

# Gradient descent: implementálás

Kormányos Andor

Komplex Rendszerek Fizikája Tanszék

2023. november 14.

# Legmeredekebb ereszkedés gépi tanulásban<sup>1</sup>

## Stratégia

- iteratív algoritmus
- válasszunk egy tetszőleges kezdőértéket a  $a_0$ ,  $a_1$ -nak
- repeat until convergence {  
 $a_j := a_j - \alpha \frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)})$   
}

Néhány gyakorlati kérdés az implemetációval kapcsolatban

- hogyan frissítsük az  $a_j$  értékeket?
- hogyan számoljuk a költségfüggvény deriváltjait?
- hogyan határozzuk meg az  $\alpha$  learning rate-t?
- hogyan figyeljük a konvergenciát?

---

<sup>1</sup>Andrew Ng, Coursera alapján

# Egyidejű frissítés (simultaneous update)

**Példa:** két paraméter,  $a_0$ ,  $a_1$  . Hogyan frissítsük  $\frac{\partial}{\partial a_j} J(\mathbf{a})$  értékét?

Előkészület:

- $a_0$  és  $a_1$  egy  $\mathbf{a\_vec}[i]$  vektor elemeiként tárolhatjuk
- hasonlóan  $\frac{\partial}{\partial a_0} J(\mathbf{a})$  és  $\frac{\partial}{\partial a_1} J(\mathbf{a})$  értékét egy  $\text{derivJ}[i]$  vektor elemeiként

# Egyidejű frissítés (simultaneous update)

**Példa:** két paraméter,  $a_0$ ,  $a_1$  . Hogyan frissítsük  $\frac{\partial}{\partial a_j} J(\mathbf{a})$  értékét?

Előkészület:

- $a_0$  és  $a_1$  egy  $\mathbf{a\_vec}[i]$  vektor elemeiként tárolhatjuk
- hasonlóan  $\frac{\partial}{\partial a_0} J(\mathbf{a})$  és  $\frac{\partial}{\partial a_1} J(\mathbf{a})$  értékét egy  $\mathbf{derivJ}[i]$  vektor elemeiként

## Egyidejű frissítés

```
1  do
2  {
3  deriv_calc(derivJ, a_vec, data, params);
4  //calculate the derivative
5
6  a_vec[0]=a_vec[0]-alpha*derivJ[0];
7  a_vec[1]=a_vec[1]-alpha*derivJ[1];
8
9  J=fun_calc(a_vec, data, params);
10 //calculate J to monitor convergence
11 }
12 while (...) //convergence achieved
```

# Egyidejű frissítés (simultaneous update)

**Példa:** két paraméter,  $a_0$ ,  $a_1$ . Hogyan frissítsük  $\frac{\partial}{\partial a_j} J(\mathbf{a})$  értékét?

Előkészület:

- $a_0$  és  $a_1$  egy  $\mathbf{a\_vec}[i]$  vektor elemeiként tárolhatjuk
- hasonlóan  $\frac{\partial}{\partial a_0} J(\mathbf{a})$  és  $\frac{\partial}{\partial a_1} J(\mathbf{a})$  értékét egy  $\mathbf{derivJ}[i]$  vektor elemeiként

## Egyidejű frissítés

```
1 do
2 {
3 deriv_calc(derivJ, a_vec, data, params);
4 //calculate the derivative
5
6 a_vec[0]=a_vec[0]-alpha*derivJ[0];
7 a_vec[1]=a_vec[1]-alpha*derivJ[1];
8
9 J=fun_calc(a_vec, data, params);
10 //calculate J to monitor convergence
11 }
12 while (...) //convergence achieved
```

## Nem egyidejű frissítés

```
1 do
2 {
3 derivJ[0]=deriv0_calc(a_vec, data, params);
4 a_vec[0]=a_vec[0]-alpha*derivJ[0];
5 //a_0 updated
6
7 derivJ[1]=deriv1_calc(a_vec, data, params);
8 //uses the updated a_vec[0]
9 // and the old a_vec[1]
10 a_vec[1]=a_vec[1]-alpha*derivJ[1];
11
12 J=fun_calc(a_vec, data, params);
13 //calculate J to minitor convergence
14 }
15 while (...) //convergence achieved
```

# Egyidejű frissítés (simultaneous update)

**Példa:** két paraméter,  $a_0$ ,  $a_1$ . Hogyan frissítsük  $\frac{\partial}{\partial a_j} J(\mathbf{a})$  értékét?

Előkészület:

- $a_0$  és  $a_1$  egy  $\mathbf{a\_vec}[i]$  vektor elemeiként tárolhatjuk
- hasonlóan  $\frac{\partial}{\partial a_0} J(\mathbf{a})$  és  $\frac{\partial}{\partial a_1} J(\mathbf{a})$  értékét egy  $\mathbf{derivJ}[i]$  vektor elemeiként

## Egyidejű frissítés

```
1 do
2 {
3 deriv_calc(derivJ, a_vec, data, params);
4 //calculate the derivative
5
6 a_vec[0]=a_vec[0]-alpha*derivJ[0];
7 a_vec[1]=a_vec[1]-alpha*derivJ[1];
8
9 J=fun_calc(a_vec, data, params);
10 //calculate J to monitor convergence
11 }
12 while (...) //convergence achieved
```

## Nem egyidejű frissítés

```
1 do
2 {
3 derivJ[0]=deriv0_calc(a_vec, data, params);
4 a_vec[0]=a_vec[0]-alpha*derivJ[0];
5 //a_0 updated
6
7 derivJ[1]=deriv1_calc(a_vec, data, params);
8 //uses the updated a_vec[0]
9 // and the old a_vec[1]
10 a_vec[1]=a_vec[1]-alpha*derivJ[1];
11
12 J=fun_calc(a_vec, data, params);
13 //calculate J to minitor convergence
14 }
15 while (...) //convergence achieved
```

Az **egyidejű frissítést** (simultaneous update) használjuk!

# Hogyan számoljuk a költségfüggvény deriváltjait?

Két széleskörben alkalmazott eljárás:

- batch gradient descent
- stochastic gradient descent

# Batch gradient descent

**Példa:** legkisebb négyzetek módszere

- költségfüggvény  $J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - h(\mathbf{x}^{(i)}; \mathbf{a}))^2$
- korábban már látott példa: egy  $h(x; \mathbf{a}) = a_0 + a_1x$  függvényt szeretnénk illeszteni
- jelölés:  $h(\mathbf{x}; \mathbf{a}) = a_0x_0 + a_1x_1$ ,  $x_0 = 1.0$
- parciális deriváltak (egzakt eredmény):

$$\frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(\mathbf{x}^{(i)}; \mathbf{a})) x_j^{(i)}$$



# Batch gradient descent

**Példa:** legkisebb négyzetek módszere

- költségfüggvény  $J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - h(\mathbf{x}^{(i)}; \mathbf{a}))^2$
- korábban már látott példa: egy  $h(x; \mathbf{a}) = a_0 + a_1x$  függvényt szeretnénk illeszteni
- jelölés:  $h(\mathbf{x}; \mathbf{a}) = a_0x_0 + a_1x_1$ ,  $x_0 = 1.0$
- parciális deriváltak (egzakt eredmény):

$$\frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(\mathbf{x}^{(i)}; \mathbf{a})) x_j^{(i)}$$

- a fenti képlet közvetlenül beprogramozható
- az összes  $N$  adatot használjuk a parciális deriváltak kiszámolására
- ez a **batch gradient descent**

# Stochastic gradient descent

Parciális deriváltak (egzakt):

$$\frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(\mathbf{x}^{(i)}; \mathbf{a})) x_j^{(i)}$$

- ha  $N$  nagyon nagy (akár  $\sim 10^6 - 10^7$ ), akkor a  $\sum_{i=1}^N (\dots)$  összegek kiszámítása minden iterációs lépésben sok időt vesz el

## Sztochasztikus számolás

- minden iterációs lépésben válasszunk véletlenszerűen egy  $(\mathbf{x}^{(k)}, y^{(k)})$  adatot
- a  $\sum_{i=1}^N (\dots)$  összeget egyetlen taggal “közelítjük”

$$\frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) \approx (y^{(k)} - h(\mathbf{x}^{(k)}; \mathbf{a})) x_j^{(k)}$$

- ez a **stochastic gradient descent**

# Mini batch gradient descent

- a batch gradient descent és a stochastic gradient descent közötti átmeneti megoldás
- minden iterációs lépésben válasszunk véletlenszerűen  $P \ll N$  darab  $(\mathbf{x}^{(k)}, y^{(k)})$  adatot
- $P$  a “mini-batch” nagysága
- a derivált számolása:

$$\frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)}) \approx \frac{1}{P} \sum_{k=1}^P \left( y^{(k)} - h(\mathbf{x}^{(k)}; \mathbf{a}) \right) x_j^{(k)}$$

# Learning rate, konvergencia figyelése

Paraméterek iteratív meghatározása:

$$a_j := a_j - \alpha \frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)})$$

Hogyan válasszuk az  $\alpha$  learning rate-t ?

# Learning rate, konvergencia figyelése

Paraméterek iteratív meghatározása:

$$a_j := a_j - \alpha \frac{\partial}{\partial a_j} J(\mathbf{a}; \mathbf{x}^{(i)}, y^{(i)})$$

Hogyan válasszuk az  $\alpha$  learning rate-t ?

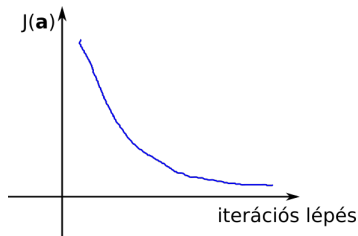
## batch gradient descent

- jelöljük  $\mathbf{a}^{(p)}$ -vel a  $p$ -ik iteráció eredményét
- ha  $\alpha$  elég kicsi, akkor  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$  minden iteráció során csökken
- viszont, ha  $\alpha$  túl kicsi, akkor lassú lesz a konvergencia
- ha  $\alpha$  nem elég kicsi, akkor  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$  növekszik vagy esetleg oszcillál

# Learning rate, konvergencia figyelése

Eljárás **batch gradient descent**-re

- válasszunk egy  $\alpha$ -t
- minden iterációs lépésben értékeljük ki  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$ -t
- ábrázoljuk  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$ -t az iterációs lépések függvényében

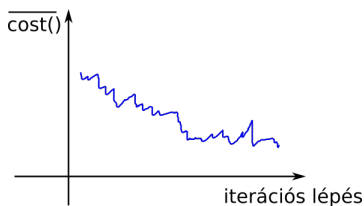


- ha  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$  csökken, akkor  $\alpha$  nem túl nagy
- ha  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$  nem csökken, akkor válasszunk egy kisebb  $\alpha$ -t
- ha  $J(\mathbf{a}^{(p)}; \mathbf{x}^{(i)}, y^{(i)})$  nem csökken bizonyos iterációs szám fölött  $\Rightarrow$  konvergencia

# Learning rate, konvergencia figyelése

Eljárás **stochastic gradient descent**-re

- minden iterációs lépésben, mielőtt még frissítenénk az  $\mathbf{a}^{(p)}$ -t, számítsuk ki a  $\text{cost}(y^{(k)}, h(\mathbf{a}^{(p)}; \mathbf{x}^{(k)}))$
- 100 v 500 stb lépésenként írjuk ki a  $\text{cost}(y^{(k)}, h(\mathbf{a}^{(p)}; \mathbf{x}^{(k)}))$ -k **átlagát**



- ha  $\alpha$  nagysága megfelelő, akkor az átlag csökkenő trendet mutat

# Learning rate változtatása

Az eddigiek során

- az  $\alpha$  learning rate-t állandónak tartottuk végig az iteráció során
- minden  $a_j$  esetén ugyanazt az  $\alpha$ -t használtuk

Ezen is lehet javítani:

- az iteráció során figyeljük a konvergenciát és valamilyen ütem szerint csökkentjük az  $\alpha$ -t
- a különböző  $a_j$  esetén különböző  $\alpha_j$ -t használunk

Lásd pl [ADAM](#) módszer