

# IDA ELŐADÁS II.

Bolgár Bence

2014. október 18.

## I. Bevezetés

Az adatfúzió célja a prediktív teljesítmény javítása több (heterogén) információforrás felhasználásával. Elvárásaink az adatfúziós módszerekkel szemben:

- Használatával javuljon a prediktív teljesítmény.
- Legyen használható heterogén, akár különböző formátumú (pl. vektoriális, gráfós) adatok esetén is.
- Skálázódjon jól (lineárisan) az információforrások számával.
- Legyen lehetőség szakértői tudás integrálására.
- Legyen automatizált, könnyen használható, felhasználóbarát (pl. paraméterezés kérdése).
- Legyen hatékonyan számítható (pl. párhuzamosítható).
- Álljanak rendelkezésre matematikai-statisztikai garanciák a teljesítményt illetően.
- Kezelje a hiányos ill. zajos adatokat.

A fúzió szintje szerint három kategóriát különböztethetünk meg: adatszintű (alacsony szintű, korai), köztes, döntésszintű (magas szintű, késői) fúzió.

## II. Adatszintű fúzió

Adatszintű fúziónál adatok direkt, leírás-szintű kombinációja történik. Ennek legegyszerűbb módja a vektortér-integráció (VSI), amely a vektoriális adatok konkatenációját (egymás után fűzését) jelenti. Ekkor a kombinált adathalmaz

$$\mathcal{X} = \bigoplus_{k=1}^n \{ \mathcal{X}_k \subset \mathbb{R}^{d_k} \}$$

formában írható, ahol  $\mathcal{X}_k$  jelöli az egyes adatforrásokat. Az

$$\mathbf{x} \in \mathcal{X} \cong \mathcal{X}_1 \times \dots \times \mathcal{X}_n \subset \mathbb{R}^{d_1 \dots d_n}$$

kombinált minta a következőképpen áll elő:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix},$$

ahol  $\mathbf{x}_i \in \mathcal{X}_i$ , az  $\mathbf{x}$  minta reprezentációja az  $i$ . adatforrásban. A módszer jellemzői:

- Előnyök: egyszerű, hatékonyan számítható, jól skálázódik, automatizált.
- Hátrányok: nincs lehetőség szakértői tudás integrálására, csak vektoriális adatokra alkalmazható, hiányos adatokat és zajt nem kezel.

## III. Köztes fúzió

A köztes fúzió az adatoknak egy átmeneti reprezentációját használja fel. Ebben a szakaszban néhány példát mutatunk a „közös nyelvre” – a teljesség igénye nélkül.

### III.A. Valószínűségi fúzió

A valószínűségi fúzióra láttunk példát az előző előadásban (modellek kombinálása naiv bayesi alapon); ide tartoznak még a később tárgyalt Bayes-hálók, mixture modellek stb. Megjegyzendő, hogy ezek közül a legtöbb döntésszintű adatfúzióra is alkalmas.

Előnyök:

- Lehetőség *a priori* ismeretek normatív kombinációjára.
- Bizonytalanság, hiány közvetlen kezelése.
- Kimenet: jól értelmezhető valószínűségi állítások.
- Évtizedes tapasztalatok, évszázadok (!) óta finomított matematikai háttér, elméleti garanciák.
- Egyre hatékonyabb algoritmusok.

Hátrányok:

- Nagy számításigény.
- Általában kedvezőtlen skálázódás.
- Nehézkes az ismeretek transzformációja a valószínűségek nyelvére.

### III.B. Hálózatos megközelítések

Rendkívül elterjedt egyes szakterületeken:

- Bioinformatika. A tudás itt több (kapcsolt!) szinten is megjelenik: génszabályozási, fehérje-fehérje interakciós, metabolikus, jelátviteli, betegség–betegség stb. hálózatok. Ezen heterogén tudás fúziója az elmúlt években a bioinformatika egyik legintenzívebben kutatott területévé vált.
- Telekommunikáció. Itt gondolhatunk például az Internet szerkezetének feltárását, hatékony kereső- és ajánlóalgoritmusok implementációját, reklámozási stratégiák kialakítását célzó kutatásokra, amelyeknek gyakran központi eleme az adatfúzió.
- Kapcsolati hálók. A kapcsolati hálók integratív elemzése a legkülönbözőbb területeket érinti: társadalomtudományok, epidemiológia, közgazdaságtan stb.

A hálózatos adatok fúziójára a fenti szakterületeken rengeteg *ad hoc* megoldást dolgoztak ki, amelyek szám-bavétele szinte lehetetlen vállalkozás – standard, „bevált” módszerek ugyanis egyelőre nem léteznek. Gyakran alkalmazzák a hálózatok valamilyen direkt integrációját (pl. tranzitív lezárás), visszavezetést optimalizációs feladatokra, hálózatillesztést, valószínűség-alapú eljárásokat stb.

Előnyök:

- Intenzív kutatás.
- Stabil matematikai alapok (pl. gráfelmélet) és empirikus eredmények.

Hátrányok:

- Konstrukció nem triviális.
- Szakértői tudás bevitele nem triviális.
- Bizonytalanság és hiány kezelése gyakran nehézkes.

### III.C. Kernel fúzió

A Multiple Kernel Learning (MKL) eljárások a hagyományos kernel gépek kiterjesztései, ahol a fúzió a kernel mátrixok szintjén történik. A konkrét módszerek előtt nézzük a közös előnyöket és hátrányokat.

Előnyök:

- Rendszerint jó prediktív teljesítmény.
- Heterogén információforrások jól integrálhatók (amennyiben tudunk hasonlóságokat származtatni).

- Lineáris skálázódás az információforrások (kernelek) számával.
- Lehetőség szakértői tudás integrálására (pl. kernel függvények tervezésével).
- Automatizált algoritmusok.
- Hatékonyan számítható, jól skálázódik nagy mennyiségű adatra is, gyakran párhuzamosítható.
- Matematikai-statisztikai garanciák.

Hátrányok:

- Paraméterezés általában nem triviális.
- Bizonytalanság és hiány kezelése nehézkes.

### III.C.1. Statikus kombináció

A kernelek kombinációjának legegyszerűbb módszerei a „statikus” lineáris vagy nemlineáris kombinációk:

- Kernelek uniform átlagolása:  $\mathbf{K} = \frac{1}{n} \sum_{k=1}^n \mathbf{K}_k$ .
- Kernelek súlyozott átlagolása:  $\mathbf{K} = \sum_{k=1}^n d_k \mathbf{K}_k$ , ahol  $\|\mathbf{d}\| = 1$  valamilyen  $\|\cdot\|$  normában.
- Kernelek szorzása:  $\mathbf{K} = \prod_{k=1}^n \mathbf{K}_k$ .
- Kernelek Hadamard-szorzata:  $\mathbf{K} = \mathbf{K}_1 \circ \mathbf{K}_2 \circ \dots \circ \mathbf{K}_n$ , ahol  $(\mathbf{A} \circ \mathbf{B})_{ij} = A_{ij} B_{ij}$ .

És még számos más eljárás. Egyszerűségük ellenére általában jó prediktív teljesítményt nyújtanak.

### III.C.2. Adaptív kombináció

A kernel gép tanítása és az adatfúzió egy lépésben is történhet, ha az optimalizációba a lineáris kombináció súlyainak tanulását is beépítjük. Ennek előnye, hogy a fúzió ekkor „adaptívan”, mindig az aktuális feladatot figyelembe véve történik. A hagyományos SVM primál feladatát a következőképpen módosítjuk:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \mathbf{d}} \quad & \frac{1}{2} \sum_k \|\mathbf{w}_k\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i \left( \sum_k \mathbf{w}_k^T \left( \sqrt{d_k} \phi_k(\mathbf{x}_i) \right) + b \right) \geq 1 - \xi_i, \\ & \|\mathbf{d}\|_p = 1, \quad \xi_i \geq 0, \quad d_k \geq 0. \end{aligned}$$

A fenti feladatban a  $\mathbf{w} \leftarrow \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_n \end{bmatrix}$  és a  $\phi(\mathbf{x}_i) \leftarrow \begin{bmatrix} \sqrt{d_1} \phi_1(\mathbf{x}_i) \\ \vdots \\ \sqrt{d_n} \phi_n(\mathbf{x}_i) \end{bmatrix}$  helyettesítéssel éltünk. Vegyük észre, hogy ez

nem más, mint a vektortér-integráció azzal a különbséggel, hogy az egyes reprezentációk meg vannak szorozva egy  $\sqrt{d_k}$  súllyal! A hagyományos SVM-hez nagyon hasonló gondolatmenetet követve, valamint a skalárszorzat tulajdonságait felhasználva a  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$  összefüggés a következőképpen módosul:

$$\begin{bmatrix} \sqrt{d_1} \phi_1(\mathbf{x}_i) \\ \vdots \\ \sqrt{d_n} \phi_n(\mathbf{x}_i) \end{bmatrix}^T \begin{bmatrix} \sqrt{d_1} \phi_1(\mathbf{x}_j) \\ \vdots \\ \sqrt{d_n} \phi_n(\mathbf{x}_j) \end{bmatrix} = \sum_k d_k k_k(\mathbf{x}_i, \mathbf{x}_j),$$

azaz a kernelek lineáris kombinációjához jutunk, ahol az algoritmus egyben az optimális súlyozást is meg fogja adni. Külön figyelmet érdekel a  $\|\mathbf{d}\|_p = 1$  kényszer – enélkül a  $d_k$  súlyokat minden határon túl növelve a célfüggvény tetszőlegesen csökkenthető lenne. Az  $\|\mathbf{a}\|_p = (\sum_l a_l^p)^{1/p}$  normát  $L_p$ -normának nevezzük, és  $p$  különböző értékeire más-más jellegű regularizációt kapunk:  $p < 2$  esetén kevés súly fog magas értéket kapni (ritka kombináció, „legjobb” kernelek kiválasztása), nagyobb  $p$  esetén egyenletesebb lesz a súlyok eloszlása.

### III.C.3. Kernel alignment

Tekintsünk egy hagyományos kétosztályos osztályozási feladatot, azaz legyen  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^P$ , ahol  $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d$ ,  $y_i \in \mathcal{Y} = \{-1, +1\}$ . Könnyen megmutatható, hogy ekkor az  $\mathbf{y}\mathbf{y}^T$  szorzat egy ideális kernel, amely minden tanítómintát helyesen osztályoz:

$$\begin{aligned} y_i = y_j &\Leftrightarrow (\mathbf{y}\mathbf{y}^T)_{ij} = 1, \\ y_i \neq y_j &\Leftrightarrow (\mathbf{y}\mathbf{y}^T)_{ij} = -1. \end{aligned}$$

[Gondolhatunk arra, hogy a kernel függvény felfogható úgy is, mint egy hasonlóságfüggvény]. Legyen most adott  $n$  darab különböző kernelünk,  $\mathbf{K}_1, \dots, \mathbf{K}_n$ , amelyek különböző információforrásokból származnak (bázis-kernel). Keressük azt a  $\mathbf{K}$  pozitív szemidefinit kernelt, amely a lehető „legközelebb” van a báziskernelnek, valamint az ideális kernelhez:

$$\begin{aligned} \min_{\mathbf{K}} \quad & L(\mathbf{K}, \mathbf{K}_1, \dots, \mathbf{K}_n, \mathbf{y}\mathbf{y}^T) \\ \text{s.t.} \quad & \mathbf{K} \succeq 0, \end{aligned}$$

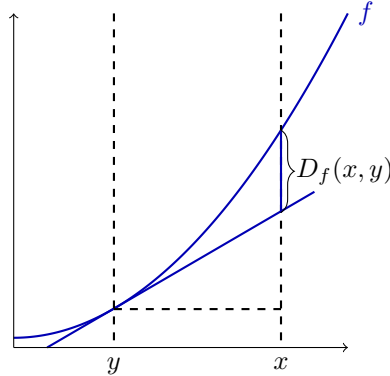
ahol  $L$  valamilyen veszteségfüggvény. A továbbiakban az egyszerűség kedvéért egy kernelen mutatjuk be a veszteségfüggvényeket. Mátrix-közelségi problémáknál gyakran használt függvény a Frobenius-norma:

$$\begin{aligned} \min_{\mathbf{K}} \quad & \|\mathbf{K} - \mathbf{y}\mathbf{y}^T\|_F \\ \text{s.t.} \quad & \mathbf{K} \succeq 0, \end{aligned}$$

ahol  $\|A\|_F = \sqrt{\sum_i \sum_j |A_{ij}|^2}$ . A Frobenius-norma hátránya, hogy nem garantálja a kapott mátrix pozitív szemidefinit tulajdonságát. Vegyük hát a Frobenius-norma egy általánosítását, amelyet Bregman-divergencia néven ismerünk. A Bregman-divergencia definíciója valós értékű függvényekre:

$$D_f(x, y) := f(x) - f(y) - f'(y)(x - y),$$

ahol  $f : \mathbb{R} \rightarrow \mathbb{R}$ , azaz személetesen a következőt jelenti:



Megjegyezzük, hogy az így definiált divergencia általánosságban nem szimmetrikus és nem teljesíti a háromszög-egyenlőtlenséget. A Bregman-divergencia triviálisan kiterjeszhető a többdimenziós esetre, azaz legyen most  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , és

$$D_f(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle.$$

**1. Példa.** Legyen  $f(\mathbf{x}) = \|\mathbf{x}\|^2$ . Ekkor  $D_f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{y}, \mathbf{y} \rangle - \langle 2\mathbf{y}, \mathbf{x} - \mathbf{y} \rangle = \|\mathbf{x} - \mathbf{y}\|^2$ , ami éppen az euklideszi távolság.

Végül definiáljuk a Bregman-divergenciát mátrixokra is,  $f : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ .

$$D_f(\mathbf{X}, \mathbf{Y}) := f(\mathbf{X}) - f(\mathbf{Y}) - \langle \nabla f(\mathbf{Y}), (\mathbf{X} - \mathbf{Y}) \rangle_F,$$

ahol  $\langle \cdot, \cdot \rangle_F$  a Frobenius-szorzat, azaz  $\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}\mathbf{B}^T) = \sum_i \sum_j A_{ij}B_{ij}$ .

**2. Példa.** Könnyen belátható, hogy  $f(\mathbf{X}) = \|\mathbf{X}\|_F^2 \rightsquigarrow D_f(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\|_F^2$ .

**3. Példa.** Legyen most  $f(\mathbf{X}) = -\log \det \mathbf{X}$ . Ekkor

$$\begin{aligned} D_f(\mathbf{X}, \mathbf{Y}) &= -\log \det \mathbf{X} + \log \det \mathbf{Y} - \langle (\det \mathbf{Y})^{-1} \text{adj } \mathbf{Y}, \mathbf{X} - \mathbf{Y} \rangle_F \\ &= -\log \det \mathbf{X}\mathbf{Y}^{-1} - n + \text{tr}(\mathbf{X}\mathbf{Y}^{-1}), \end{aligned}$$

ahol  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}$ .

A 3. példában látott divergenciát LogDet divergenciának is nevezik. A kernel alignment eljárások szempontjából igen fontosak az alábbi tulajdonságai:

1. Garantálja a pozitív szemidefinit tulajdonságot.
2. Megőrzi a rangot.
3. Első argumentumában konvex.

Az 1. tulajdonság miatt az optimalizációból elhagyhatjuk a  $\mathbf{K} \succeq 0$  kényszert. A 2. tulajdonság lehetővé teszi az ún. alacsony rangú approximációt, amellyel hatékonyan kezelhetők nagy méretű adathalmazok is, míg a 3. tulajdonság az optimalizációs stratégia kidolgozását könnyíti meg.

#### IV. Döntésszintű fúzió

A döntésszintű fúzió során a forrásonként külön-külön elvégzett elemzések (pl. osztályozás, regresszió) eredményeit egyesítjük. A korábban vett boosting, MoE, szavazásos eljárások stb. variánsai mind felfoghatók késői adatfúziós eljárásként. Itt csak egy speciális esetet, a sorrendi fúziót fogjuk megvizsgálni.

Előnyök:

- Rugalmas, gyakorlatilag bármi kombinálható; forrásonként akár különböző algoritmusokat is használhatunk.
- Előbbi egyben a szakértői tudás bevitelét is segíti.

Hátrányok:

- Könnyen számításigényessé válhat.
- A fúziót nem (közvetlenül) az adat vezérli, mivel az a döntések szintjén egyáltalán nem jelenik meg.

Sorrendi fúziót gyakran használnak információ-visszakeresési feladatokban, azaz pl. webes keresőmotorokban, ajánló rendszerekben, vagy a bio- és kemoinformatika területén (génprioritizálás, virtuális screening). A sorrendekben gyakran nem csupán egy rendezési reláció áll rendelkezésre, hanem egzakt pontszámok is, amelyeket egyfajta relevanciaként is értelmezhetünk (ti. a lekérdezés szempontjából).

Néhány egyszerű stratégia a sorrendi fúzióra:

- **Sum rank.** Adott entitás rangja a különböző sorrendekben elért rangjainak összege alapján számolható.
- **Sum score.** Adott entitás rangja a különböző sorrendekben található pontszámainak összege alapján számolható.
- **Pareto ranking.** Adott entitás rangja az alapján alakul, hogy hány entitás ért el nála jobb rangot *minden* sorrendben.
- **Rank vote.** Minden sorrend szavaz az első  $k$  elemére. A végső sorrend a kapott szavazatok alapján alakul.
- **Parallel selection.** Minden sorrendből kiválasztjuk a soron következő elemet és beillesztjük a végső sorrendbe. Ha olyat választunk, ami már bekerült, akkor helyette a következőt választjuk.

A döntetleneket a sum rank vagy sum score módszerekkel szokás feloldani.

**4. Példa.** Tekintsük az alábbi sorrendeket:

<i>A</i>	10	<i>B</i>	10	<i>A</i>	8
<i>B</i>	8	<i>A</i>	6.5	<i>D</i>	7
<i>C</i>	7	<i>C</i>	3	<i>B</i>	7
<i>D</i>	6.5	<i>D</i>	2	<i>C</i>	6
<i>E</i>	6	<i>E</i>	1	<i>E</i>	5

Az egyes módszerekkel kialakuló kombinált sorrendek:

- *Sum rank*:  $A(4), B(6), C(10), D(10), E(15)$ .
- *Sum score*:  $B(25), A(24.5), C(16), D(15.5), E(12)$ .
- *Pareto ranking*:  $A(0), B(0), D(1), C(2), E(4)$ .
- *Rank vote* ( $k = 3$ ):  $A(3), B(3), C(2), D(1), E(0)$ .
- *Rank vote* ( $k = 2$ ):  $A(3), B(2), D(1), C(0), E(0)$ .
- *Parallel selection* (balról jobbra):  $A, B, D, C, E$ .