

Taal- en informatietechnologie

Week 4
Informatie-extractie

23 september 2003

10/27/03

ABVL, Krahmer & Van den Bosch

1

Instru teki ennätystuloksen

Sairaalateknologiayritys Instrumentarium teki liikevaihdolla ja -voitolla mitattuna satavuotisen historiansa parhaan tuloksen. Liikevaihto kasvoi kymmenen prosenttia. Valtaosan liikevaihdosta tuova anestesia ja tehohoito kasvoi peräti 14 prosenttia.

10/27/03

ABVL, Krahmer & Van den Bosch

2

Nato-lähde: Turkki-umpikujaan tulossa uusi ehdotus

Naton neuvosto lykkäsi myöhemmäksi kokoustaan, jossa sen tulisi ratkaista umpikujaan ajautunut kiista Turkin avustuspyyntöä. Nato-lähde vahvisti Helsingin Sanomille tiistaina päivällä, että pääsihteeri **George Robertson** valmistelee uutta esitystä asiassa.

10/27/03

ABVL, Krahmer & Van den Bosch

3

Rond vier uur kwam Lieve Kate halen, even trap geschuurd, en dan naar de volleybal. Was zeer leuk, het ging niet helemaal zo goed, maar toch gewonnen. Achteraf een beetje blijven plakken, altijd best leuk. En toen ik thuis was, nog een goed anderhalf uur zitten schurer/krabben aan de trap. Tof, stom, vermoeiend werk. Morgen nog doen.

10/27/03

ABVL, Krahmer & Van den Bosch

4

VS woedend over veto in NAVO

De Verenigde Staten hebben woedend gereageerd op de blokkade voor NAVO-hulp aan Turkije die Frankrijk, Duitsland en België gisteren hebben opgeworpen. Volgens de Amerikaanse ambassadeur bij de NAVO, Nick Burns, kampt de alliantie nu met een 'geloofwaardigheidscrisis' vanwege het veto van de drie.

10/27/03

ABVL, Krahmer & Van den Bosch

5

Tekst vs Informatie

- Tekst *is niet* informatie!
- Informatie **zit in** tekst
- Informatie is **taalonafhankelijk**
- Informatie uit tekst halen betekent **de taal kennen**
 - grammatica
 - conventies

10/27/03

ABVL, Krahmer & Van den Bosch

6

Artikel uit 1946

Ettelijke zelfmoordpogingen werden verijdeld

Aanslag op Göring e.s. onthuld

— Het Berlin van nu is niet meer dat van weleken tijd. Het is nu een stad die door de Britten als een van de vier grote machten wordt beschouwd. Het is nu een stad die door de Amerikanen wordt beschouwd als een van de vier grote machten. Het is nu een stad die door de Sovjet-Unie wordt beschouwd als een van de vier grote machten. Het is nu een stad die door de Verenigde Naties wordt beschouwd als een van de vier grote machten.

10/27/03

ABU, Krahmer & Van den Bosch

7

1915

ALGEMEEN NIEUWS.

Dierfict door een politiekagent. Gisteren is de 14-jarige agent van politie geboort, verbleef zich in het schoolgebouw van de Kerkstraat. Hij kwam gisterenavond op de school in de school van de Kerkstraat. Hij kwam gisterenavond op de school in de school van de Kerkstraat. Hij kwam gisterenavond op de school in de school van de Kerkstraat.

10/27/03

ABU, Krahmer & Van den Bosch

8

1785

ITALIEN.

ROMEN, den 12 Oktober. In den nacht van den 11ten dezer, gevoeld men, omtrent 4 uren in den morgen, drie nieuwe Aardschuddingen, ongelijk sterker, dan die van den 2den dezer maan. De Schokken wien 20 geweldig, en de Huisen kraakten en verreesen, dat het grootste gedeelte der Bewoonders daarvan verliet, om zig op Pleinen, en in Tuinen te retireeren. Men heeft bereids publieke Gebouwen ingefield tot afweering van deze Plag, waarop men des te meer hoop heeft, doordie'er tans een heilzame regen valt, die de op een gepakte Dampen, door de droogte zedert de maand Maart, kan doen verdwynen, aan welke men de ontvanning der zwavelagtige stoffen onder den grond voornamlyk toefschryft.

10/27/03

ABU, Krahmer & Van den Bosch

9

1732

ITALIEN.

Romen den 19 July. Gisteren-morgen is de zaak van den Heer Saraini, door de Vergadering de Nonnullis afgehandeld; en men verzeekerd, dat hy van zyn Prelaatschap gedegradeerd is, en 10 jaren in 't kasteel van St. Angelo zitting zal. Genna den 19 July. Zaterdag-avond ging een van onze gewapende barken wederom zeyl na Corfica; van waerzondag nog 5 transportcheepen met 200 muylzels en eenige Keyzerlyke Officieren, overquamen. In 't begin van de week arriveerde ook van Bonifacio onze Commiffaris-Generael de Graef Gentile, en gisteren-morgen de Prins van Wurtenberg, die met alle teekenen van eerbij 'taen land stipten, door twee Heeren Gedeputeerden ontfangen, en in 't Carmeliter-klooster geleyd wierd, alwaer een logement voor zyn Hoogheyt gereed gemaakt was.

10/27/03

ABU, Krahmer & Van den Bosch

10

1655

SPANGIEN.

Adix den 27 Novem: 1655. Den Commandeur Sibion de Wilde heeft in zee veroverd en alhier opgedracht een treffelijk Curcksch Schip/ gemonteert met 32 stukke/ op hebvende 240 Curcken/ en by de 40 Chiften Slaven. De Curcken sullen/ tot goedmakinge van de Onkosten en hoort tot de zee van de Officieren ende gemeene Maets/ verhoort worden. De Chiften Slaven sijn met een Hollands Schip/ komende van Venetien/ en gaende naer Amsterdam/ op Eerghieren van hier nae Luyck geflyt. Den Hollantschen Vice Amirael de Kuprer is teghenwoordigh 10 Ologgh-schepen ende 2 Achten sterck.

10/27/03

ABU, Krahmer & Van den Bosch

11

Information Extraction

- (Jackson & Moulinier, 3.1)
- Information retrieval = vind documenten relevant voor query
- Information extraction = extraheer specifieke informatie-“snippets” uit vrije tekst
- Gebruikt oppervlakkige NLP

10/27/03

ABU, Krahmer & Van den Bosch

12

Achterliggend principe: analogie

- Conventions worden uitgedrukt door
 - Functiewoorden en frasen, in conventionele patronen
 - Inhoud op conventionele plaatsen in die patronen
- Analogieprincipe (De Saussure, 1916)
 - A staat tot B als A' tot B'
 - Wanneer A en B een vergelijkbare structuur bevatten, dan bevatten A en B hetzelfde type informatie

10/27/03

ABNL, Krahmer & Van den Bosch

13

Vereenvoudigingen

- Specifieke informatie?
 - Informatie binnen een bepaald domein
 - Terrorisme
 - Financiële transacties
 - Personele wijzigingen
 - Vacatures
- Oppervlakkige NLP?
 - "oppervlakkig ontleden", zoals op middelbare school
 - Herkennen van
 - eigennamen (personen, locaties, organisaties)
 - patronen

10/27/03

ABNL, Krahmer & Van den Bosch

14

Conventies in domeinen

- Patronen in berichtgeving

The Oceanport said its fourth-quarter profit was \$1.7 million.
Gateway said its fourth-quarter profit was 12 US cents.
Kaiser said its first-quarter loss was exacerbated by \$7 million.
Dan River Inc. said its second-quarter loss was \$6.1 million.
IBM said its first-quarter profit was \$1.2 billion.
Compaq said its first-quarter profit was 12 cents per share.
Cable operator Comcast said its second-quarter profit was four cents a share.
Chinadotcom said its third-quarter loss was \$20.5 million.

- [Company] said its [nth] quarter [loss/profit] was [\$ (per/a share)]

10/27/03

ABNL, Krahmer & Van den Bosch

15

Named entity recognition

- Wie, wat, waar?
- Tilburg corpus: 120 M woorden uit nieuwsartikelen
- 83% namen correct herkend, 88% persoonsnamen
- Oplossing: patronen

burgemeester van [GEMEENTE]
, zei [PERSOON].
de aandelen van [BEDRIJF]

10/27/03

ABNL, Krahmer & Van den Bosch

16

Named entity recognition (2)

- Beginpunt:
 - 53k namen van het web geplukt (41k pers., 5k loc., 7k orgs.)
 - >3M namen in 120M woorden krantenartikelen
- Cyclus:
 1. Zoek patronen bij niet-ambigue namen
 2. Zoek nieuwe namen in meest zekere patronen
 3. Voeg namen toe, ga terug naar 1

10/27/03

ABNL (Buchholz & Van den Bosch, 1999)

Commerciële IE

- Succesverhaal: <http://www.flipdog.com>
- IE-fabriek:
 - Webcrawling zoekt naar pagina's met vacatures (text classification)
 - Permanent 500,000 in stock
 - Informatie-extractie: type baan, bedrijf, locatie, opleidingsniveau, salaris, ...
 - Opslag in en retrieval uit web-toegankelijke database
 - Wekelijks verversen: nieuwe erin, verlopen vacatures eruit

10/27/03

ABNL, Krahmer & Van den Bosch

18

Informatie-extractie uit vacatures

- Wederom conventies:

Job type at top, in bold
"System analyst"
Candidate profile:
"must possess a [BS] degree in [chemistry]."
Salary offer
"Salary: [\$, \$\$, competitive]"
Contact info
"Contact: [person, email, phone]"

- Conventies kunnen automatisch geleerd worden op basis van voldoende voorbeelden

10/27/03

ABVI, Krahmer & Van den Bosch

19

Bonusopdracht 1

- "Precisie van search engines"
- Google, Altavista, Alltheweb
 - <http://www.google.com>
 - <http://www.altavista.com>
 - <http://www.alltheweb.com>

10/27/03

ABVI, Krahmer & Van den Bosch

20

Opdracht (vervolg)

- Google:
 - Grote coverage (meer dan 3 miljard webpagina's)
 - Page ranking, hubs/authorities
 - Standaard AND, maar zie *Advanced search*:
 - AND (all of the words ...)
 - OR (at least one of the words ...)
 - NOT (without the words ...)
 - Exact phrase

10/27/03

ABVI, Krahmer & Van den Bosch

21

Opdracht (vervolg)

- Altavista:
 - Minder grote coverage
 - Meer opties onder *More precision*
 - AND, OR, NOT, Exact phrase

10/27/03

ABVI, Krahmer & Van den Bosch

22

Opdracht (vervolg)

- Alltheweb:
 - Dekt bijna evenveel als Google
 - AND, OR, NOT, Exact phrase
 - Boolean queries
 - "RANK": "pac rank man" levert pagina's met "pac" op, en zo mogelijk ook met "man"

10/27/03

ABVI, Krahmer & Van den Bosch

23

Opdracht (vervolg)

- Doel: door het doen van queries een zo hoog mogelijke *precision* halen op tekstclassificatie.
- Twee classificaties:
 - Nederlandse pagina met CV (wel/niet)
 - Nederlandse pagina met cursusinformatie (wel/niet)

10/27/03

ABVI, Krahmer & Van den Bosch

24

Opdracht (vervolg)

- Voorbeeldqueries:
 - (gehuwd OR ongehuwd) AND opleiding AND werkervaring
 - docent AND (literatuur OR syllabus)
- Verzin zelf minimaal drie queries per taak

10/27/03

ABV, Kraemer & Van den Bosch

25

Opdracht (vervolg)

- Per query/taak/search engine: kijk naar de eerste 20 links
- Bereken *precision*: van de 20 gevonden pagina's, hoeveel zijn er daadwerkelijk CVs / cursuspagina's?
- Als er meer dan alleen CV / cursusinformatie op staat, dan telt dat als een FOUT.

10/27/03

ABV, Kraemer & Van den Bosch

26

Opdracht (vervolg)

- Rapporteer dus:
 - Voor Google, Altavista en Alltheweb
 - Precisie met minimaal 3 queries (dezelfde voor de 3 search engines)
 - Bediscussieer verschillen tussen
 - Search engines
 - Queries
 - Let op:
 - Gebruik drie hele verschillende queries

10/27/03

ABV, Kraemer & Van den Bosch

27

Opdracht (vervolg)

- Rapportage naar keuze in papieren of elektronische vorm
- Inleveren bij Antal van den Bosch / antalb@uvt.nl
- Met 1 of 2 personen
- Deadline maandag 13 oktober
- Bespreking op 14 oktober, college Text classification

10/27/03

ABV, Kraemer & Van den Bosch

28