

**Raport Badawczy**

**RB/26/2016**

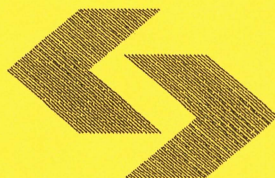
**Research Report**

**Towards handling uncertainty  
in prognostic scenarios:  
advanced learning  
from the past**

**P. Żebrowski,  
M. Jonas, J. Jarnicka**

**Instytut Badań Systemowych  
Polska Akademia Nauk**

**Systems Research Institute  
Polish Academy of Sciences**



**POLSKA AKADEMIA NAUK**

**Instytut Badań Systemowych**

ul. Newelska 6

01-447 Warszawa

tel.: (+48) (22) 3810100

fax: (+48) (22) 3810105

Kierownik Zakładu zgłaszający pracę:  
Prof. dr hab. inż. Zbigniew Nahorski

Warszawa 2016

SYSTEMS RESEARCH INSTITUTE  
POLISH ACADEMY OF SCIENCES

Piotr Żebrowski, Matthias Jonas,  
Jolanta Jarnicka

**Towards Handling Uncertainty in  
Prognostic Scenarios:  
Advanced Learning from the Past**

WARSAW 2016

# Contents

- 1. Introduction ..... 1
  - 1.1. Scientific context of the project ..... 1
  - 1.2. Motivation: problems with judging the credibility of predictions ..... 2
  - 1.3. Objectives and scope of the report..... 4
  - 1.4. Structure of the report..... 4
- 2. Learning in a controlled prognostic context ..... 5
  - 2.1. Generic notion of the explainable outreach of the data ..... 6
  - 2.2. Prognostic learning procedure..... 6
  - 2.3. Applying the prognostic learning procedure and interpretation of its results ..... 8
  - 2.4. PL versus forecasting with use of time series analysis..... 10
- 3. Regression – based construction of the EO ..... 12
  - 3.1. Analysis of historical patterns in learning phase with use of polynomial regression . 12
  - 3.2. Construction of the EO..... 14
  - 3.3. Procedure of prognostic learning based on regression method..... 15
- 4. Assessment of prognostic learning performance in the controlled conditions. Monte Carlo experiments ..... 17
  - 4.1. Method of generating the synthetic data ..... 17
  - 4.2. Description of experiments on synthetic data ..... 18
  - 4.3. Results ..... 20
    - 4.3.1. Data following a linear trend..... 21
    - 4.3.2. Data following a 4<sup>th</sup> order polynomial trend..... 24
    - 4.3.3. Data following exponential trend ..... 31
    - 4.3.4. Data following logarithmic trend ..... 37
    - 4.3.5. Data following periodic trend ..... 43
  - 4.4. Conclusions ..... 49
- 5. Real-life case studies ..... 52
  - 5.1. Global CO<sub>2</sub> emissions from technosphere..... 52
  - 5.2. Concentration of CO<sub>2</sub> in the atmosphere ..... 56

5.3. Conclusions .....	60
6. Outlook.....	61
7. Summary .....	63
8. Acronyms.....	64
9. Literature.....	64
Appendix: Nonparametric kernel-based regression .....	66

## Abstract

In this report we introduce the paradigm of learning from the past which is realized in a controlled prognostic context. It is a data-driven exploratory approach to assessing the limits to credibility of any expectations about the future system's behavior which are based on a time series of a historical observations of the analyzed system. Such horizon of the credible expectations is derived as the length of explainable outreach of the data, i.e. the spatio-temporal extent for which, in lieu of the knowledge contained in the historical observations, we may have a justified belief to contain future system's observations. Explainable outreach is of practical interest to the stakeholders since it allows to assess the credibility of scenarios produced by models of the analyzed system. It also indicates the scale of measures required to overcome the system's inertia. In this report we propose a method of learning in a controlled prognostic context which is based on polynomial regression technique. A polynomial regression model is used to grasp the system's dynamic revealed by the sample of historical observations, while the explainable outreach is constructed around the extrapolated regression function. The proposed learning method was tested on various sets of synthetic data in order to identify its strengths and weaknesses, formulate guidelines for its practical application. We also demonstrate how it can be used in context of earth system sciences by applying it to derive the explainable outreach of historical anthropogenic CO<sub>2</sub> emissions and atmospheric CO<sub>2</sub> concentrations. We arrived at conclusion that the most robust method of building the explainable outreach is based on linear regression. However, such explainable outreach of the analyzed data sets (representing credible expectations based on extrapolation of linear trend) is rather short.

## **Acknowledgments**

The authors are grateful to the Earth Systems Sciences [ESS] Research Program of the Austrian Academy of Sciences [OeAW] for financing this research.

## About the Authors

**Piotr Żebrowski** joined IIASA's Advanced Systems Analysis (ASA) Program as a research assistant in February 2015. His current research focus is on diagnostic uncertainty of greenhouse gas inventories, uncertainty propagation in climate models and on retrospective learning.

**Matthias Jonas** is a senior research scholar with IIASA's Advanced Systems Analysis (ASA) Program. His interests are in environmental science, and in the development of systems analytical models and tools to address issues of global, universal and regional change, including surprises, and their potential implications for decision and policymakers.

**Jolanta Jarnicka** is a researcher in the Systems Research Institute of the Polish Academy of Sciences. Her speciality is probability and statistics, in particular nonparametric statistical methods, data analysis, and mathematical modelling.



# Towards Handling Uncertainty in Prognostic Scenarios: Advanced Learning from the Past

Piotr Żebrowski<sup>1</sup>, Matthias Jonas<sup>1</sup>, Jolanta Jarnicka<sup>2</sup>

<sup>1</sup> IIASA, Advanced Systems Analysis Program

<sup>2</sup> Systems Research Institute of the Polish Academy of Sciences

## 1. Introduction

### 1.1. Scientific context of the project

The problem of uncertainty and horizons of credibility<sup>1</sup> of predictions of future behavior of Earth – climate system attracts growing interest as a consequence of the increasing demand of incorporating information about future climate into planning and decision making (e.g., IPCC 2007: FAQ 1.2, FAQ 8.1; NSF 2012; IPCC 2013: Box 11.1; Otto et al. 2015). Numerous scientific institutions, including IIASA, use a variety of complex integrated assessment models to generate a great number of prognostic scenarios in order to identify policy options and effectiveness of different measures for mitigating the climate change. Modelers make huge efforts trying to ensure the credibility of their scenarios and gauging their uncertainty, e.g., by carrying out sensitivity tests or inter-model comparisons under standardized conditions. In particular, multi-model-scenario exercises are becoming increasingly popular (e.g., Meinshausen *et al.* 2009). Nevertheless, such efforts are not entirely convincing and judging the credibility of climate model projections remains a notorious and unresolved question (cf. Otto et al. 2015).

In contrast to these model-related issues we propose adopting an alternative, data-driven perspective of looking at the limits to applying our current understanding of the Earth system for predicting its future behavior. We seek to assess these limits by answering the following questions:

- (1) *Given the data reflecting a system and their diagnostic uncertainty can we deduce the **explainable outreach**<sup>2</sup> of these data, which express our understanding of the prevailing patterns of system's behavior and their typical duration?*

---

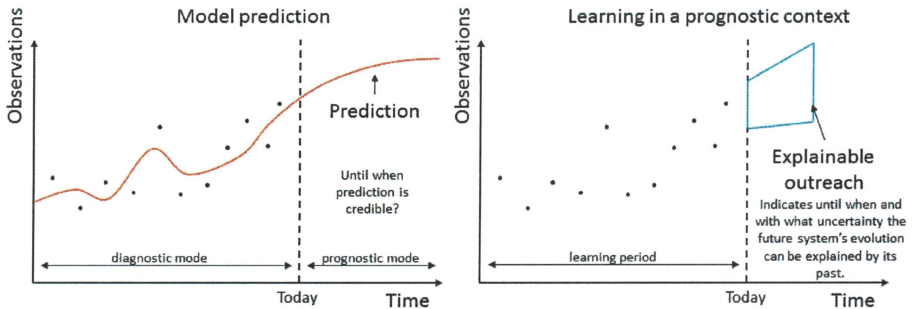
<sup>1</sup> Credibility of predictions is understood as our expectations (predictions) of its performance (Otto et al. 2015)

<sup>2</sup> The region – both in terms of time horizon and the range of plausible future values – within which we may have justifiable belief based on the past system's behaviour, that it will contain future trajectory of the process' evolution.

and

(2) Can the explainable outreach be used for assessing limits of credibility of predictions?

In order to answer these questions, we develop and apply a new (to our knowledge) exploratory method, which we call **learning in a controlled prognostic context**<sup>3</sup>. Its main idea is **to learn about the nature of the analyzed system from its past**: we use a part of the historical observations of the system to understand its basic dynamic and formulate our expectations about its future evolution (expressed as the explainable outreach) and then test these expectations against the remaining part of the sample. Such way of testing the limits of our understanding of the system based on partial and uncertain knowledge (carried by a finite set of possibly imprecise<sup>4</sup> observations) may inform us about the likely time horizon within which our expectations about its future evolution may be considered plausible in lieu of the available historical data. Therefore, the proposed method belongs to the realm of **data analysis, NOT modelling**. The difference between learning in a controlled prognostic context and modelling is explained by Figure 1.)



**Figure 1. Model prediction vs. learning in a prognostic context.** Left panel: Model prediction. A model is calibrated against historical data (diagnostic mode) before making a prediction, e.g. by extrapolating the historical trend into the future or generating a scenario pathway (prognostic mode). Modelers typically do not (or cannot) indicate until when a model prediction is in accordance with the systems past (i.e. is credible). Right panel: Learning in a prognostic context. Given the historical data the system's dynamics can be grasped and the data's explainable outreach be constructed. The explainable outreach specifies both spatial and temporal extent beyond which we cannot explain our system anymore in accordance with its past. The purpose of deriving explainable outreach directly from the data is to indicate limits of predictability of the model which we built to reflect the underlying system.

## 1.2. Motivation: problems with judging the credibility of predictions

Credibility of predictions is one of the central problems of statistical modelling. A variety of well-established statistical methods - such as regression models and machine learning techniques (Hastie *et al.* 2009, Murphy 2012) or time series analysis techniques (Brockwell & Davis 2002) - aim at predicting responses of the analyzed system in yet

<sup>3</sup> For simplicity, we call it also a **prognostic learning (PL)** method. Description of the method together with explanation of its name is provided in Chapter 2.

<sup>4</sup> We assume that the data are accurate (i.e. no systematic bias of the system's observations).

unobserved states<sup>5</sup>. Predictions are typically expressed in terms of a regression function (i.e. conditional expected value of the system's response given the value of the explanatory state variable). Quality of predictions, usually understood as expected prediction error, can be controlled<sup>6</sup> provided that the state in which we wish to make a prediction lays **within** the range of the data sample on which the analysis is based. However, analogous error control is formally unavailable for predictions of the system's responses in states laying **beyond** the range of the data sample (i.e. in conditions which may be significantly different than those to which historical observations correspond).

Similar problems haunt also the modelling community. Their common and apparently unavoidable practice is to extrapolate the current understanding of the system (e.g., discovered trends or relationships) **beyond** the range of historical data sample in order to predict its future behavior, possibly in yet unobserved states. For example, this approach was employed in the study of Meinshausen *et al.* 2009 aiming at prediction of level of global warming in the future, when GHG concentrations in the atmosphere will be at the levels without precedence in (recent) history. However, making such predictions by extrapolating the observed trends beyond the range of sample is problematic. Unless one assumes that observed process is in some sense stationary (which may be a too strong assumption, e.g. in presence of varying exogenous forcing) one lose the control over the quality of predictions, whose errors may rapidly increase the more the further away from the sample of historical observations one moves. Typically, modelers try to assess credibility of predictions either (1) by providing uncertainty ranges for the predictions<sup>7</sup>; (2) by means of sensitivity analyses<sup>8</sup>; or (3) by exploring the range of possible futures by means of selected scenario pathways (in particular in the case of computationally expensive models). Unfortunately, these methods are not entirely convincing due to a certain degree of arbitrariness in their application (e.g., assumed distributions of parameters underlying Monte Carlo methods or the choice of storylines for scenario pathways). More importantly, they **do no indicate the time horizon within which a model predictions remain in accordance with the system's past**<sup>9</sup>.

The paradigm of learning in a controlled prognostic context offers at least partial solution to these problems. It is a data analysis method designed to control the growing uncertainty of our expectations about the system's evolution in the immediate future. Moreover, this approach may provide a model-independent indicator of the time range within which the projections of the model may be judged credible in lieu of the past system behavior.

---

<sup>5</sup> i.e. in conditions not covered by the available data (out-of-sample predictions)

<sup>6</sup> The upper bands for probability of occurring large prediction errors are available and depend on the complexity of the statistical model and the length of the data sample.

<sup>7</sup> Assuming suitable probability distributions for values of exogenous parameters of the model they may be derived analytically or by means of Monte Carlo simulations.

<sup>8</sup> In this case possible correlations between exogenous parameters of the model are typically ignored. Changes in model responses are usually analysed by varying values of one of the parameters while keeping the rest constant.

<sup>9</sup> By "remaining in accordance with the system's past" we mean that predicted future trajectory of the system's evolution exhibits behavior similar to this observed in the past, such as the level of "system's inertia" or the type of dynamics. Note that this is weaker notion than stationarity of the process.

### 1.3. Objectives and scope of the report

Objectives of this report are following: (1) to introduce the generic paradigm of the **learning in a controlled prognostic context** allowing to assess the **explainable outreach**, i.e. the region – specified in terms of time horizon and the range of plausible future values (uncertainty) – within which we may have justifiable belief (based on historical observations) that it will contain future trajectory of the system’s evolution; (2) to propose a way (based on regression techniques) of implementing the PL paradigm; and (3) to demonstrate its usefulness in analysis of the real data samples relevant to understanding the Earth climate system (e.g. anthropogenic CO<sub>2</sub> emissions and atmospheric CO<sub>2</sub> concentrations).

The paradigm of learning in the controlled prognostic context is applicable both to: (1) univariate regression – like problems in which one is interested in the form of dependence of one quantity characterizing a system (response variable) on another quantity (called independent variable) which represents the state of the system or its forcing; and (2) analysis of the data forming a time series – in which case the time is treated as the independent variable.

In this report we restrict ourselves to analysis of the time series type of data only, i.e. to case (2). The reason for that is two-fold. Firstly, in context of time series “predicting beyond the range of sample” means “forecasting or predicting the future” which facilitates understanding of the idea of explainable outreach. Secondly, time series perspective is relevant both in context of prognostic modelling and in context of understanding the relevant earth systems processes (such as abovementioned CO<sub>2</sub> emissions or CO<sub>2</sub> concentrations). Hence, from now on (unless stated otherwise), all considered data samples will be assumed to consist of pairs  $(t, x_t)$ , where  $x_t$  denotes the value of the observable describing the system of interest which was recorded at time  $t$ . We will call this observable a system’s state variable<sup>10</sup>.

### 1.4. Structure of the report

In Chapter 2 we introduce the concept of learning in a controlled prognostic context. There we give a definition of the explainable outreach of the data, which is a central notion of the proposed methodology. Next, we formulate a generic procedure of learning in a controlled prognostic context and discuss how it should be applied and how to interpret its results. We conclude Chapter 2 with explaining characteristic aspects of the proposed approach vis-à-vis standard methods of time series analysis.

In Chapter 3 we propose a way of implementing the generic procedure of learning in a controlled prognostic context. Namely, we show how prognostic learning can be operationalized with use of the polynomial regression technique. We discuss how to define the shape of explainable outreach and how to determine its length. We summarize Chapter 3 with formulation of the regression-based procedure of prognostic learning.

The next two chapters are devoted to analysis of the performance of the proposed method. In Chapter 4 we present insights following from the experiments on various synthetic data sets. The purpose of these experiments is to identify strengths and

---

<sup>10</sup> or simply state variable

weaknesses of the proposed method and to formulate guidelines for its application in analysis of the real-life data. In Chapter 5 we test these insights in practice by applying the method to determine explainable outreach of the time series representing anthropogenic CO<sub>2</sub> emissions and atmospheric CO<sub>2</sub> concentrations.

We conclude this report with summary and outlook for future research followed by the Appendix in which we present yet another way of implementing the prognostic learning method – this time based on non-parametric regression techniques. We also demonstrate the potential of this variant of prognostic learning method by applying it to the abovementioned real-life time series.

## 2. Learning in a controlled prognostic context

In this chapter we present the notion of learning in a controlled prognostic context, which for the sake of brevity we also call prognostic learning (PL). Broadly speaking the purpose of this method is to indicate both the typical length of time intervals over which the trends observed in the historical data sample persist, and the level of uncertainty in grasping these trends.

Prognostic learning can be classified as a method of exploratory data analysis. Its aim is not to find a formal statistical model which can be used for testing hypothesis about the historical data sample and making predictions for the future. Instead, PL method offers a semi-formal first-order description of the system's dynamics and its "inertia"<sup>11</sup> exhibited by the system over the period in which the data sample was collected. This "inertia" is a critical factor determining the limits to credibility of predictions of the system's behavior<sup>12</sup>.

As such, the PL method informs us solely about the system's behavior in the past. However, in this report we demonstrate that it is also useful in context of expressing expectations about the immediate future of the system. Rationale for this approach is provided by the observation that patterns in the system's behavior in the relatively recent past are also likely to occur in the nearby future. Therefore, the findings of the PL method, which, in essence, concerns only the past of the system, can also be informative about its nearby future. Note that the requirement for this line of thinking to be valid is just that the nature of the system itself or its external forcing do not change too rapidly over time. This is considerably weaker requirement than stationarity of the system usually assumed by the formal statistical modelling methods<sup>13</sup>.

---

<sup>11</sup> Understood as a system's memory - a typical period within which the system does not undergo a significant change of its dynamics (e.g., average time horizon within which system exhibits linear dynamics with constant slope).

<sup>12</sup> For example, if a system has underwent a sudden and unexpected changes of its dynamics in the past it has a low "inertia". In this case any long term prediction of the future system's behaviour is not very credible.

<sup>13</sup> Some sort of stationarity is required by statistical models applied for making predictions of the future system's behaviour. That way they avoid the question of the credibility of such predictions – their uncertainty may be growing in time but, due to stationarity, the dynamics of the system does not change in any limited time horizon. On the contrary, PL method aims at indicating time horizon within which the system's behaviour is sufficiently well described – thus assumptions are significantly weaker. Cf. Table 2 for further discussion.

It is also important to note the fact that PL method is data – driven (i.e. is based only on the sample of historical observations) implies also that it adopts a conservative view of the system. Namely, it cannot anticipate systemic surprises and behaviors which had not occurred in the period over which the sample of historical observations was collected.

### 2.1. Generic notion of the explainable outreach of the data

The core idea of the PL approach is to **deduce directly from the data** their **explainable outreach (EO)**, i.e. the spatial and temporal extent beyond which we cannot explain the considered system only by the available knowledge about its past. The explainable outreach is characterized by four key attributes: (i) the instant of time in which it begins; (ii) diagnostic uncertainty of the state variable describing the system in this initial moment (defining the initial opening of explainable outreach); (iii) increase of prognostic uncertainty in time; and (iv) temporal extent (quantifying the time in the future beyond which the system’s behavior cannot be shown anymore to be in accordance with its past behavior).

Explainable outreach can be seen as a region in product space of time and the space to which values of the observations belong (i.e. real line). This region is induced by our understanding of the system (for example expressed in form of trend function). Its spatial boundaries are given by uncertainties related to the projection of our understanding of the system into future (e.g. prediction bands<sup>14</sup> centered around the extrapolated trend – to continue example), while its temporal extent is characterized by the moment in which this projection starts and the time horizon within which the uncertainty region covers the trajectory of the system.

Obviously, different hypotheses about the type of trend the system follows will result in different explainable outreaches. Some of them may be very long and wide (system’s behavior is described robustly but very imprecisely) or short and narrow (when our understanding of the system is quite precise but only locally correct). One would prefer the EO to be as long and at the same time as narrow as possible.

Comparison of different EOs derived for the same sample may be facilitated by a score assigning a numeric value to the combination of EO attributes (i) – (iv). For example, one could use the following

$$\text{Score of EO} = \frac{\text{Length of temporal extent of EO}}{\text{Width of EO at its end}}$$

Such score increases as the length of EO increases or its width decreases, thus one would prefer EO for which this score is the highest.

### 2.2. Prognostic learning procedure

Notice that an explainable outreach as defined above expresses our expectations about the consequent system’s behavior from a certain fixed instance of time on. Due to data

---

<sup>14</sup> For each instant of time prediction bands give the range which is expected to contain with predefined probability (called confidence level) an observation taken at that time. In contrast, confidence bands give range within which we expect to cover a true expected value of an observation. In this report we prefer to use prediction bands since we want to test our understanding of the system with individual data points.

variability and possible imprecision of our understanding of the system, an EO starting in another instance of time may have a different shape and length. Therefore, to gain some understanding of patterns of system's behavior it is insufficient to look at just one EO. One should rather derive this understanding from a sequence of consecutive EOs resulting from a learning procedure.

Below we provide a generic procedure of learning in a controlled prognostic context given the learning sample  $X_0, \dots, X_T$  of observations of the analyzed system collected over the period  $[0, T]$ :

1. Choose a suitable set of hypotheses (e.g., a family of regression functions) about the rules governing the system behavior and the minimal number  $k$  of data points required to select the one which represents the system best.
2. Choose the initial length  $\tau = k$  of the subsample  $X_0, \dots, X_\tau$ , which we call the learning block (LB)
3. Choose the hypothesis which reflects the system's behavior best in the learning block  $X_0, \dots, X_\tau$  (e.g., estimate parameters of the regression function) and quantify its uncertainty (e.g., with use of prediction bands)
4. Find the EO starting at point  $\tau$ . To determine the shape of the EO calculate the uncertainty region  $R \subset [\tau, \infty) \times \mathbb{R}$  spanned by the prediction of the future unfolding of the system based on hypothesis chosen in in point 3 and its uncertainty. To determine the length of the EO project the remainder of the data  $X_{\tau+1}, \dots, X_T$ , which we call testing block (TB), onto region  $R$  and find the largest  $H$  such that<sup>15</sup>

$$\forall \tau < t \leq \tau + H \quad (t, X_t) \in R$$

If  $H < T - \tau$  then length of the EO starting at point  $\tau$  is set to  $H$ ; else it is set to  $\infty$ .

5. If  $\tau < T$  then set  $\tau = \tau + 1$  and go to step 3; else end procedure.

The above procedure explains the meaning of the name "learning in a controlled prognostic context": we learn about the patterns of the past system behavior (step 3.) and then test this knowledge applying it in a prognostic mode in the controlled context of the remainder of the data sample (step 4.).

Assessment of the temporal extent of the EO,  $H$ , from step 4 of the learning procedure requires a discussion. It is either finite (not longer than the historical sample itself) or set to infinity. In the first case finite time horizon of the EO indicates limits within which we can sufficiently well predict system's evolution after time  $\tau$  by means of the method selected in step 1 to grasp the system's dynamics in the learning block  $X_0, \dots, X_\tau$ . In other words, it indicates the limits to credibility of predictions of the system's behaviour after time  $\tau$  based on our understanding of the system's dynamics given the knowledge carried by the subsample  $X_0, \dots, X_\tau$ . On the other hand, infinite time horizon indicates that we are unable to falsify this understanding of the system's behaviour with use of the testing part of the sample  $X_{\tau+1}, \dots, X_T$  (i.e. we have no ground to reject our

---

<sup>15</sup> If the hypothesis about the system's behaviour is formulated in terms of a regression model, the requirement that all points between time  $\tau$  and  $\tau + H$  belong to  $R$  may be relaxed – only a sufficient portion of these points should fall into  $R$ .

hypothesis about the system's nature). There are two possible reasons for such situation to occur: either our understanding of the system is exceptionally good or the testing sample is too short to provide evidence against it<sup>16</sup>. This indicates an important constraint of the PL approach (indeed, of any data - driven method), namely that data resources (length of sample of historical observations) set limits to the level of detail<sup>17</sup> with which we wish to describe analyzed system.

### 2.3. Applying the prognostic learning procedure and interpretation of its results

Learning in a controlled prognostic context is essentially a model – independent paradigm of exploratory data analysis. By this we mean that it does not presuppose any particular model which reflects our *a priori* knowledge<sup>18</sup> or belief about the analyzed system, and which may be calibrated on the sample of historical observations and then used for making predictions. On the contrary, PL approach is purely data – driven: we explore a sufficiently broad family of alternative methods of describing the system's behavior (e.g., different types of regressions) by running a PL procedure (cf. section 2.2) for each of them and then select the one which yields the best explainable outreaches.

After completing this task, we obtain a sequence of explainable outreaches indexed by their starting moments  $\tau = k, k + 1, \dots, T$ . Technically, this inform us how credible our predictions based on partial knowledge about the system<sup>19</sup> were over the time interval  $[0, T]$ . In particular it provides no confirmed (tested) information about the explainable outreach starting at time  $T$  and expressing our expectations about the immediate future of the system. This cannot be done formally without additional and restrictive assumptions (e.g., stationarity of the system), however, such exercise still may be informative. If only the behavior of the EOs over the period  $[0, T]$  was regular enough (i.e., EOs have comparable scores, implying similar lengths and widths) and the last  $\tau$  for which EO has finite length is sufficiently close to  $T$  we may attempt to extrapolate the characteristics of (tested) EOs to formulate expectations about likely shape and temporal extent of the (untested) EO starting at time  $T$ .

In principle, the results of PL method give us insight into system's "inertia". Such information may be useful for decision makers trying to influence future behavior of the system (e.g., mitigate global warming by implementing certain policies). Firstly, it indicates likely directions of future system evolution under "business as usual" conditions<sup>20</sup> which is a reference point with respect to which any policy is formulated. Moreover, it indicates the time horizon within which we may have some confidence in quality of predictions based on our understanding of the system. Secondly, it indicates

---

<sup>16</sup> Falsifying a good hypothesis may require a very long testing sample. In the extreme (but very unlikely) case, when we perfectly understand our system (i.e. know the process generating data – both in the past and in future) we wouldn't be able to falsify it with use of any test sample of finite length.

<sup>17</sup> Understood as complexity of the hypothesis about the system's dynamics.

<sup>18</sup> Additional knowledge (e.g. about a particular type of dynamics the system follows) obtained beforehand from some other source than the learning sample  $X_0, \dots, X_T$ .

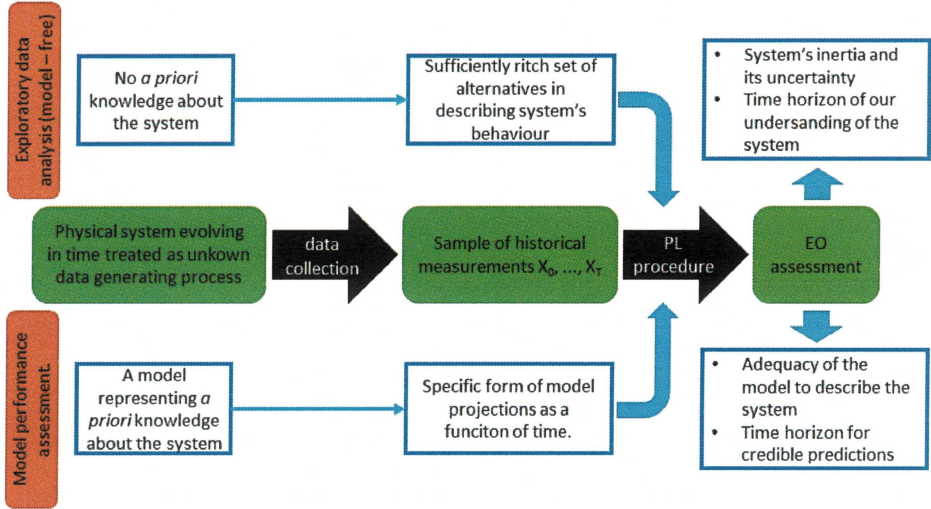
<sup>19</sup> i.e. knowledge carried by learning blocks  $X_0, \dots, X_\tau, \tau < T$ .

<sup>20</sup> i.e. in situation when current dynamics of the process and external forcing / policies / measures will not change.



the strength of the measures needed to overcome the system's inertia and to shift its future evolution towards the desirable path<sup>21</sup>.

PL methodology may also be applied to assess scenarios produced by a particular model of the system of interest. If a scenario falls out of the EO before its end it means that the model predicts a change in the system's dynamics (w.r.t. its past behavior). If so, then modeler should explain what is the reason for that, e.g. what significant changes the system is expected undergo under that scenario. If the future trajectory under "business as usual" scenario falls outside the EO it may indicate inadequacy of the model to describe the system of interest.



**Figure 2. Two modes of applying the learning in a controlled prognostic context paradigm.** In exploratory data analysis mode the selection of the best method to represent system's behavior and construct EO is purely data driven without use of any *a priori* knowledge. EO indicates the inertia of the system and uncertainty and time horizon of our understanding of the system. In model assessment mode a model-specific form of a trend function is fed into the PL procedure in order to assess model's ability to accurately describe the system and to quantify limits to its predictions. (This mode is not considered in this report).

We also speculate that a modification of the PL method may be applied to assess a particular model and its projections even more directly. If it is possible to express the model prediction as a function of time of a certain form (dependent on initial conditions and values of exogenous parameters) and calculate a region spanned by the projection and its uncertainty one can use this function directly in the prognostic learning procedure (see section 2.2). Then resulting EOs could indicate the time horizon within which the model is sufficiently adequate to describe the system's evolution. However, this generic approach would require designing of a model-specific implementation of

<sup>21</sup> If the system's trajectory under a scenario corresponding to introduction of certain policy stays within the explainable outreach it indicates that the effectiveness of such policy remains uncertain within the time horizon of this EO.

the PL procedure in order to make it operational. This modification of PL approach has not yet been tested and will not be covered in this report.

#### **2.4. PL versus forecasting with use of time series analysis**

The variant of the prognostic learning paradigm discussed in this report treats the data forming a time series. It is, however, quite different from the commonly used time series analysis (TSA) methodology. While PL trades only approximate understanding of the behavior of the data itself for ability to indicate limits to this understanding and generality of the method, the TSA strives for complete understanding of the data generating process and applying this knowledge for making predictions.

Typically, TSA is based on decomposition of the time series into deterministic component (functional trend, seasonal component, oscillations) and stochastic part. The deterministic part can be estimated from the data with use of broad range of various techniques (such as regressions, curve fitting, smoothing methods, wavelet analysis etc.) The overarching goal is to estimate the deterministic part so that it fits the data as close as possible while its extrapolation properties are a lower priority concern. Nature of the stochastic part is inferred from the behavior of residuals (i.e. remaining part after removing the estimated deterministic component from the data). This is usually done by fitting a suitable time series model (such as ARIMA or GARCH).

Obviously, the estimate of the deterministic component of the time series significantly influences the behavior of residuals and thus the statistical model of the stochastic part. As the latter may be quite complex and difficult to estimate (e.g., due to scarcity of the data resources w.r.t. number of parameters in the model) the problem of estimation of the deterministic component is somewhat subordinate to analysis of residuals. Estimate of the deterministic part is expected to produce residuals for which the statistical model is as simple as possible. The literature of the subject puts much more emphasis on the statistical models of the residuals, typically assuming that the deterministic component of analyzed time series has already been removed with use of some suitable technique (e.g., Brockwell & Davis 2002).

Once the time series is described in terms of deterministic function of time and statistical model of residuals one may use this knowledge for making forecasts. In order to do so the deterministic trend is extrapolated and the behavior of the stochastic part (i.e. residuals) is either determined theoretically (e.g., prediction bands obtained under stationarity assumptions) or simulated (using the statistical model of residuals). However, such forecasts should be considered with caution. Technical problems may arise due to incorrect structure of the model of stochastic part and/or bad extrapolation properties of the function describing deterministic component (such as instability due to uncertainty in estimated values of function parameters). Some techniques of describing the deterministic part such as smoothing splines even rule out the possibility of extrapolation. Moreover, when making forecasts the description of the analyzed time series (i.e. deterministic function plus statistical model of residuals) are treated as the true process generating data which will never change. As a result, indicator of a time horizon within which the predictions are credible cannot be derived from TSA methodology.

We conclude this section with Table 1 summarizing differences between PL method and TSA.

**Table 1.** Prognostic learning versus time series analysis.

	Learning in a controlled prognostic context	Time series analysis	
		Deterministic component	Stochastic component
<b>Approach</b>	Data – driven exploratory analysis. Emphasis on striking balance between approximate understanding of the system and ability to indicate the limits to this understanding.	Inferring the data generating process. Emphasis on statistical model of the stochastic component, while estimate of deterministic component is to yield desired statistical properties of the residuals.	
<b>Assumptions</b>	No systemic surprises (behaviors unobserved in the past will not happen in the future)	Particular form of trend function.	Particular form of the dependence structure / model of residuals. Usually also normality and weak stationarity of residuals is required.
<b>Principle</b>	<b>Optimization of the EO.</b> Selecting the type of trend generating the longest and narrowest EO.	Fitting a function minimizing <b>in-sample error</b> .	<b>Estimation from the data values of the model parameters that minimize expected forecast error.</b>
<b>Measure of performance</b>	Score of the explainable outreach	Typically sum of squared errors or mean squared error	Typically expected mean squared error
<b>Predictions</b>	Data-driven model describes the system only approximately correctly and uncertainty of predictions inevitably grows in time. <b>The method does not strive for perfect predictions. It aims to understand their limits.</b>	Within the range of observed sample the fitted function is interpreted as expected value of observations. Extrapolation of fitted function beyond the range of sample may be interpreted in the same way but there is no possibility for controlling the error of predictions with use of such extrapolation.	Future behavior of the stochastic component (typically expressed in form of prediction or confidence bands) is derived from the statistical model of residuals either theoretically (usually under assumption of stationarity) or by means of simulations utilizing model structure.
<b>Time horizon within which forecasts are supposed to be reliable</b>	Expected length of the EO based on the assessment of the results of the prognostic learning procedure.	Unknown. Fitted model of the time series (i.e. estimated deterministic component and statistical model of the stochastic part) is treated as the true data generating process and as such universally correct.	

<b>Sources of uncertainty</b>	(1) Diagnostic uncertainty (measurements errors) reflected by initial opening of the EO; and (2) prognostic uncertainty which grows into the future reflected by the shape of the EO	(1) Uncertainty in the form of the function describing deterministic component; and (2) uncertainty in the parameter estimates.	(1) Uncertainty in estimate of deterministic component defining residuals; (2) uncertainty of structure of model of residuals; and (3) uncertainty of estimates of model parameters.
-------------------------------	--	---	--

### 3. Regression – based construction of the EO

In this chapter we propose a practical method of implementing the generic paradigm of learning in a controlled prognostic context presented in Chapter 2. Making this generic notion operational requires addressing the following problems:

1. Grasping the behavior of the data from the learning block and quantifying the diagnostic uncertainty in order to specify direction and initial width of the EO.
2. Defining the shape of the explainable outreach (i.e. its spatial boundaries).
3. Determining the length of the EO by testing it against the data from the testing block.

Below we propose a solution to these questions which is based on the regression techniques.

**Ad 1.** The trend in the data is grasped by means of a regression function fitted to the points from the learning block. For each moment  $t$  belonging to the learning block the value of regression function at that moment is interpreted as the expected value of the observation taken at time  $t$ . The extrapolation of the regression function defines the main axis around which the EO is constructed. The diagnostic uncertainty is expressed as standard deviation of residuals (i.e. differences between the regression function and the actual observations) and defines the initial width of the explainable outreach.

**Ad 2.** The shape of the EO (i.e. its upper and lower band) is given by extrapolation of the prediction bands calculated for the regression model fitted to the learning block.

**Ad 3.** Given the shape of the EO its length is determined by projecting remainder of the learning sample (i.e. testing block) onto it. The moment in which the EO ends is defined as the earliest moment for which the position of the testing points with respect to the EO starts to be very unlikely if the regression model fitted over the LB is correct and true also beyond its range.

The details of the proposed solution depend on the specific regression technique to be applied. In the remainder of this section we give these details for the prognostic learning procedure based on the polynomial regression. In Appendix B we present an alternative PL procedure based on local linear regression method.

#### 3.1. Analysis of historical patterns in learning phase with use of polynomial regression

The polynomial regression is a widely used parametric technique of data analysis. Its popularity comes from the fact that it is a relatively simple and straightforward generalization of the classic linear regression method as well as from the flexibility of

the family of polynomial regression functions<sup>22</sup>. It is also a popular technique of estimating the deterministic part of a time series (Brockwell & Davies 2002).

In order to grasp the historical trend in the learning block we use a model of polynomial regression of order  $p$

$$x(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \dots + \alpha_p t^p + \varepsilon_t$$

where  $x(t) = X_t$  is a value of the observation taken at time  $t$  and the noise term  $\varepsilon_t$  is normally distributed with zero mean and standard deviation  $\sigma$ . Moreover, we assume that  $\varepsilon_t$ ,  $t = 0, 1, 2, \dots$ , are independent and identically distributed.

Let the learning block contain  $n$  observations taken in times  $t_1, \dots, t_n$ . We estimate parameters of the regression function

$$\hat{x}(t) = a_0 + a_1 t + a_2 t^2 + \dots + a_p t^p$$

with use of the ordinary least squares (OLS) method (Wolberg 2006: chapt. 2). The uncertainty of the fitted regression function at time  $t$  is then given by formula

$$s_x(t) = \sqrt{\frac{\sum_{i=1}^n (\hat{x}(t_i) - x(t_i))^2}{n - (p + 1)} \sum_{j=1}^{p+1} \sum_{k=1}^{p+1} t^{j+k-2} [C^{-1}]_{j,k}}$$

where  $[C^{-1}]_{j,k}$  is the entry at the cross-section of the  $j$ -th row and  $k$ -th column in the inverse of matrix

$$C = \left[ \sum_{i=1}^n t_i^{j+k-2} \right]_{\substack{j=1, \dots, p+1 \\ k=1, \dots, p+1}}$$

The diagnostic uncertainty over the learning block is assumed to be constant and is estimated as a standard deviation of the model residuals

$$s_r = \sqrt{\frac{\sum_{i=1}^n (\hat{x}(t_i) - x(t_i))^2}{n - (p + 1)}}$$

Upper and lower prediction bands at the confidence level  $(1 - \alpha)$  for the observations taken at time  $t$  are then given by the formulas

$$f_{up}(t) = \hat{x}(t) + t_{n-(p+1), 1-\alpha/2} \sqrt{s_x(t)^2 + s_r^2}$$

and

$$f_{low}(t) = \hat{x}(t) - t_{n-(p+1), 1-\alpha/2} \sqrt{s_x(t)^2 + s_r^2}$$

respectively, where  $t_{n-(p+1), 1-\alpha/2}$  is  $(1 - \alpha/2)$  quantile of the t-Student distribution with  $n - (p + 1)$  degrees of freedom. Notice that parameter  $\alpha$  regulates the width of the prediction bands (the lower the  $\alpha$  the wider the prediction bands). Observe also that distance between prediction bands, i.e.  $f_{up}(t) - f_{low}(t)$ , increase with  $p$ -th power of  $t$ .

---

<sup>22</sup> Indeed, any continuous trend in the data can be locally approximated with arbitrary precision by a polynomial of sufficiently high order.

### 3.2. Construction of the EO

The explainable outreach starts at time  $\tau = t_n$ , i.e. the moment in which the last observation of the learning block was taken. The EO is built around the extrapolated polynomial trend fitted to the data in the learning block, that is around  $\hat{x}(t), t \geq \tau$ . Its initial width is defined as  $f_{up}(\tau) - f_{low}(\tau)$  and is determined by the diagnostic uncertainty  $s_\tau$ . The shape of the EO, i.e. its upper and lower band are given by functions  $f_{up}(t)$  and  $f_{low}(t)$  for  $t > \tau$ , that is the prediction bands for regression model extrapolated beyond the learning block.

Notice that in order to define the initial width and the shape of the EO only the information about the system's behavior in the learning block is needed. However, to determine its temporal extent (time horizon) additional knowledge carried by the remainder of the learning sample (testing block) is required. This remaining subsample is used to determine until when our expectations about the future system's evolution after time  $\tau$  represented by the EO (given only the knowledge contained by the learning block) are in accordance with the actual evolution of the system after that time.

To explain how we determine the moment in which the EO cease to be in accordance with the actual system's evolution let us assume for a while that we know the evolution of the analyzed process only up to the moment  $\tau$  and the  $m$  remaining points in the testing block  $(t_1, X_1), \dots, (t_m, X_m)$ ,  $t_1 = \tau$ ,  $t_m = T$ , are unknown. In addition, let us define an auxiliary sequence of random variables

$$E_k = \begin{cases} 0 & \text{if } X_k \notin [f_{low}(t_k), f_{up}(t_k)] \\ 1 & \text{if } X_k \in [f_{low}(t_k), f_{up}(t_k)] \end{cases}$$

where  $(t_1, X_1), \dots, (t_m, X_m)$  are the yet unknown points from the testing block.

Now observe that if the regression model fitted to the learning block correctly describes the evolution of the analyzed process then also the points from the testing block should follow this model. If that is so, then by definition of the prediction bands at the confidence level  $(1 - \alpha)$  the probability that the future observation taken at time  $t \geq \tau$  will fall into interval  $[f_{low}(t), f_{up}(t)]$  is equal to  $(1 - \alpha)$ . Thus  $E_k = 1$  with probability  $(1 - \alpha)$  and  $E_k = 0$  with probability  $\alpha$ . In other words, all  $E_k, k = 1, \dots, m$  follows the Bernoulli distribution with parameter  $(1 - \alpha)$ <sup>23</sup>. Moreover, if the regression model fitted to the learning block is correct also for the observations in testing block, then these observations are independent. Therefore, all  $E_k, k = 1, \dots, m$  are not only identically distributed but also mutually independent. As a consequence, for each  $k = 1, \dots, m$ , a random variable

$$S_k = \sum_{i=1}^k E_i$$

---

<sup>23</sup> Random variable  $X$  follows the Bernoulli distribution with parameter  $p$  if  $P(X = 1) = p = 1 - P(X = 0)$ . Random variable  $X$  is the outcome of a so called Bernoulli trial, i.e. a random experiment with only two possible results: success (coded as 1) which occurs with probability  $p$  or failure (coded as 0) which happens with probability  $(1 - p)$ .

has a binomial distribution  $B(k, (1 - \alpha))^{24}$ .  $S_k$  may be interpreted as the number of points among the first  $k$  points of the testing block which falls into the prediction bands.

In order to determine the length of the EO<sup>25</sup> we confront our expectations based on fitted regression model about the distribution of future observations (formulated above) with the actual observations from the testing block, denoted by  $(t_1, x_1), \dots, (t_m, x_m)$ . Let  $e_1, \dots, e_m$  be the actual values of the random variables  $E_1, \dots, E_m$  and let for each  $1 \leq k \leq m$

$$s_k = \sum_{i=1}^k e_i$$

be the actual number of points among the first  $k$  points of the TB which fall into the prediction bands. Recall that if our regression model is true,  $s_k$  should follow the binomial distribution  $B(k, (1 - \alpha))$ . This is the key observation allowing us to find the temporal extent of the EO. Namely we set the end of the EO to be the first moment  $t_k$  for which actual value of  $s_k$  is an unlikely outcome given our understanding of the past of the process (represented by the fitted regression model). The observed value  $s_k$  is considered unlikely if the joint probability of such event and all not more probable (i.e. all events in which from the first  $k$  points of the TB only  $s_k$  of them or less fall into the prediction bands) is less than some suitably selected low threshold  $p_0$ . For the sake of consistency, we use  $p_0 = \alpha$ .

To summarize the above argument we present the algorithm for finding the length of the EO:

1. Select threshold  $p_0$  (e.g. equal to  $\alpha$ ) and set  $k = 1$ .
2. Calculate  $s_k$  (i.e. the number of points among the first  $k$  points of the TB which fall into the prediction bands).
3. Let  $F_{k,(1-\alpha)}$  be the cumulative distribution function of the binomial distribution  $B(k, (1 - \alpha))$ . If  $F_{k,(1-\alpha)}(s_k) < p_0$  then we set the end of the EO to the moment  $t_{k-1}$ , its length  $H$  to  $k - 1$  and we stop the algorithm.
4. If  $k = m$  (i.e. testing block is exhausted) then we cannot determine the end point of the EO. We stop the algorithm and set EO length  $H$  to  $\infty$ .
5. Set  $k = k + 1$  and go to point 2.

### 3.3. Procedure of prognostic learning based on regression method

To wrap up the present chapter, below we provide the procedure for prognostic learning based on the regression techniques presented above. It is a method – specific version of the generic PL procedure formulated in Section 2.2.

---

<sup>24</sup> Binomial distribution  $B(n, p)$  is a distribution of a number of successes in the  $n$  independent Bernoulli trials with probability of success  $p$ .

<sup>25</sup> i.e. the time horizon within which we have no reason not to believe that the actual observations are in agreement with our understanding of the system's dynamics based on its past.

1. Choose the regression technique (e.g., polynomial regression of certain order) which will be used to grasp the data behavior in the learning block.
2. Choose the initial length  $k$  of the learning block  $X_0, \dots, X_\tau$ ,  $\tau = k$ . (Note that  $k$  should be large enough with respect to the complexity of selected type of regression function in order to ensure good estimates of the trend function parameters and to prevent overfitting<sup>26</sup>.)
3. Fit the regression model to the learning block  $X_{\tau-k}, \dots, X_\tau$ .
4. Construct the EO starting at time  $\tau$  following the guidelines presented in Section 3.2 and determine its length  $H$ .
5. If  $\tau < T$  set  $\tau = \tau + 1$  and go to step 3. In the opposite case end the procedure.

Notice that in step 3 we ignore a part of the learning block  $X_0, \dots, X_\tau$  discarding all but last  $k$  points. In effect, at each stage of the learning procedure we fit a regression model to the data points falling into a window of fixed length  $k$ , which we move along the learning sample in course of the learning procedure. We call this version of PL method “rolling window”. Using window of fixed length is advantageous in two ways. First, it allows for easier comparison of EOs at different stages of the PL procedure, since width of each EO is determined not only by the uncertainty of the regression model but also by the number of points used for fitting this model. If this number is fixed, the widths of EOs depends only on appropriateness of regression model to grasp the data behavior in corresponding learning blocks. Secondly, by using only  $k$  last points from each learning block makes the method more responsive to the local behaviour of the data, acknowledging that the recent data points are more relevant to the direction of the EO than the points from the beginning of the learning sample. Throughout this report the “rolling window” learning procedure will be used<sup>27</sup>.

We conclude this chapter by emphasizing that the formulas for estimates of prognostic diagnostic and prognostic uncertainty as well as for prediction bands defining the shape of the EO given in Section 3.1 are applicable exclusively to polynomial regression. However, method of constructing the EO described in Section 3.2 and prognostic learning procedure given in Section 3.3 are readily applicable to any type of regression method for which the prediction bands can be calculated and extrapolated beyond the range of the LB. (Note, however, that the assumption on independence of residuals of the fitted regression model must be satisfied). For example, these sections are immediately applicable to the prognostic learning procedure based on non-parametric regression (as demonstrated in the Appendix).

---

<sup>26</sup> i.e. situation, in which flexible trend function is not sufficiently constrained by short sample of data points and too closely mimics the random layout of the data points. Overfitting has strong negative impact on the quality of model predictions.

<sup>27</sup> Another version of the PL method which at each stage makes use of the whole learning block is equally easy to implement as the “rolling window” procedure (in step 3 of the procedure one only needs to fit a model to all points  $X_0, \dots, X_\tau$  instead of the last  $k$  ones). We call this version “expanding”. It is useful when we want to check whether the selected regression model is able to correctly grasp the system’s dynamics over the whole period covered by the learning sample. This method is also used in the Appendix where we employ nonparametric regression techniques to grasp the behaviour of the data in the learning block. As these methods use only local information (i.e. regression curve is determined only by the nearby points, not the whole sample) the effect of increasing length of LBs on the EO (especially in its width) is negligible.



## 4. Assessment of prognostic learning performance in the controlled conditions. Monte Carlo experiments

Before we apply the prognostic learning procedure based on the polynomial regression (described in the previous chapter) in analysis of the real-life data we first test its performance under controlled conditions, that is, we conduct Monte Carlo experiments by repetitively running PL method on synthetic data sets.

Having full knowledge about the true trend in the synthetic data and the control over the strength of noise disturbing that trend allow us to clearly identify strengths and weaknesses of the PL method as well as their reasons. This enables us to draw useful conclusions and to formulate guidelines for applying the PL method in analysis of the real-life data.

By choosing to work with synthetic data we overcome a problem of data scarcity, which often occurs when working with real-life data. Real-life data sample is often too short to support application of PL method of higher order<sup>28</sup>, whereas synthetic data sample may be of any desired and suitable length. In addition, we may always afford to have additional sample used exclusively for testing our expectations about the length of the EO starting at the end of the learning sample. Moreover, we can generate multiple independent data samples following the same fixed deterministic trend and compare the performance of the PL method applied to each of them. This gives us the ability to study the stability of the method. In addition, we can repetitively compare the predicted and actual lengths of the EO starting at the end of the learning sample in order to test the extent to which we can use the insight given by the PL method about the dynamics of the observed system to inform us about its immediate future.

In the present chapter we describe the method which we use to generate synthetic data samples used for testing the PL method in controlled conditions, purpose and setup of performed numerical experiments and their results. We conclude this chapter with some general observations and guidelines of applying the prognostic learning procedure based on the polynomial regression.

### 4.1. Method of generating the synthetic data

The synthetic data samples were generated in the following way:

1. We choose the length of the sample  $N$ . For simplicity we assume that  $t_k = k$ ,  $1 \leq k \leq N$ , where  $t_k$  denote times for which synthetic observations are generated.
2. We choose a suitable trend function  $f$  which synthetic data will follow.
3. We choose the strength of the noise with which we disturb the true trend  $f$ . This strength is defined by the standard deviation  $\sigma$  of the noise, which we express as

---

<sup>28</sup> Learning block required for good estimation of parameters of higher order polynomial trend may be of comparable length as the whole learning sample leaving too few points for meaningful testing of the explainable outreach

a percentage of the width of range of the trend function values<sup>29</sup> (for example,  $\sigma = 0.01 \times (\max_{1 \leq k \leq N} f(t_k) - \min_{1 \leq k \leq N} f(t_k))$ ).

4. We generate a synthetic sample  $(t_k, x_k)$ ,  $1 \leq k \leq N$ , by setting  $x_k = f(t_k) + \varepsilon_k$ , where  $\varepsilon_1, \dots, \varepsilon_N$  is a sequence of independent random variables following normal distribution of zero mean and standard deviation  $\sigma$ .

In Section 4.3 we present result of running the PL method on five different synthetic data sets. Two of them follow polynomial trends which belong to the family of regression functions used in the employed regression method. Namely these are: the linear trend and the 4<sup>th</sup> order polynomial trend. They were selected in order to test the performance of the PL method on trends of low (linear) and high (4<sup>th</sup> order polynomial) complexity in nearly ideal conditions<sup>30</sup>, where polynomial regression may give an unbiased<sup>31</sup> model fit.

The remaining three synthetic data sets do not follow trends of polynomial type, thus allowing us to test the performance of the PL method in situations where the employed regression technique is not able to reproduce the true trend in the data (i.e. it provides only a biased estimate of the true trend). Moreover, they are intended to mimic the types of behavior often encountered in the real-life data. Namely, considered synthetic samples follow: exponential trend (increasing trend whose rate of increase accelerates), logarithmic trend (increasing but with decreasing slope) and sinusoidal with long period of oscillations, comparable with the length of the sample (to mimic a situation when apparent local trends in the historical data are in fact results of slow, long-term oscillations).

Before we present the actual results of applying the PL method on the abovementioned synthetic data samples, in the following section we describe the setup and details of performed experiments.

## 4.2. Description of experiments on synthetic data

The numerical experiments we perform for each of the abovementioned types of synthetic data involve multiple Monte Carlo runs of the “rolling window” variant of the polynomial regression based PL procedure. Each of the experiments corresponds to a fixed combination of value of order of the method (i.e. the degree of polynomial used in the regression model), level of noise and length of the learning block.

Objectives of these experiments are two-fold. Firstly, we want to identify situations (i.e. patterns in the local behavior of the data forming the learning block and the strength of the noise) in which the proposed method of prognostic learning presents its strengths or performs poorly. Secondly, we investigate the influence of order of the PL method,

---

<sup>29</sup> Expressing the strength of noise in relation to the range of the true trend function instead of in absolute terms allows us for easy comparison of different types of synthetic data samples.

<sup>30</sup> In principle, in noiseless conditions it would be possible to determine both past and future behaviour of the data given only relatively few points in the LB.

<sup>31</sup> We say that estimator is unbiased if its expected value is equal to the estimated quantity. In case of regression methods, we say that fitted trend  $\hat{f}$  is unbiased estimate of true trend  $f$  if  $E(\hat{f}(t)) = f(t)$  for all  $t$  within the range (period) of the sample. Fitted regression model is necessarily biased if the true trend does not belong to the family of considered regression functions.

strength of the noise and length of the learning block on the performance of the PL method.

In addition, we explore the reliability of predictions of future EO lengths both in-sample (i.e. using the actual EO lengths<sup>32</sup> in stages up to the present one in order to predict the EO length in the next stage of the PL procedure) as well as out-of-sample (i.e. using EO lengths calculated for all stages of the PL procedure in order to predict the length of the EO starting at the end of the learning sample on which the PL procedure was run). In both cases predictions are made by fitting the linear function (with use of the OLS method) to all available (finite) values of the past EO lengths and then extrapolating it to the future point of interest<sup>33</sup>.

Notice that in-sample predictions may be compared against the actual EO lengths calculated in due course of the learning procedure. Testing prediction of EO length out-of-sample in similar way, however, requires additional testing sample back-to-back with the learning sample used in the PL procedure. Obtaining such sample is not a problem for the synthetic data – one can easily generate it.

Observe also that for a single learning sample and corresponding additional testing sample one can only get one pair of predicted and actual EO lengths starting at the end of the learning sample. However, both values may be to large extent random and only one such pair is not very informative. Much more information carries their joint distribution. Working with synthetic data allows us to easily obtain an empirical estimate of such distribution by means of repetitive Monte Carlo simulations.

Below we describe the procedure that each of experiments follows:

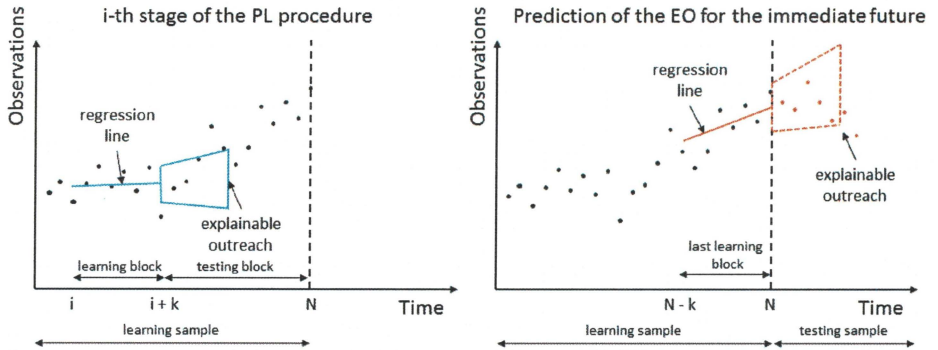
1. Select the functional trend which the synthetic data sample will follow. Choose the length  $N$  of the learning sample and the strength of the noise.
2. Select the order of the PL method and the length of the learning block (window)  $k$  to be used.
3. Select the number of repetitions of the experiment  $M$ .
4. Set the current iteration (Monte Carlo run) number  $i$  to 1.
5. Generate the synthetic data sample of length  $2N$  (cf. Section 4.1). Use the first  $N$  points as a learning sample for PL procedure and the remaining data as the additional testing sample to be used exclusively for determining the actual length of the EO starting at the end of the learning sample.
6. Run the “rolling window” prognostic learning procedure on the learning sample generated in step 5. At each stage of the procedure check the fulfilment of assumptions of the polynomial regression model fitted to the learning block and record the score of the EO, its actual length and the predicted EO length for this stage given the EO lengths for previous stages (cf. Figure 3, left panel).

---

<sup>32</sup> Actual EO length is the length of the EO determined with use of data from the testing block. In contrast, predicted EO length is just our (untested) expectation about the length based on the knowledge of actual lengths of EOs from previous stages of the learning procedure.

<sup>33</sup> This is just one, straightforward but possibly crude way of making such predictions. Application of some more subtle methods (e.g., time series model) may improve reliability of such predictions. This will be tested in course of future research.

7. After the PL procedure is complete use the calculated EO lengths (in-sample) to predict the length of the EO starting at the end of learning sample (out-of-sample).
8. In order to test the predicted length of the EO starting at the end of learning sample (cf. step 7) calculate the actual length of the EO starting at the end of this sample. To do so, take the learning block consisting of the last  $k$  points of the learning sample, fit a regression model to it and extrapolate the prediction bands to determine the shape of the EO. To find its length use the data from the additional testing sample (cf. Figure 3, right panel).
9. If  $i < M$  then set  $i = i + 1$  and go to step 5. Otherwise end the experiment.



**Figure 3. Schematic picture of the Monte Carlo experiment.** **Left panel:** One stage of the prognostic learning procedure with “rolling window” of length  $k$ . Regression model is fitted to the data forming a learning block  $[i, i + 1, \dots, i + k]$ . Prediction bands for this model define the shape of EO starting at  $i + k$ . Actual length of the EO is determined with use of the data from the testing block. **Right panel:** Determining the actual length of the EO starting at the end of the learning sample (prediction for the immediate future). The direction and shape of the EO is given by the last  $k$  points from the learning sample (last learning block). Since there are no points left in the testing sample to form a testing block, the actual length of the out-of-sample EO is determined with use of the additional testing sample.

With use of the insights gathered by performing abovementioned experiments we formulate the guidelines for selecting the order of the method and length of the LB yielding optimal performance of the PL method. By this we mean:

- (1) Satisfactory level of fulfilment of the assumptions of regression model fitted to each learning block.
- (2) As long and narrow EOs calculated at different stages of PL method as possible (i.e., with high score - cf. Section 2.1). Stable behavior of EO lengths at different stages of the PL procedure is desirable.
- (3) Ideally, good reliability of predictions of the EO lengths (both in-sample and out-of-sample).

### 4.3. Results

In this section we present the results of five sets of Monte Carlo experiments on five different types of synthetic data. This allows us to assess usefulness of the proposed

methods of prognostic learning under controlled conditions. In each set of experiments we investigate the influence of: (1) order of the method, (2) length of the learning block and (3) level of noise on the performance of prognostic learning, by varying these parameters. Below we present results only for Monte Carlo runs of the PL methods on synthetic data with low level of noise<sup>34</sup>. For each considered order of method the optimal length of the learning block is presented. General conclusions about marginal influence of each of the three abovementioned factors on the performance of PL method are presented in Section 4.4.

#### 4.3.1. Data following a linear trend

We begin our analysis of performance of the PL method with testing it in the simplest possible setting, i.e. on the synthetic noisy data following a linear trend. Such type of trend in the data is easily detected and robustly estimated with use of OLS technique, even for relatively short samples. Hence even the simplest linear regression model fitted to the data in (any) learning block not only accurately represents the in-sample data behavior but also correctly grasps the dynamic governing the whole sample. Figure 4 depicts an exemplary synthetic sample following the linear trend which are used in the set of Monte Carlo experiments, parameters of which are outlined in Table 2.

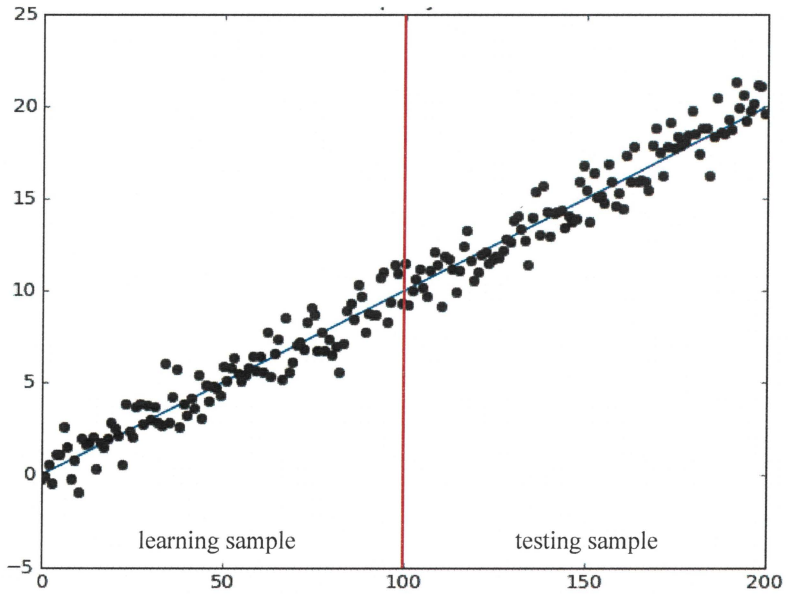
As one might have expected, the 1<sup>st</sup> order PL method is able to accurately approximate the true trend in the data, even with use of short learning blocks of 30 points – see Figure 5. However, ability to correctly estimate the true trend results in that for majority of stages of the learning procedure EOs have infinite (undefined) lengths (cf. Figure 6: infinite EO lengths do not appear on the plot, finite lengths occur sporadically). This is due to the fact that the exact grasp of the true trend in the whole sample given only information contained in the learning block is in this case equivalent to obtaining a precise model of the data generating process, which holds also beyond the learning block. As a consequence, we cannot falsify our understanding of the process based on the data from learning block with use of the testing block (i.e. part of learning sample which follows the learning block), and thus EO is infinite. Since most of the EOs in-sample are of infinite length we are also unable to formulate expectations about the limits to extrapolating our understanding of the process beyond the learning sample (i.e. the length of EO starting at the end of learning sample).

**Table 2. Experiments setup.**

<b>True trend formula</b>	$f(t) = 0.1 \times t$
<b>Length of the synthetic data sample</b>	200 points
<b>Length of the learning sample</b>	100 points
<b>Order of PL method</b>	1, 2
<b>Length of the learning blocks</b>	30, 40

<sup>34</sup> Results of Monte Carlo runs on data with higher level of noise are used to formulate general conclusions about the influence of the strength of noise on the PL method.

Strength of the noise <sup>35</sup>	0.05
Number of Monte Carlo runs for each parameter combination	40



**Figure 4.** Exemplary data (black dots) following a linear trend  $f(t) = 0.1 \times t$  (blue line). Standard deviation of noise  $\sigma = 0.05 \times (\max f - \min f)$ .

<sup>35</sup> Expressed as fraction of the range of the true trend (cf. Section 4.1)

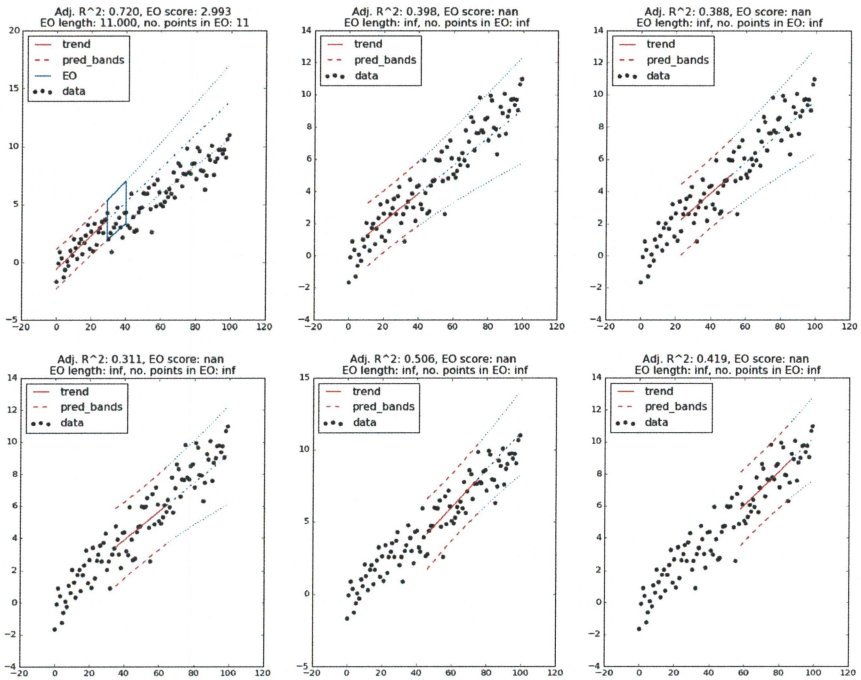


Figure 5. Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 30 points.

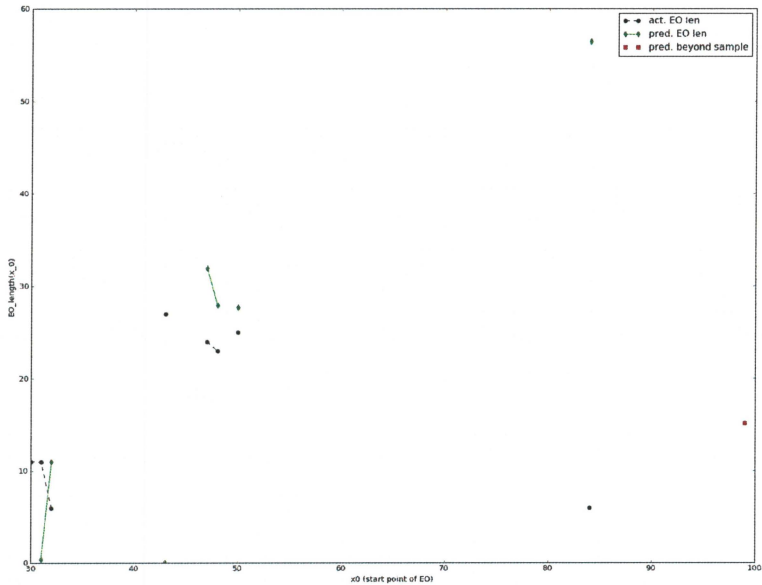


Figure 6. Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 30 points. Correlation between actual and predicted EO lengths is

-0.175. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.

In the case of noisy data following a linear trend the use of higher order PL methods (using trend functions more complex than the true linear trend) is not advisable. We demonstrate it on the example of 2<sup>nd</sup> order PL procedure. As one can see on Figure 7, prediction bands for the 2<sup>nd</sup> order polynomial regression diverge much faster than analogous prediction bands for linear regression. As a result, most often the EOs obtained in the process of 2<sup>nd</sup> order PL procedure have infinite lengths. Moreover, more flexible 2<sup>nd</sup> order polynomial model is more visibly susceptible to influence of noise in the data, and thus producing less certain and robust, often ill-directed projections. Therefore, any EO of finite length obtained with use of the 2<sup>nd</sup> order method is unreliable as it is most likely ill-directed and overly wide.

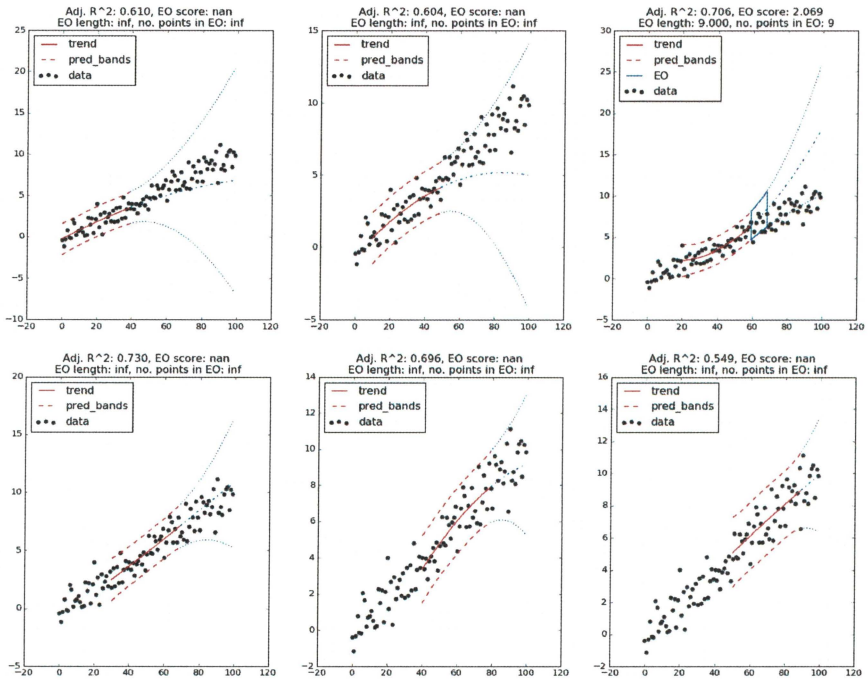


Figure 7. Six exemplary stages of the 2<sup>nd</sup> order PL procedure with learning block length of 40 points.

#### 4.3.2. Data following a 4<sup>th</sup> order polynomial trend

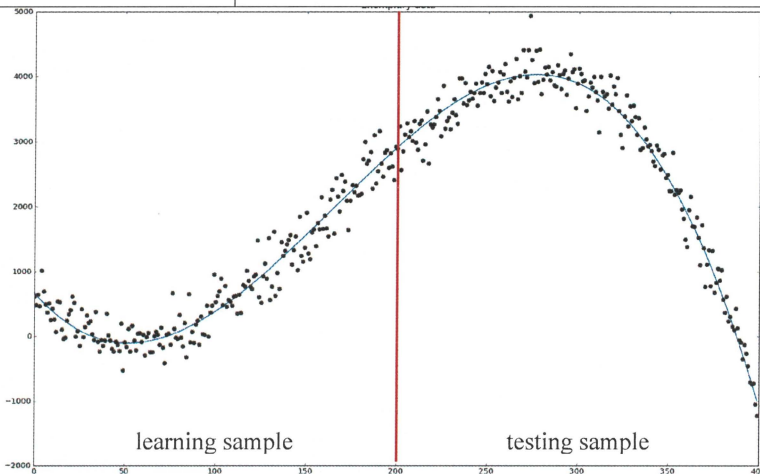
In the next set of experiments we analyze the performance of prognostic learning method applied to the noisy data following the trend of higher complexity. Method of polynomial regression is in principle able to provide an unbiased estimate of such trend.



In Table 3 we gather the parameters of these experiments. Figure 8 shows exemplary synthetic data sample used in these experiments.

**Table 3. Experiments setup. 4<sup>th</sup> order polynomial trend.**

True trend formula	$f(t) = (0.001 \times (t - 50))^4 - (0.09 \times (t - 50))^3 + (0.5 \times (t - 50))^2 - t - 50$
Length of the synthetic data sample	400 points
Length of the learning sample	200 points
Order of PL method	1, 2, 3, 4
Length of the learning blocks	20, 30, 40, 50, 60
Strength of the noise	0.01, 0.05, 0.1
Number of Monte Carlo runs for each parameter combination	40



**Figure 8.** Exemplary data (black dots) following 4<sup>th</sup> order polynomial trend (blue line) given by formula  $f(t) = (0.001 \times (t - 50))^4 - (0.09 \times (t - 50))^3 + (0.5 \times (t - 50))^2 - t - 50$ . Standard deviation of the noise  $\sigma = 0.05 \times (\max f - \min f)$ .

Table 4 presents the results obtained for the synthetic data with low level of noise<sup>36</sup> (i.e. 0.01 of width of the trend function range). For each order of the PL method the optimal learning block length is used.

**Table 4. Choices of the LB lengths for different orders of the PL method yielding the best results of experiments on data following 4<sup>th</sup> order polynomial trend.**

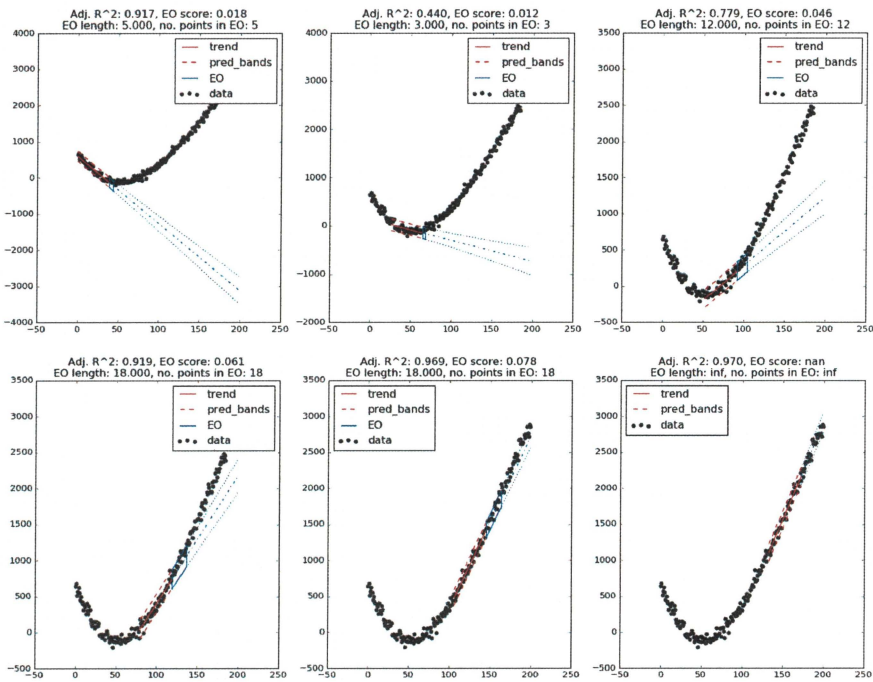
Method order	LB length	Noise level	Regression assumptions	EO Scores	EO lengths	Correlation : actual vs. predicted EO lengths (in sample)	Actual EO lengths (out-of-sample)	Predicted EO lengths (out-of-sample)	Correlation : actual vs. predicted EO lengths (out-of-sample)

<sup>36</sup> For stronger noises the performance of the PL method deteriorates, which to certain extent may be compensated by increasing the length of the learning block.

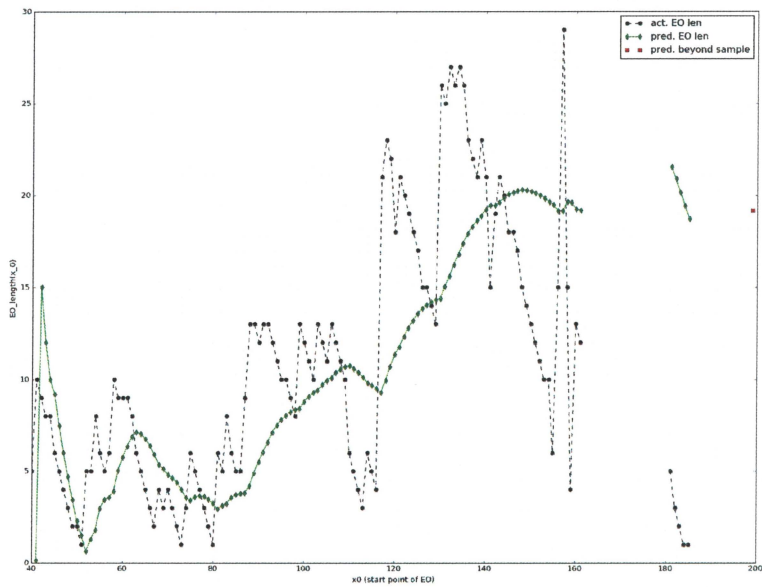
							)		sample)
1	40	0.01	Ok	0.01 -0.08	Slightly increasing Average: 15 (1 - 30)	0.54	Mode 25 [6 - 37]	Mode 18 [12 - 33]	0.2 (finite EO length in 40 out of 40 runs)
2	50	0.01	Ok	0.03 - 0.08	Oscillating decreasing 130 to 0	0.63	Flat Mode below 50 [0-180]	Left skew Mode 0 [0 - 40]	0.09 (finite EO length in 38 out of 40 runs)
3	40	0.01	acceptable (possible autocorrelation of residuals)	Up to 0.03, mostly undefined	Oscillating [2 - 10] few outliers up to 18	-0.05	[3 - 14]	[0 - 10]	0.09 (finite EO length in 7 out of 40 runs)

4	50	0.01	Ok	Up to 0.02, mostly undefined	Oscillating [1 - 15] outlier at 48	-0.09	[1-19], mostly below 6	[0 - 19], mostly below 5	-0.39 (finite EO length in 10 out of 40 runs)
---	----	------	----	------------------------------	------------------------------------	-------	------------------------	--------------------------	--

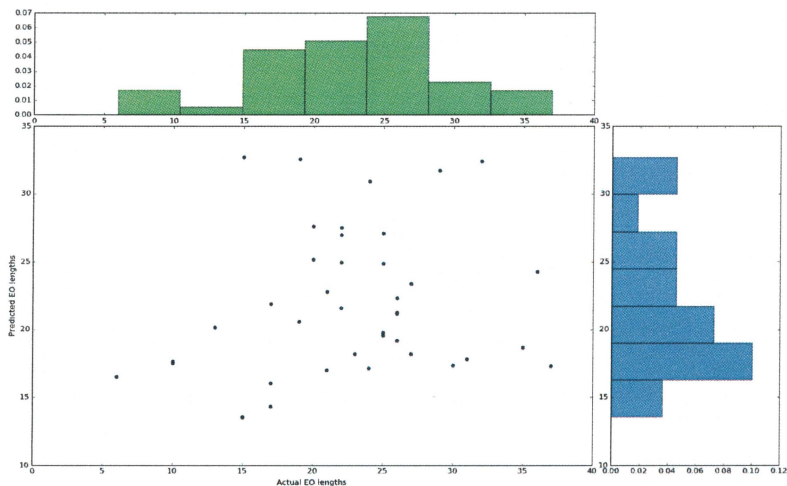
Surprisingly, the best performance is achieved for the variant of prognostic learning method which employ 1<sup>st</sup> order regression over short learning blocks (just 40 points). Figure 9 illustrates six exemplary stages of such prognostic learning procedure. This optimal combination of the order of method and the length of LB yields relatively stable behavior of the EO lengths with not too strong oscillations around slightly increasing trend (cf. Figure 10). The ranges of the actual and predicted lengths of the EO starting at the end of learning sample are in good agreement, although the correlation between these lengths is weak (see Figure 11). Notice also that all EO lengths are not longer than the learning block.



**Figure 9.** Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 40 points. In regions where the curvature of the true trend is significant linear model does not fit well to the data in the learning block and the actual lengths of the EO are low. In regions, where the true trend has approximately constant slope the PL method performs well.

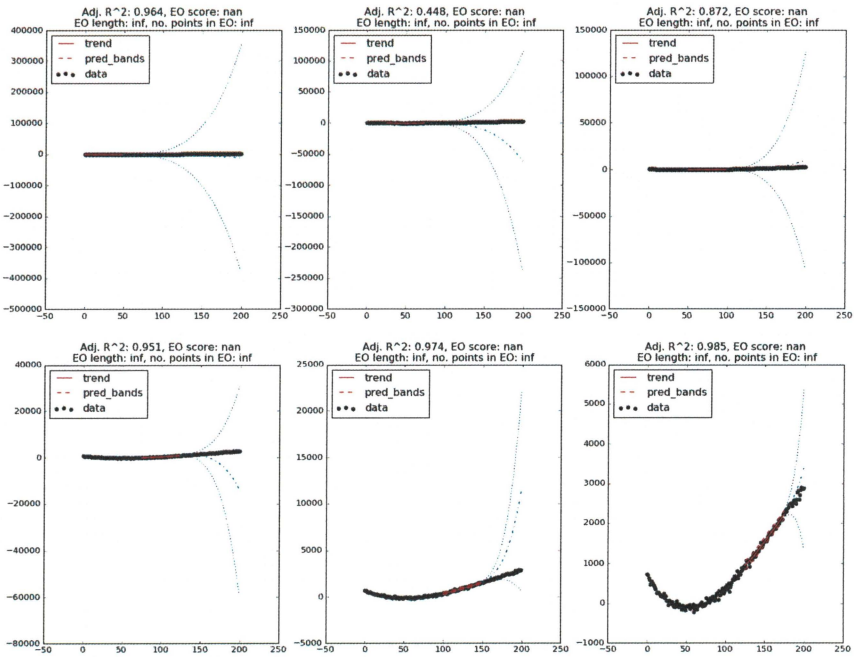


**Figure 10.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 40 points. Correlation between actual and predicted EO lengths is 0.537. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.

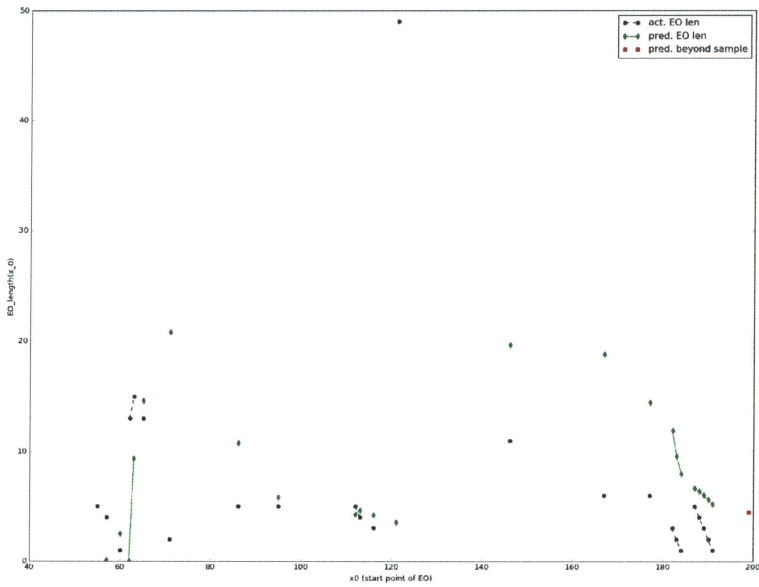


**Figure 11.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 40 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 40. Histograms approximate marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is 0.195.

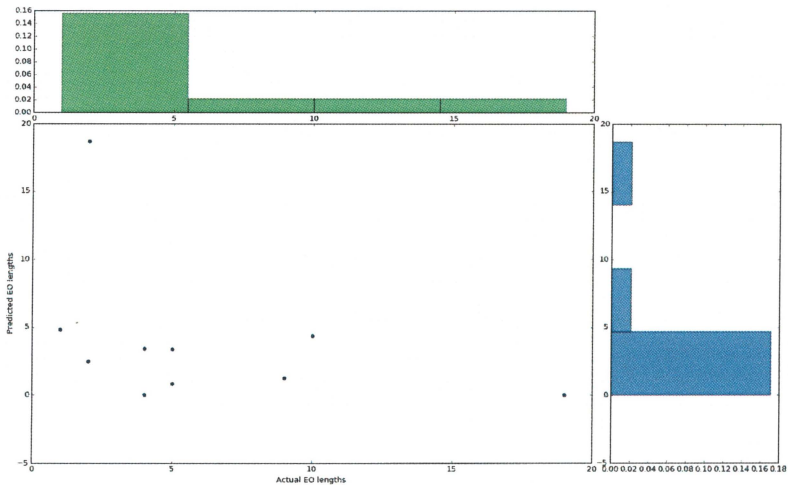
Equally surprising is a relatively poor performance of the 4<sup>th</sup> order PL method. Polynomial trends of degree 4 fitted to learnings block of length 50 grasp the behavior of the data better than linear trends. However, predictions of future behavior of the data by extrapolating 4<sup>th</sup> order polynomial regression functions are highly uncertain. This is caused by their high flexibility, which within the learning block is forced to minimize distance from the data points, but beyond it, when it is unconstrained, it may strongly deviate from the actual trend. This high uncertainty is represented by quickly diverging prediction bands. As a result, for most of the stages of the PL procedure we cannot determine the length of the EO because extremely wide prediction bands cover all points in the testing block (cf. Figures 12 and 13). This phenomenon has also a strong impact on both predicted and actual lengths of the EO starting at the end of the learning sample. Although ranges of the actual and predicted lengths are in very good agreement, only few cases in which these lengths are finite undermine meaningfulness of Monte Carlo experiments (cf. Figure 14).



**Figure 12.** Six exemplary stages of the 4<sup>th</sup> order PL procedure with learning block length of 50 points. Notice that often extrapolated trend deviates substantially from the actual data in the testing sample. High uncertainty of these predictions is exhibited by quickly diverging prediction bands.



**Figure 13.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 4<sup>th</sup> order PL procedure with learning block length of 50 points. Correlation between actual and predicted EO lengths is  $-0.086$ . The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.



**Figure 14.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 10 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 40. Histograms approximates

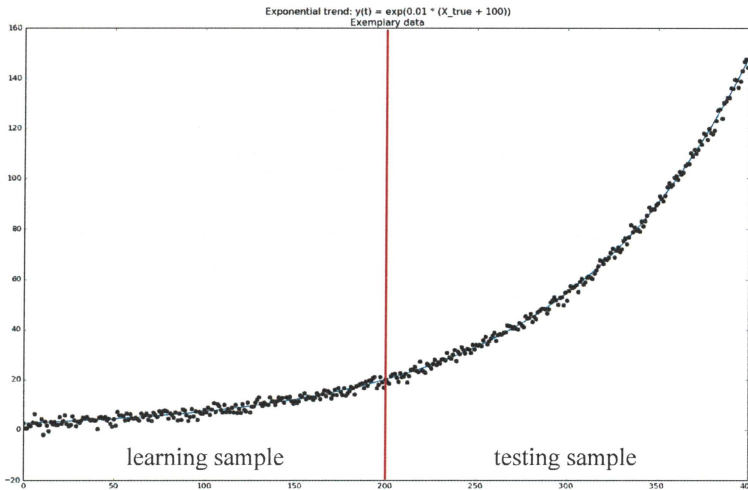
marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is -0.387.

### 4.3.3. Data following exponential trend

In this set of experiments we analyze the performance of the prognostic learning method applied to the noisy data following a commonly occurring type of trend not belonging to the family of polynomials. Although it is not possible to model the data following exponential trend with any polynomial in the long run, it is possible to achieve the satisfactory local approximation with use of polynomial function of sufficiently high order. Hence, the PL method grasping the local<sup>37</sup> behavior of the data with polynomial regression model is expected to be applicable also in this case. In Table 5 we gather the parameters of Monte Carlo experiments on synthetic exponential data. Figure 15 shows exemplary synthetic data sample used in these experiments.

**Table 5. Experiments setup. Exponential trend.**

<b>True trend formula</b>	$f(t) = \exp(0.01 \times (t + 100))$
<b>Length of the synthetic data sample</b>	400 points
<b>Length of the learning sample</b>	200 points
<b>Order of PL method</b>	1, 2, 3,
<b>Length of the learning blocks</b>	20, 30, 40, 50
<b>Strength of the noise<sup>38</sup></b>	0.001, 0.005, 0.01
<b>Number of Monte Carlo runs for each parameter combination</b>	50



<sup>37</sup> i.e. only within relatively short learning block

<sup>38</sup> Expressed as the fraction of trend function range width – cf. Section 4.1.

**Figure 15.** Exemplary data (black dots) following exponential trend (blue line) given by formula  $f(t) = \exp(0.01 \times (t + 100))$ . Standard deviation of noise  $\sigma = 0.01 \times (\max f - \min f)$ .

Table 6 gathers the results obtained for the synthetic data with low level of noise<sup>39</sup> (i.e., 0.001 of width of the trend function range). For each order of the PL method the optimal learning block length is used.

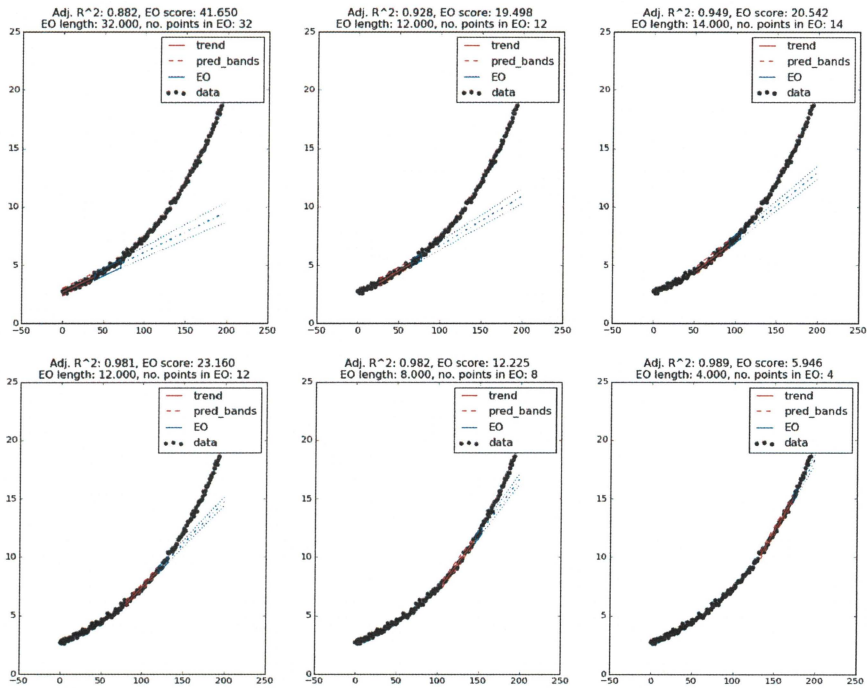
**Table 6. Choices of the LB lengths for different orders of the PL method yielding the best results of experiments for synthetic data following exponential trend.**

Method order	LB length	Noise level	Regression assumptions	EO Scores	EO lengths	Correlation: actual vs. predicted EO lengths (in sample)	Actual EO lengths (out-of-sample)	Predicted EO lengths (out-of-sample)	Correlation: actual vs. predicted EO lengths (out-of-sample)
1	40	0.001	Ok (possible autocorrelation of residuals)	Oscillating, gradually decreasing [42 to 2]	Oscillating, decreasing [30 to 1]	0.75	Flat [1 – 10]	Flat [0 – 5]	-0.03 (finite EO length in 50 out of 50 runs)
2	40	0.001	Ok	Oscillating below 20, mostly undefined	Oscillating, slight decrease [35 to 1], few outliers up to 80	0.34	Flat [0 – 200]	Left skew [0 – 30] Mode 0	0.27 (finite EO length in 50 out of 50 runs)
3	50	0.001	Ok	Oscillating below 11.2, mostly undefined	Decreasing [20 to 3] Outliers up to 75	0.02	Left skew [4 – 80] Majority below 20	[0 – 8]	0.26 (finite EO length in 10 out of 50 runs)

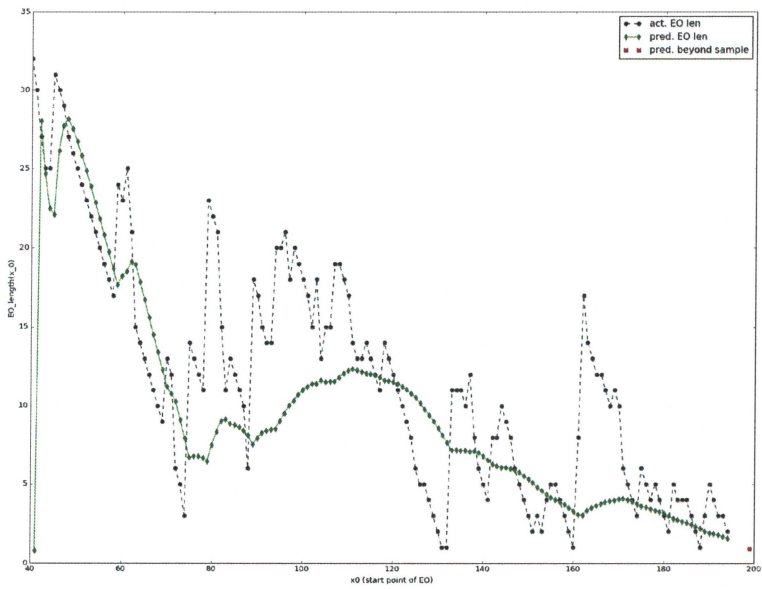
The best performance is achieved for the 1<sup>st</sup> order PL method using short learning blocks (of just 40 points). Six exemplary stages of such prognostic learning procedure are visualized on Figure 16. For the initial stages of PL procedure EOs are relatively long (due to initially small changes in the slope of exponential trend), but getting short in course of the procedure (as the increase of exponential trend accelerates) – cf. Figure 17. The ranges of the actual and predicted lengths of the EO starting at the end of learning sample are comparable (see Figure 18). The range of values of predicted EO lengths is narrower than the range of actual EO lengths, which means that expected EO length is likely to underestimate the actual EO length. However, they are virtually uncorrelated.

<sup>39</sup> For stronger noises the performance of the PL method deteriorates, which to certain extent may be compensated by increasing the length of the learning block.

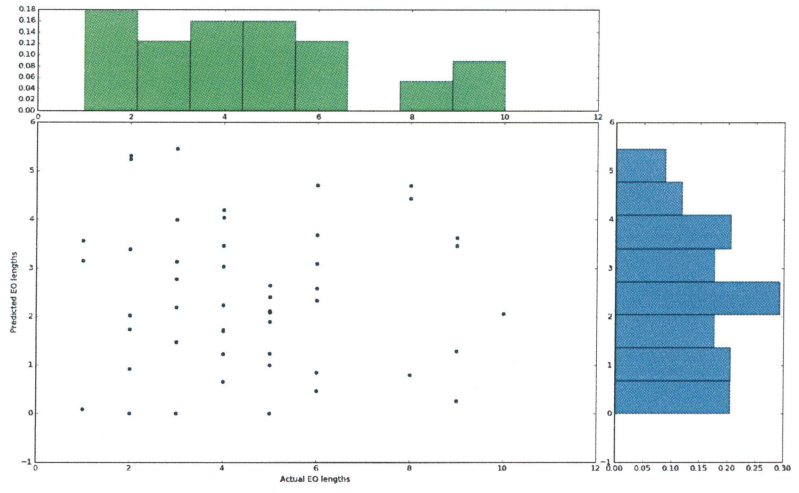




**Figure 16.** Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 40 points. For initial stages of the PL procedure lengths of the EO are comparable with the length of the learning block. This is due to slow initial increase of the exponential trend. As this increase begins to accelerate in further stages the EO lengths get shorter.



**Figure 17.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 40 points. Correlation between actual and predicted EO lengths is 0.746. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.



**Figure 18.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 50 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 50. Histograms approximates

marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is -0.032.

Higher order polynomials are much better in approximating the exponential trend, yet the performance of higher order PL methods is worse than for the one based on linear regression. We discuss it on the example of the 2<sup>nd</sup> order polynomial method. Fitted quadratic trends extrapolated beyond corresponding learning blocks always increase slower than true exponential trend (yet quicker than linear trends). However, prediction bands are usually wide enough to cover all the data points in the testing block. As a result, for most of the stages of the PL procedure we cannot determine the length of the EO (cf. Figures 19 and 20). Distribution of the predicted lengths of EO starting at the end of learning sample is strongly skewed to the left and has much narrower support than the relatively flat distribution of the actual EO lengths at the end of the learning sample (see Figure 21). Thus predicted EO length is likely to heavily underestimate the actual length, while correlation between them is weak.

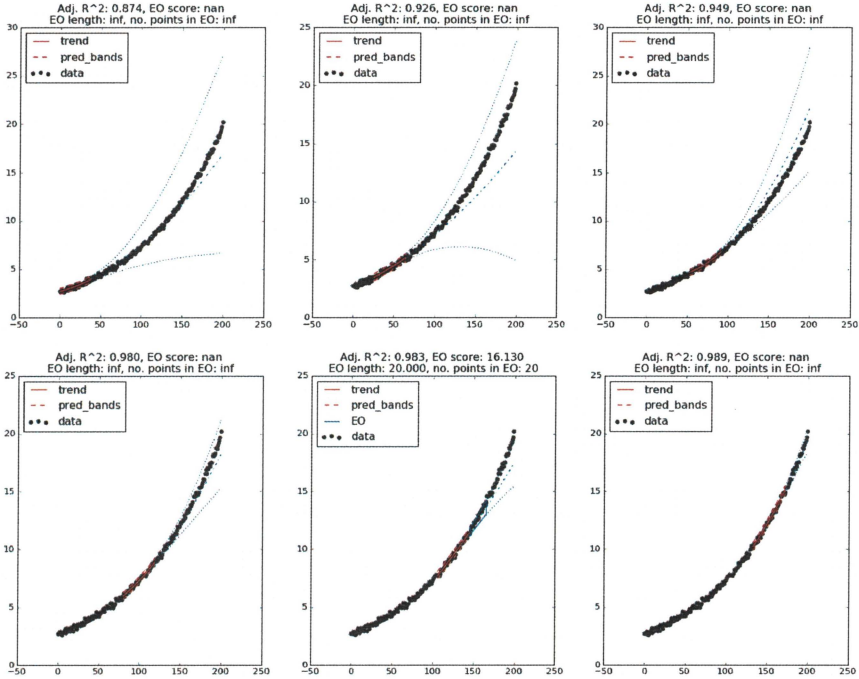
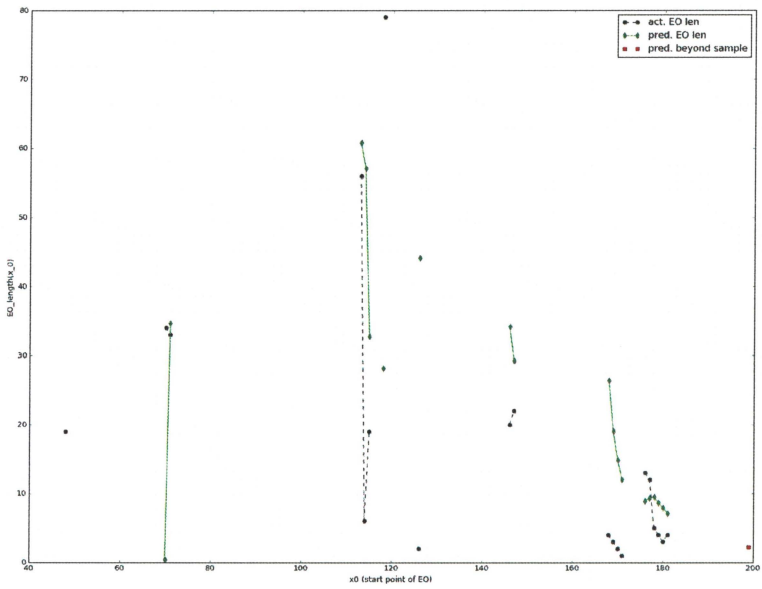
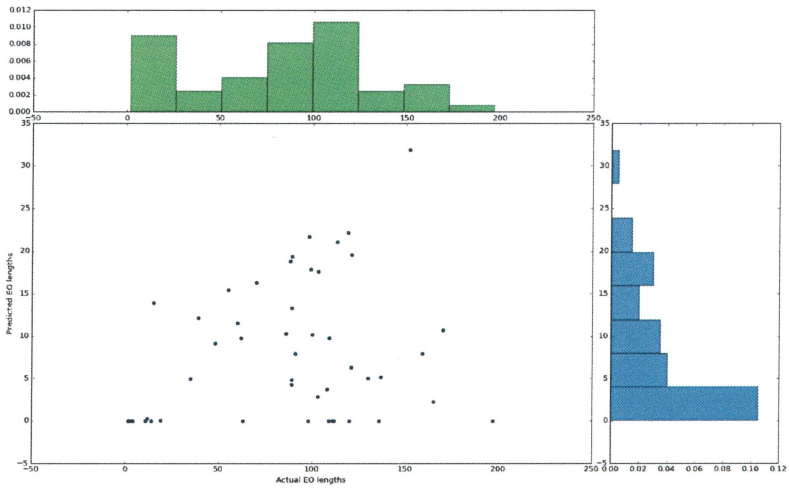


Figure 19. Six exemplary stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points.



**Figure 20.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points. Correlation between actual and predicted EO lengths is 0.335. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that majority of the EO lengths (both actual and predicted) are no longer than the length of the learning block.



**Figure 21.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 50 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 50. Histograms approximates

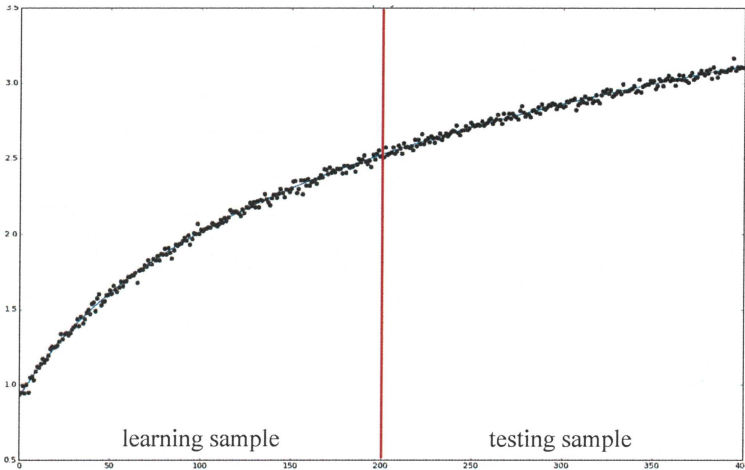
marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is 0.286.

#### 4.3.4. Data following logarithmic trend

Now we examine the performance of the prognostic learning method on the synthetic data following an increasing but decelerating trend - exemplified by logarithmic trend. Such trend, often encountered in real life data, cannot be approximated well by any polynomial in the long run, however, satisfactory local (i.e. for relatively short subsample) agreement may be achieved. This is the rationale for applying PL method for such type of data. In Table 7 we gather the parameters of Monte Carlo experiments on synthetic logarithmic data. Figure 22 shows exemplary synthetic data sample used in these experiments.

**Table 7. Experiments setup. Logarithmic trend.**

<b>True trend formula</b>	$f(t) = \log(0.05 \times (t + 50))$
<b>Length of the synthetic data sample</b>	400 points
<b>Length of the learning sample</b>	200 points
<b>Order of PL method</b>	1, 2, 3,
<b>Length of the learning blocks</b>	20, 30, 40, 50
<b>Strength of the noise<sup>40</sup></b>	0.01, 0.025, 0.05
<b>Number of Monte Carlo runs for each parameter combination</b>	50



<sup>40</sup> Expressed as the fraction of trend function range width – cf. Section 4.1.

**Figure 22.** Exemplary data (black dots) following logarithmic trend (blue line) given by formula  $f(t) = \log(0.05 \times (t + 50))$ . Standard deviation of noise  $\sigma = 0.01 \times (\max f - \min f)$ .

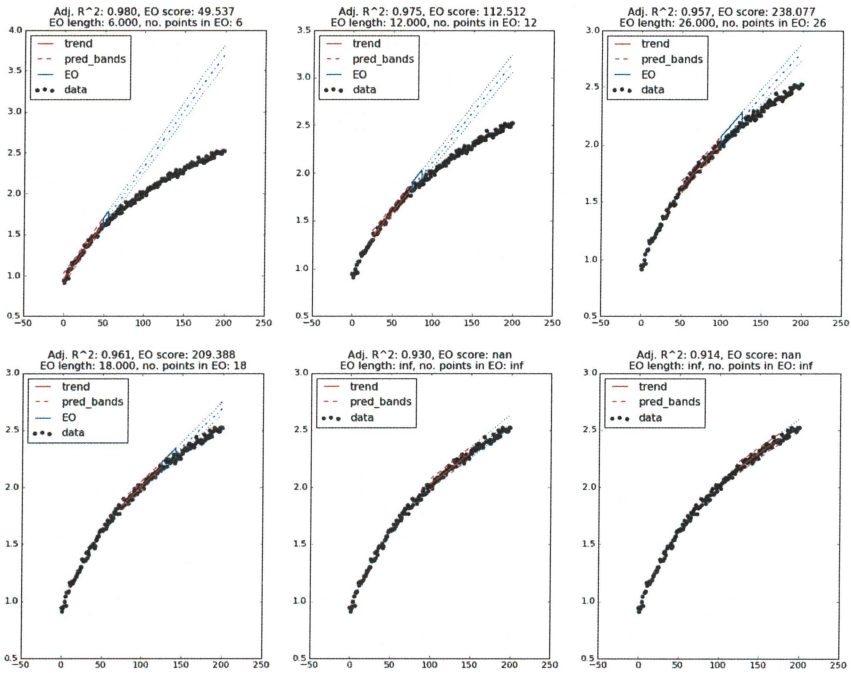
Table 8 summarizes the results obtained for the synthetic data with low level of noise<sup>41</sup> (i.e. 0.01 of width of the trend function range). For each order of the PL method the optimal learning block length is used.

**Table 8. Choices of the LB lengths for different orders of the PL method yielding the best results of experiments on synthetic data following a logarithmic trend.**

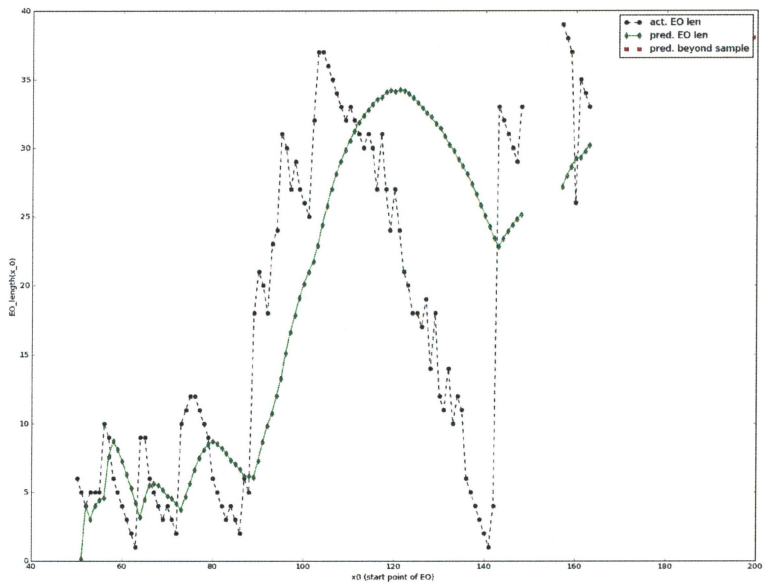
Method order	LB length	Noise level	Regression assumptions	EO Scores	EO lengths	Correlation : actual vs. predicted EO lengths (in sample)	Actual EO lengths (out-of-sample)	Predicted EO lengths (out-of-sample)	Correlation : actual vs. predicted EO lengths (out-of-sample)
1	50	0.01	Ok	Oscillating, below 405, often undefined	Oscillating, max increasing to 40	0.62	[0 – 110] Mode 40	[15 – 50] Mode 30	0.14 (finite EO length in 50 out of 50 runs)
2	50	0.01	Ok	Oscillating [20 – 160], mostly undefined	Oscillating, decreasing [120 to 1]	0.63	[3 – 26]	Left skew [0 – 23] Mode 0	0.66 (finite EO length in 7 out of 50 runs)
3	50	0.01	Ok (occasionally autocorrelation of residuals)	Oscillating [10 – 67], mostly undefined	Oscillating below 15, diminishing outliers (max 30)	0.5	[3 – 26]	[1 – 11]	-0.26 (finite EO length in 7 out of 50 runs)

As in previous sets of experiments, the best performance is achieved for the 1<sup>st</sup> order PL method – this time using slightly longer learning blocks of 50 points. Six exemplary stages of such prognostic learning procedure are shown on Figure 23. EOs calculated for the initial stages of PL procedure are short due to quickly decelerating trend at the beginning of learning sample. Slower rate of decrease of slope of logarithmic trend in the further part of the learning sample results in longer EOs for later stages of the procedure (cf. Figure 24). Notice also that the range of all (finite) lengths of EOs in-sample is narrower than the learning block. This is also the case for predicted lengths of the EO starting at the end of learning sample (see Figure 25). However, actual lengths of EO starting at the end of learning sample are significantly longer, while the correlation between actual and predicted lengths is weak.

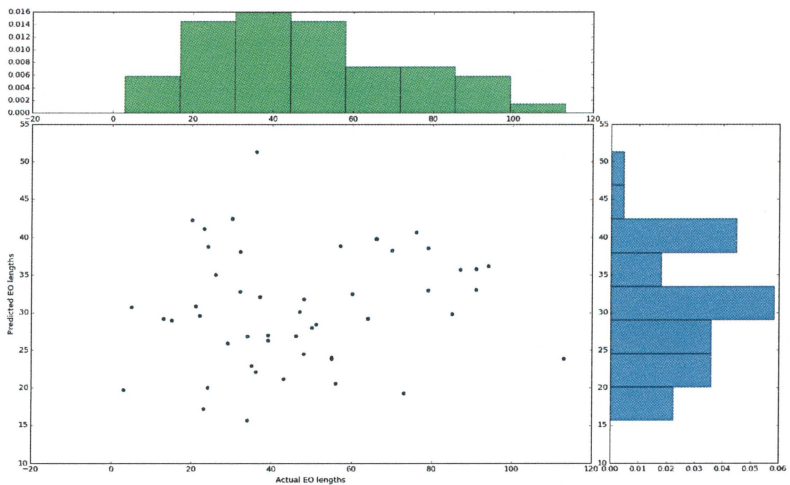
<sup>41</sup> For stronger noises the performance of the PL method deteriorates, which to certain extent may be compensated by increasing the length of the learning block.



**Figure 23.** Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 50 points. For initial stages of the PL procedure lengths of the EO are short due to initially quick decrease of slope of logarithmic trend. As this decrease begin to decelerate in further stages the EOs are getting longer.



**Figure 24.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 50 points. Correlation between actual and predicted EO lengths is 0.619. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.



**Figure 25.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 50 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 50. Histograms approximate



marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is 0.144.

PL methods based on higher order polynomial regressions perform worse than 1<sup>st</sup> order method when applied to the data following the logarithmic trend (or one of similar shape). We discuss it on the example of 2<sup>nd</sup> order polynomial method. The deviations from testing data of fitted quadratic trends extrapolated beyond the corresponding learning blocks increase faster than the analogous deviations of extrapolated linear trends. In addition to this often strong miss-direction of extrapolated higher order trends, their prediction bands diverge much faster than for prediction bands of linear models – see Figure 26. As a result, for the majority of stages of PL procedure the EOs have infinite (undefined) length (cf. Figure 27). Also the actual length of the EO starting at the end of the learning sample is infinite for the most of the Monte Carlo runs – making any analysis of joint behavior of predicted and actual lengths of the EO out-of-sample virtually impossible (cf. Figure 28).

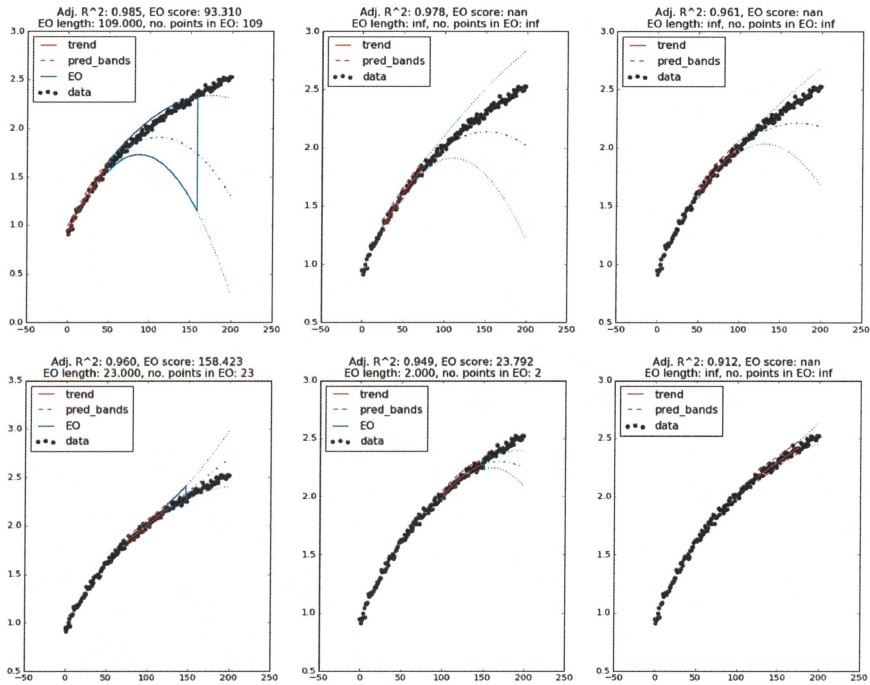
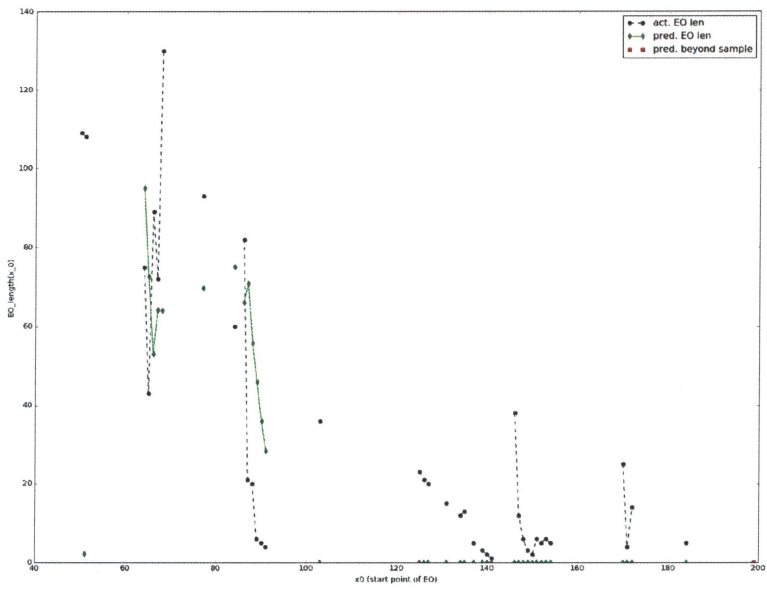
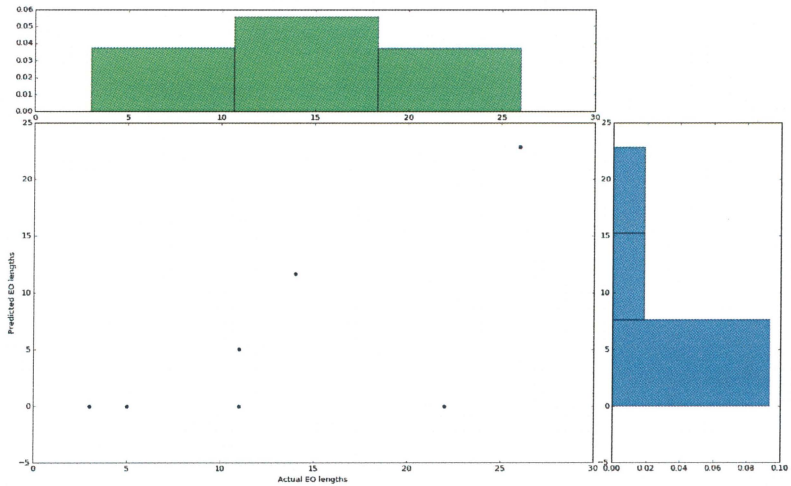


Figure 26. Six exemplary stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points.



**Figure 27.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points. Correlation between actual and predicted EO lengths is 0.628. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots).



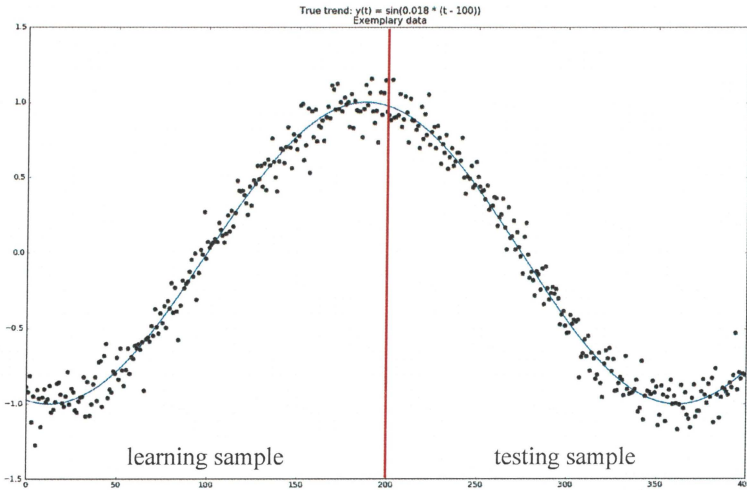
**Figure 28.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 7 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 50. Histograms approximates marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is 0.664.

### 4.3.5. Data following periodic trend

In the last set of experiments we investigate the usefulness of the prognostic learning method for analysis of the data following a sinusoidal trend with period comparable to the length of learning sample. Within short time intervals (i.e. comparable in length to the learning block) such data may appear to follow a clear non-periodic trend, which may be locally approximated by a polynomial. By applying PL method based on polynomial regression we want to understand the limits of such local approximations. Table 9 outlines the setup of Monte Carlo experiments on synthetic data following a periodic trend. Figure 29 exhibits exemplary synthetic data sample used in these experiments.

**Table 9. Experiments setup. Exponential trend.**

True trend formula	$f(t) = \sin(0.018 \times (t - 100))$
Length of the synthetic data sample	400 points
Length of the learning sample	200 points
Order of PL method	1, 2, 3,
Length of the learning blocks	20, 30, 40, 50
Strength of the noise <sup>42</sup>	0.01, 0.05, 0.1
Number of Monte Carlo runs for each parameter combination	50



<sup>42</sup> Expressed as the fraction of trend function range width – cf. Section 4.1.

**Figure 29.** Exemplary data (black dots) following sinusoidal trend with long period (blue line) given by formula  $f(t) = \sin(0.018 \times (t - 100))$ . Standard deviation of noise  $\sigma = 0.01 \times (\max f - \min f)$ .

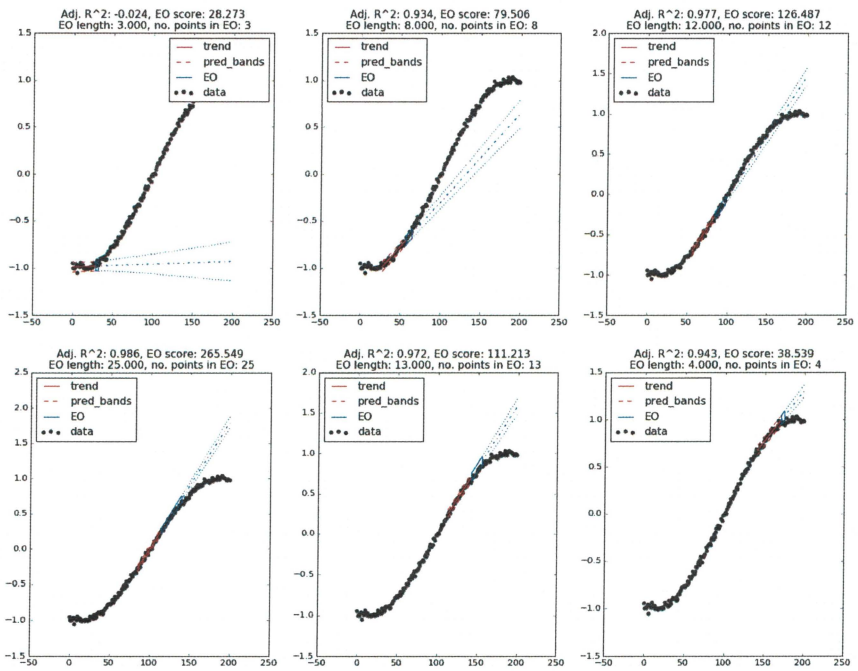
Table 10 summarizes the results of experiments performed with use of the synthetic data with low level of noise<sup>43</sup> (i.e. 0.01 of width of the trend function range). For each order of the PL method the optimal learning block length is used.

**Table 10. Results of experiments for optimal choices of LB lengths in case of synthetic data following a periodic trend.**

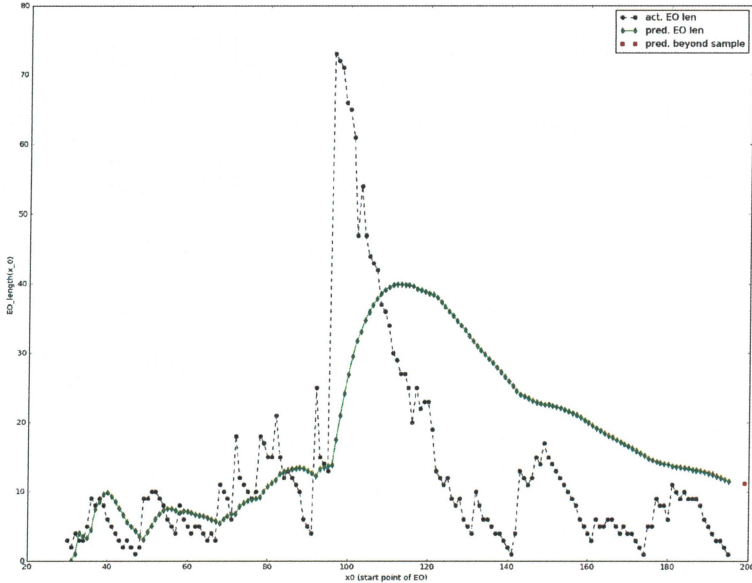
Method order	LB length	Noise level	Regression assumptions	EO Scores	EO lengths	Correlation : actual vs. predicted EO lengths (in sample)	Actual EO lengths (out-of-sample)	Predicted EO lengths (out-of-sample)	Correlation : actual vs. predicted EO lengths (out-of-sample)
1	30	0.01	Ok	Slowly oscillating, increasing to 395, then gradually decreasing to 10	Oscillating, increasing [1 - 70] then decreasing to 1. Most of the time below 20	0.43	Flat [1 - 11]	[7 - 14] Mode 11	-0.04 (finite EO length in 50 out of 50 runs)
2	50	0.01	Ok	Oscillating below 200, slightly increasing	Oscillating below 40, slightly decreasing, outliers up to 60	0.3	[0 - 150] Mode 100	[0 - 24]	0.14 (finite EO length in 50 out of 50 runs)
3	50	0.01	Ok (occasionally autocorrelation of residuals)	Oscillating [10 - 68], mostly undefined	Oscillating below 20, gradually decreasing outliers up to 40	0.53	[3 - 25]	[1 - 11]	-0.1 (finite EO length in 8 out of 50 runs)

As for the previous sets of experiments, the best performance is achieved for the 1<sup>st</sup> order PL method using short learning blocks (of just 30 points). Figure 30 shows six exemplary stages of such prognostic learning procedure. For stages of the PL method whose learning blocks are close to the bending points of the true trend the EO lengths are relatively short with respect to the length of the learning block. However, EOs are much longer when corresponding learning blocks coincide with regions in which the true trend is nearly linear – see Figure 31. The predicted EO lengths out-of-sample may be slightly over-optimistic – the range of estimated lengths is shifted to the right with respect to the range of actual lengths of the EO starting at the end of the learning sample (cf. Figure 32). Moreover, predicted and actual EO lengths are virtually uncorrelated. Note however, that they are shorter than the length of learning blocks used in the PL procedure.

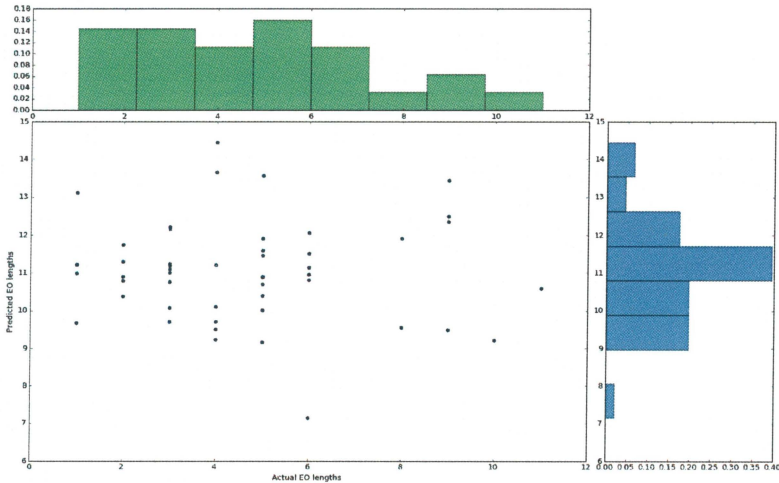
<sup>43</sup> For stronger noises the performance of the PL method deteriorates, which to certain extent may be compensated by increasing the length of the learning block.



**Figure 30.** Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 30 points. EOs are relatively short in cases when corresponding learning blocks are close to the bending points of the true trend and long otherwise.



**Figure 31.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 30 points. Correlation between actual and predicted EO lengths is 0.434. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots).



**Figure 32.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 50 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 50. Histograms approximate marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is -0.037.

Higher order polynomials are better suited to grasp the local behavior of the data in the learning blocks than the linear functions, especially when the LB is in the vicinity of bending points of the true trend (see Figure 33). This results in longer (in comparison to the 1<sup>st</sup> order method) EOs for the stages of PL procedure when the learning block coincide with intervals in which curvature of the true trend is significant – cf. Figure 34. Nevertheless, the EO scores are worse than for the 1<sup>st</sup> order PL method. This is due the fact that the prediction bands (defining the shape - and thus score - of the EO) for higher order polynomial regressions diverge faster than for linear regression. Moreover, flexibility of higher order polynomial trends is not particularly advantageous when predicting the length of the EO starting at the end of the learning sample – the predicted EO lengths grossly underestimate the actual EO lengths while their correlation is weak (see Figure 35).

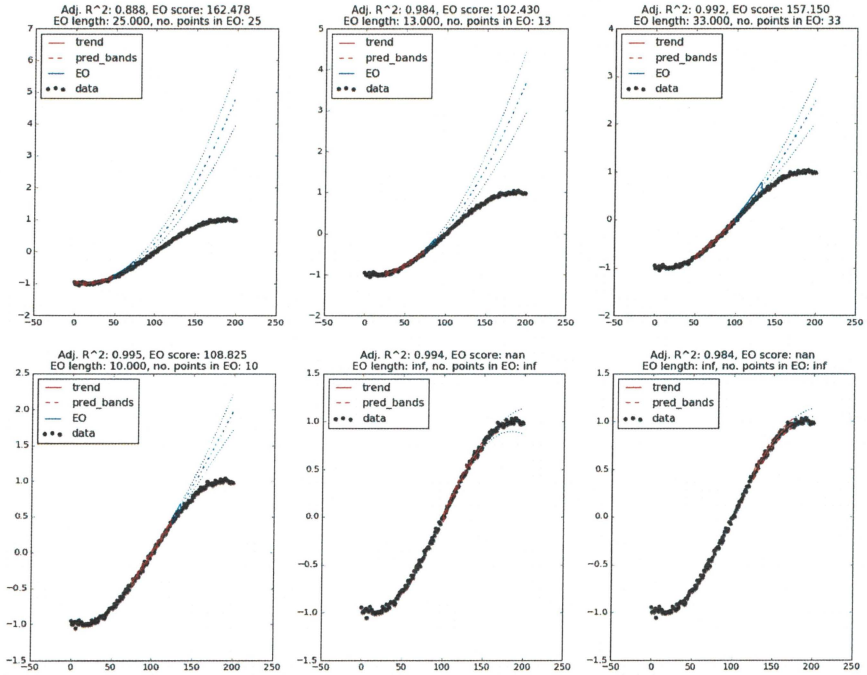
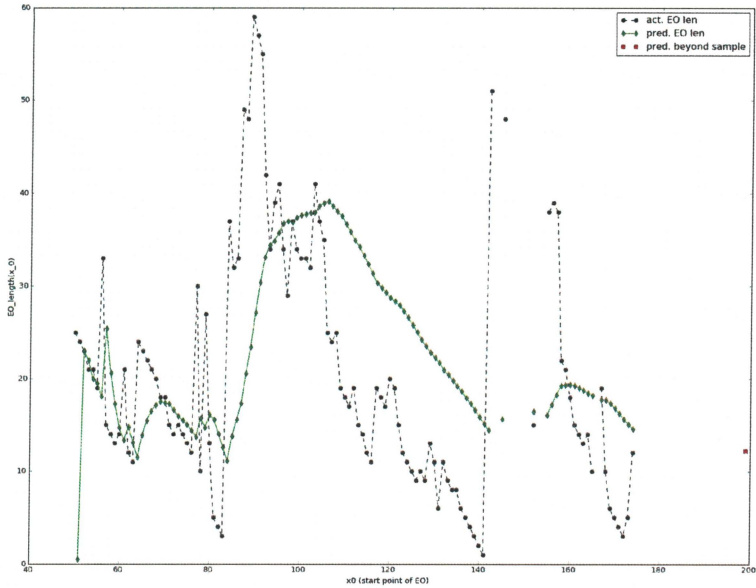
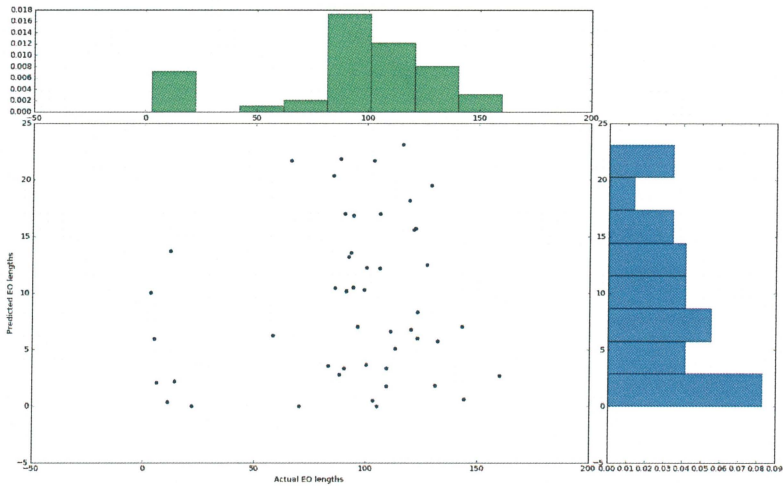


Figure 33. Six exemplary stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points.



**Figure 34.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points. Correlation between actual and predicted EO lengths is 0.309. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that majority of the EO lengths (both actual and predicted) are no longer than the length of the learning block.



**Figure 35.** Estimate of joint distribution of actual and predicted lengths of the EO starting at the end of the learning sample. Each of 50 points on the scatter plot represents the result of one Monte Carlo run resulting in finite actual EO length. Total number of Monte Carlo runs is 50. Histograms approximates marginal distributions of actual EO lengths (green) and predicted EO lengths (blue). Their correlation is 0.140.



## 4.4. Conclusions

In this section we present some general conclusions on the performance of prognostic learning method based on polynomial regression following from the results of experiments on synthetic data sets described in previous two sections.

We begin with analysis of the impact of complexity of the class of regression functions (i.e. order of polynomials) used in the PL method. This factor appears to be the most important for the performance of the prognostic learning. **With increasing complexity:**

- Fulfilment of regression method assumptions do not change significantly, however assumptions violations may be slightly more frequent.
- EO scores in principle decrease. This is due to the fact that the speed of divergence of the prediction bands - and thus width of EO – is of the same order as the polynomial trend used in the underlying regression model. In addition, the number of stages of PL procedure for which EO scores are undefined (i.e. cases for which EOs have infinite length) usually increase.
- Actual EO lengths (in-sample) - if finite - in general decrease. Clear tendencies, such as often observed decrease of the EO lengths for consecutive stages of the 1<sup>st</sup> order PL procedure, gradually change to oscillations around relatively stable level.
- Correlation between actual and predicted EO lengths (in-sample) is typically getting weaker. This correlation is relatively strong in presence of clear monotonic trend in lengths of consecutive EOs obtained in course of the learning procedure. This is most often the case for the 1<sup>st</sup> order method. As such tendencies in EO lengths change to oscillations typical for higher order methods, this correlation gets weaker.
- Actual out-of-sample EO lengths (which are determined with use of the additional testing sample back to back with the learning sample) typically decrease. This effect is especially clear for the upper limits (maximums) of the observed ranges of finite EO lengths. Moreover, for higher order methods EOs of infinite (undefined) lengths are predominant.
- Predicted out-of-sample EO lengths in principle decrease. Moreover, regardless the order of method, the range of predicted EO lengths usually lays within (or at least significantly overlaps with) the range of the actual EO lengths. Thus, at least on average, predicted EO lengths out of sample underestimate the actual ones. However, the correlation between actual and predicted EO lengths is typically weak, often negative and in principle not very reliable for higher order methods (due to EOs being predominantly infinite).

**Increasing level of noise** in the data has in principle a negative impact on the performance of prognostic learning. The most apparent effect is the deterioration of EO scores due to the fact that higher level of noise stipulates wider EOs.

Optimal length of the learning block is closely related to the order of the method used. It should not be too short or overly long (we discuss the choice of optimal LB length in the further part of this section). Therefore, it is difficult to discriminate marginal impact of increasing the length of the LB – what is too short for one method may be too long

for the other. The clearest effect one sees for the EO scores. They may slightly improve, as longer LB allows for better estimation of parameters of the regression function (lower variance of estimates of regression function parameters).

Based on the experiments on synthetic data described in the previous section, we formulate the following observations about the **1<sup>st</sup> order method of prognostic learning**:

- Any true trend and any data behavior can be locally approximated by a line. This local approximation is relatively robust to the level of noise. As a consequence, the ill-directed EOs (if appear) are the result of inability of the linear model fitted to the LB to follow quickly changing true trend rather than result of noisy conditions.
- Bias<sup>44</sup> – variance trade-off: 1<sup>st</sup> order method is biased – it looks only for linear trends in the data and cannot grasp well strongly non-linear trends. This bias may be negligible when the true trend is slowly varying, but can be significant in presence of curved true trend in the data. This bias however is balanced by the relatively low variance of predictions made with use of the linear regression model, i.e. slowly (at least slower than for higher order methods) diverging prediction bands determining the width (and thus the score) of the EO.
- This has two significant practical consequences:
  - If the true trend is linear then 1<sup>st</sup> order method is optimal (prognostic uncertainty is the lowest possible).
  - If the true trend is non-linear then predictions made by extrapolating the linear trend fitted to the LB will eventually be wrong, thus EO will **almost always have a finite length, usually not greater than the optimal length of the learning block**. Then the length of the EO informs us about **safe lower band for the time horizon within which treating the dynamics of the data as linear is a good approximation**.
- Optimal length of the LB (and thus of the learning sample) is the lowest for the 1<sup>st</sup> order PL method. This is important for the applicability of the PL methodology, since in practice the data scarcity is a common problem.

Conclusions for the **higher order PL methods** are slightly different:

- Bias – variance trade-off: any continuous true trend in the data over a specified interval may be well approximated with use of a polynomial of sufficiently high order. This ability of higher order polynomials to closely follow the data sample reduces the bias of the method. However, in noisy conditions uncertainty in estimates of the parameters of the polynomial regression model fitted to the data

---

<sup>44</sup> Here the term “bias” refers to the method. It means that  $E(\hat{f}(t)) \neq E(X_t)$  for some  $t$  within the range (period) of the sample, where  $\hat{f}$  denotes the estimate of the true trend. It is not a systematic (measurement) error of analysed data.

in a learning block almost always results in high variance of predictions beyond the range of the LB (represented by quickly diverging prediction bands).

- This has two significant practical consequences:
  - Quickly growing uncertainty of predictions made by extrapolation of the fitted polynomial trend beyond the range of the learning block makes usefulness of such predictions questionable.
  - More importantly, due to flexibility of higher order polynomial trend and quickly diverging prediction bands **in most of the cases** (stages of the PL procedure) **EO length is infinite**. Indeed, it is finite in cases when the extrapolated polynomial trend around which EO is constructed was so ill-directed that this was not offset by quickly diverging prediction bands. Thus, results of higher order PL methods should be treated somewhat differently and with more suspicion than the results of the 1<sup>st</sup> order method.
- Required length of the LB is considerably higher than for the 1<sup>st</sup> order method. Longer LB is needed to prevent overfitting – situation in which the fit of the flexible polynomial trend may be strongly impacted by the random noise. This further reduce the usefulness of higher order PL methods in analysis of the relatively short real-life data sets.

We conclude this chapter with formulation of a few **rules of thumb for applying the PL method**:

1. 1<sup>st</sup> order method should be preferred over the higher order methods.
2. The stronger the noise the longer the LB required and the more difficult it is to use the higher order methods.
3. The higher the order of method the longer the LB required. In any case there should be at least 10 points in the LB per each parameter of the regression model to be estimated.
4. Given the data and the order of the PL method one should follow the following guidelines when selecting the **optimal length of the LB**:
  - a. Choose the LB length for which EO score is the highest (or slightly longer).
  - b. Choose the LB length for which EO length exhibits stable behavior in course of the PL procedure (oscillating with few small outliers) or when trends in the behavior of the EO lengths is changing (e.g., from clear decrease of EO length in course of the PL method to oscillations around certain level or when increasing tendency appears in the oscillations).
  - c. Choose the LB length for which correlation between actual and predicted EO lengths in-sample is relatively strong and positive.

Ideally these criteria should be fulfilled simultaneously. Such choice of the optimal LB length usually coincides with a good behavior of the predicted length of the EO starting

at the end of the learning sample (i.e. good overlap of the ranges of actual and predicted EO lengths and relatively strong correlation between them).

## 5. Real-life case studies

In the present chapter we test the applicability of the prognostic learning method in determining the limits of our understanding of the dynamics of the real-life data (i.e. their explainable outreach). In finding the optimal parameters of the PL method we draw on the insights of the previous chapter.

As the case examples we choose two data sets reflecting the dynamics of two processes of fundamental importance for our understanding of the impact of humans on the climate, namely the anthropogenic CO<sub>2</sub> emissions and the increase of CO<sub>2</sub> concentration in the atmosphere. Knowledge about the dynamics of these process is also necessary to run integrated assessment models (such as IMAGE<sup>45</sup>). Hence estimation of temporal limits of our understanding of these dynamics may also shed some light on the time horizons beyond which projections of the abovementioned IAMs may be unreliable.

The data sets we use contain the annual global CO<sub>2</sub> emissions from technosphere<sup>46</sup> (i.e. resulting from fossil fuel burning and cement production) and the annual average concentration of the CO<sub>2</sub> in the atmosphere measured at the Mauna Loa station<sup>47</sup>. As the CO<sub>2</sub> concentrations are influenced by the anthropogenic CO<sub>2</sub> emissions the analyzed data sets cover the same period, namely years 1959 – 2011.

### 5.1. Global CO<sub>2</sub> emissions from technosphere

In case of the anthropogenic CO<sub>2</sub> emissions the best performance is achieved for the 1<sup>st</sup> order PL method with learning blocks of length of 25 points (which is roughly half the size of the learning sample). This is consistent with our observations following from the experiments on synthetic data – for them 1<sup>st</sup> order PL method was also the best choice. The optimal length of the learning block was chosen according to the guidelines provided at the end of the previous chapter. Exemplary stages of the optimal PL procedure are presented on Figure 36. As one can see, the data follow roughly linear trend<sup>48</sup>, although three segments of slightly different slopes are clearly visible. These segments are of comparable length as the learning blocks used in the learning procedure. Hence, two types of configurations of the learning block with respect to the abovementioned segments are possible – and each of these constellations has a negative impact on the length of the explainable outreach. If the learning block strongly overlaps with one of these segments, then the linear model describes the data in the learning data well. However, the explainable outreach representing the expected future behavior of emissions is then compared against the data in the testing block which follows a different regime (i.e. increase of different slope) than the data in the learning block. As a

---

<sup>45</sup> For brief synopsis of the IMAGE model see e.g., [http://unfccc.int/adaptation/naïrobi\\_work\\_programme/knowledge\\_resources\\_and\\_publications/items/539\\_6.php](http://unfccc.int/adaptation/naïrobi_work_programme/knowledge_resources_and_publications/items/539_6.php)

<sup>46</sup> Source: CDIAC [http://cdiac.ornl.gov/trends/emis/overview\\_2011.html](http://cdiac.ornl.gov/trends/emis/overview_2011.html)

<sup>47</sup> Source: NOAA <http://www.esrl.noaa.gov/gmd/ccgg/trends/full.html>

<sup>48</sup> In broader perspective the overall trend in CO<sub>2</sub> emissions over the last 200 years is approximately exponential, but steep growth over the last six decades alone is roughly linear.

consequence, the EO is relatively short. The other possibility is that moment of regime change lays well within the learning block. This renders the linear model less suitable to represent the data behavior within the learning block and thus in increase of autocorrelation of model residuals. Such strong violation of the PL method assumptions results in shorter EO. Analysis of both actual and predicted lengths of the EOs for different stages of 1<sup>st</sup> order the learning procedure – cf. Figure 36 – confirms these observations. It shows that in principle one should not expect the EO to be much longer than about five points<sup>49</sup>, while very short EOs for some of the stages of the learning procedure indicates that analyzed process occasionally undergo sudden regime changes.

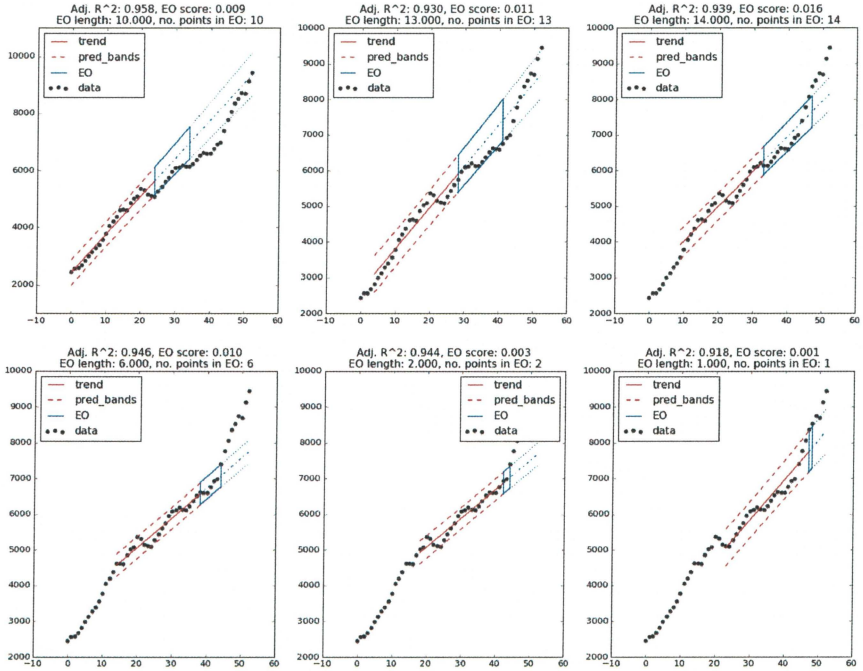
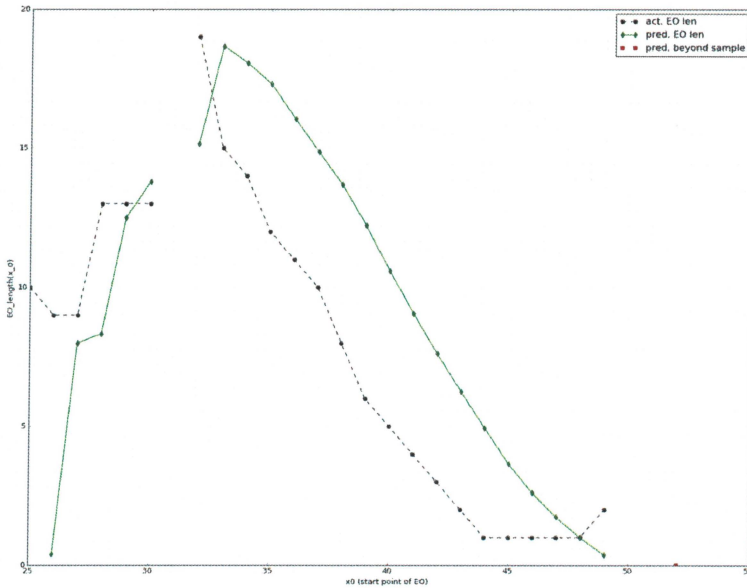


Figure 36. Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 25 points.

<sup>49</sup> Note that the EOs are shorter than the used learning blocks. This is in agreement with what we have observed for the synthetic data sets (c.f. Chapter 4).



**Figure 37.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 25 points. Correlation between actual and predicted EO lengths is 0.777. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.

Higher order PL procedures do not yield better results. As they require longer learning blocks, at each stage of the PL procedure a LB contains the moment of regime (slope of local trend) change. Although polynomial trends are more flexible than the linear trend, they too are unable grasp slight but sudden regime changes – as demonstrated on the example of the 2<sup>nd</sup> order PL method (cf. Figure 38). As a result, the EOs constructed with use of the 2<sup>nd</sup> order method are only wider (since prediction bands for higher order polynomial regression more rapidly than for linear case) but not longer – see Figure 39.

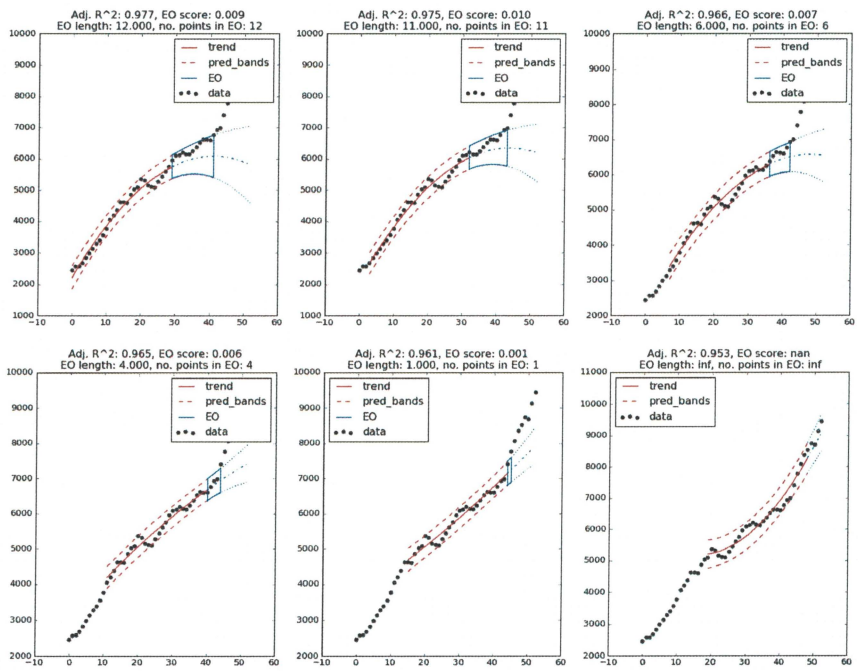
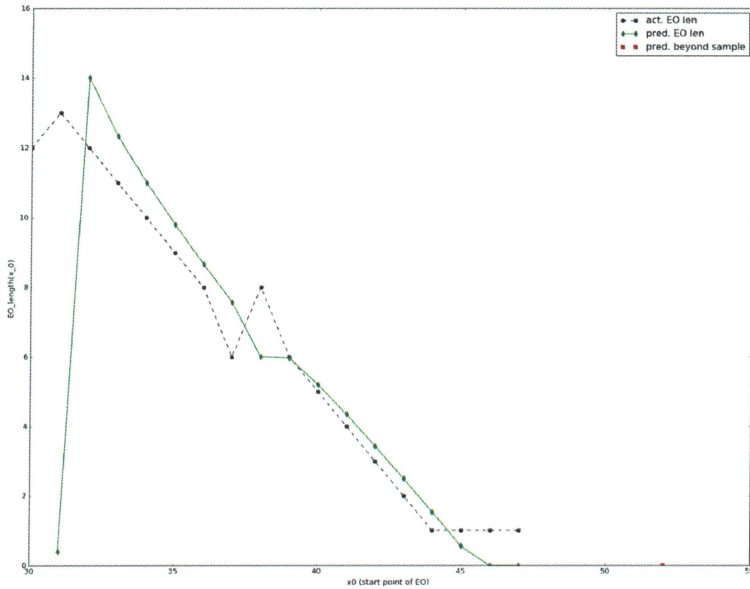


Figure 38. Six exemplary stages of the 2<sup>nd</sup> order PL procedure with learning block length of 30 points.



**Figure 39.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 2<sup>nd</sup> order PL procedure with learning block length of 25 points. Correlation between actual and predicted EO lengths is 0.713. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.

## 5.2. Concentration of CO<sub>2</sub> in the atmosphere

Time evolution of the CO<sub>2</sub> concentrations is smooth (in comparison to time evolution of anthropogenic CO<sub>2</sub> emissions) and follow a clear, exponential-like deterministic trend. The analyzed sample resembles the synthetic data with low level of noise following an exponential trend which we have analyzed in Chapter 4. Similarly to that case, the 1<sup>st</sup> order prognostic learning method proves to be the best choice among PL methods based on polynomial regressions. The optimal length of the learning block in this case is 20 points. As one can see on Figure 40 the EOs constructed with use of this method are narrow (due to low variance of the residuals for the linear models fitted to the learning blocks) but relatively short. Indeed, for most of the PL procedure stages the EOs are not longer than three points (cf. Figure 41). This is caused by the curvature of the trend the data clearly follows.



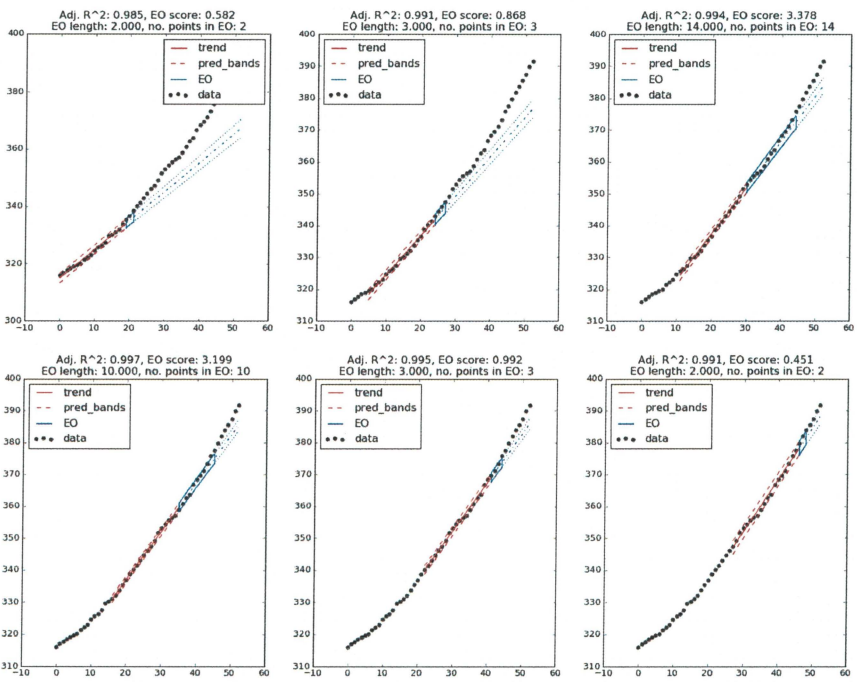
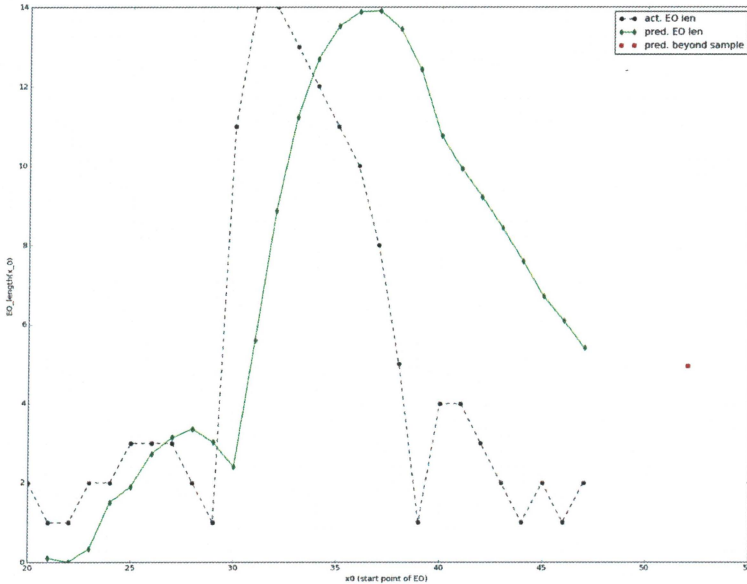


Figure 40. Six exemplary stages of the 1<sup>st</sup> order PL procedure with learning block length of 20 points.



**Figure 41.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 1<sup>st</sup> order PL procedure with learning block length of 20 points. Correlation between actual and predicted EO lengths is 0.461. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.

Quadratic trends are more suitable to approximate data following a curved trend (cf. Figure 42). However, in case of atmospheric CO<sub>2</sub> concentrations applying the 2<sup>nd</sup> order method does not result in longer EO. Indeed, although EOs constructed around quadratic trend have curved shape and are narrower than those for the 1<sup>st</sup> order method, they are still unable to follow the true trend in the long run (see Figure 43).

Applying the 3<sup>rd</sup> (or higher) order PL method to the data is not feasible as the minimal length of learning block for those methods is comparable to the size of the whole learning sample.

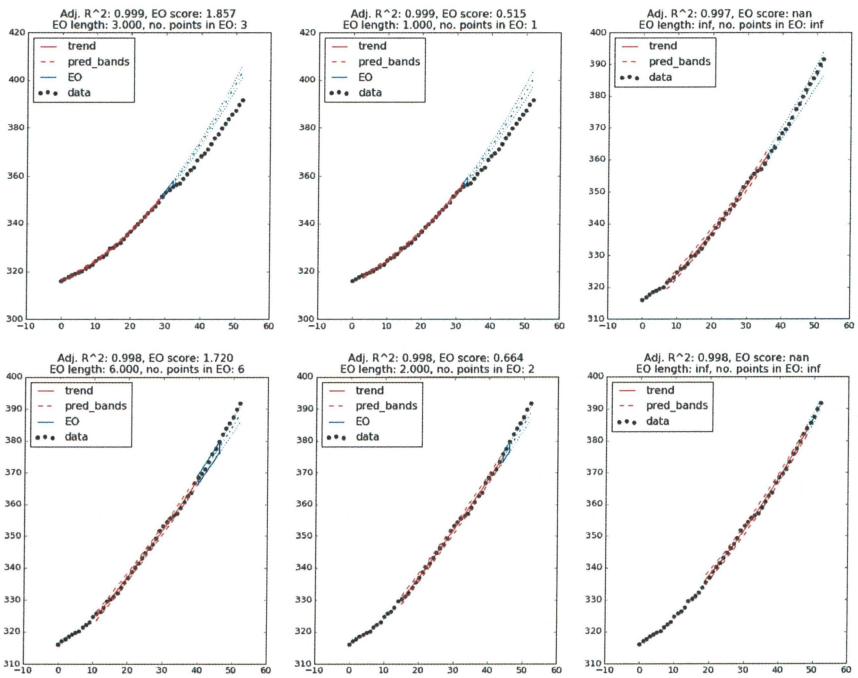
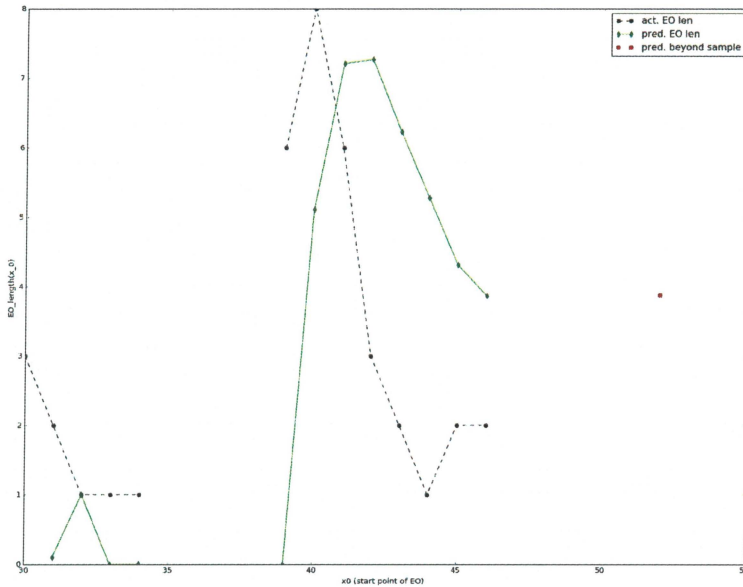


Figure 42. Six exemplary stages of the 2<sup>nd</sup> order PL procedure with learning block length of 30 points.



**Figure 43.** Actual (black dots) and predicted (green diamonds) EO lengths for all stages of the 2<sup>nd</sup> order PL procedure with learning block length of 50 points. Correlation between actual and predicted EO lengths is 0.302. The red square marks the predicted length of the EO starting at the end of testing sample. Prediction is based on all finite actual EO lengths calculated in the learning procedure (i.e. all of the black dots). Notice that all of the EO lengths (both actual and predicted) are no longer than the length of the learning block.

### 5.3. Conclusions

The temporal dynamics of both considered processes (i.e. anthropogenic CO<sub>2</sub> emissions and CO<sub>2</sub> concentrations in the atmosphere) are essentially nonlinear. The typical time horizons within which linear predictions of the behavior of upcoming data are credible is indicated by the lengths of the EOs obtained in result of applying the 1<sup>st</sup> order PL method. These limits for credible linear predictions are rather short.

For anthropogenic CO<sub>2</sub> emissions it is at most 15 points (years) but linear predictions for the immediate future are expected to be credible in much shorter time horizon. This is due to the fact that linear regression model employed in the learning procedure is not able to grasp or anticipate regime changes (i.e. sudden changes of slope).

More regular behavior of the atmospheric CO<sub>2</sub> concentrations result in slightly better, yet still short horizons for credibility of linear approximation of the process' dynamics – typical length of the EOs for the 1<sup>st</sup> order PL method is 2 to 6 points (years).

Approximations of the local dynamics of the considered processes by polynomial regression functions of higher orders are better in comparison to linear ones. However, predictions made by extrapolations of such trends are more uncertain and thus it is often impossible to assess their credibility by means of explainable outreach.

Finally, it is important to emphasize that limits of credibility assessed by means of the 1<sup>st</sup> order PL method should be treated just as the lower bound for period within which our understanding of the system's past may be used for making reliable predictions. In principle, some other method than polynomial regression may be more suitable to explain the data behavior. PL procedure based on such a method would most likely yield better (i.e. longer but still relatively narrow) EOs and thus improving the lower bounds for horizons of credibility which have established in this chapter.

## 6. Outlook

The research presented in this report is a feasibility study on the notions of prognostic learning and explainable outreach of the data. As such, it pursues the two objectives: (1) to frame the idea of the prognostic learning and place it in a broader context of earth system sciences; and (2) to develop and implement a prognostic learning procedure allowing to test the PL concept in practice.

While realizing the first objective we have restricted ourselves to analyze the data forming a time series and describing the temporal evolution of the analyzed system. Our main effort was directed to detecting the system's dynamics (i.e. the deterministic part of the analyzed time series) represented by the prevailing trend and to understanding the relation between the uncertainty of estimates of this trend and the credibility of our expectations about the future system's behavior based on projections this trend.

Understanding the temporal dynamics of the system and indicating the extent of credible predictions based on this understanding is just a first step in development of the paradigm of learning in a controlled prognostic context. However, the proposed PL method concentrates on grasping the temporal dynamics revealed by a single time series (using the time as the only explanatory variable) while hiding explicit dependence of the system on external forcing. For example, anthropogenic CO<sub>2</sub> emissions exhibit a roughly linear temporal dynamics over the last five decades (cf. Section 5) but they also strongly depend on the trends and disturbances of the global economy (such as energy crises in the 1970s, economic collapse of the soviet bloc in the 1990s or increased consumption in developing countries in recent years). We envisage a modification of the PL method by introducing additional explanatory variable(s) representing external forcing of the system (in context of anthropogenic CO<sub>2</sub> emissions this could be for example GDP) or dependence on some additional factors (e.g., carbon intensity of production processes). We speculate that explicit use of additional explanatory variables in the PL method will result in longer horizon of credible predictions (i.e. longer EOs).

Another challenge related to the objective (1) is to demonstrate the ability of the PL method to support a modelling exercise by realizing the "model performance assessment" track (cf. Figure 2) for a suitably selected climate or integrated assessment model.

Pursuing objective (2) we have proposed a way of implementing the prognostic learning concept which is based on the ordinary least squares (OLS) polynomial regression technique. This regression method was selected for its simplicity and relatively good performance. However, results presented in Sections 4.3 and 5 indicate the need for development of analogous versions of the PL method based on regressions using other parametric trends (e.g. exponential or power functions).

Moreover, we expect that the performance of the PL method based on the higher order polynomials may be improved by application of the regularization techniques (Hastie 2009, Murphy 2012). In principle, regularization penalizes the trend functions which are overly “wiggly”. It would allow to strike a balance between the flexibility of the high order polynomials and the robustness of predictions based on their extrapolations. We speculate that this would result in longer and not too much wider EOs than those obtained for the 1<sup>st</sup> order PL method.

Another way of improving the regression-based PL is to replace the OLS polynomial regressions with some more robust methods of fitting the trend, such as ridge regression or support vector regression (Hastie 2009, Murphy 2012) or nonparametric regressions (Wasserman 2006). Some preliminary results obtained with use of PL method based on selected nonparametric regression techniques are presented in the Appendix. This direction of future research is particularly interesting for the following reasons: (1) nonparametric methods do not confine us to any specific class of regression functions; (2) nonparametric offers a promising link between memory of the system (grasped by means of bandwidth parameter determining how many previous data points influences the present one) and the explainable outreach (defined as extrapolated prediction bands) and (3) Flexibility of the nonparametric regression curve results in longer (yet equally robust) EOs than the ones obtained with use of OLS linear regression.

Notice that the PL method presented in Chapter 3 relies heavily on assumption of independence of points in the learning sample<sup>50</sup>. However, by making such assumption (which we do deliberately for the sake of simplicity) we ignore the fact that the patterns of behavior of the stochastic part (such as autocorrelation structure of residuals) may also be of a significant importance. Simply assuming that the stochastic part is just an uncorrelated noise may result in underperformance of the EO<sup>51</sup>. In the future research we plan to address this problem by modifying the construction of the EO to account for the autocorrelation structure of the data.

Prognostic learning techniques discussed in this report grasp the dynamics of the system of interest by means of a regression function. Yet, trend functions are not the only way of expressing the patterns of the data behavior. Therefore, alternative<sup>52</sup> approaches to learning in a controlled prognostic are conceivable. For example, the techniques of granular computing such as quantization or clusterization (Pedrycz 2013) may be employed to grasp the patterns of data behavior. These techniques are based on assigning each of the data points to one member of a discrete collection of classes (called also information granules) in order to reduce the level of detail which may blur more fundamental features of the data (which are represented by these classes). The patterns in data behavior may then be expressed as transition rules from one information granule to the other, or more broadly by transition probabilities, i.e. likelihoods of

---

<sup>50</sup> It is required by both the OLS method of fitting a regression function to the data and by the way we determine the length of the EO (cf. Section 3.2).

<sup>51</sup> Recall that we decide to end the EO in the first moment for which layout of observations in period between the end of the learning block and this moment is unlikely under assumption that the extrapolated regression function fitted to the learning block is also a good approximate of the true trend in the testing block and the observations in the testing block are independent. However, if the observations were correlated then encountered layout of points might be not so unlikely and the actual EO length should be greater.

<sup>52</sup> i.e. alternative to the regression-based method presented in this report.

observation taken at certain time to belong to a certain information granule given the class into which the previous observation falls. This approach is currently being explored (Puchkova et al).

## 7. Summary

In this report we introduce the paradigm of learning in a controlled prognostic context. It is a data-driven exploratory approach to assessing the limits to credibility of any expectations about the future system's behavior which are based on a time series of historical observations of the analyzed system. The aim of the proposed method is to indicate the typical length of time intervals over which the trends in the historical data sample persist as well as the level of uncertainty in grasping these trends.

The key idea of the learning in a controlled prognostic context is to deduce directly from the data their explainable outreach, i.e. the spatio-temporal extent for which, in lieu of the knowledge contained in the historical observations, we may have a justified belief to contain future system's observations. The length of such explainable outreach indicates the time horizon within which predictions based on our current understanding of the system are credible. The initial width of the EO reflects the diagnostic uncertainty inherent to our system's perception, while the shape of the EO informs us about the strength of measures required to overcome the system's inertia.

We propose a method of constructing the explainable outreach based on the polynomial regression technique. The data sample is split into two parts: the learning block and the testing block. The dynamics of the system in the period covered by the learning block is grasped by means of a polynomial regression model and the explainable outreach expressing our expectations about the system's evolution beyond the learning block is constructed by extrapolating the prediction bands of the fitted regression model. These prediction bands represent both our expectations about the future system's dynamic and its uncertainty. The explainable outreach is then tested against the remainder of the data (i.e. testing block) in order to indicate the time horizon within which predictions based on the fitted regression model are believed to be credible.

We also propose a prognostic learning procedure which supports (with use of the score of explainable outreach) selection of the most appropriate type of regression model to represent the system's dynamic. In addition, the PL procedure allows also to derive an indicator of the typical length of the time interval, within which predictions made with use of such regression model match the actual future observations sufficiently well (i.e. are credible).

The proposed prognostic learning method was tested on various sets of synthetic data in order to identify its strengths and weaknesses, formulate guidelines for optimal selection of the method parameters (order of the polynomial regression and the length of the learning block) and check how useful the proposed construction of the EO may be in informing us about the immediate future of the observed system. We also indicate how the prognostic learning method can be applied in the context of earth system sciences applying it to analyze historical anthropogenic CO<sub>2</sub> emissions and atmospheric CO<sub>2</sub> concentrations. We conclude that the most robust of the analyzed methods is the one based on linear regression. However, EOs obtained with use of this method and expressing horizons within which linear projections are credible are rather short.

## 8. Acronyms

EO	Explainable outreach
GHG	Greenhouse gases
LB	Learning block (part of the learning sample to which regression model is fitted)
OLS	Ordinary least squares method of fitting a regression function to the data
PL	Learning in a controlled prognostic context (prognostic learning for short)
TB	Testing block (part of the learning sample used to test the EO in order to determine its length)
TSA	Time series analysis (statistical techniques of analysis of time series)

## 9. Literature

Brockwell, P.J., Davis, R.A. (2002): Introduction to Time series and Forecasting, Second Edition. Springer, ISBN 0-387-95351-5

Hastie, T., Tibshirani, R., Friedman, J. (2009): The Elements of Statistical Learning. Data Mining, Inference and Prediction. Springer, ISBN 978-0-387-84858-7

IPCC (2007: FAQ 1.2): What is the Relationship between Climate Change and Weather? In: *Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change* [S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K.B. Averyt, M. Tignor and H.L. Miller (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 104–105.

IPCC (2007: FAQ 8.1): How Reliable Are the Models used to Make Projections of Future Climate Change? In: *Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change* [S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K.B. Averyt, M. Tignor and H.L. Miller (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 600–601.

IPCC (2013: Box 11.1): Climate Simulation, Projection, Predictability and Prediction. In: *Climate Change 2013: The Physical Science Basis. Contribution of Working*



*Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* [T.F. Stocker, D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 959–961.

Meinshausen M., N. Meinshausen, W. Hare, S.C.B. Raper, K. Frieler, R. Knutti, D.J. Frame, M.R. Allen (2009): Greenhouse-gas emission targets for limiting global warming to 2 °C. *Nature*, **458**(7242), 1158–1162; doi: 10.1038/nature08017.

Murphy, K.P. (2012): Machine learning. A Probabilistic Perspective. MIT press, ISBN: 9780262018029

NSF (2012): Decadal and Regional Climate Prediction using Earth System Models (EaSM). National Science Foundation, Arlington VA, USA; Solicitation: <http://www.nsf.gov/pubs/2012/nsf12522/nsf12522.pdf>; FAQs: <http://www.nsf.gov/pubs/2012/nsf12029/nsf12029.jsp>.

Otto, F.E.L., C.A.T. Ferro, T.E. Fricker and E.B. Suckling (2015): On judging the credibility of climate predictions. *Clim. Change*, **132**(1–2), 47–60, doi 10.1007/s10584-013-0813-5.

Pedrycz, W. (2013): Granular computing. Analysis and design of Intelligent systems, CRC press, ISBN 9781439886816.

Puchkova, A., A. Kryazhimskiy, E. Rovenskaya, M. Jonas and P. Żebrowski (2016): Cells (working title). Manuscript, International Institute for Applied Systems Analysis, Laxenburg, Austria (Manuscript under preparation for submission to a scientific journal).

Wolberg, J. (2006): Data Analysis Using the Method of Least Squares. Springer, ISBN 978-3-540-31720-3

Wasserman, L. (2006): All of Nonparametric Statistics, Springer, ISBN 978-0-387-30623-0

## Appendix: Nonparametric kernel-based regression

Nonparametric regression is an alternative to conventional parametric methods. It can be used when we do not want to be limited to the predetermined form of the estimated regression function, when we need to relax some assumptions from the regression analysis, while maintaining a good estimate, or simply when the nature of the data analysed does not allow the selection of a reasonable model.

Among known methods of nonparametric regression (Wasserman 2006), (Härdle 1990), (Fan 1992), (Green & Silverman 1994), (Györfi et al. 2002), e.g. local averaging, regression and smoothing splines (Rice & Rosenblatt 1981,1983), (Stone 1994), (Eubank 1999), wavelets (Nason 1996), (Johnstone & Silverman 1997), (Wang 1996), or orthogonal series (Green & Silverman 1994), the *kernel estimation* is especially noteworthy. It belongs to popular *smoothing techniques* (Simonoff 1996), (Silverman 1986), etc., that allow for estimation even in the case of complicated relationships between explanatory and response variables.

This Appendix is dedicated to application of the prognostic learning method to nonparametric kernel-based regression in real-life case studies from Chapter 5:

- (1) Global CO<sub>2</sub> emissions from technosphere,
- (2) Concentration of the CO<sub>2</sub> in the atmosphere.

### A.1 Kernel functions

Kernel estimation (Wasserman 2006), (Green & Silverman 1994), (Hart 1991), etc. is an extension of local averaging and involves the use of the so-called *kernel function*  $K$ , being nonnegative, symmetric, square integrable, and satisfying the conditions

$$\int_{-\infty}^{+\infty} K(t)dt = 1, \quad \int_{-\infty}^{+\infty} tK(t)dt = 0, \quad \text{and} \quad \int_{-\infty}^{+\infty} t^2K(t) dt < \infty.$$

Given these characteristics the specific choice of a kernel function is not critical. One can take any symmetric probability density function (PDF) of a continuous random variable with zero mean and finite variance<sup>53</sup>.

The most popular choices of kernel functions (Figure A.1) are the *Gaussian* (normal) *kernel* (i.e. PDF of the standard normal distribution), and a few kernels with compact support, like rectangular (uniform), tricube, or the Epanechnikov kernel.

---

<sup>53</sup> The choice of the kernel  $K$  may slightly affect the asymptotic properties of the kernel estimator. For results in finite samples, the difference is negligible.

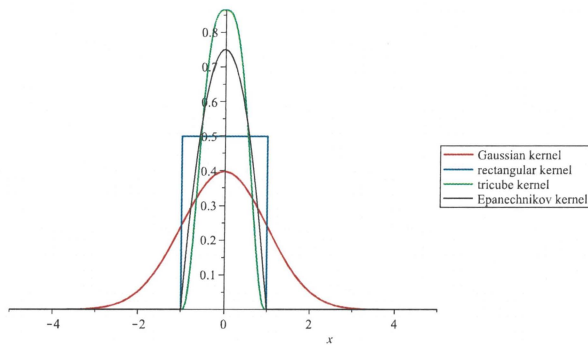


Figure A.1. Four most popular kernel functions.

## A.2 Kernel-based regression methods

Kernel regression has been known for many years and various *kernel estimators* (KE) have been used since then. The most important are (see Table A.1 for overview):

- Nadaraya-Watson KE (Nadaraya 1964), (Watson 1964),
- k-nearest neighbours KE and its modifications (Wasserman 2006),
- Priestley-Chao KE (Priestley & Chao 1972),
- Gasser-Müller KE (Gasser & Müller 1984),
- Local polynomial regression, in particular local linear KE (Li & Racine 2004), (Ruppert & Wand 1994), (Fan & Gijbels 1997).

Some of them have also been considered and analysed in the case of *time series data* or correlated errors (Hart 1991), (Opsomer et al. 2001), (Altman 1990) etc. In Section A.4 two kernel estimators are used: the Nadaraya-Watson KE (NWKE) - mostly because of its simplicity in applications, and the local linear KE (LLKE) - because of its properties and good results, even for small samples.

Each of the aforementioned KEs (except the local polynomial KE) can be considered a *linear smoother* of the form

$$\hat{f}(x) = \sum_{i=1}^n l_i(x) Y_i \quad (A.1)$$

where  $(x_1, Y_1), (x_2, Y_2), \dots, (x_n, Y_n)$ , denote the bivariate data, corresponding to continuous random variables  $x$  and  $Y$ ,

$$\hat{Y}_i = \hat{f}(x_i) + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

and residuals  $\varepsilon_i, i=1, 2, \dots, n$ , are assumed independent<sup>54</sup> and normally distributed, with zero mean and standard deviation  $\sigma > 0$ .<sup>55</sup>

<sup>54</sup> For some kernel-based methods, the independence assumption can be relaxed, in particular, when applying KE to time series data (Section A.3).

<sup>55</sup> In general, standard deviation  $\sigma$  does not need to be constant. Sometimes  $\sigma(x) > 0$ , is considered instead.

Functions  $l_i(x)$ ,  $i=1,2,\dots,n$ , satisfy

$$\sum_{i=1}^n l_i(x) = 1$$

and take various forms, depending on the estimator considered (Table A.1).

**Table A.1.** Overview of the most popular kernel regression estimators. The methods further used in this Appendix are marked in green colour.

KE	$l_i(x)$ in (A.1)	Properties & Remarks
Nadaraya-Watson (NWKE)	$l_i(x) = \frac{K\left(\frac{x-x_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right)}$	<ul style="list-style-type: none"> <li>- local constant estimator</li> <li>- can be adopted for (discrete) time series case</li> <li>- several 'rule of thumb' for selection of bandwidth <math>h</math></li> <li>- biased (design bias and strong boundary bias)</li> <li>- require large samples</li> </ul>
k-nearest neighbours (weighted) (k-NNKE)	$l_i(x) = \frac{K\left(\frac{x-x_i}{R}\right)}{\frac{1}{n}\sum_{j=1}^n K\left(\frac{x-x_j}{R}\right)}$ <p>where <math>R</math> denotes the distance between <math>x</math> and its k-nearest neighbour;</p>	<ul style="list-style-type: none"> <li>- for rectangular kernel, it reduces to NWKE</li> <li>- <math>k = 2nhf(x)</math>, where <math>f</math> denotes the PDF of the explanatory variable</li> <li>- biased (both design and boundary bias)</li> <li>- various modifications and simplifications; various weights</li> <li>- require large samples</li> </ul>
Priestley-Chao (PCKE)	$l_i(x) = \frac{x_i - x_{i-1}}{h} K\left(\frac{x-x_i}{h}\right)$	<ul style="list-style-type: none"> <li>- applicable to compactly supported data (rescaling option, with good results)</li> <li>- requires kernel function with compact support</li> <li>- no design bias, but strong boundary bias</li> <li>- require large samples</li> </ul>
Gasser-Müller (GMKE)	$l_i(x) = \frac{1}{h} \int_{v_{i-1}}^{v_i} K\left(\frac{x-u}{h}\right) du$ <p>where <math>x_i \leq v_i \leq x_{i+1}</math></p>	<ul style="list-style-type: none"> <li>- continuous version of PCKE</li> <li>- partition <math>\{v_i\}</math>, <math>i=1,\dots,n-1</math> required</li> <li>- applicable to compactly supported data (rescaling option with good results)</li> <li>- requires kernel function with compact support</li> <li>- no design bias, but boundary bias</li> <li>- require large samples</li> </ul>
Local linear (LLKE)	$l_i(x) = \frac{b_i(x)}{\sum_{j=1}^n b_j(x)}$ <p>where</p> $b_i(x) = K\left(\frac{x-x_i}{h}\right) (S_{n,2}(x) - (x_i - x)S_{n,1}(x))$ $S_{n,j}(x) = \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) (x_i - x)^j.$	<ul style="list-style-type: none"> <li>- particular case of local polynomial regression</li> <li>- local linear smoother</li> <li>- can be adopted for (discrete) time series case</li> <li>- no boundary nor design bias</li> <li>- require large samples, although thanks to good local fit, better results for smaller samples</li> </ul>
Local polynomial KE	Estimate locally (at a point $x$ ) that polynomial of degree $p$ , which approximates $F(x)$ in a small neighbourhood of the point $x$ , in the best way.	<ul style="list-style-type: none"> <li>- becomes NWKE for <math>p=0</math>, and LLKE for <math>p=1</math></li> <li>- in general cannot be represented as a linear smoother given by (A.1)</li> <li>- no boundary nor design bias</li> <li>- require large samples, although thanks to good local fit, reasonable results for smaller samples;</li> <li>- for larger <math>p</math> requires larger samples</li> </ul>

### A.2.1 The problem with bandwidth selection

Weights  $l_i(x)$ ,  $i=1, \dots, n$ , in formula (A.1) depend on kernel function  $K$ , and a *smoothing parameter*  $h > 0$  (also called a *bandwidth*)<sup>56</sup>, such that

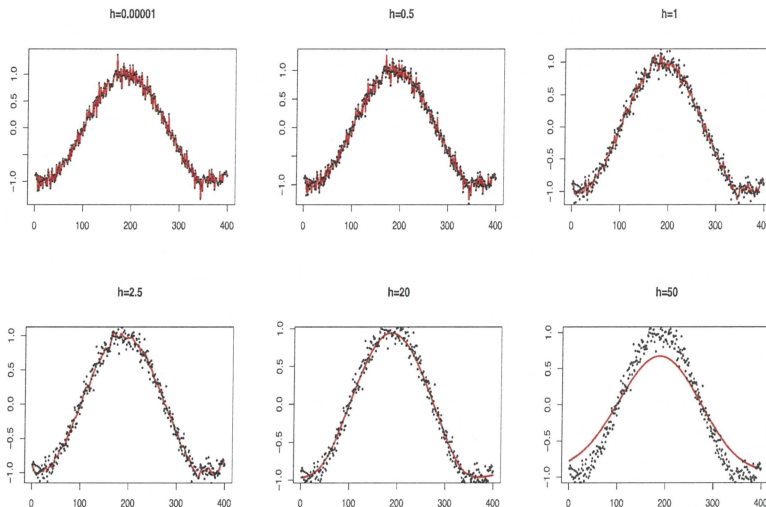
$$h \rightarrow 0 \text{ but } nh \rightarrow \infty, \quad \text{as } n \rightarrow \infty.$$

The choice of optimal value of the smoothing parameter is crucial<sup>57</sup> and corresponds to a problem of finding the 'golden mean', by minimizing the *mean squared error* (MSE), being the sum of squared bias<sup>58</sup> and sampling variance

$$MSE(\hat{f}(x)) = \text{bias}(\hat{f}(x))^2 + \text{Var}(\hat{f}(x)),$$

or its asymptotic and integrated versions.

Varying the bandwidth  $h$  controls the smoothness of the estimated regression function. Larger  $h$  results in a smoother curve, but sometimes with a worse fit and hence a larger variance. Smaller  $h$  in turn means a better fit, with smaller variance, may however cause a greater bias (see Figure A.2). Too large  $h$  means therefore *oversmoothing* (possibly failing to reflect the character of the data analysed), while too small – *undersmoothing*.



**Figure A.2.** Varying the smoothing parameter: examples of the NWKEs fitted to the data following sinusoidal trend (from Section 4.3.5) given by  $g(x) = \sin(0.018 \times (x - 100))$ , with standard deviation of noise  $\sigma = 0.01 \times (\max g - \min g)$ , where  $n=400$ , for various values of  $h$ , and using the Gaussian kernel.

<sup>56</sup> There are also methods involving variable bandwidths. Here, we focus on methods with fixed bandwidth.

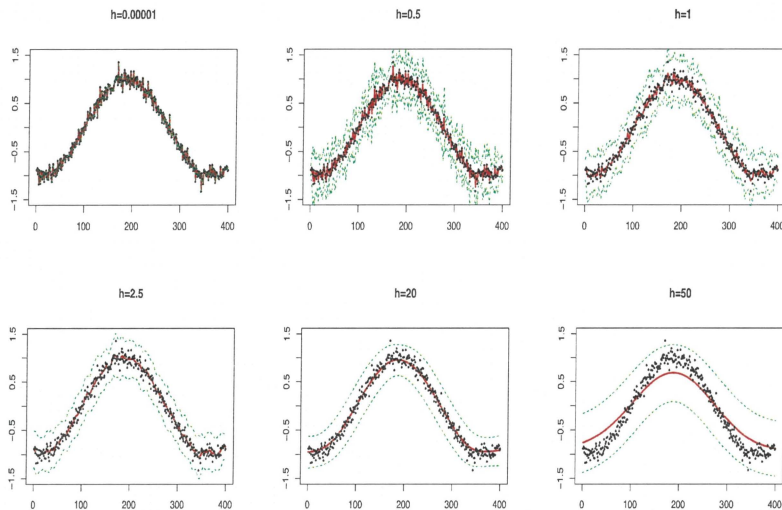
<sup>57</sup> See e.g. (Wasserman 2006), (Simonoff 1996), etc.

<sup>58</sup>  $\text{bias}(\hat{f}(x)) = E(\hat{f}(x)) - f(x)$

The shape of  $\hat{f}(x)$  changes for various values of  $h$ . The plots in the first row illustrate the case, when the smoothing parameter is too small. The variance in that case is really small, which results in a good fit, but it is done for the price of undersmoothed and fluctuating regression curve. The sample is relatively large ( $n=400$ ), so the ‘noisy’ shape of the estimator is connected with overfit. Increasing  $h$ , gives smoother  $\hat{f}(x)$ , what can be noticed for  $h=2.5$  and  $20$ . The plot in the lower right corner of Figure A.2 illustrates the evident underfit (resulting in large variance) - the curve is oversmoothed and does not grasp the behaviour of the data.

It is worth noting that, despite the problem with bandwidth selection, even the simple NWKE approximates the regression function fairly well. Despite the almost ten-fold difference between the values of  $h$ , the two figures bottom left, look satisfactorily. To assess which of them really performs better, one can look at confidence or prediction intervals (the latter works better in this regard, due to more emphasis on standard error).

Since the degree of smoothing corresponds to the variance of  $\hat{f}(x)$ , it also affects the width of prediction intervals<sup>59</sup>. Oversmoothing causes too wide intervals (interpreted as large uncertainty of results), while undersmoothing - too narrow (Figure A.3).



**Figure A.3.** Varying the smoothing parameter and illustrating its impact on 95% prediction intervals (green dashed lines): examples of the NWKEs fitted (red solid lines) to the data following sinusoidal trend (from Section 4.3.5) given by  $g(x) = \sin(0.018 \times (x - 100))$ , with standard deviation of noise  $\sigma = 0.01 \times (\max g - \min g)$ , where  $n=400$ , for various values of  $h$ , and using the Gaussian kernel.

In general,  $h$  depends on the sample size  $n$ , and asymptotically  $h \propto n^{-\frac{1}{5}}$ . The formulas for optimal  $h$ , are different, for various kernel methods. For instance, the

<sup>59</sup> see Section A.2.2

optimal value of the smoothing parameter<sup>60</sup> in the case of the NWKE satisfies the following formula<sup>61</sup>

$$h = \left( \frac{\sigma^2 \int_{-\infty}^{+\infty} K^2(x) dx \int_{-\infty}^{+\infty} f(x)^{-1} dx}{n \int_{-\infty}^{+\infty} x^2 K(x) dx \int_{-\infty}^{+\infty} \left( \hat{r}''(x) + \hat{r}'(x) \frac{f'(x)}{f(x)} \right)^2 dx} \right)^{\frac{1}{5}} \quad (\text{A.2})$$

while for the LLKE<sup>62</sup>

$$h = \left( \frac{\sigma^2 \int_{-\infty}^{+\infty} K^2(x) dx \int_{-\infty}^{+\infty} f(x)^{-1} dx}{n \int_{-\infty}^{+\infty} x^2 K(x) dx \int_{-\infty}^{+\infty} (\hat{r}''(x))^2 dx} \right)^{\frac{1}{5}}. \quad (\text{A.3})$$

The values  $\int_{-\infty}^{+\infty} K^2(x) dx$  and  $\int_{-\infty}^{+\infty} x^2 K(x) dx$  depend on the kernel used. For the Gaussian kernel  $\int_{-\infty}^{+\infty} K^2(x) dx \cong 0.28$ , while the latter one represents the variance of the standard normal distribution, i.e.  $\int_{-\infty}^{+\infty} x^2 K(x) dx = 1$ . But formulas (A.2) and (A.3) also involve unknown regression function  $\hat{r}(x)$ , that needs to be estimated, unknown variance  $\sigma^2$ , as well as  $f(x)$ , i.e. the PDF of the explanatory variable. The methods to estimate them, depend on problem requirements, the data to be analysed, and on the KE considered. In particular, for the LLKE or the GMKE,  $\sigma^2$  can be estimated by (asymptotically unbiased) estimator of the form (Gajek & Kaluszka 1993)

$$\hat{\sigma}^2 = \frac{1}{6(n-2)} \sum_{i=1}^{n-2} (Y_{i+2} - 2Y_{i+1} + Y_i)^2$$

For the NWKE, also much simpler

$$\hat{\sigma}^2 = \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (Y_{i+1} - Y_i)^2$$

can be used. However, both formulas work well mostly for large samples.

The density function of the explanatory variable can be estimated using nonparametric methods, like kernel density estimation (Silverman 1986), or (less often) parametric (e.g. MLE, provided that, we have additional information on that variable and its distribution). In complicated cases, semiparametric methods can also be used e.g. (Jarnicka 2009). To estimate  $\hat{r}''(x)$  and  $\int_{-\infty}^{+\infty} (\hat{r}''(x))^2 dx$  additional information on the data is required, since the latter one corresponds to the curvature of the estimated regression curve, or approximation by the curvature of some known curve can be used. Similarly the term  $\hat{r}'(x) \frac{f'(x)}{f(x)}$ , which is responsible for the bias.

For some estimators, like the NWKE, there exist a few ‘rules of thumbs’ for finding reasonable value of  $h$ , working well in most cases, especially for large samples (but being less useful, when applied to time series data or in the case of correlated errors).

Moreover, the smoothing parameter can also be chosen by the *cross-validation* (CV) criterion<sup>63</sup>

<sup>60</sup> see e.g. (Wasserman 2006), (Green & Silverman 1994), etc.

<sup>61</sup> The term  $\hat{r}'(x) \frac{f'(x)}{f(x)}$  in (A.2) denotes the *design bias*, typical for the NWKE (it is not present for the LLKE).

<sup>62</sup> see e.g. (Ruppert & Wand 1994), (Fan & Gijbels 1997), etc.

<sup>63</sup> see e.g. (Wasserman 2006), etc.

$$CV(h) = \sum_{i=1}^n (Y_i - \hat{r}(x_i))^2 \theta(z(x_i)),$$

where

$$z(x_i) = \frac{K(0)}{\sum_{j=1, j \neq i}^n K\left(\frac{x_i - x_j}{h}\right)}.$$

The penalizing function  $\theta(\cdot)$  takes various forms, e.g.,  $\theta(z) = \frac{1}{(1-z)^2}$ , (generalized CV), and  $\theta(z) = e^{2z}$  (AIC – Akaike’s Information Criterion), etc. and ensures various properties (i.e. possibly small bias or variance)<sup>64</sup>. The values of  $h$  obtained using the CV criteria are usually close to the MSE optimal ones. The problem starts with violation of independence assumption on residuals, as correlation may decrease the bandwidth indicated by the CV criterion, so the curve obtained is undersmoothed (Opsomer et al. 2001), (De Brabanter et al 2011), etc.

### A.2.2 100%(1- $\alpha$ )-Prediction Intervals

Choosing the right  $h$  is of great importance, due to the expected estimation result. Since this is a compromise between minimizing the variation of the KE and its bias, when choosing a bandwidth, we can put emphasis on that of them, which is more important for a particular application. In this report, we focus primarily on the variance (analysing it, but not trying to make it as small as it gets, as it may affect the EO), which determines the prediction intervals.

According to the Central Limit Theorem (CLT), regression estimates  $\hat{r}(x)$  in (A.1) have an asymptotic normal distribution

$$\frac{\hat{r}(x) - \text{bias}(\hat{r}(x))}{\sqrt{\text{Var}(\hat{r}(x))}} \rightarrow N(0,1)$$

Assuming no bias, the asymptotic 100%(1- $\alpha$ ) - prediction interval is of the form

$$\hat{r}(x) \pm z_{1-\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{r}(x)) + \hat{\sigma}(x)^2}$$

where  $z_{1-\frac{\alpha}{2}}$  denotes the  $(1 - \frac{\alpha}{2})$ th quantile of the standard normal distribution. For in-sample points, i.e. for points from the LB,  $\hat{r}(x)$  denotes the KE, and  $\hat{\sigma}^2(x)$  an estimate of the variance of residuals (corresponding to the standard error), while for new observations  $x^*$ ,  $\hat{r}(x^*)$  denotes the prediction at  $x^*$ , and  $\hat{\sigma}(x^*)^2$  prediction error. For the NWKE and the LLKE the variance is asymptotically equal

$$\text{Var}(\hat{r}(x)) \approx \frac{\hat{\sigma}^2(x) \int_{-\infty}^{+\infty} K^2(t) dt}{nhf(x)},$$

which gives the in-sample *prediction bands (PB)* of the form

$$\hat{r}(x_i) \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{\sigma}^2(x_i) \int_{-\infty}^{+\infty} K^2(t) dt}{nhf(x_i)} + \hat{\sigma}(x_i)^2} \quad (\text{A.4})$$

<sup>64</sup> This refers to finite samples, as they all guarantee the same asymptotic properties.

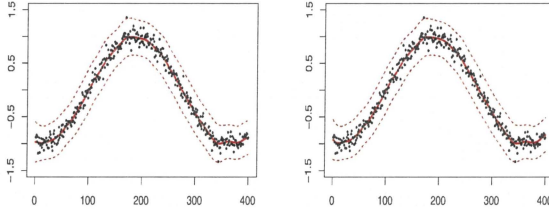


and

$$\hat{f}(x^*) \pm z_{1-\frac{\alpha}{2}} \sqrt{\sum_{i=1}^n \frac{\hat{\sigma}^2(x_i) \int_{-\infty}^{+\infty} K^2(t) dt}{nhf(x^*)} + \hat{\sigma}(x^*)^2} \quad (A.5)$$

for a new observation  $x^*$  (Green & Silverman 1994).

Formula (A.4) was used to construct the prediction intervals in Figure A.3. It is worth mentioning that, the approximately optimal value of the smoothing parameter is  $h=7.72$ , while for the LLKE applied to the same data  $h=8.06$  (see Figure A.4 for 95% in-sample PBs). Formula (A.5) will in turn be used to construct the EO in the procedure described in Section 3.2.



**Figure A.4.** 95% in-sample (LB) prediction bands (dashed) for the NWKE (left) and the LLKE (right) with the Gaussian kernel and approximately optimal bandwidths  $h=7.72$  and  $h=8.06$  for NWKE and LLKE respectively; Thanks to a large sample (LB=400, data set from Figure A.2 and A.3) and independent observations the results are almost identical. The residual standard error is equal  $0.094$  and  $0.093$  for NWKE and LLKE respectively.

### A.3 Kernel estimation of time series data

In this section we focus on the *time series case*, where the time points are fixed and equally spaced. Following notation from Section 3.1, let the learning block (LB) contain  $n$  observations  $X_1, X_2, \dots, X_n$ , taken at the time points  $t_1, \dots, t_n$  where  $t_i = \frac{i}{n}$ ,  $i=1, \dots, n$ .

Consider

$$\hat{x}(t) = \hat{f}(t) + \varepsilon_t,$$

where  $x(t) = X_t$  is a value of the observation taken at time  $t$ , and the noise term  $\varepsilon_t$  is normally distributed with zero mean and standard deviation  $\sigma > 0$ .<sup>65</sup> We assume that

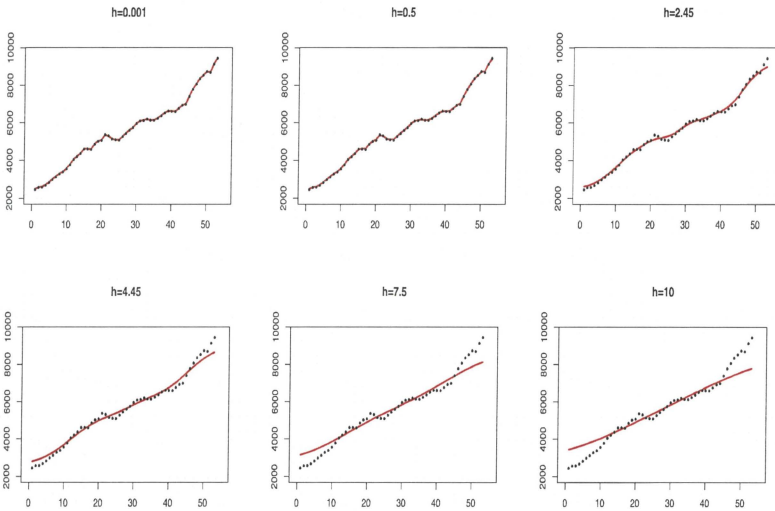
<sup>65</sup> Assumptions on residuals, when compared to parametric regression techniques can be relaxed. Two scenarios are considered in the literature: (1) allowing non-normal distribution, but ensuring covariance stationarity and possibly weak correlation ((Brabanten et al. 2011), (Opsomer et al. 2001)), or (2) ensuring normality and analyzing correlation structure e.g.(Li & Li, 2009). Both lead to problems with appropriate bandwidth selection, the second one however allows for asymptotically better results, in particular in view of predictions and the EO.

residuals  $\varepsilon_t$ ,  $t = 0, 1, 2, \dots$ , are correlated and their correlation decreases in inverse proportion to the distance between them<sup>66</sup>.

When analysing a time series, one has to deal with specific nature of the data, resulting in a need for modifications in optimal bandwidth selection methods. Moreover, the problem with applying the kernel methods to time series data is also connected with discrete distribution of explanatory variable  $t$  (discrete time), which has to be approximated by a continuous estimate.

### A.3.1 Bandwidth selection in the time series case

The problem of optimal bandwidth selection, described and illustrated in Section A.2, becomes now more visible. The time points are equally spaced, and what is more important, the data points (and hence the residuals) are correlated, so the shape of the estimated regression function changes considerably, when varying the smoothing parameter (see Figures A.5 (NWKE) and A.6 (LLKE) for examples).



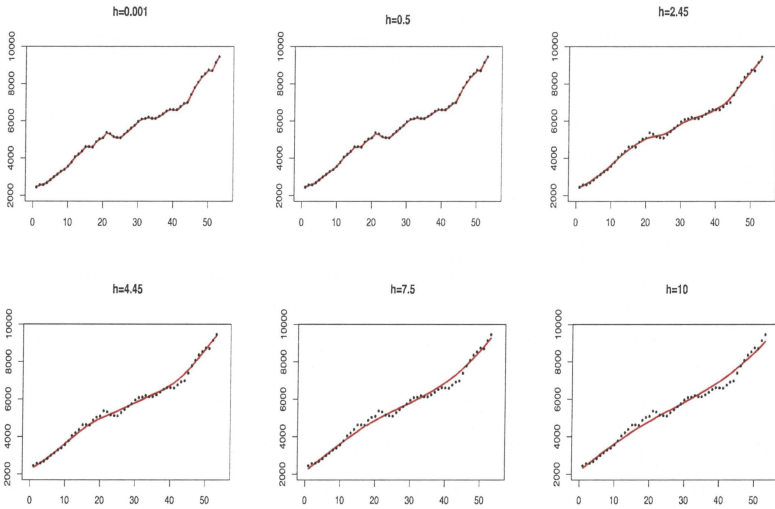
**Figure A.5.** Varying the smoothing parameter: examples of the NWKEs fitted to the data on global CO<sub>2</sub> emissions from technosphere ( $n=53$ ) for various values of  $h$ , and the Gaussian kernel.

The NWKE is fitted to the data on global CO<sub>2</sub> emissions from technosphere. To illustrate problems with finding the optimal bandwidth for the time series case, we take the whole sample, consisting of  $n=53$  data points and consider six exemplary values of  $h$ .

<sup>66</sup> This assumption corresponds to the condition  $\text{Corr}(\varepsilon_{t_i}, \varepsilon_{t_j}) = \rho(t_i - t_j)$ , based on unknown stationary correlation function  $\rho(\cdot)$ . This allows for correlation decaying, when  $n \rightarrow \infty$ , and hence better results for large samples. We will not however be interested in analysing the correlation structure in detail, using only ‘independence-like’ approximations.

It is easy to observe, that values  $h=4.45$ ,  $7.5$ , and  $10$  are too large, resulting in oversmoothing, which means that, in practice only a central part of the data is estimated, and the result is rather poor. On the other hand,  $h=0.001$  is too small, showing a perfect fit, with no visible uncertainty. Both  $h=0.5$  and  $2.45$  seem to be quite good. The first one, seems to better grasp the behaviour of the data, the latter one however leaves more room for possible improvement (and may be better in view of the EO).

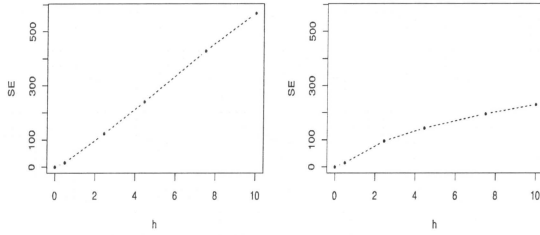
You may notice that, in four of the six examples given, we have to deal with the boundary bias, characteristic for the NWKE. It can significantly affect the length of the EO, since it cannot be overcome by slightly larger smoothing, and greater variance. Therefore, for the EO analysis, also the LLKE is used, as it is free from the boundary bias. For comparison, in Figure A.6, the LLKE is fitted to the same data series, using the Gaussian kernel, and taking the same exemplary values of  $h$ .



**Figure A.6.** Varying the smoothing parameter: examples of the LLKEs fitted to the data on global CO<sub>2</sub> emissions from technosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.

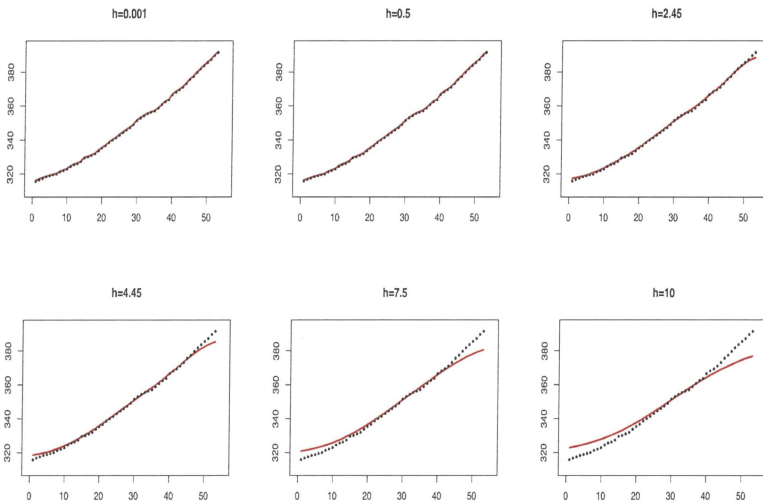
It is easy to observe that, the application of the LLKE (Figure A.6) gives better results than in the case of the NWKE (Figure A.5). This is primarily related to the lack of boundary bias. Due to the fact that, the estimator is fitted to the data locally, even in the case when smoothing parameter  $h$  is too large (e.g. for  $h=4.5$  or  $7.5$ ) the LLKE seems to properly grasp the general shape of the estimated relationship.

This is also reflected in the variation of the standard error in those cases (Figure A.7), as the standard error (SE) increases much faster in the case of the NWKE.

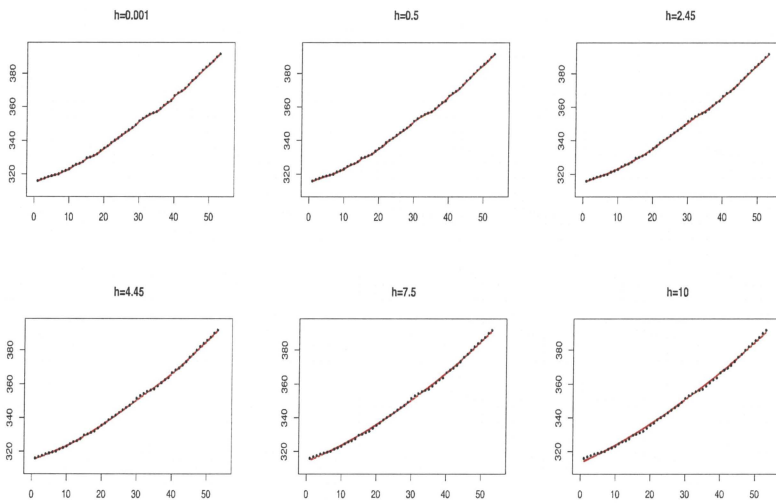


**Figure A.7.** Relationship between the smoothing parameter and the standard error for the NWKE (left) and the LLKE (right) considered in Figures A.5 and A.6.

The results obtained are connected with the type of data. The same analysis conducted for the second data set from Chapter 5 i.e. concentration of the CO<sub>2</sub> in the atmosphere, shows slightly different results (Figures A.8 and A.9).

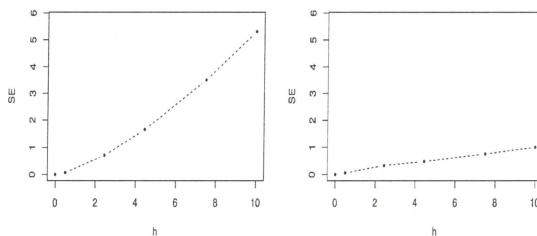


**Figure A.8.** Varying the smoothing parameter: examples of the NWKEs fitted to the data on concentration of the CO<sub>2</sub> in the atmosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.



**Figure A.9.** Varying the smoothing parameter: examples of the LLKEs fitted to the data on concentration of the CO<sub>2</sub> in the atmosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.

Although varying the smoothing parameter changes the results, thanks to the linear trend in the data, the KEs used to estimate the regression function seem to work well. As previously (Figures A.5 and A.6) the LLKE performs better, but the difference is not as evident as for the emission data. The main reason is the scale of the standard errors. The comparison of standard errors shows that, the results of the NWKE are better (Figure A.10), i.e. the standard errors for the LLKE are smaller and the difference is significant, as presented in Figure A.7.



**Figure A.10.** Relationship between the smoothing parameter and the standard error for the NWKE (left) and the LLKE (right) considered in Figures A.8 and A.9.

Since in the case of time series data, the smoothing parameter cannot be chosen by the CV criterion (usually correlation causes oversmoothing (Opsomer et al. 2001)), formulas (A.2) and (A.3) should be used.

To estimate unknown factors in (A.2) and (A.3), some additional assumptions are required.

- As a kernel function  $K$ , we take the Gaussian kernel, so

$$\int_{-\infty}^{+\infty} K^2(x)dx \cong 0.28, \quad \int_{-\infty}^{+\infty} x^2 K(x)dx = 1.$$

- The explanatory variable has discrete uniform distribution, and can therefore be roughly approximated by its continuous version. In particular, the PDF of the uniform distribution is nonzero only on  $[a, b]$ . For simplicity, the factor related to that PDF is constant and can therefore be omitted. To estimate the PDF of the explanatory variable in PB, we use kernel density estimation with the bandwidth chosen by the Silverman's rule of thumb  $h = \left(\frac{1.06\sigma}{n}\right)^{\frac{1}{5}}$  (Silverman 1986).
- For simplicity, we assume that, the unknown regression function is close to a straight line. The factor  $\int_{-\infty}^{+\infty} (\hat{r}''(x))^2 dx$  is constant and can also be omitted.
- The variance  $\hat{\sigma}^2$  is assumed constant, and is estimated by

$$\hat{\sigma}^2 = \frac{1}{6(n-2)} \sum_{i=1}^{n-2} (Y_{i+2} - 2Y_{i+1} + Y_i)^2 \quad (A.6)$$

Therefore, in Section A.4, to find the bandwidth  $h$ , we use the following rule of thumb

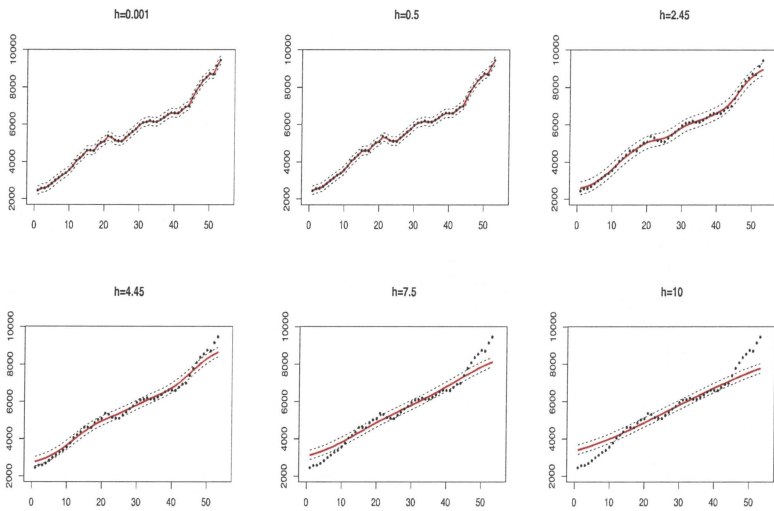
$$h = \left(\frac{\hat{\sigma}^2 0.28}{n}\right)^{\frac{1}{5}} \quad (A.7)$$

That corresponds to known rules of thumbs for NWKE (Green & Silverman 1994), and is used for both NWKE and the LLKE. In this case formula (A.7) corresponds rather to the optimal bandwidth for the LLKE (no design bias factor), but assuming no bias in the NWKE and approximating  $h$  by the same formula, as for the LLKE leads to a slight oversmoothing (and hence that assumption becomes reasonable).

### A.3.2 In-sample prediction bands – the time series case

Given the time series data, the independence assumption is not satisfied and, in general, some asymptotic properties of the KE may not be satisfied (Hart 1991). However, for some cases of correlation structure, in particular, assuming the correlation decays in inverse proportion to the distance between observations (Opsomer et al. 2001), or for the AR correlation structure (Li & Li 2009), asymptotic properties of the KE are close the ones that hold in the independent case. Moreover, generalized version of the CLT, indicates the asymptotic normal distribution, which allows for the use of formulas (A.4) and (A.5) to find the asymptotic prediction bands.

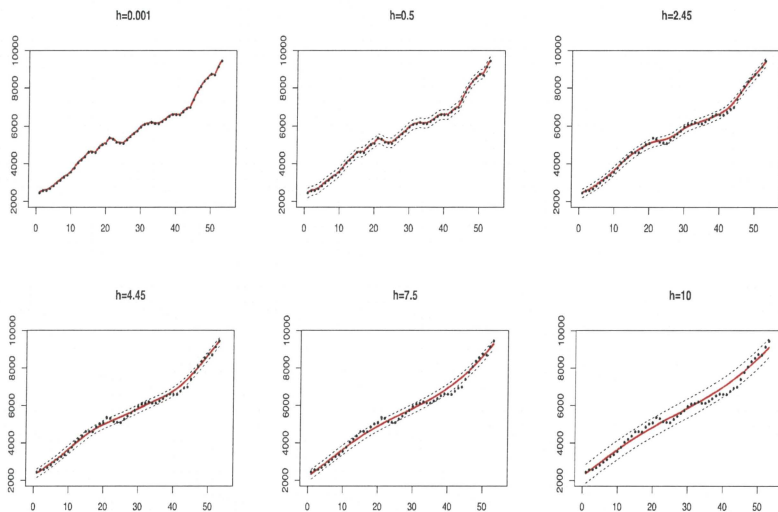
The construction of the PBs is connected with the choice of the smoothing parameter. Adding 95% prediction bands helps in illustrating differences between the results obtained in Section A.3.1 for various values of  $h$ .



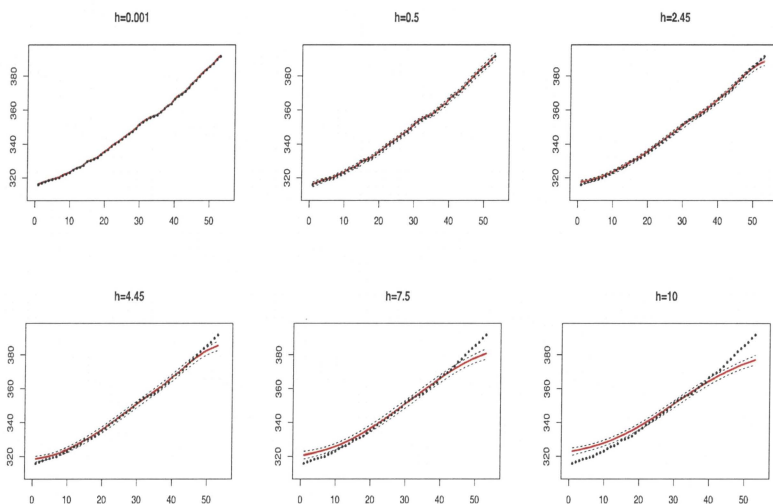
**Figure A.10.** Varying the smoothing parameter and illustrating its impact on the variance in terms of 95% prediction bands (black dashed lines): examples of the NWKEs fitted to the data on global CO<sub>2</sub> emissions from technosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.

For  $h=0.001$ , prediction bands do not cover all the data points depicted, since the variance of the estimated regression function is too small and the prediction interval too narrow. Values  $h=0.5$  and  $2.45$  provide different results – the latter appears to be slightly too large, increasing the variance and causing the wider prediction interval. For  $h=10$ , the regression estimate is obviously oversmoothed. The shape of the data is not properly reflected, and despite the large variance only few data points fall within the prediction bands<sup>67</sup>.

<sup>67</sup> That effect is partly connected with boundary bias of the NWKE.

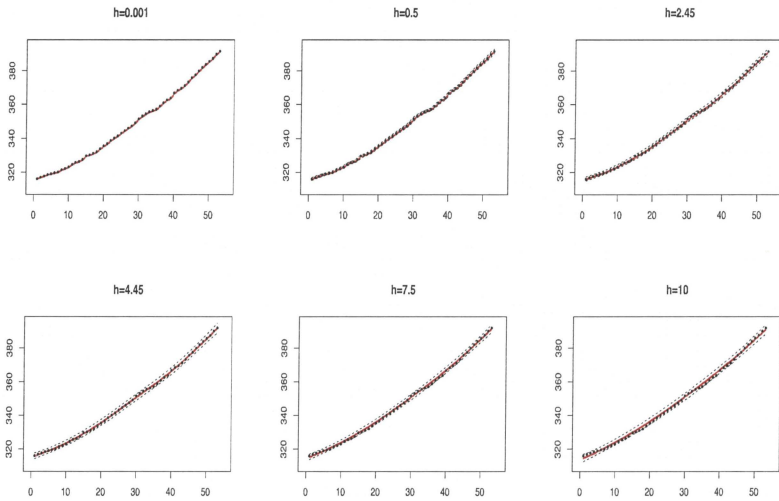


**Figure A.11** Varying the smoothing parameter and illustrating its impact on the variance in terms of 95% prediction bands (black dashed lines): examples of the LLKEs fitted to the data on global CO<sub>2</sub> emissions from technosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.



**Figure A.12** Varying the smoothing parameter and illustrating its impact on the variance in terms of 95% prediction bands (black dashed lines): examples of the NWKEs fitted to the data on concentration of the CO<sub>2</sub> in the atmosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.





**Figure A.13** Varying the smoothing parameter and illustrating its impact on the variance in terms of 95% prediction bands (black dashed lines); examples of the LLKEs fitted to the data on concentration of the  $\text{CO}_2$  in the atmosphere ( $n=53$ ) for various values of  $h$ , with Gaussian kernel.

#### A.4 Real-life case studies

The methods of prognostic learning from Chapter 3, in particular the procedure for assessing the EO, presented in Section 3.3, are applied to real-life case studies, considered in Chapter 5: (1) global  $\text{CO}_2$  emissions from technosphere, and (2) concentration of the  $\text{CO}_2$  in the atmosphere.

The ability of learning is tested in terms of the EO (described in Section 3.2) for both aforementioned kernel regression estimators: the LLKE and much simpler NWKE.

The most important problem when using nonparametric methods is their requirement of a large sample size. Each of nonparametric methods (including kernel regression) depends on the sample size in a different way. Due to the asymptotic properties of kernel estimators, the sample should be sufficiently large, although it is difficult to specify the threshold above which the results will be good. Conducted analyses and simulations (Wasserman 2006), (Green & Silverman 1994) indicate that, this depends on the type of data, in particular on their distribution. Also, correlation of data (as in the time series case) requires a larger number of test points (Opsomer et al. 2001), (Hart 1991). It can therefore be expected, that for LBs of 25 or slightly more training points, the results may not be satisfactory, which in some way will influence the EO.

### A.4.1. Procedure for analysing the EO, in the case of the kernel regression

To test the ability of prognostic learning, the following procedure is considered. Given the sample of  $n = n_1 + n_2$  data points, we perform the following steps.

**Step 1.** We take the learning block (LB) of  $n_1$  data points.

- The unknown variance of residuals is estimated by (A.6)
- The smoothing parameter is found by (A.7)
- The NWKE and the LLKE are used.
- The model assumptions are verified.
- The in-sample 95% prediction bands are found for both NWKE and LLKE, using (A.4)

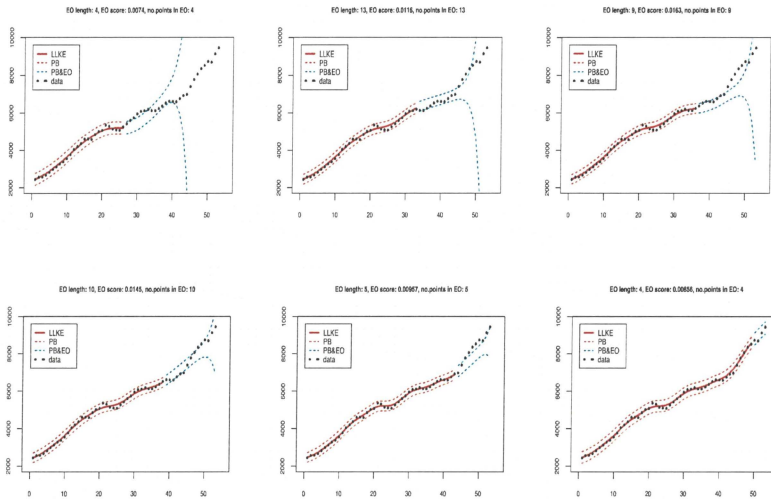
**Step 2.** We take the testing sample of  $n_2$  data points.

- The out-of-sample 95% prediction bands are found for both the NWKE and LLKE, using (A.5)
- The length and the score of the EO are found, using the procedure described in Section 3.2.

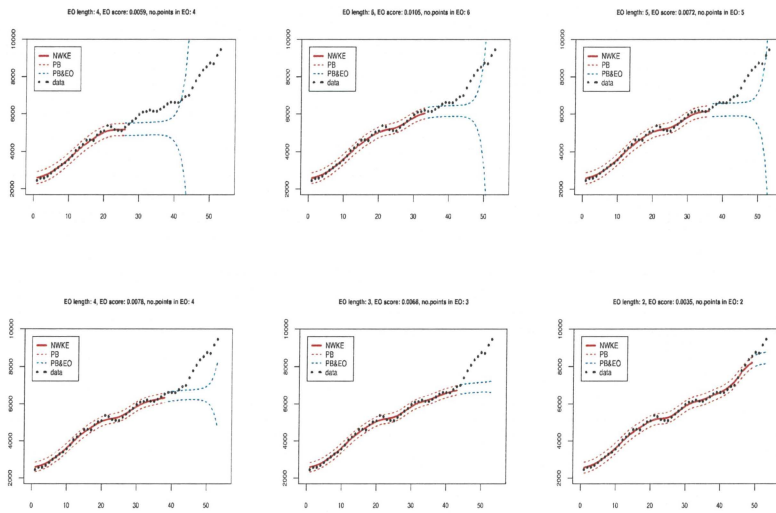
**Step 3.** We increase the LB by one and repeat Step 1 and Step 2.

### A.4.2. Global CO<sub>2</sub> emissions from technosphere

The procedure described in Section A.4.1 is applied, starting with  $n_1 = 25$ . The six exemplary stages are presented in Figures A.14 (for the LLKE) and A.15 (NWKE).



**Figure A.14.** Six exemplary stages of the PL procedure (LB lengths: 26, 33, 36, 38, 43, 49): the LLKE using the Gaussian kernel.



**Figure A.15.** Six exemplary stages of the PL procedure (LB lengths: 26, 33, 36, 38, 43, 49): the NWKE using the Gaussian kernel

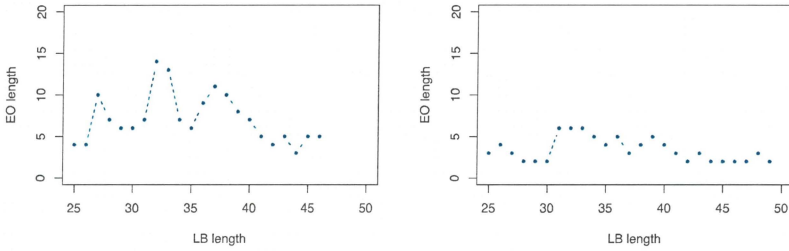
For both estimators the 95% out-of-sample predictions bands for the shortest LBs open quite fast<sup>68</sup>, but the PBs for the NWKE, in particular for the shortest LBs, seem to stabilize at first, increasing rapidly after a few out-of-sample points. This is connected with the boundary bias of the NWKE, present in particular for small samples.

The PBs for the LLKE better reflect the estimated relationship between explanatory and response variables, which also results in the longer EO. The prediction intervals for the NWKE are wider, which is connected with the greater standard errors, and results in lower EO scores (Figure A.17).

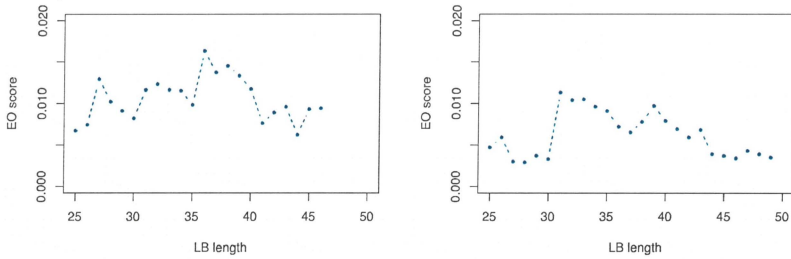
As oppose to the EO lengths presented in Figure 37, as a results of using parametric, linear regression, no decreasing trend can be observed, for  $LB > 30$ . The EO lengths decrease and increase, for the LLKE having peaks at  $LB=32$  (local maximum), 34 (local minimum), 37 (max), and then 42 (min), 43 (max), 44 (min) and 47 (max). For  $LB > 48$ , all the remaining data points are within the PBs, giving the infinite length.

It is worth mentioning, that in spite of differences in the EO lengths, the results obtained using both estimators show similar monotonic behaviour (Figure A.16). Similar effect can be observed for the EO scores (Figure A.17). This means that the EO depends on the data. Since in the case of the LLKE standard errors are smaller than for the NWKE, the prediction intervals for LLKE are narrower, and the data type affects the EO outcome stronger.

<sup>68</sup> This is connected with prediction errors, increasing really fast. The in-sample errors behaviour is completely different (Figures A.10 and A.11), as they seem to be constant.



**Figure A.16.** The EO length as a function of the LB, in the case of the LLKE (left) and the NWKE (right).



**Figure A.17.** The EO score as a function of the LB, in the case of the LLKE (left) and the NWKE (right).

The comparison of the results for the LLKE and NWKE is presented in Table A.2. The conducted analysis shows that, the LLKE performs better, giving longer EOs - between 4 and 14 data points (Figure A.16).

Moreover, starting with the LB=47, all the remaining data points fall within the PBs. The resulted EO lengths for the NWKE are in turn more stable, giving values between 2 and 6.

**Table A.2** Prognostic learning – comparison of the LLKE and NWKE results, when applied to the data on CO<sub>2</sub> emissions from technosphere.

Results		LLKE	NWKE
<b>EO</b>	max length	finite: 14 (for LB=32) ∞ (for LB≥47)	finite: 6 (for LB=31, 32, and 33) ∞ for LB≥50
	min length	4 (for LB=25, 26, 42 and 44)	2 (for LB=28, 29, 30, 41, 44-47, and 49)
	∞	for LB≥47 all tested data points fall within the PBs	for LB≥50 all tested data points fall within the PBs
	score	0.0062 – 0.0163 for LB<47 ∞ for LB≥47	0.0029 – 0.0113 for LB<50 ∞ for LB≥47
<b>Residuals</b>	normality	$\varepsilon_t$ normally distributed (Shapiro-Wilk test, $p$ -values>0.2)	$\varepsilon_t$ normally distributed (Shapiro-Wilk test, $p$ -values>0.1)
	zero mean	ok (t-test, $p$ -values>0.2)	ok (t-test, $p$ -values>0.2)
	correlation	autocorrelation at lag 1 and 2, (ACF, Box-Pierce test)	autocorrelation up to lag 5 or 6 (ACF, Box-Pierce test)

*A.4.2 Concentration of the CO<sub>2</sub> in the atmosphere.*

Now the procedure described in Section A.4.1 is applied to the second data set. As previously, we start with  $n_1 = 25$  and then increase the LB length by one. The six exemplary stages of the procedure are presented in Figures A.18 (for the LLKE) and A.19 (for the NWKE).

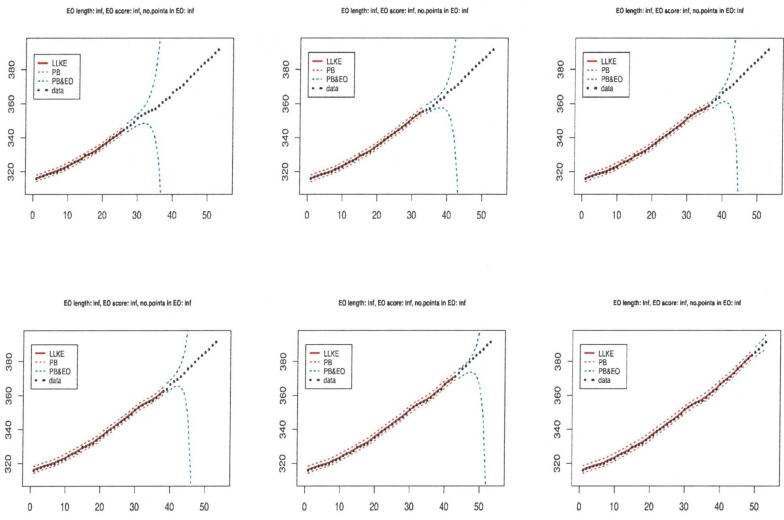


Figure A.18. Six exemplary stages of the PL procedure (LB lengths: 26, 33, 36, 38, 43, 49): the LLKE using the Gaussian kernel

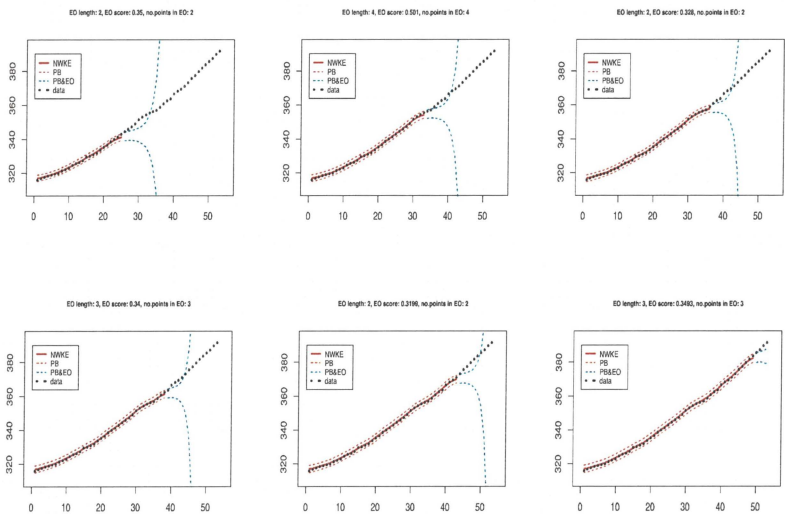
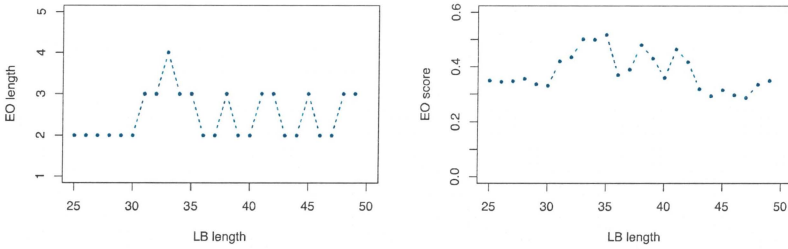


Figure A.19. Six exemplary stages of the PL procedure (LB lengths: 26, 33, 36, 38, 43, 49): the NWKE using the Gaussian kernel



**Figure A.20.** The EO length (left) and EO score (right) as a function of the LB, in the case of the NWKE.

The comparison of the results for the LLKE and NWKE is presented in Table A.3. The conducted analysis shows that, the prognostic learning method based on LLKE fails to establish length of the EO - due to quickly diverging PB all testing points fall within them and resulting EO lengths are infinite (i.e. undefined).

NWKE method on the other hand performs poorly. This is caused by the boundary bias resulting in horizontal EO while the testing points continues to follow an increasing trend.

**Table A.3** Prognostic learning – comparison of the LLKE and NWKE results, when applied to the data on concentration of the CO<sub>2</sub> in the atmosphere.

Results		LLKE	NWKE
EO	max length	$\infty$ (all tested points fall within the PBs)	4 (for LB=33)
	min length	no finite EO length	2 (for LB=25-31,36-37, 39-40, 43-44, and 49)
	$\infty$	for LB $\geq$ 25 all tested data points fall within the PBs	for LB $\geq$ 50 all tested data points fall within the PBs
	score	$\infty$ (no finite EO score)	finite: 0.287 – 0.517 or $\infty$ (for LB $\geq$ 50)
Residuals	normality	$\varepsilon_t$ normally distributed (Shapiro-Wilk test, $p$ -values $>$ 0.1)	$\varepsilon_t$ normally distributed (Shapiro-Wilk test, $p$ -values $>$ 0.09)
	zero mean	ok (t-test, $p$ -values $>$ 0.2)	ok (t-test, $p$ -values $>$ 0.2)
	correlation	autocorrelation at most at lag 1 or none (ACF, Box-Pierce test)	autocorrelation at lag 1 (at most 2) or none (ACF, Box-Pierce test)

## A.5 Conclusions

Analysis of the performance of the prognostic learning method based on nonparametric regression applied to real-life data sets of anthropogenic CO<sub>2</sub> emissions and atmospheric CO<sub>2</sub> concentrations leads to the following conclusions:

- The use of the LLKE regression performs better than the NWKE. Since it does not exhibit the boundary bias it has smaller prediction errors. This results in longer prediction errors.
- The method based on nonparametric regression easily adapts to the data behaviour, reflecting fluctuations and peaks (for CO<sub>2</sub> emissions data set) while being more stable for data exhibiting regular behaviour (as for CO<sub>2</sub> concentrations data set).
- Autocorrelation of residuals (more pronounced for NKWE method than for LLKE method) has a negative impact on the performance of prognostic learning procedure, i.e. results in shorter EOs.

## Literature

- Altman N.S., *Kernel Smoothing of Data with Correlated Errors*, J. Amer. Statist. Assoc., 1990, 85, 749-759.
- K. De Brabanter, J. De Brabanter, J. A. Bart De Moor, *Kernel Regression in the Presence of Correlated Errors*, J. Machine Learn. Research 12 (2011), 1955-1976.
- Eubank R.L., *Nonparametric regression and spline smoothing*, Marcel Dekker Inc., New York, 1999.
- Fan J., *Design-adaptive Nonparametric Regression*, Journal of the American Statistical Association, Vol. 87, 1992.
- Fan J., Gijbels I., *Local Polynomial Modelling and Its Applications*, Chapman & Hall, London, 1997.
- L. Gajek, M. Kałuszka, *Wnioskowanie statystyczne: modele i metody*, Wydawnictwa Naukowo-Techniczne, Warszawa, 1993.
- Gasser T., Müller H.G., *Estimating Regression Functions and Their Derivatives by the Kernel Method*, Scand. J. Statist., 1984, 11:171-185.
- Green P.J., Silverman B.W., *Nonparametric Regression and Generalized Linear Models: a Roughness Penalty Approach*, Chapman & Hall, London, 1994.
- Györfi L., Kohler M., Krzyżak A., Walk H., *A Distribution-Free Theory of Nonparametric Regression*, Springer, New York, 2002.
- Hart J.D., *Kernel regression estimation with time series errors*, J. Royal Statist. Soc. B. 1991, 53(1):251-259.
- Härdle W., *Applied Nonparametric Regression*, Cambridge University Press, 1990.
- Johnstone I., Silverman B.W., *Wavelet threshold estimators for data with correlated noise*, J. Royal Statist. Soc., B., 1997, 59, 319-351.



- Jarnicka J., *Multivariate kernel density estimation with a parametric support*,  
Opuscula Math. **29**, no. 1 (2009), 41-55.
- Nadaraya E.A., *On estimating regression*, Theory Prob. Appl. 1964, **9**(1): 141-142.
- Nason G.P., *Wavelet shrinkage using cross-validation*, J. Royal Statist. Soc. B, 1996,  
**58**, 463-479.
- Opsomer J., Wang Y., Yang Y., *Nonparametric Regression with Correlated Errors*,  
Statist. Sci. 2001, **16**(2): 134-153.
- Priestley M.B., Chao M.T., *Non-parametric function fitting*, J. Royal Statist. Soc. B,  
1972, **34**: 385-392.
- Rice J., Rosenblatt M., *Integrated mean squared error of a smoothing spline*, J. Approx.  
Theory, 1981, **33**, 353-369.
- Rice J., Rosenblatt M., *Smoothing splines: regression, derivatives and deconvolution*,  
Ann. Statist. 1983, **11**, 141-156.
- Ruppert D., Wand M.P., *Multivariate Locally Weighted Least Squares Regression*, The  
Annals of Statistics, 1994, Vol. 22, p. 1346-1370.
- Silverman B.W., *Density Estimation for Statistics and Data Analysis*, Champan & Hall,  
New York, 1986.
- Simonoff J.S., *Smoothing Methods in Statistics*, Springer, 1996.
- Stone C.J., *The use of polynomial splines and their tensor products in multivariate  
function estimation*, Ann. Statist., 1994, **22**, 118-184.
- Wang Y., *Function estimation via wavelet shrinkage for long-memory data*. Ann.  
Statist. **24**, 1996, 466-484.
- Wasserman L., *All of Nonparametric Statistics*, Springer Texts in Statistics, New York,  
2006.
- Watson G.S., *Smooth regression analysis*, Sankhya Ser. A, 1964, **26**(4): 359-372.

### **Acronyms**

- ACF – autocorrelation function  
 AR - autoregression  
 CLT – central limit theorem  
 CV – cross-validation  
 GMKE – Gasser-Müller kernel estimator  
 KE – kernel estimator  
 k-NNKE – k-nearest neighbour kernel estimator  
 LLKE – Local linear kernel estimator  
 MLE – maximum likelihood estimation  
 MSE – mean squared error  
 NWKE – Nadaraya-Watson kernel estimator  
 PB – prediction bands  
 PCKE – Priestley-Chao kernel estimator  
 PDF – probability density function  
 SE – standard error (i.e. residual standard error)





