

# Többrétegű Kerceptron\*

Szabó Zoltán, Lőrincz András

Információs Rendszerek Tanszék, Informatika Kar

Eötvös Loránd Tudományegyetem

Pázmány Péter sétány 1/C

Budapest, 1117

e-mail: szzoli@cs.elte.hu, andras.lorincz@elte.hu

## Kivonat

Többrétegű Perceptronokba (MLP) Támasztó Vektor Gépeket (SVM) ágyazva többrétegű SVM hálókat konstruálunk. Az összekapcsolt approximációs forma az SVM-ek általánosító képességét és az MLP-k rejtett rétegéből adódó kombinatorikus tulajdonságot egyaránt kihasználhatja. A hálózatot Többrétegű Kerceptron (MLK) hálózatnak nevezzük. Az MLK rendelkezik hibavisszaterjesztésen alapuló hangolási eljárással, amit jelen munkában bemutatunk. Négyzetes költségfüggvényre – regularizációs lehetőségekkel – hangolási szabályt származtatunk. Megközelítésünk egy további tulajdonsága, hogy az ún. *kernel trükk* segítségével az MLK-hoz tartozó számítások a duális térben kivitelezhetők.

## 1. Bevezetés

A Többrétegű Perceptronokat (MLP) és a Támasztó Vektor Gépeket (SVM) széles körben tanulmányozták az irodalomban. Kiváló áttekintést ad a témában [3, 2]. Munkánkban az SVM-eket többrétegű formára terjesztjük ki, és a kapott rendszerre hibavisszaterjesztésen alapuló hangolási szabályt vezetünk le. Az ún. kernel trükk alkalmazásával, a problémát skaláris szorzat segítségével tárgyaljuk. Más, ugyanezt a trükköt használó eljárások leírása megtalálható a [4, 6, 7] hivatkozásokban.

## 2. A hálózat felépítése

### 2.1. Jelölések

Különböző betűtípussal jelöljük a számokat ( $a$ ), a vektorokat ( $\mathbf{a}$ ), és a mátrixokat ( $\mathbf{A}$ ).  $\mathbf{A}^T$  az  $\mathbf{A}$  mátrix transzponáltja. Az  $\mathbf{a}$  vektor egy  $a$  komponenssel való

---

\* Alkalmazott Matematikai Lapok, 24:209-222, 2007.

kibővítését  $[\mathbf{a}; a]$ -ként írjuk.  $\mathbb{R}$  szimbolizálja a valós számokat.  $\|\cdot\|_2$  jelöli az  $E$  Euklideszi térbeli  $\langle \cdot, \cdot \rangle$  skaláris szorzat által indukált  $L_2$  normát, azaz  $\|\mathbf{e}\|_2 = \sqrt{\langle \mathbf{e}, \mathbf{e} \rangle}$  ( $\mathbf{e} \in E$ ).

## 2.2. Építőelemek

### 2.2.1. SVM

Az SVM-ek gyakran használt approximációs eszközök [9, 10, 8, 6, 5].  $\{\mathbf{x}(t), d(t)\}_{t=1..T}$  mintapárokat közelítenek, ahol az  $\mathbf{x}(t)$  input az  $\mathcal{X}$  *input térből* származik és  $d(t) \in \mathbb{R}$ . A közelítés lineáris, egy alkalmas  $\mathcal{H}$  térben. Ebbe a térbe a

$$\varphi : \mathbf{x} \in \mathcal{X} \rightarrow \mathcal{H} \quad (1)$$

hozzárendelés képezi le az  $\mathbf{x}(t)$  inputokat.  $\varphi(\mathbf{x})$ -et az  $\mathbf{x}$  input *reprezentációjaként* interpretálhatjuk. Az SVM közelítés a

$$f_{\mathbf{w}} : \mathbf{x} \in \mathcal{X} \mapsto \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle_{\mathcal{H}} \quad (\mathbf{w} \in \mathcal{H}) \quad (2)$$

formájú. Formálisan, az SVM feladat

$$\min_{\mathbf{w}} H[\mathbf{w}] := C \cdot \sum_{t=1}^T V[d(t), f_{\mathbf{w}}(\mathbf{x}(t))] + \frac{1}{2} \|\mathbf{w}\|_{\mathcal{H}}^2 \quad (C > 0), \quad (3)$$

ahol  $V[\cdot, \cdot]$  az ún. *veszteségfüggvény*, amely lehet kvadratikus,  $\epsilon$ -érzékeny, de más formákat is szoktak használni [4]. Röviden, az SVM-ek regularizált lineáris approximátorok [2].

Az explicit  $\varphi$  leképezés helyett, a  $\mathcal{H}$  tér egy  $k$  kernel segítségével is leírható,  $\mathcal{H} = \mathcal{H}(k)$  [11], ahol  $\varphi(\mathbf{x}) = k(\cdot, \mathbf{x})$ . A  $k$  kernel a

$$\langle f(\cdot), k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad (\mathbf{x} \in \mathcal{X}, \forall f \in \mathcal{H}), \quad (4)$$

reprodukáló tulajdonsággal rendelkeznek [1, 11] és  $\mathcal{H}$ -t Reprodukáló Kernel Hilbert Térnek (RKHS) nevezzük. Tehát, tetszőleges  $f \in \mathcal{H}$  RKHS-beli függvény  $k(\cdot, \mathbf{x})$  kernellel való skaláris szorzata az  $\mathbf{x}$  pontbeli kiértékelésnek felel meg. A skaláris szorzat a  $\mathcal{H}$  térben implicit módon számolható a kernel segítségével

$$k(\mathbf{u}, \mathbf{v}) = \langle \varphi(\mathbf{u}), \varphi(\mathbf{v}) \rangle_{\mathcal{H}} \quad (\mathbf{u}, \mathbf{v} \in \mathcal{X}). \quad (5)$$

Így, a  $\mathbf{w} = \sum_{j=1}^N \alpha_j \cdot \varphi(\mathbf{z}_j)$  ( $\alpha_j \in \mathbb{R}, \mathbf{z}_j \in \mathcal{X}$ ) választással

$$f_{\mathbf{w}}(\mathbf{x}) = \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle_{\mathcal{H}} = \sum_{j=1}^N \alpha_j \cdot \langle \varphi(\mathbf{z}_j), \varphi(\mathbf{x}) \rangle_{\mathcal{H}} = \sum_{j=1}^N \alpha_j \cdot k(\mathbf{z}_j, \mathbf{x}). \quad (6)$$

Tehát, az  $f_{\mathbf{w}}$  függvény az  $\alpha_j$  együtthatók, a  $\mathbf{z}_j$  minták és a  $k$  kernel segítségével a  $\varphi(\mathbf{x})$  reprezentáció explicit felhasználása nélkül kiértékelhető. Ez a fogás a *kernel trükk*.

### 2.2.2. MLP

Az MLP neurális hálózat többrétegű, minden réteg egy

$$\mathbf{x} \mapsto \mathbf{g}(\mathbf{W} \cdot \mathbf{x}) \quad (7)$$

formájú nem-lineáris leképezést valósít meg. Itt  $\mathbf{g}$  egy differenciálható nem-lineáris függvény. Az MLP feladatban úgy hangoljuk az egyes rétegek  $\mathbf{W}$  mátrixait, hogy a hálózat az  $\{\mathbf{x}(t), \mathbf{d}(t)\}$  input-output minta pároknak megfelelő leképezést közelítse. Formálisan, célunk a

$$\varepsilon^2(t) := \|\mathbf{d}(t) - \mathbf{y}(t)\|_2^2 \rightarrow \min_{\mathbf{w}_1, \mathbf{w}_2, \dots}, \quad (8)$$

kvadratikusan minimalizálása, ahol  $\mathbf{y}(t)$  jelöli a hálózat  $t$  időpontbeli kimenetét. Az MLP feladatot oldja meg a jól-ismert *visszaterjesztési algoritmus*.

### 2.3. Az MLK hálózat

Egy általános MLP réteg által megvalósított leképezés [lásd a (7) egyenletet] a

$$\mathbf{x} \mapsto \mathbf{g} \left( \begin{bmatrix} \vdots \\ \langle \mathbf{w}_i, \mathbf{x} \rangle \\ \vdots \end{bmatrix} \right), \quad (9)$$

formájú, ahol  $\mathbf{w}_i^T$  a  $\mathbf{W}$  mátrix  $i$ . sorát jelöli. SVM illeszthető az MLP-be, ha a hálózat egy általános rétege a

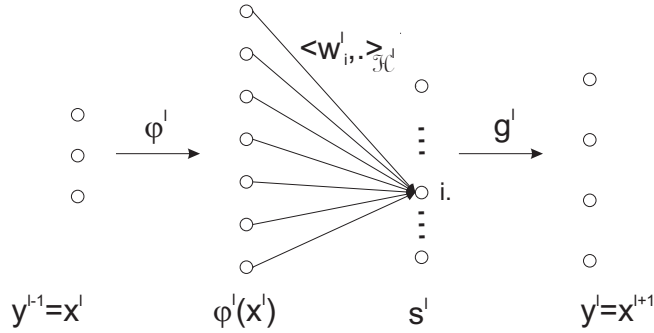
$$\mathbf{x} \mapsto \mathbf{g} \left( \begin{bmatrix} \langle \mathbf{w}_1, \varphi(\mathbf{x}) \rangle_{\mathcal{H}} \\ \vdots \\ \langle \mathbf{w}_N, \varphi(\mathbf{x}) \rangle_{\mathcal{H}} \end{bmatrix} \right) \quad (10)$$

hozzárendelést valósítja meg.<sup>1</sup>

Egy ilyen rétegekből felépített hálózatot Többrétegű Kerceptronnak (MLK) fogunk hívni, lásd az 1. ábrát. Minden egyes réteg ( $\mathbf{x}^l$ ) inputját az előző réteg ( $\mathbf{y}^{l-1}$ ) outputja szolgáltatja. A 0-adik réteg a külvilág, ami az MLK első rétege számára szolgáltatja a bemenetet.  $\mathbf{x}^l = \mathbf{y}^{l-1} \in \mathbb{R}^{N_1^l}$ , ahol  $N_1^l$  az  $l$ -edik réteg dimenziója. Az  $l$ -edik réteg  $\mathbf{x}^l$  inputja a  $\varphi^l$  leképezésen átesve a  $\mathbf{w}_i^l$  súlyokkal szorzódik. Ez a két lépcsős eljárás implicit módon elvégezhető a  $k^l$  kerneleket és a  $\mathbf{w}_i^l$ -k kifejtéseit használva. Az adódó  $\mathbf{s}^l \in \mathbb{R}^{N_s^l}$  vektorra hat a  $\mathbf{g}^l$  nem-lineáris, differenciálható függvény. Ennek a nem-lineáris függvénynek a kimenete a következő réteg bemenete, azaz  $\mathbf{x}^{l+1}$ . Az utolsó ( $L$ .) réteg outputját – azaz a hálózat outputját –  $\mathbf{y}$  jelöli.  $\mathbf{y}^l = \mathbf{x}^{l+1} \in \mathbb{R}^{N_o^l}$ , és az  $l$ -edik réteg kimenetének dimenziója  $N_o^l$ .

A következőkben megmutatjuk, hogy (i) az MLK-k is rendelkeznek visszaterjesztési szabállyal, ami (ii) csak kernelek segítségével is megadható, és így a számítások a duális térben kivitelezhetőek.

<sup>1</sup> Az egyszerűség kedvéért válasszuk az  $\mathcal{X}$  input teret a véges dimenziós Euklideszi térnek, azaz  $\mathbb{R}^n$ -nek.



1. ábra. Az MLK hálózat  $l$ -edik rétege,  $l = 1, 2, \dots, L$ . Minden egyes réteg inputját ( $\mathbf{x}^l$ ) az előző réteg outputja adja ( $\mathbf{y}^{l-1}$ ). A 0-adik réteg a külvilág, ami az MLK első réteg számára szolgáltatja a bemenetet. Az  $l$ -edik réteg  $\mathbf{x}^l$  inputja a  $\varphi^l$  leképezésen esik át, majd a réteg  $\mathbf{w}_i^l$  súlyaival szorozódik skalárisan a  $\mathcal{H}^l = \mathcal{H}^l(k^l)$  RKHS-ben. Az adódó  $\mathbf{s}^l$  vektorra hat a  $g^l$  differenciálható nem-linearitás. Ezen nem-lineáris függvény kimenete a következő réteg bemenete,  $\mathbf{x}^{l+1}$ . A hálózat kimenete az utolsó réteg kimenete.

### 3. Az MLK visszaterjesztési eljárás

Egy kicsit általánosabb, regularizációs tagokat is tartalmazó feladat a

$$c(t) := \varepsilon^2(t) + r(t) \longrightarrow \min_{\{\mathcal{H}^l \ni \mathbf{w}_i^l: l=1, \dots, L; i=1, \dots, N_S^l\}}, \quad (11)$$

probléma, ahol  $\varepsilon^2(t) = \|\mathbf{d}(t) - \mathbf{y}(t)\|_2^2$  és  $r(t) = \sum_{l=1}^L \sum_{i=1}^{N_S^l} \lambda_i^l \cdot \|\mathbf{w}_i^l(t)\|_{\mathcal{H}^l}^2$  ( $\lambda_i^l \geq 0$ ) a költségfüggvény approximációs és regularizációs tagjai, és  $\mathbf{y}(t)$  jelöli a hálózat  $t$ -edik inputra adott kimenetét. A  $\lambda_i^l$  paraméterek szabályozzák az approximáció és regularizáció közötti arányt.  $\lambda_i^l = 0$ -ra a legjobb közelítést keressük, mint az MLP feladatban [(8) egyenlet].  $\lambda_i^l$  értékeket növelve, az approximáció simasága nő.

A fenti jelölésekkel a következő állítások igazolhatók.

**1. Tétel** (explicit eset). *Tegyük fel, hogy az  $\mathbf{x} \mapsto \langle \mathbf{w}, \varphi^l(\mathbf{x}) \rangle_{\mathcal{H}^l}$  és a  $g^l$  függvények differenciálhatók ( $l = 1, \dots, L$ ). Ekkor visszaterjesztési szabály származtatható az MLK-ra, ha a költségfüggvény*

$$c(t) = \varepsilon^2(t) + \sum_{l=1}^L \sum_{i=1}^{N_S^l} \lambda_i^l \cdot \|\mathbf{w}_i^l(t)\|_{\mathcal{H}^l}^2 \quad (\lambda_i^l \geq 0) \quad (12)$$

alakú.

**2. Tétel** (implicit eset). *Tegyük fel, hogy az alábbiak teljesülnek:*

1. *Differenciálhatósági megkötés: A  $k^l$  kernelek mindkét változójukban, illetve a  $g^l$  függvények differenciálhatók ( $l = 1, \dots, L$ ).*

2. *Kifejtési tulajdonság: A hálózat kezdeti  $\mathbf{w}_i^l(1)$  súlyai egy adott*

$$\mathcal{H}^l \ni \mathbf{w}_i^l(1) = \sum_{j=1}^{N_S^l(1)} \alpha_{i,j}^l(1) \cdot \varphi^l(\mathbf{z}_{i,j}^l(1)) \quad (l = 1, \dots, L; i = 1, \dots, N_S^l) \quad (13)$$

*típusú kifejtéssel, duális reprezentációval rendelkeznek.*

*Ekkor létezik visszaterjesztési eljárás az MLK hálózatra, feltéve, hogy a költségfüggvény*

$$c(t) = \varepsilon^2(t) + \sum_{l=1}^L \sum_{i=1}^{N_S^l} \lambda_i^l \cdot \|\mathbf{w}_i^l(t)\|_{\mathcal{H}^l}^2 \quad (\lambda_i^l \geq 0) \quad (14)$$

*formájú. Az eljárás megőrzi a (13) tulajdonságot, ami így a behangolt hálózatra is fennáll. Az algoritmus implicit, abban az értelemben, hogy a duális térben realizálható.*

Az MLK visszaterjesztési eljárások pszeudokódjai az 1. és a 2. táblázatban található. Az algoritmusok levezetését, mind az explicit mind az implicit esetre, a következő alfejezetben megadjuk.

Az MLK visszaterjesztési eljárások szemléletesen (párhuzamosan lásd az 1 és a 2. táblázatokat):

1. a  $\delta^l(t)$  visszaterjesztett hiba  $\delta^L(t)$ -ből indulva egy hátráló rekurzióval fejlődik a  $\frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]}$  deriválton keresztül.
2. a  $\frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]}$  kifejezés a  $\varphi^{l+1}$  leképezés, vagy implicit módon a  $k^{l+1}$  kernel segítségével határozható meg.
3.  $\mathbf{w}$ -k hangolásában két tényező játszik szerepet:
  - (a) *felejtés* valósul meg a  $\mathbf{w}_i^l$  súlyok  $(1 - 2\mu_i^l(t) \cdot \lambda_i^l)$ -szeres szorzása által, ahol  $\lambda_i^l$  a regularizációs együtttható.
  - (b) *adaptáció* jelenik meg a visszaterjesztett hibán keresztül. Az  $l$ -edik réteg súlyait az  $\mathbf{x}^l(t)$  reprezentáció, azaz az aktuális inputnak az  $l$ -edik rétegre leképezett értéke állítja úgy, hogy a hangolást a visszaterjesztett hiba súlyozza.

### 3.1. Az MLK visszaterjesztési eljárások levezetése

Először a  $\frac{d[c(t)]}{d[\mathbf{w}_i^l(t)]}$  gradienst származtatjuk. Utána a gradienst a legmeredekebb lejtő módszerbe ágyazzuk.<sup>2</sup> A  $c(t)$  hiba két tagból áll, approximációs és regularizációs tagból:

$$c(t) = \varepsilon^2(t) + r(t). \quad (15)$$

<sup>2</sup> A legmeredekebb lejtő módszerét használjuk ötletünk bemutatásához. Más, ettől eltérő gradiens alapú technikák szintén szóba jöhetnek. Például, a momentum módszer illetve a konjugált gradiens eljárások is rendelkeznek előnyös tulajdonságokkal.

1. táblázat. Az explicit MLK visszaterjesztési algoritmus pszeudokódja

|                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>Algoritmus bemenete</b><br/> mintapontok: <math>\{\mathbf{x}(t), \mathbf{d}(t)\}_{t=1, \dots, T}</math><br/> költségfüggvény: <math>\lambda_i^l \geq 0</math> (<math>l = 1, \dots, L; i = 1, \dots, N_S^l</math>)<br/> tanulási ráták: <math>\mu_i^l(t) &gt; 0</math> (<math>l = 1, \dots, L; i = 1, \dots, N_S^l; t = 1, \dots, T</math>)</p> <p><b>Hálózat inicializációja</b><br/> méretek: <math>L</math> (rétegek száma), <math>N_1^l, N_S^l, N_o^l</math> (<math>l = 1, \dots, L</math>)<br/> súlyok: <math>\mathbf{w}_i^l(1)</math> (<math>l = 1, \dots, L; i = 1, \dots, N_S^l</math>)</p> <p><b>Számítás kezdete</b><br/> <b>Aktuális input <math>\mathbf{x}(t)</math></b><br/> <b>Előreterjesztés</b><br/> <math>\mathbf{x}^l(t)</math> (<math>l = 2, \dots, L+1</math>), <math>\mathbf{s}^l(t)</math> (<math>l = 2, \dots, L</math>)<sup>a</sup><br/> <b>Hiba visszaterjesztése</b><br/> <math>l = L</math><br/> while <math>l \geq 1</math><br/> if (<math>l = L</math>)<br/> <math>\delta^L(t) = 2 \cdot [\mathbf{y}(t) - \mathbf{d}(t)]^T \cdot (\mathbf{g}^L)'(\mathbf{s}^L(t))</math><br/> else<br/> <math display="block">\frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} = \begin{bmatrix} \vdots \\ \frac{d[\langle \mathbf{w}_i^{l+1}(t), \varphi^{l+1}(\mathbf{u}) \rangle_{\mathcal{H}^{l+1}}]}{d[\mathbf{u}]} \Big _{\mathbf{u}=\mathbf{x}^{l+1}(t)} \\ \vdots \end{bmatrix} \cdot [(\mathbf{g}^l)'(\mathbf{s}^l(t))]^b</math><br/> <math display="block">\delta^l(t) = \delta^{l+1}(t) \cdot \frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]}</math><br/> <b>Súlyok frissítése</b><br/> for <math>\forall i: 1 \leq i \leq N_S^l</math><br/> <math display="block">\mathbf{w}_i^l(t+1) = (1 - 2\mu_i^l(t) \cdot \lambda_i^l) \cdot \mathbf{w}_i^l(t) - \mu_i^l(t) \cdot \delta_i^l(t) \cdot \varphi^l(\mathbf{x}^l(t))</math><br/> <math>l = l - 1</math></p> <p><b>Számítás vége</b></p> |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

<sup>a</sup> Így a hálózat kimenete, azaz  $\mathbf{y}(t) = \mathbf{x}^{L+1}(t)$  is kiszámítható.

<sup>b</sup> Itt:  $i = 1, \dots, N_S^{l+1}$ .

2. táblázat. Az implicit MLK visszaterjesztési algoritmus pszeudokódja

|                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>Algoritmus bemenete</b><br/> mintapontok: <math>\{\mathbf{x}(t), \mathbf{d}(t)\}_{t=1, \dots, T}, T</math><br/> költségfüggvény: <math>\lambda_i^l \geq 0</math> (<math>l = 1, \dots, L; i = 1, \dots, N_S^l</math>)<br/> tanulási ráták: <math>\mu_i^l(t) &gt; 0</math> (<math>l = 1, \dots, L; i = 1, \dots, N_S^l; t = 1, \dots, T</math>)</p> <p><b>Hálózat inicializációja</b><br/> méretek: <math>L</math> (rétegek száma), <math>N_1^l, N_S^l, N_o^l</math> (<math>l = 1, \dots, L</math>)<br/> súlyok: <math>\mathbf{w}_i^l(1)</math>-kifejtések (<math>l = 1, \dots, L; i = 1, \dots, N_S^l</math>)<br/> együtthatók: <math>\alpha_i^l(1) \in \mathbb{R}^{N_i^l(1)}</math><br/> ősök: <math>\mathbf{z}_{i,j}^l(1)</math>, ahol <math>j = 1, \dots, N_i^l(1)</math></p> <p><b>Számítás kezdete</b><br/> <b>Aktuális input <math>\mathbf{x}(t)</math></b><br/> <b>Előreterjesztés</b><br/> <math>\mathbf{x}^l(t)</math> (<math>l = 2, \dots, L + 1</math>), <math>\mathbf{s}^l(t)</math> (<math>l = 2, \dots, L</math>)<sup>a</sup><br/> <b>Hiba visszaterjesztése</b><br/> <math>l = L</math><br/> while <math>l \geq 1</math><br/> if (<math>l = L</math>)<br/> <math>\delta^L(t) = 2 \cdot [\mathbf{y}(t) - \mathbf{d}(t)]^T \cdot (\mathbf{g}^L)'(\mathbf{s}^L(t))</math><br/> else<br/> <math display="block">\frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} = \begin{bmatrix} \vdots \\ \sum_{j=1}^{N_i^{l+1}(t)} \alpha_{ij}^{l+1}(t) \cdot [k^{l+1}]'_y(\mathbf{z}_{ij}^{l+1}(t), \mathbf{x}^{l+1}(t)) \\ \vdots \end{bmatrix} \cdot [(\mathbf{g}^l)'(\mathbf{s}^l(t))]^b</math><br/> <math display="block">\delta^l(t) = \delta^{l+1}(t) \cdot \frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]}</math><br/> <b>Súlyok frissítése</b><br/> for <math>\forall i: 1 \leq i \leq N_S^l</math><br/> <math>N_i^l(t+1) = N_i^l(t) + 1</math><br/> <math>\alpha_i^l(t+1) = [(1 - 2\mu_i^l(t) \cdot \lambda_i^l) \cdot \alpha_i^l(t); -\mu_i^l(t) \cdot \delta_i^l(t)]</math><br/> <math>\mathbf{z}_{i,j}^l(t+1) = \mathbf{z}_{i,j}^l(t)</math> (<math>j = 1, \dots, N_i^l(t)</math>)<br/> <math>\mathbf{z}_{i,j}^l(t+1) = \mathbf{x}^l(t)</math> (<math>j = N_i^l(t+1)</math>)<br/> <math>l = l - 1</math></p> <p><b>Számítás vége</b></p> |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

<sup>a</sup> Így a hálózat kimenete, azaz  $\mathbf{y}(t) = \mathbf{x}^{L+1}(t)$  is kiszámítható.

<sup>b</sup>  $i = 1, \dots, N_S^{l+1}$ .  $(k^l)'_y$  jelöli a  $k^l$  kernel második argumentuma szerint vett deriváltját.

### 3.1.1. Az approximációs tag gradiense

Először néhány MLK felépítéséből adódó alapösszefüggést sorolunk fel. Az egyszerűség kedvéért a továbbiakban a  $t$  indexet elhagyjuk [precízen:  $\mathbf{x}^l = \mathbf{x}^l(t)$ ,  $\mathbf{y}^l = \mathbf{y}^l(t)$ ,  $\mathbf{s}^l = \mathbf{s}^l(t)$ ,  $\mathbf{w}_i^l = \mathbf{w}_i^l(t)$ ].

$$\mathbf{x}^l = \mathbf{y}^{l-1} \in \mathbb{R}^{N_i^l} \quad (l = 1, \dots, L+1) \quad (16)$$

$$\mathbf{x}^{l+1} = \mathbf{g}^l(\mathbf{s}^l) \quad (l = 1, \dots, L) \quad (17)$$

$$\mathbf{s}^l = \begin{bmatrix} \langle \mathbf{w}_1^l, \boldsymbol{\varphi}^l(\mathbf{x}^l) \rangle_{\mathcal{H}^l} \\ \vdots \\ \langle \mathbf{w}_i^l, \boldsymbol{\varphi}^l(\mathbf{x}^l) \rangle_{\mathcal{H}^l} \\ \vdots \end{bmatrix} \quad (l = 1, \dots, L; i = 1, \dots, N_S^l) \quad (18)$$

$$= \begin{bmatrix} \langle \mathbf{w}_1^l, \boldsymbol{\varphi}^l(\mathbf{g}^{l-1}(\mathbf{s}^{l-1})) \rangle_{\mathcal{H}^l} \\ \vdots \\ \langle \mathbf{w}_i^l, \boldsymbol{\varphi}^l(\mathbf{g}^{l-1}(\mathbf{s}^{l-1})) \rangle_{\mathcal{H}^l} \\ \vdots \end{bmatrix} \quad (l = 2, \dots, L; i = 1, \dots, N_S^l) \quad (19)$$

$$\mathbf{s}^{l+1} = \begin{bmatrix} \langle \mathbf{w}_1^{l+1}, \boldsymbol{\varphi}^{l+1}(\mathbf{g}^l(\mathbf{s}^l)) \rangle_{\mathcal{H}^{l+1}} \\ \vdots \\ \langle \mathbf{w}_i^{l+1}, \boldsymbol{\varphi}^{l+1}(\mathbf{g}^l(\mathbf{s}^l)) \rangle_{\mathcal{H}^{l+1}} \\ \vdots \end{bmatrix} \quad (20)$$

$(l = 1, \dots, L-1; i = 1, \dots, N_S^{l+1})$

Az  $l$ -edik réteg visszaterjesztett hibáját definiáljuk a

$$\boldsymbol{\delta}^l(t) := \frac{d[\varepsilon^2(t)]}{d[\mathbf{s}^l(t)]} \quad (l = 1, \dots, L) \quad (21)$$

módon. Speciálisan, az utolsó rétegre:

$$\boldsymbol{\delta}^L(t) = \frac{d[\varepsilon^2(t)]}{d[\mathbf{s}^L(t)]} = \frac{d \left[ \|\mathbf{d}(t) - \mathbf{g}^L(\mathbf{s}^L(t))\|_2^2 \right]}{d[\mathbf{s}^L(t)]} \quad (22)$$

$$= 2 \cdot [\mathbf{g}^L(\mathbf{s}^L(t)) - \mathbf{d}(t)]^T \cdot (\mathbf{g}^L)'(\mathbf{s}^L(t)) \quad (23)$$

$$= 2 \cdot [\mathbf{y}(t) - \mathbf{d}(t)]^T \cdot (\mathbf{g}^L)'(\mathbf{s}^L(t)). \quad (24)$$

Itt először a láncszabályt, majd a vektorokra érvényes

$$\frac{d[\|\mathbf{d} - \mathbf{y}\|_2^2]}{d\mathbf{y}} = 2(\mathbf{y} - \mathbf{d})^T \quad (25)$$

összefüggést használtuk ki, végül beillesztettük az MLK szerkezetéből adódó

$$\mathbf{y}(t) = \mathbf{g}^L(\mathbf{s}^L(t)) \quad (26)$$



azonosságot.

A

$$\frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} \quad (l = 1, \dots, L-1) \quad (27)$$

kifejezés a (20) egyenlet segítségével számolható. Elégséges a

$$\frac{d[\langle \mathbf{w}, \varphi(\mathbf{g}(\mathbf{s})) \rangle_{\mathcal{H}}]}{d[\mathbf{s}]} \quad (28)$$

alakú kifejezéseket tekinteniünk, abból a teljes derivált „kirakható”. (28) értéke az alábbi lemma alkalmazásával megadható.

**1. Lemma.** *Legyen  $\mathbf{w} \in \mathcal{H} = \mathcal{H}(k)$  egy RKHS-beli pont. Tegyük fel, hogy*

1. *A  $k$  kernel mindkét argumentuma szerint differenciálható és jelölje  $k'_y$  a kernel második argumentuma szerint vett deriváltját.*

2. *Implicit esetben feltételezzük még, hogy  $\mathbf{w}$  véges sok  $\mathbf{z}_i$  pont  $\mathcal{H}$ -beli reprezentációjának képterében fekszik. Azaz*

$$\mathbf{w} \in \text{Im}(\varphi(\mathbf{z}_1), \varphi(\mathbf{z}_2), \dots, \varphi(\mathbf{z}_N)) \subseteq \mathcal{H}. \quad (29)$$

*Legyen ez a kifejtés  $\mathbf{w} = \sum_{j=1}^N \alpha_j \cdot \varphi(\mathbf{z}_j)$ , ahol  $\alpha_j \in \mathbb{R}$ .*

*Ekkor:*

1. *Explicit eset:*

$$\frac{d[\langle \mathbf{w}, \varphi(\mathbf{g}(\mathbf{s})) \rangle_{\mathcal{H}}]}{d[\mathbf{s}]} = \frac{d[\langle \mathbf{w}, \varphi(\mathbf{u}) \rangle_{\mathcal{H}}]}{d[\mathbf{u}]} \Big|_{\mathbf{u}=\mathbf{g}(\mathbf{s})} \cdot \mathbf{g}'(\mathbf{s}) \quad (30)$$

2. *Implicit eset:*

$$\frac{d[\langle \mathbf{w}, \varphi(\mathbf{g}(\mathbf{s})) \rangle_{\mathcal{H}}]}{d[\mathbf{s}]} = \sum_{j=1}^N \alpha_j \cdot k'_y(\mathbf{z}_j, \mathbf{g}(\mathbf{s})) \cdot \mathbf{g}'(\mathbf{s}) \quad (31)$$

*Bizonyítás.*

1. Explicit eset: az állítás adódik a láncszabályból.

2. Implicit eset:

$$\frac{d[\langle \mathbf{w}, \varphi(\mathbf{g}(\mathbf{s})) \rangle_{\mathcal{H}}]}{d[\mathbf{s}]} = \frac{d[\langle \sum_j \alpha_j \cdot \varphi(\mathbf{z}_j), \varphi(\mathbf{g}(\mathbf{s})) \rangle_{\mathcal{H}}]}{d[\mathbf{s}]} \quad (32)$$

$$= \frac{d[\sum_j \alpha_j \cdot \langle \varphi(\mathbf{z}_j), \varphi(\mathbf{g}(\mathbf{s})) \rangle_{\mathcal{H}}]}{d[\mathbf{s}]} \quad (33)$$

$$= \frac{d[\sum_j \alpha_j \cdot k(\mathbf{z}_j, \mathbf{g}(\mathbf{s}))]}{d[\mathbf{s}]} \quad (34)$$

$$= \sum_j \alpha_j \cdot k'_y(\mathbf{z}_j, \mathbf{g}(\mathbf{s})) \cdot \mathbf{g}'(\mathbf{s}). \quad (35)$$

Az első egyenletben beírtuk  $\mathbf{w}$  kifejtését, majd kihasználtuk a skaláris szorzat linearitását. Ezután a

$$k(\mathbf{u}, \mathbf{v}) = \langle \boldsymbol{\varphi}(\mathbf{u}), \boldsymbol{\varphi}(\mathbf{v}) \rangle_{\mathcal{H}} \quad (36)$$

reprezentáció és kernel közti összefüggést alkalmaztuk. Az utolsó lépés a láncszabályból adódik.

□

Folytatjuk (27) kiszámítását:

1. Explicit eset: Az előző lemma szerint

$$\frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} = \left[ \begin{array}{c} \vdots \\ \left. \frac{d[\langle \mathbf{w}_i^{l+1}(t), \boldsymbol{\varphi}^{l+1}(\mathbf{u}) \rangle_{\mathcal{H}^{l+1}}]}{d[\mathbf{u}]} \right|_{\mathbf{u}=\mathbf{g}^l(\mathbf{s}^l(t))} \\ \vdots \end{array} \right] \cdot (\mathbf{g}^l)'(\mathbf{s}^l(t)) \quad (37)$$

$$= \left[ \begin{array}{c} \vdots \\ \left. \frac{d[\langle \mathbf{w}_i^{l+1}(t), \boldsymbol{\varphi}^{l+1}(\mathbf{u}) \rangle_{\mathcal{H}^{l+1}}]}{d[\mathbf{u}]} \right|_{\mathbf{u}=\mathbf{x}^{l+1}(t)} \\ \vdots \end{array} \right] \cdot [(\mathbf{g}^l)'(\mathbf{s}^l(t))] \quad (38)$$

$$(l = 1, \dots, L-1; i = 1, \dots, N_S^{l+1}).$$

A második egyenlőségnél (i) kihasználtuk a (17) azonosságot és (ii) kiemeltük a  $(\mathbf{g}^l)'(\mathbf{s}^l(t))$  tagot mátrixok szorzásának megfelelően.

2. Implicit eset:  $\mathbf{w}_i^{l+1}(t)$ -kre fennáll a (13) kifejtési tulajdonság. Ez kezdetben feltevésünk volt. A 3.1.3 alfejezetben, látni fogjuk, hogy ez a tulajdonság az iterációk során „öröklődik”. Így

$$\mathbf{w}_i^{l+1}(t) = \sum_{j=1}^{N_i^{l+1}(t)} \alpha_{ij}^{l+1}(t) \cdot \boldsymbol{\varphi}^{l+1}(\mathbf{z}_{ij}^{l+1}(t)) \quad (l = 1, \dots, L-1; i = 1, \dots, N_S^{l+1}) \quad (39)$$

és a kívánt (27) derivált a lemma alkalmazásával

$$\begin{aligned}
& \frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} = \\
& = \begin{bmatrix} \vdots \\ \sum_{j=1}^{N_i^{l+1}(t)} \alpha_{ij}^{l+1}(t) \cdot [k^{l+1}]'_y(\mathbf{z}_{ij}^{l+1}(t), \mathbf{g}^l(\mathbf{s}^l(t))) \cdot (\mathbf{g}^l)'(\mathbf{s}^l(t)) \\ \vdots \end{bmatrix} \quad (40) \\
& = \begin{bmatrix} \vdots \\ \sum_{j=1}^{N_i^{l+1}(t)} \alpha_{ij}^{l+1}(t) \cdot [k^{l+1}]'_y(\mathbf{z}_{ij}^{l+1}(t), \mathbf{x}^{l+1}(t)) \\ \vdots \end{bmatrix} \cdot [(\mathbf{g}^l)'(\mathbf{s}^l(t))] \quad (41) \\
& \quad (l = 1, \dots, L-1; i = 1, \dots, N_S^{l+1}).
\end{aligned}$$

A második egyenlőségénél kihasználtuk a (17) azonosságot. A  $(\mathbf{g}^l)'(\mathbf{s}^l(t))$  mátrix tagot mátrixok szorzásának megfelelően kiemeltük.

Láncszabály és  $\delta^{l+1}(t)$  definíciója alapján

$$\delta^l(t) = \frac{d[\varepsilon^2(t)]}{d[\mathbf{s}^l(t)]} = \frac{d[\varepsilon^2(t)]}{d[\mathbf{s}^{l+1}(t)]} \cdot \frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} = \delta^{l+1}(t) \cdot \frac{d[\mathbf{s}^{l+1}(t)]}{d[\mathbf{s}^l(t)]} \quad (l = 1, \dots, L-1). \quad (42)$$

Ismét alkalmazva a láncszabályt,  $\delta^l(t)$  és  $\mathbf{s}^l(t)$  definíciója szerint

$$\frac{d[\varepsilon^2(t)]}{d[\mathbf{w}_i^l(t)]} = \frac{d[\varepsilon^2(t)]}{d[s_i^l(t)]} \cdot \frac{d[s_i^l(t)]}{d[\mathbf{w}_i^l(t)]} = \delta_i^l(t) \cdot \varphi^l(\mathbf{x}^l(t)) \quad (l = 1, \dots, L; i = 1, \dots, N_S^l), \quad (43)$$

ami a kívánt derivált. Figyeljük meg, hogy a derivált a  $\delta_i^l(t)$  szám és az aktuális  $\mathbf{x}(t)$  input  $l$ . rétegére eső  $\mathbf{x}^l(t)$  lenyomatának  $\varphi^l(\mathbf{x}^l(t))$  reprezentációjával kifejezhető.

### 3.1.2. Regularizációs tag

Ez a tag egyszerűen megadható:

$$\frac{d[r(t)]}{d[\mathbf{w}_i^l(t)]} = \frac{d \left[ \sum_{l=1}^L \sum_{i=1}^{N_S^l} \lambda_i^l \cdot \|\mathbf{w}_i^l(t)\|_{\mathcal{H}^l}^2 \right]}{d[\mathbf{w}_i^l(t)]} = 2\lambda_i^l \cdot \mathbf{w}_i^l(t) \quad (l = 1, \dots, L; i = 1, \dots, N_S^l). \quad (44)$$

Vegyük észre, hogy a derivált az aktuális  $\mathbf{w}_i^l(t)$  súlyok skalárszorosa. Ezen forma szerint implicit hangolási szabály adható.

### 3.1.3. Költség tag

Használva a

$$\frac{d[c(t)]}{d[\mathbf{w}_i^l(t)]} = \frac{d[\varepsilon^2(t)]}{d[\mathbf{w}_i^l(t)]} + \frac{d[r(t)]}{d[\mathbf{w}_i^l(t)]} \quad (l = 1, \dots, L; i = 1, \dots, N_S^l) \quad (45)$$

összefüggést, és az approximációs illetve regularizációs tagokra kapott eredményeinket [(43) és (44) egyenlet] a

$$\mathbf{w}_i^l(t+1) = \mathbf{w}_i^l(t) - \mu_i^l(t) \cdot \frac{d[c(t)]}{d[\mathbf{w}_i^l(t)]} \quad (l = 1, \dots, L; i = 1, \dots, N_S^l) \quad (46)$$

legmeredekebb lejtő szabályban adódik, hogy

$$\mathbf{w}_i^l(t+1) = \mathbf{w}_i^l(t) - \mu_i^l(t) \cdot (\delta_i^l(t) \cdot \boldsymbol{\varphi}^l(\mathbf{x}^l(t)) + 2\lambda_i^l \cdot \mathbf{w}_i^l(t)) \quad (47)$$

$$= (1 - 2\mu_i^l(t) \cdot \lambda_i^l) \cdot \mathbf{w}_i^l(t) - \mu_i^l(t) \cdot \delta_i^l(t) \cdot \boldsymbol{\varphi}^l(\mathbf{x}^l(t)) \quad (48)$$

$$(l = 1, \dots, L; i = 1, \dots, N_S^l).$$

Ugyanez duális formában

$$\boldsymbol{\alpha}_i^l(t+1) = [(1 - 2\mu_i^l(t) \cdot \lambda_i^l) \cdot \boldsymbol{\alpha}_i^l(t); -\mu_i^l(t) \cdot \delta_i^l(t)] \quad (l = 1, \dots, L; i = 1, \dots, N_S^l). \quad (49)$$

Így a hálózat súlyvektorainak kifejtési tulajdonsága [(13) egyenlet] az iterációk során öröklődik. Speciálisan, a számítás végeztével kapott  $\mathbf{w}_i^l$  paraméterekre is fennáll. Összefoglalva, MLK-ra létezik visszaterjesztési eljárás. A levezetett explicit és implicit eljárásokat az 1. és a 2. táblázat foglalja össze.

## 4. Konklúziók

Új többrétegű modell, a Többrétegű Kerceptron (MLK) elméleti leírásával foglalkozunk. Ez a hálózat egyesítheti a Többrétegű Perceptron (MLP) és a Támasztó Vektor Gépek (SVM) előnyeit: (i) Súlyai hangolhatók és a hangolás regularizációs elvek mentén is megtehető. (ii) MLK-ban kernelek használata lehetséges. (iii) Az MLK hálózat behangolt súlyai segítségével a hálózat kimenete gyorsan számolható. (iv) Az MLK *rejtett rétegekkel* rendelkezik és így képes lehet az SVM-ek adta partíciónálásokat kombinálni. A megközelítés különböző adatbázisokon adódó előnyei és hátrányai jövőbeni kutatásaink tárgyát képezi.

## Hivatkozások

- [1] N. Aronszajn: Theory of Reproducing Kernels. 68. évf. (1950), *Trans. of Am. Math. Soc.*, 337–404. p.
- [2] T. Evgeniou – M. Pontil – T. Poggio: Regularization Networks and Support Vector Machines. 13. évf. (2000) 1. sz., *Advances in Computational Mathematics*, 1–50. p.

- [3] S. Haykin: *Neural Networks*. New Jersey, USA, 1999, Prentice Hall.
- [4] R. Herbrich: *Learning Kernel Classifiers*. 2002, MIT Press.
- [5] K.-R. Müller–A. Smola–G. Rätsch–B. Schölkopf–J. Kohlmorgen–V. Vapnik: Predicting Time Series with Support Vector Machines. In *Advances in Kernel Methods* (konferenciaanyag). 1999, MIT Press, 243–254. p.
- [6] B. Schölkopf–A.J. Smola: *Learning with Kernels*. Cambridge, MA, 2002, MIT Press.
- [7] J. Shawe-Taylor–N. Cristianini: *Kernel Methods for Pattern Analysis*. 2004, Cambridge University Press.
- [8] V. Vapnik–S. Golowich–A. Smola: *Support Vector Method for Function Estimation, Regression Estimation and Signal Processing*. Vol. 9. köt. Neural information processing systems. kiad. 1997, MIT Press, Cambridge, MA.
- [9] V.N. Vapnik: *The Nature of Statistical Learning Theory*. 1995, Springer-Verlag New York, Inc.
- [10] V.N. Vapnik: *Statistical Learning Theory*. 1998, Wiley, Chichester, GB.
- [11] G. Wahba: Support Vector Machines, Reproducing Kernel Hilbert Spaces, and Randomized GACV. In *Advances in Kernel Methods* (konferenciaanyag). 1999, MIT Press, 69–88. p.