

Hľadanie podobností v textoch ľudových

▶ piesní

Michal Mižák

Vedúci práce: doc. RNDr. Stanislav Krajči PhD

Oživenie pamäti

- ▶ Vložím/ nájdem pieseň alebo úryvok textu
- ▶ Dostanem zoznam podobných piesní resp. piesní, ktoré v nejakej forme úryvok obsahujú
- ▶ Štatistické vyhodnotenie podobnosti
- ▶ „Mergovací“ systém
- ▶ API jednoducho využiteľné v iných aplikáciách

Využitie

- ▶ Od dediny ku dedine, od domu k domu pretvorený/ skomolený text
- ▶ Pri folklórnom výskume sa väčšinou piesne získavajú od starších ľudí - spievajú tak, ako si spomenú
 - ▶ Síce autentickosť, ale pamätníci to môžu opraviť
- ▶ Pri rozumnej databáze porovnávanie regionálnych rozdielov v piesňach
- ▶ Priestor pre odborníkov a historikov na „skonzistentnenie“ záznamov o piesňach
 - ▶ Nesprávny zápis hlások v nárečí

Využitie

- ▶ [: **Dzifče** počarovne, **ňezal'up** še do mne :]
[: ja chlapec vandrovni, co ci budze zo mne :]

[: Ja chlapec vandrovni, zo šireho poľa :]
[: co ci budze zo mne, frajirečko moja? :]
- ▶ [: **Dzifče** počarovne, **nezal'ub** se do mne :]
[: ja chlapec vandrovny, co ci budze zo mne :]

[: Ja chlapec vandrovny, zo šireho poľa :]
[: co ci budze zo mne, frajirečko moja? :]
- ▶ [: **Dzivče** počarovne, **nepopatraj** na mne :]
[: ja chlapec vandrovny, co ci budze zo mne :]

[: Ja chlapec vandrovny zo šireho poľa :]
[: co ci budze zo mne draha duša moja? :]

Seen

Ako?

- ▶ 4 hlavné časti
 - ▶ Príprava databázy
 - ▶ Analýza frekvencie slov v konkrétnom tvare
 - ▶ Zjednocovanie rôznych tvarov slova do jedného
 - ▶ Technologická implementácia UI

Príprava databázy

- ▶ Scraping a rôzne formy skriptovania

- ▶ Dáta bolo treba niekde zohnať

- ▶ Viktor Gliganič - primáš (huslista) Ľudovej hudby FS Zemplín

- ▶ Piesne372

- ▶ Iné zdroje v stave ToDo

☑ DONE!

Budeme ignoranti

- ▶ Opakovačky
 - ▶ Štandardná štruktúra ľudovky
 - ▶ [: Sloha :] [: Refrén :]
 - ▶ Pre samotnú pieseň nie je opakovaný text refrénu/ slohy zaujímavý
 - ▶ Zátvorky a špeciálne znaky (teda `{[/:.,“]}`) tým pádom môžeme ignorovať
- ▶ Čísla ignorujeme tiež

DONE!

Manuálna práca, ktorá nikoho nezaujíma

- ▶ „ ... “
 - ▶ Dopísaný text, manuálne prejdienie časti databázy
 - ▶ Čiastočne zatriedenie piesní podľa regiónov, tónin (dur, mol) a charakteru (valčík, polka...)

DONE!

Technológie - minulý rok

- ▶ Programovací jazyk Java
- ▶ Databáza
 - ▶ SQL alebo MongoDB
- ▶ User interface
 - ▶ Angular 2, JSF

and the
winner is...



Java™

{JSON}
JavaScript Object Notation



JavaServer™ Faces
JSF

Zatiaľ statické, jednoduché dáta

- Jednoduché na údržbu

Knižnica Jackson

- Rýchla serializácia/
deserializácia objektov

„Glass cannon“

- Rýchle, odľahčené
- Ťažšie sa udržuje konzistentnosť dát

JSON

Zatiaľ statické,
jednoduché dáta

FALSE

Ťažšie sa udržiava
konzistentnosť
dát

TRUE

„Oh, my sweet summer child...“

and the
winner is...



Java™




JavaServer™ Faces
JSF

Prečo nie?

Pre mňa nová
technológia, mám
sa s kým poradiť

JSF



To Do...



To Be Continued...



Nejaké...

To Do...



JavaServer™ Faces
JSF

To Do...

Čo vlastne kódim?

1. Vložím text
2. Spustím algoritmus
3. (?) Priebežne počas iterácie konfiguráciami algoritmu vypisujem výsledky
4. Vrátim percentuálnu podobnosť textov

AlgorithmConfiguration

- ▶ Abstraktná trieda rozšírená pre každý typ algoritmu
- ▶ Immutable - pre prípad, že ostane čas na paralelizáciu („Oh, my sweet summer child...“)
- ▶ Jej „deti“ obsahujú inštančné premenné definujúce beh algoritmu

AlgorithmConfiguration

- ▶ Zatiaľ pracujem len na „Vector space algorithm“
- ▶ Idea - text si rozdelím na vektory, kde dimenzia je „slovo“ a hodnota jeho váha
- ▶ Vektory dvoch textov algebraicky porovnáam

Na pripomenutie

wat wat	wat wat2	wat2 wat2	wat2 wat
2	1	1	1

```
@Immutable
public class VectorAlgorithmConfiguration extends AlgorithmConfiguration {

    private TermScheme termScheme;
    /**
     * In case term scheme is more complex
     */
    private Integer termDimension;

    private TermWeightType termWeightType;
    private VectorInclusion vectorInclusion;
    private TermComparisonAlgorithm termComparisonAlgorithm;
    private Double tolerance;
    private VectorComparisonAlgorithm vectorComparisonAlgorithm;
```

TermScheme - N-gramy

- ▶ Poradie?

- ▶ Rozšírenie databázy o n-tice slov - n-gramy

DONE!

TermScheme - zovšeobecnenie

- ▶ „Ja chlapec vandrovny, zo šireho poľa“
 - ▶ Dvojice
 - ▶ <Ja, vandrovny>, <chlapec, zo>, <vandrovny, šireho>
 - ▶ <Ja, poľa>
 - ▶ Trojice
 - ▶ <Ja, chlapec, zo>, <chlapec, vandrovny, šireho>
 - ▶ Mnoho ďalších...

DONE!

Term weighting

- ▶ „Slová“ v texte sú pre nás rôzne zaujímavé
 - ▶ Tf - term frequency
 - ▶ Idf - inverse frequency
 - ▶ Tf-idf (term frequency-inverse document frequency)
 - ▶ „A ja taka čarna jak čarna čarnica...“
 - ▶ „A ja taka dzivočka cingi lingi bom...“
 - ▶ Mnoho ďalších...

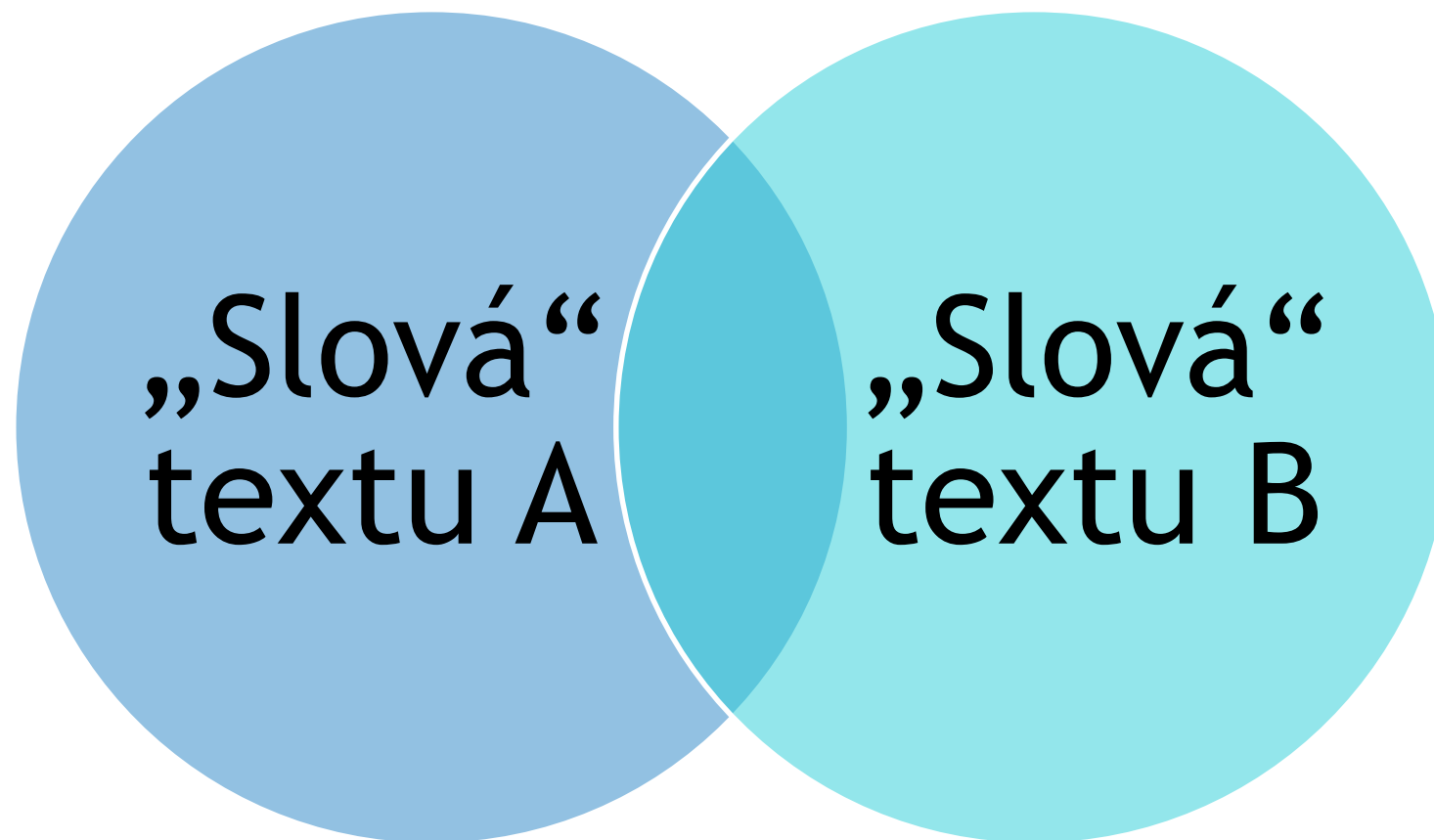


To Be Continued...

Vector inclusion

Poznámka - pracujeme s
dimenziami vektora

Vector inclusion



Vector inclusion

- ▶ Vektor 1 = $\langle A, B, C, D \rangle$
- ▶ Vektor 2 = $\langle A, B, D, E \rangle$

Vector inclusion

▶ A forma - C vo vektore B priradíme 0

▶ $\langle A, B, C, D \rangle$

▶ B forma - E vo vektore A priradíme 0

▶ $\langle A, B, D, E \rangle$

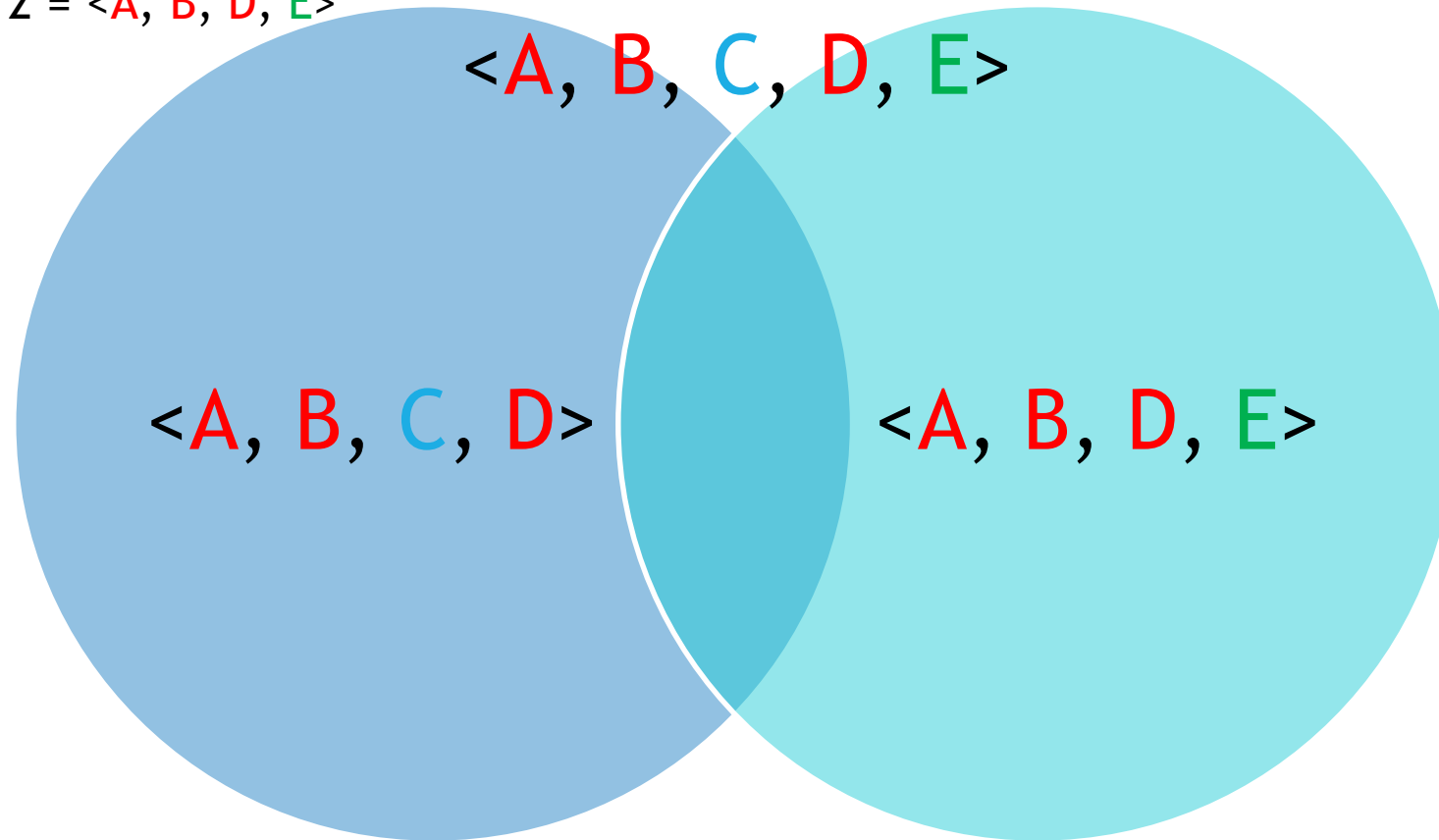
▶ AB forma

▶ $\langle A, B, C, D, E \rangle$

Vector inclusion

Vektor 1 = $\langle A, B, C, D \rangle$

Vektor 2 = $\langle A, B, D, E \rangle$



Krásne formátovanie, viem

Všetky „slová“

„Slová“
textu 1

„Slová“
textu 2

To Be
Continued...

Term comparison algorithm

- ▶ Pri formátovaní vektorov je potrebné „slová“ porovnať
- ▶ Pre základné vyhľadávanie bez podobnosti slov stačí naivný algoritmus

☑ DONE!


```
public class NaiveTermComparator implements ITermComparator {  
  
    @Override  
    public double compare(Term t1, Term t2) {  
        return t1.getLyricsFragment().equals(t2.getLyricsFragment()) ? 1 : 0;  
    }  
  
}
```

Term comparison algorithm + Vector Inclusion

- ▶ Keď už hľadáme podobnosť piesní, tak prečo nie aj „slov“ 😊
- ▶ Musíme zlučovať stĺpce

Term comparison algorithm + Vector Inclusion

- ▶ Vektor 1 = $\langle A, B, C, D \rangle$
- ▶ Vektor 2 = $\langle A, B, D, E \rangle$
 - ▶ Pozn. toto sú dimenzie vektora, nie hodnoty

Term comparison algorithm + Vector Inclusion

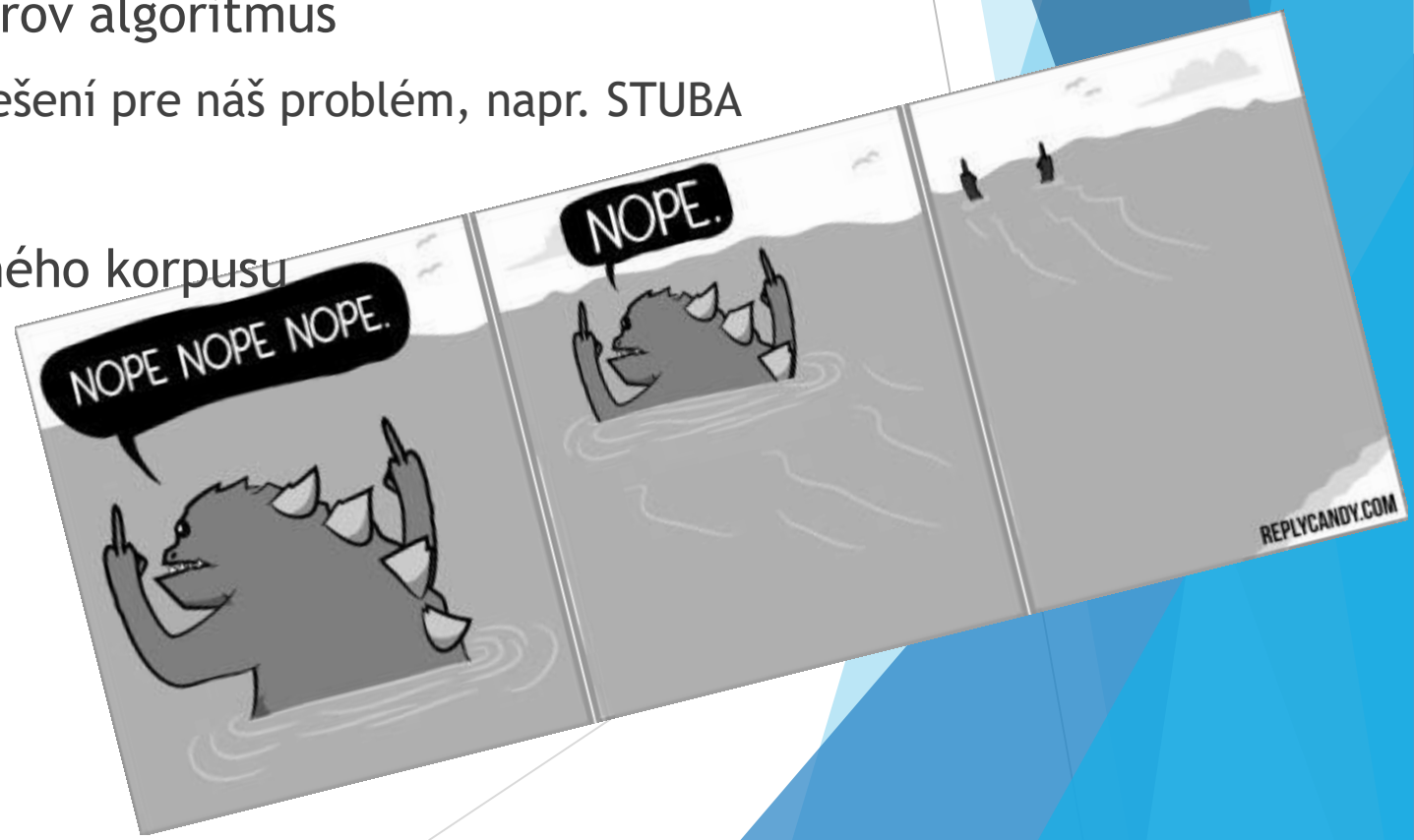
- ▶ Predpokladajme, že **C** a **E** sú v tolerancii podobnosti (definovaná v alg. config.)
- ▶ Dostaneme
 - ▶ Vektor 1 = $\langle \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D} \rangle$
 - ▶ Vektor 2 = $\langle \mathbf{A}, \mathbf{B}, \mathbf{E}, \mathbf{D} \rangle$



To Be
Continued...

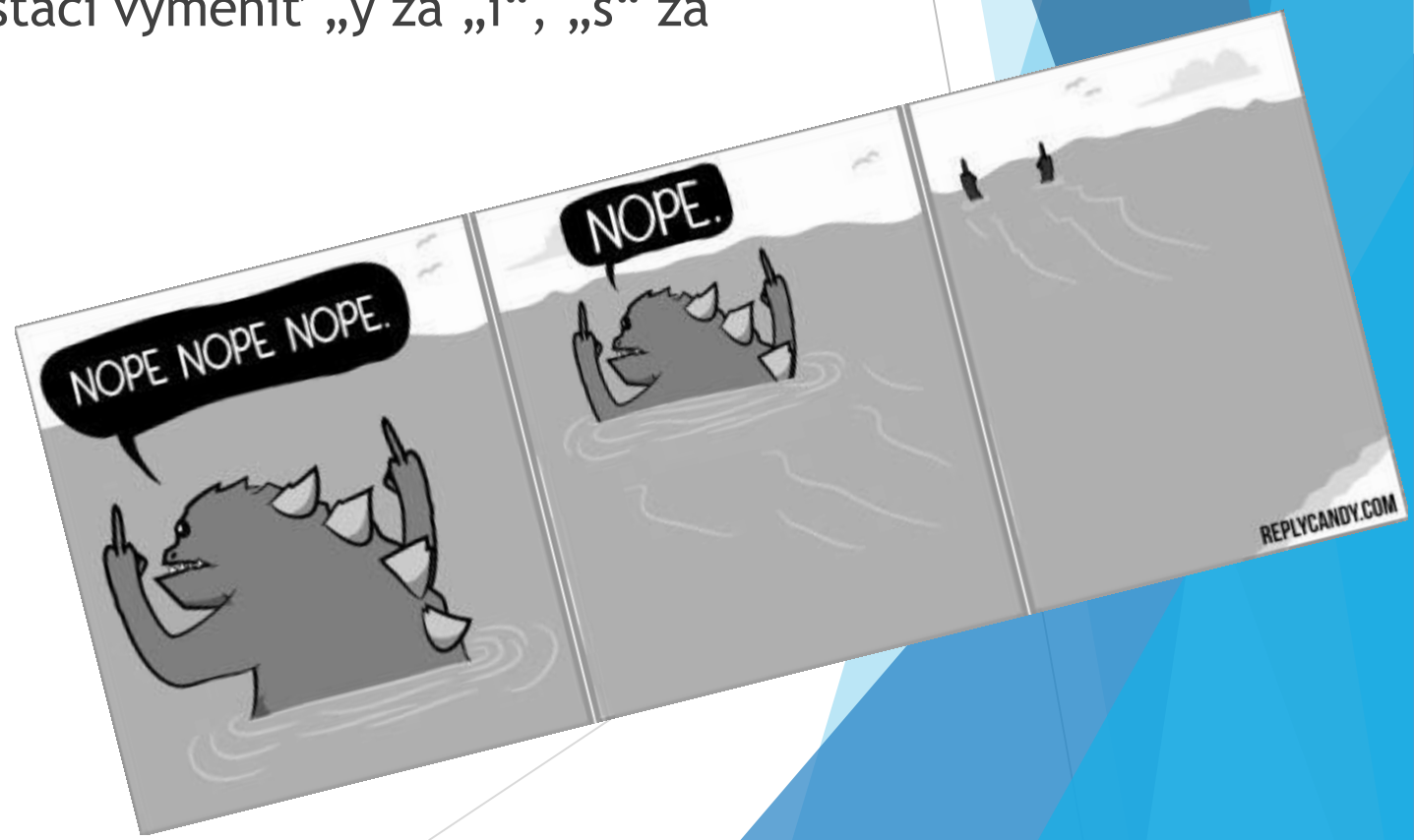
Term similarity

- ▶ Hľadanie základu slova - zlučovanie stĺpcov
 - ▶ Lemmatization, stemming - porterov algoritmus
 - ▶ Analýza efektivity existujúcich riešení pre náš problém, napr. STUBA
 - ▶ Ohýbanie riešení na ľudový jazyk
 - ▶ Korpus nárečí Slovenského národného korpusu



Term similarity

- ▶ Tvaroslovník - vytváranie vlastného ľudového „tvaroslovníka“
 - ▶ Podobnosti sú často veľmi veľké, stačí vymeniť „y“ za „i“, „š“ za „ś“, „h“ za „g“ a pod.
 - ▶ Generovanie koncoviek



Term similarity

- ▶ Levenshteinova vzdialenosť

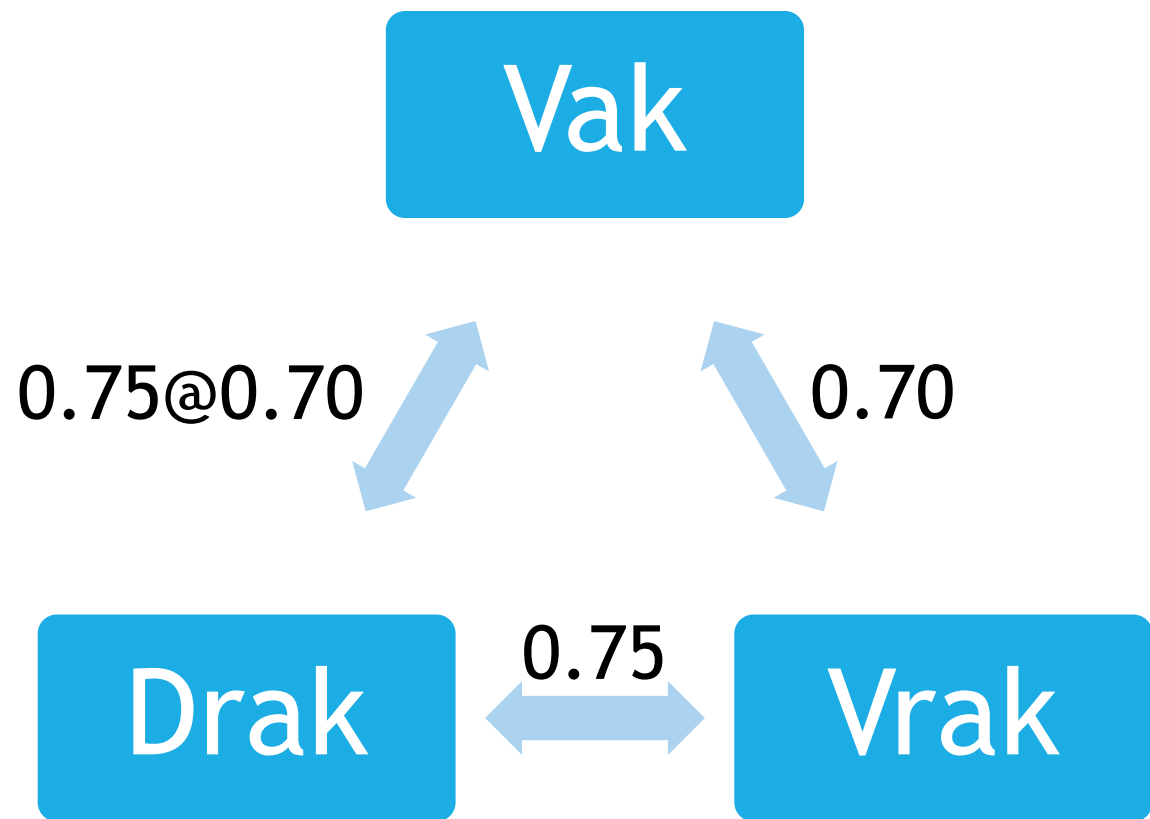
1. **kitten** → **sitten** (substitution of "s" for "k")
2. **sitten** → **sittin** (substitution of "i" for "e")
3. **sittin** → **sitting** (insertion of "g" at the end).

- ▶ Mnoho ďalších...

Indirect term similarity

- ▶ Zatiaľ len nápad, ktorý treba doladiť
- ▶ Idea - „slová“ si nemusia byť podobné priamo, ale cez nejakého suseda
- ▶ Ilustrácia - číslo naznačuje, nakoľko sú si „slová“ podobné v danom porovnávacom algoritme
- ▶ @ - nejaká vhodná operácia, napríklad súčin

Indirect term similarity



Indirect term similarity

- ▶ Dostaneme teda vzorec
- ▶ $P_i(a, b) = \max(P(a, b), \max_{c \in \text{Vsetky "Slova"}} (P_i(a, c) @ P_i(b, c)))$, kde $c \in \text{Vsetky "Slova"}$
- ▶ Kde P_i je Indirect term similarity, P je Term similarity a a, b sú „slová“

Indirect term similarity

- ▶ Prípadne len
- ▶ $P_i(a, b) = \max(P(a, b), \max_{c \in \text{Vsetky "Slova"}}(P(a, c) @ P(b, c)))$
- ▶ Alebo inak doladiť „mieru nepriamosti“

Vector comparison algorithm

- ▶ Vzdialenosť dvoch vektorov
 - ▶ Kosínusová miera

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}, w$$

☑ DONE!

$$|Q \cap D|$$

Simple matching (coordination level match)

$$2 \frac{|Q \cap D|}{|Q| + |D|}$$

Dice's Coefficient

$$\frac{|Q \cap D|}{|Q \cup D|}$$

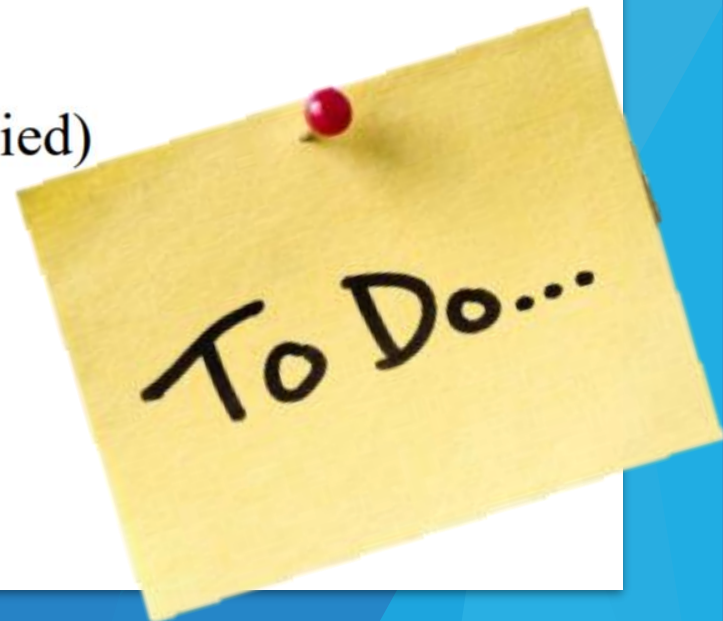
Jaccard's Coefficient

$$\frac{|Q \cap D|}{|Q|^{1/2} \times |D|^{1/2}}$$

Cosine Coefficient (what we studied)

$$\frac{|Q \cap D|}{\min(|Q|, |D|)}$$

Overlap Coefficient



Zdroje

- ▶ Texty ľudových piesní
 - ▶ Rôzne stránky
 - ▶ Kontakty u nadšencov
- ▶ <https://en.wikipedia.org/wiki/N-gram>
- ▶ https://en.wikipedia.org/wiki/Levenshtein_distance
- ▶ <https://courses.cs.washington.edu/courses/cse573/12sp/lectures/17-ir.pdf>
- ▶ <https://nlp.stanford.edu/IR-book/>



To Be Continued...

Ďakujem za pozornosť

