

Poluautomatizirana  
selekcija varijabli u  
prediktivnoj analizi

# Multicom

- **Glavna područja ekspertize:**

- Data Mining
- Obračun i naplata (**Billing**)
- Upravljanje matičnim podacima (**MDM**)
- Skladišta podataka (**DWH**) i Poslovna Inteligencija (**BI**)
- **B2B**
- Upravljanje korisničkim procesima (**CRM**)



# Prediktivna analiza

predviđanje budućnosti?



Cross-sell / up-sell  
Otkrivanje prijevara  
Churn

- Poznatiji alati:
- Mathematica
  - Matlab
  - Oracle Advanced Analytics
  - Orange
  - R
  - RapidMiner
  - SAP
  - SAS

# Prediktivna analiza pomoću R-a

R je Open Source jezik i okolina za statističke proračune i grafiku

Stvoren 1994 kao alternativa SAS-u i SPSS-u

Preko 2 milijuna R korisnika u svijetu

Tisuće open source paketa na CRAN mreži

CRAN – Comprehensive R Archive Network

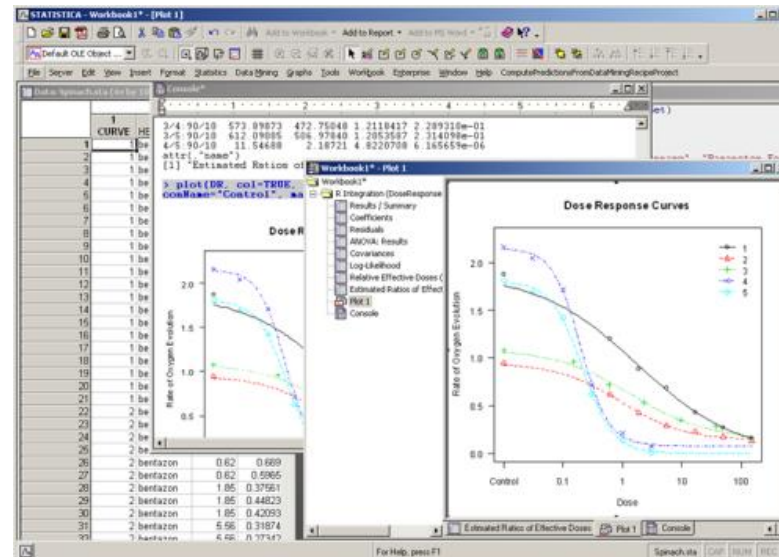


CRAN  
[Mirrors](#)  
[What's new?](#)  
[Task Views](#)  
[Search](#)  
 About R  
[R Homepage](#)  
[The R Journal](#)  
 Software  
[R Sources](#)  
[R Binaries](#)  
[Packages](#)  
[Other](#)  
 Documentation  
[Manuals](#)  
[FAQs](#)  
[Contributed](#)

CRAN Task Views

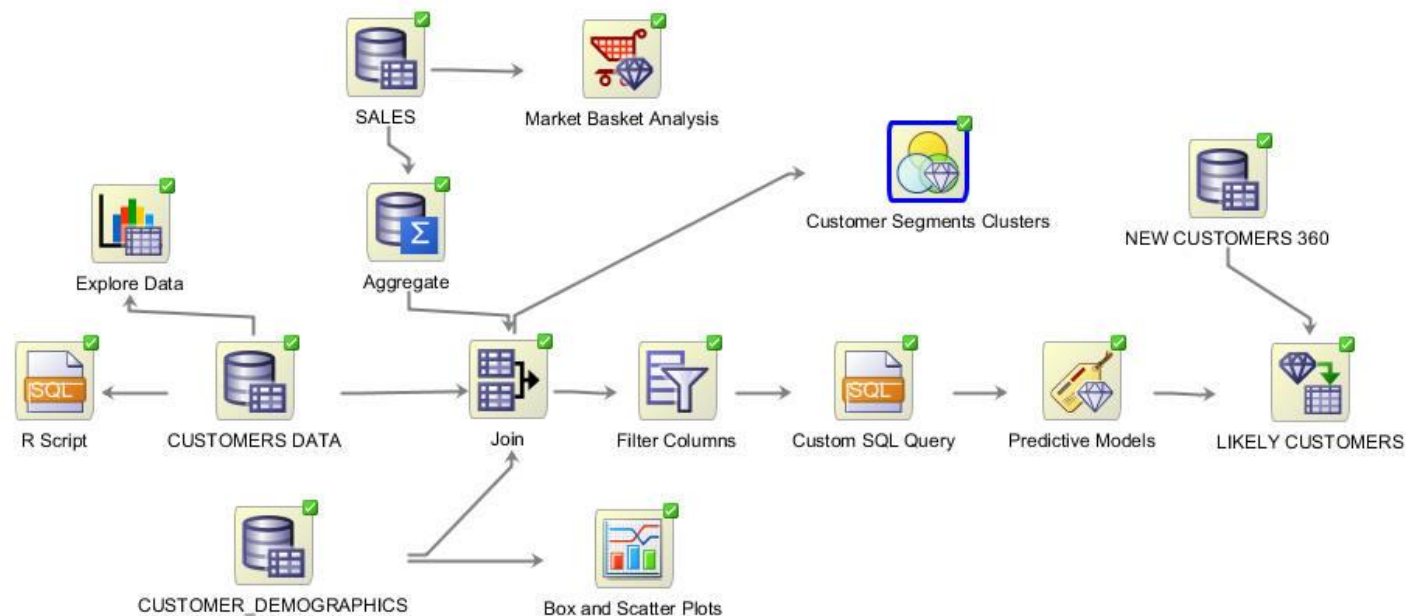
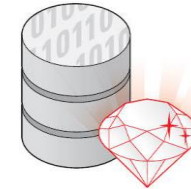
- [Bayesian Inference](#)
- [Chemometrics and Computational Physics](#)
- [Clinical Trial Design, Monitoring, and Analysis](#)
- [Cluster Analysis & Finite Mixture Models](#)
- [Differential Equations](#)
- [Probability Distributions](#)
- [Computational Econometrics](#)
- [Analysis of Ecological and Environmental Data](#)
- [Design of Experiments \(DoE\) & Analysis of Experimental Data](#)
- [Empirical Finance](#)
- [Statistical Genetics](#)
- [Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization](#)
- [High-Performance and Parallel Computing with R](#)
- [Machine Learning & Statistical Learning](#)
- [Medical Image Analysis](#)
- [Meta-Analysis](#)
- [Multivariate Statistics](#)
- [Natural Language Processing](#)
- [Numerical Mathematics](#)
- [Official Statistics & Survey Methodology](#)
- [Optimization and Mathematical Programming](#)
- [Analysis of Pharmacokinetic Data](#)
- [Phylogenetics, Especially Comparative Methods](#)
- [Psychometric Models and Methods](#)

- Statistika
- Grafika
- Data mining
- R Enterprise



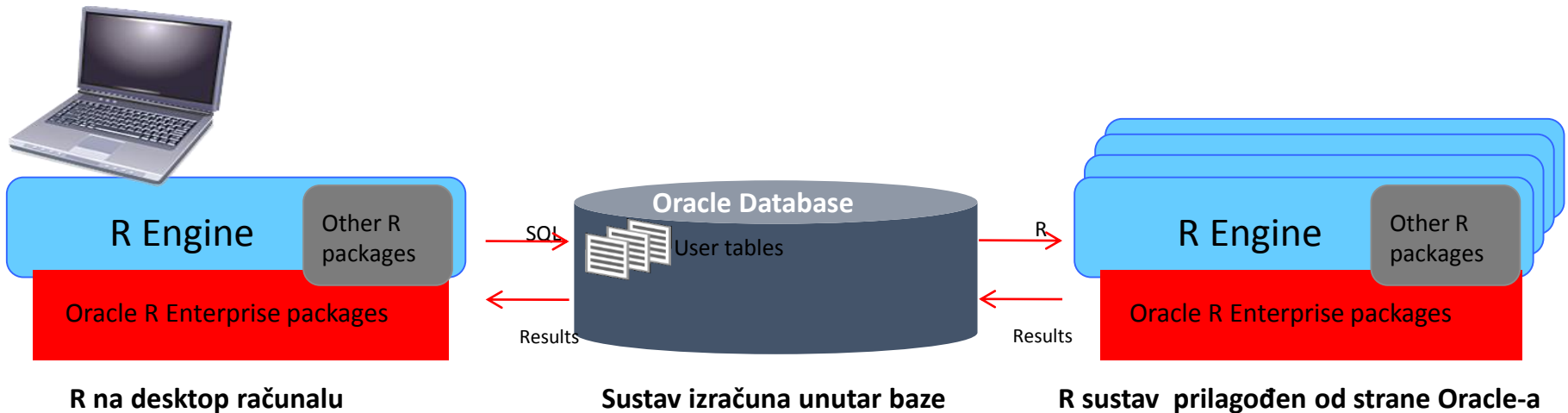
# Prediktivna analiza pomoću ODM-a

- Skalabilna unutar baze prediktivna analiza
- Jednostavno za korištenje
- Brzi i pouzdani rezultati



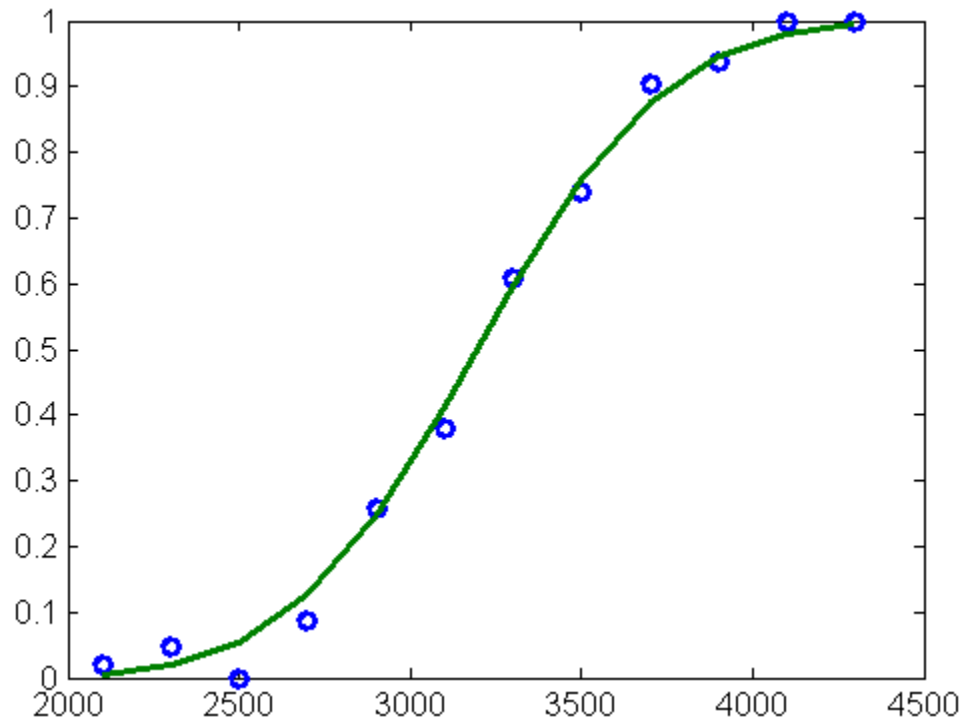
# Oracle R Enterprise (SAVRŠENI SPOJ?)

Integrira open source programski jezik R unutar Oracle baze podataka

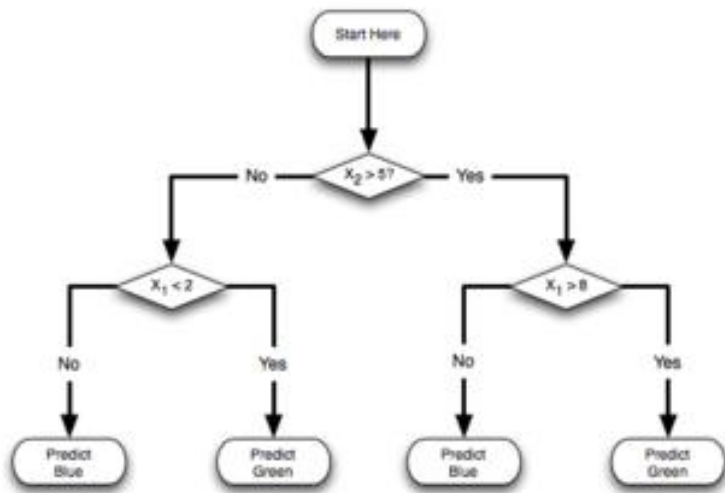


# Koristljivi algoritmi

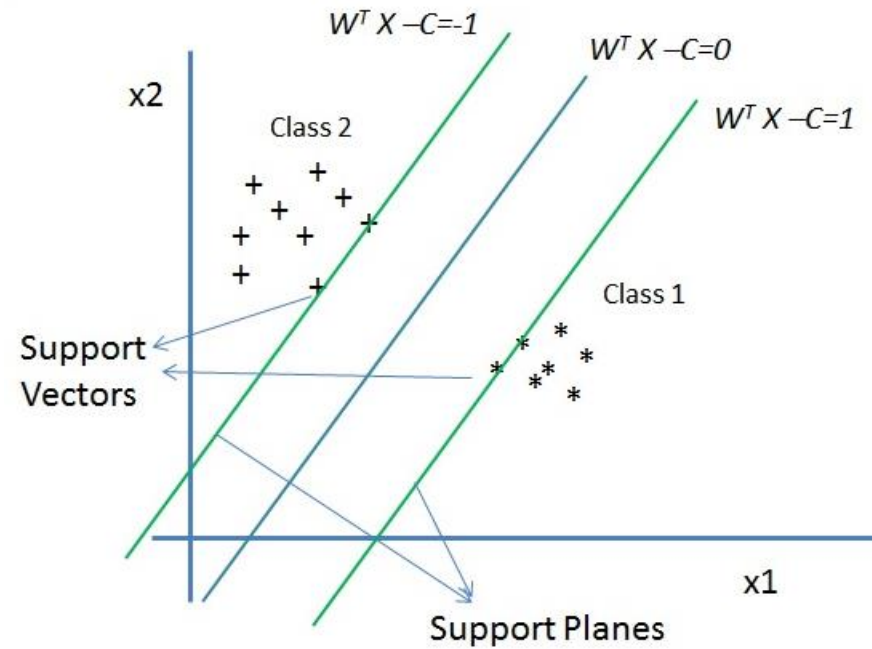
- GLM – general linear model



# Decision tree



# Support vector machine





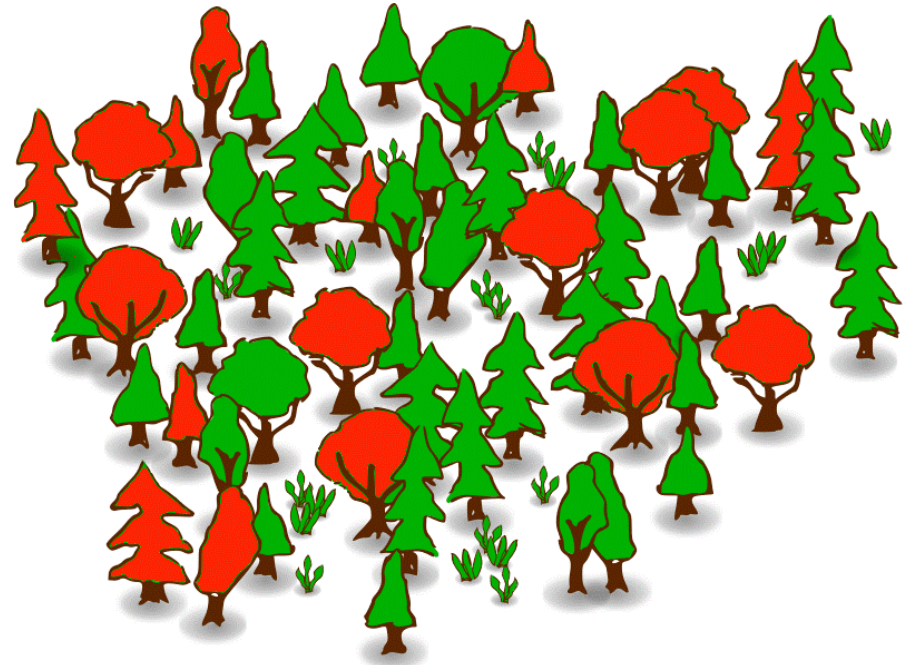
# Naive Bayes

$$\begin{aligned} \log p(C_k|\mathbf{x}) &\propto \log \left( p(C_k) \prod_{i=1}^n p_{ki}^{x_i} \right) \\ &= \log p(C_k) + \sum_{i=1}^n x_i \cdot \log p_{ki} \\ &= b + \mathbf{w}_k^\top \mathbf{x} \end{aligned}$$

$$b = \log p(C_k)$$

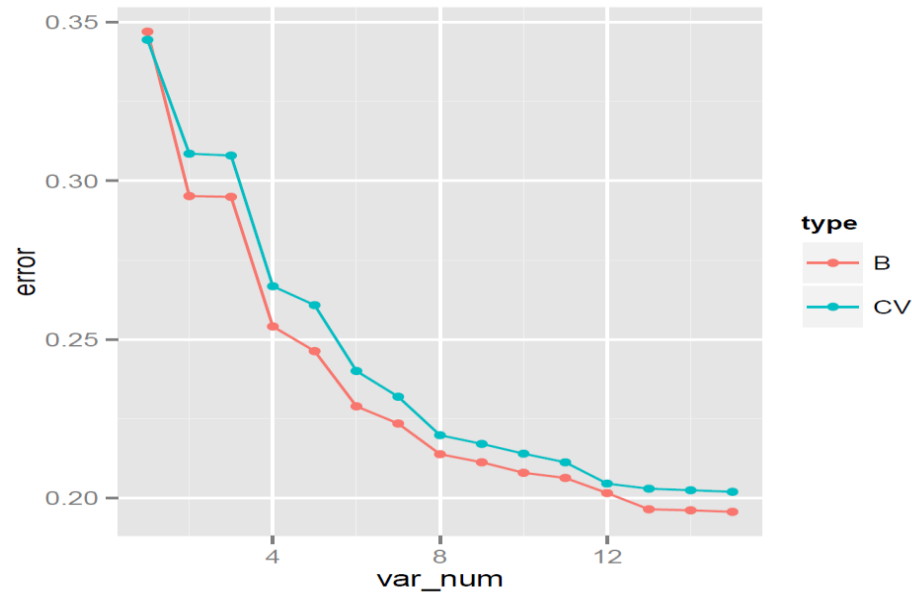
$$w_{ki} = \log p_{ki}$$

# Random forest?



# m.builder rješenje

- Sampliranje (random, stratified)
- Attribute importnace
- Iterativni odabir varijabli
- izbor algoritma
- Oracle R Enterprise – pouzdano i brzo rješenje



# Ulazne vrijednosti

```
R> results <- ore.doEval(FUN.NAME = "m.VarSelect",
  model = "GLM",
  build = "TABLE",
  attribute_importance=TRUE,
  tgt = "TGT",
  columns_to_skip=c("USER_ID","SERIAL_NUMBER"),
  sample_type="ST",
  c_per=0.3,
  ore.connect=TRUE )
```

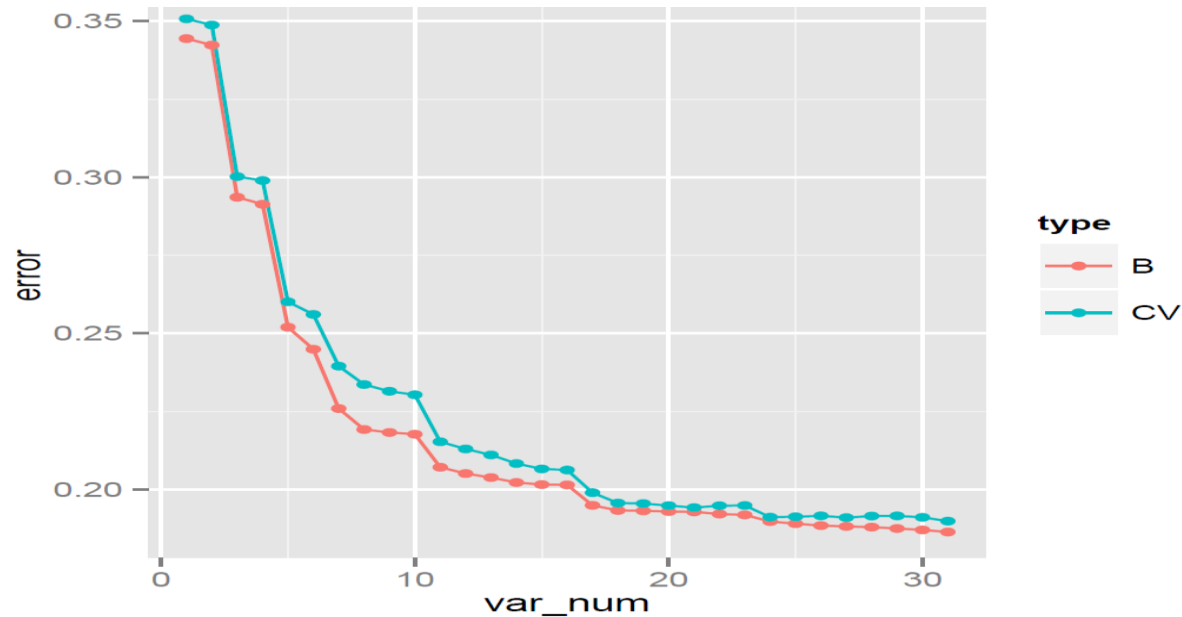
<b>build</b>	Set na kojem se gradi model za testiranje korisnosti varijabli
<b>model</b>	Izbor vrste modela: GLM (general linear model), DT (decision tree), NB (Naive Bayes) or SVM (support vector machine)
<b>tgt</b>	Ime stupca koji označava ciljnu varijablu
<b>cross_valid</b>	Ime seta za krosvalidaciju
<b>attribute_importance</b>	Izbor dali da se uzimaju u obzir samo varijable koji imaju attribute importance veći od nule (TRUE/FALSE)
<b>columns_to_skip</b>	Koje varijable bi se trebale preskočiti u procesu odabiranja
<b>c_per</b>	Koliko veliki dio BUILD seta treba uzeti za krosvalidaciju, ukoliko set za krosvalidaciju nije dezigniran ( $0 \leq c\_per \leq 1$ ).
<b>sample_type</b>	Koja vrsta sempliranja bi se trebala koristiti pri uzimanj krosvalidacijskog seta iz build seta: ST (stratified) or R (random sampling).

# Izlazne vrijednosti

<b>sample_size</b>	Veličina BUILD seta
<b>var_number</b>	Broj varijable po redu koja je automatski predviđena da će biti korisna u prediktivnom modelu ( $\Delta AUC > 0$ , attribute importance $> 0$ )
<b>var_number_full</b>	Broj varijable po redu koja je testirana (attribute importance $> 0$ )
<b>added_variable</b>	Ime varijable pridodane za zadnju iteraciju prediktivnog modela ( $\Delta AUC > 0$ , attribute importance $> 0$ )
<b>added_variable_full</b>	Ime varijable pridodane za zadnju iteraciju prediktivnog modela (AUC $> 0$ )
<b>R2_build</b>	R na kvadrat modela na BUILD setu
<b>R2_test</b>	R na kvadrat modela na krosvalidacijskom setu
<b>auc_build</b>	Površina ispod ROC krivulje za trenutni prediktivni model na BUILD setu ( $\Delta AUC > 0$ , attribute importance $> 0$ )
<b>auc_build_full</b>	Površina ispod ROC krivulje za trenutni prediktivni model na BUILD setu sa restrikcijom za attribute importance $> 0$ .
<b>auc_test</b>	Površina ispod ROC krivulje za trenutni prediktivni model na krosvalidacijskom setu
<b>m_start</b>	Vrijeme pokretanja gradnje zadnjeg prediktivnog modela

# Rezultati

- Značajna ušteda vremena
- Pojednostavljanje prediktivnog modela
- Lakša interpretacija modela
- Smanjena mogućnost „overfittinga“



...Hvala!  
Pitanja?

[Ivan.osman@multicom.hr](mailto:Ivan.osman@multicom.hr)

