

Upravljanje poslovnim podacima



PROČIŠĆAVANJE PODATAKA
ETL

PROF. DRAŽENA GAŠPAR

03.11.2015.

PREZENTACIJE IZVORI PODATAKA



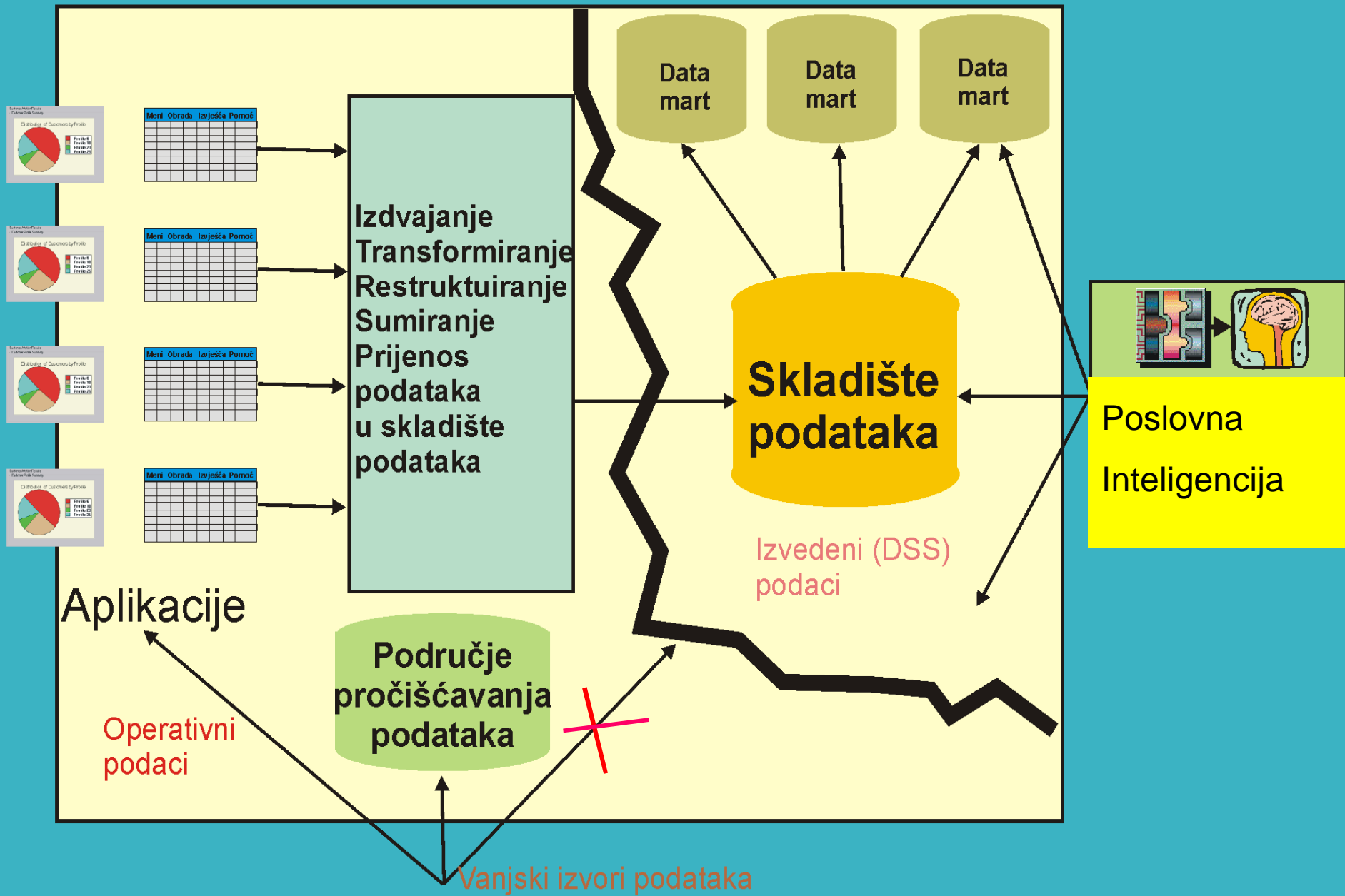
Max. 5 minuta



PITANJA ZA DISKUSIJU

- Što se podrazumijeva pod pojmom ETL?
- Uz koju tehnologiju/je se veže pojam ETL?
- Koja je osnovna svrha ETL-a?
- Što su to "prljavi podaci"? Primjer!
- Koji su osnovni koraci ETL-a?
- Načini učitavanja/ažuriranja podataka?
- Koje su značajke kvalitetnog podatka?
- Gdje ste tražili informacije o ovome?

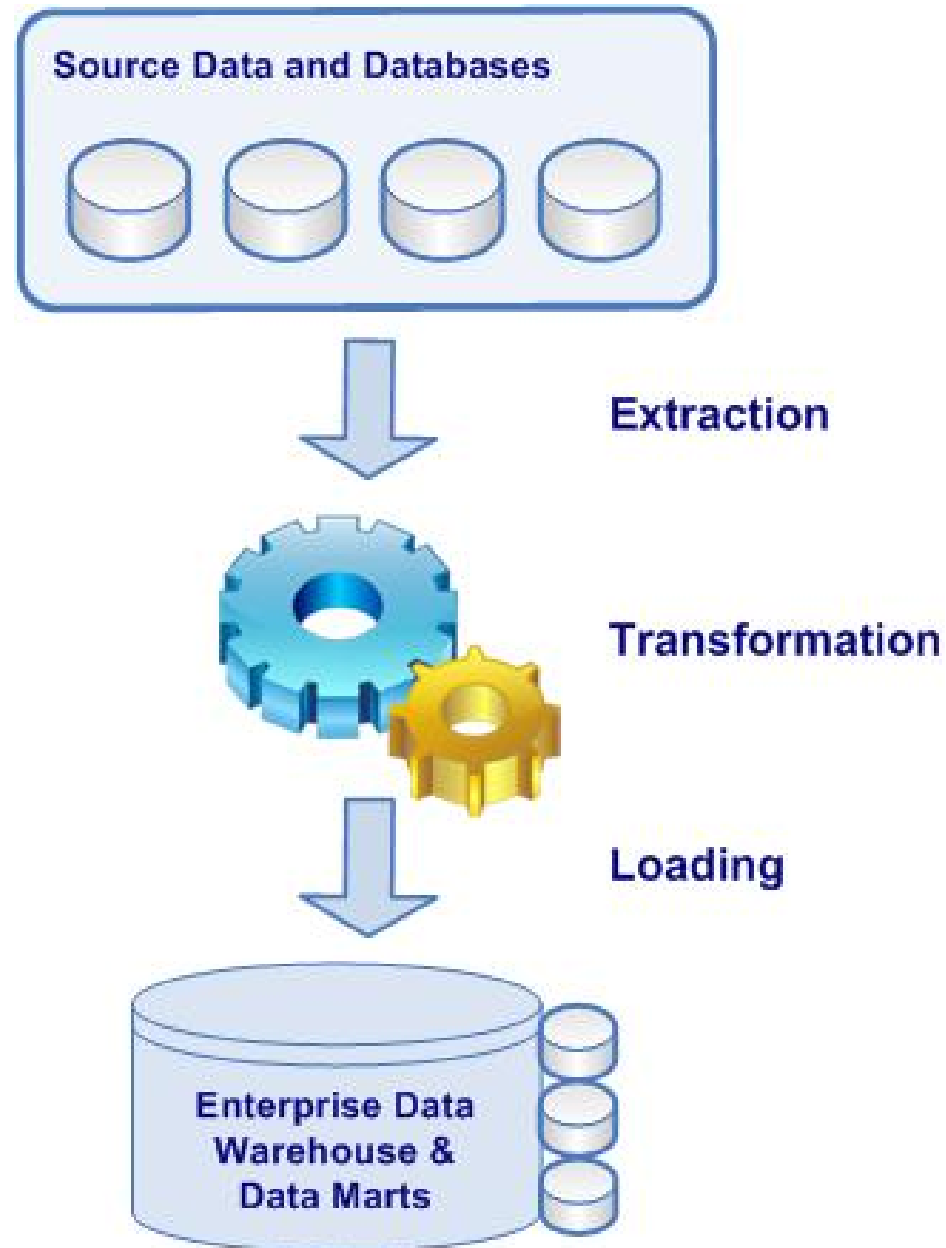


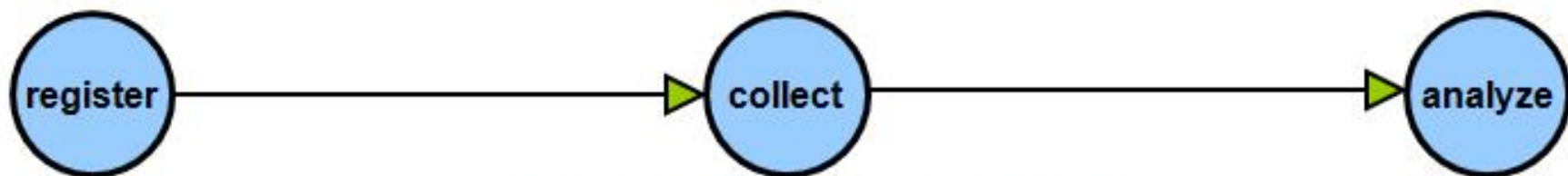


ETL

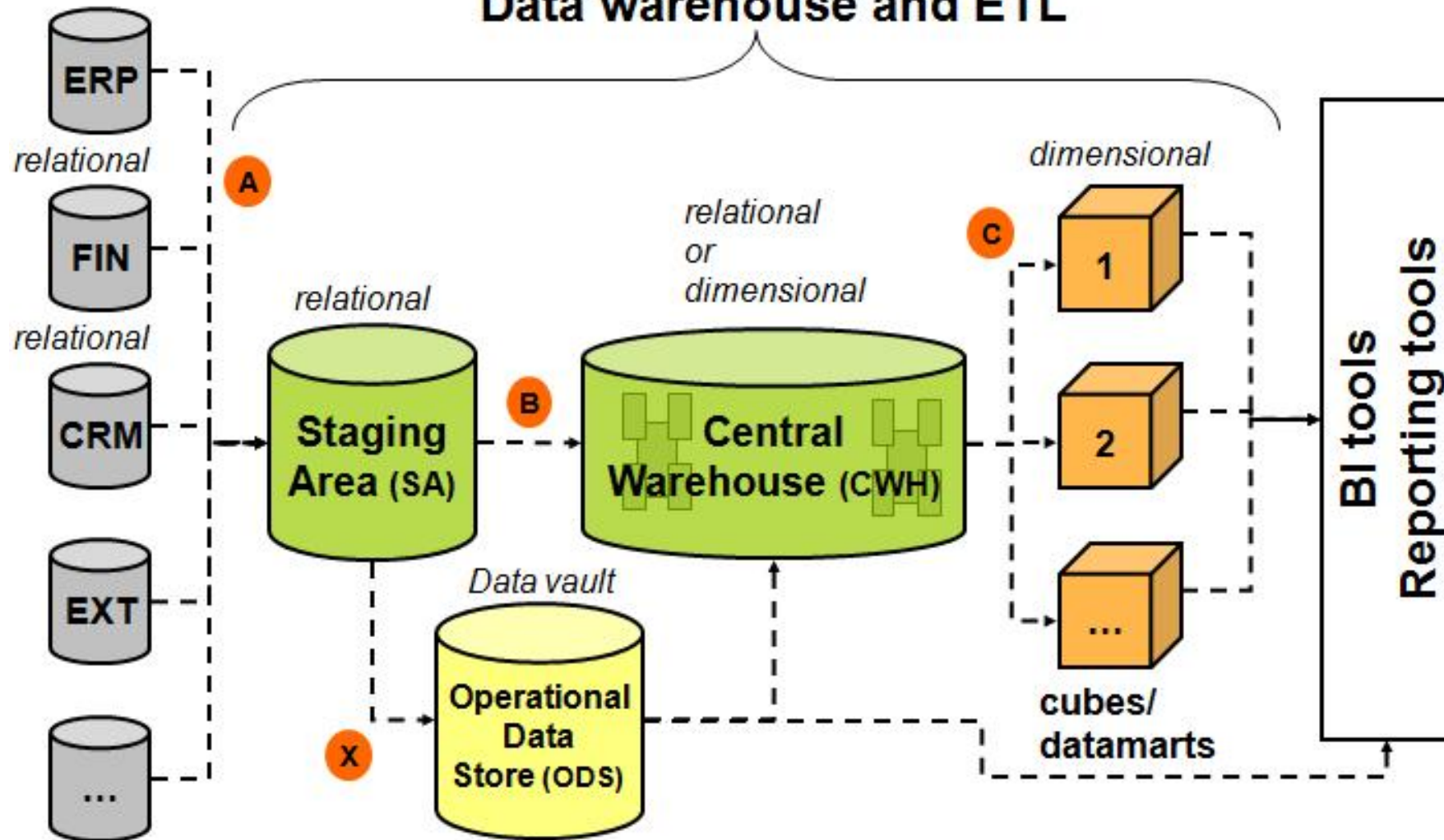


ETL

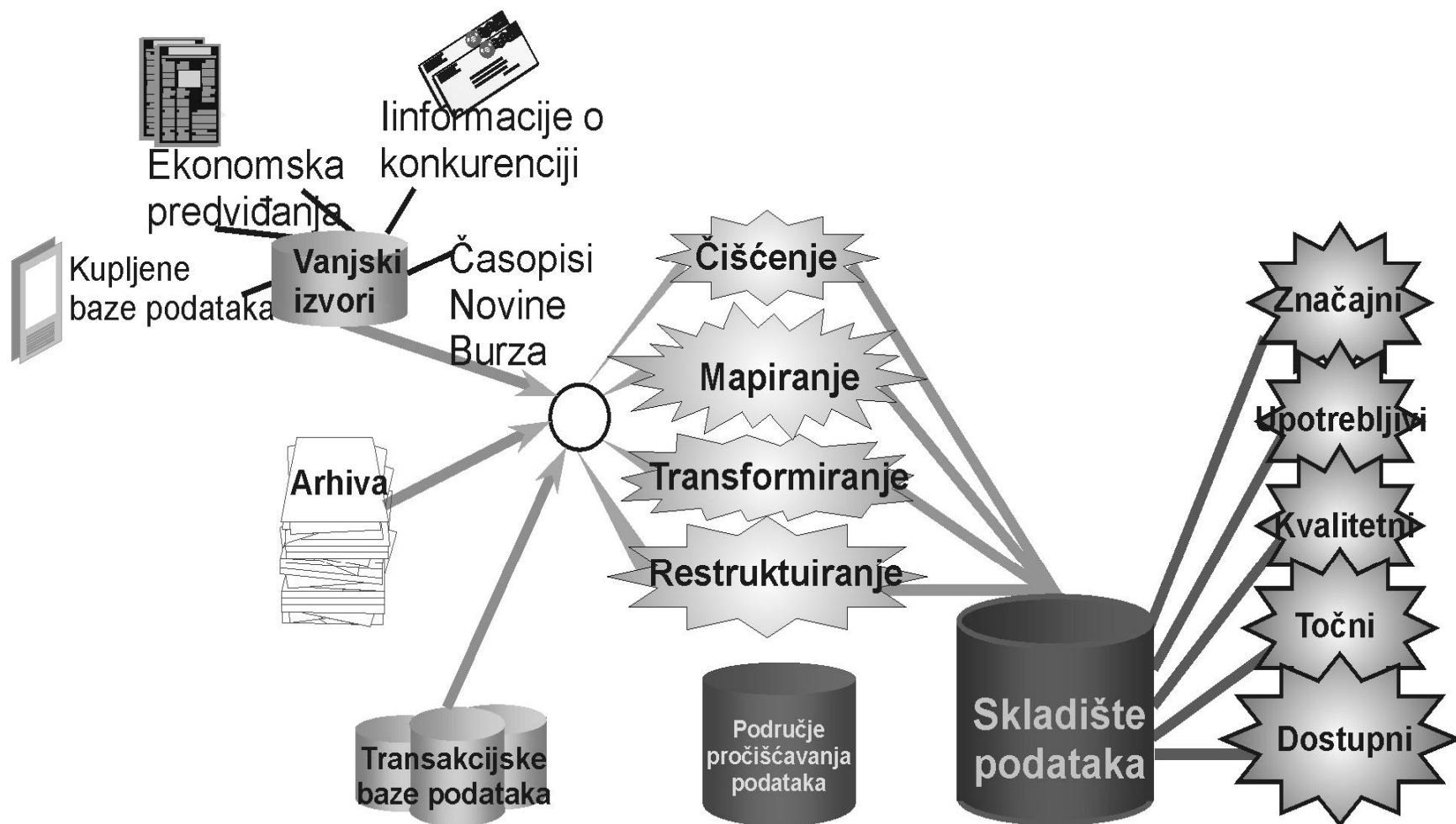




Data warehouse and ETL



PODRUČJE PROČIŠĆAVANJA PODATAKA



ZAŠTO JE NUŽNO PROČIŠĆAVANJE PODATAKA?

- Podaci u realnosti su “prljavi”
 - nepotpuni: nedostaju vrijednosti atributa, neostaju određeni atributi ili sadrže samo agregirane podatke
 - npr. zanimanje=“ ”
 - sa šumovima: sadrže greške ili vrijednosti izvan granica
 - npr. Plaća=“-10”
 - nekonzistenti: sadrže neusklađenosti u kodovima ili nazivima/imenima
 - npr. Godine=“42” Datum rođenja=“03/07/1997”
 - npr. bilo je rangiranje “1,2,3”, sada je “A, B, C”
 - npr. neusklađenost duplih slogova



PROBLEMI S INTEGRITETOM PODATAKA

- Ista osoba, različito napisano ime
Dražena, Dražana, Draženka ...
- Više načina označavanja naziva kompanije
 - Hera, Hera d.o.o., Hera – SW kompanija
- Uporaba različitih naziva
Mumbai, Bombay
- Različite šifre generirane od strane različitih aplikacija za istog kupca
- U obvezna polja unešen znak blank, . i sl.
- Pogrešna šifra proizvoda unešena na POS-u
 - Ručni unosi dovode do grešaka
 - “u slučaju problema koristiti 99999999”



ZAŠTO SU PODACI "PRLJAVI"?

- Nepotpuno podaci mogu proizići iz
 - "Nije primjenjivo" kao podatkovna vrijednost kod prikupljanja
 - Razlika u razmišljanju u vrijeme kada su podaci prikupljeni i kada se analiziraju.
 - Ljudski/hardware/software problemi
- Šumovi (netočne vrijednosti) podataka mogu proizići iz
 - Neispravni uređaji prikupljanja podataka
 - Ljudske ili računalne greške na unosu podataka
 - Greške u prijenosu podataka
- Nekonzistentnost podataka može proizići iz
 - Različiti izvori podataka
 - Narušavanje funkcijskih ovisnosti (npr. Izmjene povezanih podataka)
- Dupli slogovi također traže čišćenje



ZAŠTO JE NUŽNO PROČIŠĆAVANJE PODATAKA?

- Bez kvalitetnih podataka nema ni kvalitetnih rezultata analize!
 - Kvalitetne odluke morju biti temeljene na kvalitetnim podacima
 - Npr. Duplicirani ili nedostajući podaci mogu uzrokovati pogrešnu ili varljivu statistiku.
 - Skladište podataka treba konzistentnu integraciju kvalitetnih podataka
 - Proces izdvajanja, čišćenja i transformiranja podataka čini najveći dio posla izgradnje skladišta podataka



VIŠEDIMENZIJSKA MJERA KVALITETE PODATAKA

- Opće prihvaćeni višedimenzijski pogled:
 - Točnost
 - Potpunost
 - Dosljednost (konzistentnost)
 - Pravovremenost
 - Vjerodostojnost
 - Dodatna vrijednost
 - Interpretativnost
 - Pristup
- Opće kategorija:
 - Značajnost, odgovara kontekstu, reprezentativan, pristupačan



OSNOVNI ZADACI

○ Čišćenje podataka

- Popunjavanje vrijednosti koje nedostaju, rješavanje šumova, identificiranje ili uklanjanje "nepodobnih" i razrješavanje nekonzistentnosti

○ Integriranje podataka

- Integriranje višestrukih baza podataka, podatkovnih kocki ili datoteka

○ Transformiranje podataka

- Normalizacija i agregacija

○ Smanjivanje (redukcija) podataka

- Smanjivanje obima podatka uz zadržavanje istih ili sličnih analitičkih rezultata

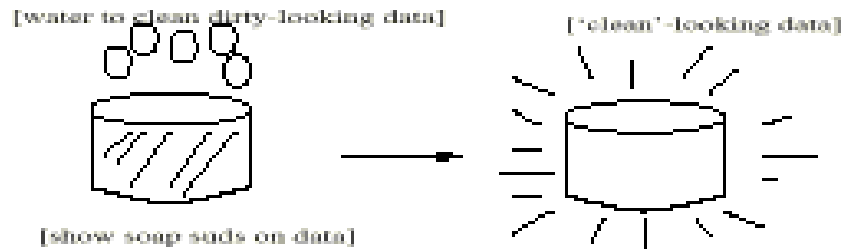
○ Diskretizacija podataka

- Dio smanjivanja (redukcije) podataka, ali posebice bitno za numeričke podatke

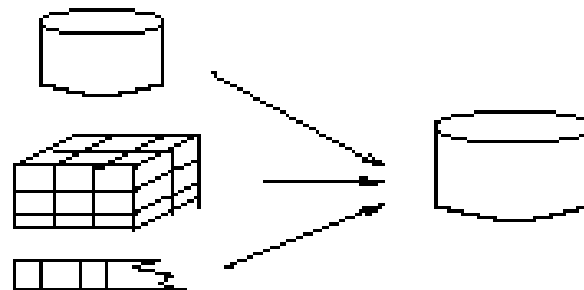


PRIMJERI OSNOVNIH ZADATAKA

Data Cleaning



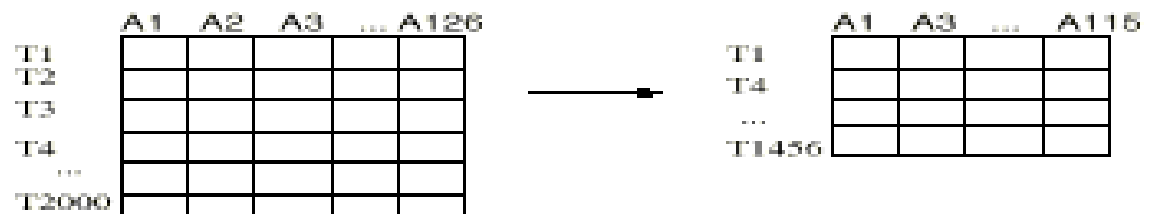
Data Integration



Data Transformation

-2, 32, 100, 59, 48 → -0.02, 0.32, 1.00, 0.59, 0.48

Data Reduction



ČIŠĆENJE PODATAKA

- Rješavanje
 - Nepotpunih podataka
 - Šumova u podacima
 - Nekonzistentnosti



NEPOTPUNI (NEDOSTAJUĆI) PODACI

- Podaci nisu raspoloživi
 - Npr. Mnoge n-torke nemaju zabilježene vrijednosti za neke od atributa (null polja)
- Podaci mogu nedostajati zbog:
 - Grešaka na opremi
 - Nekonzistentnosti s drugim podacima, pa su zbog toga izbrisani
 - Podaci nisu unešeni zbog nerazumijevanja
 - Određeni podaci nisu smatrani bitnim u vrijeme unosa
 - Ne prati se povijest promjena nad podacima



ŠTO RADITI S NEPOTPUNIM PODACIMA?

- Ignorirati n-torku: obično se radi kada nedostaje oznaka klase (kada se radi klasifikacija)—nije efektivno kada postotak nedostajućih vrijednosti značajno varira po atributima
- Ručno popunjavanje nedostajućih vrijednosti: zamorno + neizvodljivo?
- Automatsko popunjavanje s:
 - Općom konstantom: npr. “nepoznato”, nova klasa?!
 - Srednja vrijednost atributa
 - Srednja vrijednost atributa za sve koji pripadaju istoj klasi: mudriji pristup
 - Najvjerojatnija vrijednost: zaključak temeljen na Bayesian formuli ili stablu odlučivanja

ŠUMOMI U PODACIMA (NOISY DATA)

- Šum: random greška ili promjena (odstupanje) u mjerenju
- Netočne vrijednosti atributa zbog:
 - Manjkavih instrumenata prikupljanja podataka
 - Problema na unosu podataka
 - Problema pri slanju podataka
 - Tehnoloških ograničenja
 - Nekonzistentnosti u imenovanju
- Drugi problemi koji traže čišćenje podataka su:
 - Dupli slogovi
 - Nepotpuni podaci
 - Nekonzistentni podaci



ŠTO RADITI SA ŠUMOVIMA U PODACIMA?

- Binning (grupiranje)
 - Prvo se podaci sortiraju i podijele u (jednake frekvencije) grupe
 - Onda se radi "gladenje" (izravnavanje) grupa pomoću srednje vrijednosti, medijane, graničnih vrijednosti i sl.
- Regresija
 - Izravnavanje pomoću regresijskih funkcija
- Klasteriranje
 - Otkrivanje i otklanjanje nepodobnih članova grupe
- Kombiniranje računalne i ljudske kontrole
 - Otkrivanje sumnjivih vrijednosti i provjera od strane ljudi (npr. Posebno kod nepodobnih članova grupe)

ČIŠĆENJE PODATAKA KAO PROCES

- Otkrivanje proturječnosti u podacima
 - Uporaba metapodataka (npr. domene, rang, ovisnost, distribucija)
 - Provjera prekoračenja polja
 - Provjera pravila jedinstvenosti, dosljednosti i nul vrijednosti
 - Uporaba komercijalnih alata
 - Čišćenje podataka (data scrubbing): uporaba jednostavnog znanja o domenama (npr. Poštanski broj, spell-check) za otkrivanje grešaka i korekcije
 - Revizija podataka (data auditing): kroz analizu podataka kako bi se otkrila pravila i veze za pronalaženje onih koji krše pravila (npr. Korelacija i klasteriranje za pronalaženje nepodobnih)
- Migriranje i integriranje podataka
 - Alati za migriranje podataka: omogućavaju specificiranje transformacija
 - ETL alati: omogućavaju korisnicima specificiranje transformacija kroz GUI
- Integriranje dva procesa
 - Iterativnog i interaktivnog



INTEGRIRANJE PODATAKA

- Integriranje podataka:
 - Kombiniranje podataka iz višestrukih izvora u koherentnu bazu
- Shema integracija: npr. $A.cust-id \equiv B.cust-\#$
 - Integriranje metapodataka iz različitih izvora
- Problem identificiranja entiteta:
 - Identificiranje stvarnih entiteta iz višestrukih izvora podataka, npr. Bill Clinton = William Clinton
- Otkrivanje i rješavanje konfliktnih vrijednosti podataka
 - Za isti stvarni entitet, vrijednosti atributa su različite za različite izvore
 - Mogući razlozi: različito predstavljanje, različite skale, npr. Metrička vs Britanska



REDUNDANCIJA I INTEGRIRANJE PODATAKA

- Redundantni podaci se javljaju najčešće kada se radi integriranje iz više baza podataka
 - *Identifikacija objekta*: Isti atribut ili objekt mogu imati različita imena u različitim bazama
 - *Izvedeni podaci*: Jedan atribut može biti "izvedeni" atribut u drugoj tablici, npr. Godišnji prihod
- Redundantni atributi se mogu otkriti korelacijskom i analizom kovarijance
- Pažljivo integriranje podataka iz više izvora može pomoći u smanjenju/izbjegavanju redundancije i nekonzistentnosti, te poboljšati brzinu ručarenja i kvalitetu podataka



TRANSFORMIRANJE PODATAKA

- Izravnjavanje (smoothing): uklanjanje šuma u podacima
- Agregiranje: zbrajanje, izrada podatkovne kocke
- Generalizacija: koncept “penjanja” po hijerarhiji
- Normalizacija: promjena veličine kako bi se uklopilo unutar manjeg, specifičnog ranga
 - min-max normalizacija
 - z-score normalizacija
 - normalization pomoću decimalnog skaliranja
- Konstruiranje atributa/osobina
 - Novi atributi nastaju iz postojećih

TRANSFORMIRANJE PODATAKA: NORMALIZACIJA

- Min-max normalizacija: za $[new_min_A, new_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$$

- Npr. Neka se prihod s ranga \$12,000 do \$98,000 normalizira na [0.0, 1.0]. Onda se \$73,000 mapira na $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- Z-score normalizacija (μ : srednja vrijednost, σ : standardna devijacija):

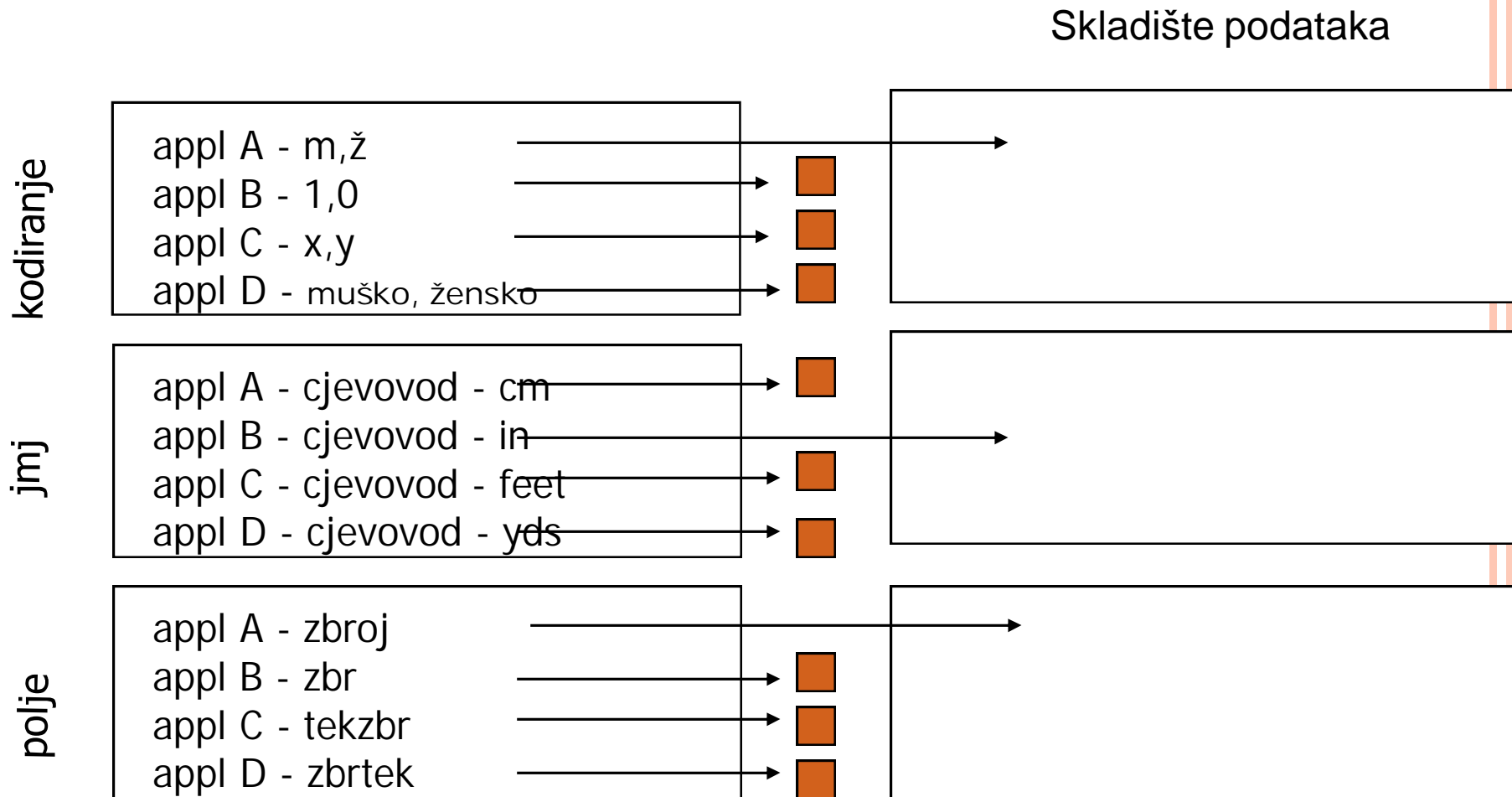
$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Npr. Neka je $\mu = 54,000$, $\sigma = 16,000$. Slijedi $\frac{73,600 - 54,000}{16,000} = 1.225$

- Normalizacija pomoću decimalnog skaliranja

$$v' = \frac{v}{10^j} \quad \text{Gdje je } j \text{ najmanji cio broj takav da je } \text{Max}(|v'|) < 1$$

TRANSFORMIRANJE PODATAKA - PRIMJER



STRATEGIJE REDUKCIJE PODATAKA

- Redukcija podataka: dobiti reducirano predstavljanje skupa podataka koji je znatno manji po veličini ali daje iste (ili skoro iste) analitičke rezultate
- Zašto redukcija podataka? — Baze/skladišta podataka mogu pohranjivati terabyte podataka. Složena analiza podataka može trajati vrlo dugo ako se izvršava na potpunom skupu podataka.
- Strategije redukcije podataka
 - **Dimenzijska redukcija**, npr. Uklanjanje nevažnih atributa
 - Elementarne transformacije
 - Principal Components Analysis (PCA)
 - Odabir podskupa mogućnosti, kreiranje mogućnosti
 - **Smanjenje veličine** (neki to zovu i redukcija podataka)
 - Regresija i Log-Linear Models
 - Histogrami, klasteriranje, uzorci
 - Gregiranje podatkovnih kocki
 - **Kompresija podataka**



PODRUČJE PROČIŠĆAVANJA PODATAKA

ETL (Extraction Transformation Loading)

- Otkrivanje promjena u izvornim podacima potrebnim za skladište podataka;
- Izdvajanje podataka iz izvornih sustava;
- Čišćenje i transformiranje podataka;
- Restrukturiranje ključeva podataka;
- Indeksiranje podataka;
- Sumiranje podataka;
- Održavanje metapodataka;
- Učitavanje podataka u skladište podataka.



POJMOVI TRANSFORMIRANJA PODATAKA

- Izdvajanje
- Prilagodba
- Ribanje
- Miješanje
- Kućanstva
- Obogaćivanje
- Procjenjivanje
- Učitavanje
- Validacija
- Ažuriranje



POJMOVI TRANSFORMIRANJA PODATAKA

○ Izdvajanje

- Izdvaja podatke iz operativnih izvora u “as is” statusu

○ Prilagodba

- Konverzija tipova podataka iz izvornih u krajnje baze (skladište podataka)



DATA TRANSFORMATION TERMS

○ Kućanstva

- Identificiranje svih članova kućanstva (koji žive na istoj adresi)
- Osigurava da se samo jedna poštanska pošiljka šalje kućanstvu
- Može rezultirati značajnim uštedama papira, poštarine ...



DATA TRANSFORMATION TERMS

- Obogaćivanje
 - Korištenje podataka iz vanjskih izvora kako bi se obogatili operativni podaci.
- Procjenjivanje
 - Izračun vjerojatnosti događaja
npr. Vjerojatnost da će kupac kupiti novi proizvod,
promijeniti marku proizvoda



UČITAVANJE (LOAD)

- Nakon izdvajanja, ribanja, čišćenja, validiranja itd. potrebno je učitati podatke u skladište podataka
- Otvorena pitanja
 - Ogromne količine podataka koje treba učitati
 - Kratko vrijeme kada skladište podataka može biti off line (često ne ni noću - web)
 - Kada praviti indekse i zbrojne tablice
 - Dozvoliti administratoru sustava nadzor, prekid, nastavak, promjenu stope učitavnja
 - Skladan oporavak – nastavak nakon ispada sustava tamo gdje se stalo bez gubitka integriteta podataka



TEHNIKE UČITAVANJA

- Korištenje SQL-a za dodavanje ili unos novih podataka
 - Slog u određenom vremenu
 - Dovodi do random disk I/O
- Korištenje batch učitavanja



TAKSONOMIJA UČITAVANJA

- Inkrementalni naspram potpunog učitavanja
- Online naspram Offline učitavanja



AŽURIRANJE (OSVJEŽAVANJE)

- Propagira ažuriranja nad izvornim podacima na skladište podataka
- Otvorena pitanja:
 - Kada osvježiti (refresh)
 - Kako osvježiti – tehnike ažuriranja



KADA OSVJEŽITI?

- periodično (npr. svaku večer, svaki tjedan) ili nakon značajnih događaja
- Za svako ažuriranje: nije zajamčeno sve dok DW ne zatraži ažuran podatak
- Politika osvježavanja postavljena od strane administratora a bazirana na korisničkim potrebama i prometu
- Moguće različite politike za različite izvore podataka



TEHNIKE OSVJEŽAVANJA

- Potpuno izdvajanje iz osnovnih tablica
 - Čita čitavu izvornu tablicu: preskupo
 - Možda jedini izbor za nasljeđene sustave



KAKO OTKRITI PROMJENE

- Kreirati snapshot log tablicu za bilježenje id-ijeva ažuriranih redaka izvornih podataka i timestamp-ova
- Otkrivanje promjena pomoću:
 - Definiranja “after row” okidača (triggers) za ažuriranje snapshot loga kada se promijeni izvorna tablica
 - Korištenje regularnih transakcijskih logova za otkrivanje promjena u izvornim podacima



IZDVAJANJE PODATAKA I ČIŠĆENJE

- Izdvajanje podataka iz postojećih operativnih i nasljeđenih podataka
- Otvorena pitanja:
 - Izvori podataka za DW
 - Kvaliteta izvornih podataka
 - Miješanje različitih izvora podataka
 - Transformiranje podataka
 - Kako propagirati ažuriranja (na izvornim podacima) u skladište podataka
 - Terabytes podataka za učitavanje



ZA SLJEDEĆE PREDAVANJE

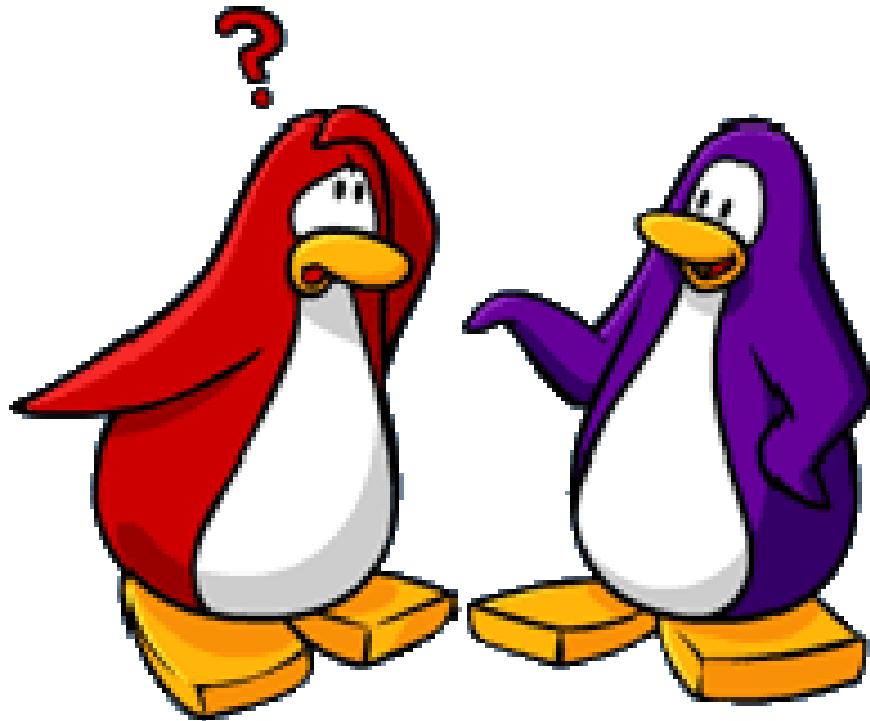
○ Datum: 10.11.2015.

1. Prezentirati:

- Opis problema
- ER model podataka
- Use case dijagram
- Dijagram aktivnosti

Prezentaciju i MS Word dokument s detaljnim opisom i svim dijagramima poslati prof. Gašpar na mail NAJDALJE DO 09.11.2015., 23:59





Questions..