# Random Forest for Environmental Data Mining

Michael Leuenberger, Carmen D. Vega Orozco
and Mikhaïl Kanevski

CRET - FGSE
University of Lausanne

June 25, 2013

**Introduction**
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

# Environmental Data Mining



## Data Mining

- Search for relevant patterns for decision making

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

# Motivation

**Why Random Forest ?**

- Random Forest can deal with high dimensional database

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

## Motivation

**Why Random Forest ?**

- Random Forest can deal with high dimensional database
- It is a powerful non-linear machine learning algorithm which can deal with the high complexity of phenomena

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
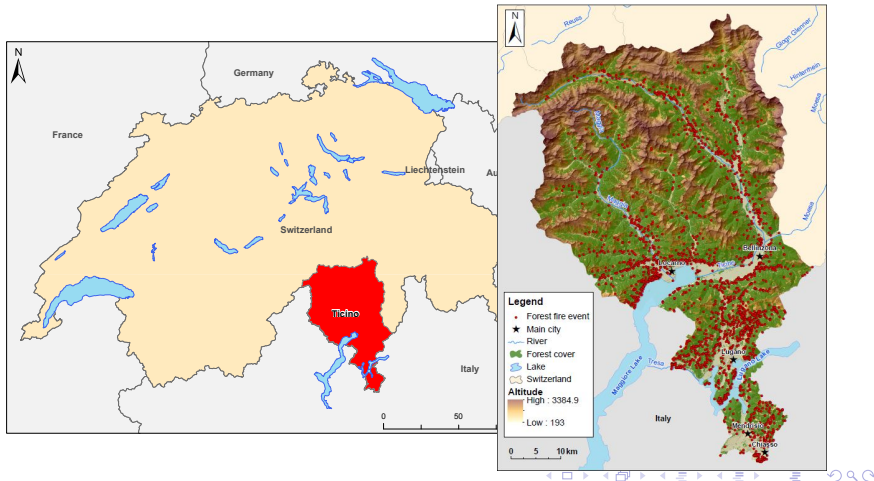Case study : Forest Fires in a Random Forest

# Motivation

**Why Random Forest ?**

- Random Forest can deal with high dimensional database
- It is a powerful non-linear machine learning algorithm which can deal with the high complexity of phenomena
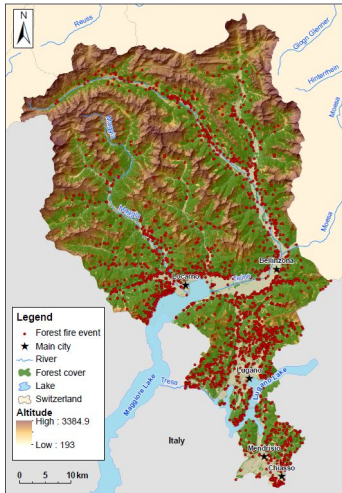- It provides directly a measure for both errors and variable importances

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

This study is focused on the 2'224 anthropogenic forest fires - ignition points - that have occured from 1969 to 2008. (*Canton Ticino, Swiss Alps, WSL-CH dataset*)

**Introduction**
Method and Methodology
Results
Conclusion

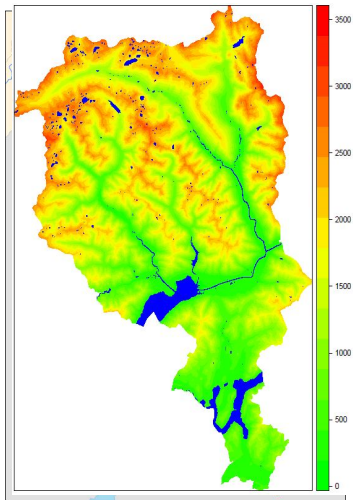Environmental Data Mining
Case study : Forest Fires in a Random Forest

This study is focused on the 2'224 anthropogenic forest fires - ignition points - that have occured from 1969 to 2008. (*Canton Ticino, Swiss Alps, WSL-CH dataset*)

Introduction
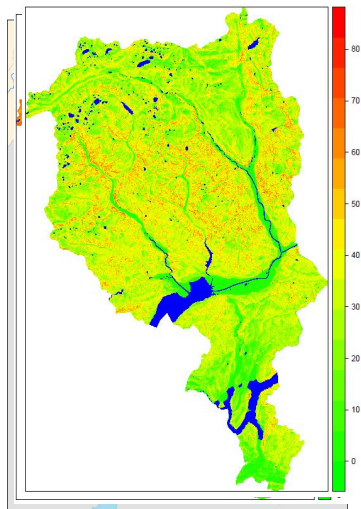Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

- Altitude

**Introduction**
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

- Altitude
- Slope

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

- Altitude
- Slope
- North aspect

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

- Altitude
- Slope
- North aspect
- West aspect

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

- Altitude
- Slope
- North aspect
- West aspect
- Dist. Streets

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest



- Altitude
- Slope
- North aspect
- West aspect
- Dist. Streets
- Dist. Hightens

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

- Altitude
- Slope
- North aspect
- West aspect
- Dist. Streets
- Dist. Hightens
- Dist. Railways

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest



- Altitude
- Slope
- North aspect
- West aspect
- Dist. Streets
- Dist. Hightens
- Dist. Railways
- Dist. Buildings

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
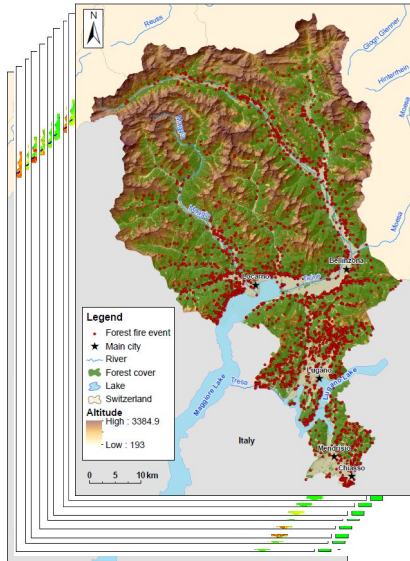Case study : Forest Fires in a Random Forest
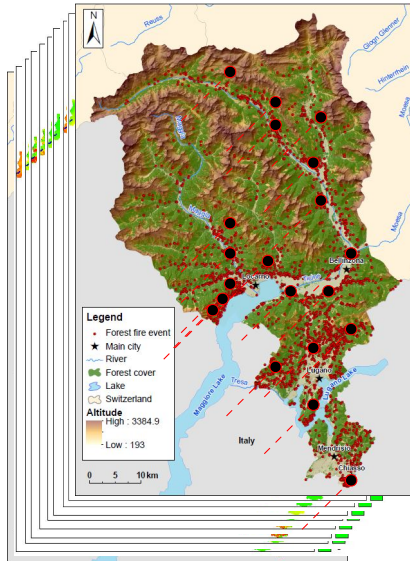
- Altitude
- Slope
- North aspect
- West aspect
- Dist. Streets
- Dist. Hightens
- Dist. Railways
- Dist. Buildings

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

Introduction
Method and Methodology
Results
Conclusion

Environmental Data Mining
Case study : Forest Fires in a Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Data

|  | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data1 | ● | ● | ● | ● | ● | ● |
| data2 | ● | ● | ● | ● | ● | ● |
| data3 | ● | ● | ● | ● | ● | ● |
| data4 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data6 | ● | ● | ● | ● | ● | ● |
| data7 | ● | ● | ● | ● | ● | ● |
| data8 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data10 | ● | ● | ● | ● | ● | ● |
| data11 | ● | ● | ● | ● | ● | ● |
| data12 | ● | ● | ● | ● | ● | ● |
| data13 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data15 | ● | ● | ● | ● | ● | ● |
| ... | ● | ● | ● | ● | ● | ● |

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

# Random Forest

Bootstrapping



| | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data1 | ● | ● | ● | ● | ● | ● |
| data2 | ● | ● | ● | ● | ● | ● |
| data3 | ● | ● | ● | ● | ● | ● |
| data4 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data6 | ● | ● | ● | ● | ● | ● |
| data7 | ● | ● | ● | ● | ● | ● |
| data8 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data10 | ● | ● | ● | ● | ● | ● |
| data11 | ● | ● | ● | ● | ● | ● |
| data12 | ● | ● | ● | ● | ● | ● |
| data13 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data15 | ● | ● | ● | ● | ● | ● |
| ... | ● | ● | ● | ● | ● | ● |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Bootstrapping



| | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data1 | • | • | • | • | • | • |
| data2 | • | • | • | • | • | • |
| data3 | • | • | • | • | • | • |
| data4 | • | • | • | • | • | • |
| data5 | • | • | • | • | • | • |
| data6 | • | • | • | • | • | • |
| data7 | • | • | • | • | • | • |
| data8 | • | • | • | • | • | • |
| data9 | • | • | • | • | • | • |
| data10 | • | • | • | • | • | • |
| data11 | • | • | • | • | • | • |
| data12 | • | • | • | • | • | • |
| data13 | • | • | • | • | • | • |
| data14 | • | • | • | • | • | • |
| data15 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

# Random Forest



$$Z_b$$

|        | V1 | V2 | V3 | V4 | V5 | ... |
|--------|----|----|----|----|----|-----|
| data2  | •  | •  | •  | •  | •  | •   |
| data4  | •  | •  | •  | •  | •  | •   |
| data5  | •  | •  | •  | •  | •  | •   |
| data5  | •  | •  | •  | •  | •  | •   |
| data7  | •  | •  | •  | •  | •  | •   |
| data9  | •  | •  | •  | •  | •  | •   |
| data9  | •  | •  | •  | •  | •  | •   |
| data10 | •  | •  | •  | •  | •  | •   |
| data11 | •  | •  | •  | •  | •  | •   |
| data14 | •  | •  | •  | •  | •  | •   |
| data14 | •  | •  | •  | •  | •  | •   |
| data16 | •  | •  | •  | •  | •  | •   |
| data17 | •  | •  | •  | •  | •  | •   |
| data19 | •  | •  | •  | •  | •  | •   |
| data19 | •  | •  | •  | •  | •  | •   |
| ...    | •  | •  | •  | •  | •  | •   |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest



$Z_b$

|        | V1 | V2 | V3 | V4 | V5 | ... |
|--------|----|----|----|----|----|-----|
| data2  | ●  | ●  | ●  | ●  | ●  | ●   |
| data4  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data7  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data10 | ●  | ●  | ●  | ●  | ●  | ●   |
| data11 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data16 | ●  | ●  | ●  | ●  | ●  | ●   |
| data17 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| ...    | ●  | ●  | ●  | ●  | ●  | ●   |

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

# Random Forest



$Z_b$

|  | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data2 | ● | ● | ● | ● | ● | ● |
| data4 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data7 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data10 | ● | ● | ● | ● | ● | ● |
| data11 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data16 | ● | ● | ● | ● | ● | ● |
| data17 | ● | ● | ● | ● | ● | ● |
| data19 | ● | ● | ● | ● | ● | ● |
| data19 | ● | ● | ● | ● | ● | ● |
| ... | ● | ● | ● | ● | ● | ● |

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

# Random Forest



$$Z_b$$

|        | V1 | V2 | V3 | V4 | V5 | ... |
|--------|----|----|----|----|----|-----|
| data2  | ●  | ●  | ●  | ●  | ●  | ●   |
| data4  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data7  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data10 | ●  | ●  | ●  | ●  | ●  | ●   |
| data11 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data16 | ●  | ●  | ●  | ●  | ●  | ●   |
| data17 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| ...    | ●  | ●  | ●  | ●  | ●  | ●   |

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

# Random Forest



$$Z_b$$

|        | V1 | V2 | V3 | V4 | V5 | ... |
|--------|----|----|----|----|----|-----|
| data2  | ●  | ●  | ●  | ●  | ●  | ●   |
| data4  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data7  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data10 | ●  | ●  | ●  | ●  | ●  | ●   |
| data11 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data16 | ●  | ●  | ●  | ●  | ●  | ●   |
| data17 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| ...    | ●  | ●  | ●  | ●  | ●  | ●   |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest



$Z_b$

|        | V1 | V2 | V3 | V4 | V5 | ... |
|--------|----|----|----|----|----|-----|
| data2  | •  | •  | •  | •  | •  | •   |
| data4  | •  | •  | •  | •  | •  | •   |
| data5  | •  | •  | •  | •  | •  | •   |
| data5  | •  | •  | •  | •  | •  | •   |
| data7  | •  | •  | •  | •  | •  | •   |
| data9  | •  | •  | •  | •  | •  | •   |
| data9  | •  | •  | •  | •  | •  | •   |
| data10 | •  | •  | •  | •  | •  | •   |
| data11 | •  | •  | •  | •  | •  | •   |
| data14 | •  | •  | •  | •  | •  | •   |
| data14 | •  | •  | •  | •  | •  | •   |
| data16 | •  | •  | •  | •  | •  | •   |
| data17 | •  | •  | •  | •  | •  | •   |
| data19 | •  | •  | •  | •  | •  | •   |
| data19 | •  | •  | •  | •  | •  | •   |
| ...    | •  | •  | •  | •  | •  | •   |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest



$Z_b$

|  | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data2 | ● | ● | ● | ● | ● | ● |
| data4 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data7 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data10 | ● | ● | ● | ● | ● | ● |
| data11 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data16 | ● | ● | ● | ● | ● | ● |
| data17 | ● | ● | ● | ● | ● | ● |
| data19 | ● | ● | ● | ● | ● | ● |
| data19 | ● | ● | ● | ● | ● | ● |
| ... | ● | ● | ● | ● | ● | ● |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest



$Z_b$

|        | V1 | V2 | V3 | V4 | V5 | ... |
|--------|----|----|----|----|----|-----|
| data2  | ●  | ●  | ●  | ●  | ●  | ●   |
| data4  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data5  | ●  | ●  | ●  | ●  | ●  | ●   |
| data7  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data9  | ●  | ●  | ●  | ●  | ●  | ●   |
| data10 | ●  | ●  | ●  | ●  | ●  | ●   |
| data11 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data14 | ●  | ●  | ●  | ●  | ●  | ●   |
| data16 | ●  | ●  | ●  | ●  | ●  | ●   |
| data17 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| data19 | ●  | ●  | ●  | ●  | ●  | ●   |
| ...    | ●  | ●  | ●  | ●  | ●  | ●   |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest



$Z_b$

|  | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data2 | ● | ● | ● | ● | ● | ● |
| data4 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data5 | ● | ● | ● | ● | ● | ● |
| data7 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data9 | ● | ● | ● | ● | ● | ● |
| data10 | ● | ● | ● | ● | ● | ● |
| data11 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data14 | ● | ● | ● | ● | ● | ● |
| data16 | ● | ● | ● | ● | ● | ● |
| data17 | ● | ● | ● | ● | ● | ● |
| data19 | ● | ● | ● | ● | ● | ● |
| data19 | ● | ● | ● | ● | ● | ● |
| ... | ● | ● | ● | ● | ● | ● |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# Random Forest

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

# Random Forest



**Bootstrapping**
**Variables selection**

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

## Random Forest



**Bootstrapping**
**Variables selection**

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
Properties

## Random Forest



**Bootstrapping**
**Variables selection**

**Random**
**Forest**

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# OOB (Out-Of-Bag)

Bootstrapping



| | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data1 | • | • | • | • | • | • |
| data2 | • | • | • | • | • | • |
| data3 | • | • | • | • | • | • |
| data4 | • | • | • | • | • | • |
| data5 | • | • | • | • | • | • |
| data6 | • | • | • | • | • | • |
| data7 | • | • | • | • | • | • |
| data8 | • | • | • | • | • | • |
| data9 | • | • | • | • | • | • |
| data10 | • | • | • | • | • | • |
| data11 | • | • | • | • | • | • |
| data12 | • | • | • | • | • | • |
| data13 | • | • | • | • | • | • |
| data14 | • | • | • | • | • | • |
| data15 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

# OOB (Out-Of-Bag)

$OOB_b$ $Z_b$

| | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data1 | • | • | • | • | • | • |
| data3 | • | • | • | • | • | • |
| data6 | • | • | • | • | • | • |
| data8 | • | • | • | • | • | • |
| data12 | • | • | • | • | • | • |
| data13 | • | • | • | • | • | • |
| data15 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

| | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| data2 | • | • | • | • | • | • |
| data4 | • | • | • | • | • | • |
| data5 | • | • | • | • | • | • |
| data5 | • | • | • | • | • | • |
| data7 | • | • | • | • | • | • |
| data9 | • | • | • | • | • | • |
| data9 | • | • | • | • | • | • |
| data10 | • | • | • | • | • | • |
| data11 | • | • | • | • | • | • |
| data14 | • | • | • | • | • | • |
| data14 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
**Properties**

## Measures in Random Forest

### OOB error estimation

For each data we compute error over the trees where she appear OOB ($E_i$). By averaging we obtain an estimation of the global error.

$$OBB_{error} = \frac{1}{N} \sum_{i=1}^{N} E_i$$

### Variable Importance Measurement

For each tree we compute the error obtained by the OOB data ($E_b$). Then we recompute the error after permuting all values in one variable ($m$) for all OOB data ($E_b^m$). By averaging over all the trees we obtain :

$$I_m = \frac{1}{B} \sum_{b \in B} E_b^m - E_b$$

Introduction
Method and Methodology
Results
Conclusion

Random Forest
Properties

## Illustrations

**Error assessment**

| Scn | X, Y (variables) | MSE test set (sd) | MSE oob (sd) |
|-----|------------------|-------------------|--------------|
| A1 | Yes (all variables) | 1.309 ($\pm$0.0026) | 1.299 ($\pm$0.0043) |
| A3 | Yes (just 8 best) | 1.334 ($\pm$0.0041) | 1.316 ($\pm$0.0050) |
| B1 | No (all variables) | 1.316 ($\pm$0.0041) | 1.333 ($\pm$0.0039) |
| B3 | No (just 5 best) | 1.363 ($\pm$0.0046) | 1.381 ($\pm$0.0032) |

- MSE test $\approx$ MSE oob
- $MSE_{A1}, ..., MSE_{B3} \in [1.309, 1.363]$
- Large decrease in number of variables $\implies$ slight increase in the error

Introduction
**Method and Methodology**
Results
Conclusion

Random Forest
**Properties**

# Illustrations

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

Introduction
Method and Methodology
**Results**
Conclusion

**Prediction**
Susceptibility map

|            | V1 | V2 | V3 | ... |
|------------|:--:|:--:|:--:|:---:|
| Newdata1   | ●  | ●  | ●  | ●   |
| Newdata2   | ●  | ●  | ●  | ●   |
| Newdata3   | ●  | ●  | ●  | ●   |
| Newdata4   | ●  | ●  | ●  | ●   |
| Newdata5   | ●  | ●  | ●  | ●   |
| Newdata6   | ●  | ●  | ●  | ●   |
| Newdata7   | ●  | ●  | ●  | ●   |
| Newdata8   | ●  | ●  | ●  | ●   |
| ...        | ●  | ●  | ●  | ●   |

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Prediction



|          | V1 | V2 | V3 | V4 | V5 | ... |
|----------|----|----|----|----|----|-----|
| Newdata1 | •  | •  | •  | •  | •  | •   |
| Newdata2 | •  | •  | •  | •  | •  | •   |
| Newdata3 | •  | •  | •  | •  | •  | •   |
| Newdata4 | •  | •  | •  | •  | •  | •   |
| Newdata5 | •  | •  | •  | •  | •  | •   |
| Newdata6 | •  | •  | •  | •  | •  | •   |
| Newdata7 | •  | •  | •  | •  | •  | •   |
| ...      | •  | •  | •  | •  | •  | •   |

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Prediction



| | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| Newdata1 | • | • | • | • | • | • |
| Newdata2 | • | • | • | • | • | • |
| Newdata3 | • | • | • | • | • | • |
| Newdata4 | • | • | • | • | • | • |
| Newdata5 | • | • | • | • | • | • |
| Newdata6 | • | • | • | • | • | • |
| Newdata7 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Prediction



|  | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| Newdata1 | • | • | • | • | • | • |
| Newdata2 | • | • | • | • | • | • |
| Newdata3 | • | • | • | • | • | • |
| Newdata4 | • | • | • | • | • | • |
| Newdata5 | • | • | • | • | • | • |
| Newdata6 | • | • | • | • | • | • |
| Newdata7 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Prediction

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Prediction

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Prediction



|  | V1 | V2 | V3 | V4 | V5 | ... |
|---|---|---|---|---|---|---|
| Newdata1 | • | • | • | • | • | • |
| Newdata2 | • | • | • | • | • | • |
| Newdata3 | • | • | • | • | • | • |
| Newdata4 | • | • | • | • | • | • |
| Newdata5 | • | • | • | • | • | • |
| Newdata6 | • | • | • | • | • | • |
| Newdata7 | • | • | • | • | • | • |
| ... | • | • | • | • | • | • |

Introduction
Method and Methodology
Results
Conclusion

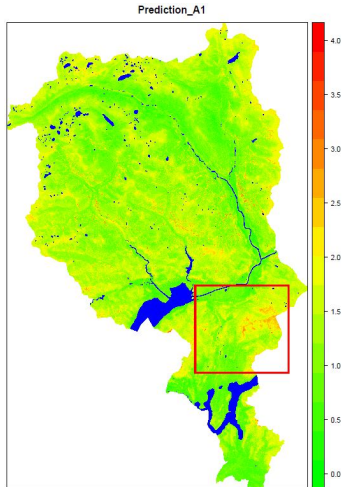Prediction
Susceptibility map

# Scenario A1



Prediction_A1

**15 Variables :**

- Altitude (DEM)
- X,Y Coordinates
- Slope
- North, West aspect
- Dist Streets, Pathways
- Dist Railways, Highways
- Dist Hightens, Building
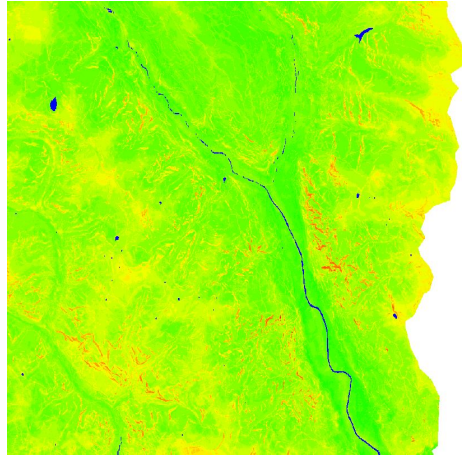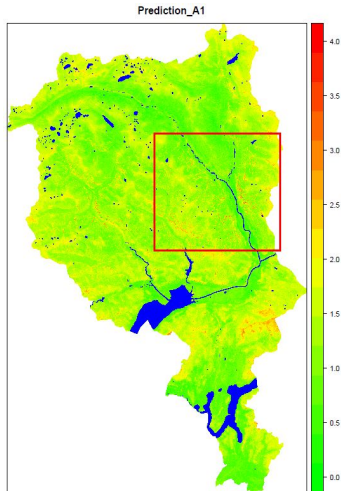- Dist Forest, Vineyard
- Primary surface

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Scenario A1



Prediction_A1

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

## Scenario A1



Prediction_A1

Introduction
Method and Methodology
**Results**
Conclusion

Prediction
Susceptibility map

# Scenario A1

# Conclusion

**Future research**

- Further analysis to optimize hyper-parameters

# Conclusion

**Future research**

- Further analysis to optimize hyper-parameters
- Construction of larger input space

# Conclusion

**Future research**

- Further analysis to optimize hyper-parameters
- Construction of larger input space
- Application of other ensemble learning methods and one-class classification methods

# Conclusion

**Future research**

- Further analysis to optimize hyper-parameters
- Construction of larger input space
- Application of other ensemble learning methods and one-class classification methods
- Application to other case studies : Landslide, Permafrost, Avalanche,...

**Thank you for your attention !**