each row that

$$\sum_{j=1}^n \left| \frac{a_{ij}}{a_{ii}} \right| = \frac{1}{|a_{ii}|} \sum_{j=1}^n |a_{ij}| < 1 \quad i = 1, 2, \dots, n .$$

It follows that \mathbf{A} is invertible.

THE NUMERICAL SOLUTION OF LINEAR SYSTEMS

The Gauss Elimination Method.

EXAMPLE: For given 4 by 4 matrix \mathbf{A} and vector $\mathbf{f} \in \mathbb{R}^4$,

$$\mathbf{A} = \begin{pmatrix} 1 & -2 & -1 & 2 \\ 2 & 0 & 1 & 2 \\ 2 & 0 & 4 & 1 \\ 1 & 6 & 1 & 2 \end{pmatrix}, \quad \mathbf{f} = \begin{pmatrix} -2 \\ 5 \\ 7 \\ 16 \end{pmatrix},$$

solve

$$\mathbf{Ax} = \mathbf{f},$$

for

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}.$$

$$\begin{pmatrix} 1 & -2 & -1 & 2 \\ 2 & 0 & 1 & 2 \\ 2 & 0 & 4 & 1 \\ 1 & 6 & 1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} -2 \\ 5 \\ 7 \\ 16 \end{pmatrix} \quad \begin{array}{l} \text{subtract } 2 \times \text{row 1 from row 2} \\ \text{subtract } 2 \times \text{row 1 from row 3} \\ \text{subtract } 1 \times \text{row 1 from row 4} \end{array}$$

$$\begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 4 & 6 & -3 \\ 0 & 8 & 2 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} -2 \\ 9 \\ 11 \\ 18 \end{pmatrix} \quad \begin{array}{l} \text{subtract } 1 \times \text{row 2 from row 3} \\ \text{subtract } 2 \times \text{row 2 from row 4} \end{array}$$

$$\begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 0 & \mathbf{3} & -1 \\ 0 & 0 & -4 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} -2 \\ 9 \\ 2 \\ 0 \end{pmatrix} \quad \text{subtract } -\frac{4}{3} \times \text{row 3 from row 4}$$

$$\begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 0 & 3 & -1 \\ 0 & 0 & 0 & \mathbf{8/3} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} -2 \\ 9 \\ 2 \\ 8/3 \end{pmatrix}$$

The bold-face numbers at the top left of each submatrix are the *pivots* :

$$\begin{pmatrix} \mathbf{1} & -2 & -1 & 2 \\ 0 & \mathbf{4} & 3 & -2 \\ 0 & 0 & \mathbf{3} & -1 \\ 0 & 0 & 0 & \mathbf{8/3} \end{pmatrix}$$

The final matrix is an *upper triangular* matrix.

The upper triangular system

$$\begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 0 & 3 & -1 \\ 0 & 0 & 0 & \mathbf{8/3} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} -2 \\ 9 \\ 2 \\ 8/3 \end{pmatrix},$$

can be solved by *backsubstitution* :

$$x_4 = (8/3)/(8/3) = 1,$$

$$x_3 = [2 - (-1)1]/3 = 1,$$

$$x_2 = [9 - (-2)1 - (3)1]/4 = 2,$$

$$x_1 = [-2 - (2)1 - (-1)1 - (-2)2]/1 = 1.$$

(Of course, actual computer computations use *floating point arithmetic*.)

Operation Count.

Using Gauss elimination for general n by n matrices, *counting multiplications and divisions only* (and treating these as equivalent).

(i) *Triangularization* (illustrated for $n = 4$) :

$$\begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{pmatrix} \Rightarrow \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \circ & \star & \star & \star \\ \circ & \circ & \star & \star \\ \circ & \circ & \circ & \star \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \bullet \\ \star \\ \star \\ \star \end{pmatrix}$$

$$(n+1)(n-1) + n(n-2) + \dots + (3)(1)$$

$$= \sum_{k=1}^{n-1} k(k+2) = \sum_{k=1}^{n-1} k^2 + 2 \sum_{k=1}^{n-1} k$$

$$= \frac{(n-1)n(2n-1)}{6} + n(n-1) = \frac{n(n-1)(2n+5)}{6}. \quad (\text{Check!})$$

(ii) *Backsubstitution* :

$$\begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \circ & \star & \star & \star \\ \circ & \circ & \star & \star \\ \circ & \circ & \circ & \star \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \bullet \\ \star \\ \star \\ \star \end{pmatrix}$$

$$1 + 2 + \dots + n = \frac{n(n+1)}{2}.$$

Taking the total of triangularization and backsubstitution we obtain

$$\frac{n(n-1)(2n+5)}{6} + \frac{n(n+1)}{2} = \frac{n^3}{3} + n^2 - \frac{n}{3}.$$

EXAMPLES:

if $n = 10$, then the total is 430,

if $n = 100$, then the total is 343 430,

if $n = 1000$, then the total is 336 333 430.

For large values of n the *dominant term* in the total operation count is $n^3/3$.

Reconsider the Gauss elimination procedure for solving the system

$$\mathbf{Ax} = \mathbf{f} ,$$

given by

$$\begin{pmatrix} 1 & -2 & -1 & 2 \\ 2 & 0 & 1 & 2 \\ 2 & 0 & 4 & 1 \\ 1 & 6 & 1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} -2 \\ 5 \\ 7 \\ 16 \end{pmatrix} .$$

- At each step *retain* the equation that *cancel*s the operation performed.

In this example this leads to :

$$\begin{array}{cccc}
\mathbf{I} & \mathbf{A} & \mathbf{x} & \mathbf{I} & \mathbf{f} \\
\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} & \begin{pmatrix} 1 & -2 & -1 & 2 \\ 2 & 0 & 1 & 2 \\ 2 & 0 & 4 & 1 \\ 1 & 6 & 1 & 2 \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} & = & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} & \begin{pmatrix} -2 \\ 5 \\ 7 \\ 16 \end{pmatrix} \\
\begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 2 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{pmatrix} & \begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 4 & 6 & -3 \\ 0 & 8 & 2 & 0 \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} & = & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 2 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{pmatrix} & \begin{pmatrix} -2 \\ 9 \\ 11 \\ 18 \end{pmatrix} \\
\begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 2 & 1 & 1 & 0 \\ 1 & 2 & 0 & 1 \end{pmatrix} & \begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 0 & 3 & -1 \\ 0 & 0 & -4 & 4 \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} & = & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 2 & 1 & 1 & 0 \\ 1 & 2 & 0 & 1 \end{pmatrix} & \begin{pmatrix} -2 \\ 9 \\ 2 \\ 0 \end{pmatrix} \\
\begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 2 & 1 & 1 & 0 \\ 1 & 2 & -\frac{4}{3} & 1 \end{pmatrix} & \begin{pmatrix} 1 & -2 & -1 & 2 \\ 0 & 4 & 3 & -2 \\ 0 & 0 & 3 & -1 \\ 0 & 0 & 0 & \frac{8}{3} \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} & = & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 2 & 1 & 1 & 0 \\ 1 & 2 & -\frac{4}{3} & 1 \end{pmatrix} & \begin{pmatrix} -2 \\ 9 \\ 2 \\ \frac{8}{3} \end{pmatrix} \\
\mathbf{L} & \mathbf{U} & \mathbf{x} & \mathbf{L} & \mathbf{g}
\end{array}$$

Thus the Gauss elimination method generates an *LU-decomposition* of \mathbf{A} :

$$\mathbf{A} = \mathbf{L} \mathbf{U} .$$

Here \mathbf{L} is lower triangular, and \mathbf{U} is upper triangular.

The below-diagonal entries of \mathbf{L} are the *multipliers* used in the elimination.

In addition we have

$$\mathbf{L} \mathbf{g} = \mathbf{f} .$$

Furthermore, $\mathbf{L} \mathbf{U} \mathbf{x} = \mathbf{L} \mathbf{g}$.

Since \mathbf{L} is nonsingular, we therefore also have

$$\mathbf{U} \mathbf{x} = \mathbf{g} .$$

Using the LU-decomposition for multiple right hand sides.

Suppose we want to solve

$$\mathbf{Ax}^{(k)} = \mathbf{f}^{(k)} ,$$

with fixed \mathbf{A} , but for *multiple right hand side vectors*

$$\mathbf{f}^{(k)} , \quad k = 1, 2, \dots, m .$$

Algorithm :

(i) Determine the LU-decomposition of \mathbf{A} .

(ii) Solve
$$\begin{cases} \mathbf{Lg}^{(k)} = \mathbf{f}^{(k)} , \\ \mathbf{Ux}^{(k)} = \mathbf{g}^{(k)} , \end{cases} \quad k = 1, 2, \dots, m .$$

Note that the decomposition need only be computed once.

Operation Count.

Multiplications and divisions for an n by n system with m right hand sides :

Step (i) (**LU**-decomposition) :

$$\begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \Rightarrow \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \circ & \star & \star & \star \\ \circ & \circ & \star & \star \\ \circ & \circ & \circ & \star \end{pmatrix}$$

$$n(n-1) + (n-1)(n-2) + \cdots + (2)(1)$$

$$= \sum_{k=1}^{n-1} k(k+1) = \sum_{k=1}^{n-1} k^2 + \sum_{k=1}^{n-1} k$$

$$= \frac{(n-1)n(2n-1)}{6} + \frac{n(n-1)}{2} = \frac{n^3}{3} - \frac{n}{3}. \quad (\text{Check!})$$

$$\begin{pmatrix} 1 & \circ & \circ & \circ \\ \star & 1 & \circ & \circ \\ \star & \star & 1 & \circ \\ \star & \star & \star & 1 \end{pmatrix} \begin{pmatrix} \mathbf{g} \\ g_1 \\ g_2 \\ g_3 \\ g_4 \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{pmatrix}, \quad \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \circ & \star & \star & \star \\ \circ & \circ & \star & \star \\ \circ & \circ & \circ & \star \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \mathbf{g} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{pmatrix}$$

Step (ii) :

$$\begin{cases} \mathbf{L}\mathbf{g}^{(k)} = \mathbf{f}^{(k)} : m(1 + 2 + \cdots + (n - 1)) , \\ \mathbf{U}\mathbf{x}^{(k)} = \mathbf{g}^{(k)} : m(1 + 2 + \cdots + n) . \end{cases} \quad k = 1, 2, \cdots, m .$$

Total Step (ii) : mn^2 (Check!).

The total of Steps (i) and (ii) is therefore

$$\frac{n^3}{3} + mn^2 - \frac{n}{3} .$$

NOTE: For m small and n large the *dominant term* remains $n^3/3$.

Tridiagonal systems.

For *tridiagonal systems* of linear equations

$$\begin{pmatrix} b_1 & c_1 & & & & \\ a_2 & b_2 & c_2 & & & \\ & a_3 & b_3 & c_3 & & \\ & & \cdot & \cdot & \cdot & \\ & & & a_{n-1} & b_{n-1} & c_{n-1} \\ & & & & a_n & b_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \cdot \\ x_{n-1} \\ x_n \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ \cdot \\ f_{n-1} \\ f_n \end{pmatrix},$$

Gauss elimination reduces to a simple algorithm :

$$\left. \begin{aligned} \beta_1 &= b_1, & g_1 &= f_1, \\ \gamma_k &= a_k / \beta_{k-1}, \\ \beta_k &= b_k - \gamma_k c_{k-1}, \\ g_k &= f_k - \gamma_k g_{k-1}, \end{aligned} \right\} \quad k = 2, 3, \dots, n.$$

This transform the tridiagonal system into the upper-triangular form

$$\begin{pmatrix} \beta_1 & c_1 & & & & & \\ & \beta_2 & c_2 & & & & \\ & & \beta_3 & c_3 & & & \\ & & & \cdot & \cdot & & \\ & & & & \beta_{n-1} & c_{n-1} & \\ & & & & & \beta_n & \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \cdot \\ x_{n-1} \\ x_n \end{pmatrix} = \begin{pmatrix} g_1 \\ g_2 \\ g_3 \\ \cdot \\ g_{n-1} \\ g_n \end{pmatrix} .$$

The backsubstitution algorithm now becomes

$$x_n = \frac{g_n}{\beta_n} ,$$

$$x_k = \frac{g_k - c_k x_{k+1}}{\beta_k} , \quad k = n - 1 , n - 2 , \dots , 1 .$$

Inverses.

The *inverse* of a n by n matrix \mathbf{A} is defined to be a matrix \mathbf{A}^{-1} such that

$$\mathbf{A} (\mathbf{A}^{-1}) = (\mathbf{A}^{-1}) \mathbf{A} = \mathbf{I},$$

where

$$\mathbf{I} \equiv \begin{pmatrix} 1 & 0 & \cdot & 0 \\ 0 & 1 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & 1 \end{pmatrix} \quad (\text{the } \textit{identity matrix}).$$

\mathbf{A} is *invertible* if and only if

$$\det \mathbf{A} \neq 0.$$

The inverse is then unique.

To compute \mathbf{A}^{-1} we can solve

$$\mathbf{A} (\mathbf{A}^{-1}) = \mathbf{I},$$

which is of the form

$$\begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} \\ c_{11} & c_{12} & c_{13} & c_{14} \\ c_{21} & c_{22} & c_{23} & c_{24} \\ c_{31} & c_{32} & c_{33} & c_{34} \\ c_{41} & c_{42} & c_{43} & c_{44} \end{pmatrix} = \begin{pmatrix} \mathbf{I} \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

This corresponds to solving a linear system with n right hand sides.

Using the earlier formula, the number of multiplications and divisions is

$$\frac{n^3}{3} + (n) n^2 - \frac{n}{3} = \frac{4n^3}{3} - \frac{n}{3}.$$

But we can *omit some operations*, because the right hand side vectors, *i.e.*, the columns of \mathbf{I} , are special.

In particular, multiplications by 0 or 1 can be omitted.

The total number of multiplications that can be omitted is seen to be

$$\begin{aligned}
 (n)(n-1) + (n-1)(n-2) + \cdots + (2)(1) &= \sum_{k=1}^{n-1} k(k+1) \\
 &= \sum_{k=1}^{n-1} k^2 + \sum_{k=1}^{n-1} k = \frac{(n-1)n(2n-1)}{6} + \frac{n(n-1)}{2} \\
 &= \frac{n^3}{3} - \frac{n}{3}. \quad (\text{Check!})
 \end{aligned}$$

Thus to find \mathbf{A}^{-1} we need only

$$\left(\frac{4n^3}{3} - \frac{n}{3}\right) - \left(\frac{n^3}{3} - \frac{n}{3}\right) = n^3 \quad \text{operations.}$$

REMARK:

To solve a n by n linear system

$$\mathbf{Ax}^{(k)} = \mathbf{f}^{(k)} ,$$

with m right hand sides, takes

$$\frac{n^3}{3} + mn^2 - \frac{n}{3} \quad \text{operations ,}$$

as derived earlier for the **LU**-decomposition algorithm.

One can also find the solution vectors by computing \mathbf{A}^{-1} , and setting

$$\mathbf{x}^{(k)} = \mathbf{A}^{-1} \mathbf{f}^{(k)} .$$

But this takes

$$n^3 + mn^2 \quad \text{operations .}$$

Thus this method is *always less efficient* , no matter how big n is !

EXERCISES:

- Compute the **LU**-decomposition of the tridiagonal matrix

$$\mathbf{T}_4 = \begin{pmatrix} 3 & 1 & 0 & 0 \\ 1 & 3 & 1 & 0 \\ 0 & 1 & 3 & 1 \\ 0 & 0 & 1 & 3 \end{pmatrix} .$$

Let $\mathbf{f} = (4, 5, 5, 4)^T$. Using the matrices \mathbf{L} and \mathbf{U} , solve $\mathbf{L}\mathbf{g} = \mathbf{f}$, followed by $\mathbf{U}\mathbf{x} = \mathbf{g}$. After having computed the vector \mathbf{x} in this way, check your answer by verifying that \mathbf{x} satisfies the equation $\mathbf{T}_4\mathbf{x} = \mathbf{f}$.

- How many multiplications and divisions are needed to compute the **LU**-decomposition of the specific tridiagonal matrix $\mathbf{T}_n = \text{diag}[1, 3, 1]$ as a function of n ? Make sure not to count unnecessary operations.
- If the **LU**-decomposition of this n by n tridiagonal matrix takes 0.01 second on a given computer if $n = 10^5$, then how much time could it take if $n = 10^9$?

EXERCISES:

- Suppose the \mathbf{LU} decomposition of a matrix \mathbf{A} is given by

$$\mathbf{L} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{U} = \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix}.$$

Using *only* \mathbf{L} , \mathbf{U} and \mathbf{f} , (that is, without explicitly determining \mathbf{A}), solve $\mathbf{Ax} = \mathbf{f}$, when $\mathbf{f} = (6, 9, 10)^T$.

- Suppose that solving a general n by n linear system of the form $\mathbf{Ax} = \mathbf{f}$ by Gauss elimination takes 10 seconds on a given computer if $n = 1000$.

Estimate how much time it will take to solve a 1000 by 1000 system $\mathbf{Lg} = \mathbf{f}$, followed by solving $\mathbf{Ux} = \mathbf{g}$, where \mathbf{L} is lower triangular with 1's along its main diagonal, and where \mathbf{U} is upper triangular?

Thus you may assume that \mathbf{L} and \mathbf{U} have already been computed.

EXERCISES:

- Suppose that multiplying two general n by n matrices takes 3 seconds on a given computer, if $n = 1000$.

Estimate how much time it will take to compute the **LU**-decomposition of such a matrix.

- Suppose that solving a general system of linear equations of dimension 1000 takes 10 seconds on a given computer.

Estimate how much time it will take to solve a *tridiagonal* linear system of dimension 10^6 on that computer.

- How many *divisions* are needed for **LU**-decomposition of an n by n *tridiagonal* matrix (not counting multiplications and additions)?
- How many *divisions* are needed for **LU**-decomposition of an n by n *general* matrix (not counting multiplications and additions)?

- Suppose multiplying two general n by n matrices takes 3 seconds on a given computer, if $n = 1000$. Estimate how much time it will take to compute the **LU**-decomposition of such a matrix.

SOLUTION : Multiplying two n by n matrices takes n^3 multiplications, while **LU**-decomposition takes approximately $n^3/3$ multiplications and divisions. Thus **LU**-decomposition of a matrix of dimension $n = 1000$ can be estimated to take one second.

- Suppose solving a general system of linear equations of dimension 1000 takes 10 seconds on a given computer. Estimate how much time it will take to solve a *tridiagonal* system of dimension 10^6 on that computer.

SOLUTION : Solving a system of n equations takes approximately $n^3/3$ multiplications and divisions, while solving a tridiagonal system takes approximately $5n$ multiplications and divisions. Thus the answer is

$$\frac{5 \cdot 10^6}{1000^3/3} \cdot 10 = \frac{15 \cdot 10^6}{10^9} \cdot 10 = 0.15 \text{ seconds .}$$

Practical Considerations.

- *Memory reduction.*

In an implementation of the **LU** decomposition algorithm, the multipliers can be stored in the lower triangular part of the original matrix **A**.

In the earlier example, with

$$A = \begin{pmatrix} 1 & -2 & -1 & 2 \\ 2 & 0 & 1 & 2 \\ 2 & 0 & 4 & 1 \\ 1 & 6 & 1 & 2 \end{pmatrix},$$

this function would return the matrix :

$$\begin{pmatrix} \mathbf{1} & -2 & -1 & 2 \\ 2 & \mathbf{4} & 3 & -2 \\ 2 & 1 & \mathbf{3} & -1 \\ 1 & 2 & -4/3 & \mathbf{8/3} \end{pmatrix}.$$

- *Row interchanges.*

Gauss elimination will *fail* for the matrix

$$\begin{pmatrix} 0 & 2 & 1 \\ 1 & 1 & 2 \\ 2 & 3 & -1 \end{pmatrix},$$

since the first pivot is zero.

A *division by zero* will occur when the first multiplier is computed !

The *remedy* is to interchange rows to get

$$\begin{pmatrix} 1 & 1 & 2 \\ 0 & 2 & 1 \\ 2 & 3 & -1 \end{pmatrix}$$

Several such *interchanges* may be needed during Gauss elimination.

- *Loss of accuracy.*

More generally, loss of accuracy may occur when there are large multipliers.

For example solve

$$\begin{pmatrix} 0.0000001 & 1 \\ & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} ,$$

on a “six-decimal-digit computer”.

(The exact solution is $x_1 \approx 1$, $x_2 \approx 1$.)

Note that the multiplier is 10,000,000 .

“*Six-decimal-digit computer*” :

Assume all arithmetic operations are performed to infinite precision, but then truncated to six decimal digits (plus exponent).

Thus, for example,

$$-100/3$$

is stored as

$$-3.33333E + 01$$

$$\begin{pmatrix} 0.0000001 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

(a) Elimination gives :

$$\begin{pmatrix} 1.000000\text{E} - 07 & 1.000000\text{E} + 00 \\ 0 & -9999999\text{E} + 01 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1.000000\text{E} + 00 \\ -9999999\text{E} + 01 \end{pmatrix} .$$

(b) Backsubstitution gives :

$$x_2 = 1.000000\text{E} + 00 , \quad x_1 = 0.000000\text{E} + 00 .$$

Clearly this result is very bad !

Again, the *remedy* is to interchange rows :

$$\begin{pmatrix} 1 & 1 \\ 0.0000001 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} .$$

Now the multiplier is only $1.000000\text{E} - 07$, and we obtain :

(a) Elimination :

$$\begin{pmatrix} 1.000000\text{E} + 00 & 1.000000\text{E} + 00 \\ 0 & .999999\text{E} + 00 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2.000000\text{E} + 00 \\ .999999\text{E} + 00 \end{pmatrix} .$$

(b) Backsubstitution : $x_2 = 1.000000\text{E} + 00$, $x_1 = 1.000000\text{E} + 00$.

This solution is accurate !

Gauss Elimination with pivoting.

One variant of the Gauss elimination procedure that avoids loss of accuracy due to large multipliers is called

“Gauss elimination with *partial pivoting* (or *row pivoting*)”.

In this modified algorithm *rows are interchanged each time a pivot element is sought*, so that the pivot is as large as possible.

All multipliers are then necessarily less than 1 in magnitude.

EXAMPLE:

$$\begin{pmatrix} 2 & 2 & 1 \\ 1 & 0 & 1 \\ 4 & 1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 5 \\ 2 \\ 7 \end{pmatrix} \quad \textit{interchange row 1 and 3}$$

$$\begin{pmatrix} 4 & 1 & 2 \\ 1 & 0 & 1 \\ 2 & 2 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 7 \\ 2 \\ 5 \end{pmatrix} \quad \begin{array}{l} \textit{subtract } \frac{1}{4} \textit{ row 1 from row 2} \\ \textit{subtract } \frac{2}{4} \textit{ row 1 from row 3} \end{array}$$

$$\begin{pmatrix} 4 & 1 & 2 \\ 0 & -1/4 & 1/2 \\ 0 & 3/2 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 7 \\ 1/4 \\ 3/2 \end{pmatrix} \quad \textit{interchange row 2 and 3}$$

$$\begin{pmatrix} 4 & 1 & 2 \\ 0 & 3/2 & 0 \\ 0 & -1/4 & 1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 7 \\ 3/2 \\ 1/4 \end{pmatrix} \quad \textit{subtract } \frac{-1}{6} \textit{ row 2 from row 3}$$

$$\begin{pmatrix} 4 & 1 & 2 \\ 0 & 3/2 & 0 \\ 0 & 0 & 1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 7 \\ 3/2 \\ 1/2 \end{pmatrix} \quad \begin{array}{l} \textit{backsubstitution : } x_1 = 1 \\ \textit{backsubstitution : } x_2 = 1 \\ \textit{backsubstitution : } x_3 = 1 \end{array}$$

EXERCISES:

- Use Gauss elimination with *row pivoting* to solve

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 8 & 11 \\ 3 & 22 & 35 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 10 \end{pmatrix} .$$

- Suppose that Gauss elimination with *row pivoting* is used to solve the tridiagonal system

$$\mathbf{T}_4 = \begin{pmatrix} 2 & 1 & 0 & 0 \\ 1 & 2 & 1 & 0 \\ 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 4 \\ 8 \\ 12 \\ 11 \end{pmatrix} .$$

Are any rows actually interchanged?

Can you also answer this question for *general* $\mathbf{T}_n = \text{diag}[1, 2, 1]$?

Error Analysis.

Suppose we want to solve

$$\mathbf{Ax} = \mathbf{f},$$

where the n by n matrix \mathbf{A} is nonsingular.

Assume that an error is made in the right hand side, *i.e.*, instead we solve

$$\mathbf{Ay} = \mathbf{f} + \mathbf{r}.$$

What will be the error $\| \mathbf{y} - \mathbf{x} \|$ in the solution ?

Subtract

$$\mathbf{Ax} = \mathbf{f} ,$$

from

$$\mathbf{Ay} = \mathbf{f} + \mathbf{r} ,$$

to get

$$\mathbf{A}(\mathbf{y} - \mathbf{x}) = \mathbf{r} .$$

Thus

$$\mathbf{y} - \mathbf{x} = \mathbf{A}^{-1}\mathbf{r} ,$$

so that

$$\| \mathbf{y} - \mathbf{x} \| = \| \mathbf{A}^{-1}\mathbf{r} \| \leq \| \mathbf{A}^{-1} \| \| \mathbf{r} \| .$$

Hence *if* $\| \mathbf{A}^{-1} \|$ *is large* then a *small perturbation* \mathbf{r} of the right hand side \mathbf{f} may lead to a *large change* in the solution .

For example,

$$\begin{pmatrix} 1 & -1.001 \\ 2.001 & -2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2.001 \\ 4.001 \end{pmatrix} ,$$

has exact solution

$$x_1 = 1 , \quad x_2 = -1 .$$

Suppose instead we solve

$$\begin{pmatrix} 1 & -1.001 \\ 2.001 & -2 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 2.000 \\ 4.002 \end{pmatrix} .$$

The exact solution of this system is

$$y_1 = 2 , \quad y_2 = 0 .$$

Note that the change in the right hand side has norm

$$\| \mathbf{r} \|_{\infty} = 0.001 .$$

Also note that the change in the solution is much larger, namely,

$$\| \mathbf{x} - \mathbf{y} \|_{\infty} = 1 .$$

In this example

$$\mathbf{A}^{-1} \approx \begin{pmatrix} -666.44 & 333.55 \\ -666.77 & 333.22 \end{pmatrix} .$$

Hence

$$\| \mathbf{A}^{-1} \|_{\infty} \approx 1000 , \quad \text{whereas} \quad \| \mathbf{A} \|_{\infty} \approx 4 .$$

Errors always occur in floating point computations due to finite word length.

For example, on a “six digit computer” $\frac{1}{3}$ is represented by .333333E+00 .

Such errors occur in both right hand side and matrix.

Suppose we want to solve

$$\mathbf{Ax} = \mathbf{f} ,$$

but instead solve the perturbed system

$$(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{f} + \mathbf{r} .$$

Here the perturbation \mathbf{E} is also a n by n matrix.

Theorem. Consider

$$\mathbf{Ax} = \mathbf{f}, \quad \text{and} \quad (\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{f} + \mathbf{r}.$$

Assume that

\mathbf{A} is nonsingular, and that $\|\mathbf{E}\| < \frac{1}{\|\mathbf{A}^{-1}\|}$, i.e., $\|\mathbf{A}^{-1}\| \|\mathbf{E}\| < 1$.

Then

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\text{cond}(\mathbf{A})}{1 - \|\mathbf{A}^{-1}\| \|\mathbf{E}\|} \left(\frac{\|\mathbf{r}\|}{\|\mathbf{f}\|} + \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|} \right),$$

where

$$\text{cond}(\mathbf{A}) \equiv \|\mathbf{A}^{-1}\| \|\mathbf{A}\|,$$

is called the *condition number* of \mathbf{A} .

PROOF:

First write

$$\mathbf{A} + \mathbf{E} = \mathbf{A} (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E}) .$$

Now, using the assumption,

$$\| \mathbf{A}^{-1} \mathbf{E} \| \leq \| \mathbf{A}^{-1} \| \| \mathbf{E} \| < 1 .$$

Hence by the Banach Lemma $(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})$ is nonsingular and

$$\| (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1} \| \leq \frac{1}{1 - \| \mathbf{A}^{-1}\mathbf{E} \|} \leq \frac{1}{1 - \| \mathbf{A}^{-1} \| \| \mathbf{E} \|} .$$

Next

$$(\mathbf{A} + \mathbf{E}) \mathbf{y} = \mathbf{f} + \mathbf{r} ,$$

implies

$$(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E}) \mathbf{y} = \mathbf{A}^{-1} (\mathbf{f} + \mathbf{r}) ,$$

so that

$$\mathbf{y} = (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1} \mathbf{A}^{-1} (\mathbf{f} + \mathbf{r}) .$$

Then

$$\begin{aligned} \mathbf{y} - \mathbf{x} &= (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1} \left(\mathbf{A}^{-1}(\mathbf{f} + \mathbf{r}) - (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})\mathbf{x} \right) \\ &= (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1} \left(\mathbf{x} + \mathbf{A}^{-1}\mathbf{r} - \mathbf{x} - \mathbf{A}^{-1}\mathbf{E}\mathbf{x} \right) \\ &= (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1} \mathbf{A}^{-1} (\mathbf{r} - \mathbf{E}\mathbf{x}) . \end{aligned}$$

Finally,

$$\begin{aligned}
\frac{\| \mathbf{y} - \mathbf{x} \|}{\| \mathbf{x} \|} &\leq \frac{\| (\mathbf{I} + \mathbf{A}^{-1} \mathbf{E})^{-1} \| \quad \| \mathbf{A}^{-1} \| \quad \left(\| \mathbf{r} \| + \| \mathbf{E} \| \| \mathbf{x} \| \right)}{\| \mathbf{x} \|} \\
&\leq \frac{\| \mathbf{A}^{-1} \|}{1 - \| \mathbf{A}^{-1} \| \| \mathbf{E} \|} \left(\frac{\| \mathbf{r} \|}{\| \mathbf{x} \|} + \| \mathbf{E} \| \right) \\
&= \frac{\| \mathbf{A}^{-1} \| \quad \| \mathbf{A} \|}{1 - \| \mathbf{A}^{-1} \| \| \mathbf{E} \|} \left(\frac{\| \mathbf{r} \|}{\| \mathbf{A} \| \| \mathbf{x} \|} + \frac{\| \mathbf{E} \|}{\| \mathbf{A} \|} \right) \\
&\leq \frac{\| \mathbf{A}^{-1} \| \quad \| \mathbf{A} \|}{1 - \| \mathbf{A}^{-1} \| \| \mathbf{E} \|} \left(\frac{\| \mathbf{r} \|}{\| \mathbf{f} \|} + \frac{\| \mathbf{E} \|}{\| \mathbf{A} \|} \right) .
\end{aligned}$$

The last step uses the fact that

$$\| \mathbf{f} \| = \| \mathbf{A} \mathbf{x} \| \leq \| \mathbf{A} \| \| \mathbf{x} \| . \quad \text{QED!}$$

From the above theorem we can conclude that :

If $\text{cond}(\mathbf{A})$ is large, then the relative error $\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|}$ can be large .

Note, however, that $\text{cond}(\mathbf{A})$ is never less than 1 because

$$1 = \|\mathbf{I}\| = \|\mathbf{A}^{-1}\mathbf{A}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \equiv \text{cond}(\mathbf{A}) ,$$

in any induced matrix norm.

A matrix having a large condition number is said to be *ill-conditioned*.

For example the 2 by 2 matrix

$$\begin{pmatrix} 1 & -1.001 \\ 2.001 & -2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix},$$

considered earlier, with inverse

$$\mathbf{A}^{-1} \approx \begin{pmatrix} -666.44 & 333.55 \\ -666.77 & 333.22 \end{pmatrix},$$

has condition number

$$\text{cond}(\mathbf{A}) \approx (4)(1000) = 4000,$$

in the matrix infinity norm.

We may say that this matrix is ill-conditioned.

Solving ill-conditioned systems numerically, if they can be solved at all on a given computer, normally requires pivoting.

However, systems that need pivoting need not be ill-conditioned.

Reconsider, for example, the 2 by 2 matrix

$$\mathbf{A} = \begin{pmatrix} .0000001 & 1 \\ 1 & 1 \end{pmatrix}, \quad \text{with} \quad \mathbf{A}^{-1} \approx \begin{pmatrix} -1 & 1 \\ 1 & 0 \end{pmatrix} \quad (\text{Check!}),$$

for which the condition number is approximately 4 using the infinity norm.

But solving a linear system with the above \mathbf{A} as coefficient matrix requires pivoting, at least on a six (decimal) digit computer.

SOLUTION: Write $\mathbf{S}_{n-1} = \mathbf{D}_{n-1} (\mathbf{I}_{n-1} + \mathbf{B}_{n-1})$, where

$$\mathbf{D}_{n-1} = \begin{pmatrix} 2(h_1 + h_2) & & & & & \\ & 2(h_2 + h_3) & & & & \\ & & 2(h_3 + h_4) & & & \\ & & & \cdot & & \cdot \\ & & & & \cdot & \cdot \\ & & & & & 2(h_{n-1} + h_n) \end{pmatrix}$$

and

$$\mathbf{B}_{n-1} = \begin{pmatrix} 0 & \frac{h_2}{2(h_1+h_2)} & & & & \\ \frac{h_2}{2(h_2+h_3)} & 0 & \frac{h_3}{2(h_2+h_3)} & & & \\ & \frac{h_3}{2(h_3+h_4)} & 0 & \frac{h_4}{2(h_3+h_4)} & & \\ & & \cdot & \cdot & \cdot & \\ & & & \frac{h_{n-1}}{2(h_{n-1}+h_n)} & 0 & \end{pmatrix}.$$

$\|\mathbf{B}_{n-1}\|_\infty = \frac{1}{2} < 1$, so $\mathbf{I}_{n-1} + \mathbf{B}_{n-1}$, and hence \mathbf{S}_{n-1} , is invertible, with

$$\|\mathbf{S}_{n-1}^{-1}\|_\infty \leq \frac{\|\mathbf{D}_{n-1}^{-1}\|_\infty}{1 - \|\mathbf{B}_{n-1}\|_\infty} \leq \frac{2}{2 \min_i \{h_i + h_{i+1}\}},$$

and

$$\text{cond}(\mathbf{S}_{n-1}) = \|\mathbf{S}_{n-1}\|_\infty \|\mathbf{S}_{n-1}^{-1}\|_\infty = \frac{3 \max_i \{h_i + h_{i+1}\}}{\min_i \{h_i + h_{i+1}\}},$$

which can be arbitrarily large.

EXERCISES:

- Consider the n by n matrix :

$$\mathbf{A}_n = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \cdot & \cdot & \cdot & \cdot \\ 1 & 1 & \cdots & 1 \end{pmatrix},$$

and let \mathbf{I}_n be the n by n identity matrix.

For what ϵ does the Banach Lemma guarantee $\mathbf{I}_n + \epsilon\mathbf{A}_n$ is invertible?

Also determine a bound on $\text{cond}(\mathbf{I}_n + \epsilon\mathbf{A}_n)^{-1}$ in this case.

- Use the Banach Lemma to prove that the five-diagonal matrix

$$\mathbf{T}_n = \text{diag}[1, 1, 5, 1, 1],$$

is invertible for all $n \geq 1$.

Derive an upper bound on $\text{cond}(\mathbf{T}_n)$ using the matrix infinity norm.

EXERCISES:

For each of the following statements, state whether it is true or false. If true then explain why; if false then give a counterexample.

- A condition number of 10^6 is large.
- All large matrices are ill-conditioned.
- All ill-conditioned matrices have small determinants.
- Only ill-conditioned matrices require pivoting.
- If pivoting is needed then the matrix is ill-conditioned.
- The condition number of a matrix is never smaller than 1.
- Tridiagonal matrices are never ill-conditioned.

EXERCISES:

For each of the following statements about matrices, say whether it is true or false. Explain your answer.

- Two n by n matrices can be multiplied using n^2 multiplications.
- **LU**-decomposition of the n by n tridiagonal matrix $\text{diag}[1, 2, 1]$ can be done using only $n - 1$ divisions and zero multiplications.
- **LU**-decomposition of a general n by n tridiagonal matrix requires $2n - 2$ multiplications and divisions.
- The n by n tridiagonal matrix $\mathbf{T}_n = \text{diag}[1, 2 + 1/n, 1]$ is nonsingular for any positive integer n .

EXERCISES:

For each of the following statements about matrices, say whether it is true or false. Explain your answer.

- For large n , the **LU**-decomposition of a general n by n matrix requires approximately $n^3/3$ multiplications and divisions.
- The inverse of a general, nonsingular n by n matrix can be computed using n^3 multiplications and divisions.
- If **D** is a diagonal matrix (*i.e.*, its entries d_{ij} are zero if $i \neq j$), then

$$\| \mathbf{D} \|_1 = \| \mathbf{D} \|_2 = \| \mathbf{D} \|_\infty .$$

- If $\| \mathbf{A}^{-1} \|$ is large then $\text{cond}(\mathbf{A})$ is large.

THE NUMERICAL SOLUTION OF NONLINEAR EQUATIONS

Introduction.

For a system of n *linear* equations in n unknowns

$$\mathbf{Ax} = \mathbf{f} ,$$

where $\mathbf{x}, \mathbf{f} \in \mathbb{R}^n$, and \mathbf{A} an n by n matrix , we have these possibilities :

- (i) \mathbf{A} is nonsingular : In this case there is a unique solution.
- (ii) \mathbf{A} is singular : There are no solutions or infinitely many. (Examples?)

Usually only case (i) is of interest.

The solution can be computed in a finite number of steps by Gauss Elimination (with pivoting, if necessary).

We can write a system of n *nonlinear* equations in n unknowns as

$$\mathbf{G}(\mathbf{x}) = \mathbf{0} ,$$

where

$$\mathbf{x} , \mathbf{0} \in \mathbb{R}^n ,$$

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T ,$$

and where \mathbf{G} is a vector-valued function of \mathbf{x} having n component functions,

$$\mathbf{G}(\mathbf{x}) = (g_1(\mathbf{x}) , g_2(\mathbf{x}) , \dots , g_n(\mathbf{x}))^T .$$

EXAMPLES (of possible situations) :

○ $x^2 - 1 = 0$ has two solutions : $x = 1$, $x = -1$.

○ $x^2 + 1 = 0$ has two solutions : $x = i$, $x = -i$, $(i \equiv \sqrt{-1})$.

○ The system
$$\begin{cases} 2x_1 - x_2 = 0 , \\ x_1^3 + x_1 - x_2 = 0 , \end{cases}$$
 has three solution pairs,

namely, $(x_1, x_2) = (0, 0)$, $(1, 2)$, $(-1, -2)$.

○ The system
$$\begin{cases} x_1x_2 - 1 = 0 , \\ x_1x_2 - 2 = 0 , \end{cases}$$
 has no solutions.

○ The system
$$\begin{cases} e^{x_1 - x_2} - 1 = 0 , \\ x_1 - x_2 = 0 , \end{cases}$$
 has infinitely many solution pairs.

REMARKS:

For nonlinear equations :

- There can be 0, 1, 2, 3, 4, \dots ∞ solutions.
- A solution can usually not be computed in a finite number of steps.
- Instead, *iterative methods* will be used.
- We will not consider the case of a *continuum of solutions*.

Some Methods for Scalar Nonlinear Equations.

Consider the *scalar* equation (one equation, one unknown)

$$g(x) = 0 ,$$

and let x^* denote a *zero* (or *root*) of this equation.

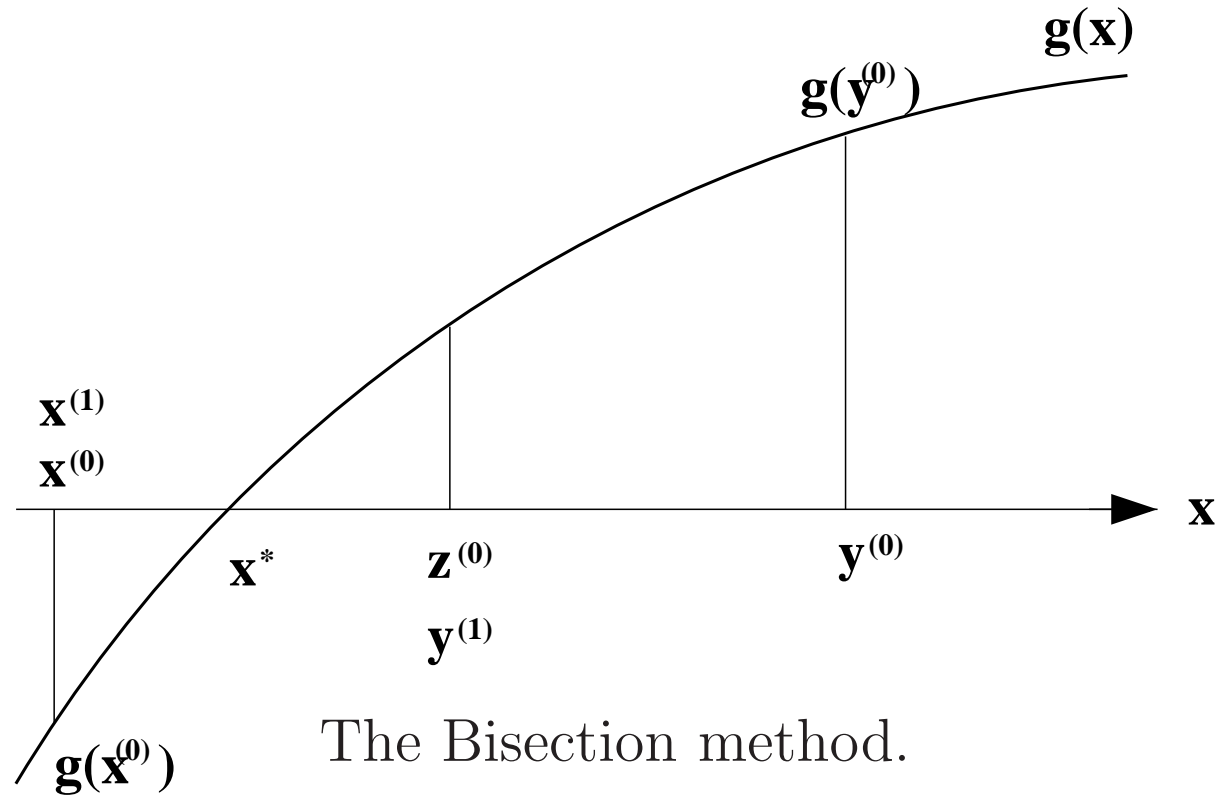
The Bisection Method.

This method requires two initial points :

$$x^{(0)} , y^{(0)} \quad \text{with} \quad g(x^{(0)}) < 0 \quad \text{and} \quad g(y^{(0)}) > 0 .$$

Algorithm ($k = 0, 1, 2, \dots$) :

- Set $z^{(k)} = \frac{1}{2}(x^{(k)} + y^{(k)})$,
- If $g(z^{(k)}) < 0$ set $x^{(k+1)} = z^{(k)}$, $y^{(k+1)} = y^{(k)}$,
- If $g(z^{(k)}) > 0$ set $x^{(k+1)} = x^{(k)}$, $y^{(k+1)} = z^{(k)}$.



The bisection method works if $g(x)$ is continuous in the interval $[x^{(0)}, y^{(0)}]$.

In fact we have

$$|x^{(k)} - x^*| \leq \frac{1}{2^k} |x^{(0)} - y^{(0)}| .$$

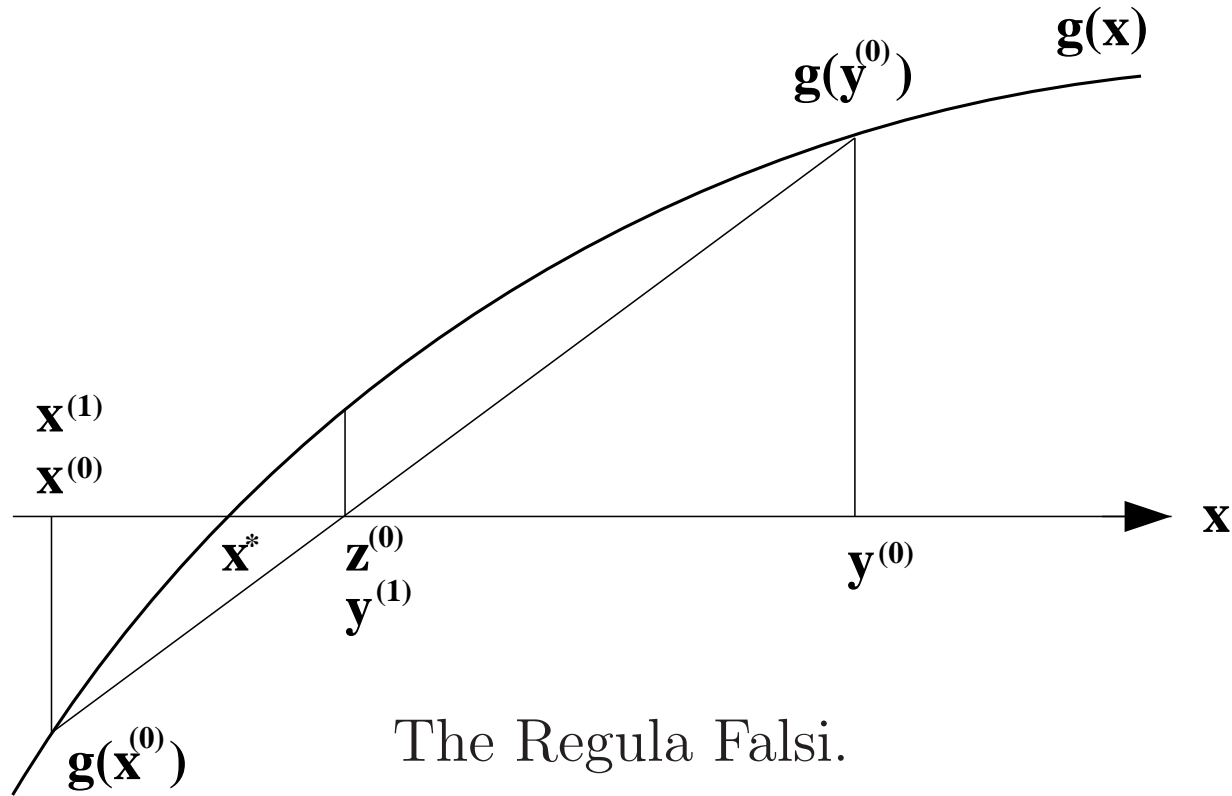
The method does not readily generalize to systems of nonlinear equations.

The Regula Falsi.

This method is similar to the bisection method.

However, in each step we now let

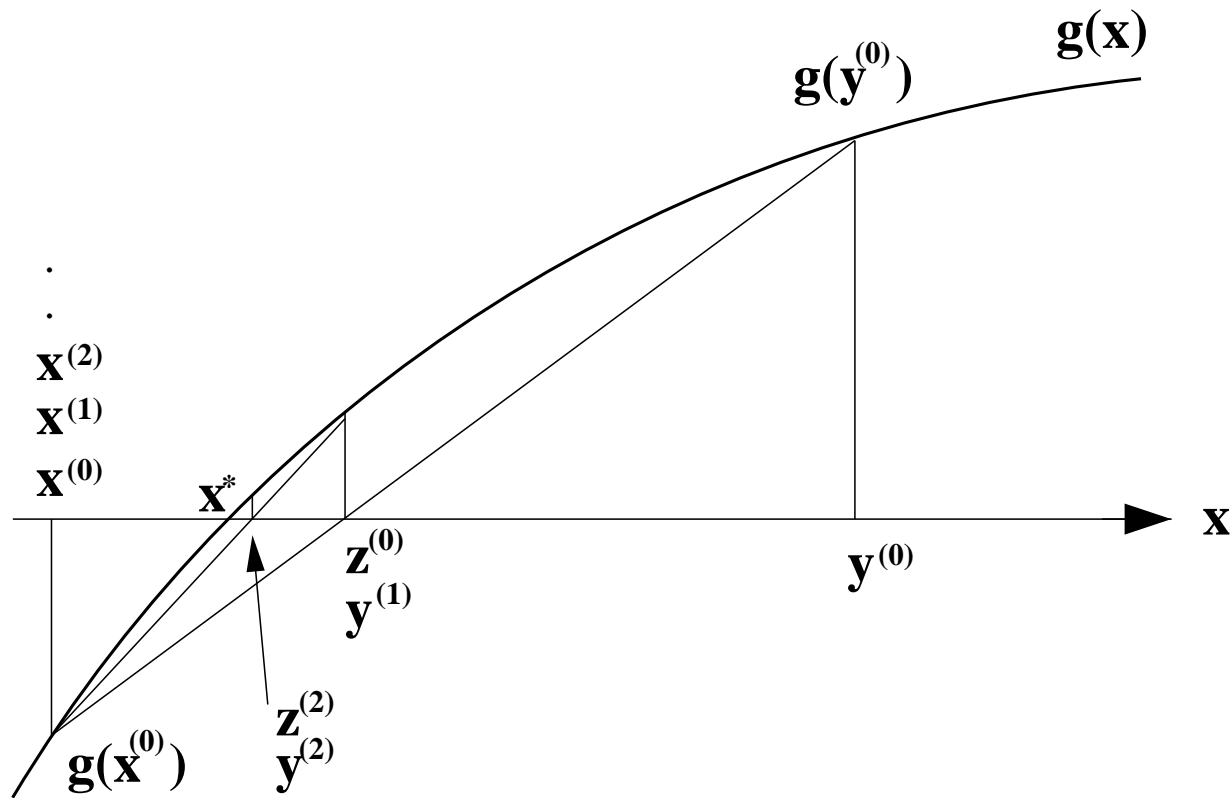
$$z^{(k)} = \frac{x^{(k)}g(y^{(k)}) - y^{(k)}g(x^{(k)})}{g(y^{(k)}) - g(x^{(k)})}.$$



$$z^{(k)} = \frac{x^{(k)}g(y^{(k)}) - y^{(k)}g(x^{(k)})}{g(y^{(k)}) - g(x^{(k)})}.$$

$z^{(k)}$ is the zero of the line connecting $(x^{(k)}, g(x^{(k)}))$ to $(y^{(k)}, g(y^{(k)}))$. (Check!)

Unlike the bisection method, not both $x^{(k)}$ and $y^{(k)}$ need converge to x^* :



The Regula Falsi : Nonconvergence of $x^{(k)}$ to x^* .

The Regula Falsi does not readily generalize to nonlinear systems.

Newton's Method.

Let $x^{(0)}$ be an initial guess for a zero x^* of $g(x) = 0$.

The line $p_0(x)$ that satisfies

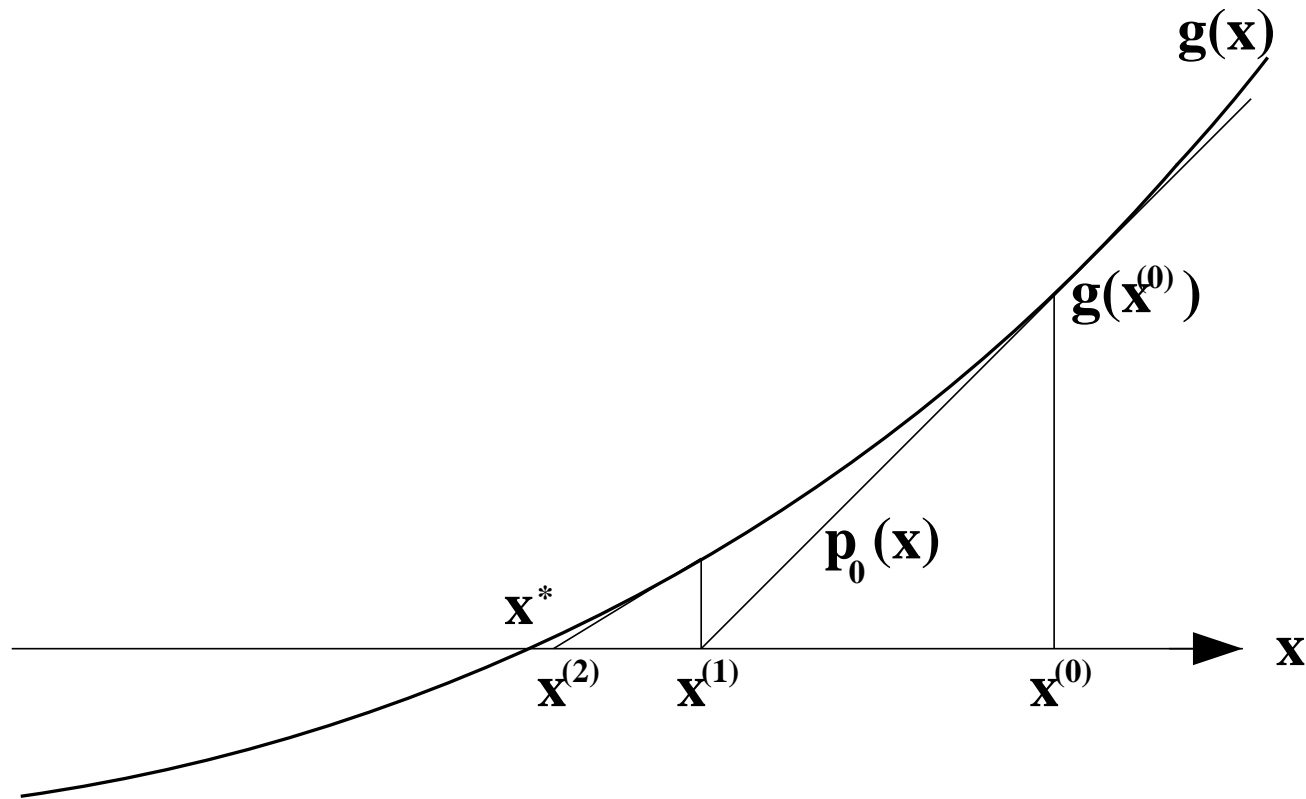
$$p_0(x^{(0)}) = g(x^{(0)})$$

and

$$p'_0(x^{(0)}) = g'(x^{(0)}) ,$$

is given by

$$p_0(x) = g(x^{(0)}) + (x - x^{(0)}) g'(x^{(0)}) .$$



Newton's method

If g is sufficiently smooth and if $x^{(0)}$ is close to x^* then we expect the zero

$$x^{(1)} = x^{(0)} - \frac{g(x^{(0)})}{g'(x^{(0)})}, \quad (\text{Check!})$$

of $p_0(x)$ to be a better approximation to x^* than $x^{(0)}$.

This procedure may now be repeated for the point $x^{(1)}$.

The general algorithm for Newton's method can therefore be written as

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(k)})}, \quad k = 0, 1, 2, \dots .$$

Later we show that Newton's method converges to a zero x^* of $g(x) = 0$ if

- g has two continuous derivatives near x^* ,
- $g'(x^*) \neq 0$,
- $x^{(0)}$ is sufficiently close to x^* .

EXAMPLE:

Use Newton's method to compute the square root of 2.

Note that this square root is a zero of $g(x) \equiv x^2 - 2$.

Thus the Newton iteration is given by

$$x^{(k+1)} = x^{(k)} - \frac{(x^{(k)})^2 - 2}{2x^{(k)}} = \frac{(x^{(k)})^2 + 2}{2x^{(k)}} .$$

With $x^{(0)} = 1.5$, we get

$$x^{(1)} = 1.41666, \quad x^{(2)} = 1.414215, \quad \textit{etc.}$$

Newton's method generalizes to systems of nonlinear equations.

This extension will be considered later.

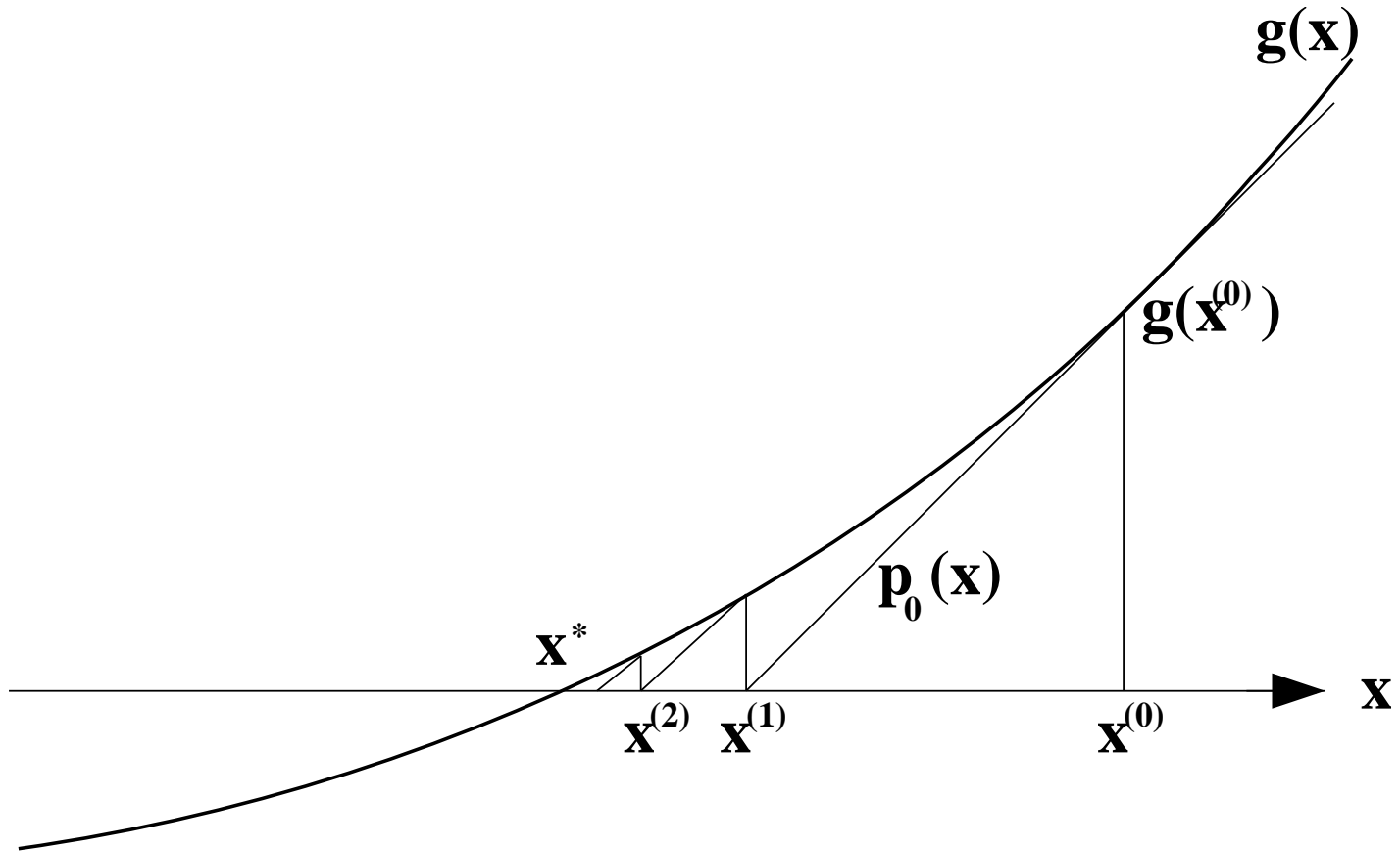
The Chord Method.

This method is similar to Newton's method.

The only difference is that $g'(x)$ is always evaluated at the initial point $x^{(0)}$.

Thus the algorithm is

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(0)})} .$$



The Chord method

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(0)})} .$$

Compared to Newton's method :

- The Chord method takes fewer arithmetic operations per iteration.
- The two methods converge under essentially the same conditions.
- The Chord method needs more iterations for a prescribed accuracy.

EXAMPLE:

With $x^{(0)} = 1.5$ the Chord method for solving $x^2 - 2 = 0$ takes the form

$$x^{(k+1)} = x^{(k)} - \frac{(x^{(k)})^2 - 2}{3} .$$

The first few iterations are

$$x^{(1)} = 1.416666 , \quad x^{(2)} = 1.414351 , \quad x^{(3)} = 1.414221 .$$

EXERCISES:

- Show how to use *Newton's method* to compute the cube root of 2. Carry out the first few iterations, using $x^{(0)} = 0.6$.
- Show how to use *the Chord method* to compute the cube root of 2. Carry out the first few iterations, using $x^{(0)} = 0.6$.

- Consider the equation

$$\sin(x) = 1/x .$$

Show the graphs of $\sin(x)$ and $1/x$ in one diagram. How many solutions are there to this equation ? Write down Newton's method for finding a solution. Carry out the first few iterations with $x^{(0)} = \pi/2$.

- Consider the equation $\sin(x) = e^{-x}$. Draw the functions $\sin(x)$ and e^{-x} in one graph. How many solutions are there to the above equation ? Show how one can use Newton's method to find a solution of the equation. Carry out the first few Newton iterations, using $x^{(0)} = 0$.

Newton's Method for Systems of Nonlinear Equations.

First reconsider Newton's method for scalar equations

$$g(x) = 0 .$$

Given $x^{(k)}$, we set $x^{(k+1)}$ to be the zero of

$$p_k(x) \equiv g(x^{(k)}) + (x - x^{(k)}) g'(x^{(k)}) .$$

REMARK:

- Here $p_k(x)$ is the tangent line to $g(x)$ at $x = x^{(k)}$, *i.e.*,
- $p_k(x)$ is the *linear approximation* to $g(x)$ at $x = x^{(k)}$, *i.e.*,
- $p_k(x)$ is the *degree 1 Taylor polynomial* of $g(x)$ at $x = x^{(k)}$.

Similarly for systems of the form

$$\mathbf{G}(\mathbf{x}) = \mathbf{0} ,$$

we have the linear approximation

$$\mathbf{P}_k(\mathbf{x}) \equiv \mathbf{G}(\mathbf{x}^{(k)}) + \mathbf{G}'(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)})$$

of $\mathbf{G}(\mathbf{x})$ about $\mathbf{x} = \mathbf{x}^{(k)}$.

Here $\mathbf{G}'(\mathbf{x}^{(k)})$ is the *Jacobian matrix* of $\mathbf{G}(\mathbf{x})$ at $\mathbf{x} = \mathbf{x}^{(k)}$.

Analogous to the scalar case we let $\mathbf{x}^{(k+1)}$ be the zero of $\mathbf{P}_k(\mathbf{x}) = \mathbf{0}$.

Thus $\mathbf{x}^{(k+1)}$ is the solution of the linear system

$$\mathbf{G}'(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = -\mathbf{G}(\mathbf{x}^{(k)}) ,$$

that is, we can get $\mathbf{x}^{(k+1)}$ by first solving

$$(1) \quad \mathbf{G}'(\mathbf{x}^{(k)}) \Delta\mathbf{x}^{(k)} = -\mathbf{G}(\mathbf{x}^{(k)}) ,$$

and then setting

$$(2) \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta\mathbf{x}^{(k)} .$$

EXAMPLE:

Use Newton's method to solve the system

$$x_1^2 x_2 - 1 = 0 ,$$

$$x_2 - x_1^4 = 0 .$$

Here

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} , \quad \mathbf{G}(\mathbf{x}) = \begin{pmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} g_1(x_1, x_2) \\ g_2(x_1, x_2) \end{pmatrix} = \begin{pmatrix} x_1^2 x_2 - 1 \\ x_2 - x_1^4 \end{pmatrix} .$$

The Jacobian matrix in this example is given by

$$\mathbf{G}'(\mathbf{x}) \equiv \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} \end{pmatrix} = \begin{pmatrix} 2x_1 x_2 & x_1^2 \\ -4x_1^3 & 1 \end{pmatrix} .$$

Hence Newton's method for this problem takes the form

$$(1) \quad \begin{pmatrix} 2x_1^{(k)}x_2^{(k)} & (x_1^{(k)})^2 \\ -4(x_1^{(k)})^3 & 1 \end{pmatrix} \begin{pmatrix} \Delta x_1^{(k)} \\ \Delta x_2^{(k)} \end{pmatrix} = \begin{pmatrix} 1 - (x_1^{(k)})^2x_2^{(k)} \\ (x_1^{(k)})^4 - x_2^{(k)} \end{pmatrix},$$

$$(2) \quad \begin{cases} x_1^{(k+1)} = x_1^{(k)} + \Delta x_1^{(k)}, \\ x_2^{(k+1)} = x_2^{(k)} + \Delta x_2^{(k)}, \end{cases}$$

for $k = 0, 1, 2, \dots$.

Thus for each iteration two linear equations in two unknowns must be solved.

With the initial guess

$$x_1^{(0)} = 2 , \quad x_2^{(0)} = 2 ,$$

the first iteration consists of solving

$$\begin{pmatrix} 8 & 4 \\ -32 & 1 \end{pmatrix} \begin{pmatrix} \Delta x_1^{(0)} \\ \Delta x_2^{(0)} \end{pmatrix} = \begin{pmatrix} -7 \\ 14 \end{pmatrix} ,$$

which gives

$$\Delta x_1^{(0)} = -0.463 , \quad \Delta x_2^{(0)} = -0.823 ,$$

and then setting

$$x_1^{(1)} = x_1^{(0)} + \Delta x_1^{(0)} = 1.537 ,$$

$$x_2^{(1)} = x_2^{(0)} + \Delta x_2^{(0)} = 1.177 .$$

After a second iteration what will $x_1^{(2)}$ and $x_2^{(2)}$ be ?

EXERCISES:

- Describe in detail how Newton's method can be used to compute solutions (x_1, x_2) of the system of two nonlinear equations

$$x_1^2 + x_2^2 - 1 = 0 ,$$

$$x_2 - e^{x_1} = 0 .$$

- Describe in detail how Newton's method can be used to compute a solution (x_1, x_2, x_3) of the system of three nonlinear equations

$$x_1^2 + x_2^2 + x_3^2 - 1 = 0 ,$$

$$x_3 - e^{x_1} = 0 ,$$

$$x_3 - e^{x_2} = 0 .$$

Residual Correction.

Suppose we use Newton's method for a *linear system* $\mathbf{Ax} = \mathbf{f}$,

that is, we let

$$\mathbf{G}(\mathbf{x}) \equiv \mathbf{Ax} - \mathbf{f} .$$

Then

$$\mathbf{G}'(\mathbf{x}) = \mathbf{A} ,$$

so that Newton's method

$$\mathbf{G}'(\mathbf{x}^{(k)}) \Delta \mathbf{x}^{(k)} = -\mathbf{G}(\mathbf{x}^{(k)}) ,$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)} ,$$

becomes

$$\mathbf{A} \Delta \mathbf{x}^{(k)} = -(\mathbf{Ax}^{(k)} - \mathbf{f}) ,$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)} .$$

$$\mathbf{A} \Delta \mathbf{x}^{(k)} = - (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{f}) ,$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)} .$$

NOTE:

- the Jacobian needs to be LU-decomposed only once.
- With *exact arithmetic* , the *exact solution* is found in *one iteration*:

$$\Delta \mathbf{x}^{(0)} = - \mathbf{A}^{-1} (\mathbf{A} \mathbf{x}^{(0)} - \mathbf{f}) = - \mathbf{x}^{(0)} + \mathbf{x} ,$$

so that

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \Delta \mathbf{x}^{(0)} = \mathbf{x}^{(0)} - \mathbf{x}^{(0)} + \mathbf{x} = \mathbf{x} .$$

$$\mathbf{A} \Delta \mathbf{x}^{(k)} = - (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{f}) ,$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)} .$$

REMARKS:

- For *inexact arithmetic* , this iteration is called *residual correction* .
- Residual correction can improve the accuracy of the solution of $\mathbf{A} \mathbf{x} = \mathbf{f}$.
- Residual correction is valuable for mildly ill-conditioned linear systems.
- The “*residual*” $\mathbf{A} \mathbf{x}^{(k)} - \mathbf{f}$ should be computed with high precision.

Convergence Analysis for Scalar Equations.

Most iterative methods for solving a scalar equation $g(x) = 0$ can be written

$$x^{(k+1)} = f(x^{(k)}) , \quad k = 0, 1, 2, \dots, \quad x^{(0)} \text{ given .}$$

For example, in Newton's method

$$f(x) = x - \frac{g(x)}{g'(x)} ,$$

and in the Chord method

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

Sometimes the iteration $x^{(k+1)} = x^{(k)} - g(x^{(k)})$ also works. In this method

$$f(x) = x - g(x) .$$

Iterations of the form

$$x^{(k+1)} = f(x^{(k)}) , \quad k = 0, 1, 2, \dots, \quad x^{(0)} \text{ given} ,$$

also arise independently, *e.g.*, as models of “population growth” :

$$x^{(k+1)} = c x^{(k)} , \quad k = 0, 1, 2, \dots ,$$

models exponential population growth.

QUESTION: What happens to the sequence $x^{(k)}$, $k = 0, 1, 2, \dots$,

- when $x^{(0)} > 0$ and $c > 1$?
- when $x^{(0)} > 0$ and $c < 1$?

The iteration

$$x^{(k+1)} = c x^{(k)} (1 - x^{(k)}), \quad k = 0, 1, 2, \dots,$$

known as the *logistic equation*, models population growth when there are limited resources.

QUESTION: For $0 < x^{(0)} < 1$:

What happens to the sequence $x^{(k)}$, $k = 0, 1, 2, \dots$,

- when $0 \leq c < 1$?
- when $1 \leq c < 2$?
- when $2 \leq c < 3$?
- when $3 \leq c \leq 4$?

In general, an iteration of the form

$$x^{(k+1)} = f(x^{(k)}) , \quad k = 0, 1, 2, \dots ,$$

is often called a *fixed point iteration* .

Suppose the sequence $x^{(k)}$, $k = 0, 1, 2, \dots$, converges, *i.e.*,

$$\lim_{k \rightarrow \infty} x^{(k)} = x^* .$$

Then x^* satisfies the equation

$$x = f(x) ,$$

(assuming that f is continuous near x^*).

We call x^* a *fixed point* of f .

EXAMPLE:

In Newton's method

$$f(x) = x - \frac{g(x)}{g'(x)} .$$

Thus a fixed point x^* satisfies

$$x^* = x^* - \frac{g(x^*)}{g'(x^*)} ,$$

that is,

$$g(x^*) = 0 ,$$

(assuming that $g'(x^*) \neq 0$.)

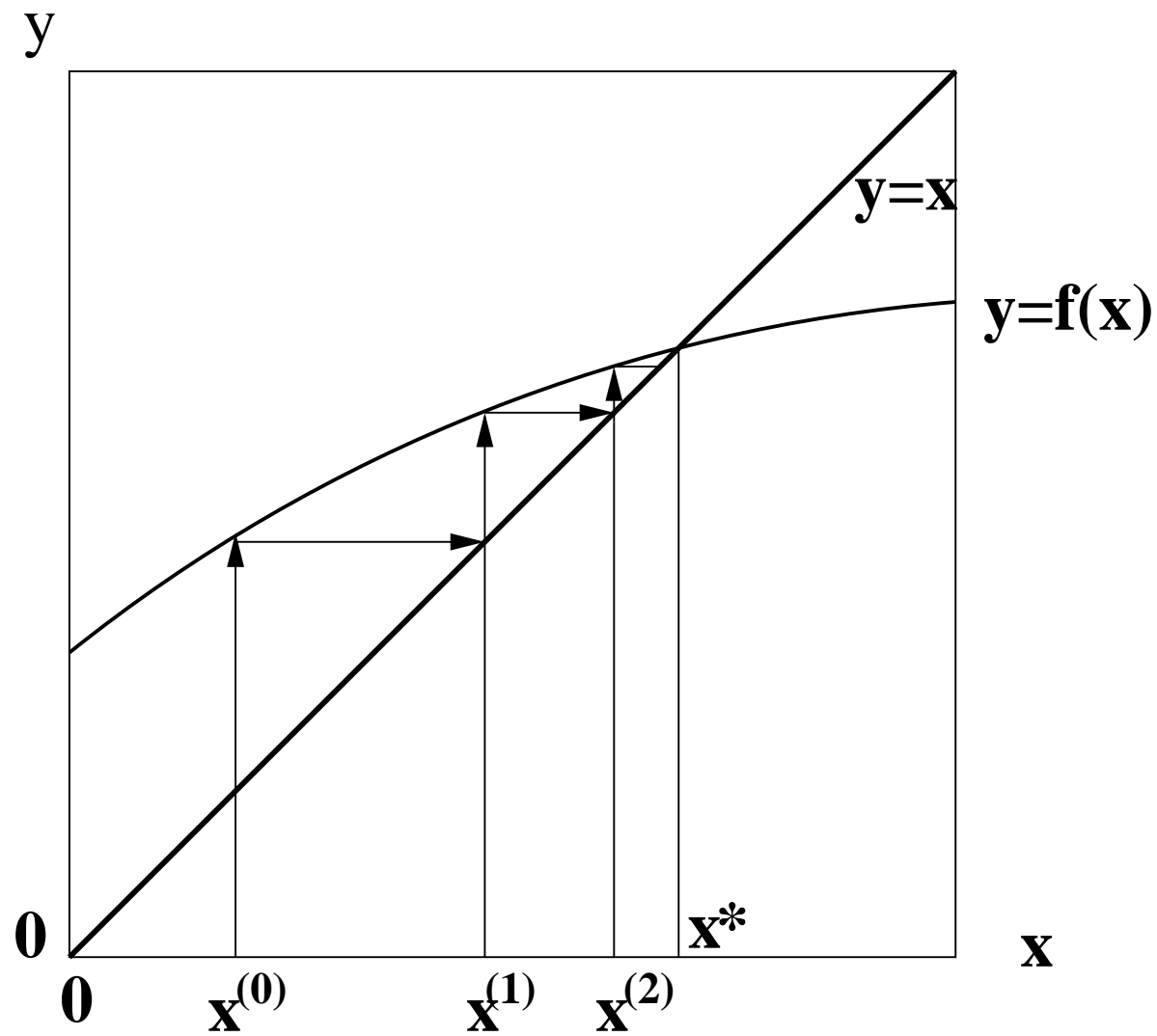
Thus x^* is a solution of $g(x) = 0$.

Assuming that f has a fixed point, when does the fixed point iteration

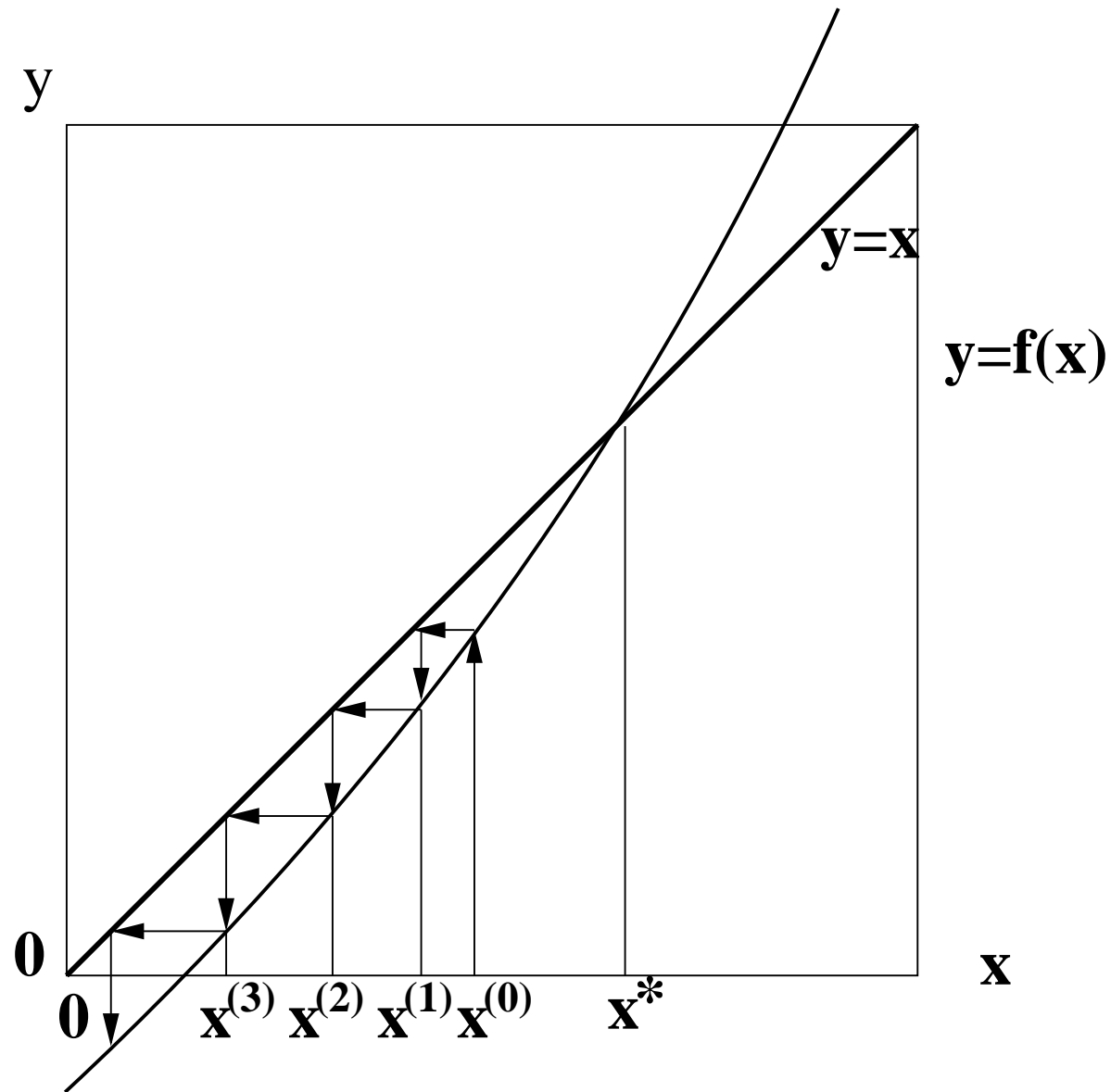
$$x^{(k+1)} = f(x^{(k)}) ,$$

converge ?

The answer is suggested in the following two diagrams :



A convergent fixed point iteration.



A divergent fixed point iteration.

THEOREM:

Let $f'(x)$ be continuous near a fixed point x^* of $f(x)$, and assume that

$$| f'(x^*) | < 1 .$$

Then the fixed point iteration

$$x^{(k+1)} = f(x^{(k)}) , \quad k = 0, 1, 2, \dots ,$$

converges to x^* , whenever the initial guess $x^{(0)}$ is sufficiently close to x^* .

PROOF:

Let $\alpha \equiv |f'(x^*)|$.

Then $\alpha < 1$.

Choose β such that $\alpha < \beta < 1$.

Then, for some $\epsilon > 0$, there exists an interval

$$I_\epsilon \equiv [x^* - \epsilon, x^* + \epsilon],$$

such that

$$|f'(x)| \leq \beta \quad \text{in} \quad I_\epsilon,$$

(because f' is continuous near x^*).

Let $x^{(0)} \in I_\epsilon$.

By Taylor's Theorem (or by the Mean Value Theorem)

$$x^{(1)} \equiv f(x^{(0)}) = f(x^*) + (x^{(0)} - x^*) f'(\eta_0) ,$$

for some η_0 between $x^{(0)}$ and x^* .

Since $f(x^*) = x^*$ it follows that

$$|x^{(1)} - x^*| = |(x^{(0)} - x^*) f'(\eta_0)| = |x^{(0)} - x^*| |f'(\eta_0)| \leq \beta |x^{(0)} - x^*| .$$

Thus $x^{(1)} \in I_\epsilon$, (because $0 < \beta < 1$) .

Again by Taylor's Theorem (or the Mean Value Theorem)

$$x^{(2)} \equiv f(x^{(1)}) = f(x^*) + (x^{(1)} - x^*) f'(\eta_1) ,$$

for some η_1 between $x^{(1)}$ and x^* .

Hence

$$| x^{(2)} - x^* | \leq \beta | x^{(1)} - x^* | \leq \beta^2 | x^{(0)} - x^* | .$$

Thus $x^{(2)} \in I_\epsilon$, (because $0 < \beta < 1$) .

Proceeding in this manner we find

$$| x^{(k)} - x^* | \leq \beta^k | x^{(0)} - x^* | .$$

Since $0 < \beta < 1$ this implies that

$$x^{(k)} \rightarrow x^* \quad \text{as} \quad k \rightarrow \infty .$$

COROLLARY: Let

$$I_\epsilon \equiv [x^* - \epsilon , x^* + \epsilon] ,$$

and assume that for some $\epsilon > 0$ we have :

- $f(x)$ has a fixed point $x^* \in I_\epsilon$,
- $f(x)$ is a smooth function in I_ϵ ,
- $| f'(x) | < 1$ everywhere in I_ϵ ,
- $x^{(0)} \in I_\epsilon$.

Then the fixed point iteration

$$x^{(k+1)} = f(x^{(k)}) , \quad k = 0, 1, 2, \dots .$$

converges to x^* .

PROOF: This follows from the proof of the Theorem.

COROLLARY:

If

- x^* is a zero of $g(x) = 0$,
- $g(x)$ has two continuous derivatives near x^* ,
- $g'(x^*) \neq 0$,
- $x^{(0)}$ is sufficiently close to x^* .

then Newton's method for solving $g(x) = 0$, *i.e.*,

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(k)})} ,$$

converges to x^* .

PROOF:

In Newton's method

$$f(x) = x - \frac{g(x)}{g'(x)} .$$

Hence

$$f'(x^*) = 1 - \frac{g'(x^*)^2 - g(x^*) g''(x^*)}{g'(x^*)^2} = \frac{g(x^*) g''(x^*)}{g'(x^*)^2} = 0 .$$

Therefore, certainly $|f'(x^*)| < 1$.

EXAMPLE:

The fixed points of the logistic equation,

$$x^{(k+1)} = c x^{(k)} (1 - x^{(k)}) ,$$

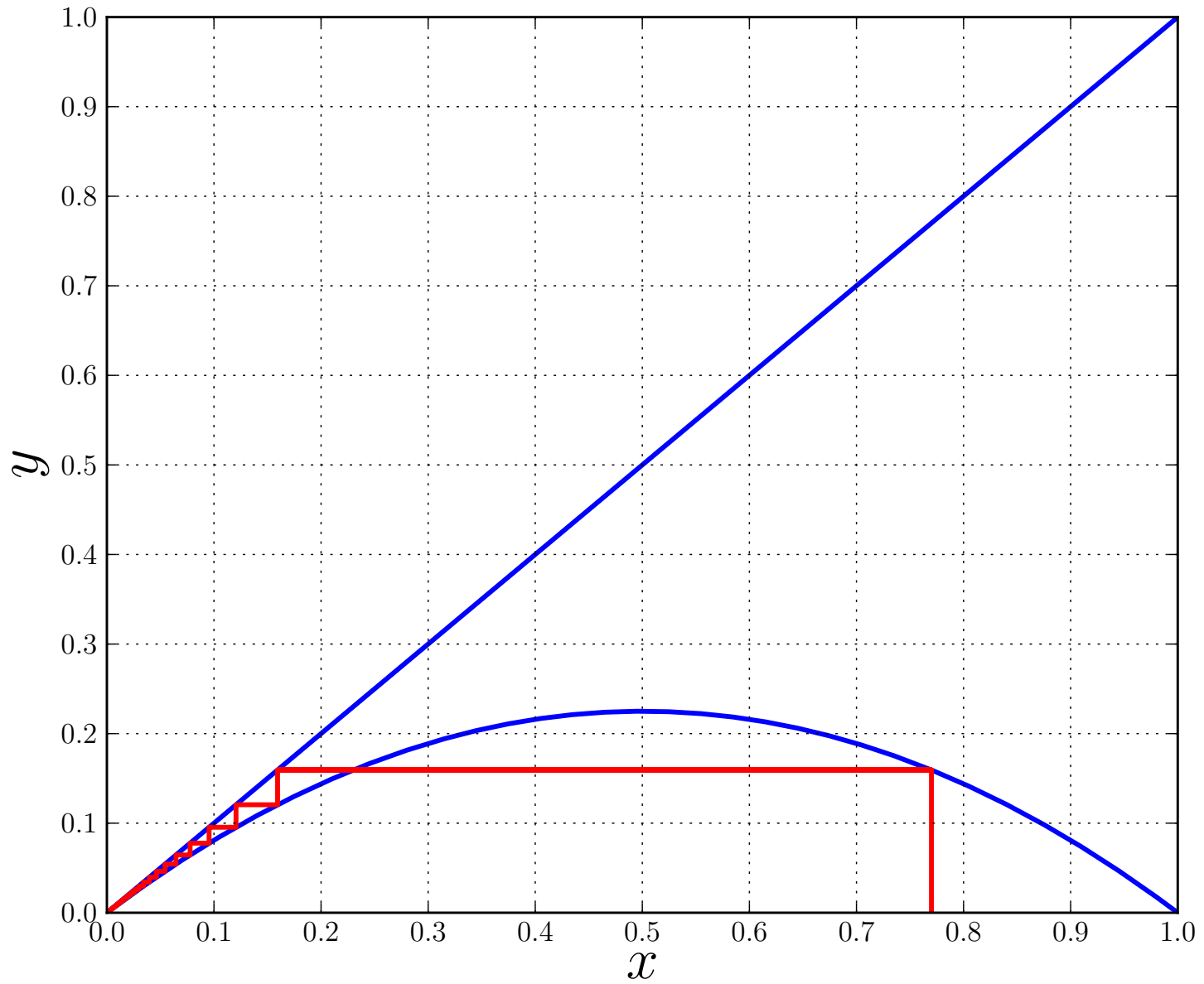
satisfy

$$x^* = c x^* (1 - x^*) .$$

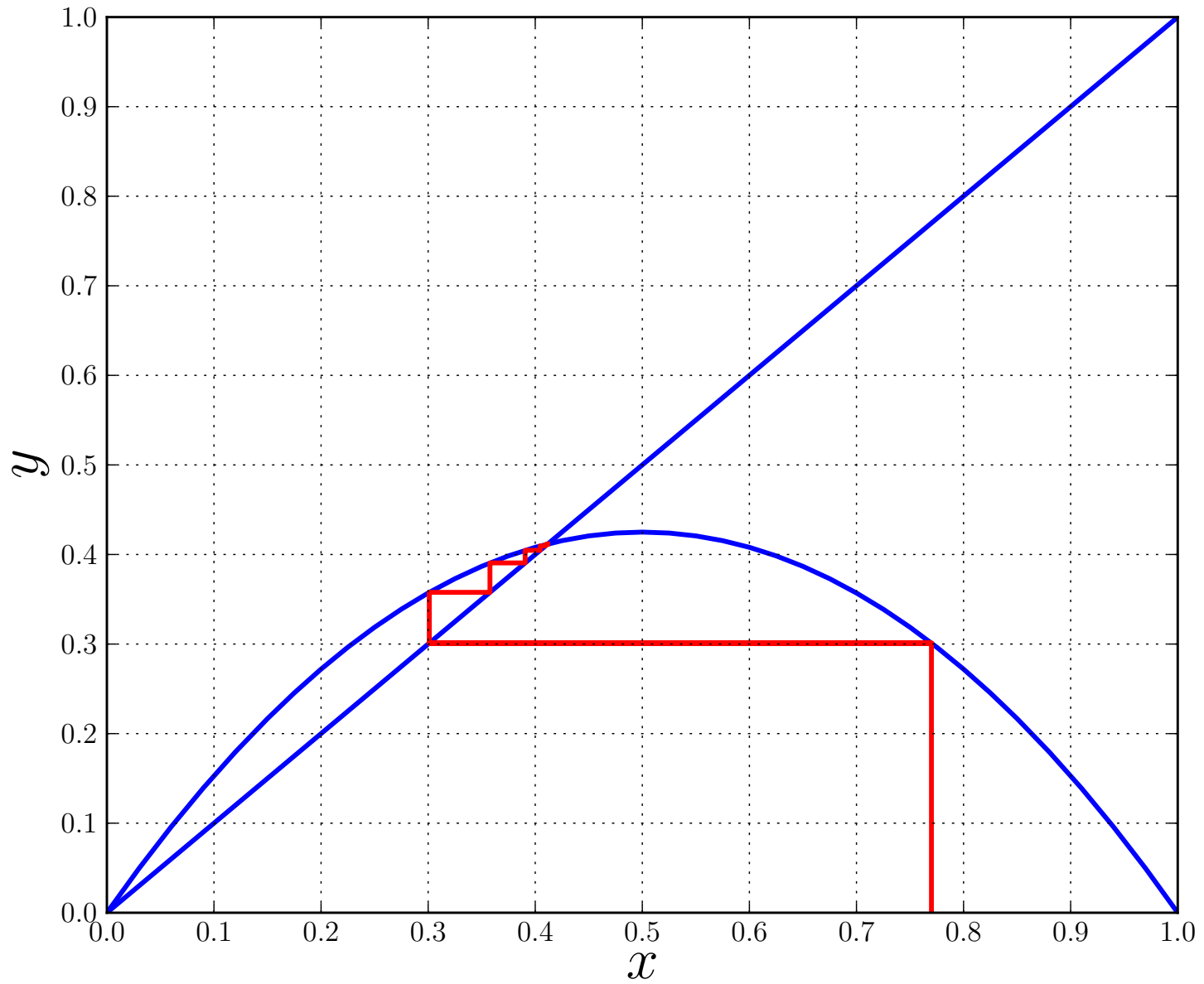
We see that the fixed points are given by

$$x^* = 0 , \quad \text{and} \quad x^* = 1 - \frac{1}{c} .$$

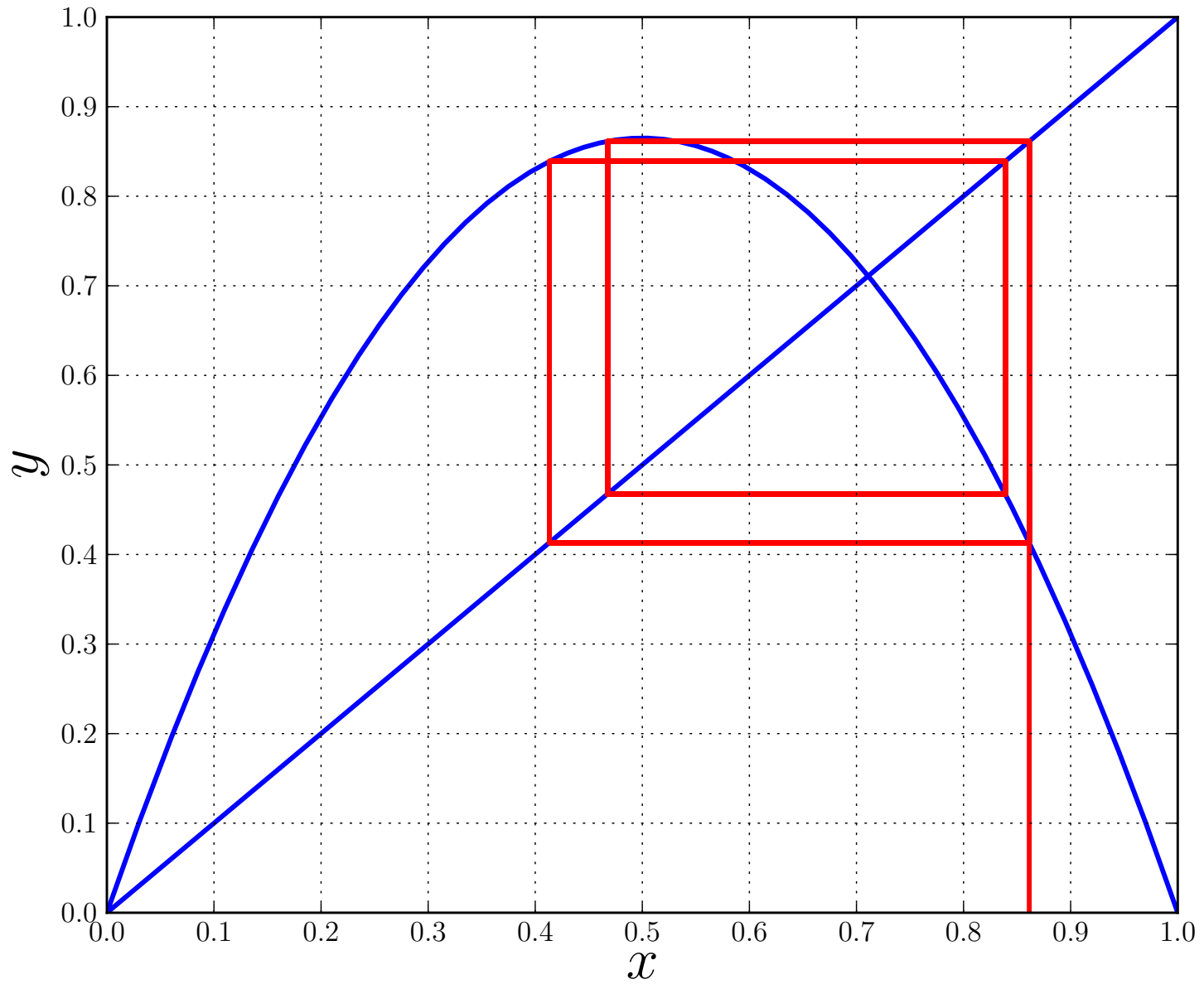
EXERCISE: Determine, for all values of c , $(0 < c \leq 4)$, whether these fixed points are *attracting* ($| f'(x^*) | < 1$) , or *repelling* ($| f'(x^*) | > 1$) .



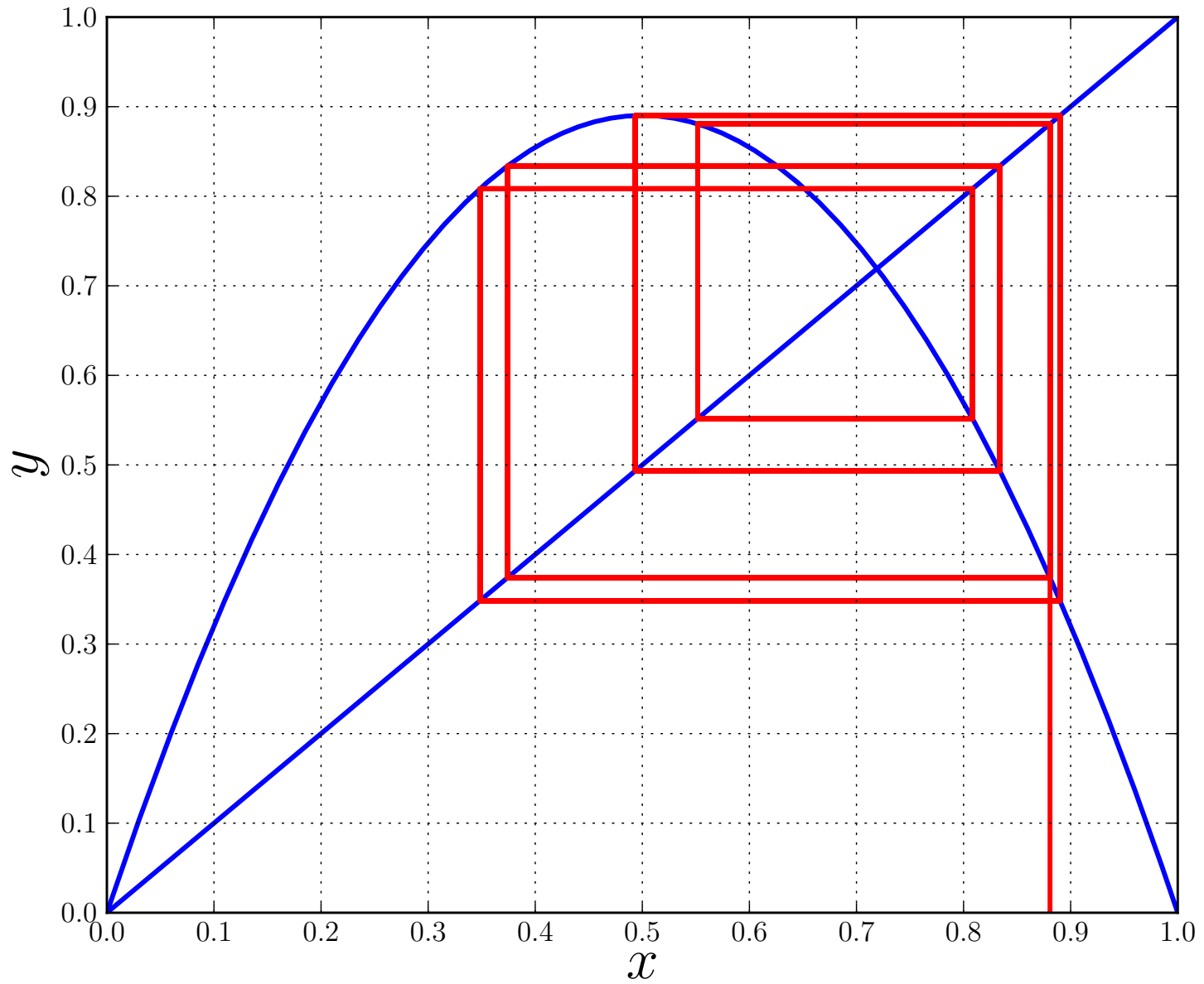
The logistic equation : $c = 0.9$.



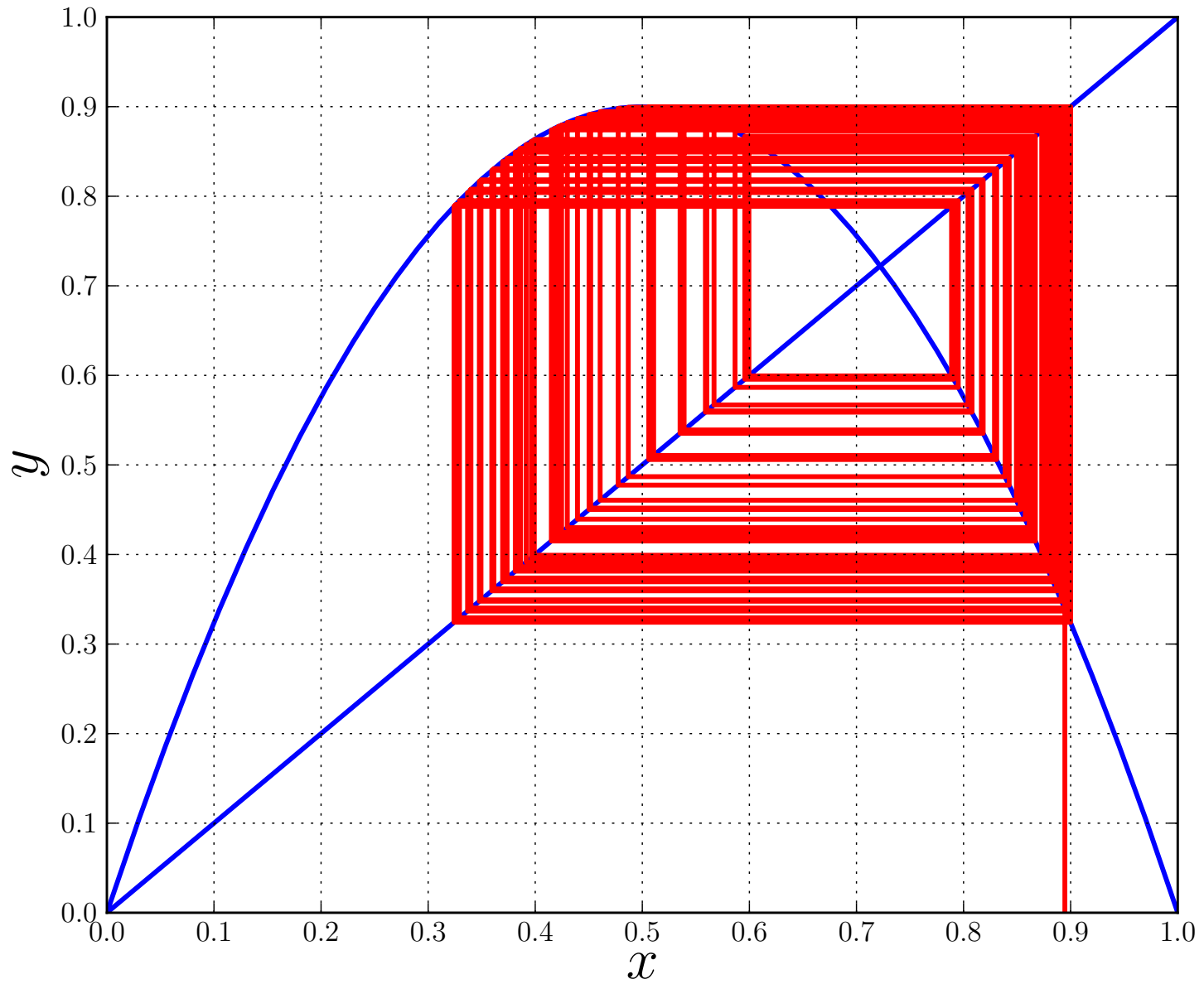
The logistic equation : $c = 1.7$.



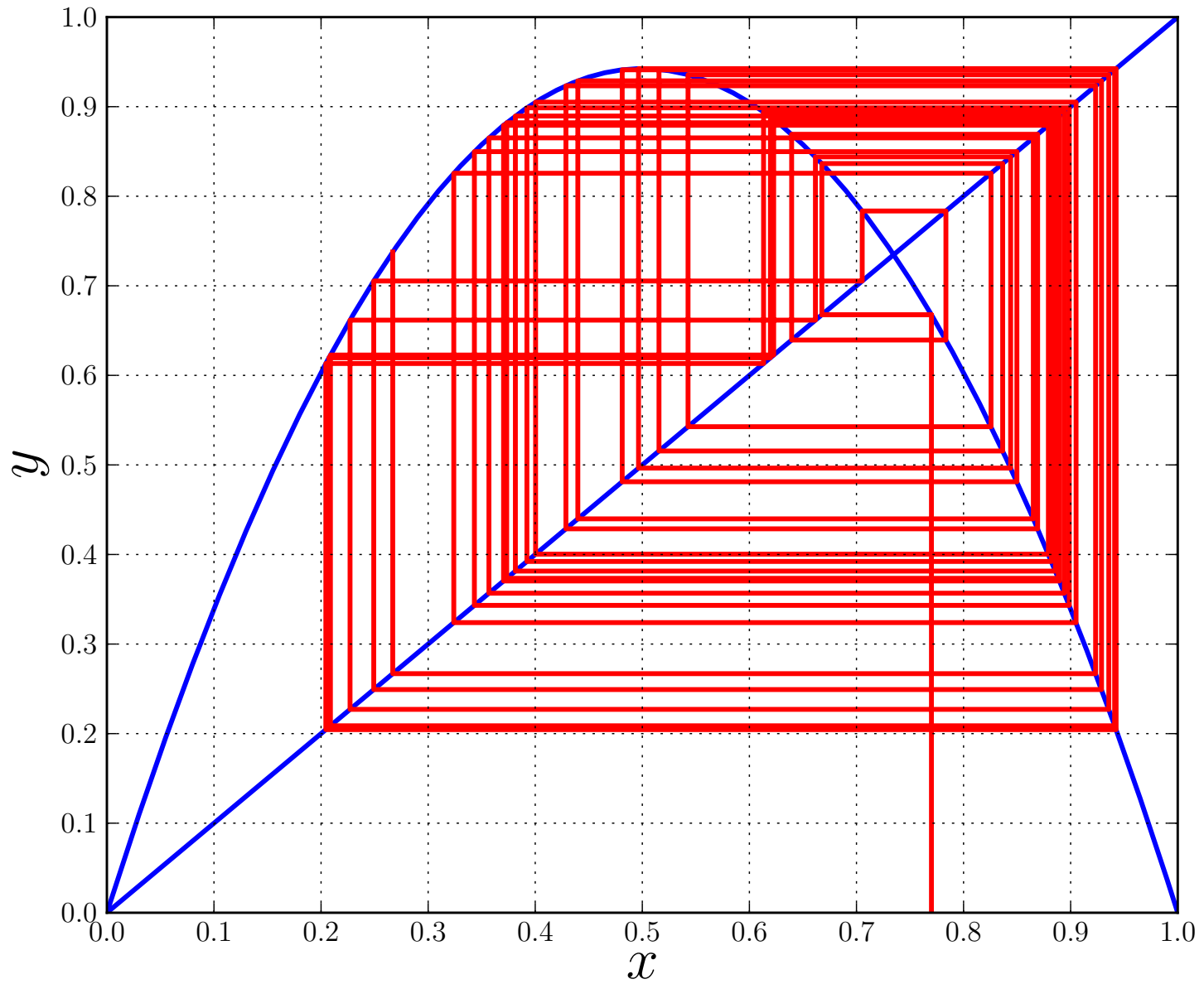
The logistic equation : $c = 3.46$.



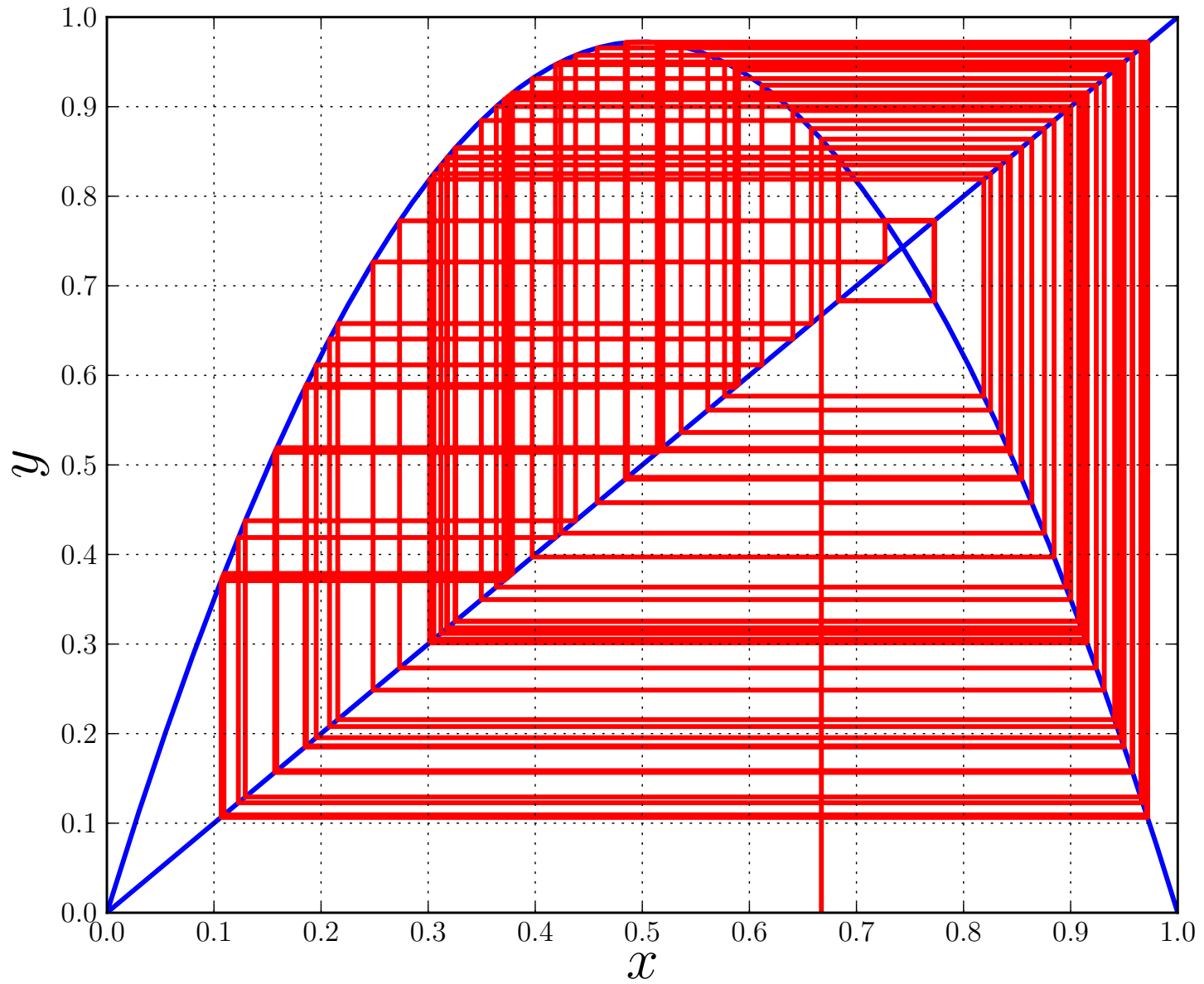
The logistic equation : $c = 3.561$.



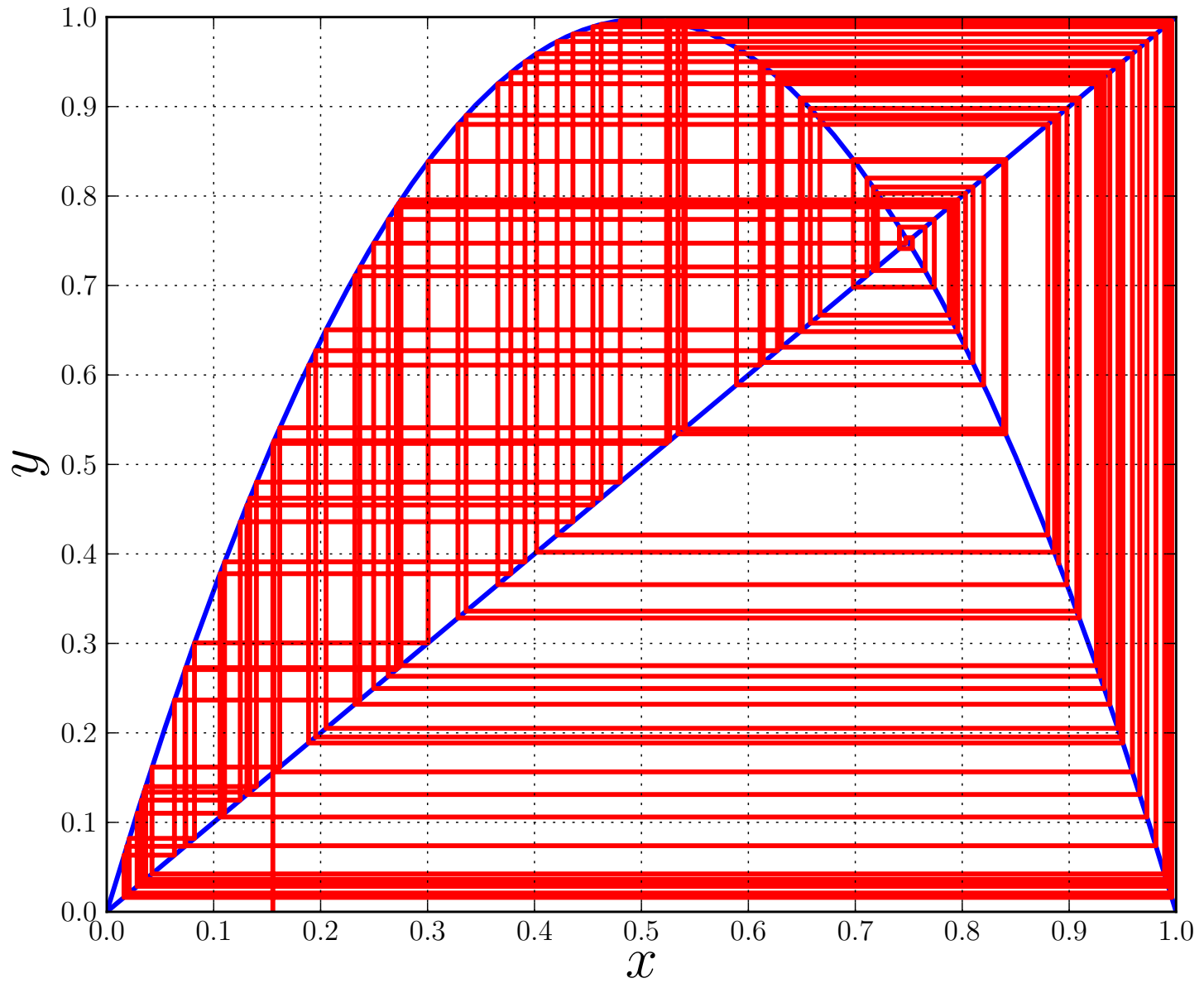
The logistic equation : $c = 3.6$.



The logistic equation : $c = 3.77$.



The logistic equation : $c = 3.89$.



The logistic equation : $c = 3.99$.

If a fixed point iteration

$$x^{(k+1)} = f(x^{(k)}) ,$$

converges to a fixed point x^* of $f(x)$, then *how fast* does it converge ?

To answer this we let

$$e_k \equiv |x^{(k)} - x^*| .$$

Thus e_k is the error after the k th iteration.

We can now show the following :

THEOREM:

Let $f'(x)$ be continuous near x^* with $|f'(x^*)| < 1$.

Assume that $x^{(0)}$ is sufficiently close to x^* , so that the fixed point iteration

$$x^{(k+1)} = f(x^{(k)})$$

converges to x^* .

Then
$$\lim_{k \rightarrow \infty} \frac{e_{k+1}}{e_k} = |f'(x^*)| \quad (\text{linear convergence}).$$

(The value of $|f'(x^*)|$ is then called *the rate of convergence*.)

If in addition $f'(x^*) = 0$ and if $f''(x)$ is continuous near x^* then

$$\lim_{k \rightarrow \infty} \frac{e_{k+1}}{e_k^2} = \frac{1}{2} |f''(x^*)| \quad (\text{quadratic convergence}).$$

PROOF:

Case : $f'(x^*) \neq 0$:

$$\begin{aligned} e_{k+1} &= |x^{(k+1)} - x^*| \\ &= |f(x^{(k)}) - x^*| \\ &= |f(x^*) + (x^{(k)} - x^*) f'(\eta_k) - x^*| \\ &= |x^{(k)} - x^*| |f'(\eta_k)| \\ &= e_k |f'(\eta_k)| , \end{aligned}$$

where η_k is some point between $x^{(k)}$ and x^* .

Hence

$$\lim_{k \rightarrow \infty} \frac{e_{k+1}}{e_k} = \lim_{k \rightarrow \infty} |f'(\eta_k)| = |f'(x^*)| .$$

Case : $f'(x^*) = 0$:

$$\begin{aligned} e_{k+1} &= |x^{(k+1)} - x^*| = |f(x^{(k)}) - x^*| \\ &= |f(x^*) + (x^{(k)} - x^*)f'(x^*) + \frac{1}{2}(x^{(k)} - x^*)^2 f''(\eta_k) - x^*| \\ &= \frac{1}{2}e_k^2 |f''(\eta_k)|. \end{aligned}$$

where η_k is some point between $x^{(k)}$ and x^* ,

Thus

$$\lim_{k \rightarrow \infty} \frac{e_{k+1}}{e_k^2} = \lim_{k \rightarrow \infty} \frac{1}{2} |f''(\eta_k)| = \frac{1}{2} |f''(x^*)| \quad \text{QED!}$$

COROLLARY:

If

- $g(x)$ has three continuous derivatives near a zero x^* of $g(x) = 0$,
- $g'(x^*) \neq 0$,
- $x^{(0)}$ is sufficiently close to x^* ,

then Newton's method for solving $g(x) = 0$ *converges quadratically* .

PROOF: In Newton's method

$$f(x) = x - \frac{g(x)}{g'(x)} ,$$

and we have already shown that

$$f'(x^*) = 0 .$$

EXAMPLE:

Newton's method for computing $\sqrt{2}$, *i.e.*, for computing a zero of

$$g(x) = x^2 - 2,$$

is given by

$$x^{(k+1)} = x^{(k)} - \frac{(x^{(k)})^2 - 2}{2x^{(k)}},$$

that is,

$$x^{(k+1)} = \frac{(x^{(k)})^2 + 2}{2x^{(k)}},$$

that is,

$$x^{(k+1)} = f(x^{(k)}),$$

where

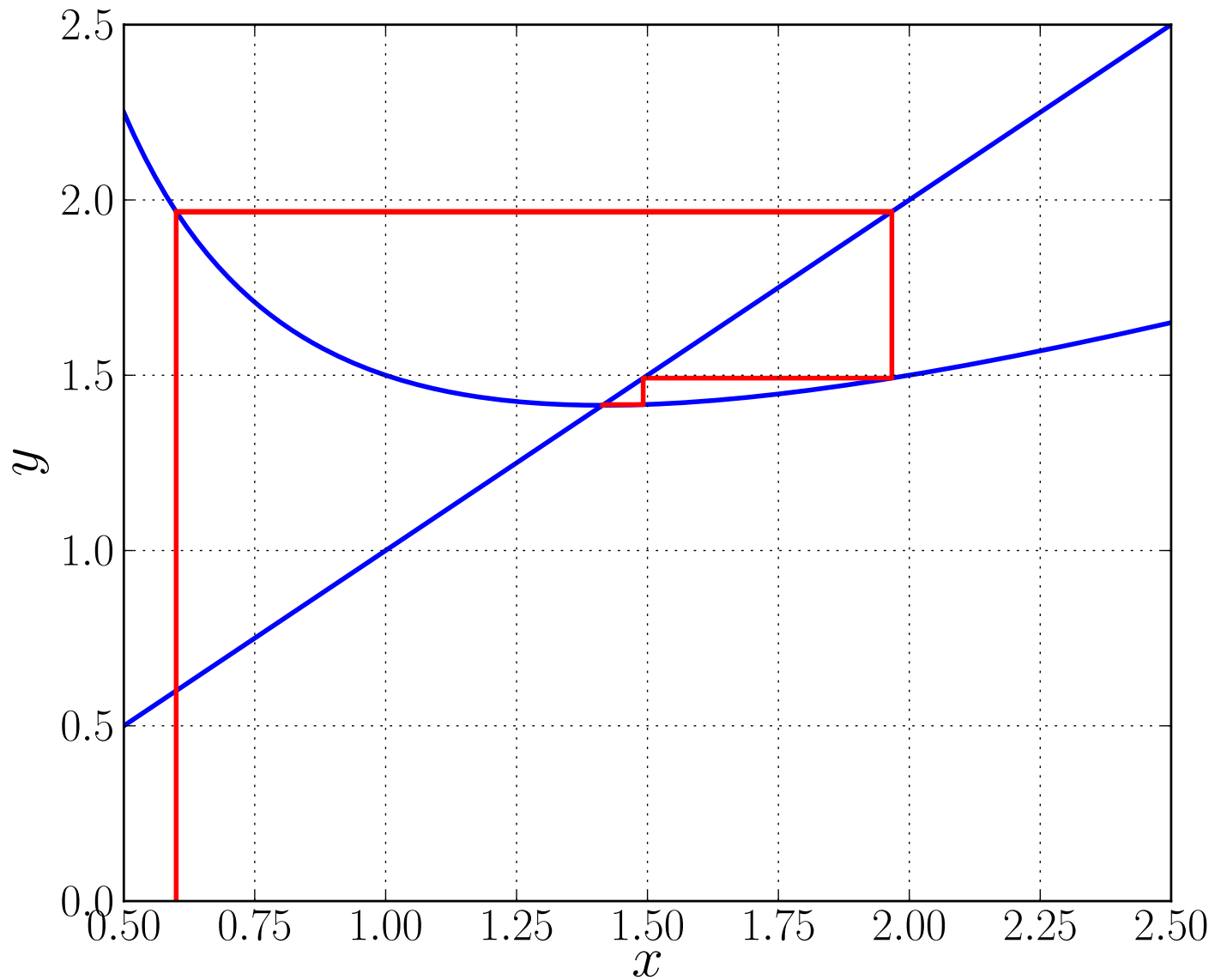
$$f(x) = \frac{x^2 + 2}{2x}.$$

$$x^{(k+1)} = f(x^{(k)}) , \quad \text{where} \quad f(x) = \frac{x^2 + 2}{2x} .$$

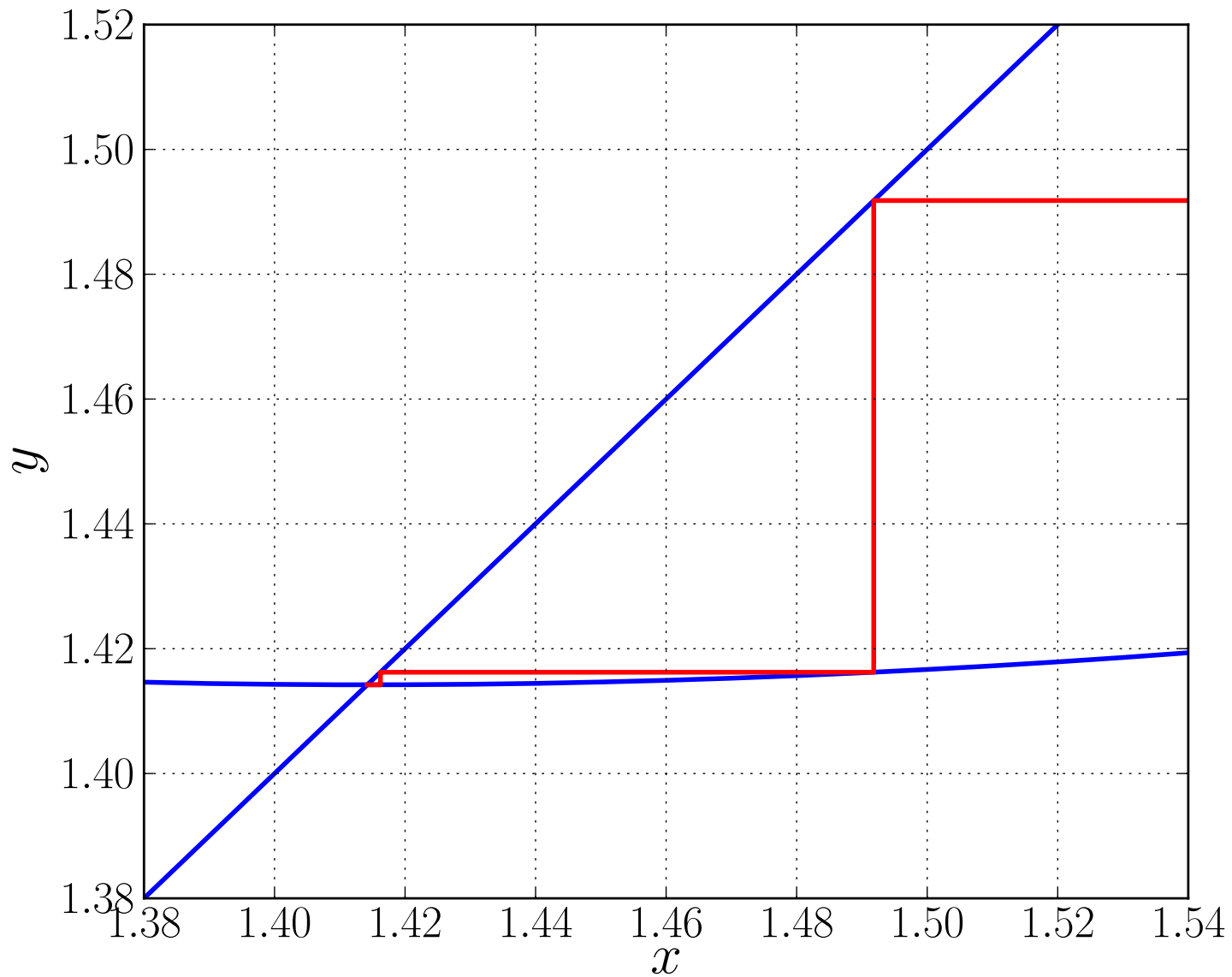
We observe that :

- The fixed points of f are $x^* = +\sqrt{2}$ and $x^* = -\sqrt{2}$.
- $f'(x^*) = 0$. (Hence quadratic convergence).
- $f(x) \rightarrow \infty$ as $x \downarrow 0$, and $f(x) \rightarrow -\infty$ as $x \uparrow 0$. (Vertical asymptotes.)
- $f(x) \approx x/2$ as $|x| \rightarrow \infty$.

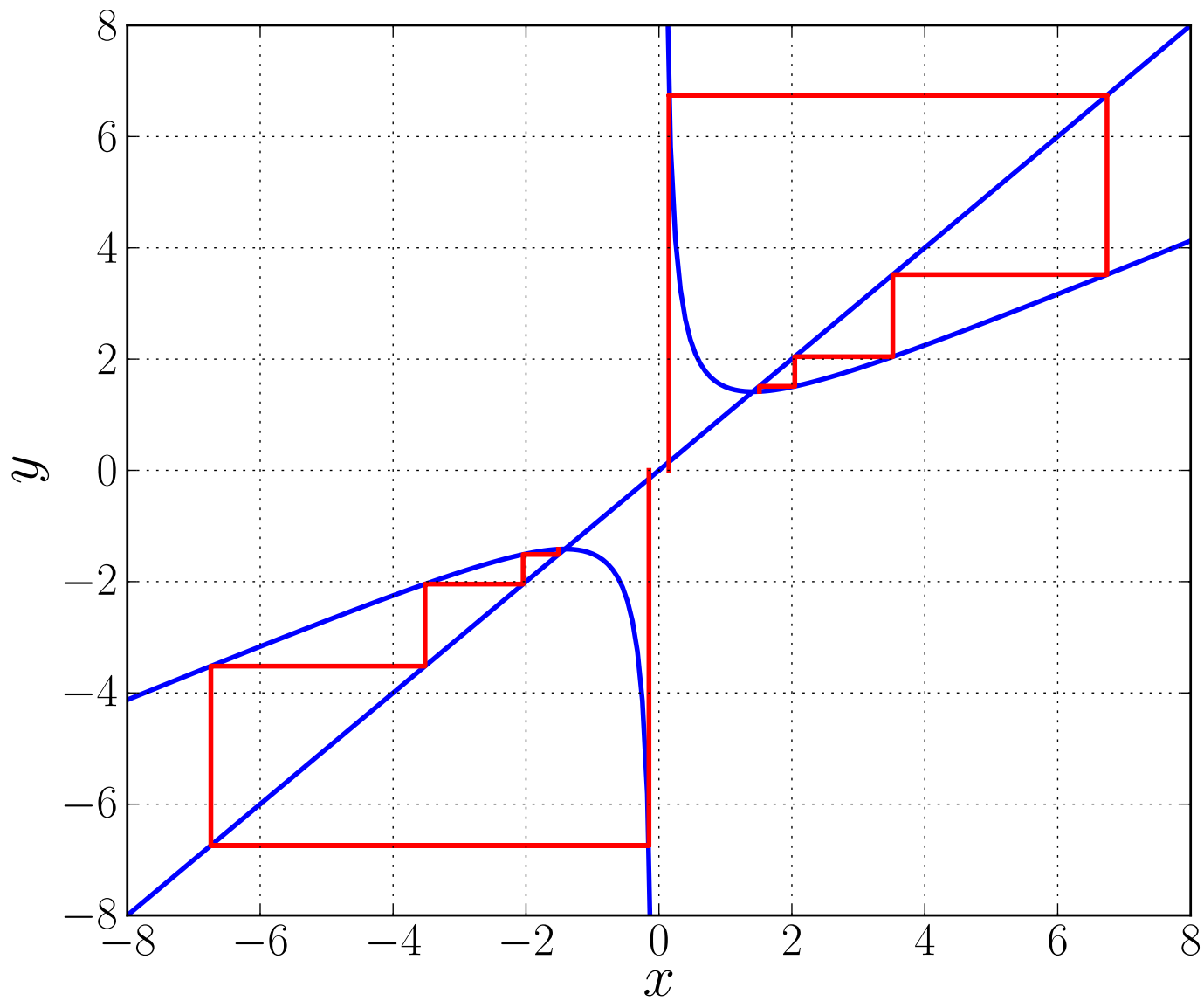
(Check all details!)



Newton's Method for $\sqrt{2}$ as a fixed point iteration



Newton's Method for $\sqrt{2}$ as a fixed point iteration (blow-up)



Newton's Method for $\pm\sqrt{2}$ as a fixed point iteration

From the last graph we see that :

- The iteration converges to $x^* = +\sqrt{2}$ for any $x^{(0)} > 0$.
- The iteration converges to $x^* = -\sqrt{2}$ for any $x^{(0)} < 0$.

(Check this!)

EXAMPLE:

Consider the fixed point iteration

$$x^{(k+1)} = x^{(k)} - \gamma g(x^{(k)}) ,$$

for computing a zero of $g(x) = 0$.

Indeed, a fixed point x^* satisfies

$$x^* = x^* - \gamma g(x^*) ,$$

that is,

$$g(x^*) = 0 . \quad (\text{Assuming } \gamma \neq 0 .)$$

In this example

$$f(x) = x - \gamma g(x) .$$

A fixed point x^* is attracting if

$$| f'(x^*) | < 1 ,$$

i.e., if

$$| 1 - \gamma g'(x^*) | < 1 ,$$

i.e., if

$$-1 < 1 - \gamma g'(x^*) < 1 ,$$

i.e., if

$$-2 < -\gamma g'(x^*) < 0 ,$$

i.e., if

$$0 < \gamma g'(x^*) < 2 .$$

The convergence is quadratic if

$$f'(x^*) = 1 - \gamma g'(x^*) = 0 ,$$

that is, if

$$\gamma = \hat{\gamma} \equiv \frac{1}{g'(x^*)} .$$

Now x^* is unknown beforehand and therefore $\hat{\gamma}$ is also unknown.

However, after the k th iteration an approximation to $\hat{\gamma}$ is given by

$$\hat{\gamma} \approx \gamma_k \equiv \frac{1}{g'(x^{(k)})} .$$

This leads to the iteration

$$x^{(k+1)} = x^{(k)} - \gamma_k g(x^{(k)}) ,$$

where $\gamma_k = 1/g'(x^{(k)})$,

i.e., we have rediscovered Newton's method !

EXERCISES:

- If the following fixed point iteration converges, then what number will it converge to? Is the convergence quadratic?

$$x^{(k+1)} = f(x^{(k)}) , \quad \text{where} \quad f(x) = \frac{2x^3 + 3}{3x^2} .$$

- Analytically determine all fixed points of

$$x^{(k+1)} = 2(x^{(k)})^2 - 2x^{(k)} + 1 , \quad k = 0, 1, 2, \dots .$$

Are the fixed points attracting or repelling? If attracting, then is the convergence linear or quadratic? Also draw a graphical interpretation.

- Analytically determine all fixed points of $x^{(k+1)} = 2x^{(k)}(1 - x^{(k)})$. Are these attracting or repelling? If attracting then is the convergence linear or quadratic? Also give a graphical interpretation.

EXERCISES:

- Give a graphical interpretation of the fixed point iteration.

$$x^{(k+1)} = \sin(x^{(k)}) .$$

What are the fixed points? Does the derivative test give conclusive evidence whether the fixed point $x = 0$ is attracting or repelling? Based on the graphical interpretation, can one conclude whether $x = 0$ is attracting or repelling?

- Consider the fixed point iteration $x^{(k+1)} = f(x^{(k)})$, where

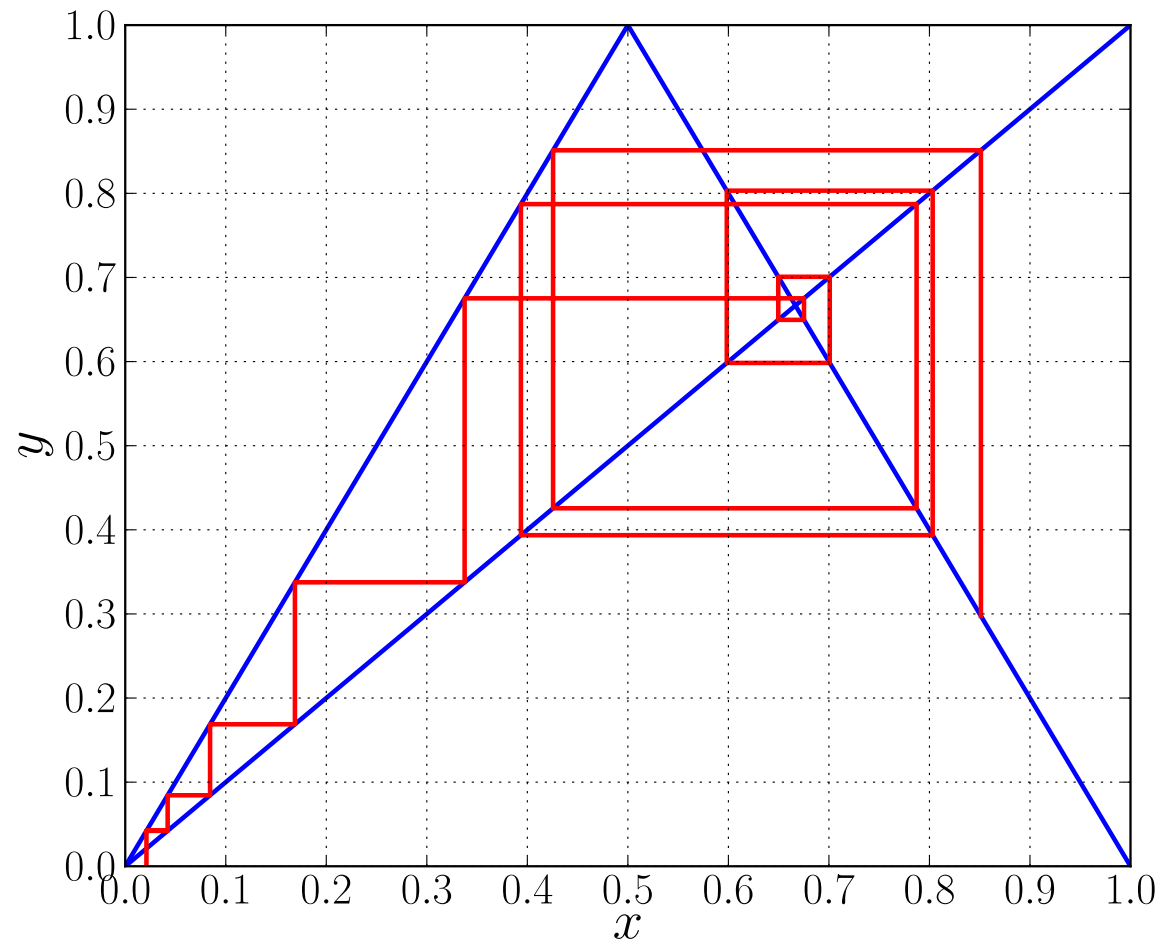
$$f(x) = \begin{cases} 2x , & \text{when } x \leq \frac{1}{2} , \\ 2(1 - x) , & \text{when } x > \frac{1}{2} . \end{cases}$$

Give an accurate graphical interpretation in the interval $[0, 1]$, with $x^{(0)} \approx 0.1$, showing enough iterations to illustrate the behavior of this fixed point iteration. Analytically determine all fixed points, and for each fixed point determine whether it is attracting or repelling.

$$f(x) = \begin{cases} 2x, & \text{when } x \leq \frac{1}{2}, \\ 2(1-x), & \text{when } x > \frac{1}{2}. \end{cases}$$

SOLUTION:

The fixed points are $x = 0$ and $x = \frac{2}{3}$. Both are repelling since $|f'(x)| = 2$.



The first 15 iterations, with $x^{(0)} = 0.0211$.

EXERCISES:

- Show how to use the Chord method to compute the cube root of 5.

Carry out the first two iterations of the Chord method, using $x^{(0)} = 2$.

Analytically determine all fixed points of this Chord iteration.

For each fixed point, determine whether it is attracting or repelling.

- Draw the graph of $g(x) = x^3 - 2$, clearly showing its zero.

Write down Newton's method for finding a zero of g , and simplify the expression for the Newton iteration as much as possible.

Will Newton's method converge if the initial guess is sufficiently close?

If yes, then what is the rate of convergence?

Will Newton's method converge for *any positive* initial guess ?

Will Newton's method converge for *negative* initial guesses ?

SOLUTION: Newton's method for the *cube root* of 2 has the form

$$x^{(k+1)} = f(x^{(k)}), \quad \text{where} \quad f(x) = x - \frac{g(x)}{g'(x)} = x - \frac{x^3 - 2}{3x^2} = \frac{2(x^3 + 1)}{3x^2}.$$

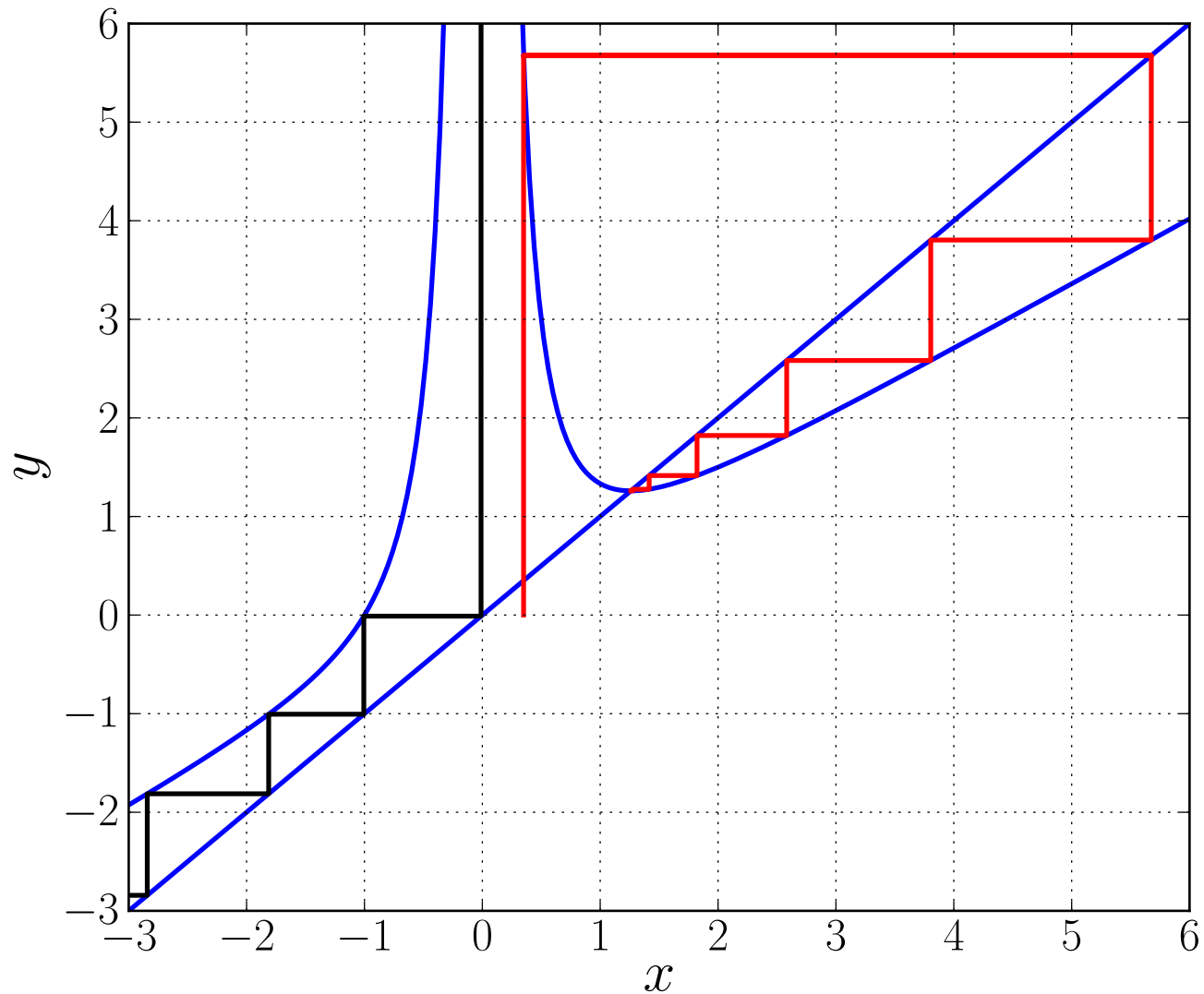
Analytically we find that $f(x)$:

- has a vertical asymptote at $x = 0$
- approaches the line $y = \frac{2}{3}x$ when $|x| \rightarrow \infty$
- has (of course!) a fixed point at $x^* = \sqrt[3]{2}$
- $|f'(x^*)| = 0$, so convergence is quadratic (once $x^{(k)}$ is close to x^*)

Graphically we see that :

- the iteration converges for any $x^{(0)} > 0$
- the iteration converges for “most” negative $x^{(0)}$, except for a countably infinite such $x^{(0)}$, namely those for which $x^{(k)} = 0$ for some k .

SOLUTION: continued ...



Newton's method for the cube root of 2, with a converging iteration (red) having $x^{(0)} = 0.35$, and the iteration (black) with $x^{(k)} = 0$ for some k .

EXERCISES:

- Suppose you enter any number on a calculator and then keep pushing the *cosine* button. (Assume the calculator is in “radian-mode”.)

What will happen in the limit to the result shown in the display?

Give a full mathematical explanation and a graphical interpretation.

Do the same for $\sin(x)$ and $\tan(x)$.

- Consider the fixed point iteration

$$x^{(k+1)} = x^{(k)} (1 - x^{(k)}) .$$

Does the derivative test give conclusive evidence whether or not the fixed point $x = 0$ is attracting?

Give a careful graphical interpretation of this fixed point iteration.

What can you say about the convergence of the fixed point iteration?

EXERCISES:

- Consider the fixed point iteration

$$x^{(k+1)} = \frac{1}{\sqrt{x^{(k)}}}, \quad k = 0, 1, 2, \dots .$$

Give a careful graphical interpretation of this fixed point iteration.

Determine all fixed points and whether they are attracting or repelling.

Does the iteration converge for all positive initial points $x^{(0)}$?

- Consider the fixed point iteration

$$x^{(k+1)} = f(x^{(k)}), \quad \text{where} \quad f(x) = \frac{1 + x^2}{1 + x} .$$

Determine all fixed points and whether they are attracting or repelling.

If attracting determine whether the convergence is linear or quadratic.

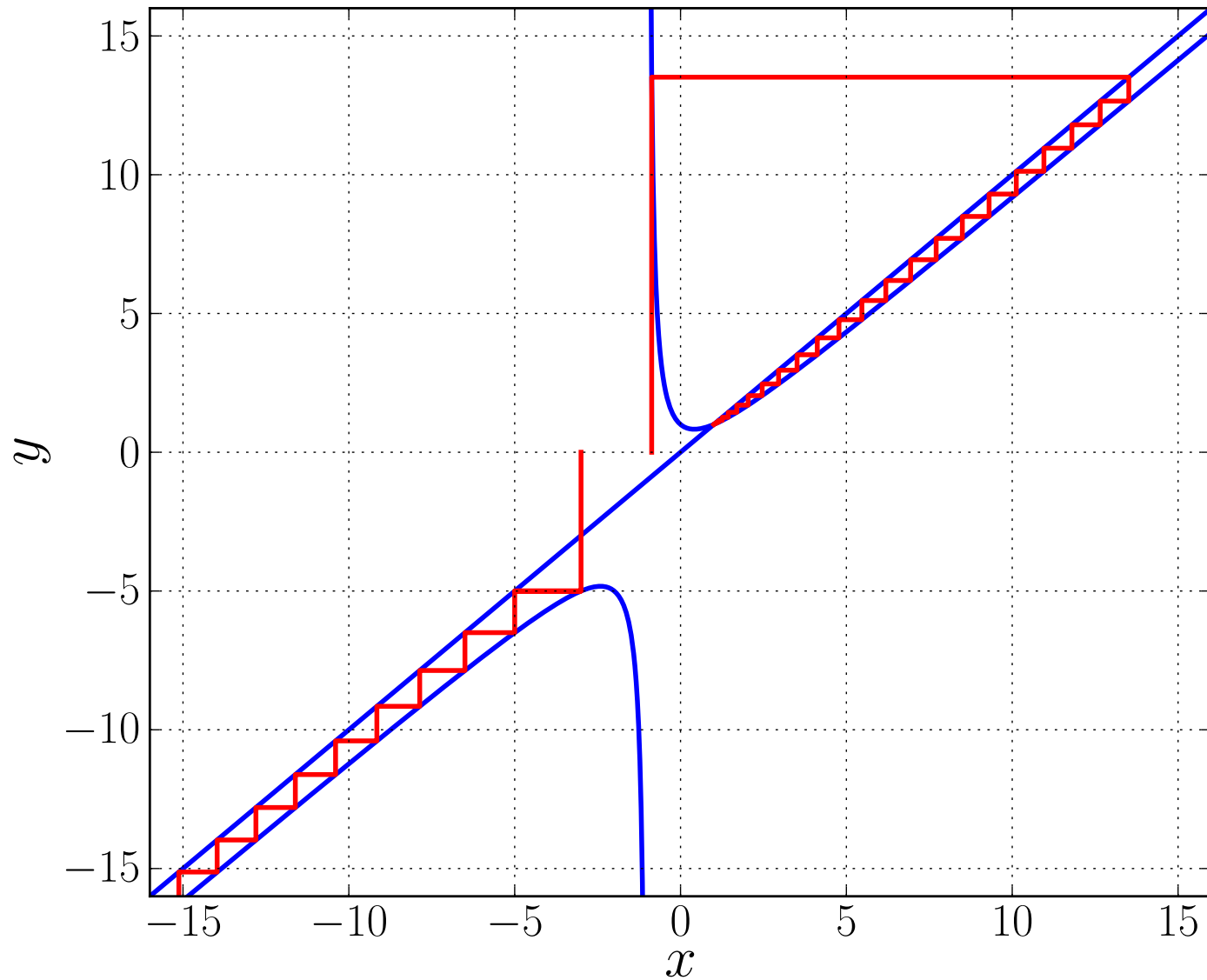
Give a graphical interpretation of the first few iterations, with $x^{(0)} = 2$.

SOLUTION: $x^{(k+1)} = f(x^{(k)})$, where $f(x) = \frac{1+x^2}{1+x}$.

Analytically we find that $f(x)$:

- has a vertical asymptote at $x = -1$
- approaches the line $y = x$ when $|x| \rightarrow \infty$
- has a fixed point at $x^* = 1$
- $|f'(x^*)| = \frac{1}{2} < 1$, so x^* is attracting
- graphically we see the iteration converges for any $x^{(0)} > -1$
- graphically we also see the iteration diverges for any $x^{(0)} < -1$
- the divergence for $x^{(0)} < -1$ is slow as $x^{(k)}$ becomes more negative

SOLUTION: continued ...



The fixed point iteration $x^{(k+1)} = f(x^{(k)})$, where $f(x) = \frac{1+x^2}{1+x}$.

Convergence Analysis for Systems.

Again most iterative methods for solving

$$\mathbf{G}(\mathbf{x}) = \mathbf{0} ,$$

can be written as

$$\mathbf{x}^{(k+1)} = \mathbf{F}(\mathbf{x}^{(k)}) , \quad k =, 1, 2, \dots ,$$

where the function \mathbf{F} should be chosen such that

\mathbf{x}^* is a root of $\mathbf{G}(\mathbf{x}) = \mathbf{0}$ if and only if \mathbf{x}^* is a fixed point of \mathbf{F} .

EXAMPLE:

Newton's method for systems is

$$\mathbf{G}'(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = -\mathbf{G}(\mathbf{x}^{(k)}) .$$

Thus

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{G}'(\mathbf{x}^{(k)})^{-1} \mathbf{G}(\mathbf{x}^{(k)}) ,$$

assuming $\mathbf{G}'(\mathbf{x})^{-1}$ to exist near $\mathbf{x} = \mathbf{x}^*$.

So here

$$\mathbf{F}(\mathbf{x}) = \mathbf{x} - \mathbf{G}'(\mathbf{x})^{-1} \mathbf{G}(\mathbf{x}) .$$

Fixed point iterations also arise as models of physical processes, where they are often called *difference equations* or *discrete dynamical systems*.

EXAMPLE :

The equations

$$x_1^{(k+1)} = \lambda x_1^{(k)}(1 - x_1^{(k)}) - c_1 x_1^{(k)} x_2^{(k)},$$

$$x_2^{(k+1)} = c_2 x_2^{(k)} + c_1 x_1^{(k)} x_2^{(k)},$$

model a “*predator-prey*” system, where, for example,

$x_1^{(k)}$ denotes the biomass of “fish” in year k ,

and

$x_2^{(k)}$ denotes the biomass of “sharks” in year k ,

and where λ , c_1 , and c_2 are constants.

Derivatives :

For scalar functions

$$f'(x) \equiv \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h},$$

or, equivalently, $f'(x)$ is the number such that

$$\frac{f(x+h) - f(x) - f'(x)h}{h} \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

If $f'(x)$ exists then f is said to be *differentiable* at x .

Similarly for vector valued functions $\mathbf{F}(\mathbf{x})$ we say that \mathbf{F} is *differentiable* at \mathbf{x} if there exists a *matrix* $\mathbf{F}'(\mathbf{x})$ such that

$$\frac{\|\mathbf{F}(\mathbf{x} + \mathbf{h}) - \mathbf{F}(\mathbf{x}) - \mathbf{F}'(\mathbf{x})\mathbf{h}\|}{\|\mathbf{h}\|} \rightarrow 0 \quad \text{as} \quad \|\mathbf{h}\| \rightarrow 0 .$$

The matrix $\mathbf{F}'(\mathbf{x})$ is the *Jacobian matrix* introduced earlier :

If $\mathbf{F}(\mathbf{x})$ has component functions

$$(f^1(\mathbf{x}), f^2(\mathbf{x}), \dots, f^n(\mathbf{x}))^T ,$$

and if

$$\mathbf{x} \equiv (x_1, x_2, \dots, x_n)^T ,$$

then

$$\{\mathbf{F}'(\mathbf{x})\}_{i,j} \equiv \frac{\partial f^i}{\partial x_j} .$$

THEOREM:

Let $\mathbf{F}'(\mathbf{x})$ be continuous near a fixed point \mathbf{x}^* of $\mathbf{F}(\mathbf{x})$ and

$$\| \mathbf{F}'(\mathbf{x}^*) \| < 1 ,$$

in some induced matrix norm.

Then the fixed point iteration

$$\mathbf{x}^{(k+1)} = \mathbf{F}(\mathbf{x}^{(k)}) , \quad k = 0, 1, 2, \dots ,$$

converges to \mathbf{x}^* whenever the initial guess $\mathbf{x}^{(0)}$ is sufficiently close to \mathbf{x}^* .

REMARK :

It can be shown that sufficient condition for

$$\| \mathbf{F}'(\mathbf{x}^*) \| < 1 ,$$

in some matrix norm, is that

$$\text{spr}(\mathbf{F}'(\mathbf{x}^*)) < 1 .$$

Here

$\text{spr}(\mathbf{F}'(\mathbf{x}^*))$ is the *spectral radius* of $\mathbf{F}'(\mathbf{x}^*)$,

that is, the size (complex absolute value) of the largest eigenvalue of $\mathbf{F}'(\mathbf{x}^*)$.

(Note that eigenvalues may be complex numbers.)

PROOF (of the Theorem):

(Similar to the proof of the scalar case.)

Let $\alpha \equiv \| \mathbf{F}'(\mathbf{x}^*) \|$. Then $\alpha < 1$.

By definition of \mathbf{F}' , given any $\epsilon > 0$, in particular

$$\epsilon \equiv \frac{1 - \alpha}{2},$$

there exists a $\delta > 0$ such that

$$\| \mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}^*) - \mathbf{F}'(\mathbf{x}^*) (\mathbf{x} - \mathbf{x}^*) \| \leq \epsilon \| \mathbf{x} - \mathbf{x}^* \| ,$$

whenever $\mathbf{x} \in \mathcal{B}_\delta(\mathbf{x}^*)$.

Here

$$\mathcal{B}_\delta(\mathbf{x}^*) \equiv \{ \mathbf{x} : \| \mathbf{x} - \mathbf{x}^* \| \leq \delta \} .$$

Let $\mathbf{x}^{(0)} \in \mathcal{B}_\delta(\mathbf{x}^*)$. Then

$$\begin{aligned} \|\mathbf{x}^{(1)} - \mathbf{x}^*\| &= \|\mathbf{F}(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^*)\| \\ &= \|\mathbf{F}(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^*) - \mathbf{F}'(\mathbf{x}^*)(\mathbf{x}^{(0)} - \mathbf{x}^*) + \mathbf{F}'(\mathbf{x}^*)(\mathbf{x}^{(0)} - \mathbf{x}^*)\| \\ &\leq \|\mathbf{F}(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^*) - \mathbf{F}'(\mathbf{x}^*)(\mathbf{x}^{(0)} - \mathbf{x}^*)\| + \|\mathbf{F}'(\mathbf{x}^*)(\mathbf{x}^{(0)} - \mathbf{x}^*)\| \\ &\leq \epsilon \|\mathbf{x}^{(0)} - \mathbf{x}^*\| + \alpha \|\mathbf{x}^{(0)} - \mathbf{x}^*\| \\ &\leq (\epsilon + \alpha) \|\mathbf{x}^{(0)} - \mathbf{x}^*\| \\ &= \left(\frac{1 - \alpha}{2} + \alpha\right) \|\mathbf{x}^{(0)} - \mathbf{x}^*\| \\ &= \frac{1 + \alpha}{2} \|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \beta \delta, \end{aligned}$$

where $\beta \equiv \frac{1 + \alpha}{2} < 1$.

Thus $\mathbf{x}^{(1)} \in \mathcal{B}_\delta(\mathbf{x}^*)$.

Since $\mathbf{x}^{(1)} \in \mathcal{B}_\delta(\mathbf{x}^*)$, we also have

$$\|\mathbf{x}^{(2)} - \mathbf{x}^*\| \leq \beta \|\mathbf{x}^{(1)} - \mathbf{x}^*\| \leq \beta^2 \delta .$$

Continuing in this fashion, we find

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \beta^k \delta .$$

Thus, since $\beta < 1$, we see that $\mathbf{x}^{(k)}$ converges to \mathbf{x}^* . QED!

EXAMPLE: In Newton's method

$$\mathbf{F}(\mathbf{x}) = \mathbf{x} - (\mathbf{G}'(\mathbf{x}))^{-1} \mathbf{G}(\mathbf{x}) .$$

Hence

$$\mathbf{G}'(\mathbf{x})\mathbf{F}(\mathbf{x}) = \mathbf{G}'(\mathbf{x})\mathbf{x} - \mathbf{G}(\mathbf{x}) .$$

$$\Rightarrow \mathbf{G}''(\mathbf{x})\mathbf{F}(\mathbf{x}) + \mathbf{G}'(\mathbf{x})\mathbf{F}'(\mathbf{x}) = \mathbf{G}''(\mathbf{x})\mathbf{x} + \mathbf{G}'(\mathbf{x}) - \mathbf{G}'(\mathbf{x}) = \mathbf{G}''(\mathbf{x})\mathbf{x}$$

$$\Rightarrow \mathbf{G}'(\mathbf{x})\mathbf{F}'(\mathbf{x}) = \mathbf{G}''(\mathbf{x})(\mathbf{x} - \mathbf{F}(\mathbf{x})) = \mathbf{G}''(\mathbf{x})(\mathbf{G}'(\mathbf{x}))^{-1}\mathbf{G}(\mathbf{x})$$

$$\Rightarrow \mathbf{F}'(\mathbf{x}) = (\mathbf{G}'(\mathbf{x}))^{-1}\mathbf{G}''(\mathbf{x})(\mathbf{G}'(\mathbf{x}))^{-1}\mathbf{G}(\mathbf{x})$$

$$\Rightarrow \mathbf{F}'(\mathbf{x}^*) = (\mathbf{G}'(\mathbf{x}^*))^{-1}\mathbf{G}''(\mathbf{x}^*)(\mathbf{G}'(\mathbf{x}^*))^{-1}\mathbf{G}(\mathbf{x}^*) = \mathbf{O} \text{ (zero matrix) .}$$

because $\mathbf{G}(\mathbf{x}^*) = \mathbf{0}$.

So $\|\mathbf{F}'(\mathbf{x}^*)\| = 0$, and therefore certainly $\|\mathbf{F}'(\mathbf{x}^*)\| < 1$.

Thus if

- $\mathbf{G}''(\mathbf{x})$ is continuous near \mathbf{x}^* ,
- $(\mathbf{G}'(\mathbf{x}^*))^{-1}$ exists ,
- $\mathbf{x}^{(0)}$ is sufficiently close to \mathbf{x}^* ,

then Newton's method converges.

Again this convergence can be shown to be *quadratic* , *i.e.*,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2} \leq C , \quad \text{for some constant } C.$$

EXERCISE:

- Consider the fixed point iteration

$$x_1^{(k+1)} = \lambda x_1^{(k)} (1 - x_1^{(k)}) - 0.2 x_1^{(k)} x_2^{(k)} ,$$

$$x_2^{(k+1)} = 0.9 x_2^{(k)} + 0.2 x_1^{(k)} x_2^{(k)} .$$

This is a “*predator-prey*” model, where, for example, $x_1^{(k)}$ denotes the biomass of “fish” and $x_2^{(k)}$ denotes the biomass of “sharks” in year k .

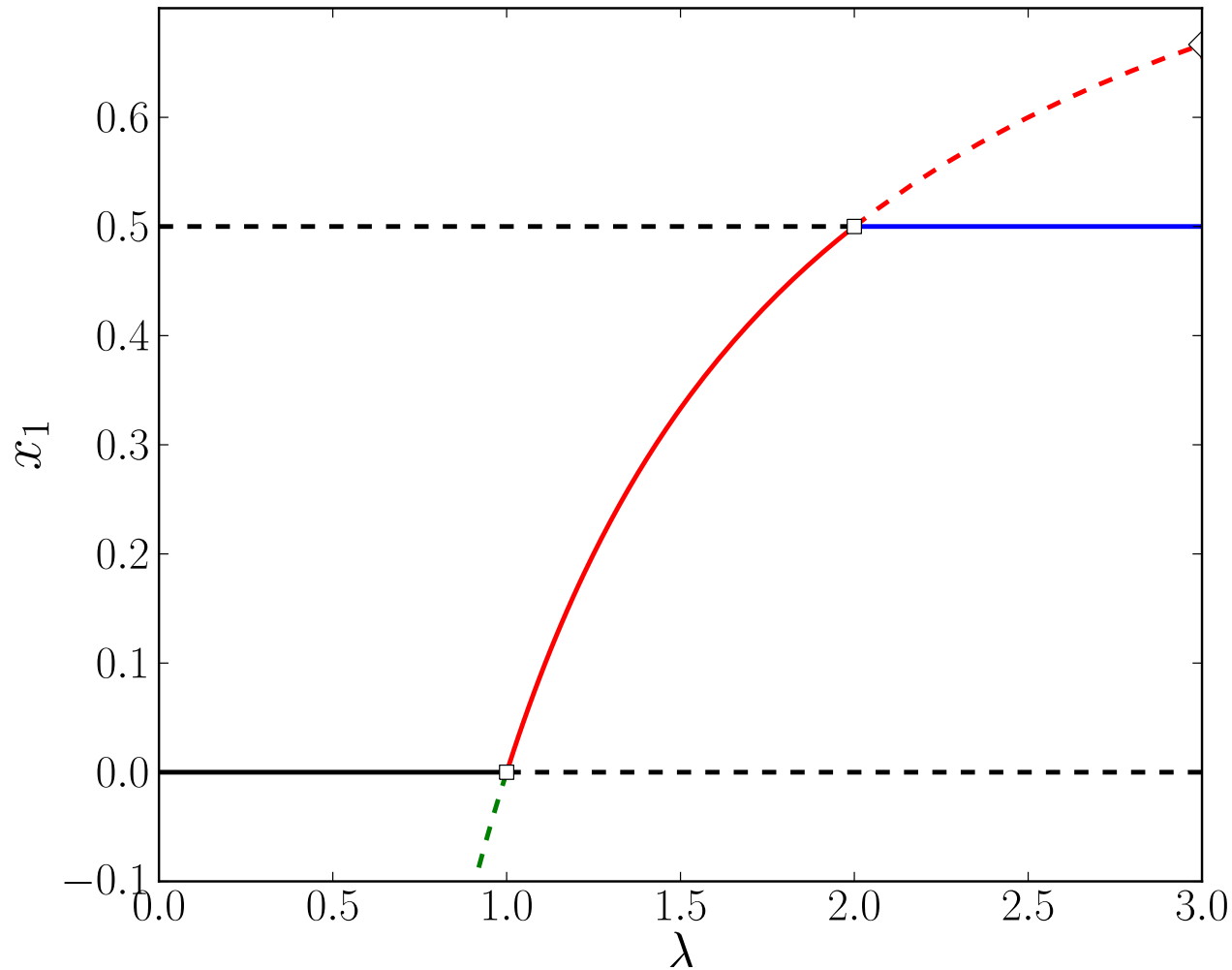
Numerically determine the long-time behavior of $x_1^{(k)}$ and $x_2^{(k)}$ for the following values of λ :

$$\lambda = 0.5, 1.0, 1.5, 2.0, 2.5 ,$$

taking, for example, $x_1^{(0)} = 0.1$ and $x_2^{(0)} = 0.1$.

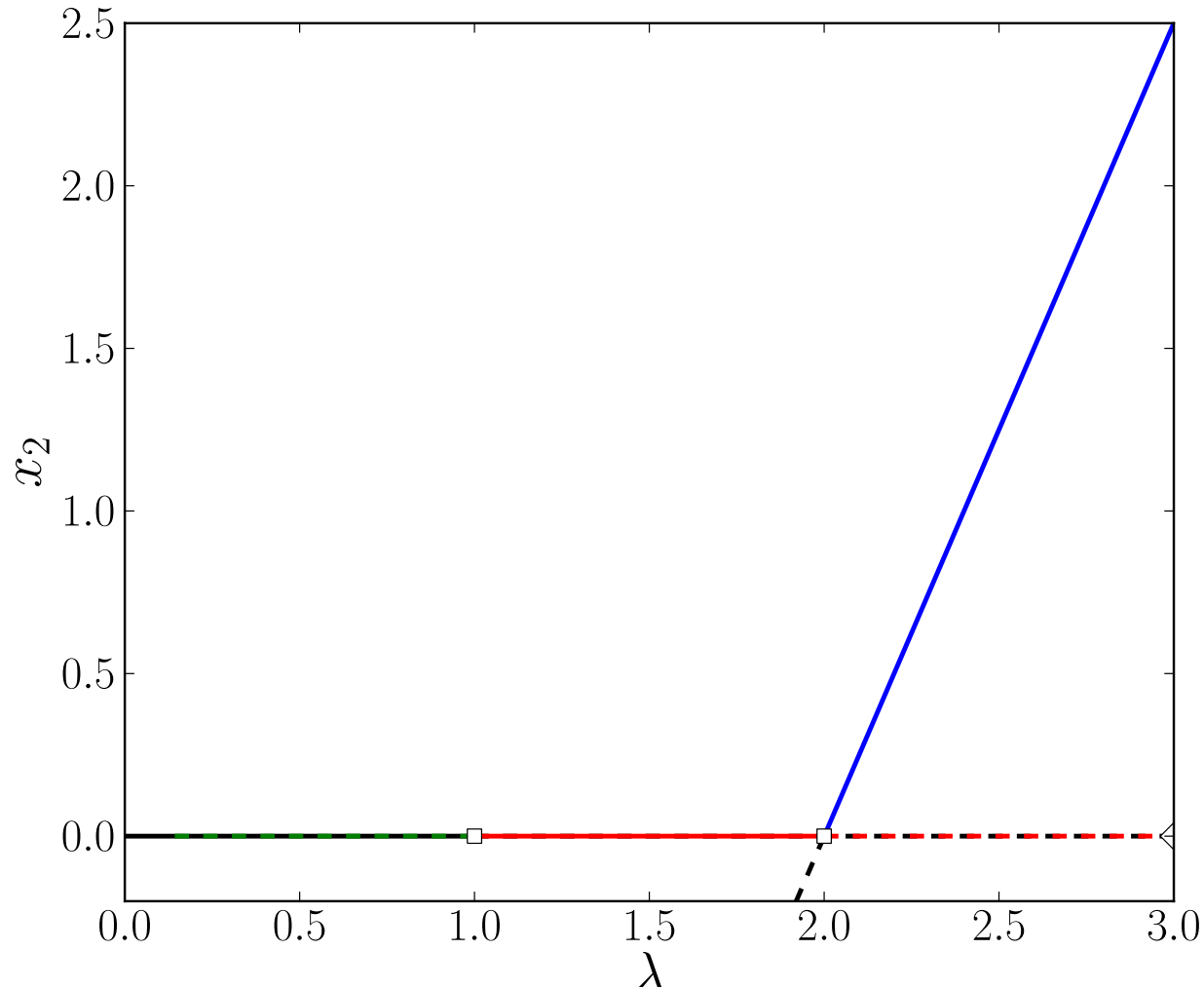
What can you say *analytically* about the fixed points of this system?

SOLUTION:



The value of x_1 along families of fixed points of the predator-prey system. Solid/dashed curves represent stable/unstable fixed points, respectively. Black: $x_1 = x_2 = 0$. Red: $x_1 \neq 0, x_2 = 0$. Blue: $x_1 = 0.5, x_2 \neq 0$.

SOLUTION: (continued ...)



The value of x_2 along families of fixed points of the predator-prey system. Solid/dashed curves represent stable/unstable fixed points, respectively. Black: $x_1 = x_2 = 0$. Red: $x_1 \neq 0, x_2 = 0$. Blue: $x_1 = 0.5, x_2 \neq 0$.

THE APPROXIMATION OF FUNCTIONS.

Function Norms.

To measure how well a given function $f \in \mathbb{C}[a, b]$ is approximated by another function we need a quantity called *function norm*.

Examples of these are :

$$\begin{aligned}\| f \|_1 &\equiv \int_a^b | f(x) | dx , \\ \| f \|_2 &\equiv \left\{ \int_a^b f(x)^2 dx \right\}^{\frac{1}{2}} , \\ \| f \|_\infty &\equiv \max_{[a,b]} | f(x) | .\end{aligned}$$

Note the similarity of these norms to the corresponding vector norms.

A function norm is required to satisfy :

$$(i) \quad \| f \| \geq 0, \quad \forall f \in \mathbb{C}[a, b], \quad \| f \| = 0 \text{ iff } f \equiv 0,$$

$$(ii) \quad \| \alpha f \| = |\alpha| \| f \|, \quad \forall \alpha \in \mathbb{R}, \quad \forall f \in \mathbb{C}[a, b],$$

$$(iii) \quad \| f + g \| \leq \| f \| + \| g \|, \quad \forall f, g \in \mathbb{C}[a, b].$$

All of the norms above satisfy these requirements. (Check!)

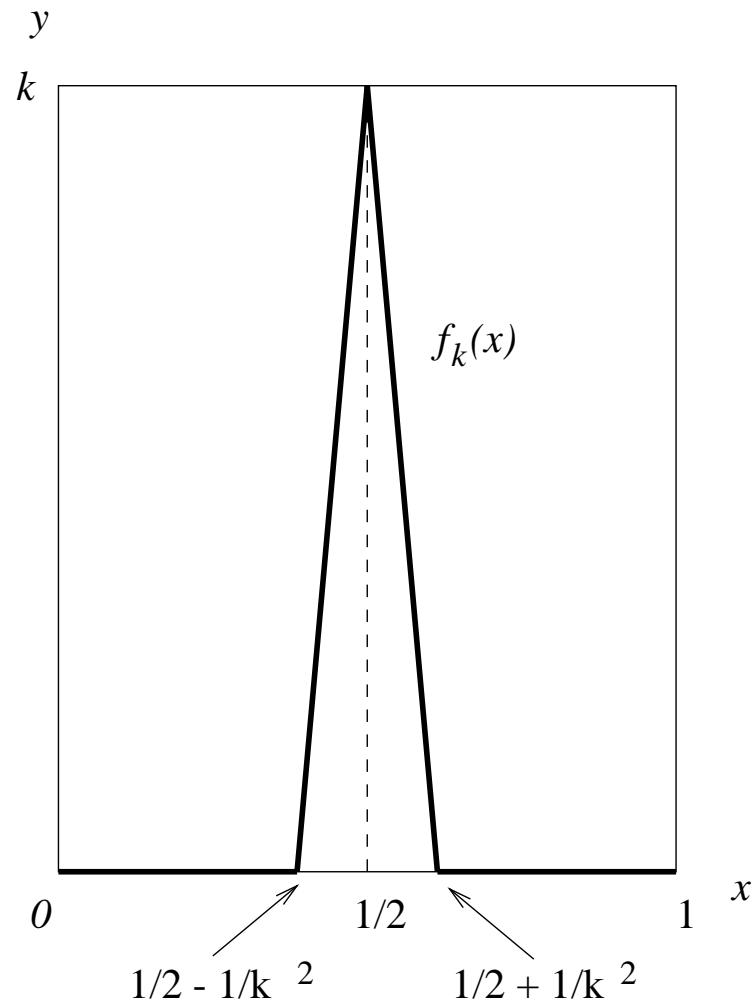
For example

$$\begin{aligned}\|f + g\|_1 &= \int_a^b |f(x) + g(x)| \, dx \\ &\leq \int_a^b |f(x)| + |g(x)| \, dx \\ &= \int_a^b |f(x)| \, dx + \int_a^b |g(x)| \, dx \\ &= \|f\|_1 + \|g\|_1 .\end{aligned}$$

(For the $\|\cdot\|_2$ we shall verify (iii) later.)

REMARK: If a function is “small” in a given function norm then it need not be small in another norm.

For example, consider $f_k(x)$, $k = 2, 3, \dots$, as shown below :



Then

$$\| f_k \|_\infty = k \rightarrow \infty \quad \text{as } k \rightarrow \infty ,$$

while

$$\| f_k \|_1 = \int_0^1 | f_k(x) | \, dx = \frac{1}{k} \rightarrow 0 \quad \text{as } k \rightarrow \infty ,$$

and

$$\| f_k \|_2 = \left\{ \int_0^1 f_k(x)^2 \, dx \right\}^{\frac{1}{2}} = \sqrt{2/3} \quad (\text{Check!}) .$$

EXAMPLE:

Approximate

$$f(x) = x^3$$

by

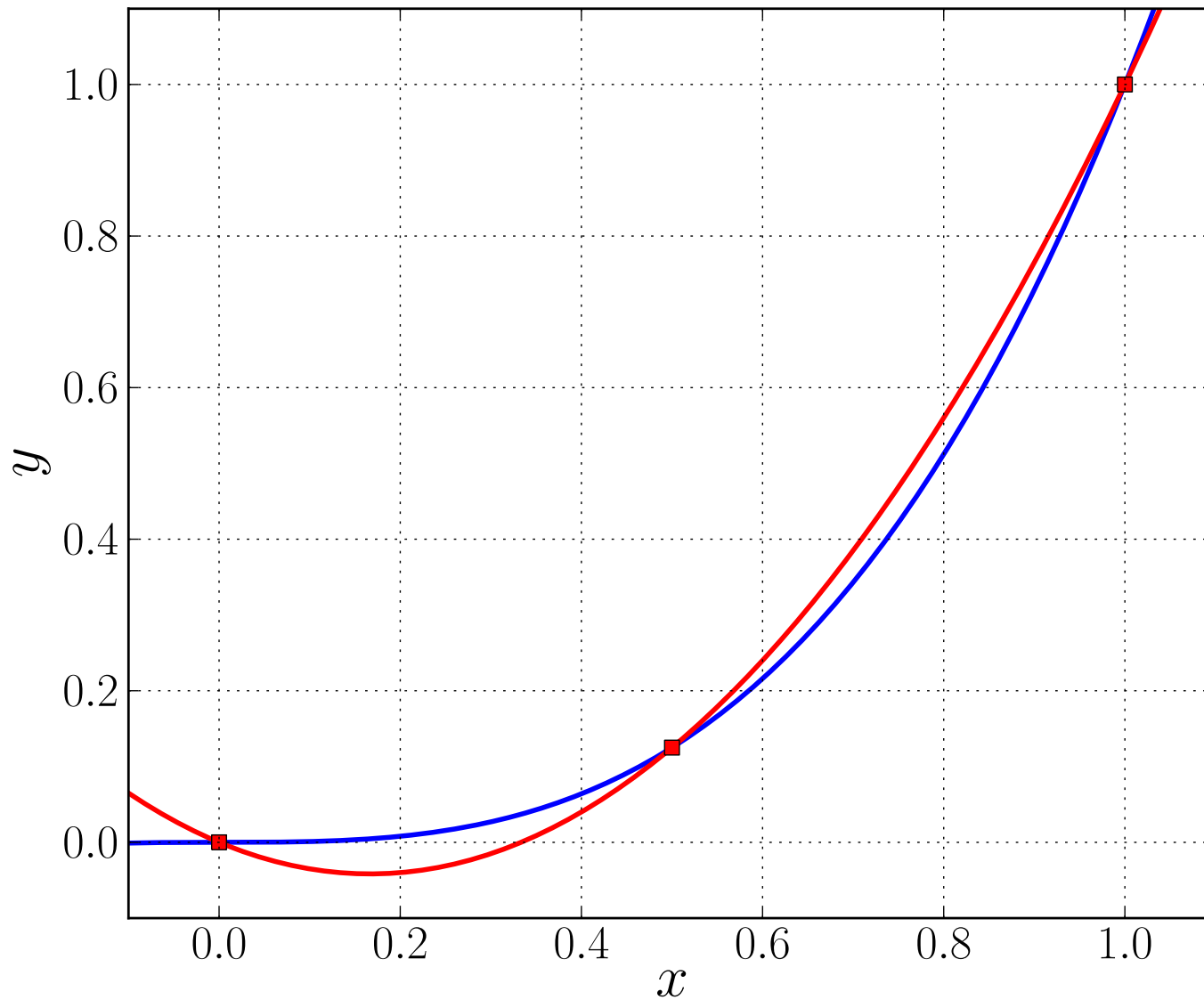
$$p(x) = \frac{3}{2}x^2 - \frac{1}{2}x,$$

on the interval $[0, 1]$.

Then

$$p(x) = f(x) \quad \text{for } x = 0, \frac{1}{2}, 1,$$

that is, $p(x)$ *interpolates* $f(x)$ at these points.



Graph of $f(x) = x^3$ (blue) and its interpolant $p(x) = \frac{3}{2}x^2 - \frac{1}{2}x$ (red) .

A measure of “how close” f and p are is then given by, for example,

$$\| f - p \|_2 = \left\{ \int_0^1 (f(x) - p(x))^2 dx \right\}^{\frac{1}{2}} .$$

We find that

$$\| f - p \|_2 = \frac{\sqrt{210}}{420} \approx 0.0345. \quad (\text{Check!})$$

The Lagrange Interpolation Polynomial.

Let f be a function defined on $[a, b]$.

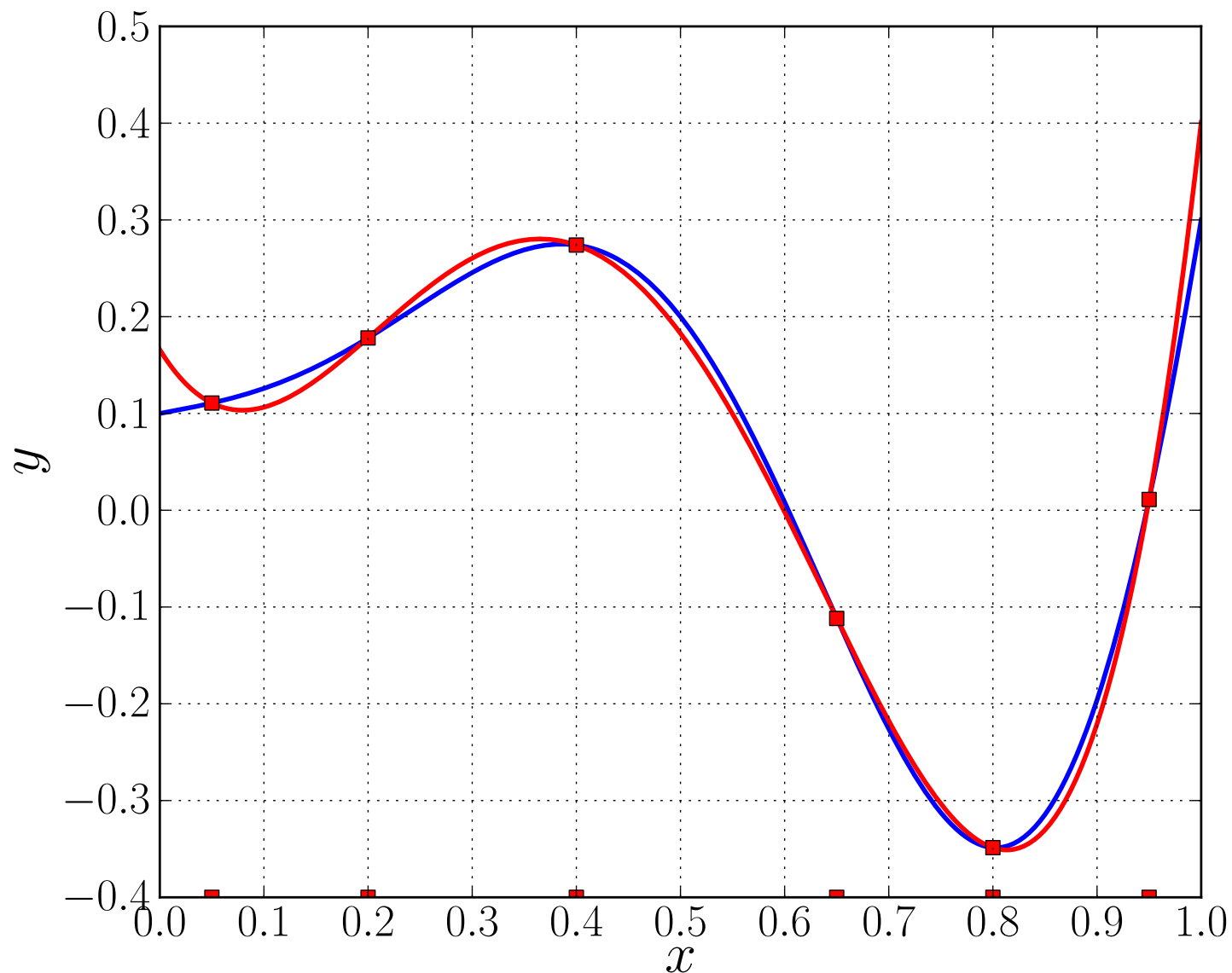
Let \mathbb{P}_n denote all polynomials of degree less than or equal to n .

Given points $\{x_k\}_{k=0}^n$ with

$$a \leq x_0 < x_1 < \cdots < x_n \leq b ,$$

we want to find $p \in \mathbb{P}_n$ such that

$$p(x_k) = f(x_k) , \quad k = 0, 1, \dots, n .$$



Graph of $f(x) = \frac{1}{10} + \frac{1}{5}x + x^2 \sin(2\pi x)$ (blue)
 and its Lagrange interpolant $p(x) \in \mathbb{P}_5$ (red)
 at six interpolation points ($n = 5$).

The following *questions* arise :

(i) Is $p(x)$ uniquely defined ?

(ii) How well does p approximate f ?

(iii) Does the approximation get better as $n \rightarrow \infty$?

To answer the above questions let

$$l_i(x) \equiv \prod_{k=0, k \neq i}^n \frac{(x - x_k)}{(x_i - x_k)}, \quad i = 0, 1, \dots, n,$$

be the *Lagrange interpolating coefficients*, or *Lagrange basis functions*.

Then each $l_i \in \mathbb{P}_n$. (Check!)

For example if $n = 2$ we have

$$\ell_0(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)},$$

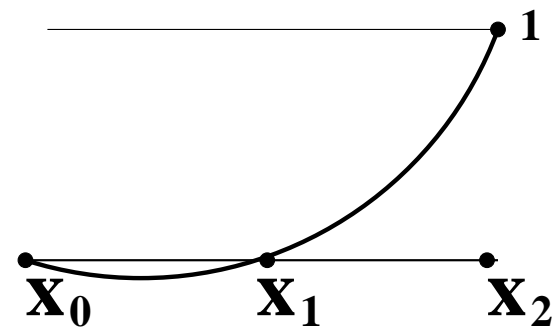
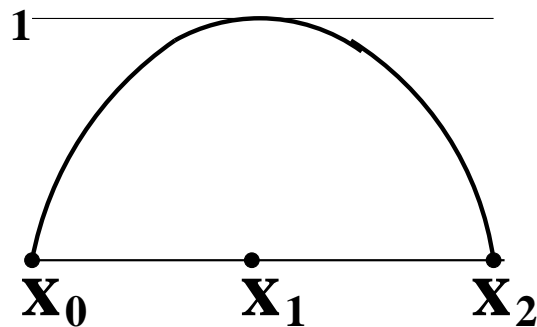
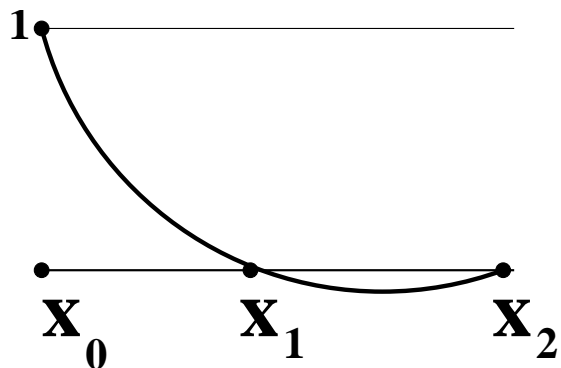
$$\ell_1(x) = \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)},$$

and

$$\ell_2(x) = \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}.$$

Note that $\ell_i \in \mathbb{P}_2$, $i = 0, 1, 2$, and that

$$\ell_i(x_k) = \begin{cases} 0 & \text{if } k \neq i, \\ 1 & \text{if } k = i. \end{cases}$$



$$l_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)}, \quad l_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)}, \quad l_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)}.$$

Lagrange basis functions (case $n = 2$).

Now given $f(x)$ let

$$p(x) = \sum_{k=0}^n f(x_k) \ell_k(x) .$$

Then $p \in \mathbb{P}_n$

and

$$p(x_i) = \sum_{k=0}^n f(x_k) \ell_k(x_i) = f(x_i) ,$$

that is, $p(x)$ *interpolates* $f(x)$ at the points x_0, x_1, \dots, x_n .

THEOREM: Let $f(x)$ be defined on $[a, b]$ and let

$$a \leq x_0 < x_1 < \cdots < x_n \leq b .$$

Then there is a unique polynomial $p \in \mathbb{P}_n$ that interpolates f at the $\{x_k\}_{k=0}^n$.

PROOF:

We have already demonstrated the existence of $p(x)$.

Suppose $q \in \mathbb{P}_n$ also interpolates f at the points $\{x_k\}_{k=0}^n$.

Let $r(x) \equiv p(x) - q(x)$.

Then $r \in \mathbb{P}_n$ and $r(x_k) = 0$, $k = 0, 1, \cdots, n$.

But $r \in \mathbb{P}_n$ can have at most n zeroes, unless $r(x) \equiv 0$.

Hence $r(x) \equiv 0$.

Thus $p(x) \equiv q(x)$, *i.e.*, p is unique.

QED!

EXAMPLE: Let $f(x) = e^x$.

Given $f(0) = 1$, $f(1) = 2.71828$, $f(2) = 7.38905$, we want to approximate $f(1.5)$ by polynomial interpolation at $x = 0, 1, 2$.

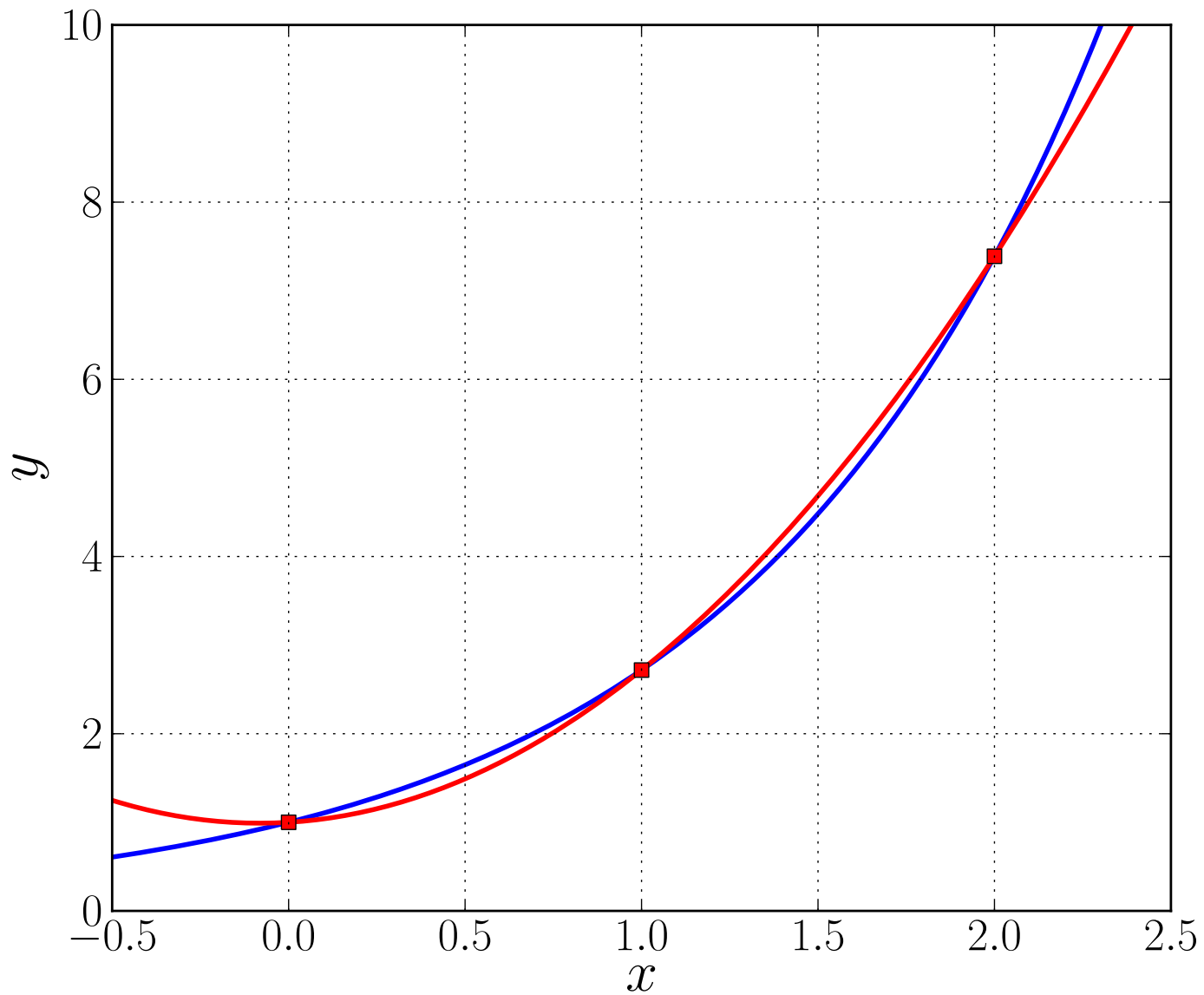
Here

$$\begin{aligned} \ell_0(1.5) &= \frac{(1.5 - 1)(1.5 - 2)}{(0 - 1)(0 - 2)} = -\frac{1}{8}, \\ \ell_1(1.5) &= \frac{(1.5 - 0)(1.5 - 2)}{(1 - 0)(1 - 2)} = \frac{6}{8}, \\ \ell_2(1.5) &= \frac{(1.5 - 0)(1.5 - 1)}{(2 - 0)(2 - 1)} = \frac{3}{8}, \end{aligned}$$

so that

$$\begin{aligned} p(1.5) &= f(0) \ell_0(1.5) + f(1) \ell_1(1.5) + f(2) \ell_2(1.5) \\ &= (1) \left(-\frac{1}{8}\right) + (2.71828) \left(\frac{6}{8}\right) + (7.38905) \left(\frac{3}{8}\right) = 4.68460. \end{aligned}$$

The exact value is $f(1.5) = e^{1.5} = 4.48168$.



Graph of $f(x) = e^x$ (blue) and its Lagrange interpolant $p(x) \in \mathbb{P}_2$ (red).

THE LAGRANGE INTERPOLATION THEOREM:

Let

$$x_0 < x_1 < \cdots < x_n, \quad \text{and let } x \in \mathbb{R}.$$

Define

$$a \equiv \min\{x_0, x\} \quad \text{and} \quad b \equiv \max\{x_n, x\}.$$

Assume that $f \in \mathbb{C}^{n+1}[a, b]$.

Let $p \in \mathbb{P}_n$ be the unique polynomial that interpolates $f(x)$ at $\{x_k\}_{k=0}^n$.

Then

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{k=0}^n (x - x_k),$$

for some point $\xi \equiv \xi(x) \in [a, b]$.

PROOF:

If $x = x_k$ for some k then the formula is clearly valid.

So assume that $x \neq x_k$, for $k = 0, 1, \dots, n$.

Let

$$w(z) \equiv \prod_{k=0}^n (z - x_k) \quad \text{and} \quad c(x) \equiv \frac{f(x) - p(x)}{w(x)} .$$

Then $c(x)$ is well defined since $w(x) \neq 0$.

We want to show that

$$c(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} . \quad (\text{Why?})$$

Consider

$$F(z) \equiv f(z) - p(z) - w(z) c(x) .$$

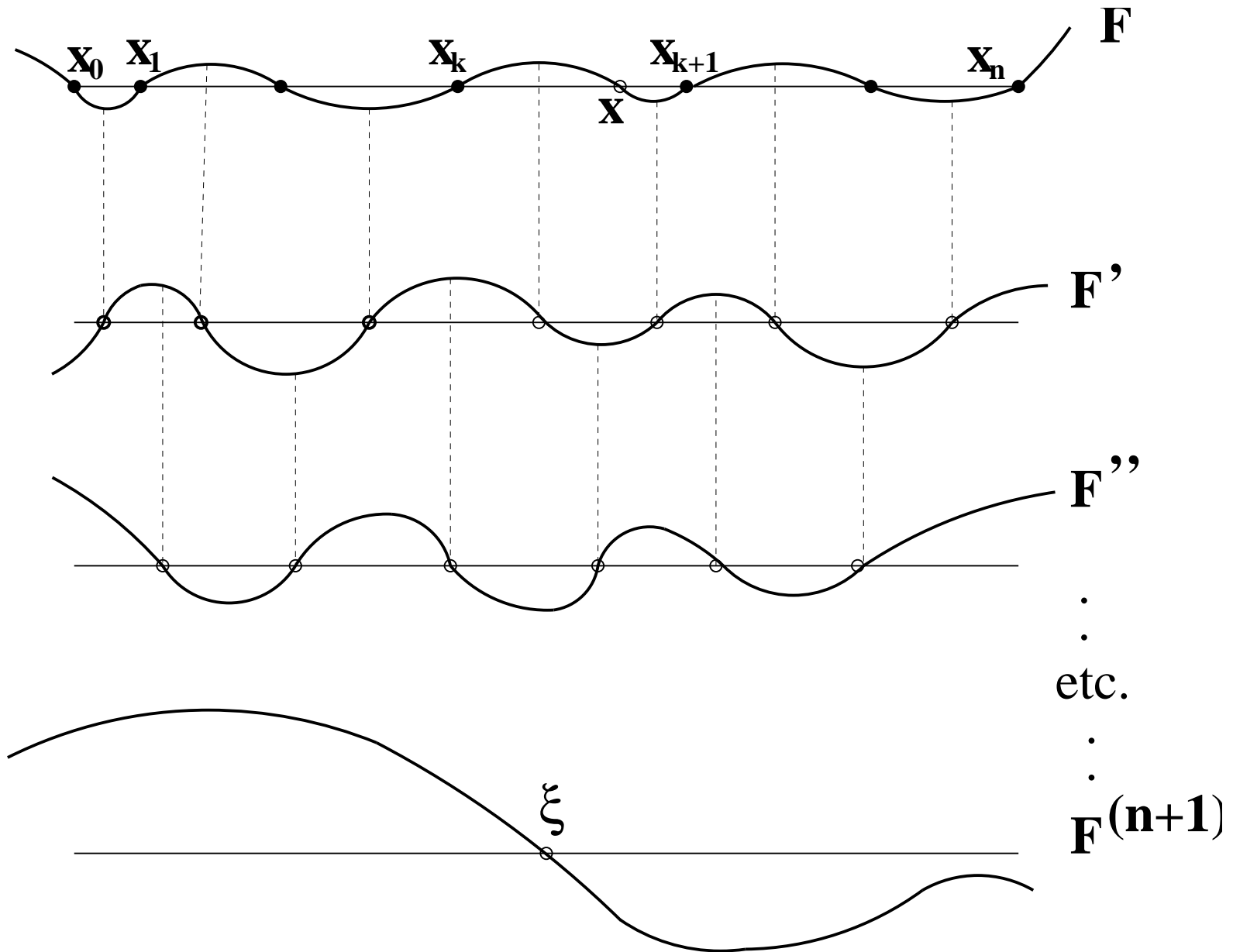
Then

$$F(x_k) = f(x_k) - p(x_k) - w(x_k) c(x) = 0, \quad k = 0, 1, \dots, n ,$$

and

$$F(x) = f(x) - p(x) - w(x) \frac{f(x) - p(x)}{w(x)} = 0 .$$

Thus $F(z)$ has (at least) $n + 2$ distinct zeroes in $[a, b]$.



The zeroes of $F(x)$ and its derivatives.

Hence, by Rolle's Theorem, $F'(z)$ has $n + 1$ distinct zeroes in $[a, b]$,

$F''(z)$ has n distinct zeroes in $[a, b]$,

$F'''(z)$ has $n - 1$ distinct zeroes in $[a, b]$, *etc.*

We find that $F^{(n+1)}(z)$ has (at least) one zero in $[a, b]$, say,

$$F^{(n+1)}(\xi) = 0, \quad \xi \in [a, b].$$

But

$$F^{(n+1)}(z) = f^{(n+1)}(z) - p^{(n+1)}(z) - w^{(n+1)}(z) c(x).$$

Hence

$$F^{(n+1)}(\xi) = f^{(n+1)}(\xi) - (n + 1)! c(x) = 0.$$

It follows that

$$c(x) = \frac{f^{(n+1)}(\xi)}{(n + 1)!}.$$

EXAMPLE: In the last example we had

$$n = 2 , \quad f(x) = e^x , \quad x_0 = 0 , \quad x_1 = 1 , \quad x_2 = 2 ,$$

and we computed the value of $p(x)$ at $x = 1.5$.

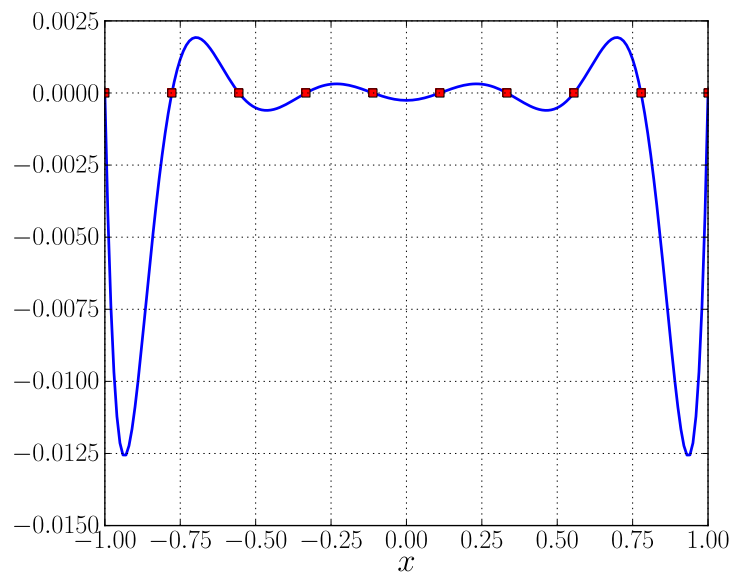
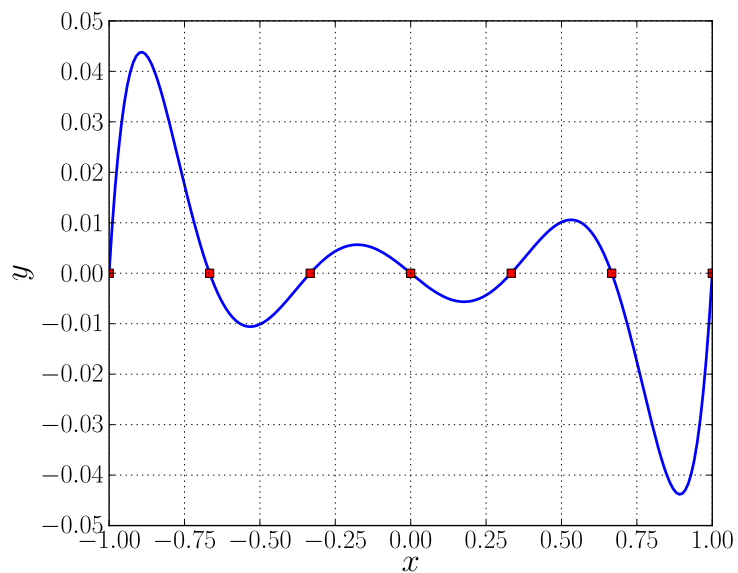
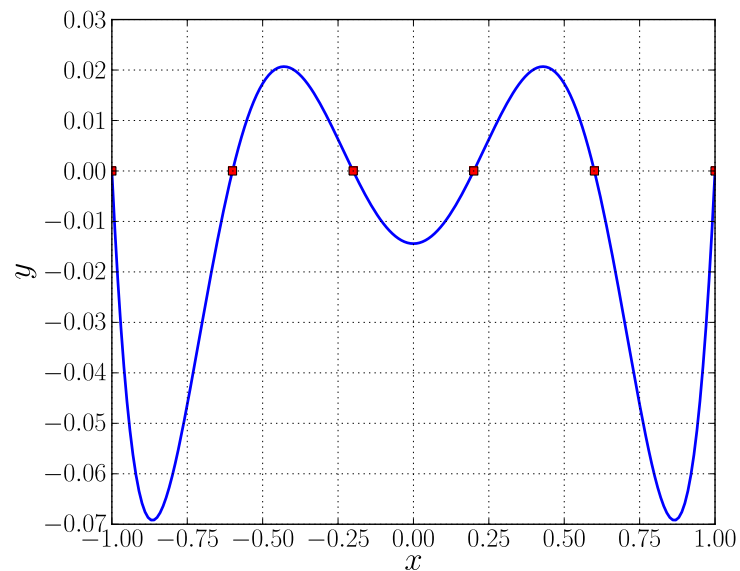
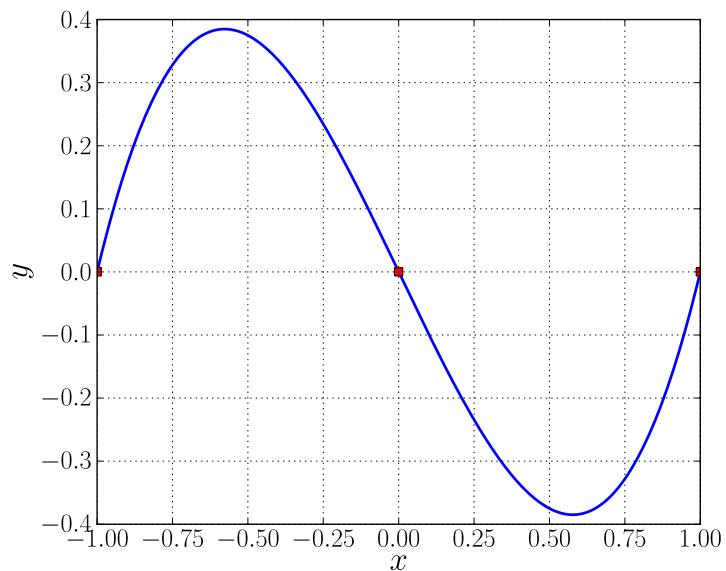
By the Theorem

$$f(x) - p(x) = \frac{f^{(3)}(\xi)}{3!} (x - 0)(x - 1)(x - 2) , \quad \xi \in [0, 2] .$$

Since $f^{(3)}(\xi) \leq e^2 < 7.4$ we find that

$$| f(1.5) - p(1.5) | < \frac{7.4}{6} (1.5) (0.5) (0.5) = \frac{7.4}{16} < 0.47 .$$

The actual error is $| p(1.5) - e^{1.5} | \approx | 4.68460 - 4.48168 | \approx 0.2 .$



The graph of $w_{n+1}(x) = \prod_{k=0}^n (x - x_k)$ for equally spaced interpolation points in the interval $[-1, 1]$, for the cases $n + 1 = 3, 6, 7, 10$.

n	max	n	max	n	max	n	max
1	1.00000	5	0.06918	9	0.01256	13	0.00278
2	0.38490	6	0.04382	10	0.00853	14	0.00193
3	0.19749	7	0.02845	11	0.00584	15	0.00134
4	0.11348	8	0.01877	12	0.00400	16	0.00095

The maximum value of $|w_{n+1}(x)|$ in the interval $[-1, 1]$ for the case of $n + 1$ equally spaced interpolation points .

EXERCISES:

- Consider the polynomial $p_n(x)$ of degree n (or less) that interpolates $f(x) = \sin(x)$ at $n + 1$ distinct points in $[-1, 1]$. Write down the general error formula for $|\sin(x) - p_n(x)|$. For distinct, but otherwise *arbitrary interpolation points*, how big should n be to guarantee that the maximum interpolation error in $[-1, 1]$ is less than 10^{-2} ?
- Also answer the above question for *equally spaced interpolation points* in $[-1, 1]$, using the Table on the preceding page .
- Also answer the above questions for the case of $f(x) = e^x$ in $[-1, 1]$.
- Consider the problem of interpolating a smooth function $f(x)$ at two points, $x_0 = -h/2$ and $x_1 = h/2$, by a polynomial $p \in \mathbb{P}_3$ such that
$$p(x_0) = f(x_0), \quad p'(x_0) = f'(x_0), \quad p(x_1) = f(x_1), \quad p'(x_1) = f'(x_1).$$
Prove that this interpolation problem has one and only one solution.

○ Consider the problem of interpolating a smooth function $f(x)$ at two points, $x_0 = -h/2$ and $x_1 = h/2$, by a polynomial $p \in \mathbb{P}_3$ such that

$$p(x_0) = f(x_0), \quad p'(x_0) = f'(x_0), \quad p(x_1) = f(x_1), \quad p'(x_1) = f'(x_1).$$

Prove that this interpolation problem has one and only one solution.

SOLUTION: Write $p(x) = c_0 + c_1x + c_2x^2 + c_3x^3$. Then we have

$$p\left(-\frac{h}{2}\right) = c_0 - c_1\frac{h}{2} + c_2\left(\frac{h}{2}\right)^2 - c_3\left(\frac{h}{2}\right)^3 = f\left(-\frac{h}{2}\right),$$

$$p\left(\frac{h}{2}\right) = c_0 + c_1\frac{h}{2} + c_2\left(\frac{h}{2}\right)^2 + c_3\left(\frac{h}{2}\right)^3 = f\left(\frac{h}{2}\right),$$

$$p'\left(-\frac{h}{2}\right) = c_1 - 2\frac{h}{2}c_2 + 3\left(\frac{h}{2}\right)^2c_3 = f'\left(-\frac{h}{2}\right),$$

$$p'\left(\frac{h}{2}\right) = c_1 + 2\frac{h}{2}c_2 + 3\left(\frac{h}{2}\right)^2c_3 = f'\left(\frac{h}{2}\right).$$

SOLUTION: continued \dots :

In matrix form

$$\begin{pmatrix} 1 & -\frac{h}{2} & \frac{h^2}{4} & -\frac{h^3}{8} \\ 1 & \frac{h}{2} & \frac{h^2}{4} & \frac{h^3}{8} \\ 0 & 1 & -h & \frac{3h^2}{4} \\ 0 & 1 & h & \frac{3h^2}{4} \end{pmatrix} \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ f'_0 \\ f'_1 \end{pmatrix},$$

where the matrix can be transformed to upper-triangular form as follows:

$$\begin{pmatrix} 1 & -\frac{h}{2} & \frac{h^2}{4} & -\frac{h^3}{8} \\ 1 & \frac{h}{2} & \frac{h^2}{4} & \frac{h^3}{8} \\ 0 & 1 & -h & \frac{3h^2}{4} \\ 0 & 1 & h & \frac{3h^2}{4} \end{pmatrix} \rightarrow \begin{pmatrix} 1 & -\frac{h}{2} & \frac{h^2}{4} & -\frac{h^3}{8} \\ 0 & h & 0 & \frac{h^3}{4} \\ 0 & 1 & -h & \frac{3h^2}{4} \\ 0 & 1 & h & \frac{3h^2}{4} \end{pmatrix}$$

$$\rightarrow \begin{pmatrix} 1 & -\frac{h}{2} & \frac{h^2}{4} & -\frac{h^3}{8} \\ 0 & h & 0 & \frac{h^3}{4} \\ 0 & 0 & -h & \frac{h^2}{2} \\ 0 & 0 & h & \frac{h^2}{2} \end{pmatrix} \rightarrow \begin{pmatrix} 1 & -\frac{h}{2} & \frac{h^2}{4} & -\frac{h^3}{8} \\ 0 & h & 0 & \frac{h^3}{4} \\ 0 & 0 & -h & \frac{h^2}{2} \\ 0 & 0 & 0 & h^2 \end{pmatrix}$$

The determinant of the upper-triangular matrix is $-h^4$, which is not zero. Thus the system is uniquely solvable.

Chebyshev Polynomials.

From the preceding Theorem it follows that

$$\| f - p \|_{\infty} \leq \frac{1}{(n+1)!} \| f^{(n+1)} \|_{\infty} \| w_{n+1} \|_{\infty} ,$$

where

$$w_{n+1}(x) \equiv \prod_{k=0}^n (x - x_k) ,$$

and where

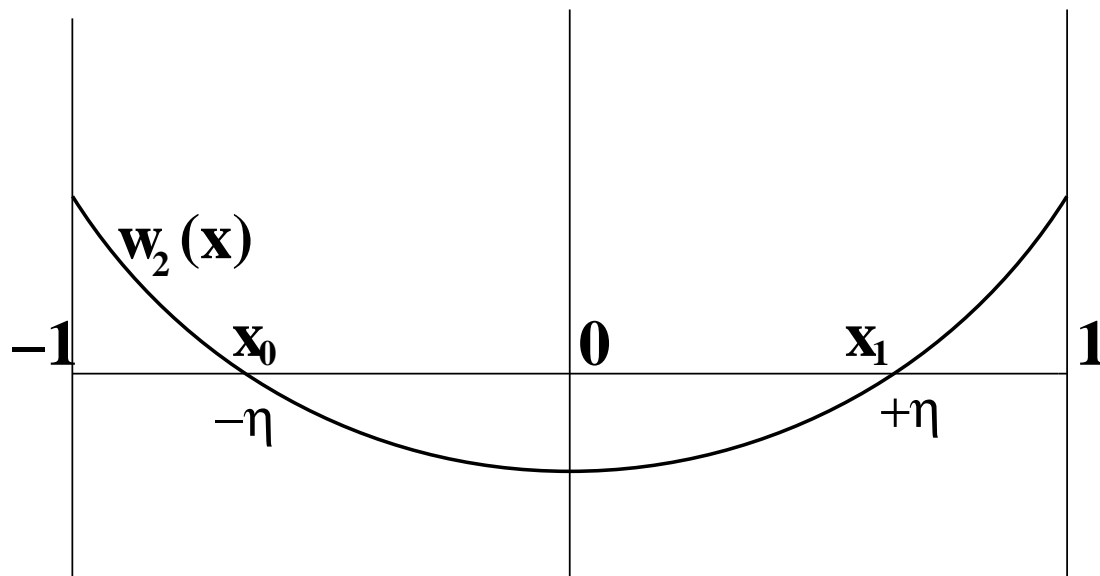
$$\| w_{n+1} \|_{\infty} \equiv \max_{[a,b]} | w_{n+1}(x) | .$$

REMARKS:

- It does *not* follow that $\| f - p \|_{\infty} \rightarrow 0$ as $n \rightarrow \infty$.
- There are examples where $\| f - p \|_{\infty} \rightarrow \infty$ as $n \rightarrow \infty$.
- For given n , can we choose $\{x_k\}_{k=0}^n$ so that $\| w_{n+1} \|_{\infty}$ is minimized ?

EXAMPLE:

Let $n = 1$ and place the x_0 and x_1 symmetrically in $[-1, 1]$:



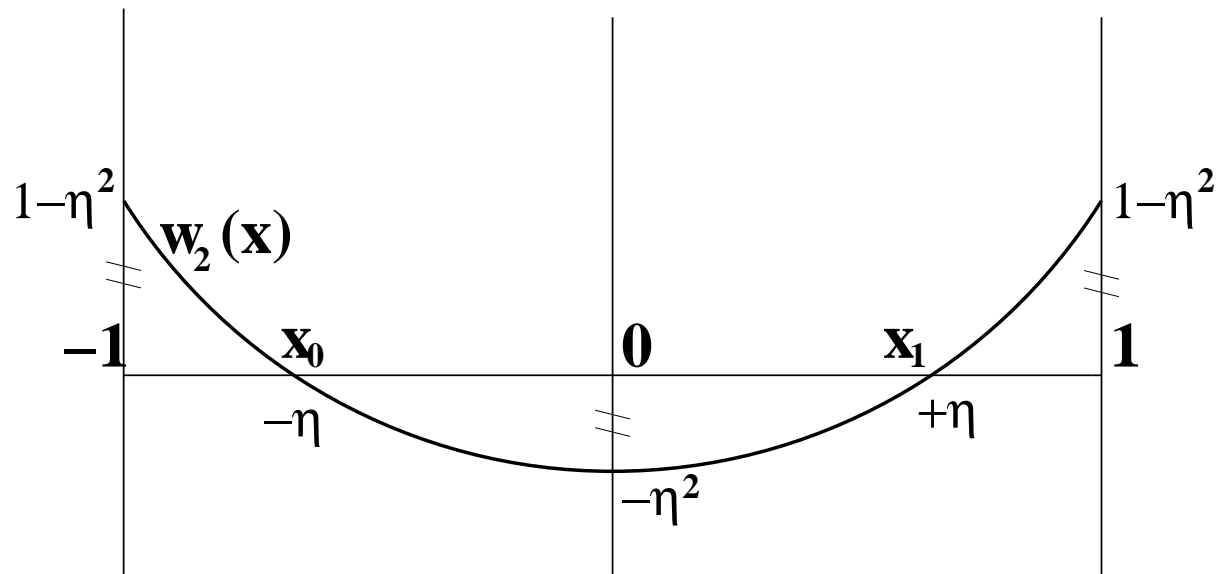
Then

$$w_2(x) = (x - x_0)(x - x_1) = (x + \eta)(x - \eta) = x^2 - \eta^2 .$$

We want to choose η such that

$$\| w_2 \|_{\infty} \equiv \max_{[-1,1]} | w_2(x) |$$

is minimized.



At the critical point : $w_2(0) = -\eta^2$.

At the endpoints : $w_2(-1) = w_2(1) = 1 - \eta^2$.

We see that $\| w_2 \|_\infty$ is minimized if we take η such that

$$| w_2(-1) | = | w_2(0) | = | w_2(1) | , \quad \text{i.e., if } \eta^2 = 1 - \eta^2 ,$$

Thus $\eta = \frac{1}{2}\sqrt{2}$ and $\| w_2 \|_\infty = \frac{1}{2}$.

In general, the points

$$\{x_k\}_{k=0}^n \text{ that minimize } \|w_{n+1}\|_\infty \text{ on } [-1, 1]$$

are the zeroes of the *Chebyshev Polynomial* T_{n+1} of degree $n + 1$.

These polynomials are defined as

$$T_k(x) \equiv \cos(k \cos^{-1}(x)) , \quad k = 0, 1, 2, \dots, \quad x \in [-1, 1] .$$

The T_k are indeed polynomials :

First $T_0(x) \equiv 1$ and $T_1(x) = x$.

Also

$$T_{k+1}(x) = 2x T_k(x) - T_{k-1}(x) .$$

$$T_{k+1}(x) = 2x T_k(x) - T_{k-1}(x)$$

To derive this recurrence formula we use the identity

$$\cos((k+1)\theta) + \cos((k-1)\theta) = 2 \cos(k\theta) \cos(\theta) .$$

which we rewrite as

$$\cos((k+1)\theta) = 2 \cos(k\theta) \cos(\theta) - \cos((k-1)\theta) .$$

so that, taking $\theta = \cos^{-1}(x)$, we have

$$\begin{aligned} T_{k+1}(x) &= \cos((k+1) \cos^{-1}(x)) \\ &= 2 \cos(k \cos^{-1}(x)) \cos(\cos^{-1}(x)) - \cos((k-1) \cos^{-1}(x)) \\ &= 2x T_k(x) - T_{k-1}(x) . \end{aligned}$$

Thus

$$T_2(x) = 2x T_1(x) - T_0(x) = 2x^2 - 1 ,$$

$$T_3(x) = 2x T_2(x) - T_1(x) = 4x^3 - 3x ,$$

$$T_4(x) = 2x T_3(x) - T_2(x) = 8x^4 - 8x^2 + 1 ,$$

·
·
·

$$T_{n+1}(x) = 2^n x^{n+1} + \dots .$$

THE CHEBYSHEV THEOREM:

Let

$$w_{n+1}(x) \equiv \prod_{k=0}^n (x - x_k) .$$

Then for fixed n the quantity

$$\| w_{n+1} \|_{\infty} \equiv \max_{[-1,1]} | w_{n+1}(x) |$$

is minimized if the points $\{x_k\}_{k=0}^n$ are the zeroes of $T_{n+1}(x)$.

For these points the value of $\| w_{n+1} \|_{\infty}$ is equal to 2^{-n} .

PROOF:

$$T_{n+1}(x) = \cos((n+1) \cos^{-1}(x)) = 0 ,$$

if

$$(n+1) \cos^{-1}(x) = (2k+1) \frac{\pi}{2} ,$$

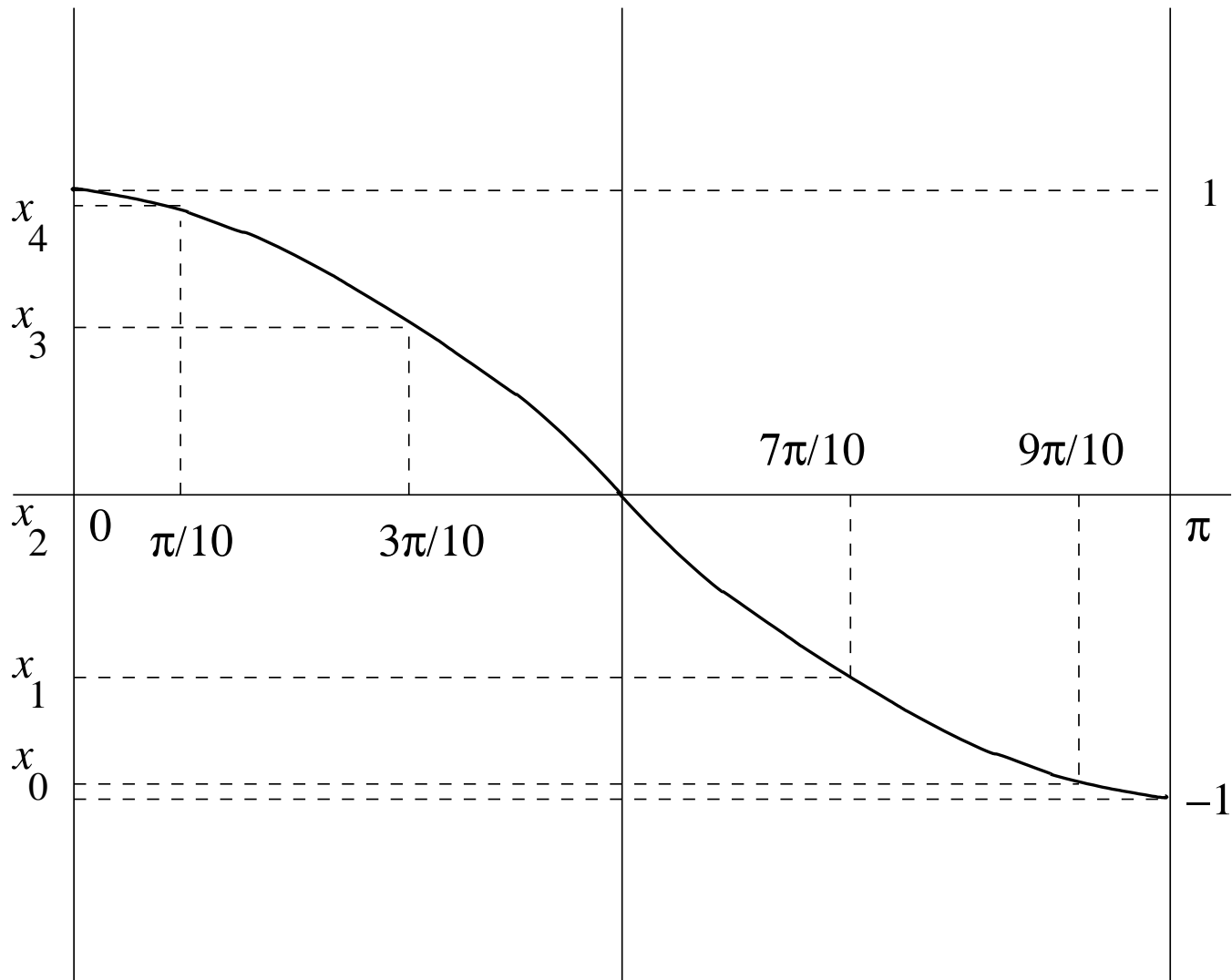
i.e., $T_{n+1}(x) = 0$ if

$$x = \cos\left(\frac{2k+1}{2(n+1)}\pi\right) , \quad k = 0, 1, 2, \dots, n .$$

Hence the zeroes of $T_{n+1}(x)$ lie indeed in $[-1, 1]$.

There are $n+1$ such zeroes.

$$x = \cos\left(\frac{2k+1}{2(n+1)}\pi\right), \quad k = 0, 1, 2, \dots, n.$$



The Chebyshev points x_k , ($k = 0, 1, 2, \dots, n$), for the case $n = 4$.

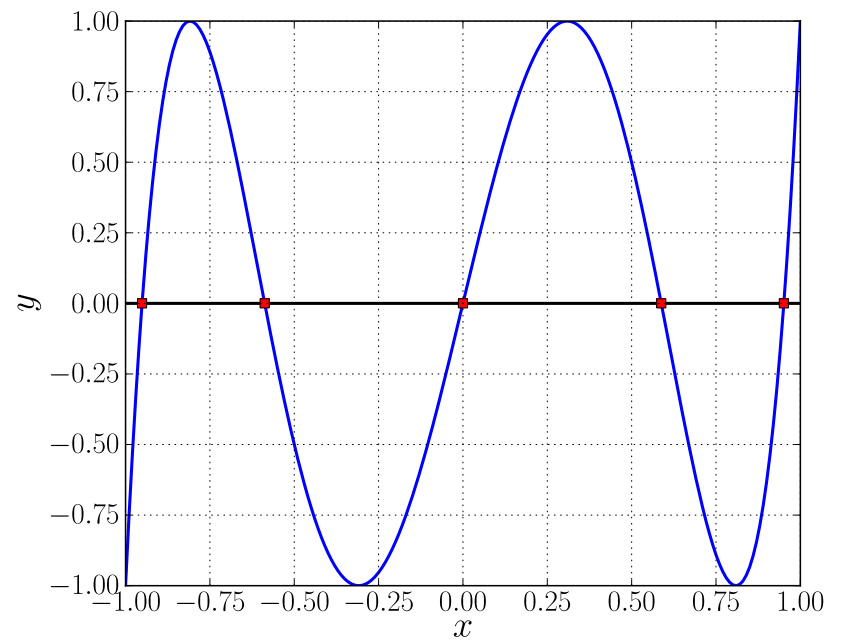
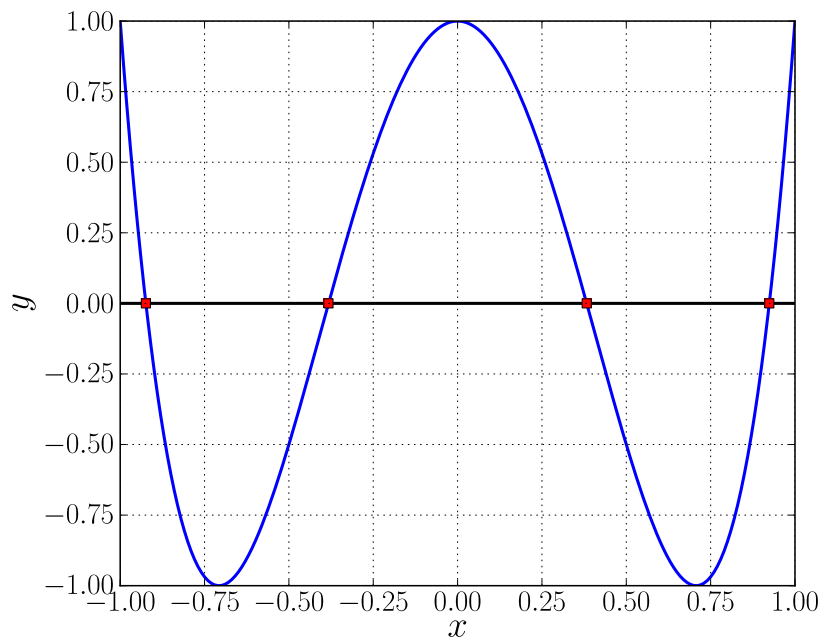
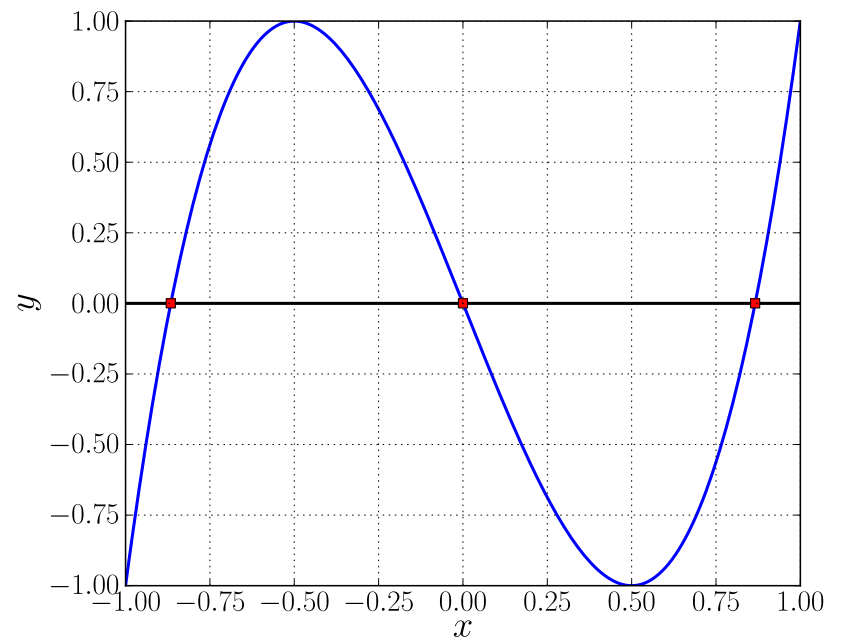
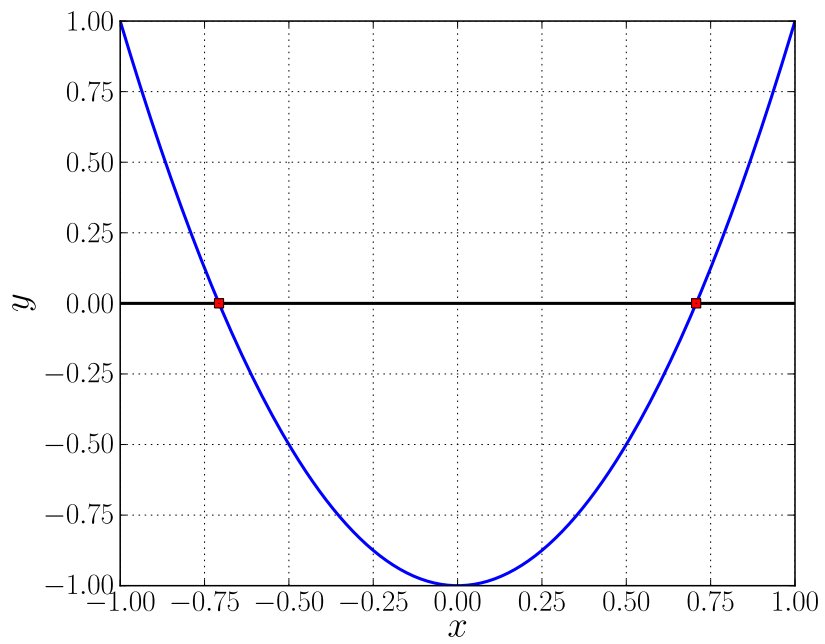
$$T_{n+1}(x) = \cos((n+1) \cos^{-1}(x)) = 0 ,$$

$$T_{n+1}(x) = \pm 1 \text{ if } (n+1) \cos^{-1}(x) = k\pi ,$$

that is, if,

$$x = \cos\left(\frac{k}{n+1}\pi\right) , \quad k = 0, 1, 2, \dots, n+1 .$$

We can now draw the *graph* of T_{n+1} :



The graph of T_n for the cases $n = 2, 3, 4, 5$.

Recall that that

$$T_{n+1}(x) = 2^n x^{n+1} + \dots .$$

Thus we can also write

$$T_{n+1}(x) = 2^n (x - x_0) (x - x_1) \cdots (x - x_n) ,$$

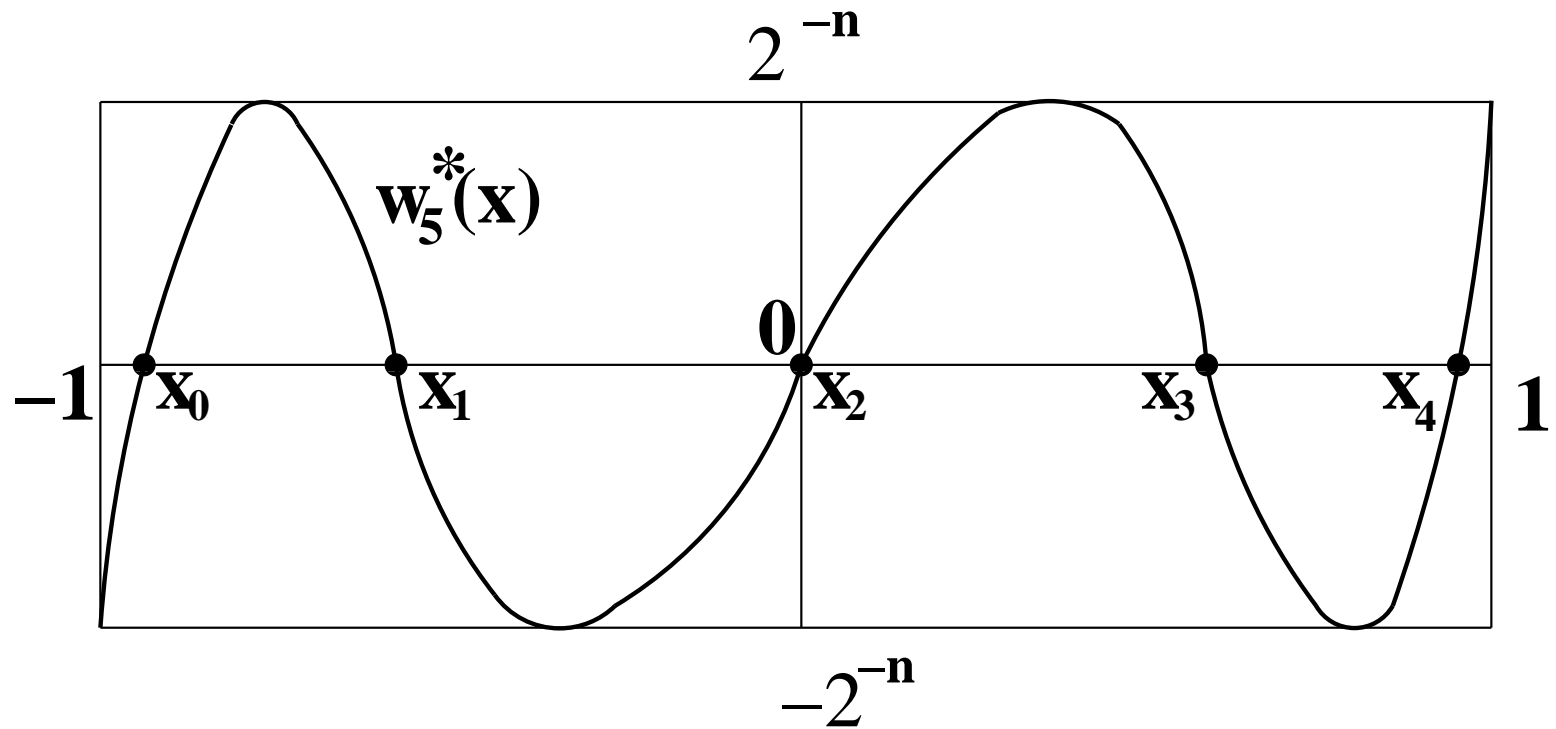
where the x_k are the zeroes of $T_{n+1}(x)$.

Let

$$w_{n+1}^*(x) \equiv 2^{-n} T_{n+1}(x) = (x - x_0) (x - x_1) \cdots (x - x_n) ,$$

Then

$$\| w_{n+1}^* \|_\infty = \| 2^{-n} T_{n+1} \|_\infty = 2^{-n} \| T_{n+1} \|_\infty = 2^{-n} .$$



The graph of $w_{n+1}^* = \prod_{k=0}^n (x - x_k)$ for the case $n = 4$.

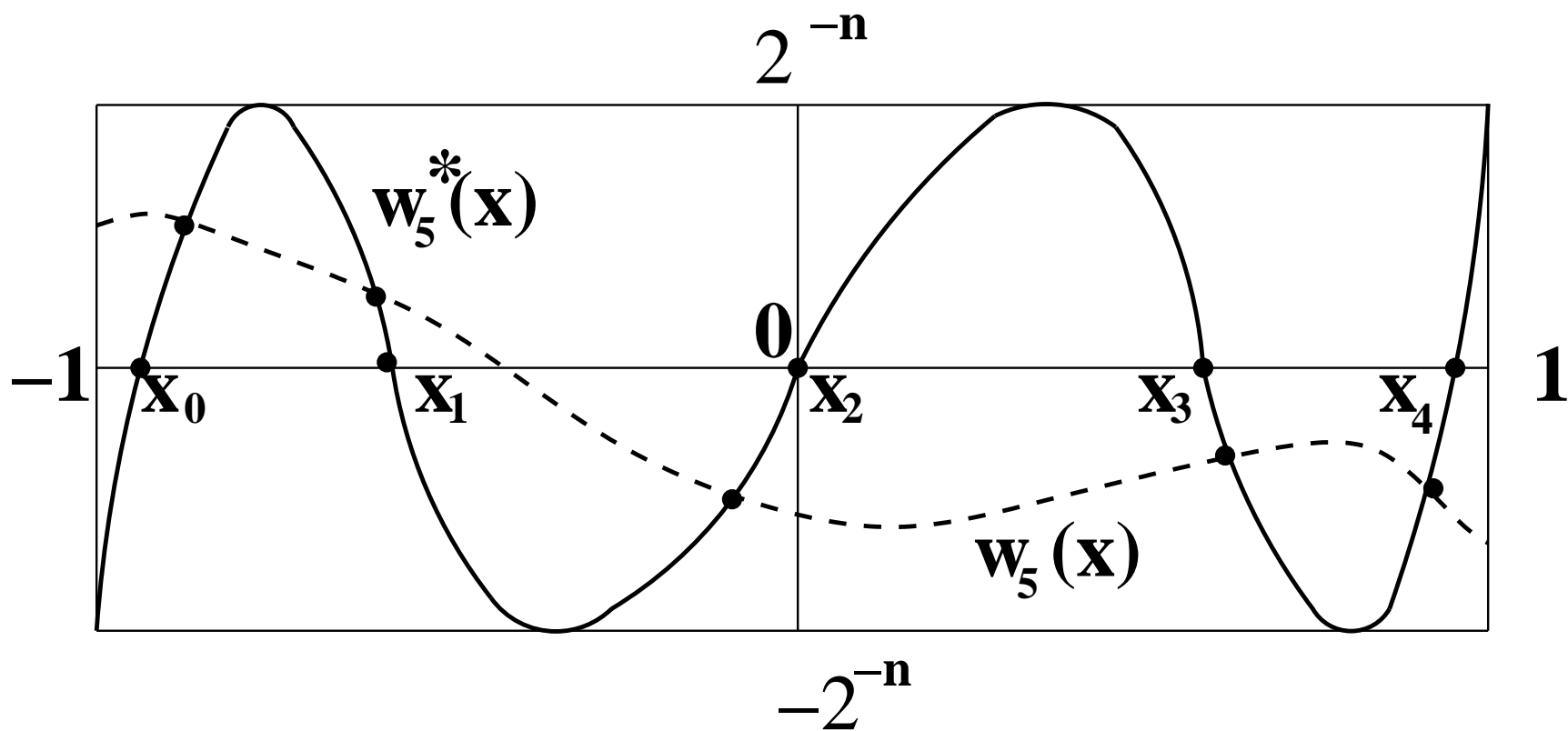
Claim :

There does not exist $w \in \mathbb{P}_{n+1}$, with leading coefficient 1, such that

$$\|w\|_{\infty} < 2^{-n}.$$

Suppose there does exist a $w_{n+1} \in \mathbb{P}_{n+1}$, with leading coefficient 1, such that

$$\|w_{n+1}\|_{\infty} < \|w_{n+1}^*\|_{\infty} = 2^{-n}.$$



Then w_{n+1} must intersect w_{n+1}^* at least $n + 1$ times in $[-1, 1]$.

Thus $(w_{n+1} - w_{n+1}^*)$ has $n + 1$ zeroes in $[-1, 1]$.

But

$$(w_{n+1} - w_{n+1}^*) \in \mathbb{P}_n$$

since both w_{n+1} and w_{n+1}^* have leading coefficient 1.

Hence $w_{n+1} - w_{n+1}^* \equiv 0$.

Thus $w_{n+1} = w_{n+1}^*$. QED!

n	uniform	Chebyshev	n	uniform	Chebyshev
1	1.00000	0.50000	5	0.06918	0.03125
2	0.38490	0.25000	6	0.04382	0.01563
3	0.19749	0.12500	7	0.02845	0.00782
4	0.11348	0.06250	8	0.01877	0.00391

The maximum of $|w_{n+1}(x)|$ in the interval $[-1, 1]$ for uniformly spaced and for Chebyshev points .

EXAMPLE:

Let $f(x) = e^x$ on $[-1, 1]$ and take $n = 2$.

$T_3(x) = 4x^3 - 3x$ has zeroes

$$x_0 = -\frac{1}{2}\sqrt{3}, \quad x_1 = 0, \quad x_2 = \frac{1}{2}\sqrt{3}.$$

Approximate $f(0.5)$ by polynomial interpolation at x_0, x_1, x_2 :

$$\ell_0(0.5) = \frac{(0.5 - x_1)(0.5 - x_2)}{(x_0 - x_1)(x_0 - x_2)} = \frac{1 - \sqrt{3}}{6},$$

$$\ell_1(0.5) = \frac{(0.5 - x_0)(0.5 - x_2)}{(x_1 - x_0)(x_1 - x_2)} = \frac{4}{6},$$

$$\ell_2(0.5) = \frac{(0.5 - x_0)(0.5 - x_1)}{(x_2 - x_0)(x_2 - x_1)} = \frac{1 + \sqrt{3}}{6}.$$

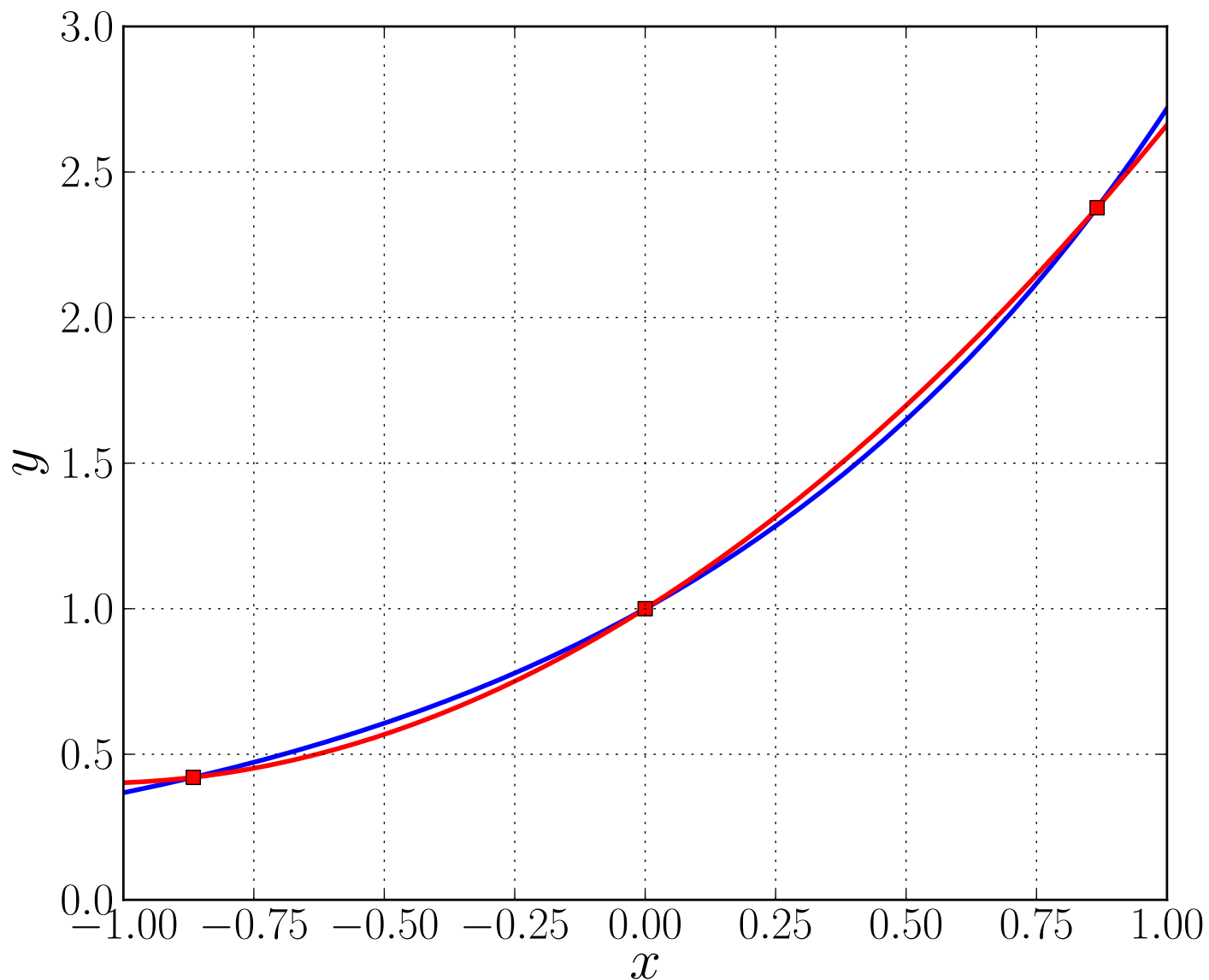
Thus

$$\begin{aligned} p(0.5) &= f(x_0) l_0(0.5) + f(x_1) l_1(0.5) + f(x_2) l_2(0.5) \\ &= e^{(-0.5\sqrt{3})} \frac{(1 - \sqrt{3})}{6} + e^0 \left(\frac{4}{6}\right) + e^{(0.5\sqrt{3})} \frac{(1 + \sqrt{3})}{6} \\ &\approx 1.697 . \end{aligned}$$

The exact value is $e^{0.5} = 1.648\dots$.

Thus the exact absolute error is

$$| e^{0.5} - p(0.5) | \approx | 1.648 - 1.697 | = 0.049 .$$



Graph of $f(x) = e^x$ (blue) on the interval $[-1, 1]$,
and its Lagrange interpolating polynomial $p(x) \in \mathbb{P}_2$ (red)
at three Chebyshev interpolation points ($n = 2$).

EXAMPLE: More generally, if we interpolate

$$f(x) = e^x \text{ by } p \in \mathbb{P}_n \text{ at } n + 1 \text{ Chebyshev points in } [-1, 1],$$

then for $x \in [-1, 1]$ we have

$$|e^x - p(x)| = \left| \frac{f^{(n+1)}(\xi)}{(n+1)!} w_{n+1}(x) \right|,$$

where $\xi \equiv \xi(x) \in [-1, 1]$, and where

$$w_{n+1}(x) = \prod_{k=0}^n (x - x_k),$$

Thus

$$\begin{aligned} \max_{x \in [-1, 1]} |e^x - p(x)| &\leq \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} \|w_{n+1}\|_\infty \\ &\leq \frac{e}{(n+1)!} \|w_{n+1}\|_\infty \\ &= \frac{e}{(n+1)!} 2^{-n}. \end{aligned}$$

REMARK:

Let f be a sufficiently smooth function.

Let p_U be the polynomial that interpolates f at $n+1$ uniformly spaced points.

Let p_C denote the polynomial that interpolates f at $n+1$ Chebyshev points.

Although the Theorem does not *guarantee* that

$$\| p_C - f \|_\infty \leq \| p_U - f \|_\infty ,$$

this inequality is “usually” valid.

EXERCISES:

- Suppose $p \in \mathbb{P}^n$ interpolates $\sin(x)$ at $n + 1$ distinct points in $[-1, 1]$. For the case of *Chebyshev points*, how big should n be for the error to be less than 10^{-4} everywhere in $[-1, 1]$?
- Suppose that $p \in \mathbb{P}^n$ interpolates e^x at $n + 1$ distinct points in $[-1, 1]$. For the case of *Chebyshev points*, how big should n be for the error to be less than 10^{-4} everywhere in $[-1, 1]$?
- Suppose $p \in \mathbb{P}^n$ interpolates $\sin(x)$ at $n + 1$ distinct points in $[0, \pi]$. For the case of *Chebyshev points*, how big should n be for the maximum interpolation error in $[0, \pi]$ to be less than 10^{-4} ?

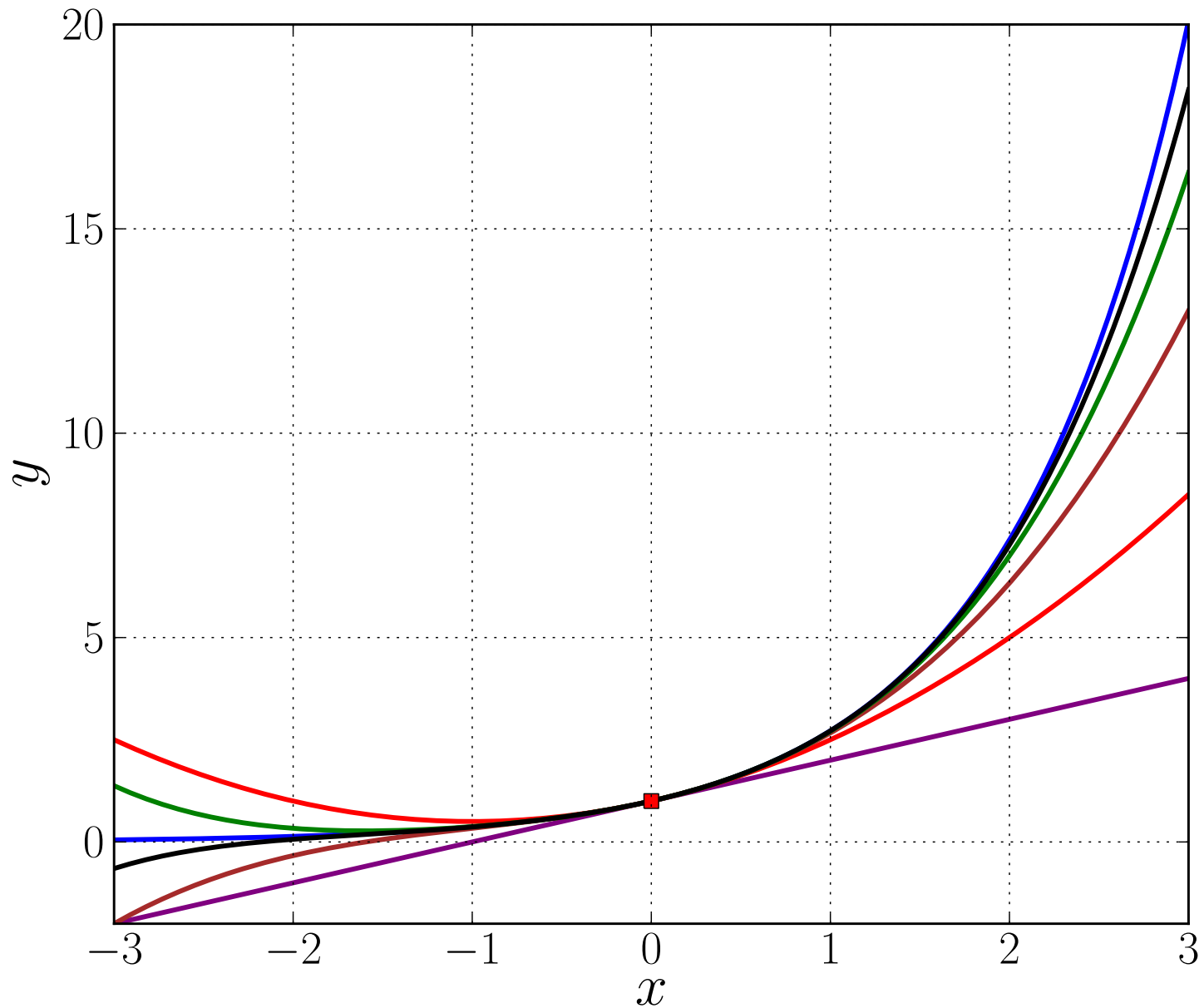
The Taylor Polynomial.

Let $f \in C^n[a, b]$.

Let \mathbb{P}_n denote all polynomials of degree less than or equal to n .

Given the point $x_0 \in [a, b]$, we want to find $p \in \mathbb{P}_n$ such that

$$p^{(k)}(x_0) = f^{(k)}(x_0) , \quad k = 0, 1, \dots, n .$$



The function e^x (blue) and its Taylor polynomials $p_k(x)$ about $x_0 = 0$:
 $k = 1$: purple, $k = 2$: red, $k = 3$: brown, $k = 4$: green, $k = 5$: black .

As for Lagrange interpolation, we have the following *questions* :

- Is the polynomial $p(x) \in \mathbb{P}_n$ uniquely defined ?
- How well does p approximate f ?
- Does the approximation get better as $n \rightarrow \infty$?

Existence :

$$p(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k .$$

Clearly

$$p^{(k)}(x_0) = f^{(k)}(x_0) , \quad k = 0, 1, \dots, n . \quad (\text{Check!})$$

DEFINITION :

$p(x)$ is called the *Taylor polynomial* of degree n for $f(x)$ about the point x_0 .

TAYLOR'S THEOREM:

Let $f \in C^{n+1}[a, b]$, $x_0 \in [a, b]$.

Let $p(x) \in \mathbb{P}_n$ be the Taylor polynomial for f about the point x_0 , *i.e.*,

$$p(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k .$$

Then, for $x \in [a, b]$,

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x - x_0)^{n+1} ,$$

for some point $\xi = \xi(x)$ that lies between x_0 and x .

DEFINITION: $R_n(x) \equiv \frac{f^{(n+1)}(\xi)}{(n+1)!} (x - x_0)^{n+1}$

is called the *Taylor remainder* .

PROOF: (Exercise!)

The steps are similar to those in the Lagrange Interpolation Theorem :

- First show that the Theorem holds if $x = x_0$.
- Next assume x is arbitrary, but $x \neq x_0$. (Consider x as fixed.)

- Define

$$c(x) = \frac{f(x) - p(x)}{(x - x_0)^{n+1}} .$$

- Define $F(z) = f(z) - p(z) - c(x) (z - x_0)^{n+1}$.
- Show that $F^{(k)}(x_0) = 0$, $k = 0, 1, \dots, n$, and that $F(x) = 0$.
- Give a qualitative graph of $F(z)$.

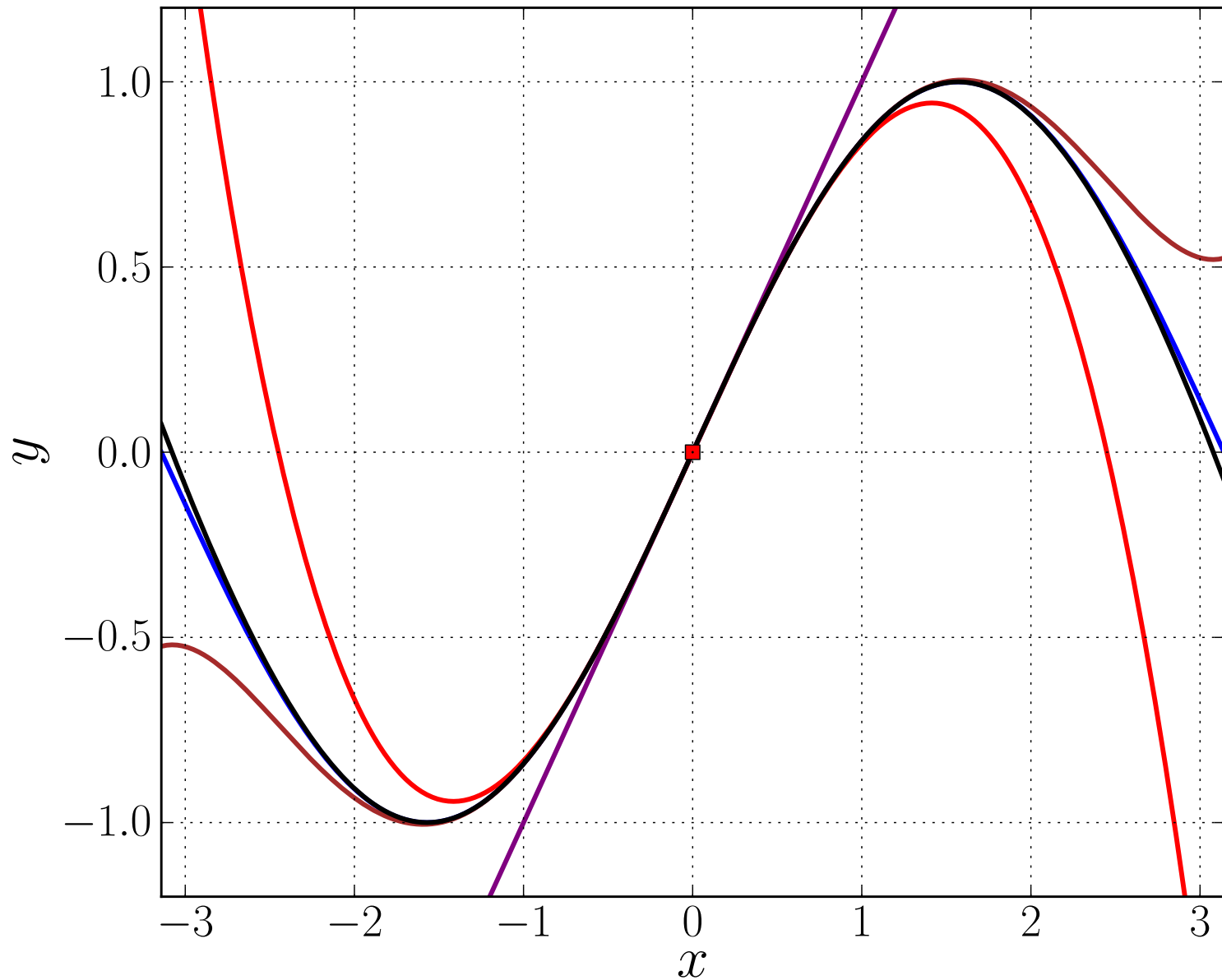
- Show $F'(\xi_0) = 0$ for some ξ_0 between x_0 and x . Graph $F'(z)$.
- Show $F''(\xi_1) = 0$ for some ξ_1 between x_0 and ξ_0 . Graph $F''(z)$.
- *etc.*
- Show that $F^{(n+1)}(\xi_n) = 0$ for some ξ_n between x_0 and ξ_{n-1} .
- From this derive that

$$c(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!}, \quad (\xi = \xi_n).$$

- Show how Taylor's Theorem follows from this last step. QED!

EXERCISES:

- Write down the Taylor polynomials $p_n(x)$ of degree n (or less) for $f(x) = e^x$ about the point $x_0 = 0$, for each of the following cases: $n = 1, 2, 3, 4$.
- How big should n be so that $|e^x - p_n(x)| < 10^{-4}$ everywhere in the interval $[-1, 1]$?
- Do the same for $f(x) = \sin(x)$ in $[0, 1]$ about the point $x_0 = 0$.
- Do the same for $f(x) = \ln(x)$ in $[\frac{1}{2}, \frac{3}{2}]$ about the point $x_0 = 1$.



Graph of the function $\sin(x)$ (blue) and its Taylor polynomials $p_k(x)$ about $x_0 = 0$: $k = 1$: purple, $k = 3$: red, $k = 5$: brown, $k = 7$: black .

Local Polynomial Interpolation.

Let $f \in \mathbb{C}^{n+1}[a, b]$.

Let $p \in \mathbb{P}_n$ interpolate f at $n + 1$ distinct points $\{x_k\}_{k=0}^n$ in $[a, b]$.

Does

$$\|f - p\|_{\infty} \rightarrow 0 \quad \text{as } n \rightarrow \infty ?$$

The answer is often NO !

For example if

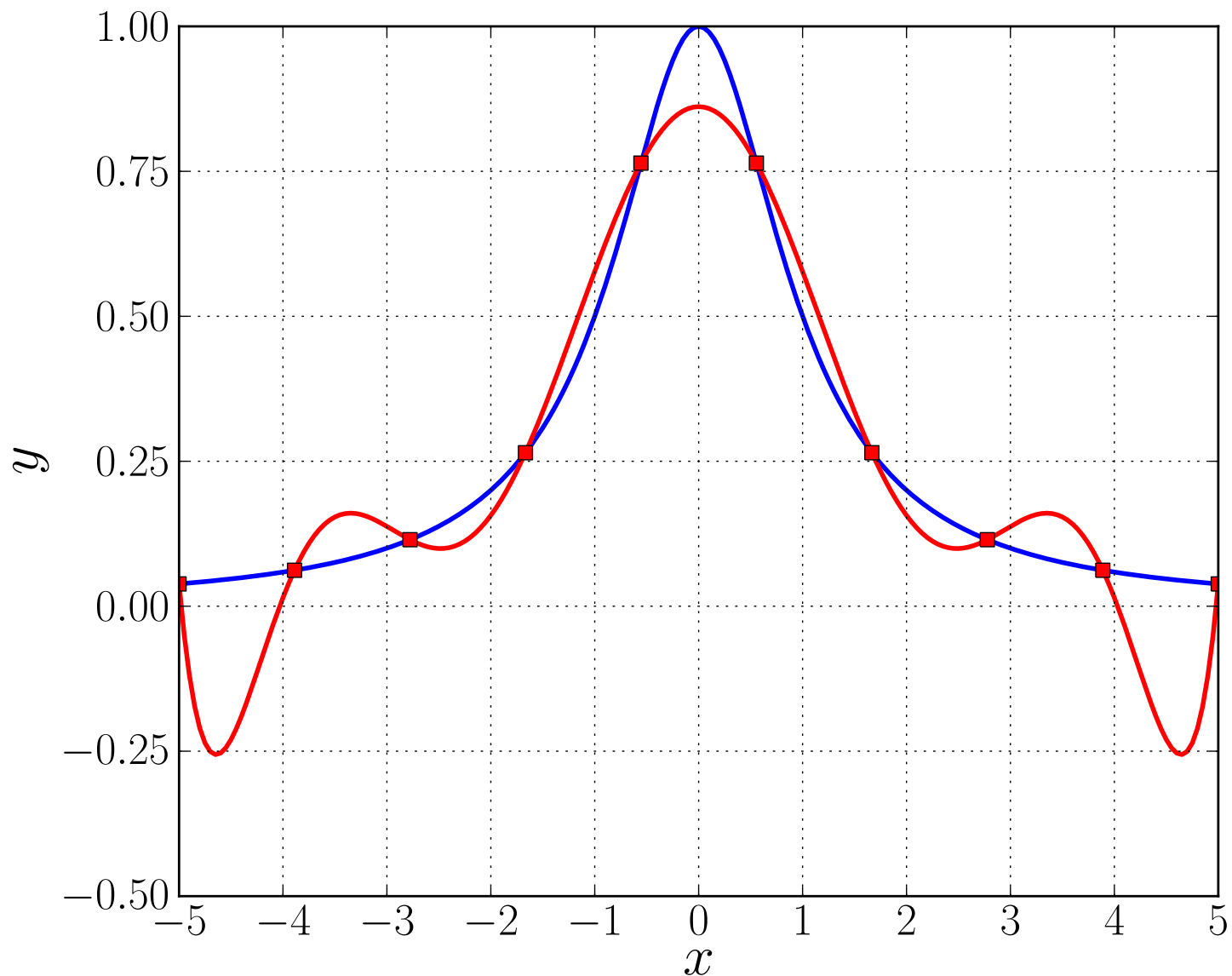
$$f(x) = \frac{1}{1+x^2} \quad \text{on} \quad [-5, 5],$$

and if $p \in \mathbb{P}_n$ interpolates f at the $n+1$ equally spaced points $\{x_k\}_{k=0}^n$ with

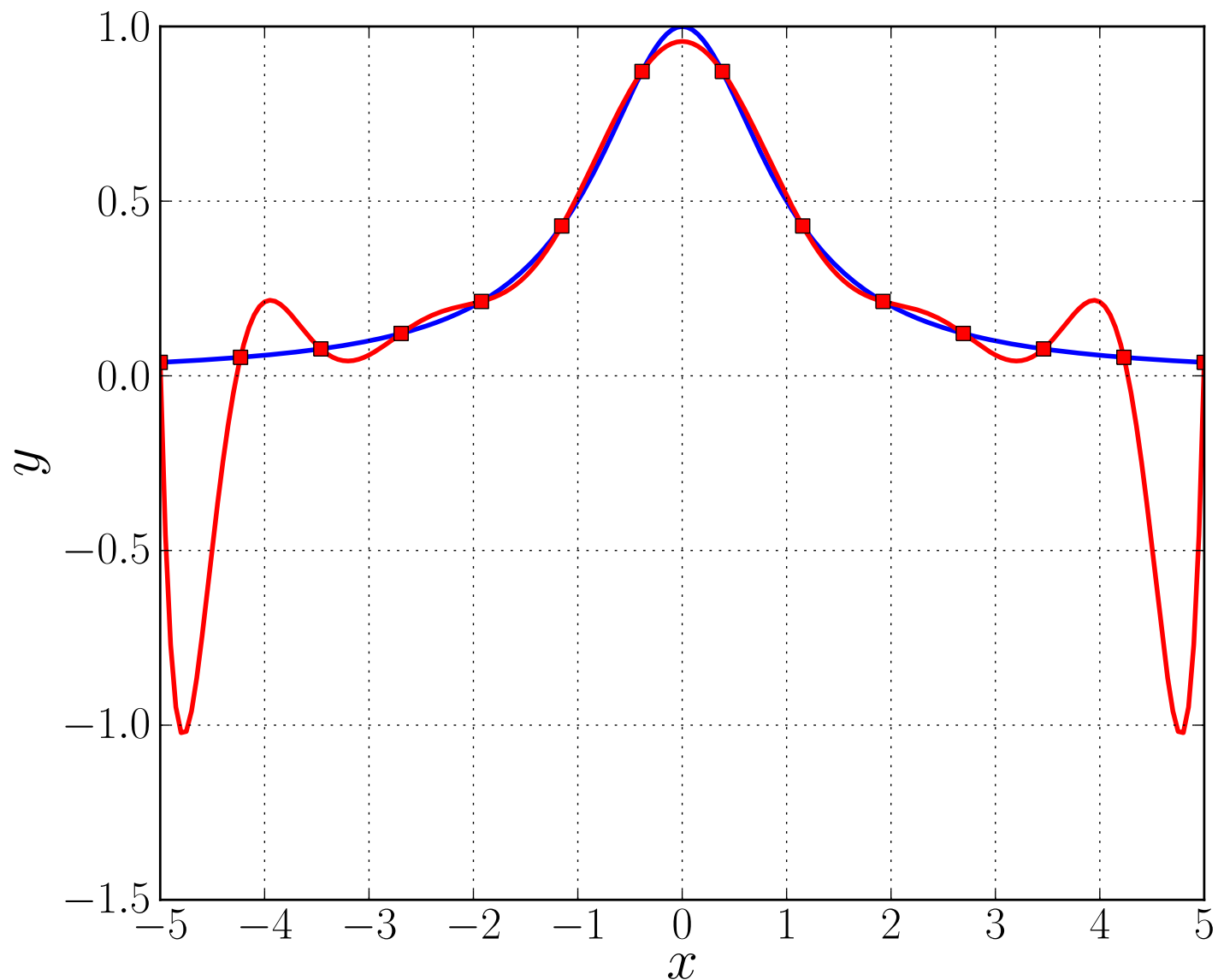
$$x_k = -5 + k \Delta x, \quad k = 0, 1, \dots, n, \quad \Delta x = \frac{10}{n},$$

then it can be shown that

$$\|f - p\|_{\infty} \rightarrow \infty \quad \text{as} \quad n \rightarrow \infty.$$



Graph of $f(x) = \frac{1}{1+x^2}$ on the interval $[-5, 5]$
 and its Lagrange interpolant $p(x) \in \mathbb{P}_9$ (red)
 at ten equally spaced interpolation points ($n = 9$).



Graph of $f(x) = \frac{1}{1+x^2}$ on the interval $[-5, 5]$
 and its Lagrange interpolant $p(x) \in \mathbb{P}_{13}$ (red)
 at fourteen equally spaced interpolation points ($n = 13$).

Conclusion :

- Interpolating a function by a polynomial of high degree is not a good idea.

Alternative :

- Interpolate the function *locally* by polynomials of relatively low degree .

For given integer N let

$$h \equiv \frac{b - a}{N},$$

and partition $[a, b]$ into

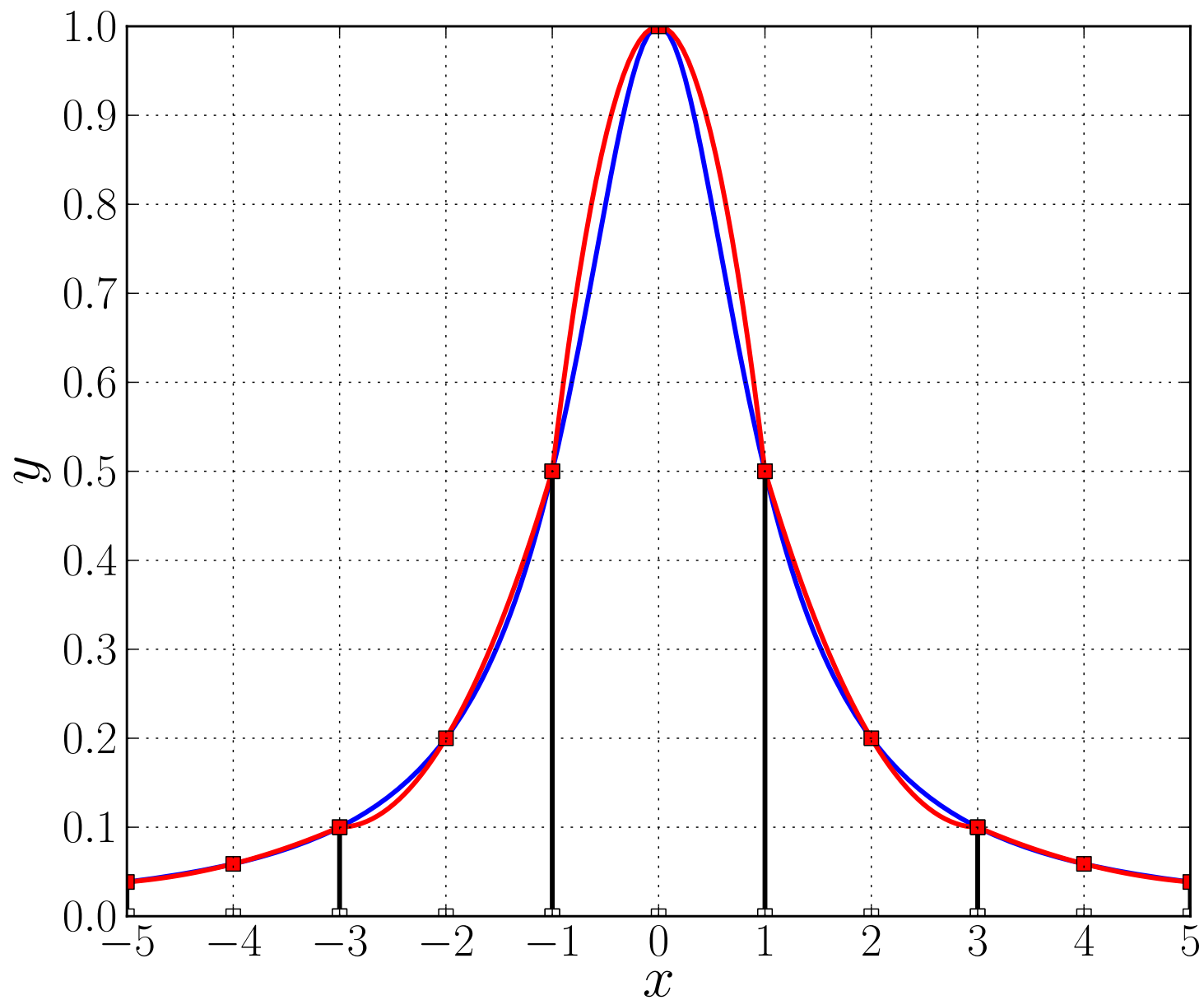
$$a = t_0 < t_1 < \cdots < t_N = b,$$

where

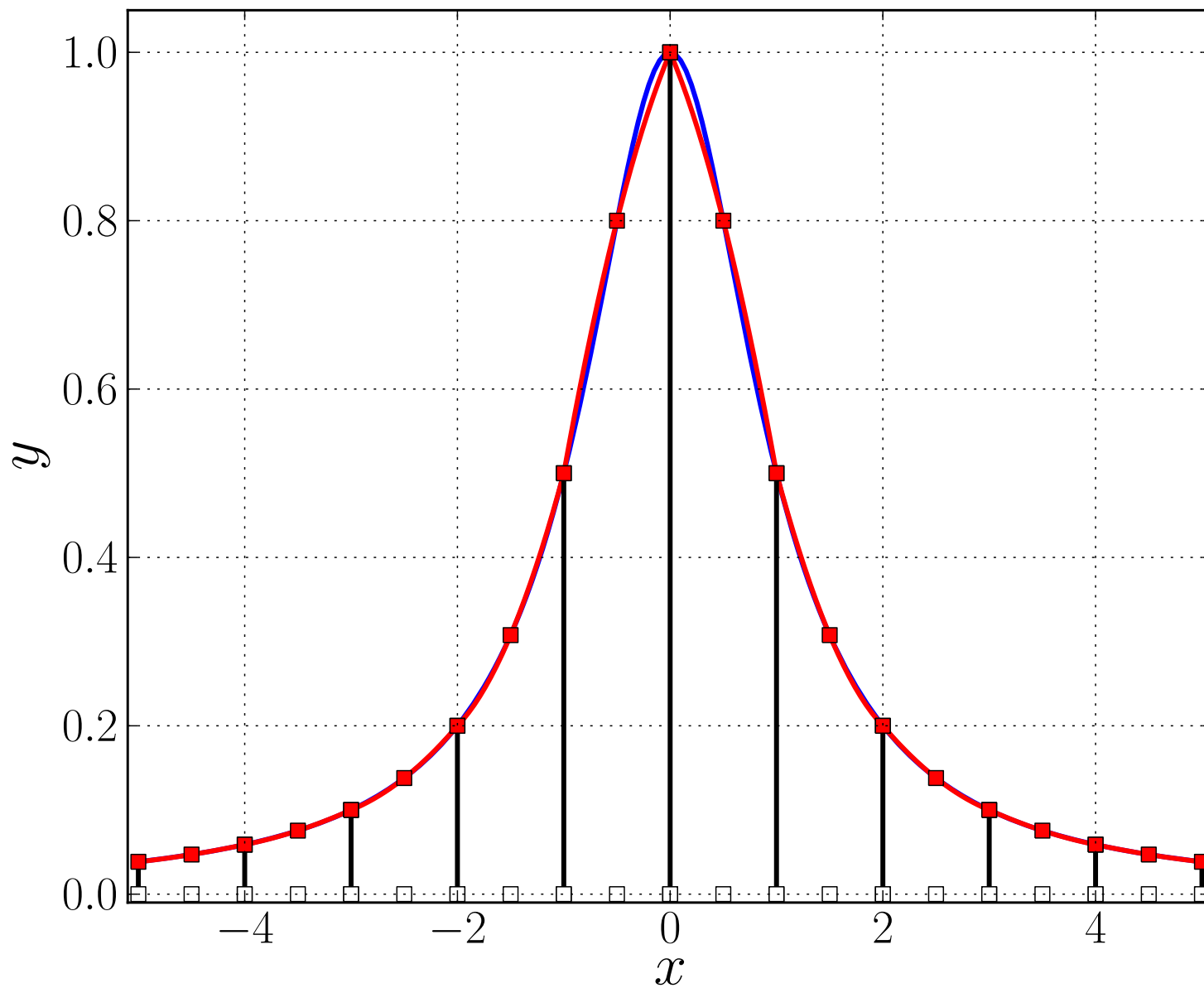
$$t_j = a + jh, \quad j = 0, 1, \dots, N.$$

In each subinterval $[t_{j-1}, t_j]$:

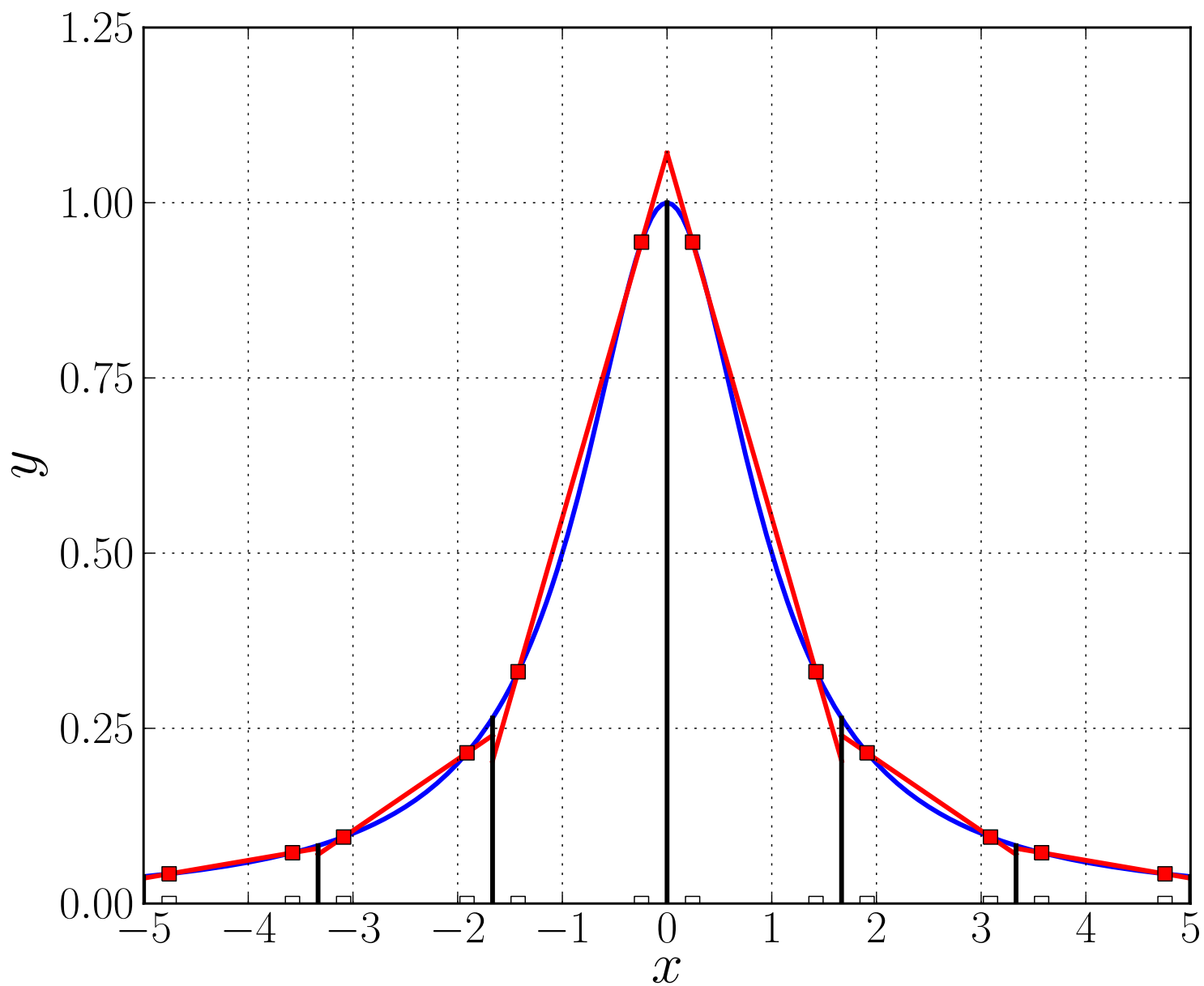
interpolate f by a *local polynomial* $p_j \in \mathbb{P}_n$ at $n+1$ distinct points $\{x_{j,i}\}_{i=0}^n$.



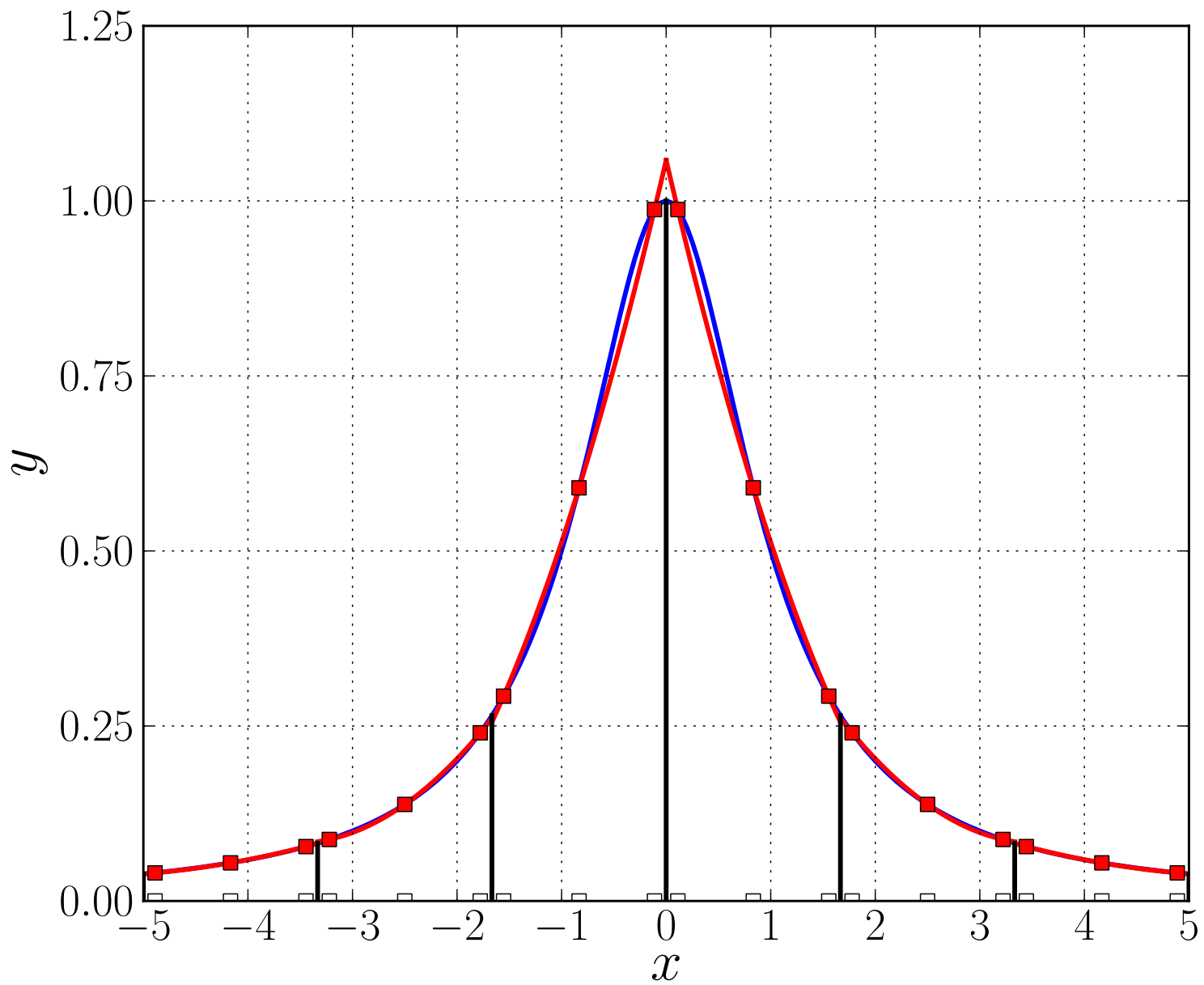
Local polynomial interpolation of $f(x) = \frac{1}{1+x^2}$ at 3 points in 5 intervals.



Local polynomial interpolation of $f(x) = \frac{1}{1+x^2}$ at 3 points in 10 intervals.



Local polynomial interpolation of $f(x) = \frac{1}{1+x^2}$ at 2 Chebyshev points.



Local polynomial interpolation of $f(x) = \frac{1}{1+x^2}$ at 3 Chebyshev points.

By the Lagrange Interpolation Theorem

$$\max_{[t_{j-1}, t_j]} | f(x) - p_j(x) | \leq \frac{1}{(n+1)!} \max_{[t_{j-1}, t_j]} | f^{(n+1)}(x) | \max_{[t_{j-1}, t_j]} | w_{n+1}(x) |$$

where

$$w_{n+1}(x) = \prod_{i=0}^n (x - x_{j,i}) , \quad h \equiv t_j - t_{j-1} = \frac{b-a}{N} .$$

For uniformly spaced interpolation points in $[-1, 1]$ the value of $C_n \equiv \max_{[-1,1]} | w_{n+1}(x) |$ is listed in the Table on Page 183, while the Table on Page 200 also shows these maxima for Chebyshev interpolation points.

A scaling argument shows that for *uniformly spaced local interpolation points*

$$\max_{[t_{j-1}, t_j]} | w_{n+1}(x) | \leq \left(\frac{h}{2}\right)^{n+1} C_n ,$$

while for *local Chebyshev interpolation points* we have

$$\max_{[t_{j-1}, t_j]} | w_{n+1}(x) | \leq \left(\frac{h}{2}\right)^{n+1} 2^{-n} .$$

NOTE:

- *Keeping n fixed*, p_j converges to f as $h \rightarrow 0$, (*i.e.*, as $N \rightarrow \infty$).
- To get more accuracy, *increase* N , keeping the degree n *fixed*.

EXAMPLE: If we approximate

$$f(x) = \cos(x) \quad \text{on} \quad \left[0, \frac{\pi}{2} \right],$$

by local interpolation at 3 *equally spaced local interpolation points*

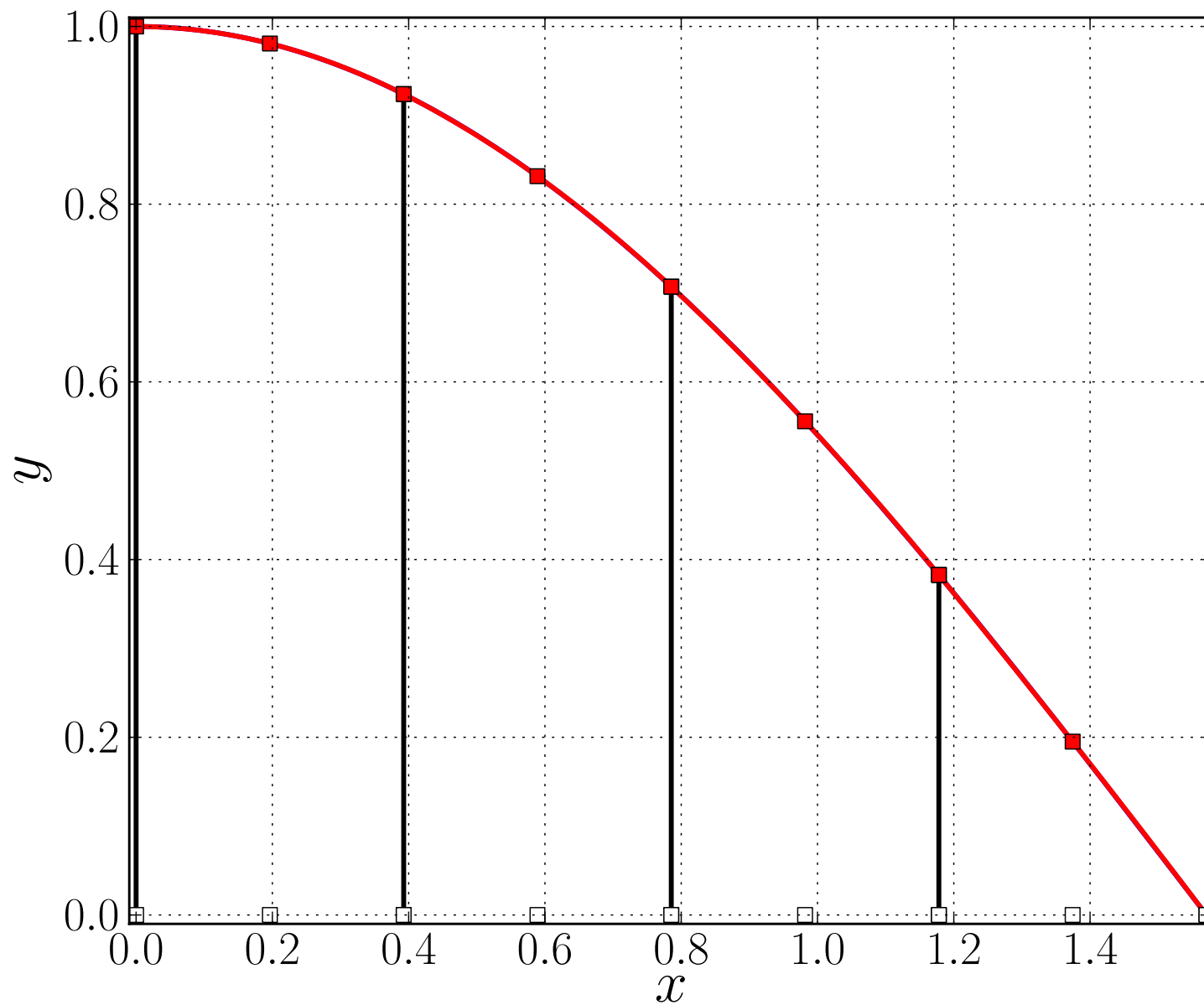
$$x_{j,0} = t_{j-1}, \quad x_{j,1} = \frac{t_{j-1} + t_j}{2}, \quad x_{j,2} = t_j,$$

then $n = 2$, $h = \pi/(2N)$, and, using the Table on Page 183,

$$\max_{[t_{j-1}, t_j]} | f(x) - p_j(x) | \leq \frac{\| f^{(3)} \|_\infty}{3!} \left(\frac{h}{2}\right)^3 C_2 \leq \frac{1}{6} \frac{h^3}{8} 0.3849 .$$

Specifically, if $N = 4$ (four intervals), then $h = \pi/8$, so that

$$\max_{[t_{j-1}, t_j]} | f(x) - p_j(x) | \leq \frac{1}{6} \frac{1}{8} \left(\frac{\pi}{8}\right)^3 0.3849 = 0.000486.$$



Local polynomial interpolation at 3 points in 4 intervals.

EXERCISES:

If we approximate a function $f(x)$ on a given interval by local interpolation with cubic polynomials, then how many intervals of equal size are needed to ensure that the maximum error is less than 10^{-4} ? Answer this question for each of the following cases:

- $f(x) = \sin(x)$ on $[0, 2\pi]$, with arbitrary local interpolation points.
- $f(x) = \sin(x)$ on $[0, 2\pi]$, with equally spaced local points.
- $f(x) = \sin(x)$ on $[0, 2\pi]$, with local Chebyshev interpolation points.
- $f(x) = e^x$ on $[-1, 1]$, with equally spaced local points.
- $f(x) = e^x$ on $[-1, 1]$, with local Chebyshev interpolation points.

Numerical Differentiation.

Numerical differentiation formulas can be derived from local interpolation :

Let $p \in \mathbb{P}_n$ interpolate $f(x)$ at points $\{x_i\}_{i=0}^n$. Thus

$$p(x) = \sum_{i=0}^n f(x_i) \ell_i(x) ,$$

where

$$\ell_i(x) = \prod_{k=0, k \neq i}^n \frac{(x - x_k)}{(x_i - x_k)} .$$

For $m \leq n$ we can approximate $f^{(m)}(x)$ by

$$f^{(m)}(x) \approx p^{(m)}(x) = \sum_{i=0}^n f(x_i) \ell_i^{(m)}(x) .$$

EXAMPLE:

Consider the case

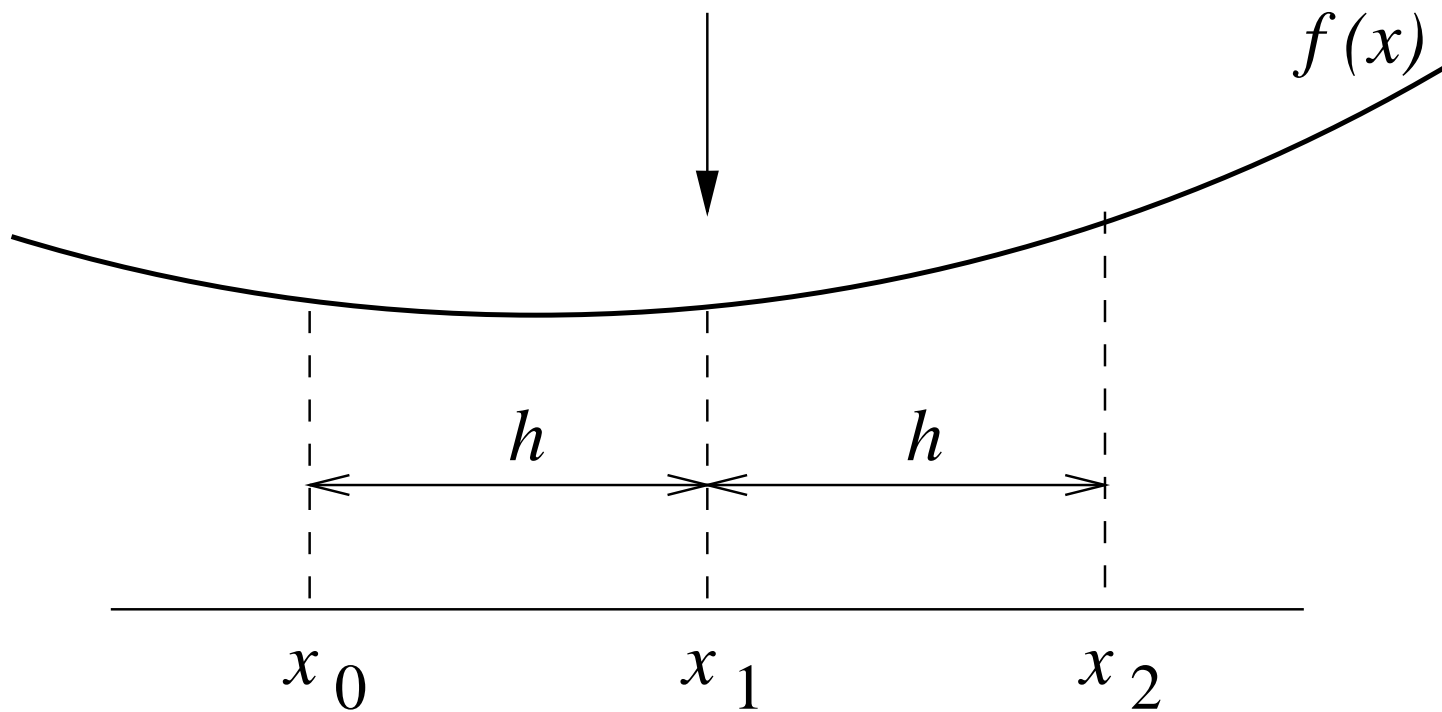
$$n = 2, \quad m = 2, \quad x = 0,$$

for the *reference interval* $[-h, h]$:

$$x_0 = -h, \quad x_1 = 0, \quad x_2 = h.$$

Thus we want to approximate $f''(x_1)$ in terms of

$$f_0, \quad f_1, \quad \text{and} \quad f_2, \quad (f_i \equiv f(x_i)).$$



In this case

$$f''(x_1) \approx p''(x_1) = f_0 \ell_0''(x_1) + f_1 \ell_1''(x_1) + f_2 \ell_2''(x_1) .$$

$$f''(x_1) \approx f_0 \ell_0''(x_1) + f_1 \ell_1''(x_1) + f_2 \ell_2''(x_1) .$$

Here

$$l_0(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} ,$$

so that

$$\ell_0''(x) = \frac{2}{(x_0 - x_1)(x_0 - x_2)} = \frac{1}{h^2} .$$

In particular,

$$\ell_0''(x_1) = \frac{1}{h^2} .$$

Similarly

$$\ell_1''(x_1) = -\frac{2}{h^2} , \quad \ell_2''(x_1) = \frac{1}{h^2} . \quad (\text{Check!})$$

Hence

$$f''(x_1) \approx \frac{f_0 - 2f_1 + f_2}{h^2} .$$

To derive an optimal error bound we use Taylor's Theorem :

$$\begin{aligned}
 & \frac{f_2 - 2f_1 + f_0}{h^2} - f_1'' \\
 = & \frac{1}{h^2} \left(f_1 + hf_1' + \frac{h^2}{2}f_1'' + \frac{h^3}{6}f_1''' + \frac{h^4}{24}f_1''''(\zeta_1) \right. \\
 & \quad - 2f_1 \\
 & \quad \left. + f_1 - hf_1' + \frac{h^2}{2}f_1'' - \frac{h^3}{6}f_1''' + \frac{h^4}{24}f_1''''(\zeta_2) \right) - f_1'' \\
 = & \frac{h^2}{24} \left(f_1''''(\zeta_1) + f_1''''(\zeta_2) \right) = \frac{h^2}{12} f_1''''(\eta) ,
 \end{aligned}$$

where $\eta \in (x_0, x_2)$.

(In the last step we used the Intermediate Value Theorem.)

EXAMPLE: With $n = 4$, $m = 2$, and $x = x_2$, and reference interval

$$x_0 = -2h, \quad x_1 = -h, \quad x_2 = 0, \quad x_3 = h, \quad x_4 = 2h,$$

we have

$$f''(x_2) \approx \sum_{i=0}^4 f_i \ell_i''(x_2).$$

Here

$$\begin{aligned} l_0(x) &= \frac{(x - x_1)(x - x_2)(x - x_3)(x - x_4)}{(x_0 - x_1)(x_0 - x_2)(x_0 - x_3)(x_0 - x_4)} \\ &= \frac{(x - x_1)(x - x_2)(x - x_3)(x - x_4)}{24h^4}. \end{aligned}$$

Differentiating, and setting x equal to x_2 , we find

$$\ell_0''(x_2) = \frac{-1}{12h^2}. \quad (\text{Check!})$$

Similarly

$$\ell_1''(x_2) = \frac{16}{12h^2}, \quad \ell_2''(x_2) = \frac{-30}{12h^2}, \quad \ell_3''(x_2) = \frac{16}{12h^2}, \quad \ell_4''(x_2) = \frac{-1}{12h^2} .$$

(Check!)

Hence we have the *five point finite difference approximation*

$$f''(x_2) \approx \frac{-f_0 + 16f_1 - 30f_2 + 16f_3 - f_4}{12h^2} .$$

By Taylor expansion one can show that the leading error term is

$$\frac{h^4 f^{(6)}(x_2)}{90} . \quad (\text{Check!})$$

We say that the *order of accuracy* of this approximation is equal to 4 .

EXERCISES:

- Derive a formula for the error in the approximation formula

$$f''(0) \approx \frac{f(2h) - 2f(h) + f(0)}{h^2} .$$

What is the order of accuracy?

- Do the same for $f''(0) \approx \frac{f(h) - 2f(0) + f(-h)}{h^2} .$

- Derive the approximation formula

$$f'(0) \approx \frac{-3f(0) + 4f(h) - f(2h)}{2h} .$$

and determine the order of accuracy.

○ Derive the approximation formula

$$f'(0) \approx \frac{-3f(0) + 4f(h) - f(2h)}{2h},$$

and determine the order of accuracy.

SOLUTION: Using the notation $f_0 = f(0)$, $f_1 = f(h)$, and $f_2 = f(2h)$, we have the Lagrange interpolating polynomial

$$p(x) = f_0 \frac{(x-h)(x-2h)}{(0-h)(0-2h)} + f_1 \frac{(x-0)(x-2h)}{(h-0)(h-2h)} + f_2 \frac{(x-0)(x-h)}{(2h-0)(2h-h)},$$

so that

$$p'(x) = f_0 \frac{(x-h) + (x-2h)}{(0-h)(0-2h)} + f_1 \frac{(x-0) + (x-2h)}{(h-0)(h-2h)} + f_2 \frac{(x-0) + (x-h)}{(2h-0)(2h-h)},$$

and

$$p'(0) = f_0 \frac{(0-h) + (0-2h)}{(0-h)(0-2h)} + f_1 \frac{(0-0) + (0-2h)}{(h-0)(h-2h)} + f_2 \frac{(0-0) + (0-h)}{(2h-0)(2h-h)},$$

from which

$$f'(0) \approx p'(0) = f_0 \frac{-3h}{2h^2} + f_1 \frac{-2h}{-h^2} + f_2 \frac{-h}{2h^2} = \frac{-3f_0 + 4f_1 - f_2}{2h}.$$

SOLUTION: continued \dots :

The order of accuracy is determined by Taylor expansion:

$$\begin{aligned} & \frac{-3f_0 + 4f_1 - f_2}{2h} - f'_0 \\ = & \frac{1}{2h} \left(-3f_0 \right. \\ & \quad \left. + 4 \left[f_0 + hf'_0 + \frac{h^2}{2}f''_0 + \frac{h^3}{6}f'''_0 + \dots \right] \right. \\ & \quad \left. - \left[f_0 + 2hf'_0 + \frac{(2h)^2}{2}f''_0 + \frac{(2h)^3}{6}f'''_0 + \dots \right] \right) - f'_0 \\ = & \frac{1}{2h} \left(\frac{4h^3}{6}f'''_0 - \frac{(2h)^3}{6}f'''_0 + \dots \right) \\ = & -\frac{h^2}{3}f'''_0 + \text{higher order terms .} \end{aligned}$$

Thus the formula is of *second order accuracy* .

EXERCISES:

- For the reference interval $[0, 3h]$, give complete details on the derivation of the four weights in the numerical differentiation formula

$$f'(0) \approx \frac{-11f(0) + 18f(h) - 9f(2h) + 2f(3h)}{6h}.$$

Also determine the the leading error term and the order of accuracy. What is the *order of accuracy* of this formula?

- For the reference interval $[-3h/2, 3h/2]$, give complete details on the derivation of the the weights in the numerical differentiation formula

$$f'''(0) \approx \frac{-f(-3h/2) + 3f(-h/2) - 3f(h/2) + f(3h/2)}{h^3}.$$

Also determine the the leading error term and the order of accuracy. What is the *order of accuracy* of this formula?

Best Approximation in the $\|\cdot\|_2$.

Introductory Example : Best approximation in \mathbb{R}^3 .

Recall (from Linear Algebra) :

- A *vector* $\mathbf{x} \in \mathbb{R}^3$ is an ordered set of three numbers, $\mathbf{x} = (x_1, x_2, x_3)^T$.
- We can think of \mathbf{x} as a *point* or an *arrow*.
- The dot product or *inner product* of two vectors \mathbf{x} and \mathbf{y} is defined as

$$\langle \mathbf{x}, \mathbf{y} \rangle \equiv x_1 y_1 + x_2 y_2 + x_3 y_3 .$$

- The length or *norm* of a vector is defined in terms of the inner product :

$$\| \mathbf{x} \|_2 \equiv \langle \mathbf{x}, \mathbf{x} \rangle^{\frac{1}{2}} = \sqrt{x_1^2 + x_2^2 + x_3^2} .$$

- Then $\| \mathbf{x}_1 - \mathbf{x}_2 \|_2$ denotes the distance between \mathbf{x}_1 and \mathbf{x}_2 .

- Two vectors are *perpendicular* if $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = 0$.

Let

$$\mathbf{e}_1 \equiv (1, 0, 0)^T, \quad \mathbf{e}_2 \equiv (0, 1, 0)^T, \quad \text{and} \quad \mathbf{e}_3 \equiv (0, 0, 1)^T.$$

The set $\{\mathbf{e}_k\}_{k=1}^3$ is a *basis* of \mathbb{R}^3 .

This basis is *orthogonal* because

$$\langle \mathbf{e}_1, \mathbf{e}_2 \rangle = \langle \mathbf{e}_1, \mathbf{e}_3 \rangle = \langle \mathbf{e}_2, \mathbf{e}_3 \rangle = 0,$$

and *normal* since

$$\|\mathbf{e}_1\|_2 = \|\mathbf{e}_2\|_2 = \|\mathbf{e}_3\|_2 = 1,$$

i.e., the basis is *orthonormal*.

Let \mathcal{S}_2 denote the x_1, x_2 -plane .

Then

$$\mathcal{S}_2 = \text{Span}\{\mathbf{e}_1, \mathbf{e}_2\} .$$

\mathcal{S}_2 is a 2-dimensional *subspace* of \mathbb{R}^3 .

Suppose we want to find the *best approximation*

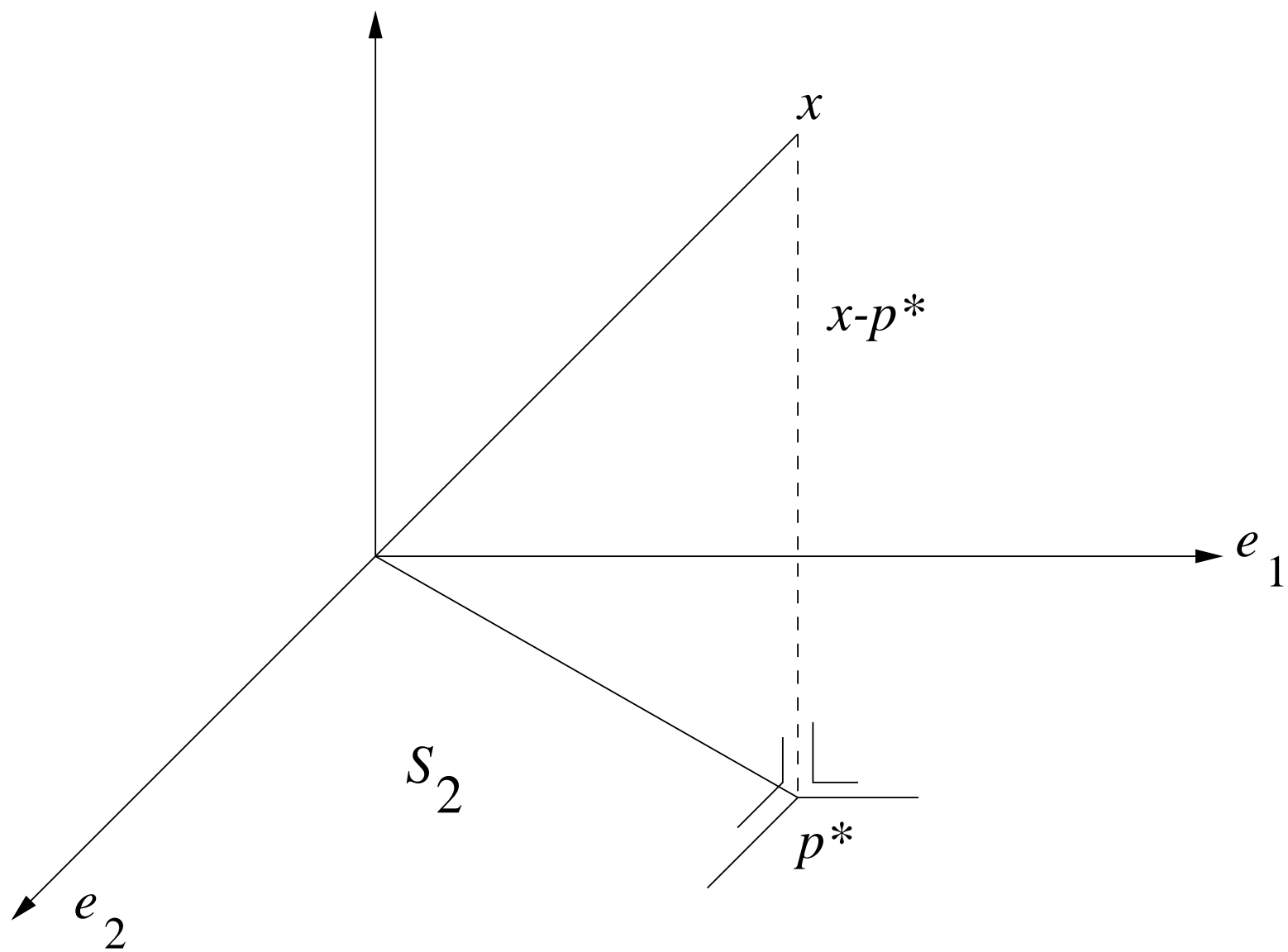
$$\mathbf{p}^* \in \mathcal{S}_2 ,$$

to a given vector $\mathbf{x} \in \mathbb{R}^3$.

Thus we want $\mathbf{p}^* \in \mathcal{S}_2$ that minimizes

$$\| \mathbf{x} - \mathbf{p} \|_2 ,$$

over all $\mathbf{p} \in \mathcal{S}_2$.



Geometrically we see that $\| \mathbf{x} - \mathbf{p} \|_2$ is minimized if and only if

$$(\mathbf{x} - \mathbf{p}) \perp \mathcal{S}_2 ,$$

i.e., if and only if

$$\langle \mathbf{x} - \mathbf{p} , \mathbf{e}_1 \rangle = 0 , \quad \text{and} \quad \langle \mathbf{x} - \mathbf{p} , \mathbf{e}_2 \rangle = 0 ,$$

i.e., if and only if

$$\langle \mathbf{x} , \mathbf{e}_1 \rangle = \langle \mathbf{p} , \mathbf{e}_1 \rangle , \quad \text{and} \quad \langle \mathbf{x} , \mathbf{e}_2 \rangle = \langle \mathbf{p} , \mathbf{e}_2 \rangle .$$

Since $\mathbf{p} \in \mathcal{S}_2$ we have

$$\mathbf{p} = c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 ,$$

for certain constants c_1 and c_2 .

Thus $\| \mathbf{x} - \mathbf{p} \|_2$ is minimized if and only if

$$\langle \mathbf{x}, \mathbf{e}_1 \rangle = \langle c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 , \mathbf{e}_1 \rangle = c_1 ,$$

$$\langle \mathbf{x}, \mathbf{e}_2 \rangle = \langle c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 , \mathbf{e}_2 \rangle = c_2 .$$

Hence

$$\mathbf{p}^* = c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 ,$$

with

$$c_1 = \langle \mathbf{x}, \mathbf{e}_1 \rangle \quad \text{and} \quad c_2 = \langle \mathbf{x}, \mathbf{e}_2 \rangle .$$

Best Approximation in General.

Let \mathbf{X} be a (possibly ∞ -dimensional) real vector space, with an inner product satisfying :

for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbf{X}$, and for all $\alpha \in \mathbb{R}$:

$$i) \quad \langle \mathbf{x}, \mathbf{x} \rangle \geq 0 , \quad \langle \mathbf{x}, \mathbf{x} \rangle = 0 \text{ only if } \mathbf{x} = \mathbf{0} ,$$

$$ii) \quad \langle \alpha \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \alpha \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle ,$$

$$iii) \quad \langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle ,$$

$$iv) \quad \langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle .$$

THEOREM:

Let \mathbf{X} be a vector space with an inner product satisfying the properties above.

Then

$$\| \mathbf{x} \| \equiv \langle \mathbf{x}, \mathbf{x} \rangle^{\frac{1}{2}},$$

defines a norm on \mathbf{X} .

PROOF: We must show that $\| \cdot \|$ satisfies the usual properties :

i) Clearly $\| \mathbf{x} \| \geq 0$, and $\| \mathbf{x} \| = 0$ only if $\mathbf{x} = \mathbf{0}$.

ii) $\| \alpha \mathbf{x} \| = \langle \alpha \mathbf{x}, \alpha \mathbf{x} \rangle^{\frac{1}{2}} = (\alpha^2 \langle \mathbf{x}, \mathbf{x} \rangle)^{\frac{1}{2}} = | \alpha | \| \mathbf{x} \|$.

iii) The triangle inequality is also satisfied :

Let

$$\alpha \equiv \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{y}, \mathbf{y} \rangle} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{y}\|^2}, \quad \text{where } \mathbf{x}, \mathbf{y} \in \mathbf{X}.$$

Then

$$\begin{aligned} 0 &\leq \|\mathbf{x} - \alpha\mathbf{y}\|^2 = \langle \mathbf{x} - \alpha\mathbf{y}, \mathbf{x} - \alpha\mathbf{y} \rangle \\ &= \|\mathbf{x}\|^2 - 2\alpha\langle \mathbf{x}, \mathbf{y} \rangle + \alpha^2\|\mathbf{y}\|^2 \\ &= \|\mathbf{x}\|^2 - 2\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{y}\|^2}\langle \mathbf{x}, \mathbf{y} \rangle + \frac{\langle \mathbf{x}, \mathbf{y} \rangle^2}{\|\mathbf{y}\|^4}\|\mathbf{y}\|^2 \\ &= \|\mathbf{x}\|^2 - \frac{\langle \mathbf{x}, \mathbf{y} \rangle^2}{\|\mathbf{y}\|^2}. \end{aligned}$$

Hence

$$\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \|\mathbf{x}\|^2 \|\mathbf{y}\|^2,$$

or

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\| \quad (\text{Cauchy - Schwartz Inequality}).$$

Now

$$\begin{aligned}\| \mathbf{x} + \mathbf{y} \|^2 &= \langle \mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y} \rangle \\ &= \| \mathbf{x} \|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \| \mathbf{y} \|^2 \\ &\leq \| \mathbf{x} \|^2 + 2 | \langle \mathbf{x}, \mathbf{y} \rangle | + \| \mathbf{y} \|^2 \\ &\leq \| \mathbf{x} \|^2 + 2 \| \mathbf{x} \| \| \mathbf{y} \| + \| \mathbf{y} \|^2 \\ &= (\| \mathbf{x} \| + \| \mathbf{y} \|)^2 .\end{aligned}$$

Hence

$$\| \mathbf{x} + \mathbf{y} \| \leq \| \mathbf{x} \| + \| \mathbf{y} \| . \quad (\textit{Triangle Inequality}) \quad \text{QED!}$$

Suppose $\{\mathbf{e}_k\}_{k=1}^n$ is an orthonormal set of vectors in \mathbf{X} , *i.e.*,

$$\langle \mathbf{e}_l, \mathbf{e}_k \rangle = \begin{cases} 0, & \text{if } l \neq k, \\ 1, & \text{if } l = k. \end{cases}$$

Let $\mathcal{S}_n \subset \mathbf{X}$ be defined by

$$\mathcal{S}_n = \text{Span}\{\mathbf{e}_k\}_{k=1}^n.$$

We want the best approximation $\mathbf{p}^* \in \mathcal{S}_n$ to a given vector $\mathbf{x} \in \mathbf{X}$.

Thus we want to find $\mathbf{p}^* \in \mathcal{S}_n$ that minimizes $\|\mathbf{x} - \mathbf{p}\|$ over all $\mathbf{p} \in \mathcal{S}_n$.

THEOREM:

The best approximation $\mathbf{p}^* \in \mathcal{S}_n$ to $\mathbf{x} \in \mathbf{X}$ is given by

$$\mathbf{p}^* = \sum_{k=1}^n c_k \mathbf{e}_k ,$$

where the *Fourier Coefficients* c_k , ($k = 1, 2, \dots, n$), are given by

$$c_k = \frac{\langle \mathbf{x}, \mathbf{e}_k \rangle}{\langle \mathbf{e}_k, \mathbf{e}_k \rangle} , \quad \text{if the basis is orthogonal ,}$$

and

$$c_k = \langle \mathbf{x}, \mathbf{e}_k \rangle , \quad \text{if the basis is orthonormal .}$$

PROOF: Let

$$F(c_1, c_2, \dots, c_n) \equiv \left\| \mathbf{x} - \sum_{k=1}^n c_k \mathbf{e}_k \right\|^2 .$$

Thus we want to find the $\{c_k\}_{k=1}^n$ that minimize F .

Now

$$\begin{aligned} F(c_1, c_2, \dots, c_n) &= \left\langle \mathbf{x} - \sum_{k=1}^n c_k \mathbf{e}_k, \mathbf{x} - \sum_{k=1}^n c_k \mathbf{e}_k \right\rangle \\ &= \langle \mathbf{x}, \mathbf{x} \rangle - 2 \left\langle \sum_{k=1}^n c_k \mathbf{e}_k, \mathbf{x} \right\rangle + \left\langle \sum_{k=1}^n c_k \mathbf{e}_k, \sum_{k=1}^n c_k \mathbf{e}_k \right\rangle \\ &= \|\mathbf{x}\|^2 - 2 \sum_{k=1}^n c_k \langle \mathbf{x}, \mathbf{e}_k \rangle + \sum_{k=1}^n c_k^2 \langle \mathbf{e}_k, \mathbf{e}_k \rangle . \end{aligned}$$

For F to be minimized we must have

$$\frac{\partial F}{\partial c_\ell} = 0, \quad \ell = 1, 2, \dots, n .$$

We had

$$F(c_1, c_2, \dots, c_n) = \|\mathbf{x}\|^2 - 2 \sum_{k=1}^n c_k \langle \mathbf{x}, \mathbf{e}_k \rangle + \sum_{k=1}^n c_k^2 \langle \mathbf{e}_k, \mathbf{e}_k \rangle .$$

Setting $\frac{\partial F}{\partial c_\ell} = 0$ gives

$$-2\langle \mathbf{x}, \mathbf{e}_\ell \rangle + 2c_\ell \langle \mathbf{e}_\ell, \mathbf{e}_\ell \rangle = 0 .$$

Hence, for $\ell = 1, 2, \dots, n$, we have

$$c_\ell = \frac{\langle \mathbf{x}, \mathbf{e}_\ell \rangle}{\langle \mathbf{e}_\ell, \mathbf{e}_\ell \rangle} ,$$

$$c_\ell = \langle \mathbf{x}, \mathbf{e}_\ell \rangle , \quad \text{if the basis is orthonormal .}$$

QED!

REMARKS:

- The proof uses the fact that \mathbf{X} is an *inner product space* , with norm defined in terms of the inner product.

- In normed vector spaces *without* inner product, *e.g.*,

$$\mathbb{C}[0, 1] \quad \text{with} \quad \|\cdot\|_{\infty} ,$$

it is *more difficult* to find a best approximation.

Gram-Schmidt Orthogonalization.

To *construct*

an orthogonal basis $\{\mathbf{e}_k\}_{k=1}^n$ of a subspace \mathcal{S}_n of \mathbf{X} ,

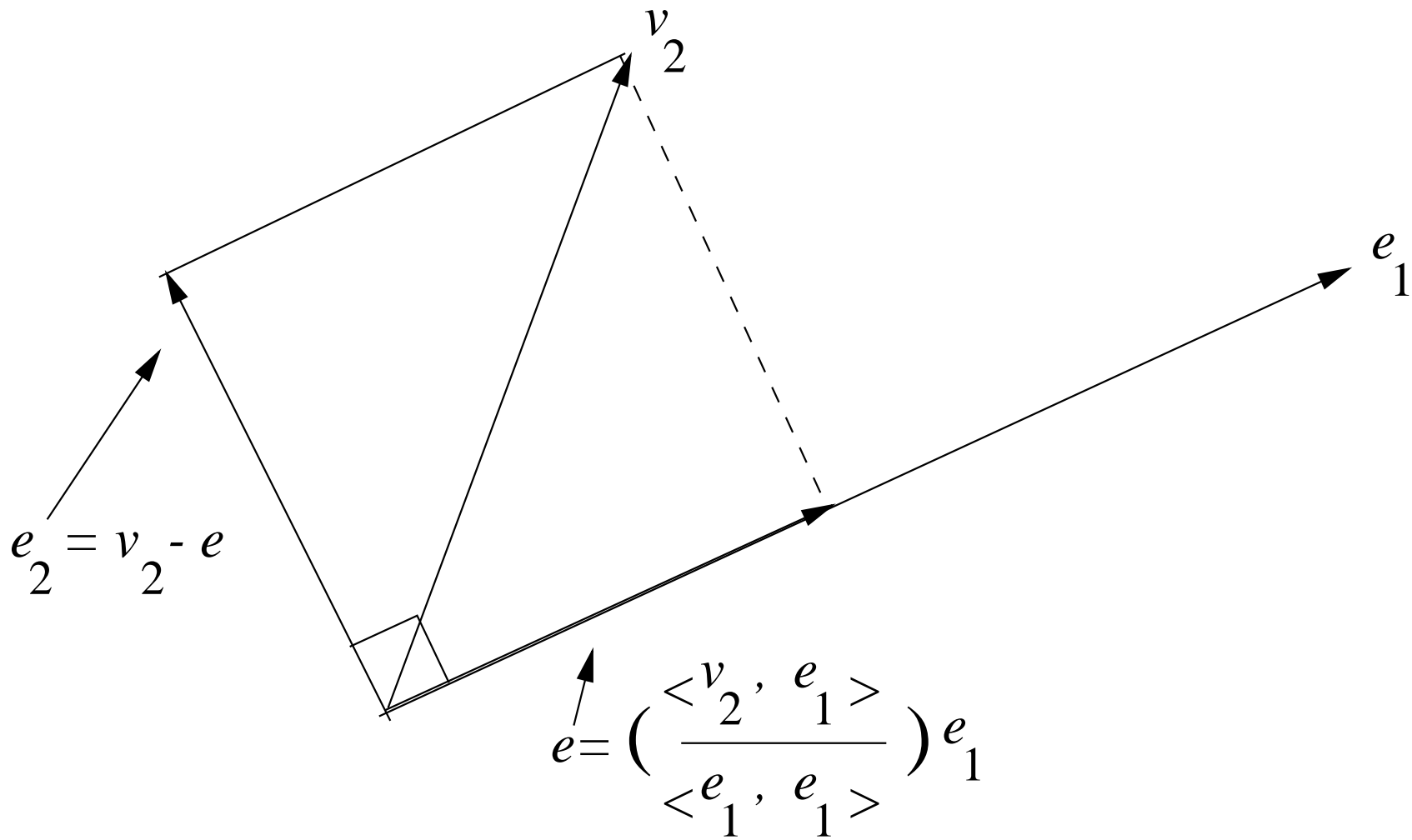
we have the *Gram-Schmidt Orthogonalization Procedure* :

- Take any nonzero $\mathbf{e}_1 \in \mathcal{S}_n$.
- Choose any $\mathbf{v}_2 \in \mathcal{S}_n$ that is linearly independent from \mathbf{e}_1 .
- Set

$$\mathbf{e}_2 = \mathbf{v}_2 - \frac{\langle \mathbf{v}_2, \mathbf{e}_1 \rangle}{\langle \mathbf{e}_1, \mathbf{e}_1 \rangle} \mathbf{e}_1 .$$

Then

$$\langle \mathbf{e}_2, \mathbf{e}_1 \rangle = 0 . \quad (\text{Check!})$$



Inductively, suppose we have mutually orthogonal $\{\mathbf{e}_k\}_{k=1}^{m-1}$, ($m \leq n$).

○ Choose $\mathbf{v}_m \in \mathcal{S}_n$ linearly independent from the $\{\mathbf{e}_k\}_{k=1}^{m-1}$.

○ Set

$$\mathbf{e}_m = \mathbf{v}_m - \sum_{k=1}^{m-1} \frac{\langle \mathbf{v}_m, \mathbf{e}_k \rangle}{\langle \mathbf{e}_k, \mathbf{e}_k \rangle} \mathbf{e}_k .$$

Then

$$\langle \mathbf{e}_m, \mathbf{e}_\ell \rangle = 0 . \quad \ell = 1, 2, \dots, m-1 . \quad (\text{Check!})$$

An orthonormal basis can be obtained by normalizing :

$$\hat{\mathbf{e}}_k = \frac{\mathbf{e}_k}{\|\mathbf{e}_k\|} , \quad k = 1, 2, \dots, n .$$

Best Approximation in a Function Space.

We now apply the general results the special case where

$$\mathbf{X} = \mathbb{C}[-1, 1] ,$$

with inner product

$$\langle f, g \rangle \equiv \int_{-1}^1 f(x) g(x) dx .$$

This definition satisfies all conditions an inner product must satisfy. (Check!)

Hence, from the Theorem it follows that

$$\| f \|_2 \equiv \langle f, f \rangle^{\frac{1}{2}} = \left(\int_{-1}^1 f(x)^2 dx \right)^{\frac{1}{2}} ,$$

is a norm on $\mathbb{C}[-1, 1]$.

Suppose we want to find $p^* \in \mathbb{P}_n$ that best approximates a given function

$$f \in \mathbb{C}[-1, 1] ,$$

in the $\| \cdot \|_2$.

Here \mathbb{P}_n is the $(n + 1)$ -dimensional subspace of $\mathbb{C}[-1, 1]$ consisting of all polynomials of degree less than or equal to n .

By the Theorem we have

$$p^*(x) = \sum_{k=0}^n c_k e_k(x) ,$$

where

$$c_k = \frac{\langle f, e_k \rangle}{\langle e_k, e_k \rangle} = \frac{\int_{-1}^1 f(x) e_k(x) dx}{\int_{-1}^1 e_k^2(x) dx} , \quad k = 0, 1, \dots, n ,$$

and where the $\{e_k\}_{k=0}^n$ denote the first $n + 1$ orthogonal polynomials.

Use the Gram-Schmidt procedure to construct an *orthogonal basis* of \mathbb{P}_n :

(These basis polynomials are called the *Legendre polynomials*.)

Take $e_0(x) \equiv 1$ and $v_1(x) = x$.

Then

$$\frac{\langle v_1, e_0 \rangle}{\langle e_0, e_0 \rangle} = \frac{\int_{-1}^1 x \, dx}{\int_{-1}^1 1^2 \, dx} = 0 .$$

Hence

$$e_1(x) = v_1(x) - 0 \cdot e_0(x) = x .$$

Take $v_2(x) = x^2$.

Then

$$\frac{\langle v_2, e_0 \rangle}{\langle e_0, e_0 \rangle} = \frac{\int_{-1}^1 x^2 dx}{\int_{-1}^1 1^2 dx} = \frac{1}{3},$$

and

$$\frac{\langle v_2, e_1 \rangle}{\langle e_1, e_1 \rangle} = \frac{\int_{-1}^1 x^3 dx}{\int_{-1}^1 x^2 dx} = 0.$$

Hence

$$e_2(x) = v_2(x) - \frac{1}{3} e_0(x) - 0 \cdot e_1(x) = x^2 - \frac{1}{3}.$$

Take $v_3(x) = x^3$. Then

$$\frac{\langle v_3, e_0 \rangle}{\langle e_0, e_0 \rangle} = \frac{\int_{-1}^1 x^3 dx}{\int_{-1}^1 1^2 dx} = 0,$$

and

$$\frac{\langle v_3, e_1 \rangle}{\langle e_1, e_1 \rangle} = \frac{\int_{-1}^1 x^4 dx}{\int_{-1}^1 x^2 dx} = \frac{3}{5},$$

and

$$\frac{\langle v_3, e_2 \rangle}{\langle e_2, e_2 \rangle} = \frac{\int_{-1}^1 x^3 \left(x^2 - \frac{1}{3}\right) dx}{\int_{-1}^1 \left(x^2 - \frac{1}{3}\right)^2 dx} = 0.$$

Hence

$$e_3(x) = v_3(x) - 0 \cdot e_0(x) - \frac{3}{5} e_1(x) - 0 \cdot e_2(x) = x^3 - \frac{3}{5}x.$$

etc.

EXAMPLE:

The polynomial $p^* \in \mathbb{P}_2$ that best approximates

$$f(x) = e^x, \quad \text{on} \quad [-1, 1], \quad \text{in} \quad \|\cdot\|_2,$$

is given by

$$p^*(x) = c_0 e_0(x) + c_1 e_1(x) + c_2 e_2(x),$$

where

$$c_0 = \frac{\langle f, e_0 \rangle}{\langle e_0, e_0 \rangle}, \quad c_1 = \frac{\langle f, e_1 \rangle}{\langle e_1, e_1 \rangle}, \quad c_2 = \frac{\langle f, e_2 \rangle}{\langle e_2, e_2 \rangle}.$$

We find that

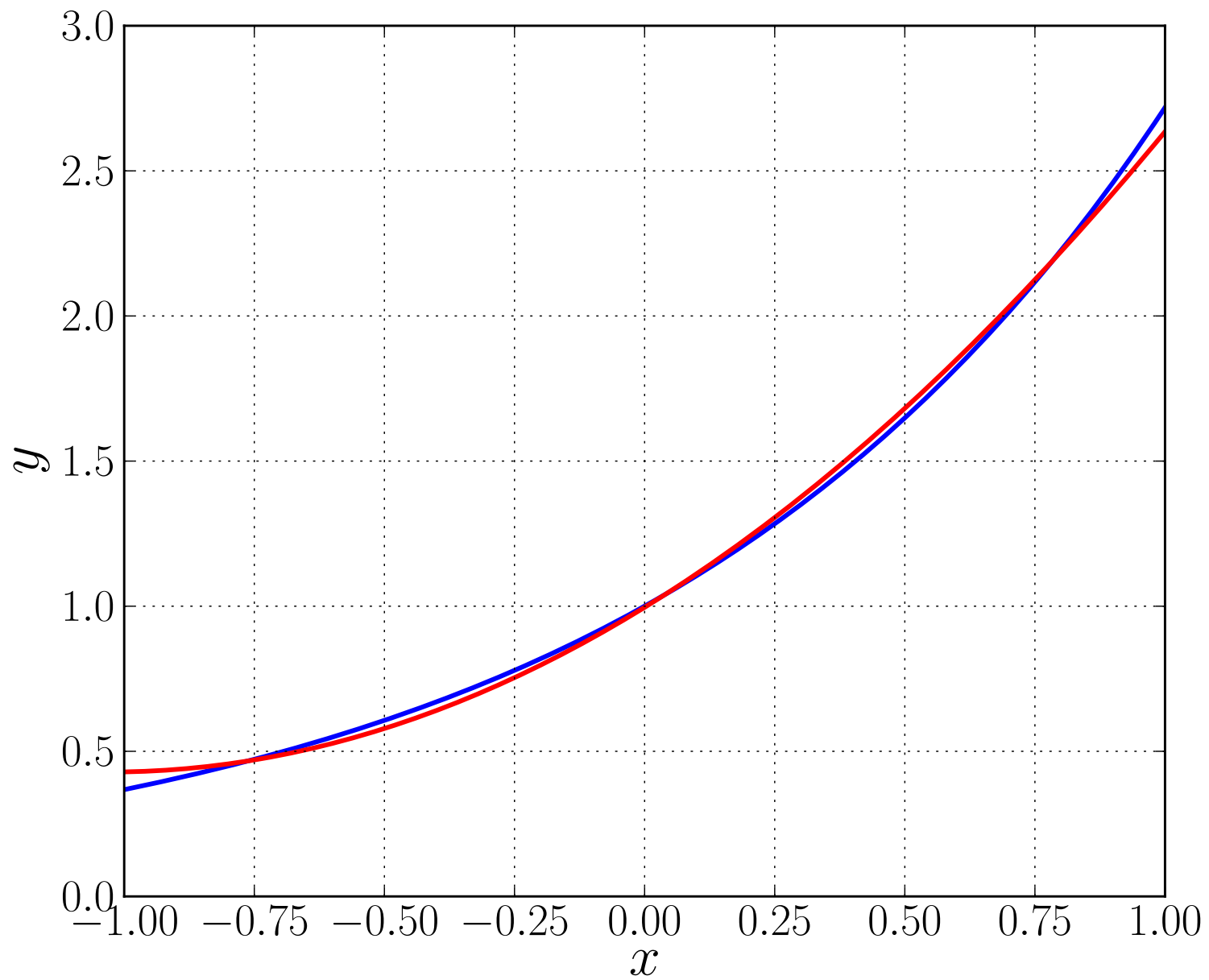
$$c_0 = \frac{\langle f, e_0 \rangle}{\langle e_0, e_0 \rangle} = \frac{\int_{-1}^1 e^x dx}{\int_{-1}^1 1^2 dx} = \frac{1}{2} \left(e - \frac{1}{e} \right) = 1.175 ,$$

$$c_1 = \frac{\langle f, e_1 \rangle}{\langle e_1, e_1 \rangle} = \frac{\int_{-1}^1 e^x x dx}{\int_{-1}^1 x^2 dx} = \frac{3}{2} (x - 1) e^x \Big|_{-1}^1 = 1.103 ,$$

$$c_2 = \frac{\langle f, e_2 \rangle}{\langle e_2, e_2 \rangle} = \frac{\int_{-1}^1 e^x \left(x^2 - \frac{1}{3} \right) dx}{\int_{-1}^1 \left(x^2 - \frac{1}{3} \right)^2 dx} = \frac{45}{8} \left(x^2 - 2x + \frac{5}{3} \right) e^x \Big|_{-1}^1 = 0.536 .$$

Therefore

$$\begin{aligned} p^*(x) &= 1.175 (1) + 1.103 (x) + 0.536 \left(x^2 - \frac{1}{3} \right) \\ &= 0.536 x^2 + 1.103 x + 0.996 . \quad (\text{Check!}) \end{aligned}$$



Best approximation of $f(x) = e^x$ in $[-1, 1]$ by a polynomial $p \in \mathbb{P}_2$.

EXERCISES:

- Use the Gram-Schmidt procedure to construct an orthogonal basis of the polynomial space \mathbb{P}_4 on the interval $[-1, 1]$, by deriving $e_4(x)$, given $e_0(x) = 1$, $e_1(x) = x$, $e_2(x) = x^2 - \frac{1}{3}$, and $e_3(x) = x^3 - \frac{3}{5}x$.
- Use the Gram-Schmidt procedure to construct an orthogonal basis of the linear space $\text{Span}\{1, x^2, x^4\}$ for the interval $[-1, 1]$. Determine the best approximation in the $\|\cdot\|_2$ to $f(x) = x^6$.
- Use the Gram-Schmidt procedure to construct an orthogonal basis of the linear space $\text{Span}\{1, x, x^3\}$ for the interval $[-1, 1]$. Determine the best approximation in the $\|\cdot\|_2$ to $f(x) = x^5$.
- Show that the functions $e_0(x) \equiv 1$, $e_1(x) = \sin(x)$, $e_2(x) = \cos(x)$, $e_3(x) = \sin(2x)$, $e_4(x) = \cos(2x)$, are mutually orthogonal with respect to the inner product $\langle f, g \rangle = \int_0^{2\pi} f(x)g(x) dx$. Also show how one can determine the coefficients c_k , $k = 0, 1, 2, 3, 4$, of the trigonometric polynomial $p(x) = c_0 + c_1 \sin(x) + c_2 \cos(x) + c_3 \sin(2x) + c_4 \cos(2x)$ that minimizes $\int_0^{2\pi} (p(x) - f(x))^2 dx$, when $f(x) = e^x$.

NUMERICAL INTEGRATION

- Many definite integrals, *e.g.*,

$$\int_0^1 e^{(x^2)} dx ,$$

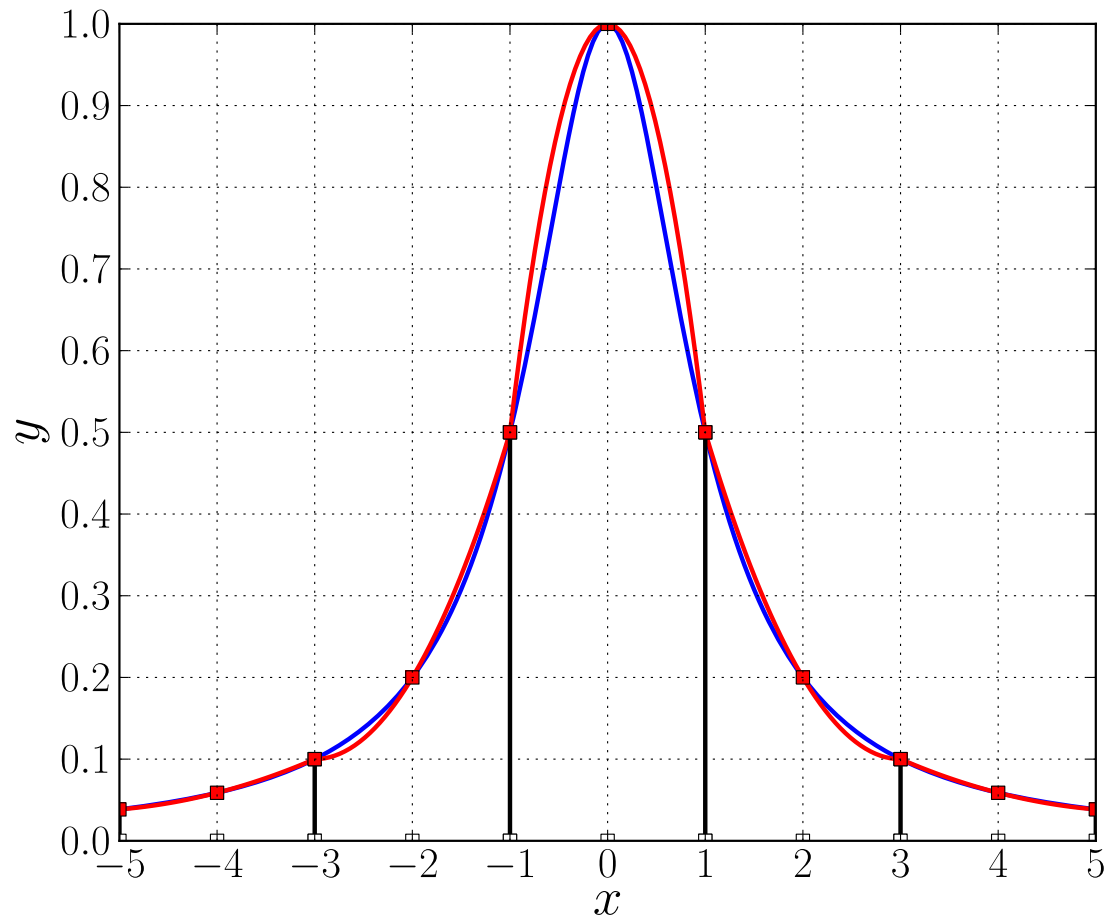
are difficult or impossible to evaluate analytically.

- In such cases we can use *numerical integration*.
- There are many *numerical integration formulas* (or *quadrature formulas*).

Most formulas are based on integrating local interpolating polynomials of f :

$$\int_a^b f(x) dx \approx \sum_{j=1}^N \int_{t_{j-1}}^{t_j} p_j(x) dx ,$$

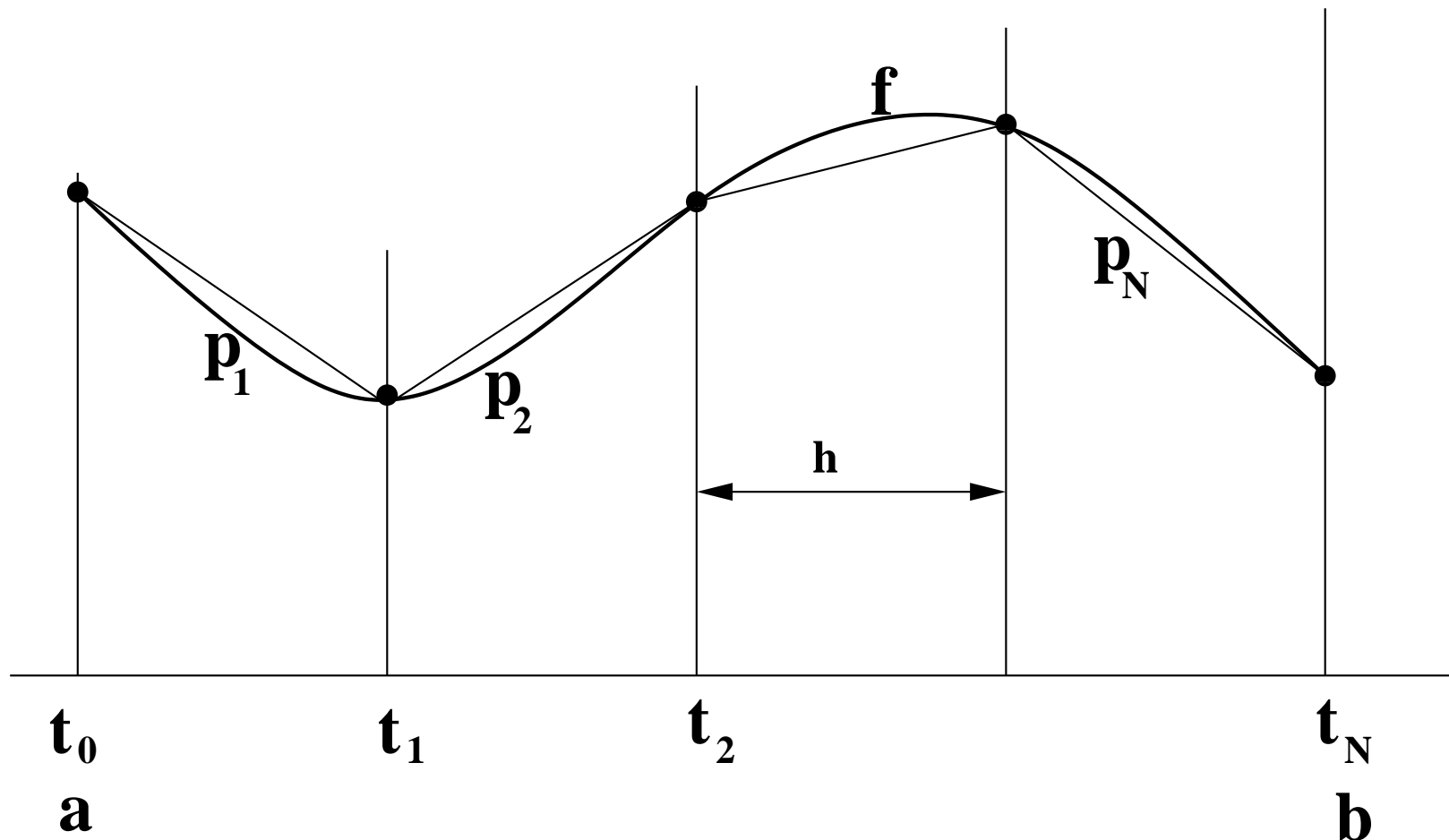
where $p_j \in \mathbb{P}_n$ interpolates f at $n + 1$ points in $[t_{j-1}, t_j]$.



The Trapezoidal Rule.

If $n = 1$, and if $p_j \in \mathbb{P}_1$ interpolates f at t_{j-1} and t_j , then

$$\int_{t_{j-1}}^{t_j} p_j(x) dx = \frac{h}{2} (f_{j-1} + f_j), \quad (\text{local integration formula}).$$



The *composite integration formula* then becomes

$$\begin{aligned}\int_a^b f(x) dx &\approx \sum_{j=1}^N \int_{t_{j-1}}^{t_j} p_j(x) dx \\ &= \sum_{j=1}^N \frac{h}{2} (f_{j-1} + f_j) \\ &= h \left(\frac{1}{2}f_0 + f_1 + \cdots + f_{N-1} + \frac{1}{2}f_N \right),\end{aligned}$$

where $f_j \equiv f(t_j)$.

This is the well-known *Trapezoidal Rule*.

In general

$$p_j(x) = \sum_{i=0}^n f(x_{ji}) \ell_{ji}(x) ,$$

where

$$\ell_{ji}(x) = \prod_{k=0, k \neq i}^n \frac{x - x_{jk}}{x_{ji} - x_{jk}} .$$

Thus we have the approximation

$$\int_{t_{j-1}}^{t_j} f(x) dx \approx \int_{t_{j-1}}^{t_j} p_j(x) dx = \sum_{i=0}^n f(x_{ji}) \int_{t_{j-1}}^{t_j} \ell_{ji}(x) dx .$$

The integrals

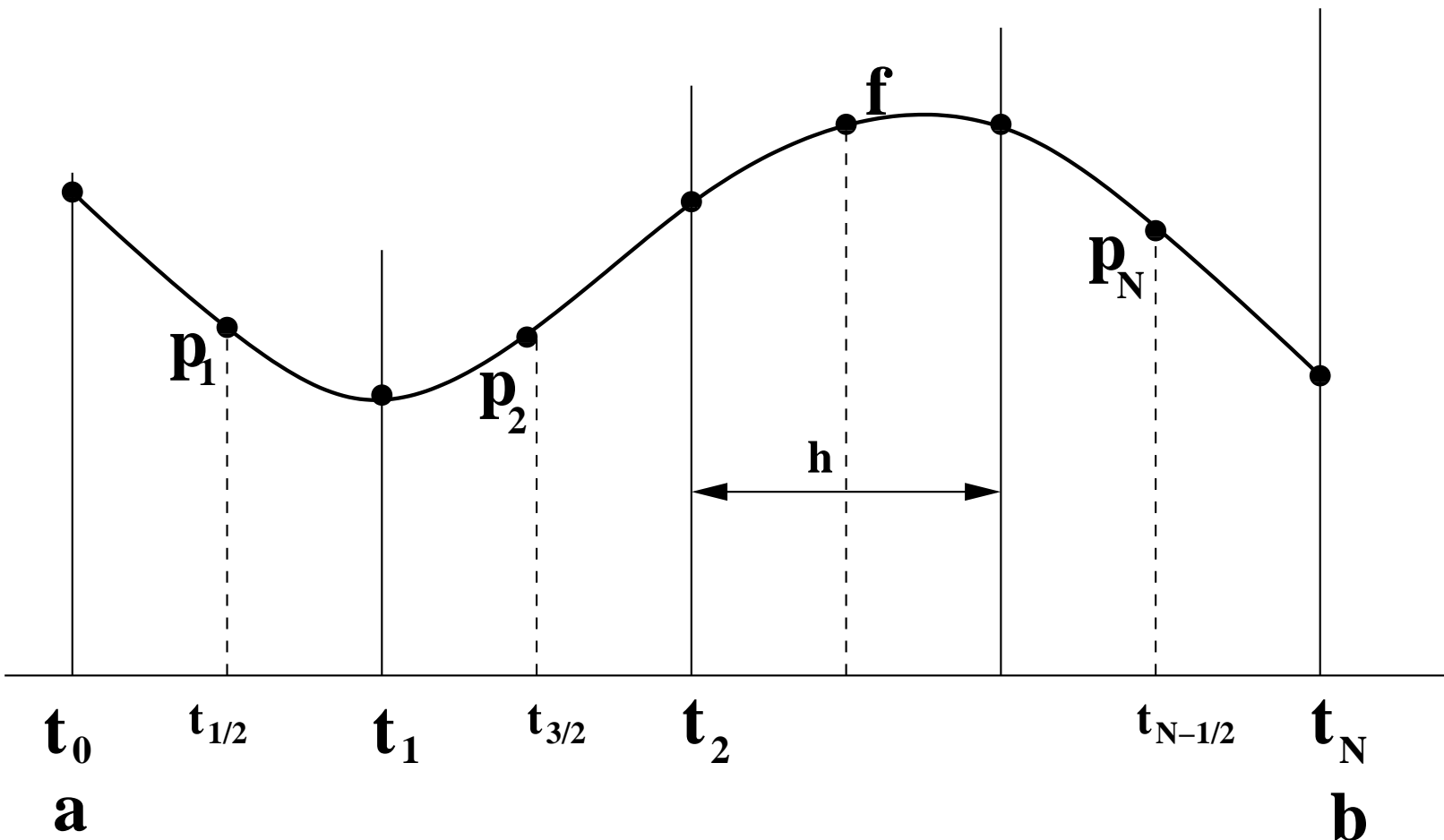
$$\int_{t_{j-1}}^{t_j} \ell_{ji}(x) dx ,$$

are called the *weights* in the local integration formula.

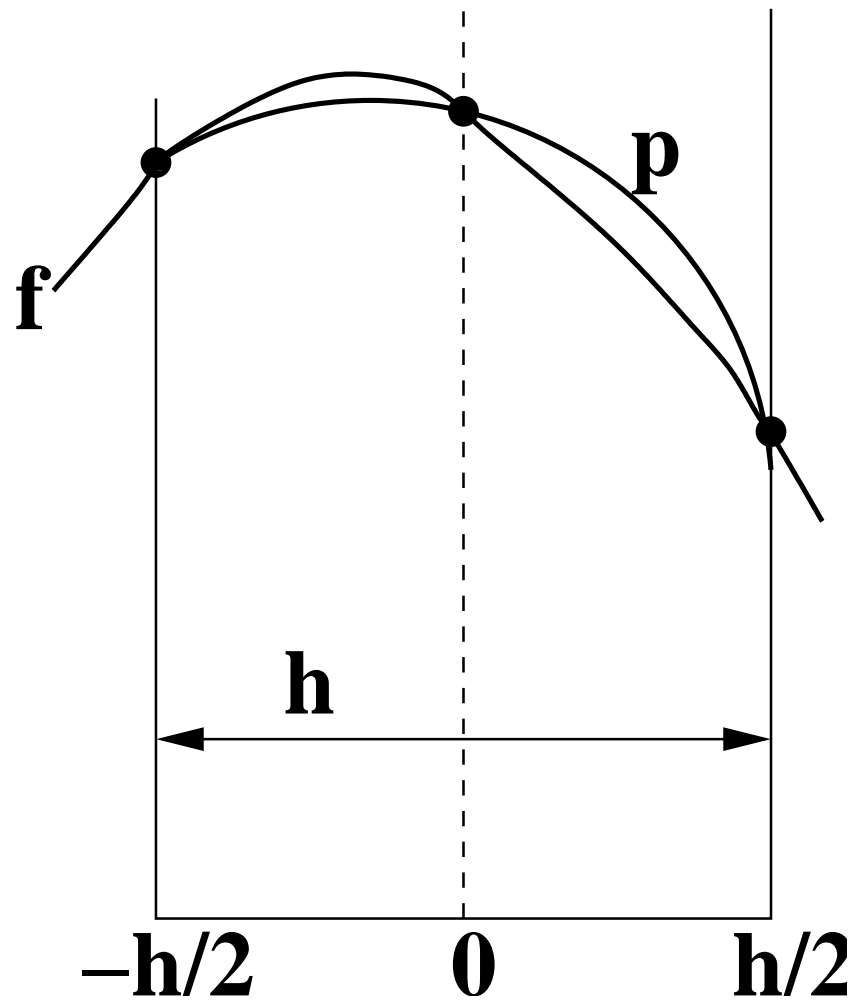
Simpson's Rule.

Let $n = 2$, and in each subinterval $[t_{j-1}, t_j]$ choose the interpolation points

$$t_{j-1}, \quad t_{j-\frac{1}{2}} \equiv \frac{1}{2}(t_{j-1} + t_j), \quad \text{and} \quad t_j.$$



It is convenient to derive the weights for the *reference interval* $[-h/2, h/2]$:



The weights are

$$\int_{-h/2}^{h/2} \frac{(x - 0) (x - \frac{h}{2})}{(-\frac{h}{2} - 0) (-\frac{h}{2} - \frac{h}{2})} dx = \frac{h}{6},$$

$$\int_{-h/2}^{h/2} \frac{(x + \frac{h}{2}) (x - \frac{h}{2})}{(0 + \frac{h}{2}) (0 - \frac{h}{2})} dx = \frac{4h}{6},$$

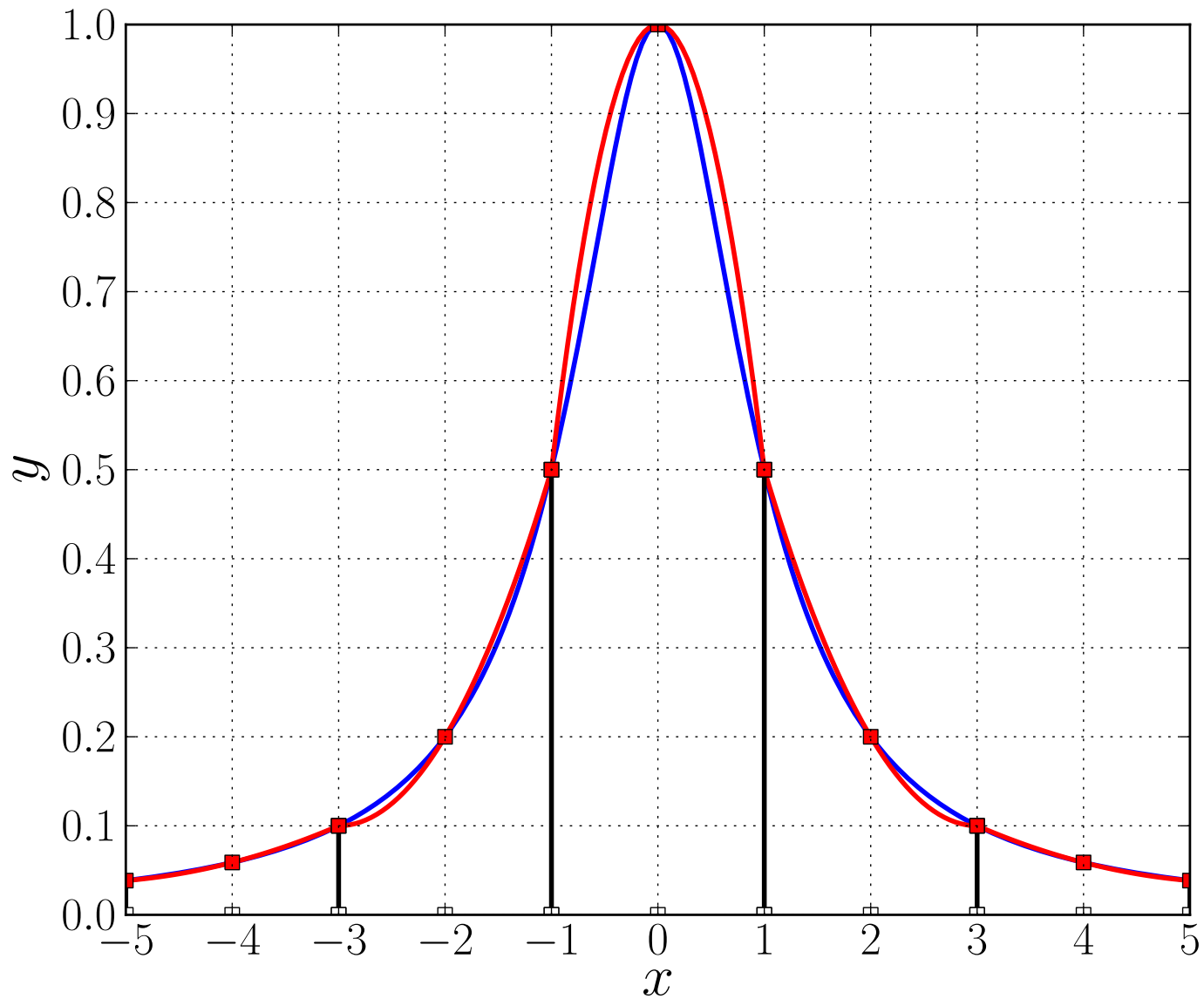
$$\int_{-h/2}^{h/2} \frac{(x + \frac{h}{2}) (x - 0)}{(\frac{h}{2} + \frac{h}{2}) (\frac{h}{2} - 0)} dx = \frac{h}{6}.$$

(Check!)

With uniformly spaced $\{t_j\}_{j=0}^N$, the composite integration formula becomes

$$\int_a^b f(x) dx \approx \sum_{j=1}^N \frac{h}{6} (f_{j-1} + 4f_{j-\frac{1}{2}} + f_j)$$
$$= \frac{h}{6} (f_0 + 4f_{\frac{1}{2}} + 2f_1 + 4f_{1\frac{1}{2}} + 2f_2 + \cdots + 2f_{N-1} + 4f_{N-\frac{1}{2}} + f_N) .$$

This formula is known as *Simpson's Rule*.



The local polynomials (red) in Simpson's Rule for numerically integrating $f(x) = \frac{1}{1+x^2}$ (blue) .

THEOREM:

The error in the composite integration formula based on local polynomial interpolation at $n + 1$ *arbitrary* distinct points satisfies the estimate

$$\left| \int_a^b f(x) dx - \sum_{j=1}^N \int_{t_{j-1}}^{t_j} p_j(x) dx \right| \leq \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} h^{n+1} (b-a),$$

where

$$h = \frac{b-a}{N}.$$

PROOF:

The *local error* is

$$\begin{aligned} \left| \int_{t_{j-1}}^{t_j} f(x) \, dx - \int_{t_{j-1}}^{t_j} p_j(x) \, dx \right| &= \left| \int_{t_{j-1}}^{t_j} f(x) - p_j(x) \, dx \right| \\ &= \left| \int_{t_{j-1}}^{t_j} \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \prod_{i=0}^n (x - x_{ji}) \, dx \right| \\ &\leq |t_j - t_{j-1}| \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} h^{n+1} \\ &= \frac{\|f^{(n+1)}\|_\infty h^{n+2}}{(n+1)!} . \end{aligned}$$

The *error in the composite formula*, using $n + 1$ *arbitrary* distinct points in each subinterval, is now easily determined as follows:

$$\begin{aligned}
 \left| \int_a^b f(x) \, dx - \sum_{j=1}^N \int_{t_{j-1}}^{t_j} p_j(x) \, dx \right| &= \left| \sum_{j=1}^N \int_{t_{j-1}}^{t_j} f(x) - p_j(x) \, dx \right| \\
 &\leq \sum_{j=1}^N \frac{\| f^{(n+1)} \|_{\infty} h^{n+2}}{(n+1)!} \\
 &= \frac{N \| f^{(n+1)} \|_{\infty} h^{n+2}}{(n+1)!} \\
 &= \frac{\| f^{(n+1)} \|_{\infty} h^{n+1} (b-a)}{(n+1)!} .
 \end{aligned}$$

The last step uses the fact that

$$h = \frac{b-a}{N} , \quad \text{i.e. ,} \quad N = \frac{b-a}{h} . \quad \text{QED!}$$

REMARKS:

- The estimate

$$\max_{x \in [t_{j-1}, t_j]} \left| \prod_{i=0}^n (x - x_{ji}) \right| \leq h^{n+1},$$

can be improved depending on the choice of the interpolation points.

- The Theorem states that *order of accuracy* is at least $n + 1$.
- We say that the method is $\mathcal{O}(h^{n+1})$.
- The actual order may be higher.
- Very high order of accuracy is possible for special interpolation points.

EXAMPLES:

For the Trapezoidal Rule ($n = 1$), the Theorem gives the error bound

$$\frac{h^2}{2} \|f''\|_{\infty} (b - a) .$$

Indeed the Trapezoidal Rule is $\mathcal{O}(h^2)$.

For Simpsons Rule ($n = 2$), the Theorem gives

$$\frac{h^3}{6} \|f^{(3)}\|_{\infty} (b - a) .$$

The actual order of Simpson's Rule is higher, namely, $\mathcal{O}(h^4)$.

EXAMPLE:

Taylor expand for the precise local error in Simpson's Rule :

$$\begin{aligned}
 & \int_{-h/2}^{h/2} f(x) \, dx - \frac{h}{6} \left(f\left(-\frac{h}{2}\right) + 4f(0) + f\left(\frac{h}{2}\right) \right) \\
 &= \int_{-h/2}^{h/2} f_0 + x f'_0 + \frac{x^2}{2} f''_0 + \frac{x^3}{6} f'''_0 + \frac{x^4}{24} f''''_0 + \dots \, dx \\
 &- \frac{h}{6} \left(f_0 - \left(\frac{h}{2}\right) f'_0 + \frac{1}{2} \left(\frac{h}{2}\right)^2 f''_0 - \frac{1}{6} \left(\frac{h}{2}\right)^3 f'''_0 + \frac{1}{24} \left(\frac{h}{2}\right)^4 f''''_0 + \dots \right. \\
 &\quad + 4f_0 \\
 &\quad \left. + f_0 + \left(\frac{h}{2}\right) f'_0 + \frac{1}{2} \left(\frac{h}{2}\right)^2 f''_0 + \frac{1}{6} \left(\frac{h}{2}\right)^3 f'''_0 + \frac{1}{24} \left(\frac{h}{2}\right)^4 f''''_0 + \dots \right)
 \end{aligned}$$

where $f_0 \equiv f(0)$, *etc.*

$$\begin{aligned}
&= \left(x f_0 + \frac{x^2}{2} f_0' + \frac{x^3}{6} f_0'' + \frac{x^4}{24} f_0''' + \frac{x^5}{120} f_0'''' + \dots \right) \Big|_{-h/2}^{h/2} \\
&\quad - \frac{h}{6} \left(6 f_0 + \frac{h^2}{4} f_0'' + \frac{h^4}{192} f_0'''' + \dots \right) \\
&= \left(h f_0 + \frac{(h/2)^3}{3} f_0'' + \frac{(h/2)^5}{60} f_0'''' + \dots \right) \\
&\quad - \left(h f_0 + \frac{h^3}{24} f_0'' + \frac{h^5}{1152} f_0'''' + \dots \right) \\
&= -\frac{h^5}{2880} f_0'''' + \text{higher order terms.}
\end{aligned}$$

Thus the leading error term of the composite Simpson's Rule is bounded by

$$\frac{h^4}{2880} \| f'''' \|_{\infty} (b - a) .$$

EXERCISE:

- The *Local Midpoint Rule*, for numerically integrating a function $f(x)$ over the reference interval $[-h/2, h/2]$, is given by

$$\int_{-h/2}^{h/2} f(x) dx \approx hf(0) .$$

Use Taylor expansion to determine the error in this local formula.

Write down the formula for the *Composite Midpoint Rule* for integrating $f(x)$ over a general interval $[a, b]$.

Derive an error formula for the composite formula.

How big must N be for the global error to be less than 10^{-6} , when integrating $f(x) = \sin(x)$ over the interval $[0, 1]$?

EXERCISE:

- The *local Trapezoidal Rule* for the reference interval $[-h/2, h/2]$ is

$$\int_{-h/2}^{h/2} f(x) dx \approx \frac{h}{2} \left[f(-h/2) + f(h/2) \right].$$

Use Taylor expansions to derive the local error formula.

Let $h = (b - a)/N$ and $x_k = a + k h$, for $k = 0, 1, 2, 3, \dots, N$.

Then the *composite Trapezoidal Rule* is given by

$$\int_a^b f(x) dx \approx \frac{h}{2} \left[f(x_0) + 2f(x_1) + 2f(x_2) + \dots + 2f(x_{N-1}) + f(x_N) \right].$$

Based on the local error, derive an upper bound on the global error.

How big must N be for the global error to be less than 10^{-6} , when integrating $f(x) = \sin(x)$ over the interval $[0, 1]$?

THE GAUSS QUADRATURE THEOREM:

If in each subinterval $[t_{j-1}, t_j]$ the interpolation points $\{x_{ji}\}_{i=0}^n$ are taken as

the zeroes of the $(n + 1)$ st orthogonal polynomial $e_{n+1}(x)$,

(relative to $[t_{j-1}, t_j]$) ,

then the the composite integration formula is $\mathcal{O}(h^{2n+2})$.

REMARKS:

- Such integration formulas are known as *Gauss Quadrature Formulas*.
- The points $\{x_{ji}\}_{i=0}^n$ are the *Gauss points* .
- The order improves from $\mathcal{O}(h^{n+1})$ to $\mathcal{O}(h^{2n+2})$ for Gauss points.

EXAMPLE: The case $n = 1$:

Relative to the interval $[-1, 1]$, the second degree orthogonal polynomial is

$$e_2(x) = x^2 - \frac{1}{3} .$$

The two Gauss points relative to $[-1, 1]$ are the zeroes of $e_2(x)$, *i.e.*,

$$x_0 = -\frac{\sqrt{3}}{3} \quad \text{and} \quad x_1 = \frac{\sqrt{3}}{3} .$$

Relative to $[t_{j-1}, t_j]$ the Gauss points are

$$x_{j,0} = t_{j-\frac{1}{2}} - h\frac{\sqrt{3}}{6} , \quad \text{and} \quad x_{j,1} = t_{j-\frac{1}{2}} + h\frac{\sqrt{3}}{6} ,$$

where $t_{j-\frac{1}{2}} \equiv \frac{1}{2}(t_{j-1} + t_j)$.

Relative to the *reference interval* $I_h \equiv [-h/2, h/2]$ the Gauss points are

$$x_0 = -\frac{h\sqrt{3}}{6} \quad \text{and} \quad x_1 = \frac{h\sqrt{3}}{6},$$

with interpolating polynomial

$$p(x) = f(x_0) l_0(x) + f(x_1) l_1(x),$$

where

$$l_0(x) = \frac{x - x_1}{x_0 - x_1} = \frac{x - h\sqrt{3}/6}{-h\sqrt{3}/3},$$

and

$$l_1(x) = \frac{x - x_0}{x_1 - x_0} = \frac{x + h\sqrt{3}/6}{h\sqrt{3}/3}.$$

The local integration formula is

$$\int_{-h/2}^{h/2} f(x) dx \approx f(x_0) \int_{-h/2}^{h/2} \ell_0(x) dx + f(x_1) \int_{-h/2}^{h/2} \ell_1(x) dx .$$

Integrating $\ell_0(x)$ and $\ell_1(x)$, we find

$$\int_{-h/2}^{h/2} f(x) dx \approx \frac{h}{2} f(x_0) + \frac{h}{2} f(x_1) .$$

Hence the composite two point Gauss quadrature formula is given by

$$\int_a^b f(x) dx \approx \frac{h}{2} \sum_{j=1}^N [f(x_{j,0}) + f(x_{j,1})] .$$

By the Theorem this integration formula is $\mathcal{O}(h^4)$.

PROOF (of the Gauss Quadrature Theorem.)

The local error for the reference interval $I_h \equiv [-h/2, h/2]$ is

$$\int_{-h/2}^{h/2} f(x) - p(x) \, dx ,$$

where $p \in \mathbb{P}_n$ interpolates $f(x)$ at Gauss points $\{x_i\}_{i=0}^n$ (relative to I_h).

By the Lagrange Interpolation Theorem

$$\begin{aligned} \int_{-h/2}^{h/2} f(x) - p(x) \, dx &= \int_{-h/2}^{h/2} \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \prod_{i=0}^n (x - x_i) \, dx \\ &= \int_{-h/2}^{h/2} c(x) e_{n+1}(x) \, dx , \end{aligned}$$

where

$$c(x) \equiv \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \quad \text{and} \quad e_{n+1}(x) = \prod_{i=0}^n (x - x_i) .$$

Note that e_{n+1} is the $(n+1)$ *st* orthogonal polynomial (relative to I_h).

FACT: If $f(x)$ is very smooth then $c(x)$ has $n + 1$ continuous derivatives.

Thus we can Taylor expand :

$$c(x) = \sum_{k=0}^n \frac{x^k}{k!} c^{(k)}(0) + \frac{x^{n+1}}{(n+1)!} c^{(n+1)}(\eta(x)) .$$

Call the remainder $r(x)$ and use the fact that each summation term is in \mathbb{P}_n :

$$c(x) = \sum_{k=0}^n c_k e_k(x) + r(x) ,$$

where e_k is the k^{th} orthogonal polynomial relative to I_h .

(Recall that the $\{e_k\}_{k=0}^n$ form an orthogonal basis of \mathbb{P}_n .)

We have

$$\int_{-h/2}^{h/2} f(x) - p(x) dx = \int_{-h/2}^{h/2} c(x) e_{n+1}(x) dx ,$$

and

$$c(x) = \sum_{k=0}^n c_k e_k(x) + r(x) .$$

It follows that

$$\begin{aligned} \left| \int_{-h/2}^{h/2} f(x) - p(x) dx \right| &= \left| \int_{-h/2}^{h/2} \left[\sum_{k=0}^n c_k e_k(x) + r(x) \right] e_{n+1}(x) dx \right| \\ &= \left| \sum_{k=0}^n c_k \int_{-h/2}^{h/2} e_k(x) e_{n+1}(x) dx + \int_{-h/2}^{h/2} r(x) e_{n+1}(x) dx \right| . \end{aligned}$$

Note that all terms in the summation term are zero by orthogonality, so that

$$\begin{aligned}
\left| \int_{-h/2}^{h/2} f(x) - p(x) dx \right| &= \left| \int_{-h/2}^{h/2} r(x) e_{n+1}(x) dx \right| \\
&= \left| \int_{-h/2}^{h/2} \frac{x^{n+1}}{(n+1)!} c^{(n+1)}(\eta(x)) \prod_{i=0}^n (x - x_i) dx \right| \\
&\leq h \max_{x \in I_h} \left| \frac{x^{n+1}}{(n+1)!} c^{(n+1)}(\eta(x)) \prod_{i=0}^n (x - x_i) \right| \\
&\leq h \frac{(h/2)^{n+1}}{(n+1)!} \max_{x \in I_h} |c^{(n+1)}(x)| h^{n+1} \\
&= \frac{h^{2n+3}}{2^{n+1}(n+1)!} \max_{x \in I_h} |c^{(n+1)}(x)|.
\end{aligned}$$

Hence the local integration formula is $\mathcal{O}(h^{2n+3})$.

As before, this implies that the composite formula is $\mathcal{O}(h^{2n+2})$. QED!

EXERCISE:

- Give complete details on the derivation of the the *local* 3-point Gauss integration formula. Also write down the *composite* 3-point Gauss formula for integrating a function $f(x)$ over a general interval $[a, b]$.
- Are the following True or False for any sufficiently smooth $f(x)$?
 - The order of accuracy of a general *composite* $(n + 1)$ -point integration formula for $f(x)$ is at least $\mathcal{O}(h^{n+1})$.
 - The order of accuracy of the *composite* $(n+1)$ -point Gauss formula for integrating $f(x)$ is $\mathcal{O}(h^{2n+4})$.
 - The order of accuracy of the composite 2-Point Gauss formula is the same as the order of accuracy of the composite Simpson formula.

DISCRETE LEAST SQUARES APPROXIMATION

We have solved the *continuous least squares problem* :

Given $f(x)$ on $[-1, 1]$, find a polynomial $p(x) \in \mathbb{P}_n$ that minimizes

$$\| p - f \|_2^2 \equiv \int_{-1}^1 [p(x) - f(x)]^2 dx .$$

Next we solve the *discrete least squares problem* :

Given a set of discrete data points

$$\{ (x_i, y_i) \}_{i=1}^N ,$$

find $p \in \mathbb{P}_n$ such that

$$e_L \equiv \sum_{i=1}^N [p(x_i) - y_i]^2$$

is minimized.

More generally, find

$$p(x) = \sum_{i=0}^n a_i \phi_i(x) ,$$

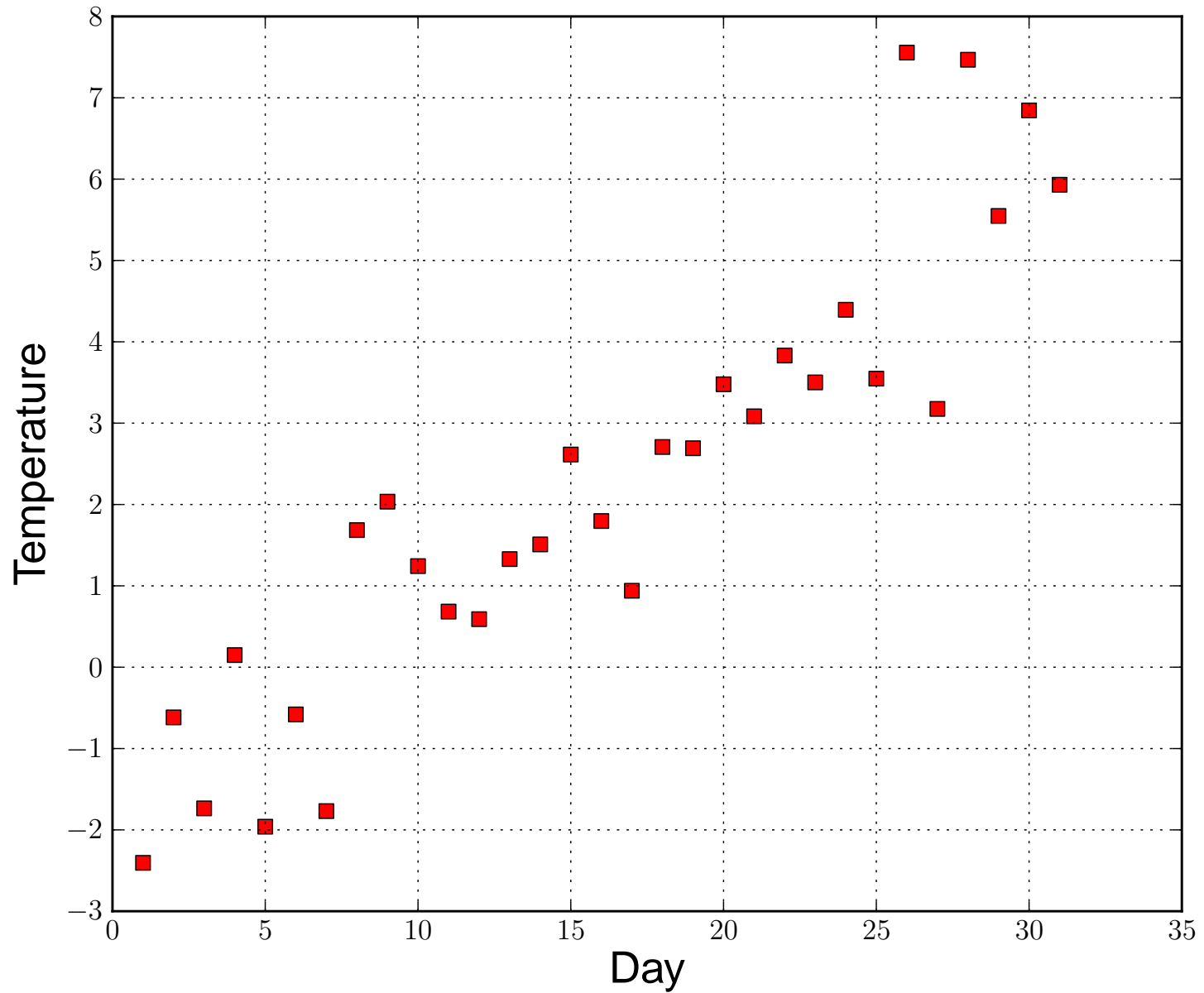
(not necessarily a polynomial), such that e_L is minimized.

Linear Least Squares

- Suppose we have data on the *daily high temperature* in March.
- For each day we compute the *average* high temperature.
- Each average is taken over *a number of years* .
- The (fictitious) data are given in the Table below.

1	-2.4	2	-0.6	3	-1.7	4	0.1	5	-2.0	6	-0.6	7	-1.8
8	1.7	9	2.0	10	1.2	11	0.7	12	0.6	13	1.3	14	1.5
15	2.6	16	1.8	17	0.9	18	2.7	19	2.7	20	3.5	21	3.1
22	3.8	23	3.5	24	4.4	25	3.5	26	7.6	27	3.2	28	7.5
29	5.5	30	6.8	31	5.9								

Average daily high temperature in Montreal in March



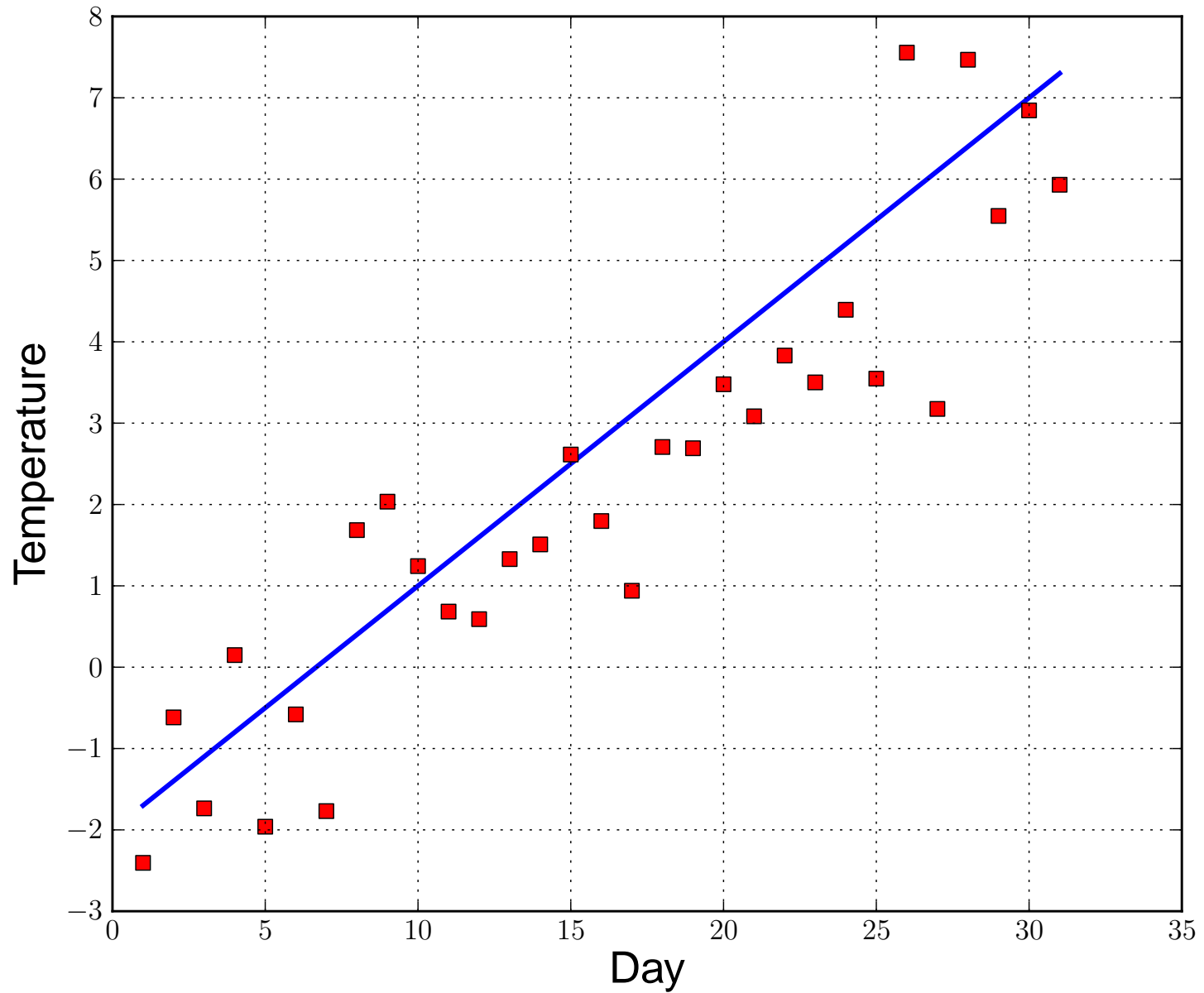
Average daily high temperature in Montreal in March

Suppose that:

- We believe these temperatures basically increase *linearly* .
- Thus we believe in a relation

$$T_k = c_1 + c_2 k , \quad k = 1, 2, \dots, 31 .$$

- The *deviations* from linearity come from *random influences* .
- These random influences can be due to *many factors* .
- We want to determine "*the best*" linear approximation.



Average daily high temperatures, with a *linear approximation* .

- There are many ways to determine such a linear approximation.
- Often used is the *least squares method*.
- This method determines c_1 and c_2 that *minimize*

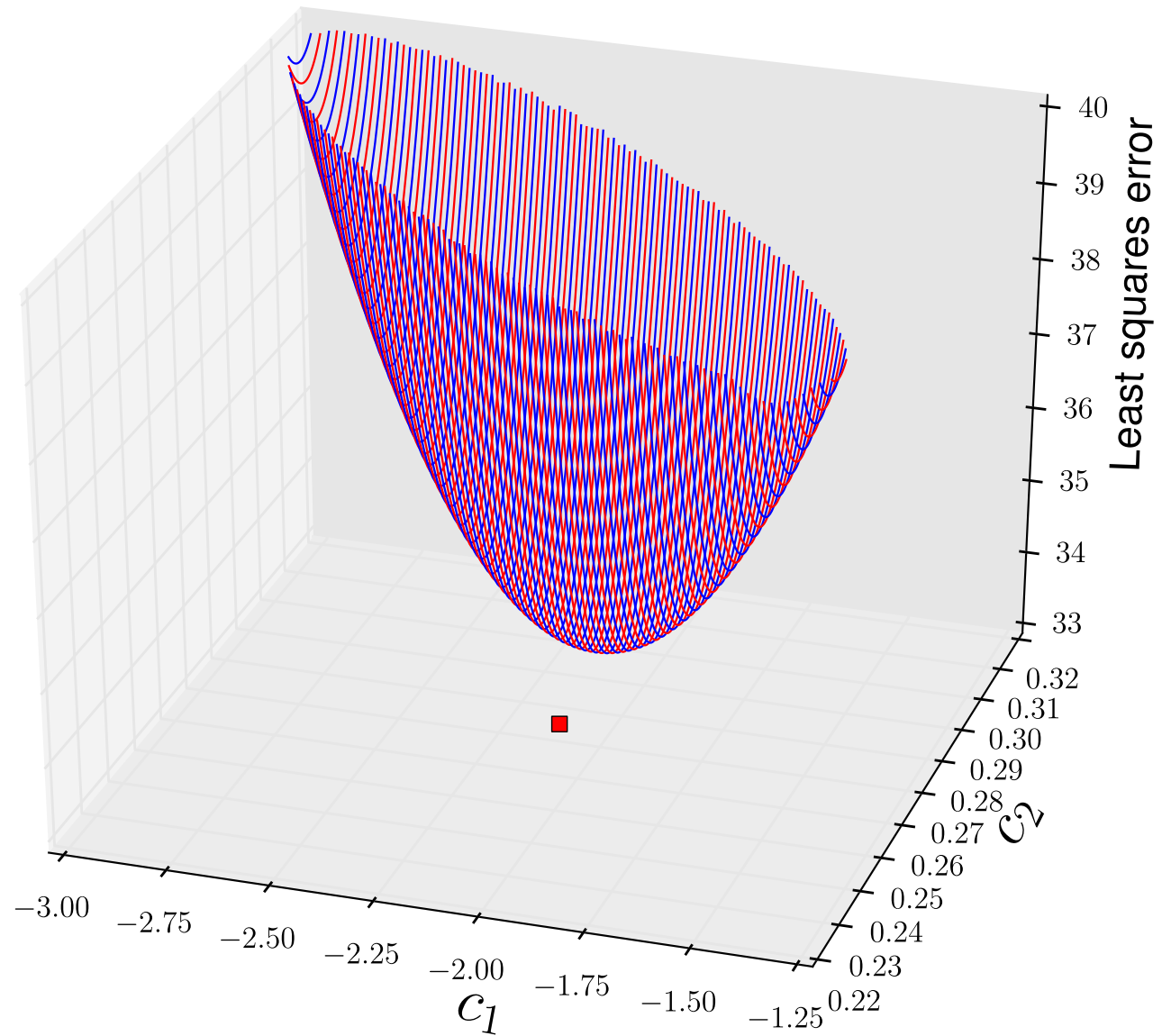
$$\sum_{k=1}^N (T_k - (c_1 + c_2 x_k))^2 ,$$

where, in our example, $N = 31$ and $x_k = k$.

- To do so set the *partial derivatives* w.r.t. c_1 and c_2 to zero:

$$\text{w.r.t. } c_1 : \quad -2 \sum_{k=1}^N (T_k - (c_1 + c_2 x_k)) = 0 ,$$

$$\text{w.r.t. } c_2 : \quad -2 \sum_{k=1}^N x_k (T_k - (c_1 + c_2 x_k)) = 0 .$$



The least squares error versus c_1 and c_2 .

From setting the partial derivatives to zero, we have

$$\sum_{k=1}^N (T_k - (c_1 + c_2 x_k)) = 0 \quad , \quad \sum_{k=1}^N x_k (T_k - (c_1 + c_2 x_k)) = 0 .$$

Solving these two equations for c_1 and c_2 gives

$$c_2 = \frac{\sum_{k=1}^N x_k T_k - \bar{x} \sum_{k=1}^N T_k}{\sum_{k=1}^N x_k^2 - N \bar{x}^2} ,$$

and

$$c_1 = \bar{T} - c_2 \bar{x} ,$$

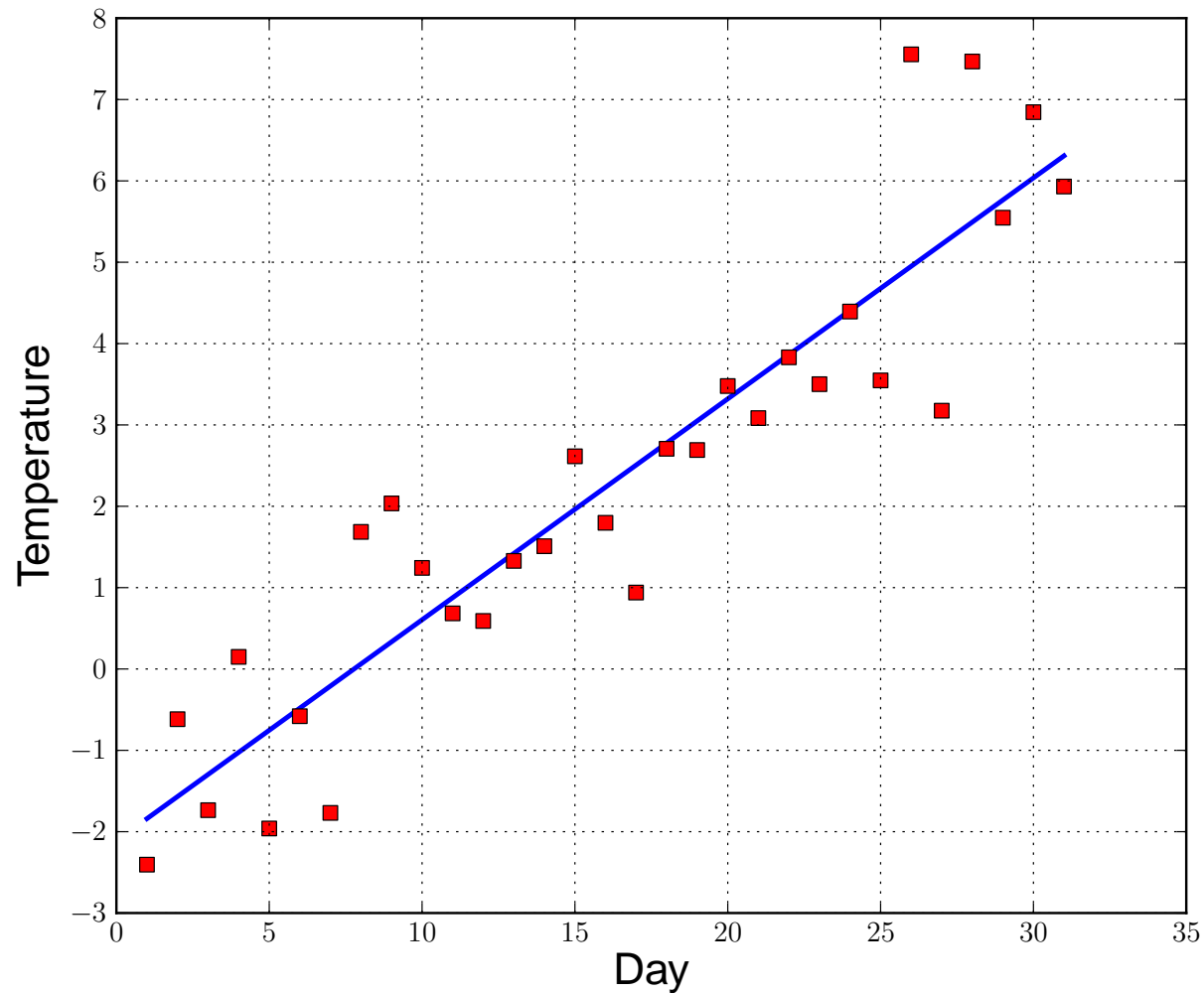
where

$$\bar{x} = \frac{1}{N} \sum_{k=1}^N x_k \quad , \quad \bar{T} = \frac{1}{N} \sum_{k=1}^N T_k .$$

EXERCISE: Check these formulas !

EXAMPLE: For our "March temperatures" example, we find

$$c_1 = -2.111 \quad \text{and} \quad c_2 = 0.272 .$$



Average daily high temperatures, with linear *least squares approximation* .

General Least Squares

Given discrete data points

$$\{ (x_i, y_i) \}_{i=1}^N ,$$

find the coefficients c_k of the function

$$p(x) \equiv \sum_{k=1}^n c_k \phi_k(x) ,$$

that *minimize* the *least squares error*

$$E_L \equiv \sum_{i=1}^N (p(x_i) - y_i)^2$$

EXAMPLES:

- $p(x) = c_1 + c_2 x .$ (Already done !)
- $p(x) = c_1 + c_2 x + c_3 x^2 .$ (Quadratic approximation)

For any vector $\mathbf{x} \in \mathbb{R}^N$ we have

$$\|\mathbf{x}\|_2^2 \equiv \mathbf{x}^T \mathbf{x} \equiv \sum_{i=1}^N x_i^2. \quad (T \text{ denotes } \textit{transpose}).$$

Then

$$\begin{aligned} E_L &\equiv \sum_{i=1}^N [p(x_i) - y_i]^2 = \left\| \begin{pmatrix} p(x_1) \\ \vdots \\ p(x_N) \end{pmatrix} - \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \right\|_2^2 \\ &= \left\| \begin{pmatrix} \sum_{i=1}^n c_i \phi_i(x_1) \\ \vdots \\ \sum_{i=1}^n c_i \phi_i(x_N) \end{pmatrix} - \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \right\|_2^2 \\ &= \left\| \begin{pmatrix} \phi_1(x_1) & \cdot & \phi_n(x_1) \\ \vdots & & \vdots \\ \phi_1(x_N) & \cdot & \phi_n(x_N) \end{pmatrix} \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} - \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \right\|_2^2 \equiv \|\mathbf{A}\mathbf{c} - \mathbf{y}\|_2^2. \end{aligned}$$

THEOREM

For the least squares error E_L to be *minimized* we must have

$$\mathbf{A}^T \mathbf{A} \mathbf{c} = \mathbf{A}^T \mathbf{y} .$$

PROOF:

$$\begin{aligned} E_L &= \|\mathbf{A}\mathbf{c} - \mathbf{y}\|_2^2 \\ &= (\mathbf{A}\mathbf{c} - \mathbf{y})^T (\mathbf{A}\mathbf{c} - \mathbf{y}) \\ &= (\mathbf{A}\mathbf{c})^T \mathbf{A}\mathbf{c} - (\mathbf{A}\mathbf{c})^T \mathbf{y} - \mathbf{y}^T \mathbf{A}\mathbf{c} + \mathbf{y}^T \mathbf{y} \\ &= \mathbf{c}^T \mathbf{A}^T \mathbf{A}\mathbf{c} - \mathbf{c}^T \mathbf{A}^T \mathbf{y} - \mathbf{y}^T \mathbf{A}\mathbf{c} + \mathbf{y}^T \mathbf{y} . \end{aligned}$$

PROOF: continued ...

We had

$$E_L = \mathbf{c}^T \mathbf{A}^T \mathbf{A} \mathbf{c} - \mathbf{c}^T \mathbf{A}^T \mathbf{y} - \mathbf{y}^T \mathbf{A} \mathbf{c} + \mathbf{y}^T \mathbf{y} .$$

For a *minimum* we need

$$\frac{\partial E_L}{\partial \mathbf{c}} = 0, \quad i.e., \quad \frac{\partial E_L}{\partial c_i} = 0, \quad i = 0, 1, \dots, n ,$$

which gives

$$\mathbf{c}^T \mathbf{A}^T \mathbf{A} + (\mathbf{A}^T \mathbf{A} \mathbf{c})^T - (\mathbf{A}^T \mathbf{y})^T - \mathbf{y}^T \mathbf{A} = 0, \quad (\text{Check !})$$

i.e.,

$$2\mathbf{c}^T \mathbf{A}^T \mathbf{A} - 2\mathbf{y}^T \mathbf{A} = 0 ,$$

or

$$\mathbf{c}^T \mathbf{A}^T \mathbf{A} = \mathbf{y}^T \mathbf{A} .$$

Transposing gives

$$\mathbf{A}^T \mathbf{A} \mathbf{c} = \mathbf{A}^T \mathbf{y} . \quad \text{QED}$$

EXAMPLE: Given the data points

$$\{ (x_i, y_i) \}_{i=1}^4 = \{ (0, 1), (1, 3), (2, 2), (4, 3) \},$$

find the coefficients c_1 and c_2 of $p(x) = c_1 + c_2x$,
that minimize

$$E_L \equiv \sum_{i=1}^4 [(c_1 + c_2x_i) - y_i]^2.$$

SOLUTION: Here $N = 4$, $n = 2$, $\phi_1(x) = 1$, $\phi_2(x) = x$.

Use the Theorem:

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 4 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 4 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \\ 2 \\ 3 \end{pmatrix},$$

or

$$\begin{pmatrix} 4 & 7 \\ 7 & 21 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} 9 \\ 19 \end{pmatrix},$$

with solution $c_1 = 1.6$ and $c_2 = 0.371429$.

EXAMPLE: Given the same data points, find the coefficients of

that minimize
$$p(x) = c_1 + c_2x + c_3x^2 ,$$

$$E_L \equiv \sum_{i=1}^4 [(c_1 + c_2 x_i + c_3 x_i^2) - y_i]^2 .$$

SOLUTION: Here

$$N = 4 \quad , \quad n = 3 \quad , \quad \phi_1(x) = 1 \quad , \quad \phi_2(x) = x \quad , \quad \phi_3(x) = x^2 .$$

Use the Theorem:

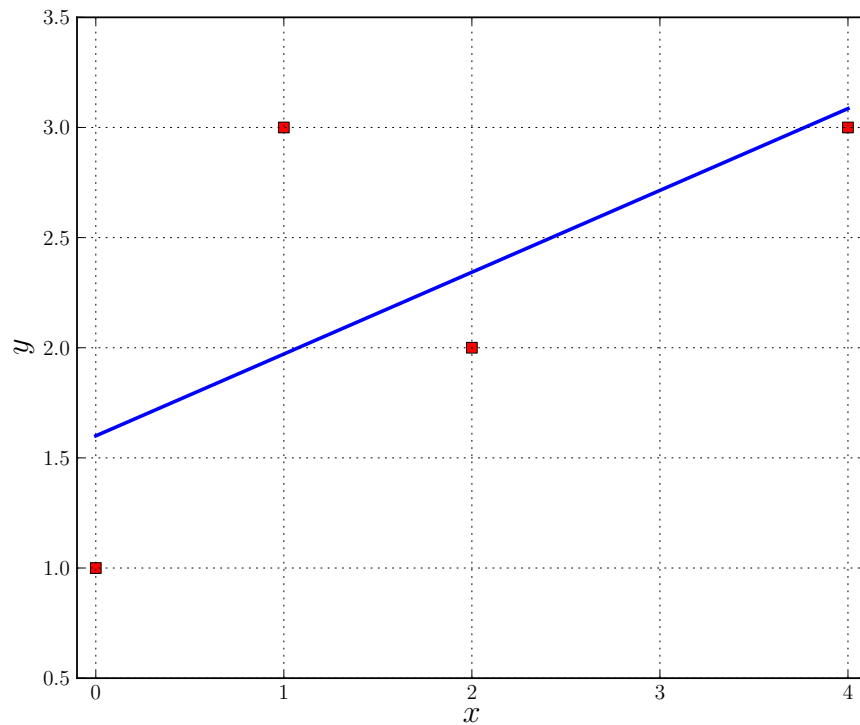
$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 4 \\ 0 & 1 & 4 & 16 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 4 & 16 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ c_3 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 4 \\ 0 & 1 & 4 & 16 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \\ 2 \\ 3 \end{pmatrix} ,$$

or

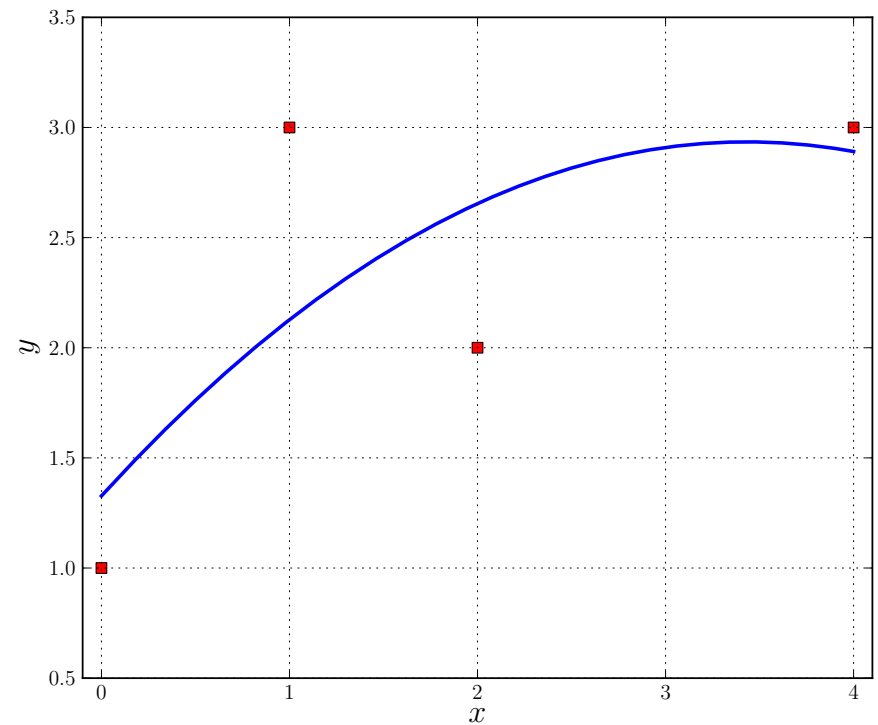
$$\begin{pmatrix} 4 & 7 & 21 \\ 7 & 21 & 73 \\ 21 & 73 & 273 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ c_3 \end{pmatrix} = \begin{pmatrix} 9 \\ 19 \\ 59 \end{pmatrix} ,$$

with solution $c_1 = 1.32727$, $c_2 = 0.936364$, $c_3 = -0.136364$.

The least squares approximations from the preceding two examples:



$$p(x) = c_1 + c_2x$$



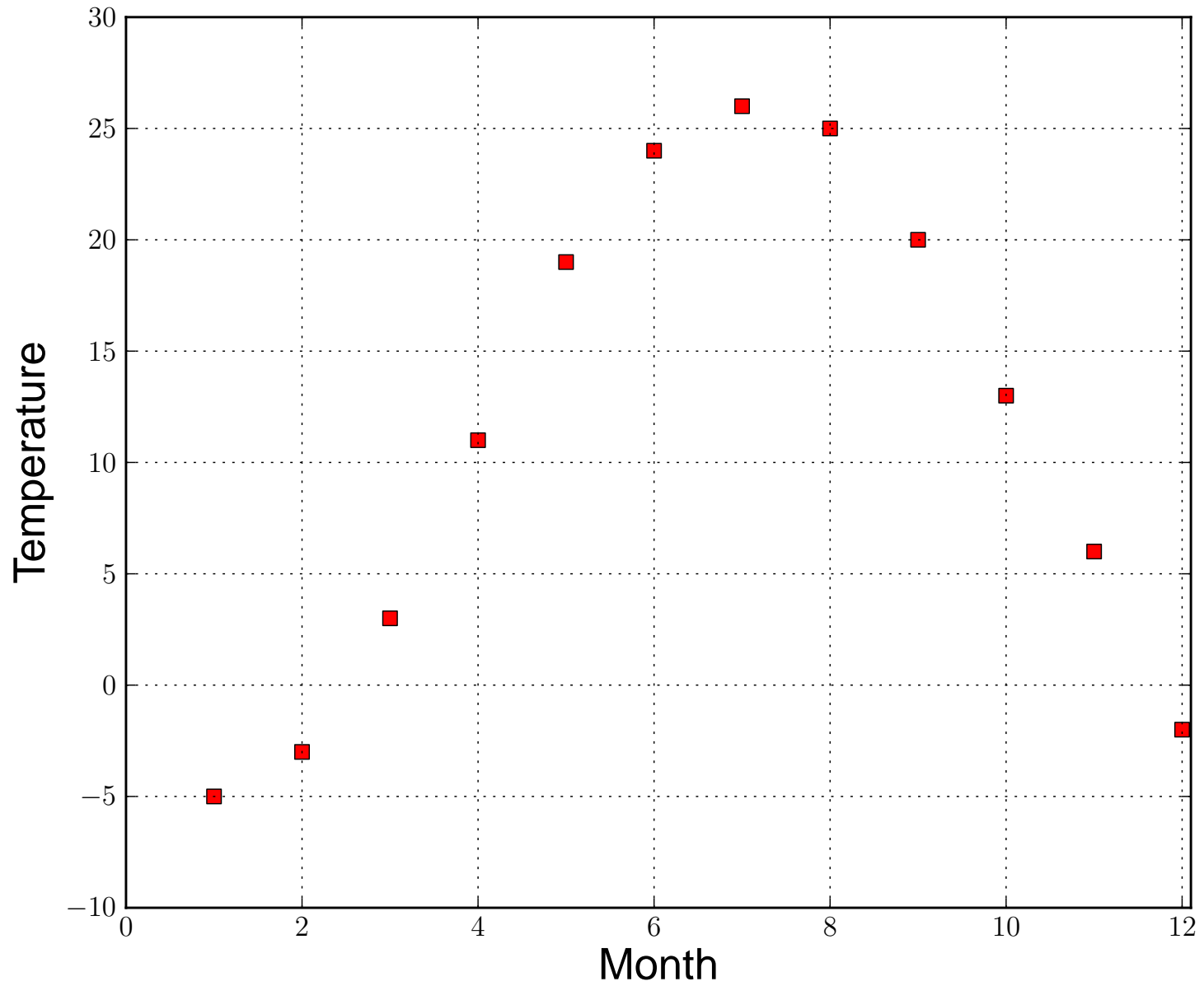
$$p(x) = c_1 + c_2x + c_3x^2$$

EXAMPLE: From actual data:

The average daily high temperatures in Montreal (by month) are:

January	-5
February	-3
March	3
April	11
May	19
June	24
July	26
August	25
September	20
October	13
November	6
December	-2

Source: <http://weather.uk.msn.com>



Average daily high temperature in Montreal (by month).

EXAMPLE: continued ...

The graph suggests using a 3-term *least squares approximation*

$$p(x) = c_1 \phi_1(x) + c_2 \phi_2(x) + c_3 \phi_3(x) ,$$

of the form

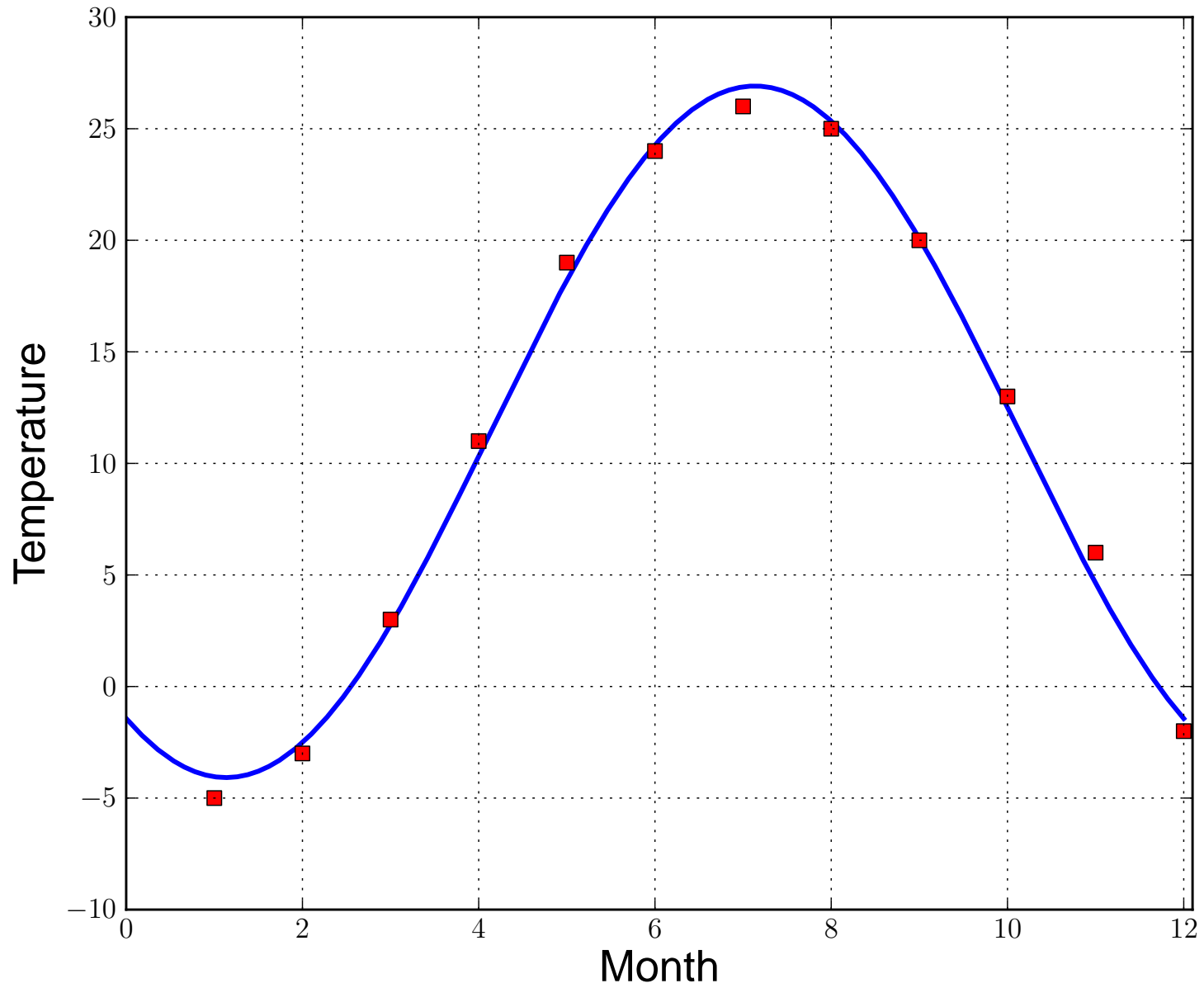
$$p(x) = c_1 + c_2 \sin\left(\frac{\pi x}{6}\right) + c_3 \cos\left(\frac{\pi x}{6}\right) .$$

QUESTIONS:

- Why include $\phi_2(x) = \sin\left(\frac{\pi x}{6}\right)$?
- Why is the argument $\frac{\pi x}{6}$?
- Why include the constant term $\phi_1(x) = c_1$?
- Why include $\phi_3(x) = \cos\left(\frac{\pi x}{6}\right)$?

In this example we find the least squares coefficients

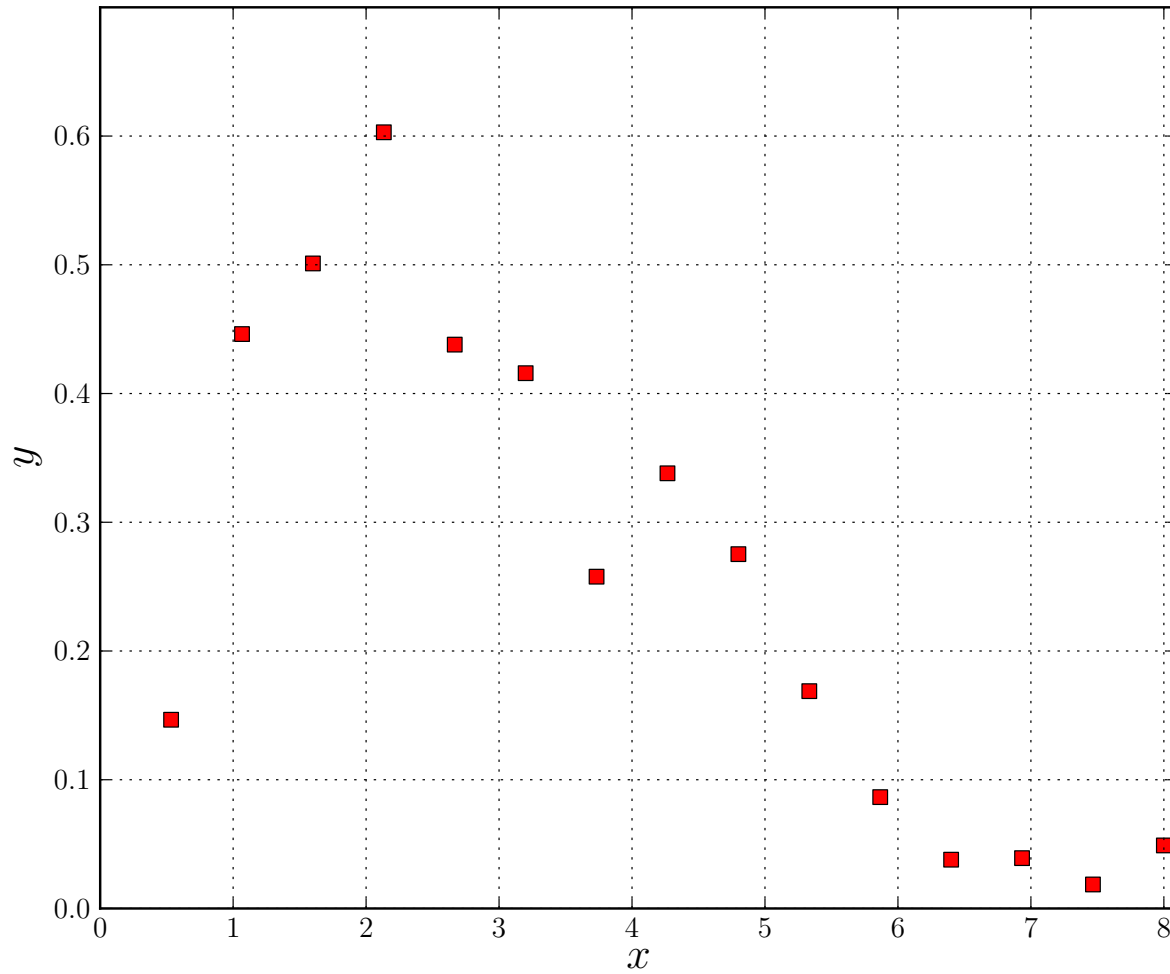
$$c_1 = 11.4 \quad , \quad c_2 = -8.66 \quad , \quad c_3 = -12.8 .$$



Least squares fit of average daily high temperatures.

EXAMPLE:

Consider the following *experimental data* :



EXAMPLE: continued ...

Suppose we are given that:

- These data contain "*noise*".
- The underlying physical process is understood.
- The *functional dependence* is *known* to have the form

$$y = c_1 x^{c_2} e^{-c_3 x} .$$

- The values of c_1 , c_2 , c_3 are *not* known.

EXAMPLE: continued ...

The functional relationship has the form

$$y = c_1 x^{c_2} e^{-c_3 x} .$$

Note that:

- The unknown coefficients c_1 , c_2 , c_3 appear *nonlinearly* !
- This gives *nonlinear equations* for c_1 , c_2 , c_3 !
- Such problems are more *difficult* to solve !
- What to do ?

EXAMPLE: continued ...

Fortunately, in this example we can take the *logarithm* :

$$\log y = \log c_1 + c_2 \log x - c_3 x .$$

This gives a *linear* relationship

$$\log y = \hat{c}_1 \phi_1(x) + c_2 \phi_2(x) + c_3 \phi_3(x) ,$$

where

$$\hat{c}_1 = \log c_1 .$$

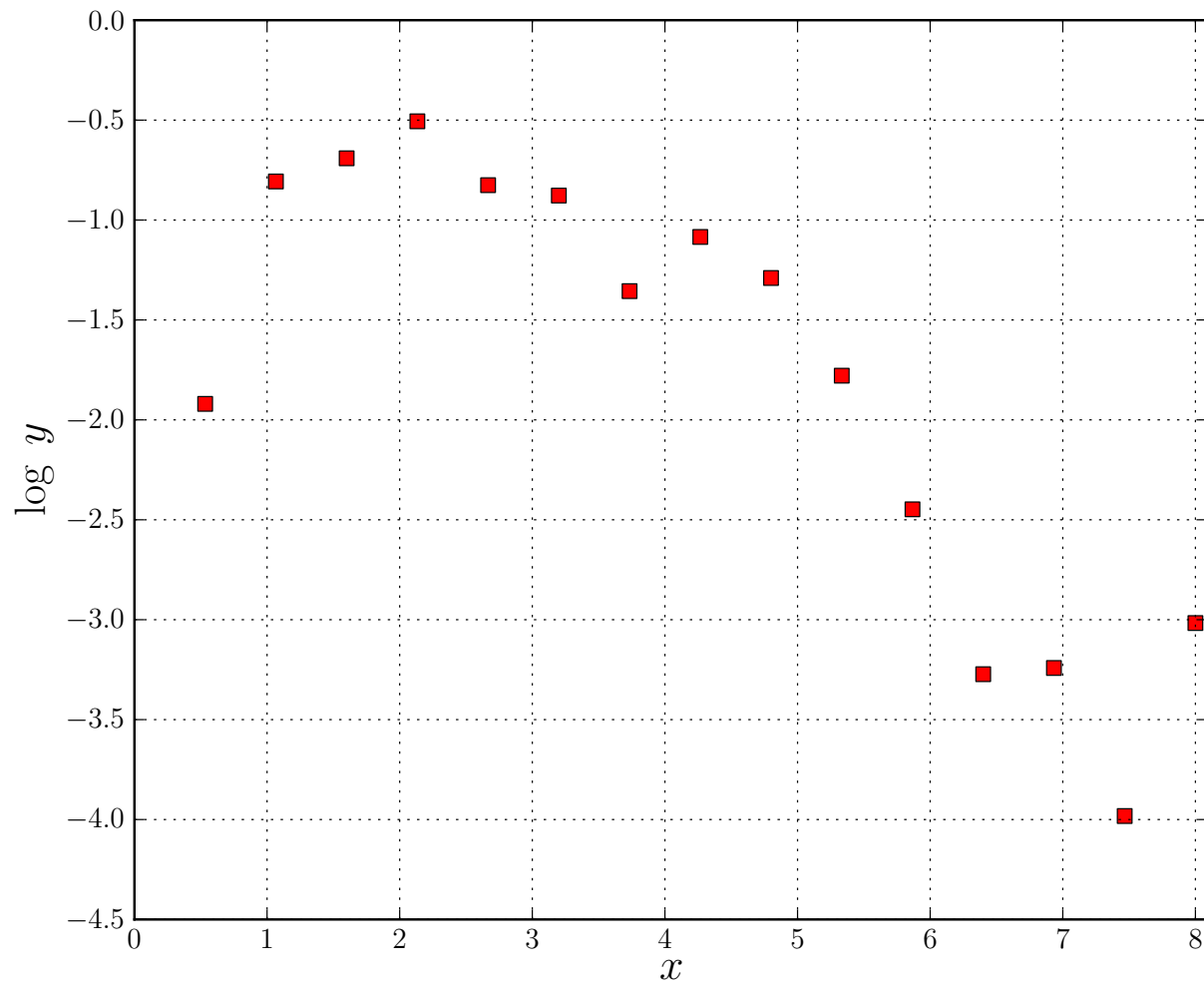
and

$$\phi_1(x) = 1 \quad , \quad \phi_2(x) = \log x \quad , \quad \phi_3(x) = -x .$$

Thus

- We can now use regular least squares.
- We first need to take the logarithm of the data.

EXAMPLE: continued ...



The logarithm of the original y -values versus x .

EXAMPLE: continued ...

We had

$$y = c_1 x^{c_2} e^{-c_3 x} ,$$

and

$$\log y = \hat{c}_1 \phi_1(x) + c_2 \phi_2(x) + c_3 \phi_3(x) ,$$

with

$$\phi_1(x) = 1 \quad , \quad \phi_2(x) = \log x \quad , \quad \phi_3(x) = -x \quad ,$$

and

$$\hat{c}_1 = \log c_1 .$$

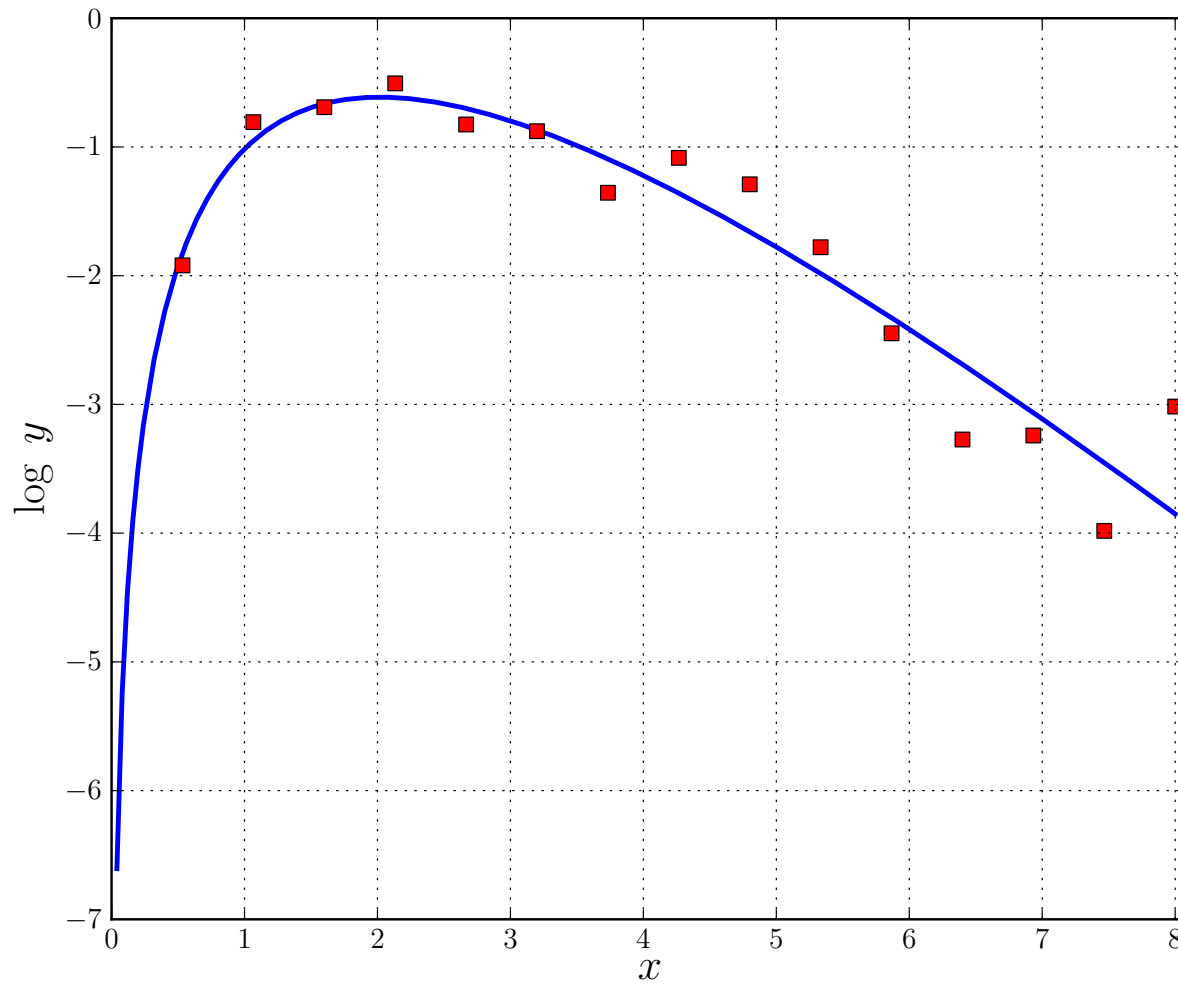
We find the following least squares values of the coefficients:

$$\hat{c}_1 = -0.00473 \quad , \quad c_2 = 2.04 \quad , \quad c_3 = 1.01 \quad ,$$

and

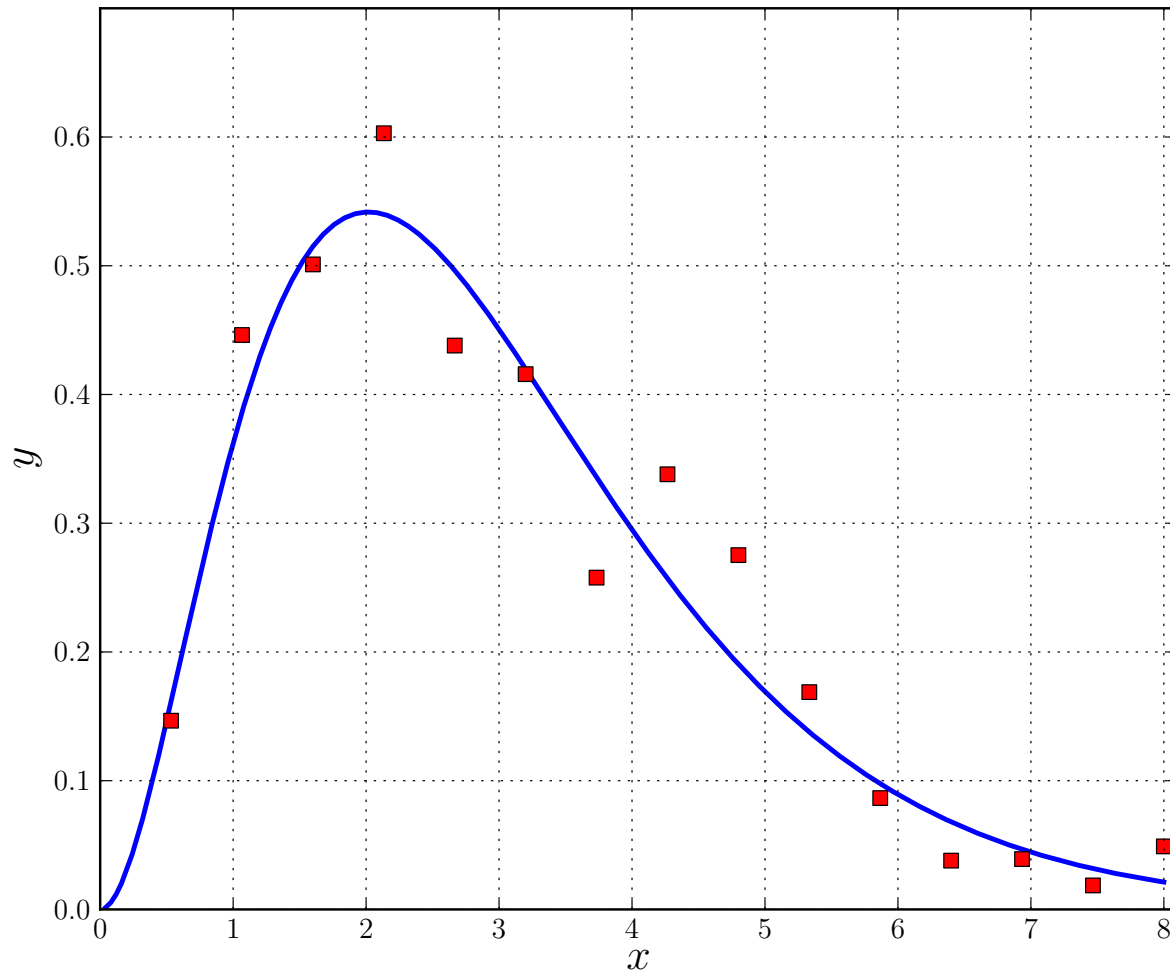
$$c_1 = e^{\hat{c}_1} = 0.995 .$$

EXAMPLE: continued ...



The least squares approximation of the transformed data.

EXAMPLE: continued ...



The least squares approximation shown in the original data.

EXERCISES:

- Compute the discrete least squares approximation of the form $p(x) = c_0 + c_1x + c_2x^2$ to the data $\{(0, 2), (1, 1), (2, 1), (3, 3)\}$.
- Compute the discrete least squares approximation of the form $p(x) = c_0 + c_1x + c_2 \frac{1}{x}$ to the data $\{(1, 5), (2, 3), (3, 2), (4, 3)\}$.
- Derive a formula in terms of N and n for the number of multiplications and divisions needed to solve the linear discrete least squares system

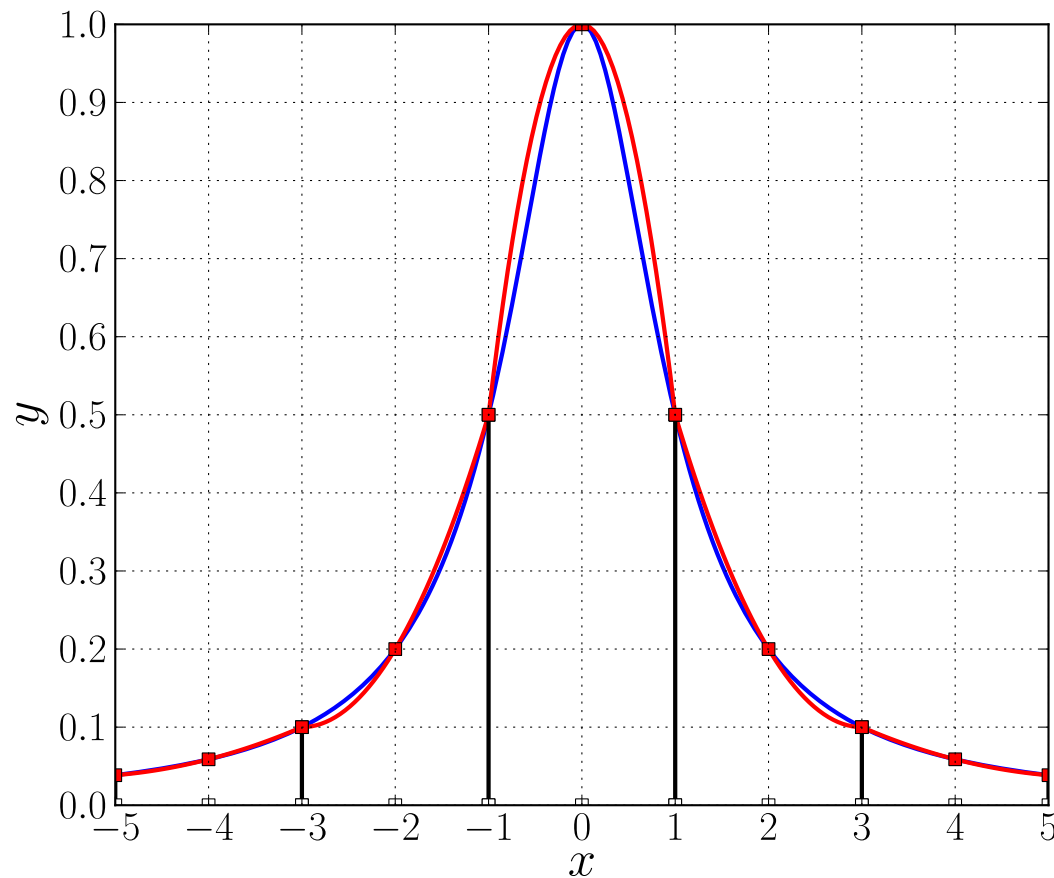
$$\mathbf{A}^T \mathbf{A} \mathbf{c} = \mathbf{A}^T \mathbf{y} ,$$

for $\mathbf{c} \in \mathbb{R}^n$, given the N by n matrix \mathbf{A} and the vector $\mathbf{y} \in \mathbb{R}^N$. Here \mathbf{A}^T denotes the transpose of \mathbf{A} . What is the total number of multiplications and divisions in terms of N for the special case $n = 2$?

SMOOTH INTERPOLATION BY PIECEWISE POLYNOMIALS

We have already discussed *local* (or *piecewise*) polynomial interpolation:

in each subinterval $[t_{j-1}, t_j]$ of the interval $[a, b]$ a given function f is interpolated by a polynomial $p_j \in \mathbb{P}_n$, at interpolation points $\{x_{j,i}\}_{i=0}^n$:



REMARKS:

- The collection $\{p_j\}_{j=1}^N$ defines a function $p(t)$ on $[a, b]$.
- $p(t)$ is generally *not smooth*, (not continuously differentiable) on $[a, b]$.
- In fact, in general $p(t)$ will not even be continuous, unless

$$x_{j,0} = t_{j-1} \quad \text{and} \quad x_{j,n} = t_j ,$$

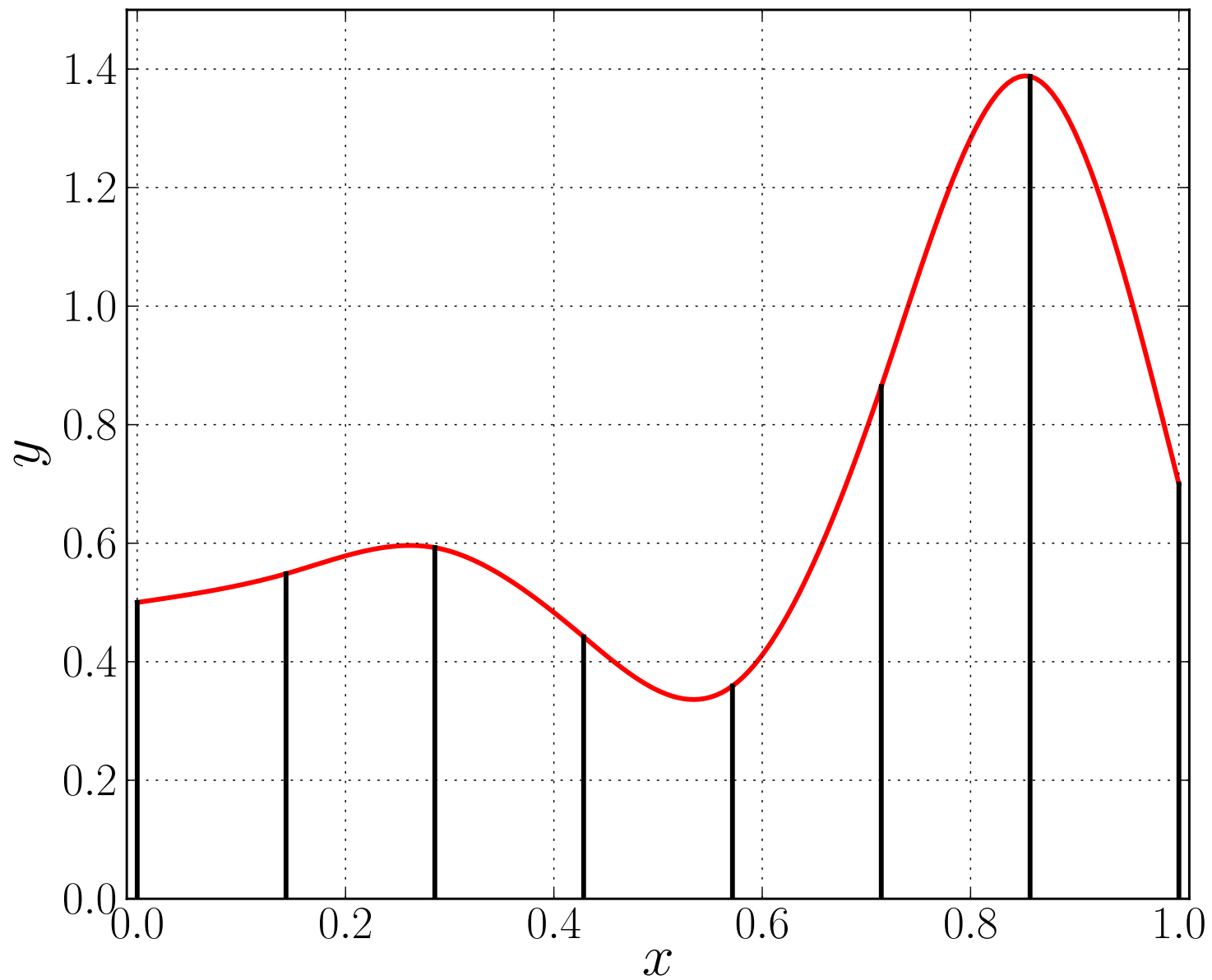
i.e., unless the first and last interpolation point in each subinterval coincide with the left and right end point of the interval, respectively,

REMARKS:

- Sometimes a *smooth interpolant* is wanted.
- These can also be constructed using piecewise polynomials.
- One class of smooth piecewise polynomial are called *cubic splines*.
- Cubic splines are piecewise polynomial functions

$$p(t) \in \mathbb{C}^2[a, b] ,$$

for which each component p_j is in \mathbb{P}_3 .



The cubic spline that interpolates the indicated data points.

Cubic Spline Interpolation.

Given $f(t)$ defined on $[a, b]$ we seek a function $p(t)$ satisfying :

- $p \in \mathbb{C}^2[a, b]$,
 - The restriction p_j of p to $[t_{j-1}, t_j]$ lies in \mathbb{P}_3 ,
 - $p(t_j) = f(t_j)$, $j = 0, 1, \dots, N$,
 - $p''(t_0) = 0$, $p''(t_N) = 0$.
-
- There are other possible choices for • .
 - With the above choice of • a spline is called the *natural cubic spline*.
 - We may also have *discrete data points* (t_j, f_j) , $j = 0, 1, \dots, N$.

This spline is "formally well defined", because

the total number of unknowns is $4N$,

(since each p_j is defined by four coefficients),

which is matched by the number of equations :

continuity equations	$3(N - 1)$
interpolation equations	$N + 1$
end point conditions	2
	<hr/>
Total	$4N$

REMARKS:

- In practice we do not solve these $4N$ equations to find the spline.
- Often we want the values of the spline at a large number of points, whereas the actual number of data points

$$\{ (t_j , f_j) \}_{j=0}^N$$

is relatively small.

- For this purpose we derive a more efficient algorithm below.

Consider the interval $[t_{j-1}, t_j]$ of size h_j .

To simplify notation take the interval $[t_0, t_1]$ of size h_1 .

Corresponding to this interval we have a polynomial $p \in \mathbb{P}_3$.

We can write

$$\begin{aligned} p(t_0) &= p_0, & p(t_1) &= p_1, \\ p''(t_0) &= p_0'', & p''(t_1) &= p_1''. \end{aligned}$$

These four equations uniquely define $p \in \mathbb{P}_3$ in terms of the values

$$p_0, \quad p_1, \quad p_0'', \quad p_1''.$$

In fact, for the interval $[t_0, t_1]$, one finds the polynomial

$$p_1(t) = \frac{p_0''}{6h_1} (t_1-t)^3 + \frac{p_1''}{6h_1} (t-t_0)^3 + \left(\frac{p_1}{h_1} - \frac{p_1''h_1}{6}\right) (t-t_0) + \left(\frac{p_0}{h_1} - \frac{p_0''h_1}{6}\right) (t_1-t).$$

Indeed, $p_1 \in \mathbb{P}_3$, and

$$p(t_0) = p_0, \quad p(t_1) = p_1,$$

$$p''(t_0) = p_0'', \quad p''(t_1) = p_1''.$$

Similarly, for the interval $[t_1, t_2]$, one finds the polynomial

$$p_2(t) = \frac{p_1''}{6h_2} (t_2-t)^3 + \frac{p_2''}{6h_2} (t-t_1)^3 + \left(\frac{p_2}{h_2} - \frac{p_2''h_2}{6}\right) (t-t_1) + \left(\frac{p_1}{h_2} - \frac{p_1''h_2}{6}\right) (t_2-t).$$

EXERCISE: Derive the formulas given above.

By construction the local polynomials p_1 and p_2 connect continuously at t_1 .

By construction the second derivatives also connect continuously.

However, the first derivatives must also match :

$$p_1'(t_1) = p_2'(t_1) .$$

This requirement leads to the *consistency relation*

$$h_1 p_0'' + 2(h_1 + h_2) p_1'' + h_2 p_2'' = 6 \left(\frac{p_2 - p_1}{h_2} - \frac{p_1 - p_0}{h_1} \right) .$$

EXERCISE: Derive this formula.

For consecutive intervals $[t_{j-1}, t_j]$ and $[t_j, t_{j+1}]$, the *consistency relation* is

$$h_j p''_{j-1} + 2(h_j + h_{j+1}) p''_j + h_{j+1} p''_{j+1} = 6 \left(\frac{p_{j+1} - p_j}{h_{j+1}} - \frac{p_j - p_{j-1}}{h_j} \right),$$

where

$$h_j \equiv t_j - t_{j-1} \quad \text{and} \quad h_{j+1} \equiv t_{j+1} - t_j.$$

We have *one* such equation for *each interior mesh point*.

To *interpolate the data points* $\{(t_j, f_j)\}_{j=0}^N$, we have

$$p_j = f_j, \quad j = 0, 1, \dots, N.$$

Furthermore we have the *natural spline endpoint conditions*

$$p''_0 = p''_N = 0.$$

REMARKS:

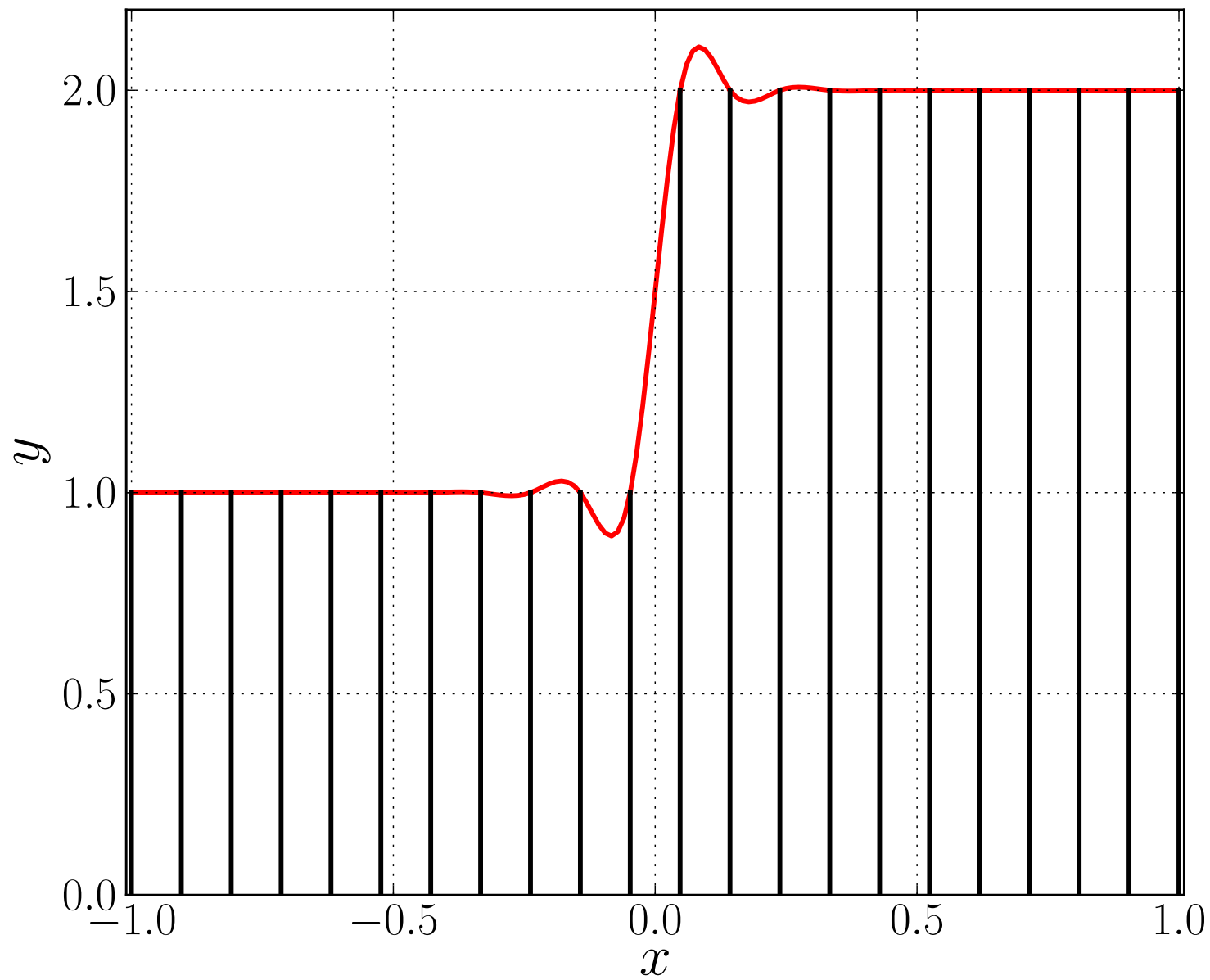
- In each row the diagonal entry is bigger than the sum of the other entries.
- Such a matrix is called *diagonally dominant*.
- By the Banach Lemma this matrix is nonsingular. (Check!)
- Thus we can compute the p_j'' using the tridiagonal algorithm.

Thereafter evaluate each local polynomial with the formula

$$p_j(t) = \frac{p_{j-1}''}{6h_j} (t_j - t)^3 + \frac{p_j''}{6h_j} (t - t_{j-1})^3 \\ + \left(\frac{p_j}{h_j} - \frac{p_j'' h_j}{6} \right) (t - t_{j-1}) + \left(\frac{p_{j-1}}{h_j} - \frac{p_{j-1}'' h_j}{6} \right) (t_j - t) .$$

REMARKS:

- The smoothness makes the component polynomials *interdependent*.
- One can not determine each component polynomial individually.
- As seen above, a tridiagonal system must be solved.
- This interdependence can lead to unwanted *oscillations*.



The cubic spline that interpolates the indicated data points.

NUMERICAL METHODS FOR INITIAL VALUE PROBLEMS

Here we discuss some *basic concepts* that arise in the numerical solution of *initial value problems (IVPs)* in *ordinary differential equations (ODEs)*.

Consider the *first order* IVP

$$\mathbf{u}'(t) = \mathbf{f}(\mathbf{u}(t)), \quad \text{for } t \geq 0,$$

with *initial conditions*

$$\mathbf{u}(0) = \mathbf{u}_0.$$

Here \mathbf{u} , $\mathbf{f}(\cdot) \in \mathbb{R}^n$.

Many *higher order ODEs* can be rewritten as *first order systems* .

EXAMPLE:

$$\mathbf{u}'' = \mathbf{g}(\mathbf{u}(t) , \mathbf{u}'(t)) , \quad \mathbf{u} , \mathbf{g}(\cdot) \in \mathbb{R}^n .$$

with *initial conditions*

$$\mathbf{u}(0) = \mathbf{u}_0 ,$$

$$\mathbf{u}'(0) = \mathbf{v}_0 ,$$

can be *rewritten* as

$$\mathbf{u}'(t) = \mathbf{v}(t) , \quad \mathbf{u} , \mathbf{v} \in \mathbb{R}^n ,$$

$$\mathbf{v}'(t) = \mathbf{g}(\mathbf{u}(t) , \mathbf{v}(t)) , \quad \mathbf{v} , \mathbf{g}(\cdot) \in \mathbb{R}^n ,$$

with *initial conditions*

$$\mathbf{u}(0) = \mathbf{u}_0 ,$$

$$\mathbf{v}(0) = \mathbf{v}_0 .$$

EXAMPLE:

The equations of *motion of a satellite* in an *Earth-Moon* - like system are :

$$x'' = 2y' + x - (1 - \mu)(x + \mu)r_1^{-3} - \mu(x - 1 + \mu)r_2^{-3} ,$$

$$y'' = -2x' + y - (1 - \mu)yr_1^{-3} - \mu yr_2^{-3} ,$$

$$z'' = -(1 - \mu)zr_1^{-3} - \mu zr_2^{-3} ,$$

where

$$r_1 = \sqrt{(x + \mu)^2 + y^2 + z^2} , \quad r_2 = \sqrt{(x - 1 + \mu)^2 + y^2 + z^2} .$$

Rewritten as a *first order system* :

$$x' = v_x ,$$

$$y' = v_y ,$$

$$z' = v_z ,$$

$$v'_x = 2v_y + x - (1 - \mu)(x + \mu)r_1^{-3} - \mu(x - 1 + \mu)r_2^{-3} ,$$

$$v'_y = -2v_x + y - (1 - \mu)yr_1^{-3} - \mu yr_2^{-3} ,$$

$$v'_z = -(1 - \mu)zr_1^{-3} - \mu zr_2^{-3} ,$$

with $r_1 = \sqrt{(x + \mu)^2 + y^2 + z^2}$ and $r_2 = \sqrt{(x - 1 + \mu)^2 + y^2 + z^2}$.

This system *is of the form*

$$\mathbf{u}'(t) = \mathbf{f}(\mathbf{u}(t)) , \quad \text{with } \textit{initial condition} \quad \mathbf{u}(0) = \mathbf{u}_0 .$$

Here μ is the *mass ratio*, i.e.,

$$\mu \equiv \frac{m_2}{(m_1 + m_2)},$$

where m_1 is the mass of the *larger body*, and m_2 of the *smaller body*.

Examples :

$$\mu \approx 0.01215 \quad \text{for the Earth Moon system,}$$

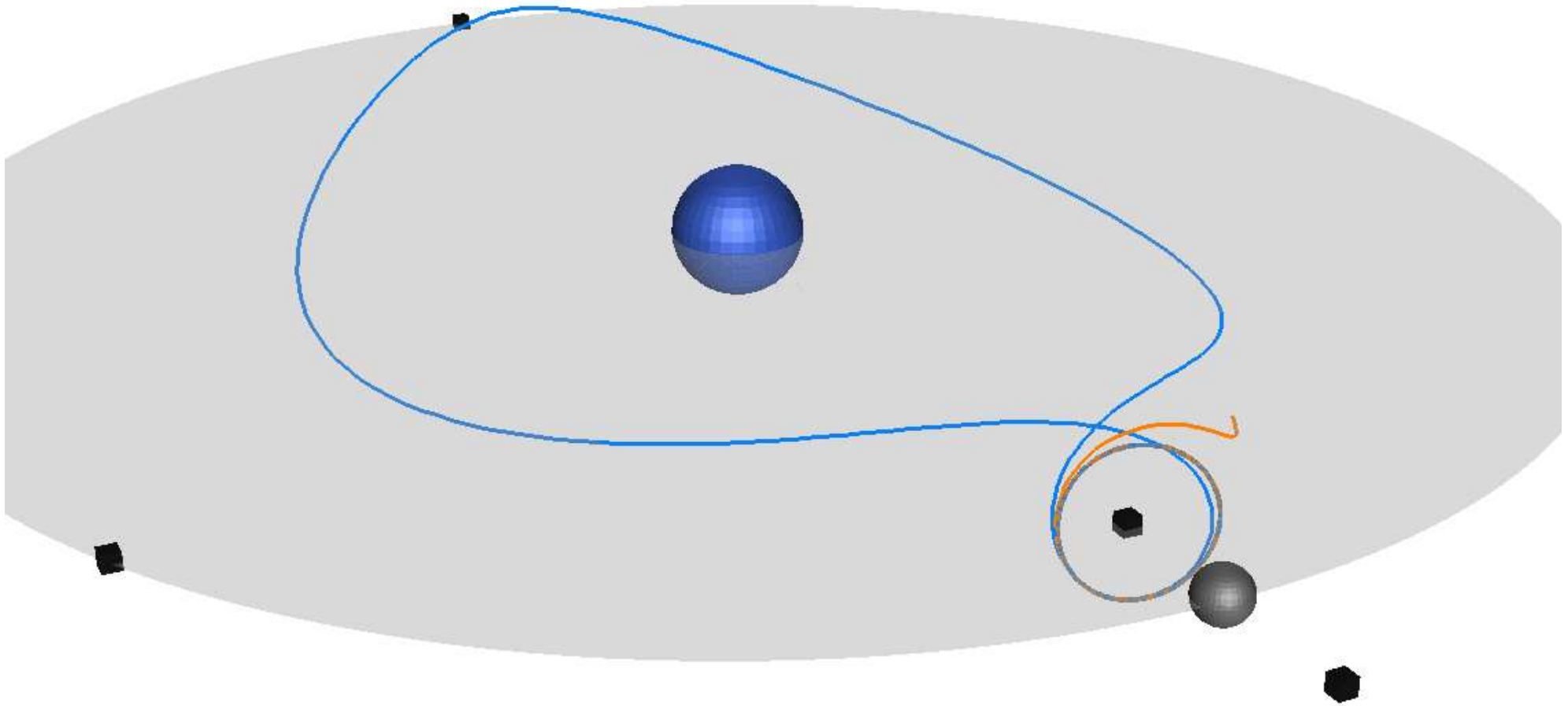
$$\mu \approx 9.53 \cdot 10^{-4} \quad \text{for the Sun Jupiter system,}$$

$$\mu \approx 3.0 \cdot 10^{-6} \quad \text{for the Sun Earth system.}$$

The *variables are scaled* such that

- the *distance between the two bodies* is 1 ,
- the *sum of their masses* is 1 .

The *larger body* is located at $(-\mu, 0, 0)$, and *the smaller body* at $(1-\mu, 0, 0)$.



A trajectory connecting a periodic “Halo orbit” to itself.

Numerical Methods.

Let

$$t_j \equiv j \Delta t, \quad j = 0, 1, 2, \dots .$$

Below we give several *basic numerical methods* for solving the IVP

$$\mathbf{u}'(t) = \mathbf{f}(\mathbf{u}(t)), \quad \mathbf{u}, \mathbf{f}(\cdot) \in \mathbb{R}^n .$$

$$\mathbf{u}(0) = \mathbf{u}_0 .$$

We use the *notation*

$$\mathbf{u}(t_j) = \text{the } \textit{exact solution} \text{ of the ODE at time } t_j ,$$

$$\mathbf{u}_j = \text{the } \textit{numerical solution} \text{ at time } t_j .$$

Euler's Method :

Using the *first order accurate* numerical differentiation formula

$$\frac{\mathbf{u}(t_{j+1}) - \mathbf{u}(t_j)}{\Delta t} \approx \mathbf{u}'(t_j) = \mathbf{f}(\mathbf{u}(t_j)) ,$$

we have

$$\mathbf{u}_{j+1} = \mathbf{u}_j + \Delta t \mathbf{f}(\mathbf{u}_j) , \quad j = 0, 1, 2, \dots ,$$

(*explicit, one-step* , $\mathcal{O}(\Delta t)$) .

(Check the order of accuracy!)

The Trapezoidal Method :

Using the *second order accurate* approximation formula

$$\frac{\mathbf{u}(t_{j+1}) - \mathbf{u}(t_j)}{\Delta t} \approx \frac{\mathbf{u}'(t_j) + \mathbf{u}'(t_{j+1})}{2} = \frac{\mathbf{f}(\mathbf{u}(t_j)) + \mathbf{f}(\mathbf{u}(t_{j+1}))}{2},$$

we have

$$\mathbf{u}_{j+1} = \mathbf{u}_j + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_j) + \mathbf{f}(\mathbf{u}_{j+1})], \quad j = 0, 1, 2, \dots,$$

(*implicit, one-step* , $\mathcal{O}(\Delta t^2)$) .

(Check the order of accuracy!)

NOTE: In each time-step a *nonlinear system* must be solved !

A Two-Step (Three-Point) Backward Differentiation Formula (BDF) :

Using the *second order accurate* approximation formula

$$\frac{3 \mathbf{u}(t_{j+1}) - 4 \mathbf{u}(t_j) + \mathbf{u}(t_{j-1}))}{2\Delta t} \approx \mathbf{u}'(t_{j+1}) = \mathbf{f}(\mathbf{u}(t_{j+1})), \quad (\text{Check!})$$

we have

$$\mathbf{u}_{j+1} = \frac{4}{3} \mathbf{u}_j - \frac{1}{3} \mathbf{u}_{j-1} + \frac{2\Delta t}{3} \mathbf{f}(\mathbf{u}_{j+1}), \quad j = 1, 2, \dots,$$

(*implicit, two-step*, $\mathcal{O}(\Delta t^2)$).

NOTE: In each time-step a *nonlinear system* must be solved!

A Two-Step (Three-Point) Forward Differentiation Formula :

Using the *second order accurate* approximation formula

$$\frac{-\mathbf{u}(t_{j+1}) + 4\mathbf{u}(t_j) - 3\mathbf{u}(t_{j-1}))}{2\Delta t} \approx \mathbf{u}'(t_{j-1}) = \mathbf{f}(\mathbf{u}(t_{j-1})), \quad (\text{Check!})$$

we have

$$\mathbf{u}_{j+1} = 4\mathbf{u}_j - 3\mathbf{u}_{j-1} - 2\Delta t \mathbf{f}(\mathbf{u}_{j-1}), \quad j = 1, 2, \dots,$$

(*explicit, two-step, $\mathcal{O}(\Delta t^2)$*).

(We will show that *this method is useless* !)

The Improved Euler Method :

$$\hat{\mathbf{u}}_{j+1} = \mathbf{u}_j + \Delta t \mathbf{f}(\mathbf{u}_j) ,$$

$$\mathbf{u}_{j+1} = \mathbf{u}_j + \frac{\Delta t}{2} [\mathbf{f}(\mathbf{u}_j) + \mathbf{f}(\hat{\mathbf{u}}_{j+1})] ,$$

for $j = 0, 1, 2, \dots$.

(*explicit, one-step, $\mathcal{O}(\Delta t^2)$*) .

An Explicit 4th order accurate Runge-Kutta Method :

$$\mathbf{k}_1 = \mathbf{f}(\mathbf{u}_j) ,$$

$$\mathbf{k}_2 = \mathbf{f}\left(\mathbf{u}_j + \frac{\Delta t}{2} \mathbf{k}_1\right) ,$$

$$\mathbf{k}_3 = \mathbf{f}\left(\mathbf{u}_j + \frac{\Delta t}{2} \mathbf{k}_2\right) ,$$

$$\mathbf{k}_4 = \mathbf{f}(\mathbf{u}_j + \Delta t \mathbf{k}_3) ,$$

$$\mathbf{u}_{j+1} = \mathbf{u}_j + \frac{\Delta t}{6} \{ \mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4 \} ,$$

for $j = 0, 1, 2, \dots$.

(*explicit, one-step, $\mathcal{O}(\Delta t^4)$*) .

The *order of accuracy* of a local formula can often be found by *Taylor expansion*.

EXAMPLE:

For the *two-step BDF* we have the *local discretization error*

$$\begin{aligned}
 \tau_j &\equiv \frac{1}{\Delta t} \left\{ \frac{3}{2}u(t_{j+1}) - 2u(t_j) + \frac{1}{2}u(t_{j-1}) \right\} - u'(t_{j+1}) \\
 &= \frac{1}{\Delta t} \left\{ \frac{3}{2} u(t_{j+1}) \right. \\
 &\quad \left. - 2 \left[u(t_{j+1}) - \Delta t u'(t_{j+1}) + \frac{\Delta t^2}{2} u''(t_{j+1}) - \frac{\Delta t^3}{6} u'''(t_{j+1}) + \cdots \right] \right. \\
 &\quad \left. + \frac{1}{2} \left[u(t_{j+1}) - 2\Delta t u'(t_{j+1}) + \frac{(2\Delta t)^2}{2} u''(t_{j+1}) - \frac{(2\Delta t)^3}{6} u'''(t_{j+1}) + \cdots \right] \right\} \\
 &\quad - u'(t_{j+1}) \\
 &= -\frac{1}{3} \Delta t^2 u'''(t_{j+1}) + \text{higher order terms.}
 \end{aligned}$$

The accuracy of this method is of *order 2*.

Stability of Numerical Approximations.

The very simple *model equation*

$$u'(t) = 0, \quad u(0) = u_0, \quad u, 0 \in \mathbb{R},$$

has *solution*

$$u(t) = u_0, \quad (\text{constant}).$$

A *general m-step approximation* is of the form

$$\alpha_m u_{j+1} + \alpha_{m-1} u_j \cdots + \alpha_0 u_{j+1-m} = 0.$$

Assume that

$$u_0 \text{ is given,}$$

and (if $m > 1$) that

$$u_1, u_2, \cdots, u_{m-1} \text{ are computed by another method,}$$

e.g., by a one-step method of the same order of accuracy.

General m-step approximation of $u'(t) = 0$, $u(0) = u_0$:

$$\alpha_m u_{j+1} + \alpha_{m-1} u_j \cdots + \alpha_0 u_{j+1-m} = 0 .$$

EXAMPLES:

(1)	$u_{j+1} - u_j = 0$,	u_0 given	Euler, Trapezoidal
(2)	$3u_{j+1} - 4u_j + u_{j-1} = 0$	u_0, u_1 given	Backward Differentiation
(3)	$-u_{j+1} + 4u_j - 3u_{j-1} = 0$	u_0, u_1 given	Forward Differentiation

The *difference equations*

$$\alpha_m u_{j+1} + \alpha_{m-1} u_j + \cdots + \alpha_0 u_{j+1-m} = 0 ,$$

can be solved explicitly :

Try *solutions of the form* $u_j = z^j$.

Then we have

$$\alpha_m z^{j+1} + \alpha_{m-1} z^j + \cdots + \alpha_0 z^{j+1-m} = 0 ,$$

or, multiplying through by z^{m-j-1}

$$\alpha_m z^m + \alpha_{m-1} z^{m-1} + \cdots + \alpha_1 z + \alpha_0 = 0 .$$

This is the *Characteristic Equation* of the difference equation.

Difference equation :

$$\alpha_m u_{j+1} + \alpha_{m-1} u_j + \cdots + \alpha_0 u_{j+1-m} = 0 .$$

Characteristic Equation :

$$\alpha_m z^m + \alpha_{m-1} z^{m-1} + \cdots + \alpha_1 z + \alpha_0 = 0 .$$

○ If $\alpha_m \neq 0$, then the characteristic equation has m *roots* $\{z_k\}_{k=1}^m$.

○ For simplicity we assume the roots are *distinct* .

○ The *general solution* of the difference equation is then

$$u_j = \gamma_1 z_1^j + \gamma_2 z_2^j + \cdots + \gamma_m z_m^j .$$

○ The coefficients γ_k are determined by the *initial data* $u_0, u_1, \cdots, u_{m-1}$.

FACT :

If the characteristic equation has *one or more zeroes* z_k with $|z_k| > 1$ then the numerical method is *unstable* .

In such a case the u_j can become *arbitrarily large* in a fixed time interval by taking Δt sufficiently small.

THEOREM:

A necessary condition for *numerical stability* of a multistep method is that the *characteristic equation*

$$\alpha_m z^{j+1} + \alpha_{m-1} z^j + \cdots + \alpha_0 = 0 ,$$

have *no zeroes outside the unit circle* .

EXAMPLES:

	Formula	Char. Eqn.	Roots	Stability
(1)	$u_{j+1} - u_j = 0$	$z - 1 = 0$	$z = 1$	Stable
(2)	$3u_{j+1} - 4u_j + u_{j-1} = 0$	$3z^2 - 4z + 1 = 0$	$z = 1, \frac{1}{3}$	Stable
(3)	$-u_{j+1} + 4u_j - 3u_{j-1} = 0$	$-z^2 + 4z - 3 = 0$	$z = 1, 3$	Unstable

Consider the last two examples in more detail :

Case (2) : Here the *general solution* is

$$u_j = \gamma_1 (1)^j + \gamma_2 \left(\frac{1}{3}\right)^j .$$

The *initial data* are u_0 and u_1 , so that

$$\gamma_1 + \gamma_2 = u_0 ,$$

$$\gamma_1 + \frac{1}{3} \gamma_2 = u_1 ,$$

from which

$$\gamma_1 = \frac{3}{2} u_1 - \frac{1}{2} u_0 ,$$

$$\gamma_2 = \frac{3}{2} u_0 - \frac{3}{2} u_1 .$$

Hence

$$u_j = \left(\frac{3}{2} u_1 - \frac{1}{2} u_0\right) + \left(\frac{3}{2} u_0 - \frac{3}{2} u_1\right) \left(\frac{1}{3}\right)^j .$$

If

$$u_1 = u_0 ,$$

then we see that

$$u_j = u_0 , \quad \text{for all } j .$$

Moreover, if

$$u_1 = u_0 + \epsilon ,$$

then

$$u_j = u_0 + \frac{3}{2} \epsilon - \frac{3\epsilon}{2} \left(\frac{1}{3}\right)^j ,$$

so that u_j *stays close* to u_0 if ϵ is small.

Case (3) : Here the *general solution* is

$$u_j = \gamma_1 (1)^j + \gamma_2 (3)^j ,$$

and using the initial data

$$\gamma_1 + \gamma_2 = u_0 ,$$

$$\gamma_1 + 3 \gamma_2 = u_1 ,$$

from which

$$\gamma_1 = \frac{3}{2} u_0 - \frac{1}{2} u_1 \quad , \quad \gamma_2 = \frac{1}{2} u_1 - \frac{1}{2} u_0 .$$

Hence

$$u_j = \left(\frac{3}{2} u_0 - \frac{1}{2} u_1 \right) + \left(\frac{1}{2} u_1 - \frac{1}{2} u_0 \right) (3)^j .$$

Again, if $u_1 = u_0$ then $u_j = u_0$ for all j .

But if $u_1 = u_0 + \epsilon$ then $u_j = u_0 - \frac{1}{2} \epsilon + \frac{1}{2} \epsilon 3^j$.

Hence u_j becomes *arbitrarily large in finite time* by taking small Δt !

THEOREM:

If the local approximation is accurate, and if the zeroes $\{z_k\}_{k=1}^m$ of the *characteristic equation*

$$\alpha_m z^m + \alpha_{m-1} z^{m-1} + \cdots + \alpha_1 z + \alpha_0 = 0 .$$

satisfy

$$|z_k| \leq 1 , \quad \text{and} \quad |z_k| = 1 \Rightarrow z_k \text{ is simple} ,$$

then the *method is stable* and

$$u_j \rightarrow u(t_j) \quad \text{as} \quad \Delta t \rightarrow 0 .$$

PROOF: Omitted.

Stiff Differential Equations.

There are ODEs for which most *explicit* difference approximations require Δt to be *very small* before one gets the convergence guaranteed by the theorem.

To investigate this, we use the *model equation*

$$u'(t) = \lambda u(t) , \quad t \geq 0 ,$$

with

$$u(0) = u_0 ,$$

where λ is a constant. (We allow λ to be *complex*.)

The *solution* is

$$u(t) = e^{\lambda t} u_0 .$$

Consider the case where

$$\operatorname{Re}(\lambda) \ll 0 ,$$

i.e., λ has *large negative real part* .

Then the *exact solution* of

$$u'(t) = \lambda u(t) ,$$

namely, $u(t) = e^{\lambda t} u_0$, *decays very quickly* as $t \rightarrow \infty$.

The *numerical solution* u_j again has the form

$$u_j = \sum_{k=1}^m \gamma_k z_k^j ,$$

and we certainly *don't want* u_j *to increase* as $j \rightarrow \infty$!

Thus we *we don't want any* z_k *outside the unit disk* in the complex plane.

However, for many difference formulas

Δt must be very small

in order that

all z_k , $k = 1, \dots, m$, are inside the unit disk .

Thus problems with

$Re(\lambda) \ll 0$, (“*Stiff Problems*”),

need *special difference approximations* .

More generally the problem

$$\mathbf{u}'(t) = \mathbf{f}(\mathbf{u}(t)) ,$$

$$\mathbf{u}(0) = \mathbf{u}_0 ,$$

is called *stiff* if the Jacobian

$$\mathbf{f}_u(\mathbf{u}(t)) ,$$

has one or more *eigenvalues* $\lambda_i = \lambda_i(t)$, with

$$\operatorname{Re}(\lambda_i) \ll 0 .$$

NOTE: Since eigenvalues can be complex *we allow λ to be complex* .

EXAMPLES:

We will approximate

$$u'(t) = \lambda u(t) ,$$

by *various discretization formulas* and determine the values of $\Delta t\lambda$ in the complex plane for which *the solution of the difference formula also decays* .

Assume

$$\Delta t > 0 , \quad \text{and} \quad \text{Re}(\lambda) < 0 .$$

Then $\Delta t\lambda$ always lies *in the negative half plane* , *i.e.*,

$$\text{Re}(\Delta t\lambda) < 0 .$$

Explicit Euler.

Applying Euler's explicit formula to the model equation

$$u'(t) = \lambda u(t) ,$$

we get the *difference equation*

$$\frac{1}{\Delta t} (u_{j+1} - u_j) = \lambda u_j , \quad i.e., \quad u_{j+1} = (1 + \Delta t \lambda) u_j .$$

Looking for solutions of the form $u_j = z^j$ gives the *characteristic equation*

$$z - (1 + \Delta t \lambda) = 0 , \quad \text{with zero} \quad z = 1 + \Delta t \lambda .$$

The *explicit solution* of the difference equation is

$$u_j = \gamma (1 + \Delta t \lambda)^j = u_0 (1 + \Delta t \lambda)^j .$$

For *explicit Euler* we just found that

$$z = 1 + \Delta t \lambda .$$

which represents a function (or *map*)

from the complex $\Delta t \lambda$ -plane *to* the complex z -plane .

We want to know:

which part of the $\Delta t \lambda$ -plane *is mapped into the unit circle* $|z| \leq 1$.

For such values of $\Delta t \lambda$ the difference formula will give *decaying solutions* .

NOTE: By $\Delta t \lambda$ -plane we mean all values of the product $\Delta t \lambda$,

where λ is complex, and where $\Delta t \in \mathbb{R}$, with $\Delta t > 0$.

For *explicit Euler* we found that

$$z = 1 + \Delta t \lambda ,$$

Now $|z| = 1$ if

$$|1 + \Delta t \lambda| = 1 ,$$

that is, if

$$1 + \Delta t \lambda = e^{i\theta} ,$$

that is, if

$$\Delta t \lambda = -1 + e^{i\theta} ,$$

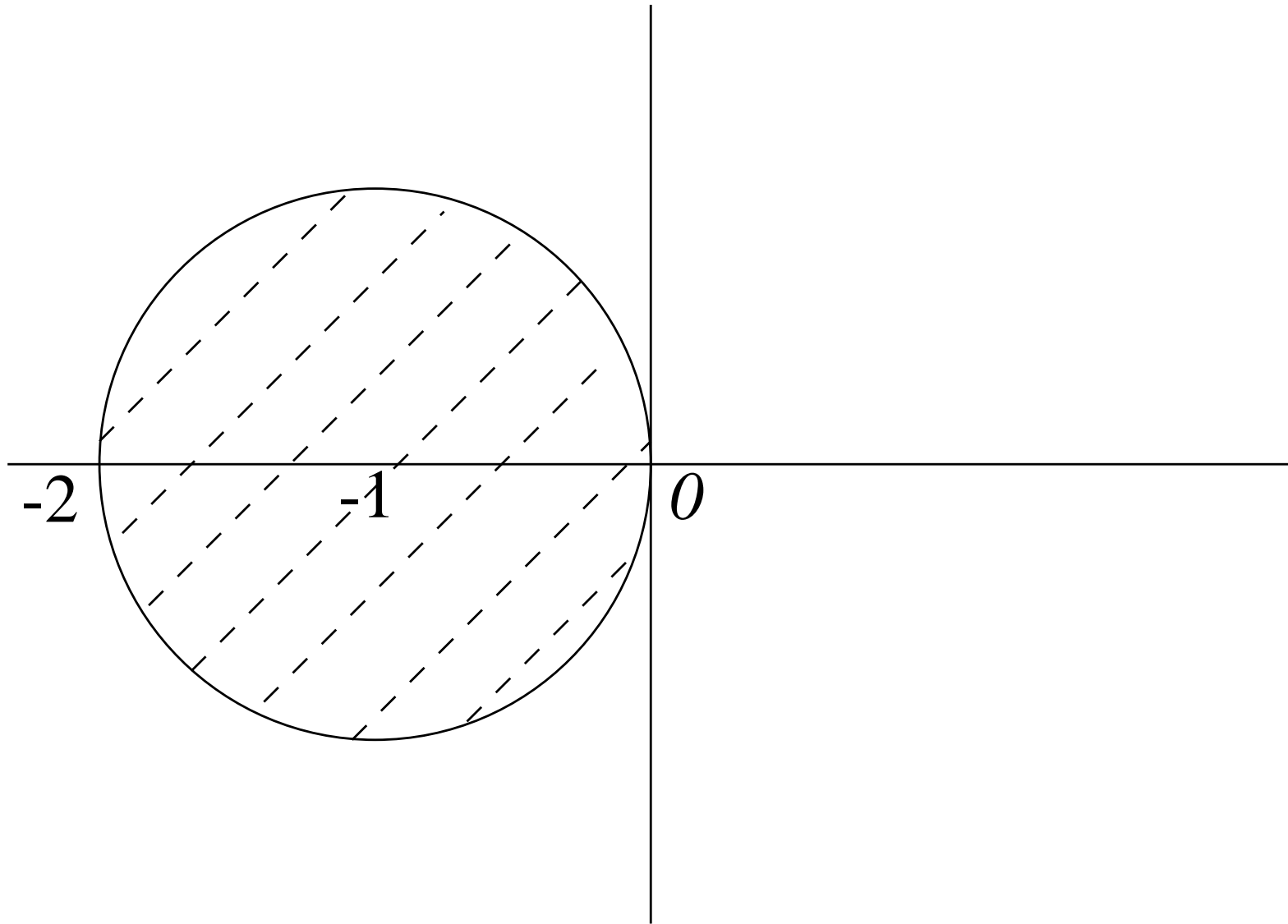
i.e., if $\Delta t \lambda$ is on the *circle of radius* 1 centered at -1 in the $\Delta t \lambda$ -plane.

(Check *the inside of this circle is mapped into the unit circle in the z -plane* !)

Thus the *region of stability* of the explicit Euler method is :

the *disk of radius* 1 centered at -1 in the $\Delta t \lambda$ -plane.

Complex $(\Delta t \cdot \lambda)$ -plane



Stability region of the Explicit Euler method

Consider a *specific example* :

Take $\lambda = -10^6$.

Then

$$u(t) = e^{(-10^6 t)} u_0 ,$$

which *decays very rapidly* for increasing t !

However, *for u_j to decay*, one must take

$$(\Delta t \lambda) > -2 ,$$

that is,

$$\Delta t < 2 \cdot 10^{-6} !$$

Thus the explicit Euler method is *useless* for stiff equations !

Implicit Euler.

The *difference formula* is

$$\frac{1}{\Delta t} (u_{j+1} - u_j) = \lambda u_{j+1} ,$$

i.e.,

$$u_{j+1} = \frac{1}{1 - \Delta t \lambda} u_j .$$

The *characteristic equation*

$$z - \frac{1}{1 - \Delta t \lambda} = 0 ,$$

has *zero*

$$z = \frac{1}{1 - \Delta t \lambda} .$$

$$\text{Implicit Euler : } z = \frac{1}{1 - \Delta t \lambda} .$$

Now

$$z = e^{i\theta} \quad \text{if} \quad \Delta t \lambda = 1 - e^{-i\theta} ,$$

that is,

if $\Delta t \lambda$ is on the circle of radius 1 centered at +1 in the complex $\Delta t \lambda$ -plane.

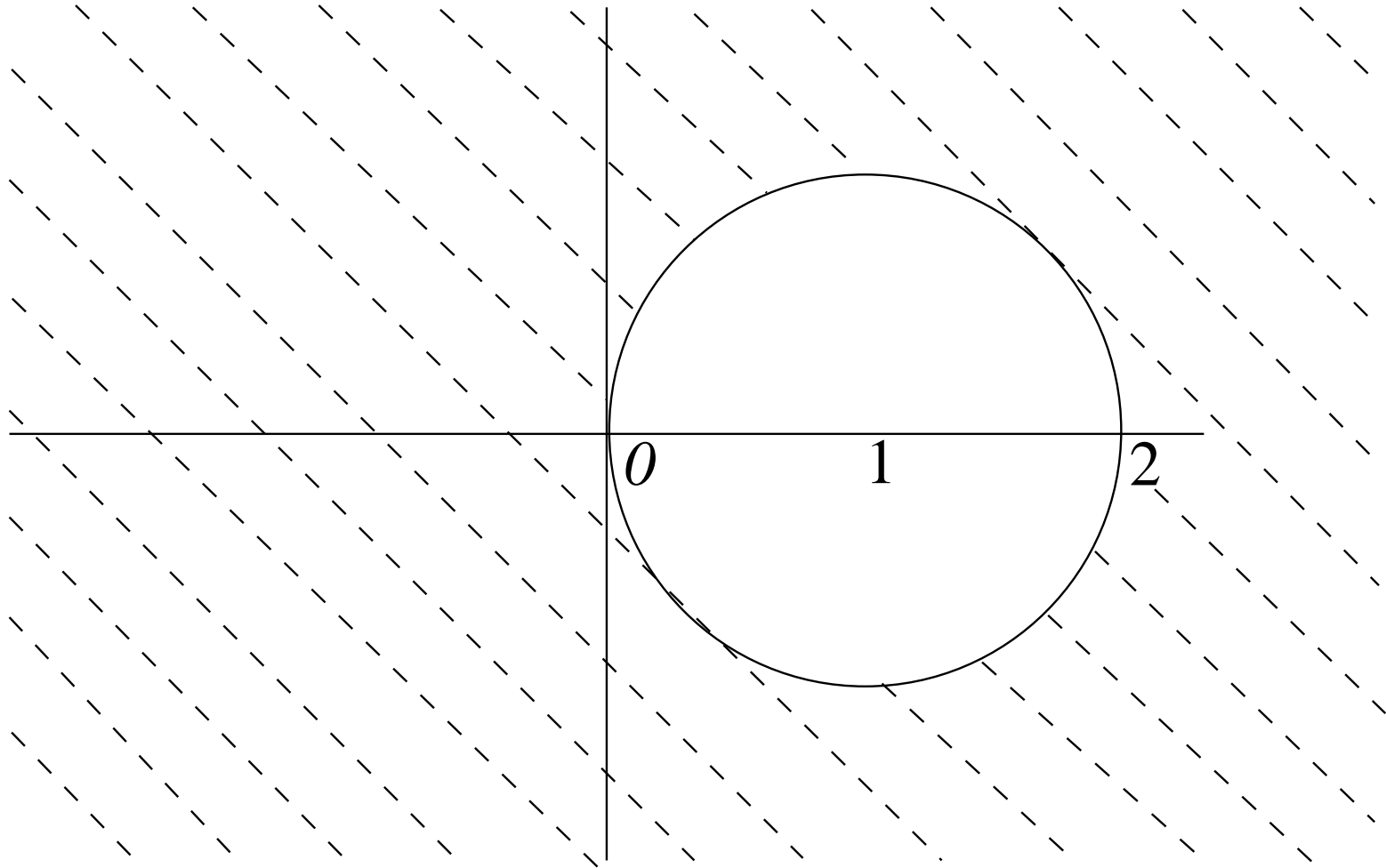
The *outside of this circle* is *mapped into the unit circle* in the z -plane.

Thus :

- u_j *decays* for all Δt .
- The step size Δt is *not* determined by stability requirements.
- Instead Δt is determined by accuracy requirements.
- Hence the implicit Euler method is *appropriate for stiff problems*.

(However its accuracy is not very high : only $\mathcal{O}(\Delta t)$.)

Complex $(\Delta t \cdot \lambda)$ -plane



Stability region of the Implicit Euler method

Trapezoidal Method.

When applied to $u'(t) = \lambda u(t)$, the Trapezoidal Method gives

$$\frac{1}{\Delta t} (u_{j+1} - u_j) = \frac{1}{2} \lambda (u_j + u_{j+1}) .$$

Thus the *characteristic equation* is

$$\left(1 - \frac{1}{2}\Delta t\lambda\right) z - \left(1 + \frac{1}{2}\Delta t\lambda\right) = 0 , \quad \text{with zero } z = \frac{1 + \frac{1}{2}\Delta t\lambda}{1 - \frac{1}{2}\Delta t\lambda} .$$

This time we find that $z = e^{i\theta}$ if

$$\Delta t\lambda = 2 \left(\frac{z - 1}{z + 1}\right) = 2 \left(\frac{e^{i\theta} - 1}{e^{i\theta} + 1}\right) = 2i \tan\left(\frac{\theta}{2}\right) .$$

The *region of stability* is now precisely the *entire negative half plane*.

Thus, $z \leq 1$ if and only if $Re(\Delta t\lambda) < 0$, which is *very desirable*.

A *disadvantage* is that the decay rate becomes smaller when

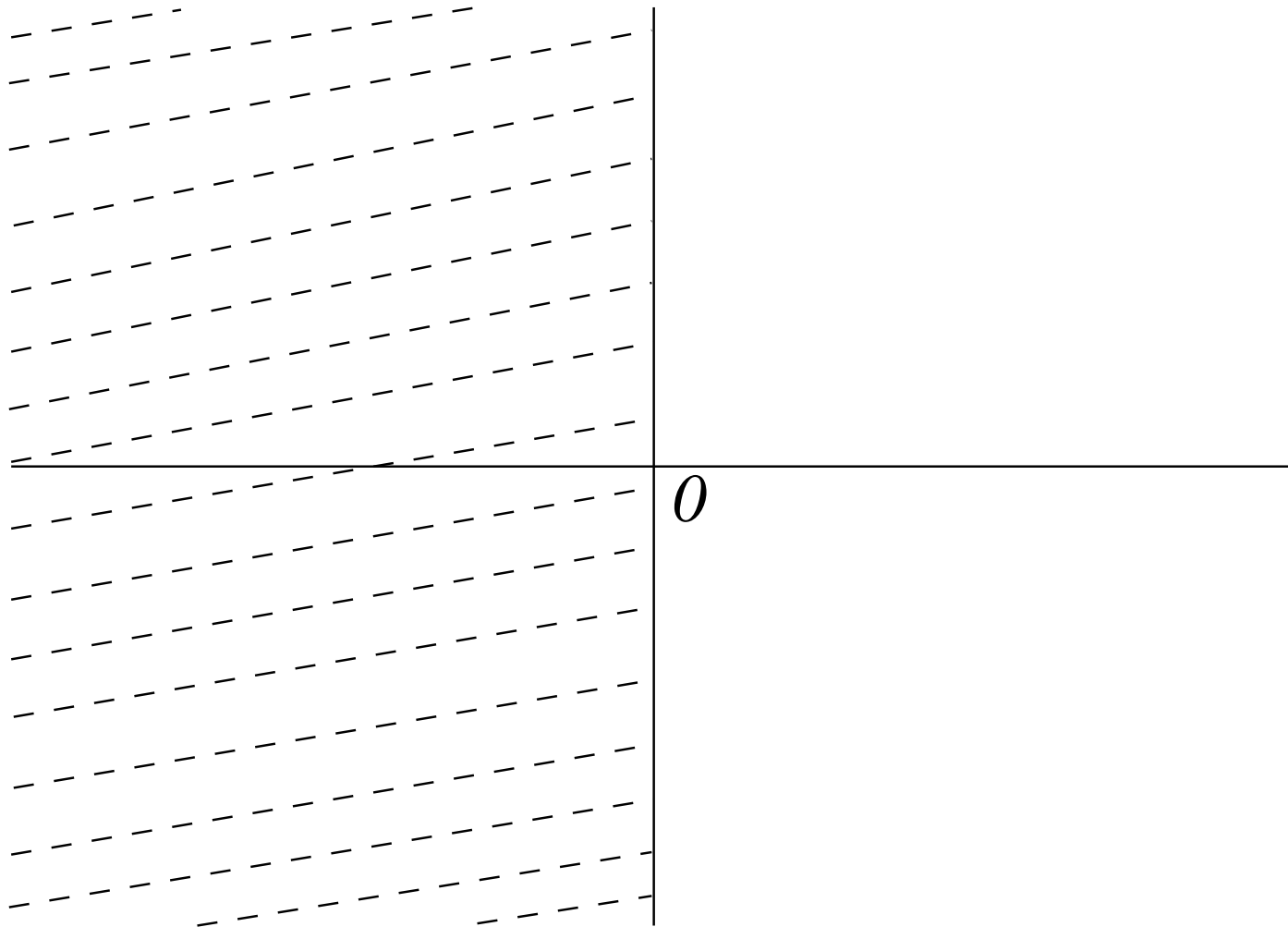
$$\operatorname{Re}(\lambda) \rightarrow -\infty ,$$

contrary to the decay rate of the solution of the differential equation.

In fact (thinking of Δt as fixed) we have

$$\lim_{\lambda \rightarrow -\infty} z(\lambda) = \lim_{\lambda \rightarrow -\infty} \frac{1 + \frac{1}{2}\Delta t\lambda}{1 - \frac{1}{2}\Delta t\lambda} = -1 .$$

Complex $(\Delta t \cdot \lambda)$ -plane



Stability region of the Trapezoidal method

Backward Differentiation Formulas (BDF).

For the differential equation $u'(t) = f(u(t))$ the BDF take the form

$$\frac{1}{\Delta t} \sum_{i=0}^m \alpha_i u_{j+1-i} = f(u_{j+1}) .$$

The $\{\alpha_i\}_{i=0}^m$ are chosen so the order is as high as possible, namely, $\mathcal{O}(\Delta t^m)$.

These formulas follow from the *numerical differentiation formulas* that approximate $u'(t_{j+1})$ in terms of

$$u(t_{j+1}) , u(t_j) , \dots , u(t_{j+1-m}) .$$

All of these methods are *implicit* .

The choice $m = 1$ gives the implicit Euler method.

Let \mathcal{S}_m denote the *stability region* of the m -step BDF.

About \mathcal{S}_m one can show the following :

$m = 1, 2 :$

\mathcal{S}_m *contains the negative half plane* .

These methods are called *A-stable* .

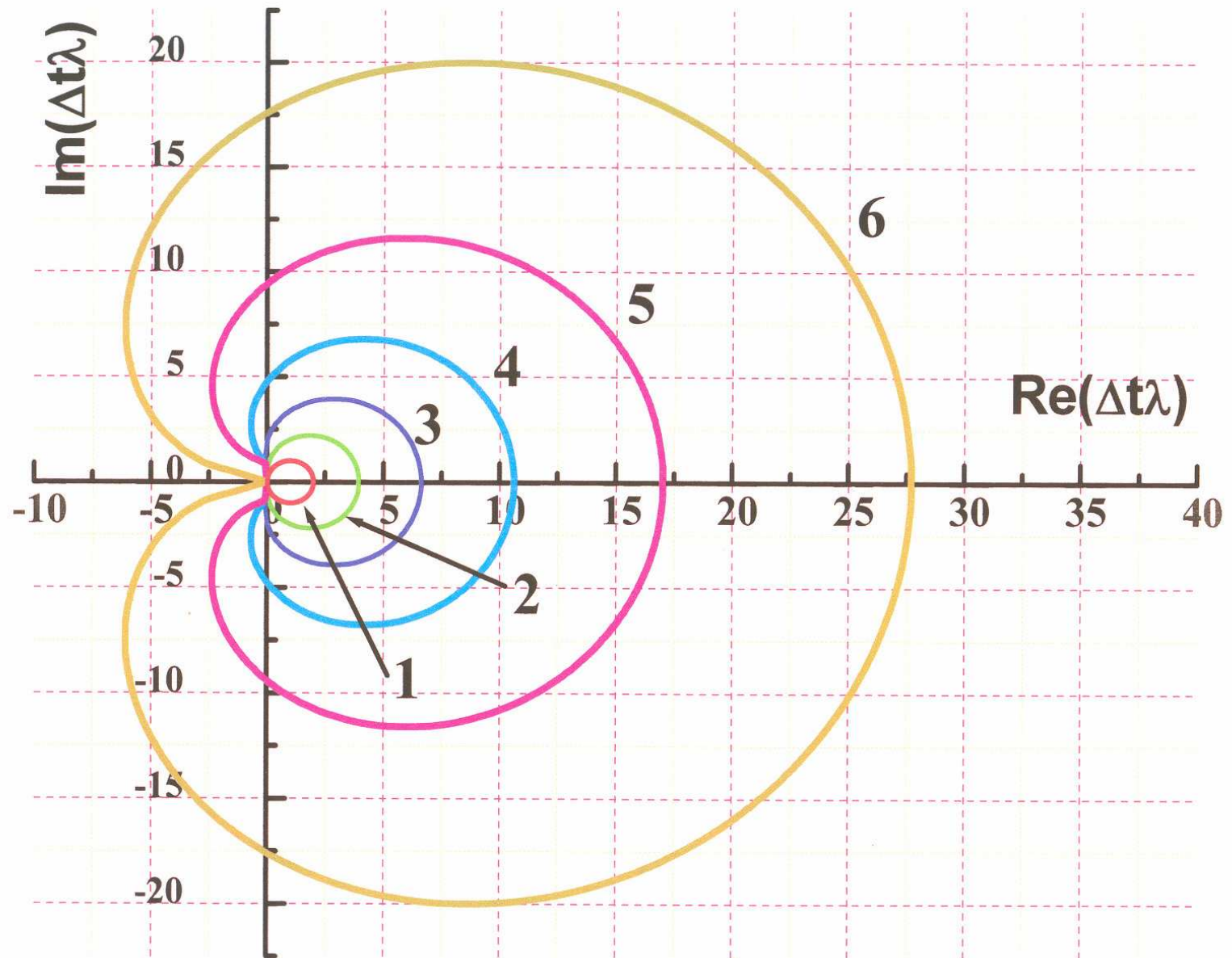
$m = 3, 4, 5, 6 :$

\mathcal{S}_m *contains the negative axis* , but not the entire negative half plane.

These methods are called *A(α)-stable* .

$m \geq 7 :$

These methods are *unstable* , even for solving $u'(t) = 0$!



Stability region of Backward Differentiation Formulas.

Collocation at 2 Gauss Points.

The 2-point *Gauss collocation method* for taking a time step for the IVP

$$u'(t) = f(u(t)) , \quad u(0) = u_0 ,$$

is defined by finding a *local polynomial* $p \in \mathbb{P}_2$ that satisfies

$$p(t_j) = u_j ,$$

and

$$p'(x_{j,i}) = f(p(x_{j,i})) , \quad i = 1, 2 , \quad (\textit{collocation}) ,$$

where

$$x_{j,i} = \frac{t_j + t_{j+1}}{2} \pm \frac{\Delta t \sqrt{3}}{6} ,$$

and then setting

$$u_{j+1} = p(t_{j+1}) .$$

Applied to the *model equation* $u'(t) = \lambda u(t)$ this gives

$$u_{j+1} = \frac{1 + \Delta t \lambda + \frac{1}{12}(\Delta t \lambda)^2}{1 - \Delta t \lambda + \frac{1}{12}(\Delta t \lambda)^2} u_j \equiv z(\Delta t \lambda) u_j .$$

It can be shown that the *stability region*

$$\mathcal{S} \equiv \left\{ \Delta t \lambda : \frac{1 + \Delta t \lambda + \frac{1}{12} (\Delta t \lambda)^2}{1 - \Delta t \lambda + \frac{1}{12} (\Delta t \lambda)^2} \leq 1 \right\},$$

is the *entire negative half plane*.

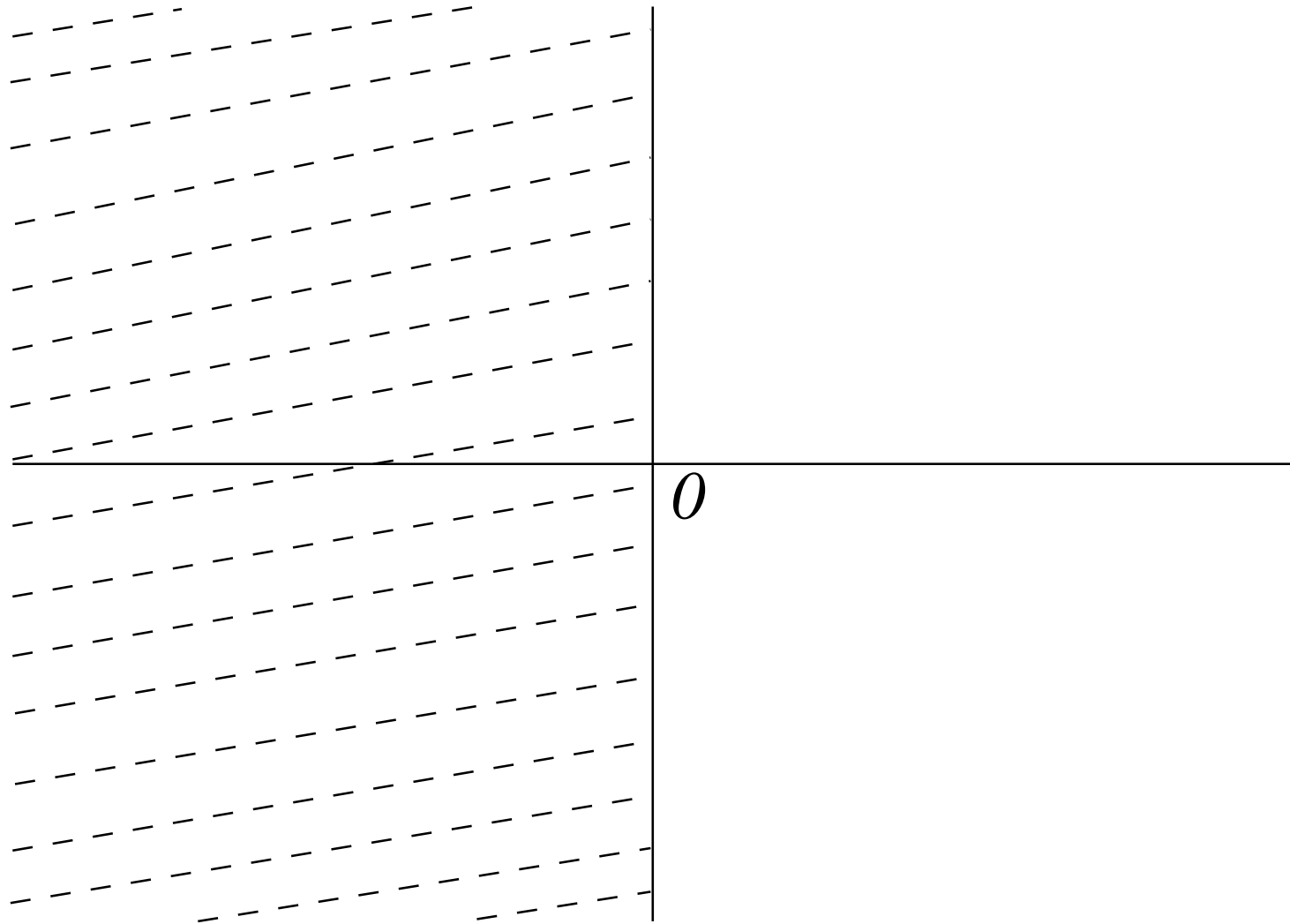
All Gauss collocation methods have this property and thus *are A-stable*.

However,

$$\lim_{\lambda \rightarrow -\infty} z(\Delta t \lambda) = 1,$$

so that the methods lead to *slow decay* for *stiff problems*.

Complex $(\Delta t \cdot \lambda)$ -plane



Stability region of the collocation method

BOUNDARY VALUE PROBLEMS IN ODEs

EXAMPLE: The *boundary value problem* (BVP)

$$y''(x) - y(x) = -5 \sin(2x) , \quad x \in [0, \pi] ,$$

$$y(0) = 0 , \quad y(\pi) = 0 ,$$

has the exact (and unique) *solution*

$$y(x) = \sin(2x) .$$

- This BVP is a simple example of problems from science and engineering.
- Usually it is difficult or impossible to find an exact solution.
- In such cases numerical techniques can be used.

Partition $[0, \pi]$ into a *grid* or *mesh* :

$$0 = x_0 < x_1 < x_2 < \cdots < x_N = \pi ,$$

where

$$x_j = jh , \quad (j = 0, 1, 2, \cdots, N) , \quad h = \frac{\pi}{N} .$$

We want to find approximations u_j to $y(x_j)$, $j = 0, 1, 2, \cdots, N$.

A *finite difference approximation* to $y''(x_j)$ is given by

$$y''(x_j) \approx \frac{\frac{y_{j+1} - y_j}{h} - \frac{y_j - y_{j-1}}{h}}{h} = \frac{y_{j+1} - 2y_j + y_{j-1}}{h^2} ,$$

where

$$y_j \equiv y(x_j) .$$

We want to find approximations u_j to $y(x_j)$, $j = 0, 1, 2, \dots, N$.

The u_j are computed by solving the *finite difference equations*:

$$\begin{aligned} u_0 &= 0, \\ \frac{u_2 - 2u_1 + u_0}{h^2} - u_1 &= -5 \sin(2x_1), \\ \frac{u_3 - 2u_2 + u_1}{h^2} - u_2 &= -5 \sin(2x_2), \\ &\cdot \\ &\cdot \\ &\cdot \\ \frac{u_N - 2u_{N-1} + u_{N-2}}{h^2} - u_{N-1} &= -5 \sin(2x_{N-1}), \\ u_N &= 0. \end{aligned}$$

We found that :

The *finite difference equations* can be written in *matrix form* as

$$\mathbf{A}_h \mathbf{u}_h = \mathbf{f}_h ,$$

where

$$\mathbf{A}_h = \text{diag} \left[\frac{1}{h^2} , - \left(1 + \frac{2}{h^2} \right) , \frac{1}{h^2} \right] ,$$

$$\mathbf{u}_h \equiv (u_1 , u_2 , \dots , u_{N-1})^T ,$$

and

$$\mathbf{f}_h \equiv -5 (\sin(2x_1) , \sin(2x_2) , \dots , \sin(2x_{N-1}))^T .$$

QUESTIONS:

- How to *solve the linear systems efficiently* , especially when N is large ?
- How to *approximate derivatives* and find the *error* in the approximation ?
- What is *the actual error* after solving the system,

i.e. , what is

$$\max_j | u_j - y(x_j) | \quad ?$$

(assuming exact arithmetic)

- How to *solve the linear systems efficiently* , especially when N is large ?

ANSWER :

The matrix is *tridiagonal* .

Thus the linear system can be solved by the specialized Gauss elimination algorithm for tridiagonal systems.

- How to *approximate derivatives* and find the *error* in the approximation ?

ANSWER : As done earlier, the *local discretization error* is

$$\begin{aligned}
 \tau_j &\equiv \frac{y_{j+1} - 2y_j + y_{j-1}}{h^2} - y_j'' \\
 &= \frac{1}{h^2} \left(y_j + hy_j' + \frac{h^2}{2}y_j'' + \frac{h^3}{6}y_j''' + \frac{h^4}{24}y_j''''(\zeta_1) \right. \\
 &\quad \left. - 2y_j \right. \\
 &\quad \left. + y_j - hy_j' + \frac{h^2}{2}y_j'' - \frac{h^3}{6}y_j''' + \frac{h^4}{24}y_j''''(\zeta_2) \right) - y_j'' \\
 &= \frac{h^2}{24} \left(y_j''''(\zeta_1) + y_j''''(\zeta_2) \right) \\
 &= \frac{h^2}{12} y_j''''(\eta_j) , \quad \text{for some } \eta_j \in (x_{j-1}, x_{j+1}) ,
 \end{aligned}$$

using Taylor and Intermediate Value Theorem, assuming y'''' is continuous.

We found that

$$\tau_j \equiv \frac{y_{j+1} - 2y_j + y_{j-1}}{h^2} - y_j'' = \frac{h^2}{12} y''''(\eta_j) .$$

In our BVP, we have

$$y(x) = \sin(2x) , \quad \text{and} \quad y''''(x) = 16 \sin(2x) .$$

Thus $|y''''(x)| \leq 16$, and

$$|\tau_j| \leq \frac{16}{12} h^2 = \frac{4}{3} h^2 , \quad j = 1, 2, \dots, N - 1 .$$

- What is the *actual error* after solving the system ?

i.e., what is

$$\max_j | u_j - y(x_j) | \text{ ?}$$

ANSWER :

For this, we will use the *Banach Lemma* .

We already showed that

$$|\tau_j| \equiv \left| \frac{(y_{j-1} - 2y_j + y_{j+1}))}{h^2} - y_j'' \right| \leq \frac{4h^2}{3}, \quad j = 1, 2, \dots, N-1.$$

Now

$$\begin{aligned} & \frac{1}{h^2}y_{j-1} - \left(1 + \frac{2}{h^2}\right)y_j + \frac{1}{h^2}y_{j+1} \\ &= \frac{(y_{j+1} - 2y_j + y_{j-1}))}{h^2} - y_j \\ &= y_j'' + \tau_j - y_j \\ &= \tau_j - 5 \sin(2x_j). \end{aligned}$$

Thus if we define

$$\mathbf{y}_h \equiv (y_1, y_2, \dots, y_{N-1})^T, \quad \text{and} \quad \boldsymbol{\tau}_h \equiv (\tau_1, \tau_2, \dots, \tau_{N-1})^T,$$

then

$$\mathbf{A}_h \mathbf{y}_h = \boldsymbol{\tau}_h + \mathbf{f}_h.$$

We found that

$$\mathbf{A}_h \mathbf{y}_h = \boldsymbol{\tau}_h + \mathbf{f}_h .$$

Since

$$\mathbf{A}_h \mathbf{u}_h = \mathbf{f}_h ,$$

it follows from subtraction that

$$\mathbf{A}_h (\mathbf{y}_h - \mathbf{u}_h) = \boldsymbol{\tau}_h .$$

We found that

$$\mathbf{A}_h(\mathbf{y}_h - \mathbf{u}_h) = \boldsymbol{\tau}_h .$$

Thus if we can show that \mathbf{A}_h has an *inverse* and that

$$\| \mathbf{A}_h^{-1} \|_{\infty} \leq K ,$$

for some constant K that does not depend on h , then

$$\begin{aligned} \| \mathbf{y}_h - \mathbf{u}_h \|_{\infty} &= \| \mathbf{A}_h^{-1} \boldsymbol{\tau}_h \|_{\infty} \\ &\leq \| \mathbf{A}_h^{-1} \|_{\infty} \| \boldsymbol{\tau}_h \|_{\infty} \\ &\leq K \frac{4h^2}{3} . \end{aligned}$$

Now

$$\begin{aligned}
 \mathbf{A}_h &= \begin{pmatrix} -1 - \frac{2}{h^2} & \frac{1}{h^2} & & & \\ \frac{1}{h^2} & -1 - \frac{2}{h^2} & \frac{1}{h^2} & & \\ & \cdot & \cdot & \cdot & \\ & & \frac{1}{h^2} & -1 - \frac{2}{h^2} & \frac{1}{h^2} \\ & & & \frac{1}{h^2} & -1 - \frac{2}{h^2} \end{pmatrix} \\
 &= -\frac{h^2 + 2}{h^2} \mathbf{I}_h + \frac{1}{h^2} \begin{pmatrix} 0 & 1 & & & \\ 1 & 0 & 1 & & \\ & \cdot & \cdot & \cdot & \\ & & & 1 & 0 \end{pmatrix} \\
 &= -\frac{h^2 + 2}{h^2} \left[\mathbf{I}_h - \frac{h^2}{h^2 + 2} \frac{1}{h^2} \begin{pmatrix} 0 & 1 & & & \\ 1 & 0 & 1 & & \\ & \cdot & \cdot & \cdot & \\ & & & 1 & 0 \end{pmatrix} \right] \\
 &= -\frac{h^2 + 2}{h^2} \left[\mathbf{I}_h - \frac{1}{h^2 + 2} \begin{pmatrix} 0 & 1 & & & \\ 1 & 0 & 1 & & \\ & \cdot & \cdot & \cdot & \\ & & & 1 & 0 \end{pmatrix} \right].
 \end{aligned}$$

We have

$$\begin{aligned} \mathbf{A}_h &= -\frac{h^2 + 2}{h^2} \left[\mathbf{I}_h - \frac{1}{h^2 + 2} \begin{pmatrix} 0 & 1 & & \\ 1 & 0 & 1 & \\ & \cdot & \cdot & \cdot \\ & & 1 & 0 \end{pmatrix} \right] \\ &= -\frac{h^2 + 2}{h^2} (\mathbf{I}_h + \mathbf{B}_h), \end{aligned}$$

where \mathbf{I}_h is the identity matrix and

$$\mathbf{B}_h \equiv \frac{-1}{h^2 + 2} \begin{pmatrix} 0 & 1 & & \\ 1 & 0 & 1 & \\ & \cdot & \cdot & \cdot \\ & & 1 & 0 \end{pmatrix}.$$

Since

$$\|\mathbf{B}_h\|_\infty = \frac{2}{h^2 + 2} < 1,$$

it follows by the *Banach Lemma* that $(\mathbf{I}_h + \mathbf{B}_h)^{-1}$ exists and that

$$\|(\mathbf{I}_h + \mathbf{B}_h)^{-1}\|_\infty \leq \frac{1}{1 - \frac{2}{h^2 + 2}} = \frac{h^2 + 2}{h^2}.$$

We have

$$\mathbf{A}_h = -\frac{h^2 + 2}{h^2} (\mathbf{I}_h + \mathbf{B}_h) ,$$

and

$$\| (\mathbf{I}_h + \mathbf{B}_h)^{-1} \|_\infty \leq \frac{h^2 + 2}{h^2} .$$

Hence

$$\| \mathbf{A}_h^{-1} \|_\infty = \left\| \frac{-h^2}{h^2 + 2} (\mathbf{I}_h + \mathbf{B}_h)^{-1} \right\|_\infty \leq \frac{h^2}{h^2 + 2} \frac{h^2 + 2}{h^2} = 1 .$$

Thus $K = 1$, and

$$\| \mathbf{y}_h - \mathbf{u}_h \|_\infty \leq \frac{4h^2}{3} .$$

A Nonlinear Boundary Value Problem.

Consider the *Gelfand-Bratu* problem

$$u''(x) + \lambda e^{u(x)} = 0, \quad x \in [0, 1],$$

$$u(0) = 0, \quad u(1) = 0, \quad \lambda \text{ is a parameter,}$$

and its *finite difference approximation*

$$g_1(\mathbf{u}) \equiv \frac{u_2 - 2u_1}{h^2} + \lambda e^{u_1} = 0,$$

$$g_2(\mathbf{u}) \equiv \frac{u_3 - 2u_2 + u_1}{h^2} + \lambda e^{u_2} = 0,$$

.

.

$$g_{N-2}(\mathbf{u}) \equiv \frac{u_{N-1} - 2u_{N-2} + u_{N-3}}{h^2} + \lambda e^{u_{N-2}} = 0,$$

$$g_{N-1}(\mathbf{u}) \equiv \frac{-2u_{N-1} + u_{N-2}}{h^2} + \lambda e^{u_{N-1}} = 0,$$

where $\mathbf{u} \equiv (u_1, u_2, \dots, u_{N-1})^T$.

If we let

$$\mathbf{G}(\mathbf{u}) \equiv (g_1(\mathbf{u}), g_2(\mathbf{u}), \dots, g_{N-1}(\mathbf{u}))^T,$$

and

$$\mathbf{0} \equiv (0, 0, \dots, 0)^T \in \mathbf{R}^{N-1},$$

then these equations can be compactly written as

$$\mathbf{G}(\mathbf{u}) = \mathbf{0}.$$

The Jacobian matrix is an $N - 1$ by $N - 1$ tridiagonal matrix :

$$\mathbf{G}'(\mathbf{u}) = \begin{pmatrix} -\frac{2}{h^2} + \lambda e^{u_1} & & & & \\ \frac{1}{h^2} & -\frac{2}{h^2} + \lambda e^{u_2} & & & \\ & \cdot & \frac{1}{h^2} & & \\ & & \cdot & \cdot & \\ & & \frac{1}{h^2} & -\frac{2}{h^2} + \lambda e^{u_{N-1}} & \end{pmatrix}.$$

Each Newton iteration for solving the nonlinear system

$$\mathbf{G}(\mathbf{u}) = \mathbf{0} ,$$

consists of solving the tridiagonal system

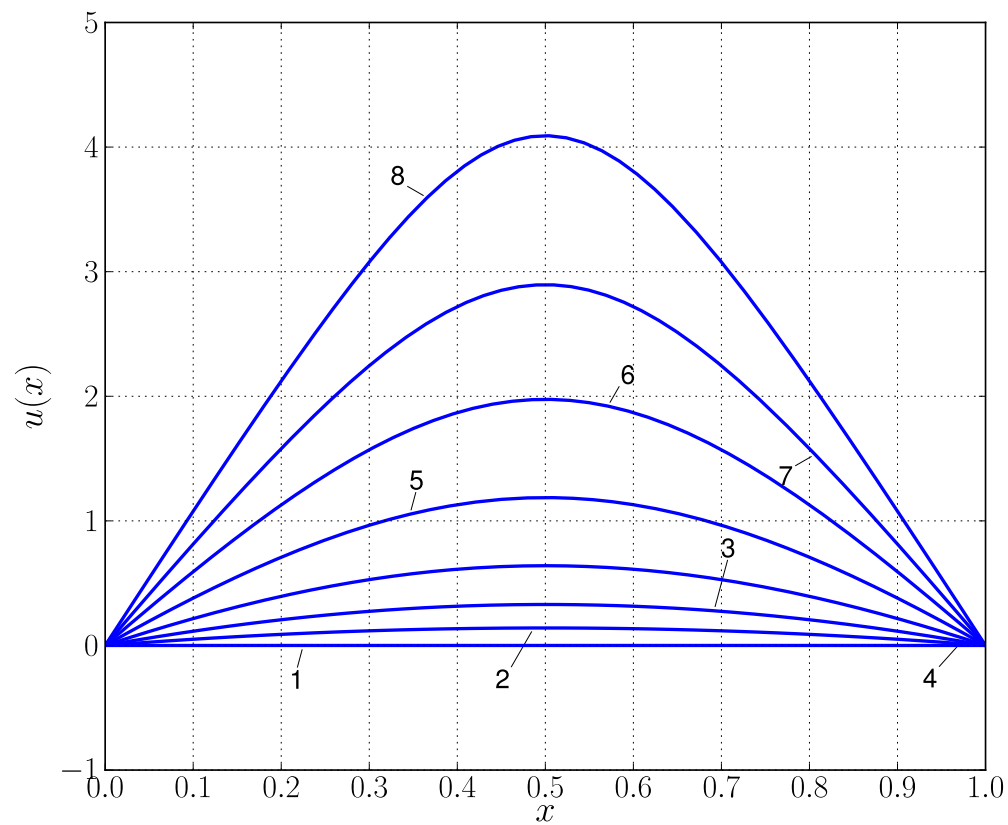
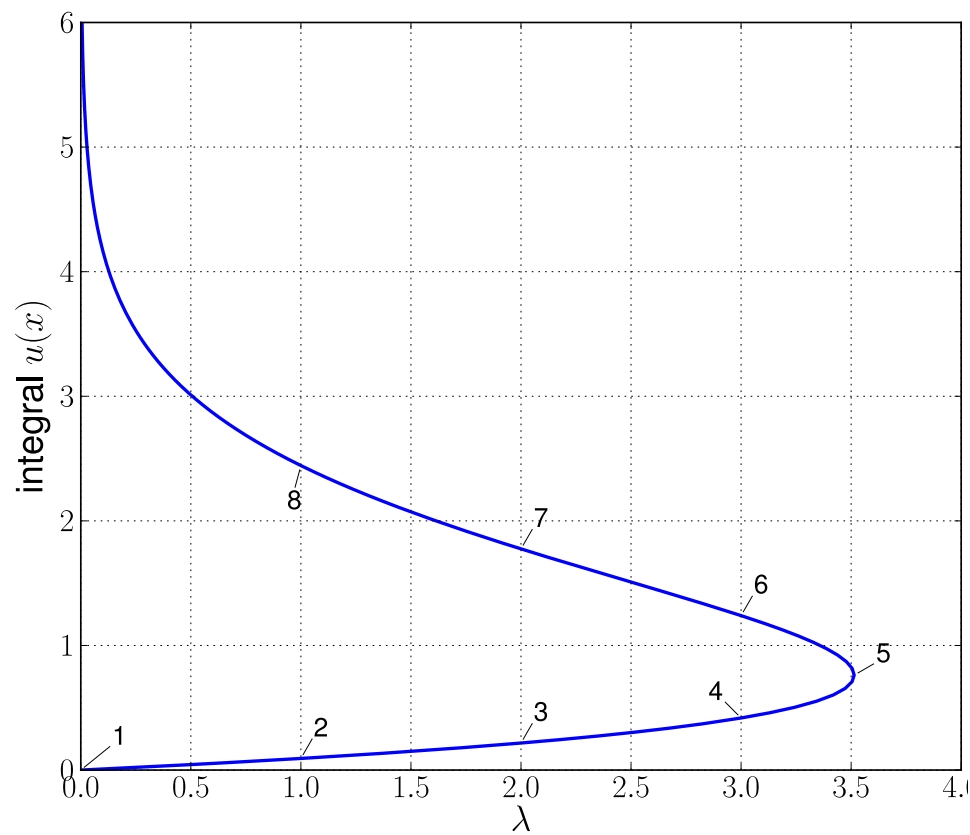
$$\begin{pmatrix} -\frac{2}{h^2} + \lambda e^{u_1^{(k)}} & & & & \\ & \frac{1}{h^2} & & & \\ & \frac{1}{h^2} & -\frac{2}{h^2} + \lambda e^{u_2^{(k)}} & & \\ & & \cdot & \frac{1}{h^2} & \\ & & \frac{1}{h^2} & & -\frac{2}{h^2} + \lambda e^{u_{N-1}^{(k)}} \end{pmatrix} \Delta \mathbf{u}^{(k)} = -\mathbf{G}(\mathbf{u}^{(k)}) ,$$

where

$$\Delta \mathbf{u}^{(k)} \equiv (\Delta u_1^{(k)} , \Delta u_2^{(k)} , \dots , \Delta u_{N-1}^{(k)})^T ,$$

and updating

$$\mathbf{u}^{(k+1)} = \mathbf{u}^{(k)} + \Delta \mathbf{u}^{(k)} .$$



Solutions of the Gelfand-Bratu equations for different values of λ .

DIFFUSION PROBLEMS

Here we look at some simple *parabolic partial differential equations*.

The simplest is the *linear diffusion equation* or *heat equation* :

$$u_t(x, t) = u_{xx}(x, t) , \quad x \in [0, 1] , \quad t \geq 0 ,$$

$$u(x, 0) = g(x) ,$$

$$u(0, t) = u(1, t) = 0 .$$

This equation governs, for example, the temperature in an insulated rod of which the endpoints are kept at the constant temperature zero, and in which the initial temperature distribution is $g(x)$.

First *discretize in space* :

$$u'_j(t) = \frac{u_{j-1}(t) - 2u_j(t) + u_{j+1}(t)}{\Delta x^2} ,$$

$$u_j(0) = g(x_j) ,$$

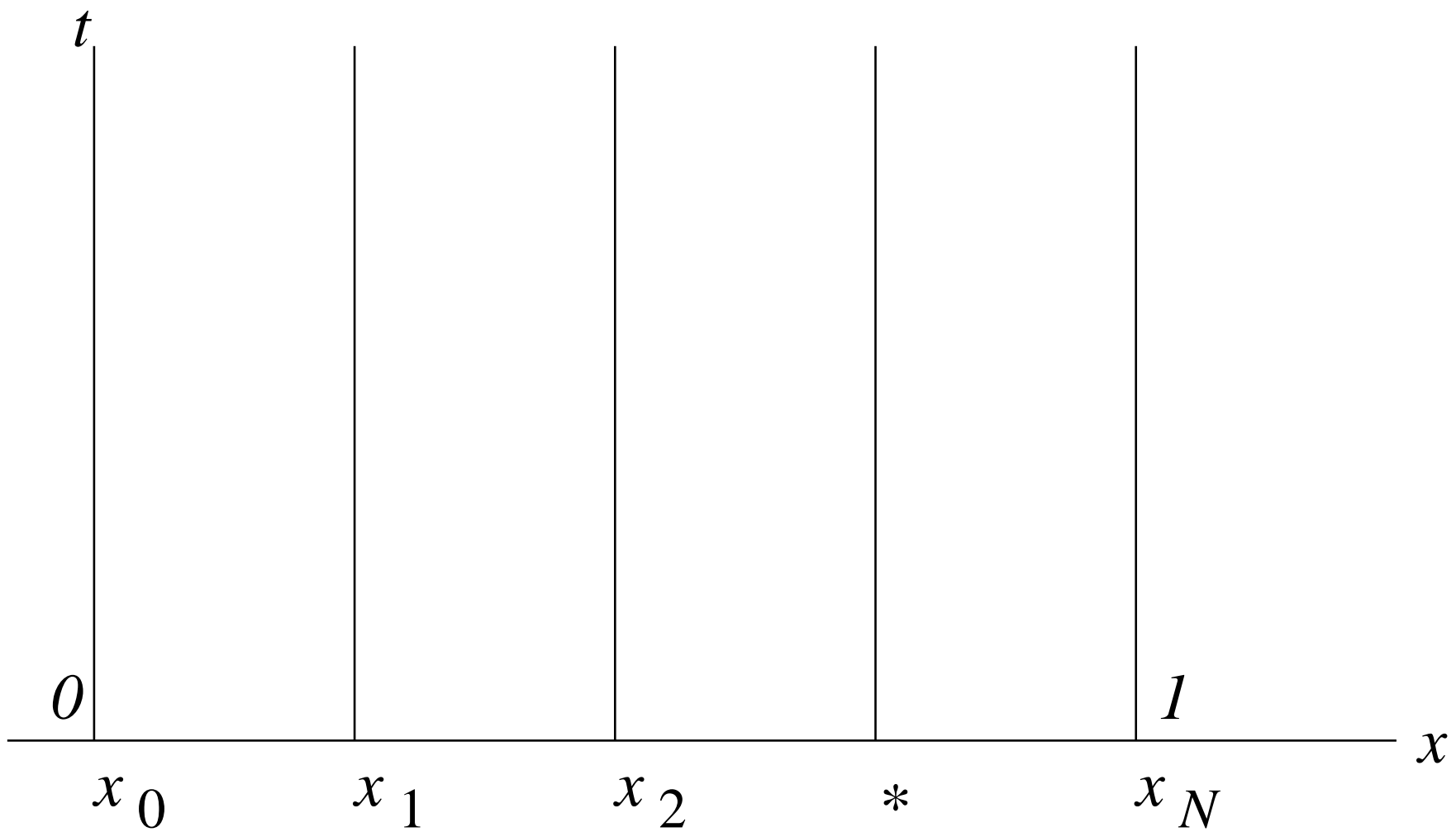
$$u_0(t) = u_N(t) = 0 ,$$

where we have introduced the *notation*

$$u_j(t) \equiv u(x_j, t) ,$$

and where ' denotes differentiation with respect to t .

These space-discretized equations represents a system of $N - 1$ coupled *ordinary* differential equations.



In *matrix-vector notation* we can write the space-discretized equations as

$$\mathbf{u}'(t) = \frac{1}{\Delta x^2} \mathbf{D} \mathbf{u}(t) ,$$

where

$$\mathbf{D} \equiv \begin{pmatrix} -2 & 1 & & & \\ 1 & -2 & 1 & & \\ & \cdot & \cdot & \cdot & \\ & & 1 & -2 & 1 \\ & & & 1 & -2 \end{pmatrix} ,$$

and

$$\mathbf{u} \equiv \begin{pmatrix} u_1 \\ u_2 \\ \cdot \\ u_{N-2} \\ u_{N-1} \end{pmatrix} .$$

Now *discretize in time* :

Often used is the *Trapezoidal rule* :

$$\frac{\mathbf{u}^{k+1} - \mathbf{u}^k}{\Delta t} = \frac{1}{2\Delta x^2} \mathbf{D} \{ \mathbf{u}^{k+1} + \mathbf{u}^k \} ,$$

where

$$\mathbf{u}^k \equiv \begin{pmatrix} u_1^k \\ u_2^k \\ \cdot \\ u_{N-1}^k \end{pmatrix} ,$$

and

u_j^k approximates $u(x_j, t^k)$.

Assume that the solution has been computed up to time t^k .

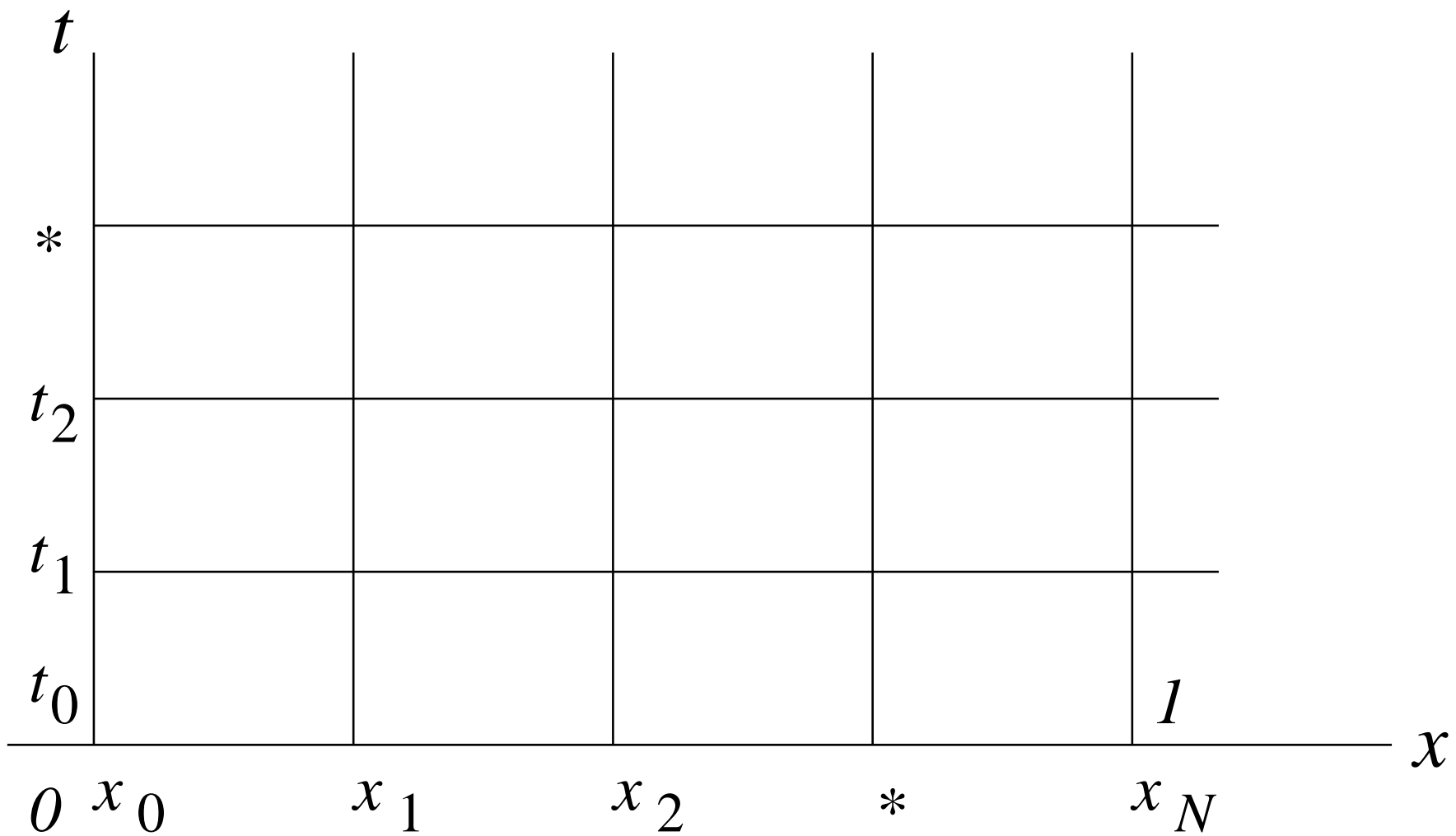
Thus \mathbf{u}^k *is known*, and we want to *solve for* \mathbf{u}^{k+1} .

Rewrite the above equation as

$$\left(I - \frac{\Delta t}{2\Delta x^2} \mathbf{D}\right) \mathbf{u}^{k+1} = \left(I + \frac{\Delta t}{2\Delta x^2} \mathbf{D}\right) \mathbf{u}^k .$$

Thus to take a step in time we have to solve a *tridiagonal* linear system.

This method is also known as the *Crank-Nicolson scheme*.



REMARK :

- We can also use an *explicit* method in time, for example explicit Euler.
- But this can be a *bad choice* because the ODE system is *stiff*.
- The time step Δt may have to be *very small* to have stability.

For the simple ODE (2) we can *demonstrate the stiffness* analytically.

In fact, we can explicitly *compute the eigenvalues* of the matrix

$$\frac{1}{\Delta x^2} \mathbf{D} ,$$

as follows:

An *eigenvalue-eigenvector* pair λ, \mathbf{v} satisfies

$$\frac{1}{\Delta x^2} \mathbf{D} \mathbf{v} = \lambda \mathbf{v} ,$$

that is,

$$\frac{1}{\Delta x^2} (v_{j-1} - 2v_j + v_{j+1}) = \lambda v_j , \quad j = 1, 2, \dots, N-1 , \quad v_0 = v_N = 0 .$$

We had the *difference equation*

$$\frac{1}{\Delta x^2}(v_{j-1} - 2v_j + v_{j+1}) = \lambda v_j .$$

Try a solution of the form $v_j = z^j$.

This gives the *characteristic equation*

$$z^2 - (2 + \Delta x^2 \lambda)z + 1 = 0 ,$$

or

$$\lambda = \frac{z + z^{-1} - 2}{\Delta x^2} .$$

This equation has *zeroes*

$$z = z_1 \quad \text{and} \quad z = z_1^{-1} .$$

The *general solution* of the difference equation then has the form

$$v_j = c_1 z_1^j + c_2 z_1^{-j} .$$

From the *first boundary condition* we have

$$v_0 = 0 \quad \Rightarrow \quad c_1 + c_2 = 0 .$$

Thus we can take

$$c_1 = c \quad \text{and} \quad c_2 = -c .$$

Then

$$v_j = c (z_1^j - z_1^{-j}) .$$

From the *second boundary condition* we now find

$$v_N = 0 \quad \Rightarrow \quad c (z_1^N - z_1^{-N}) = 0 ,$$

from which

$$z_1^{2N} = 1 \quad \Rightarrow \quad z_1 = e^{\frac{k2\pi i}{2N}} .$$

The *eigenvalues* are therefore

$$\begin{aligned}\lambda_k &= \frac{z + z^{-1} - 2}{\Delta x^2} \\ &= \frac{e^{\frac{k2\pi i}{2N}} + e^{-\frac{k2\pi i}{2N}} - 2}{\Delta x^2} \\ &= \frac{2 \cos\left(\frac{k2\pi}{2N}\right) - 2}{\Delta x^2} \\ &= \frac{2\left(\cos\left(\frac{k2\pi}{2N}\right) - 1\right)}{\Delta x^2} \\ &= -\frac{4}{\Delta x^2} \sin^2\left(\frac{k\pi}{2N}\right), \quad k = 1, 2, \dots, N - 1.\end{aligned}$$

The eigenvalue with *largest negative real part* is

$$\lambda_{N-1} = -\frac{4}{\Delta x^2} \sin^2\left(\frac{(N-1)\pi}{2N}\right),$$

which for large N behaves like

$$\lambda_{N-1} \approx \lambda^* \equiv -\frac{4}{\Delta x^2}.$$

Thus the system is *stiff* if Δx is small .

For example, to make the explicit Euler method stable we need to take the timestep Δt so that $\Delta t\lambda^*$ lies in the circle of radius 1 centered at -1 , *i.e.*, we must take

$$\Delta t < \frac{1}{2}\Delta x^2.$$

Using explicit Euler is often *not a good idea* .

Nonlinear Diffusion Equations.

An example of a nonlinear diffusion equation is *the Fisher equation*

$$u_t(x, t) = u_{xx}(x, t) + \lambda u(x, t) (1 - u(x, t)) ,$$

for

$$x \in [0, 1] , \quad t \geq 0 ,$$

with

$$u(x, 0) = g(x) , \quad u(0, t) = u(1, t) = 0 .$$

This is a simple model of *population growth* with *diffusion* and with *maximal sustainable population* 1.

Another example is *the time-dependent Gelfand-Bratu equation*

$$u_t(x, t) = u_{xx}(x, t) + \lambda e^{u(x, t)},$$

for

$$x \in [0, 1], \quad t \geq 0,$$

with

$$u(x, 0) = g(x), \quad u(0, t) = u(1, t) = 0,$$

for which we have already considered the *stationary equations*

$$u_{xx}(x) + \lambda e^{u(x)} = 0, \quad x \in [0, 1],$$

with

$$u(0) = u(1) = 0.$$

We illustrate the numerical solution procedure for the *general equation*

$$u_t(x, t) = u_{xx}(x, t) + f(u(x, t)) , \quad x \in [0, 1] , \quad t \geq 0 ,$$

$$u(x, 0) = g(x) ,$$

$$u(0, t) = u(1, t) = 0 ,$$

where

$$f(u) = \lambda u (1 - u) \quad \text{for the } \textit{Fisher equation} ,$$

and

$$f(u) = \lambda e^u \quad \text{for the } \textit{Gelfand-Bratu equation} .$$

We approximate this equation as follows :

First *discretize in space* to get a *system of ODEs* :

$$u'_j(t) = \frac{u_{j-1}(t) - 2u_j(t) + u_{j+1}(t)}{\Delta x^2} + f(u_j(t)) ,$$

for $j = 1, 2, \dots, N - 1$, with

$$u_j(0) = g(x_j) ,$$

$$u_0(t) = u_N(t) = 0 .$$

Then *discretize in time* using *Implicit Euler* :

$$\frac{u_j^{k+1} - u_j^k}{\Delta t} = \frac{u_{j-1}^{k+1} - 2u_j^{k+1} + u_{j+1}^{k+1}}{\Delta x^2} + f(u_j^{k+1}) .$$

Rewrite these equations as

$$F_j^{k+1} \equiv u_j^{k+1} - u_j^k - \frac{\Delta t}{\Delta x^2} (u_{j-1}^{k+1} - 2u_j^{k+1} + u_{j+1}^{k+1}) - \Delta t f(u_j^{k+1}) = 0 ,$$

for $j = 1, 2, \dots, N - 1$,

with

$$u_0^{k+1} = 0 \quad \text{and} \quad u_N^{k+1} = 0 .$$

We can assume that the solution has been computed up to time t^k ,

i.e., the u_j^k *are known* and we must *solve for* the u_j^{k+1} .

Since *the equations are nonlinear* we use Newton's method.

As *initial approximation* to u_j^{k+1} in Newton's method use

$$(u_j^{k+1})^{(0)} = u_j^k, \quad j = 1, 2, \dots, N-1.$$

Each *Newton iteration* then consists of solving a linear *tridiagonal system*

$$\mathbf{T}^{k+1,(\nu)} \Delta \mathbf{u}^{k+1,(\nu)} = -\mathbf{F}^{k+1,(\nu)},$$

where

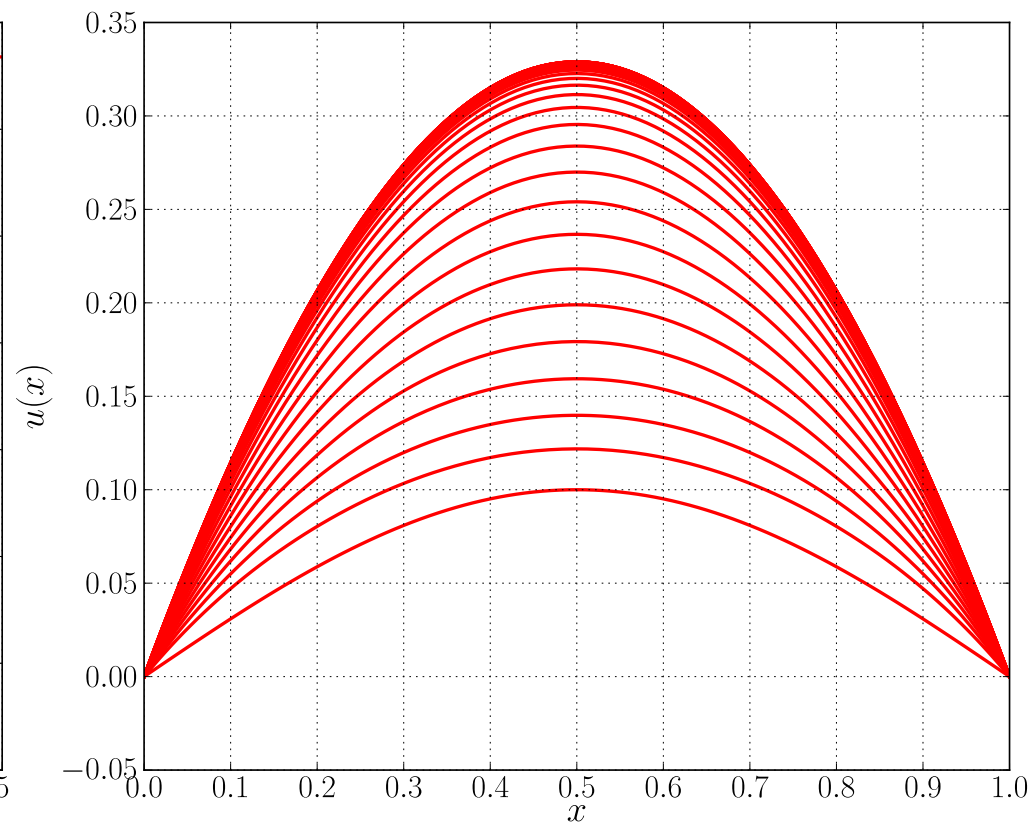
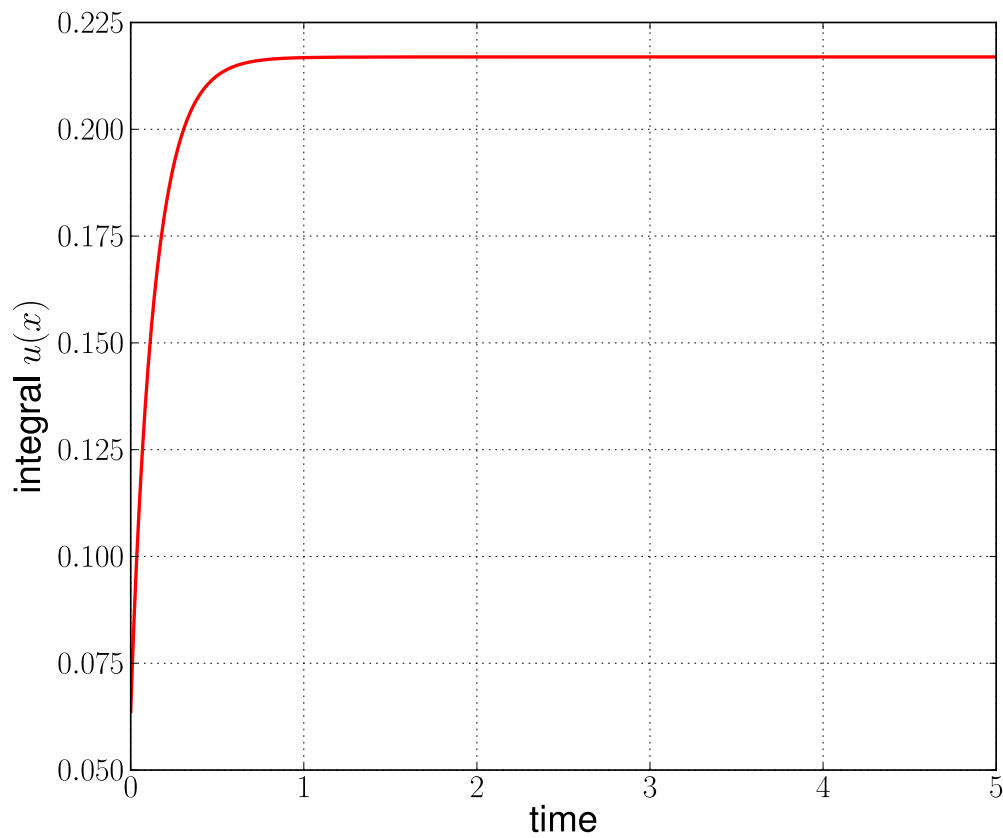
$$\mathbf{T}^{k+1,(\nu)} = \begin{pmatrix} 1 + 2\frac{\Delta t}{\Delta x^2} - \Delta t f_u(u_1^{k+1,(\nu)}) & -\frac{\Delta t}{\Delta x^2} & & & \\ -\frac{\Delta t}{\Delta x^2} & 1 + 2\frac{\Delta t}{\Delta x^2} - \Delta t f_u(u_2^{k+1,(\nu)}) & & & \\ \cdot & & \cdot & & \\ \cdot & & & \cdot & \\ -\frac{\Delta t}{\Delta x^2} & & & & 1 + 2\frac{\Delta t}{\Delta x^2} - \Delta t f_u(u_{N-1}^{k+1,(\nu)}) \end{pmatrix}$$

and

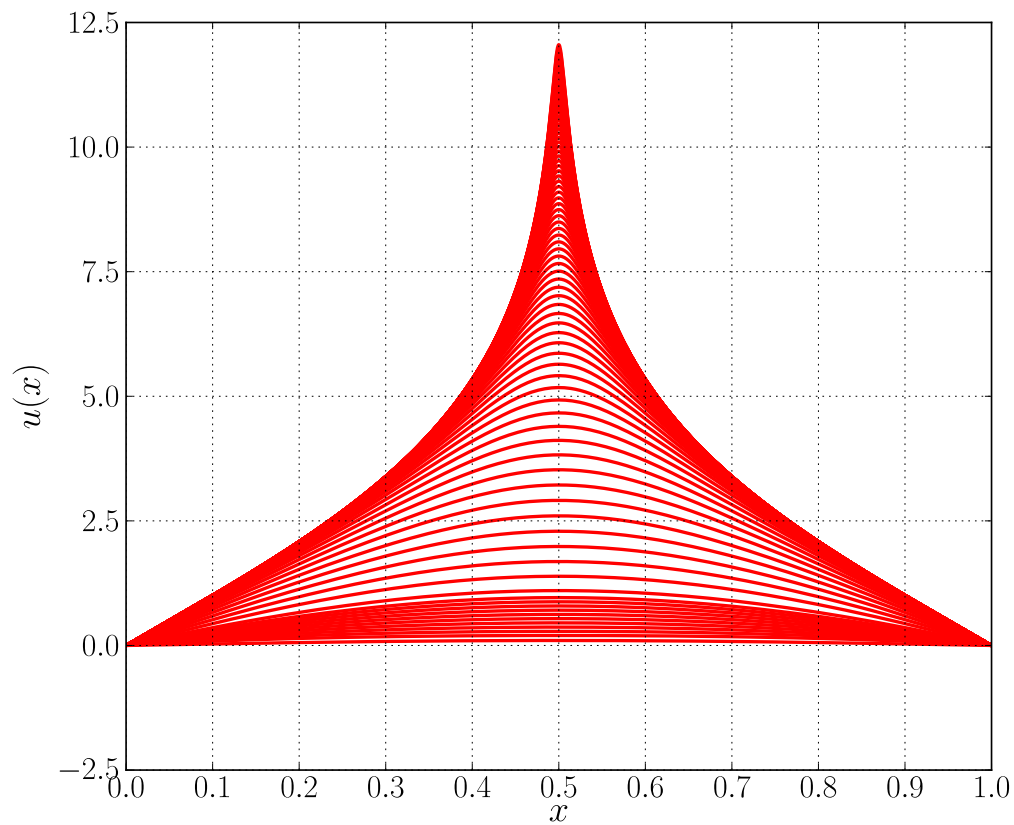
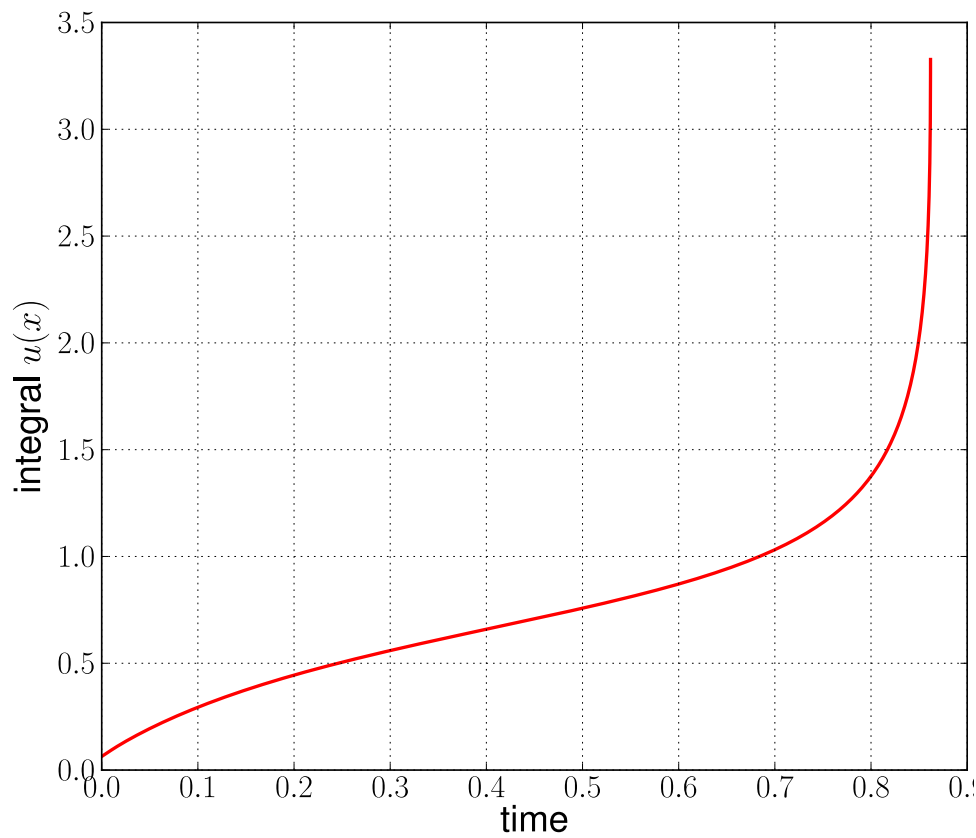
$$\Delta \mathbf{u}^{k+1,(\nu)} = \begin{pmatrix} \Delta u_1^{k+1,(\nu)} \\ \Delta u_2^{k+1,(\nu)} \\ \cdot \\ \Delta u_{N-1}^{k+1,(\nu)} \end{pmatrix}, \quad \mathbf{F}^{k+1,(\nu)} = \begin{pmatrix} F_1^{k+1,(\nu)} \\ F_2^{k+1,(\nu)} \\ \cdot \\ F_{N-1}^{k+1,(\nu)} \end{pmatrix}.$$

Then set the *next approximation* to the solution at time $t = t^{k+1}$ equal to

$$\mathbf{u}^{k+1,(\nu+1)} = \mathbf{u}^{k+1,(\nu)} + \Delta \mathbf{u}^{k+1,(\nu)}.$$



Time-evolution of solutions of the Gelfand-Bratu equations for $\lambda = 2$.



Time-evolution of solutions of the Gelfand-Bratu equations for $\lambda = 4$.