


# Active Output Selection Strategies for Multiple Learning Regression Models

Adrian Prochaska<sup>1</sup> <sup>a</sup>, Julien Pillas<sup>1</sup> and Bernard Bäker<sup>2</sup>

<sup>1</sup>Mercedes-Benz AG, 71059 Sindelfingen, Germany

<sup>2</sup>TU Dresden, Chair of Vehicle Mechatronics, 01062 Dresden, Germany  
{adrian.prochaska, julien.pillas}@daimler.com, bernard.baeker@tu-dresden.de

Keywords: Gaussian Processes, Active Learning, Regression, Active Output Selection, Drivability Calibration

Abstract: Active learning shows promise to decrease test bench time for model-based drivability calibration. This paper presents a new strategy for *active output selection*, which suits the needs of calibration tasks. The strategy is actively learning multiple outputs in the same input space. It chooses the output model with the highest cross-validation error as leading. The presented method is applied to three different toy examples with noise in a real world range and to a benchmark dataset. The results are analyzed and compared to other existing strategies. In a best case scenario, the presented strategy is able to decrease the number of points by up to 30 % compared to a sequential space-filling design while outperforming other existing active learning strategies. The results are promising but also show that the algorithm has to be improved to increase robustness for noisy environments. Further research will focus on improving the algorithm and applying it to a real-world example.

## 1 Introduction

Active learning methods – sometimes called *online design of experiments* or *optimal experimental design* – increase the capabilities of algorithms taking part in test design and execution (Cohn, 1996). They reduce the required number of measurements significantly, while guaranteeing adequate model qualities (Klein et al., 2013). However, most methods aim at optimally identifying only one model. In most real-world applications, there are not one but multiple outputs. That leaves the test engineer with a question: Should all models be learned sequentially or simultaneously? And if they learn simultaneously, how to decide which model is the leading one? Drivability calibration applications can be further distinguished from other active learning tasks because

- the goal is to identify all measured outputs equally well and
- pulling one query reveals the values of all outputs of interest.


(Dursun et al., 2015) showed a comparison of a sequential and a round-robin learning strategy for a drivability calibration task. To the authors knowledge, no other publication analyses more sophisticated strategies for multiple learning regression models, which

follow the conditions described above. This paper proposes a new concept of learning strategy, which decides on the leading output by evaluating a cross validation error. This new strategy is compared to other existing strategies. Multiple toy examples are used to create a noisy but reproducible test environment with different complexities. At last, the strategy is also applied to a benchmark dataset.

The paper is structured in six sections. Section 2 of this paper introduces previous works in context of active learning in general and in particular for regression tasks. Section 3 focuses on describing the specialties of active learning in the calibration context. A new active learning task called *active output selection* (AOS) is introduced there. Section 4 describes the analyzed approaches. Furthermore, a new approach for AOS is presented. The approaches are evaluated using a toy example and a benchmark dataset. Experimental details and a discussion of results are shown in section 5. At the end, section 6 concludes the results and presents fields of possible future works.

## 2 Previous work

The field of active learning is a growing branch of the very present machine learning domain. It is

<sup>a</sup>  <https://orcid.org/0000-0003-2707-1266>

also referred to as *optimal experimental design* (Cohn, 1996). (Settles, 2009) shows a broad overview of the current state of the art in this discipline and gives an outlook to multiple possible future work fields. Recent methodological advances in the scientific community mainly focused on classification problems. The main application domains are speech recognition and text information extraction (Settles, 2009).

While regression tasks in the context of active learning have not been as popular, the methodological development is relevant as well. (Sugiyama and Rubens, 2008) propose an approach which actively learns multiple models for the same task and picks the best one to query new points. (Cai et al., 2013) introduced an approach which uses expected model change maximization (EMCM) to improve the active learning progress for gradient boosted decision trees, which was later extended to choose a set of informative queries and to gaussian process regression models (GPs) by (Cai et al., 2017). (Park and Kim, 2020) propose a learning algorithm based on the EMCM, which handles outliers more robustly than before. Those publications focus on new criteria for single output regression models to improve the active learning process. (Zhang et al., 2016) present a learning algorithm for multiple-output gaussian processes (MOGP) which outperforms multiple single-output gaussian processes (SOGP). However, this publication focuses on improving the prediction accuracy of one target output with the help of several correlated auxiliary outputs. The experiments indicate that a global consideration is beneficial.

There were also advances in active learning for automotive calibration tasks for which the identification of multiple process outputs in the same experiment is more relevant to the application. (Klein et al., 2013) applied a design of experiments for hierarchical local model trees (HiLoMoT-DoE), which was presented by (Hartmann and Nelles, 2013), successfully to an engine calibration task. They presented two application examples with two outputs each and five respectively seven inputs. The two outputs were modeled with a sequential strategy, which identifies an output model completely before moving to the next one (Klein et al., 2013).

(Dursun et al., 2015) applied the HiLoMoT-DoE active learning algorithm to a drivability calibration example characterized by multiple static regression tasks with identical input spaces. They analyzed the sequential strategy already shown by (Klein et al., 2013) and compared it to a round-robin strategy, which switches the leading model after each iteration/measurement (Dursun et al., 2015). The authors show that the round-robin strategy outperforms offline methods and the

online sequential strategy in this experiment. It might indicate, that round-robin is preferably used in general, but further experiments are necessary. Since then, no efforts have been made to analyze active learning strategies for multiple outputs.

### 3 Problem definition

The analyses of this paper are motivated by the field of model-based drivability calibration. For this application, an active learning algorithm learns a number of  $M$  different outputs, which are possibly non-correlated. Their models are equally important for succeeding optimizations, so the goal is to identify adequate models for all outputs. The input dimensions of all models are the same. Querying a new instance corresponds to conducting a measurement on powertrain test benches. Therefore, a measurement point is cost-sensitive, which is inherent to active learning problems. Contrary to other applications, every single measurement provides values for all  $M$  outputs<sup>1</sup>. Tasks of simultaneously learning  $M > 1$  process outputs with equal priority and multi-output measurements are not known in the scientific community. In the following, they are referred to as *active output selection* (AOS). All measured outputs contain to some extent noise. The signal-to-noise-ratio  $\text{SNR}_m$  of model  $m$  is the ratio between the range of all measurements  $y_m$  and the standard deviation  $\sigma_N$  of normally distributed noise:  $\text{SNR}_m = \frac{\max(y_m) - \min(y_m)}{\sigma_N}$ . The SNR for drivability criteria lies approximately in a range of (7 .. 100) and can be different for each criterion.

For applications on a test bench, conducting a measurement is timely more expensive than the evaluation of code. This is why the performance of code is not crucial in this context and is only discussed openly in this paper instead of analyzing it systematically.

### 4 Active output selection strategies

This paper analyzes strategies for AOS tasks with  $M > 1$  regression models. In this paper, each of those  $M$  process outputs is modeled with a GP since they handle noise in the range of vehicle calibration tasks very robustly (Tietze, 2015). The leading process output defines the placement of the query in each iteration. A simple maximum variance strategy is deployed as

<sup>1</sup>This is in contrast to e. g. geostatistics, where measuring any individual output, even at the same place (i. e. model inputs), has its own costs (Zhang et al., 2016).

active learning algorithm: A new query  $x_m^*$  in the input space  $\mathbb{X}$  is placed at that point, where the output variance is maximal.

$$x_m^* = \arg \max_{x_m^* \in \mathbb{X}} (\hat{\sigma}_m^2(x^*)) \quad (1)$$

This approach was presented by (MacKay, 1992) for general active learning purposes and applied and evaluated on GPs e. g. by (Seo et al., 2000) or (Pasolli and Melgani, 2011). The implementation of such a learning strategy is straightforward for GPs since the output variance at each input point is directly calculated in the model. Equation (2) and eq. (3) show the calculations of the predicted mean  $\hat{y}$  and output variance  $\hat{\sigma}^2$  of a GP.  $\hat{\mathbf{k}}$  is the vector of covariances  $k(X, \hat{x})$  between the measured training points  $X$  and a single test point  $\hat{x}$ ,  $K = K(X, X)$  are the covariances of  $X$  and  $\mathbf{y}$  contains the observations under noise with variance  $\sigma_n^2$  (Rasmussen and Williams, 2008).

$$\hat{y}(\hat{x}) = \hat{\mathbf{k}}^T (K + \sigma_n^2 I)^{-1} \mathbf{y} \quad (2)$$

$$\hat{\sigma}^2(\hat{x}) = k(\hat{x}, \hat{x}) - \hat{\mathbf{k}}^T (K + \sigma_n^2 I)^{-1} \hat{\mathbf{k}} \quad (3)$$

Depending on the AOS strategy the leadership of the learning process is chosen differently. In the following, three already existing and one new active learning strategy (CVH) as well as a passive sequential design are described. All of those strategies are empirically analyzed in section 5.

**sequential strategy (SQ)** After measuring a set of initial points, the first process output is leading. When the desired model accuracy or the maximum number of points is reached, the next model places measurements and is identified. This procedure is repeated until the criteria for all  $M$  models are fulfilled. The maximum number of measurements for every  $i$ -th model is calculated as follows:

$$p_{m,\max} = \frac{p_{\max} - p_{\text{init}}}{M} \quad (4)$$

An advantage of SQ is, that it identifies only one model each iteration. Depending on the complexity and noise of all process outputs, the order of leading models might influence the performance of this strategy.

**round-robin strategy (RR)** This strategy changes the leading model after each measurement. Models that have reached the desired model quality are not leading any longer. An advantage of round-robin is, that the order of process outputs only has a very small influence on planning the measurements, since the models are switched with every step. Therefore, this strategy should be more suited to handle tasks where the outputs have different complexities. RR also identifies only one model each iteration.

**global strategy (G)** This strategy chooses that query  $x^*$ , which maximizes the sum of output variances.

$$x^* = \arg \max_{x^* \in \mathbb{X}} \left( \sum_{m=1}^M w_m \hat{\sigma}_m^2(x^*) \right) \quad (5)$$

This is a weighted compromise between all models with the weights being  $w_m = 1$ . G identifies all  $M$  models each iteration and is therefore computationally more expensive than SQ and RR.

**CV<sub>10, high</sub> strategy (CVH)** Algorithm 1 shows the CV<sub>10, high</sub> strategy. In the beginning, CVH plans the queries of an initial set of points and conducts the measurements. Afterwards, CVH identifies the models of all outputs in each iteration. Additionally, the model errors are calculated. In this case, a model error is expressed using the normalized root mean squared  $K$ -fold cross-validation-error  $CV_K$  with  $K = 10$ . Equation (6) shows the general form of  $CV_K$  of the  $m$ -th model with the predictions  $\hat{y}_{m,i}^{-\kappa(i)}$  of the  $m$ -th model being identified without measurements of set  $\kappa: \{1, \dots, N\} \mapsto \{1, \dots, K\}$

$$CV_{K,m} = \sqrt{\frac{\sum_{i=1}^N (y_{m,i} - \hat{y}_{m,i}^{-\kappa(i)})^2}{\max(y_m) - \min(y_m)}} \quad (6)$$

The usage of another accuracy or error criterion is possible, but  $CV_K$  is well-comparable between models. For stability reasons,  $CV_{10}$  is filtered with a digital moving average filter, which reduces the influence of fluctuations during runtime. In every following iteration, the output with the highest model error is leading the learning process. This output is assumed to benefit the most from being in leadership of learning.

---

**Algorithm 1** CVH active output selection strategy.

---

- 1: **repeat**
  - 2:   **if** no initial points have been carried out **then**
  - 3:     plan queries of initial points
  - 4:   **else**
  - 5:     find model with the highest filtered cross-validation error
  - 6:     calculate next query
  - 7:     conduct measurements on planned queries
  - 8:   **for all** models **do**
  - 9:     update model
  - 10:    assess cross-validation error
  - 11:    filter the cross-validation error
  - 12: **until** maximum number of points or desired model quality is reached
- 

Using  $CV_{10}$  obliges identifying each of the  $M$  models for 10 times in each iteration. Compared to the

other strategies, this results in a higher computational effort than the previously presented methods. However this argument is not crucial for drivability calibration tasks, as the measurements itself take a lot longer than calculating the succeeding query. Since  $CV_{10}$  also increases with higher noise, this strategy might be prone to one process output with significantly larger noise than the others. Its model cannot reach a model error as low as those of the other outputs; after reaching the minimum possible  $CV_{10}$  the model will not benefit from actively planning points anymore. To the authors' knowledge, this strategy has not been presented or analyzed in any other publication.

**sequential space-filling strategy (passive, SF)** Instead of a random sequential set of queries, the authors choose including a passive but sequential design as baseline method to verify the benefits of those AOS strategies. This kind of design is derived from an s-optimal (space-filling) experimental design, which is preferred over a random set of points in drivability calibration applications. A sequential method additionally enables a fair comparison on whether an active learning strategy is truly beneficial over a passive one. After an initial set of measurements, the next point is always placed in a maximin-way which maximizes the minimum Mahalanobis-distances  $d_{\min}^2$  between a huge set of candidate points and the already measured points.

$$d_{\min}(\hat{x}) = \min \| \hat{x} - X \| \quad (7)$$

$$x^* = \arg \max_{x^* \in \mathbb{X}} (d_{\min}(\hat{x})) \quad (8)$$

Because points are planned sequentially, this design does not exactly result in a test design which is optimally space-filling for the current number of points. However, it is an easy way to be close to this optimality during a sequential design where the number of points is not predefined.

Due to the characteristics of the AOS strategies described above, the following hypothesis are tested with the experiments:

1. The non-heuristic CVH is in many cases beneficial but also has drawbacks concerning high noise-induced generalization error.
2. RR is robust in all use cases but can be outperformed by CVH.
3. The active learning strategies perform significantly better than a SF.

<sup>2</sup>The Mahalanobis distance for uncorrelated data in a range between 0 and 1 is identical to the Euclidian distance.

## 5 Experiments

The application of the presented learning strategies in the field of drivability calibration is designed for the use on a test bench. However, typical static drivability criteria have a signal-to-noise-ratio SNR of (7 .. 100).

Figure 1 demonstrates the influence of noise in that range in a toy example. It shows the  $NRMSE_{\text{val}}$ -values over the number of measurements  $n_{\text{Meas}}$  of 3 learning procedures of a space-filling design for the same example. The normalized root mean squared

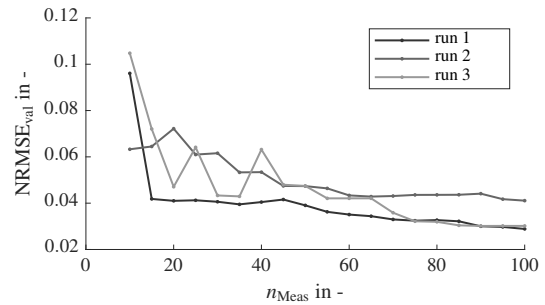


Figure 1:  $NRMSE_{\text{val}}$  for 3 runs of a generic process with a space-filling design and different noise observations. This figure shows the influence of noise on  $NRMSE_{\text{val}}$  for a space-filling strategy.

error of the validation points  $NRMSE_{\text{val}}$  is calculated according to eq. (9).

$$NRMSE_{\text{val},m} = \sqrt{\frac{\sum_{i=1}^{N_{\text{val}}} (y_{m,i,\text{val}} - \hat{y}_{m,i,\text{val}})^2}{\max(y_m) - \min(y_m)}} \quad (9)$$

The only difference between those runs are the locations of initial points and the noise observations, which are characterized by  $SNR = 12$ . This amount of noise leads to completely different generalization errors. Those results propose that to compare learning strategies for noisy environments, an experiment has to be repeated multiple times. For comparing the results of different learning strategies, not only the mean but also the standard deviation of  $NRMSE_{\text{val}}$  is relevant. The conduction of numerous tests for example on a powertrain or engine test bench is time- and cost-expensive and therefore not practicable. Furthermore, the test conditions would be slightly different every time which makes a direct comparison difficult. The authors' goal in this publication is to compare the described learning strategies in a way that is reproducible and representative for applications in vehicle calibration. That is why three different toy examples are chosen for comparison here. They use the possibility of computer generated noise to be reset to the same starting point of the random-number-generator.

## 5.1 Toy Examples

Every toy example includes three different generic processes, which are analytical multi-dimensional sigmoid or polynomial models generated by a random function generator presented in (Belz and Nelles, 2015). Their outputs are overlaid with normally distributed noise to simulate the measurement inaccuracy. Each process has two input dimensions. The setup used for comparisons consists of multiple runs. One run is understood as one single observation for those comparisons. The overlaid noise of one run is the same for each learning strategy. This is a condition that real world tests cannot fulfill – or at least with untenable effort. However, it increases comparability: Each strategy has the same initial conditions for each run. Furthermore inside one run, the chronological order of examined process outputs and the randomly created, initial points are the same for each strategy.

Another advantage of a comparison with analytical models is the knowledge about the real values from the underlying process without influences of noise. All identified models are validated with 121 gridded validation points and the  $\text{NRMSE}_{\text{val}}$  is calculated (see eq. (9)).

The hypotheses stated earlier are analyzed with 3 different toy examples. For every toy example three different generic processes are chosen. This is a realistic experience value for the number of process outputs to be modeled. The characteristics of those different setups are shown in table 1. Every learning strategy is tested for 50 times in each setup.

Table 1: Specification of the three setups. The symbol + is indicating high complexity or noise.

| setup | model type of each toy example output | complexity | noise |
|-------|---------------------------------------|------------|-------|
| 1     | sigmoid                               | +          | +     |
|       | sigmoid                               | +          | +     |
|       | sigmoid                               | +          | +     |
| 2     | polynomial                            | o          | +     |
|       | polynomial                            | o          | +     |
|       | multiple sigmoids combined with steps | ++         | +     |
| 3     | sigmoid                               | +          | ++    |
|       | sigmoid                               | +          | +     |
|       | multiple sigmoids combined with steps | ++         | +     |

For better understandability of the results the squared sum of the  $\text{NRMSE}_{\text{val}}$  of each model is used

for comparison.

$$\text{NRMSE}_{\text{val},\Sigma} = \sqrt{\sum_{m=1}^M (\text{NRMSE}_{\text{val},m})^2} \quad (10)$$

Figure 2 shows the results of setup 1. When  $n_{\text{Meas}} \lesssim 30$ , the performance of the analyzed strategies are all very similar. From that point on, the new CVH strategy performs significantly better than all other strategies. This includes a low standard deviation  $\sigma_{\text{NRMSE}_{\text{val},\Sigma}}$ . The low  $\sigma_{\text{NRMSE}_{\text{val},\Sigma}}$  shows that there is not a big difference between different runs and therefore stands for the high robustness of this method.

The mean of CVH has the lowest end value. This difference is significant compared to the results of RR and CVH, which have very similar means. Assuming that we want to quit our experiment at any time, CVH should be preferred. RR and SQ have similar mean performance, but RR inherits a lower standard deviation and is therefore more robust. Compared to the  $\text{NRMSE}_{\text{val},\Sigma}$  of SF after  $n_{\text{Meas}} = 100$ , RR and SQ reduce the number of points by 5 % whereas CVH reduces the number of points even further by 15 %. The mean performance of G is comparable to SF. However, the variance of G is higher, which ranks the performance of SF over G.

Figure 3 shows that the difference between active learning strategies is especially high in setup 2 compared to the SF design. In contrast to the previous setup, one process output has a much higher complexity than the other ones. This combination shows the benefits of CVH very clearly. When  $n_{\text{Meas}} \gtrsim 45$ , CVH performs better than all other strategies. G performs significantly worse than the other strategies. This is unlike the third hypothesis in section 4. Depending on the application and the chosen strategy, it is not always beneficial to use active learning. After  $n_{\text{Meas}} = 100$ , RR and SQ have no significant difference in results. However, the standard deviation of SQ is higher during the runs, especially for  $n_{\text{Meas}} \lesssim 45$ . The authors assume that the influence of the process output order plays a role in that. Compared to that, RR shows a more robust behavior. RR reaches the end value of SF at  $n_{\text{Meas}} = 75$ . The standard deviations of RR and CVH are both on similar levels. The CVH performs significantly better than all other strategies in this scenario. It reduces the number of measurement points of a SF by 30 % and reaches the end value of RR after  $n_{\text{Meas}} = 80$ .

Figure 4 shows the results of setup 3. Those results match the expectations that CVH performs worse in an environment where a process output with high noise and a complex one exist. For  $n_{\text{Meas}} \lesssim 80$ , the overall performance of CVH is not much worse than other strategies. However, all other strategies have a better

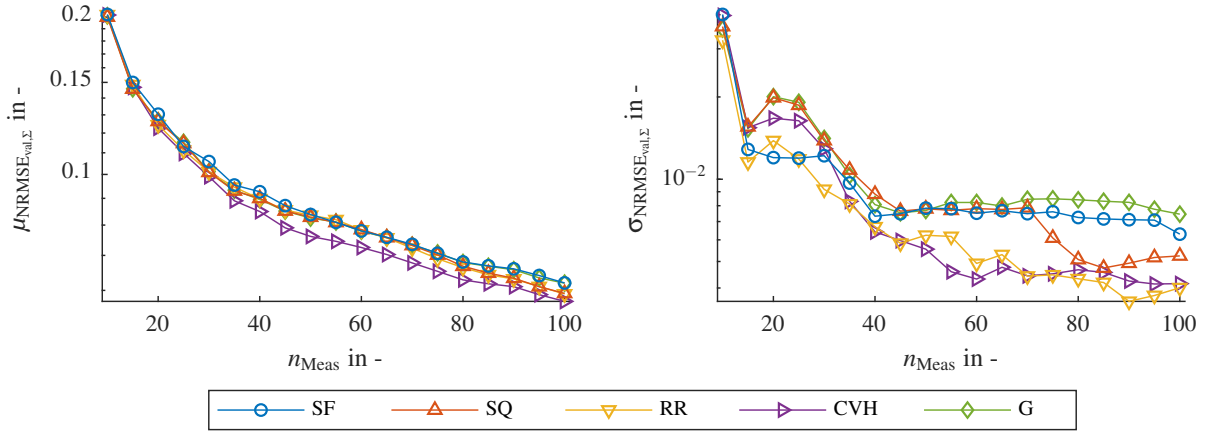


Figure 2: Mean  $\mu$  and standard deviation  $\sigma$  of the  $\text{NRMSE}_{\text{val},\Sigma}$  of setup 1 over the number of measurements  $n_{\text{Meas}}$ .

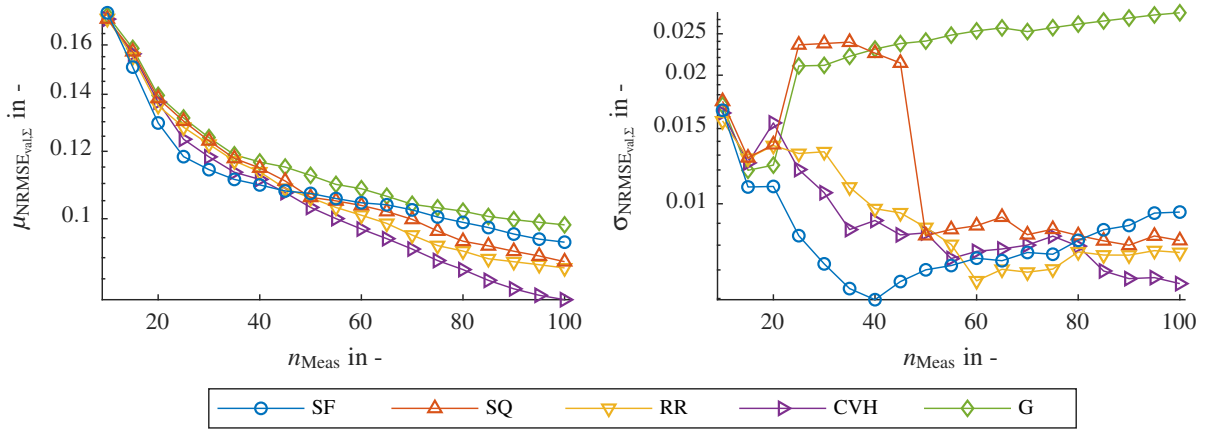


Figure 3: Mean  $\mu$  and standard deviation  $\sigma$  of the  $\text{NRMSE}_{\text{val},\Sigma}$  of setup 2 over the number of measurements  $n_{\text{Meas}}$ .

end value. In this setup, where there is one very noisy and one very complex process output, RR and SQ perform best. They outperform SF at  $n_{\text{Meas}} = 85$  and CVH already at  $n_{\text{Meas}} = 75$ . Compared to SQ, RR is rather robust in the beginning and in the end.

The experiments confirm the hypotheses stated in section 4. Only the third hypothesis turned out not to be true in all cases: SF can outperform active learning in some cases.

## 5.2 Benchmark dataset

As stated in section 3, applications of active learning in the domain of drivability calibration are rather unique concerning the conditions and goals of other existing tasks. That is, why there is no benchmark dataset that fully suits the needs of an example. However we wanted to demonstrate the practical use of such a learning strategy. This is why the jura dataset (Goovaerts, 1997), which is actually a dataset from the domain of geostatistics, is used as a benchmark dataset here. This dataset contains the concentration of 7 heavy metal

concentrations at different locations in the Swiss Jura. It is the best fitting dataset, which is also used to evaluate the learning algorithms in (Zhang et al., 2016). In contrast to that publication, we set the goal to model every of the three chosen outputs equally well. Three concentrations (Ni, Cd, Zn) are modeled as a function of the locations during every test run. The results of the AOS strategies presented in section 4 are averaged over 50 test runs.

Figure 5 shows the results of the benchmark dataset. In the beginning, SF performs significantly worse than CVH. After  $n_{\text{Meas}} \approx 100$ , there are no significant differences between those two strategies. CVH and SF outperform G, SQ and RR in the end. Throughout all measurements however, the mean of CVH is the lowest. Since experiments on a test bench might be stopped after any fixed number of measurements, the results indicate that CVH is preferably used, although it is not significantly better than SF regarding the final model accuracy.

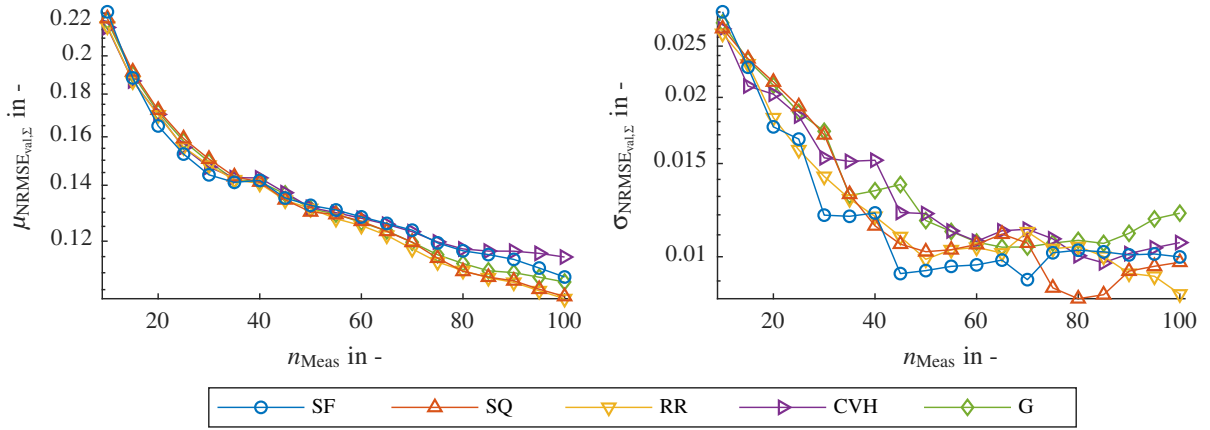


Figure 4: Mean  $\mu$  and standard deviation  $\sigma$  of the  $\text{NRMSE}_{\text{val},\Sigma}$  of setup 3 over the number of measurements  $n_{\text{Meas}}$ .

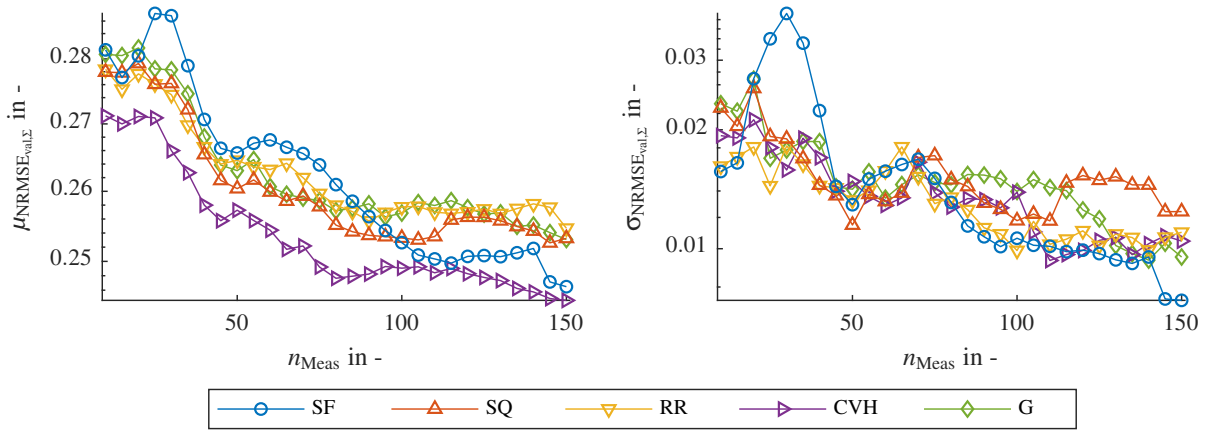


Figure 5: Mean  $\mu$  and standard deviation  $\sigma$  of the  $\text{NRMSE}_{\text{val},\Sigma}$  of the Jura Dataset over the number of measurements  $n_{\text{Meas}}$ .

## 6 Conclusion

In this paper, *active output selection*, a new task for active learning is introduced. It is characterized by identifying models of multiple process outputs with the same input dimensions. We present a new strategy (CVH) to define the leading models in an active output selection setup. The decision is based on the cross-validation errors of the identified models. The paper thoroughly analyzes the advantages and disadvantages of said strategy. The process outputs are identified using gaussian processes (GP). A simple maximum variance algorithm is chosen as active learning strategy for each individual output. The strategy is analyzed on different toy examples, which include noisy generic process outputs. The results of CVH are compared against three existing learning strategies: round-robin, sequential and global. Furthermore a passive, sequential space-filling strategy is chosen as baseline for the active learning strategies.

The results show that the presented strategy is preferably used in most real-world setups. The per-

formance and robustness are good compared to other multi-output strategies. Compared to the baseline strategy, CVH saves up to 30% of the measurements. In this setup, which has one process output with higher complexity, CVH outperforms any other active output selection strategy. In the particular case of a setup with one output with high noise and one output with high complexity, other strategies perform better than CVH. The consideration of the estimated generalization error could further improve the performance of CVH especially for those setups and will therefore be context to further investigation.

The results of a benchmark dataset confirm the good performance of CVH in the toy examples. Due to the lack of a more suitable public benchmark dataset however, a geostatistics example was chosen. A public benchmark dataset from the field of drivability calibration would facilitate comparisons and simplify further work on the subject. Another future research will apply the presented strategy to setups with different numbers of input dimensions. Moreover the application of those strategies on different modeling and single model ac-

tive learning approaches is promising.

## REFERENCES

- Belz, J. and Nelles, O. (2015). Proposal for a function generator and extrapolation analysis. In *2015 International Symposium on Innovations in Intelligent Systems and Applications (INISTA)*. IEEE.
- Cai, W., Zhang, M., and Zhang, Y. (2017). Batch Mode Active Learning for Regression With Expected Model Change. *IEEE Transactions on Neural Networks and Learning Systems*.
- Cai, W., Zhang, Y., and Zhou, J. (2013). Maximizing Expected Model Change for Active Learning in Regression. In *2013 IEEE 13th International Conference on Data Mining*, Dallas, TX, USA. IEEE.
- Cohn, D. A. (1996). Neural Network Exploration Using Optimal Experiment Design. *Neural Networks*.
- Dursun, Y., Kirschbaum, F., Jakobi, R., Gebhardt, A., Goos, J.-C., and Rinderknecht, S. (2015). Ansatz zur adaptiven Versuchsplanung für die Längsdynamikapplikation von Fahrzeugen auf Prüfständen. In *6. Internationales Symposium für Entwicklungsmethodik*.
- Goovaerts, P. (1997). *Geostatistics for natural resources evaluation*. Oxford University Press on Demand.
- Hartmann, B. and Nelles, O. (2013). Adaptive Test Planning for the Calibration of Combustion Engines – Methodology. *Design of Experiments (DoE) in Engine Development*.
- Klein, P., Kirschbaum, F., Hartmann, B., Bogachik, Y., and Nelles, O. (2013). Adaptive Test Planning for the Calibration of Combustion Engines – Application. *Design of Experiments (DoE) in Engine Development*.
- MacKay, D. J. C. (1992). Information-Based Objective Functions for Active Data Selection. *Neural Computation*.
- Park, S. H. and Kim, S. B. (2020). Robust expected model change for active learning in regression. *Applied Intelligence*.
- Pasolli, E. and Melgani, F. (2011). Gaussian process regression within an active learning scheme. In *2011 IEEE International Geoscience and Remote Sensing Symposium*, Vancouver, BC, Canada. IEEE.
- Rasmussen, C. E. and Williams, C. K. I. (2008). *Gaussian processes for machine learning*. MIT Press, Cambridge, Mass., 3. print edition.
- Seo, S., Wallat, M., Graepel, T., and Obermayer, K. (2000). Gaussian Process Regression: Active Data Selection and Test Point Rejection. In *Mustererkennung 2000*. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Settles, B. (2009). Active Learning Literature Survey. Technical report, University of Wisconsin–Madison.
- Sugiyama, M. and Rubens, N. (2008). Active Learning with Model Selection in Linear Regression. In *Proceedings of the 2008 SIAM International Conference on Data Mining*. Society for Industrial and Applied Mathematics.
- Tietze, N. (2015). *Model-based Calibration of Engine Control Units Using Gaussian Process Regression*. PhD thesis, Technische Universität Darmstadt, Darmstadt.
- Zhang, Y., Hoang, T. N., Low, K. H., and Kankanhalli, M. (2016). Near-optimal active learning of multi-output Gaussian processes. In *Thirtieth AAAI Conference on Artificial Intelligence*.