

ВЪВЕДЕНИЕ В ИЗВЛИЧАНЕТО НА ИНФОРМАЦИЯ

или как работят търсачките



СУ, ФМИ
Катедра “Компютърна информатика”
Мария Матева

Съдържание

- 1 Вашите водещи(на упражнения)
- 2 Организационни
- 3 Предизвикателство #1: Обратно инженерство на търсачката на Google
- 4 Понятия
- 5 Обща архитектура на търсачките
- 6 Булев модел
- 7 Предварителна обработка на текст
- 8 Обърнат индекс
- 9 Въпроси и отговори

- Мария Матева



- Ясен Кипров



Завършили са “Компютърни науки” и “Изкуствен интелект” във ФМИ, СУ.

Имат опит в изграждането на специализирани търсачки, както и в задачи от машинното самообучение.

Мария в момента работи върху софтуер за аналитични решения.

Ясен е докторант и има интереси в областта на обработката на естествен език. (Той ще разкаже повече.)

Moodle

- Регистрирайте се в курса в moodle
<http://moodle.openfmi.net/>
- Курсове » Магистри, зимен семестър 2014/2015 » ИИОЗ
» Извличане на информация 14/15
- IR_2014!
- там ще намерите тази презентация! ;)
- както и всички други материали
- във форума ще поддържаме актуална организационна информация
- можете да задавате всякакви въпроси там

Провеждане

- Упражненията биват **теоретични и практически**
- Ще бъдат в **понеделник, 19:30 - 21:00, бл.2 на СУ**
- Теоретичните упражнения ще бъдат в зала "С"
- Практическите упражнения ще бъдат в зали "30X" -
очаквайте подробности
- Ще ви осведомяваме във форума своевременно
- За практическите упражнения можете да си носите
собствени лаптопи
- За теоретичните упражнения е добре да си носите нещо за
писане

Софтуер, с който ще работим

На практическите упражнения ще се запознаем с:

- GATE 8, <https://gate.ac.uk>
- Apache Lucene, <http://lucene.apache.org/core/>
- Apache Solr, <http://lucene.apache.org/solr/>
- Apache Nutch, <http://nutch.apache.org> и други кролери
- други

Оценяване

Крайното оценяване се базира на

- проект
- два теста
- домашни(по желание)
- участие в упражнения(бонуси)
- предизвикателства(бонуси)

Предизвикателство #1: Обратно инженерство на търсачката на Google

Обратно инженерство (на английски: Reverse engineering) е процесът по преоткриване и възстановяване (реконструиране) на технологичните принципи и механизми, по които е създаден определен обект, машина или програма.

Как работи търсачката на Google?

- Задачата е да се намерят интересни/неочаквани поведения в търсачката, които доказват конкретни характеристики на работата ѝ.
- Добрите предложения ще получат бонус :)



- Ще повторим предизвикателството накрая на семестъра.

Предизвикателството - пример: домат, БГ

gotvach.bg > Продукти > Зеленчуци

доматите са най-често срещаните, обичани и полезни зеленчуци в света. Според мнозина **доматите** всъщност са плод, но с оглед на ...

Манастирски ливади - Пицария Дон Домат | пица, бира и ..
dondomat.com/zavedenia.php?z=12 ▾
... бизнес без рисък - Доктор Дон Домат ресторантски услуги · Контакти Връзка с

Домати градина Вискар
домати.вискар.com/

Предизвикателство #1: Обратно инженерство на търсачката на Google

Предизвикателството - пример: tomato, БГ

The screenshot shows a Google search results page for the query "томат". The search bar at the top contains "томат". Below it, the search interface includes a magnifying glass icon, a "Вход" (Sign In) button, and a gear icon for settings. The main search area has a red box highlighting the search term "томат" in the search bar.

Below the search bar, there are several navigation links: "Мрежата" (Network), "Изображения" (Images), "Видеоклипове" (Videos), "Новини" (News), "Още" (More), and "Инструменти за търсене" (Search tools). A "Change language" button with a gear icon is also present.

The search results section starts with a message: "„Езиките“ ни помагат да предоставяме услуги си. С използването им приемате употребата на „български“ от наша страна." followed by "Научете повече" and "OK" buttons.

The first result is a link to "Tomato - Wikipedia, the free encyclopedia" with the URL en.wikipedia.org/w/index.php?title=Tomato&oldid=6450000. The snippet describes Tomato as an edible fruit and its botanical classification.

The second result is a link to "Tomato (firmware) - Wikipedia, the free encyclopedia" with the URL [en.wikipedia.org/w/index.php?title=Tomato_\(firmware\)&oldid=6450000](http://en.wikipedia.org/w/index.php?title=Tomato_(firmware)&oldid=6450000). The snippet describes Tomato as a partially free HyperVRT-based Linux core firmware distribution.

The third result is a link to "Tomato Firmware | polarcloud.com" with the URL www.polarcloud.com/tomato. The snippet describes Tomato as a small, lean and simple replacement firmware for Linksys' WRT54G/GL/GS, Buffalo WHR-054SV/WHR-HP-054 and other Broadcom-based routers.

The fourth result is a link to "Rotten Tomatoes: Movies | TV Shows | Movie Trailers ..." with the URL www.rottentomatoes.com/. The snippet describes Rotten Tomatoes as a consensus opinion of professional critics from across the nation.

A yellow callout box on the right side of the page, titled "Looking for results in English?", contains options to "Change to English", "оставане на български" (Stay in Bulgarian), and "Езикови настройки" (Language settings). It also features a large image of a tomato and a "Отговори" (Answers) button.

Предизвикателството - пример: домат, UK

DOMAT

Web Images Videos Shopping News More Search tools

About 401,000 results (0.43 seconds)

Домат — Уикипедия
[bg.wikipedia.org/w/index.php?title=%D0%Е%0%М%D0%Е%D0%BB%D0%Е%D1%82%D0%Е&oldid=1432476210](https://bg.wikipedia.org/w/index.php?title=%D0%BE%D0%BC%D0%BE%D0%BB%D0%BE%D1%82%D0%BE&oldid=1432476210) ▾ Translate this page
 Доматът (Solanum lycopersicum) е растение от семейство Картофови (Solanaceae). Родина на домата са Централна и Южна Америка, от Мексико до ...

Манастирски ливади - Пицария Дон Домат | пица, бира и ...
dondomat.com/zavedenia.php?z=12 ▾ Translate this page
 ... бизнес без риск: Доктор Дон Домат ресторантни услуги - Контакти Бързка с нас. » Начало » Нашите заведения » Дон Домат - Манастирски ливади ...

Дондуков - Пицария Дон Домат | пица, бира и добро ...
dondukov.com/zavedenia.php?z=9 ▾ Translate this page
 всяка събота и неделя през месец февруари Дон Домат - Шести Септември и Дон Домат - Дондуков Ви предлагат всички малки пици с 20% отстъпка.

Дон Домат - Манастирски ливади - Sofia, Bulgaria ...
<https://en-gb.facebook.com/DonDomatMansiLivi/> ▾ Translate this page
 Дон Домат - Манастирски ливади, Sofia, Bulgaria. 218 харесвания · 19 говорят за това · 263 бяха тук. Ресторант Дон Домат отвори врати и в кв.

Дон Домат - Дондуков - Пицария | Фейсбук | Facebook
https://en-gb.facebook.com/_/DonMat_41234876210... ▾ Translate this page
 Дон Домат - Дондуков. 867 харесвания · 30 говорят за това · 271 бяха тук. Верига пицарии Дон Домат. Пица, бира и добро настроение. Най-добрата ...

Дон Домат - Шести Септември - Sofia, Bulgaria - Пицария ...
https://en-gb.facebook.com/_/DonMat_413526160655... ▾ Translate this page
 ... Шести Септември, Sofia, Bulgaria. 368 харесвания · 9 говорят за това · 318 бяха тук. Верига пицарии Дон Домат. Пица, бира и добро настроение....

Домат - Bonduelle България
bonduelle.bg/franeme/encyclopedia-na.../domat/ ▾ Translate this page
 Доматът е вкусен и полезен зеленчук, който съдържа ценни биолепчни вещества, а също капици капици фосфор и много витамини A, B, C и F и



More images

Tomato

Vegetable

The tomato is the edible, often red fruit/berry of the nightshade Solanum lycopersicum, commonly known as a tomato plant. [Wikipedia](#)

Nutrition Facts

Tomatoes, red ▾

Amount Per 100 grams ▾

Calories 18

% Daily Value*

Total Fat 0.2 g 0%

Saturated fat 0 g 0%

Polyunsaturated fat 0.1 g

Monounsaturated fat 0.1 g

Cholesterol 0 mg 0%

Sodium 5 mg 0%

Potassium 237 mg 6%

Total Carbohydrate 3.9 g 1%

Dietary fiber 1.2 g 4%

Предизвикателството - пример: tomato, UK

tomato

Web Images Shopping Videos News More Search tools

About 152,000,000 results (0.41 seconds)

Tomato - Wikipedia, the free encyclopedia
en.wikipedia.org/wiki/Tomato ▾
 The tomato is the edible, often red fruit/berry of the nightshade Solanum lycopersicum, commonly known as a tomato plant. The species originated in the South ...
 Lycopene - Nix v. Heden - Beefsteak tomato - Tomatine

tomato
www.tomato.co.uk/ ▾
 tomato is a collective of artists, designers, musicians and writers engaged in publishing, exhibitions, live performances, advertising, architecture, fashion, public ...

Tomatoes - The World's Healthiest Foods
www.whfoods.com/genpage.php?name=foodspice&id=44 ▾
 What's New and Beneficial About Tomatoes. Did you know that tomatoes do not have to be a deep red color to be an outstanding source of lycopene? Lycopene ...

Tomato | BBC Good Food
www.bbcgoodfood.com/glossary/tomato ▾
 The British tomato season runs from June to October. In winter, you could use more canned tomatoes to save on food miles (the environmental cost of food ...

News for tomato

Rotten Tomatoes: Farmers Pay the Price for a False Food ...

 Businessweek - 7 hours ago
 Graves Williams, a farmer in Quincy, Fla., was just a few days into a six-week tomato harvest in June 2008, when the U.S. Food and Drug ...

Heinz unveils new Tomato Ketchup TV ad
 Talking Retail - 2 hours ago
 Recipe: Ricotta and Tomato Jam Jar (9.29.14)
 fox13now.com - 17 minutes ago



More images

Tomato

Vegetable

The tomato is the edible, often red fruit/berry of the nightshade Solanum lycopersicum, commonly known as a tomato plant. Wikipedia

Nutrition Facts
 Tomatoes, red ▾

Amount Per 100 grams	% Daily Value*
Calories 18	
Total Fat 0.2 g	0%
Saturated fat 0 g	0%
Polyunsaturated fat 0.1 g	
Monounsaturated fat 0 g	
Cholesterol 0 mg	0%
Sodium 5 mg	0%
Potassium 237 mg	6%
Total Carbohydrate 3.9 g	1%
Dietary fiber 1.2 g	4%

Мрежата от близки/свързани понятия

Извличане на информация(Information retrieval)

представлява придобиването на (информационни) ресурси, които отговарят на дадена информационна нужда, избрани измежду колекция от такива ресурси. Търсенето се основава на съдържанието на ресурсите или на техните метаданни.

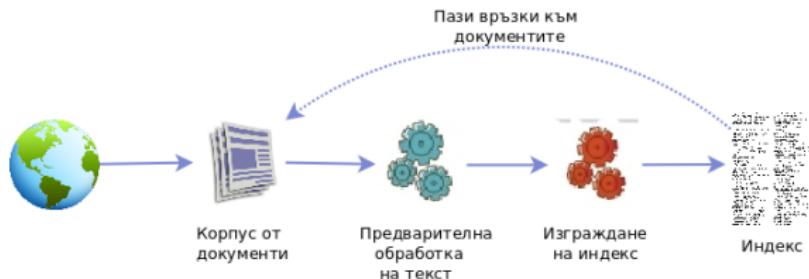
- Изкуствен интелект(Artificial Intelligence)
- Машинно самообучение(Machine Learning)
- Обработка на естествен език(Natural Language Processing)
- Анализ на текст(Text analysis; Text Mining)
- Извличане на закономерности от данни(Data mining)
- Анализ на данни(Data Analysis)
- Структурирано извлечение на информация(Information Extraction)
- Извличане на значещи обекти(Named-entity Recognition, Semantic Annotation)
- Препоръчващи системи(Recommendation Systems), и още,

Важни понятия

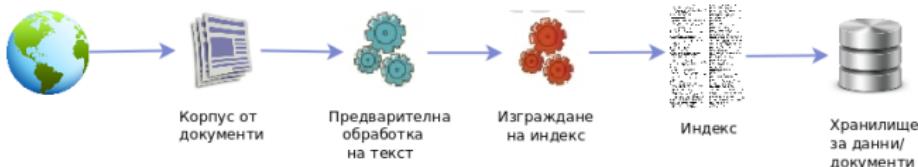
Важни (уводни) понятия в извличането на информация

- документ
- корпус
- информационна нужда
- заявка
- индекс
- резултат
- метрики за оценка на точност: precision, recall, accuracy, F1

Индексиране



ИЛИ



Извличане



Булев модел

Направете матрица на инцидентност на следните документи и думите в тях:

- $d1 = \{ \text{Иван затвори вратата.} \}$
- $d2 = \{ \text{Петър погледна навън.} \}$
- $d3 = \{ \text{Иван и Петър излязоха навън през вратата.} \}$

Изпълнете заявките: “Иван”, “навън”, “Петър” OR “вратата”, “Петър” AND “вратата”

Основни стъпки в обработката на текст:

- определяне какво е документ
- кодировка на текста (encoding)
- определяне на езика на текста
- определяне на азбуката на писане (кирилица или шльокавица?)
- лексикален анализ - извлечение на наредено множество от части на речта
- токенизация - отделяне за думи (и фрази), които по-нататък считаме за единици при анализа
- премахване на стоп-думи
- определяне на части на речта
- стеминг - определяне на корен на думата чрез прости евристични правила
- лематизация - определяне на корен на думата чрез по-сложни, морфо-синтактични правила
- синтактичен анализ - намиране на връзки на подчинение между отделните части на речта в рамките на изречението

Обърнат индекс

След първоначална обработка на текста:

- за всяка дума правим списък от документите, в които се среща
- първо обработваме всеки документ - създаваме двойки дума - ID на документ
- сортираме по думи(азбучен ред)
- сливаме повторенията, формирайки списък с ID-та на документите, съдържащи дадена дума
- докато сливаме повторенията, пъддържаме списъка със съответстващите документи

Задачка

Прилагайки стеминг, Съставете обърнат индекс на следните документи и думите в тях:

- $d1 = \{ \text{Иван затвори вратата.} \}$
- $d2 = \{ \text{Иван затвори една врата.} \}$
- $d3 = \{ \text{Иван излезе навън.} \}$
- $d4 = \{ \text{Петър погледна навън.} \}$
- $d5 = \{ \text{Иван и Петър излязоха навън през вратата.} \}$

Операции върху множества: обединение, разлика, сечение
Пропускащи списъци(skip-lists) и операции при голяма разлика
в дължината на списъците

Въведение в извлечането на информация

Въпроси и отговори

Въпроси?

Следващия път...

- **Практическо упражнение**
- Ще говорим за предварителната обработка на текст
- И GATE. Моля, изтеглете си GATE оттук:
<https://gate.ac.uk/download/>

Благодаря за вниманието! :)