

Tema nr. 8

Fie $F : \mathbb{R}^n \longrightarrow \mathbb{R}$ o funcție reală, $F(x) = F(x_1, x_2, \dots, x_n)$. Să se aproximeze un punct de minim (local sau global) al funcției F folosind metoda gradientului descendent. Să se testeze diversele metode de calcul a ratei de învățare descrise în acest text. Să se calculeze gradientul funcției F folosind formula analitică și formula aproximativă. Să se compare soluțiile obținute folosind cele două moduri de calcul a gradientului funcției F , din punctul de vedere al numărului de iterații efectuate pentru obținerea soluțiilor (pentru aceeași precizie $\epsilon > 0$). Testați toate funcțiile listate la sfârșitul acestui text.

Bonus: Aplicați algoritmi de minimizare pentru rezolvarea problemei de regresia logistică descrisă în fișierul [RL](#) (pag. 65).

Obțineți **10pt** dacă folosiți funcția de log-verosimilitate și gradientul ei deja calculate.

Obțineți **15pt** dacă construiți funcția de log-verosimilitate și gradientul ei pornind de la tabelul cu date.

$$\max\{F(x); x \in V\} = -\min\{-F(x); x \in V\} \quad ,$$

$$\operatorname{argmax}\{F(x); x \in V\} = \operatorname{argmin}\{-F(x); x \in V\}$$

Minimizarea funcțiilor

Fie $F : \mathbb{R}^n \longrightarrow \mathbb{R}$ o funcție reală de două ori derivabilă, $F \in C^2(\mathbb{R}^n)$, pentru care vrem să aproximăm soluția x^* a problemei de minimizare:

$$\min\{F(x); x \in V\} \iff F(x^*) \leq F(x) \quad \forall x \in V \quad (1)$$

unde $V = \mathbb{R}^n$ (x^* este punct de minim global) sau $V = S(\bar{x}, r)$, sfera de centru \bar{x} și rază r (punct de minim local). Se numește *punct critic* pentru funcția F , un punct \tilde{x} care este rădăcină a sistemului de ecuații:

$$\nabla F(\tilde{x}) = 0 \quad , \quad \nabla F(x) = \begin{pmatrix} \frac{\partial F}{\partial x_1}(x) \\ \vdots \\ \frac{\partial F}{\partial x_i}(x) \\ \vdots \\ \frac{\partial F}{\partial x_n}(x) \end{pmatrix} . \quad (2)$$

Se știe că pentru funcțiile de două ori derivabile, punctele de minim ale funcției F se găsesc printre punctele critice. Un punct critic este punct de minim dacă matricea hessiană este pozitiv semidefinită:

$$H(\tilde{x}) = \left(\frac{\partial^2 F}{\partial x_i \partial x_j}(\tilde{x}) \right)_{i,j=1,\dots,n}, \quad (H(\tilde{x})z, z)_{\mathbb{R}^n} \geq 0 \quad \forall z \in \mathbb{R}^n.$$

Metoda gradientului descendent

Punctul de minim al funcției F se aproximează construind un șir de vectori $\{x^k\} \subseteq \mathbb{R}^n$ care, în anumite condiții, converge la punctul de minim x^* căutat. Convergența șirului depinde de alegerea primului element al șirului, x^0 .

Elementul $k + 1$ al șirului, x^{k+1} , se construiește pornind de la elementul precedent, x^k , astfel:

$$x^{k+1} = x^k - \eta_k \nabla F(x^k), \quad k = 0, 1, \dots, \quad x^0 - \text{ales random} \quad (3)$$

Elementul η_k poartă numele de rată de învățare sau pasul iterației.

Strategii de alegere a ratei de învățare

1. $\eta_k = \eta$, $\forall k$ ($\eta = 10^{-3}, 10^{-4}, \dots$). O rată de învățare constantă prea mare poate face ca punctul de minim să nu poată fi găsit, iar o valoare prea mică pentru rata de învățare are dezavantajul unui cost de calcul mare.
2. Un mod de a rezolva problemele care apar în cazul ratei de învățare constante este de a considera o valoare variabilă, în funcție de contextul local. Metoda descrisă mai jos poartă denumirea de ajustare de tip *backtracking* a lungimii pasului/ratei de învățare (*backtracking line search*). Această metodă funcționează pentru funcții convexe.

Se considera $\beta \in (0, 1)$ fixat (de obicei se alege $\beta = 0.8$). La fiecare pas rata de învățare se calculează astfel:

$$\begin{aligned} \eta &= 1; \\ p &= 1; \\ \text{while } F(x^k - \eta \nabla F(x^k)) &> F(x^k) - \frac{\eta}{2} \|\nabla F(x^k)\|^2 \ \&\& \ p < 8 \\ \eta &= \eta \beta; \\ p &++; \end{aligned}$$

Observație importantă: Alegerea vectorului inițial, x^0 , poate determina convergența sau divergența șirului $\{x^k\}$ la punctul de minim x^* . De obicei, o alegere a datelor inițiale în vecinătatea lui x^* asigură convergența $x^k \rightarrow x^*$ pentru $k \rightarrow \infty$.

Nu este necesară memorarea întregului șir $\{x^k\}$ ci avem nevoie doar de 'ultimul' element x^{k_0} calculat. Se consideră că un element x^{k_0} aproximează punctul de minim căutat, x^* , $x^{k_0} \approx x^*$ (x^{k_0} este ultimul element al șirului care se calculează) atunci când diferența dintre două elemente succesive ale șirului devine suficient de mică, i.e.,

$$\|x^{k_0} - x^{k_0-1}\| \leq \epsilon \quad (4)$$

unde ϵ este precizia cu care vrem să aproximăm soluția x^* .

Prin urmare, o schemă posibilă de aproximare a soluției x^* este următoarea:

Schema de calcul

```
se aleg random elementele vectorului inițial  $x$  ;  
 $k = 0$  ;  
do  
  {  
    - calculează  $\nabla F(x)$  ;  
    - calculează rata de învățare  $\eta$  folosind  
      una din cele 2 metode;  
    -  $x = x - \eta \nabla F(x)$  ;  
    -  $k = k + 1$ ;  
  }  
while (  $\eta \|\nabla F(x)\| \geq \epsilon$  și  $k \leq k_{\max}$  și  
         $\eta \|\nabla F(x)\| \leq 10^{10}$  )  
if (  $\eta \|\nabla F(x)\| \leq \epsilon$  )  $x \approx x^*$  ;  
else "divergență" ; //(de încercat schimbarea datelor  
                        inițiale)
```

O valoare posibilă pentru k_{\max} este 30000 și $\epsilon > 10^{-5}$.

Pentru a calcula valoarea gradientului funcției F într-un punct oarecare se va folosi formula analitică de calcul a gradientului (funcție declarată în program) și de asemenea se va folosi următoarea formulă aproximativă:

$$\nabla F(x) \approx (G_i(x, h))_{i=1, \dots, n} \quad , \quad \frac{\partial F}{\partial x_i}(x) \approx G_i(x, h)$$

unde

$$\frac{\partial F}{\partial x_i}(x) \approx G_i(x, h) = \frac{-F_1 + 8F_2 - 8F_3 + F_4}{12h} \quad , \quad \forall i = 1, \dots, n$$

cu $h = 10^{-5}$ sau 10^{-6} (poate fi considerat parametru de intrare) și:

$$\begin{aligned} F_1 &= F(x_1, \dots, x_{i-1}, x_i + 2h, x_{i+1}, \dots, x_n), \\ F_2 &= F(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n), \\ F_3 &= F(x_1, \dots, x_{i-1}, x_i - h, x_{i+1}, \dots, x_n), \\ F_4 &= F(x_1, \dots, x_{i-1}, x_i - 2h, x_{i+1}, \dots, x_n). \end{aligned}$$

Exemple

$$l(w_0, w_1) = -\ln(1 + \exp(w_0 - w_1)) + w_0 + w_1 - \ln(1 + \exp(w_0 + w_1))$$

$$\nabla l(w_0, w_1) = \begin{pmatrix} -\sigma(w_0 - w_1) + \sigma(-w_0 - w_1) \\ \sigma(w_0 - w_1) + \sigma(-w_0 - w_1) \end{pmatrix}, \quad \sigma(z) = \frac{1}{1 + \exp(-z)}.$$

$$F(x_1, x_2) = x_1^2 + x_2^2 - 2x_1 - 4x_2 - 1, \nabla F(x_1, x_2) = \begin{pmatrix} 2x_1 - 2 \\ 2x_2 - 4 \end{pmatrix}, x_1^* = 1, x_2^* = 2.$$

$$F(x_1, x_2) = 3x_1^2 - 12x_1 + 2x_2^2 + 16x_2 - 10, \nabla F(x_1, x_2) = \begin{pmatrix} 6x_1 - 12 \\ 4x_2 + 16 \end{pmatrix}, x_1^* = 2, x_2^* = -4.$$

$$F(x_1, x_2) = x_1^2 - 4x_1x_2 + 5x_2^2 - 4x_2 + 3, \nabla F(x_1, x_2) = \begin{pmatrix} 2x_1 - 4x_2 \\ -4x_1 + 10x_2 - 4 \end{pmatrix}, x_1^* = 4, x_2^* = 2.$$

$$F(x_1, x_2) = x_1^2x_2 - 2x_1x_2^2 + 3x_1x_2 + 4, \nabla F(x_1, x_2) = \begin{pmatrix} 2x_1x_2 - 2x_2^2 + 3x_2 \\ x_1^2 - 4x_1x_2 + 3x_1 \end{pmatrix}, x_1^* = -1, x_2^* = 0.5.$$