

Hash-alapú keresések az információ-visszakeresésben

Bodó Zalán

Babeş–Bolyai Tudományegyetem
Magyar Matematika és Informatika Intézet

Információ-
visszakeresés

Hash-alapú keresések

Locality-Sensitive
Hashing (LSH)

Spektrális hashing

Spektrális klaszterezés

Lineáris spektrális
hashing

Kísérletek, eredmények

Bibliográfia

**A Magyar Tudomány Napja Erdélyben
Kolozsvár, 2012**

Tartalom

Információ-visszakeresés

Hash-alapú keresések

Locality-Sensitive Hashing (LSH)

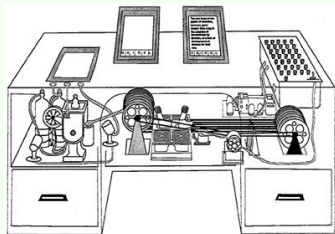
Spektrális hashing

Spektrális klaszterezés

Lineáris spektrális hashing

Kísérletek, eredmények

Bibliográfia



- ▶ **információ-visszakeresés** (*information retrieval*) = információ visszatérítése egy adathalmazból, amely releváns az adott információk igény/kérés szempontjából
- ▶ alapötlet: Vannevar Bush: *As We May Think* (1945)

"It consists of a desk, and while it can presumably be operated from a distance, it is primarily the piece of furniture at which he works. On the top are slanting translucent screens, on which material can be projected for convenient reading. There is a keyboard, and sets of buttons and levers. Otherwise it looks like an ordinary desk."

<http://www.theatlantic.com/magazine/archive/1945/07/as-we-may-think/303881/>

- ▶ az adatok nem strukturáltak; nem adatbázis-lekérdezést hajtunk végre
- ▶ általában nem *egyetlen* találatot térítünk vissza, hanem **az első k legközelebbit**, vagy akár az **összes illeszkedőt**
- ▶ a találatokat **rangsoroljuk**
- ▶ az adatok legtöbbször szövegek és képek
- ▶ **probléma**: be kell járni az **egész** adatbázist és össze kell hasonlítani a kéréssel

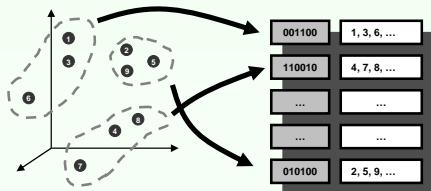
Példa: *Vektortér modell*

A szövegeket tárolhatjuk *bag-of-words*-ként, azaz szavak bag-jeként (multihalmazként) vagy ekvivalens módon vektorokként. A tér dimenzióit a (valamilyen speciális módon kiválasztott) indexelő szavak adják.

Két dokumentum hasonlóságát kiszámíthatjuk a vektorok közötti szög koszinuszaként:

$$s(\mathbf{d}_i, \mathbf{d}_j) = \frac{\mathbf{d}'_i \mathbf{d}_j}{\|\mathbf{d}_i\| \|\mathbf{d}_j\|}$$

Hash-alapú keresések



- ▶ alapötlet: a pontokat tegyük egy **hash táblába** úgy, hogy *közeli* pontok *azonos* vagy egymáshoz *közeli kulcsú* rekeszbe kerüljenek
- ▶ kulcs = bináris szekvencia
- ▶ a Hamming-távolságot “egyszerű” kiszámítani
- ▶ legyen $\mathbf{u}, \mathbf{v} \in \{0, 1\}^n$; ekkor

$$d_H(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^n (\mathbf{u} \text{ XOR } \mathbf{v})$$

Locality-Sensitive Hashing (LSH)

ALG 1 LSH

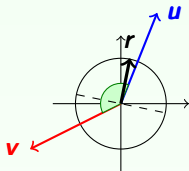
- 1: Generáljunk r darab egységnyi hosszúságú (egyenletes eloszlású) d -dimenziós vektort: $\{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_r\}$
- 2: Az r darab hash függvényünk legyen

$$h_i(\mathbf{u}) = \begin{cases} 1, & \mathbf{r}_i' \mathbf{u} \geq 0 \\ 0, & \mathbf{r}_i' \mathbf{u} < 0 \end{cases}, \quad i = 1, 2, \dots, r$$

- 3: Egy \mathbf{x} ponthoz rendelt bináris hash szekvencia

$$(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_r(\mathbf{x}))$$

Miért működik az LSH?



Biz.

- ▶ u, v – vektorok \mathbb{R}^d -ben
- ▶ r – véletlen vektor \mathbb{R}^d -ben,
 $\|r\| = 1$

Ha a hash függvény

$$h_r(u) = \begin{cases} 1, & r'u \geq 0 \\ 0, & r'u < 0 \end{cases}$$

akkor

$$P(h_r(u) = h_r(v)) = 1 - \frac{\theta(u, v)}{\pi}$$

$$\begin{aligned} &P(h_r(u) \neq h_r(v)) \\ &= P(r'u \geq 0, r'v < 0) \\ &\quad + P(r'u < 0, r'v \geq 0) \\ &= \frac{\theta(u, v)}{2\pi} + \frac{\theta(u, v)}{2\pi} \\ &= \frac{\theta(u, v)}{\pi} \end{aligned}$$



Spektrális hashing

- ▶ $\mathbf{y}_i \in \{-1, 1\}^r$, $i = 1, \dots, N$ kódszavakat keresünk a pontokhoz úgy, hogy egymáshoz hasonló pontok esetén ezen kódszavak távolsága kicsi legyen:

$$\min_{\{\mathbf{y}_i\}_{i=1, \dots, N}} \sum_{i, j=1}^N W_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|^2$$

$$\text{ú.h.} \quad \mathbf{y}_i \in \{-1, 1\}^r, \quad \sum_{i=1}^N \mathbf{y}_i = \mathbf{0}, \quad \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i \mathbf{y}_i' = \mathbf{I}$$

- ▶ 2. feltétel: a bitek "kiegyensúlyozottak" legyenek
- ▶ 3. feltétel: a bitek ne korreláljanak egymással
- ▶ a feladatot relaxáljuk:

$$\min_{\{\mathbf{y}_i\}_{i=1, \dots, N}} \text{tr}(\mathbf{Y}'\mathbf{L}\mathbf{Y})$$

$$\text{ú.h.} \quad \mathbf{Y}'\mathbf{1} = \mathbf{0}, \quad \mathbf{Y}'\mathbf{Y} = \mathbf{I}$$

- ▶ a feladat **megoldása**: $L = D - W$ első r darab (legkisebb sajátértékű) sajátvektora, $Y = (\mathbf{v}_2 \quad \mathbf{v}_3 \quad \dots \quad \mathbf{v}_{r+1})$
- ▶ az i -edik kódszó ekkor : $\mathbf{y}'_i = (v_{2i} \quad v_{3i} \quad \dots \quad v_{(r+1)i})$
- ▶ az eredeti cikkben az általánosítást úgy oldották meg, hogy feltételezték a Gauss-kernel használatát, a többdimenziós egyenletes eloszlást, és felhasználták a Laplace–Beltrami operátorok sajátfüggvényeire vonatkozó eredményeket (approximáció)

- ▶ **normalizált minimális vágás** gráfokban:

$$\frac{\sum_{i \in A, j \in \bar{A}} W_{ij}}{\sum_{i \in A, j \in V} W_{ij}} + \frac{\sum_{i \in \bar{A}, j \in \bar{A}} W_{ij}}{\sum_{i \in \bar{A}, j \in V} W_{ij}}$$

- ▶ átalakítható az alábbi numerikus optimalizálási feladatra (relaxáció):

$$\begin{aligned} \max_{\mathbf{z}} \quad & \frac{\mathbf{z}' \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2} \mathbf{z}}{\mathbf{z}' \mathbf{z}} \\ \text{ú.h.} \quad & \mathbf{z}' \mathbf{D}^{1/2} \mathbf{1} = 0 \end{aligned}$$

- ▶ kimutatható, hogy a spektrális klaszterezés egy **hipersíkkal** “vágja” el az adatokat ($\mathbf{W} = \Phi' \Phi$)

$$\begin{aligned} \max_{\mathbf{w}} \quad & \|\mathbf{w}' \Phi \mathbf{D}^{-1/2}\| \\ \text{ú.h.} \quad & \sum_{i=1}^N \mathbf{w}' \phi(\mathbf{x}_i) = 0, \quad \|\mathbf{w}\| = 1 \end{aligned}$$

- ▶ ez felírható a következőképpen:

$$\max_{\mathbf{w}} \frac{\mathbf{w}' \Phi \mathbf{D}^{-1} \Phi' \mathbf{w}}{\mathbf{w}' \mathbf{w}}$$

$$\text{ú.h.} \quad \sum_{i=1}^N \mathbf{w}' \phi(\mathbf{x}_i) = 0$$

- ▶ ennek megoldása a $\Phi \mathbf{D}^{-1} \Phi'$ második legnagyobb sajátértékű sajátvektora (\mathbf{u}_2)
- ▶ $\mathbf{D}^{-1/2} \Phi' \Phi \mathbf{D}^{-1/2}$ 2. legnagyobb sajátértékű sajátvektora \mathbf{v}_2 ; a kapcsolat a kettő között: $\mathbf{u}_2 = \Phi \mathbf{D}^{-1/2} \mathbf{v}_2 \mathbf{e}_2^{-1}$
- ▶ ezáltal egy induktív klaszterezőt nyerünk \mathbf{w} normálissal:

$$\mathbf{w} = \mathbf{u}_2 = \Phi \mathbf{D}^{-1/2} \mathbf{v}_2 \mathbf{e}_2^{-1}$$

- ▶ ekkor egy új pont klasztercímkejét a következő módon határozzuk meg:

$$f_2(\mathbf{x}) = \phi(\mathbf{x})' \mathbf{w} = \phi(\mathbf{x})' \mathbf{u}_2 = \phi(\mathbf{x})' \Phi \mathbf{D}^{-1/2} \mathbf{v}_2 \mathbf{e}_2^{-1}$$

- ▶ megmutatható, hogy

$$\begin{aligned} \text{sgn}(\Phi' \mathbf{w}) &= \text{sgn}(\Phi' \Phi \mathbf{D}^{-1/2} \mathbf{v}_2 \mathbf{e}_2^{-1}) = \text{sgn}(\mathbf{D}^{-1/2} \Phi' \Phi \mathbf{D}^{-1/2} \mathbf{v}_2 \mathbf{e}_2^{-1}) \\ &= \text{sgn}(\mathbf{v}_2 \mathbf{e}_2 \mathbf{e}_2^{-1}) = \text{sgn}(\mathbf{v}_2) \end{aligned}$$

Lineáris spektrális hashing

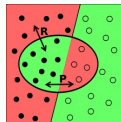
- ▶ a gond az **általánosítás** új pontokra
- ▶ **alapötlet**: használjuk fel az hipersíkkal való szétválasztást a spektrális klaszterezéstől
- ▶ cseréljük ki a Laplace-mátrixot a normalizált Laplace-mátrixra a spektrális hashingben: $D^{-1/2}LD^{-1/2}$
- ▶ felhasználva az általánosítási eredményeket a spektrális klaszterezéstől kapjuk, hogy

$$f(\mathbf{x}) = (f_2(\mathbf{x}) \quad f_3(\mathbf{x}) \quad \dots \quad f_{r+1}(\mathbf{x}))'$$

- ▶ **probléma**: $\phi(\mathbf{x})'\Phi$ kiszámításához N kernel számolásra van szükség
- ▶ **megoldás**: használjuk a pontokat az input térben \Rightarrow kernel = skalárszorzat
- ▶ ekkor kiszámítjuk az első r sajátvektort, $\mathbf{g}_k := \mathbf{u}_k$ vagy $\mathbf{g}_k := \mathbf{X}D^{-1/2}\mathbf{v}_k$, $k = 2, \dots, r + 1$
- ▶ új pont esetén a ponthoz rendelt kódszót a $\mathbf{x}'\mathbf{g}_k$, $k = 2, \dots, r + 1$ skalárszorzatokkal számítjuk ki

Kísérletek, eredmények

- ▶ tesztelés a Reuters-21578¹ és 20Newsgroups² szövegtörzseken
- ▶ szövegek feldolgozása:
 - ▶ stopszavak kiszűrése (199 db.)³
 - ▶ az első 5000 leggyakoribb szó használata dimenziókként
 - ▶ **bag-of-words** reprezentáció + **tf-idf** súlyozás
 - ▶ dokumentumvektorok normalizálása
- ▶ nem használtuk a címkéket
- ▶ a tesztalmanax minden dokumentumához megkerestük az első 50 legközelebbit a tanuló adatokból
- ▶ ezek megtalálásának sikerességét mértük **precision** és **recall** segítségével (pl. 32 bit esetén $\{0, 1, 2, 4, \dots, 32\}$ Hamming-távolságokra számítottunk recallt és precisiót)



¹<http://disi.unitn.it/moschitti/corpora.htm>

²<http://qwone.com/~jason/20Newsgroups/20news-bydate.tar.gz>

³WordNet::Similarity 2.05, <http://wn-similarity.sourceforge.net/>

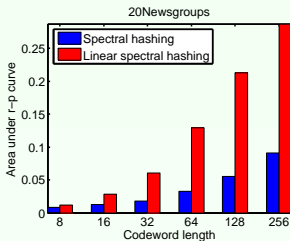
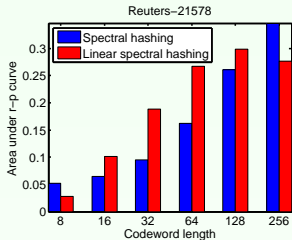


Figure: A recall-precision görbe alatti terület a hash szekvencia hosszának függvényében.

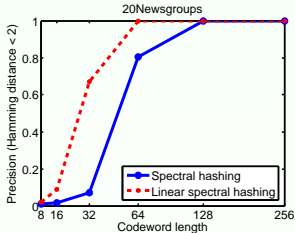
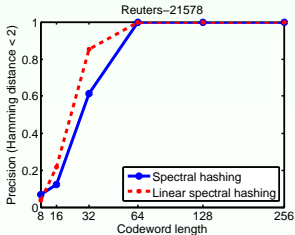


Figure: Precison eredmények Hamming-távolság < 2 esetén.

Bibliográfia

- [Weiss, Torralba, Fergus 2008] Yair Weiss, Antonio B. Torralba, and Robert Fergus. Spectral hashing. In NIPS, pages 1753–1760. MIT Press, 2008.
- [Rahimi, Recht 2004] Ali Rahimi and Ben Recht. Clustering with normalized cuts is clustering with a hyperplane. In Statistical Learning in Computer Vision, 2004.

Köszönöm a figyelmet!