# Practical Subgroup Discovery

**Janis Kalofolias**                                based on work by Mario Boley
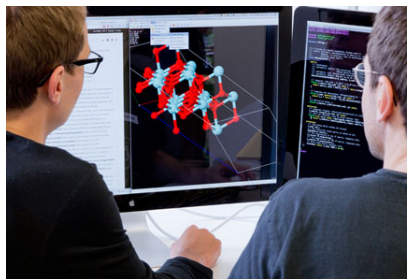
max planck institut
informatik

UNIVERSITÄT
DES
SAARLANDES

CISPA
HELMHOLTZ CENTER FOR
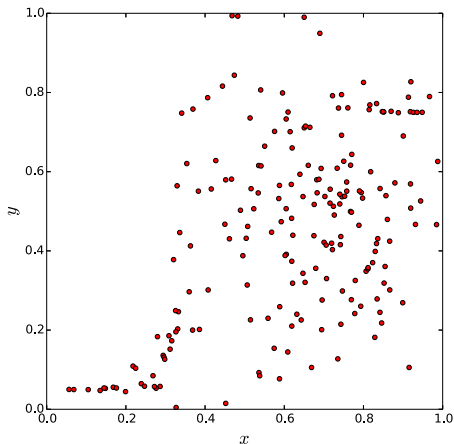INFORMATION SECURITY

# Two flavors of data science

**Predictive modelling**



**Exploratory data analysis**

# Global versus local modelling

**Given**
Population $P = \{1, \ldots, n\}$
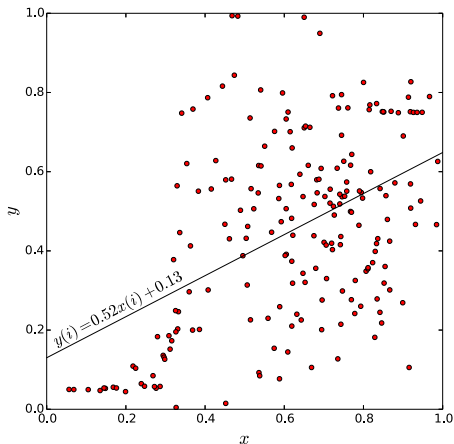Target variable $y \colon P \to \mathbb{R}$
Description variable $x \colon P \to \mathbb{R}$

**What can we tell about $y$?**

(in terms of $x$)

# Global versus local modelling

**Given**
Population $P = \{1, \ldots, n\}$
Target variable $y: P \to \mathbb{R}$
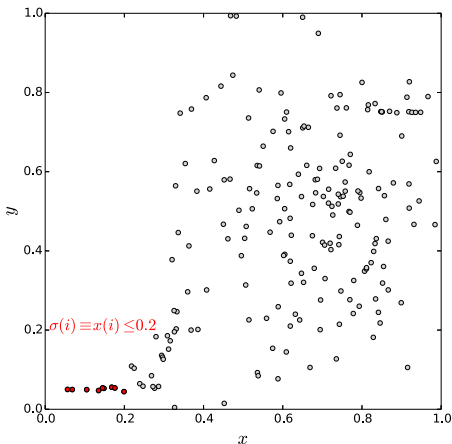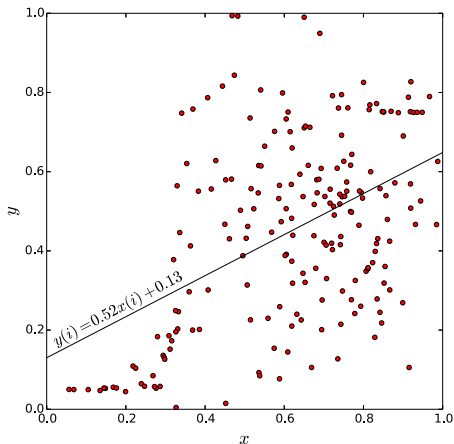Description variable $x: P \to \mathbb{R}$

**Find**
Coefficients $\alpha, \beta \in \mathbb{R}$ such that objective function

$$f(\alpha, \beta) = \frac{1}{n} \sum_{i \in P} (\alpha x(i) + \beta - y(i))^2 + \lambda \|\alpha, \beta\|_1$$
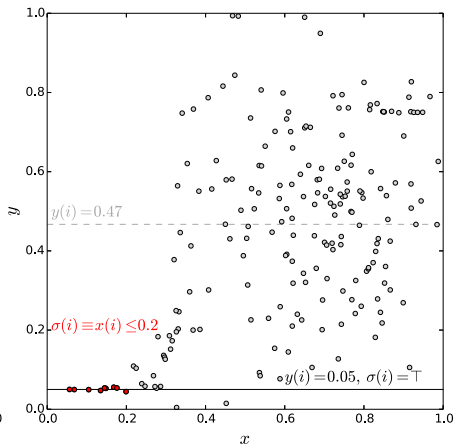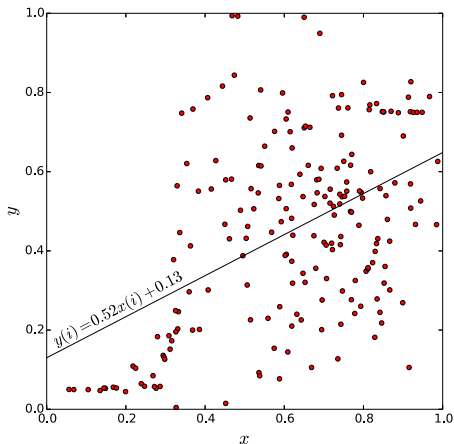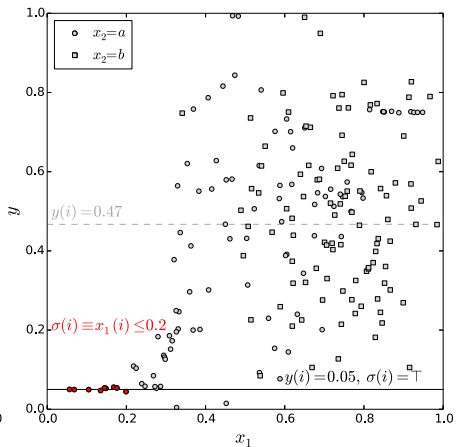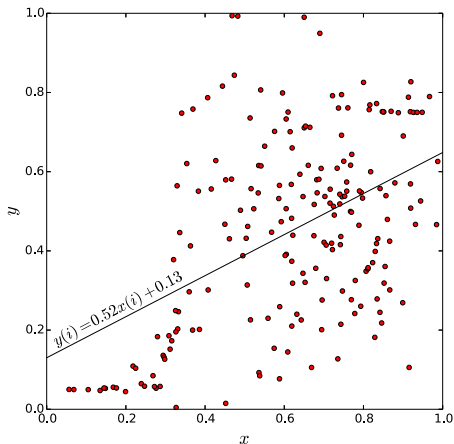
is minimal

# Global versus local modelling

$y(i) = 0.52x(i) + 0.13$

$\sigma(i) \equiv x(i) \leq 0.2$

# Global versus local modelling

$y(i) = 0.52x(i) + 0.13$

$y(i) = 0.47$

$\sigma(i) \equiv x(i) \leq 0.2$

$y(i) = 0.05, \ \sigma(i) = \top$

# Global versus local modelling

Left plot: $y(i) = 0.52 x(i) + 0.13$

Right plot legend: $x_2 = a$, $x_2 = b$

$y(i) = 0.47$

$\sigma(i) \equiv x_1(i) \leq 0.2$

$y(i) = 0.05,\ \sigma(i) = \top$

# Global versus local modelling

# Practical example

**Background**
Fraction of non-voters increased in German federal election 2009 to 0.28 (from 0.21 in 2005)

**Question**
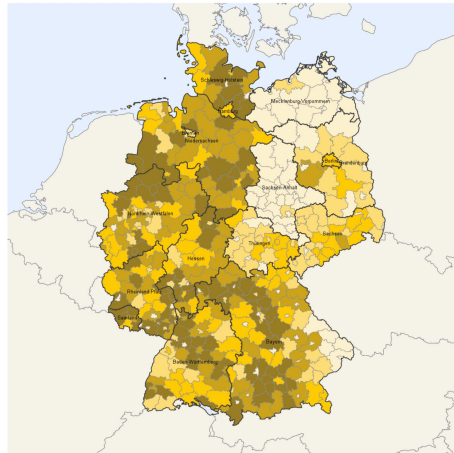"Where" did increase come from?

**Data**
Population: `admin. districts of Germany`
Target variable: `non-voter diff. 2005-2009`
Description attributes:
- Geographical (`region, state`)
- Demographic (`pop. density, highsch. degrees, …`)
- Economic (`GDP growth, web domains,…`)



©2017, Mario Boley

# Practical example

**Background**
Fraction of non-voters increased in German federal election 2009 to 0.28 (from 0.21 in 2005)

**Question**
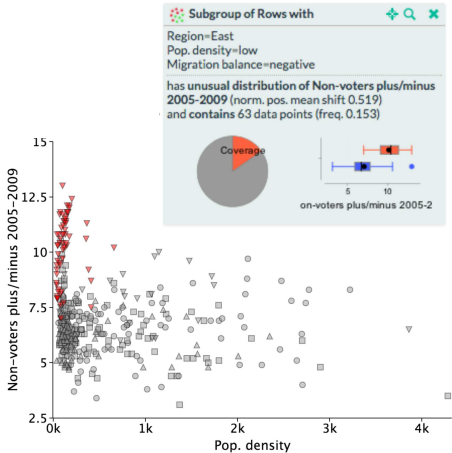"Where" did increase come from?

**Data**
Population: `admin. districts of Germany`
Target variable: `non-voter diff. 2005-2009`
Description attributes:
- Geographical (`region, state`)
- Demographic (`pop. density, highsch. degrees, …`)
- Economic (`GDP growth, web domains,…`)

- VIKAMINE

  More for computer scientists: (communities, simpler patterns)
- RealKD/Creedo `https://bitbucket.org/realKD`

  Multiple algorithms, several objectives

  From a material scientist for material scientists

## Getting RealKD/Creedo

A docker image is available

- Getting the image
  Option a) Build:
  ```
  docker build -t kalofoli/creedo-deps \
      git@github.com:kalofoli/docker.git#:creedo-deps
  docker build -t kalofoli/creedo \
      git@github.com:kalofoli/docker.git#:creedo
  ```

  Option b) Download: ($\approx 1.3$GB)
  ```
  docker create -t kalofoli/creedo
  ```

- Running docker
  ```
  docker run -it -p 8080:8080/tcp  kalofoli/creedo
  ```

# Subgroup Discovery with Creedo

- Open browser at: `http://localhost:8080/Creedo`
- login with User: default and empty password
- (Optional) Upload xarf data file
- Click Analyze

Settings



Titanic



Germany

## Automating Experiments

Using RealKD Job descriptions

```
 1   {
 2   "type" : "productWorkScheme",
 3   "id" : "octet_binaries_fdd",
 4   "workspaces" : [ {
 5     "type" : "workspaceFromXarf",
 6     "id" : "binaries",
 7     "datafile" : "octet_binaries_2.1.1.xarf"
 8   } ],
 9   "computations" : [ {
10     "type" : "functionalDependencyDiscovery",
11     "id" : "titanic_functional_pattern_discovery
            ",
12     "target" : "sign_delta_e",
13     "num_res" : 3,
14     "alpha" : 1.0
15   } ],
16   "computationTimeLimit" : 3600
17   }
```