

Optimalizáló algoritmusok

Baran Ágnes

Előadás

Sztochasztikus optimalizálás

A feladat:

Keresett

$$\min_{x \in \Omega} f(x),$$

ahol $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\Omega \subset \mathbb{R}^n$.

Sztochasztikus optimalizálásról beszélünk, ha

- az $f(x)$ függvény csak zajjal terheltén figyelhető meg, azaz csak $y(x) = f(x) + \varepsilon(x)$ ismert, ahol ε egy véletlen zaj,
- a keresési irányokat véletlenszerűen választjuk.

Determinisztikus esetben az f teljesen ismert és a keresési irányokat determinisztikusan választjuk.

Sztochasztikus optimalizálás. Egy egyszerű megoldás.

egyszerű elvétel



Ha Ω korlátos, akkor generáljunk $u_1, \dots, u_m \sim U(\Omega)$ értékeket.

Keresett $f_m^* = \min\{f(u_1), \dots, f(u_m)\}$.

Ez tart $f^* = f(x^*)$ -hoz, de nagyon lassan.

Példa. Keresett a következő függvény minimuma:

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f(x) = -(\cos(10x) - \sin(60x))^2, \quad x \in [0, 1]$$

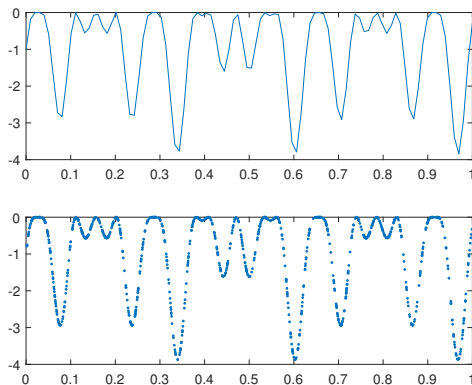
A valódi minimum: $f^* = \underline{-3.868428}$ az 0.3396384 és 0.6028403 helyeken.

Közelítések:

m	100	500	1000	5000
f_m^*	-3.8303	-3.8667	-3.8680	-3.8684

Példa.

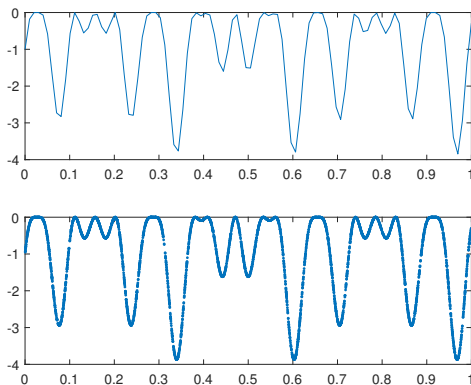
$$f(x) = -(\cos(10x) - \sin(60x))^2, \quad x \in [0, 1]$$



Az $f(x)$ és $f(u)$ függvények, ahol $m = 1000$

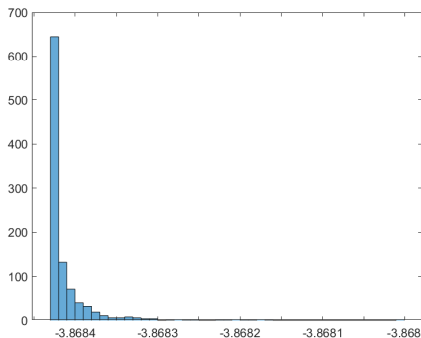
Példa.

$$f(x) = -(\cos(10x) - \sin(60x))^2, \quad x \in [0, 1]$$



Az $f(x)$ és $f(u)$ függvények, ahol $m = 5000$

$$f(x) = -(\cos(10x) - \sin(60x))^2, \quad x \in [0, 1]$$



1000 tesztfutás eredménye ($m = 5000$). A valódi minimum -3.868428 .

Az algoritmus:

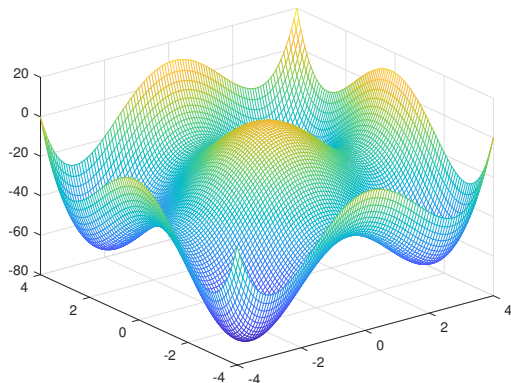
- Válasszunk egy $x_0 \in \Omega$ kezdőpontot (véletlenszerűen vagy determinisztikusan)
- $x_k \rightarrow x_{k+1}$: generáljunk egy $x_n \in \Omega$ véletlen, független pontot egy adott eloszlás segítségével. Ha $f(x_n) < f(x_k)$, akkor legyen $x_{k+1} = x_n$, egyébként $x_{k+1} = x_k$.
- Leállás: ha elértük a maximális iterációszámot (vagy elégedettek vagyunk az aktuális közelítőértékkel)

Az algoritmus elegendően általános feltételek mellett majdnem biztosan konvergál x^* -hoz, de csak alacsony dimenzióban hatékony.

Vak: az új pont generálásánál teljesen figyelmen kívül hagyja az x^* korábbi közelítéseit.

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

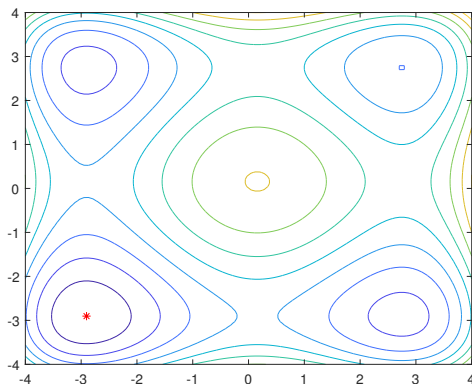
$$f(x) = (x_1^4 - 16x_1^2 + 5x_1)/2 + (x_2^4 - 16x_2^2 + 5x_2)/2;$$



$$[-4, 4] \times [-4, 4]$$

$$x^* = [-2.9035, -2.9035]$$

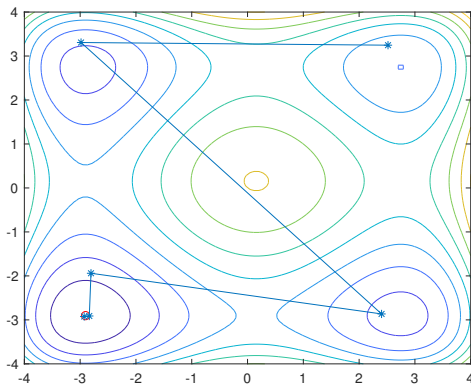
$$f(x) = (x_1^4 - 16x_1^2 + 5x_1)/2 + (x_2^4 - 16x_2^2 + 5x_2)/2;$$



$$x^* = [-2.9035, -2.9035]$$

Tesztfutások

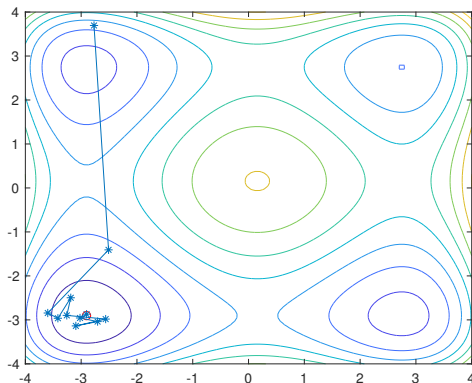
$$f(x) = (x_1^4 - 16x_1^2 + 5x_1)/2 + (x_2^4 - 16x_2^2 + 5x_2)/2;$$



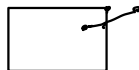
(5000 iteráció)
az elfogadható
pontos
néhány

$$x^* = [-2.9035, -2.9035], x_{opt} = [-2.9276, -2.9233] \boxed{k=6}$$

$$f(x) = (x_1^4 - 16x_1^2 + 5x_1)/2 + (x_2^4 - 16x_2^2 + 5x_2)/2;$$



$$x^* = [-2.9035, -2.9035], x_{opt} = [-2.9028, -2.8707], k = 11$$



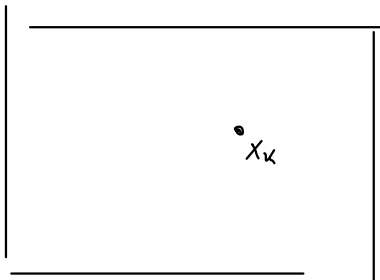
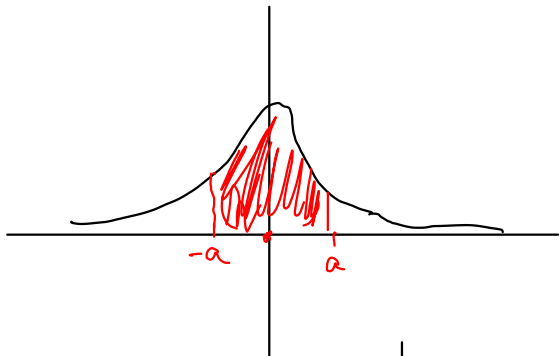
Az új pontot az eddigi legjobb közelítés egy környezetéből választjuk.

Az algoritmus:

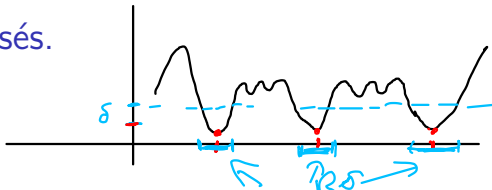
- Válasszunk egy $x_0 \in \Omega$ kezdőpontot (véletlenszerűen vagy determinisztikusan)
- $x_k \rightarrow x_{k+1}$: generáljunk egy véletlen irányt $d_k \in \mathbb{R}^n$, ha $x_k + d_k \notin \Omega$ akkor generáljunk egy új d_k irányt, ezt ismételjük (vagy válasszuk az Ω legközelebbi pontját). Ha $x_k + d_k \in \Omega$ akkor legyen $\underline{x}_n = x_k + d_k$.
- Ha $f(x_n) < f(x_k)$, akkor legyen $x_{k+1} = x_n$, egyébként $x_{k+1} = x_k$.
- Leállás: ha elértük a maximális iterációszámot (vagy elégedettek vagyunk az aktuális közelítőértékkel)

A d_k irány generálásánál egy 0 várható értékű normális eloszlást használunk. A szórásokat válasszuk az aktuális x_k komponenseinek nagyságával összhangban.

Az 1 dim normális eloszlás sűrűségf.
 $\hat{\sigma}^2$ értékét értékelni)



Lokalizált véletlen keresés.



Tétel. Legyen Ω^* az f függvény Ω feletti globális minimumhelyeinek a halmaza. Tegyük fel, hogy f folytonos, és ha $x_k + d_k \notin \Omega$ akkor az új d_k irányt véletlenszerűen választjuk. Adott $\delta > 0$ esetén legyen

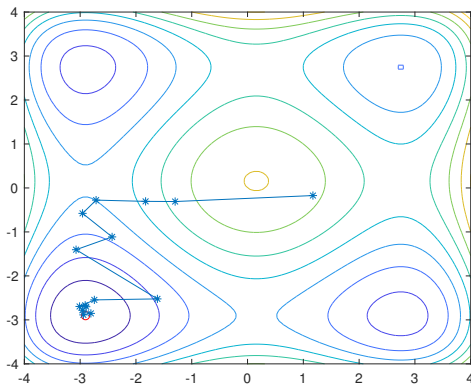
$$R_\delta = \bigcup_{x \in \Omega^*} \{x \in \Omega : |f(x) - f(x^*)| < \delta\}.$$

Ekkor ha a d_k sorozat elemei független, azonos $\mathcal{N}(0, I)$ eloszlásúak, akkor

$$\lim_{k \rightarrow \infty} P(x_k \in R_\delta) = 1.$$

Tesztfutások

$$f(x) = (x_1^4 - 16x_1^2 + 5x_1)/2 + (x_2^4 - 16x_2^2 + 5x_2)/2;$$



↙ eljósodott
pontos
név

$$x^* = [-2.9035, -2.9035], x_{opt} = [-2.9416, -2.8864], k = 15$$

Gradiens módszer

Determinisztikus esetben: $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$

Ha f nem differenciálható, vagy az f -et nem ismerjük analitikus alakban, akkor

$$x_{k+1} = x_k - \underline{\alpha_k} \hat{\nabla} f(x_k).$$

Legyen most

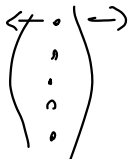
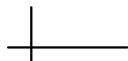
$$\hat{\nabla} f(x_k) = \begin{pmatrix} \frac{f(x_k + \beta_k \xi_1) - f(x_k - \beta_k \xi_1)}{2\beta_k} \\ \vdots \\ \frac{f(x_k + \beta_k \xi_n) - f(x_k - \beta_k \xi_n)}{2\beta_k} \end{pmatrix},$$

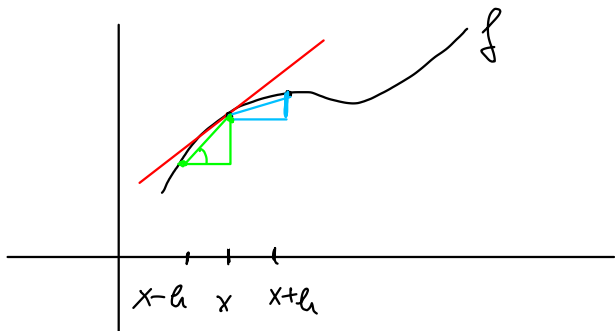
ahol ξ_i jelöli az i -edik egységvektort, és $\beta_k > 0$.

Az α_k és β_k egy tipikus megválasztása:

$$\alpha_k = \frac{\alpha}{(k+1+A)^a}, \quad \beta_k = \frac{\beta}{(k+1)^\gamma},$$

ahol A nemnegatív, a többi konstans pozitív.



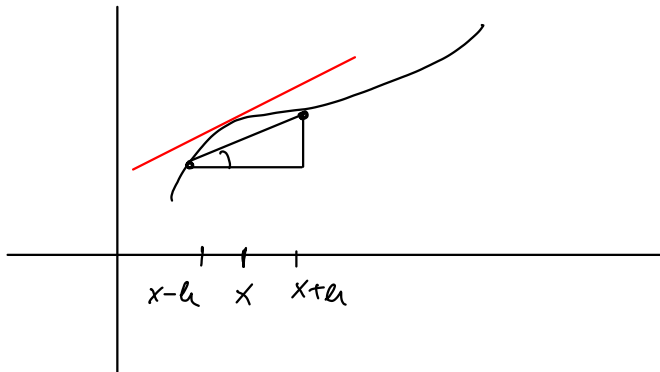


$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

haladó diff.

$$f'(x) \approx \frac{f(x) - f(x-h)}{h}$$

retrográd diff.



$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}$$

centered
diff.

A gradiens módszer egy sztochasztikus módosítása

$$\nabla f(x_k) \approx \underbrace{\frac{f(x_k + \beta_k \zeta_k) - f(x_k - \beta_k \zeta_k)}{2\beta_k}}_{\downarrow} \zeta_k = \frac{\Delta f(x_k, \beta_k \zeta_k)}{2\beta_k} \zeta_k$$

és

$$x_{k+1} = x_k - \frac{\alpha_k}{2\beta_k} \Delta f(x_k, \beta_k \zeta_k) \zeta_k,$$

ahol (β_k) egy csökkenő sorozat, míg ζ_k az $\|\zeta\| = 1$ egységkörön egyenletes eloszlású. ^{gömb}

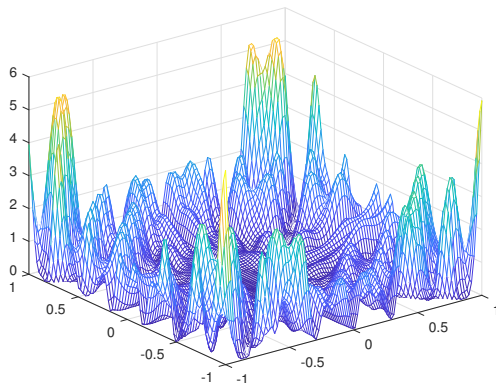
A gyakorlatban ajánlott egy skálázási lépés használata, mellyel elkerülhető a túl nagy lépés: legyen

$$g := \frac{\alpha_k}{2\beta_k} \Delta f(x_k, \beta_k \zeta_k) \zeta_k,$$

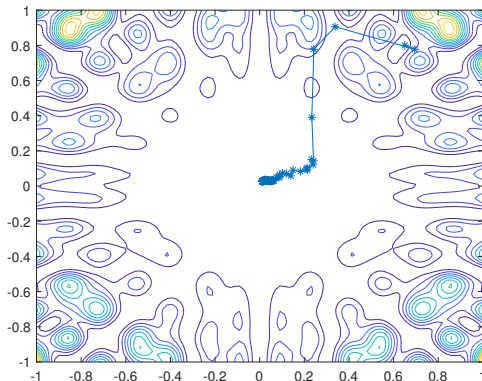
ha $\|g\| > 1$ akkor generáljunk egy új ζ_k -t.

Példa.

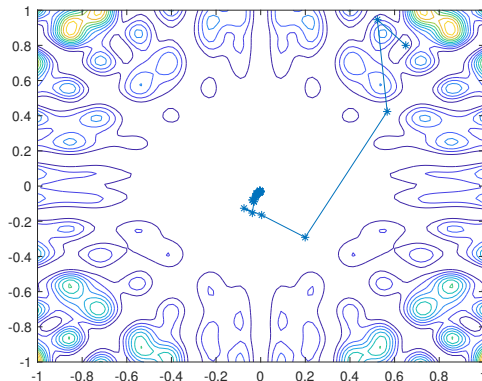
$$f(x) = (x_1 \sin(20x_2) + x_2 \sin(20x_1))^2 \cosh(\sin(10x_1)x_1) \\ + (x_1 \cos(10x_2) - x_2 \sin(10x_1))^2 \cosh(\cos(20x_2)x_2).$$



Az $f(x)$ függvény a $[-1, 1]^2$ tartomány felett. Minimumhely: $(0, 0)$.



$$\alpha_k = \frac{1}{k+1}, \beta_k = \frac{1}{(1+k)^{0.5}}, x_{opt} = [0.0109, 0.0268], k = 110$$



$$\alpha_k = \frac{1}{k+1}, \beta_k = \frac{1}{(1+k)^{0.5}}, x_{opt} = [-0.0042, 0.0260], k = 53$$

Sztochasztikus gradiens módszer

Az

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

iteráció helyett az

$$x_{k+1} = x_k - \alpha_k g(x_k, \xi_k),$$

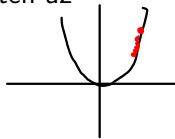
közelítést használjuk, ahol ξ_k egy valószínűségi változó, és $\mathbb{E}g(x_k, \xi_k) = \nabla f(x_k)$.

($g(x_k, \xi_k)$ a $\nabla f(x_k)$ egy torzítatlan becslése)

Gépi tanulásban (SGD: stochastic gradient descent)

Adott $\{(x^{(i)}, d^{(i)}), i = 1, \dots, M\}$ címkézett tanulóhalmaz esetén az

$$L(\Theta) = \frac{1}{M} \sum_{i=1}^M L_i(\Theta),$$

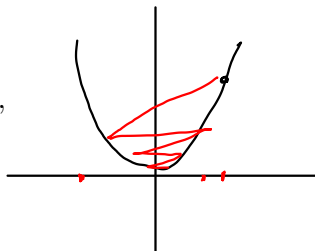


veszteségfüggvényt akarjuk minimalizálni, ahol Θ a hálózat paramétereit tartalmazó vektor, míg L_i az i -edik tanulóadatnak megfelelő veszteség.

Gradiens módszer esetén a paraméter módosítása:

$$\Theta := \Theta - \frac{\eta}{M} \sum_{i=1}^M \nabla L_i(\Theta),$$

ahol η a tanulási paraméter (a lépéshossz).



MSE

$$L(\theta) = \frac{1}{M} \sum_{i=1}^M \underbrace{(y^{(i)} - d^{(i)})^2}_{L_i(\theta)}$$

what $y^{(i)}$ or $x^{(i)}$ has a
label and what model output

Az **SGD** formula:

$$\Theta := \Theta - \eta \nabla L_i(\Theta).$$

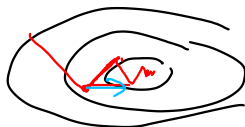
1 epoch: amikor i minden értéket felvesz $\{1, \dots, M\}$ -ből (a tanulóhalmaz elemeit véletlenszerűen keverjük minden epoch előtt).

Mini-batch SGD:

az $\{1, \dots, M\}$ indexhalmazt véletlenszerűen m elemű M_1, \dots, M_K részhalmazokra osztjuk, és

$$\Theta := \Theta - \frac{\eta}{m} \sum_{i \in M_j} \nabla L_i(\Theta).$$

1 epoch: amikor j minden értéket felvesz $\{1, \dots, K\}$ -ből.



Momentum módszer

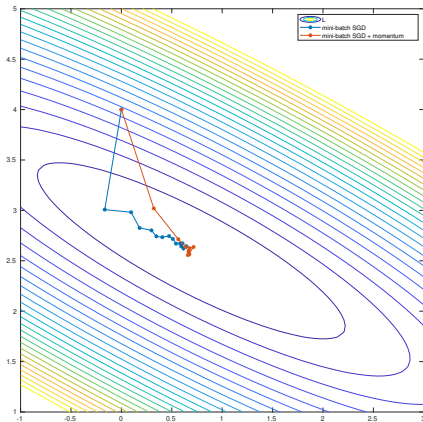
Az aktuális lépéshossz az előző lépéshossz és az aktuális gradiens lineáris kombinációja:

$$\Delta\Theta := \varepsilon \Delta\Theta - \eta G(\Theta),$$

$$\Theta := \Theta + \Delta\Theta,$$

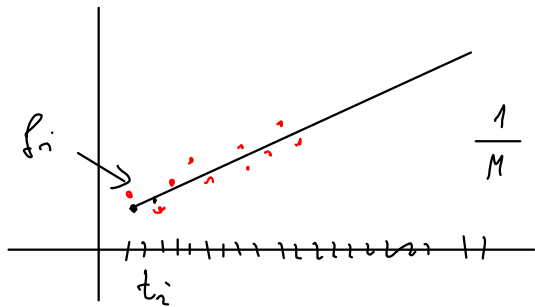
ahol $G(\Theta)$ az aktuális gradiens (a módszernek megfelelően).

“Játékpélda”



6 darab
mini-batch
30 tanulósor
⇒ epochonként
T idő múlva

Egyenes illesztése 30 generált, zajos, síkbeli pontra, a veszteségfüggvény az átlagos négyzetes hiba. A mini-batchek mérete 6, a momentum módszernél $\varepsilon = 0.9$. A tanulási paraméter állandó (nem ideális!) A közelítéseket az epochok végén ábrázoltuk.



$$y(t) = \boxed{a} + \boxed{b} \cdot t$$

$x_1 \quad x_2$

$$\frac{1}{M} \sum_{i=1}^M (y(t_i) - f_i)^2$$

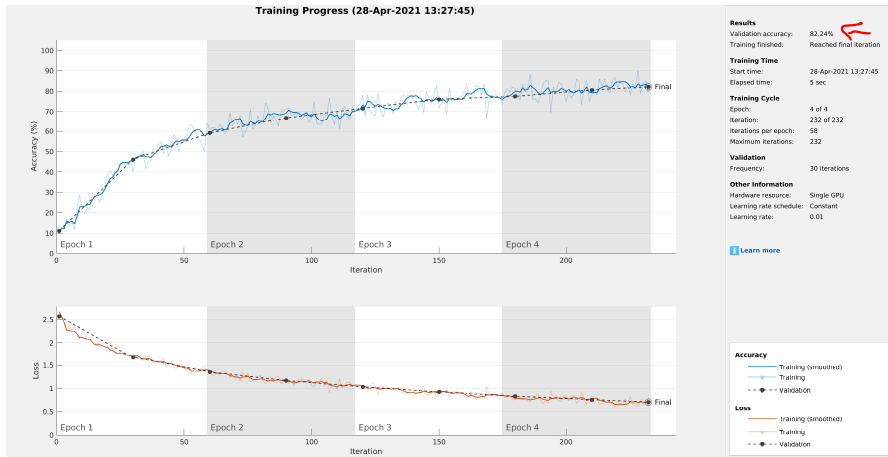
Példa

10 osztályos osztályozási feladat, 10x1000 kép, ennek 75%-ával tanítunk.

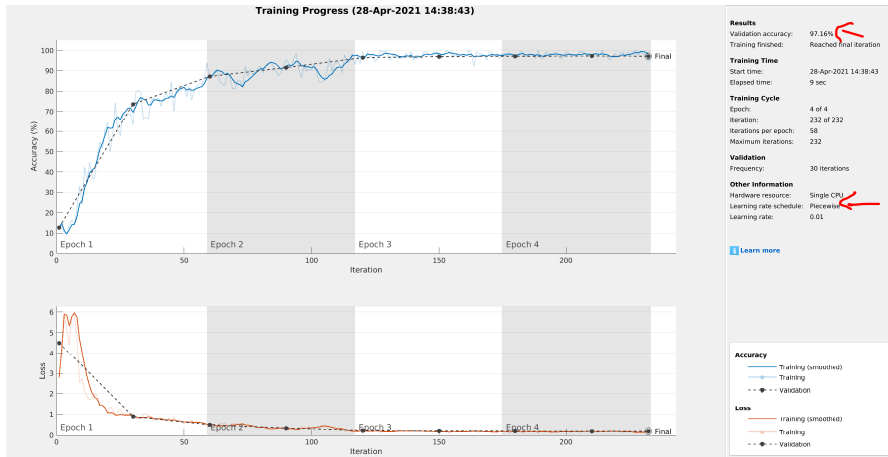
Veszteségfüggvény: keresztentrópia

- (a) mini-batch SGD, momentum nélkül, mini-batch mérete 128, tanulási paraméter állandó (0.01)
- (b) mini-batch SGD, momentum nélkül, mini-batch mérete 128, tanulási paraméter változik (0.1 két epochon át, majd 0.01)
- (c) mini-batch SGD, momentummal ($\varepsilon = 0.9$), mini-batch mérete 128, tanulási paraméter állandó (0.01)

(a) eset:



(b) eset:



(c) eset:

momentum módszer
 $\epsilon = 0.9$

