

Aurkibidea

1. Datu-Meatzaritzarako sarrera

Irrakasgia: Datu Meatzaritzat
Especialitatea: Konputazioa, hirugaren ikasmarria
Titulazioa: Informatika Ingeniaritzako Gradua
Konputazio Zientzia eta Adimen Arifitziala saila
Universidad del País Vasco - Euskal Herriko Unibertsitatea

Aurkibidea

Datu-Meatzaritzia

Definizio batzuk

- Data mining.** Aurretik ezezaguna den ezagutzaren erabilgarri eta ulergarria erauztua formatu desberdinetan gordeta dauden datu-multzo handietatik (Witten eta Frank, 2000) **Knowledge discovery in databases.** Datubaseetan ezagutza aurkitzea, datuetatik abiaturu eredu baliagarriak, berriak eta azken batean ulergarriak identifikatuak (Fayyad eta lao, 1996)

Edu-motak

- ▶ **Datuatik ezagutza** lortzea eredu konputazioalak erabiliz
 - ▶ **Eedu deskribatzaleak:** datuak deskribatzaten edo laburten dira
 - ▶ **Erregeak:** datuen portaera-ereduak erakusten dira
 - ▶ **Clustering:** homogenoak diren kasuak multzokatzen dira
 - ▶ **Eedu iragarleak:** aldagai iragarleen baloeratik abiatuz iragarri beharreko beste aldagai batzuren baliok estimatzenten dira

Datu-motak

- **Datubase erlazionak**
 - Erlazio-multzoak (taulak). Ezagari-multzo baterako (aldagaiak, Zutabeakak, erremuak) n-koleak gozdezen dira (kasuluak, erreinkakadak, erregistriak)
 - Taula moduko aurkepena: ezagari-balio
 - **Datubase espazioak**: geografia-datuak, medizina-zentzu-irudiak, garraio-sareak, ...
 - **Datubase temporalaik**: denborako une edo tarte desberdinak
 - **Datubase dokumentalak**: Objektuak, testu-dokumentuak dira, aldagueiak hitzak adierazi dituzte, edo laburpenak...
 - **Multimedia datubaseak**: irudiek, soinua, bideoa

■ **World Wide Web**, gauegian dagoren informazioariko guneenik handiena eta anitzetan

 - Edukarien analisia: web orietatuen ereduak aukitztea
 - Egizialarien analisia: hipertextuak eta URLak aztertzea
 - Erabileraaren analisia: pertsonaia, erabiltzailea, erabili...

Erlazioa beste zenbait arlorekin	Aplikazioak
<ul style="list-style-type: none"> ▶ Estatistika. Datu-meatzaritzaren "ama" ▶ Ikasketa automatikoa. Konputagailuak adibideetatik ikasten du ▶ Ereduen ezagutza. Clustering. Sailkapen-gainbegiratua ▶ Erabakiak haritzeko sistemak. Zuzendaritzari laguntzeko tresnak era sistemak ▶ Datuaren bistaraztea. Grafiken bidez datuetatik ateratako ereduak ikusten eta ulertzen laguntza ▶ Datubaseak. Datuen biltegiak. Atzipen eraginkorra ▶ Informazioaren erauzketa. Testuak. Liburutegi digitalak. Bilaketaik Interneten ▶ Konputazio paraleloa eta banatu. Datu-meatzaritzak dakaren kostu konputazionala banatu egiten da prozesamendu paraleloa eta banatuera erabiliz 	<p>7/22</p> <p>7/22</p>

Erlazioa beste zenbait arlorekin	Aplikazioak
<ul style="list-style-type: none"> ▶ Estatistika. Datu-meatzaritzaren "ama" ▶ Ikasketa automatikoa. Konputagailuak adibideetatik ikasten du ▶ Ereduen ezagutza. Clustering. Sailkapen-gainbegiratua ▶ Erabakiak haritzeko sistemak. Zuzendaritzari laguntzeko tresnak era sistemak ▶ Datuaren bistaraztea. Grafiken bidez datuetatik ateratako ereduak ikusten eta ulertzen laguntza ▶ Datubaseak. Datuen biltegiak. Atzipen eraginkorra ▶ Informazioaren erauzketa. Testuak. Liburutegi digitalak. Bilaketaik Interneten ▶ Konputazio paraleloa eta banatu. Datu-meatzaritzak dakaren kostu konputazionala banatu egiten da prozesamendu paraleloa eta banatuera erabiliz 	<p>7/22</p> <p>7/22</p>

Erlazioa beste zenbait arlorekin	Aplikazioak
<ul style="list-style-type: none"> ▶ Estatistika. Datu-meatzaritzaren "ama" ▶ Ikasketa automatikoa. Konputagailuak adibideetatik ikasten du ▶ Ereduen ezagutza. Clustering. Sailkapen-gainbegiratua ▶ Erabakiak haritzeko sistemak. Zuzendaritzari laguntzeko tresnak era sistemak ▶ Datuaren bistaraztea. Grafiken bidez datuetatik ateratako ereduak ikusten eta ulertzen laguntza ▶ Datubaseak. Datuen biltegiak. Atzipen eraginkorra ▶ Informazioaren erauzketa. Testuak. Liburutegi digitalak. Bilaketaik Interneten ▶ Konputazio paraleloa eta banatu. Datu-meatzaritzak dakaren kostu konputazionala banatu egiten da prozesamendu paraleloa eta banatuera erabiliz 	<p>9/22</p> <p>9/22</p>

Erlazioa beste zenbait arlorekin	Aplikazioak
<ul style="list-style-type: none"> ▶ Garesti gerta daitzekeen bezeroak identifikatzeara ▶ Poliza berriak kontratzen dituzten bezero-motak identifikatzeara ▶ Arrisku egoeran egon daitzezen bezeroen portera-ereduak identifikatzeara ▶ Iruzurrerako porterarek identifikatzeara 	<p>10/22</p> <p>10/22</p>

Erlazioa beste zenbait arlorekin	Aplikazioak
<ul style="list-style-type: none"> ▶ Garesti gerta daitzekeen bezeroak identifikatzeara ▶ Poliza berriak kontratzen dituzten bezero-motak identifikatzeara ▶ Arrisku egoeran egon daitzezen bezeroen portera-ereduak identifikatzeara ▶ Iruzurrerako porterarek identifikatzeara 	<p>11/22</p> <p>11/22</p>

Erlazioa beste zenbait arlorekin	Aplikazioak
<ul style="list-style-type: none"> ▶ Garesti gerta daitzekeen bezeroak identifikatzeara ▶ Poliza berriak kontratzen dituzten bezero-motak identifikatzeara ▶ Arrisku egoeran egon daitzezen bezeroen portera-ereduak identifikatzeara ▶ Iruzurrerako porterarek identifikatzeara 	<p>12/22</p> <p>12/22</p>

Aplikazioak

- Bioinformatika, bioingeniaritza**
 - ▲ Geneen bilaketa (genoma kodifikatzaren duten eremuak)
 - ▲ Proteinen bigarren mailako egituraaren iagarrapena egitea
 - ▲ Uhodeen iragarpena
 - ▲ Uren kalitatea neuritzeko ereduak sortzea

13/12/22

14/12/22

15/12/22

Aplikazioak

Beste zenbait arlotan

- ▲ Telekomunikabideak: iruzura identifikatzea
- ▲ Posta elektronikoa eta agenda personalak: postaren sailkapena eta banaketia automatikoa, spam mezuk hautematea
- ▲ Ogasuna: zerga-iruzurrak hautematea
- ▲ Web: erabiltzailen portaeraren analisia egitea
- ▲ Kirotak: datu medikuatik abiatuz lesioak jasateko arriskua identifikatzea

Aurkibidea

Datu-Meatzaritza

Knowledge Discovery from Databases (KDD): datubaseetan ezagutza aurkitzea

16/12/22

17/12/22

Datubaseetan ezagutza aurkitzea

1. Datuak bildu eta bateratu

Prozesu iteratibo eta elkarreragilearen faseak

1. Datuak bildu eta bateratu
2. Aukeratu, garbitu eta eraaldatu
3. Datu-meatzaritza egin
4. Ebaluatu eta interpretatu
5. Zabaldu eta erabili

16/12/22

17/12/22

Datubaseetan ezagutza aurkitzea

2. Aukeratu, garbitu eta eraaldatu

- ▲ Aurkitutako ezagutzaren kalitatea datu-meatzaritan erabilitako algoritmoaren mende egoeteaz gain **aztertutako datuen kalitatearen mende** dago
- ▲ Datuen portera orokorrera egokitzten ez den zenbait daturen presentzia egon daiteke (**outliers**)
- ▲ Biloak falta dituzten datuak egon daitezke (**missing values**)
- ▲ Aldagai esanguratsuak aukera daitezke (**feature selection**)
- ▲ Datubase oso handietan **kasuen auzazko aukeraketa** egin daiteke
- ▲ Aldagai jarraiatik **diskretiliza** daitezke

16/12/22

17/12/22

Datubaseetan ezagutza aurkitzea

3. Datu-meatzaritza egin

- ▲ **Datubaseak eta datuen prozesamendu tradizionala** (On-Line Transaction Processing, OLTP): eguneroko beharrak asetzeko nahiakoa dira (fakturazioa, inventarioen kontrola, ...)
- ▲ **Ertabaki estrategikoak** analisian, plangintzan eta iragarpenean oinarrituta: datuak sail desberdinetan egon daitezke
- ▲ **Jatorrizko datuak** formatu desberdinetan egoten dira
- ▲ Datubaseen batera: **datuen biltegiak** (data warehousing)
- ▲ Datuen biltegiak gomendagarriak dira informazio-kantitatea handia denean. Hala ez denean ez dira beharrezkoak (**testu-fixategiak, kalkulu-omnik, ...**)

17/12/22

18/12/22

Datubaseetan ezagutza aurkitzea

- ▲ **Aurkitutako ezagutzaren kalitatea datu-meatzaritan erabilitako algoritmoaren mende egoeteaz gain **aztertutako datuen kalitatearen mende** dago**
- ▲ Datuen portera orokorrera egokitzten ez den zenbait daturen presentzia egon daiteke (**outliers**)
- ▲ Biloak falta dituzten datuak egon daitezke (**missing values**)
- ▲ Aldagai esanguratsuak aukera daitezke (**feature selection**)
- ▲ Datubase oso handietan **kasuen auzazko aukeraketa** egin daiteke
- ▲ Aldagai jarraiatik **diskretiliza** daitezke

Datubaseetan ezagutza aurkitzea

3. Datu-Meatzaritza

- ▶ **Erodu deskribatzaleak**
 - ▶ Eregeiak
 - ▶ **Multzikatzea** (clustering): partizionala, probabilistikoa, ierarkikoa, kontzeptuala
 - ▶ **Erodu iragartzaileak**
 - ▶ Eregresioa: regresio lineala...
 - ▶ Gainbegiraturako **sailkapena**: sailkapen-zuhaitzak, K-NN, sailkatzaille Bayestarrak, erregelen indukzioa, erregresio logistikoa, sare neuronak, sailkatzaleen konbinaketa
 - ▶ Erodu ulerterrazak

Datubaseetan ezagutza aurkitzea

4. Ebaluatu eta interpretatu

- ▶ Ebaluazio teknikak: **validazio simplea** (entrenamendua + testa), ***k*-geruzatako balidazio gurutzatua**, **bootstrapping**
- ▶ Eregeiak: **estalitza, konfidantza**
- ▶ Multzikatzea (Clustering): **multzo barriko elementuen eta multzoen arteko distantzia**
- ▶ Eregresioa: **batazbesteiko errore kuartatikoa**
- ▶ Gainbegiratutako sailkapena: ondo sailkatutakoen **porzentzia, erroreen matrizea, ROC analisia**
- ▶ Erodu ulerterrazak eta interesgarriak (erabilgarriak eta berritzaleak)

Datubaseetan ezagutza aurkitzea

5. Zabaldu eta erabili

- ▶ **Zabaldu**: erakikitako ereduia erabilitzaleen artean zabaldu
 - ▶ eta erabili erabakiak hartzeko
- ▶ **Ereduaaren garapena**, neurri behar da denboran zehar:
 - ▶ Berrebaliatu
 - ▶ Berrentenatu
 - ▶ Berreraiki

Oinarrizko bibliografía

- ▶ Liburua: **Introducción a la Minería de Datos**
 - ▶ Capítulo 1: ¿Qué es la minería de datos?
 - ▶ Capítulo 2: El proceso de extracción de conocimiento
- ▶ Egileak: José Herrández Orelló, M^a José Ramírez Quintana, César Ferri Ramírez
- ▶ Departamento de Sistemas Informáticos y Computación, Universidad Politécnica de Valencia
- ▶ Argitaraztalea: Pearson Prentice Hall, 2004
- ▶ ISBN: 84-205-4091-9