

# SZÁMÍTÓGÉPES SZEMANTIKA

## VECTOR SEMANTICS AND EMBEDDINGS

---

Ferenczi Zsanett

2020. április 27.

1. Bevezetés
2. Lexikális szemantika
3. Szavak és vektorok
4. Vektorok hasonlósága
5. tf-idf
6. PPMI
7. Skip-gram
8. Szóbeágyazási modellek vizualizációja
9. A szóbeágyazások szemantikai jellemzői
10. Elfogultság, sztereotípia
11. Vektoros modellek kiértékelése

# Bevezetés

---

- a szavakat ábrázolhatjuk atomi elemekként, de így elveszítjük a kapcsolatokat közöttük
  - szavak hasonlósága fontos bizonyos feladatok esetén
  - ötlet: a szavakat próbáljuk **jelentésük** alapján kódolni, ezáltal a hasonlóság is megragadható lesz két szó között
  - 1950-es évek: **disztribúciós hipotézis**: a szinonimák gyakran hasonló környezetben fordulnak elő
  - pl. *oculist* és *eye-doctor* gyakran állnak együtt az *eye*, *examined* szavakkal
- a hasonló környezetben álló szavak hasonló jelentéssel bírnak

## Tesgüino?

A bottle of **tesgüino** is on the table

Everybody likes **tesgüino**

**Tesgüino** makes you drunk

We make **tesgüino** out of corn.

*"You shall know a word by the company it keeps!" (Firth (1957))*

## Lexikális szemantika

---

## egér (főnév)

1. nagy szemű és fülű, hegyes orrú rágcsáló...
2. számítógép kézi vezérlőeszköze...

- **lemma (címszó):** egér
- **szóalakok:** egér, egerek, egérnek, stb.
- **jelentés (word sense):** itt a két definíció adja meg az egyes jelentéseket
- **poliszémia:** egy szónak több jelentése van → jelentés egyértelműsítés (WSD)

# Jelentések közötti kapcsolatok

- **szinonímia**: egyik szó jelentése közel azonos egy másik szó jelentésével
- **hasonlóság (word similarity)**: szavak között állhat fenn, pl. *kutya* és *macska*
- **word relatedness**: szavak közötti egyéb kapcsolat
  - **szemantikai mező**: olyan szavak halmaza, amelyek egy szemantikai domént fednek le (pl. *pincér*, *étlap*, *szakács*, stb.)
  - ilyen még: **hiperonímia**, **antonímia**, **meronímia** (ld. 19. fejezet)
- **szemantikai keret, szerepek**: olyan szavak, szereplők halmaza, melyek egy eseményhez köthetők
- **konnotáció**: pl. negatív konnotáció: *szomorú*, pozitív konnotáció: *boldog*



# Vector Semantics

|            | Valence | Arousal | Dominance |
|------------|---------|---------|-----------|
| courageous | 8.05    | 5.5     | 7.38      |
| music      | 7.67    | 5.57    | 6.5       |
| heartbreak | 2.45    | 5.65    | 3.58      |
| cub        | 6.71    | 3.95    | 4.24      |
| life       | 6.68    | 5.59    | 5.89      |

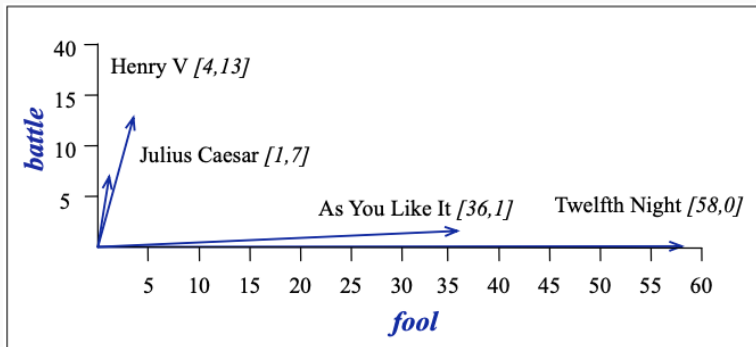
- egy szó jelentése ábrázolható pontként a térben (Osgood et al. (1957))
- hasonló környezetben előforduló szavak hasonló jelentésűek (egy szó ábrázolása a körülötte előforduló szavak számlálásával)
- ezen két megfigyelést ötvözi a vector semantics
- **szóbeágyazás**: olyan vektor, amely egy szót reprezentál

# Szavak és vektorok

---

- a vektorok általában együttes előfordulási mátrixon alapulnak (co-occurence matrix)
- pl. **term-document matrix**: minden sor egy szót jelöl, minden oszlop egy dokumentumot
- V. Henriket a [13, 89, 4, 3] vektorral lehetne ábrázolni
- a vektortér dimenziója ebben az esetben 4

|        | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|--------|----------------|---------------|---------------|---------|
| battle | 1              | 0             | 7             | 13      |
| good   | 114            | 80            | 62            | 89      |
| fool   | 36             | 58            | 1             | 4       |
| wit    | 20             | 15            | 2             | 3       |



- Shakespeare 4 darabja két dimenzióan ábrázolva (4 dokumentumvektor)
- egy term-document mátrix annyi dimenziós lenne, ahány type van a dokumentumokban ( $|V|$  sor)

# Term-term mátrix

- **term-term (word-word vagy term-context) mátrix:**  
sorok és oszlopok is szavakat jelölnek
- a dimenziója:  $|V| \times |V|$  (a  $|V|$  általában 10 000 és 50 000 közötti)
- sor: **célszó** (target word), oszlop: **kontextus (szavak)**
- az egyes cellák azt jelölik, hogy a célszó hányszor fordult elő a kontextusszó környezetében
- a környezet lehet egy dokumentum, de lehet kisebb egység is, pl. a szó körüli **ablak** (a célszótól jobbra és balra 4-4 szó)

|             | aardvark | ... | computer | data | result | pie | sugar | ... |
|-------------|----------|-----|----------|------|--------|-----|-------|-----|
| cherry      | 0        | ... | 2        | 8    | 9      | 442 | 25    |     |
| strawberry  | 0        | ... | 0        | 0    | 1      | 60  | 19    |     |
| digital     | 0        | ... | 1670     | 1683 | 85     | 5   | 4     |     |
| information | 0        | ... | 3325     | 3982 | 378    | 5   | 13    |     |

# Vektorok hasonlósága

---

# Vektorok hasonlósága

- két vektor közötti hasonlóság mérése → **skaláris szorzattal**
- nagy lesz, ha két vektor ugyanazon dimenzióinak értékei nagyok
- 0 pedig, ha egyáltalán nem hasonlók

$$\text{dot product}(\mathbf{v}, \mathbf{w}) = \mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^N v_i w_i = v_1 w_1 + v_2 w_2 + \dots + v_N w_N$$

# Vektorok hasonlósága

- a skaláris szorzat előnyben részesíti a hosszabb vektorokat, így a gyakoribb szavak magasabb értékeket kapnak, míg a kevésbé gyakori szavakhoz nehéz hasonlót találni → ez probléma
- egy megoldás: elosztjuk a vektorok hosszával → **a bezárt szög koszinusza** (0 és 1 közötti szám)
- vektor hossza:

$$|v| = \sqrt{\sum_{i=1}^N v_i^2}$$

$$\frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}||\mathbf{b}|} = \cos\theta$$



|             | pie | data | computer |
|-------------|-----|------|----------|
| cherry      | 442 | 8    | 2        |
| digital     | 5   | 1683 | 1670     |
| information | 5   | 3982 | 3325     |

$$\cos(\text{cherry}, \text{information}) = \frac{442 * 5 + 8 * 3982 + 2 * 3325}{\sqrt{442^2 + 8^2 + 2^2} \sqrt{5^2 + 3982^2 + 3325^2}} = .017$$

$$\cos(\text{digital}, \text{information}) = \frac{5 * 5 + 1683 * 3982 + 1670 * 3325}{\sqrt{5^2 + 1683^2 + 1670^2} \sqrt{5^2 + 3982^2 + 3325^2}} = .996$$

tf-idf

---

- a sokszor együtt előforduló szavak fontosabbak, mint a csak néhányszor előfordulók
- de a túl gyakori szavak nem segítenek: a `good` minden dokumentumban nagyjából ugyanannyiszor fordul elő (ld. 7. dia)
- ezt egyensúlyozni kell → **tf-idf algoritmus**
- **term frequency (tf)**:  $tf_{t,d}$  =  $t$  szó gyakorisága  $d$  dokumentumban

$$tf_{t,d} = \log_{10}(\text{count}(t, d) + 1)$$

|        | Collection Frequency | Document Frequency |
|--------|----------------------|--------------------|
| Romeo  | 113                  | 1                  |
| action | 113                  | 31                 |

- **document frequency (df)**:  $df_t$  = hány dokumentumban fordul elő  $t$  szó
- **collection frequency**: hányszor szerepel  $t$  szó az összes dokumentumban
- az **Romeo**-hoz hasonló szavaknak nagyobb súlyt ad az **idf**
- **inverse document frequency (idf)**: a kevesebb dokumentumban előforduló szavak előnyben részesítése ( $N$  = dokumentumok száma)

$$idf_t = \log_{10}\left(\frac{N}{df_t}\right)$$

- a **tf-idf** ezek szorzata:

$$w_{t,d} = tf_{t,d} \times idf_t$$

- **tf-idf model:**
  - szavak vektorokként való ábrázolása
  - annyi dimenzióval, ahány type van a dokumentumokban
  - az előfordulások **tf-idf**-fel súlyozva
  - vektorok hasonlósága: koszinusszal

|               | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|---------------|----------------|---------------|---------------|---------|
| <b>battle</b> | 0.074          | 0             | 0.22          | 0.28    |
| <b>good</b>   | 0              | 0             | 0             | 0       |
| <b>fool</b>   | 0.019          | 0.021         | 0.0036        | 0.0083  |
| <b>wit</b>    | 0.049          | 0.044         | 0.018         | 0.022   |

- dokumentumok hasonlósága:

$$d = \frac{w_1 + w_2 + \dots + w_k}{k}$$

PPMI

---

- **PMI (pointwise mutual information)**: két szó mennyivel gyakrabban fordulnak elő együtt, mint ha függetlenek lennének
- **w** - célszó, **c** - kontextus(szó)

$$PMI(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)}$$

- ez  $-\infty$  és  $+\infty$  közötti eredményt ad
- ahhoz, hogy valami kevesebbszer forduljon elő, mint várnánk, hatalmas korpusz kellene → **Positive PMI**
- minden negatív értéket 0-ra cserélünk

$$PPMI(w, c) = \max(\log_2 \frac{P(w, c)}{P(w)P(c)}, 0)$$

|                       | computer | data | result | pie | sugar | count(w) |
|-----------------------|----------|------|--------|-----|-------|----------|
| cherry                | 2        | 8    | 9      | 442 | 25    | 486      |
| strawberry            | 0        | 0    | 1      | 60  | 19    | 80       |
| digital               | 1670     | 1683 | 85     | 5   | 4     | 3447     |
| information           | 3325     | 3982 | 378    | 5   | 13    | 7703     |
| <b>count(context)</b> | 4997     | 5673 | 473    | 512 | 61    | 11716    |

$$P(\text{information}, \text{data}) = \frac{3982}{11716} = 0.3399$$

$$P(\text{information}) = \frac{7703}{11718} = 0.6575$$

$$P(\text{data}) = \frac{5673}{11716} = 0.4842$$

$$\text{ppmi}(\text{information}, \text{data}) = \log_2(0.3399 / (0.6575 * 0.4842)) = 0.0944$$

|             | computer | data | result | pie  | sugar |
|-------------|----------|------|--------|------|-------|
| cherry      | 0        | 0    | 0      | 4.38 | 3.30  |
| strawberry  | 0        | 0    | 0      | 4.10 | 5.51  |
| digital     | 0.18     | 0.01 | 0      | 0    | 0     |
| information | 0.02     | 0.09 | 0.28   | 0    | 0     |



## Skip-gram

---

# Skip-gram

- a **word2vec** egyik algoritmus (CBOW mellett)
- a vektorok eddig hosszúak és ritkák (sok elem 0) voltak
- rövidebb (50-1000) és sűrűbb vektorok → **skip-gram with negative sampling**
- eddig: az „apricot” szó környezetében mi hányszor szerepel → ehelyett egy bináris klasszifikációs feladat
- *Valószínű, hogy  $w$  szó közel lesz az „apricot” szóhoz?*

|                                                       |    |                            |          |                 |
|-------------------------------------------------------|----|----------------------------|----------|-----------------|
| ... lemon, a [tablespoon of apricot jam, a] pinch ... |    |                            |          |                 |
|                                                       | c1 | c2                         | t        | c3              |
|                                                       |    |                            |          | c4              |
| <b>positive examples +</b>                            |    | <b>negative examples -</b> |          |                 |
| t                                                     | c  | t                          | c        | t c             |
| apricot tablespoon                                    |    | apricot                    | aardvark | apricot seven   |
| apricot of                                            |    | apricot                    | my       | apricot forever |
| apricot jam                                           |    | apricot                    | where    | apricot dear    |
| apricot a                                             |    | apricot                    | coaxial  | apricot if      |

# Skip-gram

... lemon, a [tablespoon of apricot jam, a] pinch ...  
                  c1                  c2      t      c3          c4

## positive examples +

| t       | c          |
|---------|------------|
| apricot | tablespoon |
| apricot | of         |
| apricot | jam        |
| apricot | a          |

## negative examples -

| t       | c        | t       | c       |
|---------|----------|---------|---------|
| apricot | aardvark | apricot | seven   |
| apricot | my       | apricot | forever |
| apricot | where    | apricot | dear    |
| apricot | coaxial  | apricot | if      |

- logisztikus regresszióval betanítunk egy **bináris osztályozót**
- k-szor annyi negatív példát használ, mint pozitívat (itt k=2)
- a cél- és kontextusszavak hasonlóságának maximalizálása
- negatív párok hasonlóságának minimalizálása

- minden szóhoz két különböző vektort tanul meg: amikor  $t$  célszó, vagy amikor  $c$  kontextusszó
- két mátrixban vannak ezek tárolva:  $T$  **target matrix**, és  $C$  **context matrix**
- 3 lehetőség:
  - csak a  $T$ -t tartjuk meg
  - összeadjuk egy szó minden szóbeágyazását  $\rightarrow$  új,  $d$ -dimenziós vektor
  - konkaténáljuk őket  $\rightarrow$  új,  $2d$ -dimenziós vektor

# Szóbeágyazási modellek vizualizációja

---

# Vizualizáció hasonló vektorok lekérdezésével

| 0  | alma       | 1      | 63906 |
|----|------------|--------|-------|
| 1  | körte      | 0.8392 | 13339 |
| 2  | eper       | 0.8356 | 16159 |
| 3  | banán      | 0.8222 | 17732 |
| 4  | szilva     | 0.8046 | 12602 |
| 5  | őszibarack | 0.8011 | 4698  |
| 6  | uborka     | 0.7971 | 14735 |
| 7  | répa       | 0.7937 | 14107 |
| 8  | cseresznye | 0.7848 | 11676 |
| 9  | ananasz    | 0.7820 | 4827  |
| 10 | dinnye     | 0.7689 | 11428 |

| 0  | kenyerek      | 1      | 2270 |
|----|---------------|--------|------|
| 1  | zsemlék       | 0.8105 | 283  |
| 2  | péksütemények | 0.8048 | 997  |
| 3  | kekszek       | 0.7972 | 1046 |
| 4  | pékárúk       | 0.7957 | 771  |
| 5  | tészták       | 0.7881 | 2466 |
| 6  | lepények      | 0.7849 | 202  |
| 7  | kiflik        | 0.7843 | 349  |
| 8  | kalácsok      | 0.7841 | 277  |
| 9  | sonkák        | 0.7836 | 613  |
| 10 | pogácsák      | 0.7787 | 539  |

- Novák et al. (2017)

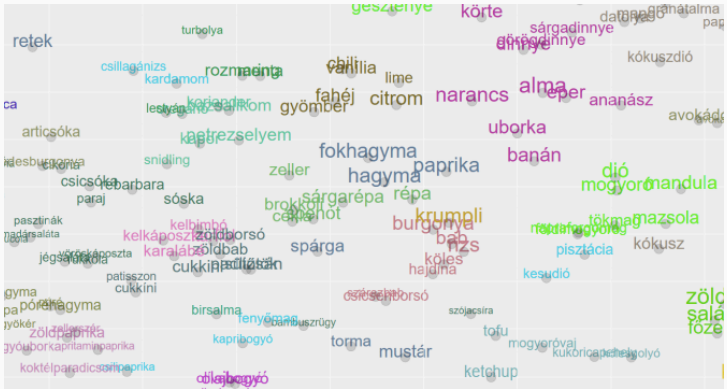
# Vizualizáció - hierarchikus klaszterezés

- a vektortérben egymáshoz közeli vektorok hierarchikus reprezentációja



# Vizualizáció t-SNE algoritmussal

- sokdimenziós vektorokat kétdimenziós térben jelenítjük meg
- megtartja az elemek közötti távolságok arányát



- Novák et al. (2017)



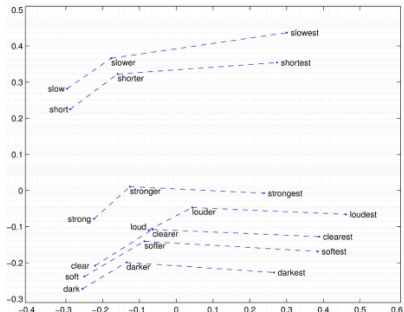
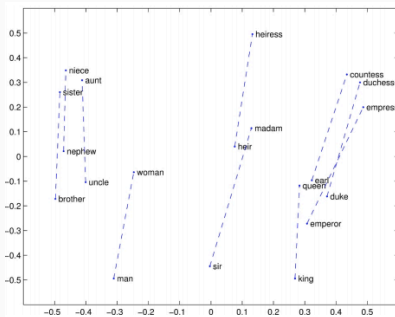
## A szóbeágyazások szemantikai jellemzői

---

# A szóbeágyazások szemantikai jellemzői

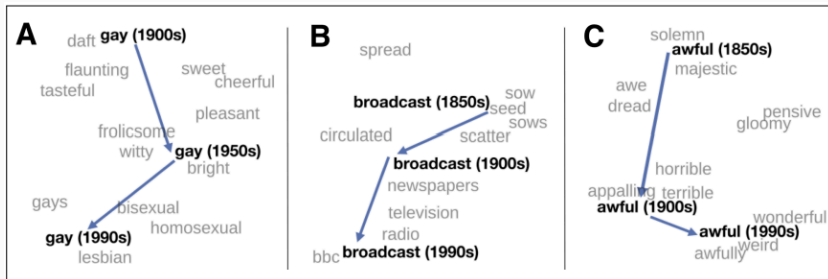
- **ablak mérete**: általában 3 és 21 közötti (1-10 mindkét oldalon)
  - ha kisebb, inkább szintaktikai hasonlóságot mutatnak (pl. szófaj megegyezik)
  - ha nagyobb, témában hasonlóak
- pl. **Hogwarts**hoz legközelibb szavak (Levy és Goldberg (2014))
  - $\pm 2$  ablakkal: Evernight, Sunnydale, Collinwood, stb.
  - $\pm 5$ -össel: Dumbledore, half-blood, Malfoy, stb.
- **first-order co-occurrence**: tipikusan közel állnak egymáshoz, pl. *wrote, book*
- **second-order co-occurrence**: hasonló szomszédjaik vannak, pl. *wrote, said*

# A szóbeágyazások szemantikai jellemzői



- **analógia**: az egyes vektorok közötti eltolások mintha valamilyen jelentéssel bírnának
- $\text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'}) = \text{vector}(\text{'queen'})$
- $\text{vector}(\text{'Paris'}) - \text{vector}(\text{'France'}) + \text{vector}(\text{'Italy'}) = \text{vector}(\text{'Rome'})$

# A szóbeágyazások szemantikai jellemzői



- **történeti szemantika:** hogyan változott a jelentése egy szónak
- egy szóhoz különböző korokban íródott korpuszokból számolunk szóbeágyazásokat
- a kontextusszavak szóbeágyazását a modern szöveg alapján készült vektortérből vesszük

## Elfogultság, sztereotípia

---

# Elfogultság, sztereotípia

- a korpuszban megtalálható sztereotípiákat, hiányosságokat is ábrázolják a szóbeágyazások
- 'man' - 'computer programmer' + 'woman' = 'homemaker'
- 'father' is to 'doctor' as 'mother' is to 'nurse'
- kimutatták azt is, hogy a kompetenciával kapcsolatos melléknevek szóbeágyazása közelebb állt a férfi(aka)t jelölő szavakhoz, 1960 óta ez a hatás csökkenni látszik
- **debiasing**: a szóbeágyazások ilyen jellegű elfogultságának enyhítése
- lehet, hogy sikerül az egyes szóbeágyazásokból kiirtani ezt, de ettől még nem szűnik meg a probléma

## Vektoros modellek kiértékelése

---

- **extrinsic** (beépítve más NLP feladatokba)
- **intrinsic**
  - hasonlóság mérése, összevetve egy gold standarddal
  - **kontextus nélkül:**
    - **WordSim-353** (0-10-es skálán 353 főnévpárt osztályoztak)
    - **SimLex-999** (melléknevek, igék, főnevek)
    - **TOEFL dataset** (80 kérdés, 4 lehetséges válasszal)
  - **kontextussal:**
    - **Stanford Contextual Word Similarity (SCWS)** dataset (2 003 szópár mondatba illesztve, ezek hasonlóságának értékelése)
    - **Word-in-Context dataset** (egy célszó két kontextusban való megadása, el kell dönten, hogy azonos jelentésben szerepel-e)
  - **analógiatesztek:**
    - "a is to b as c is to d", ahol d-t keressük
    - Athén olyan Görögországnak, mint Oslo \_\_\_\_\_-nak
    - szintaktikai: *mouse – mice, dollar – dollars*



- **GloVe**:
  - célszavakhoz valószínűséget számolunk
  - ahol két célszó hasonló valószínűségű, ott kb. 1 lesz a valószínűségek hányadosa
- **fasttext**:
  - a **word2vec** kiterjesztése
  - kezeli az ismeretlen szavakat és ezáltal a gazdag morfológiájú nyelveket
  - a szót n-gramokra bontjuk, ezekre külön-külön számolunk skip-gram szóbeágyazásokat
  - az összes alkotóelem vektorát összeadjuk, ez reprezentálja a szót

John R Firth. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*, 1957.

Eric H Huang, Richard Socher, Christopher D Manning, and Andrew Y Ng. Improving word representations via global context and multiple word prototypes. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pp. 873–882. Association for Computational Linguistics, 2012.

Daniel Jurafsky and James Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, volume 3.

Omer Levy and Yoav Goldberg. Dependency-based word embeddings. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 302–308, 2014.

Attila Novák, Borbála Siklósi, and Nóra Wenszky. Szóbeágyazási modellek vizualizációjára és böngészésére szolgáló webes felület. *XIII. Magyar Számítógépes Nyelvészeti Konferencia, Szeged, Szegedi Tudományegyetem, Informatikai Tanszékcsoport*, pp. 355–362, 2017.

Charles Egerton Osgood, George J Suci, and Percy H Tannenbaum. *The measurement of meaning*. Number 47. University of Illinois press, 1957.