

## Informaatiojärjestelmät, tietotulva ja tiedon louhinta

Hannu Toivonen

Tietojenkäsittelytieteen laitos  
Hannu.Toivonen@cs.helsinki.fi

1 TKTL\_S2004.PPT

## Informaatiotulva

- Vuoden 2000 aikana tuotettiin 3 exatavua dataa
  - kilotavu = 1024 tavua
  - megatavu =  $1024 \times 1024$  tavua
  - ...giga, tera, peta...
  - exatavu =  $1024^6$  tavua  $\approx 10^{18}$  tavua
- Datat määrä kaksinkertaistuu vuosittain
- Informaatiota esitetään moninaisissa muodoissa
  - relaatiotietokannat
  - teksti (Google-hakukone tuntee 4.3 miljardia sivua)
  - mittaus- ja lokitietokannat
  - geneettiset aineistot (ihmisen dna: 3 miljardia emäsparia)
  - ...

2 TKTL\_S2004.PPT

## Informaatiojärjestelmät

- Informaation hallinta
  - tiedon tallettaminen
  - tiedon esittäminen
  - tiedonhaku
  - tiedon analysointi
- TKTL:n "info"-linja



3 TKTL\_S2004.PPT

## Infon opetus

- Tietokantojen perusteet
- Tietokannan hallinta
- Johdatus sovellussuunnittelun
- Digitaalisen median tekniikat
- XML-metakielii
- Tutkimustiedonhallinnan peruskurssi
- Tietokannan mallinnus
- Tietokantarakenteet ja -algoritmit
- Tiedon louhinnan menetelmät
- Tiedon louhinnan erikoiskurssi
- Tietovarastot
- Tiedonhakumenetelmät
- Rakenteisten dokumenttien käsitteily
- seminaareja
- ...

} cum laude  
} laudatur

4 TKTL\_S2004.PPT

## Infon tutkimus

- Tiedon louhinta (Hannu Toivonen, Helena Ahonen-Myka, Marko Salmenkivi)
  - (tästä tarkemmin seuraavilla kalvoilla...)
- Dokumentit ja kieliteknologia (Helena Ahonen-Myka, Greger Linden)
  - rakenteisten dokumenttien hallinta, tiedonhaku, tiedon eristäminen tekstillä, tekstin louhinta
- Tietokannat (Seppo Sippu, Harri Laine)
  - tietojen mallintaminen, samanaikaisen käytön valvontan ongelmat, tietokantarakenteiden elvytys, tietokantasovellusten suunnitteluvälineet ja toteutusmenetelmät

5 TKTL\_S2004.PPT

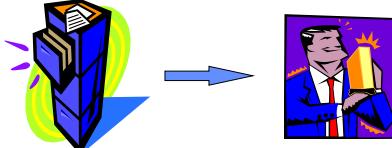
## Sisällysluettelo

- Informaatiotulva
- Informaatiojärjestelmien opetus ja tutkimus laitoksella
- Tiedon louhinta
- Esimerkkisovellus: sairausgeenien paikannus
- Esimerkkimenetelmä: assosiaatiosäännöt
- Tiedon louhinta tieteentalana
- Tiedon louhinnan tutkimus laitoksella
- Yhteenvetö

6 TKTL\_S2004.PPT

### Tiedon louhinta

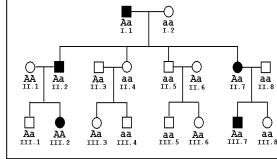
- Uuden ja hyödyllisen tiedon päätelemisen suurista datamassooista



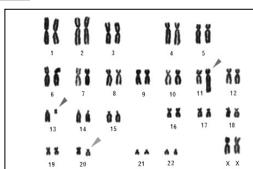
- "Moderni data-analyysi" tai "algoritminen tilastotiete"
- "Mitä data kertoisi, jos siltä osaisi kysyä oikeat kysymykset?"

7 TKTL\_S2004.PPT

### Sairausgeenien paikannus



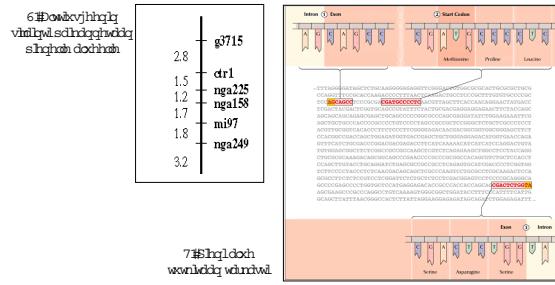
4 Bxoxesxel wdnwWhip daxl  
kdyldldq fwhfi dwxevvdlrlwd  
wrxwedydq wdoWq shulw |



5 #Ship fwfi sdlnodgihdak  
nux rvp 1#rcai doamdy  
jninglrg

8 TKTL\_S2004.PPT

### Sairausgeenien paikannus



6 Bowekvhingq  
vihnlq sdlnodgihdak  
sighnq dloohnh

g3715

2.8  
drl  
1.5  
nga225  
1.2  
nga158  
1.7  
mi97  
1.8  
nga249  
3.2

7 Sholdoh  
wrxwedydq wdoWq

- Tiedon louhinnalla voidaan havaita esim. yhteys kromosomin tietyn kohdan ja sairauden välillä

9 TKTL\_S2004.PPT

### Assosiaatiostäällöt

- Alkuperäinen ongelmatyyppi: mitä tavaroita ostetaan usein yhdessä?
- Ostoskorianalyysi
  - Jos vaippoja niin olutta (todennäköisyys 56 %, frekvenssi 12 %)

10 TKTL\_S2004.PPT

### Assosiaatiostäällöt

- 1. yleistys: mitkä asiat esiintyvät usein yhdessä? Eriisia sovelluskohteita:
  - kurssi-ilmostautumiset
    - Jos tietoliikenne ja UNIX-ohjelmointi niin C-ohjelmointi (tod.näk. 72 %, frekv. 6 %)
  - tekstidokumenttien analysointi
    - Jos "www" ja "netscape" niin "browser" ja "internet" (tod.näk. 89 %, frekv. 0.12 %)
  - geneettisten markkerit ja perinnöllinen sairaus
    - Jos "marker9" ja "marker33" ja "tupakoi" niin "sairas" tod.näk. 34 %, frekv. 8 %)

11 TKTL\_S2004.PPT

### Assosiaatiostäällöt

- Tavoitteena on kuvilla mahdollisesti mielenkiintoisia yksinkertaisia ilmiöitä
- Menetelmä tuottaa kaikki assosiaatiostäällöt, joilla frekvenssi > kynnysarvo
- Mahdollisia säätöjä on valtavasti, läpikäynti käsin olisi mahdotonta
- Joukossa voi olla yllättäväkin säätöjä
- Tiedon louhintaan liittyvä ongelma: miten autetaan käyttäjää löytämään juuri häntä kiinnostavat säällöt?
- Menetelmä on sovellusriippumaton

12 TKTL\_S2004.PPT

## Muita tiedon louhinnan muotoja

- Säännönmukaisuuksien etsintä
  - Millaiset hahmot ovat aineistossa tyypillisiä?
- Klusterointi
  - Millaisiin luonteviin ryhmiin aineiston voi jakaa?
- Luokittelu, ennustaminen
  - Miten havaintojen tietyn ominaisuuden voi ennustaa havainnon muista ominaisuuksista?
- Poikkeuksien etsiminen
  - Mitkä havainnot vaikuttavat poikkeuksellisilta?

13 TKTL\_S2004.PPT

## Tiedon louhintaprosessi

- Tiedon louhinnassa tutkitaan algoritmien lisäksi myös koko analyysiprosessia



14 TKTL\_S2004.PPT

## Tutkimusmenetelmä

- Esimerkkinä assosiaatiosäännöt
- Algoritmikehitys:
  - assosiaatiosäännöt
    - episodisäännöt (assosiaatiot tapahtumajonoissa)
    - yleinen menetelmärunko
- Teoreettinen kehitys:
  - konkreettinen ongelma (ostoskorianalyysi)
  - yleistetty ongelmatyyppi (toistuvat ilmiöt)
  - tehtävätyyppi ja ratkaisuvaihtoehtojen analyysi

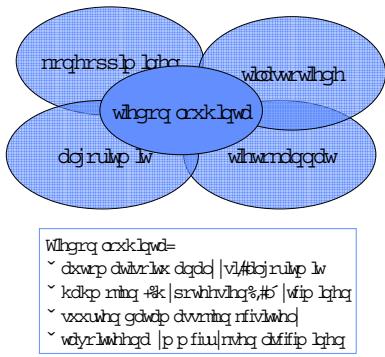
15 TKTL\_S2004.PPT

## Assosiaatiosäännöt

- 2. yleistys: mitkä hahmot esiintyvät aineistossa usein?
  - Syöte
    - r: tietokanta
    - P: suuri joukko hahmoja tai hahmojen "kieli"
    - k: yleisyyden kynnysarvo
  - Tulos
    - kaikki joukon P hahmot, joiden yleisyys ylittää kynnysarvon k tietokannassa r
  - Analyysi
    - ongelman teoreettisista ominaisuuksista seuraa, että tietty yksinkertainen algoritmi on tehtävään optimaalinen (tiettyllä oletuksilla)

16 TKTL\_S2004.PPT

## Tiedon louhinnan lähinaapurit



18 TKTL\_S2004.PPT

## Millaisista taidoista on hyötyä

- algoritmiikka
- todennäköisyyslaskenta
- tilastotiede
- tietokannat (??)
- koneoppiminen
- sovellusalueen tuntemus
  - poikkitieteellisillä taidolla iso tutkimuspotentiaali
- tiedon louhinta ei ole helppoa: jokainen ongelma vaatii luovuutta ratkaisujen kehittämisessä ja soveltamisessa

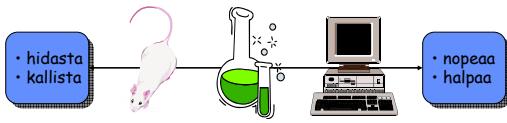
19 TKTL\_S2004.PPT

## Tiedon louhinta ja TKTL

- Informaatiojärjestelmien linja
  - geenikartoitusmenetelmät, geneettisen datan analyysi
  - ekologiset data-analyysiongelmat (mm. ilmaston rekonstruointi)
  - hahmokielet, algoritmikehitys
  - tekstien ja dokumenttirakenteiden louhinta
- mukana laitoksella toimivissa "virtuaaliorganisaatioissa"
  - FDK-huippuyksikkö (From Data to Knowledge) tiedon louhinnan ja hahmonsovitukseen "kattoprojekti"
  - HIIT/BRU
    - data-analyysi, proaktiivinen laskenta

21 TKTL\_S2004.PPT

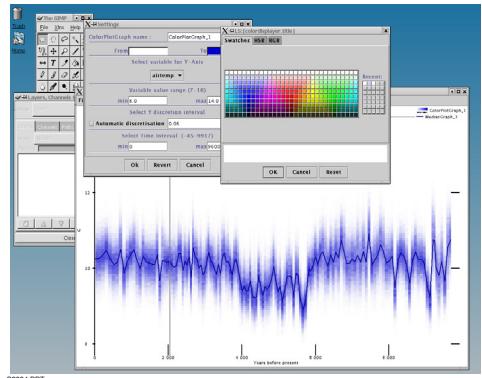
## Syövälle altistavien yhdisteiden tunnistaminen



- Kansainvälinen "haastekilpailu" tiedon louhijoille
  - järjestäjänä mm. NIH Yhdysvalloista
  - tieteellinen sokkotesti
- Mallinnus- ja ennustusongelma
- Mallien ja tulosten testaaminen ja arviointi

22 TKTL\_S2004.PPT

## Ilmaston rekonstruointi



23 TKTL\_S2004.PPT

## Yhteenveto tiedon louhinnasta

### Tiedon louhinta tieteenalana

- tuottaa ja tunnistaa erilaisia datan automaattiseen analysointiin ja kuvaleimiseen liittyviä tehtävätyyppejä tai lähestymistapoja
- analysoi ja kategorisoi niitä
- kehittää niihin tehokkaita ratkaisuja
- myös: tietosuoja ja etiikka tiedon louhinnassa

Laitoksella kansainvälisti korkealaatuista tutkimusta

Runsaasti tieteellisiä yhteistyöprojekteja

24 TKTL\_S2004.PPT