# Using Road Sensor Data for Official Statistics

## Towards a Big Data Methodology

Marco Puts, Martijn Tennekes, Piet Daas

# Road sensors

**Road sensor data (NDW)**
– Passing vehicle counts for each minute (24/7) at about 60.000 sensors in the Netherlands
– Types of sensors:
  - Induction loop
  - Camera
  - Bluetooth
– Length categories (e.g. small, medium, long vehicles)
– Large volume: approx. 230 mln records/day

# Challenges at Statistics Netherlands

**Volume**

- How to deal with large volumes of data?

**Historical time series**

- How to create a historical time series?

**Accuracy**

- Can we create accurate statistics based on this data?

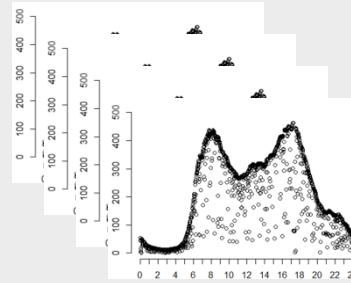**Representativity**

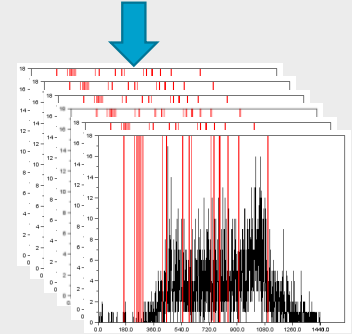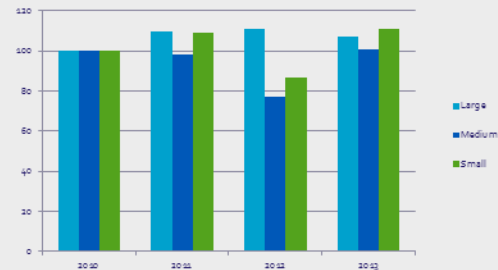- Loops are not homogeneous distributed.

# Statistical Process

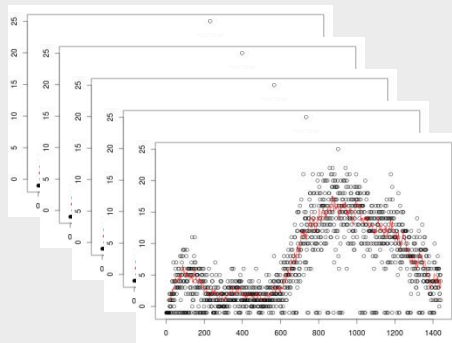# Process

Raw Data
**4 TB**

2010 - 2014

Transform
+
Select

Transformed
Data
**10 GB**

Cleaning

Clean Data
**500 MB**
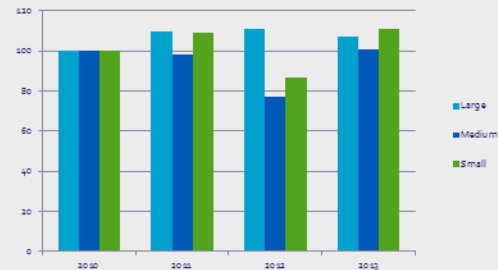
Road Network

Framing

Estimation

Traffic Index
**6 KB**

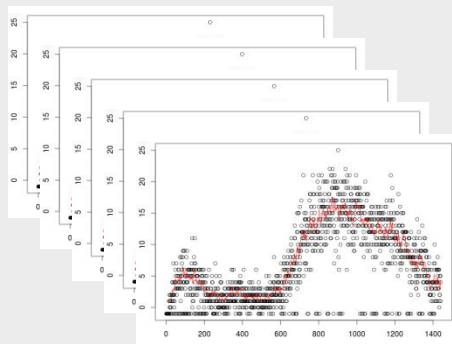# Statistical Process
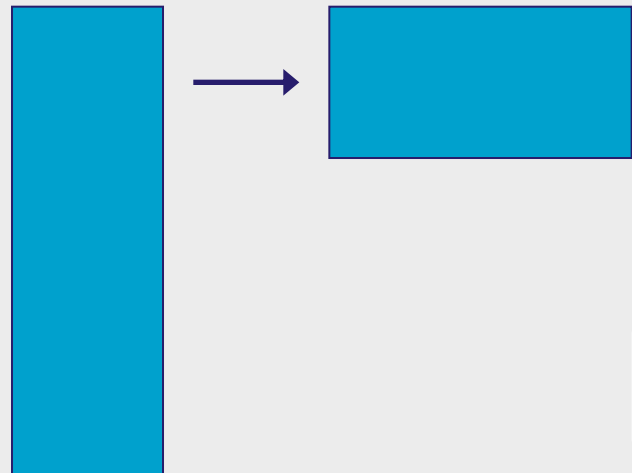
# Transform + Select

Reduce the Volume of the Data

–Select

- Only necessary variables

- Only valid data

- On the main routes (without ramps and interchanges)
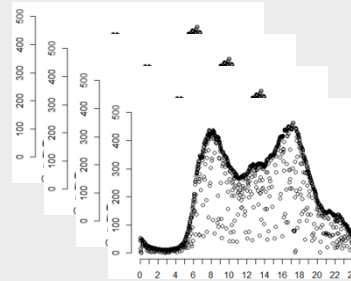
–Transform
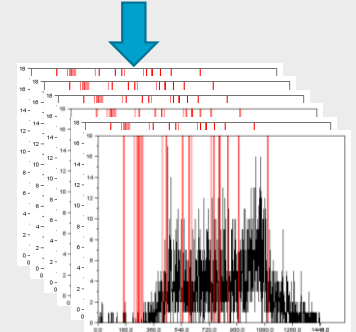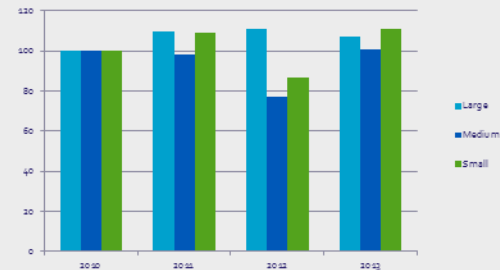
- Put one day in one record

# Statistical Process

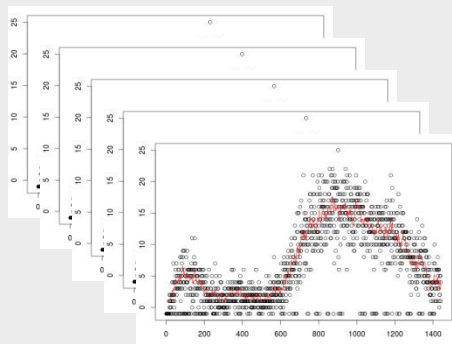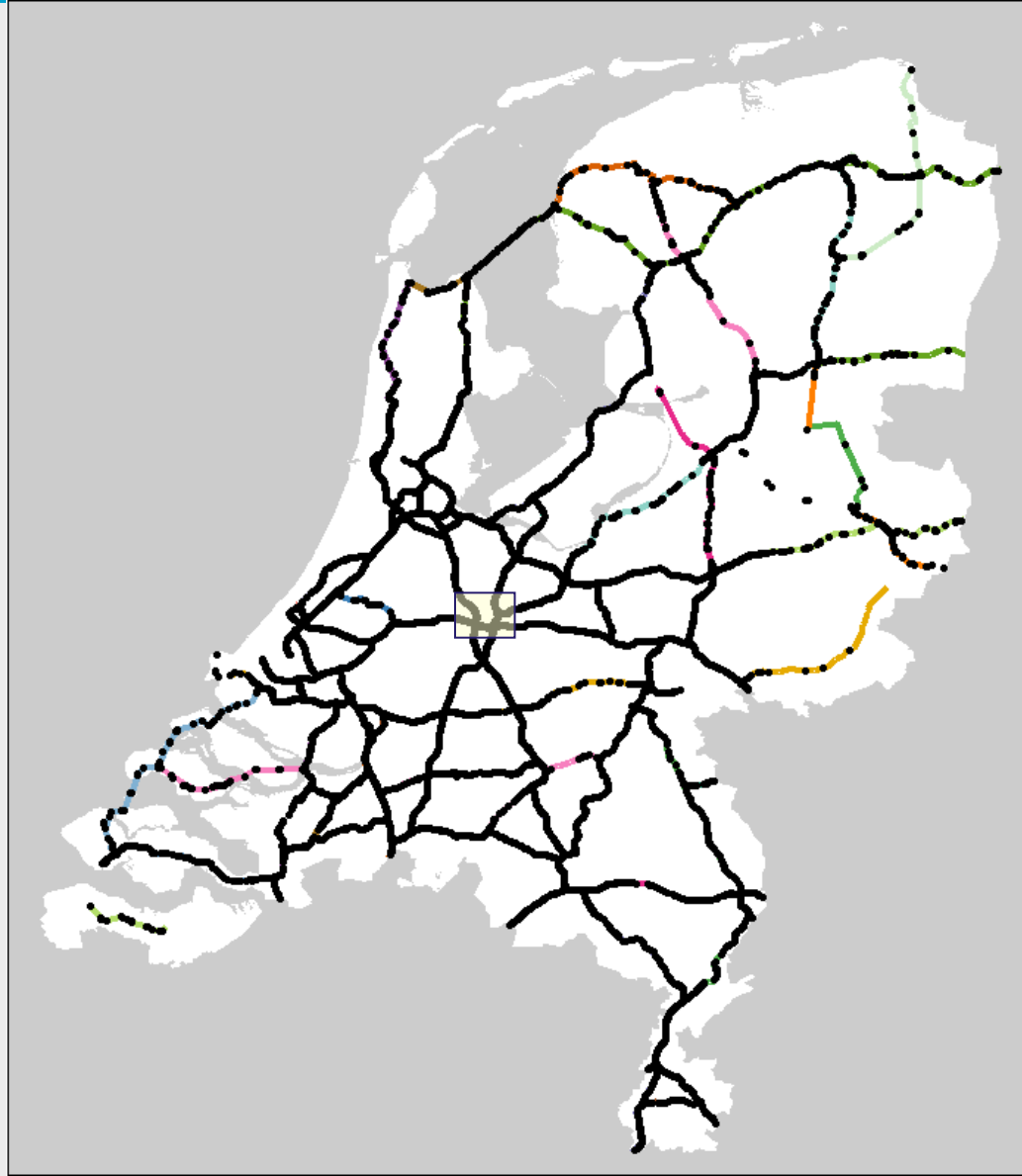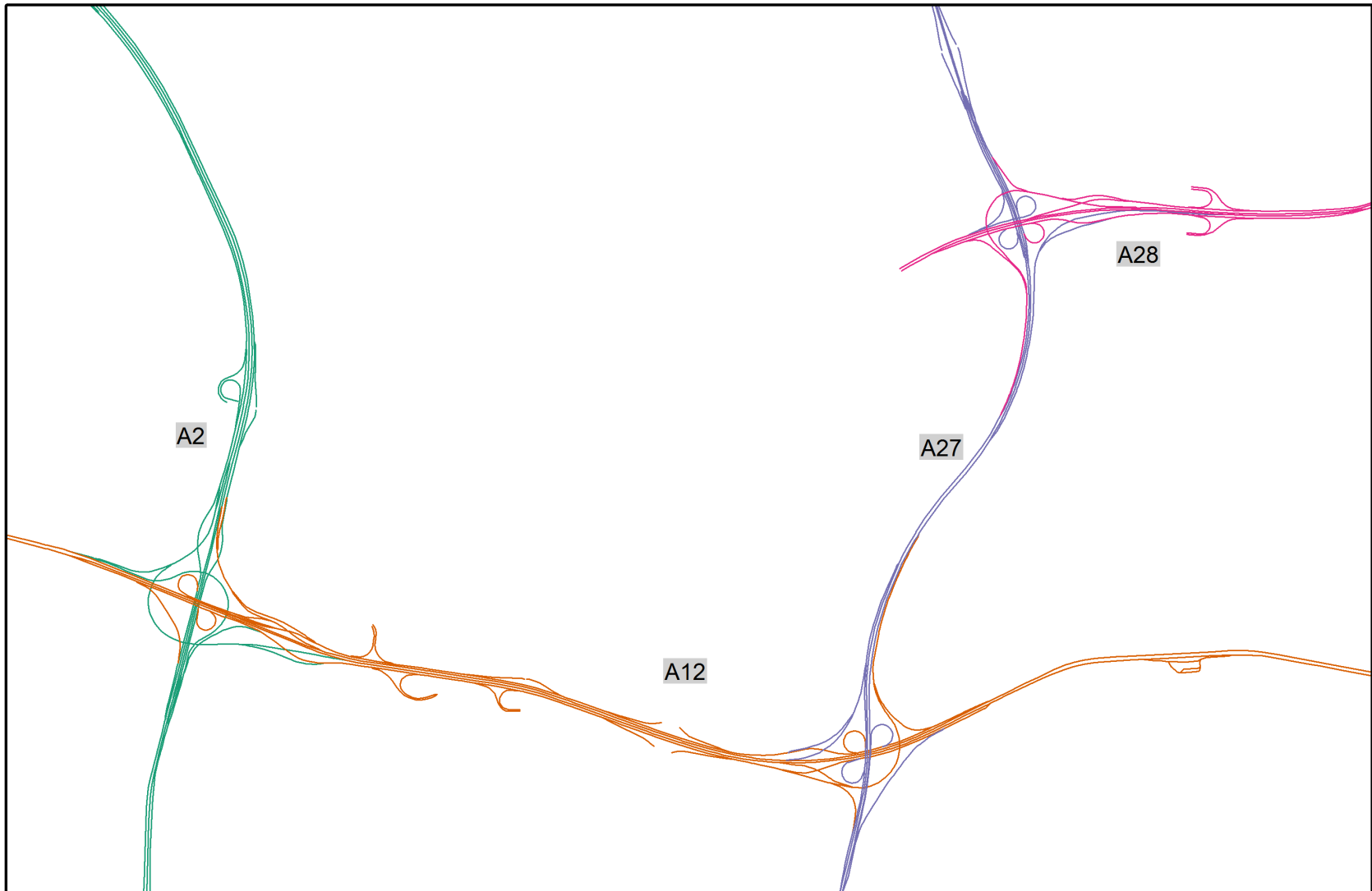# Dutch highways

# Dutch highways with road sensors

# A closer look...

# Road selection

- Dutch Highways
- Main routes (no interchange, entrance and exit ramps)
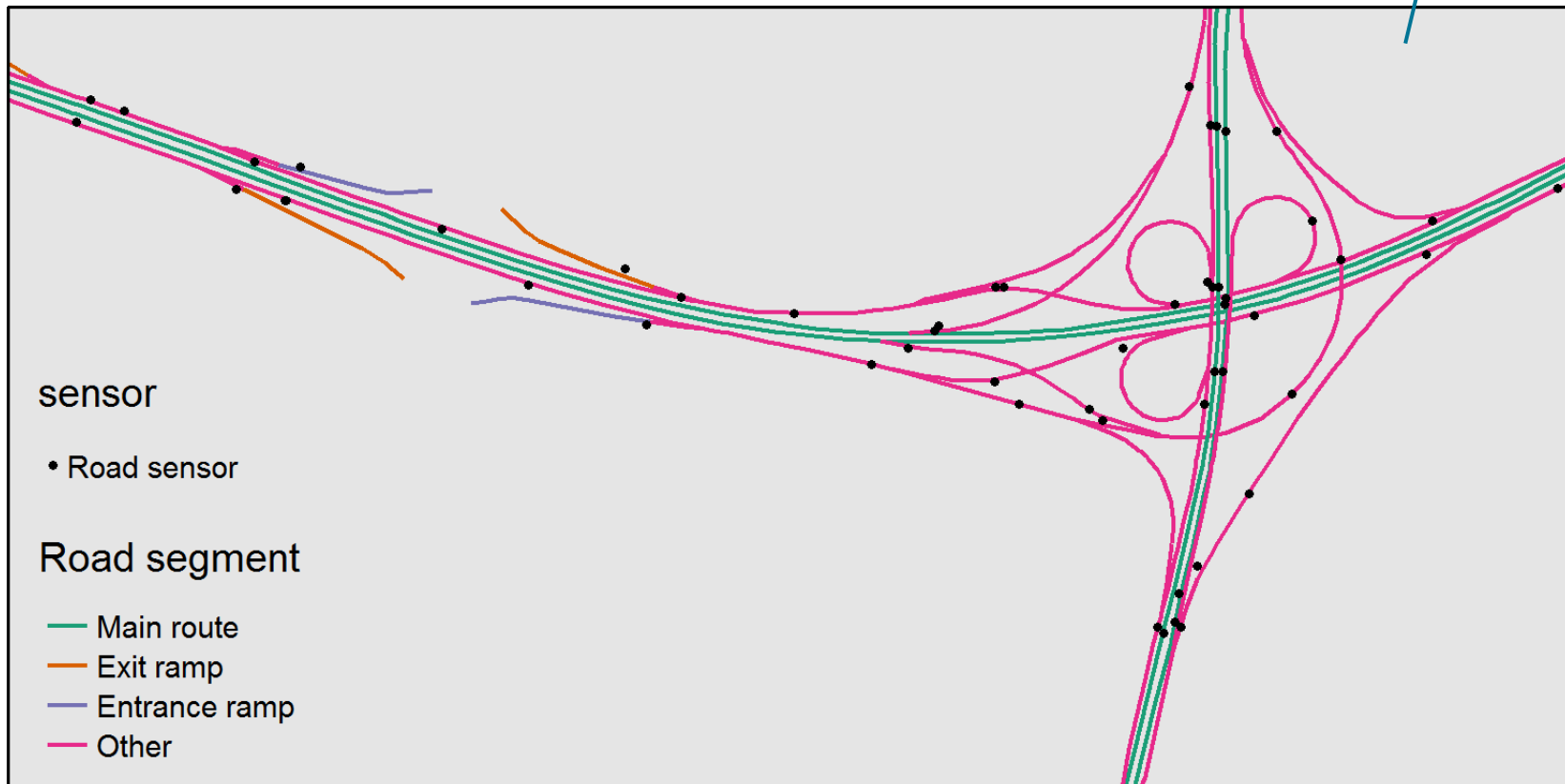


sensor

• Road sensor

Road segment

— Main route
— Exit ramp
— Entrance ramp
— Other

# Metadata input

– ESRI shape file of Dutch roads



– Road sensor metadata

| Road | Direction | Type | Lat | Long |
|------|-----------|------|---------|--------|
| A79 | West | Main | 50.8779 | 5.7502 |
| A79 | West | Main | 50.8772 | 5.7625 |
| A79 | West | Main | 50.8768 | 5.7737 |
| A79 | West | Main | 50.8747 | 5.8082 |
| A79 | West | Main | 50.8828 | 5.8650 |
| … | … | … | … | … |

# Map projection
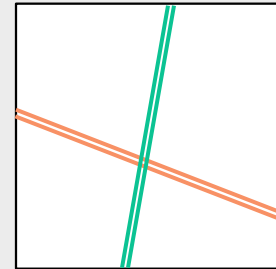


Amersfoort: projection centre

- Dutch National Grid (Rijksdriehoekstelsel)
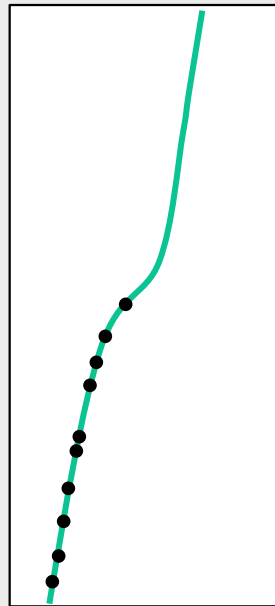- Preserves real-world distances

# Main routes
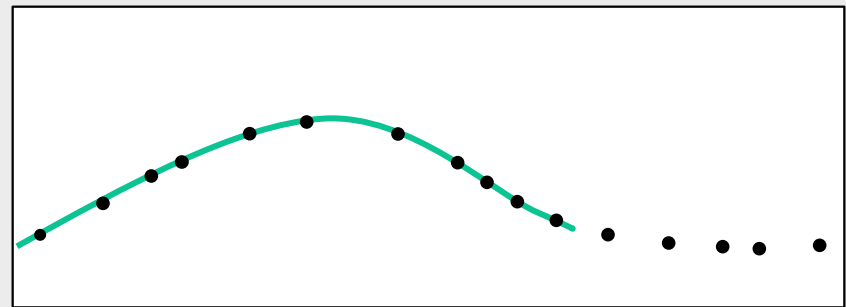


Raw shape

Simplify →

Main routes

# **Metadata inconsistencies**



No road sensors?

Where is the road?

Possible causes:
- Errors in metadata
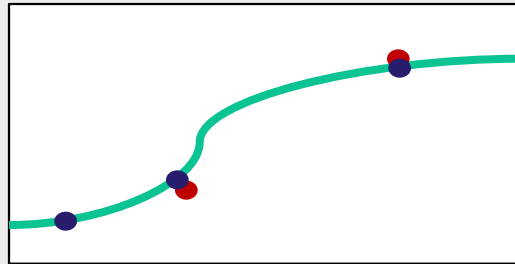- Different time references
- Different definitions

Solutions:
- Shape is leading:    Impute empty part    Remove loose road sensors
- Sensors are leading:    Cut off empty part    Extrapolate main route
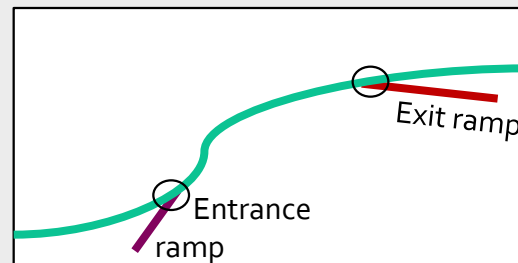
# Projections

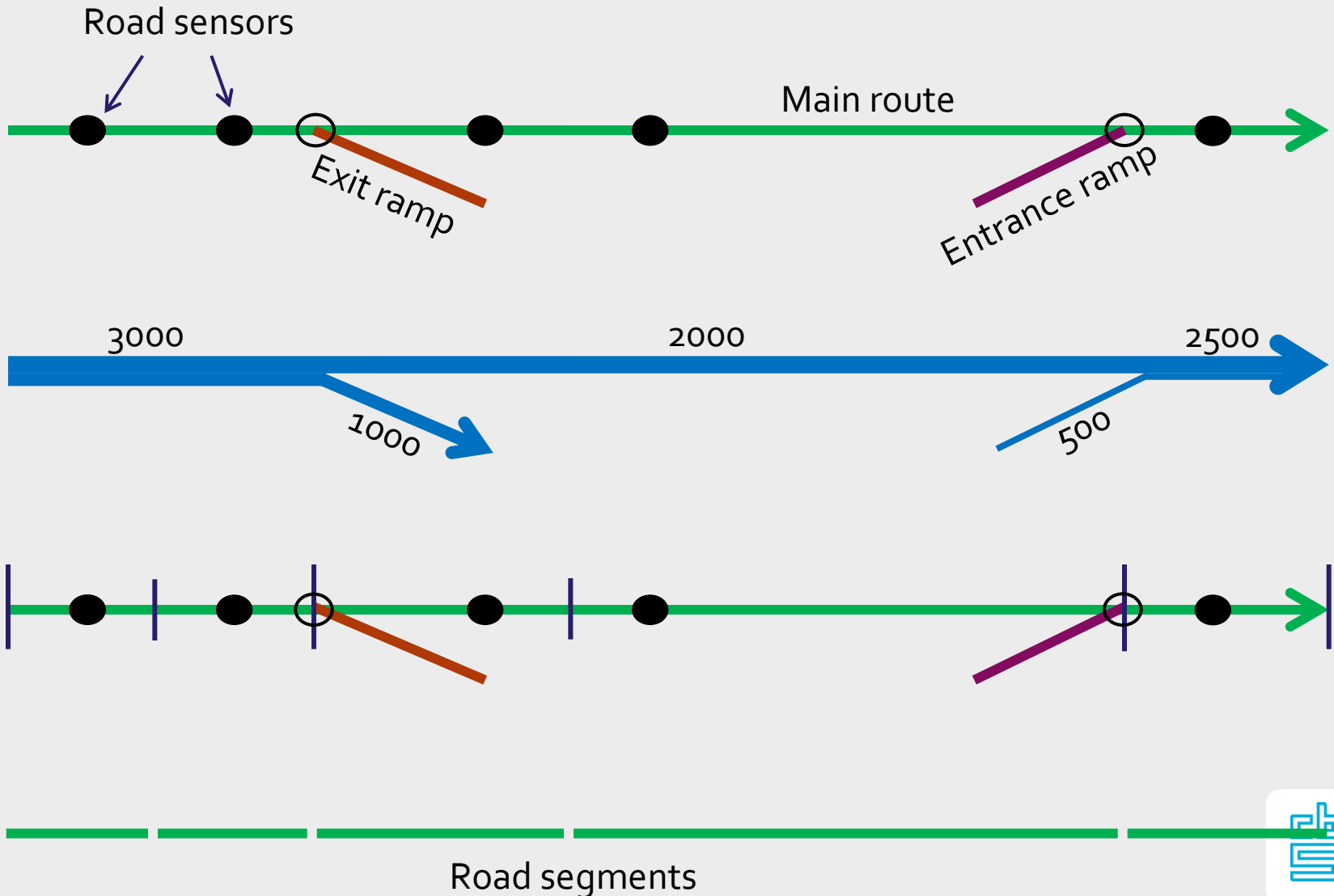– Project road sensors on main routes



– Determine points of bifurcation for all entrance and exit ramps
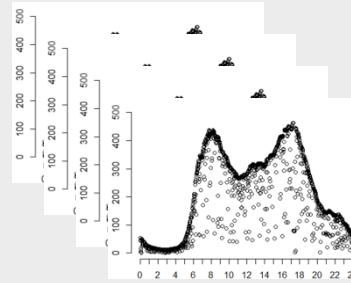
# Metadata output: road segments

Road sensors

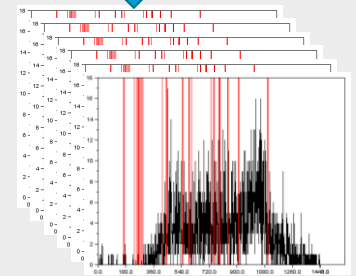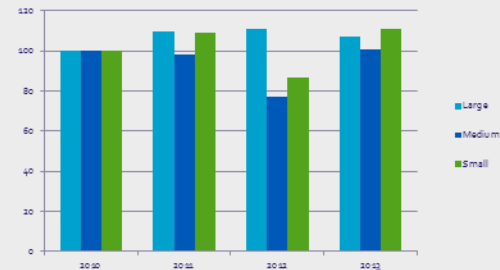Main route

Exit ramp

Entrance ramp

3000

2000

2500

1000

500

Road segments

# Statistical Process

# Cleaning the data

# Cleaning the data
# Hidden Markov Model

state

$X_1$ → $X_2$ → $X_3$ → $X_4$ →

$Y_1$   $Y_2$   $Y_3$   $Y_4$

observation

# Cleaning the Data
# Recursive Bayesian Estimation



state

$X_1$ → $X_2$ → $X_3$ → $X_4$

$Y_1$ $Y_2$ $Y_3$ $Y_4$

observation

Update

# Cleaning the Data
# Recursive Bayesian Estimation



state

observation

Prediction

24

# Cleaning the Data
# Recursive Bayesian Estimation



state

$X_1$   $X_2$   $X_3$   $X_4$

$Y_1$   $Y_2$   $Y_3$   $Y_4$

observation

Missing Data

# Cleaning the Data
# Recursive Bayesian Estimation



state

$X_1$ → $X_2$ → $X_3$ → $X_4$

$Y_1$ $Y_2$ $Y_3$ $Y_4$

observation

Smoothing

# Cleaning the Data
# Recursive Bayesian Estimation

**Assumption**:

– Arrival times of vehicles follow a Poisson Process

– Gaussian Random Walk

**Algorithm:**

– Discretization of Probability Density Function

**Advantage**:

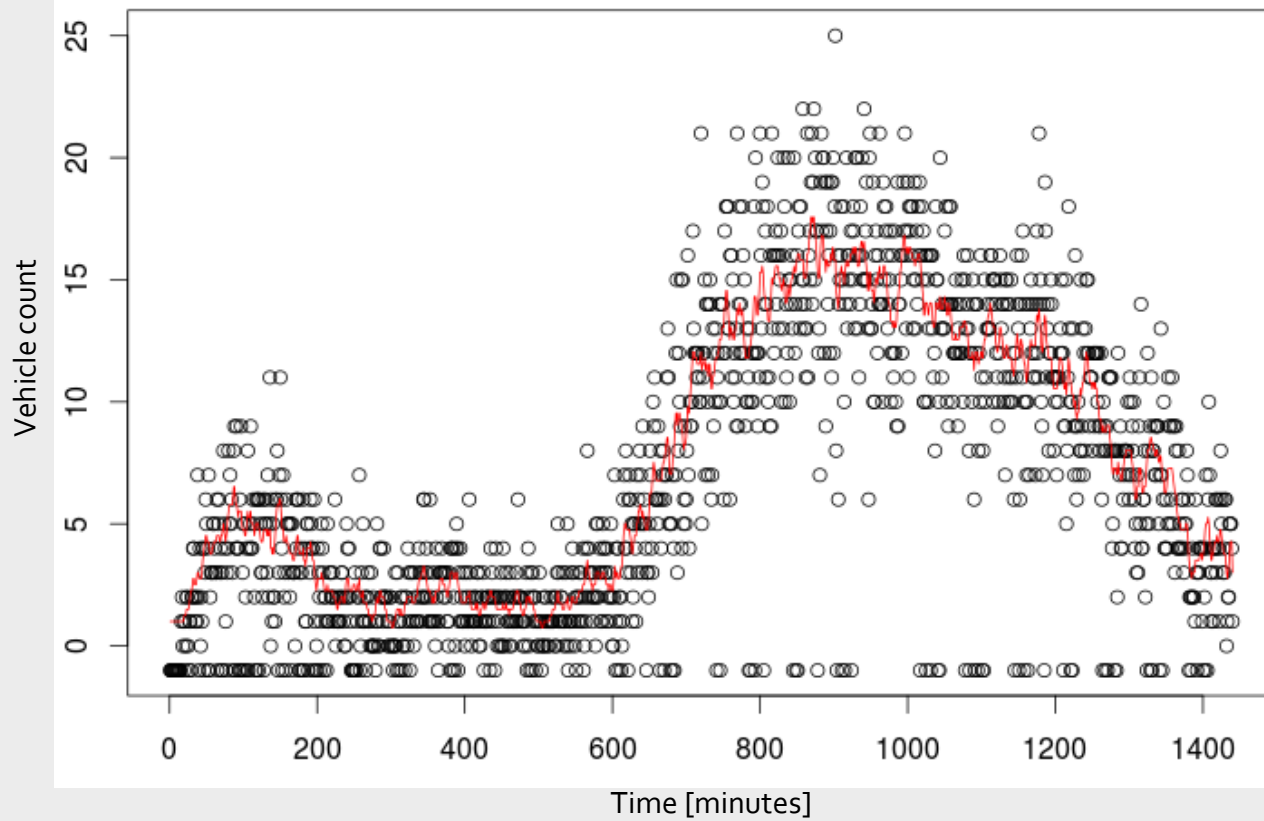– High Accuracy

**Disadvantage**:

– Slow… (due to convolutions)

# Cleaning the Data
# Speeding Up Things

**Use Fuzzy Logic**

– Discrete PDF => Membership Function
– Convolutions => Dilation operators

Vehicle count vs Time [minutes]
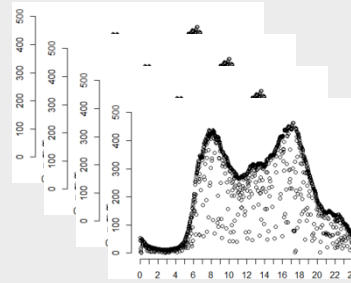
# Precision/Accuracy

The filter does not introduce extra errors:
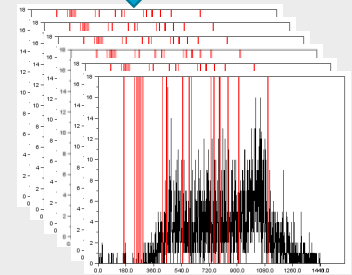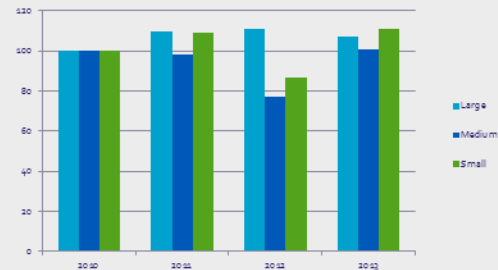
– Precision: 3.6%
– Accuracy:+0.13%

# Final Process

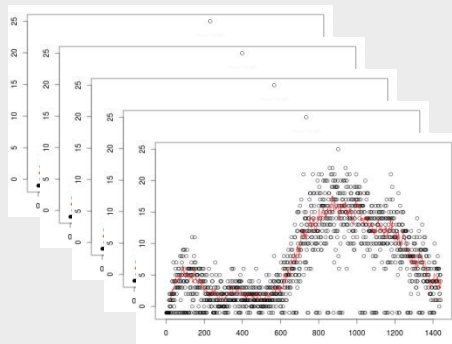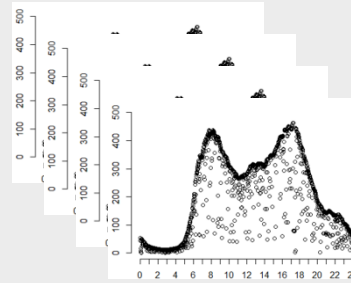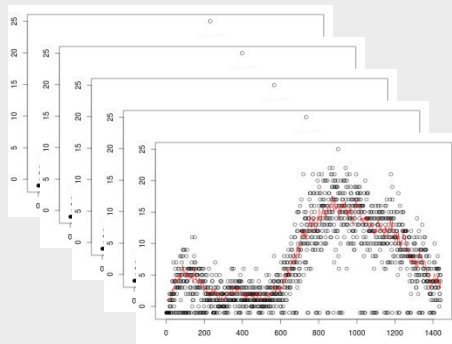# Estimation

# Statistical Process



Frame

Transform & Select

Clean

Estimate

# Questions?