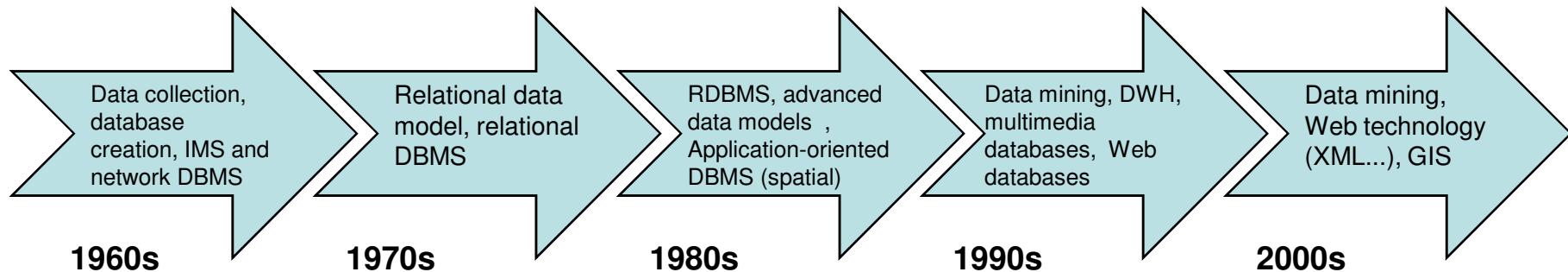


U rudnik s SQL Developerom primjena Oracle data mining alata



K.Bokulić, M.Šipek

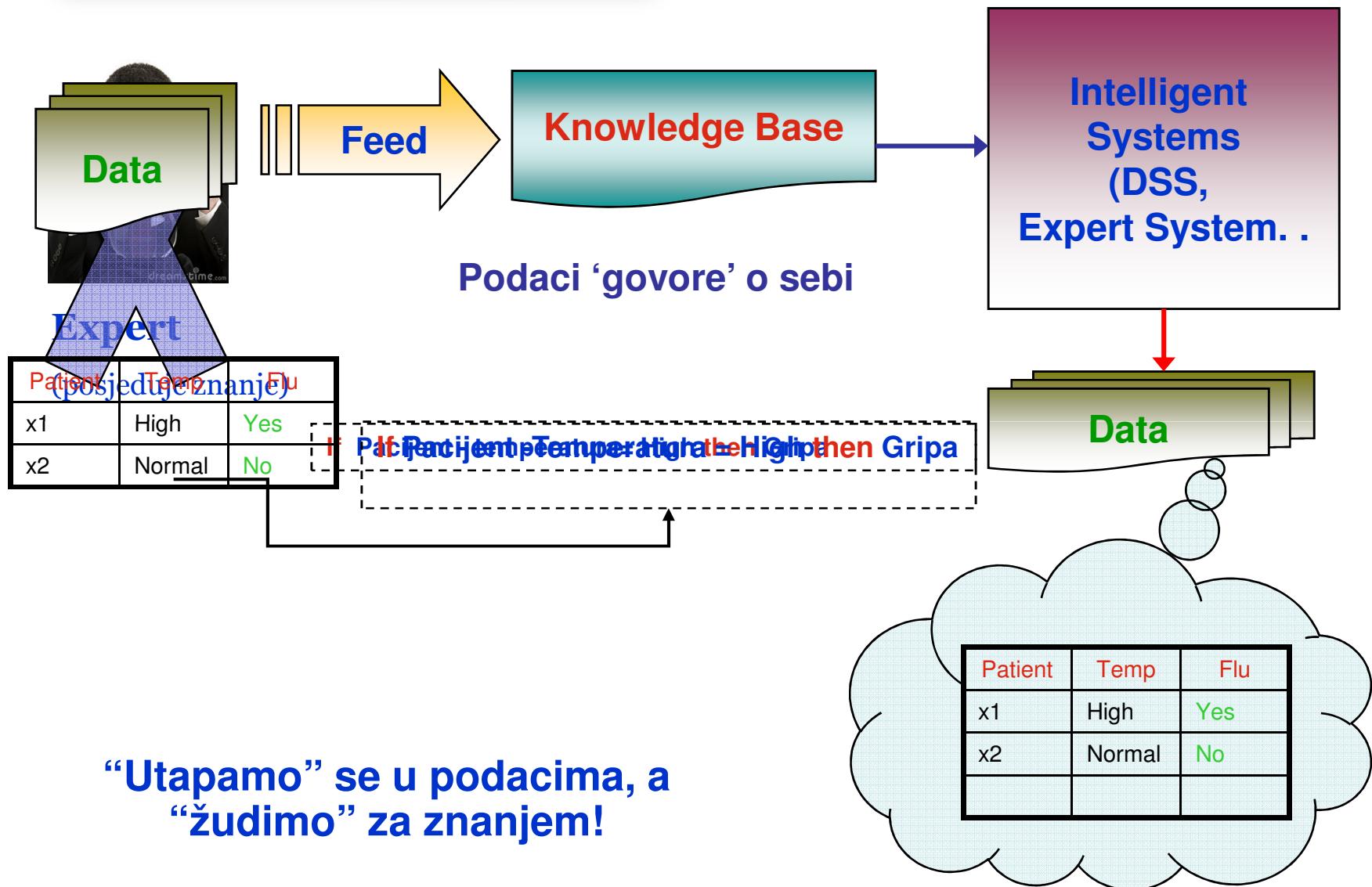
Evolution of Database Technology



- Ekstrakcija detaljnih informacija
- Daje informaciju -
Tko je kupio proizvod A u zadnja 3 mjeseca

- Otkrivanje znanja i uzoraka
- Daje uvid i predikciju -
Tko će kupiti proizvod A u iduća 3 mjeseca **i zašto?**

The New Story: Data Mining and KDD



Svjetla budućnost:

“A recent Gartner Group Advanced Technology Research Note listed data mining and artificial intelligence at the top of the five key technology areas that "will clearly have a major impact across a wide range of industries within the next 3 to 5 years.”

Gartner also listed parallel architectures and data mining as two of the top 10 new technologies in which companies will invest during the next 5 years.”

Gartner

Što je Data mining

- **Rudarenje podataka** (*Data mining*) je postupak *sortiranja, organiziranja ili grupiranja* velikog broja podataka i izvlačenje relevantnih informacija.
- Možemo ga definirati kao proces pronalaženja korisnog znanja kojim se ***otkrivaju odnosi, logičnost, pravilnost*** te općenito bilo kakve ***strukture u podacima***.
- Alternativni nazivi:
 - Knowledge discovery(mining) in databases (KDD), knowledge extraction, data / pattern analysis, data archeology, data dredging, information harvesting, business intelligence, dr.

Primjena DM-a



- **Direktni marketing** – najveća vjerojatnost odaziva
- **Izradu profila kupaca** – prilagođena ponuda
- **Segmentacija** – segmente kojima se mogu posebno prilagoditi usluge
- **Istraživanje povezanosti prodaje različitih proizvoda** – analiza kupovne košarice
- **Stimulacija kupovine drugih artikala istog poduzeća** - crossell, upsell
- **Zadržavanje klijenata, životna vrijednost i aktivacija klijenta**
- **Racionalizacija poslovanja**
- **Razni modeli rizika**- odlaska klijenata, pranja novca, kreditni rizik

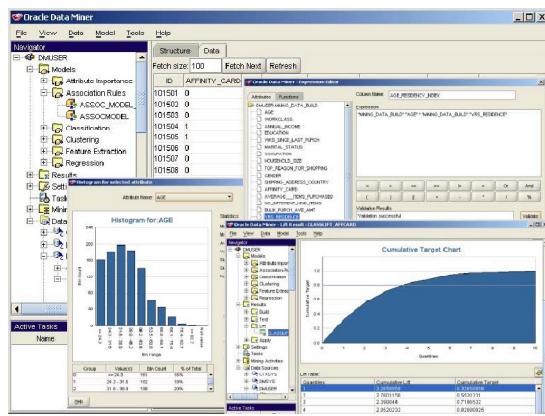
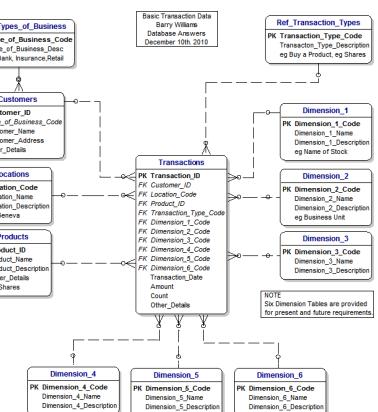
Tko koristi Data Mining

- Financijske institucije
- Telekomi
- Osiguravajuća društva
- Maloprodaja
- Servisi
- Zdravstvo
- Autoindustrija
- Proizvodnja
- Kemijska industrija

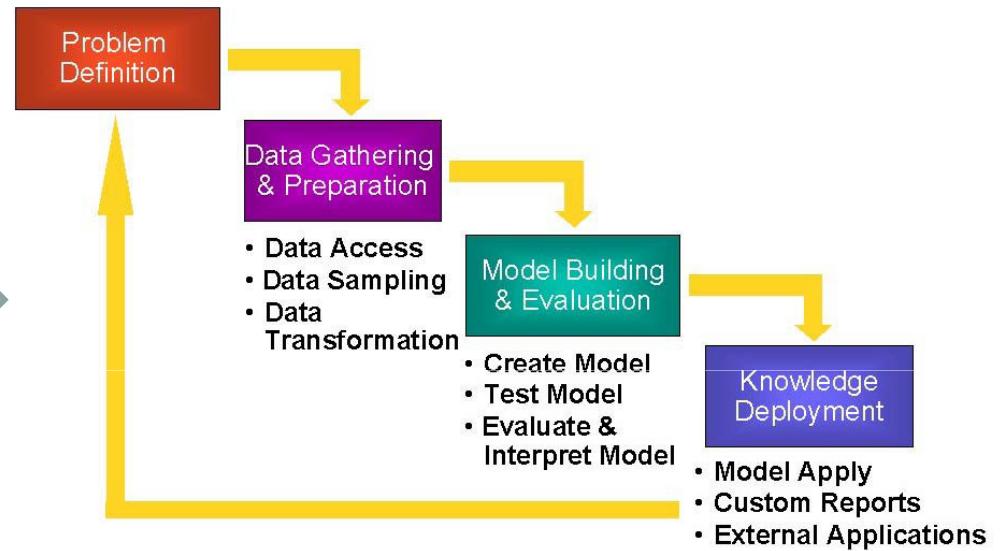


Proces data mininga

DWH model



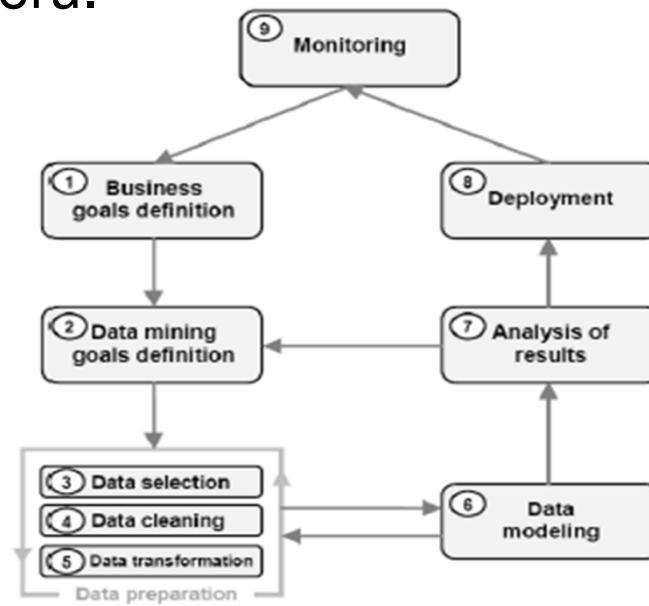
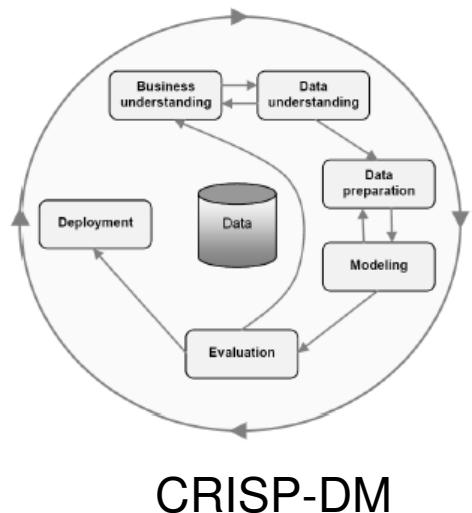
The Data Mining Process



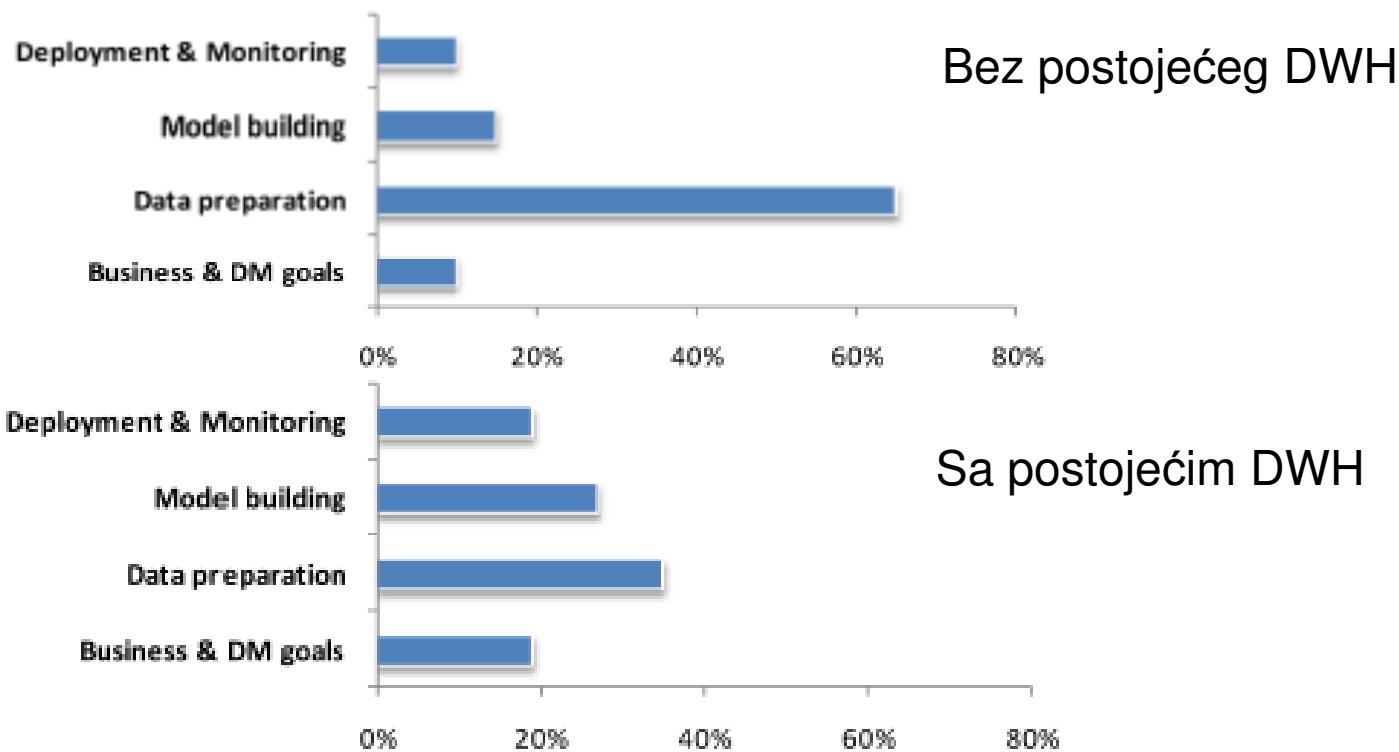
- Odabir metode ovisi o cilju analize.
- Proces rudarenja podataka je **iterativan!**
- Svaka iteracija → precizniji model

Metodologija i framework

- Često govorimo o upotrebi rudarenja podataka radi vidljivih “benefita”, međutim, metodologija koja se koristi i koraci unutar frameworka nisu uvek jasni.
- Za izradu modela korištena je **CRISP-DM** metodologija, preporučena prema Gartneru.



Utrošeno vrijeme po koracima metodologije



Tipovi DM-a



- Dva osnovna tipa rudarenja podataka:
 - **verifikacija hipoteze** – cilj je provjeriti da li je neka ideja ili dojam o važnosti odnosa među određenim podacima utemeljen ili ne
 - **otkrivanje novih znanja** – među nekim pojavama mogu postojati neki još nepoznati, a statistički važni odnosi koje čovjek ni iskustvom niti svojim intelektualnim sposobnostima ne može spoznati

Kako DM model uči?

Tipovi dubinske analize

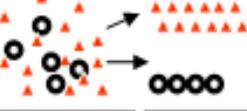
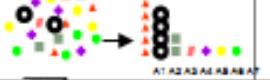
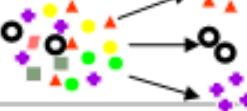
- Za otkrivanje znanja primjereni su postupci učenja iz primjera.
- Nadzirano učenje (eng. ***Supervised/directed data mining***) podrazumijeva pristup „odozgo prema dolje“ i koristi se kada analitičar zna što traži ili što želi predvidjeti. Stoga se ovaj stil često naziva **prediktivnim modeliranjem**.
- S druge je strane nenadzirano učenje (eng. ***Unsupervised / Undirected data mining***) koje se često naziva i **deskriptivnim modeliranjem** te podrazumijeva pristup „odozdo prema gore“. Problem učenja definiran je kao traženje grupa sličnih primjera (engl. *clustering*).

Informacije koje utvrđujemo

Rudarenjem je moguće utvrditi sljedeće vrste informacija:

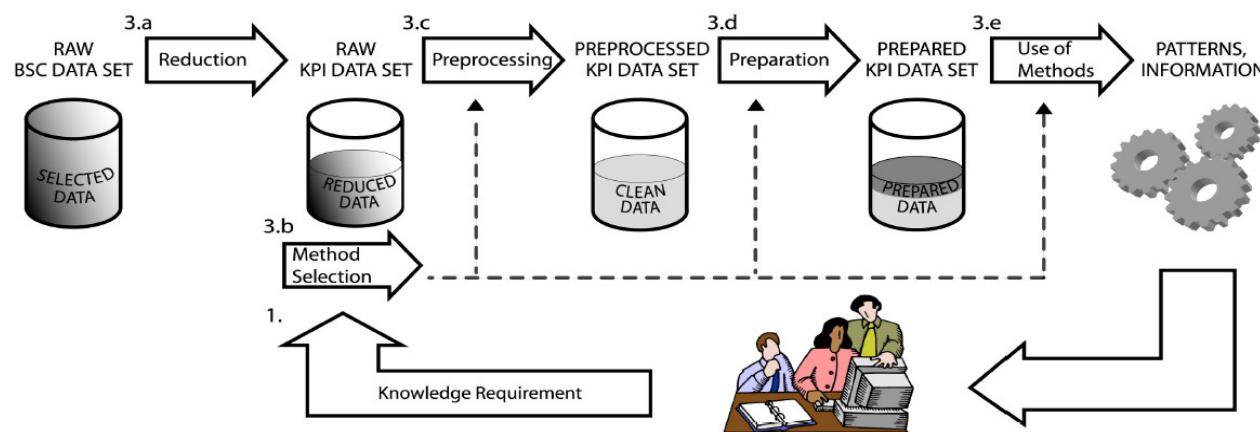
- **Klase**, (unaprijed definirane klase)
- **Klastere** (bez unaprijed zadanih klasa)
- **Asocijacije** (uvjetovanost događaja)
- **Sekvence** (događaji koji u određenoj vjerojatnosti slijede jedan za drugim)
- **Prognoze** (prognozira se budućnost iz postojećih podataka)

Modeli i algoritmi Oracle DM-a

Problem	Algorithm	Applicability
Classification 	Logistic Regression (GLM) Decision Trees Naïve Bayes Support Vector Machine	Classical statistical technique Popular / Rules / transparency Embedded app Wide / narrow data / text
Regression 	Multiple Regression (GLM) Support Vector Machine	Classical statistical technique Wide / narrow data / text
Anomaly Detection 	One Class SVM	Lack examples of target field
Attribute Importance 	Minimum Description Length (MDL)	Attribute reduction Identify useful data Reduce data noise
Association Rules 	Apriori	Market basket analysis Link analysis
Clustering 	Hierarchical K-Means Hierarchical O-Cluster	Product grouping Text mining Gene and protein analysis
Feature Extraction 	Nonnegative Matrix Factorization	Text analysis Feature reduction

Data Preprocessing - Zašto?

- Podaci su stvarnom svjetu su “prljavi”
 - nekompletni: nedostatak vrijednosti atributa, nedostatak određenih atributa od interesa ili postojanje samo agregiranih podataka
 - “bučni”: sadrže grešku ili krajnost
 - nekonzistentni: sadrže diskrepanciju u nazivu ili šifri
- Bez kvalitetnih podataka nema ni rezultata data mining-a!
 - Kvalitetne odluke temelje se na kvalitetnim podacima

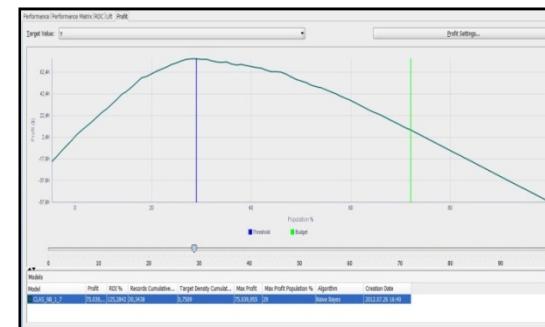


Evaluacija modela

- **Potvrda modela** - Nakon izgradnje modela, provjeravaju se rezultati i tumači signifikantnost. Stopa točnosti pronađena tijekom testiranja odnosi se samo na podatke na kojima je model izgrađen
- **Confusion matrica** - Za klasifikacijske probleme, pokazuje stvaran odnos predviđenih klasa vrijednosti. Stupci pokazuju stvarne klase, redovi pokazuju predviđene klase, a dijagonala pokazuje sva točna predviđanja

Prediction	Actual		
	Class A	Class B	Class C
Class A	45	2	3
Class B	10	38	2
Class C	4	6	40

Confusion matrices

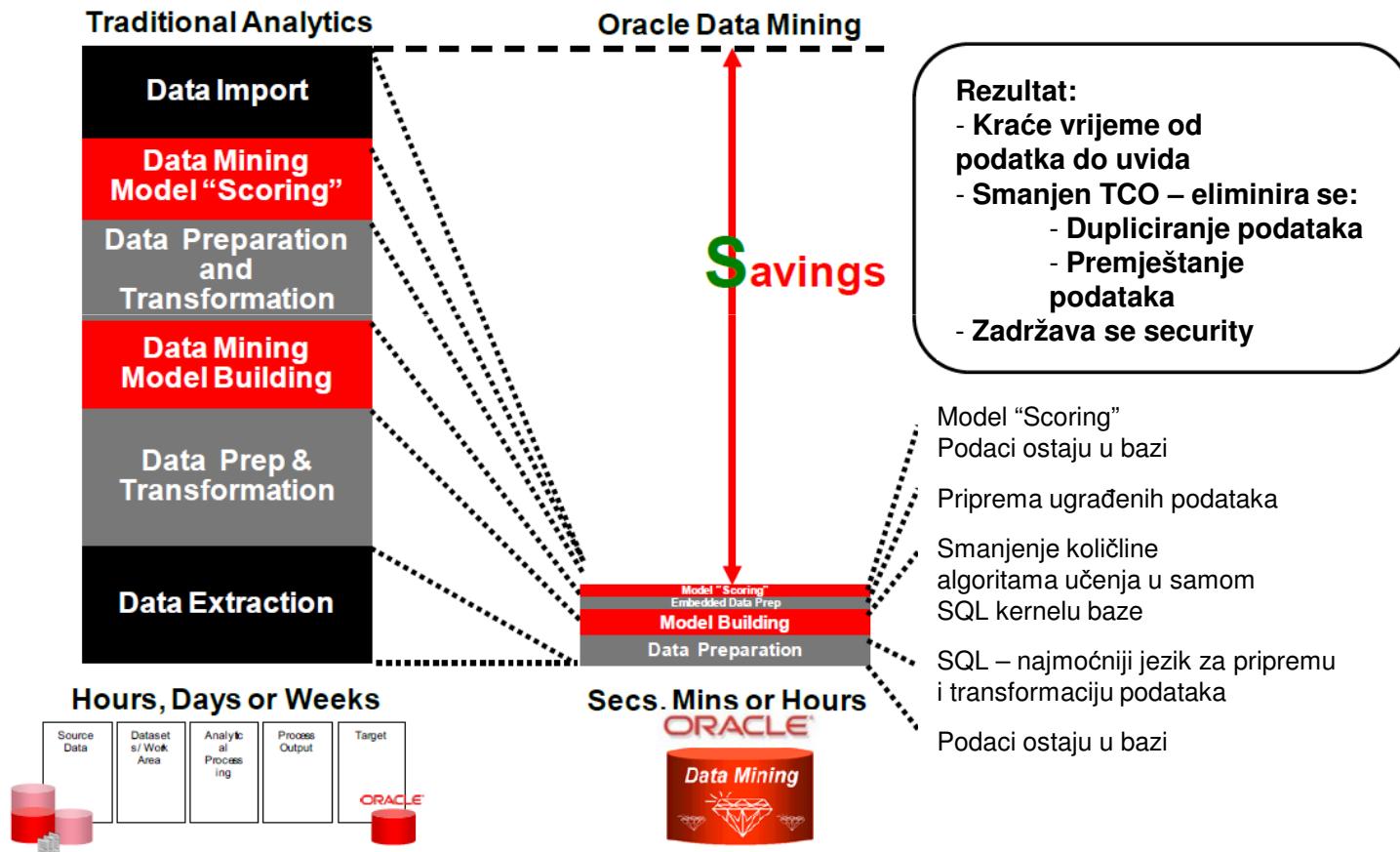


ROI

Zašto Oracle?

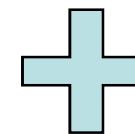
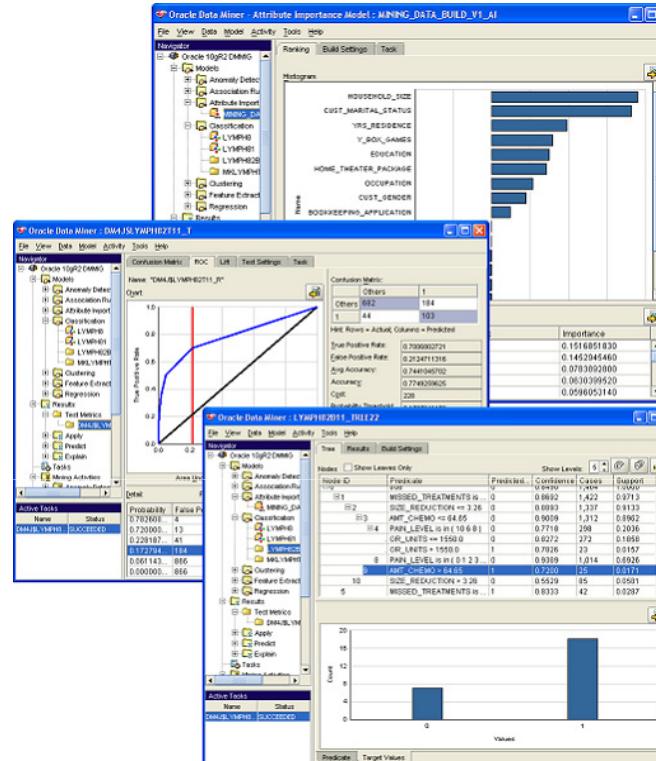
OD17
hroug

In-Database Data Mining



Oracle Data Miner

- 10 godina sa nama, od sada integrirana opcija SQL Developera
 - Obuhvaćeni svi aspekti: dohvati, transformacija, procesiranje, izrada i evaluacija modela
 - Mogućnost utjecaja na sve elemente procesa
 - Interaktivna vizualizacija i izrada izvještaja, uz set predefiniranih izvještaja i analitika
 - Alat za analitičare



Oracle Data Mining

- Automated knowledge discovery, model building and deployment
 - Domain expertise to assemble the “right” data to mine



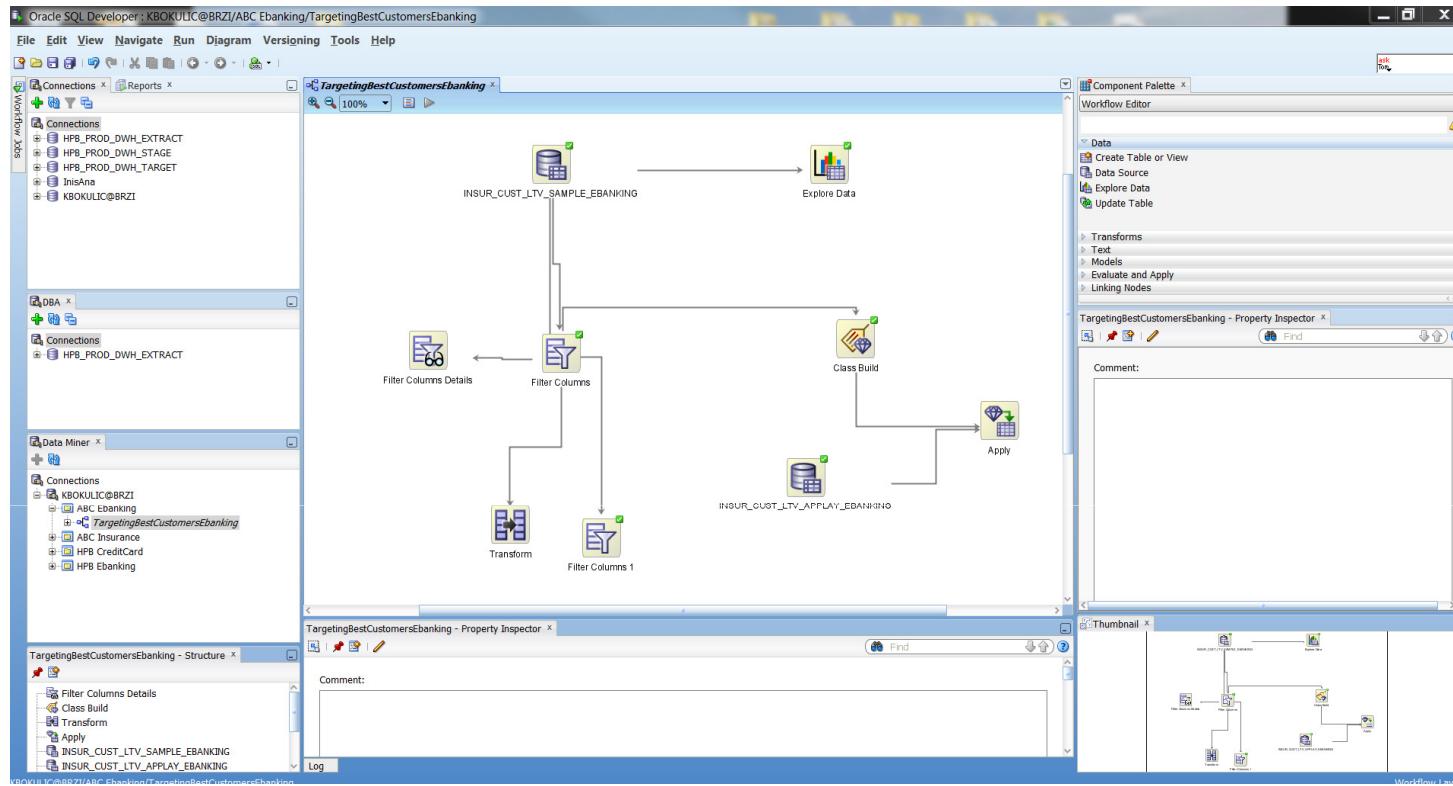
• QRM “Verbs”

- PREDICT
 - DETECT
 - CLUSTER
 - CLASSIFY
 - REGRESS
 - PROFILE
 - IDENTIFY FACTORS
 - ASSOCIATE

multicom
INFORMACIJSKI SUSTAVI

Oracle Data Miner

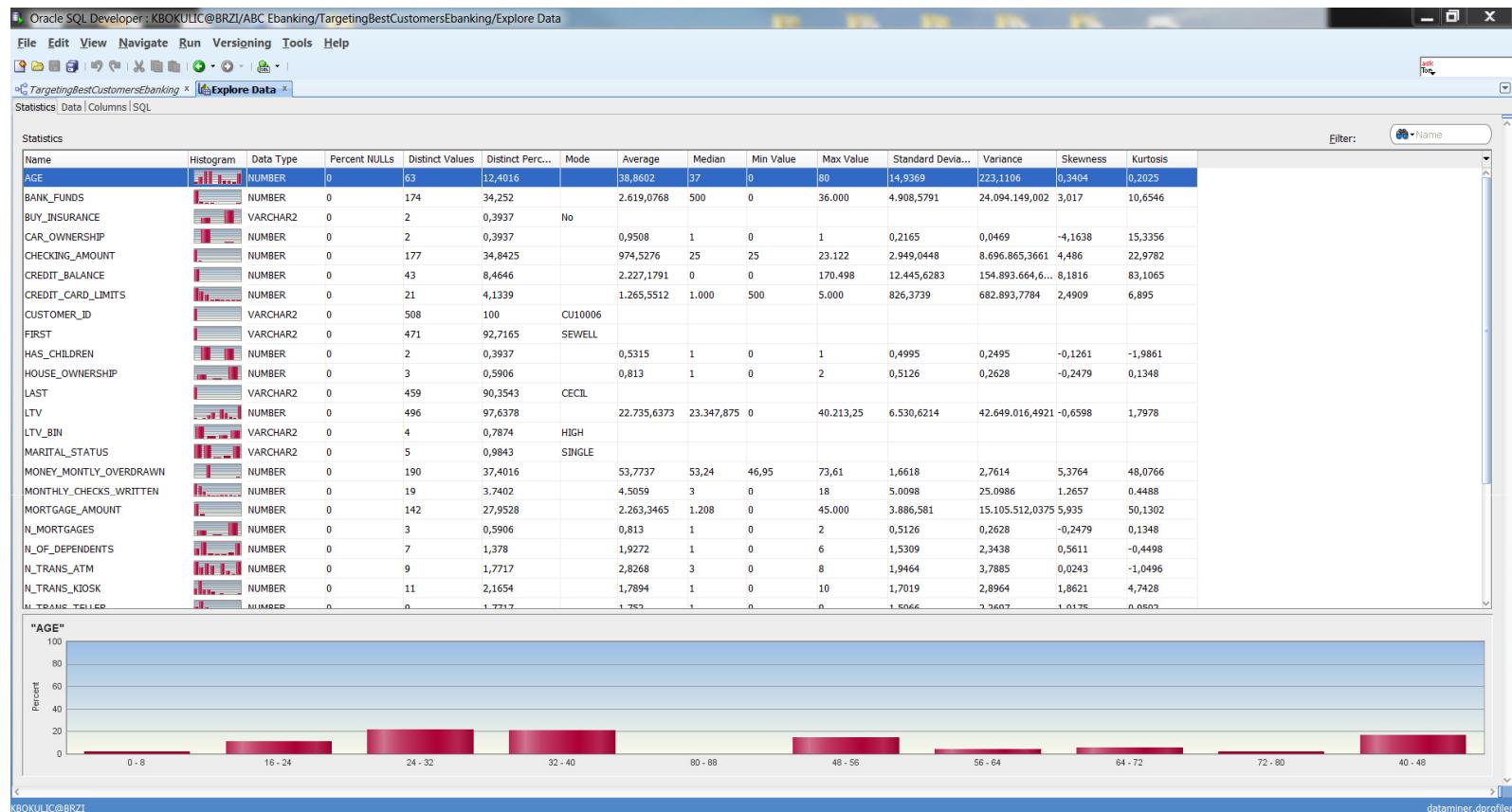
OD17
hroug



- Repozitorij unutar poznatog sučelja SQL Developera
- Izrada modela kroz work-flow
- Jednostavna rekalkulacija i remodulacija modela

Oracle Data Miner

017
hroug



- Statistička analiza podataka i uzorka
- Histogram, prosjek, mode, meidan, standardna deviacija, skewness, kurtosis ...

Oracle Data Miner

Filter Columns - Property Inspector

Name	Type	Output	Hints
AGE	NUMBER	↳	
BANK_FUNDS	NUMBER	↳	
BUY_INSURANCE	VARCHAR2	↳	Attribute passed filters
CAR_OWNERSHIP	NUMBER	↳	Attribute passed filters
CHECKING_AMOUNT	NUMBER	↳	Min importance not reached
CREDIT_BALANCE	NUMBER	↳	Min importance not reached
CREDIT_CARD_LIMITS	NUMBER	↳	Min importance not reached
CUSTOMER_ID	VARCHAR2	↳	Min importance not reached
FIRST	VARCHAR2	↳	Min importance not reached
HAS_CHILDREN	NUMBER	↳	Attribute passed filters
HOUSE_OWNERSHIP	NUMBER	↳	Attribute passed filters
LAST	VARCHAR2	↳	Min importance not reached
LTV	NUMBER	↳	Attribute passed filters
LTV_BIN	VARCHAR2	↳	Attribute passed filters
MARITAL_STATUS	VARCHAR2	↳	
MONEY_MONTLY_OVERD...	NUMBER	↳	Attribute passed filters
MONTHLY_CHECKS_WRI...	NUMBER	↳	Attribute passed filters
MORTGAGE_AMOUNT	NUMBER	↳	Attribute passed filters
NL_MORTGAGES...	NUMBER	↳	Attribute passed filters

Columns

Filters

- Data Quality
 - % Nulls less than or equal: 95
 - % Unique less than or equal: 95
 - % Constant less than or equal: 95
- Attribute Importance
 - Target: N_TRANS_WEB_BANK
 - Importance Cutoff: 0
 - Top N: 100
- Sampling (Stratified)
 - Sample Size: 2.000

TargetingBestCustomersEBanking x Filter Columns x Explore Data x

Attribute Importance Data | Columns | SQL

Target: N_TRANS_WEB_BANK

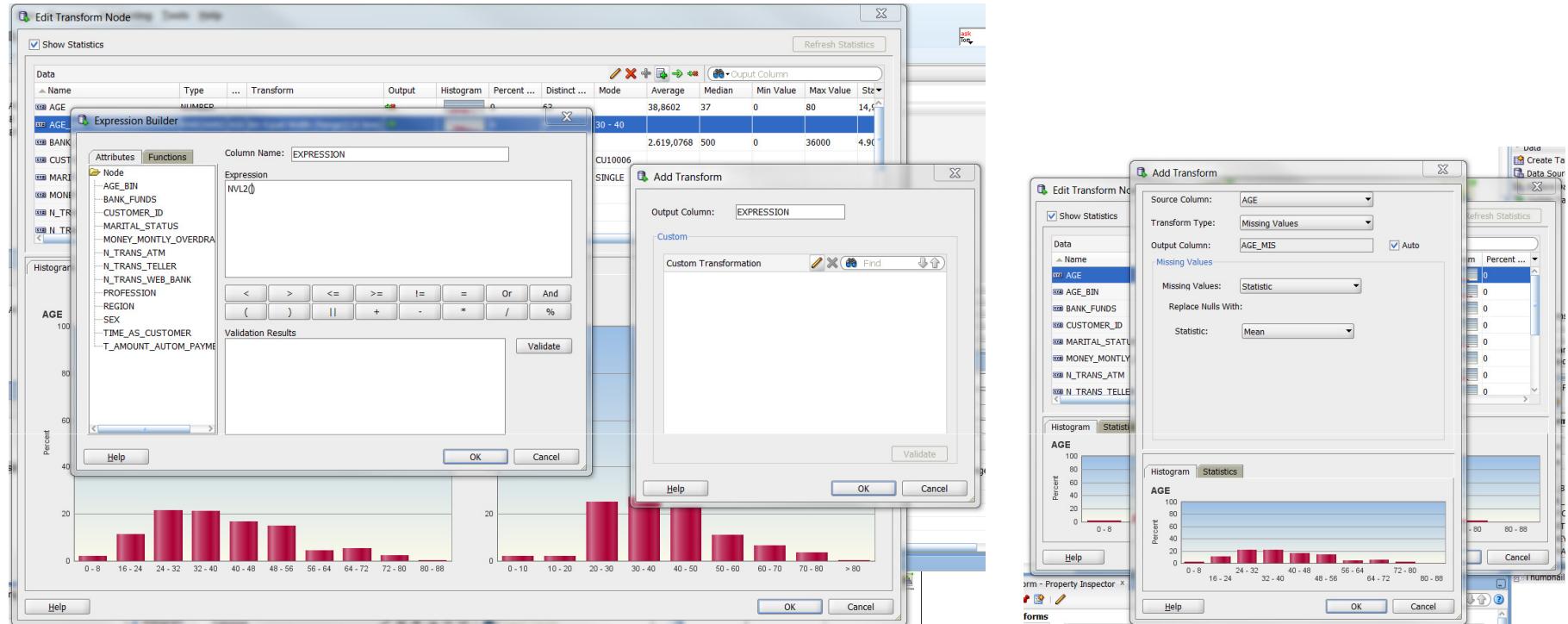
Attribute Ranking

Name	Type	Rank	Importance
HOUSE_OWNERSHIP	NUMBER	1	0,5804
MORTGAGE_AMOUNT	NUMBER	1	0,5804
N_MORTGAGES	NUMBER	1	0,5804
MARITAL_STATUS	VARCHAR2	2	0,4078
LTV	NUMBER	3	0,3229
LTV_BIN	VARCHAR2	4	0,2636
BANK_FUNDS	NUMBER	5	0,148
N_OF_DEPENDENTS	NUMBER	6	0,1215
CAR_OWNERSHIP	NUMBER	7	0,1071
TIME_AS_CUSTOMER	NUMBER	8	0,103
SEX	VARCHAR2	9	0,0854
N_TRANS_TELLER	NUMBER	10	0,0832
N_TRANS_ATM	NUMBER	11	0,0771
MONEY_MONTLY_OVERDRA...	NUMBER	12	0,0741
T_AMOUNT_AUTOM_PAYME...	NUMBER	13	0,0659
HAS_CHILDREN	NUMBER	14	0,0281
MONTHLY_CHECKS_WRITTEN	NUMBER	15	0,0253
AGE	NUMBER	16	0,0211
BUY_INSURANCE	VARCHAR2	17	0,0134
CHECKING_AMOUNT	NUMBER	18	0
CREDIT_BALANCE	NUMBER	18	0
CREDIT_CARD_LIMITS	NUMBER	18	0
CUSTOMER_ID	VARCHAR2	18	0
FIRST	VARCHAR2	18	0
LAST	VARCHAR2	18	0
N_TRANS_KIOSK	NUMBER	18	0
PROFESSION	VARCHAR2	18	0
REGION	VARCHAR2	18	0
SALARY	NUMBER	18	0
STATE	VARCHAR2	18	0

- Analiza kvalitete i prediktivnosti varijabli
- MDL algoritam

Oracle Data Miner

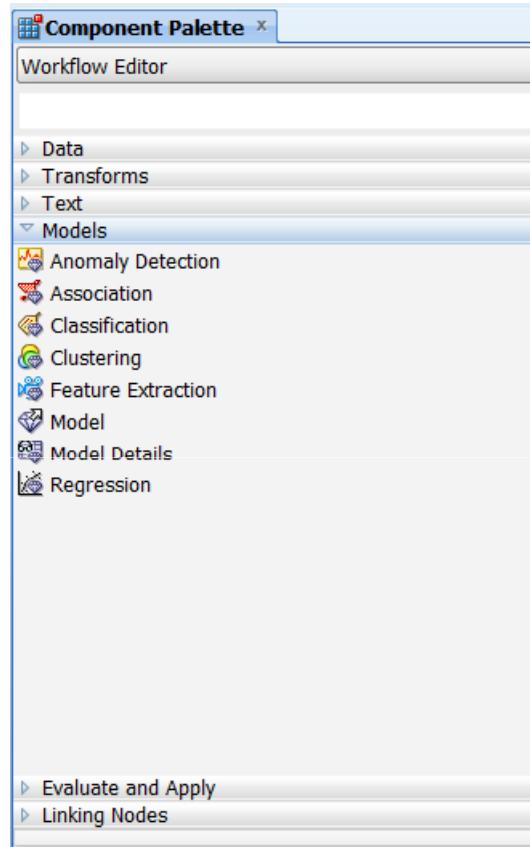
OD17
hroug



- Transformacije nad podacima, pomoću pl/sql izraza
- Popunjavanje nedostajećih podataka ...

Oracle Data Miner

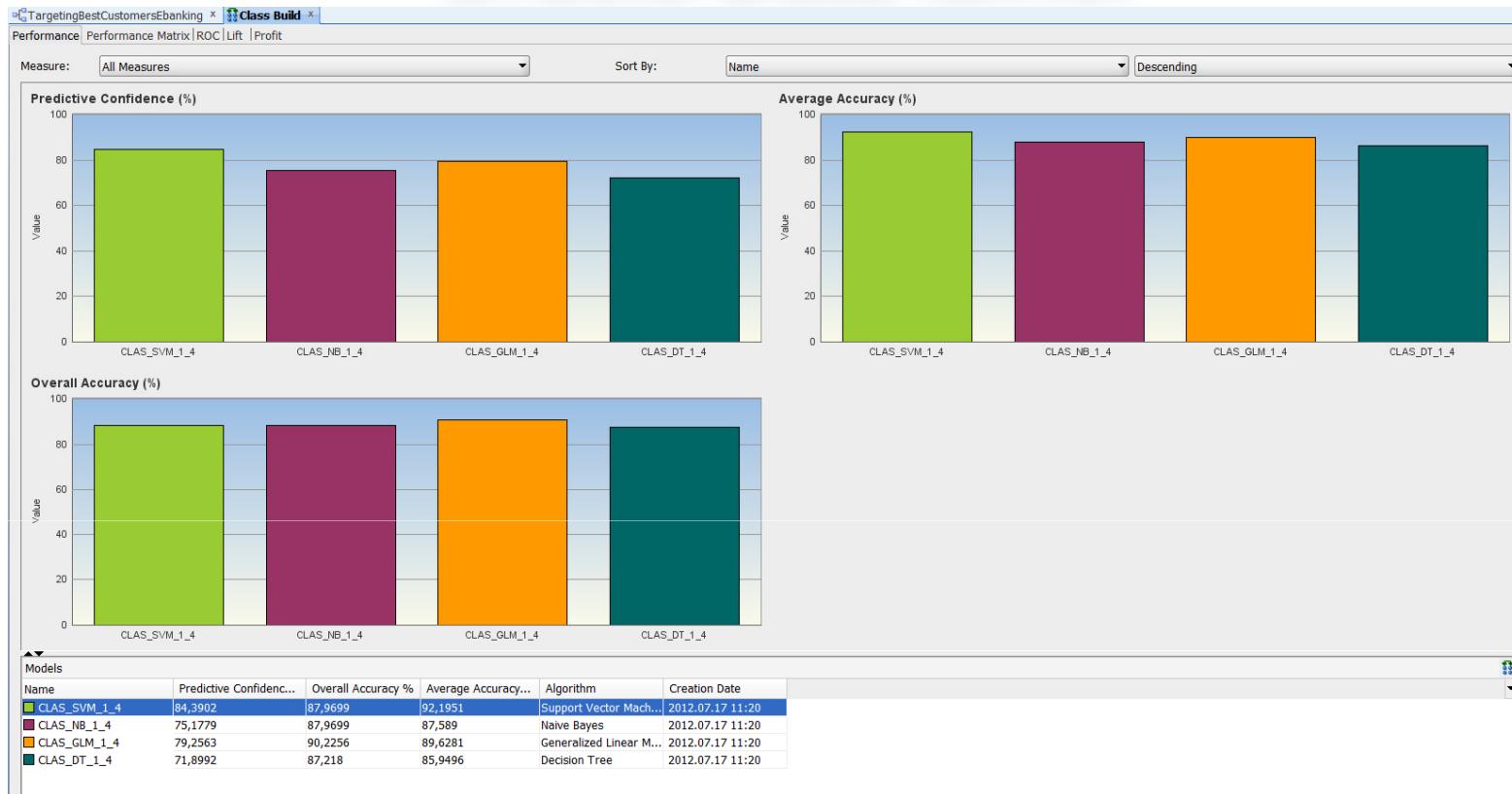
017
hroug



- Učitavanje podataka
- Transformacije i analiza podataka
- Izrada DM modela za ne strukturirane podatke
- Podržani svi industriski standardni modeli za DM
- Jednostavna primjena modela nad podacima iz različitih izvora

Oracle Data Miner

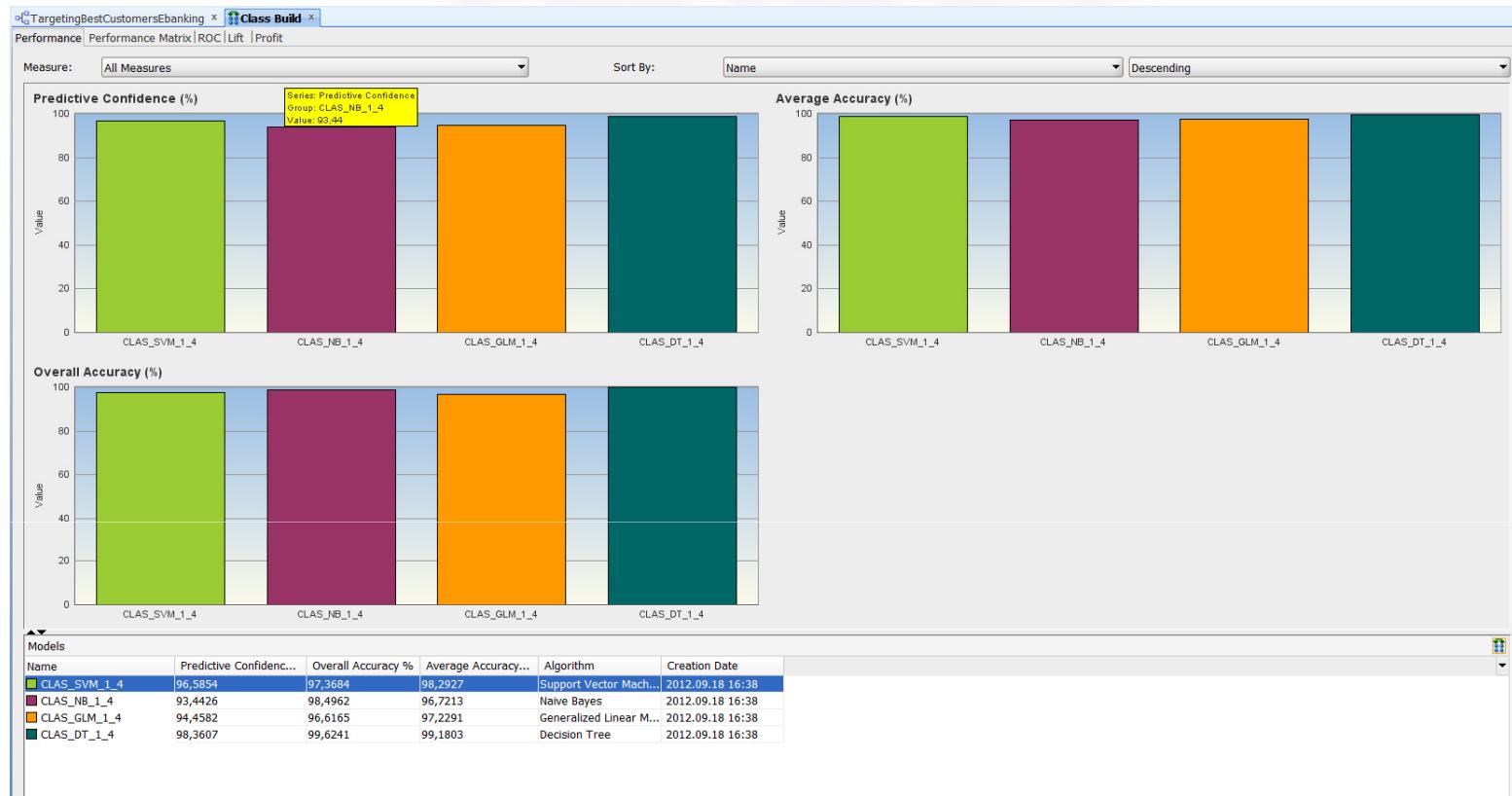
017
hroug



Rezultati modela za predikciju kupovine usluge internet bankarstva

- Istovremena izgradnja višestrukih modela različitim algoritmima
- Grafička usporedba rezultata više modela

Oracle Data Miner

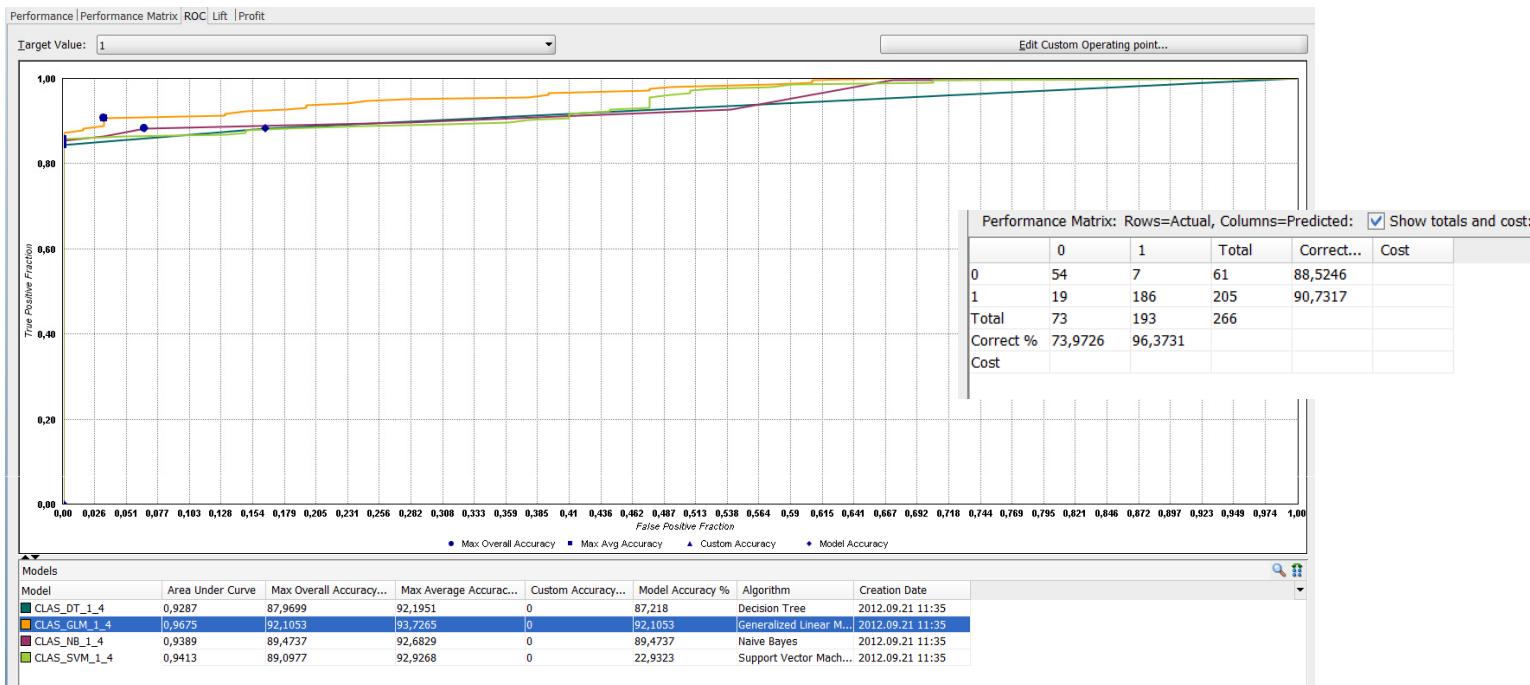


Rezultati modela koji koristi dodatnu varijablu **vlasnik_nekretnine**

- Modeli su osjetljivi na dobar odabir varijabli s kvalitetnim podacima
- Kvalitetna varijabla s dobrom prediktivnošću - bolja prediktivnost modela

Oracle Data Miner

017
hroug

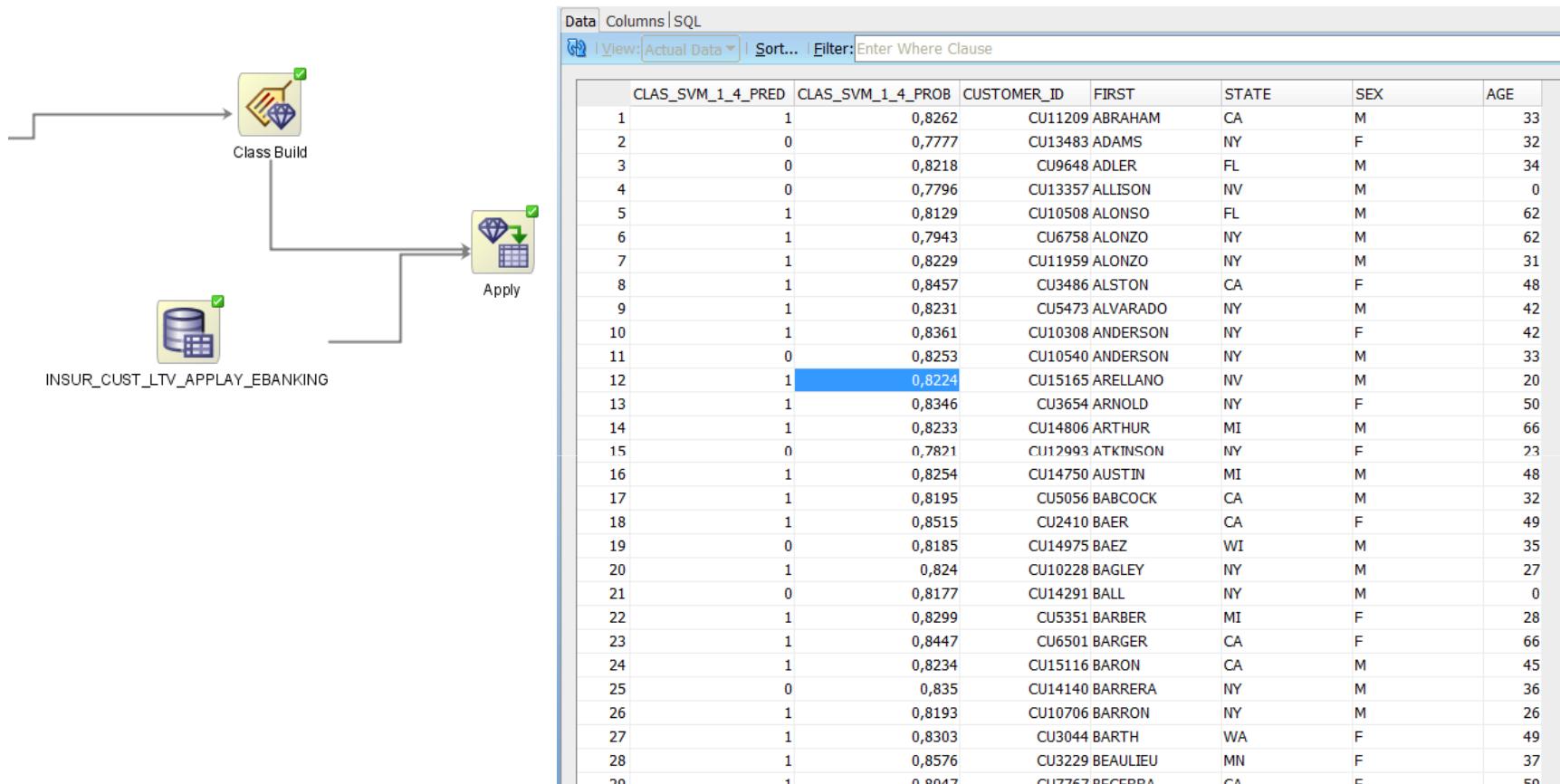


ROC krivulja i performance matrica

- ROC krivulja nam prikazuje omjer dobrih i loših predviđanja u testnom uzorku za ciljanu vrijednost

Oracle Data Miner

017
hroug

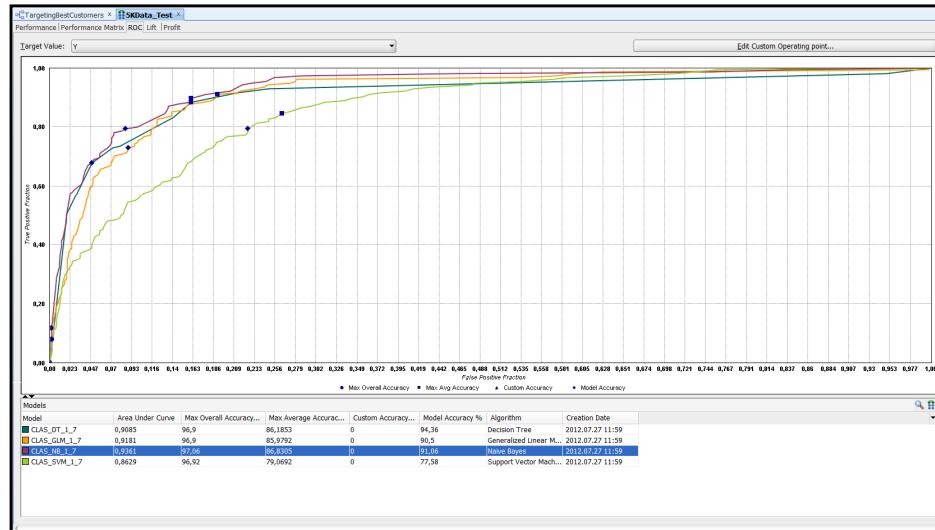


Rezultat modela s vjerojatnostima ishoda

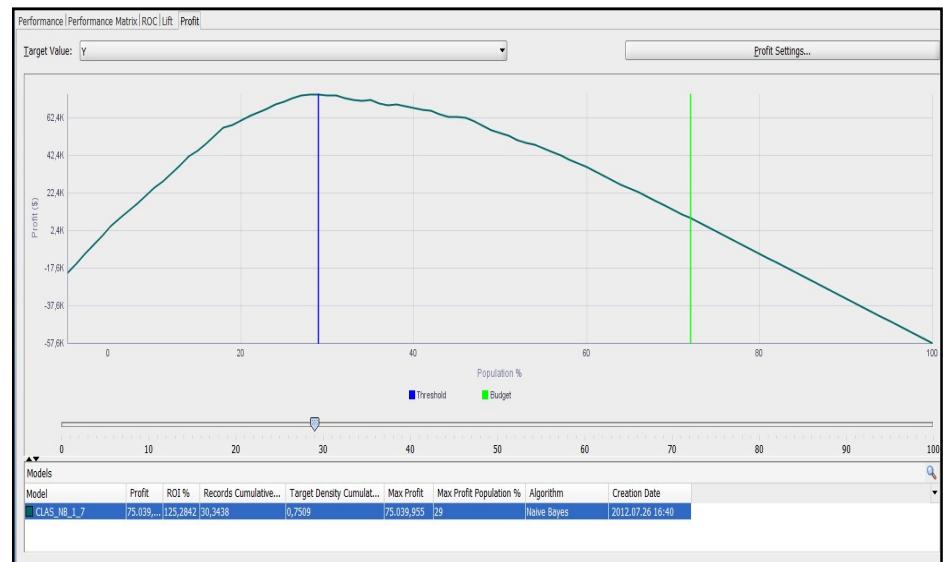
- Podaci nad kojima je primijenjen model smještaju se u tablicu u bazi

Oracle Data Miner

OD17
hroug



ROC krivulje



Projekcija profita kampanje

Primjer: stvarni rezultati modela povećanje prodaje usluge (up-sell)

Oracle Data Miner

OD17
hroug

Izazovi modeliranja:

- Kako definirati problem ?
- Koliko varijabli za model ?
- Koje varijable za model ?
- Koji algoritam ?
- Koji uzorak za trening podatke ?
- Koliko je podataka dovoljno?
- Da li je model 'pretreniran' ili nije dovoljno 'treniran' ?
- Koliko i koji podaci za test modela ?
- Koje razdoblje trening podataka ?
- Koje razdoblje za ishod modela ?



Pitanja?



Mario.Sipek@multicom-is.hr

Kresimir.Bokulic@multicom-is.hr