

Indeks pojęć (jest ich tu więcej, ale się nie zmieściły w spisie). Kolejność rzeczy na stronach jak w wykładach

Algebra/Analiza	2
Algorytm Dollittle'a	4
Algorytm Floatera-Hormanna	10
Algorytm PageRank	7
Algorytm QR	7
Algorytm Shermana-Morrisona	5
Algorytmy genetyczne	18
Aproksymacja	18
Całki wielowymiarowe	12
Całkowanie numeryczne	10
Cholesky	4
Ekstrapolacja Richardsona	11
Eliminacja Gaussa	3
Faktoryzacja	3
Housholder	4
Interpolacja	8
Interpolacja Hermite'a	9
Kryterium Akaike	19
Kwadratura adaptacyjna w dwóch wymiarach	13
Kwadratury Newtona-Cotesa	10
Kwadratury adaptacyjne	12
LDL	4
LU	4
Liniowe zagadnienie najmniejszych kwadratów	18
Metoda Broydena (wielowymiarowa metoda siecznych)	14
Metoda Gaussa-Seidla	6
Metoda Jacobiego	5
Metoda Laguerre'a	15
Metoda Levenberga-Marquardta	17
Metoda Newtona	13
Metoda Romberga	11
Metoda gradientów sprzężonych	6
Metoda najszybszego spadku	16
Metoda potęgowa	7
Metoda siecznych	13
Metoda trapezów	10
Metoda zmiennej metryki	17
Minimalizacja	15
Minimalizacja funkcji wielu zmiennych	16
Minimalizacja globalna	17
Minimalizacja gradientami sprzężonymi	16
Nieliniowe zagadnienie najmniejszych kwadratów	19
Obroty Givensa	5
Ortogonalna transformacja podobieństwa	7
Oscylacje Rungego	9
Particle Swarm Optimization	18
Pivoting (wybór elementu podstawowego)	7
QR	5
Reprezentacja liczb i błędy zaokrągleń	3
Rezolwenta	8
Równania algebraiczne	13
Różniczkowanie numeryczne	10
SOR (Successive Over-Relaxation)	6
SVD	5
Splajny (funkcja sklejana)	9
Transformacja Möbiusa	8
Układ równań liniowych	3
Układ równań nieliniowych	14
Uogólnione zagadnienie własne	8
Uwarunkowanie numeryczne	10
Wartości własne macierzy Hermitowskiej	8
Wzór interpolacyjny Lagrange'a	8
Zagadnienie własne	8

Macierz diagonalizowalna / normalna - Macierz  $A \in \mathbb{C}^{N \times N}$ , która ma N niezależnych liniowo wektorów własnych.

Dla X – macierz z wektorów własnych:  $X^{-1}AX = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_N\}$  //macierz diagonalna z  $\lambda$  na przekątnej

Macierz symetryczna –  $A = A^T$ , jej unormowane wektory własne tworzą bazę ortonormalną.

Macierz dodatnio określona –  $x^T A x > 0$ ; wszystkie główne wyznaczniki ( $1 \times 1, 2 \times 2, \dots$ ) i wartości własne są  $> 0$

Macierz osobliwa – wyznacznik równy 0.

Macierz hermitowska:  $H^\dagger = H$

Macierz unitarna:  $U^\dagger U = U U^\dagger = I$  oraz  $U^\dagger = U^{-1}$

Macierze A i B są podobne, jeśli istnieje taka nieosobliwa macierz S, że:  $B = S^{-1}AS$ . Mają one takie same widma.

Ortogonalna transformacja podobieństwa:  $B = O^T A O$ , albo  $A_k = O_k^T \dots O_1^T A_1 O_1 \dots O_k = P_k^T A P_k$ . O jest macierzą ortogonalną.

Macierz rzadka – liczba elementów niezerowych rośnie wolniej niż  $N^2$ .

Macierz diagonalna – ma zera poza diagonalą (przekątną), oznaczana:  $\text{diag}(a_i)$  lub  $[\text{diag}(a_i)]$

Macierz trójdzielna – tak jak diagonalna, ale nad i pod przekątną ma jeszcze jedną przekątną

Macierz trójkątna – ma same zera poniżej (górna) lub powyżej (dolna) diagonal

Macierz ortogonalna –  $AA^T = I$ ; macierz przekształcenia ortogonalnego; wyznacznik = 1 lub -1;  $(AB)^T = A^T + B^T$  dla AB, BA ortogonalnych //I to macierz jednostkowa

Macierz symetryczna rzeczywista – ma rzeczywiste wartości własne, a jej unormowane wektory własne tworzą bazę w  $\mathbb{R}^n$ . Jej norma równa się największemu modułowi wartości własnych. //n to wymiar macierzy

Ortogonalna macierz rzutowa –  $P^T = P, P^2 = P, P^T = P^{-1}$

Norma macierzy –  $\|A\| = \max\{\|Ax\| \div \|x\|, x \neq 0\}$

Promień spektralny macierzy –  $\rho = \sqrt{\|AA^T\|}$

$\rho(G) = \max\{|\lambda| : \exists y \neq 0 : Gy = \lambda y\}$  dla G diagonalizowalnej

Norma wektora x –  $\|x\| = \sqrt{(x_1^2 + x_2^2 + \dots)}$

Wektor unormowany – o długości jeden (normalizacja = dzielenie przez długość)

Dla  $A \in \mathbb{R}^{M \times N}$

Jądro operatora:  $\text{Ker } A = \{x \in \mathbb{R}^N : Ax = 0\}$

Zasięg operatora:  $\text{Range } A = \{y \in \mathbb{R}^M : \exists x \in \mathbb{R}^N : Ax = y\}$

Dla  $M = N$ :  $\dim(\text{Ker } A) + \dim(\text{Range } A) = N$

Biegun – punkt, w otoczeniu którego funkcja nie jest ograniczona

Szereg Taylora: 
$$f(x) = \sum_{k=0}^n \left( \frac{(x-a)^k}{k!} f^{(k)}(a) \right) + R_n(x, a)$$

Gradient funkcji: 
$$\nabla f = \left[ \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]$$

Hessjan – macierz drugich pochodnych.

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

W komputerze można reprezentować ściśle tylko liczby całkowite i wymierne ze skończonym rozwinięciem, ale tylko z pewnego zakresu. Inne liczby obarczone są błędem zaokrąglenia.

Liczba  $x$  jest poprawnie zaokrąglona na  $d$ -ej pozycji do liczby, którą oznaczamy  $x^{(d)}$  jeśli  $|\varepsilon| = |x - x^{(d)}| \leq \frac{1}{2} \cdot 10^d$ , gdzie  $\varepsilon$  to błąd zaokrąglenia.

Jesli  $\bar{x}$  jest przybliżeniem dokładnej wartości  $x$ , to  $k$ -tą cyfrę dziesiętną liczby  $\bar{x}$  nazwiemy znaczącą, jeśli  $|x - \bar{x}| \leq \frac{1}{2} \cdot 10^k$

$x = \bar{x} + \varepsilon$ , dla  $\varepsilon$  będącego dowolną liczbą w przedziale  $[-\frac{1}{2}10^d, \frac{1}{2}10^d]$ .

Błędy, a działania:

$x + y = \bar{x} + \bar{y} + \varepsilon_x + \varepsilon_y$  //błędy się sumują

$x \cdot y \approx \bar{x} \cdot \bar{y} + \bar{x}\varepsilon_y + \bar{y}\varepsilon_x$  //błędy się mnożą, mały błąd może znacznie urosnąć

$x / y \approx \bar{x} / \bar{y} + (1 / \bar{y}) \varepsilon_x + (\bar{x} / \bar{y}^2) \varepsilon_y$  //dzielenie przez względnie małe liczby powoduje znaczne zwiększenie błędu

Na komputerach pracujemy w arytmetyce ze skończoną dokładnością i dodawania nie zawsze jest przemienne.

Zagadnienie obliczenia  $\phi(x)$  jest numerycznie dobrze uwarunkowane, jeżeli niewielkie względne zmiany danych dają niewielkie względne zmiany rozwiązania. Inaczej jest źle uwarunkowane.

$\forall x, \bar{x} : \|\phi(x) - \phi(\bar{x})\|_{R^m} \div \|\phi(x)\|_{R^m} \leq \kappa \cdot \|x - \bar{x}\|_{R^n} \div \|x\|_{R^n}$  –  $\kappa$  to współczynnik uwarunkowania dla  $\phi(x)$  będącego funkcją  $R^n \rightarrow R^m$ . Mówi on jak bardzo błąd przybliżenia wpływa na błąd wyniku. Mniejszy = lepiej.

---

Układ równań liniowych to układ typu  $Ax = b$ , gdzie  $A$  jest macierzą  $N \times N$ ,  $x$  i  $b$  są wektorami.  $A$  i  $b$  są znane.  $\text{Det } A \neq 0$

Nie rozwiązuje się ich wzorami Cramera, bo są bardzo złożone numerycznie.

Współczynnik uwarunkowania rozwiązywania układu równań liniowych –  $\kappa = \|A\| \cdot \|A^{-1}\|$   
 $\kappa = \max\{\lambda_i\} \div \min\{\lambda_i\}$  dla macierzy symetrycznych i rzeczywistych. Bo  $\|A^{-1}\| = \max_i \{1 / |\lambda_i|\}$

Kolumny i wiersze macierzy układu równań można zamieniać miejscami i dodawać do siebie.

**Eliminacja Gaussa** prowadzi do utworzenia równań z trójkątną macierzą górną. Polega ona na dzieleniu/mnożeniu przez następne współczynniki i odejmowaniu równań stronami. Ma złożoność  $O(N^3)$ . Algorytm nie działa, gdy trzeba dzielić przez 0, ale można wtedy próbować zamienić miejscami równania. Wymaga znajomości kolumny wyrazów wolnych.

Układ z trójkątną macierzą górną można rozwiązać metodą back substitution, która ma złożoność  $O(N^2)$ . Polega na podstawianiu rozwiązań od dołu do góry.

Zamiast brać kolejne współczynniki, można jako następne wybierać największe co do modułu i permutować wiersze. Nazywa się to pivoting (wybrany współczynnik to pivot, czyli element podstawowy) i zwiększa stabilność algorytmu. Ma złożoność  $O(N^2)$ .

Pełny pivoting polega na wybieraniu pivota ze wszystkich kolumn i ma złożoność  $O(N^3)$ . Rzadko stosowany.

---

Układ równań można rozwiązać metodą macierzy odwrotnej, ALE TAK SIĘ NIE ROBI BO TO JEST  $O(2N^3)$ . Zapis  $A^{-1}b$  traktujemy jako znalezienie wektora  $z$ , takiego, że  $Az = b$ .

---

Faktoryzacja to przekształcenie macierzy równań (bez wyrazów wolnych) w iloczyn dwóch macierzy:  $A = Y \cdot Z$

**Faktoryzacja LU** – macierz ma postać  $A = L \cdot U$ , gdzie L jest macierzą trójkątną dolną z samymi jedynkami na diagonalu, a U jest trójkątną górną. Rozwiązanie równania wygląda:  $Ly = b$   
 $Ux = y$

Pierwsze równanie rozwiązujemy metodą forward substitution (podstawianie od góry do dołu), drugie back substitution. Koszt obu tych rozwiązań to  $O(2N^2)$ .

Algorytm Doolittle'a – służy do wyznaczania faktoryzacji LU. Jeśli mamy macierz A z elementami  $a_{ij}$ , L z  $l_{ij}$  oraz U z  $u_{ij}$  to:

W pierwszej kolumnie:  $u_{11} = a_{11}$ ,  $l_{21}u_{11} = a_{21}$ ,  $l_{31}u_{11} = a_{31}$  itd.  $// A_{ij} \equiv A_{\text{rzęd kolumna}}$

W drugiej kolumnie:  $u_{12} = a_{12}$ ,  $l_{21}u_{12} + u_{22} = a_{22}$ ,  $l_{31}u_{12} + l_{32}u_{22} = a_{32}$ ,  $l_{N1}u_{12} + l_{N2}u_{22} = a_{N2}$  itd.  
 itd., ma złożoność  $O(N^3)$ , bo jeden element obliczamy  $O(N)$ , a jest ich  $O(N^2)$ .

Algorytm Crouta – podobny do algorytmu Doolittle'a, jednak zawiera operacje Pivotingu. Złożoność się nie zmienia, przestawiamy tylko wyrazy już obliczone. Na koniec trzeba dokonać takich samych permutacji na kolumnie wyrazów wolnych jak na wierszach. Pełny Pivoting dla LU jest niemożliwy.

Przewaga LU nad eliminacją Gaussa jest taka, że wykonujemy to raz dla wielu równań z takimi samymi lewymi stronami. Faktoryzacja ta nie wymaga dodatkowej pamięci na macierze L i U.

**Faktoryzacja Cholesky'ego** – macierz ma postać  $A = CC^T$ , gdzie C (czynnik Cholesky'ego) jest macierzą trójkątną dolną, z dodatnimi elementami diagonalu. Jest mniej więcej o połowę szybsze od LU, ale macierz A musi być symetryczna i dodatnio określona oraz nie można wykorzystać pivotów. Elementy macierzy C można przechowywać w tym samym miejscu co elementy A, dzięki kolejności obliczeń.

$$l_{ss} = \sqrt{a_{ss} - \sum_{k=1}^{s-1} (l_{sk})^2} \quad \text{dla } s = 1, 2, \dots, n$$

$$l_{is} = \frac{(a_{is} - \sum_{k=1}^{s-1} l_{ik}l_{sk})}{u_{ii}} \quad \text{dla } i = s + 1, \dots, n$$

$// l_{ij}$  to elementy macierzy C  
 $// u_{ij}$  to elementy macierzy  $C^T$

Rozwiązania wyznaczamy:  $Cy = b$   
 $C^T x = y$

**Faktoryzacja LDL** – macierz ma postać  $A = LDL^T$ , gdzie L jest macierzą trójkątną dolną z samymi jedynkami na diagonalu, a D jest macierzą diagonalną o dodatnich elementach. Faktoryzacja ta wymaga takich samych założeń jak Cholesky'ego, ale do jej wyznaczenia nie trzeba pierwiastkować.

---  
 Jeśli macierz jest rzadka, TRZEBA TO UWZGLĘDNIĆ w naszym algorytmie.

Faktoryzacja LU dla macierzy trójdzielnej jest rzędu  $O(N)$ .  $// \text{Ale nie można pivotować}$

Rozkład Cholesky'ego macierzy M-diagonalnej ma czynnik Cholesky'ego M-diagonalny. Jednak może się pojawić wypełnienie "wewnątrz" pasma. Jeśli jest ryzyko wypełnienia, rozwiązujemy układ  $(PAP)^T (Px) = Pb$ , gdzie P to macierz permutacji dobrana tak, by wypełnienie było możliwie małe.  $// \text{wykorzystuje się do tego minimum degree algorithms, które są oparte o heurystykę}$

**Transformacja Householdera** – tworzymy macierz  $P = I - 2(uu^T / \|u\|^2)$ , dla u niezerowego oraz I - macierzy jednostkowej. Wtedy  $P^T = P$ ,  $P^2 = I$ ,  $P^T = P^{-1}$ ,  $PP^T = I$ , czyli jest symetryczna, rzeczywista oraz ortogonalna = ortogonalna macierz rzutowa. Gdy  $u = x \mp \|x\| \hat{e}_1$  ( $\hat{e}_1$  to pierwszy wektor jednostkowy), macierz P jest macierzą Householdera.

$Px = \pm \|x\| \hat{e}_1$ , czyli działając macierzą P na wektor x, zerujemy wszystkie jego składowe oprócz pierwszej, która przyjmuje wartość jego długości. Operacja ta jest  $O(N)$ .

**Faktoryzacja QR** – jeśli działamy transformacją Householdera na kolejne kolumny macierzy (zaczynając od elementu diagonalnego w każdej kolumnie), otrzymujemy macierz trójkątną górną. Iloczyn przekształceń  $P_1 \dots P_N$ , oznaczany  $Q^T$ , jest ortogonalny i nie musimy pamiętać kolejnych czynników. Faktoryzacja powstała w ten sposób wygląda tak:  $A = QR$  ( $R$  to wynikowa macierz trójkątna).

UWAGA: Transformację Householdera działamy na całą macierz, jednak do każdego kroku można zastąpić policzone elementy macierzą jednostkową. Mimo iż transformacja działa na całą macierz (przez co jest  $O(N^2)$ ) to zmienia tylko jedną kolumnę.

Złożoność całej faktoryzacji to  $O(N^3)$  (jest gorsza od LU). Układ rozwiązujemy:  $Rx = Q^T b$ , przez back substitution.

**Obroty Givensa** – macierz obrotów Givensa  $G(i,j)$  jest macierzą jednostkową, ale z elementami  $(i,i)$  i  $(j,j)$  zastąpionymi przez  $c = \cos \theta$ , oraz elementami  $(i,j)$  i  $(j,i)$  przez odpowiednio  $s = \sin \theta$  i  $-s$ .  $\theta$  to kąt obrotu. Macierz ta jest ortogonalna. Jeśli  $y = G(i,j)x$ , to składowe  $y$  wynoszą:

$$y_k = \begin{cases} cx_i + sx_j & k = i \\ -sx_i + cx_j & k = j \\ x_k & \text{poza tym} \end{cases}$$

$$c = \frac{x_i}{\sqrt{x_i^2 + x_j^2}}, \quad s = \frac{x_j}{\sqrt{x_i^2 + x_j^2}}$$

Jeśli chcemy, by  $y_j$  było równe 0, mamy warunek:

Obrót Givensa z tym warunkiem zeruje  $j$ -tą składową wybranego wektora, a  $i$ -ta ma wartość:  $\sqrt{x_i^2 + x_j^2}$

Jeśli mamy macierz trójdagonalną symetryczną  $A$ , możemy tak dobierać odpowiednie obroty, by zerować kolejne elementy poddiagonalne. Otrzymana macierz  $G_{N-1} \dots G_2 G_1 A = R$  jest trójdagonalna i otrzymujemy faktoryzację QR ( $Q = G_{N-1}^T \dots G_1^T$ ). Wtedy  $Rx = G_{N-1} \dots G_2 G_1 b$ , a równanie  $Rx = \tilde{b}$  możemy otrzymać w czasie liniowym, również w liniowym je rozwiązując (przez back substitution). A więc całkowita złożoność rozwiązania to  $O(N)$ .

---

Wzór Shermana-Morrisona – jeśli  $A_1 = A + uv^T$  to

$$A_1^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u}$$

**Algorytm Shermana-Morrisona** – gdy chcemy policzyć  $A_1^{-1}b$ , gdzie  $b$  jest znane, a  $A^{-1}$  łatwo policzyć, wtedy szukamy:

$$w = A_1^{-1}b = \left( A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) b$$

Żeby to zrobić, liczymy:

$$Az = b \wedge Aq = u$$

Wtedy:

$$w = z - \frac{v^T z}{1 + v^T q} q$$

Problem sprowadza się do rozwiązania dwóch równań z macierzą  $A$  i dzięki temu możemy łatwo policzyć równanie z macierzą  $A_1$ .

**Singular Value Decomposition (SVD)** –  $A = U[\text{diag}(\omega_i)]V^T$ , dla  $A, U \in \mathbb{R}^{M \times N}$ ,  $V \in \mathbb{R}^{N \times N}$  ortogonalna,  $M \geq N$ ,  $\omega_1 = 0$ . Macierz  $U$  jest kolumnowo ortogonalna. Kolumny macierzy  $V$  odpowiadające zerowym współczynnikom  $\omega_i$  stanowią bazę operatora  $A$ .  $|\det A| = \prod_i \omega_i$ .  $\Pi$  to symbol iloczynu (odpowiednik  $\Sigma$  dla sumy). Zastosowanie: dla  $Ax = b \rightarrow x = A^{-1}b = V[\text{diag}(\omega_i^{-1})]U^T b$   $\omega_i^{-1} = 0$ , gdy  $\omega_i = 0$  oraz  $b \in \text{Range}(A)$

Gdy macierz jest źle uwarunkowana, można użyć przybliżenia:  $|\omega_i| \leq \tau \Rightarrow \omega_i^{-1} = 0$ , gdzie  $\tau$  to zadana tolerancja. Jest to tzw. pseudoodwrotność ( $\tilde{A}^{-1}$ ). Jeśli  $A$  jest nadokreślona:  $\|A(\tilde{A}^{-1}b) - b\|_2 = \text{minimum}$

Metody iteracyjne – możemy uzyskać rozwiązanie w granicach błędu zaokrąglenia w skończonej liczbie kroków.

**Metoda Jacobiego**

$$x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^N a_{ij}x_j^{(k)} \right) / a_{ii} \quad //x_i^{(k)} \text{ oznacza przybliżenie w } k\text{-tym kroku.}$$

Zbieżna, gdy macierz jest silnie diagonalnie dominująca.

### Metoda Gaussa-Seidela

$$x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^N a_{ij}x_j^{(k)} \right) / a_{ii}$$

Różni się od metody Jacobiego tylko tym, że w jednym miejscu używamy wyrazów z aktualnego przybliżenia. Dzięki temu jest szybciej zbieżna, ale trudno ją zrównoleglić. Rozbieżna, gdy macierz nie jest symetryczna i dodatnio określona.

Ogólny zapis powyższych metod:  $Mx^{(k+1)} = Nx^{(k)} + b$ .  $A = M - N$  jest podziałem macierzy.  $\text{Det}(M)$  musi być różny od 0, a promień spektralny  $M^{-1}N < 0$ .

Dla Jacobiego:  $M = D$ ,  $N = -(L + U)$  //  $D$  to macierz diagonalna,  $L$  trójkątna dolna,  $U$  górna

Dla Gaussa-Seidla:  $M = D + L$ ,  $N = -U$

### Metoda successive over-relaxation

$$x_i^{(k+1)} = w \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^N a_{ij}x_j^{(k)} \right) / a_{ii} + (1 - w)x_i^{(k)},$$

$\omega$  to współczynnik relaksacji i należy go tak dobrać, by zminimalizować  $\rho(M_\omega^{-1}N\omega)$  //  $\rho$  to promień spektralny

**Metoda gradientów sprzężonych** – forma kwadratowa  $\frac{1}{2}x^T Ax - b^T x + c$  ma jedno minimum dla  $Ax = b$ , które jest w punkcie  $\nabla f = 0$  i jest równoważne rozwiązaniu układu. Można to wykorzystać.  $A$  jest symetryczna i dodatnio określona (i dobrze, gdy jest rzadka).

Algorytm:

```

r1 = b - Ax1, p1 = r1 //ten wariant metody nazywamy algebraicznym
while ||rk|| > ε
    αk = (rk^T rk) / (pk^T Apk)
    rk+1 = rk - αk Apk
    βk = (rk+1^T rk+1) / (rk^T rk)
    pk+1 = rk+1 + βk pk
    xk+1 = xk + αk pk
end

```

Koszt algorytmu to  $O(N^3)$  (bo  $N$  kroków, każdy  $N^2$ , przez  $Ap_k$ ). Dla macierzy rzadkiej, odpowiednio zaprogramowany algorytm ma złożoność  $O(M \cdot N^2)$ , gdzie  $M$  to szerokość pasma macierzy.

Metoda bierze nazwę od warunku sprzężenia, które spełniają wektory  $r_k$  (są wzajemnie prostopadłe w metryce  $A$ ).

Metoda ma problem, jeśli macierz ma duży współczynnik uwarunkowania, bo jest wtedy powolna.

Wtedy macierz prewarunkujemy. Bierzymy macierz  $C$ , która jest odwracalna, symetryczna, rzeczywista i dodatnio określona. Wtedy  $C^{-1}AC^{-1}Cx = C^{-1}b$ , czyli  $\tilde{A}\tilde{x} = \tilde{b}$ .

Ostatecznie:

```

r1 = b - Ax1
rozwiąż Mz1 = r1
p1 = z1
while ||rk|| > ε
    αk = (rk^T zk) / (pk^T Apk)
    rk+1 = rk - αk Apk
    rozwiąż Mzk+1 = rk+1
    βk = (rk+1^T zk+1) / (rk^T zk)
    pk+1 = zk+1 + βk pk
    xk+1 = xk + αk pk
end

```

Zagadnienie własne to problem szukania takich  $\lambda$ , że istnieje niezerowy wektor  $x$ , taki że  $Ax = \lambda x$ . Liczbę  $\lambda$  nazywamy wartością własną, a  $x$  – wektorem własnym.

Równanie  $Ax = \lambda x$  można zapisać również:  $(A - \lambda I)x = 0$ . Niezerowe rozwiązanie istnieje tylko, gdy  $\det(A - \lambda I) = 0$ . Jest to tak zwane równanie charakterystyczne macierzy, a zbiór jego rozwiązań to widmo macierzy.

Rozwiązywanie równania charakterystycznego jest na ogół numerycznie nieefektywne.

**Algorytm PageRank** – algorytm opracowany przez twórców Google, służący do liczenia rankingu stron WWW, który sprowadza się do zagadnienia własnego. Zakładamy, że strony można ponumerować. Popularność strony zależy od ilości linków, które do niej prowadzą oraz popularności stron linkujących. Niech  $r_i$  oznacza ranking  $i$ -tej strony. Wzór wygląda tak:

$r_i = C \cdot \sum_k r_k / l_k$ , czyli suma po stronach  $k$  linkujących do  $i$ , gdzie  $C > 0$  jest stałą proporcjonalności, a  $l_k$  oznacza liczbę linków na stronie  $k$ .

Jeśli przyjmiemy notację  $p_{ik} = \{l_k^{-1}$  gdy strona  $k$  linkuje do  $i$ ; 0 gdy nie linkuje}, to  $r_i = C \cdot \sum_k p_{ik} r_k$  idzie po wszystkich  $k$ . Wtedy możemy użyć macierzy  $P$  dla wartości  $p_{ij}$  oraz wektora  $r$  dla rankingów. Otrzymujemy:  $Pr = \lambda r$ , dla  $\lambda = 1/C$ . Wektor rankingów to wektor własny macierzy  $P$ .

**Metoda potęgowa** – mamy macierz  $A$ , która jest symetryczna.  $e_i$  oznacza wektor własny, zakładamy, że wartości własne są uporządkowane malejąco.

Wektor  $y$  posiada rozkład  $y = \sum_i \beta_i e_i$  ( $e_i$  to składowa wektora). Wtedy  $A^k y = \sum_i \beta_i \lambda_i^k e_i$ . Ponieważ założyliśmy, że  $\lambda_1$  jest największa, to będzie dominować dla dużych  $k$ , czyli prawa strona dąży do  $e_1$ .

Tak wyglądają kolejne iteracje.  $\|y_1\| = 1$ . Gdy  $y_k$  przestanie się zauważalnie zmieniać, liczymy  $\lambda_1 = \|z_k\|$ .

Tą metodą możemy znaleźć pierwszą wartość własną (największą). Jeśli chcemy znaleźć następne,  $y_k$  musi być prostopadły do każdego ze znalezionych  $z_k$ . Do algorytmu dodajemy

krok pośredni (ortogonalizację):  $\tilde{z}_k = z_k - e_1(e_1^T z_k)$ . Reortogonalizacja jest potrzebna przynajmniej co kilka kroków (przez błąd obcięcia), co zwiększa koszt numeryczny.

Jeżeli w algorytmie zamienimy  $A$  na  $A^{-1}$ , znajdziemy wartość najmniejszą. Ujemne wartości własne powodują, że po stabilizacji, współrzędne wektora własnego zmieniają znak. Metoda jest bardzo wolna dla wartości własnych bliskich co do modułu.

**Ortogonalna transformacja podobieństwa** – liczymy kolejno:  $P_0 = I$ ,  $A_1 = Q_1^T A Q_1$ ,  $P_1 = P_0 Q_1$ ,  $A_2 = Q_2^T A_1 Q_2$  itd. W końcu:  $P_s = P_{s-1} Q_s$  – macierz, której kolumnami są wektory własne macierzy  $A$ . //  $Q$  to jakaś macierz ortogonalna, taka, że  $A_{n+1}$  jest ortogonalnie podobne do  $A_n$ .

Jeśli interesują nas wartości własne, a nie wektory, nie musimy pamiętać zakumulowanych macierzy transformacji.

Jeśli zrobimy  $P_s^T A P_s$ , otrzymamy macierz z wartościami własnymi na diagonalu.

**Algorytm QR** – jeśli macierz ma faktoryzację  $QR$ , to  $RQ = Q^T A Q$ , czyli ortogonalna transformacja podobieństwa. Procedurę można iterować:  $A_n = R_{n-1} Q_{n-1} = Q_n R_n$ . Algorytm ten zachowuje symetrię, trójdagonalność i postać Hessenberga. Jeśli  $A$  jest diagonalizowalna oraz ma różne i rzeczywiste wartości własne, to iteracja dąży do macierzy trójkątnej górnej z wartościami własnymi na przekątnej.

Ponieważ algorytm jest drogi numerycznie ( $N^3$  na każdą iterację), używa się go dla macierzy rzadkich. Używając obrotów Givensa:  $A' = RQ = G_{N-1} \dots G_1 A G_1^T \dots G_{N-1}^T$ , faktoryzację macierzy można znaleźć w czasie liniowym, otrzymując macierz trójdagonalną symetryczną.

Jeżeli macierz nie jest symetryczna, faktoryzacja prowadzi do górnej macierzy Hessenberga, która jest macierzą trójkątną górną z jedną dodatkową diagonalą pod główną. Koszt jej faktoryzacji to  $O(N^2)$ .

Macierze w prostej postaci (trójdagonalnej/Hessenberga) diagonalizujemy za pomocą faktoryzacji QR.



**Wartości własne macierzy Hermitowskiej** –  $H$  jest macierzą hermitowską i  $Hu = \lambda u$ .  $H$  można rozłożyć na część rzeczywistą i urojoną:  $H = A + iB$ , wtedy  $A^T = A$ ,  $B^T = -B$ . Dla  $u = x + iy$ , mamy równanie: 
$$\begin{bmatrix} A & -B \\ B & A \end{bmatrix} = \lambda \begin{bmatrix} x \\ y \end{bmatrix}$$

Pierwsza macierz,  $\tilde{H}$ , ma 2 razy większe wymiary od  $H$ , czyli ma  $2N$  wartości własnych, a każda jest co najmniej dwukrotnie zdegenerowana.

**Rezolwenta** – jeśli liczba zespolona  $\xi$  nie należy do widma  $A$ , to macierz  $Z = (A - \xi I)^{-1}$ , będąca funkcją argumentu  $\xi$ , nazywa się rezolwentą.  $Au = \lambda u \Rightarrow (\lambda - \xi)^{-1}u = Zu$ , czyli  $u$  jest wektorem własnym rezolwenty do wartości własnej  $(\lambda - \xi)^{-1}$ .

**Transformacja Möbiusa** – jeśli  $u_1, u_2$  są wektorami własnymi, a  $\lambda_1, \lambda_2$  odpowiadającymi im wartościami własnymi oraz  $\exists \tau : \lambda_1 = \tau + \varepsilon_1, \lambda_2 = \tau + \varepsilon_2$  i  $|\varepsilon_{1,2}| \ll 1$ , rozważmy rezolwentę  $(A - \tau I)^{-1}$ . Wektory  $u_{1,2}$  są jej wektorami własnymi do wartości własnych  $\varepsilon_{1,2}^{-1}$ , a różnica między nimi jest duża. Czyli: bliskie wartości własne są dalekie w rezolwencie, przy takich samych wektorach własnych. Kosz policzenia rezolwenty rekompensuje to, że łatwo znaleźć rozseparowane wartości.

Rezolwenty można użyć również do poszukiwania wektora znanej wartości własnej. Rezolwentę  $(A - \tau I)^{-1}$  podstawiamy zamiast  $A$  do metody potęgowej ( $\tau$  to liczba  $\simeq \lambda$ ). Wynikiem jest wektor do wartości własnej  $\lambda$ . Macierz  $A - \tau I$  jest źle uwarunkowana, dlatego trzeba zwiększyć precyzję.

**Uogólnione zagadnienie własne** – jeśli w równaniu charakterystycznym zamienimy  $I$  na jakieś  $B$ , otrzymamy  $Ax = \lambda Bx$ , czyli uogólnione zagadnienie własne, z uogólnionymi wartościami i wektorami własnymi. Zakładamy, że  $B$  posiada faktoryzację  $B = CC^T$ .

Mamy równanie:  $\tilde{A}y = \lambda y$ ,  $\tilde{A} = C^{-1}A(C^T)^{-1}$ ,  $y = C^Tx$ . Macierz  $C^{-1}$  trzeba niestety znaleźć jawnie.

Uogólnione wektory własne są sprzężone względem macierzy  $B$  oraz ortogonalne względem iloczynu skalarnego generowanego przez macierz  $B$ .

Interpolacja jest gdy znamy kilka punktów funkcji i chcemy znaleźć co jest pomiędzy nimi. Można je połączyć prostymi liniami, ale taka funkcja może stwarzać problemy, m.in. nie jest gładka.

**Macierz Vandermonde'a** –  $n$  różnych punktów wyznacza wielomian  $n-1$  stopnia. Jeśli za  $x$  w wielomianie podstawimy kolejno znane punkty, to otrzymujemy układ równań, w którym współrzędne  $y$  punktów robią za wektor  $b$ . Wyznacznik macierzy to iloczyn różnic każdej pary punktów i jest niezerowy, gdy żadne z punktów się nie pokrywają. Koszt rozwiązania tego to  $O(n^3)$ .

**Wzór interpolacyjny Lagrange'a** – zamiast poprzedniego równania, możemy użyć:

$$f(x) = \sum_j l_j(x) f_j + E(x), \text{ gdzie } l_j(x) = \frac{(x - x_1) \dots (x - x_{j-1})(x - x_{j+1}) \dots (x - x_n)}{(x_j - x_1) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_n)}$$

$l_j$  jest wielomianem  $n-1$  stopnia i  $l_j(x_k) = \delta_{jk}$ . Jeżeli  $f(x)$  jest stopnia  $\leq n-1$ , to  $E(x)$ , czyli błąd interpolacji (reszta), znika tożsamościowo. Mówimy, że interpolacja ma dokładność  $n-1$ . Takie równanie rozwiązujemy w czasie  $O(n^2)$ .

//  $f_j \equiv f(x_j)$ ,  $n$  to liczba węzłów interpolacji

Alternatywnie:  $l_j(x) = \frac{p_n(x)}{(x - x_j)p'_n(x_j)}$ ,  $p_n$  jest iloczynem  $\prod (x - x_i)$ , a  $p'_n$  to jego pochodna.

Oznaczmy:  $y(x) = \sum_j l_j(x) f_j$  (wielomianowa część wzoru, bez reszty).

Funkcja  $F(z) = f(z) - y(z) - [f(x) - y(x)] \left( \frac{p_n(z)}{p_n(x)} \right)$  ma  $n+1$  miejsc zerowych równych  $x_1, \dots, x_n$ , zakładając, że jest odpowiednio gładka, stosujemy  $n$ -krotnie twierdzenie Rolle'a i okazuje się, że  $n$ -ta pochodna ma co najmniej jedno miejsce zerowe, które oznaczamy  $\xi$ .

Otrzymujemy:  $E(x) = \left( \frac{p_n(x)}{p_n(x)} \div n! \right) f^{(n)}(\xi)$ . Trudność w oszacowaniu błędu interpolacji polega na trudności w oszacowaniu wysokich pochodnych funkcji.



Oscylacje Rungego – zbyt duża liczba punktów do utworzenia wielomianu prowadzi w granicach przedziału do tzw. oscylacji Rungego, czyli wysokich amplitud między punktami. Dzieje się tak, bo narzucamy punkty wielomianom, które są sztywną strukturą.

**Interpolacja Hermite'a** – jeśli oprócz wartości funkcji w punktach znamy wartości pochodnej, można skonstruować interpolację wielomianową rzędu  $2n - 1$  (bo mamy  $2n$  warunków).

Mamy funkcję:  $y(x) = \sum_i h_i(x)f_i + \sum_i \bar{h}_i(x)f'_i + E(x)$ , gdzie  $h_i(x) = (1 - 2(x - x_i)l'_i(x_i))l_i^2(x)$  i  $\bar{h}_i(x) = (x - x_i)l_i^2(x)$ , oraz

$$E(x) = \frac{p_n^2(x)}{(2n)!} f^{(2n)}(\xi)$$

Wielomian ten zgadza się z interpolowaną funkcją oraz jej pochodną w zadanych węzłach.

**Funkcja sklejana** – "splajn" rzędu  $k$  to funkcja, która lokalnie jest wielomianem rzędu  $k$  i jest  $k-1$  krotnie różniczkowalna (czyli pochodne rzędu  $< k-1$  są ciągłe).

Splajn kubiczny to splajn trzeciego rzędu (cubic spline). Jeśli oprócz wartości funkcji znamy jakieś wartości  $\xi_i$  w węzłach, konstruujemy takie wyrażenie interpolacyjne, że będzie się zgadzać z wartościami  $f_i$  w węzłach, a druga pochodna wyrażenia będzie się zgadzać z  $\xi_i$ .

W przedziałach  $[x_j, x_{j+1}]$  konstruujemy wielomian 3-stopnia:  $y_j(x) = Af_j + Bf_{j+1} + C\xi_j + D\xi_{j+1}$ , dla

$$A = \frac{x_{j+1} - x}{x_{j+1} - x_j}, B = \frac{x - x_j}{x_{j+1} - x_j}, C = \frac{1}{6}(A^3 - A)(x_{j+1} - x_j)^2, D = \frac{1}{6}(B^3 - B)(x_{j+1} - x_j)^2$$

Jednak tak na prawdę nie znamy wartości  $\xi$ . Ale,  $\xi$  to wartości drugich pochodnych, czyli żeby istniały, pierwsze pochodne muszą być ciągłe, czyli możemy zarządzić, by pochodna  $y_j(x)$  z prawej strony równała się  $y_{j+1}(x)$  z lewej. Otrzymamy równanie:

$$\frac{x_j - x_{j-1}}{6}\xi_{j-1} + \frac{x_{j+1} - x_{j-1}}{3}\xi_j + \frac{x_{j+1} - x_j}{6}\xi_{j+1} = \frac{f_{j+1} - f_j}{x_{j+1} - x_j} - \frac{f_j - f_{j-1}}{x_j - x_{j-1}}$$

Jest to trójdiodagonalny układ równań na nieznane wartości  $\xi$ .

Ponieważ dla  $n$  węzłów mamy  $n-2$  wewnętrznych punktów zszycia, mamy układ  $n-2$  równań z  $n$  niewiadomymi, czyli trzeba podać więcej warunków. Najczęściej jest to  $\xi_1 = \xi_n = 0$ , przez co uzyskujemy naturalny splajn kubiczny. Możemy oczywiście narzucić inne warunki jeśli trzeba.

Jeśli węzły są równoodległe, równanie jest szczególnie proste:

$$\begin{bmatrix} 4 & 1 & & & & \\ 1 & 4 & 1 & & & \\ & 1 & 4 & 1 & & \\ \dots & \dots & \dots & \dots & \dots & \dots \\ & & & 1 & 4 & 1 \\ & & & & 1 & 4 \end{bmatrix} \begin{bmatrix} \xi_2 \\ \xi_3 \\ \xi_4 \\ \vdots \\ \xi_{n-2} \\ \xi_{n-1} \end{bmatrix} = \frac{6}{h^2} \begin{bmatrix} f_1 - 2f_2 + f_3 \\ f_2 - 2f_3 + f_4 \\ f_3 - 2f_4 + f_5 \\ \vdots \\ f_{n-3} - 2f_{n-2} + f_{n-1} \\ f_{n-2} - 2f_{n-1} + f_n \end{bmatrix}$$

Macierz po lewej posiada łatwą do znalezienia faktoryzację Cholesky'ego. Z prawej strony są ilorazy różnicowe z dokładnością do czynników "6".

Aby użyć splajnów w praktyce, robimy:

- 1) Rozwiązujemy układ na wartości  $\xi$ , co jest  $O(n)$  i robimy to tylko raz.
- 2) Wykorzystujemy równanie  $y_j(x)$  tyle razy, ile chcemy znaleźć wartości między węzłami. W każdym przedziale  $[x_j, x_{j+1}]$  używamy odpowiedniego wielomianu  $y_j(x)$ .

**Splajny bikubiczne** – jeśli mamy funkcję dwóch zmiennych i znamy siatkę punktów, wiersze tej siatki odpowiadają wartościom  $y$ , a kolumny  $x$ . Chcemy znaleźć wartość funkcji punktów wewnętrznych siatki:  $f(x^*, y^*)$ .

- 1) Przeprowadzamy splajn przez każdy z wierszy.  $//O(n^2)$ , bo  $n$  wierszy razy  $O(n)$
- 2) Liczymy wartości każdego z tych splajnów w punkcie  $x^*$ . Otrzymujemy  $n$  wartości w  $f(x^*, y_n)$
- 3) Przez punkty  $x^*$  przeprowadzamy splajn po  $y$  i liczymy wartość splajnu w  $(x^*, y^*)$

Zamiast wielomianami, możemy interpolować za pomocą funkcji wymiernych:

$r(x) = (P_\mu(x) \div Q_\nu(x))$ . Funkcje wymierne mogą modelować więcej różnych zachowań.

**Algorytm Floatera-Hormanna** – mamy węzły  $x_1, \dots, x_n$  i wybieramy parametr interpolacji  $0 \leq d \leq n$ . Jeśli  $p_i(x)$  jest wielomianem interpolującym na  $x_i, \dots, x_{i+d}$ , to

$$r(x) = \frac{\sum_{i=0}^{n-d} \lambda_i(x) p_i(x)}{\sum_{i=0}^{n-d} \lambda_i(x)}, \lambda_i(x) = \frac{(-1)^i}{(x - x_i) \cdots (x - x_{i+d})}$$

$r(x)$  jest gładką mieszanką wielomianów interpolacyjnych. Nie ma biegunów na osi rzeczywistej.

Alternatywnie (barycentrycznie):

$$r(x) = \frac{\sum_{k=0}^n \frac{w_k}{x - x_k} f_k}{\sum_{k=0}^n \frac{w_k}{x - x_k}}, w_k = \sum_{i \in J_k} (-1)^i \prod_{j=i, j \neq k}^{i+d} \frac{1}{x_k - x_j}$$

//Wagi  $w_k$  obliczamy tylko raz.

// $J_k = \{i \in I : k - d \leq i \leq k\}, I = \{0, 1, \dots, n - d\}$

Parametr  $d$  wystarczy brać 3, ale nieraz potrzebne jest 8. Jeśli interpolowana funkcja jest dostatecznie gładka to błąd interpolacji wynosi  $O(h_{\max}^{d+1})$  //  $h_{\max}$  to największa odległość między punktami

Algorytm ten w pewnych sytuacjach jest tak dokładny jak splajny i mniej złożony.

**Różniczkowanie numeryczne** –  $f'_j = (f_{j+1} - f_{j-1}) \div 2h$  – symetryczne przekształcenie wzoru na pochodną. Stosowanie go jest niezalecane, bo jest mało dokładny i zależy od rozmiaru kroku interpolacji (większy = gorzej).

Niesymetryczne wzory mają postać:  $(f_{j+1} - f_j) \div h$  oraz  $(f_j - f_{j-1}) \div h$  i są jeszcze gorsze.

Różniczkowanie splajnu:  $y_{j+1}' = 1/h (f_{j+2} - f_{j+1}) - 1/6 h (2\xi_{j+1} + \xi_{j+2})$  (w lewym krańcu), po przekształceniach:

$$y'(x_j) = \frac{f_{j+1} - f_{j-1}}{2h} - \frac{1}{12} h (\xi_{j+1} - \xi_{j-1}), j > 1. \text{ Wykorzystujemy znajomość drugiej pochodnej.}$$

Numerycznie wolno obliczać całki, o których wiemy, że istnieją. Kwadratura to całka przybliżona, którą otrzymuje się przez całkowanie odpowiednich wielomianów interpolacyjnych. Jeśli  $f(x) = y(x) + E(x)$  (gdzie  $y(x) = \sum_i h_i(x) f_i$  + ewentualne pochodne), a  $E(x)$  to błąd interpolacji, otrzymujemy:  $\int_a^b f(x) dx = \sum_i H_i f_i + E$ , dla  $H_i = \int_a^b h_i(x) dx$ ,  $E = \int_a^b E(x) dx$ .

Należy tak wykonać interpolację, by całki z ewentualnych pochodnych zniknęły tożsamościowo. Z całkowania wielomianów Lagrange'a otrzymuje się kwadratury Newtona-Cotesa, a z interpolacji Hermite'a – kwadratury Gaussa. Jeśli granice całkowania wielomianów Lagrange'a są węzłami interpolacji, to otrzymujemy zamknięte kwadratury Newtona-Cotesa. Błąd całkowania  $E$  powinien być proporcjonalny do pochodnej rzędu  $(n+1)$ , jednak dla parzystych wielomianów jest on o rząd wyższy (ze względu na symetrie i takie tam).

**Kwadratury Newtona-Cotesa** – opierają się na interpolacji wielomianami niskiego rzędu, bo są prostsze i jest mniej oscylacji Rungego, oraz szacowanie pochodnych wysokiego rzędu jest trudne.

4 najczęściej stosowane kwadratury: //  $\zeta \in [a, b]$

Metoda trapezów

$$\int_a^b f(x) dx \simeq \frac{b-a}{2} (f_0 + f_1)$$

$$E = -\frac{1}{12} (b-a)^3 f''(\zeta)$$

Metoda 3/8

$$\int_a^b f(x) dx \simeq \frac{b-a}{8} (f_0 + 3f_1 + 3f_2 + f_3)$$

$$E = -\frac{3}{80} \left(\frac{b-a}{3}\right)^5 f^{(4)}(\zeta)$$

Metoda Simpsona

$$\int_a^b f(x) dx \simeq \frac{b-a}{6} (f_0 + 4f_1 + f_2)$$

$$E = -\frac{1}{90} \left(\frac{b-a}{2}\right)^5 f^{(4)}(\zeta)$$

Metoda Milne'a

$$\int_a^b f(x) dx \simeq \frac{b-a}{90} (7f_0 + 32f_1 + 12f_2 + 32f_3 + 7f_4)$$

$$E = -\frac{8}{945} \left(\frac{b-a}{4}\right)^7 f^{(6)}(\zeta)$$

Wzór na kwadraturę przybliżoną podaje pole pod wykresem odpowiedniej funkcji.

Większą dokładność, niż kwadraturami wyższego poziomu, można uzyskać dzieląc przedział na mniejsze i używając niższych kwadratur – są to tak zwane kwadratury złożone. Procedurę tą można iterować. Podziały najlepiej zagęszczać tak, by węzły grubszego podziału były też węzłami gęstszego, w ten sposób można korzystać z obliczonych już wartości.

Stosowanie kwadratur złożonych dodatkowo zmniejsza błąd. Np. dla metody trapezów, błąd zmniejsza się  $n^2$  razy dla  $n$  podziałów (same błędy się sumują).

Złożony wzór trapezów (dla wielokrotnych podziałów) ma postać:  $I_N \simeq h (\frac{1}{2}f_0 + f_1 + \dots + f_N + \frac{1}{2}f_N)$ , gdzie  $h$  to długość najmniejszego podprzedziału (średnica podziału),  $f_N$  jest wartością w prawym krańcu przedziału, a  $N = 2^k$ . Widać, że przy kolejnych iteracjach trzeba pamiętać tylko sumę wartości z węzłów wyższego rzędu. Kończymy iterować, gdy uzyskamy zadaną dokładność. Jeżeli wartość całki jest bardzo mała/duża, trzeba stosować dokładność względną:

$$\frac{|I_{k+1} - I_k|}{|I_k| + \varepsilon'} < \varepsilon // 0 < \varepsilon' \ll 1, \text{ zabezpiecza przed zerem w mianowniku}$$

**Ekstrapolacja Richardsona** – jak iterujemy zagęszczanie przedziału, otrzymujemy ciąg różnych przybliżeń całki, które możemy wykorzystać. Jeśli mamy 2 całki z metody trapezów:

$$\begin{aligned} I &= I_n - \frac{(b-a)^3}{12n^2} f''(\zeta_n), \quad n \text{ i } 2n \text{ to liczba podprzedziałów. Jeśli pochodne są odpowiednio podobne} \\ I &= I_{2n} - \frac{(b-a)^3}{12(2n)^2} f''(\zeta_{2n}) \quad (\text{zakładamy, że są równe}), \text{ to po wyeliminowaniu ich mamy:} \end{aligned}$$

$$I \simeq \frac{4I_{2n} - I_n}{3}$$

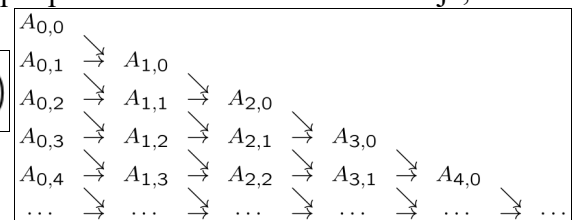
Jeśli ciąg otrzymanych przybliżeń jest monotoniczny, to znaczy, że funkcja nie zmienia bardzo gładkości, czyli założenie o stałości drugiej pochodnej jest mniej więcej dobre. Jeśli jednak ten ciąg nie jest monotoniczny, ekstrapolacja jest wątpliwa.

**Metoda Romberga** – możemy wyjść dalej poza ekstrapolację Richardsona. Korzystamy z faktu, że błąd metody trapezów zawiera jedynie parzyste potęgi średnicy przedziałów:

$$I = \int_a^b f(x) dx = h \left( \frac{1}{2}f_0 + f_1 + f_2 + \dots + f_{N-1} + \frac{1}{2}f_N \right) + \sum_{j=1}^{\infty} \alpha_j h^{2j}$$

Niech  $A_{0,k}$  oznacza przybliżenie całki uzyskane wzorem trapezów z  $2^k$  podprzedziałami. Jeśli całka istnieje, to  $\lim_{k \rightarrow \infty} A_{0,k} = I$  //może również zachodzić, gdy całki nie ma

Definiujemy:  $A_{n,k} = \frac{1}{4^n - 1} (4^n A_{n-1,k+1} - A_{n-1,k})$



Ułożone po kolei elementy przypominają macierz trójkątną dolną (rysunek ↑), przy czym musimy pamiętać tylko ostatnio policzony wiersz. Na diagonalu macierzy leżą wyrazy  $A_{k,0}$ , reszta wiersza to  $A_{i,j}$ , gdzie  $i$  to kolumna, a  $j$  to wiersz.

Ciąg  $A_{0,k}$  można przekształcić w ciąg  $A_{k,0}$ , który jest szybciej zbieżny do tej samej wartości. Robimy to za pomocą specjalnej macierzy (trójkątnej dolnej).

Ogólne postępowanie:

Obliczamy  $A_{0,k}$  za pomocą złożonej metody trapezów dla  $2^k$  podprzedziałów. Aktualnym przybliżeniem całki jest  $A_{k,0}$ . Obliczamy  $A_{0,k+1}$  korzystając ze złożonej metody trapezów dla  $2^{k+1}$  podprzedziałów. Zapęlamy cały wiersz korzystając z powyższego wzoru. Kończymy, gdy  $A_k$  i  $A_{k+1}$  różnią się mniej niż zadana tolerancja (albo jest za dużo iteracji, czyli metoda jest rozbieżna).

Metoda Romberga, w zależności od przypadku, może być dużo szybsza lub nieco wolniejsza od trapezów złożonych.

Całki nieskończone liczymy w ten sposób, że wybieramy jakąś dużą stałą  $A$  i dzielimy całkę na 2 przedziały:  $[0, A]$  oraz  $[A, \infty)$ .  $[0, A]$  liczymy numerycznie, a  $[A, \infty]$  analitycznie. Całka musi być oczywiście zbieżna oraz dobieramy  $A$  tak, by dla  $x > A$  zachodziło  $|f(x)| \leq B \cdot g(x)$ , gdzie  $g(x)$  jest funkcją, którą łatwo policzyć analitycznie.

**Kwadratury adaptacyjne** – jeśli funkcja, którą całkujemy posiada lokalne oscylacje, może być problem z dobraniem optymalnego kroku całkowania. Kwadratury adaptacyjne same lokalnie dobierają krok. W tym celu algorytm musi mieć dwa niezależne oszacowania całki, których różnica jest miarą popełnianego błędu. Szukamy przybliżenia:

$$\int_a^b f(x) dx \simeq I = \sum_{j=1}^N I_{[x_{j-1}, x_j]}$$

gdzie sumujemy numeryczne przybliżenia całki w przedziałach.  $\varepsilon_{[x_{j-1}, x_j]}$  oznacza błąd całkowania na danym przedziale, a  $\tau$  to zadany błąd maksymalny.

$$|\varepsilon_{[x_{j-1}, x_j]}| < \frac{x_j - x_{j-1}}{b - a} \tau$$

Chcemy oszacować całkę  $\int_a^b f(x) dx$  dla  $a \leq \alpha < \beta \leq b$ . Mamy aktualne przybliżenie na przedziale  $I_{[\alpha, \beta]}$ , mamy też  $I$ , czyli aktualną zakumulowaną sumę na tych przedziałach, które spełniły też warunek z  $\tau$ . Liczymy:

$$I_{[\alpha, (\alpha+\beta)/2]} \simeq \int_{\alpha}^{(\alpha+\beta)/2} f(x) dx, \quad I_{[(\alpha+\beta)/2, \beta]} \simeq \int_{(\alpha+\beta)/2}^{\beta} f(x) dx$$

Mając dwa oszacowania, czyli  $I_{[\alpha, \beta]}$  oraz  $I_{[\alpha, (\alpha+\beta)/2]} + I_{[(\alpha+\beta)/2, \beta]}$ , szacujemy błąd  $\varepsilon_{[x_{j-1}, x_j]}$ .

Jeżeli błąd nie spełnia warunku z  $\tau$ , prawy podprzedział  $[(\alpha+\beta)/2, \beta]$  odkładamy na stos i powtarzamy procedurę dla lewego podprzedziału. \*Jeśli spełnia oszacowanie, dodajemy obliczoną wartość do całki po przedziale  $[\alpha, \beta]$  i zdejmujemy ze stosu kolejny przedział, dla którego powtarzamy procedurę.

Jeśli stos się przepełni, to znaczy, że algorytm zawiódł, albo zadaliśmy za dużą dokładność (wtedy zwiększamy  $\tau$ ).

Do algorytmu zwykle stosujemy kwadratury niskiego rzędu (trapezów, Simpsona). Tego przybliżenia całki używamy w kwadraturach adaptacyjnych (dla trapezów) w \*:

$$\int_{\alpha}^{\beta} f(x) dx \simeq I_{[\alpha, (\alpha+\beta)/2]} + I_{[(\alpha+\beta)/2, \beta]} + \frac{1}{3} (I_{[\alpha, (\alpha+\beta)/2]} + I_{[(\alpha+\beta)/2, \beta]} - I_{[\alpha, \beta]})$$

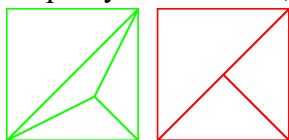
Wynika on z pewnych różnych długich przekształceń wzoru metody trapezów.

---

Dla całek wielowymiarowych o wymiarze  $> 2$ , używamy metod Monte Carlo. Dla 2-wymiarowych, tradycyjna metoda jest zwykle lepsza.

**Całkowanie po siatce prostokątnej** – licząc całkę  $I = \int_a^b dx \int_c^d dy f(x, y)$  po obszarze prostokątnym, korzystamy z metody iteracyjnej:  $I = \int_a^b g(x) dx$ , dla  $g(x) = \int_c^d f(x, y) dy$ . Stosujemy kwadratury jednowymiarowe, co wymaga  $N^2$  obliczeń. Można też liczyć w odwrotnej kolejności.

**Kwadratura adaptacyjna w dwóch wymiarach** – stosujemy, gdy obszar całkowania nie jest prostokątem. Wtedy triangulujemy obszar, czyli pokrywamy go trójkątami w taki sposób, że trójkąty albo się nie stykają, albo mają wspólny wierzchołek, albo wspólną krawędź.



Przykład **poprawnej** i **niepoprawnej** triangulacji.

Wstępnej triangulacji możemy dokonać ręcznie. Natomiast dalej można tworzyć trójkąty ze środków ciężkości (średniej współrzędnych wierzchołków) już istniejących trójkątów, co tworzy 3 nowe trójkąty. W wyższych wymiarach, pokrywamy sympleksami (np. czworościanami w  $R^3$ ).

Można zauważyć, że dla trójkątów  $f(x_1, y_1)$ ,  $f(x_2, y_2)$ ,  $f(x_3, y_3)$  wyznaczają płaszczyznę. Objętość graniastoslupa ściętego utworzonego z tych płaszczyzn może służyć jako przybliżenie całki, co jest dwuwymiarowym odpowiednikiem metody trapezów. Zagęszczanie triangulacji wymaga tylko jednego dodatkowego policzenia wartości funkcji.

Postępowanie:

Mając trójkąt  $\Omega$ , liczę przybliżenie całki ze "wzoru graniastoslupów":  $I_\Omega$ . Dzielę trójkąt na 3 potomne i liczę przybliżenia całek:  $I_{\Omega_1}$ ,  $I_{\Omega_2}$ ,  $I_{\Omega_3}$ . Miarą błędu jest:  $\epsilon_\Omega = I_\Omega - (I_{\Omega_1} + I_{\Omega_2} + I_{\Omega_3})$ .

$|\epsilon_\Omega| < \frac{S_\Omega}{S_D} \tau$  Jeśli zachodzi taki warunek  $\leftarrow$ , gdzie  $S_\Omega$  to powierzchnia trójkąta, a  $S_D$  to cała powierzchnia, za przybliżenie całki uznaję  $I_{\Omega_1} + I_{\Omega_2} + I_{\Omega_3}$ . Zdejmuję kolejny trójkąt ze stosu i powtarzam. Jeśli warunek nie został spełniony, odkładam 2 trójkąty na stos i powtarzam procedurę dla trzeciego. Jeśli stos jest pusty i nie ma więcej trójkątów, to sukces. Porażka, gdy stos się przepełni.

Procedurę tą można stosować do obszarów niewielokątowych, zagęszczając i uzyskując coraz lepsze przybliżenie.

**Równanie algebraiczne** ma postać  $f(x) = 0$ . Zakładamy, że  $f$  jest ciągła i różniczkowalna ileś razy. Miejsce zerowe jest k-krotne, jeśli funkcja zeruje się tam z k-1 pochodnymi.

**Metoda bisekcji** – wybieramy sobie dwa punkty takie, że funkcja zeruje się między nimi. Punkt pośrodku tych dwóch traktujemy jako przybliżenie miejsca zerowego. Sprawdzamy w którym z dwóch przedziałów funkcja się zeruje i powtarzamy metodę na tym przedziale. Jeśli przybliżona wartość funkcji będzie mniejsza modulem od  $\epsilon$  (jakieś zadanej dokładności) to przerywamy. Zbieżność tej metody jest liniowa. Nie działa dla miejsc zerowych o parzystej krotności.

**Metoda regula falsi** – wybieramy dwa punkty jak poprzednio i prowadzimy sieczną. Jako punkt trzeci wybieramy punkt przecięcia siecznej z OX:  $x_3 = (f(x_1)x_2 - f(x_2)x_1) / (f(x_1) - f(x_2))$ . Postępujemy jak w bisekcji.

**Metoda siecznych** – wybieramy dwa dowolne, różne punkty. Trzeci punkt to punkt przecięcia siecznej z OX, ale jako następne dwa punkty wybieramy zawsze dwa ostatnie. Metoda ta może być nieco szybciej zbieżna niż regula falsi, ale niekoniecznie musi być zbieżna do miejsca zerowego.

---

Możemy dokonać interpolacji funkcji odwrotnej. Jej wartość w zerze jest miejscem zerowym funkcji interpolowanej.

**Metoda Newtona** – kolejne iteracje wyglądają tak:  $x_{n+1} = x_n - f(x_n) / f'(x_n)$ . Wynika to z zastosowania rozwinięcia Taylora:  $f(x_0 + \delta) \simeq f(x_0) + \delta \cdot f'(x_0)$  i żądania, żeby lewa strona rozwinięcia zniknęła. Metoda może być rozbieżna, gdy zastosowana na pewien przedział  $[a, b]$  wychodzi poza niego. Może również prowadzić do cykli. Pochodną w miarę możliwości powinniśmy znać analitycznie, inaczej metoda na ogół nie ma sensu. Metoda ta jest tym bardziej efektywna, im bliżej miejsca zerowego zacznie. Dla miejsc o nieparzystej krotności ma zbieżność kwadratową, liniową dla parzystych. Można ją uogólnić na liczby zespolone, a baseny atrakcji tworzą często wtedy fraktale.



Jeśli iloraz w iteracji pomnożymy przez pewną stałą  $\alpha \in (0, 1]$ , to otrzymamy łumioną metodę Newtona, i użycie jej na kilka kroków pozwala wyjść z wielocyklu.

Dla drugiej pochodnej, iteracje wyglądają tak:

$$x_{n+1} = x_n - \frac{2f(x_n)}{f'(x_n) \pm \sqrt{[f'(x_n)]^2 - 2f(x_n)f''(x_n)}}$$

Znak dobieramy, by moduł mianownika był większy. Może prowadzić w iteraty zespolone dla rzeczywistych startów, w przeciwieństwie do normalnych iteracji.

**Metoda Halleya** – wygląda jak metoda Newtona, ale w mianowniku iteracji jest  $\sqrt{|f(x)|}$ . Po przekształceniu:

$$x_{n+1} = x_n - \frac{2f(x_n)f'(x_n)}{2[f'(x_n)]^2 - f(x_n)f''(x_n)}$$

---

Metody Newtona można użyć dla układów równań nieliniowych. Równanie takie ma postać:

$g_k(x_1, \dots, x_n) = 0$ , dla  $k = 1, \dots, n$ .  $g$  jest funkcją  $\mathbb{R}^n \rightarrow \mathbb{R}^n$ .

Iteracja metody Newtona ma postać:  $x_{k+1} = x_k - J^{-1}(x_k)g(x_k)$  //  $J$  to jakobian funkcji  $g$

Jakobian trzeba liczyć w każdym kroku, ale można to robić co kilka, dla dużych macierzy.

Minimalizacja funkcji w celu znalezienia rozwiązań układu równań to zły pomysł. Chyba, że połączy się ją z metodą Newtona:

1.  $w = 1$ . Oblicz  $\delta x$ .

2.  $x_{\text{test}} = x_i + w \delta x$ .

3. Jeśli  $G(x_{\text{test}}) < G(x_i)$ , to

(a)  $x_{i+1} = x_{\text{test}}$

(b) goto 1

4. Jeśli  $G(x_{\text{test}}) > G(x_i)$ , to

(a)  $w \rightarrow w/2$

(b) goto 2

$\delta x = -J^{-1}g$ , przerywamy, gdy  $w$  jest za małe. Jest to forma łumionej metody Newtona. Jeśli znalezione minimum nie jest globalnym (tym co nas interesuje), to zaczynamy gdzieś indziej.

**Metoda Broydena (a.k.a. wielowymiarowa metoda siecznych)** – dobra, gdy nie znamy pochodnych, albo liczenie jakobianu jest zbyt kosztowne.  $x_{i+1} = x_i - B_i^{-1}g(x_i)$ , gdzie  $B$  jest przybliżeniem jakobianu, liczonym:

$$B_{i+1} = B_i + \frac{(\Delta g_i - B_i \Delta x_i)(\Delta x_i)^T}{(\Delta x_i)^T \Delta x_i}$$

natomiast  $\Delta x = x_{i+1} - x_i$ , a  $\Delta g_i = g(x_{i+1}) - g(x_i)$ .

Do obliczenia  $B_{i+1}^{-1}g(x_{i+1})$  można skorzystać ze wzoru Shermana-Morrisona. Metoda wymaga podania  $B_1$  (np.  $J(x_1)$ ) oraz  $x_1$ .

Gdy szukamy miejsc zerowych wielomianu, interesują nas wszystkie.

Wartości wielomianu najlepiej liczyć algorytmem Hornera:  $P = a_n$ ,  $k = n$

{  $k--$ ;  $P = P \cdot z + a_k$  dopóki  $k > 0$  } //  $z$  to liczba dla której liczymy wielomian  $W(z)$

$$|\varepsilon| \simeq \frac{\sum_{k=0}^n \delta_k \tilde{z}_0^k}{|\tilde{P}'_n(\tilde{z}_0)|}$$

Oszacowanie wpływu zaburzeń współczynników na zaburzenie miejsca zerowego (zwykle nie znamy dokładnych współczynników).  $\tilde{P}_n$  oznacza zaburzony wielomian.

//  $\tilde{P}_n(\tilde{z}_0) = 0 \Rightarrow \tilde{z}_0 = z_0 + \varepsilon$  oraz  $\tilde{a}_k = a_k + \delta_k$ ,  $i \varepsilon, \delta_k \ll 0$

Zagadnienie to jest złe uwarunkowane. Przykład Wilkinsona:  $W(z) = (z+1)\dots(z+20)$ . Jeśli zaburzymy jeden współczynnik (przy wysokiej potęgze), miejsce zerowe zaburzy się 7 rzędów wielkości bardziej.

Powyższe oszacowanie załamuje się dla wielokrotnych miejsc zerowych, które w szczególności mogą się rozszczepiać. Żeby zapobiec takim różnym niedokładnościom, trzeba użyć odpowiedniej metody do szukania miejsc zerowych wielomianów oraz zmniejszać stopień wielomianu po każdym znalezionym miejscu.



**Metoda Laguerre'a** –  $P_n(z)$  to wielomian stopnia  $n$ . Mamy taką iterację:

$$z_{i+1} = z_i - \frac{n P_n(z_i)}{P'_n(z_i) \pm \sqrt{(n-1) \left( (n-1) [P'_n(z_i)]^2 - n P_n(z_i) P''_n(z_i) \right)}}$$

Znak w mianowniku wybieramy tak, by moduł mianownika był największy. Metoda jest zbieżna sześciennie do wszystkich pojedynczych miejsc zerowych, liniowo do wielokrotnych.

Może się zdarzyć, że trafimy drugi raz na to samo miejsce zerowe. Żeby tego uniknąć, znajdujemy faktoryzację:  $P_n(z) = (z - z_1)P_{n-1}(z)$ , czyli obniżamy stopień wielomianu. Jednak ze względu na to, że zaburzenia współczynników powodują duże rozbieżności, używamy miejsca zerowego  $z_2$  znalezione w wielomianie  $P_{n-1}$  do metody Laguerre'a w pierwotnym wielomianie (wygładzamy). Następnie znajdujemy faktoryzację  $P_{n-2}$  itd.

Uwagi: jeśli wielomian ma rzeczywiste współczynniki i znajdziemy zespolone miejsce zerowe, to automatycznie istnieje drugie sprzężone z nim miejsce. Jeżeli wielomian ma niewielkie, całkowite współczynniki, można wykorzystać znane twierdzenia do znalezienia wymiernych miejsc zerowych. Jeśli faktoryzując dojdziemy do wielomianu stopnia 2, liczymy jego pierwiastki ze znanego wzoru i je również wygładzamy.

**Macierz stowarzyszona** – może być użyta do szukania miejsc zerowych wielomianu. Na diagonalu same jedynki, współczynniki wielomianu zajmują ostatni rząd, z zerowym po lewej i  $a_{n-1}$  po prawej (wszystkie mają też odwrotny znak). Ma ona równanie charakterystyczne:

$x^n + a_{n-1}x_{n-1}^{n-1} + a_{n-2}x_{n-2}^{n-2} + \dots + a_1x + a_0 = 0$ . W tym przypadku,  $a_n$  równy jest 1, co nie ma wpływu na miejsca zerowe.

Minimalizacja to szukanie minimum (jednego lub więcej) jakiejś funkcji. Jeśli szukamy maksimum funkcji  $f(x)$ , to tak jak szukanie minimum  $g(x) = -f(x)$ . Jesteśmy w stanie zlokalizować minimum z dokładnością równą co najwyżej pierwiastkowi z precyzji obliczeń. Szukanie przerywamy, gdy zmiana kolejnych wartości funkcji jest mniejsza od zadanej tolerancji.

Mówimy, że punkty  $a, b, c$  ograniczają minimum, jeśli  $f(a) > f(b)$  oraz  $f(c) > f(b)$ .

Wybieramy dwa punkty. Najczęściej wyznaczają one kierunek spadku funkcji. W takim wypadku, bierzemy trzeci punkt, który jest oddalony w tym kierunku o jakiś krok. Jeśli powyższy warunek nie jest spełniony, podwajamy krok. Trzeba założyć maksymalną liczbę kroku (albo rozmiar), żeby nie iść sobie w nieskończoność.

Jak mamy już 3 punkty, które spełniają warunek, wybieramy jakiś czwarty pomiędzy nimi i posuwamy się w dół, zmieniając odpowiednie punkty (i obliczone wartości funkcji, żeby nie liczyć znowu). Warunkiem skończenia iteracji jest:  $|c - a| < \tau(|b| + |d|)$ .

**Metoda złotego podziału** – minimum szukamy w większym z dwóch przedziałów:  $[a, b]$ ,  $[b, c]$ . Żeby zminimalizować szansę na minimum w mniejszym przedziale, przedziały powinny spełniać złotą proporcję. Jeśli przedział  $[a, b]$  jest dłuższy to:  $|b - a| \div |c - a| = |c - b| \div |b - a|$ . Żeby to zapewnić, robimy: {jeśli  $|b - a| > |c - b|$  to  $d = a + \omega \cdot |b - a|$ , inaczej  $d = b + \omega \cdot |c - b|$ }.  $\omega \simeq 0.381966$ . Metoda ta jest zbieżna liniowo.

**Interpolacja paraboliczna i metoda Brenta** – stosujemy wzór:

$$d = \frac{1}{2} \cdot \frac{a^2(f_c - f_b) + b^2(f_a - f_c) + c^2(f_b - f_a)}{a(f_c - f_b) + b(f_a - f_c) + c(f_b - f_a)} // f_a \equiv f(a)$$

Czyli tworzymy parabolę między trzema punktami i bierzemy jej minimum jako  $d$ . Mogą pojawić się kłopoty, gdy początkowo funkcja nie przypomina paraboli.

W takim razie odrzucamy  $d$ , gdy nie leży w przedziale  $[a, c]$  lub rozmiar przedziału nie jest mniejszy niż połowa przedziału w przedostatniej iteracji. Jeśli odrzuciliśmy punkt, dokonujemy bisekcji na  $[a, c]$  (bisekcja, bo jest prostsza od złotego podziału, a daje podobny rezultat).

Zamiast metody Brenta, można skorzystać z pochodnych:  $d = \frac{af'_c - cf'_a}{f'_c - f'_a}$   
Jest to jednak na ogół mniej efektywne.

Jeszcze gorszą metodą jest wykorzystanie interpolacji Hermite'a: robimy wielomian 3-ego stopnia, który zgadza się z badaną funkcją i jako d bierzemy jego minimum.

---

Żeby znaleźć minima funkcji wielu zmiennych, wybieramy jakiś punkt startowy  $x_k$ , który jest aktualnym przybliżeniem minimum oraz kierunek  $p_k$ . Konstuuujemy funkcję:  $g(\alpha) = f(x_k + \alpha p_k)$ . Znajdujemy minimum funkcji  $g$ , która jest funkcją jednej zmiennej. Wtedy  $x_{k+1} = x_k + \alpha_{\min} p_k$ .

Powtarzamy, dopóki nie otrzymamy porządanej dokładności. Problem sprowadza się do odpowiedniego wyboru kierunku poszukiwań.

Kierunek powinniśmy wybierać inaczej, gdy jesteśmy blisko minimum, a inaczej gdy daleko.

**Minimalizacja po współrzędnych** – kolejnymi kierunkami są kierunki równoległe do kolejnych osi układu współrzędnych. Zły pomysł: za dużo kroków i są one niezwiązane z kształtem funkcji. To nie może być optymalne.

**Metoda najszybszego spadku** – kierujemy się minus gradientem funkcji. Jednak szukanie minimów kierunkowych jest dość kosztowne, więc kierujemy się tylko najszybszym spadkiem funkcji i tylko do okolic minimum. Używanie tej metody w okolicach minimum prowadzi do małych i licznych kroków.

---

Waunek na to, by  $f$  osiągała minimum kierunkowe (czyli minimum  $g$ ):

$$\frac{dg_k}{d\alpha} = \sum_i \frac{\partial f}{\partial x_i} (p_k)_i = \left( \nabla f|_{x=x_{\min}} \right)^T p_k = 0$$

Czyli w minimum, gradient jest prostopadły do kierunku poszukiwań (styczny do poziomicy).

Rozwinięcie minimalizowanej funkcji w szereg Taylora:  $f(x) \simeq \frac{1}{2} x^T A x - b^T x + f_0$ . A to hessjan obliczany w minimum. Macierz ta jest dodatnio określona (bo minimum) i może być symetryczna, gdy funkcja jest gładka. Czyli w pobliżu minimum, funkcja zachowuje się podobnie do formy kwadratowej.

**Gradienty sprzężone** – przy rozwinięciu w szereg Taylora, gradient funkcji w  $x_k$  wynosi:

$\nabla f|_{x=x_k} = A x_k - b$ . Zmiana gradientu w kierunku poszukiwań  $p_k$  wynosi:  $\delta(\nabla f) = \alpha A p_{k+1}$ . Żeby nie naruszyć żadnego z poprzednio wyliczonych minimów kierunkowych, zmiana gradientu musi być prostopadła do wszystkich poprzednich kierunków poszukiwań:  $p_i^T A p_j = 0, i > j$ .

**Algorytm**: Zaczynamy w  $x_1$ . Bierzemy  $r_1 = p_1 = -\nabla f|_{x_1}$ . Z punktu  $x_k$  minimalizujemy kierunkowo w kierunku  $p_k$ , otrzymując  $x_{k+1}$ . Obliczamy  $r_{k+1} = -\nabla f|_{x_{k+1}}$ . Liczymy:

W końcu:  $p_{k+1} = r_{k+1} + \beta p_k$ . Powtarzamy.

$$\beta = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

Wektor  $r_{k+1}$  jest tutaj tym samym wektorem, który otrzymalibyśmy w algebraicznej metodzie gradientów. Dzięki temu, nie potrzebujemy hessjanu do obliczeń, bo zastępujemy go minimalizacją kierunkową, a resztę kroków metodą algebraiczną. Alternatywna postać równania, która może przyspieszyć obliczenia, gdy grozi stagnacja:

$$\beta = \frac{r_{k+1}^T (r_{k+1} - r_k)}{r_k^T r_k}$$

**Metoda zmiennej metryki** – ponieważ szukając minimum, mamy rozwiązać równanie  $\nabla f = 0$ , to możemy użyć takiego rozwiązania (opartego na metodzie Newtona):

$\mathbf{x}_{i+1} = \mathbf{x}_i - \mathbf{H}_i^{-1} \nabla f|_{\mathbf{x}_i}$ , dla  $\mathbf{H}$  będącego macierzą drugich pochodnych cząstkowych w  $\mathbf{x}_i$ . Żeby dochodziło do zmniejszania się wartości funkcji, musi dodatkowo zachodzić:  $-(\mathbf{x}_{i+1} - \mathbf{x}_i)^T \mathbf{H}_i (\mathbf{x}_{i+1} - \mathbf{x}_i) < 0$  (jeśli  $\mathbf{H}$  jest dodatnio określone, mamy to spełnione automatycznie). Zamiast hessjanu, w tej metodzie używamy przybliżenia hessjanu, które jest zbliżone do tego prawdziwego, w minimum.

**Algorytm:**  $\mathbf{H}_0$  to jakaś macierz symetryczna, dodatnio określona. Rozwiązujemy  $\mathbf{H}_k \mathbf{p}_k = -\nabla f|_{\mathbf{x}_k}$ . Minimalizujemy kierunkowo  $f(\mathbf{x}_k + \alpha \mathbf{p}_k)$ , a  $\alpha_k$  to znalezione minimum kierunkowe.  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$ . Liczymy  $\mathbf{y}_k = \nabla f|_{\mathbf{x}_{k+1}} - \nabla f|_{\mathbf{x}_k}$ . Na koniec liczymy  $\mathbf{H}_{k+1}$ , np. ze wzoru BFGS:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\alpha_k \mathbf{p}_k^T \mathbf{y}_k} - \frac{\nabla f|_{\mathbf{x}_k} (\nabla f|_{\mathbf{x}_k})^T}{\mathbf{p}_k^T \mathbf{H}_k \mathbf{p}_k}$$

Kończymy, gdy  $|\mathbf{p}_k| < \varepsilon$  (zadanej tolerancji).

W tej metodzie, jeśli  $\mathbf{H}_0$  jest symetryczne i dodatnio określone, to takie pozostaje. Metoda sprawdza się dla dużej liczby wymiarów.

Dwie ostatnie metody są dobre, gdy jesteśmy blisko minimum. Gdy zaczynamy daleko, to są wolne. Dzieje się tak dlatego, że daleko do minimum, funkcja nie jest podobna do formy kwadratowej.

**Metoda Levenberga-Marquardta** – modyfikujemy hessjan (żeby był dodatnio określony również z dala od minimum):

$$\begin{aligned} \tilde{\mathbf{H}}_{ii} &= (1 + \lambda) \frac{\partial^2 f}{\partial x_i^2}, \\ \tilde{\mathbf{H}}_{ij} &= \frac{\partial^2 f}{\partial x_i \partial x_j}, \quad i \neq j \end{aligned}$$

$\lambda \geq 0$ . Zaczynamy z jakimś niewielkim  $\lambda$ , np.  $2^{-10}$ .  $\mathbf{x}_i$  to aktualne przybliżenie. Liczymy po kolei:  $\nabla f(\mathbf{x}_i)$ ,  $\tilde{\mathbf{H}}(\mathbf{x}_i)$ .  $\mathbf{x}_{\text{test}} = \mathbf{x}_i - \tilde{\mathbf{H}}^{-1}(\mathbf{x}_i) \nabla f(\mathbf{x}_i)$ .  
Jeżeli  $f(\mathbf{x}_{\text{test}}) > f(\mathbf{x}_i)$ , to zwiększamy  $\lambda$  8 (lub nie 8) razy i liczymy znowu hessjan.  
Inaczej, zmniejszamy  $\lambda$  8 razy,  $\mathbf{x}_{i+1} = \mathbf{x}_{\text{test}}$  i liczymy znowu gradient.

Daleko od minimum znowu nie minimalizujemy, a podążamy za spadkiem. Jeśli  $\lambda$  jest zbyt duże (ustalamy jakieś maksymalne), to znaczy, że jesteśmy poza basenem atrakcji i algorytm zawodzi. Jeśli  $\lambda$  jest bardzo małe, możemy przerzucić się na gradienty sprzężone albo zmienną metrykę. Gdy wybierzemy to drugie, aktualne przybliżenie hessjanu może służyć za początkowe w tamtej metodzie.

Jeśli chcemy przyspieszyć algorytm, możemy zrezygnować na kilka kroków ze zmniejszania  $\lambda$  i przeliczania  $\mathbf{H}$  na spadku. Daleko od minimum, ten algorytm zachowuje się jak metoda najszybszego spadku.

**Metoda Powella** – dobra, gdy nie możemy policzyć pochodnych (a więc i gradientów). Zaczynamy od minimalizacji po współrzędnych:  $\{\mathbf{p}_i\}_{i=1}^N = \{\mathbf{e}_i\}_{i=1}^N$ . Znajdujemy się w jakimś  $\mathbf{x}_0$ , minimalizujemy wzdłuż kierunków  $\mathbf{p}_i$ , osiągając kolejne  $\mathbf{x}_i$ . Od  $i$  do  $N-1$ ,  $\mathbf{p}_i = \mathbf{p}_{i+1}$ .  $\mathbf{p}_N = \mathbf{x}_N - \mathbf{x}_0$ . Minimalizujemy wzdłuż  $\mathbf{p}_N$  i nowy punkt oznaczamy jako  $\mathbf{x}_0$ , i robimy następną iterację. Jeśli funkcja jest rozwijalna w szereg Taylora,  $\mathbf{p}_i$  stają się sprzężone.

Wszystkie powyższe algorytmy są nieprzydatne dla minimalizacji globalnej. Nie sprawdzają się dla niej również metody zachłanne (sprawdzanie dużej ilości punktów) oraz losujące punkty do sprawdzania, szczególnie w wielu wymiarach są zbyt złożone.

Algorytmy do szukania minimów globalnych idą zawsze w dół, więc jak wejdą do basenu atrakcji minimum, to mogą nie wyjść (a nie musi to być minimum globalne). Więc robimy tak, że losowo możemy zrobić "zły" krok pod górkę (wyższe kroki są rzadsze).

**Algorytm Monte Carlo** – mamy funkcję  $f: \mathbb{R}^N \rightarrow \mathbb{R}$ . Startujemy z jakiegoś  $\mathbf{x}_0$ , liczymy  $f_0 = f(\mathbf{x}_0)$ .  $\mathbf{x}_{\min} = \mathbf{x}_0$ ,  $f_{\min} = f_0$ .  $\sigma, T > 0$  są parametrami. Zaczynamy iterację:  $\tilde{\mathbf{x}} = \mathbf{x}_k + \sigma \xi_k$  ( $\xi$  to  $N$  wymiarowa zmienna losowa gaussowska).  $\tilde{f} = f(\tilde{\mathbf{x}})$ . Jeżeli  $\tilde{f} < f_k$ , to akceptujemy położenie. Jeżeli  $\tilde{f} < f_{\min}$ , to  $\mathbf{x}_{\min} = \tilde{\mathbf{x}}$  oraz  $f_{\min} = \tilde{f}$ . Jeśli  $\tilde{f} > f_k$ , to losujemy  $z$  o rozkładzie jednorodnym i jeśli  $z < \exp(-(\tilde{f} - f_k)/T)$ , to akceptujemy nowe położenie. //  $\exp(-(f - f_k)/T)$  to kryterium boltzmannowskie

Nie ma kryterium stopu, zadajemy  $M$  iteracji i tyle robimy. Zamiast brać duże  $M$ , kilkakrotnie losujemy punkt startowy i robimy niewielką liczbę iteracji. Nie ma gwarancji, że wylosowane punkty same znajdują się w basenach atrakcji.

Kombinacja Gauss-Boltzmann powoduje, że im większ paramter T, tym łatwiej wydostać się funkcji z basenu atrakcji (czas na to potrzebny to  $\sim \exp(\text{wysokość bariery} / T)$ ). Parametry dobieramy np. metodą prób i błędów, T nie powinno być za wysokie. Metoda działa również dla problemów dyskretnych.

Algorytm ten ma charakter lokalny, bo nie wprowadza znacznych zmian w sposób sensowny i może być wolny.

**Algorytmy genetyczne** – opierają się na zasadzie survival of the fittest. Mamy jakąś populację osobników (dużo mniej niż możliwych stanów) i każdy reprezentowany jest przez chromosom. Dla każdego chromosomu można wyliczyć wartość funkcji dopasowania. Każdemu z nich przyporządkowujemy wycinek koła, proporcjonalny do jego wartości funkcji:  $p_i = f_i / \sum_j f_j$  (tzw. metoda ruletki).

Wybieramy następnie dwoje "rodziców", którzy stworzą chromosom do następnej iteracji. Dzielimy chromosomy rodziców (najczęściej po połowie) i tworzymy nowy chromosom. Następnie z niewielkim prawdopodobieństwem, możemy zmienić losową z wartości w chromosomie (zmutować). Tworzymy tyle potomków ile ma wynosić cała populacja, a operację tą można zrównoleglić.

Populację startową wybieramy losowo z puli dostępnych stanów i zadajemy jakąś maksymalną ilość epok (pokoleń). Jeśli wyczerpią się epoki, zmiany są zbyt małe, albo populacja jest w pobliżu jakiegoś stanu, przerywamy. Zapamiętujemy chromosom, który miał najlepszą wartość funkcji dopasowania. **//np. możemy mieć dodatkia funkcję i jej maksimum będzie najlepsze**

**Particle Swarm Optimization** – mamy pewne niezbyt liczne stado. Każdemu z osobników przyporządkowujemy położenie  $x_i$  oraz prędkość  $v_i$ . Początkowe położenia i prędkości są losowe (prędkości nie za duże). Dla każdego członka stada zapamiętujemy najmniejszą wartość funkcji dopasowania  $f_{\min, i}$  oraz odpowiadające jej położenie  $x_{\min, i}$ . Zapamiętujemy też globalną najmniejszą wartość funkcji dopasowania  $f_{\text{glob}}$  oraz jej położenie.  $0 < \omega < 1$ ;  $a, b > 0$  są parametrami.

Stado eksploruje przestrzeń. Każdy członek ma tendencję do lecenia jak wcześniej, orientowania się w stronę swojej najlepszej wartości oraz orientowania w stronę najlepszej globalnej. Sekwencyjnie dla wszystkich wykonujemy: losujemy liczby  $p, q$  na przedziale  $[0, 1]$ ;  $v_{i+1} = \omega v_i + a \cdot p \cdot (x_{\min, i} - x_i) + b \cdot q \cdot (x_{\text{glob}} - x_i)$ ;  $x_{i+1} = x_i + v_{i+1}$ ;  $f_{i+1} = f(x_{i+1})$ ; aktualizujemy wartości  $f, x_{\min, i}$  oraz  $f, x_{\text{glob}}$  na podstawie  $f, x_{i+1}$ .

Tendencja do podążania w stronę minimum znalezionej przez danego członka odpowiada za podążanie w dół, a w stronę globalnego, za inteligencję roju. Wadą algorytmu jest, że nie zwraca uwagi na kształt i zmienność funkcji, ale zaletą jest, że nie trzeba liczyć gradientów. Parametry można dobierać na podstawie wyników minimalizacji podobnych, ale prostszych problemów (metaoptymalizacja). Jeśli  $\omega > 1$ , to prędkości wybuchają. Żeby uniknąć sytuacji, gdy stado skupia się wokół jednego minimum (zdarza się dla wielu wymiarów), robimy kilka niezależnych stad.

Aproksymacja może być: punktowa – szukamy znanej funkcji oraz jej parametrów, która będzie przebiegać możliwie najbliżej danych punktów, oraz może być ciągła – mamy jakąś funkcję, której obliczanie jest trudne i szukamy podobnej, ale łatwiejszej.

Interpolacja to skrajny przypadek aproksymacji punktowej. W aproksymacji chcemy mieć prostą funkcję, która nie musi przechodzić przez wszystkie punkty. Zależność między  $x$  i  $y$  powinna być  $y = y(x)$ , ale problem jest z błędami pomiaru.

**Liniowe zagadnienie najmniejszych kwadratów** – każdej zmierzonej wartości  $y_i$  odpowiada wartość teoretyczna  $\tilde{y}_i$ , którą zmienna  $y$  "powinna" przybrać dla  $x$ . Przypuśćmy, że  $\tilde{y}$  jest kombinacją liniową znanych funkcji:  $\tilde{y}_i = a_1 \cdot f_1(x_1) + \dots + a_s \cdot f_s(x_s)$ . Zespół wszystkich teoretycznych wartości możemy przedstawić jako  $\tilde{\mathbf{y}} = \mathbf{A}\mathbf{p}$ , gdzie  $\mathbf{A}$  jest macierzą funkcji  $f_i x_j$ , z argumentami w rzędach i funkcjami w kolumnach, a  $\mathbf{p}$  to wektor współczynników. Problem, który chcemy rozwiązać to znalezienie najlepszego wektora  $\mathbf{p}$ .

Przykładowa macierz dla  $n$  punktów pomiarowych i wielomianu 2 stopnia, który chcemy dopasować:

$$\tilde{\mathbf{y}} = \begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \dots & \dots & \dots \\ x_n^2 & x_n & 1 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

Dopasowanie danych do wielomianu ustalonego stopnia jest zagadnieniem liniowym. Różnica między wartością teoretyczną, a pomiarową jest spowodowana błędem pomiaru:  $y_i - \tilde{y}_i = \xi_i$ . Przypuśćmy, że liczby  $\xi$  są zmiennymi losowymi z rozkładu Gaussa. Wektor wszystkich błędów pomiarowych oznaczamy  $\xi = [\xi_1, \dots, \xi_n]^T$ . Przyjmijmy też, że łącznie wszystkie błędy tworzą rozkład Gaussa o macierzy kowariancji  $G = \langle \xi \xi^T \rangle$ , gdzie  $\langle \dots \rangle$  oznacza średniowanie po realizacjach zmiennych losowych.  $G$  jest symetryczna i dodatnio określona. Dla takiej macierzy, estymator największej wiarygodności wynosi:  $Q = \frac{1}{2} \xi^T G^{-1} \xi$  i z pewnością posiada minimum. Obecność odwrótności  $G$  oznacza, że pomiary z większym błędem dają mniejszy wkład do  $Q$ .

W tym zagadnieniu, minimalizowana funkcja jest formą kwadratową w parametrach (zależy liniowo od parametrów, nie argumentów) i wiemy, że posiada jednoznaczne minimum. Aby znaleźć estymator, trzeba znaleźć taki wektor  $p$ , że forma kwadratowa  $Q$  osiąga minimum. Można to zrobić metodą zmiennej metryki albo gradientów sprzężonych, lub formalnie rozwiązując  $\nabla Q = 0$ . Otrzymujemy:  $A^T G^{-1} A p = A^T G^{-1} y$ . Członu  $A^T G^{-1}$  nie da się uprościć, bo to nie jest macierz kwadratowa. Samo równanie jest dobrze określone, bo  $A^T G^{-1} A$  jest symetryczna i dodatnio określona.

$y = A p^* + \xi$ , gdzie  $p^*$  jest zbudowany z "prawdziwych" współczynników  $a_i$ . Wtedy:  $A^T G^{-1} A (p - p^*) = A^T G^{-1} \xi$ , czyli  $\langle p \rangle = p^*$ ,  $\langle \xi \rangle = 0$ . Obliczone estymatory są przybliżeniem "prawdziwych" wartości parametrów. Macierz kowariancji estymatorów wynosi:  $C_p = \langle (p - p^*)(p - p^*)^T \rangle$ , co po przekształceniach daje:  $\langle A^T G^{-1} A \rangle^{-1}$ .

Możemy rozwiązać, aby  $y_i = \tilde{y}_i$  było ściśle spełnione, czyli musimy rozwiązać układ równań  $A p = y$ . Jest on nadokreślony i nie ma ścisłego rozwiązania, ale możemy znaleźć optymalne, w sensie najmniejszych kwadratów, rozwiązanie za pomocą SVD. Jeśli pomiary są nieskorelowane, rozwiązanie uzyskane w ten sposób jest równoważne minimalizacji formy kwadratowej  $Q$  lub rozwiązaniu  $\nabla Q = 0$ . Nadokreślony układ  $G^{-1} A p = G^{-1} y$  również daje równoważne rozwiązanie z SVD (przy uwzględnieniu wag pomiarów).

Jeśli pomiary są nieskorelowane,  $Q$  upraszcza się do:

$$Q = \sum_{i=1}^n \frac{(a_1 f_1(x_i) + a_2 f_2(x_i) + \dots + a_s f_s(x_i) - y_i)^2}{\sigma_i^2}$$

A macierz kowariancji estymatorów upraszcza się do:  $C_p = \sigma^2 (A^T A)^{-1}$ .

**Kryterium Akaike** – pozwala zbalansować najlepsze dopasowanie z najmniejszą liczbą parametrów.

$AIC = \ln Q + 2s / N$ , należy to zminimalizować, dla  $s$  będącego liczbą parametrów, a  $N$  liczbą punktów. // **AIC to skrót od Akaike Information Criterion**

Kryterium to nagradza za lepsze dopasowanie i karze za dużą liczbę parametrów.

**Nieliniowe zagadnienie najmniejszych kwadratów** – kiedy dopasowana zależność teoretyczna zależy od danych nieliniowo:  $\tilde{y}_i = f(x_i; p)$ , dla  $p$  będącego wektorem parametrów.  $f$  jest znaną funkcją o nieznanym parametrach, np.

$$y(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \bar{x})^2}{2\sigma^2}\right)$$

//parametrami są  $\bar{x}$  oraz  $\sigma^2$ , widać, że zależność jest nieliniowa

$Q$  do minimalizacji jest podobne, ale  $\xi_i = y_i - f(x_i; p)$ . Dla najczęstszego przypadku nieskorelowanych danych,  $Q$  redukuje się do:  $Q = \text{const } \frac{1}{2} \sum_i (y_i - f(x_i; p))^2$ . Jednak w nieliniowym zagadnieniu,  $Q$  nie jest formą kwadratową w parametrach.  $Q$  jest dodatnio określona, więc posiada minimum, ale niekoniecznie jedno. Należy je znaleźć numerycznie, metodą Levenberga-Marquardta, chyba, że  $\nabla_p Q = 0$  można wyjątkowo łatwo rozwiązać.

**Pseudolinearyzacja** – stosuje się ją czasami do minimalizacji  $Q$ . Jeśli  $p_n$  jest aktualnym przybliżeniem  $p$ , zakładamy, że prawdziwe wartości są małą poprawką  $\delta p$  w stosunku do  $p_n$ .  $y_i = f(x_i; p_n + \delta p)$  rozwijamy w szereg Taylora, otrzymując  $f(x_i; p_n) + [\nabla_p f | p_n]^T \delta p$ . Funkcja:  $Q = \frac{1}{2} \sum_i (y_i - f(x_i; p_n) - [\nabla_p f | p_n]^T \delta p)^2$  jest formą kwadratową. Po znalezieniu minimum  $\delta p$ , podstawiamy:  $p_{n+1} = p_n + \delta p_{\min}$  i powtarzamy.

Działa ona dość dobrze dla nieliniowej metody najmniejszych kwadratów, ale nie należy jej polecać jako ogólnej metody minimalizacji. Niewątpliwym jej zastosowaniem jest, że po znalezieniu ostatecznych wartości minimalizujących  $Q$ , za pomocą pseudolinearyzacji wokół tego punktu, znajdujemy macierz kowariancji estymatorów będącą charakterystyką liniową.



**Przybliżenie Padé** – jedno z zagadnień aproksymacji ciągłej. Przypuśćmy, że znamy przybliżenie  $g(z)$  oraz jej pochodnych do rzędu  $N$  w  $z = 0$  (jeśli w innym punkcie, to przesuwamy). Chcemy skonstruować jej przybliżenie w pewnym przedziale, by zgadzało się z nią oraz z jej pochodnymi w 0. Najprościej jest skonstruować rozwinięcie Maclaurina, ale ono jest wielomianem i załamuje się w niewielkiej nawet odległości od zera.

Lepiej użyć przybliżeń wymiernych:

$$R_{mk}(x) = \frac{P_m(x)}{Q_k(x)} = \frac{\sum_{j=0}^m a_j x^j}{\sum_{j=0}^k b_j x^j}$$

Dla  $b_0 = 0$ , a szeregiem Maclaurina aproksymowanej funkcji będzie:

$$g(x) = \sum_{j=0}^{\infty} c_j x^j$$

Przyjmijmy, że  $m + k + 1 = N + 1$ , czyli tyle, ile wyrazów zawiera szereg Maclaurina funkcji  $g(x)$  do rzędu  $N$ .

Liczmy:

$$g(x) - R_{mk}(x) = \frac{\left(\sum_{j=0}^{\infty} c_j x^j\right) \left(\sum_{j=0}^k b_j x^j\right) - \sum_{j=0}^m a_j x^j}{\sum_{j=0}^k b_j x^j}$$

$g(x)$  będzie się zgadzać z przybliżeniem wyżej, jeśli najniższy nieznikający wyraz sumy w liczniku będzie się zgadzać z  $X^{N+1}$ .

Otrzymujemy warunki, które są układem  $N+1$  równań liniowych na  $N+1$  współczynników  $a_r, b_j$ :

$$\begin{aligned} \sum_{j=0}^k c_{N-s-j} b_j &= 0 \quad s = 0, 1, \dots, N - m - 1; \quad c_j = 0 \text{ dla } j < 0 \\ \sum_{j=0}^r c_{r-j} b_j &= a_r \quad r = 0, 1, \dots, m; \quad b_j = 0 \text{ dla } j > k \end{aligned}$$

**Przybliżenie Czebyszewa** – zamiast  $R_{mk}$ , lepszym przybliżeniem są wielomiany Czebyszewa:

$$C_{mk}(x) = \frac{\sum_{j=0}^m a_j T_j(x)}{\sum_{j=0}^k b_j T_j(x)} \quad // T_j(x) = \cos(j \arccos x) \text{ dla } x \in (-1, 1)$$

Po analitycznym przedłużeniu,  $T_j(x)$  jest poza przedziałem  $(-1, 1)$  wielomianem stopnia  $j$ . Jeśli  $g(x)$  posiada rozwinięcie Czebyszewa:

$$g(x) = \frac{1}{2} c_0 + \sum_{j=1}^{\infty} c_j T_j(x)$$

To jego współczynniki obliczamy tak:

$$c_j = \frac{2}{\pi} \int_{-1}^1 \frac{g(x) T_j(x)}{\sqrt{1-x^2}} dx$$

Liczy się to kwadraturami Gaussa-Czebyszewa, których nie trzeba znać, bo nie były omówione.