

## Prezentácia projektových zadaní

Možnosť vlastnej témy, môže súvisieť s diplomkou alebo bakalárkou

Zameranie: spracovanie textových dát za účelom vyhľadávania alebo extrakcie informácií

### **DBPedia spotlight**

naprogramovať NE extraction pomocou [DBPedia spotlight](#)

### **Regex editor**

vytvorenie jednoducheho editora na regularne vyrazy ked bude moct uzivatel cez Java aplikaciju testovat regularne vyrazy na nacitanom txt subore.

Vypisanie extrahovanych udajov aj skupin (groups)

Podla moznosti integracia z Ontea nastrojom.

Integracia makier z ontea aj s prikladmi textov a spustenim na nich v description. Vytvaranie a editacia makier pomocou nastroja.

### **Štatistický prekladač**

Na základe rôznych jazykových verzií stránok vytvoriť prekladač na základe štatistických údajov.

Existuje napríklad package [Moses](#).

Alebo vytvorenie prekladača slov na základe spracovania anchor textov.

Je možné riešiť nasledovné projekty:

- prekladač slovenčina <=> angličtina na základe anchor textov
- prekladač slovenčina <=> čeština na základ hociakých textov. Pri podobných jazykoch by mohol byť menší problém s tým že nevieme rozpoznať slovné druhy (POS tagging)

### **Gazetteer**

Vytvorenie podobného gazetteera (slovníka) pre information extraction ako ma GATE.

Vlastnosti:

- lineárna zložitost' (iba jeden prechod textom)
- definovanie oddeľovača slov (tokenizatora), najlepšie po znakoch, pričom hľadanie zhody začína vždy za white space znakom
- generovanie aliasov pomocou regularneho vyrazu. Napr. v slovníku je "Meno Priezvisko" a bude hľadať aj "M. Priezvisko"

## Podpora slovenskeho vyhladávania

- analyzer ktorý rieši diakritiku (napr vyhadzuje), spellcheck aj bez diakritiky - urobiť nad lucene

## Extrakcia faktov

Extrakcia faktov zo slovenského webu alebo iných textových dokumentov. Niečo na spôsob [knowItAll](#) ale pre slovenčinu. Napríklad zo stránky FIIT vytiahnuť zoznam učiteľov, predmetov alebo študentov. Možno jednoduchý príklad. Proste vytiahnuť fakty ktore sa nedajú nájsť z jedného dokumentu. Iný príklad je napríklad vytiahnuť zoznam sklenárstiev v Bratislave.

- zoznam ľudí zo stránok ústavou SAV, automatické porovnanie so SAV.sk

## Tvorba vyhľadávača kontaktov v mobile a iné

- aplikácia na spôsob písania SMS pomocou T9. Teda niečo kde termy v indexe budú vlastne čísla.
- Može to byť vyhľadávanie kontaktov alebo menších dokumentov pomocou zadania query cez čísla reprezentujúce písmená.

## Focused Crawler

- sťahovanie dynamicky generovaných stránok tak aby sa nesťahoval rovnaký obsah (inak utriedený zoznam, verzie na wiki stránke, printová verzia stránky a podobne)
- sťahovanie toho istého obsahu v inej jazykovej verzii (tvorba corpusu pre štatistický prekladač)

## Rozpoznávanie slovných druhov

Part of Speech Tagging (POS) tagging pre slovenčinu. Založené na slovníkovom princípe z dostupných slovníkov (aspell, ispell, OpenOffice a pod.) alebo založený na štatistickom princípe. Možno pomyliť princíp OpenNLP. Netreba všetky slovné druhy ani nemusí 100 percentne fungovať.

## Fazetový a fultextový prehliadač

Kto má prístup k databáze nejakých produktov alebo nejaký rozsiahlejší web. Treba urobiť fultext toho webu s kombináciou fazetového prehliadača. (Podobne ako na amazon.com a iných)

Je to možné urobiť pomocou systému Apache Solr. Ide najmä o konfiguráciu a napojenie systému a jeho odladenie.

## Spread Activation

spread activation algoritmus (pozri wikipediou) treba naprogramovať a použiť na nejakých dátach. Napr na extrakcii z emailov.

## Tag Cloud

Generovanie Tag cloudu (pozri wikipediou) z webstránky a jej podstránok.

Teória okolo TF-IDF, stop slová, lematizácia....

## Name Entity Recognition - Machine Learning

Rozpoznávanie mien (osoby, mesta, organizácie, ...) pomocou OpenNLP

2 projekty - jeden rozpoznávanie na Slovenskom a anglickom texte, druhý trenovanie na Slovenskom

## Name Entity Recognition - Extrakcia Anchor Text

Pomocou extrakcie anchor textov liniek (text v rámci Tagu <a href>TEXT</a>) robiť Named Entity recognition a Aliasy.

## Analyzer ktorý vyhodí diakritiku

použitie vo vyhľadávачi emailov alebo fajlov na disku spolu so spell checkom.

## Advanced email search

Vyhľadavac pomocou socialnej siete-grafu extrahovaneho z emailu.

Pomocou existujuceho softveru ([Ontea](#) alebo acoma = [emailSocNet](#)) sa extrahuje graph. Nad nim urobiť vyhľadavanie pomocou spread activation ?+lucene kde budu fazety podla objavenych typov objektov.

## Porovnanie implementacii spread activation algoritmov

Z minuleho roka treba porovnat implementacie spread activation (3 projekty - Fridrich, Blazko, Marton ) a zistiť v čom sa líšia a ako je najvyhodnejšie implementovať spread activation na Jung knižnici tak aby bol algoritmus vseobecny a pouzitelny

## Vyhľadávanie s využitím anotácií (tagov)

V systémoch ako delicious.com, twitter alebo youtube uzivatelia generujú množstvo tagov ktoré v kombinácii s klasickými technikami vyhľadavanie (indexovanie) môžu priniesť lepšie výsledky.

Úlohou je vytvoriť vyhľadavач ktorý tieto tag-y využije.

## Extrakcia udalostí

extrahovanie udalostí z emailov alebo webových stránok.

- dátum, čas
- miesto
- názov udalosti