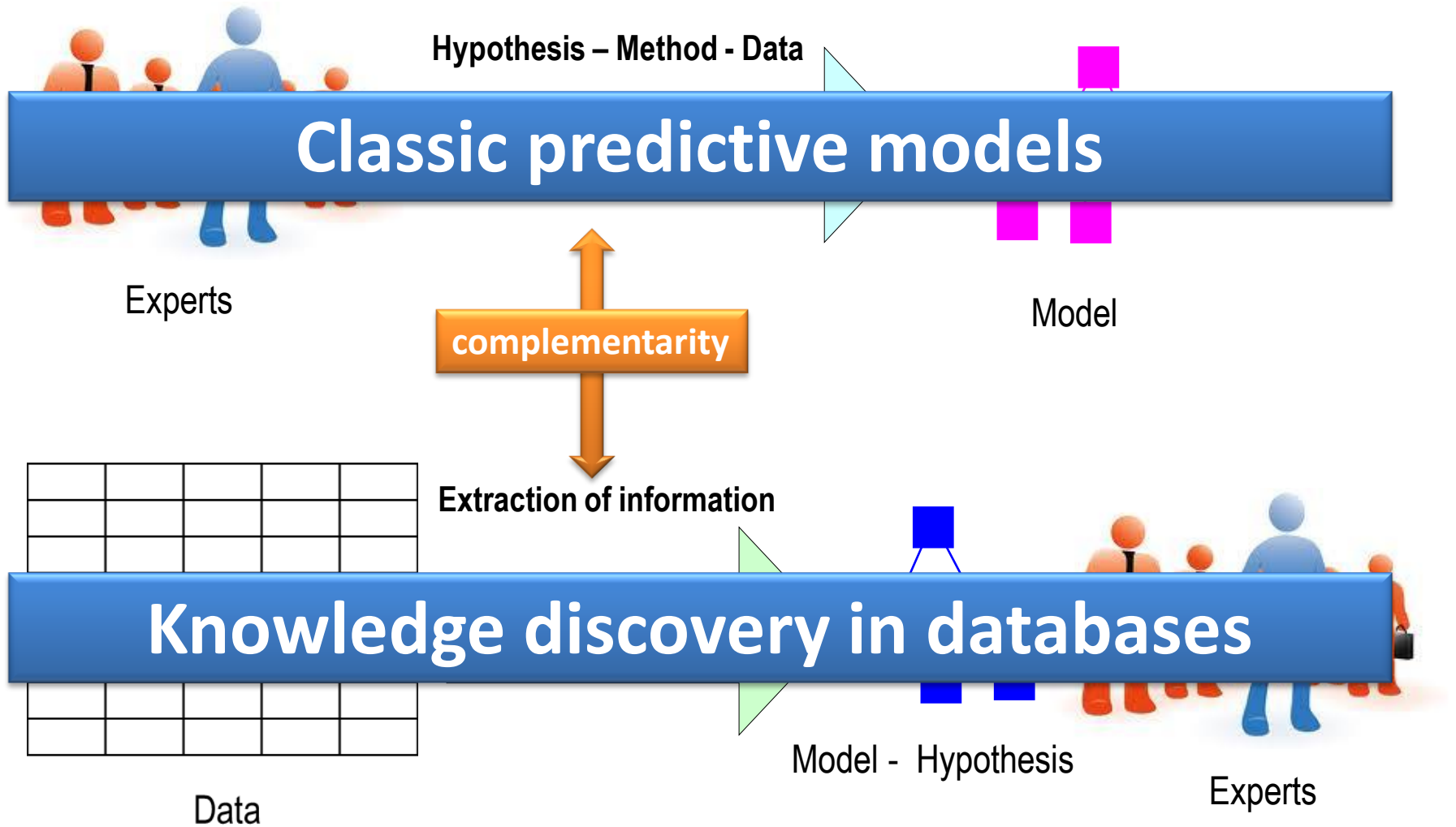


# Inductive machine learning in ecological modeling

*Marko Debeljak*

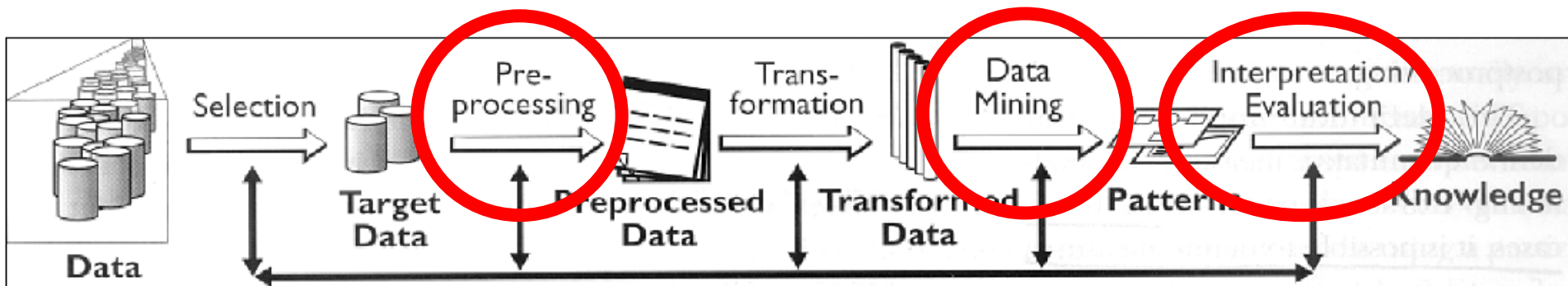
*Jožef Stefan Institute, Slovenia*



## 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**



# 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, then it is called

**predictive or supervised learning.**

## 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.

## PROPOSITIONAL data mining

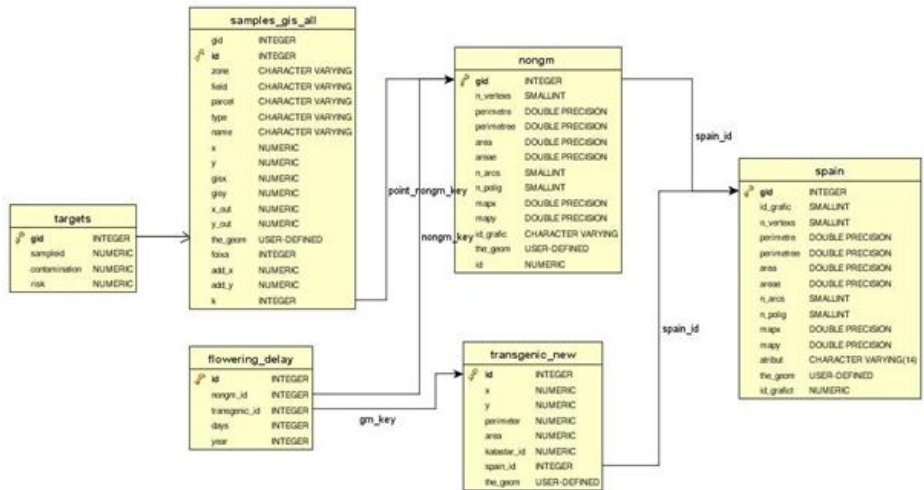
Loss of information due to aggregation

Objects	Properties of objects			
	Distance (m)	Wind direction (°)	Wind speed (m/s)	Out-crossing rate (%)
10	123	3	8	
12	88	4	7	
14	121	6	3	
18	147	2	4	
20	93	1	5	
22	115	3	1	
...	...	...	..	

Data stored in **original** tables or relations.

**No loss** of information, due to aggregation.

## RELATIONAL data mining

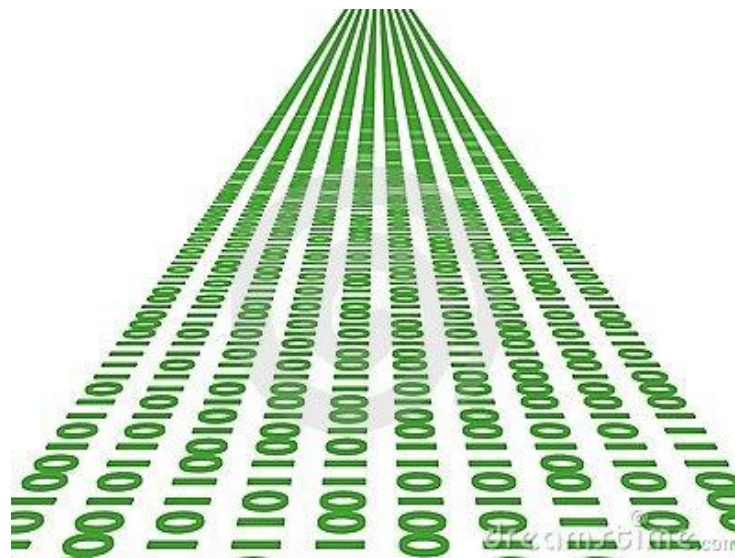


# Data

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

Each example can be propositional or relational.

## DATA STREAM mining



## 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.



## 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:

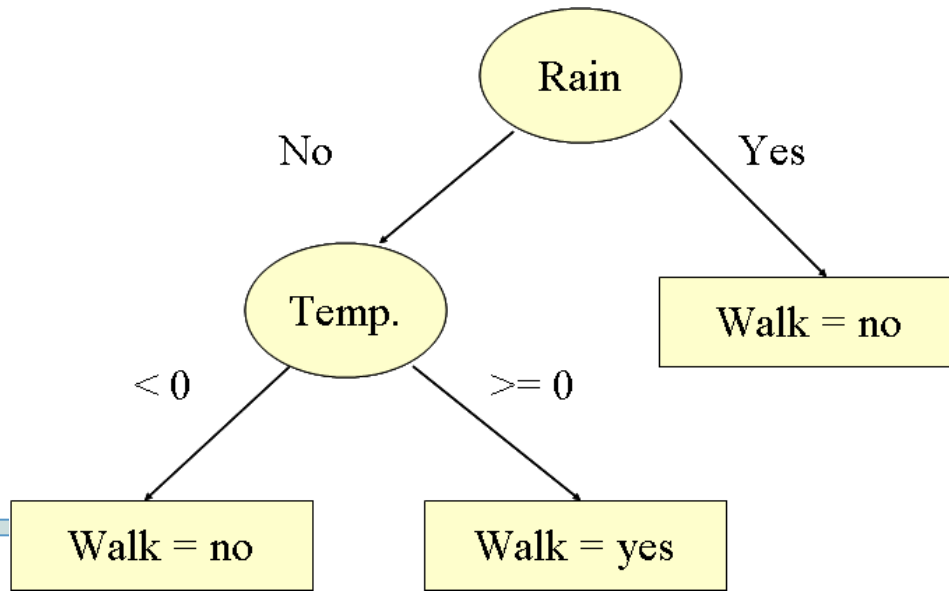
- **Difference equations**
- **Differential equations**

## 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:

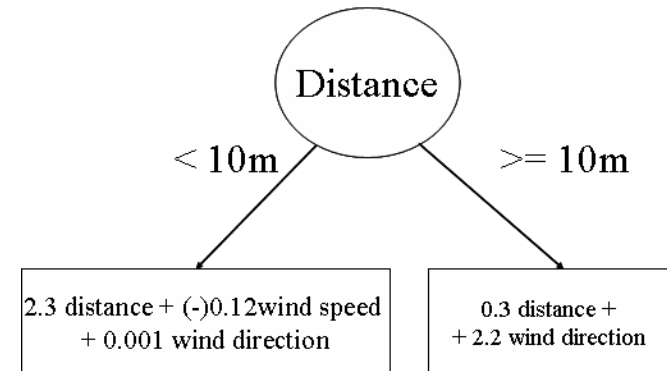
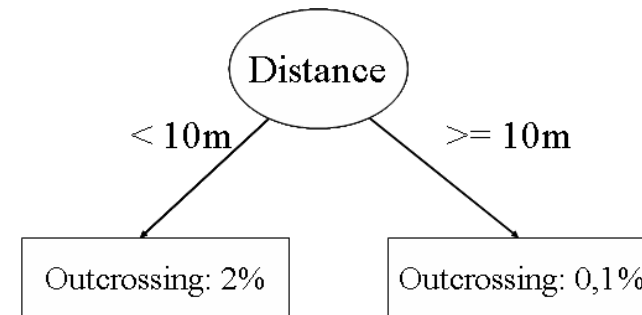
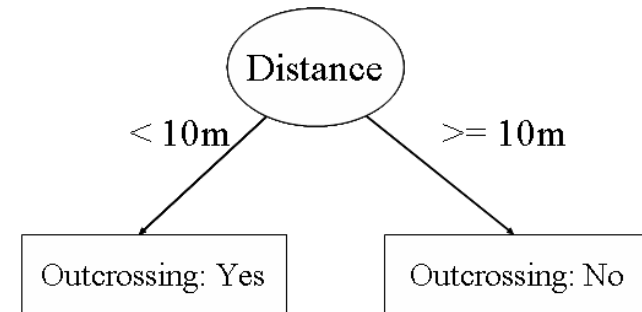
- each internal **node** contains a test on an independent variable,
- each **branch** corresponds to an outcome of the test (critical values of independent variable),
- each **leaf** gives a prediction for the value of the dependent (predicted) variable.



## 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:  

$$\text{out-crossing rate} = 12.3 \text{ distance} - 0.123 \text{ wind speed} + 0.00123 \text{ wind direction}$$



## 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, ...

## 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

**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**

# ***DATA MINING – CASE STUDIES***

## **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)**



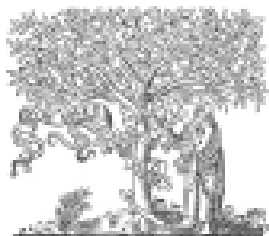
**POPULATION  
DYNAMICS**

**HABITAT  
MODELLING**

**GENE FLOW  
MODELLING**

**RISK MODELLING**

ECOLOGICAL MODELLING 215 (2008) 180–189

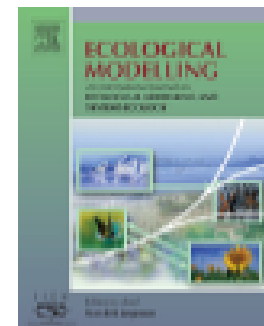


ELSEVIER

available at [www.sciencedirect.com](http://www.sciencedirect.com)



journal homepage: [www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)



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

Jana Laganis<sup>a,\*</sup>, Aleksandar Pečkov<sup>b</sup>, Marko Debeljak<sup>b</sup>

<sup>a</sup> Laboratory for Environmental Research, University of Nova Gorica, Vipavska 13, Nova Gorica, Slovenia

<sup>b</sup> Department of Knowledge Technologies, "Jožef Stefan" Institute, Jamova 39, Ljubljana, Slovenia

Problem: **Prediction of radial increment**

Type of pattern: **Algebraic equation**

Algorithm: **CIPER**

## 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; m<sup>3</sup>/y removed from the stand; *Forestry Unit Lendava*)

## Dataset

### 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:
  - minT>0°C
  - minT<-10°C
  - minT<-4°C
  - minT>25°C
  - maxT>10°C
  - maxT>25°C

- Monthly data + aggregated data (AMJ, MJJ, JJA, MJJA etc.)
- $\Sigma$ : **333** attributes; 35 years

- 52 different combinations of attributes were tested.  
 **$\Sigma$ : 124 models**

Experiment	RRSE	# eq. elements
<b>jnj3_2m</b>	<b>0,7282</b>	<b>6</b>
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

## Model jnj3\_2m:

RadialGrowthIncrement =

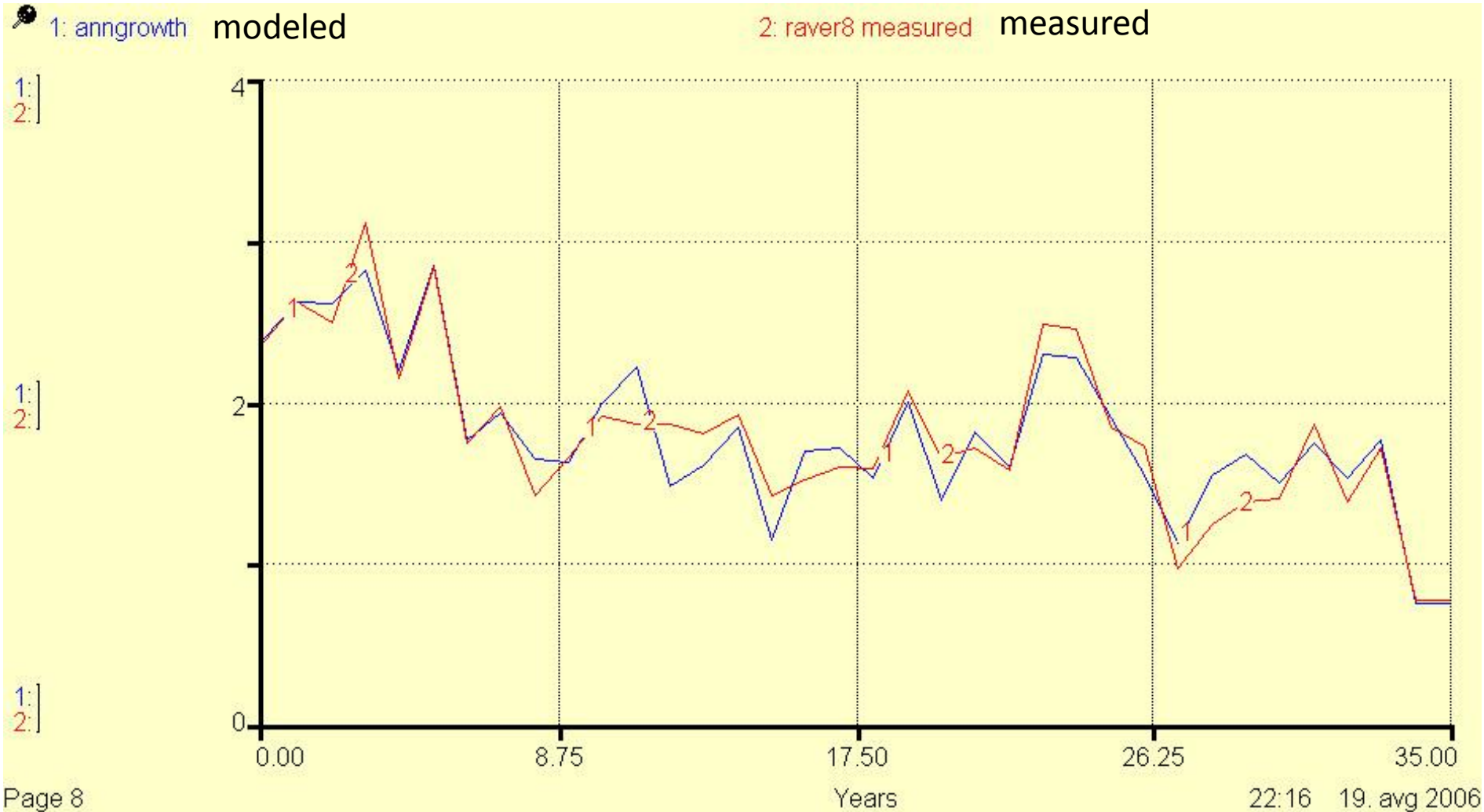
+ -0.0511025526922 minL8-10<sup>1</sup>  
+ -0.0291795197998 maxL8-10<sup>1</sup>  
+ -0.017479975134 t-sun4-7<sup>1</sup>  
+ 0.0346935385853 t-sun8-10<sup>1</sup>  
+ -1.950606536e-05 t-sun8-10<sup>2</sup>  
+ -2.01014710248 d-wf-4-7<sup>1</sup>  
+ 9.35586778387e-05 minL4-7<sup>1</sup> t-sun4-7<sup>1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

Relative Root Squared Error  
= **0.728229824611**

Correlation between average  
measured  
(r-aver8) and modeled increments:  
linear regression:  
**R<sup>2</sup> = 0.8771**

**8 out of 333 attributes**

# Algebraic equations: POPULATION DYNAMICS



Cyper - jnj3 2m measured versus modelled increments

MEDNARODNA PODIPLOMSKA ŠOLA JOŽEFA STEFANA  
JOŽEF STEFAN INTERNATIONAL POSTGRADUATE SCHOOL



Vladimir Kuzmanovski

## **Integration of expert knowledge and predictive learning: Modelling water flows in agriculture**

**Master Thesis**

*Supervisor:* Prof. Dr. Marko Debeljak

*Co-Supervisor:* Prof. Dr. Sašo Džeroski



Problem: **Prediction of drainage water**

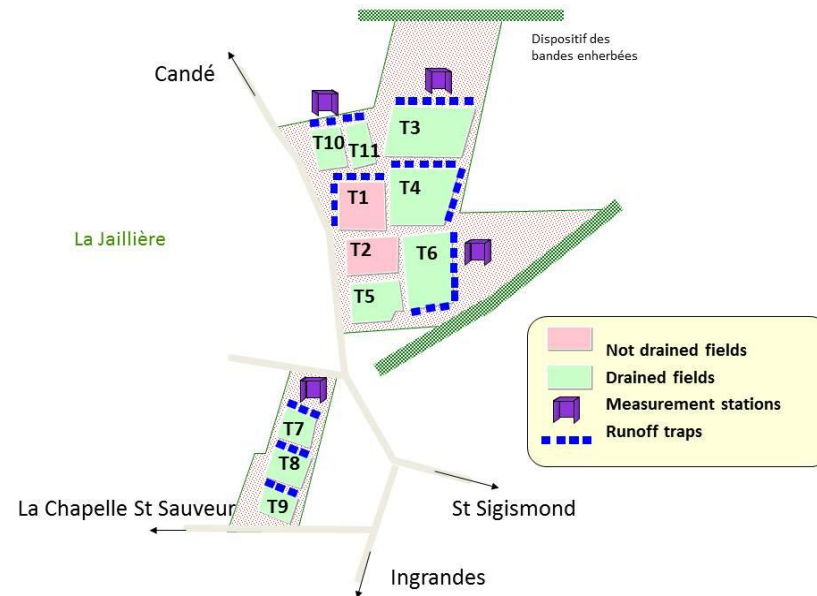
Type of pattern: **Algebraic equation**

Algorithm: **CIPER**

# PCQE Database



- Experimental site La Jaillière
- Western France
- Owned by ARVALIS
  
- Shallow silt clay soils
- 11 fields are observed
- Field size about 0.3 - 1 ha



# PCQE Database *(continued)*



- Agricultural practices
  - Fertilization
  - Irrigation
  - Phytochemical protection
  - Harvesting
  - Tillage
- Slope
- Water flow
  - Drainage
  - Runoff
- 25 campaigns (1987 - 2011)
- Campaign is defined as period starting from 01.09 and finishing on 31.08, following year



# DRAINAGE predictive model - CIPER



- Polynomial equations induced on data for a whole campaign - CIPER algorithm
- Evaluation
  - “Leave one out” approach

Fields	Test field	Std. Dev.	RMSE	RRSE	Corr. coeff. (r)
All	T3	3.187	2.1119	66.26 %	0.7855
All	T4	3.188	1.7220	54.01 %	0.8273
All	T5	3.163	2.2478	71.06 %	0.7467
All	T6	3.096	2.1784	70.36 %	0.8096
All	T7	3.229	1.3286	41.15 %	0.7812
All	T8	3.210	1.3813	43.03 %	0.7839
All	T9	3.210	1.5927	49.62 %	0.7434
All	T10	3.130	1.6672	53.26 %	0.7766
All	T11	3.108	1.6274	52.36 %	0.7841
T3	T6	3.056	2.1251	69.54 %	0.8125
T6	T3	3.600	2.0961	58.22 %	0.7745



# Predictive models



Model (All/T4)

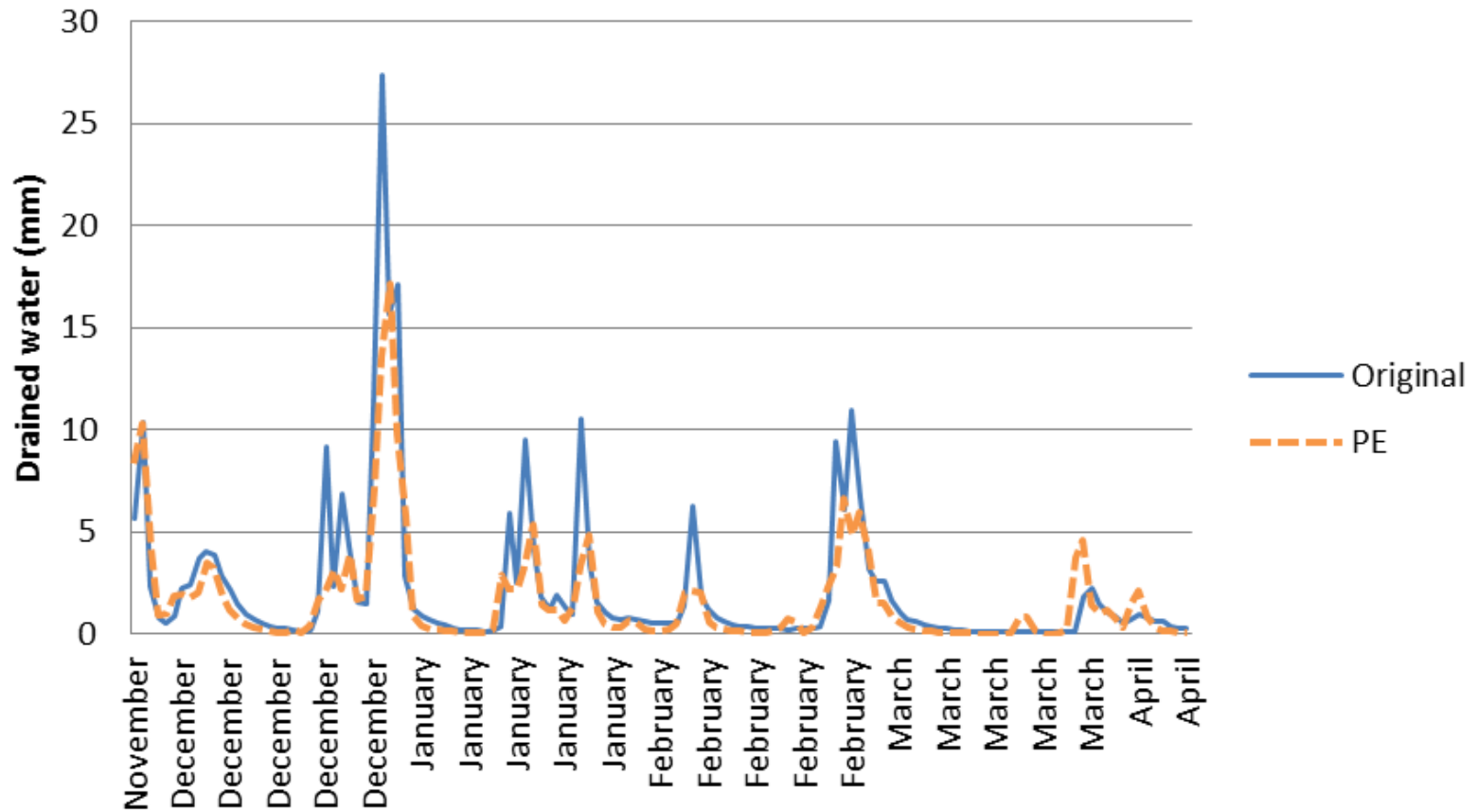
$$\begin{aligned} \text{Drainage} = & 0.0196445 * \text{RainfallA1} * \text{Temp} \\ & + 0.33246 * \text{DrainageN1} \\ & + 0.000662861 * \text{CDCoef}^2 * \text{RainfallA1}^2 * \text{Slope} \\ & + 0.00000107253 * \text{Runoff} * \text{DrainageN1} * \text{Temp}^2 * \text{RainfallA1}^2 \\ & - 0.00115983 * \text{Runoff}^2 * \text{DrainageN1} * \text{Slope}^3 \\ & - 0.00114057 * \text{Temp}^2 * \text{RainfallA1} \\ & + 0.00153725 * \text{RainfallA1}^2 \\ & + 1.63563 * \text{Runoff} * \text{Slope} \\ & - 1.90622 * \text{Runoff} \\ & - 0.0231748 * \text{Slope}^2 * \text{RainfallA1} \\ & + 0.0654042 * \text{RainfallA1} * \text{Slope} \\ & - 0.00755737 * \text{Slope}^3 * \text{CDCoef}^3 * \text{Runoff}^2 \\ & + 0.0675951 * \text{Slope} \\ & - 0.146702 \end{aligned}$$



# Predictive models



(2009/2010)



Ecological Modelling 220 (2009) 1063–1072

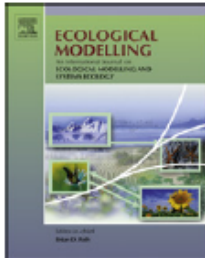


ELSEVIER

Contents lists available at ScienceDirect

## Ecological Modelling

journal homepage: [www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)



## Modelling the outcrossing between genetically modified and conventional maize with equation discovery

Aneta Ivanovska<sup>a,\*</sup>, Ljupčo Todorovski<sup>b</sup>, Marko Debeljak<sup>a</sup>, Sašo Džeroski<sup>a</sup>

<sup>a</sup> Jožef Stefan Institute, Department of Knowledge Technologies, Jamova cesta 39, SI-1000 Ljubljana, Slovenia

<sup>b</sup> University of Ljubljana, Faculty of Administration, Gosarjeva 5, SI-1000 Ljubljana, Slovenia

Problem: **Prediction of gene flow**

Type of pattern: **Constraint algebraic equation**





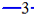
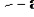
Algorithm: **LAgrange**



## Experiment design: Federal Biological Research Centre, BBA, D

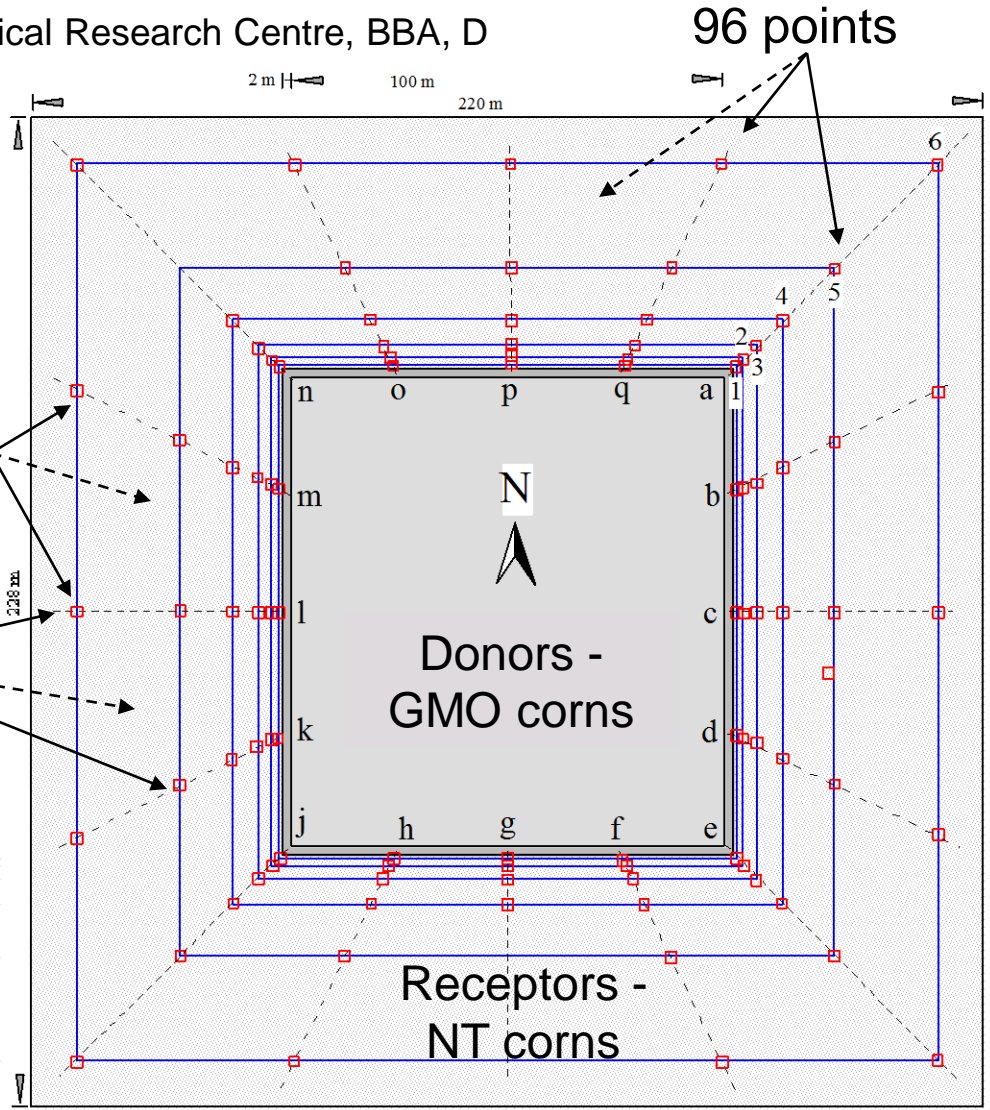
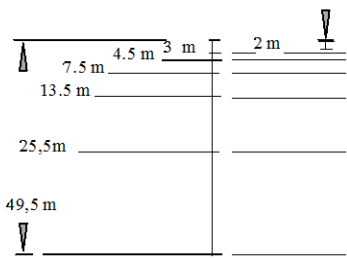
Temperature  
Relative humidity  
Wind velocity  
Wind direction

Field design 2000

-  transgenic field / donor
-  non-transgenic field / recipient
-  access paths
-  sampling point
-  system of coordinates
-  for the sampling points

60 cobs – 2500 kernels

% of outcrossing



## Table 5

The context free grammar used by Lagrange to model the outcrossing between GM and non-GM maize. The grammar specifies the space of possible equations to be considered by Lagrange.

Outcrossing  $\rightarrow$  const  $\times$  (Distance Influence <sup>$\alpha$</sup> )  $\times$  (Wind Influence <sup>$\beta$</sup> )

Distance Influence  $\rightarrow$  1

Distance Influence  $\rightarrow F$

Distance Influence  $\rightarrow F$

$F \rightarrow e^{-\text{Distance}}$

$F \rightarrow 1/\text{Distance}$

$F \rightarrow 1/\text{Distance}^2$

$F \rightarrow \text{Distance}^{-\gamma}; (0 \leq \gamma \leq 1000)$

Distance  $\rightarrow$  variable\_minDistance

Distance  $\rightarrow$  variable\_distanceCenter

Wind Influence  $\rightarrow$  1

Wind Influence  $\rightarrow$  PWind

PWind  $\rightarrow$  (PWind)  $\times$  Wind + const|const

Wind  $\rightarrow$  variable\_appropriateWindProc

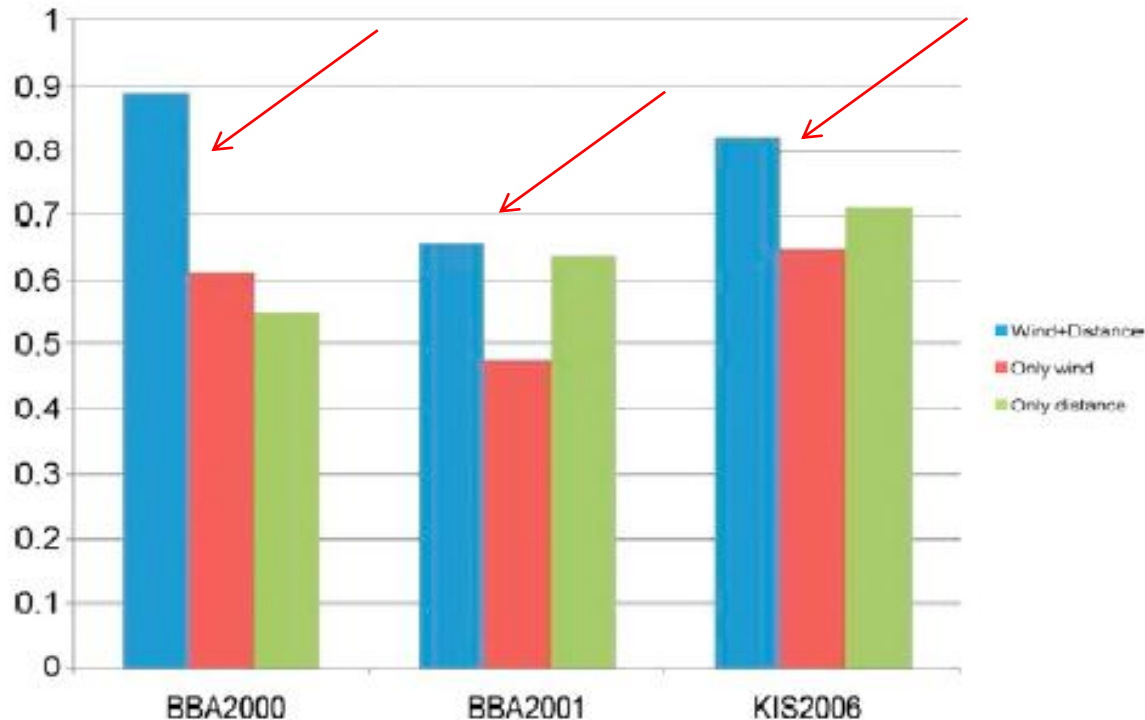
Wind  $\rightarrow$  variable\_windTunnelLength

# Algebraic equations: GENE FLOW

**Table 8**

Correlation coefficients ( $r$ ), relative mean squared error (reMSEs) and best equations of the experiments carried out on BBA2000, BBA2001, KIS2006, all BBA, and all three datasets.

	Correlation coefficient (reMSE)	Best equation
BBA2000	0.89 (0.50)	Outcrossing = $\frac{0.02}{\text{minDistance}^{1.8}} \times [0.007 \times \text{windTunnelLength}^2 \times \text{appropriateWindProc} + 602.93]$
BBA2001	0.68 (0.90)	Outcrossing = $\frac{0.01}{\text{distanceCenter} \times \text{minDistance}^2} \times [\text{windTunnelLength}^3 + \text{windTunnelLength}^2 + \text{windTunnelLength} + 1]$
KIS2006	0.83 (0.33)	Outcrossing = $\frac{531.12}{\text{distanceCenter} \times \text{minDistance}}$
BBA2000 + 2001	0.86 (0.48)	Outcrossing = $\frac{0.01}{\text{distanceCenter} \times \text{minDistance}^2} \times [\text{appropriateWindProc} \times \text{windTunnelLength}^2 + \text{windTunnelLength}^2 + \text{windTunnelLength} + 1]$
ALL	0.64 (1.52)	Outcrossing = $\frac{0.01}{\text{distanceCenter} \times \text{minDistance}^{0.7}} \times [\text{appropriateWindProc}^2 + \text{appropriateWindProc} + 1]$



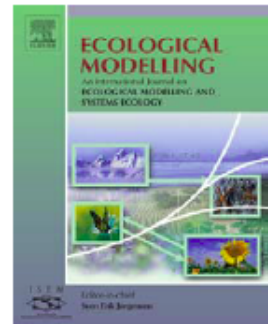
ECOLOGICAL MODELLING 194 (2006) 37–48



available at [www.sciencedirect.com](http://www.sciencedirect.com)



journal homepage: [www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)



## Automated modelling of a food web in lake Bled using measured data and a library of domain knowledge

Nataša Atanasova<sup>a,\*</sup>, Ljupčo Todorovski<sup>b</sup>, Sašo Džeroski<sup>b</sup>, Špela Rekar Remec<sup>c</sup>,  
Friedrich Recknagel<sup>d</sup>, Boris Kompare<sup>a</sup>

<sup>a</sup> Faculty of Civil and Geodetic Engineering, University of Ljubljana, Slovenia

<sup>b</sup> Jožef Stefan Institute, Slovenia

<sup>c</sup> Environmental Agency of the Republic of Slovenia, Slovenia

<sup>d</sup> University of Adelaide, Australia

Problem: **Time dependent ecosystem processes**

Type of pattern: **Differential equations**

Algorithm: **LAgранже**

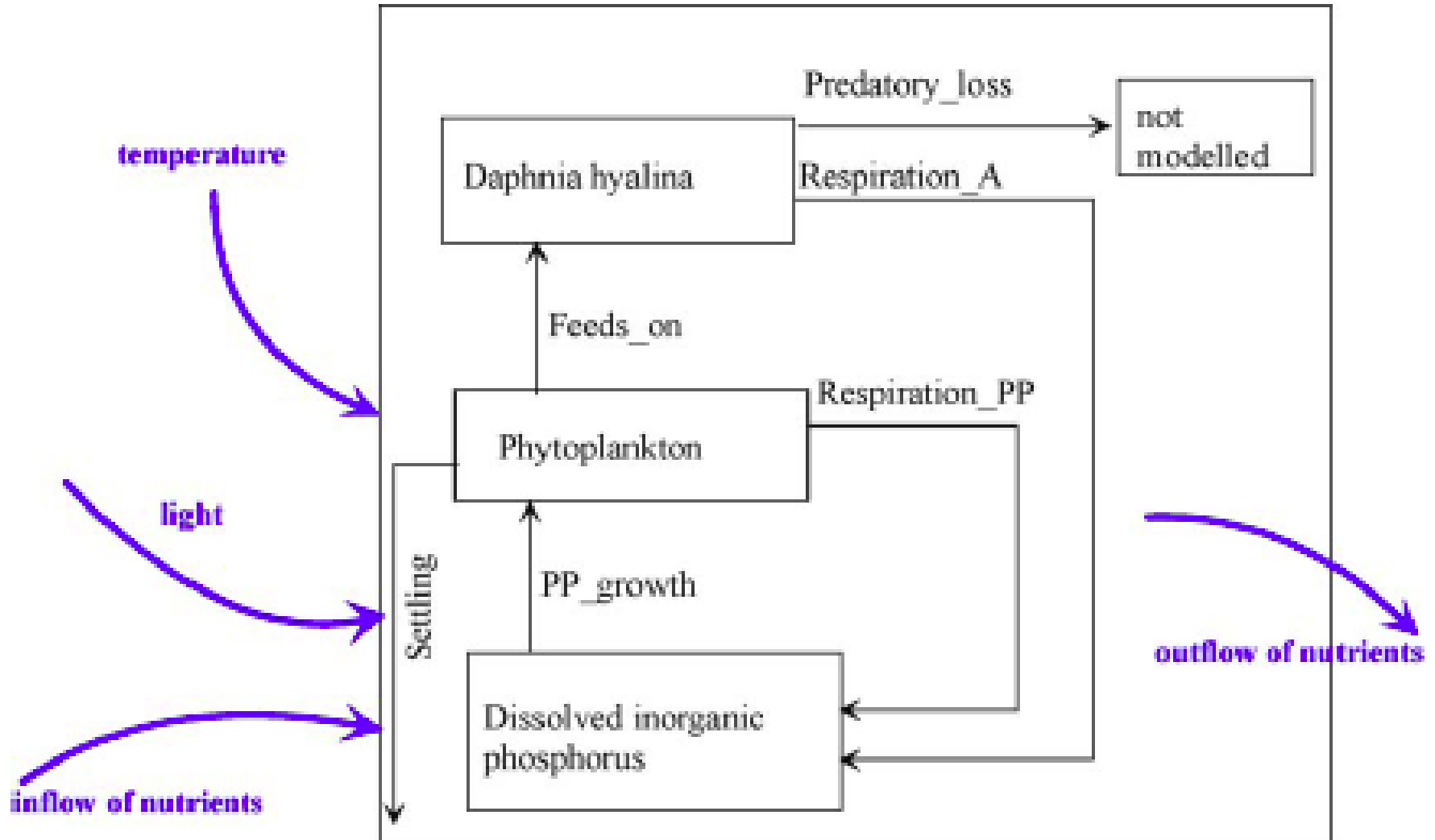


Fig. 2 – Simple conceptual model for lake Bled.

Data: 1995 to 2002

**Table 2 – Measured data (variables) in lake Bled used for model induction**

Variable name	Description	Frequency
q_krivica (m <sup>3</sup> /day)	Inflow to the lake	Daily
q_misca (m <sup>3</sup> /day)	Inflow to the lake	Daily
q_radovna (m <sup>3</sup> /day)	Inflow to the lake	Daily
q_jezemica (m <sup>3</sup> /day)	Outflow (at surface)	Daily
q_natega (m <sup>3</sup> /day)	Outflow (syphon)	Daily
ps_krivica, ps_misca, ps_radovna (mg/l)	Nutrient (orthophosphate) concentration in the inflows	Monthly
temp (°C)	Water temperature of the streams and lake	Monthly
light (J/(cm <sup>2</sup> day))	Calculated underwater light	Monthly
ps, no, silica (mg/l)	Inorganic nutrients' concentration in the lake (ps is soluble phosphorus and NO is nitrate)	Monthly
phyto (mgDW/l)	Phytoplankton biomass concentration in the lake	Monthly
daph (No. ind/ml or mgDW/l (see text))	Zooplankton ( <i>Daphnia hyalina</i> ) biomass concentration in the lake	Monthly

## Phosphorus

$$\begin{aligned}
 \frac{d(ps)}{dt} = & \underbrace{ps\_krivica \cdot \frac{q\_krivica}{7 \cdot 10^6} + ps\_misca \cdot \frac{q\_misca}{7 \cdot 10^6}}_{\text{water in-flow}} \\
 & + ps\_radovna \cdot \frac{q\_radovna}{7 \cdot 10^6} - ps \cdot \frac{q\_jezemica}{7 \cdot 10^6} \\
 & - ps \cdot \frac{q\_natega}{7 \cdot 10^6} + \underbrace{0.0022 \cdot phyto^2 \cdot 0.072 \cdot \frac{temp - 2.7}{20.4 - 2.7}}_{\text{respiration}} \\
 & + 0.07 \cdot daph \cdot 0.0026 \cdot \frac{temp}{12.3} - \underbrace{0.0023 \cdot phyto \cdot 0.21}_{\text{growth}} \\
 & \cdot \frac{ps}{ps + 0.00042} \cdot \frac{temp}{16.7} \cdot \frac{light}{170} \cdot e^{\left(1 - \frac{light}{170}\right)} \quad (10)
 \end{aligned}$$



## Phytoplankton

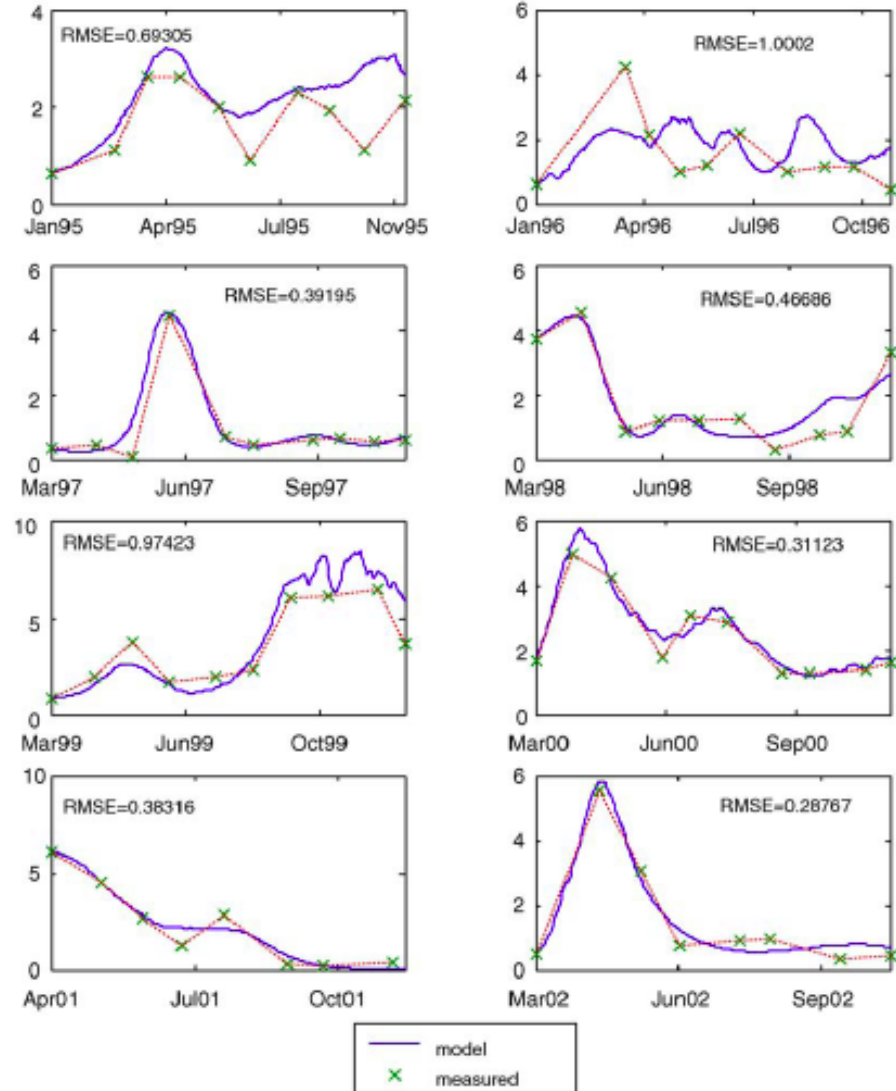
$$\frac{d(\text{phyto})}{dt} = \text{phyto} \cdot 0.21 \cdot \frac{\text{ps}}{\text{ps} + 0.00042} \cdot \frac{\text{temp}}{16.7} \cdot \frac{\text{light}}{170} \cdot e^{(1 - \frac{\text{light}}{170})} - \text{phyto}^2 \cdot 0.072 \cdot \frac{\text{temp} - 2.7}{19.7 - 2} - \text{phyto} \cdot \frac{0.5}{10} \cdot \frac{\text{temp} - 2}{18 - 4} - \text{daph} \cdot 0.5 \cdot \frac{\text{temp} - 2.6}{18 - 4} \cdot (1 - \exp(-0.58 \cdot \text{phyto})) + 0.56 \cdot \text{phyto} \quad (13)$$

growth

respiration

sedimentation

grazing



## Zooplankton

$$\frac{d(\text{daph})}{dt} = 0.14 \cdot \text{daph} \cdot 0.5 \cdot \frac{\text{temp} - 2.6}{18 - 4} \cdot (1 - \exp(-0.58 \cdot \text{phyto})) \cdot 0.56 \cdot \text{phyto}$$

Feeds on  
phytoplankton

$$- \text{daph} \cdot 0.026 \cdot \frac{\text{temp}}{12.3} - 0.01 \cdot \frac{\text{daph}^2}{0.001 + \text{daph}}$$

respiration

mortality

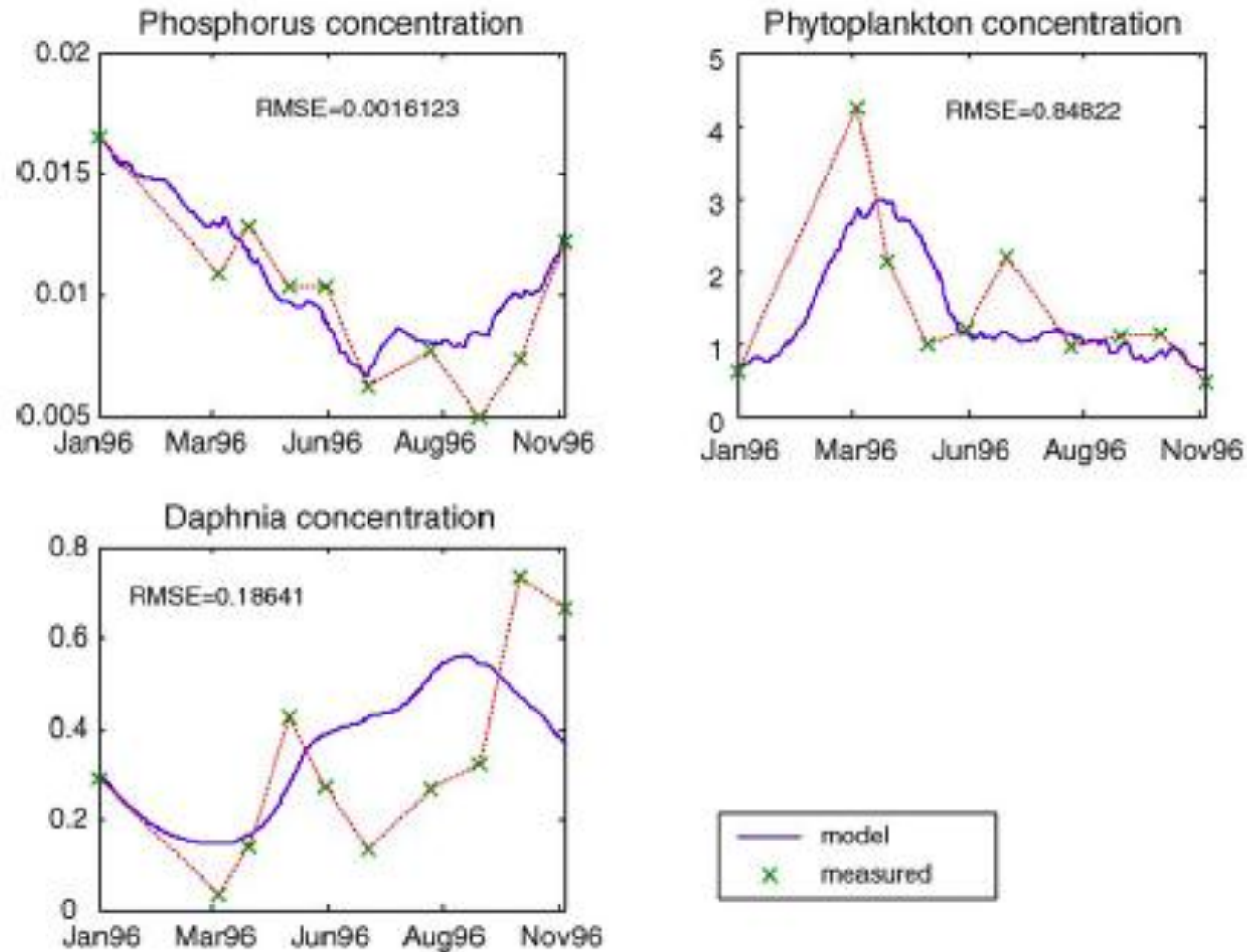


Fig. 6 – Performance of the food web model for phosphorus, phytoplankton and daphnia.

Algebraic - CIPER

Constraint algebraic- LAgamge

Differential- Lagramge

**SUITABLE for predictions,  
UNSUITABLE for interpretation**



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Ecological Modelling 170 (2003) 453–469

ECOLOGICAL  
MODELLING

[www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)

## Modeling the brown bear population in Slovenia A tool in the conservation management of a threatened species

Klemen Jerina<sup>a,\*</sup>, Marko Debeljak<sup>b</sup>, Sašo Džeroski<sup>c</sup>, Andrej Kobler<sup>d</sup>, Miha Adamič<sup>a</sup>

<sup>a</sup> Department of Forestry, Biotechnical Faculty, University of Ljubljana, Večna pot 83, 1000 Ljubljana, Slovenia

<sup>b</sup> Nova Gorica Polytechnic, School of Environmental Sciences, Vipavska 13, 5000 Nova Gorica, Slovenia

<sup>c</sup> Jozef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia

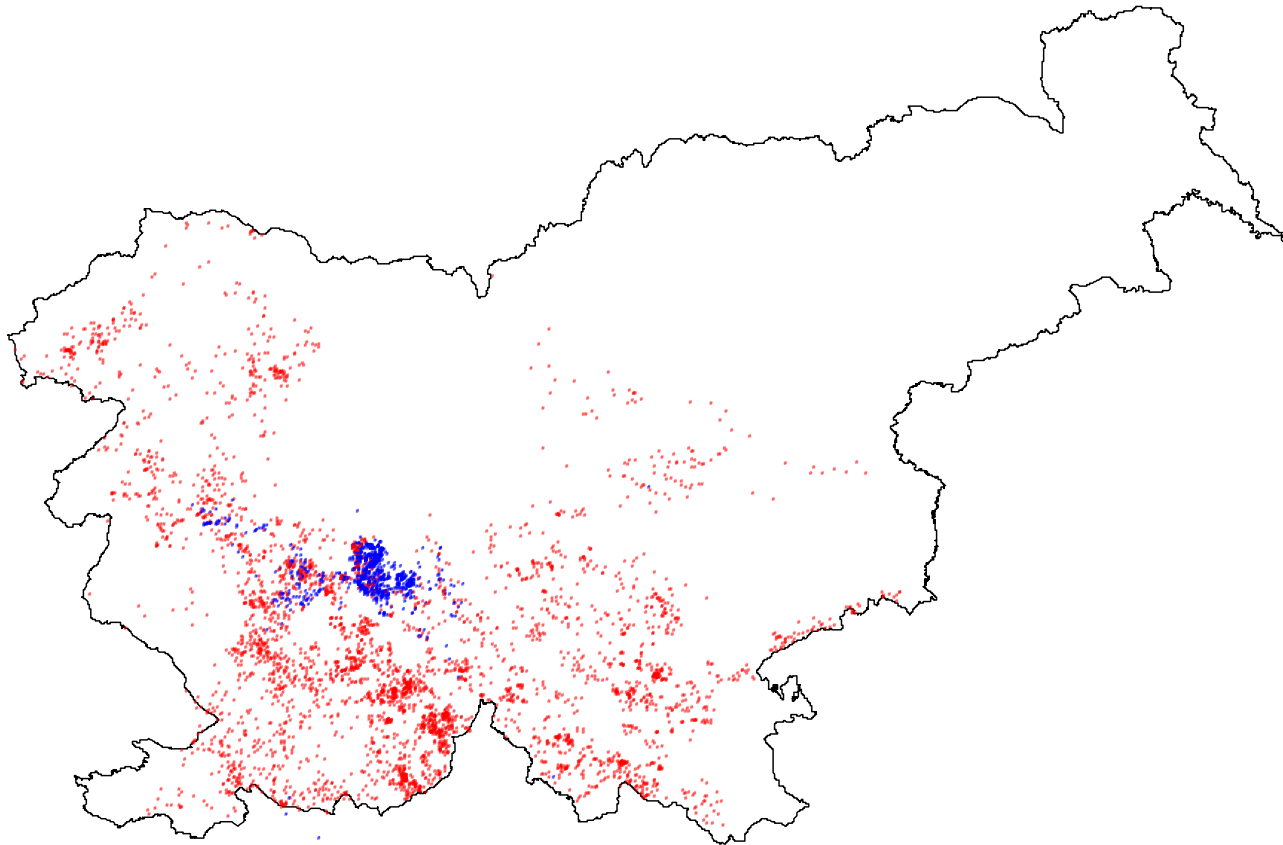
<sup>d</sup> Slovenian Forestry Institute, Večna pot 2, 1000 Ljubljana, Slovenia

Problem: **Classification of habitats** (suitable, unsuitable)

Type of pattern: **Classification decision trees**

Algorithm: **J4.8**

## Observed locations of BBs



## 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*

## Propositional dataset

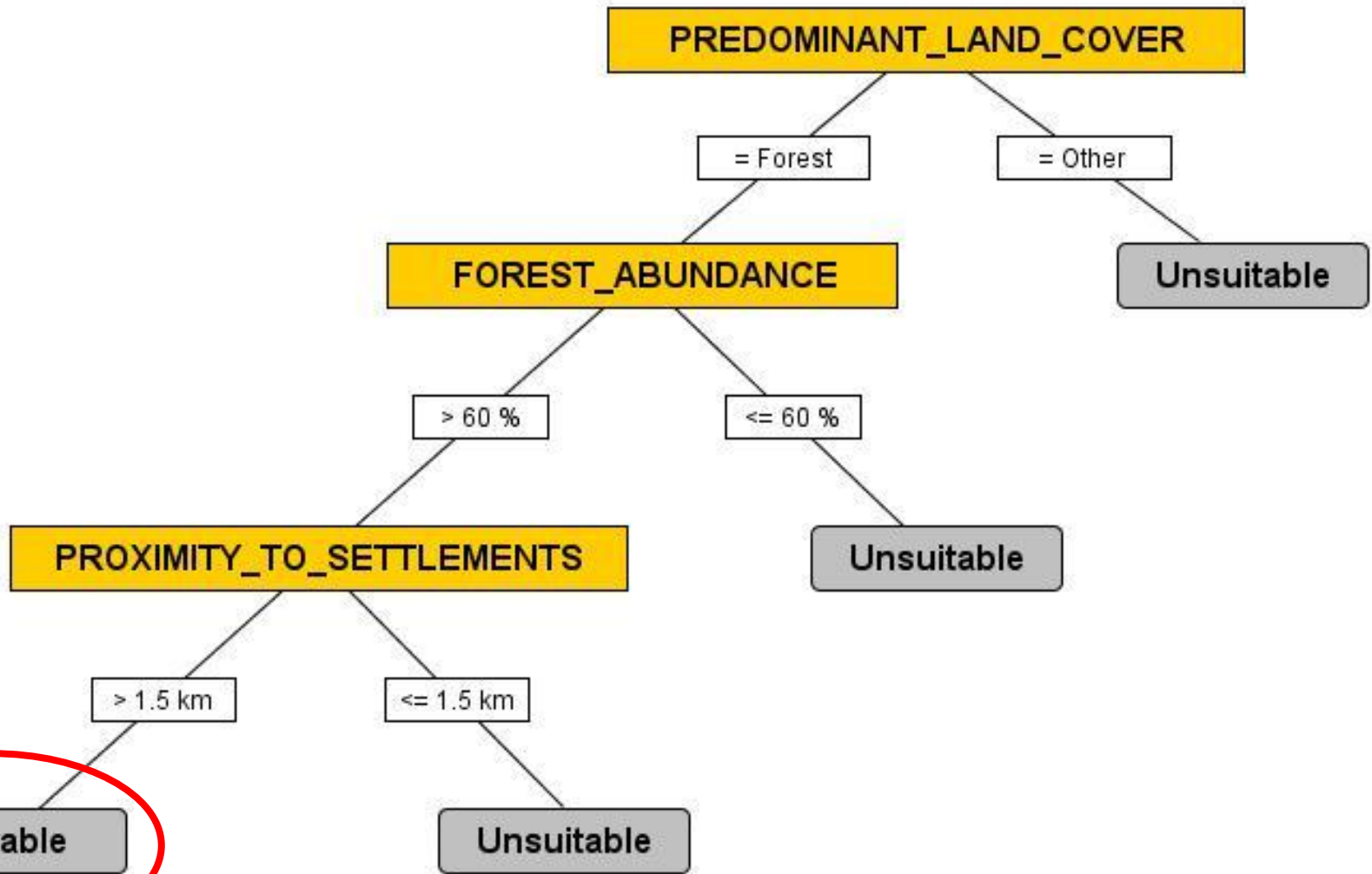
1,73,26,0,0,1,88,0,2,70,7,20,1,0,0,1,0,60,0,0,0,0,0,2,0,0,0,4123,0,0,0,0,63,211,11,11,11,83,213,213,0,0,4155.  
1,62,37,0,0,2,88,0,2,70,7,20,1,0,0,1,0,60,0,0,0,0,0,2,1,53,0,3640,0,0,0,-1347,63,211,11,11,11,83,213,213,11,89,3858.  
2,0,99,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,6,82,0,10404,0,2074,-309,48,0,0,11,11,11,83,83,83,0,20,3862.  
2,0,100,0,0,1,76,0,16,71,0,12,0,0,0,0,0,0,0,0,0,0,0,1,6,82,0,7500,0,1661,-319,-942,0,0,11,11,11,0,0,0,0,20,4088.  
1,8,91,0,0,1,52,0,59,41,0,0,0,0,0,0,0,0,4,0,0,0,0,5,1,6,82,0,6500,0,1505,-166,879,9,57,11,11,11,281,281,281,0,20,3199.  
4,3,0,86,9,0,75,0,33,67,0,0,0,0,0,0,1,2,0,0,0,0,0,1,2,54,0,0,0,465,-66,-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,-264,34,220,11,41,41,151,141,112,60,619,3897.  
1,100,0,0,0,3,52,0,86,6,3,5,9,6,7,38,40,0,0,0,0,0,0,0,1,17,64,0,8062,0,932,-603,-71,100,337,11,41,41,171,232,202,4,24,3732.

.....  
.....  
.....

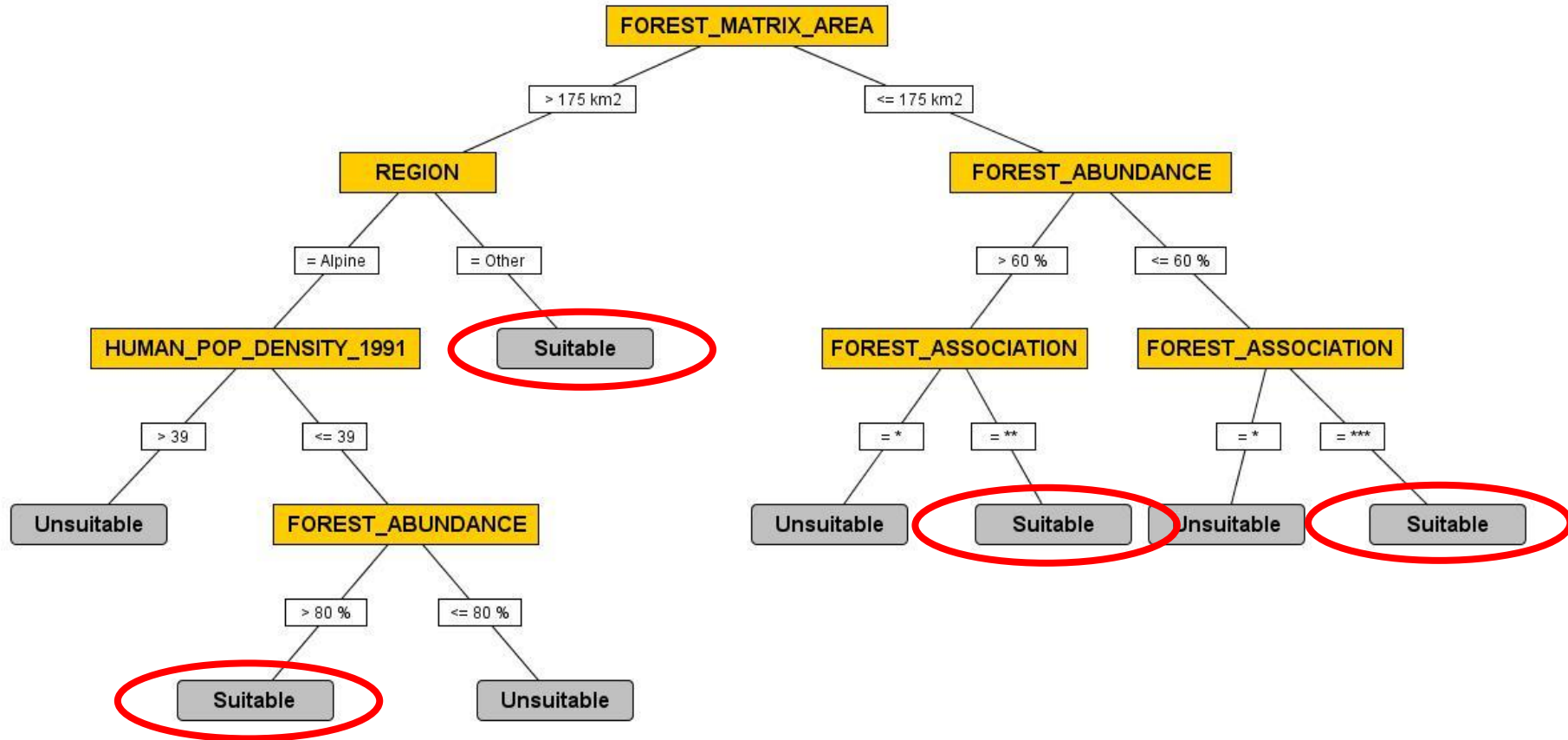
Present: 1

Absent: 0

# The model for maximal habitat



# The model for optimal habitat



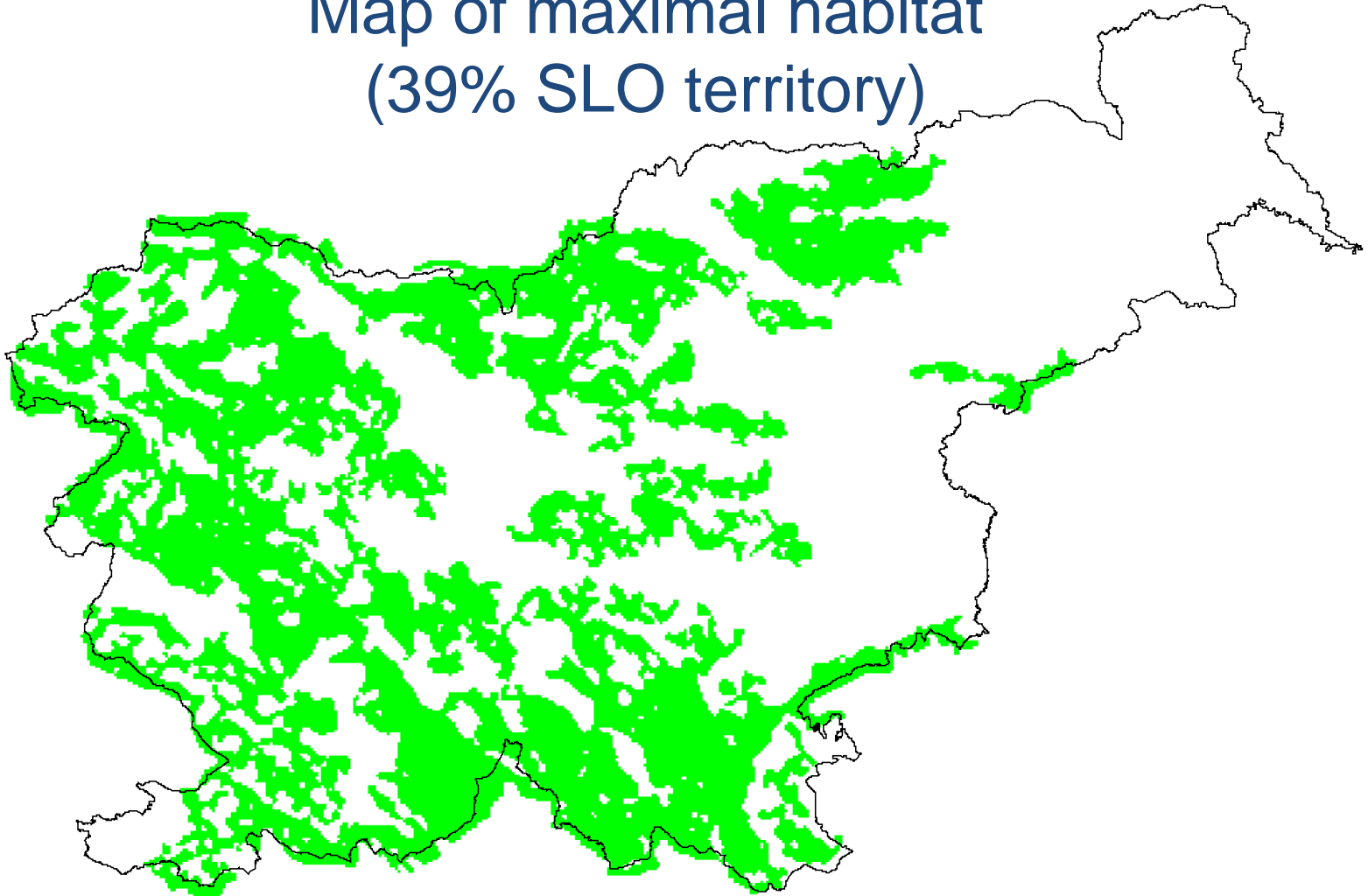
\* QUERICO ROBORI - CARPINETUM, CARICI ELATE - ALNETUM GLUTINOSAE, CARICI BRIZOIDI - ALNETUM GLUTINOSAE, ALNETUM GLUTINOSO - INCANAE, ALNETUM INCANAE, SALICI - POPULETUM, SALICETUM GR., QUERCO - CARPINETUM VAR. HACQUETIA, QUERCO - CARPINETUM VAR. LUZULA, etc ...

\*\* MELAMPYRO VULGATI - QUERCETUM, ADENOSTYLO - FAGETUM, FAGETUM SUBALPINUM, QUERCO - FAGETUM, ABIETI - FAGETUM DINARICUM

\*\*\* ADENOSTYLO - FAGETUM, FAGETUM SUBALPINUM

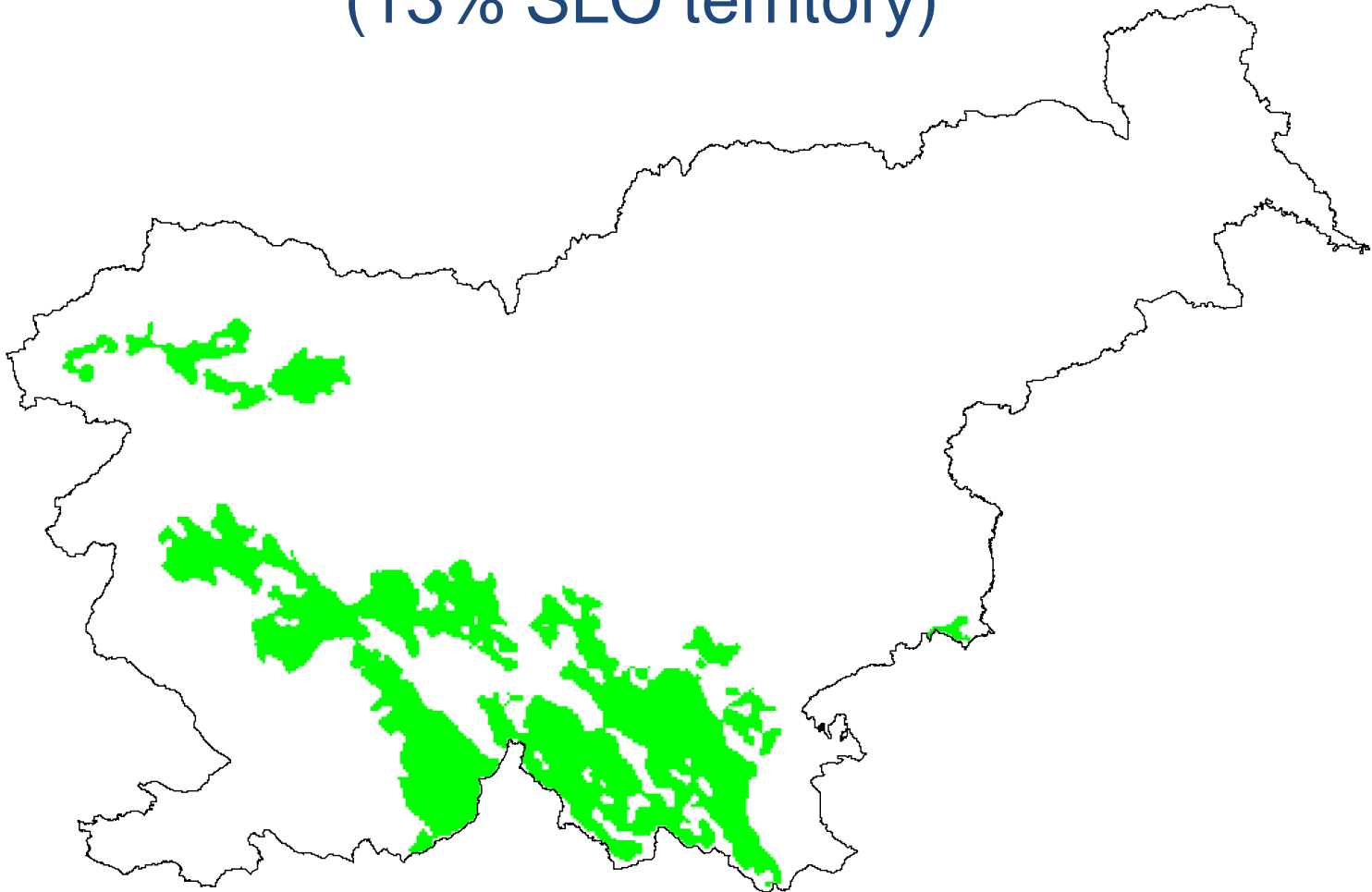
# Decision trees: *HABITAT MODELS*

Map of maximal habitat  
(39% SLO territory)



# Decision trees: *HABITAT MODELS*

Map of optimal habitat  
(13% SLO territory)



## ***Predictive models of forest development in Slovenia***

**Problem A: Prediction of time dynamics of growing stock**

Type of pattern: **Multi-target regression tree**

Algorithm: **CLUS**

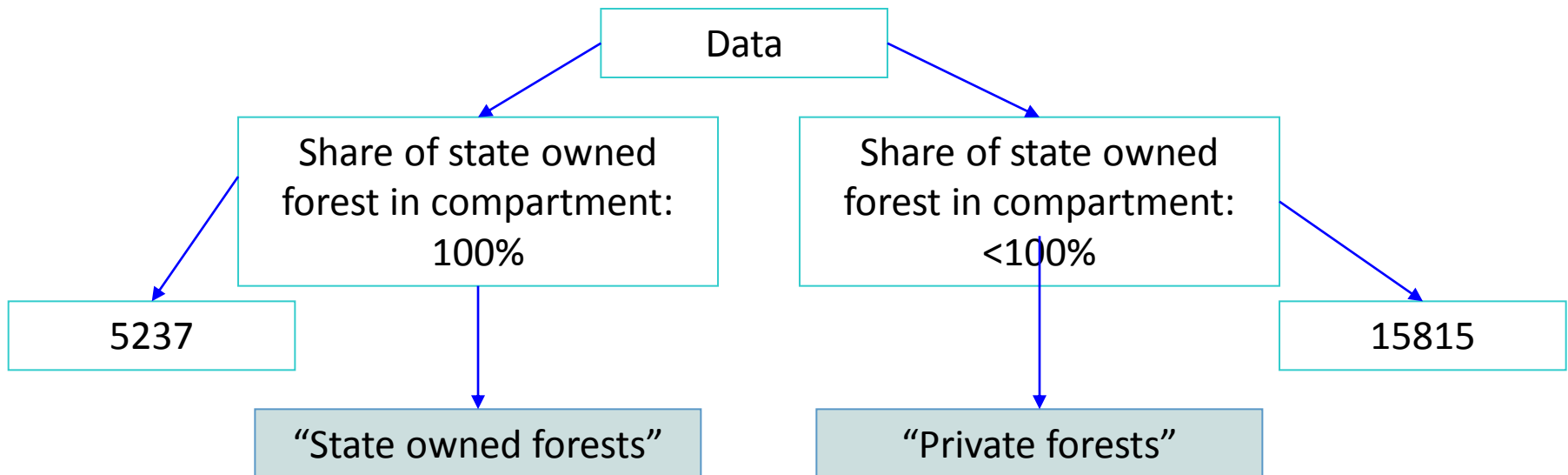
**Problem B: Prediction of growing stock**

Type of pattern: **Model tree**

Algorithm: **M5'**

## Database *Silva* 1970-2008:

- data from 1970, 1980, 1990, 2000 and 2008,
- data unit is a **permanent compartment**
- **21052 permanent compartments**





EACH  
compartment  
described by:

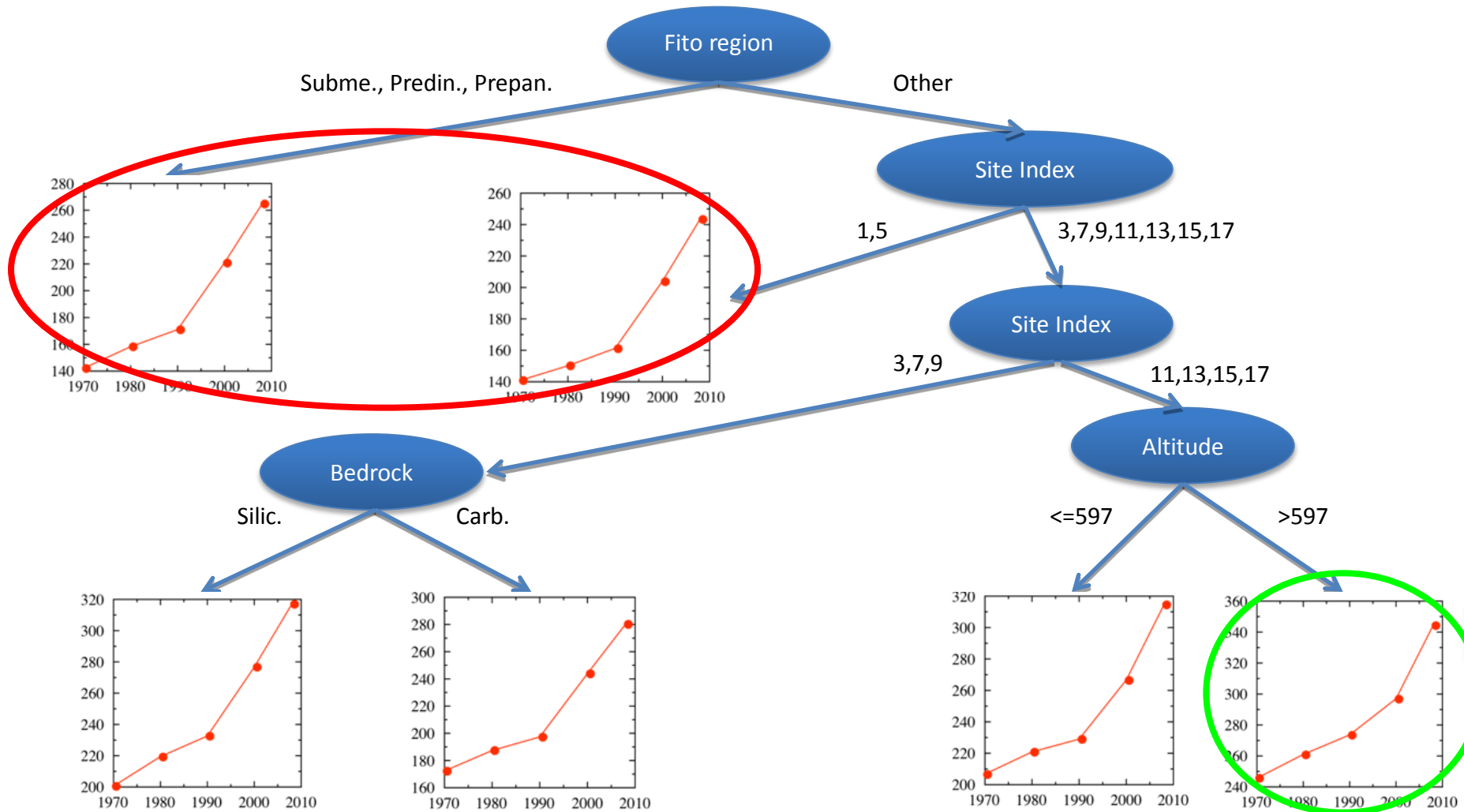
Environmental  
attributes

Forest stands  
attributes

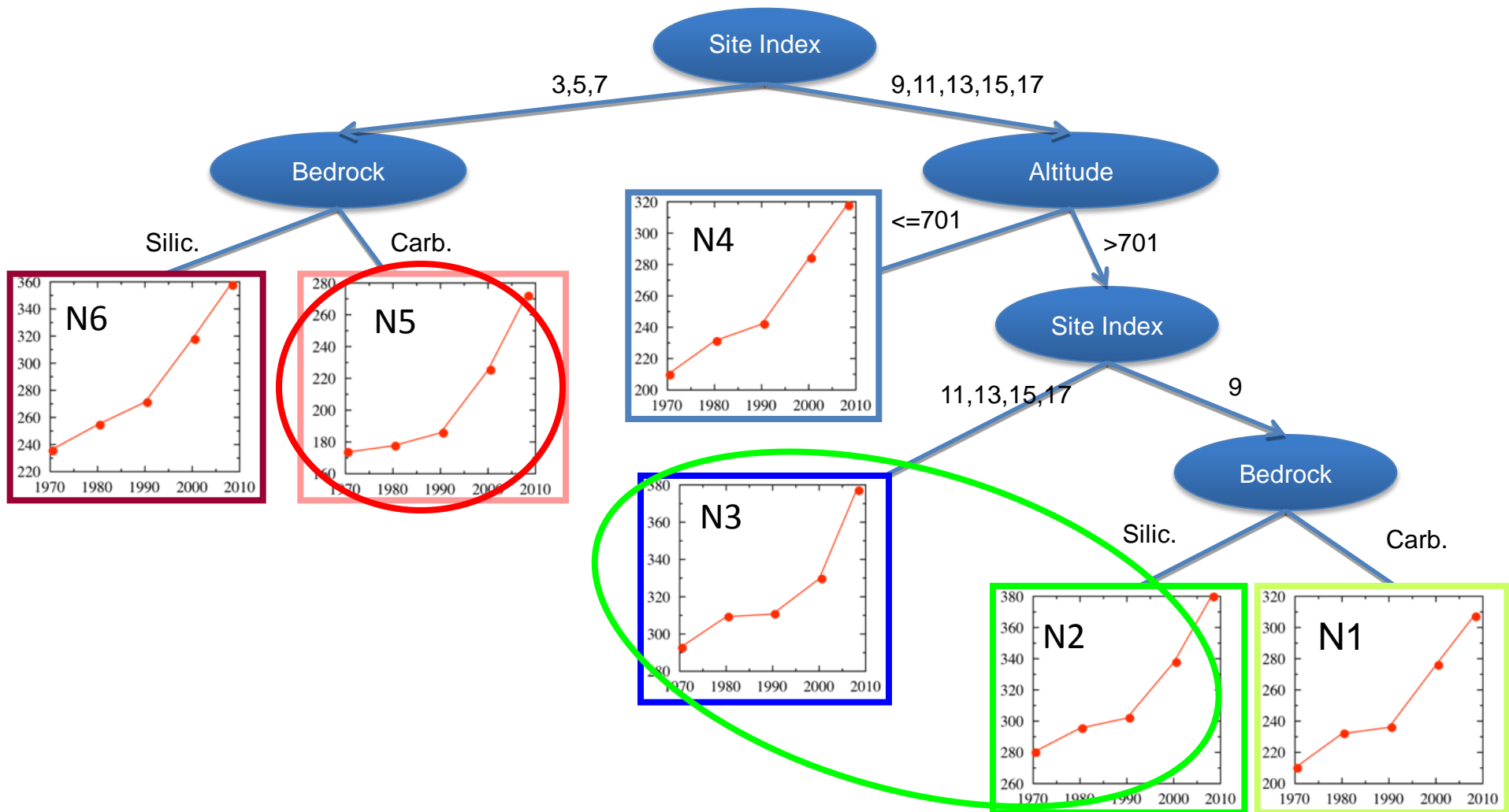
Management  
activities  
attributes

NMV\_M : elevation (m)  
NAKLON\_M: slope (%)  
EXP: aspect (ranks 0-8)  
MAT\_POFLAG: bedrock (1-carbonate, 0-silicate)  
RK: site quality (1 (the worst)- 17 (the best))  
FITOGEO\_IM: fito-geographical regions (Pre-dinaric, Alpine, Dinaric, Sub-mediterranean, Pre-panonic)  
LZSKU70: growing stock 1970 (m<sup>3</sup>/ha)  
LZSKU80: growing stock 1980 (m<sup>3</sup>/ha)  
LZSKU90: growing stock 1990 (m<sup>3</sup>/ha)  
LZSKU00: growing stock 2000 (m<sup>3</sup>/ha)  
LZSKU08: growing stock 2008 (m<sup>3</sup>/ha)  
E70: cut (etat) 1970 (m<sup>3</sup>/ha)  
E80: cut (etat) 1980 (m<sup>3</sup>/ha)  
E90: cut (etat) 1990 (m<sup>3</sup>/ha)  
E00: cut (etat) 2000 (m<sup>3</sup>/ha)

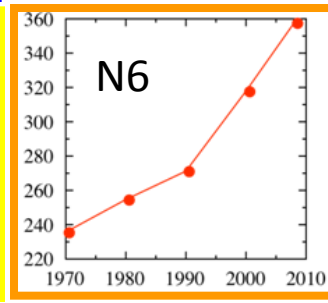
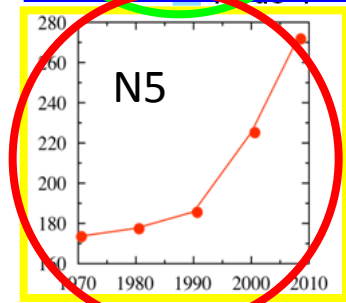
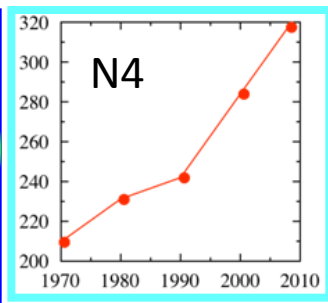
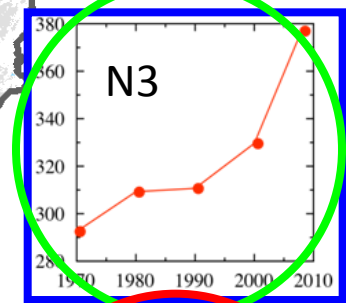
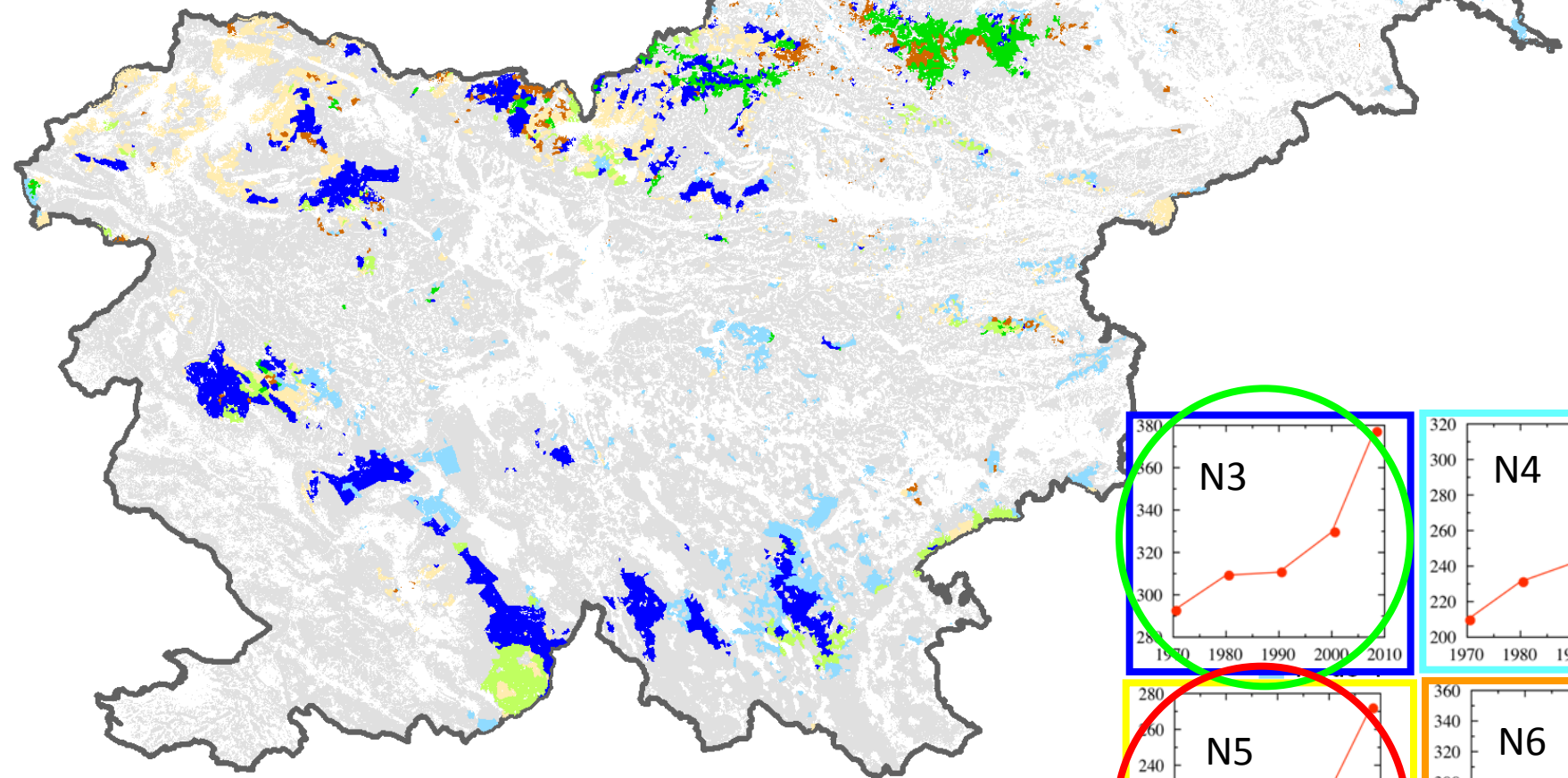
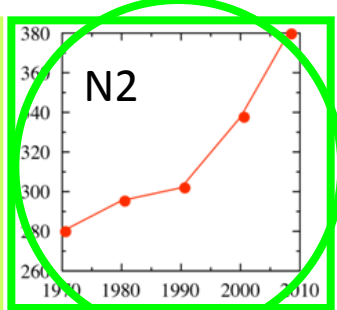
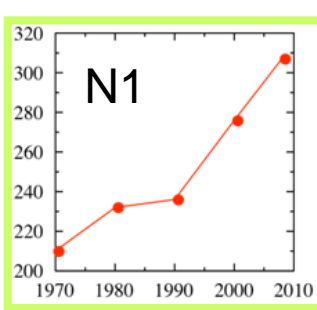
## Private forests



## State forests



# Growing stock in **state** forests



# *Prediction model for year 2018 – MODEL decision trees*

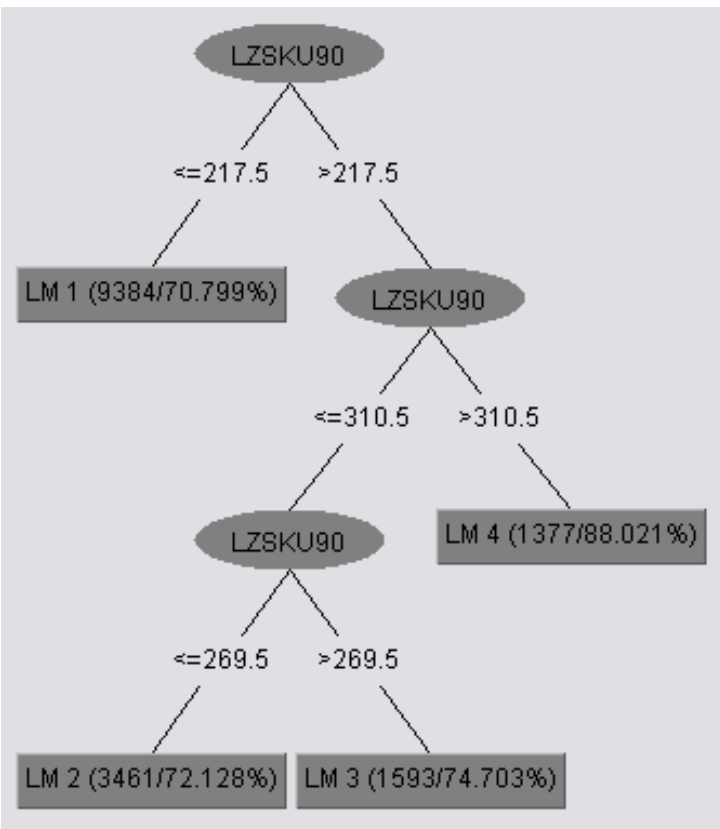
## Step 1: Verification of methodology on prediction for 2008:

- Extrapolation of the linear trends for model 2000 on 2008
- Verification on real data for 2008

## Step 2: Prediction of growing stock for 2018:

- Extrapolation of linear trend of the model for 2008 to the year 2018

## Step 1: Verification of methodology on year 2008 using extrapolation of linear trend from the model for year 2000



LM 1:  
LZSKU00 =  
-0.1233 \* NAKLON\_M  
+ 9.3081 \* MAT\_PODLAG=0  
+ 20.5809 \* RK=13,9,11,3,7,15,17  
+ 4.0773 \* RK=11,3,7,15,17  
- 2.7453 \* RK=3,7,15,17  
- 0.0116 \* RK=7,15,17  
+ 6.7562 \* FITOGEO\_IM=Predpanon,Dinarsko,Predalpsko,Alpsko  
- 6.685 \* FITOGEO\_IM=Dinarsko,Predalpsko,Alpsko  
+ 17.1923 \* FITOGEO\_IM=Predalpsko,Alpsko  
+ 0.0249 \* LZSKU70  
+ 0.0001 \* LZSKU80  
+ 0.5976 \* LZSKU90  
+ 87.0434

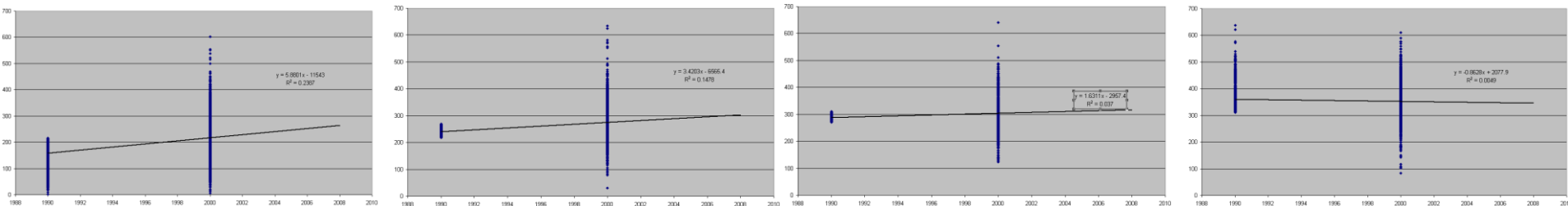
LM 2:  
LZSKU00 =  
0.0002 \* NAKLON\_M  
+ 11.9402 \* MAT\_PODLAG=0  
+ 0.0405 \* RK=13,9,11,3,7,15,17  
+ 0.0132 \* RK=11,3,7,15,17  
+ 0.0088 \* RK=3,7,15,17  
+ 0.0091 \* RK=7,15,17  
+ 0.0099 \* FITOGEO\_IM=Predpanon,Dinarsko,Predalpsko,Alpsko  
- 0.0047 \* FITOGEO\_IM=Dinarsko,Predalpsko,Alpsko  
+ 11.7842 \* FITOGEO\_IM=Predalpsko,Alpsko  
+ 0.0284 \* LZSKU70  
+ 0.0007 \* LZSKU80  
+ 0.5109 \* LZSKU90  
+ 130.2721

LM 3:  
LZSKU08 =  
0.0002 \* NAKLON\_M  
+ 0.1377 \* MAT\_PODLAG=0  
+ 0.0405 \* RK=13,9,11,3,7,15,17  
+ 0.0132 \* RK=11,3,7,15,17  
+ 0.0088 \* RK=3,7,15,17  
+ 0.0281 \* RK=7,15,17  
+ 0.0099 \* FITOGEO\_IM=Predpanon,Dinarsko,Predalpsko,Alpsko  
- 0.0047 \* FITOGEO\_IM=Dinarsko,Predalpsko,Alpsko  
+ 0.0822 \* FITOGEO\_IM=Predalpsko,Alpsko  
+ 0.0005 \* LZSKU70  
+ 0.1804 \* LZSKU80  
+ 0.0074 \* LZSKU90  
+ 254.9865

LM 4:  
LZSKU00 =  
0.0105 \* NMV\_M  
+ 0.0002 \* NAKLON\_M  
+ 0.1034 \* MAT\_PODLAG=0  
+ 0.0405 \* RK=13,9,11,3,7,15,17  
+ 0.0132 \* RK=11,3,7,15,17  
+ 0.0088 \* RK=3,7,15,17  
+ 0.0185 \* RK=7,15,17  
+ 0.0099 \* FITOGEO\_IM=Predpanon,Dinarsko,Predalpsko,Alpsko  
- 0.0047 \* FITOGEO\_IM=Dinarsko,Predalpsko,Alpsko  
+ 0.0212 \* FITOGEO\_IM=Predalpsko,Alpsko  
+ 0.0006 \* LZSKU70  
+ 0.162 \* LZSKU80  
+ 0.4448 \* LZSKU90  
+ 132.4082

2000	LM1	LM2	LM3	LM4
Average real growing stock	217.1	275.1	304.9	352.3
Average model growing stock	253.4	350.9	304.8	352.3
Mean absolute error (MAE)	61.0	93.4	43.3	50.9

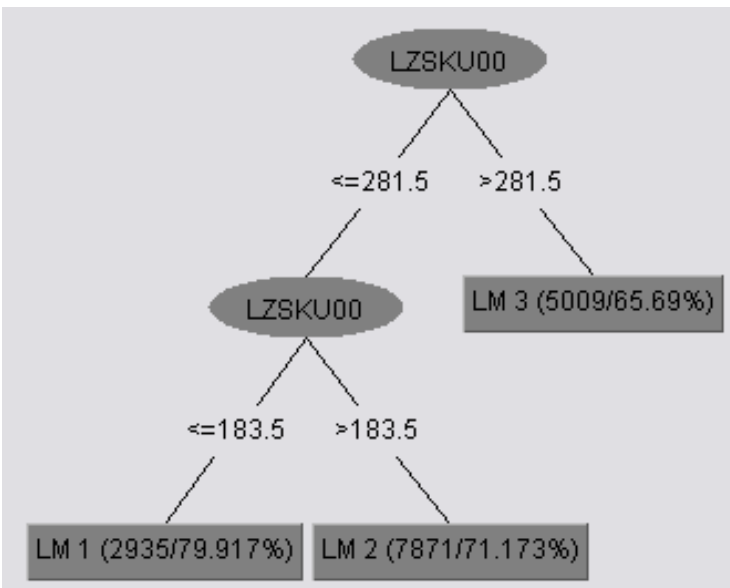
**Step 1:** - Extrapolation of linear trends from 1990-2000 to 2008  
- Verification on real data for 2008



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

# Prediction for private forests in 2018

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



LM 1:  
LZSKU08 =

$$\begin{aligned}
 & -0.0387 * NMV\_M \\
 & + 0.0845 * RK=11,3,7,15,17 \\
 & + 0.0795 * FITOGEO\_IM=Predalpsko,Alpsko \\
 & + 0.0001 * LZSKU70 \\
 & + 0.2302 * LZSKU90 \\
 & + 0.6201 * LZSKU00 \\
 & + 120.3242
 \end{aligned}$$

LM 2:  
LZSKU08 =

$$\begin{aligned}
 & 0 * NMV\_M \\
 & + 14.6206 * RK=11,3,7,15,17 \\
 & + 15.1998 * FITOGEO\_IM=Predalpsko,Alpsko \\
 & + 0.0203 * LZSKU70 \\
 & + 0.2167 * LZSKU90 \\
 & + 0.3798 * LZSKU00 \\
 & + 127.4198
 \end{aligned}$$

LM 3:  
LZSKU08 =

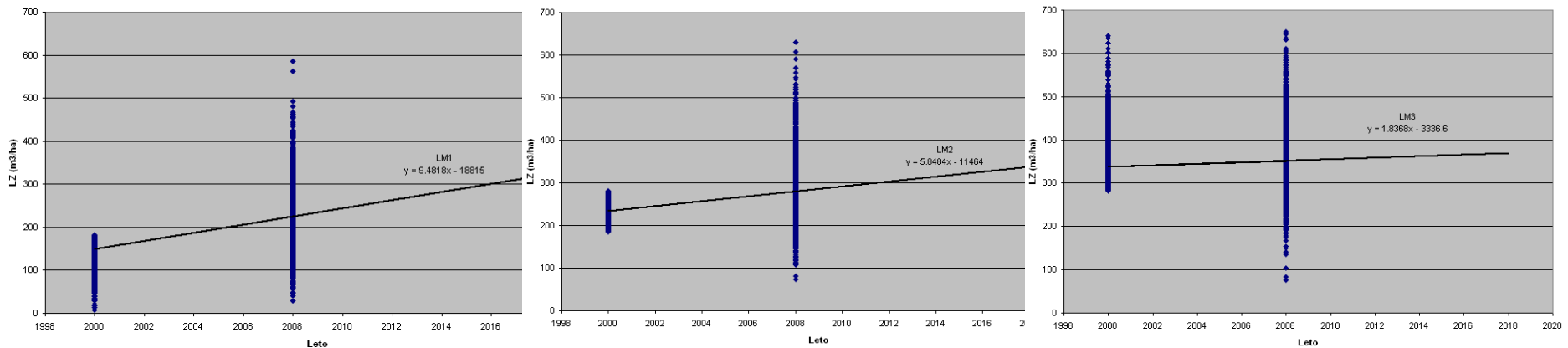
$$\begin{aligned}
 & 0.0399 * NMV\_M \\
 & + 10.8066 * RK=11,3,7,15,17 \\
 & + 0.0283 * FITOGEO\_IM=Predalpsko,Alpsko \\
 & + 0.0001 * LZSKU70 \\
 & + 0.1484 * LZSKU90 \\
 & + 0.591 * LZSKU00 \\
 & + 80.9512
 \end{aligned}$$

2008	LM1	LM2	LM3
Average real growing stock (m <sup>3</sup> /ha)	224.4	280.0	351.7
Average predicted growing stock (m <sup>3</sup> /ha)	224.4	280.0	351.8
Mean absolute error (MAE)	49.1	43.3	35.5



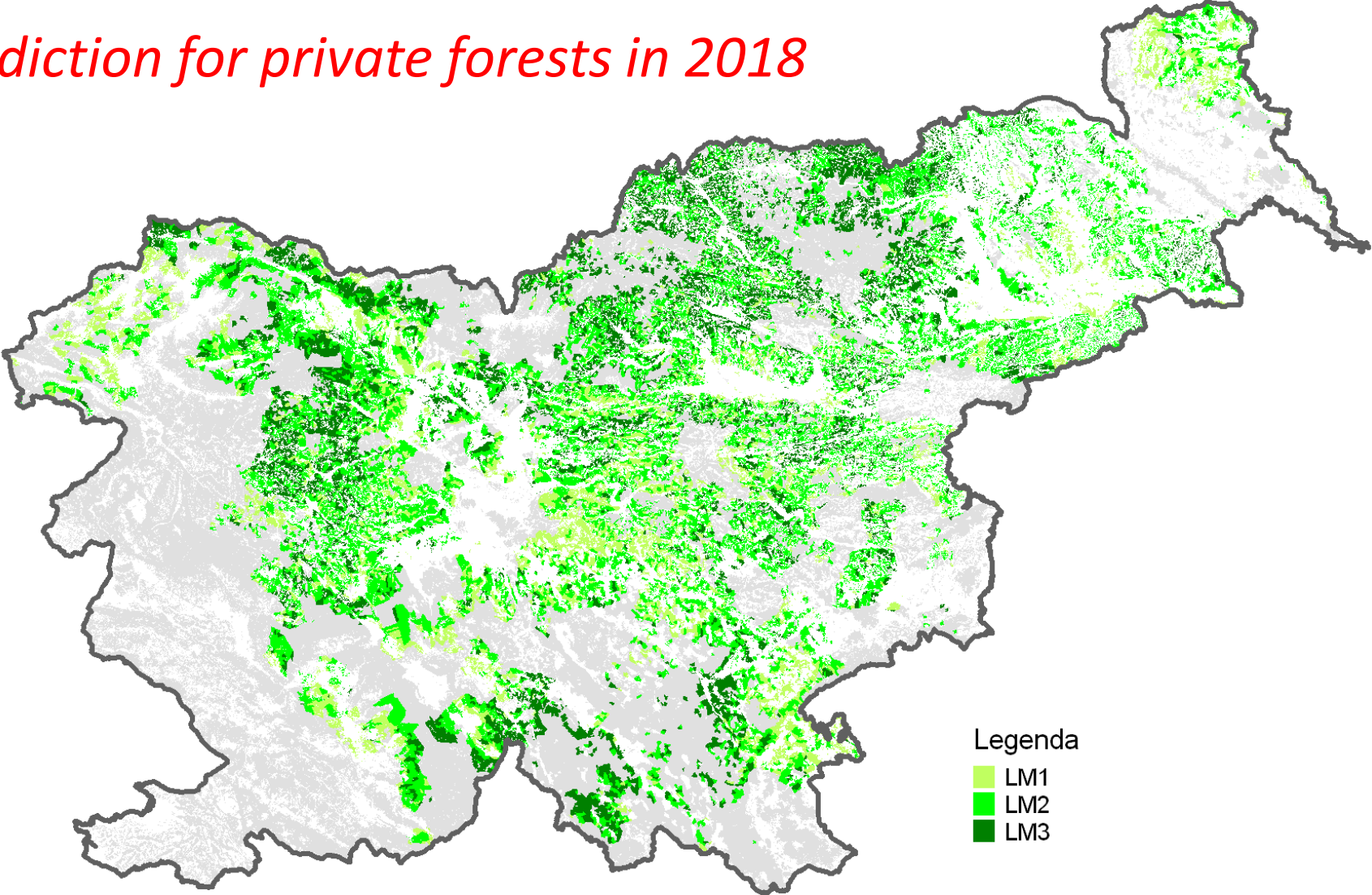


## Step 2: - Extrapolation of linear trends from 2000-2008 to the year 2018



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
<b>Predicted growing stock 2018 (m3/ha)</b>	<b>319.3</b>	<b>338.1</b>	<b>370.1</b>
<b>Index (2008=100)</b>	<b>142.3</b>	<b>120.8</b>	<b>105.2</b>

# Prediction for private forests in 2018



Legenda  
LM1  
LM2  
LM3

	LM1	LM2	LM3
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

Ecological Informatics 5 (2010) 256–266



Contents lists available at [ScienceDirect](#)

## Ecological Informatics

journal homepage: [www.elsevier.com/locate/ecolinf](http://www.elsevier.com/locate/ecolinf)



## Estimating vegetation height and canopy cover from remotely sensed data with machine learning<sup>☆</sup>

Daniela Stojanova<sup>a</sup>, Panče Panov<sup>b,\*</sup>, Valentin Gjorgjioski<sup>b</sup>, Andrej Kobler<sup>a</sup>, Sašo Džeroski<sup>b</sup>

<sup>a</sup> Slovenian Forestry Institute, Večna pot 2, SI-1000 Ljubljana, Slovenia

<sup>b</sup> Jožef Stefan Institute, Department of Knowledge Technologies, Jamova cesta 39, SI-1000 Ljubljana, Slovenia

Problem: Spatial prediction of vegetation high and canopy cover

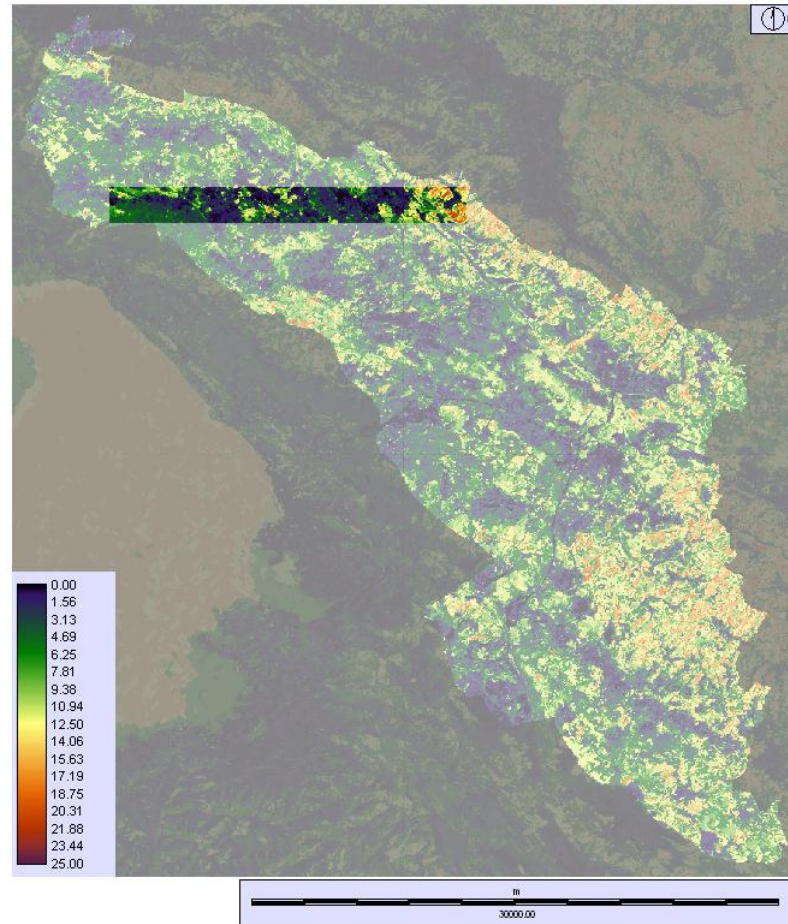
Type of pattern: : **Random Forest Multi Target  
Regression Trees**

Algorithm: **CLUS**

## 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



## 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**





# Vegetation height

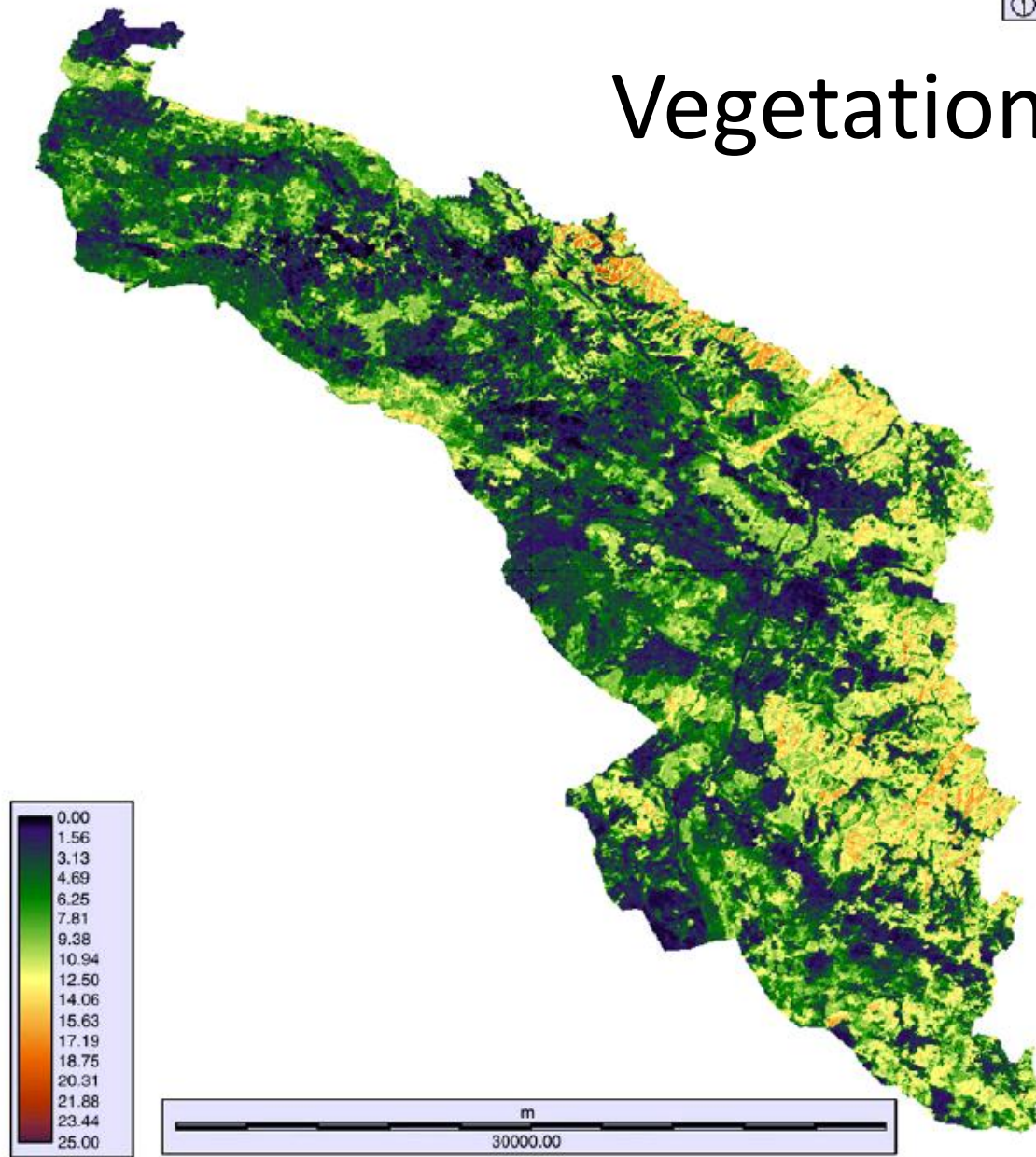


Fig. 5. Map of vegetation height for the Kras region generated by using a random forest of multi-target regression trees model. The legend shows the vegetation height in meters.





# Canopy cover

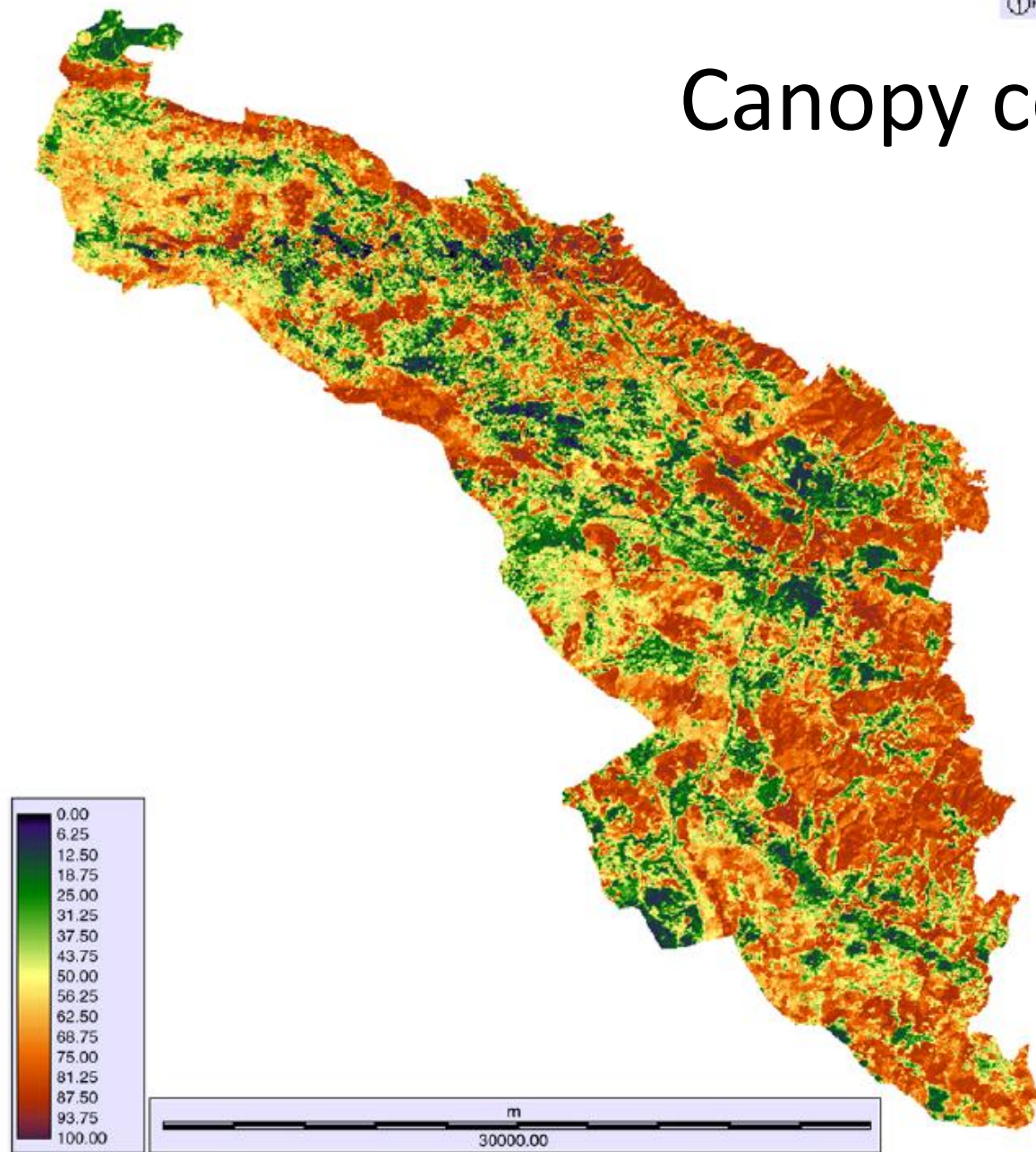


Fig. 6. Map of canopy cover for the Kras region generated by using a random forest of multi-target regression trees model. The legend shows the percentage of canopy cover.

# Data mining of time series: COMMUNITY STRUCTURE

Ecological Modelling 222 (2011) 2524–2529

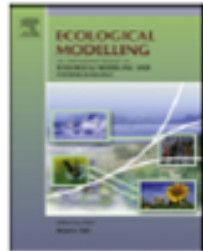


ELSEVIER

Contents lists available at ScienceDirect

Ecological Modelling

Journal homepage: [www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)



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

Marko Debeljak<sup>a,\*</sup>, Geoffrey R. Squire<sup>b</sup>, Dragi Kocev<sup>a</sup>, Cathy Hawes<sup>b</sup>, Mark W. Young<sup>b</sup>, Sašo Džeroski<sup>a</sup>

<sup>a</sup> Jozef Stefan Institute, Jamova cesta 39, 1000 Ljubljana, Slovenia

<sup>b</sup> Scottish Crop Research Institute, Invergowrie, Dundee DD2 5DA, Scotland, UK

**Problem: Temporal prediction of height and cover of crop and weeds**

Type of pattern:

**a) Multi target predictive clustering trees**

**AND**

**b) Constraint multi target predictive clustering trees**

Algorithm: **CLUS**

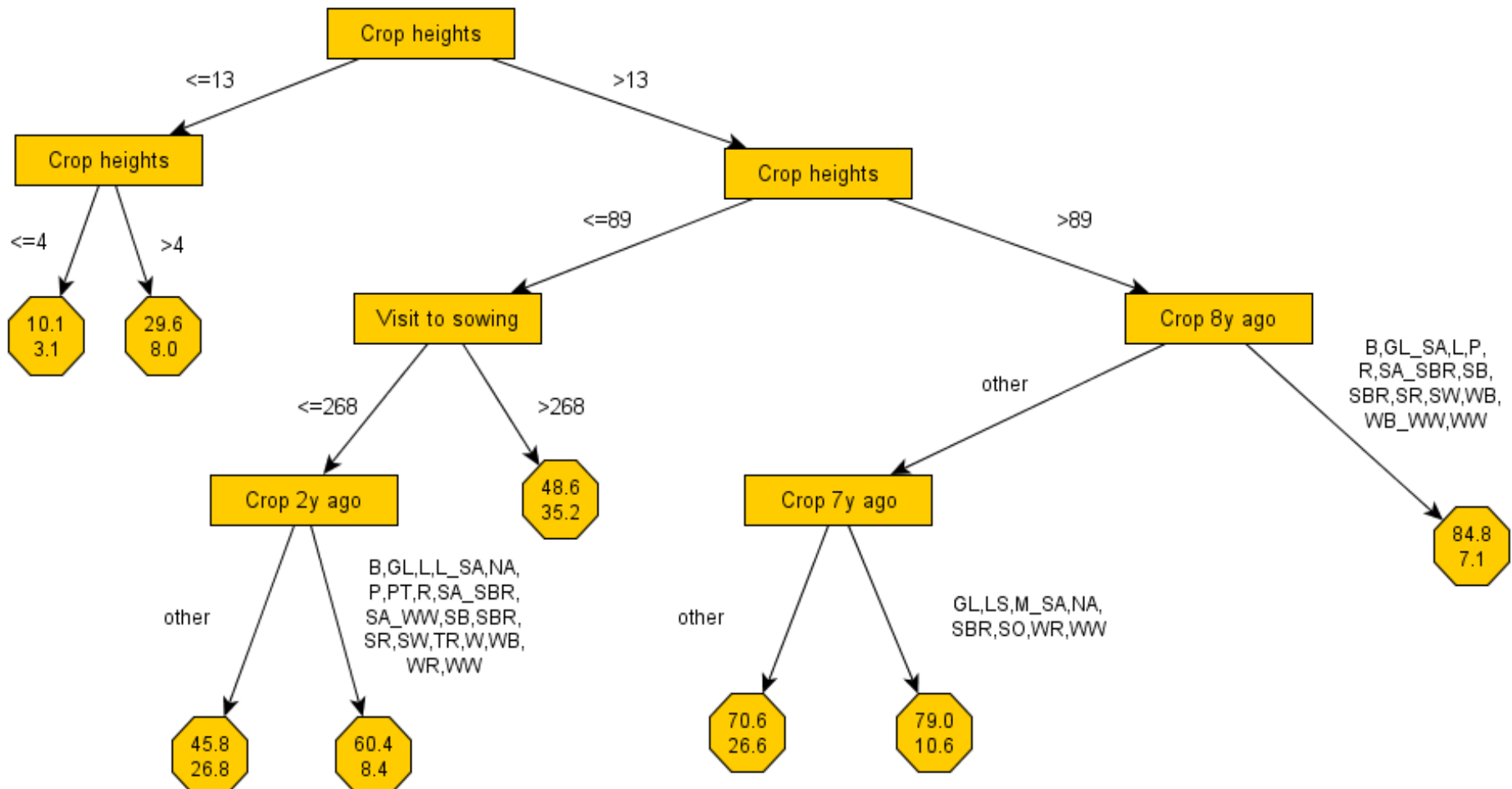
# Data

- 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)



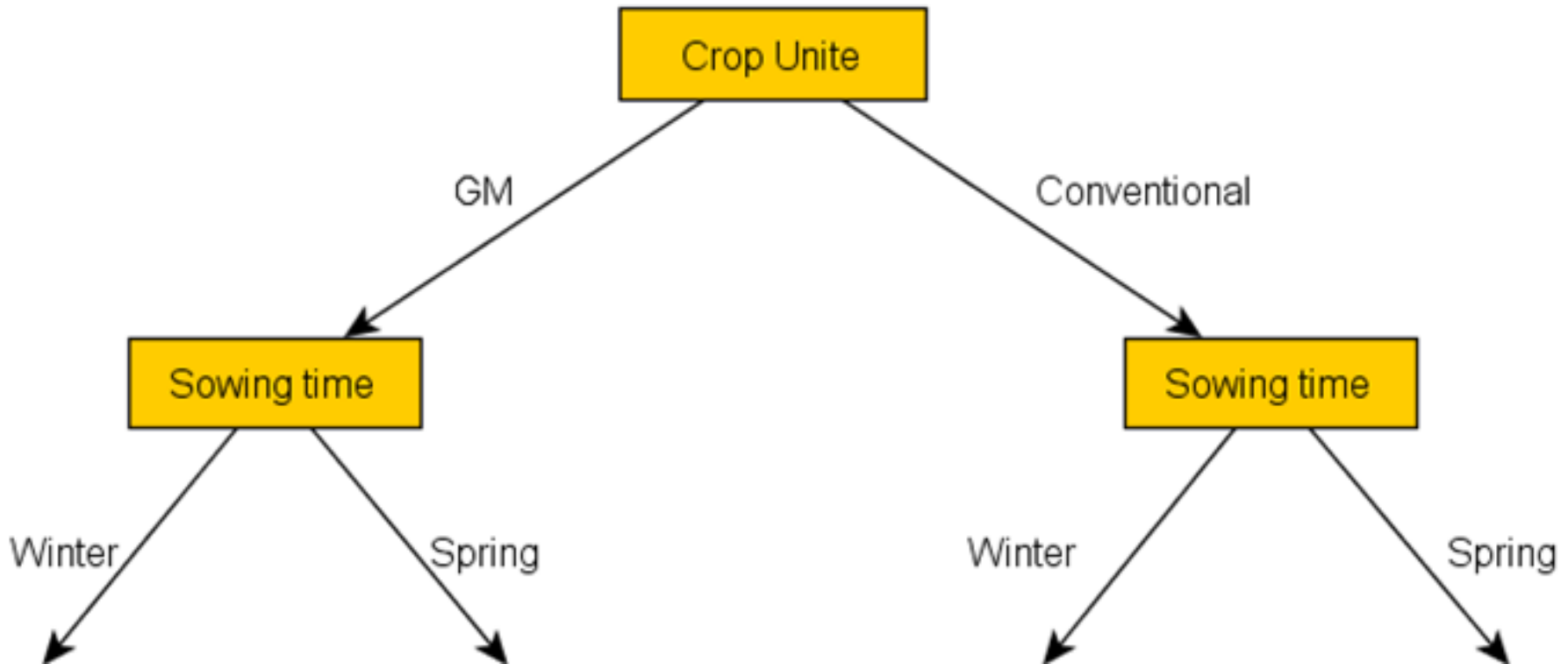
# Results scenario B: Multiple target regression tree

Target: Avg Crop Covers, Avg Weed Covers				Excluded attributes: /		
Constraints: MinInstances = 64.0; MaxSize = 15						
Predictive power:	Corr.Coeff.:	0.8513, 0.3746	RMSE:	16.504, 12.6038	RRMSE:	0.5248, 0.9301



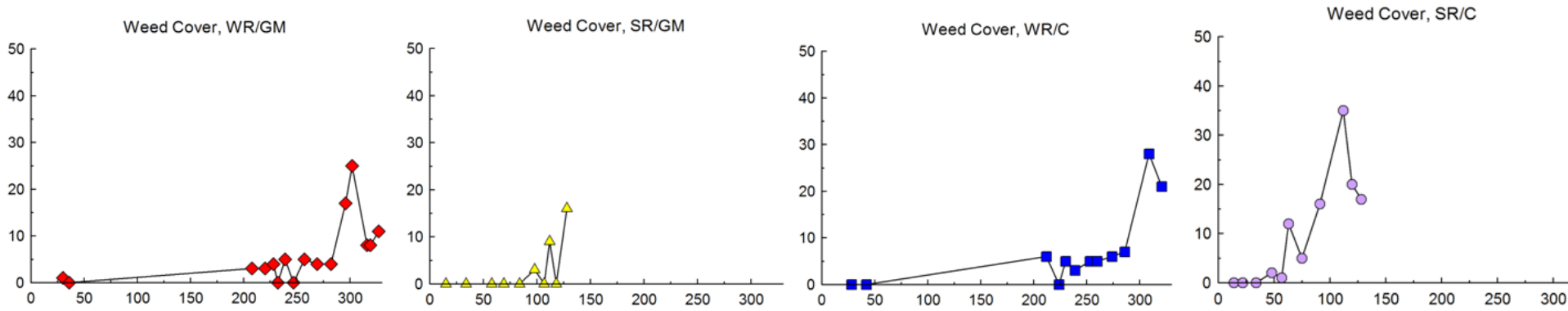
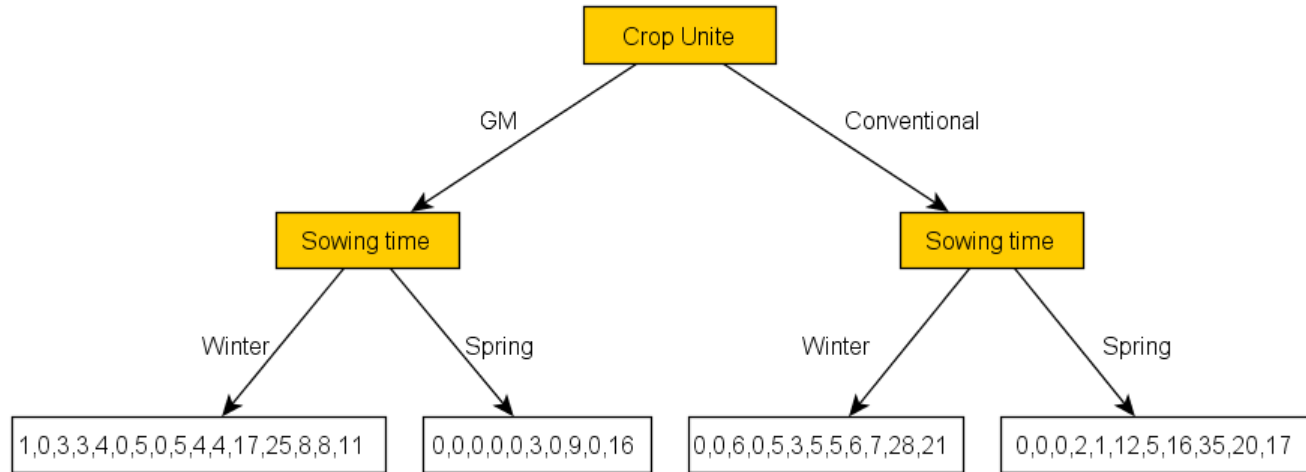
# Results scenario D: Constraint predictive clustering trees for time series including TS clusters for crop (CLUS)

## syntactic constraint



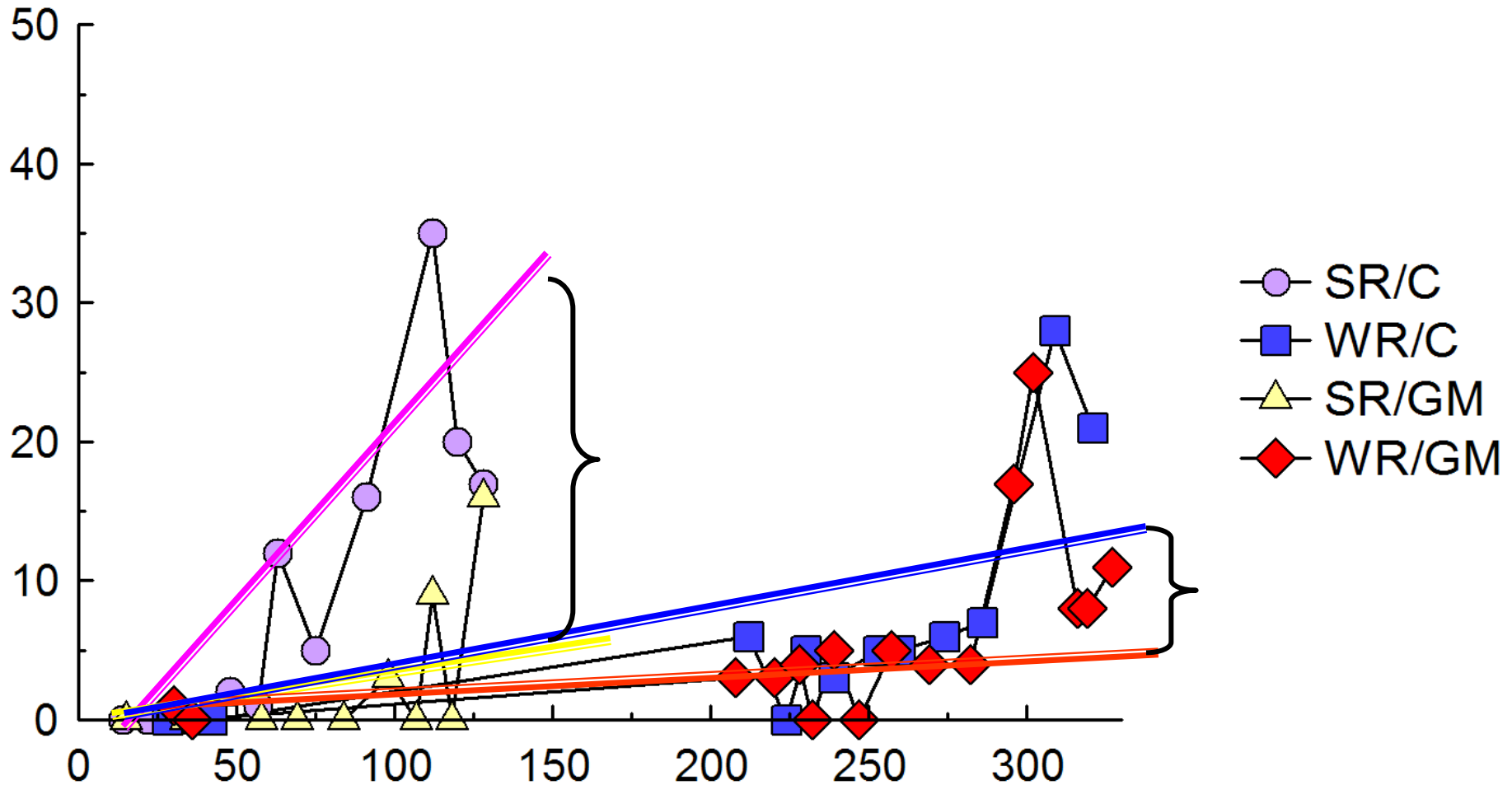
# Results scenario D: Constraint predictive clustering trees for time series including TS clusters for crop (CLUS)

Target: Avg Weed Covers (Time Series)		Scenario 3.9					
Constraints: Syntactic, MinInstances = 32							
Predictive power:	TSRMSExval:	4.98	TSRMSEtrain:	4.86	ICVtrain:	30.44	



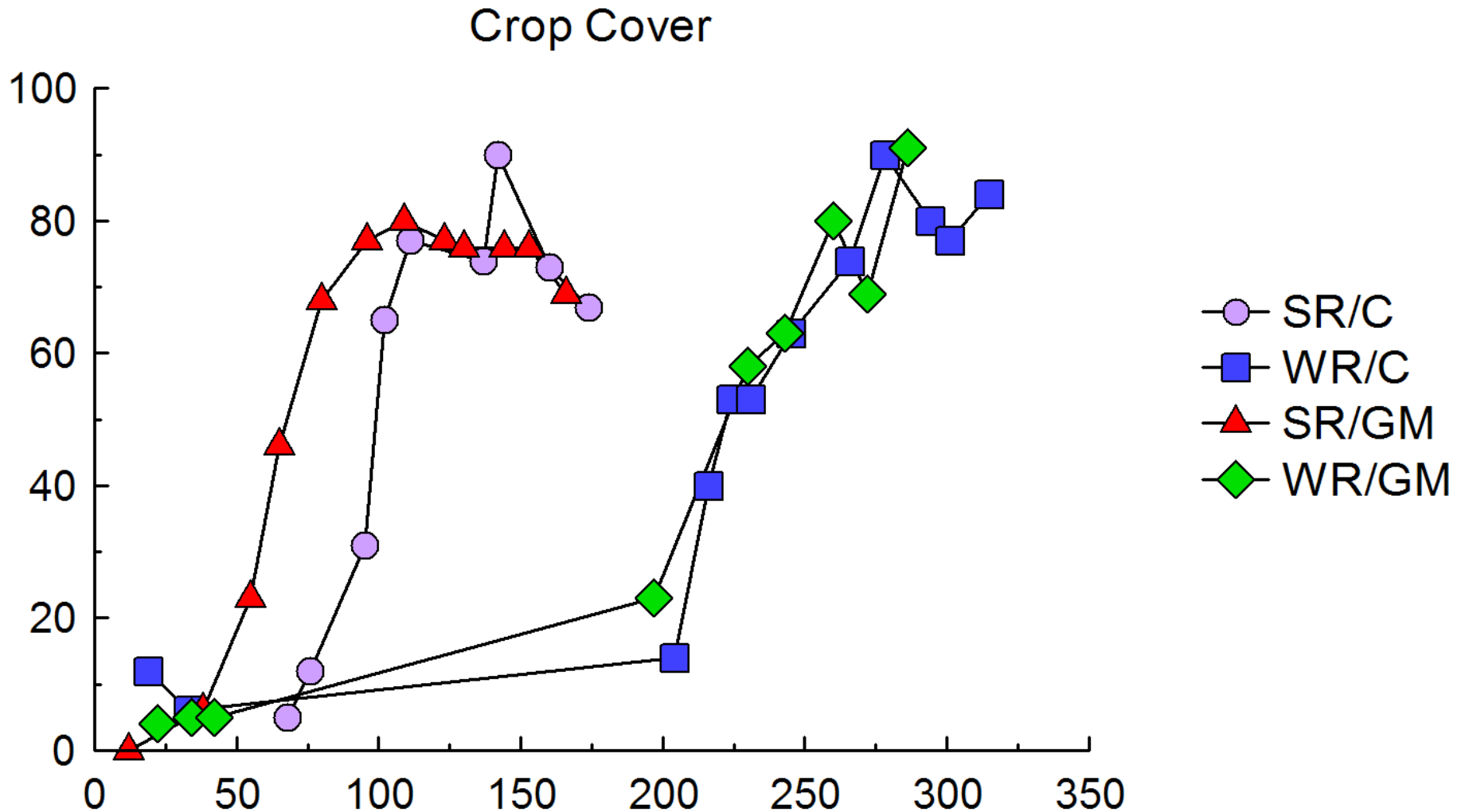
# Results scenario D: Constraint predictive clustering trees for time series including TS clusters for crop (CLUS)

Weed Cover





# Results scenario D: Constraint predictive clustering trees for time series including TS clusters for crop (CLUS)



# Relational data mining: GENE FLOW MODELLING

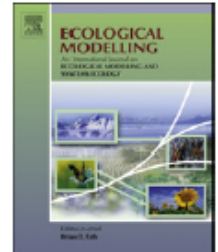
Ecological Modelling 245 (2012) 75–83



Contents lists available at SciVerse ScienceDirect

## Ecological Modelling

journal homepage: [www.elsevier.com/locate/ecolmodel](http://www.elsevier.com/locate/ecolmodel)



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

Marko Debeljak<sup>a,b,\*</sup>, Aneta Trajanov<sup>a</sup>, Daniela Stojanova<sup>a,b</sup>, Florence Leprince<sup>c</sup>, Sašo Džeroski<sup>a,b,d</sup>

<sup>a</sup> Jozef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia

<sup>b</sup> Jozef Stefan International Postgraduate School, Jamova 39, 1000 Ljubljana, Slovenia

<sup>c</sup> ARVALIS-Institut du Végétal, 21, Chemin de Pau, 64121 Montardon, France

<sup>d</sup> Centre of Excellence for Integrated Approaches in Chemistry and Biology of Proteins, Jamova 39, 1000 Ljubljana, Slovenia

Problem: **Classification of fields to above or below  
0.9% of outcrossing**

Type of pattern: **Relational classification decision tree**

Algorithm: **TILDE**

# Spatial temporal relations

**2004:**

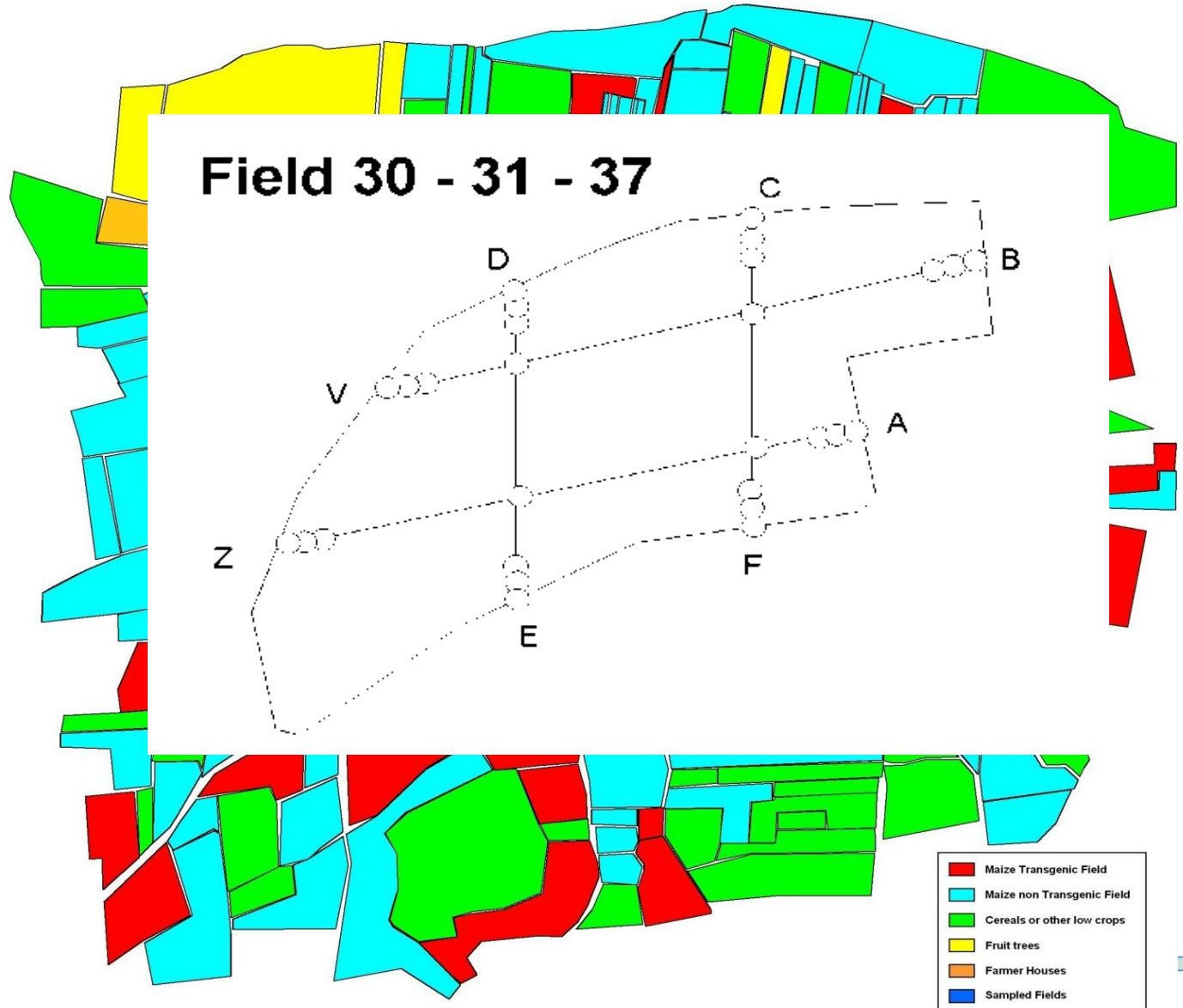
40 GM fields  
7 non-GM fields  
181 sampling points

**2005:**

17 GM fields  
4 non-GM fields  
127 sampling points

**2006:**

43 GM fields  
4 non-GM fields  
112 sampling points

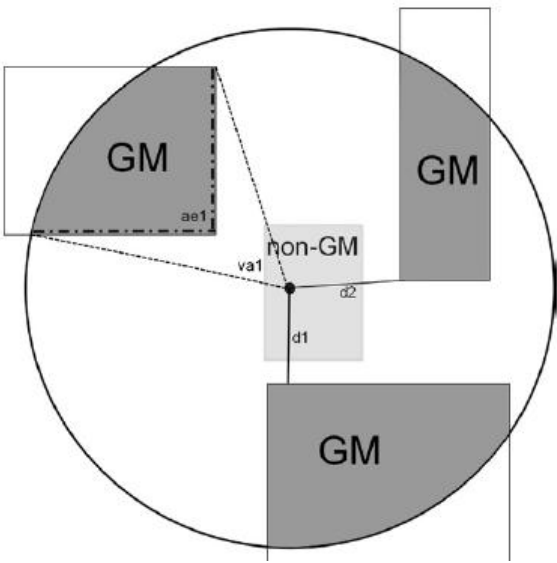


## Relational data mining: GENE FLOW MODELLING

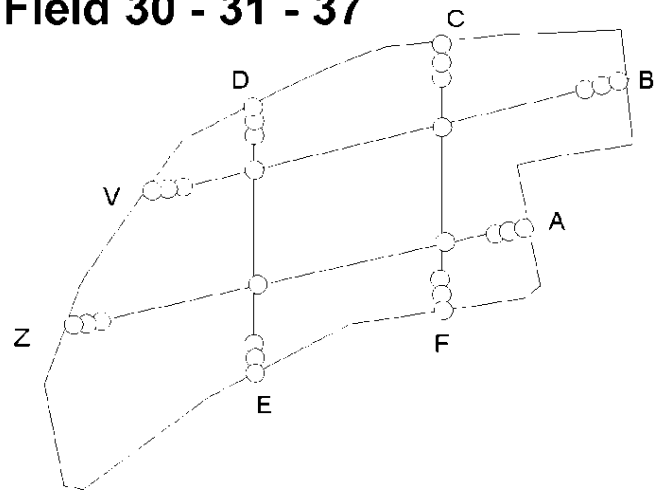
Data scattered over several **tables** or **relations**:

- 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

# Relational data mining: GENE FLOW MODELLING

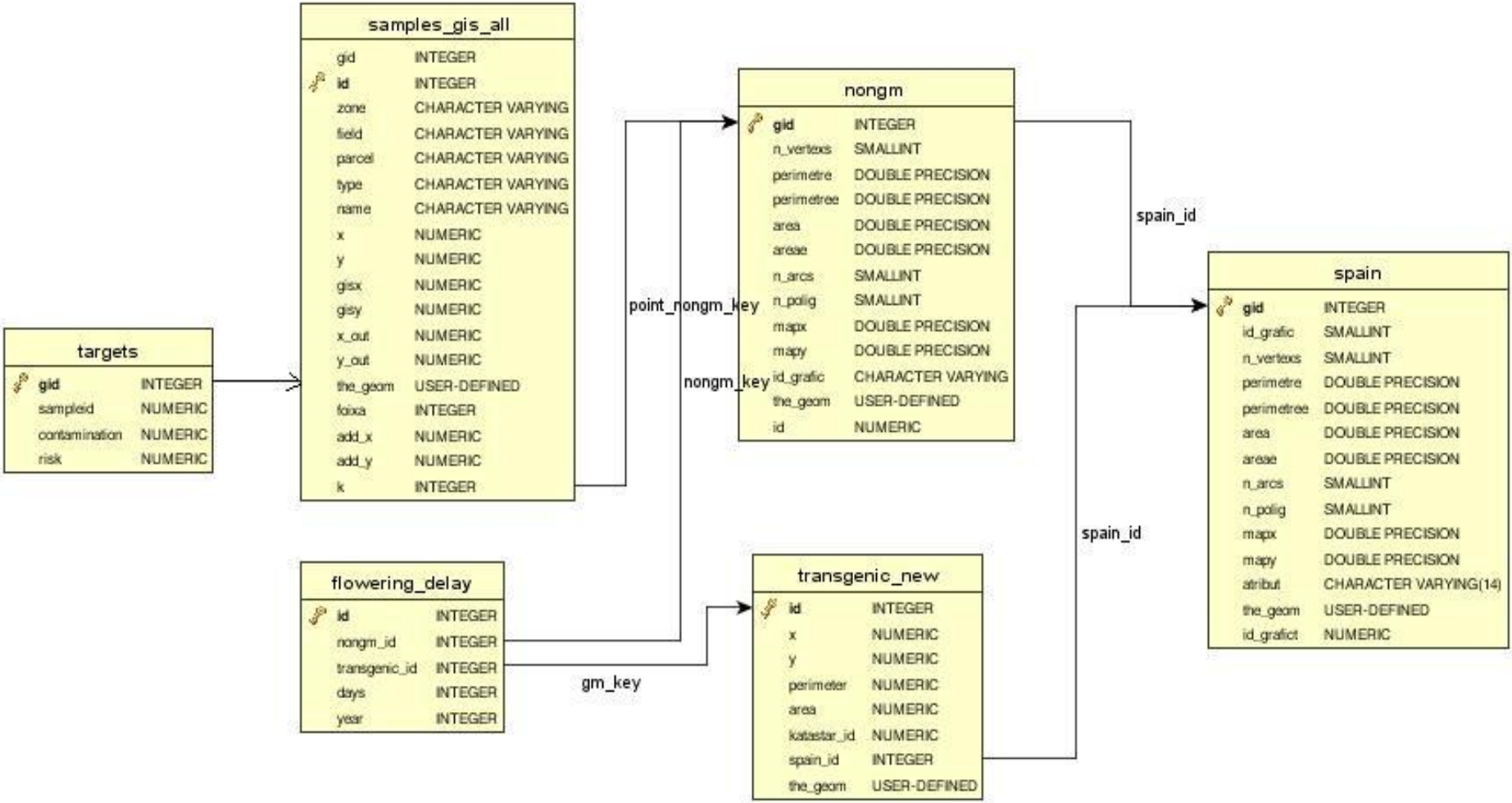


Field 30 - 31 - 37



# Relational data mining: GENE FLOW MODELLING

## Relation database system **PostGIS**



# Relational data mining – building model

---

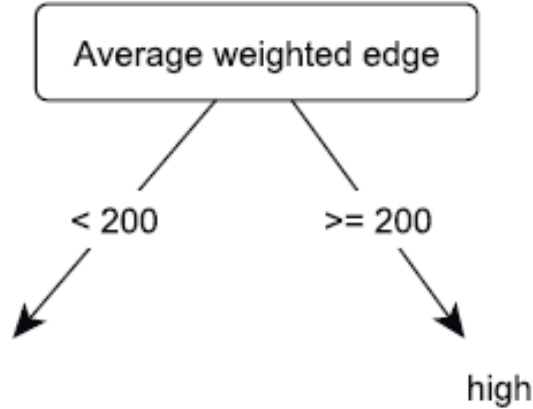
## Relation data analysis:

- Algorithm **Tilde** (Blockeel and De Raedt, 1998; De Raedt et al., 2001) => upgrade of algorithm **C.4.5** (Quinlan, 1993) for classification decision trees
- The algorithm is included in the **ACE-iiProlog** data mining system

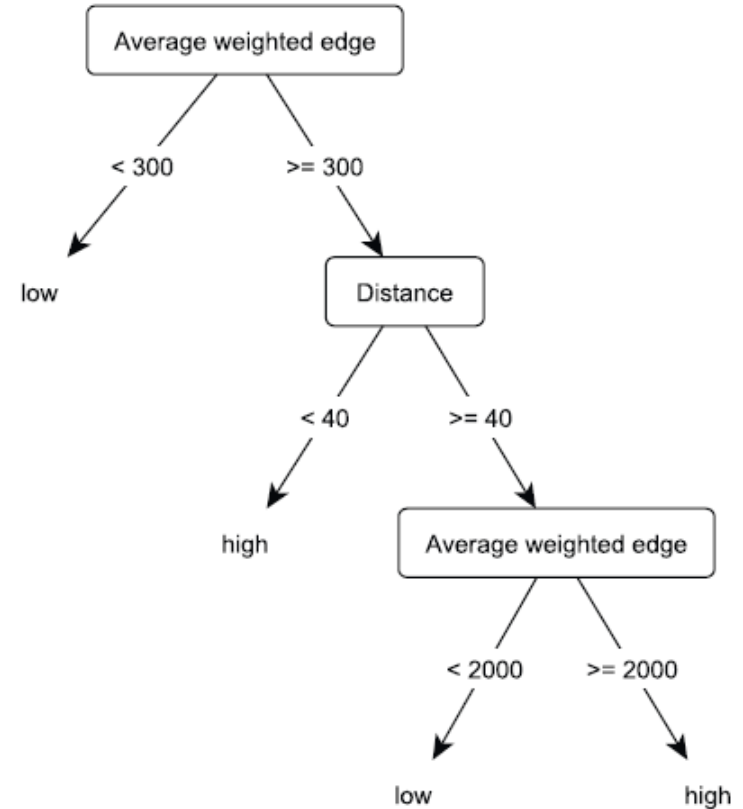


# Relational data mining – results

Threshold 0.01%

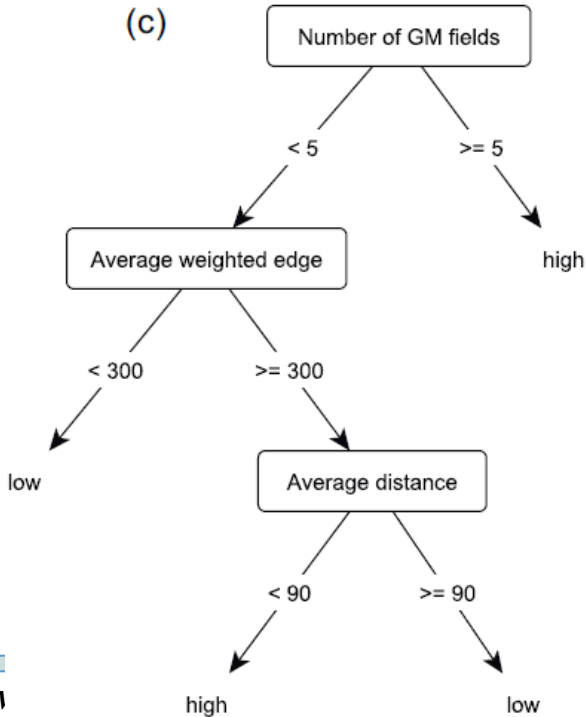


Threshold 0.45%



Threshold 0.9%

(c)



## **HABITAT MODELING OF TICK *Ixodes ricinus*, MAJOR VECTOR FOR LYME BORRELIOSIS and TICK-BORNE MENINGOENCEPHALITIS IN THE UPPER SOČA VALLEY**

Prof. dr. Marko Debeljak

Problem: **Prediction of habitat**

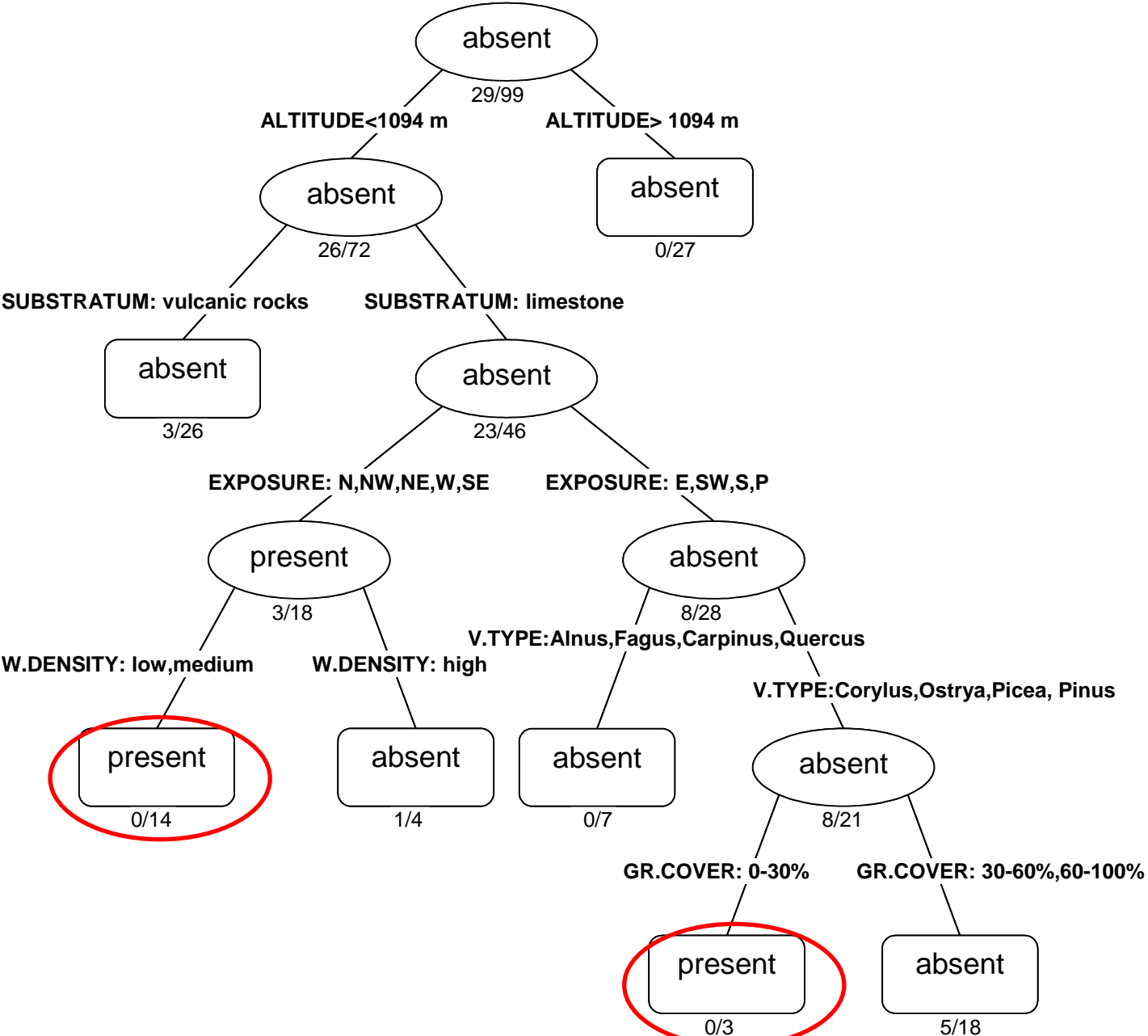
Type of pattern: **Classification decision tree**

Algorithm: **J4.8**

# Study area: Zgornje Posočje



# Classification tree - Trentino region:

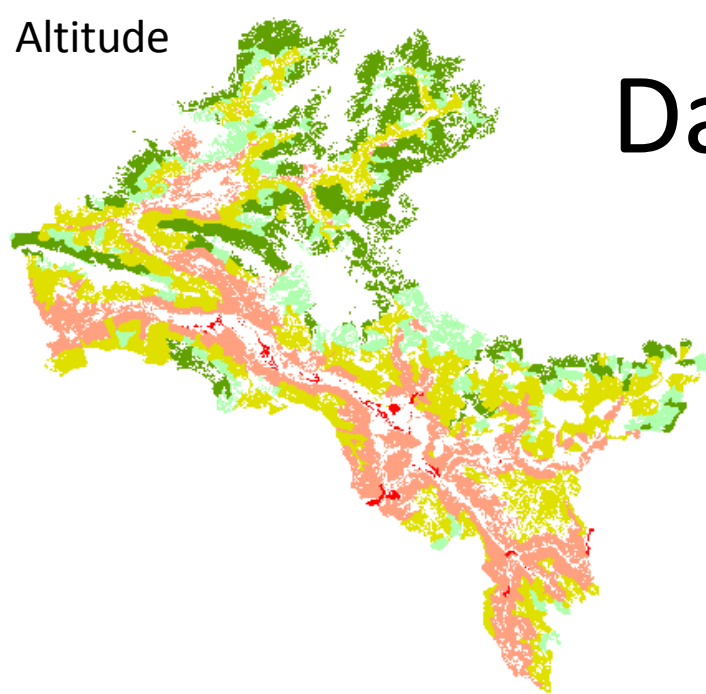




Altitude

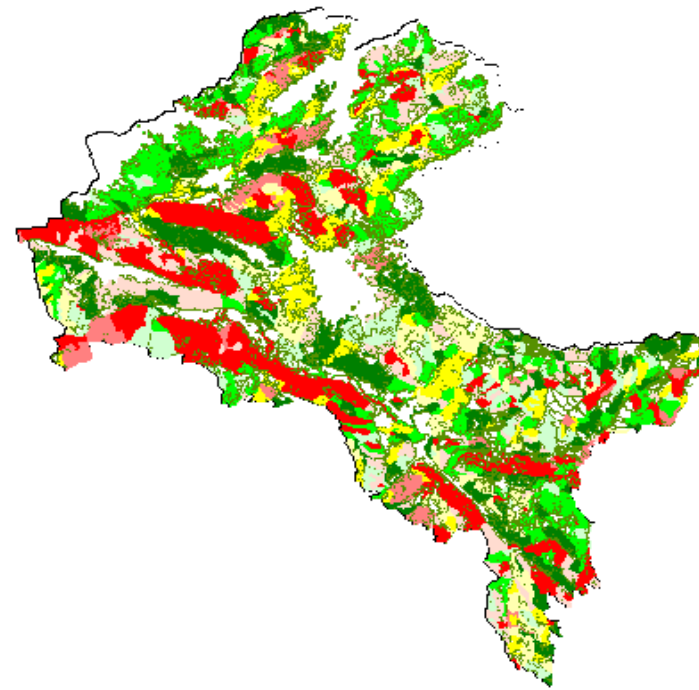
# Data

Exposition



Višinski pasovi:

0-300 m	(25)
300-600 m	(353)
600-900 m	(480)
900-1100 m	(220)
nad 1100 m	(289)

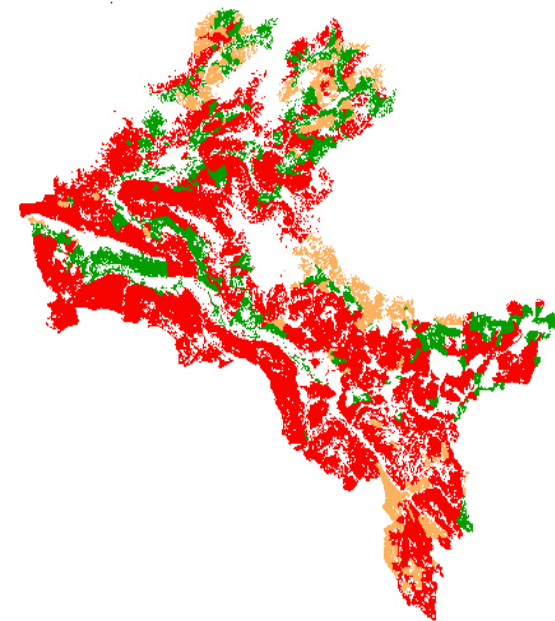


Ekspozicija površine:

0	J	(196)
1	JV	(224)
2	V	(139)
5	JZ	(109)
7	Z	(158)
10	SV	(177)
15	SZ	(77)
20	S	(287)

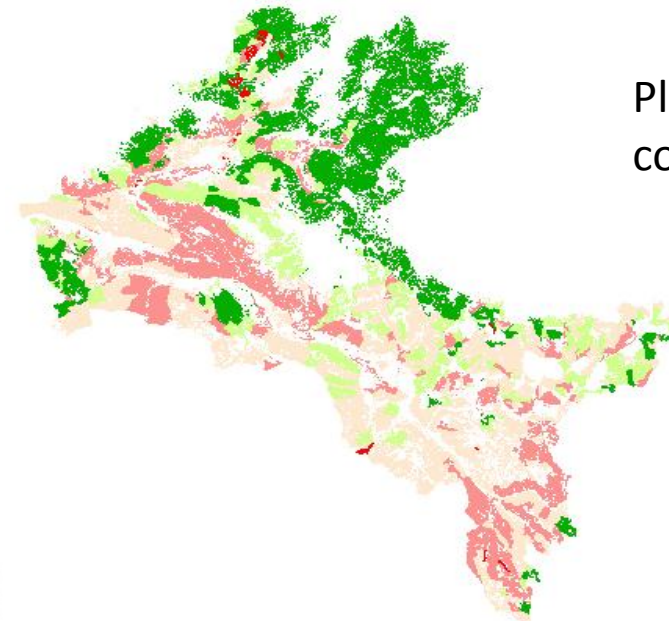
Bedrock

Plant communities



Vrste kamnin:

apnenci, apn. breče, apn. melišča, apn. naplavine	(836)
dolomitizirani apnenci	(161)
ostalo	(370)

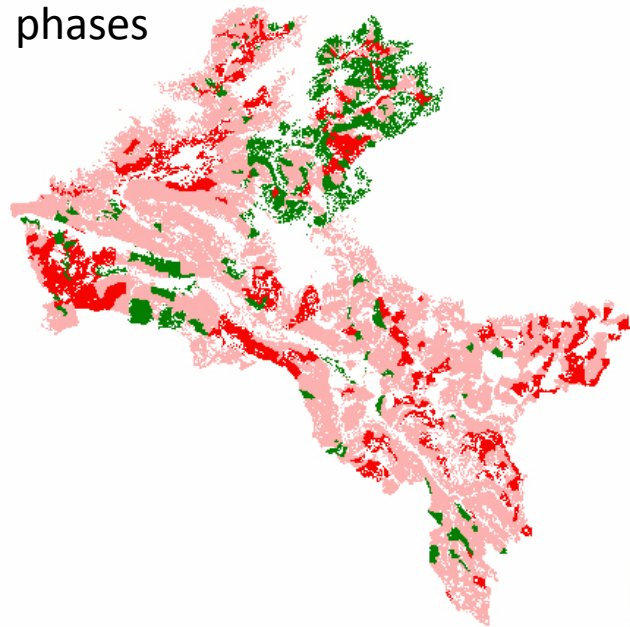


Gozdne združbe (glede na pogoje rastišča):

zelo tople lege	(19)
tople lege	(255)
zmerne lege	(426)
hladne lege	(197)
zelo hladne lege	(470)

# Forest development phases

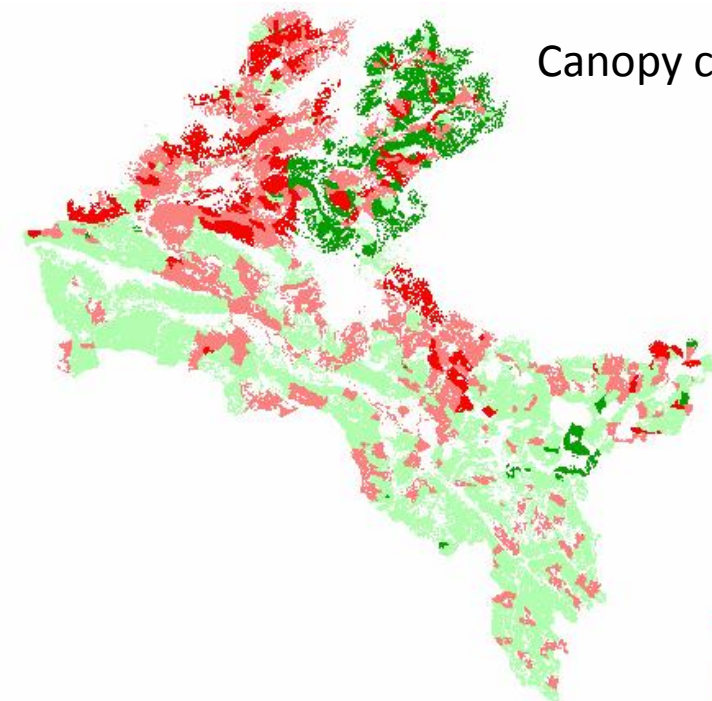
# Data



Razvojna faza:

■ mladovje, parjevec, grmičav gozd, pionirski gozd z grmišči	(262)
■ drogovnjak, dvoslojni sestoj, sestoj v obnovi	(886)
■ debeljak, tipični prebiralni sestoj	(219)

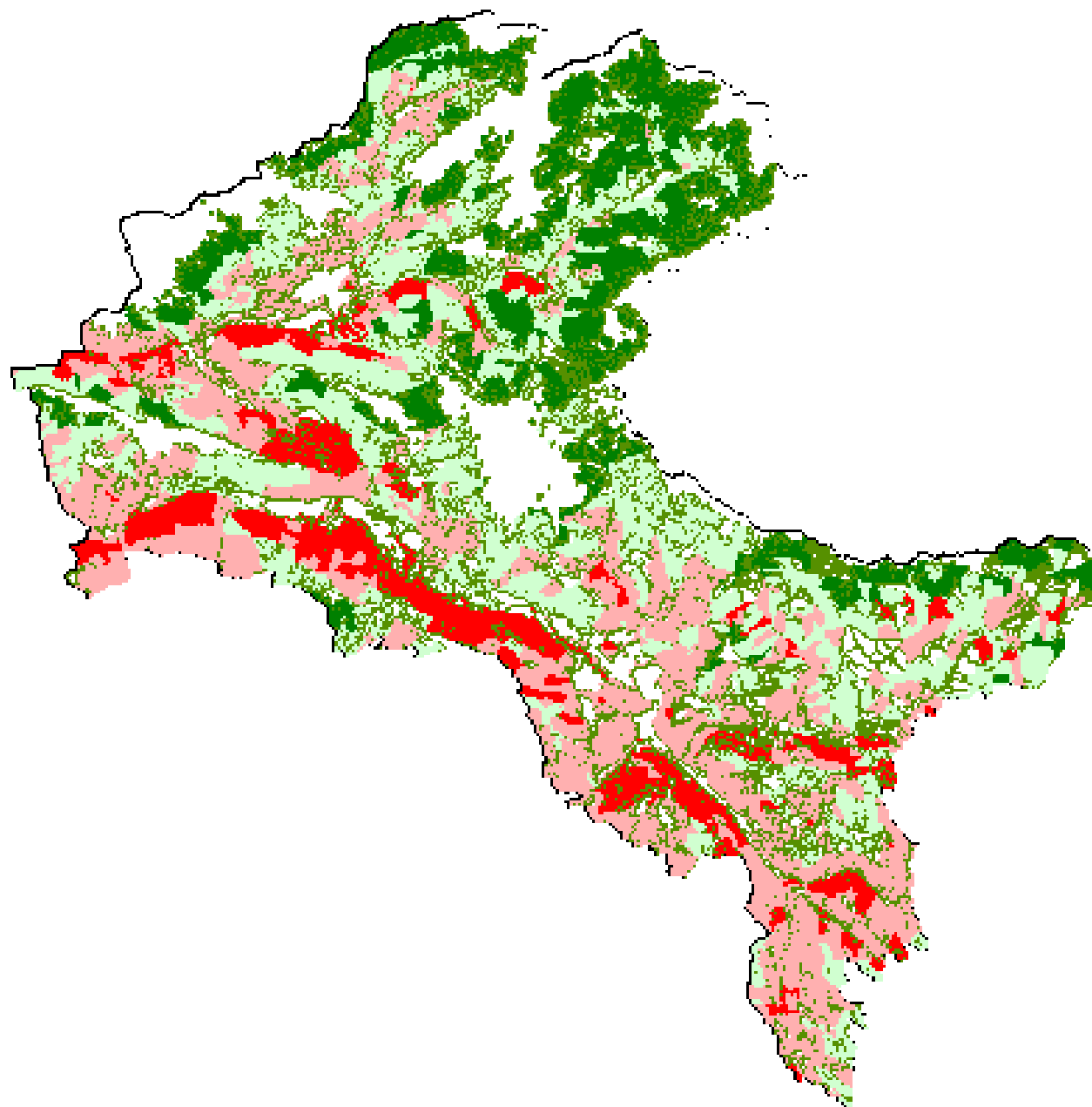
# Canopy coverage




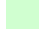

Sklep krošenj:

■ vrzelast do pretrgan	(170)
■ rahel	(401)
■ normalen	(656)
■ tesen	(140)

# Spatial probability of tick presence (tick habitat)

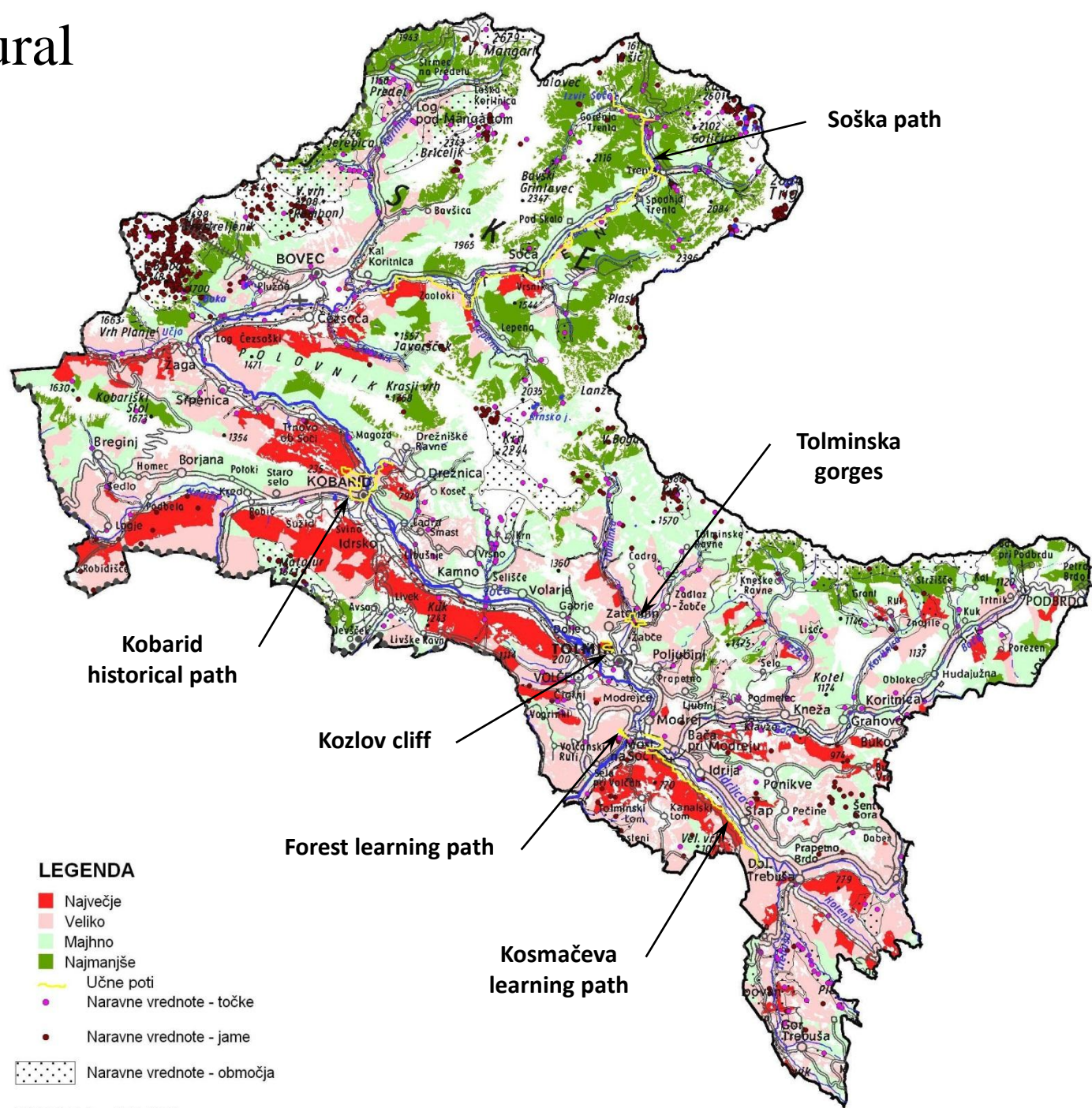


Verjetnost pojavljanja klopov:

	1 največja	(155)
	2 velika	(482)
	3 manjša	(467)
	4 najmanjša	(263)



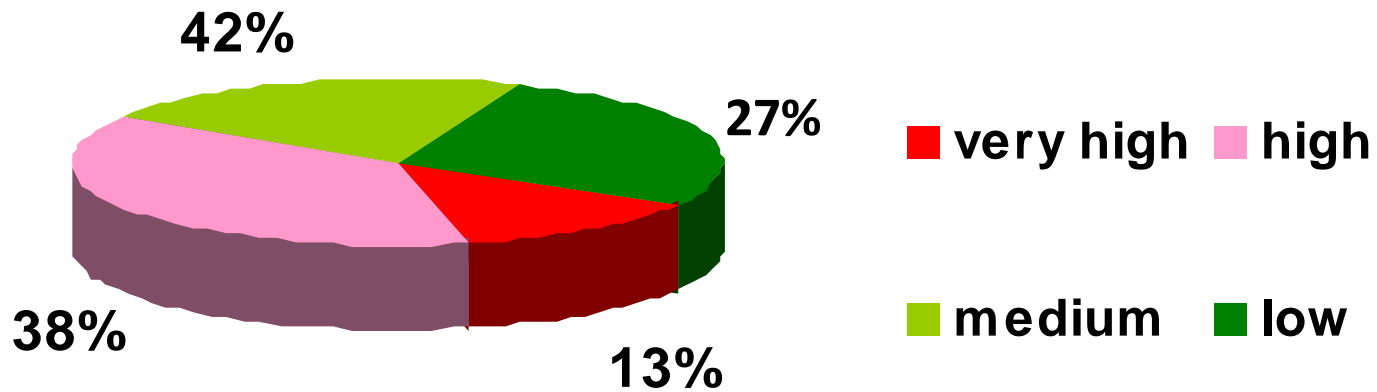
# Locations of natural heritages and exposure of their visitors to ticks



- LEGENDA**
- Največje
  - Veliko
  - Majhno
  - Najmanjše
  - Učne poti
  - Naravne vrednote - točke
  - Naravne vrednote - jame
  - Naravne vrednote - območja

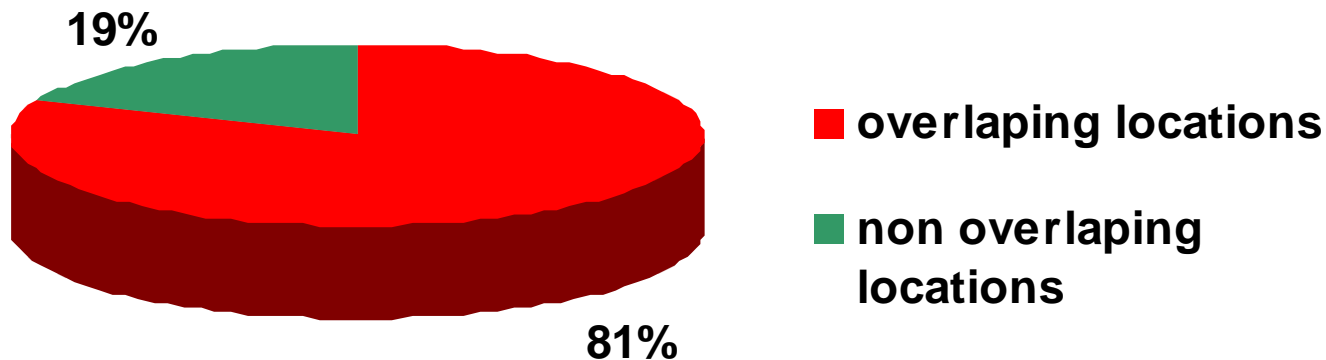
MERILO 1 : 250 000

# Distribution of **natural heritage** according to the risk of infections with tick-borne diseases in upper Soča valley



## Verification of the model:

- Overlapping with the locations of confirmed infections



- Locations with the highest confirmed infections with tick-borne diseases: Tolmin, Kobarid, Volče, Zatoimin, Idrsko, Most na Soči, Poljubinj, Borjana, Bovec, Ljubinj, Žabče, Trebuša, Vrsno.

## 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. DEBELJAK<sup>1</sup>, D. KOCEV<sup>1</sup>, W. TOWERS<sup>2</sup>, M. JONES<sup>2</sup>, B. S. GRIFFITHS<sup>3,\*</sup> & P. D. HALLETT<sup>3</sup>

<sup>1</sup>Department of Knowledge Technologies, Jožef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia, <sup>2</sup>Macaulay Institute, Craigiebuckler, Aberdeen AB15 8QH, UK, and <sup>3</sup>Scottish Crop Research Institute, Invergowrie, Dundee DD2 5DA, UK

Problem: **Prediction of soil exposure to disturbances**

Type of pattern: **Multi target regression trees** (resistance, resilience)

Algorithm: **CLUS**



# 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

## 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:  $1/C_c$ ,  
recovery from compression:  $C_e/C_c$ , overburden stress: eg,  
recovery from overburden stress after two days cycles:  
eg2dc
- **biological resilience:** heat, copper

# Multi-target regression model: RISK MODELLING

Independent variables

Major soil subgroup

Sand

Silt

Clay

pH

C

N

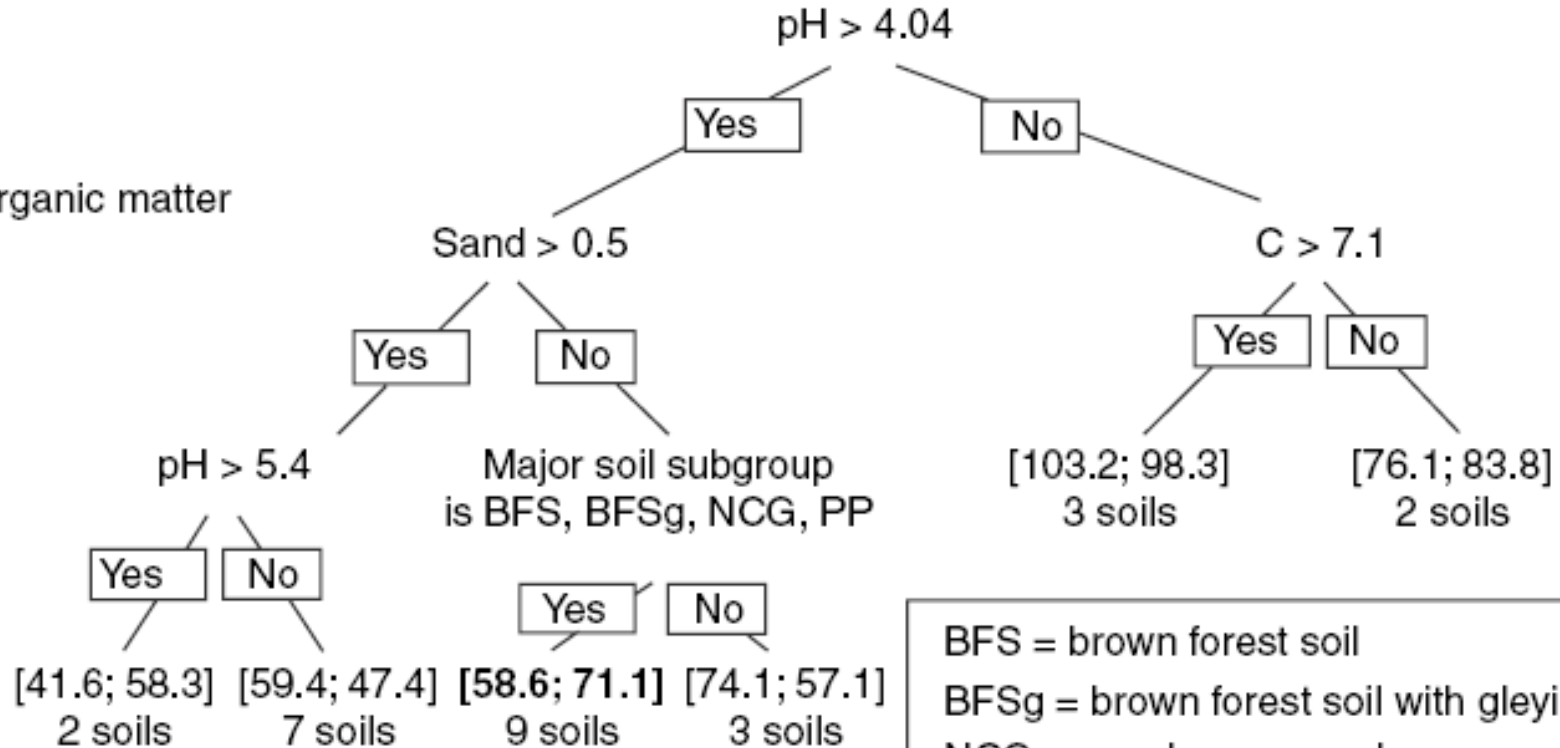
Soil organic matter

Dependent variables

Copper resistance

Copper resilience

Model 10 (see Table 1)



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



**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

- o MAJOR\_SOIL\_SUBGROUP
- o Sand
- o Silt
- o Clay
- o pH
- o C
- o N
- o SOM

### Dependent attributes

- o Heat\_resist
- o Heat\_Resil
- o Cu\_resist
- o Cu\_Resilience
- o Over\_resist\_eg
- o Over\_resil\_2dc\_eg
- o 1/Cc
- o Ce/Cc

Validated Parameters:

RMSE: [15.4825,20.8952,11.0126,18.5897,0.0809,0.2114,0.3469,0.0098]

Correlation Coefficient: [-0.3097,-0.0784,0.8062,0.4445,-0.1455,0.2886,0.355,0.8447]

C > 7.6

+--yes: [56.47,110.4,103.24,98.3,0.789667,1.359333,0.371243,0.055778]: 3

+--no: C > 4.58

+--yes: pH > 4.09

| +--yes: [48.836,93.08,70.264,64.36,0.8828,1.1358,0.626798,0.016037]: 5

| +--no: [39.8,91.733333,73.483333,79.1,0.805333,0.875,0.847719,0.040065]: 3

+--no: N > 0.14

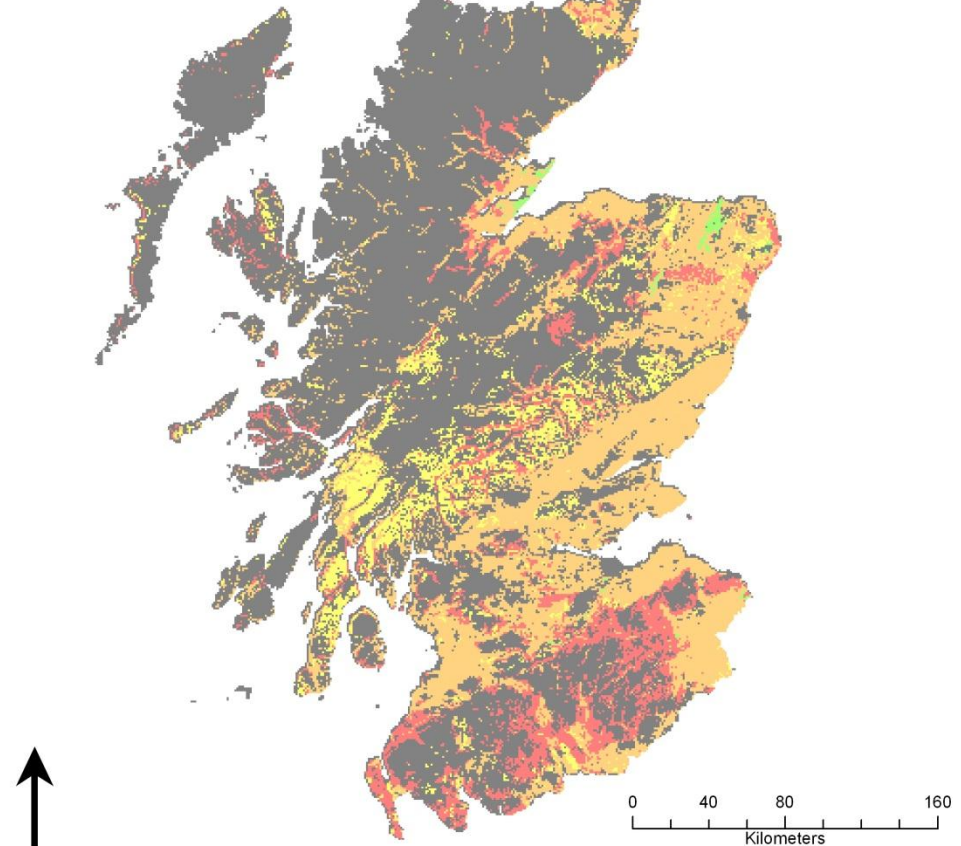
+--yes: [53.946667,86.833333,54.366667,57.35,0.8165,0.9415,0.914021,0.011523]

+--no: [69.673333,81,59.116667,59.866667,0.865,0.911,1.54998,0.02033]: 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



# Application

## Experiment 3A

### Independent Attributes

- o MAJOR\_SOIL\_SUBGROUP
- o Sand
- o Silt
- o Clay
- o pH
- o C
- o N
- o SOM

### Dependent attributes

- o Over\_resist\_eg
- o Over\_resil\_2dc\_eg

Validated Parameters:

RMSE: [0.0965,0.2644]

Correlation Coefficient: [-0.3021,0.1116]

C > 7.6

+--yes: [0.789667,1.359333]: 3

+--no: N > 0.21

+--yes: pH > 5.5

| +--yes: [0.742,0.9415]: 2

| +--no: [0.863077,1.042462]: 13

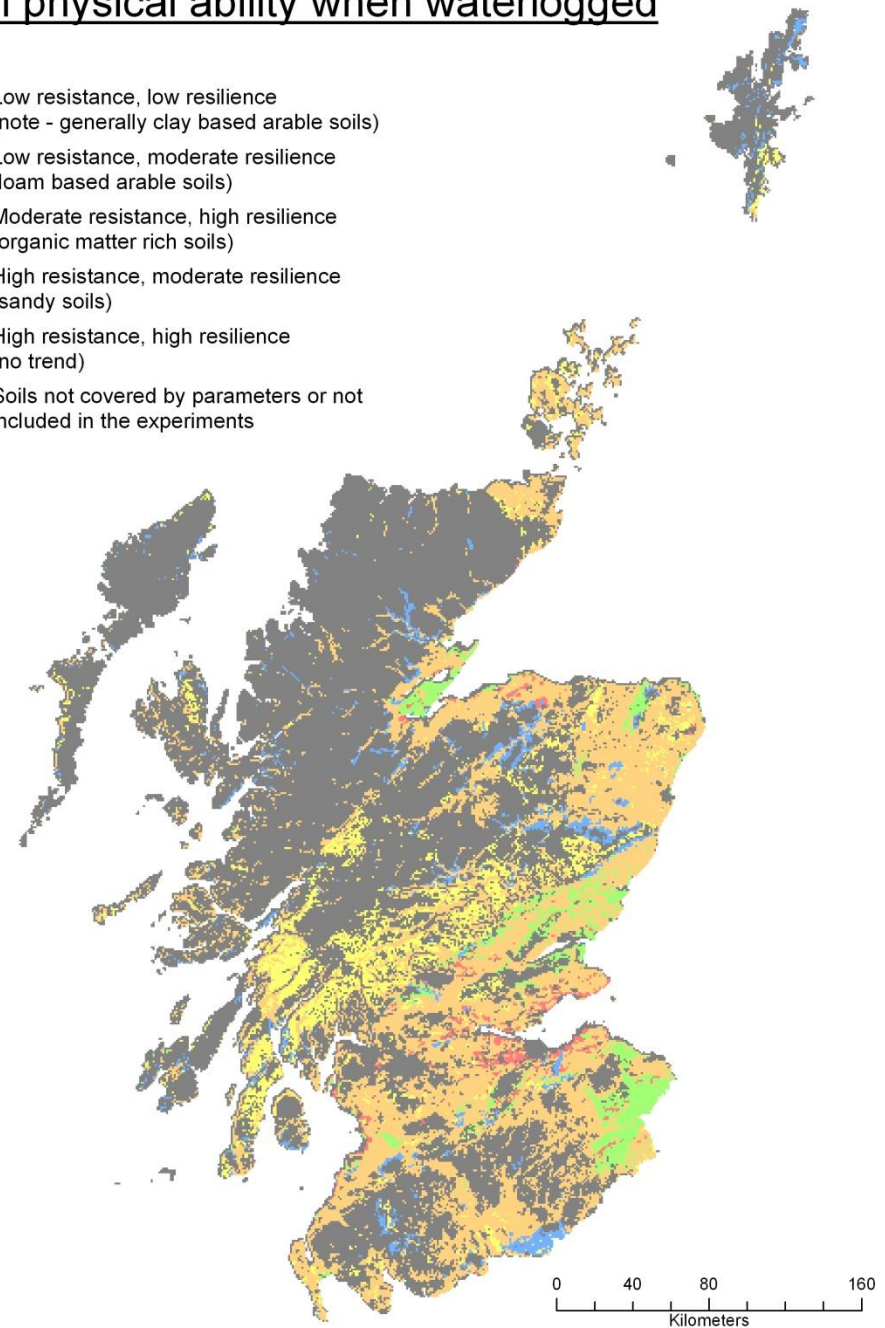
+--no: C > 2.75

+--yes: [0.739,0.788]: 2

+--no: [0.840167,0.887333]: 6

## Soil physical ability when waterlogged

- Low resistance, low resilience  
(note - generally clay based arable soils)
- Low resistance, moderate resilience  
(loam based arable soils)
- Moderate resistance, high resilience  
(organic matter rich soils)
- High resistance, moderate resilience  
(sandy soils)
- High resistance, high resilience  
(no trend)
- Soils not covered by parameters or not  
included in the experiments



# Application

## Soil physical stability to compression

### Experiment 3B

#### Independent Attributes

- o MAJOR\_SOIL\_SUBGROUP
- o Sand
- o Silt
- o Clay
- o pH
- o C
- o N
- o SOM

#### Dependent attributes

- o  $1/Cc$
- o  $Ce/Cc$

Validated Parameters:

RMSE: [0.2433,0.0164]

Correlation Coefficient: [0.708,0.4774]

$C > 5.92$

+--yes: MAJOR\_SOIL\_SUBGROUP in {Regosol,Peaty\_Podzol,Humus\_Iron\_F

| +--yes: Sand > 0.54

| | +--yes: [0.718108,0.059179]: 2

| | +--no: [0.299543,0.055084]: 2

| +--no: [0.438268,0.025836]: 2

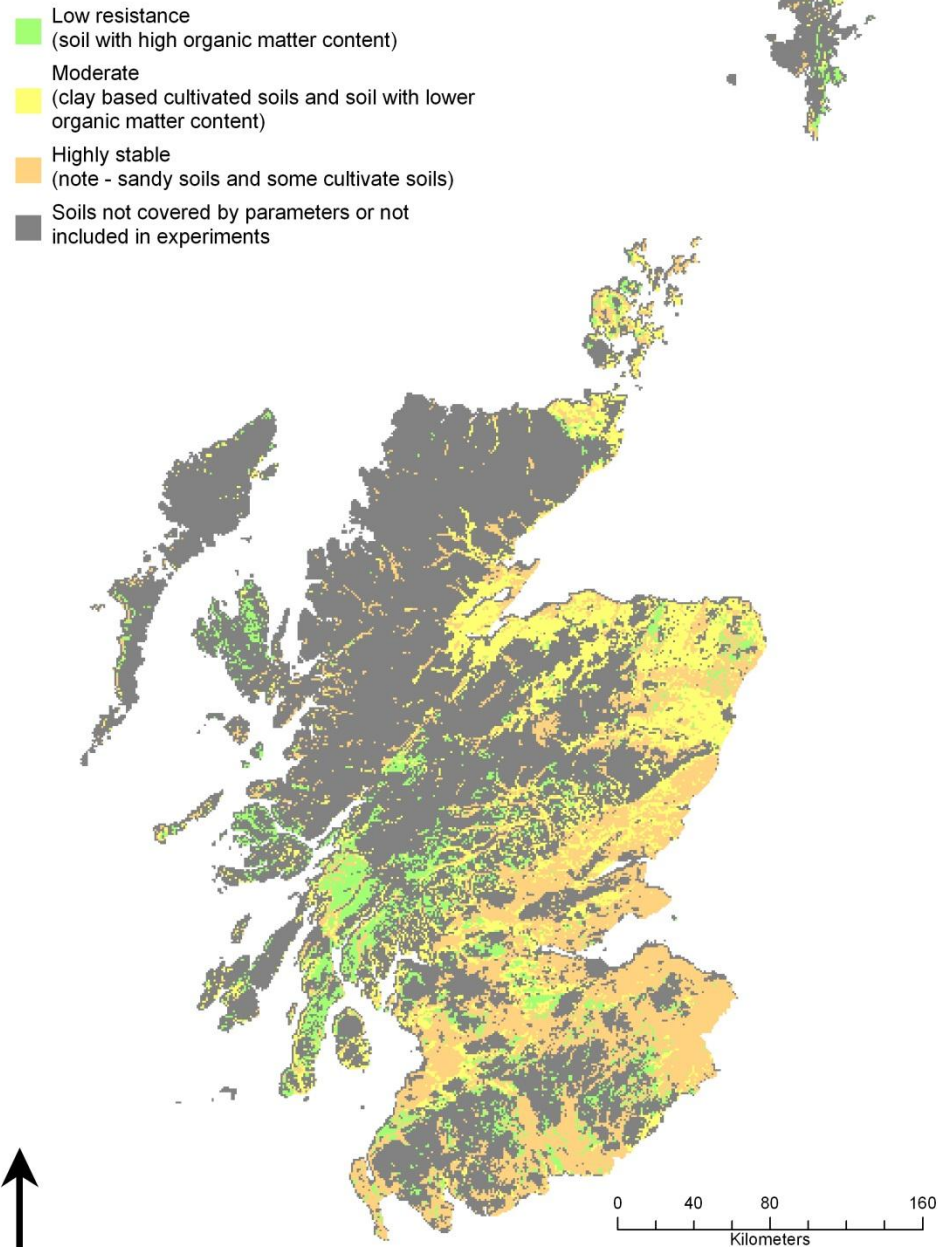
+--no: N > 0.14

+--yes: MAJOR\_SOIL\_SUBGROUP in {Brown\_Forest\_Soil,Noncalcareo

| +--yes: [0.858735,0.00968]: 13

| +--no: [0.920932,0.024988]: 4

+--no: [1.54998,0.02033]: 3



Single or multiple decision trees

Classification, regression, model trees

Propositional and relational data

Temporal and spatial data

**SUITABLE** for predictions,  
**SUITABLE** for interpretation



## 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

## 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

## 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



## 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