



Automatic measured data validation applied on hydraulic and water quality parameters in sewer systems

Nemanja Branislavljević¹

Dušan Prodanović¹

Zoran Kapelan²

¹University of Belgrade

²University of Exeter

Research challenges

Challenge 1:

Large amount of data can not be validated traditionally

Challenge 2:

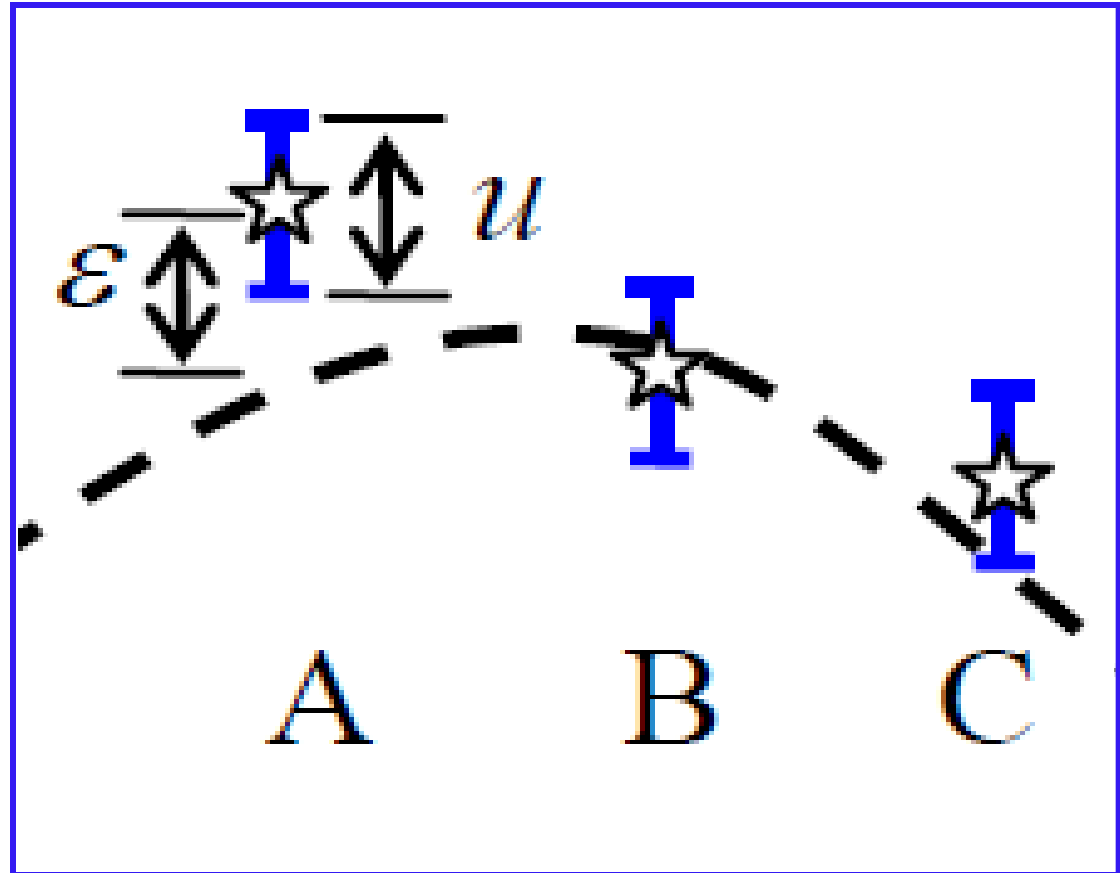
Data validation procedure have to be capable to use any kind of information available

Challenge 3:

Data validation procedure have to be independent on type or number of relations between data

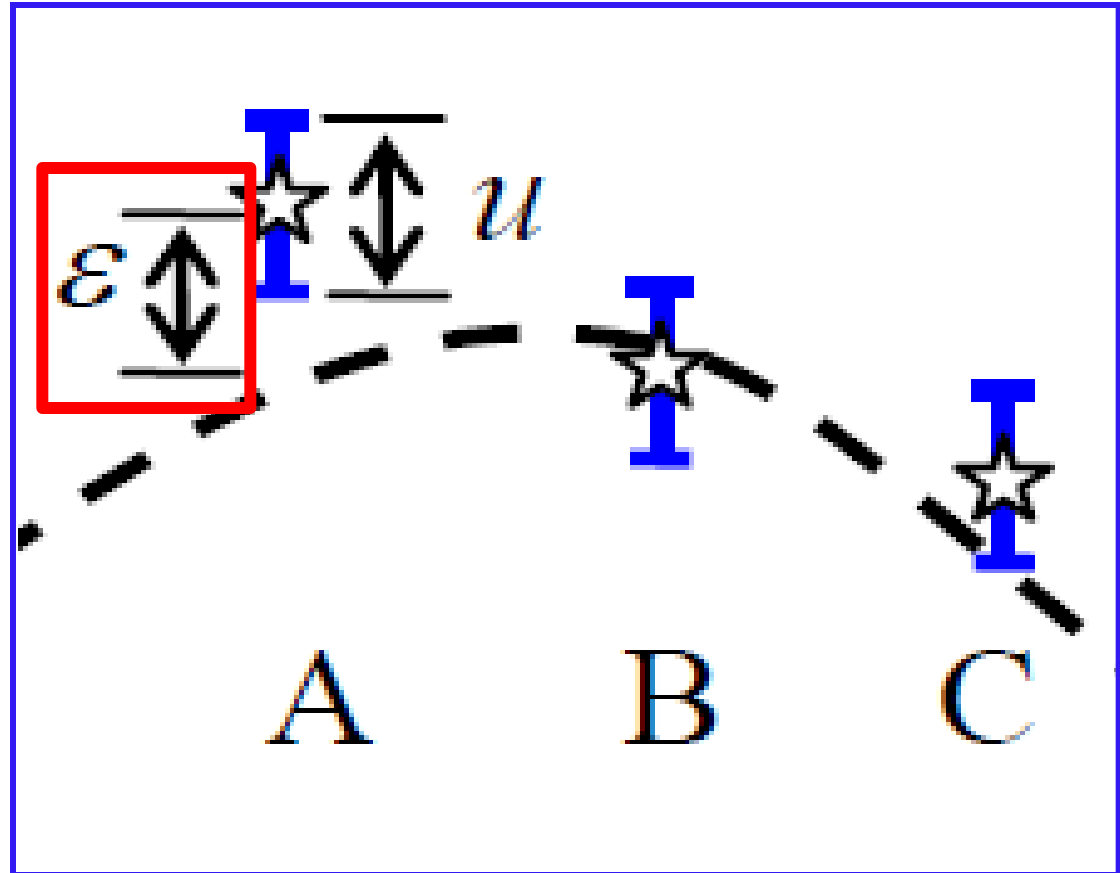
Measured data values

- Exact value (-----) – not available
- u - uncertainty
- ε - error



Measured data values

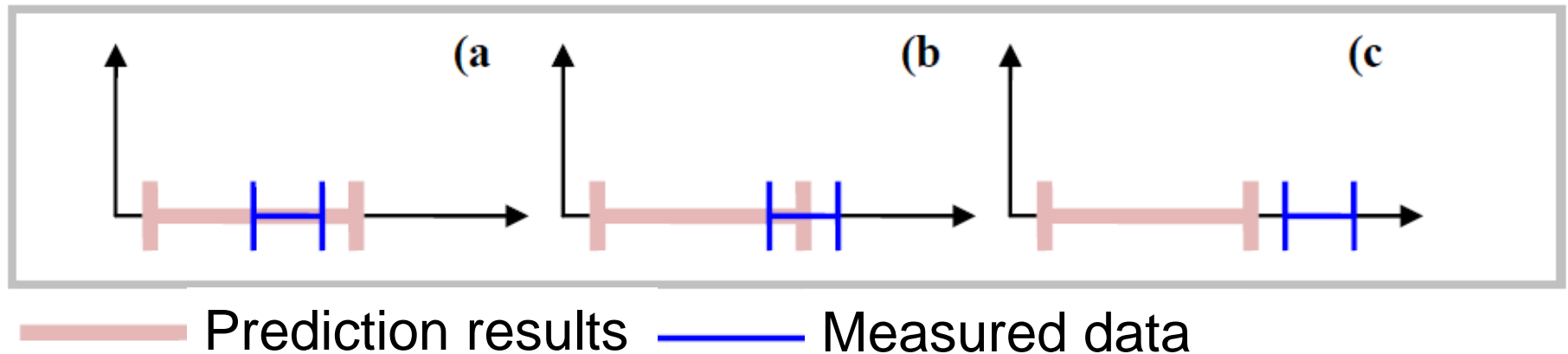
- Exact value (-----) – not available
- u - uncertainty
- ε - error



Data validation – detecting errors if they exist

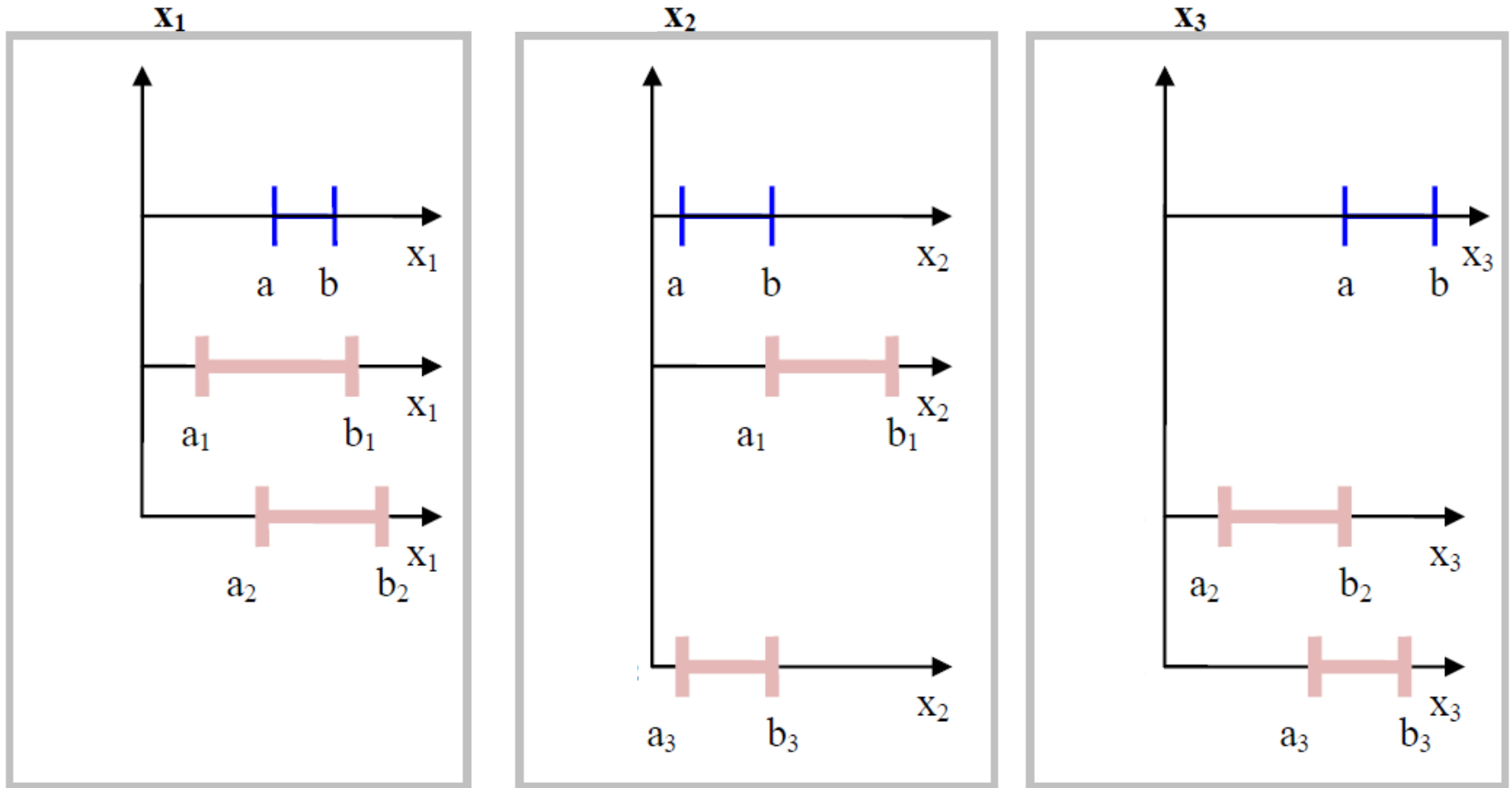
Comparing measured and predicted values

Measured and calculated values are presented with intervals



Enables detection of errors in data – **data validation**

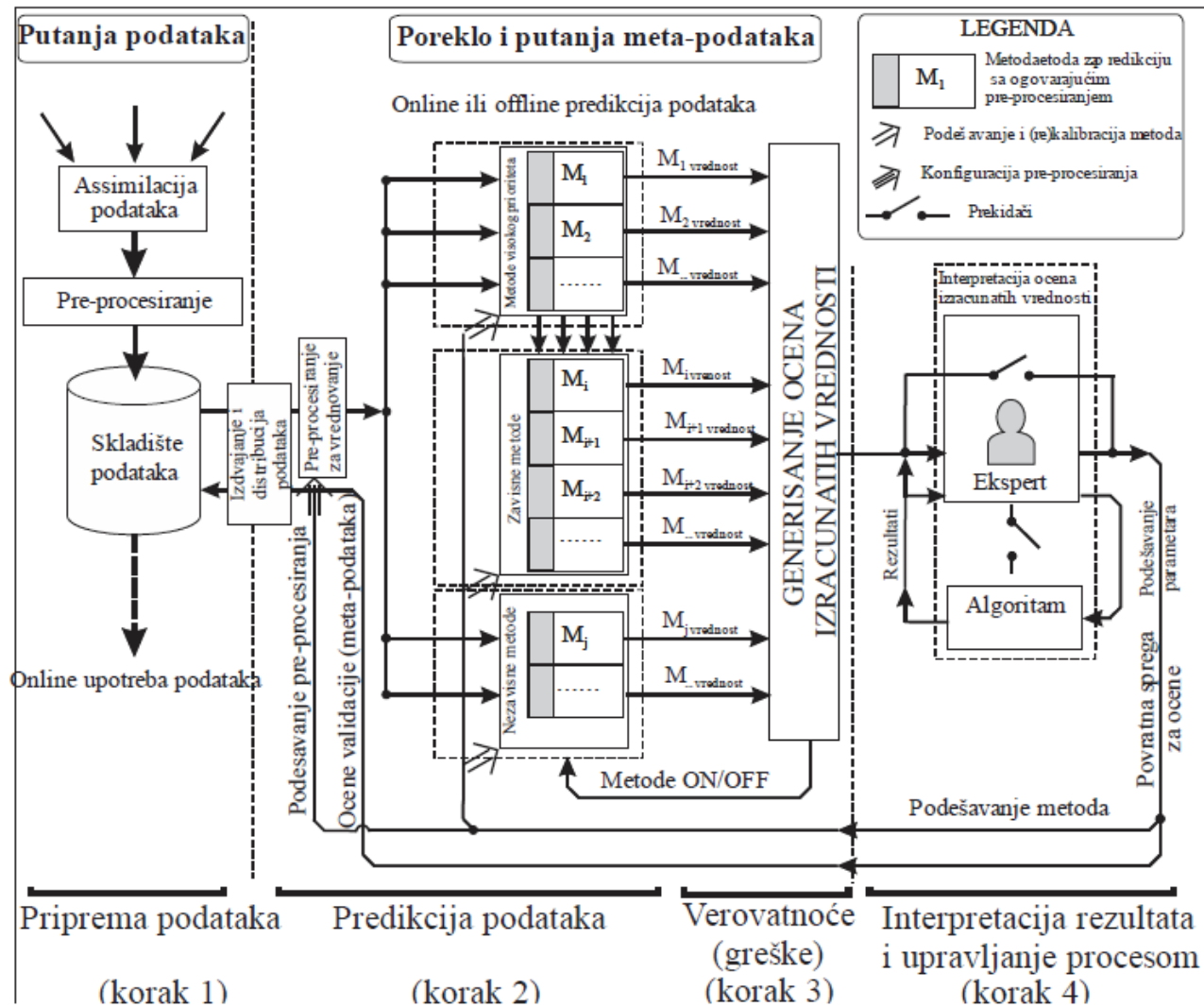
Comparing measured and predicted values -complex system-



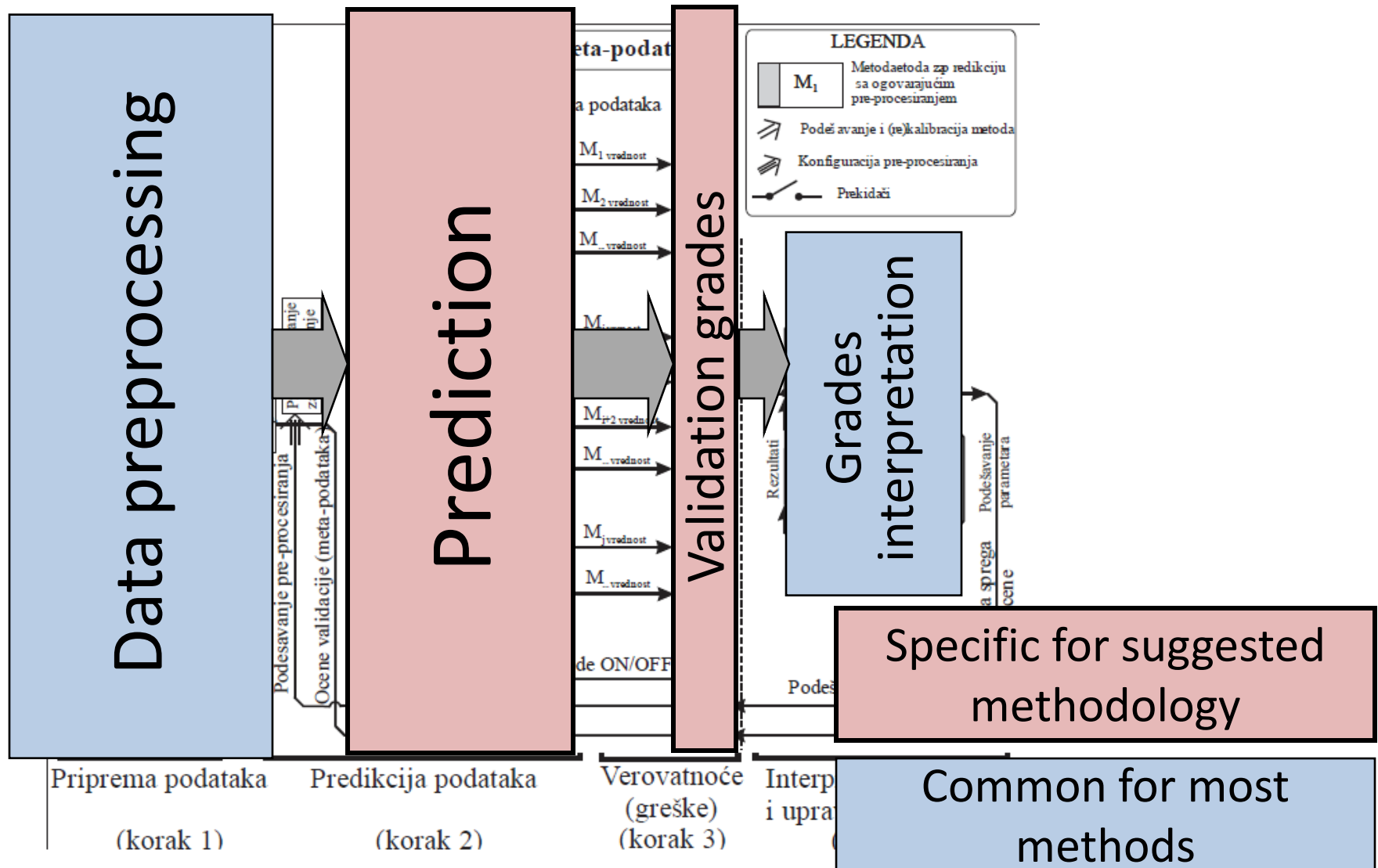
— Prediction results — Measured data

Suggested data validation methodology

sorry for this mess on the screen – will make it more readable...

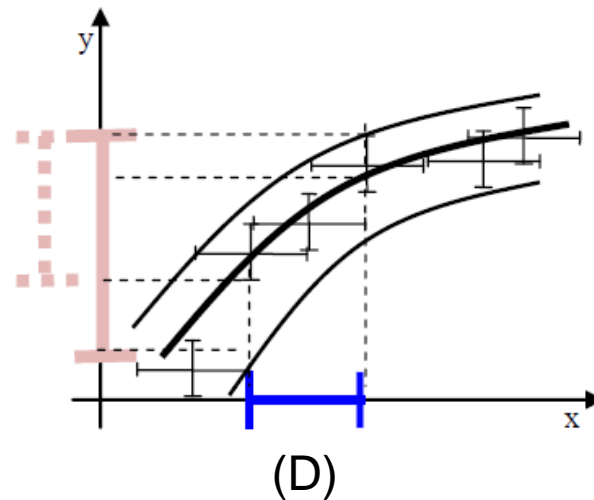
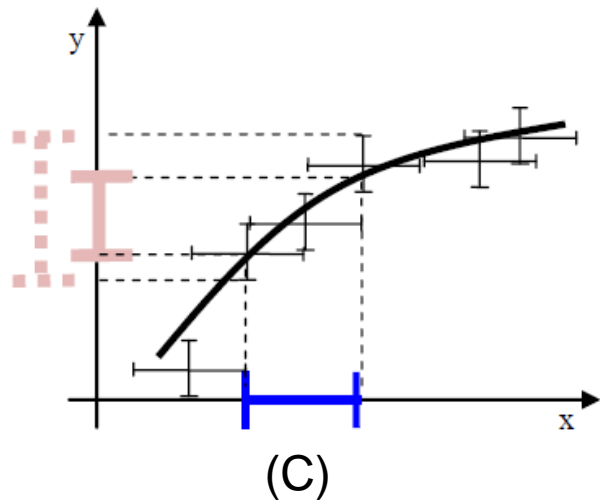
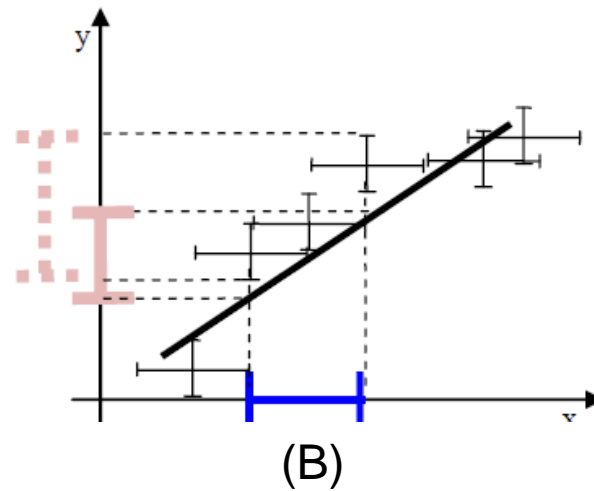
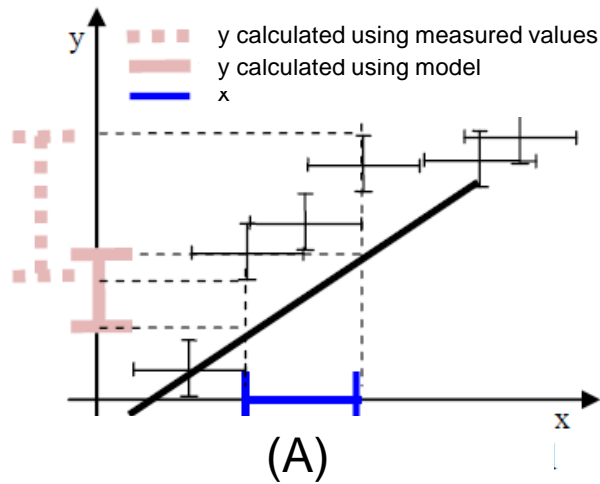


Suggested data validation methodology



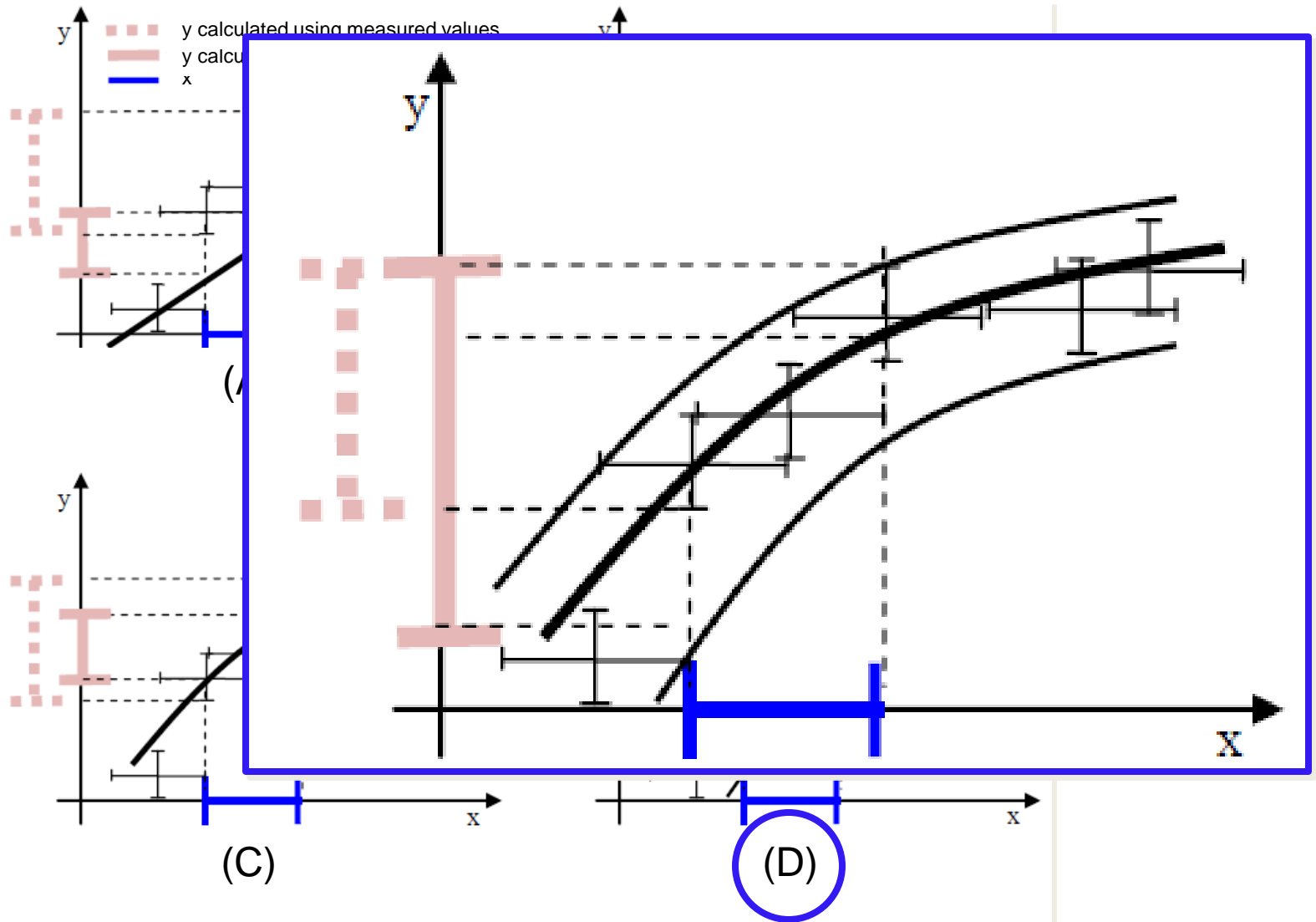
Data prediction

-relation errors and uncertainty-

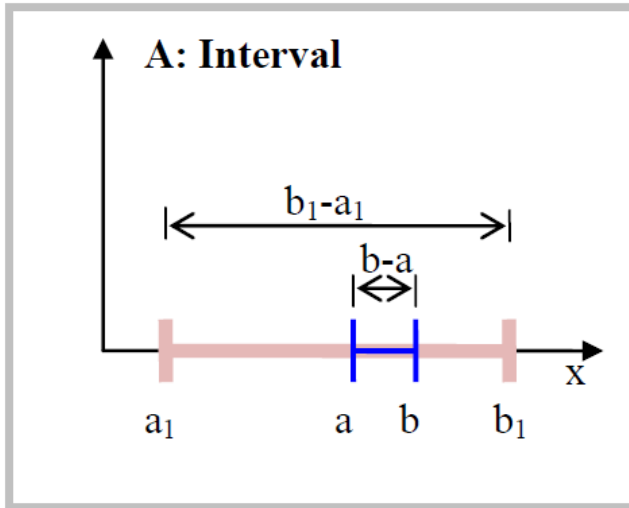


Data prediction

-relation errors and uncertainty-



Validation grades -data probability-



— Prediction results
— Measured data

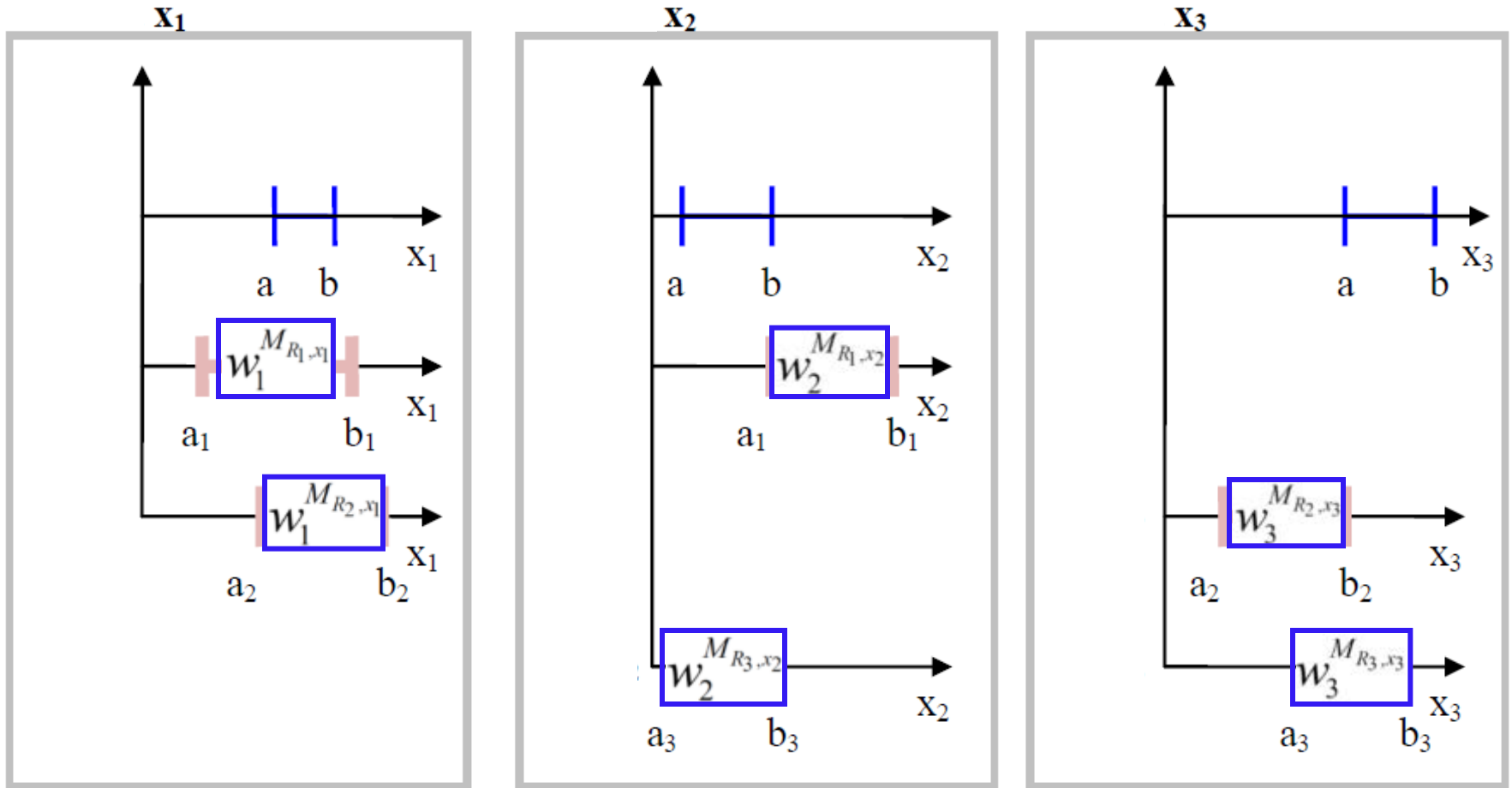
$$p([a, b] | [a_1, b_1]) = \frac{p([a, b] \cap [a_1, b_1])}{p([a_1, b_1])}$$

$$p(x_i | x_i^{M_j}, X_{x_i}^{M_j}) = p([a, b] | [a_1, b_1]) = \begin{cases} 1 - \frac{\max(0, a - a_1) + \max(0, b_1 - b)}{(b_1 - a_1)}, & a_1 \leq b \\ 0, & b_1 \geq a \end{cases}$$

$$\max(p(x_i | x_i^{M_j}, X_{x_i}^{M_j})) = \max(P([a, b] | [a_1, b_1])) = \frac{b - a}{b_1 - a_1}$$

Validation grades

-weighted average of data probabilities-



— Prediction results — Measured data

Validation grades -likelihood of calculated values-

Weights– likelihood of predicted value

$$w_i^{M_{R_j, x_i}} = p\left(x_i^{M_{R_j, x_i}} \mid X_{x_i}^{M_{R_j, x_i}}\right)$$

Maximization of a sum of likelihoods - EM algorithm

$$J = \max \sum_i \sum_j p\left(x_i^{M_{R_j, x_i}} \mid X_{x_i}^{M_{R_j, x_i}}\right)$$

E step:

$$p\left(x_i^{M_{R_j, x_i}} \mid X_{x_i}^{M_{R_j, x_i}}\right) = \sum_{X_{x_i}^{M_{R_j, x_i}}} \left(\frac{\partial M_{R_j, x_i}\left(X_{x_i}^{M_{R_j, x_i}}, \theta\right)}{\partial X_{x_i}^{M_{R_j, x_i}}} \right)^{-1} P_{x_i}$$

M step:

$$P_{x_i} = \sum_j p\left(x_i \mid x_i^{M_{R_j, x_i}}, X_{x_i}^{M_{R_j, x_i}}\right) \times \frac{p\left(x_i^{M_{R_j, x_i}} \mid X_{x_i}^{M_{R_j, x_i}}\right)}{\sum_j p\left(x_i^{M_{R_j, x_i}} \mid X_{x_i}^{M_{R_j, x_i}}\right)}$$

Validation grades -grades-

1. Not normalized

$$x_i^{grade} = \sum_j w_i^{M_{R_j, x_i}} \times p\left(x_i \mid x_i^{M_{R, x_i}}, X_{x_i}^{M_{R, x_i}}\right) \quad (\textit{between 0 and } a < 1)$$

2. Normalized

$$x_i^{grade} = p_{x_i}^{norm} = \left(\frac{\sum_{x_i} w_i^{M_j} \frac{p\left(x_i \mid x_i^{M_{R_j, x_i}}, X_{x_i}^{M_{R_j, x_i}}\right)}{\max\left(p\left(x_i \mid x_i^{M_{R_j, x_i}}, X_{x_i}^{M_{R_j, x_i}}\right)\right)}} \right) \quad (\textit{between 0 and 1})$$

Evaluating validation systems

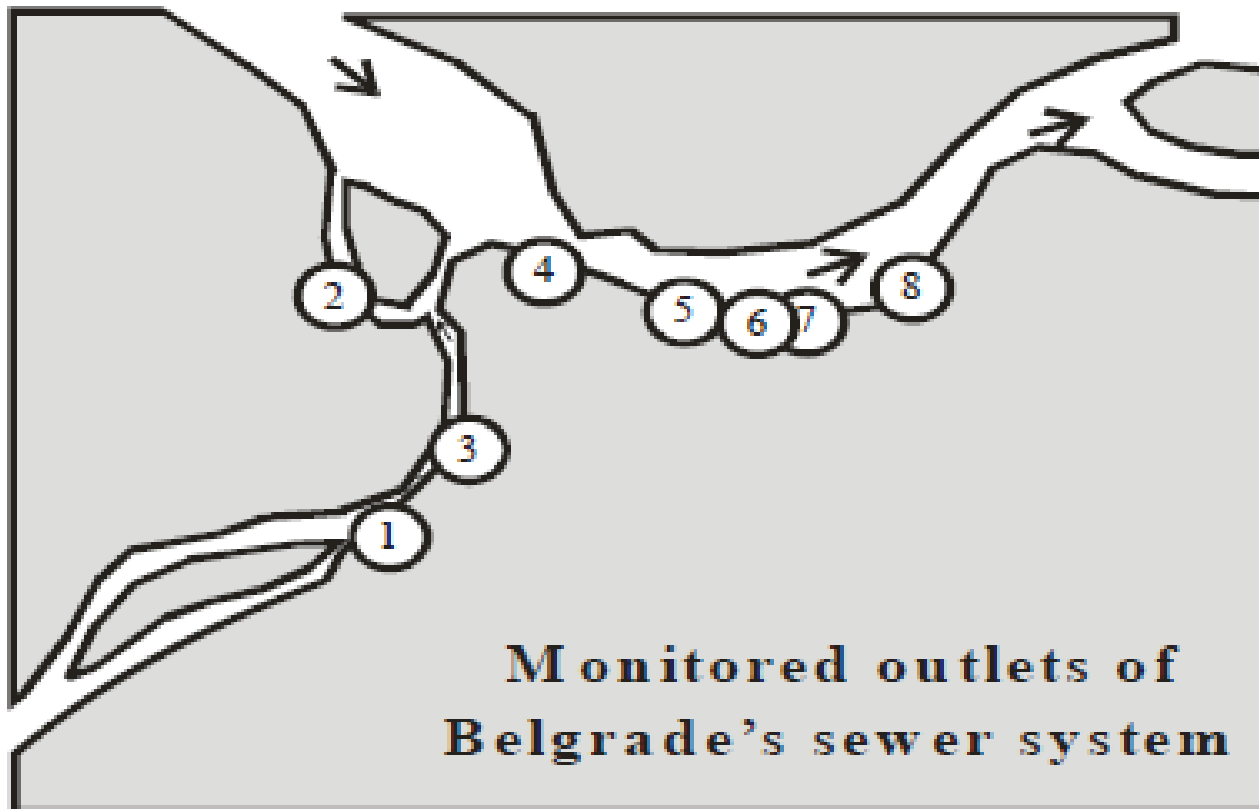
$$p = \frac{N_{registered}}{N_{anomalies} + N_{missed} + N_{registered\ nonanomalies}}$$

Comparing to Kappa index of coincidence:

- Higher penalty to missed anomalies, and
- Lower penalty to correct values interpreted as anomalies

Case study

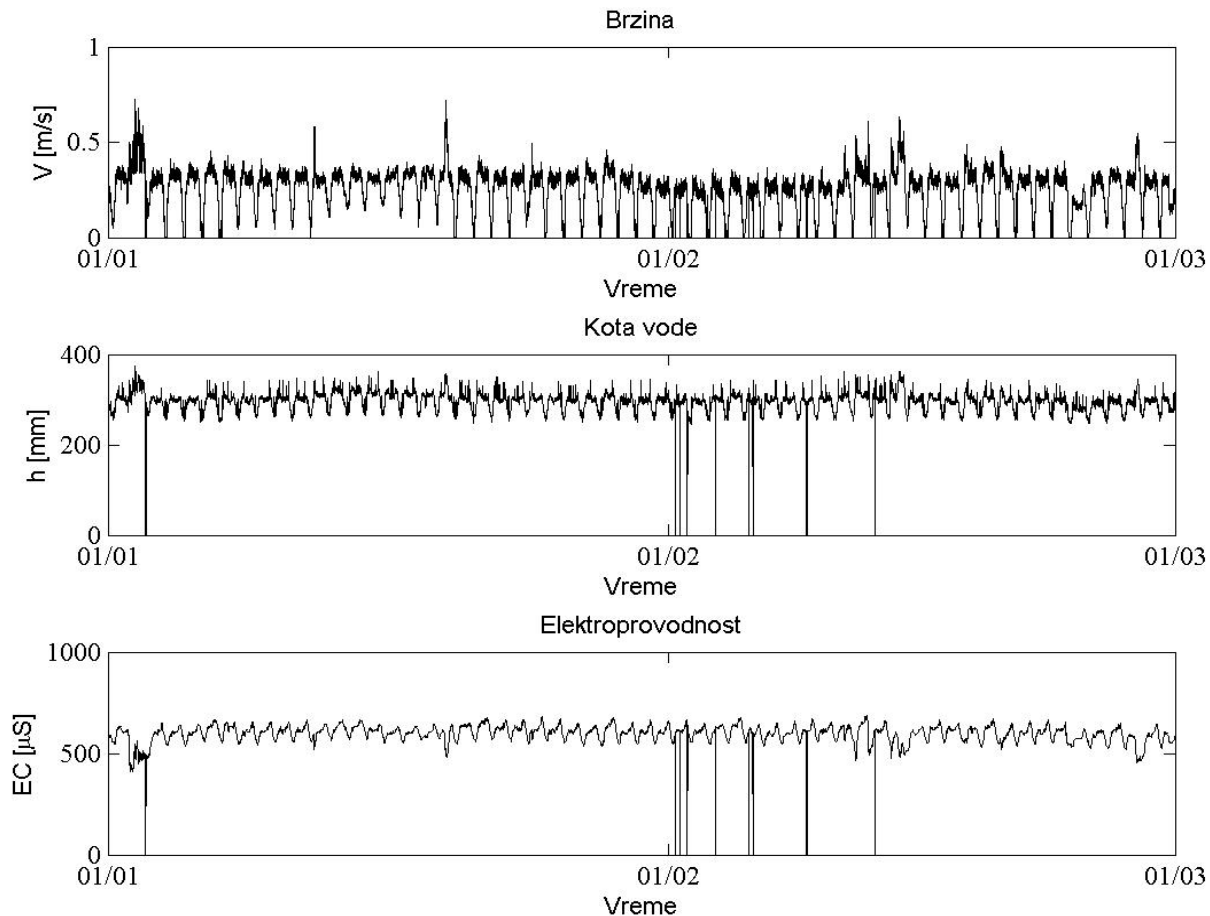
Measurements of hydraulic and water quality parameters in Belgrade sewer CSO



Measurements in Belgrade sewer CSO

Measured variables: Velocity (**V**), water depth (**h**), electro-conductivity(**EC**)

Period: 1/1/2007-1/3/2007



Relations between measured variables

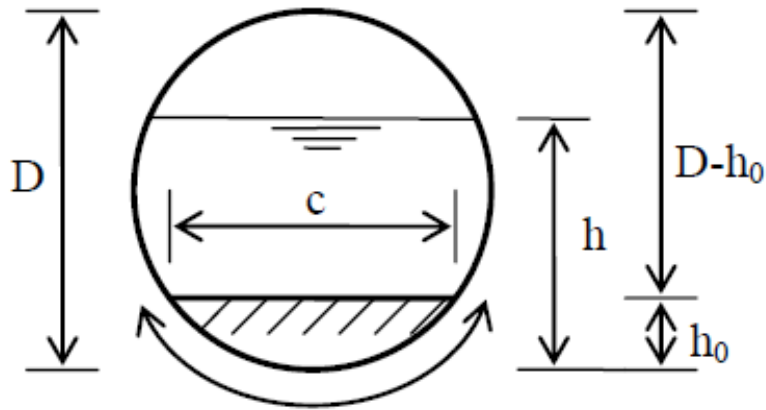
| | V | h | EC | V^{t-1} | h^{t-1} | EC^{t-1} |
|-------|--------------------------------|--------------------------------|-----------------------------------|---|---|--|
| R_1 | $V^t = f_{V^t}^{R_1}(h^t)$ | $h^t = f_{h^t}^{R_1}(V^t)$ | | $V^{t-1} = f_{V^{t-1}}^{R_1}(h^{t-1})$ | $h^{t-1} = f_{h^{t-1}}^{R_1}(V^{t-1})$ | |
| R_2 | | $h^t = f_{h^t}^{R_2}(EC^t)$ | $EC^t = f_{EC^t}^{R_2}(h^t)$ | | $h^{t-1} = f_{h^{t-1}}^{R_2}(EC^{t-1})$ | $EC^{t-1} = f_{EC^{t-1}}^{R_2}(h^{t-1})$ |
| R_3 | $V^t = f_{V^t}^{R_3}(EC^t)$ | | $EC^t = f_{EC^t}^{R_3}(V^t)$ | $V^{t-1} = f_{V^{t-1}}^{R_3}(EC^{t-1})$ | | $EC^{t-1} = f_{EC^{t-1}}^{R_3}(V^{t-1})$ |
| R_4 | $V^t = f_{V^t}^{R_4}(V^{t-1})$ | | | $V^{t-1} = f_{V^{t-1}}^{R_4}(V)$ | | |
| R_5 | | $h^t = f_{h^t}^{R_5}(h^{t-1})$ | | | $h^{t-1} = f_{h^{t-1}}^{R_5}(h)$ | |
| R_6 | | | $EC^t = f_{EC^t}^{R_6}(EC^{t-1})$ | | | $EC^{t-1} = f_{EC^{t-1}}^{R_6}(EC)$ |

R₁ – Chezy-Manning equation

R₂, R₃ – Water quality model of EC reduction

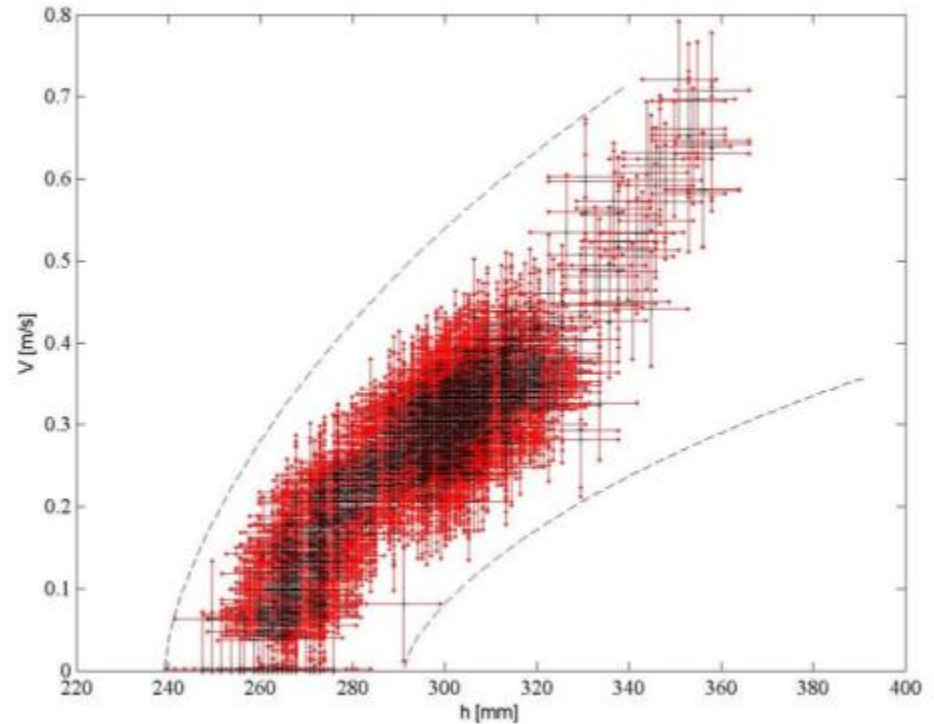
R₄, R₅, R₆ – AR(1) models

Chezy-Manning equation



$$V(h_0, C, h) = C \times R(h_0, h)^{2/3}$$

$$C = \sqrt{I_d} / n$$



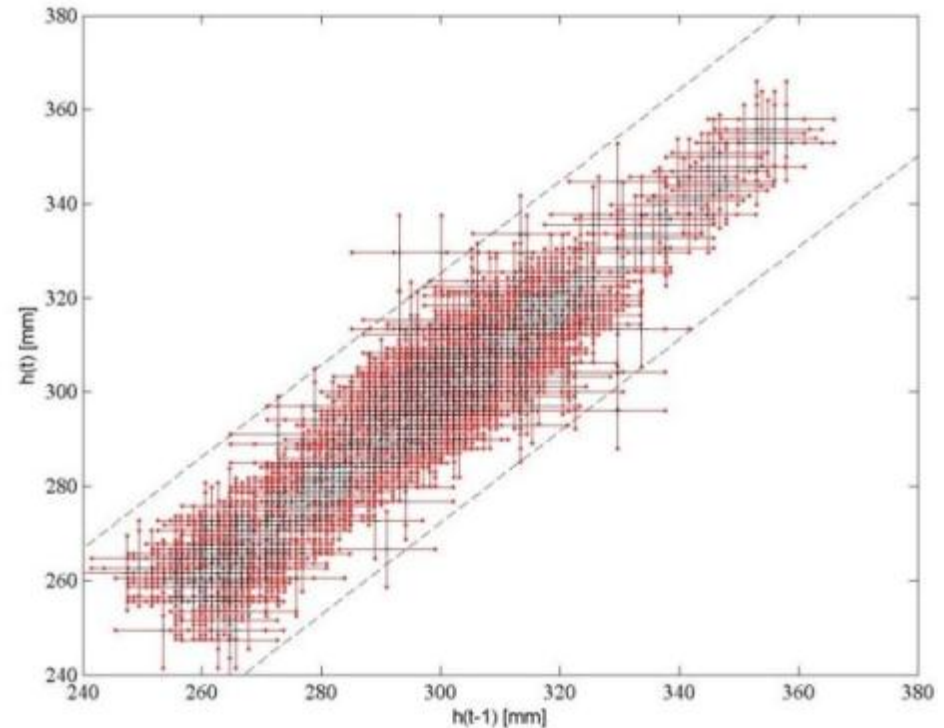
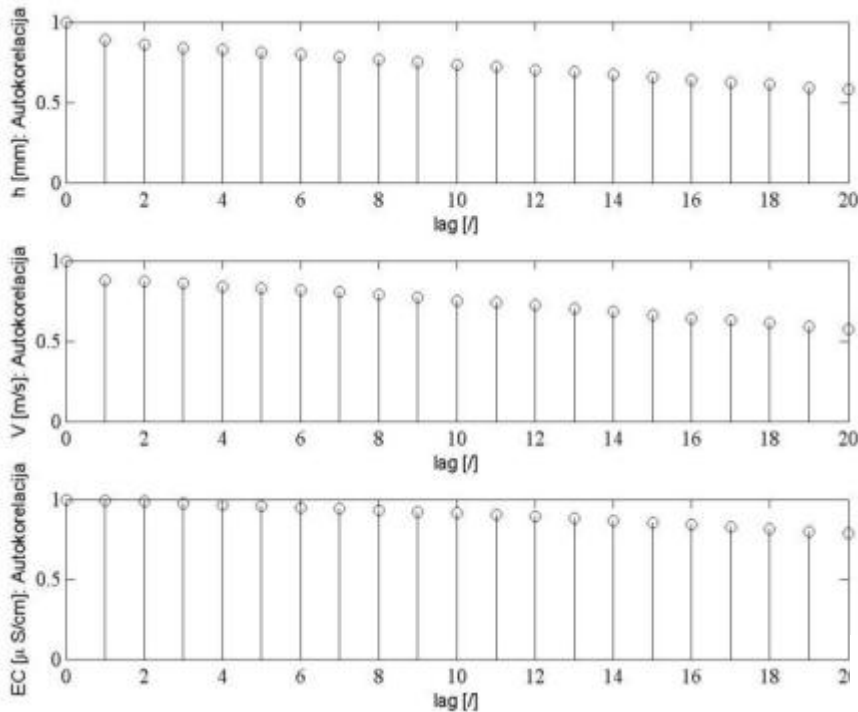
$$M_{R_1, V} : V(h) = [3.8, 1.9] \times R([239, 291], h)^{2/3}$$

AR(1) models

$$x^t = ax^{t-1} + b$$

Autocorrelation coefficient

$$R_{ff}(\tau) = \int_{-\infty}^{\infty} f(t + \tau)\bar{f}(t) dt$$

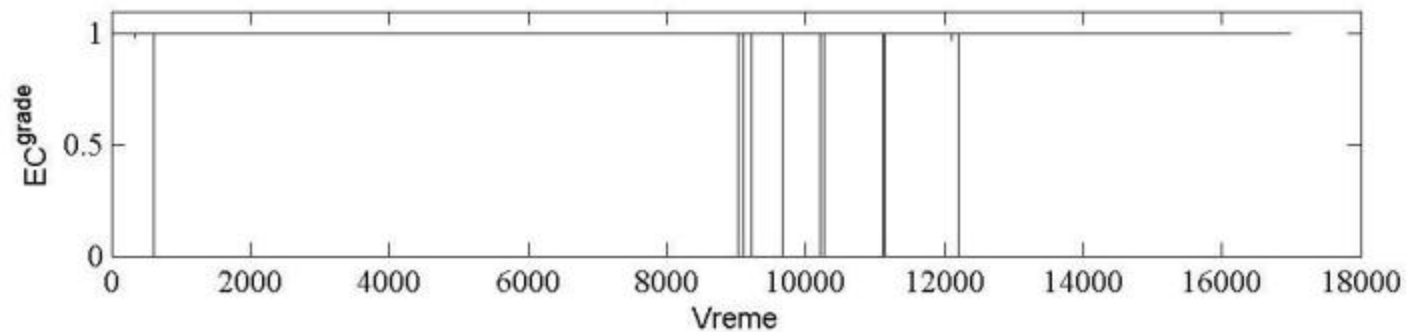
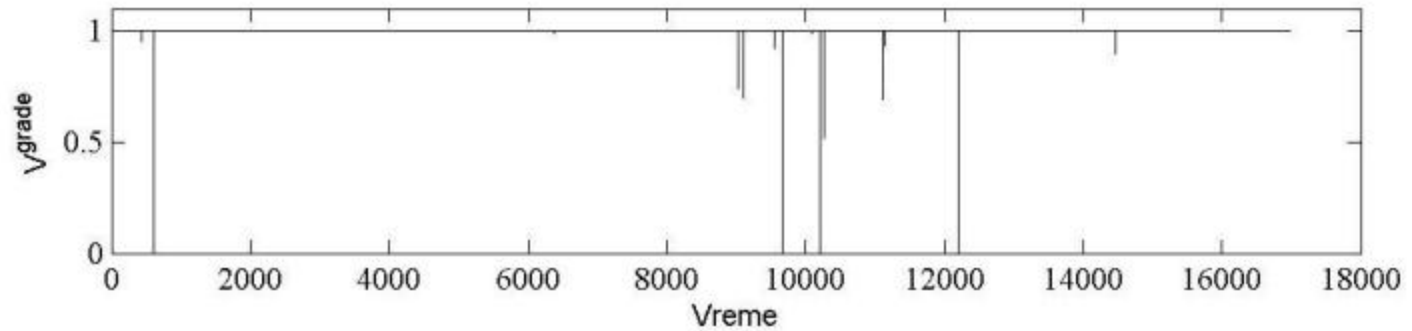
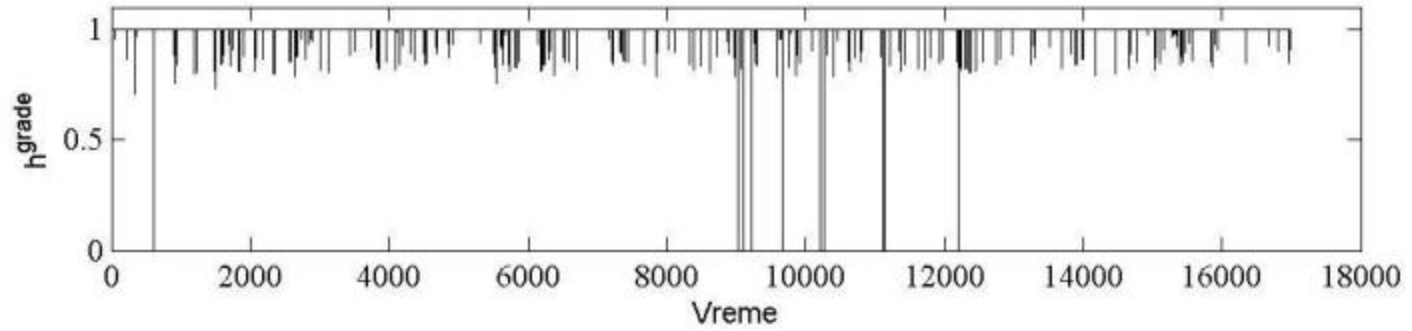


Model calibration results

| x | h | V | EC |
|--------------------------------|-----------|---------------|-----------|
| a | 0.97427 | 0.97258 | 0.99712 |
| $b = [\underline{b}, \bar{b}]$ | [-33, 52] | [-0.17, 0.19] | [-59, 49] |

Results

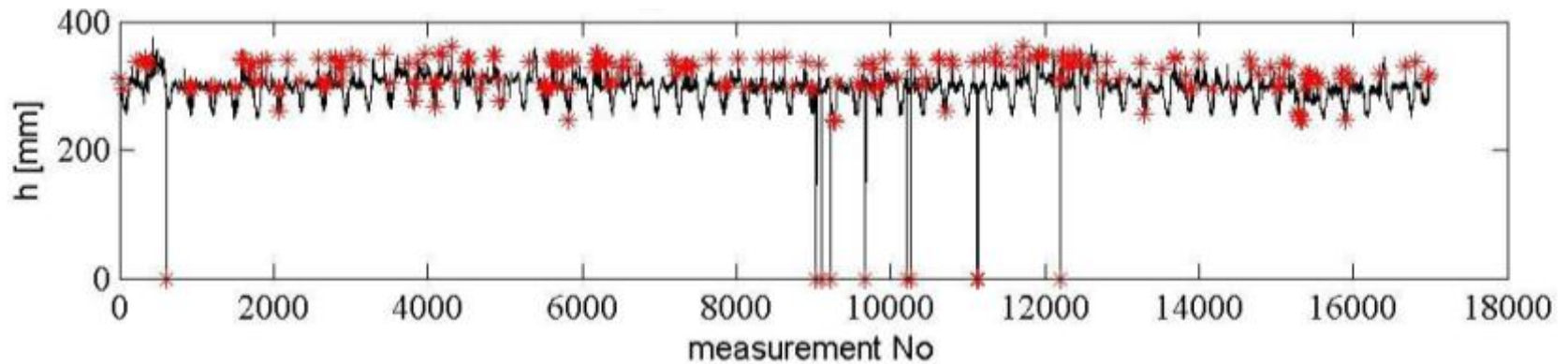
-Validation grades-



Results

-comparison with manual data validation-

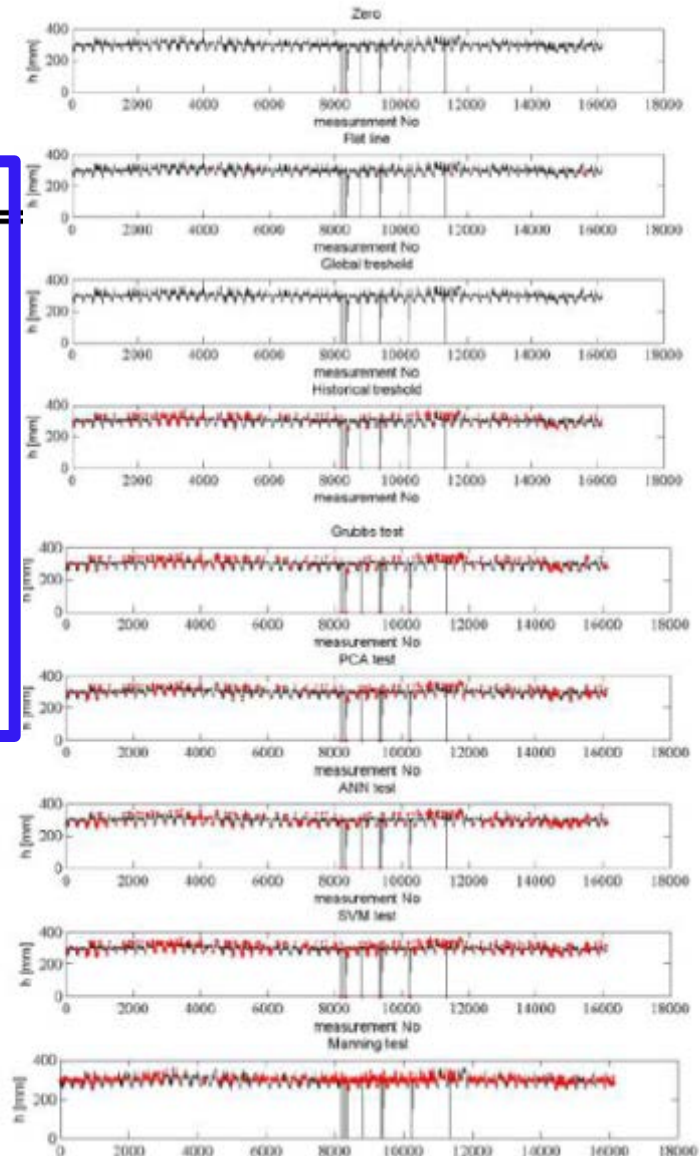
Threshold value= **0.99**



| | | | | |
|-----------------------|----------------------|--------------------|-------------------|----------|
| <i>(No anomalies)</i> | <i>(No detected)</i> | <i>(No missed)</i> | <i>(No false)</i> | <i>P</i> |
| 244 | 219 | 25 | 73 | 0.64 |

Nine validation methods

| | No anomalies | No detected | No missed | No false | p |
|-------|--------------|-------------|-----------|----------|-------|
| M_1 | 244 | 24 | 220 | 0 | 0.052 |
| M_2 | 244 | 0 | 244 | 28 | 0.000 |
| M_3 | 244 | 0 | 244 | 0 | 0.000 |
| M_4 | 244 | 202 | 42 | 220 | 0.399 |
| M_5 | 244 | 187 | 57 | 334 | 0.294 |
| M_6 | 244 | 224 | 20 | 359 | 0.360 |
| M_7 | 244 | 109 | 135 | 896 | 0.085 |
| M_8 | 244 | 237 | 11 | 483 | 0.321 |
| M_9 | 244 | 71 | 173 | 2725 | 0.023 |



Averaging validation results:

$$P = \frac{N_{registered}}{N_{anomalies} + N_{missed} + N_{registered\ nonanomalies}} = \frac{195}{244 + 49 + 58} = 0.56$$

Conclusions

- New theoretical framework for data validation
- No restrictions about number and type of additional information that may be used in the validation process
- No restrictions of type and number of relations that may be used
- The validation system may be evaluated and improved by improving the relations between the data values



Automatic measured data validation applied on hydraulic and water quality parameters in sewer systems

Nemanja Branislavljević¹

Dušan Prodanović¹

Zoran Kapelan²

¹University of Belgrade

²University of Exeter