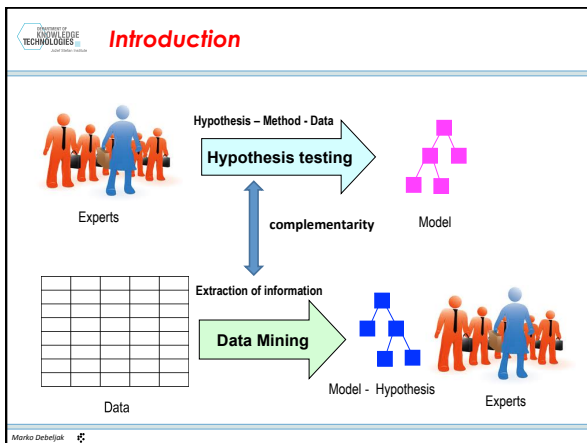


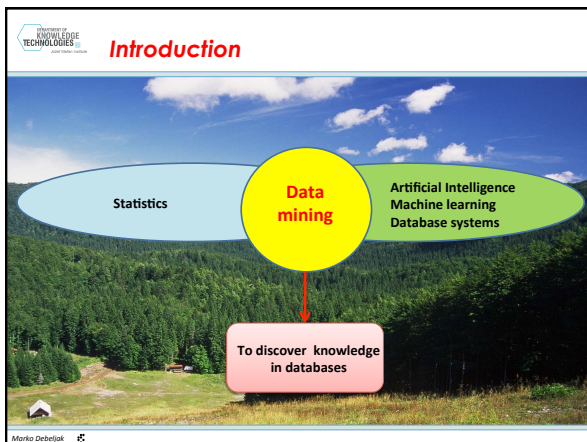
UNIVERSITY OF
KOPENICE
TECHNOLOGIES

Introduction of data mining to forest ecology and forestry

Marko Debeljak
Jožef Stefan Institute, Slovenia

Marko Debeljak





Knowledge discovery in data bases (KDD)

What is KDD?

Frawley et al., 1991: "KDD is the non-trivial **process** of identifying valid, novel, potentially useful, and ultimately understandable **patterns** in data"

The key task is the discovery of **previously unknown knowledge**

How to find patterns in data?

Data mining (DM) – **central step** in the KDD process concerned with applying computational techniques to actually find **patterns** in the data (15-25% of the effort of the overall KDD process).

- step 1: **data preprocessing (50%)**
- step 3: **evaluation of discovered patterns (25%)**

Data mining and machine learning

Data mining focuses on the discovery of **previously unknown knowledge** and integrates **machine learning**.

Machine learning focuses on **descriptions** and **prediction**, based on known properties **learned** from the training empirical data (examples) using computer algorithms.

Learning from examples is called **inductive learning**

If the goal of **inductive learning** is to obtain model that **predicts** the value of target variable from learning examples, the it is called **predictive or supervised learning.**

Data mining

The most relevant notions of data mining:

1. **Data**
2. **Patterns**
3. **Data mining algorithms**

Data

Data stored in **one flat** table.
Each example represented by a **fixed number** of attributes.

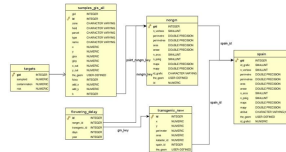
Data stored in **original** tables or relations.
No loss of information, due to aggregation.

PROPOSITIONAL data mining

Loss of information due to aggregation

RELATIONAL data mining

| Properties of objects | | | |
|-----------------------|--------------------|-------------------|-------------------------|
| Incidence (%) | Wind direction (°) | Wind speed (km/h) | Catch crossing rate (%) |
| 10 | 123 | 3 | 8 |
| 12 | 88 | 4 | 7 |
| 14 | 121 | 6 | 5 |
| 18 | 147 | 2 | 4 |
| 20 | 93 | 1 | 6 |
| 22 | 115 | 3 | 1 |
| ... | ... | ... | ... |



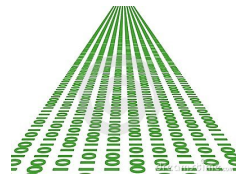
Marko Dehajuk

Data

Data are **not stored at all** but they continuously flow through algorithm.

Each example can propositional or relational.

DATA STREAM mining



Marko Dehajuk

Pattern

2. What is a pattern?

A **pattern** is defined as: "A statement (**expression**) in a given language, that describes (**relationships** among) the facts in a **subset** of the given data and is (in some sense) simpler than the enumeration of all facts in the subset" (Frawley et al. 1991, Fayyad et al. 1996).

- Classes of patterns considered in data mining
- A. **equations**,
- B. **decision trees, relational decision trees**
- C. association, classification, and regression **rules**.

Selection of the pattern type depends on the data mining task at hand.

Pattern

A. Equations
 To predict the value of a **target** (dependent) variable as a **linear or non linear combination** of the **input** (independent) variables:
 - Algebraic equations

To predict the behavior of **dynamic** systems, which change their rate over time:
 - Difference equations
 - Differential equations

Pattern


B. Decision trees
 To predict the value of **one or several target dependent variables** from the values of other **independent variables** by **decision tree**.

Decision tree has a hierarchical structure, where:
 - each internal **node** contains a test on an independent variable,
 - each **branch** corresponds to an outcome of the test (critical values of independent variable),
 - each **leaf** gives a prediction for the value of the dependent (predicted) variable.

Pattern

Decision tree is called:


- A **classification tree**: value of dependent variable in leaf is **discrete (finite set of nominal values)**: e.g., (yes, no), (spec. A, spec. B, ...)
- A **regression tree**: value of dependent variable in leaf is a **constant (infinite set of values)**: e.g., 120, 220, 312, ...
- A **model tree**: leaf contains **linear model** predicting the value of piece-wise linear function:
 $out-crossing\ rate = 12.3\ distance - 0.123\ wind\ speed + 0.00123\ wind\ direction$

 **Pattern**

C. Rules
To perform association analysis between variables discovered by **association rules**.


The **rule denotes** patterns of the form:
IF „Conjunction of conditions“ **THEN** „Conclusion.“

- For **classification rules**, the conclusion assigns one of the possible **discrete** values to the dependent variable (finite set of nominal values): e.g., (yes, no), (spec. A, spec. B, spec. D)
- For **predictive rules**, the conclusion gives a prediction for the value of the dependent variable (**infinite** set of values): e.g., 120, 220, 312, ...

 **Algorithm**

3. What is data mining algorithm?
Algorithm in general:
 - a **procedure** (a finite set of well-defined instructions) for accomplishing some task which will terminate in a defined end-stat.

Data mining algorithm:
 - a computational process for finding patterns in data

 **Data mining (DM) - algorithm**

Selection of algorithm depends on problem at hand:

1. **Equations = Linear and multiple regressions, equation discovery**
2. **Decision trees = Top/down induction of decision trees**
3. **Rules = Rule induction**

UNIVERSITY OF KNOWLEDGE TECHNOLOGIES
 INSTITUT "JOŽEF STEFAN"
 LJUBLJANA, SLOVENIJA

DATA MINING – CASE STUDIES

UNIVERSITY OF KNOWLEDGE TECHNOLOGIES
 INSTITUT "JOŽEF STEFAN"
 LJUBLJANA, SLOVENIJA

Applications – forest ecology and forestry

| | |
|----------------------------|--------------------------|
| POPULATION DYNAMICS | HABITAT MODELLING |
| GENE FLOW MODELLING | RISK MODELLING |

UNIVERSITY OF KNOWLEDGE TECHNOLOGIES
 INSTITUT "JOŽEF STEFAN"
 LJUBLJANA, SLOVENIJA

Applications – forest ecology and forestry

Propositional and relational supervised data mining:

- Simple data mining
- Data mining of time series
- Spatial data mining

1. Equations:

- Algebraic equations
- Differential equations

2. Single and multi target decision trees:

- Classification trees
- Regression trees
- Model trees (single target only)

Algebraic equations: POPULATION DYNAMICS

Algebraic equations: CIPER

ECOLOGICAL MODELLING 215 (2008) 180-189

available at www.sciencedirect.com

ScienceDirect

Journal homepage: www.elsevier.com/locate/ecolmodel

Modeling radial growth increment of black alder (*Alnus glutionsa* (L.) Gaertn.) tree

Jana Laganis^{a,*}, Aleksandar Pečkov^b, Marko Debeljak^b

^a Laboratory for Environmental Research, University of Nova Gorica, Vipavska 13, Nova Gorica, Slovenia
^b Department of Knowledge Technologies, "Jozef Stefan" Institute, Jamova 39, Ljubljana, Slovenia

Algebraic equations: POPULATION DYNAMICS

Measured radial increments:
 - 8 trees
 - 69 years old

Hydrological conditions
 (HMS Lendava, monthly data on minimal, average and maximum values)
 - Ledava River levels
 - groundwater levels

Management data
 (thinning; m³/y removed from the stand; Forestry Unit Lendava)

Meteorological conditions
 (monthly data, HMS Lendava):
 - Time of solar radiation (h),
 - precipitation (mm),
 - ET (mm)
 - Number of days with white frost
 - Number of days with snow
 - T: max, aver, min
 - Cumulative T > 0°C, > 5°C, and > 10°C
 - Number of days with:
 - min T > 0°C
 - min T < -10°C
 - min T < -4°C
 - min T > 25°C
 - max T > 10°C
 - max T > 25°C

Dataset

• Monthly data + aggregated data (AMJ, MJJ, JJA, MJJA etc.)
 • **Σ: 333** attributes; 35 years

Algebraic equations: POPULATION DYNAMICS

• 52 different combinations of attributes were tested.
Σ: 124 models

| Experiment | RRSE | # eq. elements |
|------------|---------|----------------|
| jnj3_2m | 0,7282 | 6 |
| jnj3_3s | 0,7599 | 6 |
| jnj3_1s | 0,7614 | 6 |
| jnj3_4m | 0,76455 | 3 |
| jnj2_2 | 0,7685 | 5 |
| jly_4xl | 0,7686 | 6 |

Algebraic equations: POPULATION DYNAMICS

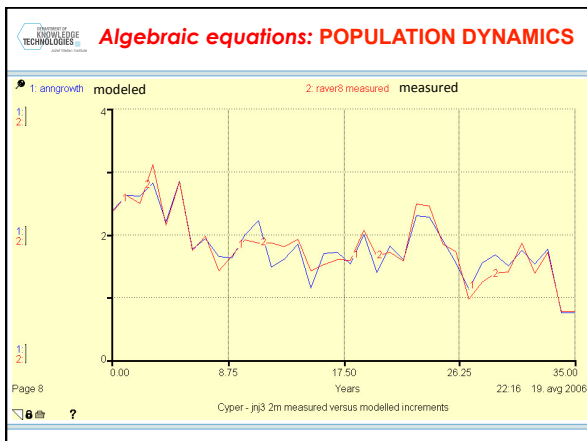
Model jnj3_2m:

RadialGrowthIncrement =
 + -0.0511025526922 minL8-10^1
 + -0.0291795197998 maxL8-10^1
 + -0.017479975134 t-sun4-7^1
 + 0.0346935385853 t-sun8-10^1
 + -1.950606536e-05 t-sun8-10^2
 + -2.01014710248 d-wf-4-7^1
 + 9.35586778387e-05 minL4-7^1 t-sun4-7^1
 + -0.000179339939732 minL4-7^1 t-sun8-10^1
 + 6.45688563611e-05 minL8-10^1 t-sun8-10^1
 + 3.06551434164e-05 maxL8-10^1 t-sun4-7^1
 + 0.00282485442386 t-sun4-7^1 d-wf-4-7^1
 + -0.00141078675225 t-sun8-10^1 d-wf-4-7^1
 + 7.91071710872

Relative Root Squared Error = **0.728229824611**

Correlation between average measured (r-aver8) and modeled increments: linear regression: **R² = 0.8771**

8 out of 333 attributes



Algebraic equations: WATER CYCLE

Prediction of drainage water : CIPER

MEĐNARODNA POPLOVNA SILA JOŽEFA ŠTEFANA
 VIŠJA ŠOLSKA VEŠTOVSKA INŽENIRSKA ŠOLA

Vladimir Kuzmanovski

Integration of expert knowledge and predictive learning: Modelling water flows in agriculture
 Master Thesis

Integracija ekspertnega znanja z napovednim učenjem: Modeliranje vodnih tokov v poljedelstvu
 Magistrsko delo

Supervisor: Prof. Dr. Marko Delžek
 Co-Supervisor: Prof. Dr. Ivo Džuranič
 Ljubljana, Slovenia, September 2012

Differential equations: COMMUNITY STRUCTURE

Phosphorus

$$\frac{d(ps)}{dt} = \underbrace{ps_krivica \cdot \frac{q_krivica}{7 \cdot 10^6} + ps_misca \cdot \frac{q_misca}{7 \cdot 10^6}}_{\text{water in-flow out-flow}} + \underbrace{ps_radovna \cdot \frac{q_radovna}{7 \cdot 10^6} - ps \cdot \frac{q_jezemica}{7 \cdot 10^6}}_{\text{respiration}} - \underbrace{ps \cdot \frac{q_natega}{7 \cdot 10^6}}_{\text{respiration}} + \underbrace{0.0022 \cdot phyto^2 - 0.072 \cdot \frac{temp - 2.7}{20.4 - 2.7}}_{\text{respiration}} + \underbrace{0.07 \cdot daph \cdot 0.0026 \cdot \frac{temp}{12.3}}_{\text{growth}} - \underbrace{0.0023 \cdot phyto \cdot 0.21}_{\text{growth}} \cdot \frac{ps}{ps + 0.00042} \cdot \frac{temp}{16.7} \cdot \frac{light}{170} \cdot e^{\left(1 - \frac{light}{170}\right)} \quad (10)$$

Differential equations: COMMUNITY STRUCTURE

Phytoplankton

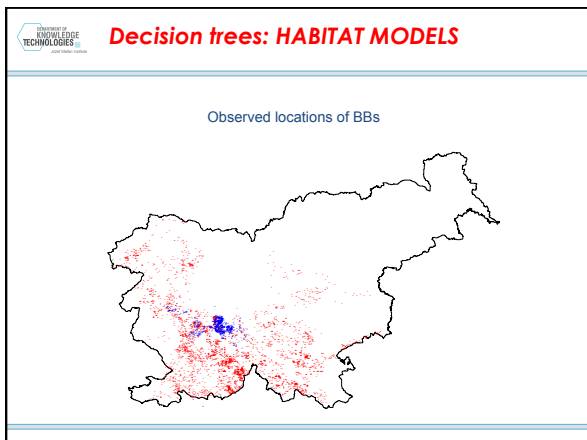
$$\frac{d(phyto)}{dt} = \underbrace{phyto \cdot 0.21}_{\text{growth}} - \underbrace{ps \cdot 0.00042}_{\text{respiration}} \cdot \frac{temp}{16.7} \cdot \frac{light}{170} \cdot e^{\left(1 - \frac{light}{170}\right)} - \underbrace{phyto^2 \cdot 0.072}_{\text{respiration}} \cdot \frac{temp - 2.7}{19.7 - 2} - \underbrace{phyto \cdot 0.5}_{\text{sedimentation}} \cdot \frac{temp - 2}{18 - 4} - \underbrace{daph \cdot 0.5}_{\text{grazing}} \cdot \frac{temp - 2.6}{18 - 4} \cdot (1 - \exp(-0.58 \cdot phyto)) - 0.56 \cdot phyto \quad (13)$$

growth
respiration
sedimentation
grazing

Differential equations: COMMUNITY STRUCTURE

Zooplankton

$$\frac{d(daph)}{dt} = \underbrace{0.14 \cdot daph \cdot 0.5 \cdot \frac{temp - 2.6}{18 - 4}}_{\text{Feeds on phytoplankton}} - \underbrace{daph \cdot 0.026 \cdot \frac{temp}{12.3}}_{\text{respiration}} - \underbrace{0.01 \cdot \frac{daph^2}{0.001 + daph}}_{\text{mortality}} \cdot (1 - \exp(-0.58 \cdot phyto)) + 0.56 \cdot phyto$$



Decision trees: HABITAT MODELS

The training dataset

- Positive examples:**
 - Locations of bear sightings (Hunting association; telemetry)
 - Females only
 - Using home-range (HR) areas instead of "raw" locations
 - Narrower HR for optimal habitat, wider for maximal
- Negative examples:**
 - Sampled from the unsuitable part of the study area
 - Stratified random sampling
 - Different land cover types equally accounted for

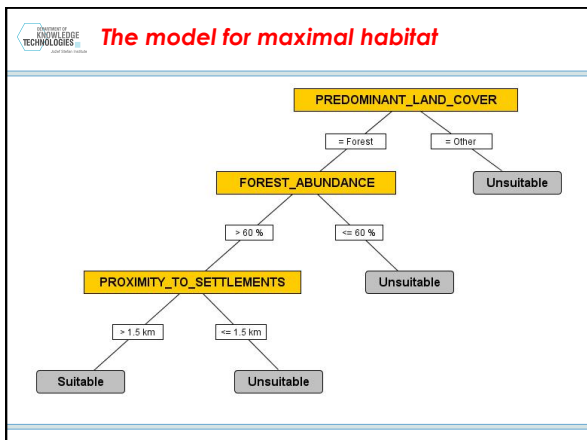
Decision trees: HABITAT MODELS

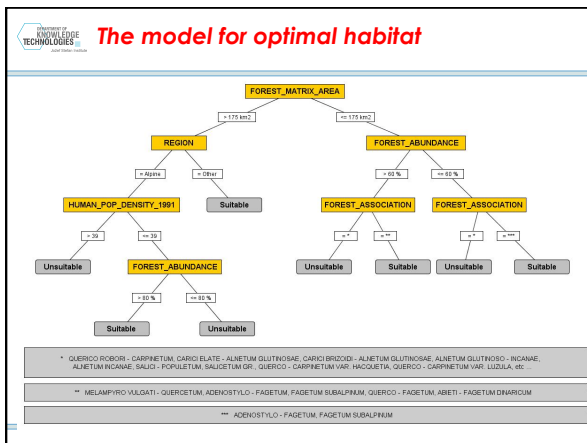
Dataset

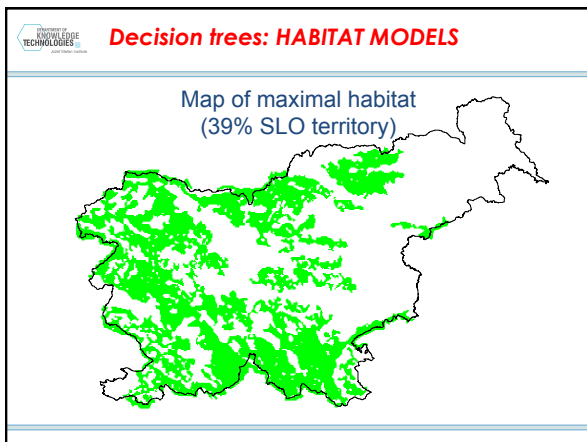
```

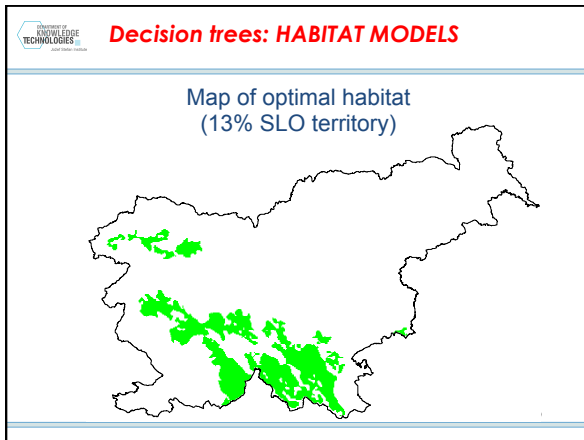
1.73,26,0.0,1.88,0.2,70,7.20,1.0,0.1,0.60,0.0,0.0,0.2,0.0,0.4123,0.0,0.0,63,211,11,11,11,83,213,213,0.0,4155
1.62,37,0.0,2.88,0.2,70,7.20,1.0,0.1,0.60,0.0,0.0,0.2,1.53,0.3640,0.0,0.0,-1347,63,211,11,11,11,83,213,213,11,89,3858
2.0,99,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,1.6,82,0.10404,0.2074,-309,48,0.0,11,11,11,83,83,83,0.20,3862
2.0,100,0.0,1.76,0.16,71,0.12,0.0,0.0,0.0,0.0,0.0,1.6,82,0.7590,0.1681,-319,-942,0.0,11,11,11,0.0,0.0,20,4088
1.8,91,0.0,1.52,0.59,41,0.0,0.0,0.0,0.4,0.0,0.0,5.1,6,82,0.6950,0.1505,-166,879,9,57,11,11,11,281,281,281,0.20,3199
4.3,0,96,9,0,75,0,33,67,0,0,0,0,0,0,1,2,0,0,0,0,0,1,2,54,0,0,0,465,-96,-191,4,225,11,31,31,41,72,272,60,619,4013
1.34,65,0.0,2.51,9,76,9,5,1,4,1,0,1,0,29,0,0,0,0,0,1,2,54,0,3000,0,841,-111,-284,34,220,11,41,41,151,141,112,60,619,3897
1,100,0,0,3,52,0,66,6,3,5,9,6,7,38,40,0,0,0,0,0,1,17,64,0,6062,0,932,-603,-71,100,3,37,11,41,41,171,232,202,4,24,3732
.....
.....
.....
    
```

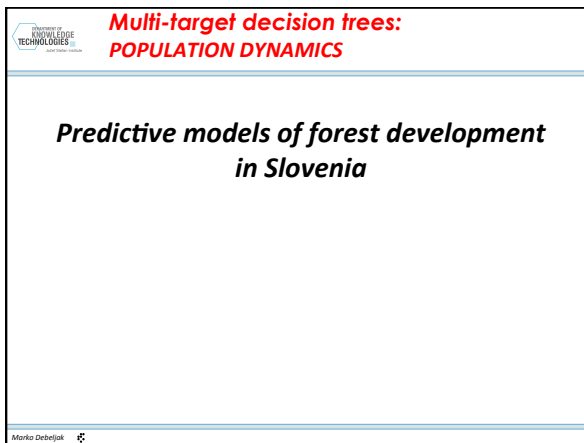
↑ Present: 1
Absent: 0

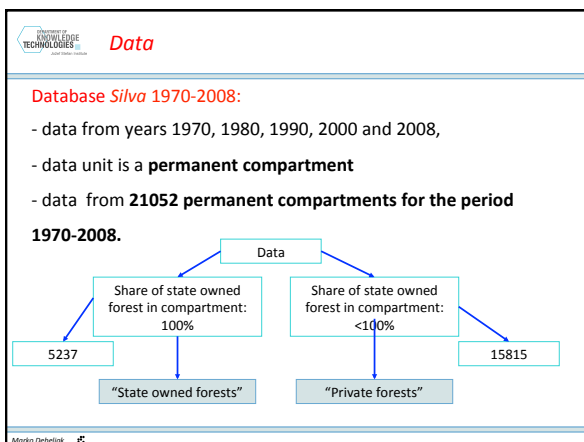


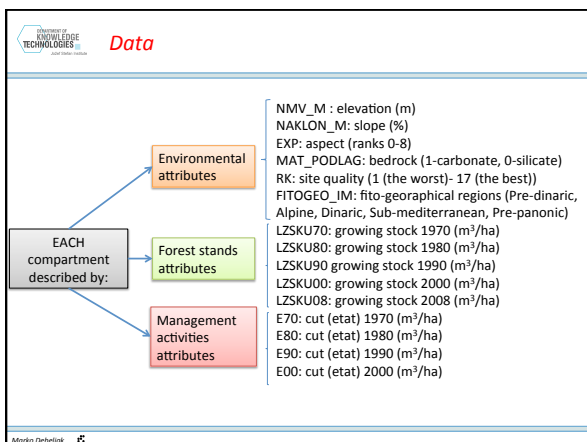


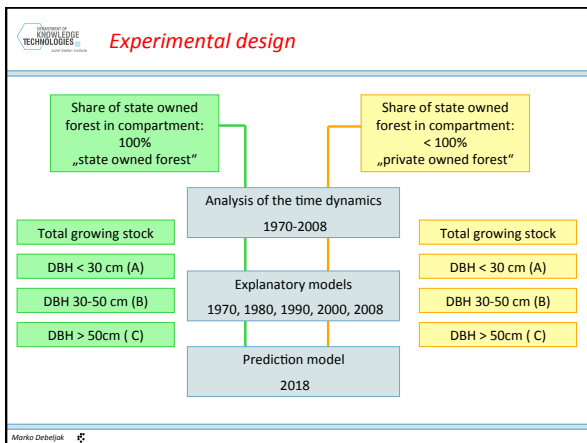


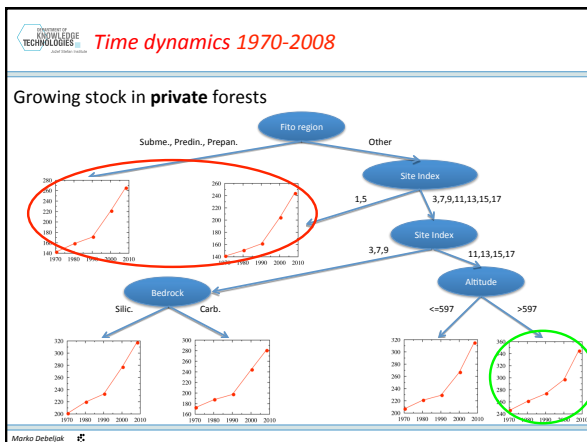












Prediction for private forests in 2018

Step 1: - Extrapolation of linear trends from 1990-2000 to 2008
- Verification on real data for 2008

| 2008 | LM1 | LM2 | LM3 | LM4 |
|---------------------------------|----------------|-----------------|-----------------|-----------------|
| Linear regression a (k, n) | 5.8801, -11543 | 3.4203, -6565.4 | 1.6311, -2957.4 | -0.8628, 2077.9 |
| Average real growing stock 2008 | 264.2 | 312.1 | 339.5 | 380.5 |
| Model prediction for 2008 | 264.2 | 302.6 | 317.8 | 345.4 |
| Difference (%) | 0.0 | -3.0 | -6.4 | -9.2 |

Marko Debeljak

Prediction for private forests in 2018

Step 2: - Prediction of wood stock for 2018 with the model for 2008

- Correlation coefficient: 0.7002
- Mean absolute error: 41.9041 m³/ha
- Root mean squared error: 55.0751 m³/ha
- Relative absolute error: 59.7019 %
- Root relative squared error: 71.3877 %

| 2008 | LM1 | LM2 | LM3 |
|---|-------|-------|-------|
| Average real growing stock (m³/ha) | 224.4 | 280.0 | 351.7 |
| Average predicted growing stock (m³/ha) | 224.4 | 280.0 | 351.7 |
| Mean absolute error (MAE) | 49.1 | 43.3 | 35.5 |

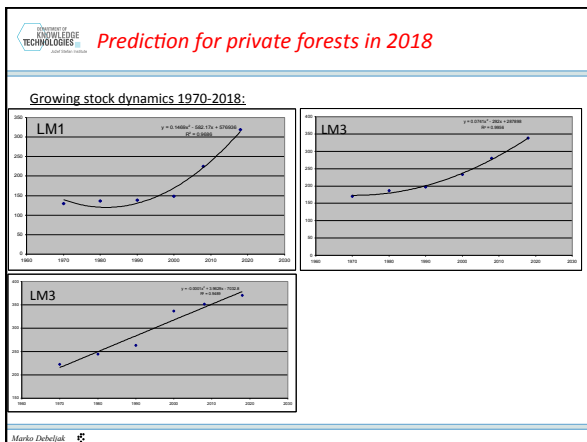
Marko Debeljak

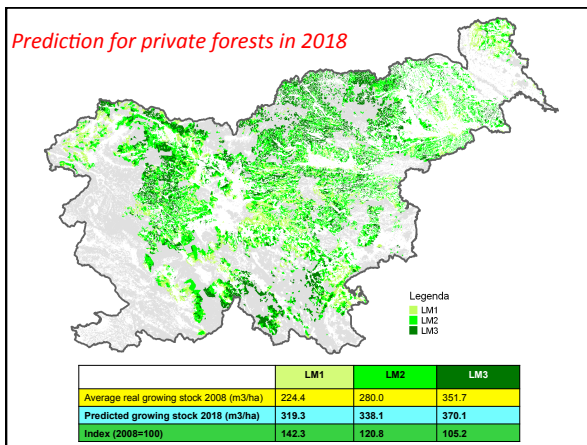
Prediction for private forests in 2018

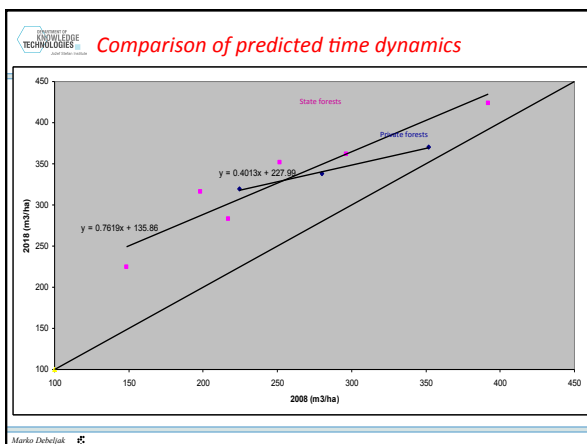
Step 2: - Extrapolation of linear trends from 2000-2008 to the year 2018

| 2018 | LM1 | LM2 | LM3 |
|---|----------------|----------------|-----------------|
| Linear regression (k, n) | 9.4818, -18815 | 5.8484, -11464 | 1.8368, -3336.6 |
| Average real growing stock 2008 (m³/ha) | 224.4 | 280.0 | 351.7 |
| Predicted growing stock 2018 (m³/ha) | 319.3 | 338.1 | 370.1 |
| Index (2008=100) | 142.3 | 120.8 | 105.2 |

Marko Debeljak







GIS data mining: COMMUNITY STRUCTURE

Ecological Informatics 5 (2010) 256–266

Contents lists available at ScienceDirect

Ecological Informatics

journal homepage: www.elsevier.com/locate/ecolinf

Estimating vegetation height and canopy cover from remotely sensed data with machine learning[☆]

Daniela Stojanova^a, Panče Panov^{b,*}, Valentin Gjorgjioski^b, Andrej Kobler^a, Sašo Džeroski^b

^a Slovenian Forestry Institute, Večna pot 2, SI-1000 Ljubljana, Slovenia
^b Jozef Stefan Institute, Department of Knowledge Technologies, Jamova cesta 39, SI-1000 Ljubljana, Slovenia

GIS data mining: COMMUNITY STRUCTURE

Data

- Locations: Kras region (Karst)
- Attributes:
 - Statistical information (max, min, avg, std) from Landsat, IRS, SPOT & aerial photographs
 - Normalized Difference Vegetation Index (NDVI)
 - Textures
 - Relief: Aspect, Slope, Elevation
- Targets (forest properties) from LiDAR data:
 - Vegetation height (H)
 - Canopy Cover (CC)

Landsat and LiDAR data

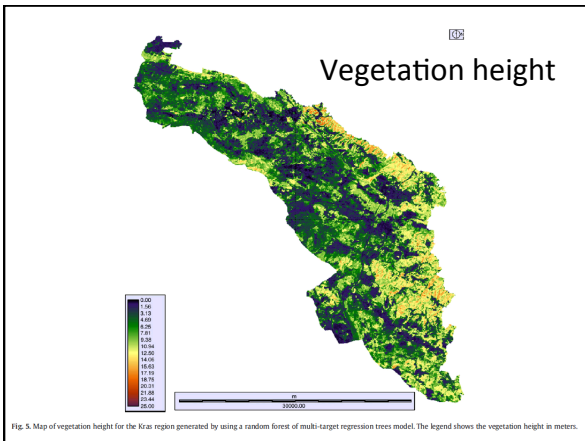
UNIVERSITY OF
SOUTH ALABAMA
TECHNOLOGICAL CENTER

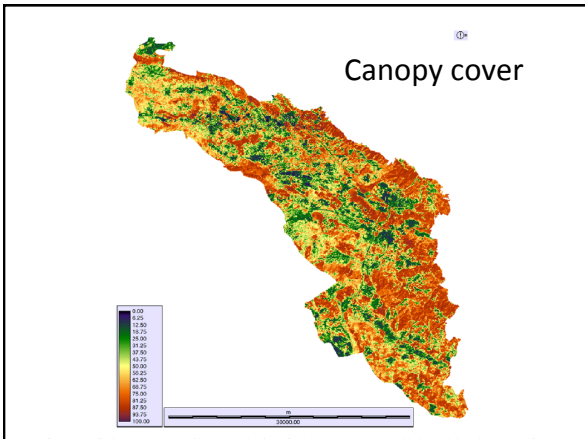
GIS data mining: COMMUNITY STRUCTURE

Machine Learning Methodology

- WEKA
 - Regression (RT) and Model (MT) trees
 - Bagging of Model Trees (BagMT)
- CLUS
 - Single Target Regression Trees (STRT)
 - Multi Target Regression Trees (MTRT)
 - Ensembles: Bagging of Model Trees (MTBG) and Random Forest (MTRF)

SELECTED: Random Forest Multi Target Regression Trees





GIS data mining: COMMUNITY STRUCTURE

- The integration of LiDAR and RS promises detailed estimation of forest parameters
- Detail forest vegetation maps can be generated and used for forest management

Data mining of time series: COMMUNITY STRUCTURE

Ecological Modelling 222 (2011) 2524–2529

Contents lists available at ScienceDirect

Ecological Modelling

Journal homepage: www.elsevier.com/locate/ecolmodel


Analysis of time series data on agroecosystem vegetation using predictive clustering trees

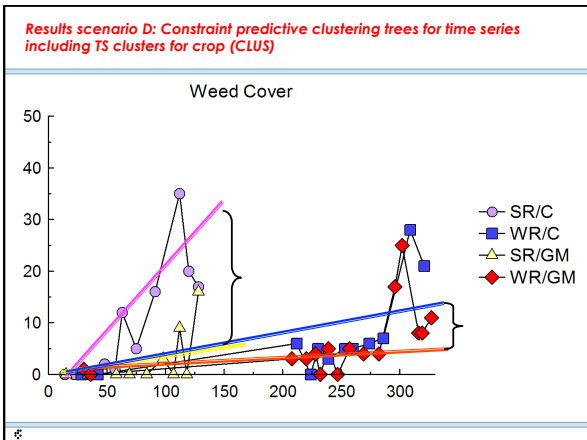
Marko Debeljak^{a,*}, Geoffrey R. Squire^b, Dragi Kocev^a, Cathy Hawes^b, Mark W. Young^b, Sašo Džeroski^a

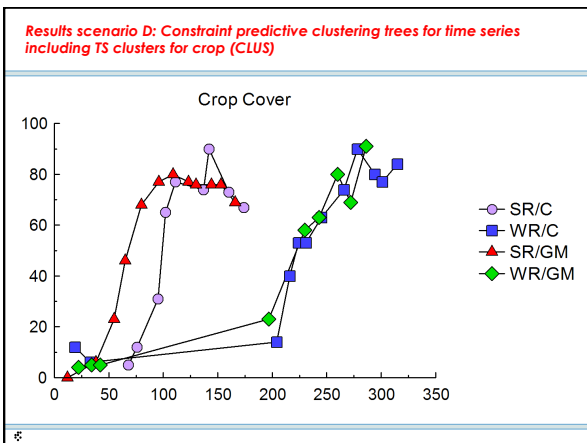
^a Jozef Stefan Institute, Jamova cesta 39, 1000 Ljubljana, Slovenia
^b Scottish Crop Research Institute, Invergowrie, Dundee DD2 9DA, Scotland, UK

Data

- 130 sites, monitoring every 7 to 14 days for 5 month (2665 samples: 1322 conventional, 1333, HT OSR observations)
- Each sample (observation) described with 65 attributes
- Original data collected by Centre for Ecology and Hydrology, Rothamsted Research and SCRI within Farm Scale Evaluation Program (2000, 2001, 2002)







Relational data mining: GENE FLOW MODELLING

Ecological Modelling 245 (2012) 75–83

Contents lists available at SciVerse ScienceDirect

Ecological Modelling

journal homepage: www.elsevier.com/locate/ecolmodel

Using relational decision trees to model out-crossing rates in a multi-field setting

Marko Debeljak^{a,b,*}, Aneta Trajanov^a, Daniela Stojanova^{a,b}, Florence Leprince^c, Sašo Džeroski^{a,b,d}

^a Jozef Stefan Institute, Jamnova 39, 1000 Ljubljana, Slovenia
^b Jozef Stefan International Postgraduate School, Jamnova 39, 1000 Ljubljana, Slovenia
^c INRA-Misc-Institut de Vegetation, 21, Chemin de Pélissier, 64121 Montfort, France
^d Centre of Excellence for Integrated Approaches in Chemistry and Biology of Proteins, Jamnova 39, 1000 Ljubljana, Slovenia

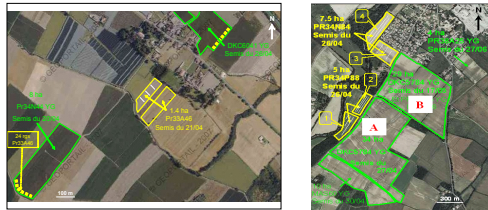
Marko Debeljak

Relational data mining: GENE FLOW MODELLING

Initial questions:

To **what extent** will GM maize grown on **Geens** genetically **interfere** with the maize on **Yelows**?

Will this interference be small enough to **allow co-existence**?



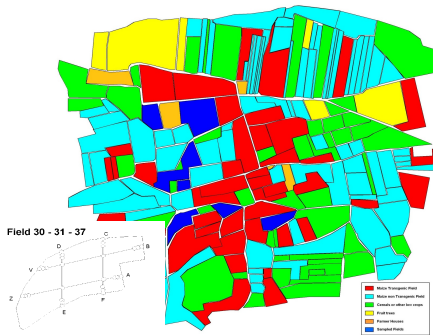
Marko Dehajuk

Spatial temporal relations

2004:
40 GM fields
7 non-GM fields
181 sampling points

2005:
17 GM fields
4 non-GM fields
127 sampling points

2006:
43 GM fields
4 non-GM fields



4 non-GM fields

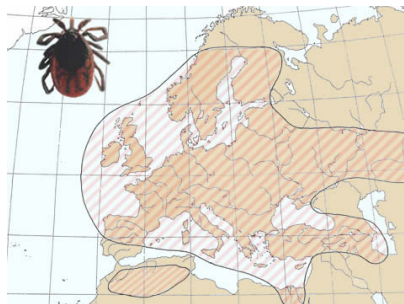
Relational data mining: GENE FLOW MODELLING

Data scattered over several **tables** or **relations**:

- A table storing **general information** on each field (e.g., area)
- A table storing the **cultivation techniques** for each field and each year
- A table storing the **relations** (e.g., distance) between fields

Marko Dehajuk

- Present in most parts of Europe
- Habitat: mixed, shadow forest with diverse tree composition with shrubs and dense vegetation cover with a lot of litter

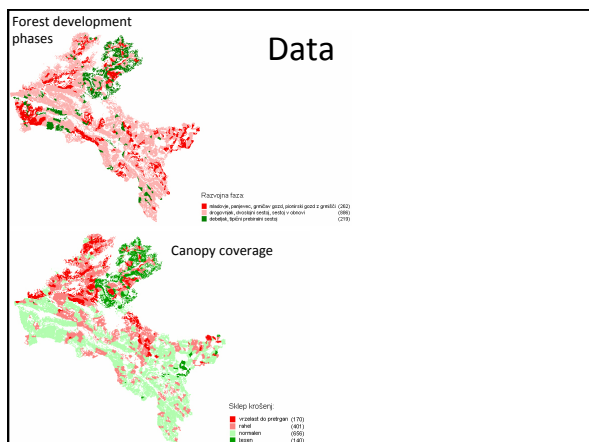


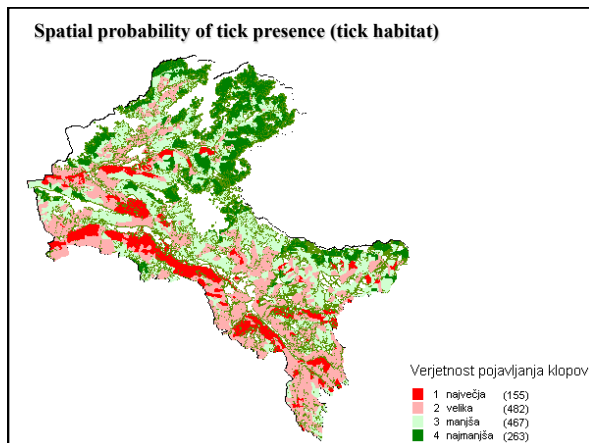
Visitors of natural heritage in upper valley of Soča are exposed to the risk of tick-borne diseases.

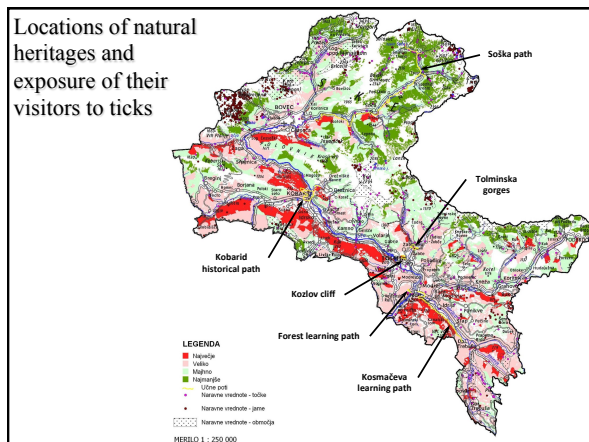
The goal: To identify areas of tick habitats where visitors can be infected with tick-borne diseases.

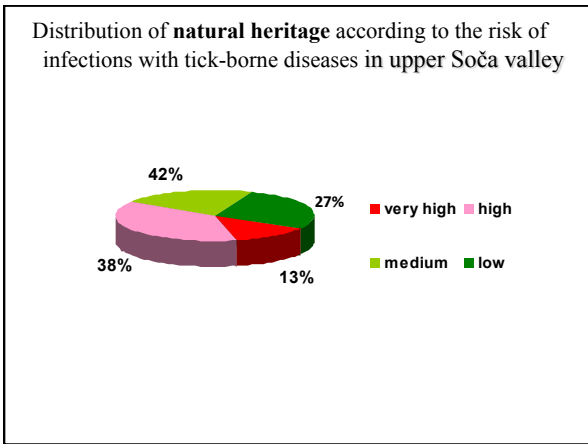
Study area: Zgornje Posočje

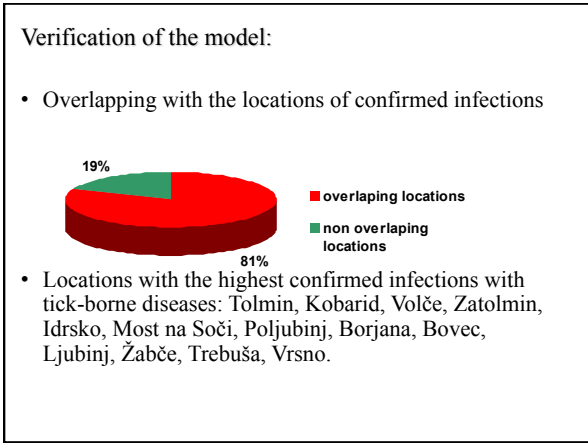















Multi-target regression model: RISK MODELLING

Soil Use and Management


Soil Use and Management, March 2009, 25, 66-77 doi: 10.1111/j.1475-2743.2009.00196.x

Potential of multi-objective models for risk-based mapping of the resilience characteristics of soils: demonstration at a national level

M. DEBELJAK¹, D. KOČEV¹, W. TOWERS², M. JONES², B. S. GRIFFITHS^{3,*} & P. D. HALLET³
¹Department of Knowledge Technologies, Jozef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia; ²Maccaulay Institute, Craigiebuckler, Aberdeen AB15 8QH, UK, and ³Scottish Crop Research Institute, Invergowrie, Dundee DD2 5DA, UK

Multi-target regression model: RISK MODELLING

The dataset: soil samples taken on 26 location throughout SCO

The dataset: The flat table of data: 26 by 18 data entries

Multi-target regression model: RISK MODELLING

The dataset:

- > **physical properties:** soil texture: sand, silt, clay
- > **chemical properties:** pH, C, N, SOM (soil organic matter)
- > **FAO soil classification:** Order and Suborder
- > **physical resilience:** resistance to compression: 1/Cc, recovery from compression: Ce/Cc, overburden stress: eg, recovery from overburden stress after two days cycles: eg2dc
- > **biological resilience:** heat, copper

Multi-target regression model: RISK MODELLING

Different scenarios and multi-target regression models have been constructed:

A model predicting the resistance and resilience of soils to copper perturbation.

Independent variables

Major soil subgroup
Sand
Silt
Clay
pH
C
N
Soil organic matter

Dependent variables

Copper resistance
Copper resilience

Model 10 (see Table 1)

```

    graph TD
      Root["pH > 4.04"] -- Yes --> Sand["Sand > 0.5"]
      Root -- No --> C["C > 7.1"]
      Sand -- Yes --> pH54["pH > 5.4"]
      Sand -- No --> MajorSubgroup["Major soil subgroup  
is BFS, BFSg, NCG, PP"]
      pH54 -- Yes --> R1["[41.6; 58.3]  
2 soils"]
      pH54 -- No --> R2["[59.4; 47.4]  
7 soils"]
      MajorSubgroup -- Yes --> R3["[58.6; 71.1]  
9 soils"]
      MajorSubgroup -- No --> R4["[74.1; 57.1]  
3 soils"]
      C -- Yes --> R5["[103.2; 98.3]  
3 soils"]
      C -- No --> R6["[76.1; 83.8]  
2 soils"]
  
```

Legend:
 BFS = brown forest soil
 BFSg = brown forest soil with gleying
 NCG = nonclacareous gley
 PP = peaty podzol

Multi-target regression model: RISK MODELLING


The increasing importance of mapping soil functions to **advise on land use and environmental management** - to make a **map of soil resilience for Scotland**.

The **models = filters** for existing GIS datasets about physical and chemical properties of Scottish soils.

Multi-target regression model: RISK MODELLING

Macaulay Institute (Aberdeen): soils data – attributes and maps:

Approximately 13 000 soil profiles held in database
 Descriptions of over 40 000 soil horizons



Application

Experiment 1

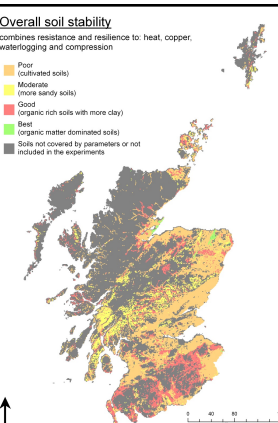
| | |
|--|--|
| Independent Attributes | Dependent attributes |
| <ul style="list-style-type: none"> o MAJOR_SOIL_SUBGROUP o Sand o SF o Clay o pH o C o N o SOM | <ul style="list-style-type: none"> o Heat_resil o Heat_Resil o Cu_resil o Cu_Resilience o Over_resil_eg o Over_Resil_Eg o I/CC o CerCC |


Validated parameters:
 RMSE: [15.4832 20.8162 11.0124 18.8897 0.0809 0.2114 0.3419 0.0098]
 Correlation Coefficient: [-0.3097 -0.0784 0.8042 0.4448 -0.1465 0.2884 0.388 0.8447]

C > 7.4
 ***res [56.47,110.4,103.2498,3.0,789667,1.389333,0.371243,0.058778]: 3
 ***res: C > 4.28
 ***res: SPT > 4.09
 | ***res [48.834,93.06,70.2444,4.34,0.8028,1.1336,0.426798,0.016037]: 3
 | ***res [39.631,203333,72.403333,79.0,0.003333,0.970,0.047719,0.040945]: 3
 ***res: N > 0.14
 ***res [58.94667,86.833333,54.246667,57.35,0.8145,0.9418,0.914021,0.011320]
 ***res: [49.473333,81.59,116667,59.866667,0.865,0.911,1.54998,0.020033]: 3

Overall soil stability
 combines resistance and resilience to heat, copper, waterlogging and compression

- o Poor (cultivated soils)
- o Moderate (more sandy soils)
- o Good (organic rich soils with more clay)
- o Best (organic matter dominated soils)
- o Soils not covered by parameters or not included in the experiments




 **Conclusions**

What can data mining do for you?


Knowledge discovered by analyzing data with DM techniques can help:

- Understand the domain studied
- Make predictions/classifications
- Support decision processes in environmental management

 **Conclusions**

What data mining cannot do for you?

- The law of information conservation (garbage-in-garbage-out)
- The knowledge we are seeking to discover has to come from the combination of data and background knowledge
- If we have very little data of very low quality and no background knowledge no form of data analysis will help

 **Conclusions**

Side-effects?

- Discovering problems with the data during analysis
 - missing values
 - erroneous values
 - inappropriately measured variables
- Identifying new opportunities
 - new problems to be addressed
 - recommendations on what data to collect and how
