

## 1.1 ELEMENTY TEORII BŁĘDÓW

Metody numeryczne dają nam możliwość rozwiązywania pewnych zagadnień w sposób przybliżony wtedy, gdy dokładne metody nie mogą być stosowane lub nie są znane.

Każde rozwiązywanie numeryczne wiąże się z popełnianiem błędów obliczeń. Błędy te mogą wynikać z różnych przyczyn.

Jedną z nich jest podawanie danych w sposób przybliżony - podczas pomiarów lub doświadczeń popełniamy nieścisłości związane np.: z dokładnością przyrządu pomiarowego.

Błędy danych mogą silnie wpływać na wyniki obliczeń, ale nie zawsze. Czasami można podać tzw.: wskaźniki uwarunkowania, które "przenoszą" błędy danych na błędy obliczeń końcowych. Podamy takie wskaźniki dla błędów funkcji jednej i wielu zmiennych.

Drugą przyczyną jest dokładność samego algorytmu stosowanego do obliczeń. Może się zdarzyć, że dokładny wzór np.: rekurencyjny nie nadaje się do obliczeń numerycznych, będziemy o nim mówić, że nie jest stabilny- podamy przykład takiego algorytmu. Nawet kolejność działań na liczbach przybliżonych może mieć znaczenie i wpływać na wynik, a niektóre działania np.: odejmowanie liczb przybliżonych bliskich daje czasami zaskakująco duże błędy.

Trzeba również pamiętać o błędach maszynowych, wynikające z reprezentacji liczb w komputerze.

Wprawdzie dzisiejsze komputery są bardzo dokładne, to jednak nakładanie się tych wszystkich błędów może dawać absurdalne wyniki. Będziemy ilustrować takie sytuacje.

## 1.2 BŁĘDY BEZWZGLĘDNE I WZGLĘDNE

Dużymi literami  $A, B, C, \dots$  będziemy oznaczać liczby dokładne, a małymi  $a, b, c \dots$  przybliżone wartości tych liczb.

**Błędem bezwzględnym** liczby przybliżonej  $a$  nazywamy wartość bezwzględną różnicy między liczbą dokładną, a jej przybliżeniem, zatem

$$\Delta a = |A - a| \quad (1.2.1)$$

Na ogół nie jest znana wartość liczby  $A$ , wtedy  $\Delta a$  możemy tylko oszacować z góry. W praktyce błędem bezwzględnym nazywamy możliwie najmniejsze oszacowanie takiej różnicy. Na ogół błąd bezwzględny możemy przyjmować jako dokładność przyrządu pomiarowego.

**Błędem względnym** nazywamy stosunek błędu bezwzględnego do wartości bezwzględnej liczby przybliżonej:

$$\delta a = \frac{\Delta a}{|a|} \quad (1.2.2)$$

dla  $a$  różnego od zera.

**Przykład 1.2.1.:** Jeśli liczba dokładna  $A=1,88$ , a chcemy podać jej przybliżenie z dokładnością do jednego miejsca po przecinku to  $a=1,9$  i błąd  $\Delta a = |A - a| = 0,02$ . Oczywiście **nie** "obcinamy" liczby dokładnej do dwóch miejsc po przecinku, tylko zaokrąglamy cyfrę 8 do cyfry 9. Gdybyśmy "obcięli" i za liczbę przybliżoną przyjęli  $a^*=1,8$  błąd bezwzględny byłby o wiele większy  $\Delta a^* = |A - a^*| = 0,08$ . Obliczymy jeszcze błędy względne obu przybliżeń:

$$\delta a = \frac{\Delta a}{|a|} = \frac{0,02}{1,9} = 0,0105 \text{ co stanowi } 1,05\% \text{ liczby przybliżonej}$$

natomiast

$$\delta a^* = \frac{\Delta a^*}{|a^*|} = \frac{0,08}{1,8} = 0,0444 \text{ co stanowi } 4,44\% \text{ liczby przybliżonej.}$$

Będziemy się posługiwać pojęciem cyfr dokładnych liczby przybliżonej. W podanym przykładzie w liczbie  $a=1,9$  wszystkie cyfry będą dokładne mimo, że w liczbie dokładnej nie występuje cyfra 9, natomiast w przybliżeniu  $a^*=1,8$  cyfra 8 nie jest dokładna mimo, że taka sama cyfra i na takim samym miejscu występuje w  $A$ . Liczba przybliżona będzie mieć wszystkie cyfry dokładne, jeśli jej błąd bezwzględny nie będzie przekraczać połowy ostatniego miejsca dziesiętnego.

**Przykład 1.2.2.:** Dany jest szereg:  $\sum_{n=1}^{\infty} (-1)^n \frac{2n}{(n+2)!}$ , obliczymy jego przybliżoną sumę,

przyjmując dokładność  $\varepsilon = 10^{-10}$ .

Sumujemy szereg naprzemienny zbieżny. Jeśli za przybliżoną sumę szeregu będziemy brać  $n$ -tą sumę częściową, to błąd bezwzględny między dokładną sumą a jej przybliżeniem nie będzie przekraczał wartości bezwzględnej pierwszego odrzuconego wyrazu czyli  $|a_{n+1}|$ . Zatem

będziemy brać tyle wyrazów, aż sąsiednie sumy będą się różnić o mniej niż podana dokładność 0,0000000001.

$$s_1 = \frac{-2}{3!} \quad s_2 = s_1 + a_2 = \frac{-2}{3!} + \frac{4}{4!} \quad s_n = s_{n-1} + (-1)^n \frac{2n}{(n+2)!}$$

i sumujemy tak długo, aż  $|s_n - s_{n-1}| < \varepsilon = 0,0000000001$ . Okazuje się, że wystarczy przesumować  $n=13$  wyrazów i wtedy przybliżona suma będzie wynosić  $s_n = -0,2072766470$ .

### 1.3 BŁĘDY FUNKCJI JEDNEJ ZMIENNEJ

**Przykład 1.3.1.:** Zmierzyliśmy długość boku sześcianu i otrzymaliśmy wynik  $x=2,3\text{cm}$ , ale naszą miarką możemy zmierzyć z dokładnością do  $0,03\text{cm}$ . Jak błąd długości boku wpłynie na błąd objętości tego sześcianu?

Mamy następujące dane: bok  $x=2,3\text{cm}$ , błąd  $\Delta x = 0,03\text{cm}$ , objętość sześcianu jest funkcją boku i wynosi  $v(x)=x^3$ . Szukamy  $\Delta v$  i  $\delta v$ , czyli błędu bezwzględnego i względnego objętości. Aby wykonać obliczenia podamy ogólne wzory na te błędy.

Rozpatrujemy funkcję jednej zmiennej  $f(x)$  i argument  $x$  jest obarczony błędem bezwzględnym  $\Delta x$  wtedy błąd bezwzględny funkcji, oznaczany przez  $\Delta f$ , równa się:

$$\Delta y = \Delta f = |f'(x)| \Delta x \quad (1.3.1)$$

gdzie pochodną we wzorze obliczamy dla wartości podanego argumentu  $x$ . Wzór ten wynika ze wzoru Taylora funkcji jednej zmiennej - nie będziemy go tutaj wyprowadzać. Z ogólnego wzoru na błąd względny możemy zapisać:

$$\delta y = \delta f = \frac{\Delta f}{|f(x)|} \quad (1.3.2)$$

Wzór ten można przekształcić, wstawiając do niego wzór na błąd bezwzględny i otrzymamy:

$$\delta y = \delta f = \frac{\Delta f}{|f(x)|} = \frac{|f'(x)| \Delta x}{|f(x)|} = \frac{|f'(x) \cdot x|}{|f(x)|} \cdot \frac{\Delta x}{|x|} = w \cdot \delta x$$

gdzie wielkość  $w = \left| \frac{f'(x) \cdot x}{f(x)} \right|$  nazywamy **wskaźnikiem uwarunkowania** i za jego pomocą możemy zapisać wzór na błąd względny funkcji:

$$\delta f = w \cdot \delta x \quad (1.3.3)$$

Z tego wzoru widać, że wskaźnik ten "przenosi" błąd względny z argumentu na funkcję.

Wróćmy do przykładu 1.3.1. Korzystając z powyższych wzorów mamy:

$$v(x) = x^3, \quad v'(x) = 3x^2, \quad x = 2,3, \quad \Delta x = 0,03$$

$$v(x) = 12,2, \quad \Delta v = |3 \cdot (2,3)^2| \cdot 0,03 = 0,5$$

$$\delta v = \frac{0,5}{(2,3)^3} = 0,039, \quad \delta x = \frac{0,03}{|2,4|} = 0,013, \quad w = 3$$

Błąd względny funkcji powiększył się 3 razy w stosunku do błędu względnego argumentu. Oczywiście wyniki są zaokrąglone. Ponieważ błąd bezwzględny objętości wynosi 0,5 nie ma sensu podawać w wyniku więcej cyfr po przecinku, nawet cyfra 2 po przecinku nie jest cyfrą

dokładną. Wyniki na błędy względne podane są z trzema cyframi po przecinku, żeby wyraźnie było widać, że błąd względny wzrósł 3 razy.

W następnym przykładzie błąd względny funkcji dla wartości  $x=1$  i  $x=-1$  rośnie do nieskończoności, a im bliższe są argumenty tych wartości tym błąd jest większy. Wiąże się to z odejmowaniem liczb przybliżonych bliskich - porównać przykład z tematu: Błędy działań arytmetycznych.

**Przykład 1.3.2.** Dana jest funkcja  $f(x) = x^2 - 1$ . Napisać wzór na wskaźnik uwarunkowania. Obliczyć go dla różnych wartości argumentu  $x$ . Obliczyć błędy względne funkcji dla różnych argumentów, biorąc za błąd względny argumentu 5% jego wartości (bezwzględnej).

Obliczymy pochodną funkcji i podamy wzory na błędy:

$$f'(x) = 2x, \quad \Delta f = |2x| \cdot \Delta x, \quad \delta x = 0.05 |x|, \quad \Delta x = \delta x |x|,$$

$$w(x) = \left| \frac{f'(x) \cdot x}{f(x)} \right| = \left| \frac{2x \cdot x}{x^2 - 1} \right| = \left| \frac{2x^2}{x^2 - 1} \right|$$

$$\delta f(x) = w(x) \cdot \delta x$$

Wartość argumentu  $x =$

Wartość wskaźnika uwarunkowania wynosi=

Wartość funkcji dla podanego argumentu=

Błąd względny funkcji dla podanego argumentu wynosi=

Dla  $x=1$  nie można obliczyć błędu względnego funkcji, jeśli  $x$  będzie bliski jedności, wskaźnik uwarunkowania będzie duży i błąd względny funkcji też będzie duży. Proszę wstawiać argumenty dalekie od 1 i bliskie np.: 1,03.

## 1.4 BŁĘDY FUNKCJI WIELU ZMIENNYCH

**Przykład 1.4.1:** Zmierzyliśmy boki prostopadłościanu i otrzymaliśmy  $x=1,2\text{cm}$ ,  $y=1,8\text{cm}$  oraz  $z=2,1\text{cm}$ . Nasz przyrząd pomiarowy ma dokładność  $0,01\text{cm}$ . Jaki popełnimy błąd bezwzględny i względny licząc pole powierzchni całkowitej tego prostopadłościanu?

Mamy dane:  $x=1,2$ ,  $y=1,8$ ,  $z=2,1$ , błędy bezwzględne przyjmujemy dla wszystkich boków  $\Delta x = \Delta y = \Delta z = \Delta = 0,01$ , wzór na pole  $p(x,y,z) = 2xy + 2xz + 2yz$ . Skorzystamy z następujących wzorów dla funkcji wielu zmiennych.

Rozważania przeprowadzimy dla funkcji 2 zmiennych, ale wszystkie wzory można uogólnić na więcej zmiennych. Dana jest funkcja  $f(x,y)$  i argumenty są obarczone błędami  $\Delta x, \Delta y$ . Wzór na błąd bezwzględny funkcji wynika, tak jak dla funkcji jednej zmiennej, ze wzoru Taylora (nie będziemy go wyprowadzać) i jest następujący:

$$\Delta f(x,y) = \left| \frac{\partial f(x,y)}{\partial x} \right| \Delta x + \left| \frac{\partial f(x,y)}{\partial y} \right| \Delta y \quad (1.4.1)$$

Pochodne cząstkowe są liczone dla tych wartości argumentów, dla których liczymy błąd.

Z ogólnego wzoru na błąd względny otrzymujemy:

$$\delta f = \frac{\Delta f}{|f(x)|} \quad (1.4.2)$$

Przekształcimy ten wzór tak, jak wzór na błąd funkcji jednej zmiennej, aby wprowadzić wskaźniki uwarunkowania.

$$\begin{aligned} \delta f &= \frac{\Delta f}{|f(x)|} = \frac{|f_x(x,y)| \Delta x + |f_y(x,y)| \Delta y}{|f(x,y)|} = \left| \frac{f_x(x,y) \cdot x}{f(x,y)} \right| \frac{\Delta x}{|x|} + \left| \frac{f_y(x,y) \cdot y}{f(x,y)} \right| \frac{\Delta y}{|y|} = \\ &= w_1 \cdot \delta x + w_2 \cdot \delta y \end{aligned} \quad (1.4.3)$$

Wielkości, które wprowadziliśmy

$$w_1 = \left| \frac{f_x(x,y) \cdot x}{f(x,y)} \right| \quad w_2 = \left| \frac{f_y(x,y) \cdot y}{f(x,y)} \right| \quad (1.4.4)$$

nazywamy **wskaźnikami uwarunkowania** odpowiednio zmiennej  $x$  i  $y$ , "przenoszą" one błędy względne argumentów na błąd względny funkcji.

Powróćmy do przykładu 1.4.1, podając jednocześnie wzory dla funkcji 3 zmiennych.

Mamy dane:  $x=1,2$ ,  $y=1,8$ ,  $z=2,1$ , błędy bezwzględne przyjmujemy dla wszystkich boków  $\Delta x = \Delta y = \Delta z = \Delta = 0,01$ , wzór na pole  $p(x,y,z) = 2xy + 2xz + 2yz$ .

$$p_x(x,y,z) = 2y + 2z, \quad p_y(x,y,z) = 2x + 2z, \quad p_z(x,y,z) = 2x + 2y$$

Wartość pola  $p=16,9\text{cm}^2$

$$\Delta p(x, y, z) = \left| \frac{\partial p(x, y, z)}{\partial x} \right| \Delta x + \left| \frac{\partial p(x, y, z)}{\partial y} \right| \Delta y + \left| \frac{\partial p(x, y, z)}{\partial z} \right| \Delta z =$$

$$= (|2 \cdot 1,8 + 2 \cdot 2,1| + |2 \cdot 1,2 + 2 \cdot 2,1| + |2 \cdot 1,2 + 2 \cdot 1,8|) \cdot 0,01 = 0,2$$

$$\delta p = \frac{\Delta p}{|p(x, y, z)|} = 0,012, \quad w_1 = 0,553, \quad w_2 = 0,702, \quad w_3 = 0,745$$

## 2.1 BŁĘDY DZIAŁAŃ ARYTMETYCZNYCH

Skorzystamy ze wzoru (1.4.1) na błąd bezwzględny funkcji dwóch zmiennych, aby wyprowadzić wzory na błędy działań arytmetycznych.

**Błąd sumy dwóch liczb przybliżonych:** Argumenty  $x$  i  $y$  są obarczone odpowiednio błędami bezwzględnymi  $\Delta x, \Delta y$ , sumę argumentów zapisujemy jako  $s(x, y) = x + y$ . Pochodne cząstkowe tej funkcji zarówno po  $x$  jak i po  $y$  są równe 1, zatem

$$\Delta s(x, y) = 1 \cdot \Delta x + 1 \cdot \Delta y = \Delta x + \Delta y \quad (2.1.1)$$

Zatem błąd bezwzględny sumy równa się sumie błędów bezwzględnych składników.

Błąd względny trzeba obliczyć z ogólnego wzoru:

$$\delta s(x, y) = \frac{\Delta s}{|s(x, y)|} = \frac{\Delta x + \Delta y}{|x + y|} \quad (2.1.2)$$

**Błąd różnicy dwóch liczb przybliżonych:** Argumenty  $x$  i  $y$  są obarczone odpowiednio błędami bezwzględnymi  $\Delta x, \Delta y$ , różnicę argumentów zapisujemy jako  $r(x, y) = x - y$ . Pochodne cząstkowe tej funkcji: po  $x$  jest równa 1, po  $y$  jest równa -1, zatem

$$\Delta r(x, y) = 1 \cdot \Delta x + |-1| \cdot \Delta y = \Delta x + \Delta y \quad (2.1.3)$$

Zatem błąd bezwzględny różnicy równa się **sumie** błędów bezwzględnych składników.

Błąd względny trzeba obliczyć z ogólnego wzoru:

$$\delta r(x, y) = \frac{\Delta r}{|r(x, y)|} = \frac{\Delta x + \Delta y}{|x - y|} \quad (2.1.4)$$

Wzór na błąd względny ma sens jeśli  $x$  jest różne od  $y$ .

**Przykład 2.1.1:** Dane są dwie liczby przybliżone  $a=0,0035$  i  $b=0,0033$ . Wszystkie cyfry tych liczb są dokładne tzn.: błędy bezwzględne tych liczb są równe 0,00005. Obliczymy błędy względne tych liczb i błąd względny różnicy  $a - b$ .

$$a = 0,0035 \quad b = 0,0033 \quad \Delta a = 0,00005 \quad \Delta b = 0,00005 \quad \delta a = \frac{\Delta a}{|a|} = 0,015 \quad \delta b = \frac{\Delta b}{|b|} = 0,015$$

$$r = a - b = 0,0001 \quad \Delta r = \Delta a + \Delta b = 0,0001 \quad \delta r = 1$$

Błędy względne składników to 1,5% dla  $a$  i 1,5% dla  $b$ , natomiast błąd względny różnicy jest bardzo duży w porównaniu z błędami składników i wynosi 100%. Ten niekorzystny efekt jest związany z odejmowaniem liczb przybliżonych bliskich. Jeśli jest możliwość zastąpienia różnicy innym działaniem, należy zastosować inny wzór, aby nie odejmować liczb przybliżonych bliskich.

**Błąd iloczynu dwóch liczb przybliżonych:** Argumenty  $x$  i  $y$  są obarczone odpowiednio błędami bezwzględnymi  $\Delta x, \Delta y$ , iloczyn argumentów zapisujemy jako  $i(x, y) = x \cdot y$ . Pochodna



cząstkowa tej funkcji: po  $x$  jest równa  $y$ , po  $y$  jest równa  $x$ , zatem

$$\Delta i(x,y) = |y| \cdot \Delta x + |x| \cdot \Delta y \quad (2.1.5)$$

Błąd względny trzeba obliczyć z ogólnego wzoru:

$$\delta i(x,y) = \frac{\Delta i}{|i(x,y)|} = \frac{|y| \cdot \Delta x + |x| \cdot \Delta y}{|x \cdot y|} = \frac{\Delta x}{|x|} + \frac{\Delta y}{|y|} = \delta x + \delta y \quad (2.1.6)$$

Zatem błąd względny iloczynu równa się **sumie** błędów względnych czynników.

Wzór na błąd względny ma sens jeśli  $x$  i  $y$  są różne od zera.

**Błąd ilorazu dwóch liczb przybliżonych:** Argumenty  $x$  i  $y$  są obarczone odpowiednio błędami bezwzględnymi  $\Delta x, \Delta y$ , iloraz argumentów zapisujemy jako  $ir(x,y) = x/y$  dla  $y$  różnego od zera.

Pochodna cząstkowa tej funkcji: po  $x$  jest równa  $\frac{1}{y}$ , po  $y$  jest równa  $-\frac{x}{y^2}$ , zatem

$$\Delta ir(x,y) = \left| \frac{1}{y} \right| \cdot \Delta x + \left| \frac{-x}{y^2} \right| \cdot \Delta y \quad (2.1.7)$$

Błąd względny trzeba obliczyć z ogólnego wzoru:

$$\delta ir(x,y) = \frac{\Delta ir}{|ir(x,y)|} = \frac{\left| \frac{1}{y} \right| \cdot \Delta x + \left| \frac{x}{y^2} \right| \cdot \Delta y}{\left| \frac{x}{y} \right|} = \frac{\Delta x}{|x|} + \frac{\Delta y}{|y|} = \delta x + \delta y \quad (2.1.8)$$

Zatem błąd względny ilorazu równa się **sumie** błędów względnych czynników.

Wzór na błąd względny ma sens jeśli  $x$  i  $y$  są różne od zera.

Na koniec tego tematu podamy przykład, który pokazuje, że sumowanie liczb za pomocą pewnych programów może nie być przemienne.

**Przykład 2.1.2.:** W programie, w którym różnica rzędu między liczbami nie może przekraczać  $10^{15}$ , obliczymy na dwa sposoby sumę liczb bardzo małych i dużej. Zmiana kolejności sumowania wpłynęła na wynik końcowy.

I sposób:

Do największej liczby  $c=26$  dodajemy po kolei liczby  $d_i = 9 \cdot 10^{-16}$  gdzie  $i = 0, 1, \dots, n$ , a liczba  $n=300000$ . Korzystamy ze wzoru rekurencyjnego  $s_i = c + d_i$  i w wyniku otrzymujemy sumę

$S=26,00000000000000$  ( 13 zer po przecinku).

II sposób:

Sumujemy po kolei liczby małe według wzoru rekurencyjnego:

$$s_0 = d_0 \quad j = 1, \dots, n \quad s_j = s_{j-1} + d_j$$

A potem dodajemy wynik do liczby dużej:  $S = s_n + c$  i w wyniku otrzymujemy  $S = 26,0000000002700$ .

Nie wszystkie wzory, matematycznie poprawne, nadają się do obliczeń numerycznych. Niektóre z nich są bardzo "czułe" na błędy danych wejściowych i bardzo małe błędy tych danych powodują olbrzymie błędy wyniku. Mówimy o tych wzorach, że są numerycznie niestabilne. Podamy tutaj, często cytowany w literaturze, przykład takiego niestabilnego algorytmu.

$$\int_0^1 x^n e^{x-1} dx = 1 - n \int_0^1 x^{n-1} e^{x-1} dx$$
$$\int_0^1 x^n e^{x-1} dx = \left| \begin{array}{ll} u = x^n & u' = nx^{n-1} \\ v' = e^{x-1} & v = e^{x-1} \end{array} \right|_0^1 = x^n e^{x-1} \Big|_0^1 - \int_0^1 nx^{n-1} e^{x-1} dx = 1 - n \int_0^1 x^{n-1} e^{x-1} dx$$
$$I_n = \int_0^1 x^n e^{x-1} dx$$
$$I_1 = \int_0^1 x e^{x-1} dx = x e^{x-1} \Big|_0^1 - \int_0^1 e^{x-1} dx = 1 - e^{x-1} \Big|_0^1 = 1 - 1 + \frac{1}{e} = \frac{1}{e}$$
$$I_n = 1 - n \cdot I_{n-1}$$

Obliczyliśmy tę całkę z powyższego wzoru rekurencyjnego dla  $n$  od 2 do 15 i dla  $n$  od 16 do 30. Jednocześnie obliczyliśmy numerycznie tę samą całkę za pomocą wzorów na całkowanie metodą Simpsona (rozdział 5) , oznaczamy te drugie wyniki przez  $c_n$ . Obok w kolumnach są różnice między wartościami  $I_n$  i  $c_n$ . Na początku te wyniki się pokrywają, ale dla pewnych  $n$  zaczynają się silnie "rozjeżdżać". To właśnie algorytm na  $I_n$  daje olbrzymie błędy. Poniżej wyniki obliczeń dla  $j=2, 3, \dots, 15$  i dla  $k=16, 17, \dots, 30$

|  |  |
|--|--|
|  |  |
|--|--|

| $c_j =$ | $I_j =$ | $ c_j - I_j  =$        | $c_k =$ | $I_k =$                | $ c_k - I_k  =$       |
|---------|---------|------------------------|---------|------------------------|-----------------------|
| 0.264   | 0.264   | 0                      | 0.056   | 0.055                  | $2.6 \cdot 10^{-4}$   |
| 0.207   | 0.207   | 0                      | 0.053   | 0.057                  | $4.421 \cdot 10^{-3}$ |
| 0.171   | 0.171   | 0                      | 0.05    | -0.029                 | 0.08                  |
| 0.146   | 0.146   | $1.471 \cdot 10^{-15}$ | 0.048   | 1.56                   | 1.512                 |
| 0.127   | 0.127   | $8.965 \cdot 10^{-15}$ | 0.046   | -30.192                | 30.238                |
| 0.112   | 0.112   | $6.263 \cdot 10^{-14}$ | 0.044   | 635.04                 | 634.997               |
| 0.101   | 0.101   | $5.011 \cdot 10^{-13}$ | 0.042   | $-1.397 \cdot 10^4$    | $1.397 \cdot 10^4$    |
| 0.092   | 0.092   | $4.51 \cdot 10^{-12}$  | 0.04    | $3.213 \cdot 10^5$     | $3.213 \cdot 10^5$    |
| 0.084   | 0.084   | $4.51 \cdot 10^{-11}$  | 0.039   | $-7.711 \cdot 10^6$    | $7.711 \cdot 10^6$    |
| 0.077   | 0.077   | $4.961 \cdot 10^{-10}$ | 0.037   | $1.928 \cdot 10^8$     | $1.928 \cdot 10^8$    |
| 0.072   | 0.072   | $5.953 \cdot 10^{-9}$  | 0.036   | $-5.012 \cdot 10^9$    | $5.012 \cdot 10^9$    |
| 0.067   | 0.067   | $7.739 \cdot 10^{-8}$  | 0.035   | $1.353 \cdot 10^{11}$  | $1.353 \cdot 10^{11}$ |
| 0.063   | 0.063   | $1.084 \cdot 10^{-6}$  | 0.033   | $-3.789 \cdot 10^{12}$ | $3.789 \cdot 10^{12}$ |
| 0.059   | 0.059   | $1.625 \cdot 10^{-5}$  | 0.032   | $1.099 \cdot 10^{14}$  | $1.099 \cdot 10^{14}$ |
|         |         |                        | 0.031   | $-3.297 \cdot 10^{15}$ | $3.297 \cdot 10^{15}$ |

W następnym przykładzie podajemy wzór, który również jest poprawny, ale nie można z niego obliczać wartości funkcji dla argumentu bliskiego zeru. Wystarczą błędy dokładności obliczeń numerycznych, a wyniki odbiegają znacznie od wartości dokładnych.

**Przykład 2.2.2.** Obliczymy pewne wartości funkcji  $f(x) = \frac{1 - \cos x}{x^2}$  i funkcji

$$g(x) = \frac{1}{2} \cdot \left( \frac{\sin^2 \frac{x}{2}}{\left(\frac{x}{2}\right)^2} \right). \text{ To jest ta sama funkcja, tylko skorzystaliśmy ze wzorów}$$

trygonometrycznych i inaczej ją przekształciliśmy. Przypominamy, że  $\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1$  i z tej

granicy wynika, że wartości funkcji w punktach bliskich zeru będą bliskie 0,5. Obliczyliśmy wartości jednej i drugiej funkcji w podanych, bliskich zeru punktach i otrzymane wyniki wskazują, że pierwszy wzór jest niestabilny w punktach bliskich zeru. Obliczenia prowadziliśmy z dokładnością do 15 miejsc po przecinku.

Wyniki:

| $x_n =$            | $f(x_n) =$ | $g(x_n) =$ |
|--------------------|------------|------------|
| $2 \cdot 10^{-3}$  | 0.500000   | 0.500000   |
| $2 \cdot 10^{-4}$  | 0.500000   | 0.500000   |
| $2 \cdot 10^{-5}$  | 0.500000   | 0.500000   |
| $2 \cdot 10^{-6}$  | 0.499989   | 0.500000   |
| $2 \cdot 10^{-7}$  | 0.499600   | 0.500000   |
| $2 \cdot 10^{-8}$  | 0.555112   | 0.500000   |
| $2 \cdot 10^{-9}$  | 0.000000   | 0.500000   |
| $2 \cdot 10^{-10}$ | 0.000000   | 0.500000   |



### 3.1 FUNKCJE INTERPOLACYJNE

Z doświadczeń lub pomiarów określiliśmy w  $n+1$  różnych punktach :  $x_0, x_1, x_2, \dots, x_n$  z przedziału  $< a, b >$  wartości funkcji  $y=f(x)$  i te wartości oznaczyliśmy przez:

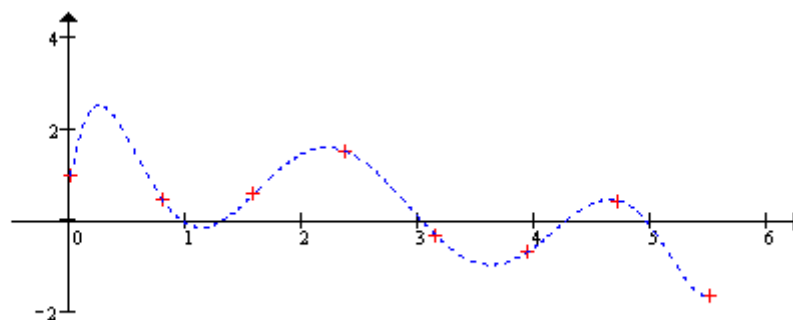
$$y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$$

Interpolacja służy do znajdowania przybliżonych wartości funkcji  $f(x)$  w dowolnym punkcie przedziału  $< a, b >$  nawet w przypadku, gdy znane jest tylko kilka wartości funkcji  $f(x)$  w tym przedziale.

Zadaniem interpolacji jest wyznaczenie funkcji  $F(x)$ , zwanej **funkcją interpolacyjną**, określonej w przedziale  $< a, b >$ , która w punktach  $x_0, x_1, x_2, \dots, x_n$ , zwanych **węzłami interpolacji**, przyjmuje wartości funkcji  $f(x)$  i w punktach poza węzłami przybliża wartość tej funkcji.

Zatem dla punktów  $x_i \quad i = 0, 1, 2, \dots, n$  w przedziale  $< a, b >$  funkcja  $F(x)$  musi spełniać  $n+1$  warunków:  $F(x_i) = y_i = f(x_i) \quad i = 0, 1, 2, \dots, n$

Wykres funkcji interpolacyjnej musi przechodzić przez punkty  $(x_i, y_i) \quad i = 0, 1, 2, \dots, n$  zaznaczone czerwonymi krzyżykami. Funkcja interpolacyjna jest narysowana niebieską przerywaną linią.



Rys3.1.1.Funkcja interpolacyjna.

Jako funkcje interpolacyjne stosuje się bardzo często:

- wielomiany algebraiczne stopnia  $n$ , oznaczmy je przez  $W_n(x)$ ,
- funkcje sklejane  $S(x)$ .

Wielomiany algebraiczne stopnia  $n$  będziemy zapisywać w znanej postaci:

$$W_n(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1} + a_n x^n$$

gdzie  $a_i \quad i = 0, 1, 2, \dots, n$  to współczynniki rzeczywiste wielomianu.

W dalszych rozważaniach właśnie tym wielomianom poświęcimy najwięcej miejsca i będziemy je wykorzystywać jako funkcje interpolacyjne.

Funkcje sklejane tzw. "splajny" będziemy omawiać w lekcji 5 .

### 3.2 WIELOMIANY INTERPOLACYJNE

Dla funkcji  $y = f(x)$ , która w  $n+1$  różnych punktach:  $x_0, x_1, x_2, \dots, x_n$  z przedziału  $\langle a, b \rangle$  przyjmuje wartości  $y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$  zbudujemy wielomian interpolacyjny algebraiczny w postaci:

$$W_n(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1} + a_n x^n \quad (3.2.1)$$

Ponieważ wielomian  $n$ -tego stopnia ma  $n+1$  niewiadomych współczynników  $a_i$   $i = 0, 1, 2, \dots, n$ , aby go jednoznacznie określić trzeba wyznaczyć  $n+1$  równań, w których te współczynniki są niewiadomymi. Inaczej mówiąc trzeba podać  $n+1$  punktów, przez które ma przechodzić wykres tego wielomianu. I tak wielomian pierwszego stopnia  $W_1(x) = a_0 + a_1 x$  graficznie przedstawia prostą, ma dwa współczynniki  $a_0, a_1$  i wiadomo, że przez dwa punkty przechodzi jedna prosta. Wielomian drugiego stopnia  $W_2(x) = a_0 + a_1 x + a_2 x^2$  ma trzy współczynniki  $a_0, a_1, a_2$  i wymaga trzech punktów, aby określić te współczynniki jednoznacznie, graficznie taki wielomian przedstawia parabolę (przez trzy punkty przechodzi jedna parabola). Ogólnie zatem prawdziwe jest twierdzenie:

**Twierdzenie o istnieniu i jednoznaczności wielomianu interpolacyjnego:** Istnieje jedyny wielomian interpolacyjny  $W_n(x)$  stopnia co najwyżej  $n$ , który w  $n+1$  różnych punktach

$x_0, x_1, x_2, \dots, x_n$  z przedziału  $\langle a, b \rangle$  pokrywa się z funkcją  $y = f(x)$ , tzn.:

$$W(x_i) = y_i = f(x_i) \quad i = 0, 1, 2, \dots, n.$$

Dowód: Dla wielomianu  $W_n(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1} + a_n x^n$  warunki

$W(x_i) = y_i = f(x_i) \quad i = 0, 1, 2, \dots, n$  sprowadzają się do rozwiązania układu  $n+1$  następujących równań liniowych z  $n+1$  niewiadomymi  $a_i \quad i = 0, 1, 2, \dots, n$ :

$$\begin{cases} a_0 + a_1 x_0 + a_2 x_0^2 + \dots + a_{n-1} x_0^{n-1} + a_n x_0^n = y_0 \\ a_0 + a_1 x_1 + a_2 x_1^2 + \dots + a_{n-1} x_1^{n-1} + a_n x_1^n = y_1 \\ \dots \\ a_0 + a_1 x_n + a_2 x_n^2 + \dots + a_{n-1} x_n^{n-1} + a_n x_n^n = y_n \end{cases} \quad (3.2.2)$$

Przypominamy, że układ równań liniowych  $(n+1) \times (n+1)$  ma jednoznaczne rozwiązania na niewiadome  $a_i \quad i = 0, 1, 2, \dots, n$ , jeśli wyznacznik tego układu jest różny od zera. Jak wygląda wyznacznik naszego układu?

$$V_n = \begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^{n-1} & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^{n-1} & x_1^n \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} & x_n^n \end{vmatrix} \quad (3.2.3)$$

Jest to wyznacznik Vandermonde'a i jest on równy iloczynowi wszystkich możliwych różnic

$x_j - x_i$  gdzie  $j > i$ . Zapisujemy ten fakt w postaci

$$V_n = \prod_{i=0, \dots, n-1} \prod_{j=1, \dots, n \atop j > i} (x_j - x_i)$$

Ponieważ węzły  $x_i$   $i = 0, 1, 2, \dots, n$  są różne to wyrazy  $x_j - x_i$  są różne od zera, zatem wyznacznik jest różny od zera i układ ma zawsze jedyne rozwiązanie. Oznacza to, że istnieje zawsze jedyny szukany wielomian interpolacyjny. Może być stopnia niższego niż  $n$ , bo współczynnik  $a_n$  może być równy zero, oczywiście jeszcze jakieś inne współczynniki mogą się zerować. (Więcej o sposobie **obliczania** wyznacznika Vandermonde'a).



Aby wyjaśnić obliczanie wyznacznika Vandermonde`a

$$V_n = \begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^{n-1} & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^{n-1} & x_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} & x_n^n \end{vmatrix} \quad V_n = \prod_{i=0, \dots, n-1} \prod_{j=1, \dots, n, j>i} (x_j - x_i)$$

obliczymy ten wyznacznik dla stopnia 2 i 3, dalsze obliczanie wynika z indukcji matematycznej - nie będziemy jej prowadzić, aby nie komplikować rozważań.

Dla  $n=2$  mamy

$$V_2 = \begin{vmatrix} 1 & x_0 \\ 1 & x_1 \end{vmatrix} = x_1 - x_0 \neq 0 \text{ bo węzły } x_0 \neq x_1.$$

Dla  $n=3$  wyznacznik obliczamy następująco:

$$\begin{aligned} V_3 &= \begin{vmatrix} 1 & x_0 & x_0^2 \\ 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \end{vmatrix} = \begin{vmatrix} 1 & x_0 & x_0^2 \\ 0 & x_1 - x_0 & x_1^2 - x_0^2 \\ 0 & x_2 - x_0 & x_2^2 - x_0^2 \end{vmatrix} = \begin{vmatrix} x_1 - x_0 & x_1^2 - x_0^2 \\ x_2 - x_0 & x_2^2 - x_0^2 \end{vmatrix} = \\ &= (x_1 - x_0)(x_2 - x_0) \begin{vmatrix} 1 & x_1 + x_0 \\ 1 & x_2 + x_0 \end{vmatrix} = (x_1 - x_0)(x_2 - x_0)(x_2 - x_1) \end{aligned}$$

Jak powstały te przekształcenia? W pierwszym kroku od drugiego wiersza i od trzeciego wiersza został odjęty wiersz pierwszy i wyzerowały się wyrazy w pierwszej kolumnie w drugim i trzecim wierszu. Dalej następuje rozwinięcie wyznacznika względem pierwszej kolumny i zostaje jeden wyznacznik drugiego stopnia. Wyciągamy z pierwszego wiersza wyrażenie  $x_1 - x_0$ , a z drugiego wiersza  $x_2 - x_0$  przed wyznacznik korzystając ze wzoru skróconego mnożenia  $a^2 - b^2 = (a - b)(a + b)$ . I znów ponieważ węzły są różne otrzymujemy wyrażenie na wyznacznik, które jest różne od zera.

Potem prowadzimy indukcję względem stopnia wyznacznika.

[powrót](#)

### 3.3 WIELOMIAN LAGRANGE`A

Poszukujemy wielomianu algebraicznego w postaci:

$$WL_n(x) = \frac{(x-x_1)(x-x_2)\cdots(x-x_n)}{(x_0-x_1)(x_0-x_2)\cdots(x_0-x_n)}y_0 + \frac{(x-x_0)(x-x_2)\cdots(x-x_n)}{(x_1-x_0)(x_1-x_2)\cdots(x_1-x_n)}y_1 + \dots +$$

$$+ \frac{(x-x_0)\cdots(x-x_{k-1})(x-x_{k+1})\cdots(x-x_n)}{(x_k-x_0)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)}y_k + \dots + \frac{(x-x_0)(x-x_1)\cdots(x-x_{n-1})}{(x_n-x_0)(x_n-x_1)\cdots(x_n-x_{n-1})}y_n,$$

Wielomian ten nosi nazwę **wielomianu interpolacyjnego Lagrange`a**.

Wielomian ten ma  $n+1$  składników, w każdym ze składników przy  $y_k$   $k=0,1,2,\dots,n$  w liczniku nie występuje czynnik  $(x-x_k)$  to znaczy w liczniku jest wielomian  $n$ -tego stopnia, a w mianowniku od węzła  $x_k$  odejmowane są wszystkie inne węzły i te różnice są pomnożone przez siebie.

Ponieważ suma wielomianów stopnia  $n$  jest wielomianem co najwyżej  $n$ -tego stopnia (jakieś wyrazy mogą się zredukować i możemy dostać wielomian niższego stopnia, ale nigdy wyższego) wielomian Lagrange`a  $WL_n(x)$  jest wielomianem co najwyżej  $n$ -tego stopnia.

Sprawdźmy, czy jest spełniony warunek  $WL_n(x_k) = y_k = f(x_k)$   $k=0,1,2,\dots,n$ .

Wstawiając za  $x$  do wielomianu  $WL_n(x)$  węzeł  $x_k$  otrzymujemy:

$$WL_n(x_k) = \frac{(x_k-x_1)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)}{(x_0-x_1)(x_0-x_2)\cdots(x_0-x_n)}y_0 + \frac{(x_k-x_0)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)}{(x_1-x_0)(x_1-x_2)\cdots(x_1-x_n)}y_1 +$$

$$+ \frac{(x_k-x_0)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)}{(x_k-x_0)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)}y_k + \dots + \frac{(x_k-x_0)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)}{(x_n-x_0)(x_n-x_1)\cdots(x_n-x_{n-1})}y_n$$

W każdym składniku oprócz tego, który stoi przy  $y_k$  jest w liczniku różnica  $(x_k-x_k)$  czyli zero, natomiast w tym składniku przy  $y_k$  w liczniku jest takie samo wyrażenie jak w mianowniku, to znaczy, że przy  $y_k$  jest współczynnik 1. Zatem

$$WL_n(x_k) = 0 \cdot y_0 + 0 \cdot y_1 + \dots + 1 \cdot y_k + \dots + 0 \cdot y_n = y_k$$

Wielomian Lagrange`a spełnia wymagania z twierdzenia o istnieniu, zatem jest szukany wielomianem interpolacyjnym. Jeśli wprowadzimy następujące oznaczenie:

$$\Phi_k(x) = \frac{(x-x_0)\cdots(x-x_{k-1})(x-x_{k+1})\cdots(x-x_n)}{(x_k-x_0)\cdots(x_k-x_{k-1})(x_k-x_{k+1})\cdots(x_k-x_n)} \quad (3.3.3)$$

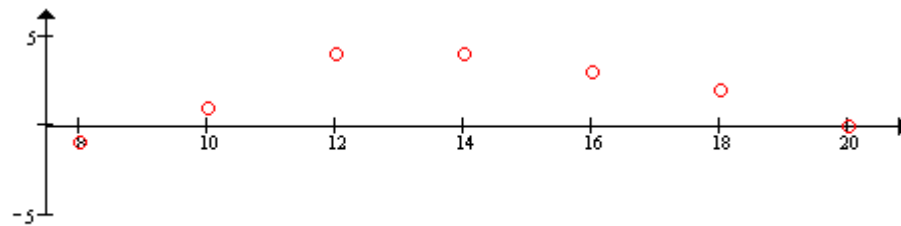
to wielomian będzie mieć postać:

$$WL_n(x) = \sum_{k=0}^n \Phi_k(x)y_k \quad (3.3.4)$$

**Przykład 3.3.1.** Zmierzyliśmy w pierwszym dniu wiosny - 21 marca 2005 roku- w Warszawie na Mokotowie temperaturę za oknem i wyniki zapisaliśmy w tabeli:

|                         |    |    |    |    |    |    |    |
|-------------------------|----|----|----|----|----|----|----|
| Godzina pomiaru         | 8  | 10 | 12 | 14 | 16 | 18 | 20 |
| Temperatura w stopniach | -1 | 1  | 4  | 4  | 3  | 2  | 0  |

Na wykresie wygląda to następująco:

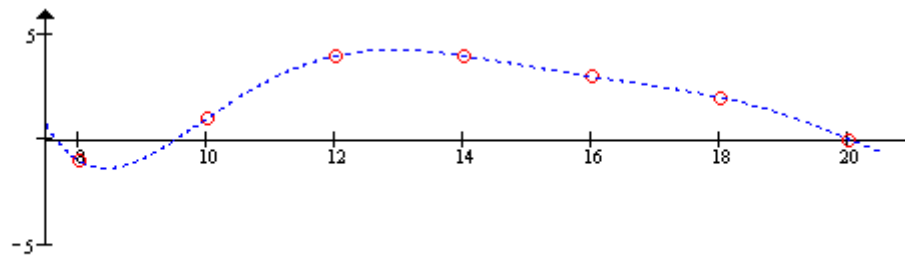


Rys3.3.1.Rozkład temperatury.

Mamy 7 pomiarów, to znaczy wartość funkcji - temperatury - jest dana w 7 równoodległych węzłach. Szukany wielomian interpolacyjny będzie zatem  $n=6$  stopnia. Po zastosowaniu wzoru na wielomian Lagrange'a dostajemy wynik:

$$W_6(x) = 1187 - \frac{7949}{15}x + \frac{22847}{240}x^2 - \frac{283}{32}x^3 + \frac{173}{384}x^4 - \frac{23}{1920}x^5 + \frac{1}{7680}x^6$$

Poniżej podajemy rysunek punktów pomiarowych i wielomianu interpolacyjnego.



Rys3.3.2.Wielomian interpolacyjny.

Możemy teraz obliczać przybliżoną temperaturę w dowolnej porze między godziną 8 a godziną 20-tą. Otrzymujemy np. że o godzinie 13<sup>30</sup> temperatura wynosiła 4,187 stopnia, o godzinie 15<sup>45</sup> wynosiła 3,119 stopnia, a o godzinie 10<sup>15</sup> miała wartość 1,52 stopnia.

**Przykład 3.3.2** Będziemy obliczać wielomian stopnia 1 dla funkcji danej w dwóch węzłach, a następnie za pomocą tego wielomianu obliczać przybliżone wartości tej funkcji w dowolnym punkcie między węzłami. (proszę wstawiać dowolne wartości dla pomiarów i dla x).

Wykonaliśmy dwa pomiary i otrzymaliśmy następujące wartości:  $x_1 =$  ,  $x_2 =$  ,

$y_1 =$  ,  $y_2 =$  .

Będziemy obliczać przybliżoną wartość funkcji w punkcie:  $x =$

Obliczymy wielomian interpolacyjny  $n=1$  stopnia, bo pomiarów jest  $n+1=2$ .

Oblicz wielomian

$W_1(x) =$    $+$    $x$

Oblicz wartość wielomianu

Wartość wielomianu dla podanego  $x$  wynosi  $W(x)$

$=$

wyczyść obliczenia

Wykresem tego wielomianu jest prosta przechodząca przez punkty  $(x_1, y_1)$ ,  $(x_2, y_2)$ .

### 3.4 WIELOMIAN NEWTONA

Inną postacią wielomianu interpolacyjnego jest **wielomian interpolacyjny Newtona**. Do zdefiniowania tego wielomianu wykorzystamy **ilorazy różnicowe**  $n$ -tego rzędu dla funkcji  $y = f(x)$ , która w  $n+1$  różnych punktach:  $x_0, x_1, x_2, \dots, x_n$  z przedziału  $\langle a, b \rangle$  przyjmuje wartości  $y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$ .

Ilorazy różnicowe I rzędu:

$$\begin{aligned} f(x_0, x_1) &= \frac{y_1 - y_0}{x_1 - x_0}, \quad f(x_1, x_2) = \frac{y_2 - y_1}{x_2 - x_1}, \quad \dots \quad f(x_{k-1}, x_k) = \frac{y_k - y_{k-1}}{x_k - x_{k-1}}, \dots \\ \dots f(x_{n-1}, x_n) &= \frac{y_n - y_{n-1}}{x_n - x_{n-1}} \end{aligned} \quad (3.4.1)$$

Ilorazy różnicowe II rzędu:

$$\begin{aligned} f(x_0, x_1, x_2) &= \frac{f(x_1, x_2) - f(x_0, x_1)}{x_2 - x_0}, \quad f(x_1, x_2, x_3) = \frac{f(x_2, x_3) - f(x_1, x_2)}{x_3 - x_1}, \dots \\ f(x_{k-2}, x_{k-1}, x_k) &= \frac{f(x_{k-1}, x_k) - f(x_{k-2}, x_{k-1})}{x_k - x_{k-2}}, \dots \\ f(x_{n-2}, x_{n-1}, x_n) &= \frac{f(x_{n-1}, x_n) - f(x_{n-2}, x_{n-1})}{x_n - x_{n-2}} \end{aligned} \quad (3.4.2)$$

I ogólnie iloraz  $m$ -tego rzędu ( $m = 2 \dots n, k = 2, \dots, n$ )

$$f(x_{k-m}, \dots, x_{k-1}, x_k) = \frac{f(x_{k-m+1}, \dots, x_k) - f(x_{k-m}, \dots, x_{k-1})}{x_k - x_{k-m}} \quad (3.4.3)$$

Iloraz  $n$ -tego rzędu jest tylko jeden i ma postać:

$$f(x_0, x_1, \dots, x_n) = \frac{f(x_1, \dots, x_n) - f(x_0, \dots, x_{n-1})}{x_n - x_0} \quad (3.4.4)$$

Zauważmy, że aby obliczyć ilorazy różnicowe rzędu  $m$ , trzeba podzielić różnice ilorazów rzędu  $m-1$  przez różnice wartości węzłów o współrzędnych różniących się o  $m$ . Jest to na pierwszy rzut oka dość skomplikowane. Zapiszemy w tabeli ilorazy różnicowe dla funkcji, dla której są dane wartości tylko w trzech punktach:

| $x_0$ | $y_0$ | Ilorazy I rzędu                             | Iloraz II rzędu  |
|-------|-------|---|--|
| $x_1$ | $y_1$ | $f(x_0, x_1) = \frac{y_1 - y_0}{x_1 - x_0}$ |  |
| $x_2$ | $y_2$ | $f(x_1, x_2) = \frac{y_2 - y_1}{x_2 - x_1}$ | $f(x_0, x_1, x_2) = \frac{f(x_1, x_2) - f(x_0, x_1)}{x_2 - x_0}$ |

**Przykład 3.4.1** : Funkcja  $f(x)$  dana jest za pomocą tabelki:

| $x_i$ | $y_i$ |
|-------|-------|
| -1    | 0     |
| 1     | 4     |
| 3     | 6     |
| 5     | 2     |

Obliczymy dla niej ilorazy różnicowe do  $n=3$  rzędu włącznie. Zauważmy, że węzłów jest  $n+1=4$ , wtedy ostatni iloraz jest  $n$ -tego rzędu - czyli trzeciego.

| $x_i$ | $y_i$ | I rzędu                  | II rzędu                                    | III rzędu   |
|-------|-------|--------------------------|---|---|
| -1    | 0     |                          |   |   |
| 1     | 4     | $\frac{4-0}{1-(-1)} = 2$ |   |   |
| 3     | 6     | $\frac{6-4}{3-1} = 1$    | $\frac{1-2}{3-(-1)} = -\frac{1}{4} = -0,25$ |   |
| 5     | 2     | $\frac{2-6}{5-3} = -2$   | $\frac{-2-1}{5-1} = -\frac{3}{4} = -0,75$   | $\frac{-0,75-(-0,25)}{5-(-1)} = \frac{-0,5}{6} = -\frac{1}{12}$ |

Możemy teraz podać postać **wielomianu interpolacyjnego Newtona** dla funkcji  $f(x)$ :

$$WN_n(x) = y_0 + f(x_0, x_1)(x - x_0) + f(x_0, x_1, x_2)(x - x_0)(x - x_1) + \dots + f(x_0, x_1, \dots, x_n)(x - x_0)(x - x_1) \cdot \dots \cdot (x - x_{n-1}) \quad (3.4.5)$$

Widać z tej postaci, że jest to wielomian co najwyżej  $n$ -tego stopnia. W ostatnim składniku jest co najwyżej  $x$  do potegi  $n$ -tej ( ale iloraz różnicowy  $n$ -tego rzędu może być równy zero).

Trudniej jest sprawdzić, że wielomian ten w węzłach pokrywa się z funkcją  $f(x)$ . Sprawdzimy ten warunek tylko dla wielomianu drugiego stopnia , który ma postać:

$$WN_2(x) = y_0 + f(x_0, x_1)(x - x_0) + f(x_0, x_1, x_2)(x - x_0)(x - x_1)$$

Dla  $x = x_0$  widać ,że

$$WN_2(x_0) = y_0 + f(x_0, x_1)(x_0 - x_0) + f(x_0, x_1, x_2)(x_0 - x_0)(x - x_1)$$

zatem  $WN_2(x_0) = y_0$

Dla  $x = x_1$  mamy

$$WN_2(x_1) = y_0 + f(x_0, x_1)(x_1 - x_0) + f(x_0, x_1, x_2)(x_1 - x_0)(x_1 - x_1)$$

Ostatni składnik znika, zostają dwa składniki, w tym drugim skorzystamy ze wzoru na iloraz różnicowy pierwszego rzędu i otrzymamy:

$$WN_2(x_1) = y_0 + f(x_0, x_1)(x_1 - x_0) = y_0 + \frac{y_1 - y_0}{x_1 - x_0}(x_1 - x_0) = y_0 + y_1 - y_0 = y_1$$

Dla  $x = x_2$  mamy

$$WN_2(x_2) = y_0 + f(x_0, x_1)(x_2 - x_0) + f(x_0, x_1, x_2)(x_2 - x_0)(x_2 - x_1)$$

i rozpiszemy iloraz drugiego rzędu w ostatnim składniku:

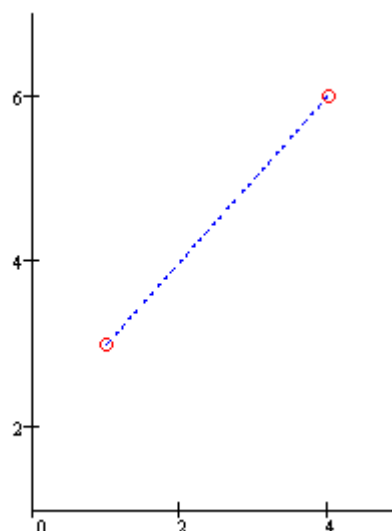
$$\begin{aligned} WN_2(x_2) &= y_0 + f(x_0, x_1)(x_2 - x_0) + \frac{f(x_1, x_2) - f(x_0, x_1)}{x_2 - x_0}(x_2 - x_0)(x_2 - x_1) = \\ &= y_0 + f(x_0, x_1)(x_2 - x_0) + f(x_1, x_2)(x_2 - x_1) - f(x_0, x_1)(x_2 - x_1) = \\ &= y_0 + f(x_0, x_1)(x_2 - x_0 - x_2 + x_1) + f(x_1, x_2)(x_2 - x_1) = \\ &= y_0 + \frac{y_1 - y_0}{x_1 - x_0}(x_1 - x_0) + \frac{y_2 - y_1}{x_2 - x_1}(x_2 - x_1) = y_0 + y_1 - y_0 + y_2 - y_1 = y_2 \end{aligned}$$

Można to sprawdzić ogólnie, że  $WN_n(x_k) = y_k = f(x_k) \quad k = 0, 1, \dots, n$ .

Powróćmy do przykładu 3.4.1, dla tej funkcji wielomian Newtona będzie następujący:

$$WN_3(x) = 0 + 2(x+1) - \frac{1}{4}(x+1)(x-1) + \frac{1}{12}(x+1)(x-1)(x-3)$$

**Przykład 3.4.2.** Zmierzyliśmy wartość funkcji w przedziale  $<1, 4>$  tylko w dwóch węzłach i otrzymaliśmy wyniki: dla  $x=1$   $y=3$  i dla  $x=4$   $y=6$ . Wielomian interpolacyjny dla tej funkcji jest prostą przechodzącą przez te punkty. Węzłów jest  $n+1=2$ , wielomian jest stopnia  $n=1$ . Potem dodaliśmy jeszcze punkt dla  $x=2$  wartość funkcji  $y=3$  i poprowadziliśmy wielomian stopnia 2. Następnie dodaliśmy jeszcze jeden pomiar i dla  $x=3$  otrzymaliśmy  $y=2$ . Teraz mamy 4 węzły w przedziale, a wielomian jest 3-stopnia. n=1 n=2 n=3



Rys3.4.1. Wielomiany interpolacyjne 1, 2 i 3 stopnia.

Wielomiany te obliczyliśmy korzystając ze wzorów na ilorazy różnicowe i wielomian Newtona:

$$w1(x) = 2 + x \quad w2(x) = 4 - \frac{3}{2}x + \frac{1}{2}x^2 \quad w3(x) = -4 + \frac{25}{2}x - \frac{13}{2}x^2 + x^3$$

**Przykład 3.4.3:** Obliczyć ilorazy różnicowe dla funkcji danej z 4 pomiarów za pomocą tabelki:  
(proszę wstawiać dowolne wartości węzłów i wartości funkcji)

| $x_i$                        | $y_i$                        |
|------------------------------|------------------------------|
| $x_0 =$ <input type="text"/> | $y_0 =$ <input type="text"/> |
| $x_1 =$ <input type="text"/> | $y_1 =$ <input type="text"/> |
| $x_2 =$ <input type="text"/> | $y_2 =$ <input type="text"/> |
| $x_3 =$ <input type="text"/> | $y_3 =$ <input type="text"/> |

Oblicz:

Ilorazy I rzędu są trzy:  $f(x_0, x_1) = \frac{y_1 - y_0}{x_1 - x_0} =$    $f(x_1, x_2) = \frac{y_2 - y_1}{x_2 - x_1} =$    
 $f(x_2, x_3) = \frac{y_3 - y_2}{x_3 - x_2} =$

Ilorazy II rzędu są dwa:  $f(x_0, x_1, x_2) = \frac{f(x_1, x_2) - f(x_0, x_1)}{x_2 - x_0} =$    
 $f(x_1, x_2, x_3) = \frac{f(x_2, x_3) - f(x_1, x_2)}{x_3 - x_1} =$

I jest jeden iloraz rzędu III:  $f(x_0, x_1, x_2, x_3) = \frac{f(x_1, x_2, x_3) - f(x_0, x_1, x_2)}{x_3 - x_0} =$

wyczyść obliczenia



## 4.1 BŁĄD INTERPOLACJI

Jeśli funkcja  $f(x)$  jest dana za pomocą tabelki, to znaczy jej wartości są wynikiem doświadczeń lub pomiarów, nie możemy określić błędu jaki popełniamy biorąc za wartość funkcji w punkcie nie będącym węzłem wartość wielomianu interpolacyjnego. Ale są przypadki gdy funkcja jest dana wzorem analitycznym  $y=f(x)$  w przedziale  $\langle a, b \rangle$ , a mimo to potrzebujemy zbudować dla niej wielomian interpolacyjny. Taka sytuacja ma miejsce przede wszystkim przy całkowaniu, o czym będziemy mówić w rozdziale dotyczącym całkowania numerycznego. W takim przypadku można obliczyć błąd interpolacji, zależy on od funkcji, a właściwie od pochodnej rzędu  $n+1$ , oraz od sposobu rozmieszczenia węzłów interpolacji. Ponieważ wzór (podajemy ten wzór bez wyprowadzania) na różnicę między funkcją interpolowaną (oczywiście taką, która ma w tym przedziale wszystkie pochodne do rzędu  $n+1$  włącznie) a wielomianem interpolacyjnym jest następujący:

$$f(x) - W_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x-x_0)(x-x_1)\dots(x-x_n) \quad (4.1.1)$$

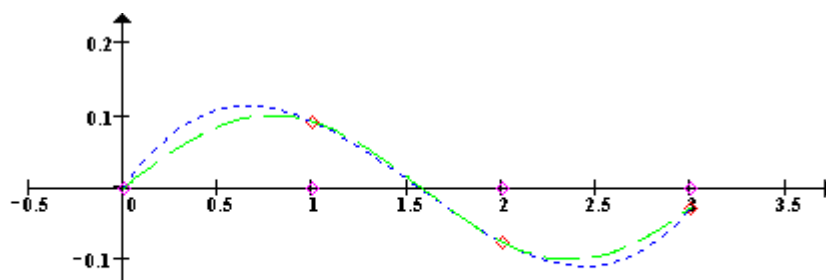
gdzie  $\xi$  jest pewnym punktem w przedziale  $\langle a, b \rangle$ , to błąd bezwzględny interpolacji można oszacować przez:

$$|f(x) - W_n(x)| \leq \frac{M_{n+1}}{(n+1)!} |(x-x_0)(x-x_1)\dots(x-x_n)| \quad (4.1.2)$$

gdzie  $M_{n+1} = \sup_{x \in \langle a, b \rangle} |f^{(n+1)}(x)|$ .

Porównajmy, dla przykładu funkcję która ma ograniczone pochodne z jej wielomianem interpolacyjnym.

**Przykład 4.1.1:** Weźmy funkcję  $f(x) = 0,1 \sin 2x$  w przedziale  $\langle 0, 3 \rangle$  i jako węzły interpolacji przyjmijmy cztery punkty równoodległe w tym przedziale 0,1,2 i 3. Funkcję i wielomian przedstawimy na wykresie:

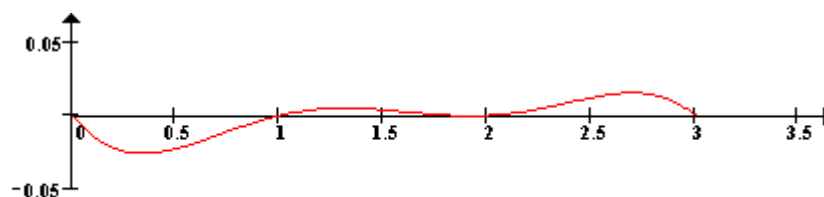


Rys4.1.1. Wykres funkcji  $f$  i jej wielomianu interpolacyjnego 3 stopnia.

Na osi  $Ox$  zaznaczone są węzły, widać, że funkcja i wielomian pokrywają się w węzłach. Funkcja jest narysowana zieloną linią, wielomian niebieską przerywaną.

Na następnym wykresie przedstawiona jest funkcja będąca różnicą  $f(x)$  i wielomianu

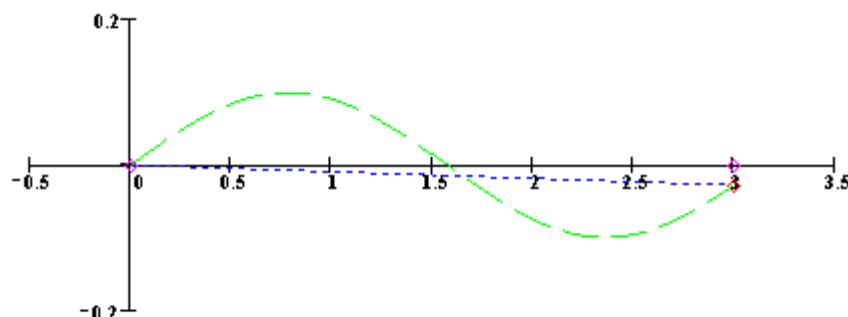
interpolacyjnego. Jak widać błąd bezwzględny nie przekracza 0,03 w rozpatrywanym przedziale.



Rys4.1.2. Wykres różnicy między funkcją  $f$  i wielomianem interpolacyjnym.

W poprzednim przykładzie obliczyliśmy wielomian 3 stopnia dla podanej funkcji, w następnym będziemy zmieniać stopień wielomianu dla tej samej funkcji.

**Przykład 4.1.2.** Będziemy rozpatrywać jeszcze raz poprzednią funkcję, a mianowicie:  $f(x) = 0,1 \sin 2x$  w przedziale  $<0,3>$ . Funkcja ta ma ograniczone pochodne w przedziale, wraz ze wzrostem ilości węzłów (bierzemy węzły równoodległe) wielomian interpolacyjny coraz lepiej będzie przybliżał daną funkcję. Proszę zmieniać stopień wielomianu  $n$ , ilość węzłów jest  $n+1$ . Na osi  $Ox$  zaznaczone są węzły, widać na rysunku, że funkcja w węzłach pokrywa się z wielomianem. Wykres funkcji narysowany jest na zielono, dla  $n=8$  wykresy na naszym rysunku prawie się pokrywają.



Rys4.1.3. Wykres funkcji  $f$  i jej wielomianów interpolacyjnych 1, 2, 3, 4 i 8 stopnia.

Dla dwóch węzłów maksymalny błąd interpolacji w tym przedziale równa się 0,107, dla trzech węzłów równa się 0,098, dla czterech 0,026, dla pięciu 0,016, a dla  $n=8$  czyli dla dziewięciu węzłów błąd nie przekracza 0,000084.

## 4.2 WĘZŁY Czebyszewa

Najczęściej stosuje się węzły równoodległe dla funkcji interpolacyjnej, które są proste w użyciu, a podczas doświadczeń można mierzyć badaną wartość funkcji co ustaloną jednostkę czasu. Jednak jak wynika ze wzoru na błąd interpolacji:

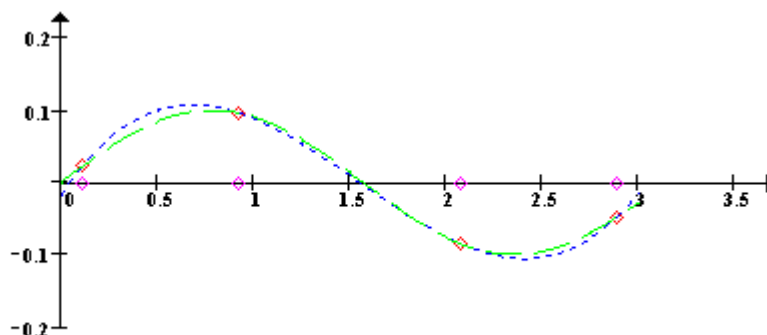
$$|f(x) - W_n(x)| \leq \frac{M_{n+1}}{(n+1)!} |(x - x_0)(x - x_1) \dots (x - x_n)| \quad (4.2.1)$$

jego wartość zależy w istotny sposób od rozmieszczenia węzłów. Okazuje się, że węzły równoodległe nie zawsze są najlepsze. Tę część błędu, zależną od węzłów minimalizują tzw.: węzły Czebyszewa, które podamy tutaj bez wyprowadzania. Jeśli szukamy w dowolnym przedziale  $< a, b >$   $n+1$  optymalnych węzłów, można je wyliczyć ze wzoru:

$$x_k = \frac{b-a}{2} \cos\left(\frac{2k+1}{2n+2}\pi\right) + \frac{b+a}{2} \quad k = 0, 1, 2, \dots, n \quad (4.2.2)$$

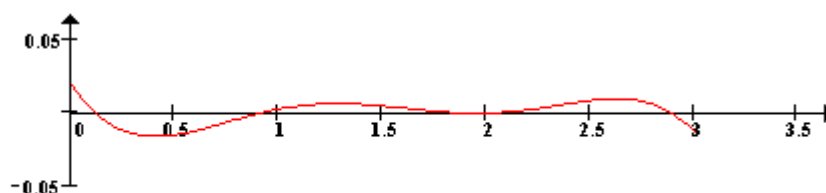
Wróćmy do przykładu 4.1.1 z funkcją  $f(x) = 0,1 \sin 2x$  w przedziale  $<0, 3>$ .

**Przykład 4.2.1:** Poprzednio dla zbudowania wielomianu interpolacyjnego braliśmy cztery węzły równoodległe : 0, 1 2 i 3. Teraz obliczymy 4 węzły Czebyszewa w tym przedziale z powyższego wzoru. Będą to : 0,114; 0,926; 2,074; 2,886 ( podajemy te wartości z dokładnością do trzech cyfr po przecinku). Te węzły są zaznaczone na osi 0x, funkcja jest narysowana zieloną linią, wielomian niebieską przerywaną:



Rys4.2.1. Wykres funkcji  $f$  i jej wielomianu interpolacyjnego 3 stopnia z 4 węzłami Czebyszewa.

Tak jak poprzednio na następnym rysunku przedstawiamy różnicę między funkcją daną, a jej wielomianem interpolacyjnym opartym na węzłach Czebyszewa. Największy błąd bezwzględny nie przekroczy 0,02 w tym przedziale ( jest on trochę mniejszy niż w poprzednim przykładzie z węzłami równoodległymi- tam było 0,03).



Rys4.2.2. Wykres różnicy między funkcją  $f$  i jej wielomianem interpolacyjnym.

### 4.3 ZBIEŻNOŚĆ PROCESÓW INTERPOLACYJNYCH

Jeszcze raz podamy wzór na błąd interpolacji dla funkcji określonej w przedziale  $\langle a, b \rangle$  i mającej w tym przedziale pochodną do rzędu  $n+1$  włącznie.

$$|f(x) - W_n(x)| \leq \frac{M_{n+1}}{(n+1)!} |(x-x_0)(x-x_1)\dots(x-x_n)| \quad (4.3.1)$$

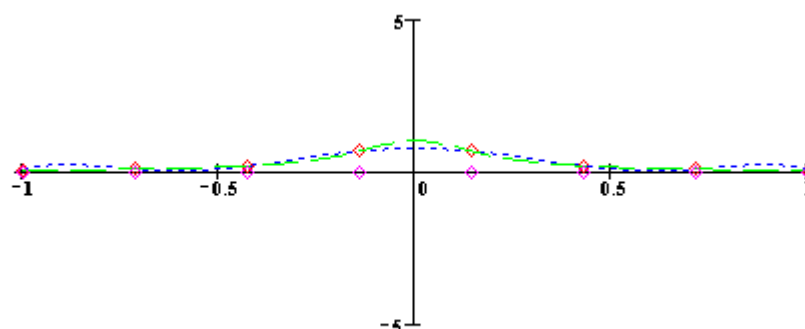
gdzie  $M_{n+1} = \sup_{x \in \langle a, b \rangle} |f^{(n+1)}(x)|$ .

Z tego wzoru wynika, że wraz ze wzrostem ilości węzłów mianownik szybko rośnie, bo jest w nim wyraz  $(n+1)!$  zatem cały ułamek winien maleć i przez to maleć powinien błąd. Ale na błąd ma wpływ wielkość ograniczająca pochodną  $n+1$  rzędu. Podamy popularny w literaturze przykład funkcji, która ma wszystkie pochodne ograniczone, ale na tyle dużej wartości, że wraz ze wzrostem ilości węzłów błąd interpolacji rośnie tzn.: "rozjeżdża" się wielomian z funkcją interpolowaną.

**Przykład 4.3.1:** Będziemy rozpatrywać funkcję:  $f(x) = \frac{1}{1+25x^2}$  w przedziale  $\langle -1, 1 \rangle$ . Na początku będziemy brać osiem węzłów równoodległych tzn.  $n=7$ . Wtedy wielomian interpolacyjny i funkcja będą zachowywać się "poprawnie", niezbyt się od siebie różnić.

Pierwszy rysunek obrazuje tę sytuację. Jeśli będziemy zwiększać ilość węzłów interpolacja będzie obarczona coraz to większym błędem.

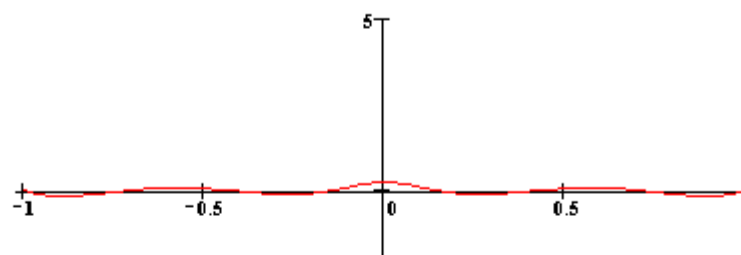
n=7 n=8 n=10 n=12



Rys4.3.1. Wykres funkcji  $f$  i jej wielomianów interpolacyjnych 7, 8, 10 i 12 stopnia.

Podobnie będzie z rysunkiem przedstawiającym różnice między funkcją interpolowaną a wielomianem interpolacyjnym o równoodległych węzłach. Te różnice dla  $n=12$  czyli dla 13 węzłów będą bardzo duże w porównaniu do wartości funkcji w tym przedziale.

n=7 n=8 n=10 n=12



Rys4.3.2. Wykres różnic między funkcją  $f$  i jej wielomianami interpolacyjnymi 7, 8, 10 i 12 stopnia.

Trzeba sobie zdać sprawę z takich faktów, że dla  $n=12$  pochodna rozpatrywanej funkcji  $n+1=13$  rzędu równa jest w punkcie 0,1 wartości  $-3,29 \cdot 10^{17}$ , a  $13! = 6,277 \cdot 10^9$ .

## 5.1 BAZA FUNKCJI SKLEJANYCH

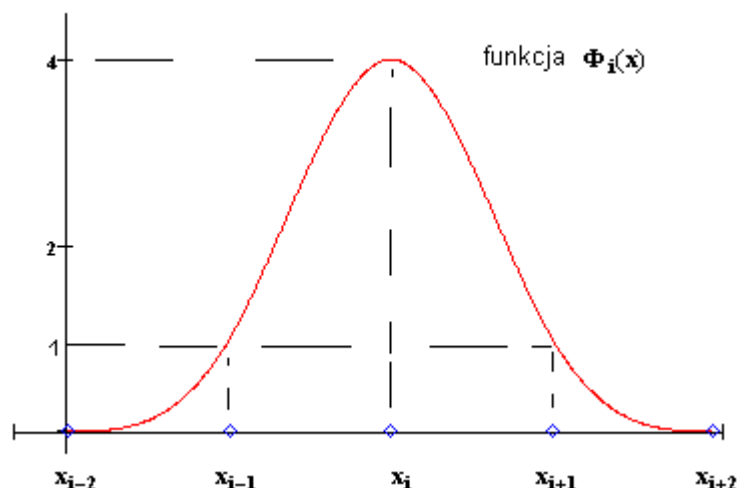
Stosując do interpolacji wielomian interpolacyjny nie możemy narzucać stopnia wielomianu, ten stopień zależy od ilości węzłów. Jeśli mamy 20 różnych węzłów ( 20 pomiarów) to wielomian interpolacyjny może być nawet 19 stopnia. Wraz ze wzrostem ilości węzłów rośnie na ogół stopień wielomianu. Natomiast stopień niżej zdefiniowanej funkcji skleianej tzw.: splajnu, nie będzie zależał od ilości węzłów.

Ograniczymy się w tym opracowaniu do splajnu 3-iego stopnia, jest on na ogół najczęściej używany do interpolacji. Będziemy rozpatrywać przedział  $\langle a, b \rangle$  i podzielimy go na  $n$  części, czyli na  $n$  podprzedziałów o długości  $h = \frac{b-a}{n}$ . Otrzymamy węzły równoodległe

$x_i = a + i \cdot h \quad i = 0, 1, \dots, n$ . **Funkcją sklejaną 3-iego stopnia** będziemy nazywać funkcję, która na każdym podprzedziale jest wielomianem 3 stopnia, ale posklejaną tak, aby była ciągła i miała pierwszą i drugą pochodną ciągłą na  $\langle a, b \rangle$ .

Aby dokładnie określić funkcję sklejaną 3-iego stopnia  $S_3(x)$  na przedziale  $\langle a, b \rangle$  określimy najpierw bazę splajnów 3-iego stopnia dla węzłów równoodległych. Jedna funkcja bazowa jest podana za pomocą bardzo skomplikowanego wzoru, ale musi spełniać powyższe wymagania, tzn.: musi być wielomianem 3-iego stopnia na każdym podprzedziale, mieć pierwszą i drugą pochodną ciągłą na  $\langle a, b \rangle$ . Funkcja bazowa o numerze  $i$ , oznaczona przez  $\Phi_i(x)$  ma w węźle o numerze  $i$  maksimum równe 4, w węzłach obok ma wartość 1, a w węzłach o numerach  $i-2$  i  $i+2$  ma wartość 0. Oto wzór i wykres takiej funkcji:

$$\Phi_i(x) = \frac{1}{h^3} \cdot \begin{cases} (x - x_{i-2})^3 & \text{dla } x \in \langle x_{i-2}, x_{i-1} \rangle \\ (x - x_{i-2})^3 - 4(x - x_{i-1})^3 & \text{dla } x \in \langle x_{i-1}, x_i \rangle \\ (x_{i+2} - x)^3 - 4(x_{i+1} - x)^3 & \text{dla } x \in \langle x_i, x_{i+1} \rangle \\ (x_{i+2} - x)^3 & \text{dla } x \in \langle x_{i+1}, x_{i+2} \rangle \\ 0 & \text{dla } x \in \mathbb{R} - \langle x_{i-2}, x_{i+2} \rangle \end{cases}$$

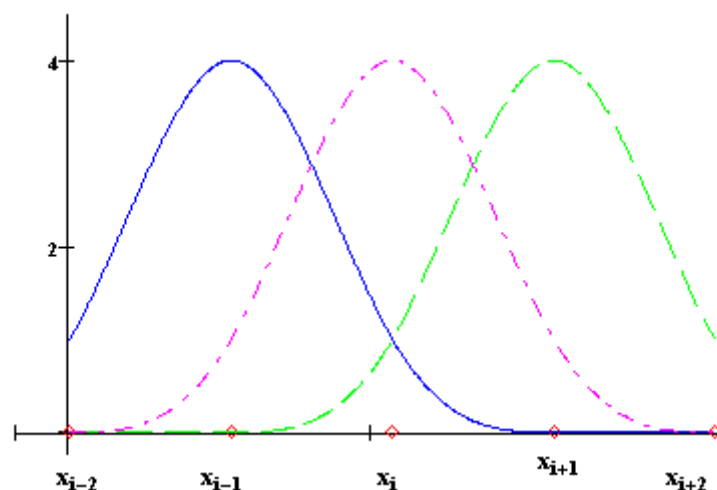


Rys5.1.1. Funkcja bazowa o numerze  $i$ , takim samym jak węzeł  $x_i$ .

Na podstawie tych funkcji bazowych będziemy określać w przedziale  $\langle a, b \rangle$  z  $n+1$  węzłami równoodległymi splajn 3-iego stopnia. Ile jest takich funkcji bazowych w przedziale  $\langle a, b \rangle$ ?

Narysujemy po kolei funkcje bazowe dla  $n=4$  tzn.: dla 5 węzłów. Narysujemy najpierw trzy  $\Phi_1(x)$ ,  $\Phi_2(x)$ ,  $\Phi_3(x)$ . Dodamy jeszcze dwie :  $\Phi_0(x)$ ,  $\Phi_4(x)$ . I mamy tyle funkcji ile jest węzłów. Ale są jeszcze dwie funkcje, które nie są równe 0 w całym przedziale, te funkcje odpowiadają węzłom, których na rysunku nie ma, jednemu o numerze wcześniejszym niż 0 i jednemu o numerze późniejszym niż 4. Te funkcje oznaczymy przez  $\Phi_{-1}(x)$ ,  $\Phi_{n+1}(x)$ . Okazuje się, że niezerowych funkcji jest o dwie więcej niż węzłów, czyli o 3 więcej niż  $n$ .

(1,2,3) (1,2,3,0,4) (wszystkie 7)



Rys5.1.2.Funkcje bazowe dla 5 węzłów.

Za pomocą tych funkcji zdefiniujemy funkcję sklejającą (splajn) 3-iego stopnia:

$$S_3(x) = c_{-1}\Phi_{-1}(x) + c_0\Phi_0(x) + c_1\Phi_1(x) + \dots c_n\Phi_n(x) + c_{n+1}\Phi_{n+1}(x) \quad (5.1.1)$$

lub w skrócie:

$$S_3(x) = \sum_{i=-1}^{n+1} c_i \Phi_i(x) \quad (5.1.2)$$

Współczynniki  $c_i$  są dowolnymi liczbami rzeczywistymi, będziemy je dobierać tak, aby splajn był funkcją interpolacyjną dla funkcji  $f(x)$ .



## 5.2 INTERPOLACJA SPLAJNAMI

Zastosujemy funkcję sklejającą  $S_3(x)$  do interpolacji funkcji  $f(x)$  danej w przedziale  $\langle a, b \rangle$ .

Dzielimy przedział na  $n$  części,  $h = \frac{b-a}{n}$ , węzły równoodległe  $x_i = a + i \cdot h \quad i = 0, 1, \dots, n$ .

Funkcja interpolacyjna musi się pokrywać w węzłach z funkcją  $f(x)$  tzn.:

$$S_3(x_i) = y_i = f(x_i) \quad i = 0, 1, 2, \dots, n \quad (5.2.1)$$

Otrzymaliśmy z tych związków  $n+1$  równań, a współczynników jest  $n+3$ , przypominamy wzór:

$$S_3(x) = c_{-1}\Phi_{-1}(x) + c_0\Phi_0(x) + c_1\Phi_1(x) + \dots + c_n\Phi_n(x) + c_{n+1}\Phi_{n+1}(x)$$

Nasz układ ma zatem dwa stopnie swobody i aby jednoznacznie wyznaczyć  $S_3(x)$  musimy mieć jeszcze dwa równania. Na ogół zadaje się wartości pochodnej funkcji  $S_3'(x)$  w punktach  $a$  i  $b$  - tzn.: zadaje się współczynniki kierunkowe stycznych pod jakimi funkcja interpolacyjna ma startować z punktu  $a$  w prawo i jak ma wpadać do  $b$  z lewej strony. Dodatkowe warunki to:  $S_3'(a^+) = \alpha$ ,  $S_3'(b^-) = \beta$ .

Z warunków (5.2.1) i z własności funkcji bazowych i ich pochodnych dostajemy układ równań:

$$c_{i-1} + 4c_i + c_{i+1} = y_i \quad i = 0, 1, \dots, n \quad (5.2.2)$$

$$\begin{aligned} -c_{-1} + c_1 &= \frac{h}{3} \cdot \alpha \\ -c_{n-1} + c_{n+1} &= \frac{h}{3} \cdot \beta \end{aligned} \quad (5.2.3)$$

Po wyliczeniu współczynników  $c_{-1}$ ,  $c_{n+1}$  z równań (5.2.3) i po wstawieniu ich do (5.2.2) otrzymujemy następujący układ  $n+1$  równań liniowych z  $n+1$  niewiadomymi  $c_0, c_1, \dots, c_n$ :

$$\begin{array}{rcl} 4c_0 + 2c_1 & & = y_0 + \frac{h}{3} \cdot \alpha \\ c_0 + 4c_1 + c_2 & & = y_1 \\ c_1 + 4c_2 + c_3 & & = y_2 \\ \dots & & \dots \\ c_{n-2} + 4c_{n-1} + c_n & & = y_{n-1} \\ 2c_{n-1} + 4c_n & & = y_n - \frac{h}{3} \cdot \beta \end{array} \quad (5.2.4)$$

Układ ten ma zawsze jedyne rozwiązanie na  $c_0, c_1, \dots, c_n$ , pozostałe 2 współczynniki obliczymy ze

wzorów:  $c_{-1} = c_1 - \frac{h}{3} \cdot \alpha$ ,  $c_{n+1} = c_{n-1} + \frac{h}{3} \cdot \beta$ .

Nie wyprowadzaliśmy układu równań, aby nie rozbudowywać tego tematu. Zainteresowanych obliczeniami odsyłamy do podanej literatury.

**Przykład 5.2.1:** Dana jest funkcja  $f(x) = x + \cos(2x)$  w przedziale  $<0, 5>$ . Znajdziemy dla niej funkcję sklejaną 3-iego stopnia dla różnej ilości węzłów równoodległych. Liczba  $n$  oznacza ilość podprzedziałów, węzłów jest  $n+1$ . Ponieważ  $f'(x) = 1 - 2\sin(2x)$  oraz  $\alpha = f'(0) = 1$ ,  $\beta = f'(5) = 2,088$  przyjmujemy, że  $S_3'(a^+) = \alpha = 1$ ,  $S_3'(b^-) = \beta = 2,088$ .

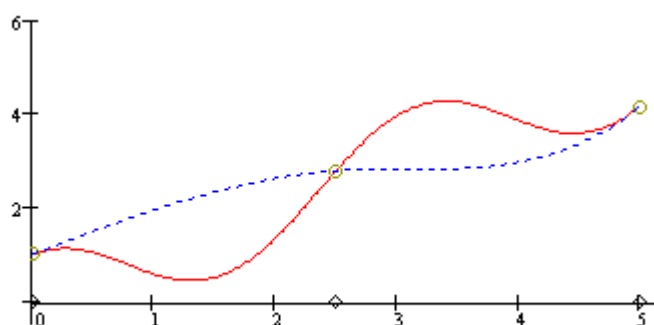
Rozwiązując układ dla  $n=2$  ( 3 węzły) dostajemy następujące współczynniki:

$$c_{-1} = -0,26; c_0 = 0,172; c_1 = 0,573; c_2 = 0,319; c_3 = 2,313$$

Narysujemy ten wykres, a później będziemy zwiększać liczbę podprzedziałów ( węzłów):

( nie podajemy współczynników następnych funkcji sklejanych, a tylko ich wykresy i błędy).

Dla



Rys5.2.1. Wykres funkcji  $f$  i interpolacyjnych splajnów dla  $n=2,3,4,5$  i 8.

Błędy dla  $n=2$  są jak widać na rysunku bardzo duże, rzędu 1,822, dla  $n=3$  maksymalny błąd interpolacji równa się 0,273, dla  $n=4$  błąd wynosi 0,306- jest większy!, dla  $n=5$  jest 0,097, a dla  $n=8$  już tylko 0,008. Na rysunku funkcje się pokrywają.

Obliczenia podaliśmy z dokładnością do trzech cyfr po przecinku.

## 6.1 APROKSYMACJA DYSKRETNA

Aproksymacja, tak jak interpolacja, służy do znajdowania przybliżonych wartości funkcji  $f(x)$  w dowolnym punkcie przedziału  $\langle a, b \rangle$ . Jednak funkcja aproksymacyjna na ogół jest inna niż funkcja interpolacyjna. W przypadku funkcji interpolacyjnej pokrywała się ona w pewnych punktach z funkcją interpolowaną, na funkcję aproksymacyjną nie będziemy narzucać takiego warunku. Będziemy od niej żądać aby była "bliska" funkcji aproksymowanej. Wyjaśnimy co będziemy uważać za "bliskość" dwóch funkcji. Ogólnie chodzi o to, aby wartości tych funkcji w pewnych wyróżnionych punktach były sobie bliskie. Założmy, podobnie jak w interpolacji, że z doświadczeń lub pomiarów określiliśmy w  $n+1$  różnych punktach :

$$x_0, x_1, x_2, \dots, x_n$$

z przedziału  $\langle a, b \rangle$  wartości funkcji  $y = f(x)$  i te wartości oznaczyliśmy przez:

$$y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$$

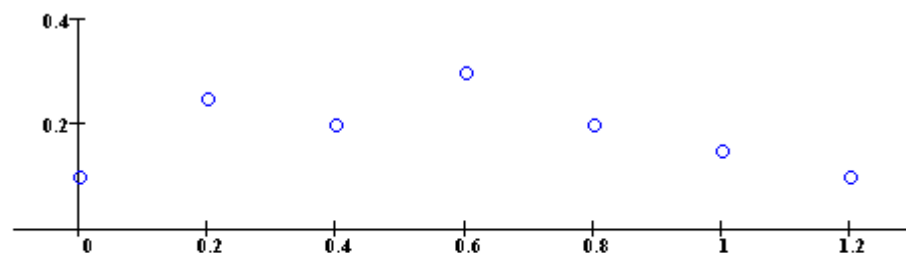
Funkcję aproksymacyjną oznaczmy przez  $F(x)$  i będziemy wymagać, aby kwadraty odległości między wartościami  $y_i$  a  $F(x_i)$  w sumie były jak najmniejsze, tzn.: aby suma

$$\sum_{i=0}^n (y_i - F(x_i))^2 \quad (6.1.1)$$

była minimalna. Taka metoda aproksymacyjna nazywana jest **metodą najmniejszych kwadratów**. Na rysunku przedstawiona jest na czerwono funkcja aproksymacyjna  $F(x)$  i są zaznaczone te odcinki (na czarno), których suma kwadratów długości ma być najmniejsza. Wartości funkcji  $F(x_i)$  oznaczone są przez  $F_i$ , wartości  $(x_i, y_i)$  są zaznaczone kółkami.

Na rysunku a) zaznaczone są węzły i wartości funkcji, na rysunku b) dochodzi jeszcze  $F(x)$ , na rysunku c) zaznaczone są dodatkowo odcinki- różnice między  $f(x)$  i  $F(x)$  w węzłach.

a) b) c)



Rys6.1.1. Wartości funkcji  $f$  w węzłach i funkcja aproksymacyjna.

## 6.2 FUNKCJE APROKSYMACYJNE

Założmy, że z doświadczeń lub pomiarów określiliśmy w  $n+1$  różnych punktach :

$$x_0, x_1, x_2, \dots, x_n$$

z przedziału  $< a, b >$  wartości funkcji  $y = f(x)$  i te wartości oznaczyliśmy przez:

$$y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$$

Będziemy rozpatrywać funkcje aproksymacyjne w różnej postaci, w szczególności wielomiany algebraiczne i wielomiany trygonometryczne. Jeśli za funkcję aproksymacyjną będziemy brać wielomian  $m$ -tego stopnia, to ten wielomian zapisywać będziemy w postaci:

$$W_m(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{m-1} x^{m-1} + a_m x^m \quad (6.2.1)$$

gdzie  $a_i \quad i = 0, 1, 2, \dots, m$  to współczynniki rzeczywiste wielomianu, które trzeba znaleźć.

Jeśli za funkcję aproksymacyjną będziemy brać  $m$ -ty wielomian trygonometryczny to będzie to następująca funkcja:

$$T_m(x) = a_0 + a_1 \cos(c \cdot x) + b_1 \sin(c \cdot x) + a_2 \cos(2c \cdot x) + b_2 \sin(2c \cdot x) + \dots + a_m \cos(mc \cdot x) + b_m \sin(mc \cdot x) \quad (6.2.2)$$

W takim  $m$ -tym wielomianie występują cosinusy i sinusy wielokrotności kąta  $cx$ , współczynnik  $c$  jest znany, niewiadome są współczynniki  $a_0, a_i, b_i \quad i = 1, 2, \dots, m$ .

Ogólnie, jeśli funkcja aproksymacyjna oparta będzie na  $m+1$  znanych niezależnych liniowo funkcjach bazowych:  $\varphi_0(x), \varphi_1(x), \dots, \varphi_m(x)$ , to będzie mieć postać:

$F(x) = a_0 \varphi_0(x) + a_1 \varphi_1(x) + \dots + a_m \varphi_m(x)$  gdzie  $a_i \quad i = 0, 1, 2, \dots, m$  są szukanymi współczynnikami.

Metoda najmniejszych kwadratów polega zatem na znalezieniu współczynników przy funkcjach bazowych takich, aby funkcja określająca sumę kwadratów odchylen  $\sum_{i=0}^n (y_i - F(x_i))^2$  była jak najmniejsza.

Oznaczmy przez:

$$H(a_0, a_1, \dots, a_m) = \sum_{i=0}^n (y_i - F(x_i))^2 = \sum_{i=0}^n (y_i - (a_0 \varphi_0(x_i) + a_1 \varphi_1(x_i) + \dots + a_m \varphi_m(x_i)))^2 \quad (6.2.3)$$

Będziemy szukać minimum tej funkcji  $m+1$  zmiennych  $a_i \quad i = 0, 1, 2, \dots, m$ . Punkty, w których funkcja wielu zmiennych przyjmuje minimum są punktami, w których zerują się pochodne cząstkowe (jeśli istnieją) tej funkcji po  $a_i \quad i = 0, 1, 2, \dots, m$ . Funkcja  $H(a)$  jest wielomianem ze względu na niewiadome  $a_i \quad i = 0, 1, 2, \dots, m$ , więc te pochodne istnieją i są ciągłe. Otrzymujemy

$$\frac{\partial H}{\partial \alpha_j} = 0 \quad j = 0, 1, 2, \dots, m$$
$$\frac{\partial H}{\partial a_j} = 2 \sum_{i=0}^n (y_i - (a_0 \varphi_0(x_i) + a_1 \varphi_1(x_i) + \dots + a_m \varphi_m(x_i))) \cdot (-\varphi_j(x_i)) = 0$$
$$a_0 \sum_{i=0}^n \varphi_0(x_i) \varphi_j(x_i) + a_1 \sum_{i=0}^n \varphi_1(x_i) \varphi_j(x_i) + \dots + a_m \sum_{i=0}^n \varphi_m(x_i) \varphi_j(x_i) = \sum_{i=0}^n y_i \varphi_j(x_i) \quad (6.2.4)$$
$$\begin{aligned} & a_0 \sum_{i=0}^n \varphi_0^2(x_i) + a_1 \sum_{i=0}^n \varphi_1(x_i) \varphi_0(x_i) + \dots + a_m \sum_{i=0}^n \varphi_m(x_i) \varphi_0(x_i) = \sum_{i=0}^n y_i \varphi_0(x_i) \\ & a_0 \sum_{i=0}^n \varphi_0(x_i) \varphi_1(x_i) + a_1 \sum_{i=0}^n \varphi_1^2(x_i) + \dots + a_m \sum_{i=0}^n \varphi_m(x_i) \varphi_1(x_i) = \sum_{i=0}^n y_i \varphi_1(x_i) \quad (6.2.5) \\ & \dots\dots\dots \\ & a_0 \sum_{i=0}^n \varphi_0(x_i) \varphi_m(x_i) + a_1 \sum_{i=0}^n \varphi_1(x_i) \varphi_m(x_i) + \dots + a_m \sum_{i=0}^n \varphi_m^2(x_i) = \sum_{i=0}^n y_i \varphi_m(x_i) \end{aligned}$$
$$M = \begin{bmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \dots & \varphi_m(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \dots & \varphi_m(x_1) \\ \dots & \dots & \dots & \dots \\ \varphi_0(x_n) & \varphi_1(x_n) & \dots & \varphi_m(x_n) \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \quad A = \begin{bmatrix} a_0 \\ a_1 \\ \dots \\ a_m \end{bmatrix} \quad (6.2.6)$$
$$M^T M \cdot A = M^T \cdot Y \quad (6.2.7)$$

2005-08-12

$a_i$   $i = 0, 1, 2, \dots, m$  istnieje, jest jedyne, zatem możemy znaleźć funkcję aproksymacyjną spełniającą narzucone warunki - minimalizacja sumy kwadratów różnic między funkcją daną a aproksymacyjną w wybranych punktach.

### Uwagi:

1. W interpolacji ilość węzłów narzucała stopień wielomianu interpolacyjnego, w aproksymacji wielomian może być stopnia 2, a ilość węzłów 100. Możemy sami sterować stopniem wielomianu. Jeśli natomiast węzłów będzie tyle ile funkcji bazowych to funkcja aproksymacyjna pokrywa się z funkcją interpolacyjną.
2. Macierz  $M$ , powyżej zdefiniowana, ma tyle wierszy ile jest węzłów, a tyle kolumn ile jest funkcji bazowych. W pierwszej kolumnie jest zerowa funkcja bazowa we wszystkich węzłach, w drugiej kolumnie następna funkcja bazowa dla wszystkich węzłów po kolei, w ostatniej kolumnie jest ostatnia funkcja bazowa dla wszystkich węzłów.
3. Macierz  $M^T M$  jest macierzą kwadratową wymiaru  $(m+1) \times (m+1)$ .

**Przykład 6.2.1:** Dla  $n+1=7$  punktów zmierzaliśmy wartości funkcji  $f(x)$  i otrzymaliśmy następujące wyniki ( w tabelce):

| $X_i$ | $Y_i$ |
|-------|-------|
| 0     | 0,1   |
| 0,2   | 0,25  |
| 0,4   | 0,2   |
| 0,6   | 0,3   |
| 0,8   | 0,2   |
| 1,0   | 0,15  |
| 1,2   | 0,1   |

Dla  $n = 6$  oraz  $i = 0, 1, \dots, n$ . Szukamy funkcji aproksymacyjnej  $F(x)$ ,

$$F(x) = a_0 + a_1 \sin x + a_2 e^x$$

to znaczy, że bazą dla tej funkcji jest układ  $\{1, \sin x, e^x\}$ .

Macierz  $M$  ma w pierwszej kolumnie same jedynki, w drugiej wartości funkcji  $\sin x$  dla wszystkich  $X_i$ , a w trzeciej kolumnie wartości funkcji  $e^x$  dla  $X_i$  i wygląda następująco (wyniki podaliśmy z dokładnością do trzech cyfr po przecinku):

$$M = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 0.199 & 1.221 \\ 1 & 0.389 & 1.492 \\ 1 & 0.565 & 1.822 \\ 1 & 0.717 & 2.226 \\ 1 & 0.841 & 2.718 \\ 1 & 0.932 & 3.32 \end{pmatrix}$$

Z układu równań  $M^T M \cdot A = M^T \cdot Y$ , w którym macierze  $M^T M$  i  $M^T Y$  są następujące:

$$M^T \cdot M = \begin{pmatrix} 7 & 3.644 & 13.799 \\ 3.644 & 2.601 & 8.831 \\ 13.799 & 8.831 & 31.403 \end{pmatrix} \quad M^T \cdot Y = \begin{pmatrix} 1.3 \\ 0.66 \\ 2.435 \end{pmatrix}$$

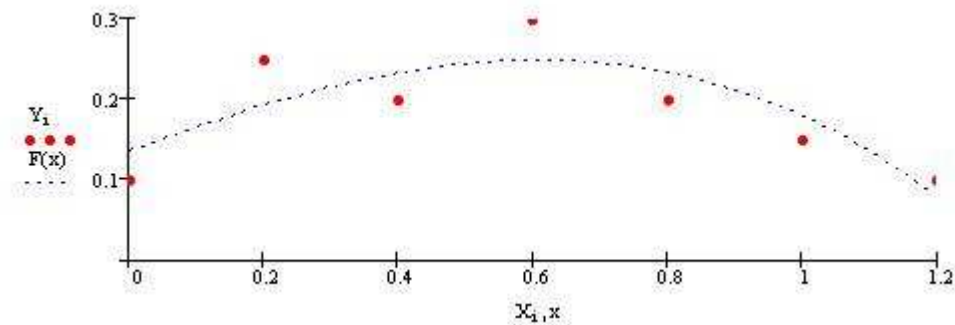
otrzymujemy współczynniki :

$$a = \begin{pmatrix} 0.39 \\ 0.572 \\ -0.255 \end{pmatrix}$$

Zatem szukana funkcja aproksymacyjna ma postać:

$$F(x) = 0,39 + 0.572 \sin x - 0,255e^x$$

Wykres:



Rys6.2.1. Wykres funkcji aproksymacyjnej  $F(x)$ .

### 6.3 APROKSYMACJA WIELOMIANAMI ALGEBRAICZNYMI

Założmy, że z doświadczeń lub pomiarów określiliśmy w  $n+1$  różnych punktach :

$$x_0, x_1, x_2, \dots, x_n$$

z przedziału  $< a, b >$  wartości funkcji  $y = f(x)$  i te wartości oznaczyliśmy przez:

$$y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$$

Będziemy rozpatrywać funkcje aproksymacyjne w postaci wielomianów algebraicznych:

$$W_m(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{m-1} x^{m-1} + a_m x^m \quad (6.3.1)$$

gdzie  $a_i$   $i = 0, 1, 2, \dots, m$  to współczynniki rzeczywiste wielomianu, które trzeba znaleźć. Bazą takiego wielomianu są funkcje:  $\{1, x, x^2, \dots, x^m\}$ . Zbudujemy macierz  $M$  dla tej bazy:

$$M = \begin{bmatrix} 1 & x_0 & \dots & x_0^m \\ 1 & x_1 & \dots & x_1^m \\ \dots & \dots & \dots & \dots \\ 1 & x_n & \dots & x_n^m \end{bmatrix} \quad (6.3.2)$$

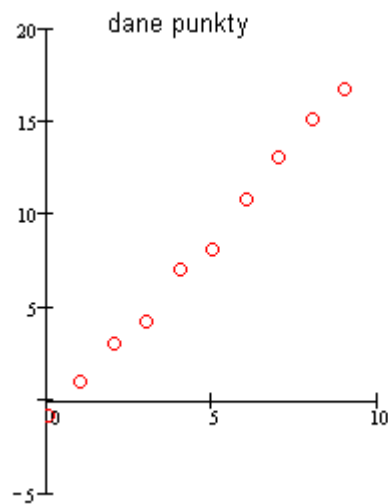
i aby znaleźć współczynniki  $a_i$   $i = 0, 1, 2, \dots, m$ , trzeba rozwiązać układ  $M^T M \cdot A = M^T \cdot Y$  z powyższą macierzą. Ponieważ macierz układu  $M^T M$  jest na ogół dla wysokich stopni wielomianów źle uwarunkowana (małe błędy danych powodują duże błędy wyników), stosuje się najczęściej aproksymację wielomianami niskich stopni tzn.:  $m=1, 2$  lub  $3$ . Różne programy numeryczne liczą wskaźniki uwarunkowania macierzy i rozwiązują układy równań liniowych. Prześledzimy tylko jeszcze raz powstawanie tego układu dla wielomianu pierwszego i drugiego stopnia.

Założmy, że z doświadczeń dostaliśmy takie wartości badanej funkcji, że punkty  $(x_i, y_i)$  ułożyły się tak, jak na wykresie:

Na rysunku a) są tylko dane punkty, na rysunku b) również wielomian interpolacyjny stopnia 1- oznaczony jako  $F(x)$ .

a) b)





Rys6.3.1. Dane i wykres funkcji liniowej- wielomianu aproksymacyjnego 1 stopnia.

Wtedy naturalnie jest stosować jako funkcję aproksymacyjną wielomian pierwszego stopnia czyli  $F(x) = a_0 + a_1 x$ . Jeśli będziemy korzystać z powyższych gotowych wzorów to dostaniemy:

$$M = \begin{bmatrix} 1 & x_0 \\ 1 & x_1 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \quad M^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_0 & x_1 & \dots & x_n \end{bmatrix} \quad M^T M = \begin{bmatrix} n+1 & \sum_{i=0}^n x_i \\ \sum_{i=0}^n x_i & \sum_{i=0}^n x_i^2 \end{bmatrix} \quad M^T Y = \begin{bmatrix} \sum_{i=0}^n y_i \\ \sum_{i=0}^n x_i y_i \end{bmatrix}$$

Zatem układ na współczynniki niewiadome  $a_0, a_1$  będzie następujący:

$$\begin{aligned} a_0(n+1) + a_1 \sum_{i=0}^n x_i &= \sum_{i=0}^n y_i \\ a_0 \sum_{i=0}^n x_i + a_1 \sum_{i=0}^n x_i^2 &= \sum_{i=0}^n x_i y_i \end{aligned} \quad (6.3.3)$$

Ten sam układ otrzymamy wracając do funkcji :

$$H(a_0, a_1) = \sum_{i=0}^n (y_i - F(x_i))^2 = \sum_{i=0}^n (y_i - (a_0 + a_1 x_i))^2$$

i obliczając jej minimum .

Obliczone pochodne cząstkowe przyrównamy do zera :

$$(6.3.4)$$

$$\frac{\partial H}{\partial a_0} = 2 \sum_{i=0}^n (y_i - (a_0 + a_1 x_i))(-1) = 0$$

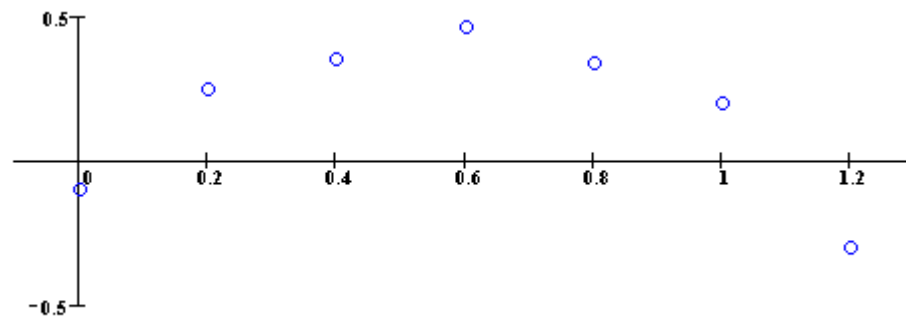
$$\frac{\partial H}{\partial a_1} = 2 \sum_{i=0}^n (y_i - (a_0 + a_1 x_i))(-x_i) = 0$$

A stąd otrzymamy ten sam układ co podany powyżej.

Jeśli punkty pomiarowe ułożą się tak jak na rysunku poniżej, to nie ma co szukać funkcji aproksymacyjnej jako wielomianu stopnia 1, tylko co najmniej stopnia 2.

Na rysunku c) są tylko dane punkty, na rysunku d) również wielomian interpolacyjny stopnia 2

c) d)



Rys6.3.2. Dane i wykres funkcji kwadratowej- wielomianu aproksymacyjnego 2 stopnia.

W tym wypadku za funkcję aproksymacyjną można przyjąć wielomian  $F(x) = a_0 + a_1 x + a_2 x^2$ .  
Wtedy

$$M = \begin{bmatrix} 1 & x_0 & x_0^2 \\ 1 & x_1 & x_1^2 \\ \dots & \dots & \dots \\ 1 & x_n & x_n^2 \end{bmatrix} \quad M^T M = \begin{bmatrix} n+1 & \sum_{i=0}^n x_i & \sum_{i=0}^n x_i^2 \\ \sum_{i=0}^n x_i & \sum_{i=0}^n x_i^2 & \sum_{i=0}^n x_i^3 \\ \sum_{i=0}^n x_i^2 & \sum_{i=0}^n x_i^3 & \sum_{i=0}^n x_i^4 \end{bmatrix} \quad M^T Y = \begin{bmatrix} \sum_{i=0}^n y_i \\ \sum_{i=0}^n x_i y_i \\ \sum_{i=0}^n x_i^2 y_i \end{bmatrix}$$

I układ na współczynniki  $a_0, a_1, a_2$  jest następujący:

$$\begin{aligned}
 a_0(n+1) + a_1 \sum_{i=0}^n x_i + a_2 \sum_{i=0}^n x_i^2 &= \sum_{i=0}^n y_i \\
 a_0 \sum_{i=0}^n x_i + a_1 \sum_{i=0}^n x_i^2 + a_2 \sum_{i=0}^n x_i^3 &= \sum_{i=0}^n x_i y_i \\
 a_0 \sum_{i=0}^n x_i^2 + a_1 \sum_{i=0}^n x_i^3 + a_2 \sum_{i=0}^n x_i^4 &= \sum_{i=0}^n x_i^2 y_i
 \end{aligned} \tag{6.3.5}$$

Taki sam układ otrzymamy, jeśli określimy funkcję:

$$H(a_0, a_1, a_2) = \sum_{i=0}^n (y_i - F(x_i))^2 = \sum_{i=0}^n (y_i - (a_0 + a_1 x_i + a_2 x_i^2))^2$$

obliczymy jej pochodne cząstkowe po  $a_0, a_1, a_2$  i przyrównamy je do zera.

**Przykład 6.3.1:** Funkcja  $f(x)$  jest dana w 7 węzłach za pomocą tabelki:

| $X_i :=$ | $Y_i :=$ |
|----------|----------|
| 0        | 0.5      |
| 0.2      | 0.4      |
| 0.4      | 0.3      |
| 0.6      | 0.3      |
| 0.8      | 0.2      |
| 1.0      | 0.15     |
| 1.2      | 0.1      |

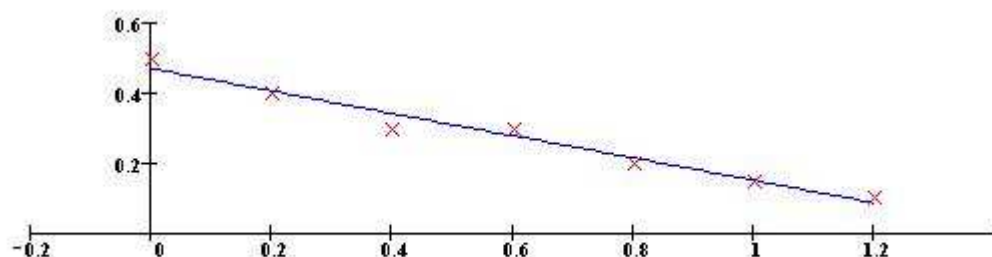
Będziemy szukać wielomianu aproksymacyjnego stopnia 1. Dla tych danych macierz  $M$  ma w pierwszej kolumnie jedynki, w drugiej węzły. Układ dwóch równań na współczynniki  $a$  jest bardzo prosty:

$$7a_0 + 4.2a_1 = 1.95$$

$$4.2a_0 + 3.64a_1 = 0.81$$

Rozwiązując ten układ dostajemy :  $a_0 = 0.471, a_1 = -0.321$

Zatem wielomian aproksymacyjny pierwszego stopnia ma postać:  $W_1(x) = 0.471 - 0.321x$



Rys6.3.3. Wykres funkcji  $W_1(x)$ .

Możemy obliczać za pomocą tego wielomianu wartość funkcji poza danymi punktami w przedziale .

Będziemy obliczać przybliżoną wartość funkcji w punkcie:  $x=$

Przybliżona wartość funkcji dla podanego  $x$  wynosi  $f(x)$

=

## 7.1 WIELOMIANY TRYGONOMETRYCZNE

Założmy, że z doświadczeń lub pomiarów określiliśmy w  $n+1$  różnych punktach :

$$x_0, x_1, x_2, \dots, x_n$$

z przedziału  $< a, b >$  wartości funkcji  $y = f(x)$  i te wartości oznaczyliśmy przez:

$$y_0 = f(x_0), y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$$

Będziemy rozpatrywać jako funkcje aproksymacyjne wielomiany trygonometryczne:

$$T_m(x) = a_0 + a_1 \cos(c \cdot x) + b_1 \sin(c \cdot x) + a_2 \cos(2c \cdot x) + b_2 \sin(2c \cdot x) + \dots + a_m \cos(mc \cdot x) + b_m \sin(mc \cdot x) \quad (7.1.1)$$

W takim  $m$ -tym wielomianie występują cosinusy i sinusy wielokrotności kąta  $cx$ , współczynnik  $c$  jest znany, niewiadome są współczynniki  $a_0, a_i, b_i$   $i = 1, 2, \dots, m$ .

Ograniczymy się do węzłów równoodległych, podzielimy przedział  $< a, b >$  na  $n$  części, otrzymamy podprzedziały o długości  $h = \frac{b-a}{n}$  i węzły  $x_i = a + i \cdot h$ ,  $i=0, 1, \dots, n$ .

Ze względu na okresowość funkcji  $\sin$  i  $\cos$  przyjmujemy  $c = \frac{\pi}{l}$ ,  $l = \frac{n+1}{2}h$ . Wtedy pierwszy wielomian trygonometryczny ma postać:

$T_1(x) = a_0 + a_1 \cos(\frac{\pi}{l}x) + b_1 \sin(\frac{\pi}{l}x)$  gdzie współczynniki  $a_0, a_1, b_1$  wyliczamy z zerowania się pochodnych cząstkowych po  $a_0, a_1, b_1$  funkcji :

$$H(a_0, a_1, b_1) = \sum \left( y_i - (a_0 + a_1 \cos(\frac{\pi}{l}x_i) + b_1 \sin(\frac{\pi}{l}x_i)) \right)^2 \quad (7.1.2)$$

lub budujemy macierz  $M$  :

$$M = \begin{bmatrix} 1 & \cos(\frac{\pi}{l}x_0) & \sin(\frac{\pi}{l}x_0) \\ 1 & \cos(\frac{\pi}{l}x_1) & \sin(\frac{\pi}{l}x_1) \\ \dots & \dots & \dots \\ 1 & \cos(\frac{\pi}{l}x_n) & \sin(\frac{\pi}{l}x_n) \end{bmatrix} \quad M^T M = \begin{bmatrix} n+1 & 0 & 0 \\ 0 & \frac{n+1}{2} & 0 \\ 0 & 0 & \frac{n+1}{2} \end{bmatrix} \quad M^T Y = \begin{bmatrix} \sum_{i=0}^n y_i \\ \sum_{i=0}^n y_i \cos(\frac{\pi}{l}x_i) \\ \sum_{i=0}^n y_i \sin(\frac{\pi}{l}x_i) \end{bmatrix}$$

W tym wypadku macierz układu  $M^T M$  jest dobrze uwarunkowana, jest macierzą diagonalną i układ ma bardzo prostą postać:

$$(n+1)a_0 = \sum_{i=0}^n y_i, \quad \frac{n+1}{2}a_1 = \sum_{i=0}^n y_i \cos\left(\frac{\pi}{l} x_i\right), \quad \frac{n+1}{2}b_1 = \sum_{i=0}^n y_i \sin\left(\frac{\pi}{l} x_i\right) \quad (7.1.3)$$

Możemy podać wzory na współczynniki  $a_0, a_1, b_1$ :

$$a_0 = \frac{1}{n+1} \sum_{i=0}^n y_i, \quad a_1 = \frac{2}{n+1} \sum_{i=0}^n y_i \cos\left(\frac{\pi}{l} x_i\right), \quad b_1 = \frac{2}{n+1} \sum_{i=0}^n y_i \sin\left(\frac{\pi}{l} x_i\right) \quad (7.1.4)$$

Dla drugiego wielomianu trygonometrycznego:

$T_2(x) = a_0 + a_1 \cos\left(\frac{\pi}{l} x\right) + b_1 \sin\left(\frac{\pi}{l} x\right) + a_2 \cos\left(2\frac{\pi}{l} x\right) + b_2 \sin\left(2\frac{\pi}{l} x\right)$  współczynniki można wyliczyć analogicznie do powyższych, dostajemy wtedy:

$$a_0 = \frac{1}{n+1} \sum_{i=0}^n y_i, \quad a_1 = \frac{2}{n+1} \sum_{i=0}^n y_i \cos\left(\frac{\pi}{l} x_i\right), \quad b_1 = \frac{2}{n+1} \sum_{i=0}^n y_i \sin\left(\frac{\pi}{l} x_i\right),$$

$$a_2 = \frac{2}{n+1} \sum_{i=0}^n y_i \cos\left(2\frac{\pi}{l} x_i\right), \quad b_2 = \frac{2}{n+1} \sum_{i=0}^n y_i \sin\left(2\frac{\pi}{l} x_i\right) \quad (7.1.5)$$

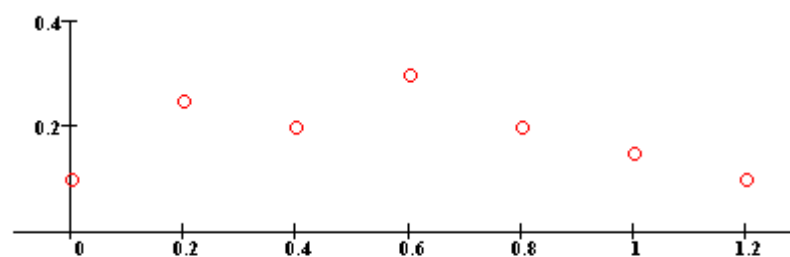
Można uogólnić to postępowanie dla wielomianów trygonometrycznych zawierających wyrazy z  $\cos$  i  $\sin$  kąta  $m$ -krotnego.

**Przykład 7.1.1:** Funkcja jest dana za pomocą tabelki:

| $X_i$ | $Y_i$ |
|-------|-------|
| 0     | 0,08  |
| 0,2   | 0,2   |
| 0,4   | 0,25  |
| 0,6   | 0,27  |
| 0,8   | 0,22  |
| 1,0   | 0,15  |
| 1,2   | 0,1   |

Dane:  $n=6, i=0,1,\dots,n, h=0,2, l=\frac{n+1}{2}h=0,7$ . Rysunki: a) dane, b) wielomian  $T_1(x)$ , c) wielomian  $T_2(x)$

a) b) c)



Rys7.1.1. Dane i dwa wielomiany aproksymacyjne trygonometryczne.

Po obliczeniu współczynników korzystając z powyższych wzorów otrzymujemy:

$$T_1(x) = 0,181 - 0,075 \cos\left(\frac{\pi}{l}x\right) + 0,056 \sin\left(\frac{\pi}{l}x\right)$$

$$T_2(x) = 0,181 - 0,075 \cos\left(\frac{\pi}{l}x\right) + 0,056 \sin\left(\frac{\pi}{l}x\right) - 0,012 \cos\left(2\frac{\pi}{l}x\right) + 0,004 \sin\left(2\frac{\pi}{l}x\right)$$

## 7.2 BŁĄD APROKSYMACJI

Czym będziemy się kierować decydując się na tę, a nie inną funkcję aproksymacyjną? Ponieważ chcemy, aby suma kwadratów odchyleń między funkcją daną a funkcją aproksymacyjną w węzłach była jak najmniejsza, możemy przyjąć dla prostoty, że najlepsza będzie ta funkcja, dla której ta suma jest jak najmniejsza. Będziemy posługiwać się wzorem na średni błąd przypadający na jeden węzeł, to znaczy

$$bl = \sqrt{\frac{\sum_{i=0}^n (y_i - F(x_i))^2}{n+1}} \quad (7.2.1)$$

We wzorze pod pierwiastkiem w liczniku jest suma kwadratów odchyleń, którą minimalizowaliśmy, w mianowniku jest ilość węzłów.

Jeśli będziemy szukać funkcji aproksymacyjnej spośród danych możliwych, wybierać będziemy tę dla której wartość powyższego błędu  $bl$  jest najmniejsza.

W przykładzie (poniżej) na wielomiany trygonometryczne te błędy są odpowiednio równe: dla podanego w poprzednim temacie wielomianu  $T_1(x)$  błąd średni wynosi 0,014, dla wielomianu  $T_2(x)$  błąd średni równa się 0,010.

**Przykład 7.2.1:** Funkcja jest dana za pomocą tabelki:

| $X_i$ | $Y_i$ |
|-------|-------|
| 0     | 0,08  |
| 0,2   | 0,2   |
| 0,4   | 0,25  |
| 0,6   | 0,27  |
| 0,8   | 0,22  |
| 1,0   | 0,15  |
| 1,2   | 0,1   |

Dane:  $n = 6$ ,  $i = 0, 1, \dots, n$ ,  $h = 0,2$ ,  $l = \frac{n+1}{2}h = 0,7$ ,

$$T_1(x) = 0,181 - 0,075 \cos\left(\frac{\pi}{l}x\right) + 0,056 \sin\left(\frac{\pi}{l}x\right)$$

$$T_2(x) = 0,181 - 0,075 \cos\left(\frac{\pi}{l}x\right) + 0,056 \sin\left(\frac{\pi}{l}x\right) - 0,012 \cos\left(2\frac{\pi}{l}x\right) + 0,004 \sin\left(2\frac{\pi}{l}x\right)$$

Wyberzemy z tych dwóch funkcji wielomian drugi, bo ma mniejszy średni błąd.

Można stosować inne kryteria doboru funkcji aproksymacyjnej, czasami stosuje się błąd średni statystyczny, ale nie będziemy komplikować rozważań i zostaniemy przy tym najprostszym



wzorze na blad.

### 7.3 APROKSYMACJA CIĄGŁA

Będziemy aproksymować funkcję ciągłą  $f(x)$  w przedziale  $\langle a, b \rangle$  funkcją  $F(x)$  postaci:

$F(x) = a_0 \varphi_0(x) + a_1 \varphi_1(x) + \dots + a_m \varphi_m(x)$  gdzie  $a_i \quad i = 0, 1, 2, \dots, m$  są szukanyymi

współczynnikami. Jeśli założymy, że funkcje bazowe  $\varphi_i(x) \quad i = 0, 1, \dots, m$  są w przedziale  $\langle a, b \rangle$

całkowalne z kwadratem ( tzn. istnieje skończona wartość całki  $\int_a^b (\varphi_i(x))^2 dx$  ) to funkcję

aproksymacyjną  $F(x)$  będziemy poszukiwać taką, aby funkcja:

$$H(a_0, a_1, \dots, a_m) = \int_a^b (f(x) - F(x))^2 dx = \int_a^b (f(x) - (a_0 \varphi_0(x) + \dots + a_m \varphi_m(x)))^2 dx \quad (7.3.1)$$

miała jak najmniejszą wartość. Podobnie jak poprzednio, gdy funkcja  $f(x)$  dana była tylko w skończonej ilości punktów, warunkiem koniecznym na minimum funkcji  $H(a_0, a_1, \dots, a_m)$  jest

zerowanie się pochodnych cząstkowych  $\frac{\partial H}{\partial a_j} = 0 \quad j = 0, 1, 2, \dots, m$ . I tak jak poprzednio układ tych

równań posiada jednoznaczne rozwiązanie na współczynniki  $a_i \quad i = 0, 1, 2, \dots, m$ , a warunek

konieczny w tym wypadku zapewnia ( ze względu na postać funkcji  $H(a_0, a_1, \dots, a_m)$  ) istnienie minimum.

**Przykład 7.3.1:** Wyznaczyć wielomian aproksymacyjny pierwszego stopnia najlepiej aproksymujący funkcję  $f(x) = \frac{1}{x}$  w przedziale  $\langle 1, 2 \rangle$ .

Funkcja aproksymacyjna ma postać  $F(x) = a_0 + a_1 x$  gdzie dwa współczynniki znajdziemy z zerowania się pochodnych funkcji:

$$H(a_0, a_1) = \int_1^2 (f(x) - (a_0 + a_1 x))^2 dx \text{ po } a_0 \text{ i po } a_1.$$

$$\frac{\partial H}{\partial a_0} = (-2) \int_1^2 (f(x) - (a_0 + a_1 x)) dx = 0$$

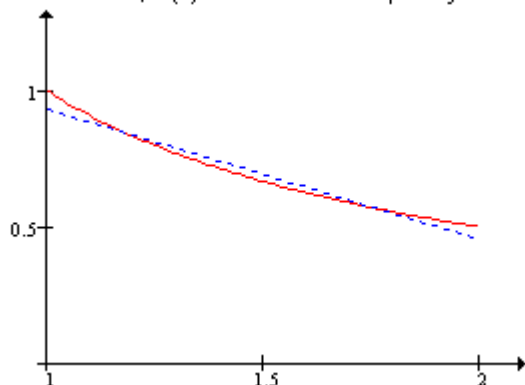
$$\frac{\partial H}{\partial a_1} = (-2) \int_1^2 (f(x) - (a_0 + a_1 x)) x dx = 0$$

Stąd otrzymujemy układ równań:

$$\begin{aligned} \int_1^2 (a_0 + a_1 x) dx &= \int_1^2 \frac{1}{x} dx \\ \int_1^2 (a_0 x + a_1 x^2) dx &= \int_1^2 \frac{1}{x} \cdot x dx \Rightarrow \begin{cases} a_0 + \frac{3}{2}a_1 = \ln 2 \\ \frac{3}{2}a_0 + \frac{7}{3}a_1 = 1 \end{cases} \end{aligned}$$

Rozwiązaniem tego układu są liczby:  $a_0 = 1,408$ ,  $a_1 = -0,477$ , zatem szukana funkcja aproksymacyjna ma postać:  $F(x) = 1,408 - 0,477 x$ .

$f(x)$  - linia czerwona,  $F(x)$  - linia niebieska przerywana



Rys7.3.1. Wykres funkcji  $f$  i funkcji aproksymacyjnej  $F$ .

## 8.1 WPROWADZENIE

Równanie nieliniowe np. równanie kwadratowe, logarytmiczne, wykładnicze, trygonometryczne, będziemy zapisywać ogólnie jako równanie postaci:

$$f(x) = 0$$

gdzie funkcja  $f$  jest funkcją nieliniową zmiennej rzeczywistej  $x$  i jest funkcją ciągłą na pewnym przedziale skończonym lub nieskończonym  $(a, b)$ . Poznane algorytmy rozwiązywania wymienionych równań dotyczą jednak pewnej wąskiej klasy funkcji  $f(x)$ , np. wielomianów stopnia nie większego niż czwarty, natomiast olbrzymiej klasy równań nieliniowych - głównie równań przestępnych - nie da się rozwiązać dokładnie. Potrzebne są metody przybliżone, które umożliwiają znalezienie pierwiastków rzeczywistych tych równań z góry podaną dokładnością. Niektóre z tych metod, oraz problemy z nimi związane, będą przedstawione w tej lekcji.

*Definicja.* Liczbę rzeczywistą  $p$ , która spełnia równanie  $f(x) = 0$  tzn. dla której  $f(p) \equiv 0$ , nazywamy pierwiastkiem rzeczywistym równania lub zerem funkcji  $f$ .

*Definicja.* Liczbę rzeczywistą  $p$  nazywamy  $k$ -krotnym pierwiastkiem równania  $f(x) = 0$  lub  $k$ -krotnym miejscem zerowym funkcji  $f$ , jeśli dla wartości  $p$  funkcja i jej pochodne do  $k-1$  rzędu włącznie przyjmują wartość zero, natomiast wartość pochodnej  $k$ -tego rzędu jest różna od zera, tzn:

$$f(p) \equiv 0, \quad f'(p) \equiv 0, \quad \dots \quad f^{(k-1)}(p) \equiv 0, \quad f^{(k)}(p) \neq 0.$$

Przedstawione poniżej przybliżone metody rozwiązywania równań nieliniowych, można stosować jedynie pod warunkiem, że znany jest pewien przedział, w którym znajduje się jeden i tylko jeden pierwiastek rzeczywisty danego równania. Taki przedział będziemy nazywać *przedziałem izolacji* dla równania  $f(x) = 0$ . Przedziały izolacji, przed przystąpieniem do rozwiązywania równań, będziemy wyznaczać graficznie.

Wybrane przybliżone metody rozwiązywania równań nieliniowych są metodami iteracyjnymi, polegającymi na budowaniu ciągu przybliżeń liczb rzeczywistych:

$$x_0, x_1, x_2, \dots, x_n, \dots$$

zbieżnego do szukanego rozwiązania - pierwiastka  $p$  - równania  $f(x) = 0$ . Oczywiście, aby ciąg  $\{x_n\}$  był zbieżny do pierwiastka  $p$  niezbędne są na ogół dodatkowe, oprócz ciągłości, założenia o funkcji  $f$ , które będą podane przy każdej metodzie osobno.

Ograniczymy się do metod jednokrokowych, polegających na znalezieniu następnego przybliżenia  $x_{n+1}$ , mając dane jego poprzednie przybliżenie  $x_n$ , oraz do metod dwukrokowych, polegających na znalezieniu następnego przybliżenia  $x_{n+1}$ , mając obliczone dwa poprzednie  $x_{n-1}$  i  $x_n$ . Formuły określające ciągi przybliżeń można zapisać ogólnie w postaci:

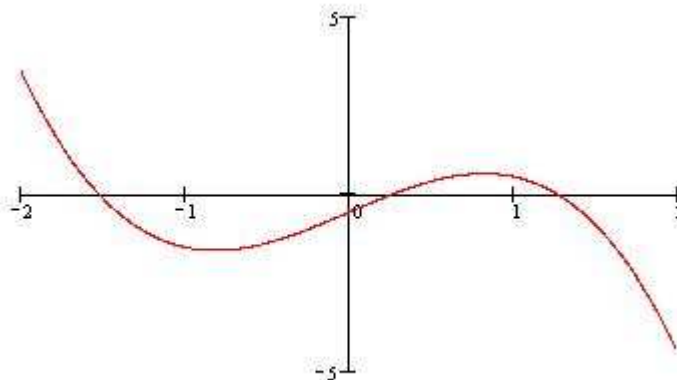
$$x_{n+1} = F(x_n) \quad \text{- dla metod jednokrokowych,}$$

$$x_{n+1} = F(x_{n-1}, x_n) \quad \text{- dla metod dwukrokowych.}$$

Będziemy korzystać z następujących twierdzeń:

*Twierdzenie.* Jeżeli funkcja  $f$  ciągła w przedziale  $[a, b]$  ma na końcach tego przedziału różne znaki tzn.  $f(a)f(b) < 0$ , to wewnątrz tego przedziału istnieje co najmniej jeden pierwiastek równania  $f(x) = 0$ .

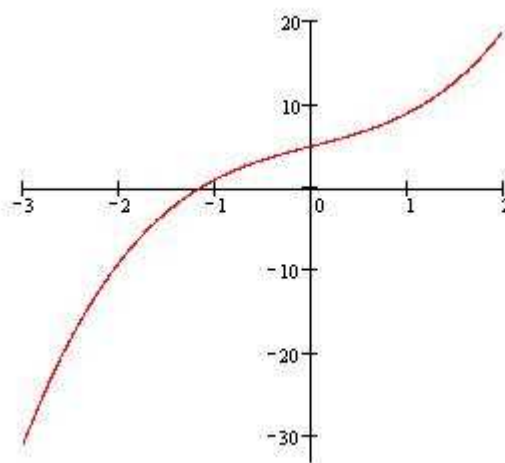
Ilustracja twierdzenia: Na rysunku funkcja jest ciągła i  $f(b) < 0$ ,  $f(a) > 0$ . Funkcja przecina oś  $Ox$  trzy razy, ma zatem trzy pierwiastki rzeczywiste w przedziale  $[a, b]$ .



Rys8.1.1. Wykres funkcji ciągłej  $f$  zmieniającej znak trzy razy w przedziale  $[-2, 2]$ .

*Twierdzenie.* Jeżeli w przedziale  $(a, b)$  istnieje pochodna  $f'(x)$  i nie zmienia znaku w tym przedziale tzn. albo jest w nim cały czas dodatnia albo ujemna, a  $f(a)f(b) < 0$ , to równanie  $f(x) = 0$  ma dokładnie jeden pierwiastek jednokrotny.

Ilustracja twierdzenia: Na rysunku funkcja jest ciągła, ma ciągłą pochodną, która jest cały czas dodatnia w  $(a, b)$ , tzn.: funkcja w  $(a, b)$  rośnie, oraz  $f(a)f(b) < 0$ , raz tylko wykres funkcji przecina się z osią  $Ox$ .



Rys8.1.2. Wykres funkcji ciągłej  $f$ , zmieniającej znak raz w przedziale  $[-3, 2]$ .

—

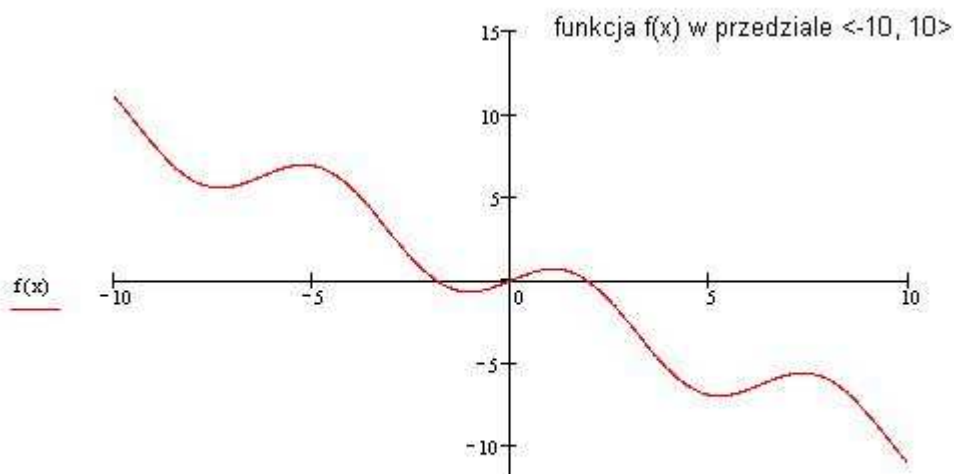
## 8.2 IZOLACJA PIERWIASTKÓW

Najprostszą metodą przekonania się czy równanie  $f(x) = 0$  posiada pierwiastki rzeczywiste jest narysowanie funkcji  $f$  i osi  $0x$ , i przekonanie się czy funkcja przecina się z osią. Można tym sposobem znaleźć przedziały, które się nie przecinają, a w których znajduje się po jednym pierwiastku danego równania. Zilustrujemy to na przykładzie:

**Przykład 8.2.1:** Znaleźć przedziały izolacji równania:  $2 \sin x - x = 0$ .

Rysujemy w kartezjańskim układzie współrzędnych funkcję  $f$  na możliwie dużym przedziale, a potem tak zawężamy ten przedział, aby nie stracić pierwiastków - aby wszystkie nadal były na wykresie. Rysunek a) przedstawia funkcję w przedziale  $<-10, 10>$ , rysunek b) w przedziale  $<-5, 5>$ , a rysunek c) w przedziale  $<-3, 3>$ .

a) b) c)



Rys8.2.1. Izolacja pierwiastków równania  $f(x)=0$ .

Jak widać na rysunku funkcja  $f(x)$  dąży do minus nieskończoności, gdy argumenty dążą do minus nieskończoności i dąży do minus nieskończoności, gdy argumenty dążą do nieskończoności. Można przyjąć, że wszystkie pierwiastki danego równania tzn. wszystkie punkty przecięcia funkcji z osią  $0x$ , są w przedziale  $<-3, 3>$ . Widzimy, że równanie ma pierwiastek  $x=0$ , co łatwo sprawdzić, drugi pierwiastek w przedziale np.  $(1,5; 2,5)$  i trzeci pierwiastek w przedziale np.  $(-2,5; -1)$ . Zatem równanie ma trzy przedziały izolacji:  $(-2,5; -1)$ ,  $(-0,5; 0,5)$  i  $(1; 2,5)$ . Oczywiście widać, że przedziały izolacji nie są wyznaczone jednoznacznie, można podać je w postaci:  $(-2,2; -1,3)$ ,  $(-0,2; 0,2)$  i  $(1,3; 2,2)$ . Więcej pierwiastków równanie nie posiada.

Inną metodą graficzną jest rysowanie równania  $f(x) = 0$  za pomocą dwóch funkcji, jeśli równanie da się przedstawić w postaci różnicy funkcji  $g(x) - h(x) = 0$ . Punkty przecięcia funkcji  $g$  i  $h$ , są pierwiastkami równania  $f(x) = 0$ .

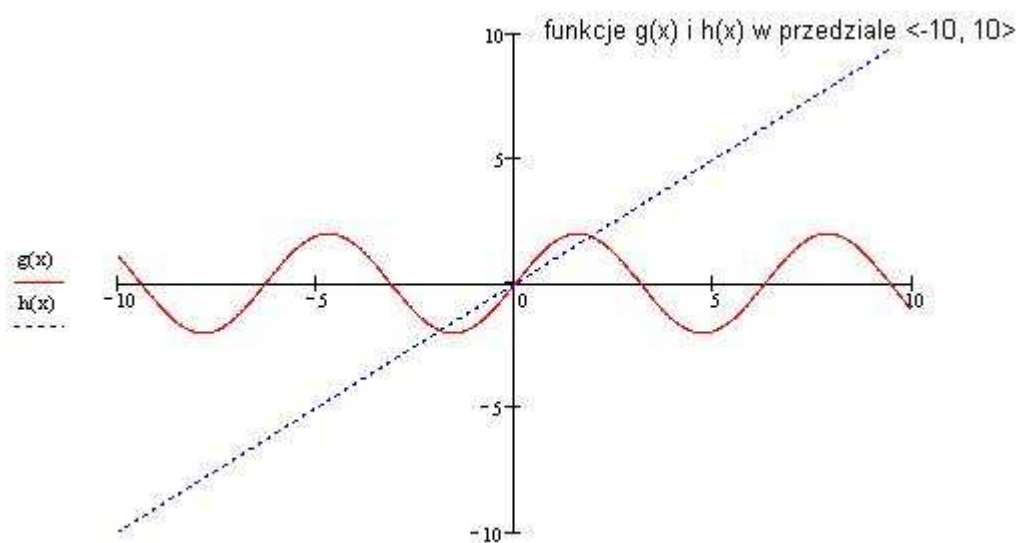
**Przykład 8.2.2:** Ponieważ rozpatrywane równanie można przedstawić jako równanie :

$$2 \sin x = x$$

to oznaczając przez  $g(x) = 2 \sin x$ , a przez  $h(x) = x$ , narysujemy te funkcje w takim samym przedziale jak poprzednio funkcję  $f$ .

Funkcje  $g$  i  $h$  przecinają się w zerze i dla tych samych  $x$ , w których funkcja  $f$  przecina się z osią  $0x$ . Funkcje te nie przecinają się w innym przedziale, bo funkcja  $g$  jest funkcją ograniczoną o wartościach w przedziale  $[-2, 2]$ , a funkcja  $h$  rośnie na prawo do nieskończoności, a na lewo maleje do minus nieskończoności. Rysunek a) przedstawia funkcje w przedziale  $<-10, 10>$ , rysunek b) w przedziale  $<-5, 5>$ , a rysunek c) w przedziale  $<-3, 3>$ .

a) b) c)



Rys8.2.2. Izolacja pierwiastków równania  $g(x)=h(x)$ .

### 8.3 UWAGI O DOKŁADNOŚCI

Istotnym problemem w metodach iteracyjnych jest decyzja, którą iterację wziąć za przybliżenie pierwiastka równania  $f(x) = 0$ , co ma decydować o zakończeniu postępowania iteracyjnego (wyboru warunku "stopu"), lub jak się da oszacować przyjęte przybliżenie w stosunku do nieznannej dokładnej wartości pierwiastka. Do każdej metody będzie ten problem rozważany osobno, tutaj podamy jak można wykorzystać następujące twierdzenie:

**Twierdzenie** . Niech  $p$  będzie dokładną, a  $x^*$  przybliżoną wartością pierwiastka równania  $f(x) = 0$ , przy czym obie te liczby znajdują się w przedziale domkniętym  $[a, b]$ . Jeśli  $f$  posiada pochodną i jeśli dla  $x$  z przedziału  $[a, b]$  zachodzi nierówność  $|f'(x)| \geq m_1 > 0$  to prawdziwe jest oszacowanie :

$$|x^* - p| \leq \frac{|f(x^*)|}{m_1}$$

*Dowód:* Stosując wzór Lagrange'a otrzymujemy:  $f(x^*) - f(p) = (x^* - p)f'(c)$  gdzie wartość  $c$  jest liczbą między  $p$  i  $x^*$ .

Ponieważ  $f(p) = 0$  i  $f'(c) \geq m_1$  to  $|f(x^*) - f(p)| = |f(x^*)| \geq m_1 |x^* - p|$  zatem

$$|x^* - p| \leq \frac{|f(x^*)|}{m_1}.$$

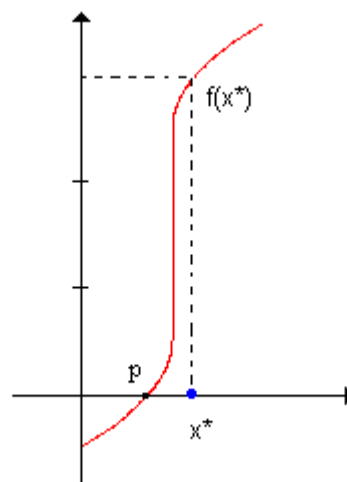
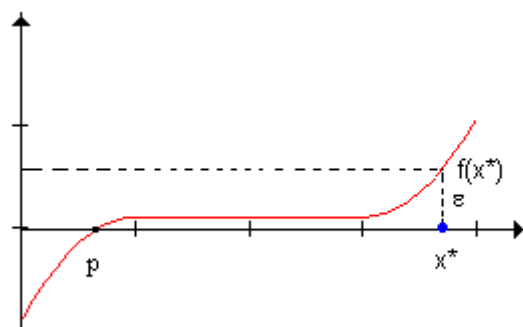
**Przykład.8.3.1.** Rozpatrzmy równanie  $x^4 - x - 1 = 0$ . Weźmy za  $x^* = 1,22$ . Oszacujemy błąd bezwzględny tego przybliżenia. Mamy  $f(x^*) = -0,0047$ . Dla  $x^{**} = 1,23$  wartość funkcji  $f(x^{**}) = 0.0588$ , dokładny pierwiastek  $p$  jest zatem w przedziale  $(1,22; 1.23)$ .

Pochodna  $f'(x) = 4x^3 - 1$  jest w tym przedziale rosnąca, a najmniejszą wartość przyjmuje dla

$$x^* = 1,22, \text{ zatem } m_1 = 4 \cdot (1,22)^3 - 1 = 6,264 \text{ więc } |x^* - p| \leq \frac{0,0047}{6,264} < 0,001.$$

*Uwaga:* Niekiedy w praktyce ocenia się dokładność przybliżenia pierwiastka  $x^*$  według tego, czy liczba  $|f(x^*)|$  jest mała, czy duża. Jeśli jest mała, to uważa się że  $x^*$  jest dobrym przybliżeniem dokładnej wartości pierwiastka  $p$  i na odwrót, jeśli  $|f(x^*)|$  jest duże to  $x^*$  zostaje uznane za złe przybliżenie. Jak widać z następujących rysunków takie podejście nie zawsze jest prawidłowe. Nie należy również zapominać, że po pomnożeniu równania  $f(x) = 0$  przez dowolną liczbę  $N$  różną od zera, otrzymujemy równanie równoważne, a liczbę  $Nf(x^*)$  można uczynić dowolnie dużą lub dowolnie małą, dzięki doborowi  $N$ .





Rys8.3.1. Sytuacja, gdy  $x^*$  nie jest bliskie  $p$ . Rys8.3.2. Sytuacja gdy  $f(x^*)$  nie jest bliskie zeru.

W dalszych rozważaniach, aby zapobiec takim opisanym wyżej sytuacjom, będziemy zakładać brak punktów przegięcia funkcji  $f$  w przedziałach izolacji, tzn, tak będziemy dobierać (zawężać) przedział izolacji, aby druga pochodna funkcji opisującej lewą stronę równania miała stały znak w rozpatrywanym przedziale.

## 8.4 RZĄD METODY

Podstawowym warunkiem, jaki powinna spełniać dana metoda jest zbieżność ciągu iteracyjnego do pierwiastka równania. Oczywiście, tym lepsza jest metoda im szybciej ciąg przybliżeń jest zbieżny do  $p$ . Szybkość zbieżności można określić za pomocą dwóch wielkości: rzędu metody i stałej asymptotycznej  $C$  błędu metody.

*Definicja.* Mówimy, że metoda jest rzędu  $r$ , jeśli istnieje granica:

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - p|}{|x_n - p|^r} = C \neq 0 \quad (8.4.1)$$

Liczbę  $C$  nazywamy stałą asymptotyczną błędu metody.

*Uwagi:* 1. Jeśli  $r = 1$  i  $C < 1$  to  $x_n$  dąży do  $p$  dla dowolnego  $x_0$  - punktu startowego.

2. Jeśli  $r > 1$  i  $x_0$  jest dostatecznie bliskie  $p$  to  $x_n$  zbiega do  $p$ .

3. Im większy jest rząd metody i im mniejsza jest stała asymptotyczna błędu, tym szybciej ciąg jest zbieżny do pierwiastka.

Zilustrujemy te uwagi na przykładzie.

**Przykład 8.4.1.** Załóżmy, że metoda jest rzędu 2, ze stałą asymptotyczną równą 5. Czyli:

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - p|}{|x_n - p|^2} = 5$$

Dla każdego  $\varepsilon > 0$  istnieje takie  $N$ , że dla  $n > N$  spełnione są nierówności:

$$5|x_n - p|^2 - \varepsilon \leq |x_{n+1} - p| \leq 5|x_n - p|^2 + \varepsilon$$

Oznacza to, że różnica pomiędzy  $(n+1)$ -szym przybliżeniem  $x_{n+1}$  i dokładnym pierwiastkiem  $p$  może przyjmować, dla wystarczająco dużych  $n$  wartości dowolnie mało różniące się od  $5|x_n - p|^2$

$$|x_{n+1} - p| \cong 5|x_n - p|^2$$

Przyjmijmy, że obliczyliśmy przybliżenie  $x_n$  pierwiastka, różniące się od niego o mniej niż 0,1 tzn:  $|x_n - p| < 0,1$ . Wówczas:

$$|x_{n+1} - p| \cong 5 \cdot (0,1)^2 = 0,05$$

$$|x_{n+2} - p| \cong 5 \cdot (0,05)^2 = 0,0025$$

$$|x_{n+3} - p| \cong 5 \cdot (0,0025)^2 = 0,00003125$$

Widzimy więc, że różnice pomiędzy kolejnymi przybliżeniami i pierwiastkiem dokładnym szybko maleją.

Założmy, że metoda jest rzędu 3, a stała asymptotyczna równa się 5. Zakładając, że  $|x_n - p| < 0,1$  otrzymujemy teraz:

$$|x_{n+1} - p| \cong 5 \cdot (0,1)^3 = 0.005$$

$$|x_{n+2} - p| \cong 5 \cdot (0,005)^3 = 0,000000625$$

Zatem metoda rzędu 3 jest wyraźnie szybsza od metody rzędu 2.

Gdyby wziąć metodę rzędu 2 ze stałą  $C=1$ , to przy takiej samej przyjętej dokładności między  $n$ -tym przybliżeniem a pierwiastkiem :  $|x_n - p| < 0,1$  dostajemy:

$$|x_{n+1} - p| \cong 0,01$$

$$|x_{n+1} - p| \cong 0,0001$$

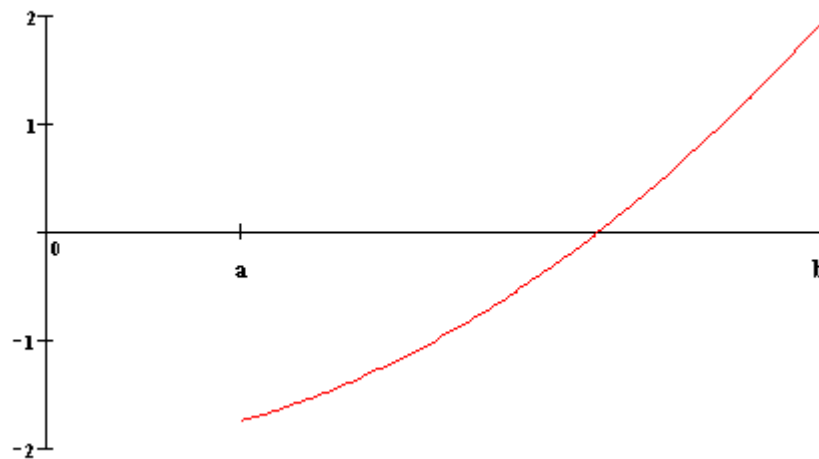
$$|x_{n+3} - p| \cong 0,00000001$$

Metoda rzędu 2 ze stałą 1 jest nieco szybsza od metody rzędu 2 ze stałą 5 , ale wolniejsza od metody rzędu 3 ze stałą większą.

## 8.5 METODA BISEKCJI

Omówimy teraz najprostsze metody znajdowania pierwiastków równania nieliniowego  $f(x) = 0$ . Podamy za każdym razem warunki wystarczające i konieczne, aby ciąg iteracyjny był zbieżny do szukanego pierwiastka. W każdej z omawianych metod rozpatrujemy przedział izolacji  $(a, b)$  w którym znajduje się aktualnie szukany pierwiastek, jeśli pierwiastków jest więcej - więcej jest przedziałów izolacji, metodę stosujemy po kolei do każdego przedziału izolacji osobno.

**Metoda bisekcji.** Założenia: Funkcja  $f$  jest funkcją ciągłą na  $[a, b]$ , oraz  $f(a)f(b) < 0$ .

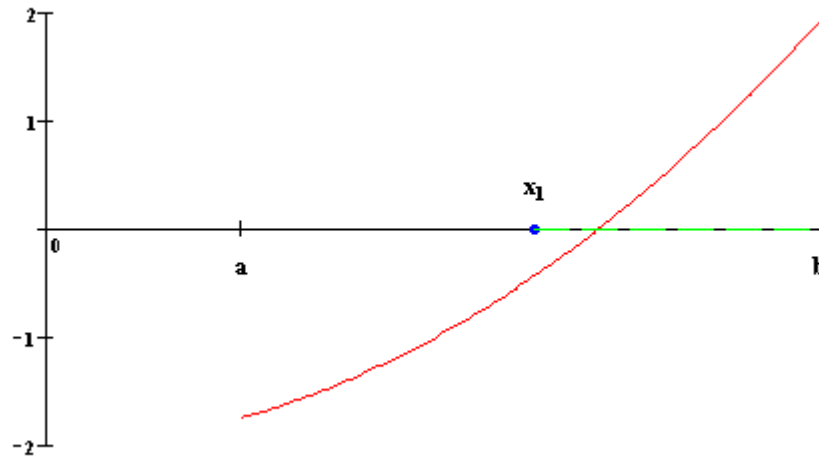


Rys8.5.1. W przedziale  $[a, b]$  jest jeden pierwiastek rzeczywisty równania  $f(x)=0$ .

*Opis metody:* Dzielimy przedział  $[a, b]$  na połowy punktem:  $x_1 = \frac{a+b}{2}$

a) b) c) d)

- a) Wybieramy przedział, w którym funkcja zmienia znak, czyli w naszym przypadku  $[x_1, b]$
- b) Dzielimy go na połowę punktem  $x_2$  i wybieramy znów przedział, w którym funkcja zmienia znak, w naszym przypadku  $[x_1, x_2]$
- c) Dzielimy go na połowę punktem  $x_3$ , wybieramy przedział  $[x_1, x_3]$ ,
- d) dzielimy znów na pół punktem  $x_4$  wybieramy przedział  $[x_4, x_3]$  itd.



Rys8.5.2. Kolejne cztery iteracje szukanego pierwiastka.

Jeśli  $f(x_1) = 0$  to  $x_1$  jest pierwiastkiem równania. Jeśli  $f(x_1)$  jest różne od zera to z otrzymanych dwóch podprzedziałów  $[a, x_1]$  i  $[x_1, b]$  wybieramy ten, w którym funkcja  $f$  zmienia znak. Z kolei ten przedział dzielimy na połowy punktem  $x_2$  i badamy wartość funkcji w  $x_2$  oraz znaki w otrzymanych podprzedziałach, wybierając do dalszych obliczeń zawsze ten, w którym funkcja zmienia znak. Otrzymujemy albo po  $n$  krokach  $f(x_n) = 0$  albo ciąg podprzedziałów takich, że  $f(x_n)f(x_{n+1}) < 0$  przy czym  $x_n, x_{n+1}$  są końcami przedziału, a jego długość

$$|x_n - x_{n+1}| < \frac{1}{2^n}(b - a).$$

Ponieważ, z konstrukcji, lewe końce przedziałów tworzą ciąg niemalejący i ograniczony z góry (przez  $p$ ), a prawe końce przedziałów tworzą ciąg nierosnący i ograniczony z dołu (przez  $p$ ), istnieje granica wspólna dla tych ciągów równa  $p$ .

Podstawową zaletą tej metody jest jej prostota i pewność, że w każdej kolejnej iteracji szukany pierwiastek leży między dwiema wartościami zmiennej  $x$ , dla których funkcja zmienia znak. Teoretycznie można uzyskać dowolną dokładność przy obliczeniach pierwiastka, stosować iterację tak długo, aż

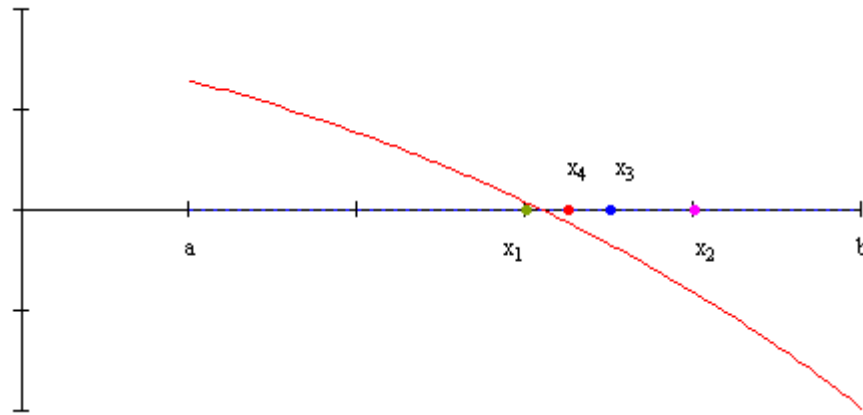
$$|x_n - x_{n+1}| < \frac{1}{2^n}(b - a) < \varepsilon \quad (8.5.1)$$

jednak przy dużej ilości kroków, błędy zaokrągleń mogą nie dopuścić do otrzymania żądanej dokładności. Metoda ta jest wolno zbieżna, bo z konstrukcji wynika, że przedziały mają za każdym razem mniejszą długość tylko o połowę. Rząd metody bisekcji jest równy 1.

#### Przykład 8.5.1: Szukamy pierwiastka wielomianu

$$f(x) = x^8 - 10x^6 + 5$$

w przedziale  $[a, b]$ , gdzie  $a = 0,8$ ,  $b = 1$  z dokładnością  $d = 0,0001$ , gdzie dokładność oznacza dla nas zakończenie iteracji jeśli  $|f(x_n)| < d$ , i jeśli funkcja nie ma w przedziale  $[a, b]$  punktów przegięcia.



Rys8.5.3. Ilustracja graficzna 4 iteracji.

Na rysunku pokazany jest wielomian w rozpatrywanym przedziale. Widać, że funkcja  $f$  nie ma punktów przegięcia w tym przedziale, i zmienia znak w  $[a, b]$ . Punkt  $x_1 = \frac{a+b}{2} = 0,9$ .

Otrzymujemy następujące wyniki, ciąg  $\{x_n\}$  dla  $j=1, \dots, 13$

|          |
|----------|
| $x_j =$  |
| 0.9      |
| 0.95     |
| 0.925    |
| 0.9125   |
| 0.90625  |
| 0.903125 |
| 0.904688 |
| 0.903906 |
| 0.903516 |
| 0.903711 |
| 0.903613 |
| 0.903662 |
| 0.903638 |

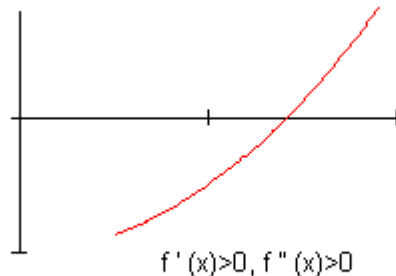
Ciąg  $\{x_n\}$  nie jest monotoniczny, oscyluje wokół pierwiastka  $p$  wielomianu  $f(x)$ .

Ostatnia iteracja daje nam pierwiastek z podaną dokładnością  $x_{13}=0,903638$  i  $f(x_{13})=0,000016$ .

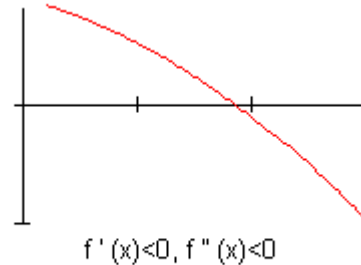
Jeśli warunkiem stopu będzie warunek  $|x_{n-2} - x_{n-1}| < \epsilon$  to dostaniemy jako pierwiastek 11-tą iterację  $x_{11} = 0,903613$  i wartość funkcji  $f : f(x_{11}) = 0,000771$ .

## 9.1 METODA SIECZNYCH

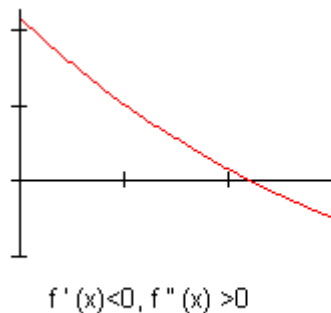
Założenia: Funkcja  $f$  jest klasy  $C^2(a,b)$ , zmienia znak w przedziale  $(a, b)$  oraz pochodne pierwsza i druga mają stały znak w rozpatrywanym przedziale. To znaczy, że w przedziale izolacji  $(a, b)$  może zachodzić któryś z czterech podanych na rysunkach przypadków:



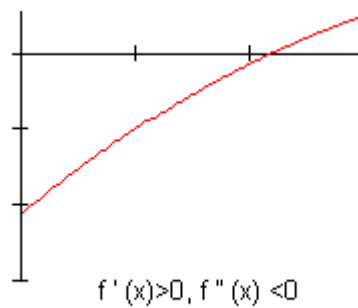
funkcja rośnie, jest wypukła



funkcja maleje, jest wklęsła



funkcja maleje, jest wypukła



funkcja rośnie, jest wklęsła

Rys9.1.1. Możliwe przypadki wykresów funkcji  $f$  w przedziale izolacji  $[a,b]$ .

Na rysunku opieramy się o przypadek, kiedy pierwsza i druga pochodna są dodatnie i startujemy z punktu  $b=x_0$  oraz z punktu  $x_1$  leżącego po lewej stronie  $b$ , ale po prawej stronie od pierwiastka.

*Opis metody:* Metoda siecznych jest metodą dwukrokową, startujemy z dwóch punktów  $x_0$  i  $x_1$  takich, że  $f(x_0)f''(x_0) > 0$ ,  $f(x_1)f''(x_1) > 0$ . Przez punkty  $(x_0, f(x_0))$  i  $(x_1, f(x_1))$  prowadzimy sieczną i przecinamy ją z osią  $Ox$ , punkt przecięcia wyznacza następną iterację  $x_2$ .

$$y - f(x_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0} (x - x_1) \quad \text{stąd} \quad x_2 = x_1 - f(x_1) \frac{x_1 - x_0}{f(x_1) - f(x_0)}$$

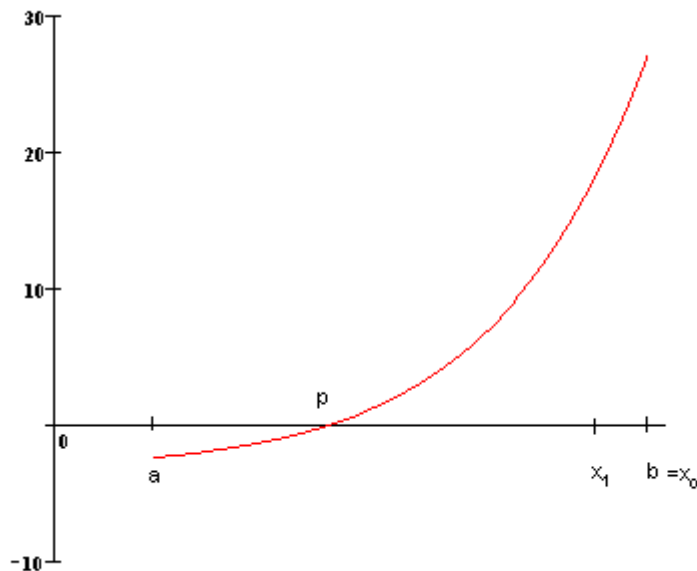
$y = 0$

Jeśli  $f(x_2) = 0$  to  $x_2$  jest pierwiastkiem, jeśli nie, przez punkty  $(x_2, f(x_2))$  i  $(x_1, f(x_1))$  prowadzimy sieczną i przecinamy ją z osią  $Ox$ , dostajemy następną iterację:

$$x_3 = x_2 - f(x_2) \frac{x_2 - x_1}{f(x_2) - f(x_1)}$$

0) a) b) c)

- a) Przez punkty  $(x_0, f(x_0))$  i  $(x_1, f(x_1))$  prowadzimy prostą i przecinamy ją z osią  $Ox$ , punkt przecięcia oznaczamy przez  $x_2$ ,
- b) Przez punkty  $(x_1, f(x_1))$  i  $(x_2, f(x_2))$  prowadzimy prostą i przecinamy ją z osią  $Ox$ , punkt przecięcia oznaczamy przez  $x_3$ ,
- c) Przez punkty  $(x_2, f(x_2))$  i  $(x_3, f(x_3))$  prowadzimy prostą i przecinamy ją z osią  $Ox$ , punkt przecięcia oznaczamy przez  $x_4$  itd.,



Rys9.1.2. Kolejne trzy iteracje w metodzie siecznych.

Postępując kolejno w wyżej opisany sposób otrzymamy wzór ogólny na ciąg iteracyjny:

$$x_{n+1} = x_n - f(x_n) \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \quad (9.1.1)$$

Jeśli nie będziemy przestrzegać spełnienia warunków na punkty startu, ciąg może być zbieżny do pierwiastka, ale nie zawsze (przykład 9.1.2).

Dla tej metody jest prawdziwe twierdzenie:

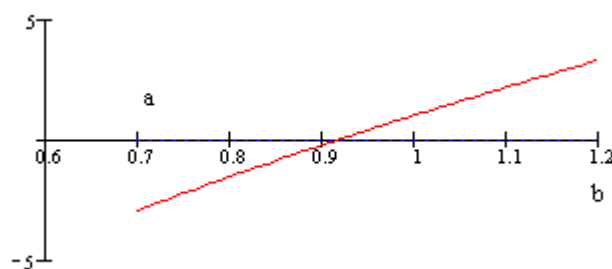
**Twierdzenie.** Jeśli w otoczeniu  $|x - p| < \delta$  pierwiastka  $p$  równania  $f(x) = 0$  funkcja  $f$  ma ciągłą drugą pochodną, a pierwsza i druga pochodna jest różna od zera w tym otoczeniu oraz przybliżenia  $x_0$  i  $x_1$  (są różne) są dostatecznie bliskie pierwiastka  $p$ , to metoda siecznych jest

zbieżna, jej rząd jest równy  $\frac{\sqrt{5}+1}{2} = 1,618\dots$ , a stała asymptotyczna błędów jest równa

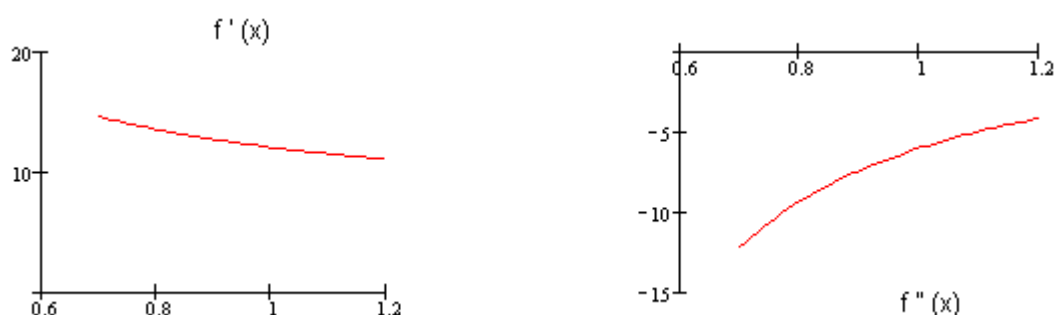
$$C = \left( \frac{|f''(p)|}{2|f'(p)|} \right)^{\frac{\sqrt{5}-1}{2}}.$$

**Przykład 9.1.1.** Rozpatrujemy równanie :  $6x + 6\ln x - 5 = 0$  w przedziale  $[0,7, 1.2]$  z dokładnością  $d = 10^{-8}$ .



Rys9.1.3. Wykres funkcji  $f$  w podanym przedziale.

Na rysunku widać, że funkcja zmienia znak w danym przedziale. Na następnych rysunkach są zilustrowane pierwsza i druga pochodna funkcji  $f$  w tym samym przedziale.

Rys9.1.4. Wykresy pochodnych funkcji  $f$ .

Widać, że pierwsza pochodna jest dodatnia w rozpatrywanym przedziale, a druga pochodna jest ujemna w  $[a, b]$ . Spełnione są założenia dla metody siecznych, pochodne są ciągłe i nie zmieniają znaku w  $[a, b]$ .

Możemy zastosować metodę siecznych, wybierając za punkty startu takie punkty, w których funkcja ma taki sam znak jak druga pochodna. Ponieważ druga pochodna jest ujemna wybieramy punkty po lewej stronie pierwiastka, w którym funkcja też jest ujemna.

$$x_0 = a, x_1 = a + 0,01$$

Za pomocą wzoru iteracyjnego na  $x_{n+1}$  dostajemy wektor iteracji, wzór jest przeliczany tak długo dopóki nie będzie osiągnięta dokładność, tzn. aż  $|f(x_n)| < d$ .

Wektor iteracji:

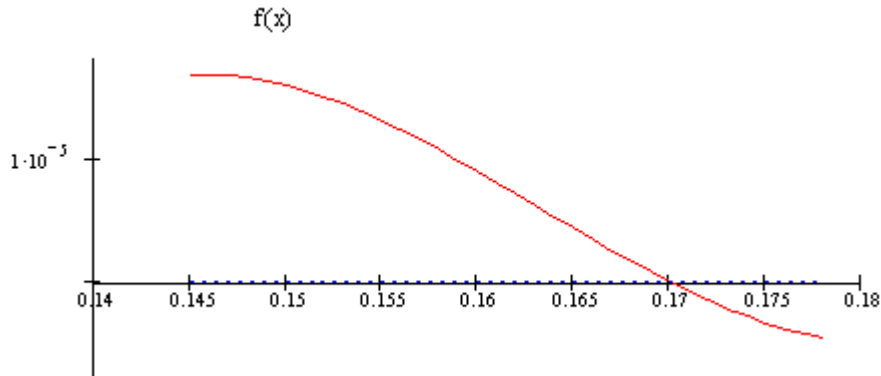
$x_j =$

|              |
|--------------|
| 0.7          |
| 0.71         |
| 0.9026114008 |
| 0.9173854594 |
| 0.9184219035 |
| 0.91842661   |
| 0.9184266114 |

Rozwiązaniem jest:  $x_6 = 0,9184266114$  dla którego  $f(x_6) = -2,309 \cdot 10^{-14}$ .

**Przykład 9.1.2.** Przykład ilustruje sytuację, w której nie są spełnione założenia przy jakich możemy stosować tę metodę. Dane jest równanie :  $x^3 - 0,49x^2 + 0,0791x - 0,004199 = 0$

Szukamy pierwiastka tego równania w przedziale  $[0,145, 0,178]$ . Pierwiastek istnieje, bo jak widać na rysunku, funkcja zmienia znak, ten pierwiastek jest blisko punktu 0,17.

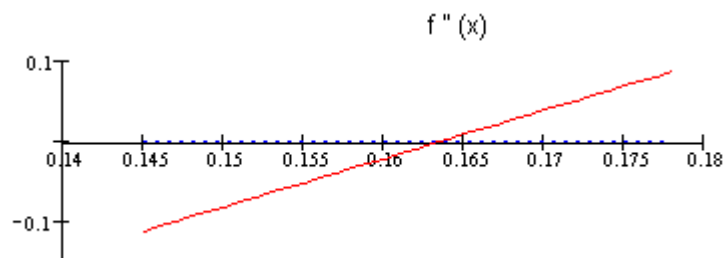


Rys9.1.5. Wykres funkcji  $f$  w podanym przedziale.

Znajdziemy ten pierwiastek metodą siecznych. Startujemy z punktów:  $x_0 = a$  i  $x_1 = a + 0,01$  i

stosujemy wzór: 
$$x_{n+1} = x_n - f(x_n) \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})}$$

Dostajemy dla 10 iteracji  $x_{10} = 0,19$  i  $f(x_{10}) = 0$ . Jednak **to nie jest pierwiastek z tego przedziału**, nasz pierwiastek był blisko punktu 0,17. Dlaczego tak się stało? Pochodna druga zmienia znak w tym przedziale, w dodatku wystartowaliśmy ze złych punktów. Ponieważ wykres drugiej pochodnej jest następujący:

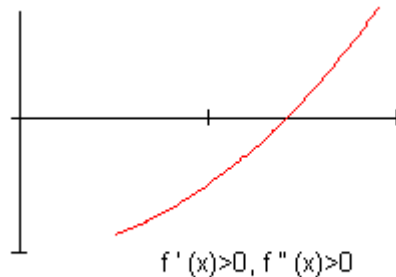


Rys9.1.6. Wykres drugiej pochodnej, która zmienia znak w rozpatrywanym przedziale.

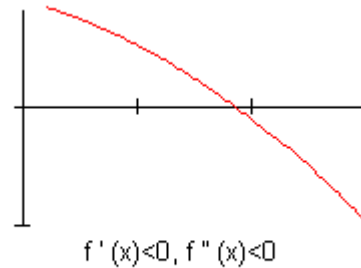
Przedział nie może być w tym wypadku taki duży, powinniśmy zmienić go na  $[0,165, 0,178]$  i sprawdzić pozostałe założenia.

## 9.2 METODA STYCZNYCH - NEWTONA

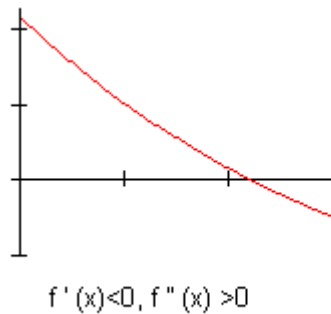
Założenia: Podobnie jak w poprzedniej metodzie założymy, że funkcja  $f$  jest klasy  $C^2(a,b)$ , zmienia znak w przedziale  $(a, b)$  oraz pochodne pierwsza i druga mają stały znak w rozpatrywanym przedziale. To znaczy, że w przedziale izolacji  $(a, b)$  może zachodzić któryś z czterech podanych na rysunkach przypadków:



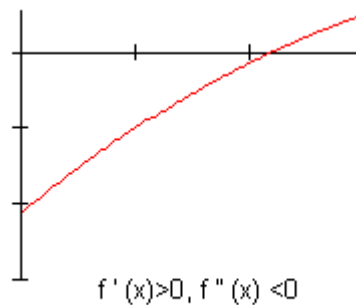
funkcja rośnie, jest wypukła



funkcja maleje, jest wklęsła



funkcja maleje, jest wypukła



funkcja rośnie, jest wklęsła

Rys9.2.1. Możliwe przypadki wykresów funkcji  $f$  w przedziale izolacji  $[a,b]$ .

*Opis metody:* Metodę opiszemy korzystając z przypadku pierwszego, kiedy pierwsza pochodna jest dodatnia (funkcja rośnie) i druga pochodna jest dodatnia (funkcja jest wypukła).

Jako punkt startu obieramy taki punkt  $x_0$ , w którym funkcja ma taki sam znak jak druga pochodna:  $f(x_0)f''(x_0) > 0$ , w naszym przypadku punkt  $b$  - ponieważ w tym punkcie funkcja jest dodatnia tak jak druga pochodna. Z punktu  $(x_0, f(x_0))$  wystawiamy styczną do krzywej  $y = f(x)$ . Równanie stycznej  $y - f(x_0) = f'(x_0)(x - x_0)$ . Przecinamy styczną z osią  $Ox$  i otrzymany punkt przecięcia jest pierwszym przybliżeniem pierwiastka.

$$y - f(x_0) = f'(x_0)(x - x_0), \quad y = 0 \quad \text{Wstawiając za } y \text{ zero a za } x \text{ wartość } x_1 \text{ otrzymujemy:}$$

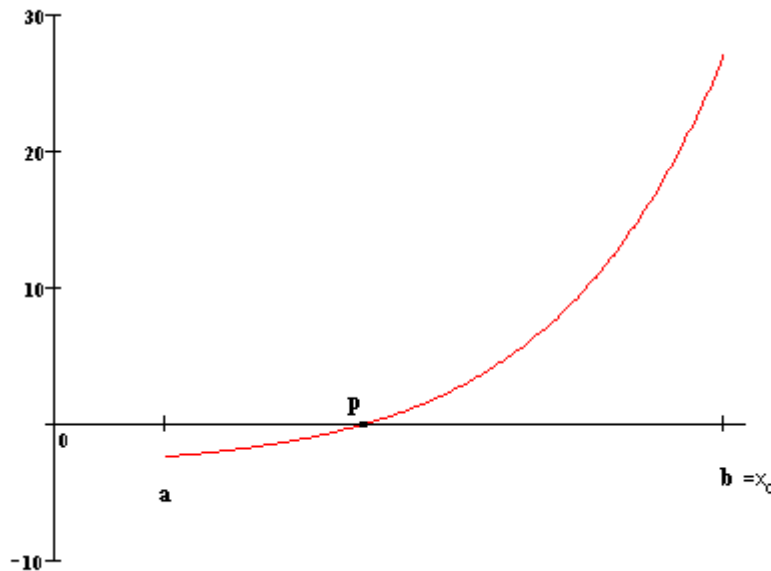
$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Jeśli  $f(x_1) = 0$  to  $x_1$  jest pierwiastkiem, jeśli nie, postępujemy analogicznie dalej, z punktu  $(x_1, f(x_1))$  wystawiamy styczną do krzywej i przecinamy ją z osią  $Ox$ :

$$y - f(x_1) = f'(x_1)(x - x_1), \quad y = 0, \quad \text{Otrzymujemy: } x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

0) a) b) c)

- a) Przez punkt  $(x_0, f(x_0))$  prowadzimy styczną do  $f(x)$  i przecinamy ją z osią  $Ox$ , punkt przecięcia oznaczamy przez  $x_1$ ,  
 b) Przez punkt  $(x_1, f(x_1))$  prowadzimy styczną i przecinamy ją z osią  $Ox$ , punkt przecięcia oznaczamy przez  $x_2$ ,  
 c) Przez punkt  $(x_2, f(x_2))$  prowadzimy styczną i przecinamy ją z osią  $Ox$ , punkt przecięcia oznaczamy przez  $x_3$  itd.,



Rys9.2.2. Kolejne trzy iteracje w metodzie stycznych.

Powtarzając w ten sposób budowanie kolejnej iteracji otrzymujemy ciąg iteracyjny  $x_n$  określony wzorem :

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (9.2.1)$$

Ciąg ten przy podanych założeniach jest zbieżny do szukanego pierwiastka  $p$ . Może się zdarzyć , że startując z innego punktu, nie spełniającego podany warunek  $f(x_0)f''(x_0) > 0$  , ciąg iteracyjny też będzie zbieżny do szukanego pierwiastka, ale bez tego warunku nie mamy gwarancji , że ciąg  $x_n$  zbiega do  $p$  (przykład 9.2.2).

Dla tej metody jest prawdziwe twierdzenie:

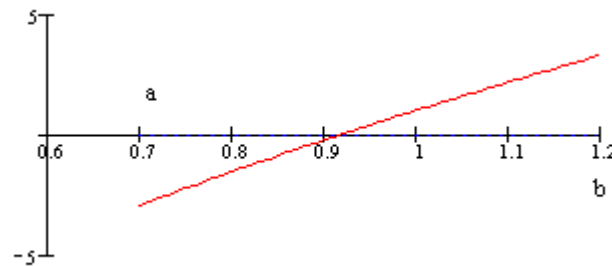
**Twierdzenie:** Jeżeli w otoczeniu  $|x - p| < \delta$  pierwiastka  $p$  równania  $f(x) = 0$  funkcja  $f$  ma ciągłą drugą pochodną oraz pierwsza i druga pochodna są różne od zera w tym otoczeniu oraz  $x_0$  leży wystarczająco blisko pierwiastka  $p$ , to metoda Newtona jest rzędu 2 ze stałą

asymptotyczną błędu  $C = \frac{|f''(p)|}{|2f'(p)|}$ .

Metoda Newtona jest szybkozbieżną metodą jednokrokową wymagającą na każdym kroku obliczania jednej wartości funkcji i jednej wartości pierwszej pochodnej.

**Przykład 9.2.1.** Rozpatrujemy równanie :  $6x + 6\ln x - 5 = 0$  w przedziale  $[0.7, 1.2]$ .

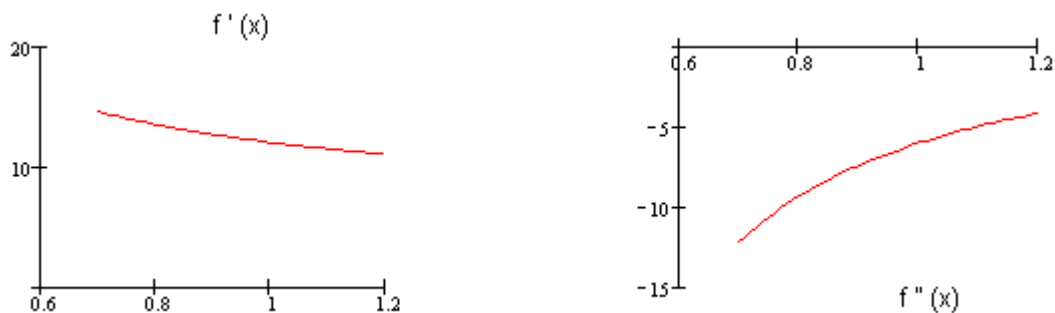
Na rysunku widać, że funkcja zmienia znak w danym przedziale.



Rys9.2.3. Wykres funkcji  $f$  w podanym przedziale.

Na następnych rysunkach są zilustrowane pierwsza i druga pochodna funkcji  $f$  w tym samym przedziale.

Widać, że pierwsza pochodna jest dodatnia w rozpatrywanym przedziale, a druga pochodna jest ujemna w  $[a, b]$ . Spełnione są założenia dla metody Newtona, pochodne są ciągłe i nie zmieniają znaku w  $[a, b]$ .



Rys9.2.4. Wykres pochodnych w rozpatrywanym przedziale.

Możemy zastosować metodę Newtona przyjmując za punkt startu ten koniec przedziału  $[a, b]$ , dla którego jest spełniony warunek  $f'(x_0)f''(x_0) > 0$ . W tym przypadku jest to punkt  $a$ , zatem  $x_0 = a$ . Dla dokładności  $\epsilon = 10^{-8}$  otrzymujemy 4 iteracje i  $x_4 = 0.9184266114$  jest przybliżonym pierwiastkiem równania oraz  $f(x_4) = 0$  (przyjmujemy za zero wszystkie liczby mniejsze niż  $10^{-15}$ ).

Wektor iteracji ma postać:  $j = 0, 1, \dots, 4$

$x_j =$

|              |
|--------------|
| 0.7          |
| 0.9017681142 |
| 0.9183466866 |
| 0.9184266096 |
| 0.9184266114 |

Ten sam przykład, dla tej samej dokładności obliczyliśmy w poprzednim temacie metodą

siecznych. Aby otrzymać żadaną dokładność trzeba było dla tamtej metody wziąć o dwie iteracje więcej. Metoda stycznych jest szybciej zbieżną metodą.

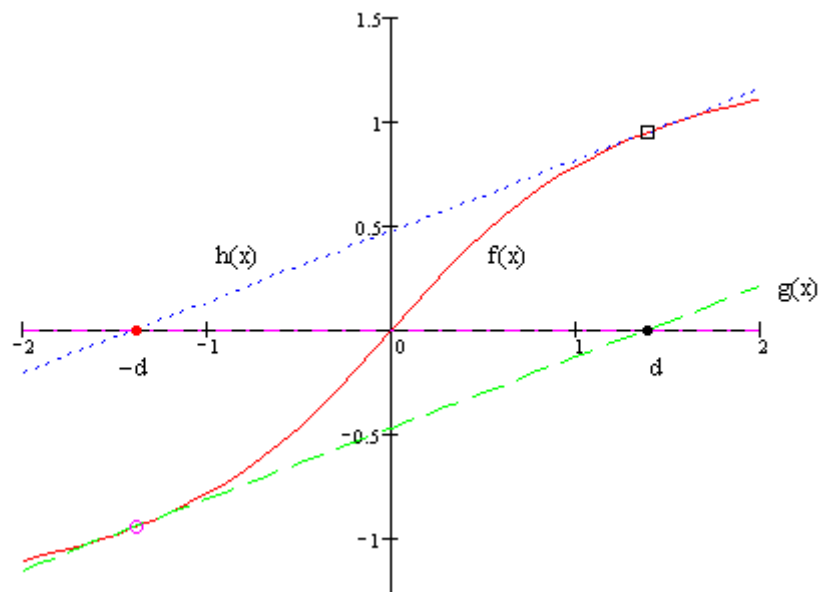
**Przykład 9.2.2.** Rozpatrujemy równanie  $\arctg x = 0$  w przedziale  $[-2, 2]$ . Jako punkt startu obieramy dokładny pierwiastek równania:

$$\arctg x - \frac{2x}{1+x^2} = 0$$

Oznaczmy ten pierwiastek przez  $d$  i w przybliżeniu równa się on 1,39125043. Będziemy stosować metodę Newtona dla równania:  $\arctg x = 0$  z punktem startowym  $x_0 = d$ .

Wystawiamy z punktu  $(d, f(d))$  styczną do krzywej  $\arctg x$ :  $g(x) = f'(d)(x - d) + f(d)$

i przecinamy ją z osią  $Ox$  wyznaczając punkt  $x_1$ . Okazuje się że punkt przecięcia będzie  $x_1 = -d$ . Jeśli z punktu  $(-d, f(-d))$  wystawimy do krzywej  $\arctg x$  styczną:  $h(x) = f'(-d)(x + d) + f(-d)$  i przetniemy ją z osią  $Ox$  dostaniemy znów punkt  $x_2 = d$ . W ten sposób metoda Newtona "zapętiła" się i ze wzoru Newtona dostajemy na zmianę punkty  $d$  i  $-d$  jako kolejne iteracje, a widać na rysunku, że pierwiastkiem równania jest  $p = 0$ .



Rys9.2.5. Wykres funkcji  $f$  i stycznych wychodzących z punktów  $(d, 0)$  i  $(-d, 0)$ .

To zapętienie wynika z tego, że druga pochodna zmienia znak w przedziale  $[-2, 2]$ , ma w zerze punkt przegięcia. Nie są zatem spełnione założenia podane do metody Newtona.

### 9.3 PIERWIASTKI WIELOKROTNE

Metody iteracyjne wymagają na ogół, aby szukany pierwiastek był pierwiastkiem jednokrotnym. Tak jest przy metodzie Newtona i metodzie siecznych. Metoda bisekcji dopuszcza pierwiastki nieparzystokrotne, przy parzystokrotnych funkcja nie zmienia znaku w przedziale izolacji. Na ogół nie znamy krotności szukanych pierwiastków.

Wprowadzamy funkcję pomocniczą  $u(x) = \frac{f(x)}{f'(x)}$  i rozwiązujemy równanie  $u(x) = 0$  zamiast równania  $f(x) = 0$ . Równanie  $u(x) = 0$  ma takie same pierwiastki jak równanie  $f(x) = 0$ , ale wszystkie są jednokrotne.

$$\text{Ponieważ : } u'(x) = \frac{f'(x) \cdot f'(x) - f(x) \cdot f''(x)}{(f'(x))^2} = 1 - \frac{f(x)}{f'(x)} \frac{f''(x)}{f'(x)} = 1 - u(x) \frac{f''(x)}{f'(x)}$$

wzory na metodę Newtona i metodę siecznych przybierają postać:

dla metody Newtona:

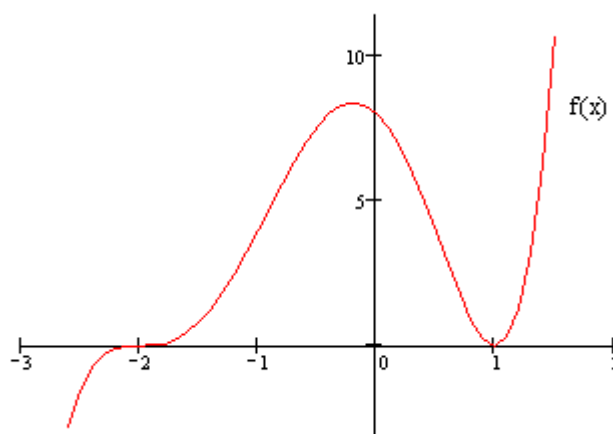
$$x_{n+1} = x_n - \frac{u(x_n)}{u'(x_n)} \quad (9.3.1)$$

$$\text{gdzie } u'(x_n) = 1 - u(x_n) \frac{f''(x_n)}{f'(x_n)}$$

dla metody siecznych:

$$x_{n+1} = x_n - u(x_n) \frac{x_n - x_{n-1}}{u(x_n) - u(x_{n-1})} \quad (9.3.2)$$

**Przykład 9.3.1:** Funkcja nieliniowa  $f(x)$  będzie wielomianem stopnia piątego mającym jeden pierwiastek dwukrotny i jeden trzykrotny.  $f(x) = (x-1)^2(x+2)^3$



Rys9.3.1. Wykres funkcji  $f$ .

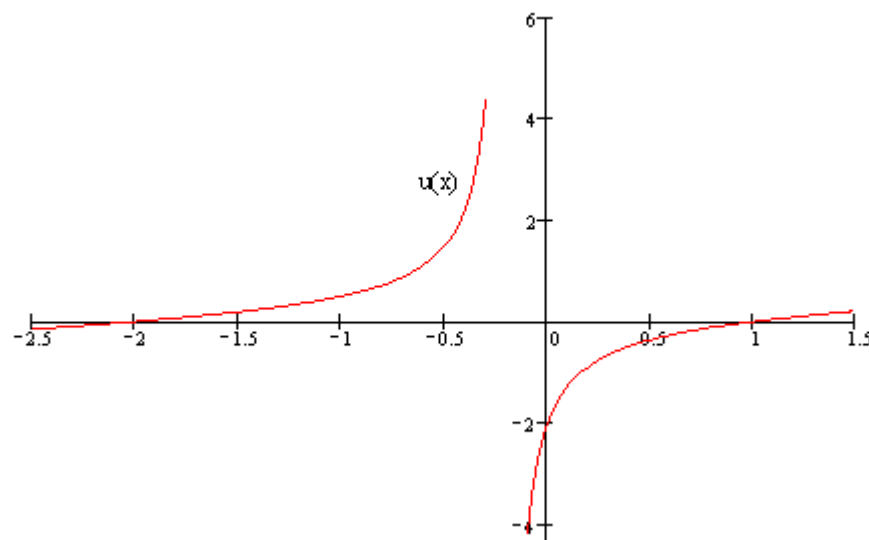
Obliczymy jej pochodną i następnie wprowadzimy funkcję  $u(x) = \frac{f'(x)}{f''(x)}$ .

$$f(x) = (x-1)^2(x+2)^3,$$

$$\begin{aligned} f'(x) &= 2(x-1)(x+2)^3 + 3(x-1)^2(x+2)^2 = (x-1)(x+2)^2(2(x+2) + 3(x-1)) = \\ &= (x-1)(x+2)^2(5x+1) \end{aligned}$$

$$u(x) = \frac{(x-1)^2(x+2)^3}{(x-1)(x+2)^2(5x+1)} = \frac{(x-1)(x+2)}{5x+1}$$

Równanie  $u(x) = 0$  ma dwa pierwiastki, takie jak funkcja  $f$  ale są już jednokrotne. Funkcja  $u$  nie jest ciągła na całej osi  $\mathbb{R}$ , ale istnieją przedziały izolacji pierwiastków, w których jest ciągła i ma ciągłe pochodne.



Rys9.3.2. Wykres funkcji  $u$ .



## 9.4 UKŁADY NIELINIOWE

Dany jest układ  $n$  równań nieliniowych z  $n$  niewiadomymi:

$$\begin{cases} f_1(x_1, x_2, \dots, x_n) = 0 \\ f_2(x_1, x_2, \dots, x_n) = 0 \\ \dots \dots \dots \\ f_n(x_1, x_2, \dots, x_n) = 0 \end{cases} \quad (9.4.1)$$

który będziemy zapisywać wektorowo :  $F(x) = 0$  , gdzie  $x \in \mathbb{R}^n$  ;  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ .

Warunki istnienia rozwiązań układu są znacznie trudniejsze do sprawdzenia, a nawet nie można sformułować jednolitego kryterium istnienia rozwiązania bez założenia szczególnych własności odwzorowania  $F$ , takich jak różniczkowalność itd. Będziemy zakładać istnienie rozwiązania układu  $F(x) = 0$  i ograniczymy się do jednej metody : poszukiwania rozwiązań metodą Newtona.

Rozpatrzmy metodę iteracyjną jednokrokową daną ogólnym wzorem :  $x^{(k)} = G(x^{(k-1)})$  i wektor początkowy  $x^{(0)}$  będziemy dobierać dostatecznie blisko rozwiązania.

*Definicja:* Niech  $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Punkt  $w$  nazywamy punktem przyciągania metody iteracyjnej , jeżeli istnieje takie otoczenie  $U$  tego punktu, że biorąc dowolny wektor początkowy  $x^{(0)}$  z tego otoczenia uzyskamy ciąg punktów  $x^{(1)}, x^{(2)}, \dots, x^{(k)}, \dots$  zbieżny do  $w$ . Największe z tych otoczeń nazywamy obszarem przyciągania punktu  $w$ .

Oznaczmy przez :

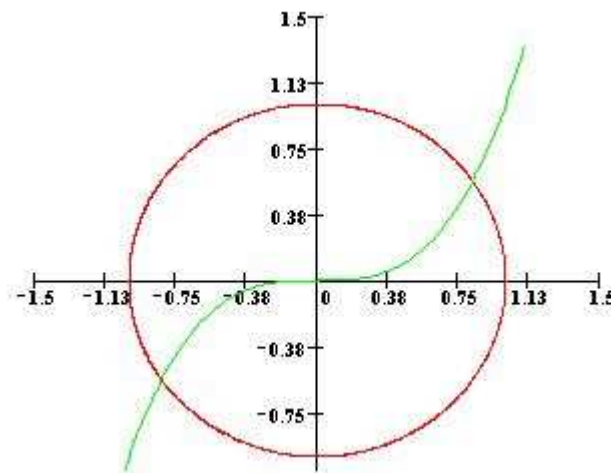
$$F(x) = \begin{bmatrix} f_1(x_1, \dots, x_n) \\ \dots \dots \dots \\ f_n(x_1, \dots, x_n) \end{bmatrix} \quad J(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix} \quad (9.4.2)$$

Jeśli funkcje  $f_i(x)$  są różniczkowalne w sposób ciągły w pewnym otoczeniu punktu  $p$  , w którym  $F(p) = 0$  , i macierz  $J(x)$  jest w tym otoczeniu nieosobliwa, to jeśli wektor startu dobierzemy odpowiednio blisko punktu  $p$  to punkt  $p$  jest punktem przyciągania metody iteracyjnej danej wzorem ( metoda Newtona):

$$x^{(n+1)} = x^{(n)} - (J(x^{(n)}))^{-1} F(x^{(n)}) \quad (9.4.3)$$

**Przykład 9.4.1:** Rozpatrujemy układ równań:

$$\begin{aligned} x^2 + y^2 - 1 &= 0 \\ x^3 - y &= 0 \end{aligned}$$



Rys9.4.1. Graficzna interpretacja układu.

Na rysunku czerwona linia opisuje pierwsze równanie, zielona drugie. Widać, że krzywe przecinają się w dwóch punktach, układ ma dwa rozwiązania. Lewą stronę pierwszego równania oznaczmy przez  $f_1(x, y)$ , lewą stronę drugiego równania oznaczmy przez  $f_2(x, y)$ . Oznaczmy przez :

$$F(x, y) = \begin{bmatrix} f_1(x, y) \\ f_2(x, y) \end{bmatrix} \text{ oraz } J(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix}$$

Znajdziemy rozwiązanie w pierwszej ćwiartce. Jako wektor startu bierzemy (odczytujemy z rysunku), wektor  $z$  ma pierwszą współrzędną  $x$  a drugą  $y$ .

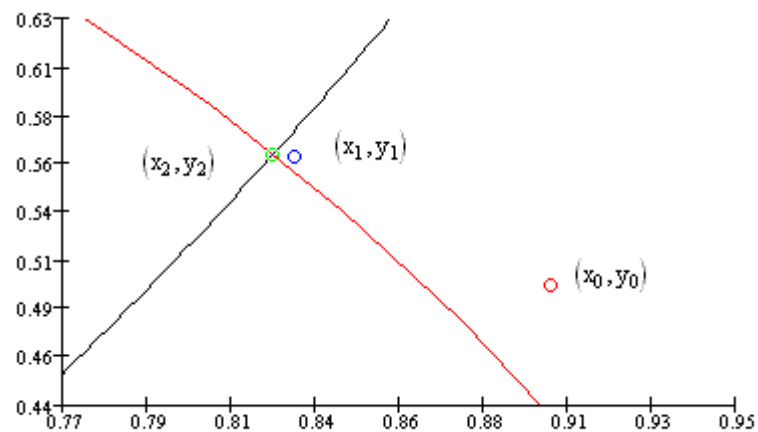
$$z^{(0)} = \begin{bmatrix} 0.9 \\ 0.5 \end{bmatrix}$$

Korzystając ze wzoru Newtona dla układów równań :  $z^{(n+1)} = z^{(n)} - (J(z^{(n)}))^{-1} F(z^{(n)})$  dostajemy dla dokładności  $d = 10^{-8}$ , tzn.  $|F(z^{(n)})| < d$  całą macierz iteracji, w kolumnach której są kolejne wektory iteracyjne:

$$z = \begin{pmatrix} 0.9 & 0.831678 & 0.826062 & 0.826031 \\ 0.5 & 0.562979 & 0.563608 & 0.563624 \end{pmatrix}$$

Rozwiązaniem jest trzecia iteracja skąd  $x = 0,826031$ , a  $y = 0,563624$ .

Na rysunku, w dużym powiększeniu, czerwone kółeczko to punkt startu, pierwsze przybliżenie to niebieskie kółeczko, drugie przybliżenie to zielone kółeczko, a rozwiązanie przybliżone, czyli trzecia iteracja pokrywa się na rysunku z drugą.



Rys9.4.2. Kolejne trzy iteracje rozwiązania układu.

Wartość funkcji wektorowej opisującej równania jest dla tego rozwiązania następująca:

$$F(z^{(3)}) = \begin{bmatrix} 1,19 \cdot 10^{-9} \\ 2,294 \cdot 10^{-9} \end{bmatrix}$$

Jeśli będziemy brać jako wektor startu wektor o współrzędnych o przeciwnych znakach

$$z^{(0)} = \begin{bmatrix} -0.9 \\ -0.5 \end{bmatrix}$$

dostaniemy symetryczne rozwiązanie  $x = -0,826031$  ,a  $y = -0,563624$ .

## 10.1 CAŁKOWANIE NUMERYCZNE

Bardzo dużo zagadnień geometrycznych, fizycznych, mechanicznych sprowadza się do obliczania całek oznaczonych funkcji jednej zmiennej. Można za pomocą tych całek liczyć np.: długości łuków, pola obszarów, pola powierzchni obrotowych, objętości brył obrotowych, masy ciał, momenty statyczne i bezwładności itd. Dokładne obliczenie tych całek wymaga znajomości funkcji pierwotnych dla funkcji podcałkowych, nie każda jednak funkcja posiada funkcję pierwotną. Zachodzi konieczność znalezienia całki oznaczonej metodą przybliżoną. Również w przypadku, gdy funkcja jest dana za pomocą pomiarów, tylko w pewnej ilości punktach, można całkować taką funkcję podanymi poniżej metodami.

Zajmiemy się w tym rozdziale obliczaniem całek oznaczonych za pomocą wielomianów interpolacyjnych. W całości:

$$\int_a^b f(x) dx$$

będziemy zastępować funkcję  $f(x)$  jej wielomianem interpolacyjnym  $n$ -tego stopnia, którego węzły  $x_0, x_1, \dots, x_n$  będą leżały w przedziale całkowania  $[a, b]$  i będą teraz węzłami całkowania. Ograniczymy się w tym opracowaniu do całkowania funkcji bez osobliwości w przedziale  $[a, b]$  tzn.: funkcji przyjmującej skończone wartości w rozpatrywanym przedziale, a przedział  $[a, b]$  jest skończony. Wynika z tych założeń, że nie będziemy się zajmować całkowaniem całek niewłaściwych.

Do szacowania błędu całkowania wykorzystamy podane już wcześniej wzory na błąd interpolacji. Błąd ten zależy od pochodnych funkcji podcałkowych i od węzłów. Będziemy rozpatrywać węzły równoodległe - dla prostoty obliczeń, a również węzły optymalne, tzn.: takie, które minimalizują tę część błędu całkowania, która zależy od węzłów.

Obliczając zatem całkę

$$\int_a^b f(x) dx$$

wstawiamy za  $f(x)$  wielomian interpolacyjny Lagrange'a  $n$ -tego stopnia  $L_n(x)$  dany wzorem:

$$L_n(x) = \sum_{k=0}^n \Phi_k(x) f(x_k) \quad (10.1.1)$$

gdzie :

$$\Phi_k(x) = \frac{(x - x_0) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_n)}{(x_k - x_0) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n)} \quad (10.1.2)$$

a  $x_0, x_1, \dots, x_n$  są węzłami w przedziale  $[a, b]$ .

Otrzymamy:

$$\begin{aligned}
 I(f) &= \int_a^b f(x) dx \cong \int_a^b L_n(x) dx = \int_a^b \sum_{k=0}^n \Phi_k(x) f(x_k) dx = \sum_{k=0}^n f(x_k) \left( \int_a^b \Phi_k(x) dx \right) = \\
 &= \sum_{k=0}^n A_k f(x_k) = S(f)
 \end{aligned}
 \tag{10.1.3}$$

gdzie  $A_k = \int_a^b \Phi_k(x) dx$ .

$A_k$  są współczynnikami zależnymi od węzłów, nie zależą od funkcji podcałkowej  $f(x)$ , łatwo je wyliczyć, bo są to całki z wielomianów w przedziale  $[a, b]$ .

Wzór na wartość  $S(f)$  będziemy nazywać kwadraturą.

Błąd przybliżenia:  $E(f) = I(f) - S(f)$ , to różnica między dokładną wartością całki  $I(f)$ , a jej wartością przybliżoną  $S(f)$ , jest to przecalkowany w przedziale  $[a, b]$  błąd interpolacji.

## 10.2 METODY PROSTE TRAPEZÓW I PARABOL

Rozpatrzmy skończony przedział  $[a, b]$  oraz równoodległe węzły  $x_0, x_1, \dots, x_n$  w tym przedziale:

$h = \frac{b-a}{n}$ ,  $x_k = a + k \cdot h$ ,  $k = 0, 1, \dots, n$ . Wtedy :

$$I(f) = \int_a^b f(x) dx \cong \sum_{k=0}^n A_k f(x_k) = S(f) \quad \text{gdzie} \quad A_k = \int_a^b \Phi_k(x) dx.$$

### Wzór prosty trapezów.

Ustalmy  $n=1$  i wtedy  $h = b - a$  i mamy dwa węzły  $x_0 = a$ ,  $x_1 = b$ . Wielomian interpolacyjny Lagrange'a jest stopnia 1 i wstawiając za funkcję  $f(x)$  wielomian  $L_1(x)$  mamy:

$$\begin{aligned} I(f) &= \int_a^b f(x) dx \cong \int_a^b L_1(x) dx = \int_a^b \left( \frac{x-x_1}{x_0-x_1} f(x_0) + \frac{x-x_0}{x_1-x_0} f(x_1) \right) dx = \\ &= f(x_0) \int_a^b \frac{x-x_1}{x_0-x_1} dx + f(x_1) \int_a^b \frac{x-x_0}{x_1-x_0} dx = \frac{h}{2} f(x_0) + \frac{h}{2} f(x_1) = S(f) \end{aligned} \quad (10.2.1)$$

Całki w ostatniej linijce wzoru łatwo obliczyć, bo są to całki oznaczone z wielomianów pierwszego stopnia.

W tym wypadku

$$A_0 = \int_a^b \frac{x-x_1}{x_0-x_1} dx = \frac{1}{x_0-x_1} \int_a^b (x-x_1) dx = \frac{1}{-h} \int_a^b (x-b) dx = \frac{1}{-h} \left( \frac{1}{2} x^2 - bx \right) \Big|_a^b = \frac{h}{2},$$

$$A_1 = \int_a^b \frac{x-x_0}{x_1-x_0} dx = \frac{1}{h} \int_a^b (x-a) dx = \frac{1}{h} \left( \frac{1}{2} x^2 - ax \right) \Big|_a^b = \frac{h}{2}.$$

Możemy wzór na przybliżoną wartość  $I(f)$  zapisać w postaci :

$$I(f) = \int_a^b f(x) dx \cong A_0 f(x_0) + A_1 f(x_1) = \frac{h}{2} (f(x_0) + f(x_1)) = \frac{h}{2} (f(a) + f(b)) \quad (10.2.2)$$

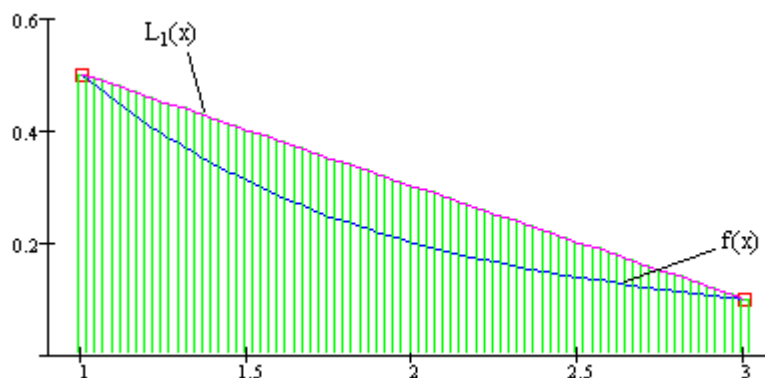
Korzystając ze wzoru na błąd interpolacji dostajemy:

$$E(f) = I(f) - S(f) = \int_a^b (I(f) - S(f)) dx = \frac{1}{2} \int_a^b (x-a)(x-b) f''(\xi) dx = -\frac{1}{12} h^3 f''(\xi^*) \quad (10.2.3)$$

gdzie  $\xi, \xi^* \in (a, b)$ .

Wzór na przybliżoną wartość całki ma prostą interpretację geometryczną. Pokażemy ją na przykładzie.

**Przykład 10.2.1:** Obliczymy przybliżoną wartość całki z funkcji  $f(x) = \frac{1}{1+x^2}$  w przedziale  $[1,3]$  za pomocą wielomianu interpolacyjnego stopnia 1. Graficznie, ponieważ dana funkcja jest dodatnia, całka z tej funkcji równa się polu pod krzywą opisaną daną funkcją. Na rysunku krzywa jest zaznaczona na niebiesko. Zamiast pola pod krzywą liczymy ze wzoru przybliżonego pole pod prostą łączącą punkty  $(a, f(a))$  i  $(b, f(b))$  (czyli pod wielomianem interpolacyjnym pierwszego stopnia) zaznaczoną na czerwono. Pole, które otrzymamy jest zakreskowane na zielono. To zielone pole jest polem trapezu, który "leży" na swojej wysokości  $h$ . I wzór dlatego nosi nazwę wzoru trapezów, a jak widać, we wzorze jest suma podstaw trapezu dzielona przez 2 i pomnożona przez wysokość.



Rys10.2.1. Interpretacja geometryczna wzoru prostego trapezów.

Po wykonaniu obliczeń dostajemy:

$$S(f) = \frac{h}{2}(f(x_0) + f(x_1)) = \frac{h}{2}(f(a) + f(b)) = \frac{2}{2}\left(\frac{1}{2} + \frac{1}{10}\right) = 0.6$$

wartość przybliżona całki  $S(f)=0.6$ , a ponieważ dokładną wartość możemy w tym wypadku podać, bo funkcja pierwotna dla funkcji podcałkowej to  $\arctg(x)$ , zatem  $I(f)=\arctg(b)-\arctg(a)=0,463648$  (z dokładnością do 6 cyfr po przecinku), to błąd bezwzględny równa się  $bl=|I(f)-S(f)|=0,136352$ , i stanowi aż 23%. Widać na rysunku, że wartości przybliżona i dokładna znacznie się różnią (pole pod funkcją i pole pod prostą).

### Wzór prosty parabol (Simpsona).

Ponieważ w podanym przykładzie wartość całki obarczona jest dużym błędem, wstawimy zamiast funkcji wielomian interpolacyjny stopnia 2.

Ustalmy  $n = 2$  i wtedy  $h = (b - a)/2$  i mamy trzy węzły  $x_0 = a$ ,  $x_1 = a + h = \frac{a+b}{2}$ ,  $x_2 = b$ .

Wielomian interpolacyjny Lagrange'a jest stopnia 2 i wstawiając za funkcję  $f(x)$  wielomian  $L_2(x)$  mamy:

$$\begin{aligned}
I(f) &= \int_a^b f(x) dx \cong \int_a^b L_2(x) dx = \\
&= \int_a^b \left( \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} f(x_0) + \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} f(x_1) + \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} f(x_2) \right) dx = \\
&= f(x_0) \int_a^b \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} dx + f(x_1) \int_a^b \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} dx + f(x_2) \int_a^b \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} dx = \\
&= \frac{h}{3} f(x_0) + \frac{4h}{3} f(x_1) + \frac{h}{3} f(x_2) = S(f)
\end{aligned} \tag{10.2.}$$

Całki w ostatniej linijce wzoru łatwo obliczyć, bo są to całki oznaczone z wielomianów drugiego stopnia.

W tym wypadku

$$\begin{aligned}
A_0 &= \int_a^b \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} dx = \frac{h}{3}, & A_1 &= \int_a^b \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} dx = \frac{4h}{3}, \\
A_2 &= \int_a^b \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} dx = \frac{h}{3}
\end{aligned}$$

Możemy wzór na przybliżoną wartość  $I(f)$  zapisać w postaci:

$$I(f) = \int_a^b f(x) dx \cong A_0 f(x_0) + A_1 f(x_1) + A_2 f(x_2) = \frac{h}{3} (f(x_0) + 4f(x_1) + f(x_2)) \tag{10.2.5}$$

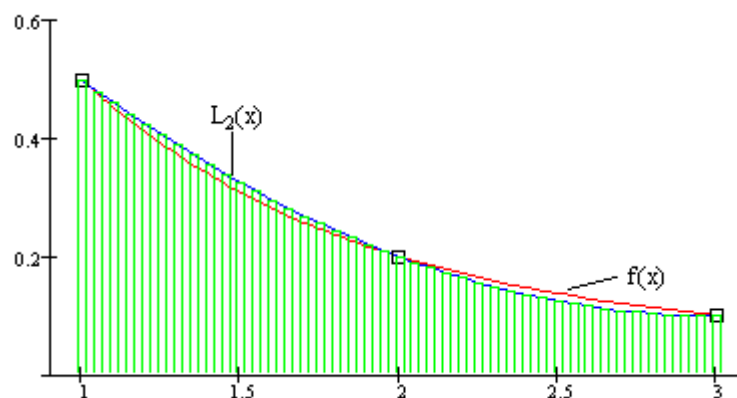
Korzystając ze wzoru na błąd interpolacji dostajemy:

$$E(f) = I(f) - S(f) = \int_a^b (I(f) - S(f)) df = -\frac{1}{90} h^5 f^{(4)}(\xi^*) \tag{10.2.6}$$

gdzie  $\xi, \xi^* \in (a, b)$ .

Powróćmy do przykładu 10.2.1. Teraz prowadzimy parabolę przez trzy punkty  $(a, f(a))$ ,  $((a+h), f(a+h))$  i  $(b, f(b))$ . Pole zakreskowane na zielono równa się polu pod parabolą (wielomianem interpolacyjnym stopnia 2) narysowaną na niebiesko, funkcja jest narysowana na czerwono.





Rys10.2.2. Interpretacja geometryczna wzoru prostego parabol.

Po obliczeniu przybliżonej wartości całki według wzoru parabol

$$\begin{aligned}
 S(f) &= A_0 f(x_0) + A_1 f(x_1) + A_2 f(x_2) = \frac{h}{3} (f(x_0) + 4f(x_1) + f(x_2)) = \\
 &= \frac{1}{3} \left( \frac{1}{2} + \frac{4 \cdot 1}{5} + \frac{1}{10} \right) = \frac{7}{15}
 \end{aligned}$$

otrzymujemy:  $S(f)=0,466667$ , błąd bezwzględny  $|I(f) - S(f)|=0,003019$ , co stanowi tylko 0,3%.

### 10.3 METODY ZŁOŻONE TRAPEZÓW I PARABOL

W poprzednim temacie rozpatrywaliśmy przybliżone całkowanie funkcji za pomocą wielomianów interpolacyjnych 1 i 2 stopnia z równoodległymi węzłami. Można by wyprowadzić również podobne wzory dla wielomianów wyższych stopni, ale okazało się, że lepiej podzielić przedział całkowania na  $m$  części i w każdym otrzymanym podprzedziale zastosować wzór prosty trapezów lub parabol, korzystając z tego faktu, że całka po przedziale  $[a, b]$  jest sumą całek po otrzymanych podprzedziałach. Zajmiemy się tym poniżej.

#### Wzór złożony trapezów.

Dzielimy przedział  $[a, b]$  na  $m$  części :  $h = \frac{b-a}{m}$ ,  $x_k = a + k \cdot h$ ,  $k = 0, 1, \dots, m$ . Otrzymamy  $m$

podprzedziałów o długości  $h$ , w każdym z podprzedziałów  $[x_i, x_{i+1}]$ ,  $i = 0, 1, \dots, m-1$  stosujemy wzór prosty trapezów:

$$\int_{x_i}^{x_{i+1}} f(x) dx \cong \frac{h}{2} (f(x_i) + f(x_{i+1})) \quad (10.3.1)$$

Sumując całki po wszystkich podprzedziałach dostajemy przybliżoną wartość całki w przedziale  $[a, b]$ :

$$\begin{aligned} I(f) &= \int_a^b f(x) dx = \int_a^{x_1} f(x) dx + \int_{x_1}^{x_2} f(x) dx + \dots + \int_{x_i}^{x_{i+1}} f(x) dx + \dots + \int_{x_{m-1}}^b f(x) dx \cong \\ &\cong \frac{h}{2} (f(x_0) + f(x_1)) + \frac{h}{2} (f(x_1) + f(x_2)) + \dots + \frac{h}{2} (f(x_i) + f(x_{i+1})) + \dots + \frac{h}{2} (f(x_{m-1}) + f(x_m)) = \\ &= \frac{h}{2} (f(x_0) + 2(f(x_1) + f(x_2) + \dots + f(x_{m-1})) + f(x_m)) = S(f) \end{aligned} \quad (10.3)$$

Jeśli przesumujemy błędy po wszystkich podprzedziałach otrzymamy;

$$E(f) = -\frac{(b-a)^3}{12m^2} f''(\xi) \quad , \quad \xi \in (a, b) \quad (10.3.3)$$

Zauważmy, że we wzorze na  $S(f)$ , w ostatniej linijce, wartości funkcji podcałkowej w skrajnych węzłach są w nawiasie wzięte z mnożnikiem 1, a w pozostałych węzłach z mnożnikiem 2. Prosta interpretację geometryczną tego faktu zilustrujemy na przykładzie, który był przeliczany w poprzednim temacie.

**Przykład 10.3.1:** Obliczymy przybliżoną wartość całki z funkcji  $f(x) = \frac{1}{1+x^2}$  w przedziale  $[1, 3]$ , dzieląc przedział na 4 części. Mamy :  $m = 4$ ,  $h = 0,5$ , zatem

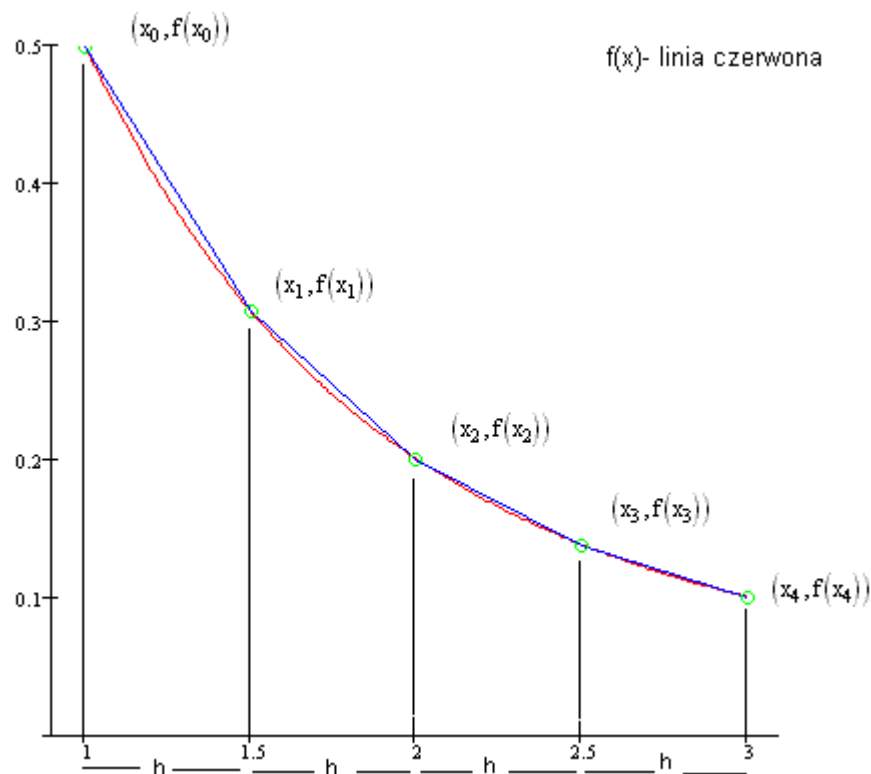
$$S(f) = \frac{0,5}{2} (f(1) + f(3) + 2(f(1,5) + f(2) + f(2,5))) =$$

$$= 0,25 \left( \frac{1}{2} + \frac{1}{10} + 2 \cdot \left( \frac{1}{1+(1,5)^2} + \frac{1}{1+2^2} + \frac{1}{1+(2,5)^2} \right) \right) = 0,473$$

Obliczając wartość całki za pomocą funkcji pierwotnej  $F(x)=\arctg(x)$ , otrzymamy  $I(f)=0,464$ , zatem błąd bezwzględny wyniesie 0,009, tzn. 0,9%.

Graficznie: trzeba obliczyć pola czterech trapezów o tej samej wysokości  $h$ , pierwszy trapez ma podstawy równe  $f(x_0)$  i  $f(x_1)$ , drugi ma podstawy  $f(x_1)$  i  $f(x_2)$ , podstawami trzeciego trapezu są  $f(x_2)$  i  $f(x_3)$  i podstawami czwartego trapezu są  $f(x_3)$  i  $f(x_4)$ .

Jak widać trzy podstawy są wspólne w tych czterech trapezach i dlatego są we wzorze wzięte podwójnie. Skrajne podstawy są uwzględnione tylko raz. Po przesumowaniu pól tych trapezów dostajemy przybliżoną wartość całki  $S(f)$ .



Rys10.3.1. Interpretacja geometryczna wzoru złożonego trapezów.

### Wzór złożony parabol (Simpsona).

Dzielimy przedział  $[a, b]$  na  $m$  części ( ale bierzemy  $m$  **parzyste**):

$h = \frac{b-a}{m}$ ,  $x_k = a + k \cdot h$ ,  $k = 0, 1, \dots, m$ . Otrzymamy  $m$  podprzedziałów o długości  $h$ , inaczej  $m/2$  podprzedziałów o długości  $2h$ , w każdym z podprzedziałów o długości  $2h$ :

$[x_i, x_{i+2}]$ ,  $i = 0, 1, \dots, m-2$  stosujemy wzór prosty parabol:

$$\int_{x_i}^{x_{i+2}} f(x) dx \cong \frac{h}{3} (f(x_i) + 4f(x_{i+1}) + f(x_{i+2})) \quad (10.3.4)$$

Sumując otrzymane całki po  $m/2$  podprzedziałach otrzymujemy:

$$\begin{aligned} I(f) &= \int_a^b f(x) dx = \int_a^{x_2} f(x) dx + \dots + \int_{x_i}^{x_{i+2}} f(x) dx + \dots + \int_{x_{m-2}}^b f(x) dx \cong \\ &\cong \frac{h}{3} (f(x_0) + 4f(x_1) + f(x_2)) + \frac{h}{3} (f(x_2) + 4f(x_3) + f(x_4)) + \dots \\ &+ \frac{h}{3} (f(x_i) + 4f(x_{i+1}) + f(x_{i+2})) + \dots + \frac{h}{3} (f(x_{m-2}) + 4f(x_{m-1}) + f(x_m)) = \\ &= \frac{h}{3} (f(x_0) + 2(f(x_2) + f(x_4)) \dots + f(x_{m-2})) + 4(f(x_1) + f(x_3) + \dots + f(x_{m-1})) + f(x_m)) = \\ &= S(f) \end{aligned} \quad (10.3.5)$$

W nawiasie w ostatniej linijce wzoru wartości funkcji w skrajnych węzłach są wzięte z mnożnikiem 1, wartości funkcji w węzłach pozostałych numerach parzystych są z mnożnikiem 2, a w węzłach o numerach nieparzystych z mnożnikiem 4.

Sumując błędy, podane w poprzednim temacie dla wzoru parabol, po  $m/2$  podprzedziałach dostajemy:

$$E(f) = -\frac{(b-a)^5}{180m^4} f^{(4)}(\xi^*), \quad \xi^* \in (a, b) \quad (10.3.6)$$

Wróćmy do przykładu 10.3.1 i obliczmy całkę metodą parabol biorąc też  $m = 4$ ,  $h = 0,5$ .

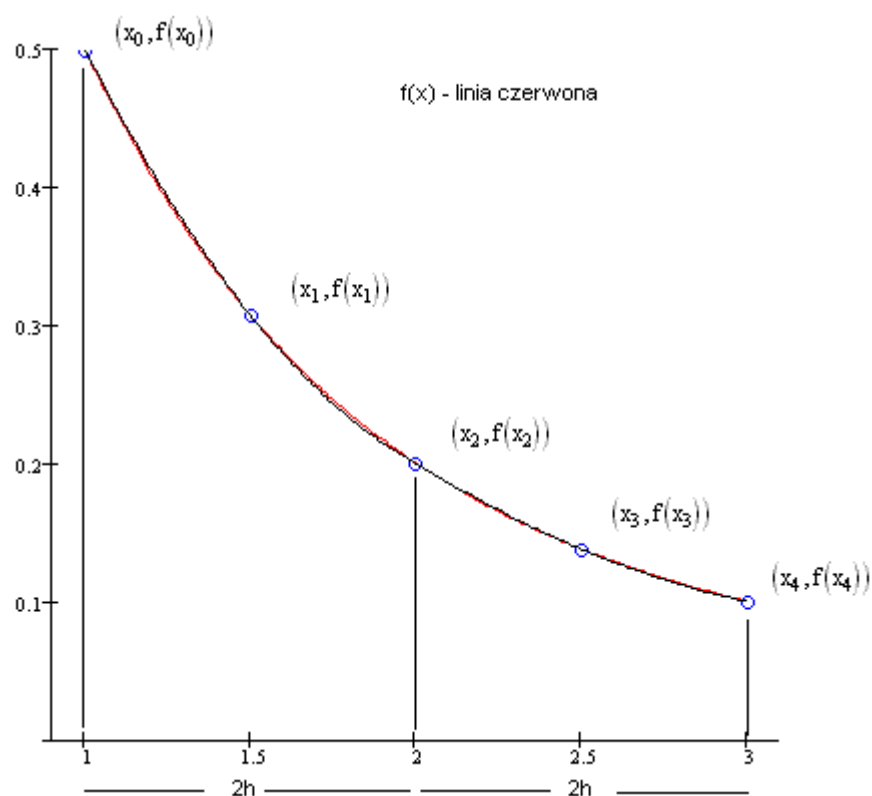
Bierzemy podprzedziały  $[1, 2]$  i  $[2, 3]$ , w każdym z nich stosujemy wzór prosty parabol.

Dostajemy:

$$\begin{aligned} S(f) &= \frac{h}{3} (f(x_0) + 4f(x_1) + f(x_2)) + \frac{h}{3} (f(x_2) + 4f(x_3) + f(x_4)) = \\ &= \frac{h}{3} (f(1) + f(3) + 2 \cdot f(2) + 4 \cdot (f(1,5) + f(2,5))) = 0,4637 \end{aligned}$$

Korzystając z funkcji pierwotnej, tak jak w przykładzie powyżej dostajemy błąd bezwzględny równy 0,0001, czyli 0,01%.

Na rysunku są zaznaczone dwie parabole (na czarno), jedna w przedziale  $[1, 2]$ , druga w przedziale  $[2, 3]$ , ale błędy są tak małe, że prawie się pokrywają się z funkcją (na czerwono). Przybliżona wartość całki to suma pól pod tymi parabolami.



Rys10.3.2. Interpretacja geometryczna wzoru złożonego parabol.

## 11.1 WĘZŁY LEGENDRE`A

Do tej pory obliczaliśmy przybliżoną wartość całki oznaczonej z funkcji  $f(x)$  zastępując ją wielomianami interpolacyjnymi z równoodległymi węzłami. Ale błąd całkowania zależy od położenia węzłów, tak jak w interpolacji, więc choć węzły równoodległe są wygodne do liczenia, nie zawsze są najlepsze. Okazuje się, że optymalnymi węzłami są pierwiastki pewnych wielomianów, które noszą nazwę wielomianów Legendre`a. Nie będziemy wprowadzać tutaj teorii wielomianów ortogonalnych, tylko podamy wartości tych pierwiastków i wartości związanych z nimi współczynników. Ponieważ pierwiastki wielomianów Legendre`a są w przedziale  $(-1,1)$ , aby z nich skorzystać, jako z węzłów, zamienimy naszą całkę po przedziale  $[a, b]$  na przedział  $[-1, 1]$  za pomocą podstawienia:  $x = \frac{b-a}{2}t + \frac{b+a}{2}$ , wtedy  $dx = \frac{b-a}{2}dt$  i otrzymujemy:

$$I(f) = \int_a^b f(x)dx = \int_{-1}^1 f\left(\frac{b-a}{2}t + \frac{b+a}{2}\right) \cdot \frac{b-a}{2} dt = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{b-a}{2}t + \frac{b+a}{2}\right) dt \quad (11.1.1)$$

Jeśli  $t_k$  dla  $k=0,1 \dots n$  będą węzłami Legendre`a, to wstawiając za funkcję  $f$  wielomian interpolacyjny stopnia  $n$  z tymi węzłami otrzymamy wzór na przybliżoną kwadraturę :

$$I(f) = \int_a^b f(x)dx = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{b-a}{2}t + \frac{b+a}{2}\right) dt \cong \frac{b-a}{2} \sum_{k=0}^n A_k f\left(\frac{b-a}{2}t_k + \frac{b+a}{2}\right) = S(f) \quad (11.1.2)$$

gdzie  $A_k = \int_{-1}^1 \Phi_k(x)dx$  nie zależą od funkcji, tylko od węzłów, a w wielomianach  $\Phi_k(x)$  też występują węzły  $t_k$ .

W każdej książce z metod numerycznych jest podana tablica węzłów i współczynników Legendre`a. Podajemy poniżej tabelkę dla wielomianów interpolacyjnych stopnia  $n=1, 2, 3$  i  $4$ .

| $n$ | $k$  | Węzły $t_k$ | Współczynniki $A_k$ |
|-----|------|-------------|---------------------|
| 1   | 0; 1 | -/+0,577350 | 1                   |
| 2   | 0; 2 | -/+0,774597 | 5/9                 |
|     | 1    | 0           | 8/9                 |
| 3   | 0; 3 | -/+0,861136 | 0,347855            |
|     | 1; 2 | -/+0,339981 | 0,652145            |
| 4   | 0; 4 | -/+0,906180 | 0,236927            |
|     | 1; 3 | -/+0,538469 | 0,478629            |
|     | 2    | 0           | 0,568889            |

**Przykład 11.1.1 :** Obliczymy przybliżoną wartość całki z funkcji  $f(x) = \frac{1}{1+x^2}$  w przedziale  $[1,3]$  za pomocą wielomianu interpolacyjnego stopnia 1 wykorzystując dwa węzły Legendre'a :

$t_0 = -0,577350, t_1 = 0,577350, A_0 = 1, A_1 = 1$ . Dla  $n=1$  wzór na  $S(f)$  ma postać :

$$I(f) = \int_a^b f(x) dx \cong \frac{b-a}{2} \left( A_0 f\left(\frac{b-a}{2}t_0 + \frac{b+a}{2}\right) + A_1 f\left(\frac{b-a}{2}t_1 + \frac{b+a}{2}\right) \right) = S(f)$$

Przeliczając węzły z przedziału  $(-1, 1)$  do przedziału  $[1, 3]$  otrzymujemy

$$x_0 = \frac{b-a}{2}t_0 + \frac{b+a}{2} = 1,42265, \quad x_1 = \frac{b-a}{2}t_1 + \frac{b+a}{2} = 2,57735$$

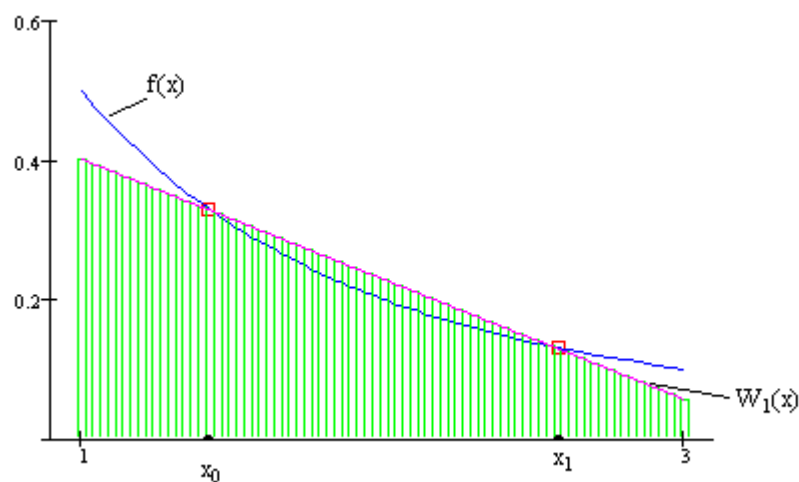
Zatem:

$$S(f) = \frac{3-1}{2} (A_0 f(x_0) + A_1 f(x_1)) = 1 \cdot (1 \cdot f(1,42265) + 1 \cdot f(2,57735)) = 0,462$$

Jak porównamy ten wynik z całką obliczoną za pomocą funkcji pierwotnej  $\arctg(x)$  dostajemy błąd równy 0,002 czyli 0,2%.

Ten przykład był przeliczany na różne sposoby, proszę porównać wynik otrzymany wzorem prostym trapezów.

Graficznie, zamiast pola pod krzywą liczymy ze wzoru przybliżonego pole pod prostą łączącą punkty  $(x_0, f(x_0)), (x_1, f(x_1))$  (czyli pod wielomianem interpolacyjnym pierwszego stopnia). Pole, które otrzymamy jest zakreskowane na zielono.



Rys11.1.1. Interpretacja graficzna wzoru z dwoma węzłami Legendre'a.

Jeśli zastosujemy wielomian stopnia 2, będą nam potrzebne trzy węzły i odpowiadające im trzy współczynniki Legendre'a :

$$t_0 = -0,774597, t_1 = 0, t_2 = 0,774597, A_0 = \frac{5}{9}, A_1 = \frac{8}{9}, A_2 = \frac{5}{9}$$

Wtedy:

$$I(f) = \int_a^b f(x) dx \cong \frac{b-a}{2} (A_0 f(x_0) + A_1 f(x_1) + A_2 f(x_2)) = S(f)$$

gdzie

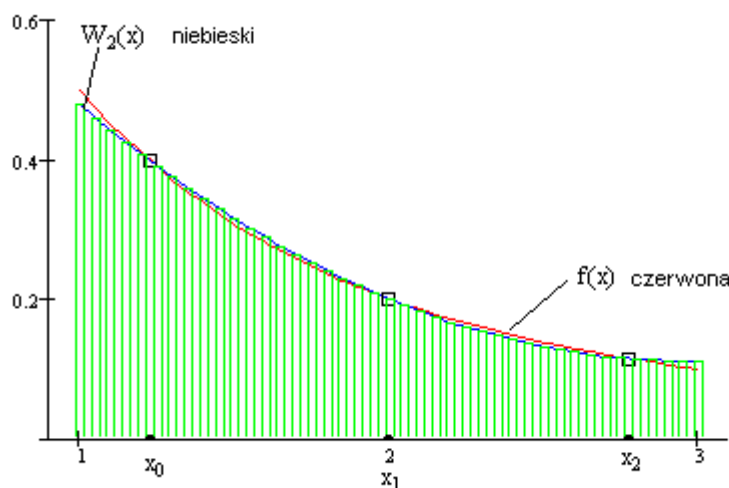
$$x_0 = \frac{b-a}{2} t_0 + \frac{b+a}{2} = 1,225403, \quad x_1 = \frac{b-a}{2} t_1 + \frac{b+a}{2} = 2, \quad x_2 = \frac{b-a}{2} t_2 + \frac{b+a}{2} = 2,774597$$

Wstawiając te wartości otrzymujemy:

$$S(f) = \frac{b-a}{2} (A_0 f(x_0) + A_1 f(x_1) + A_2 f(x_2)) = 0,4637$$

a porównując przybliżoną wartość całki z otrzymaną za pomocą funkcji pierwotnej otrzymujemy błąd 0,00008, czyli 0,008%.

Graficznie, zamiast pola pod krzywą liczymy ze wzoru przybliżonego pole pod parabolą przechodzącą przez punkty  $(x_0, f(x_0))$ ,  $(x_1, f(x_1))$ ,  $(x_2, f(x_2))$  (czyli pod wielomianem interpolacyjnym drugiego stopnia). Pole, które otrzymamy jest zakreskowane na zielono.



Rys11.1.2. Interpretacja graficzna wzoru z trzema węzłami Legendre'a.

Porównując wyniki, które otrzymaliśmy w powyższych przykładach z wynikami otrzymanymi za pomocą węzłów równoodległych widzimy, że w tym wypadku liczona tą metodą wartość całki jest bardzo bliska jej wartości "dokładnej". W tym wypadku węzły Legendre'a dają dużo lepszy wynik. Można również, analogicznie jak w przypadku wzorów trapezów i parabol wyprowadzić wzory złożone oparte na 2 lub 3 węzłach Legendre'a i wtedy dokładność jeszcze się poprawi



## 11.2 UWAGI O DOKŁADNOŚCI

**Przykład 11.2.1.** Zacniemy od tego samego przykładu, który już był wcześniej wielokrotnie przeliczany: Obliczymy przybliżoną wartość całki z funkcji  $f(x) = \frac{1}{1+x^2}$  w przedziale  $[1,3]$ , metodą złożoną trapezów, dzieląc przedział na  $m$  części.

Przypominamy, że funkcją pierwotną funkcji  $f$  jest  $\arctg(x)$ , zatem wartość "dokładna" tej całki wynosi (z dokładnością do 10 miejsc po przecinku) 0,4636476090. Podane niżej obliczenia wskazują na to, że jeśli zwiększamy ilość podprzedziałów na które dzielimy przedział  $[a, b]$ , dokładność obliczeń rośnie. Metoda trapezów jest wolno zbieżna i jak widać wymaga dużej ilości podprzedziałów, jeśli błąd ma być mały.

Podać  $m =$   oblicz: całka = , błąd =

Przeliczmy, teraz całkę z tej samej funkcji metodą złożoną parabol. Jest to metoda zdecydowanie szybsza. Proszę obserwować wyniki.

Podać **parzyste**  $m =$   oblicz: całka = , błąd =

Jeśli pojawi się w błędzie wynik typu 3.245e-7, to zapis ten oznacza liczbę, która dopiero na 7 miejscu po przecinku ma 3 potem 2,4 i 5.

W podanych przykładach obliczaliśmy całkę z funkcji, która miała funkcję pierwotną, można było określić błąd wyników, przez porównanie wartości przybliżonej z wartością dokładną. Na ogół całkujemy w sposób numeryczny funkcję taką, dla której nie istnieje funkcja pierwotna, albo trudno ją znaleźć, jednak chcemy aby nasze obliczenia nie przekraczały z góry zadanego błędu. Można to zrobić, jeśli oszacujemy w przedziale całkowania maksymalną wartość modułu drugiej pochodnej dla wzoru trapezów i czwartej pochodnej dla wzoru parabol. Wtedy możemy dobrać tak liczbę podprzedziałów, na które dzielimy przedział  $[a, b]$ , aby uzyskać żadaną dokładność.

Ponieważ wzór na błąd całkowania dla metody trapezów był następujący:

$$E(f) = -\frac{(b-a)^3}{12m^2} f''(\xi) \quad , \quad \xi \in (a,b) \quad (11.2.1)$$

to oznaczając przez  $M2 = \sup_{x \in [a,b]} |f''(x)|$ , błąd bezwzględny całkowania nie przekroczy wartości

$$\varepsilon = \left| \frac{(b-a)^3}{12m^2} \cdot M2 \right| , \text{ skąd przy danym maksymalnym dopuszczalnym błędzie możemy obliczyć}$$

$m$ , czyli ilość podprzedziałów w metodzie trapezów.

Mamy:  $m = \sqrt{\frac{(b-a)^3}{12\varepsilon} M2}$ , ale ponieważ  $m$  musi być liczbą naturalną, bierzemy za  $m$ :

$$m = \left\lceil \sqrt[4]{\frac{(b-a)^3}{12\varepsilon}} M_2 \right\rceil + 1 \quad (11.2.2)$$

gdzie nawias  $\lceil \cdot \rceil$  oznacza część całkowitą liczby  $u$ .

Dla metody parabol wzór na błąd był następujący:

$$E(f) = -\frac{(b-a)^5}{180m^4} f^{(4)}(\xi^*), \quad \xi^* \in (a, b) \quad (11.2.3)$$

Oznaczamy przez  $M_4 = \sup_{x \in [a, b]} |f^{(4)}(x)|$ , wtedy błąd bezwzględny nie przekroczy wartości

$\varepsilon = \left| \frac{(b-a)^5}{180m^4} M_4 \right|$ , zatem wzór na  $m$  jest następujący:  $m = \sqrt[4]{\frac{(b-a)^5}{180\varepsilon}} M_4$ , ale  $m$  jest naturalne i musi być **parzyste** zatem:

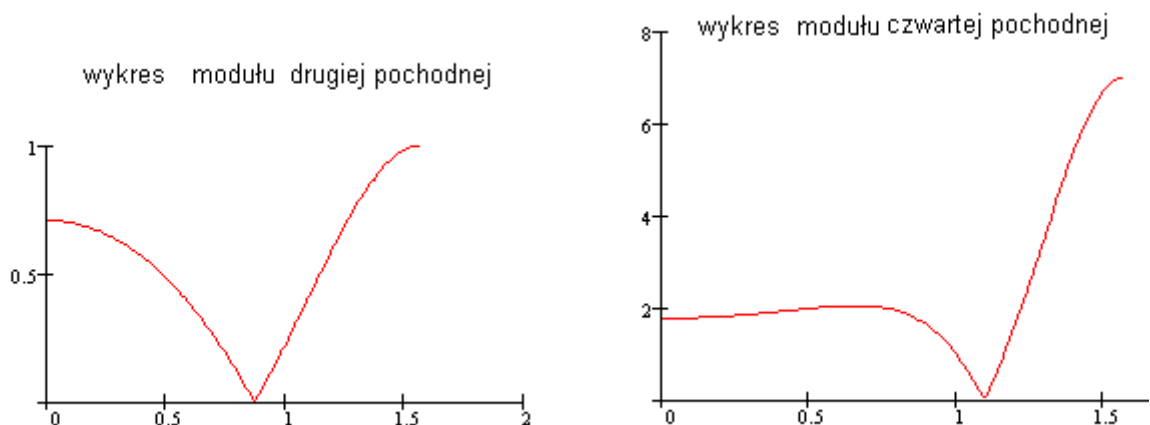
$$m = 2 \cdot \left\lceil \frac{1}{2} \cdot \sqrt[4]{\frac{(b-a)^5}{180\varepsilon}} M_4 \right\rceil + 2 \quad (11.2.4)$$

To sztuczne podzielenie pierwiastka przez 2, a później pomnożenie znów przez 2, zapewnia parzystość otrzymanego  $m$ .

**Przykład 11.2.2.** Obliczyć z dokładnością  $10^{-6}$  długość łuku krzywej  $y = \sin(x)$  w przedziale  $[0, \frac{\pi}{2}]$

Skorzystamy ze wzoru na długość łuku:  $l = \int_a^b \sqrt{1 + (y'(x))^2} dx$ , jeśli krzywa jest opisana wzorem  $y$

przedziale  $[a, b]$ . W naszym przypadku funkcją podcałkową będzie  $f(x) = \sqrt{1 + \cos^2(x)}$ , a ta fun posiada funkcji pierwotnej, zatem całkę z niej w przedziale  $[a, b]$  liczymy numerycznie. Skorzystamy metody trapezów i parabol. Oszacujemy z rysunku wartości modułu drugiej i czwartej pochodnej



Rys11.2.1. Wykresy wartości bezwzględnych pochodnych: drugiej i czwartej.

Przyjmujemy  $M_2=1$ , oraz  $M_4=7$ . Wstawiając do odpowiednich wzorów otrzymujemy: stosując metodę złożoną trapezów, aby uzyskać dokładność  $10^{-6}$ , musimy podzielić przedział  $[a, b]$  na  $m=569$  części, wtedy otrzymamy wynik  $I = 1.910099$ , a stosując metodę złożoną parabol, aby uzyskać ten sam wynik wystarczy podzielić przedział na  $m=26$  części.

---