

Online Machine Learning Techniques for Predicting Operator Performance

MASTER'S THESIS BY
AHMET ANIL PALA

Thesis Advisor: MAX HEIMEL

Thesis Supervisor: Prof. Dr. VOLKER MARKL

September 3rd, 2015



Fachgebiet Datenbanksysteme und Informationsmanagement
Technische Universität Berlin

<http://www.dima.tu-berlin.de/>

- Motivation
- Problem
- Existing Methods
- Approach
- Evaluation
- Summary

- Tuning a hardware-oblivious database involves two decision-making tasks
 - Device Selection
 - Algorithm Selection

Given:

Operator Input Data:



Operator Features:



Select:

cpu1

cpu2

gpu

fpga

⋮

Alg1

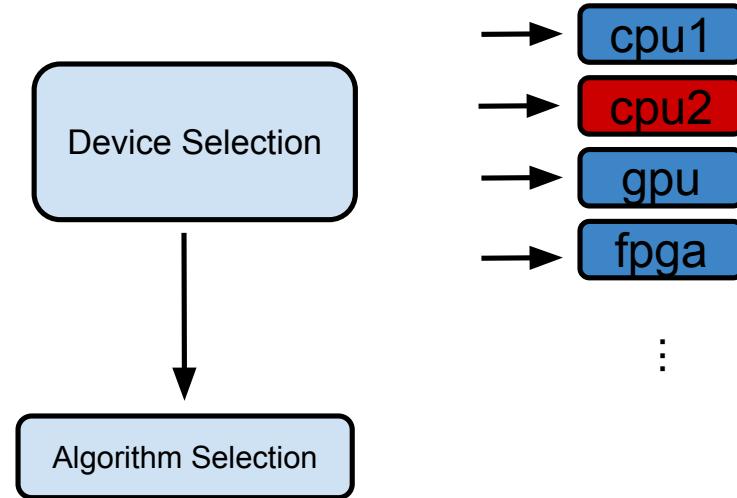
Alg2

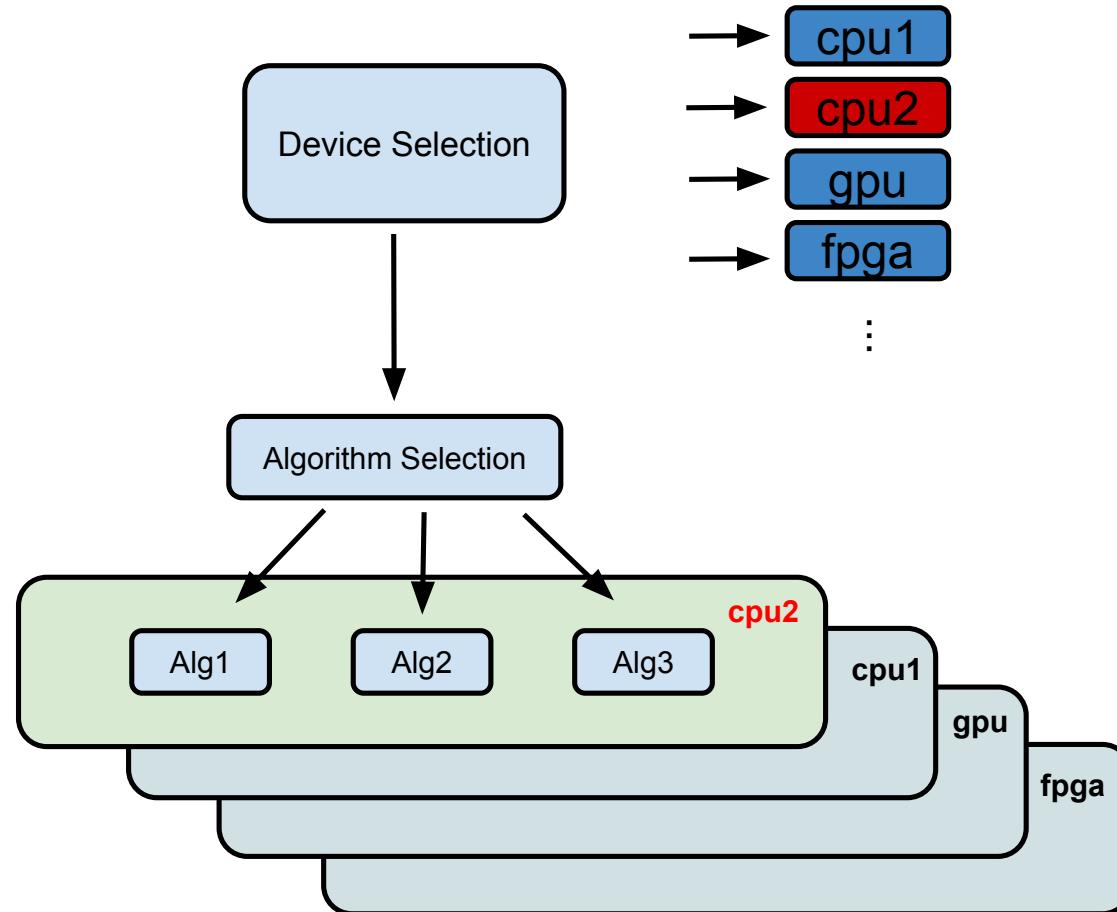
Alg3

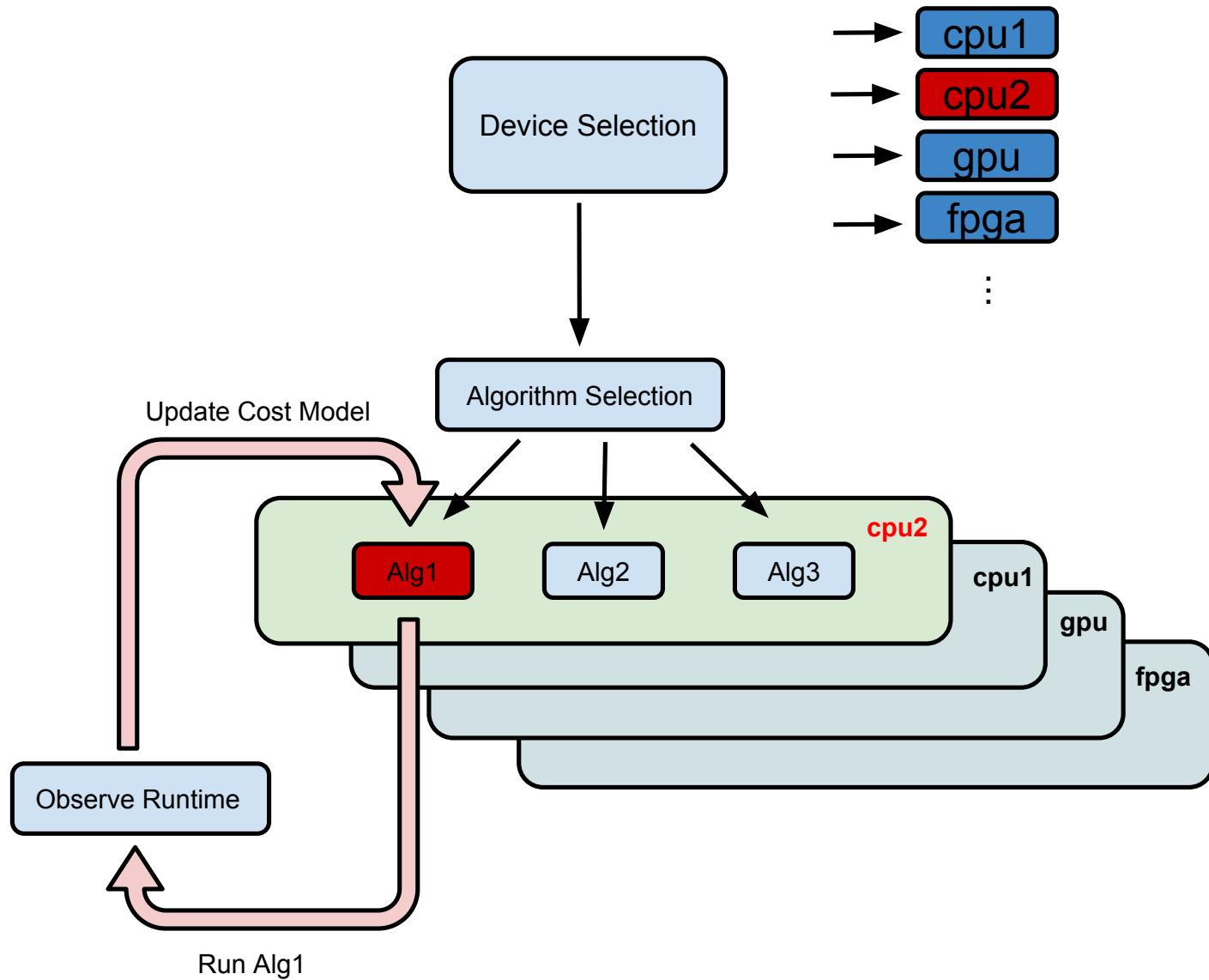
hash-based join

sort-merge join

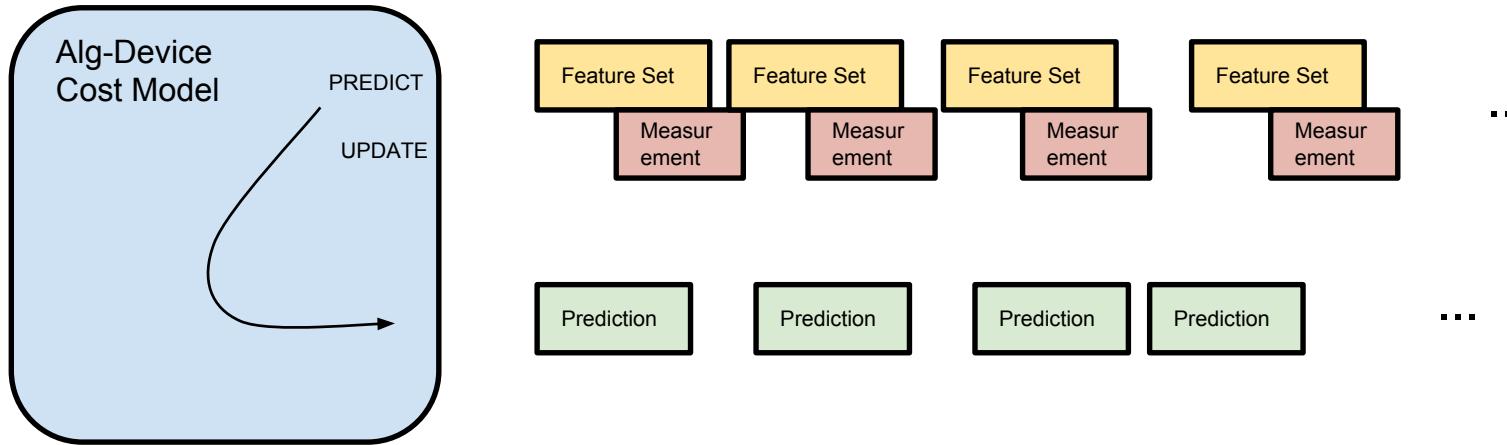
nested-loop join

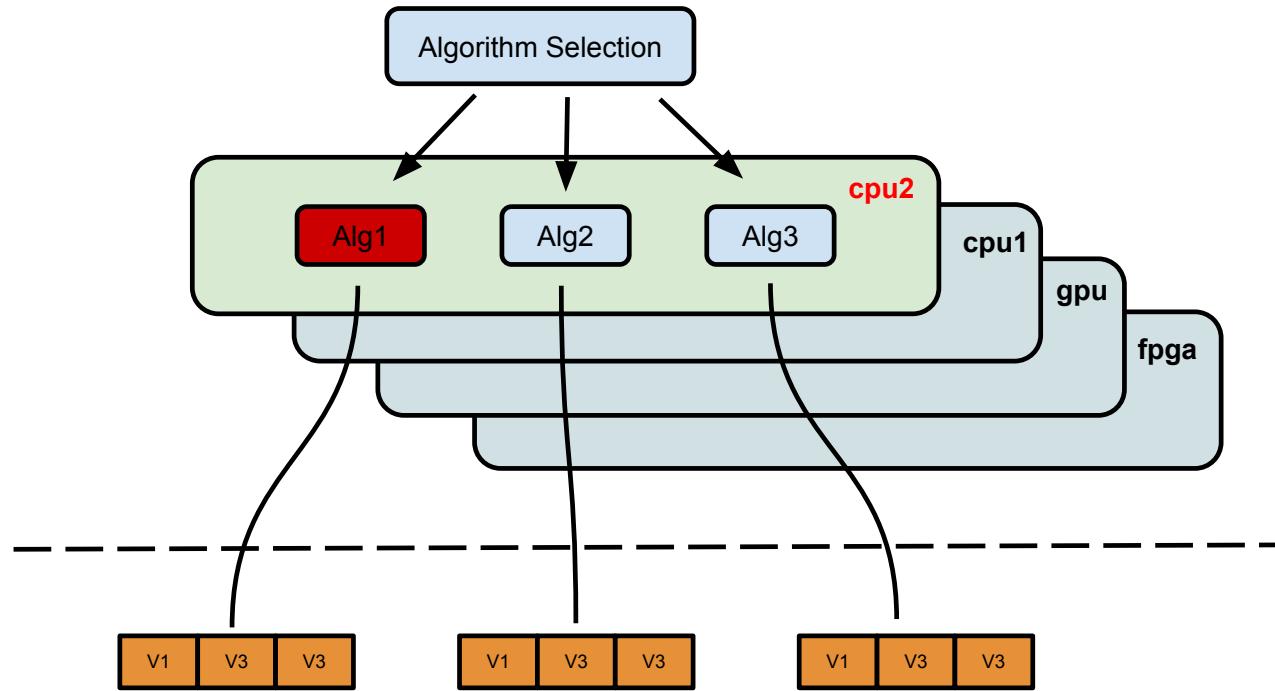






Learning Scenario





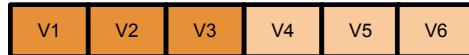
Variant Library for Alg1:

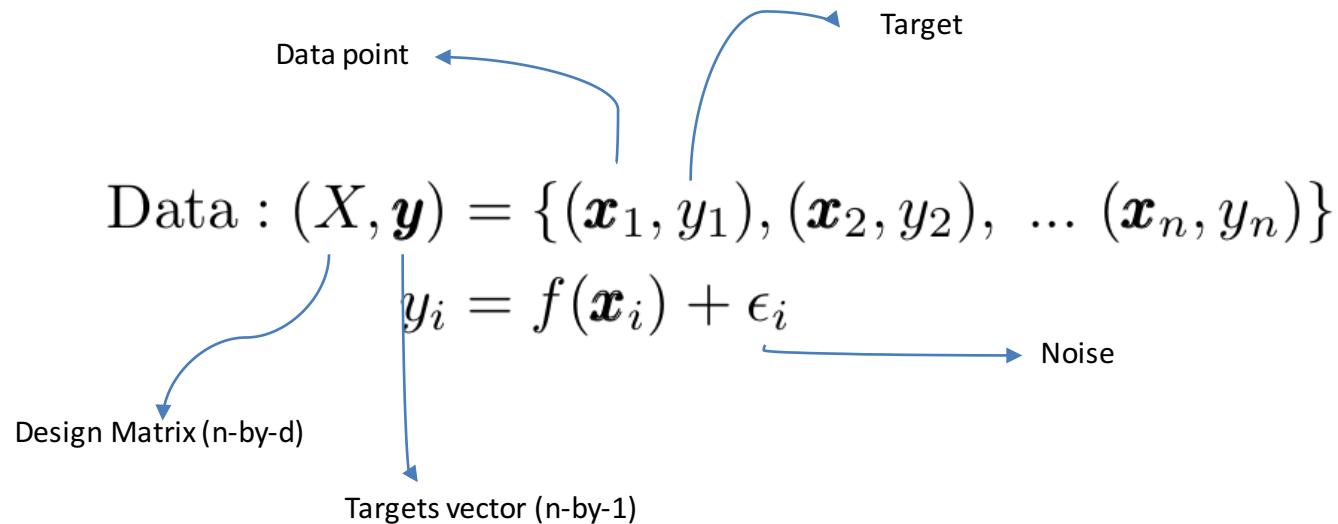


Variant Library for Alg2:



Variant Library for Alg3:





Find $y_{n+1} | (X, \mathbf{y})$

- High Accuracy
- Robust to Abrupt Concept Changes
- Robust to Measurement Noise
- Provide Estimation Bounds
- Constant Space
- Efficiency ($\leq 1\text{ms/item}$)

- Parametric models
 - MAP-method
 - MLE-method
- Non-parametric models
 - Gaussian Process Regression
 - Kernel Regression

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots$$

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) = w_1 x_1^2 + w_2 x_1 x_2 + w_3 x_1 x_3 + w_4 x_2 x_3 + \dots$$

Aim is to find \mathbf{w} !

$$\boldsymbol{w}_{MLE} = (XX^\top)^{-1} X \boldsymbol{y}$$

$$\boldsymbol{w}_{MAP} = (XX^\top + \sigma_y^2 \Sigma_w^{-2})^{-1} X \boldsymbol{y}$$

regularization term

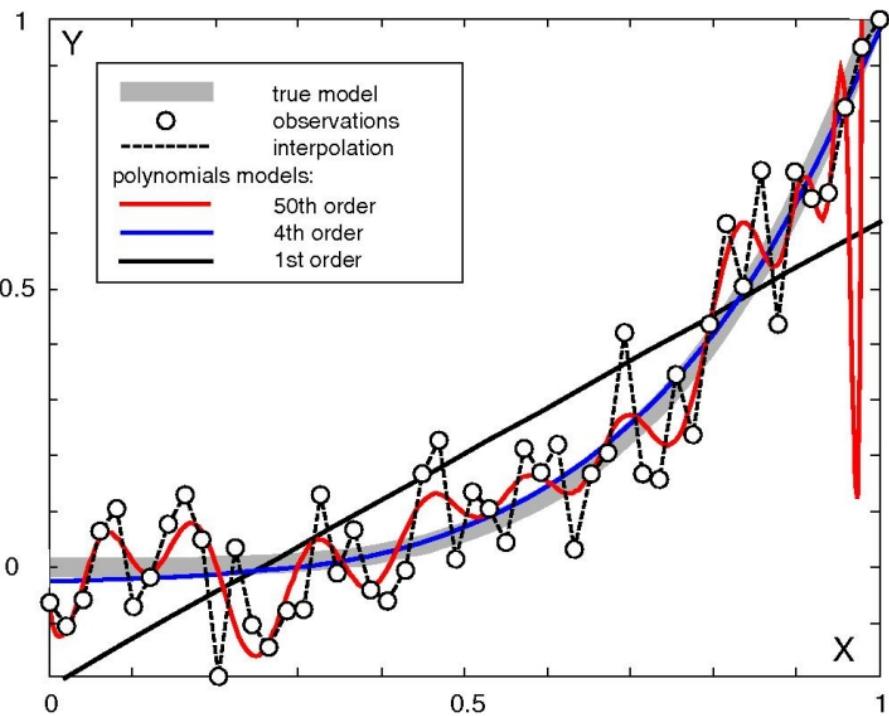
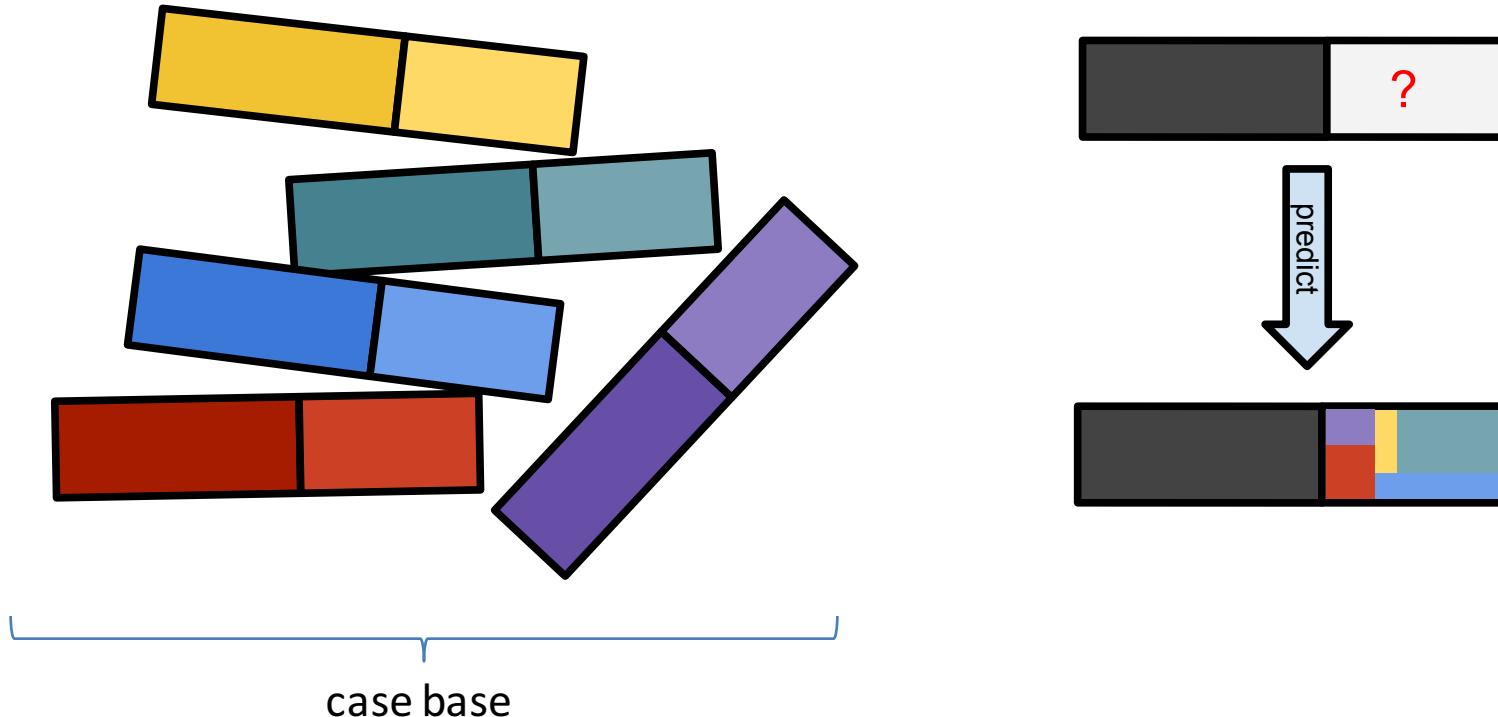
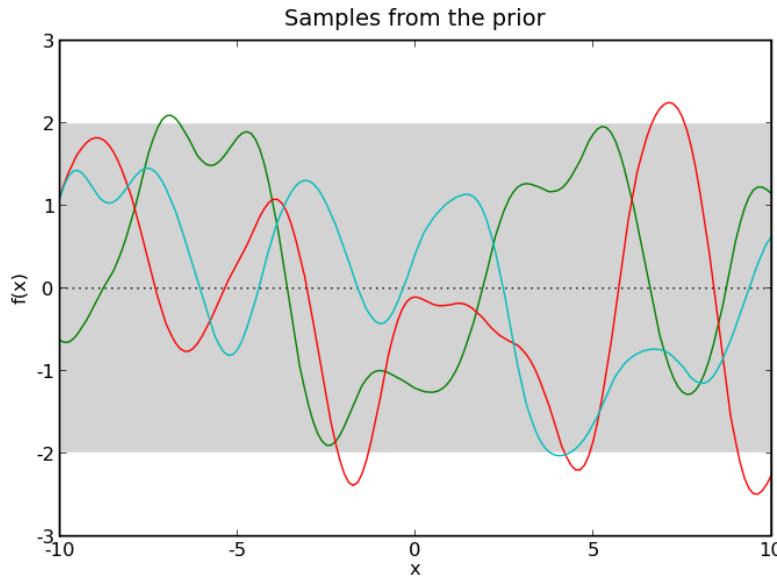


Image taken from <http://pvermees.andropov.org/noble/disc/discPaper/node19.html>

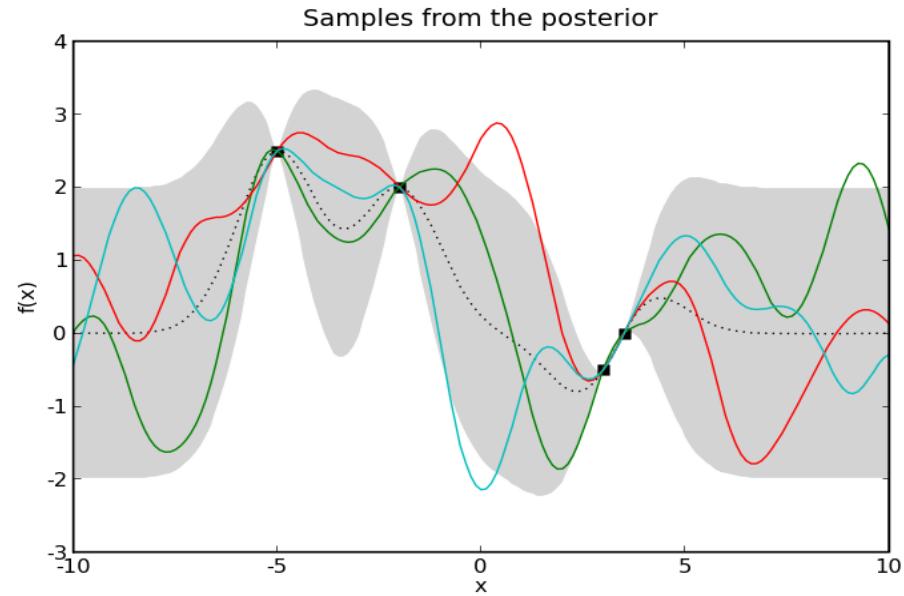
- Data-driven
- Accumulative



- All targets as a single point in Gaussian distribution
- Considering the functions that can map data points to the observed targets in the function space

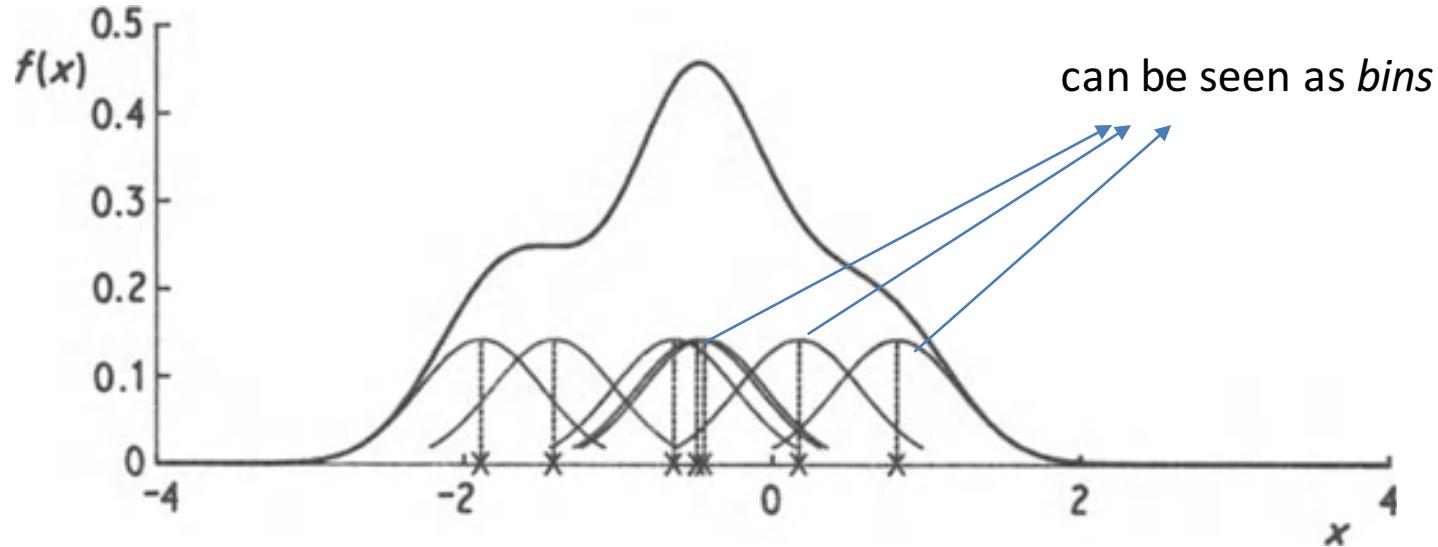


$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$



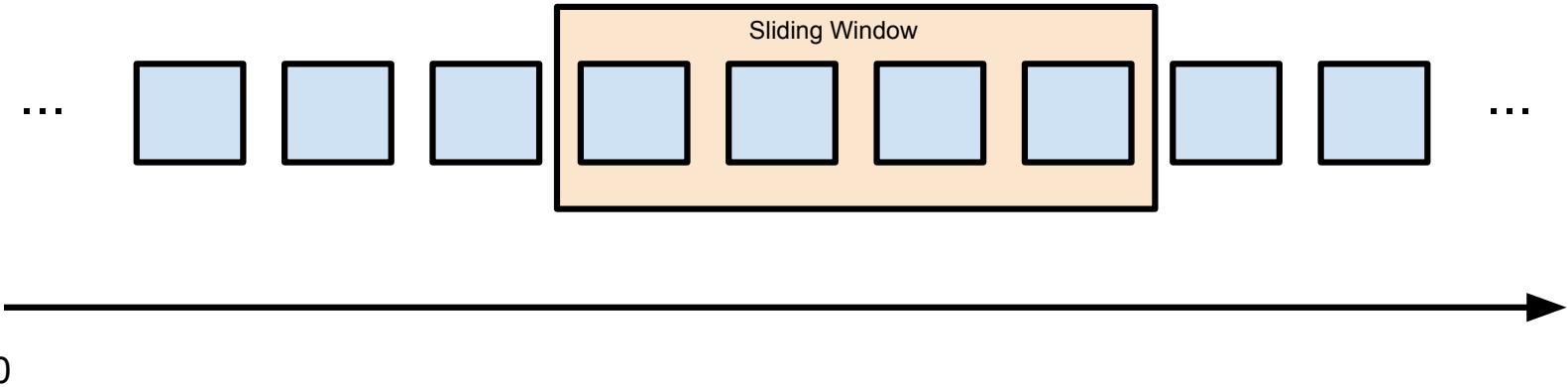
Images taken from <http://pythonhosted.org/infpy/gps.html>

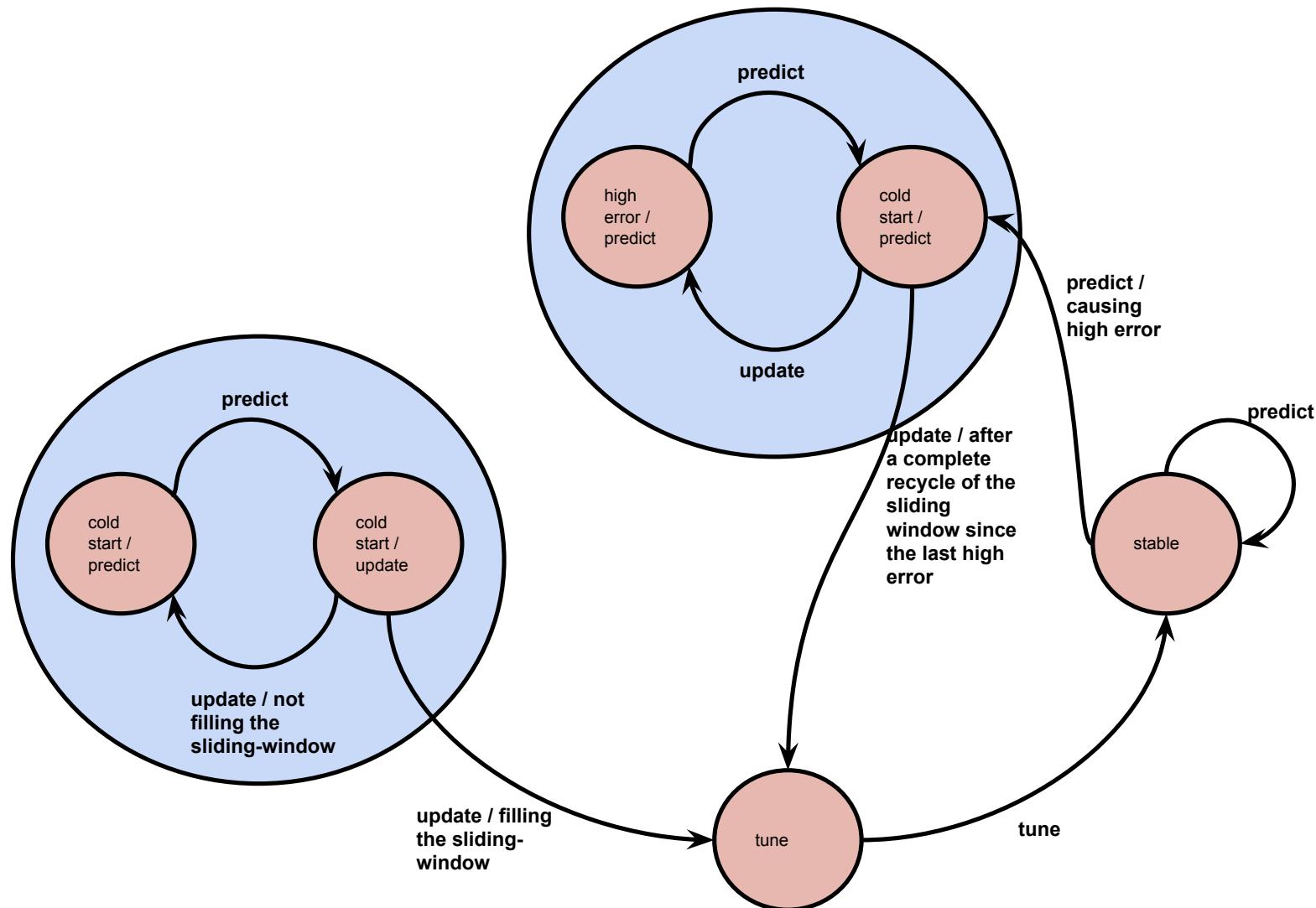
- Based on kernel-density estimation
- Generalization of k-NN method

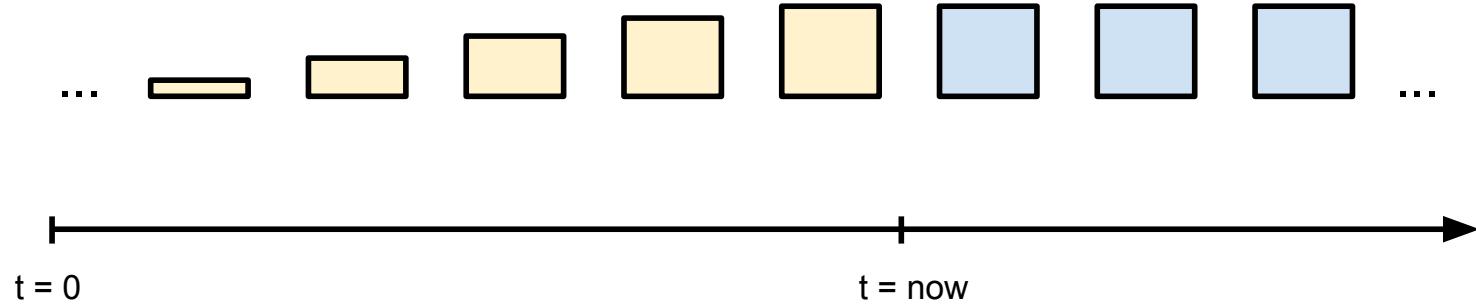


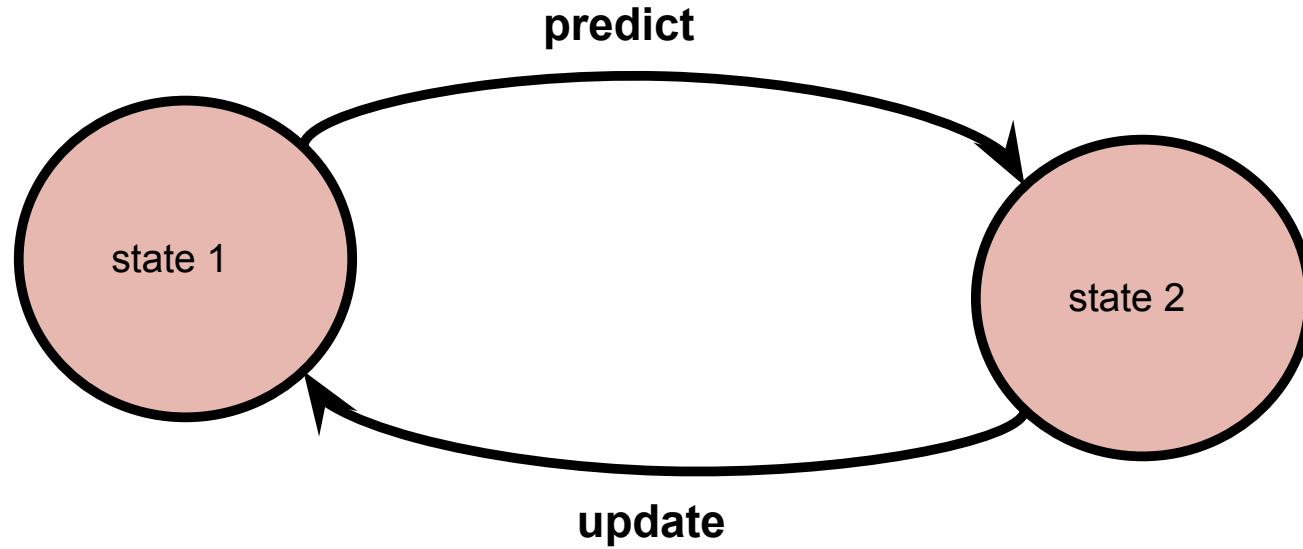
Graph taken from Silverman, Bernard W. *Density estimation for statistics and data analysis*.

- Incremental update mechanism
- Removal of obsolete examples (downdate)







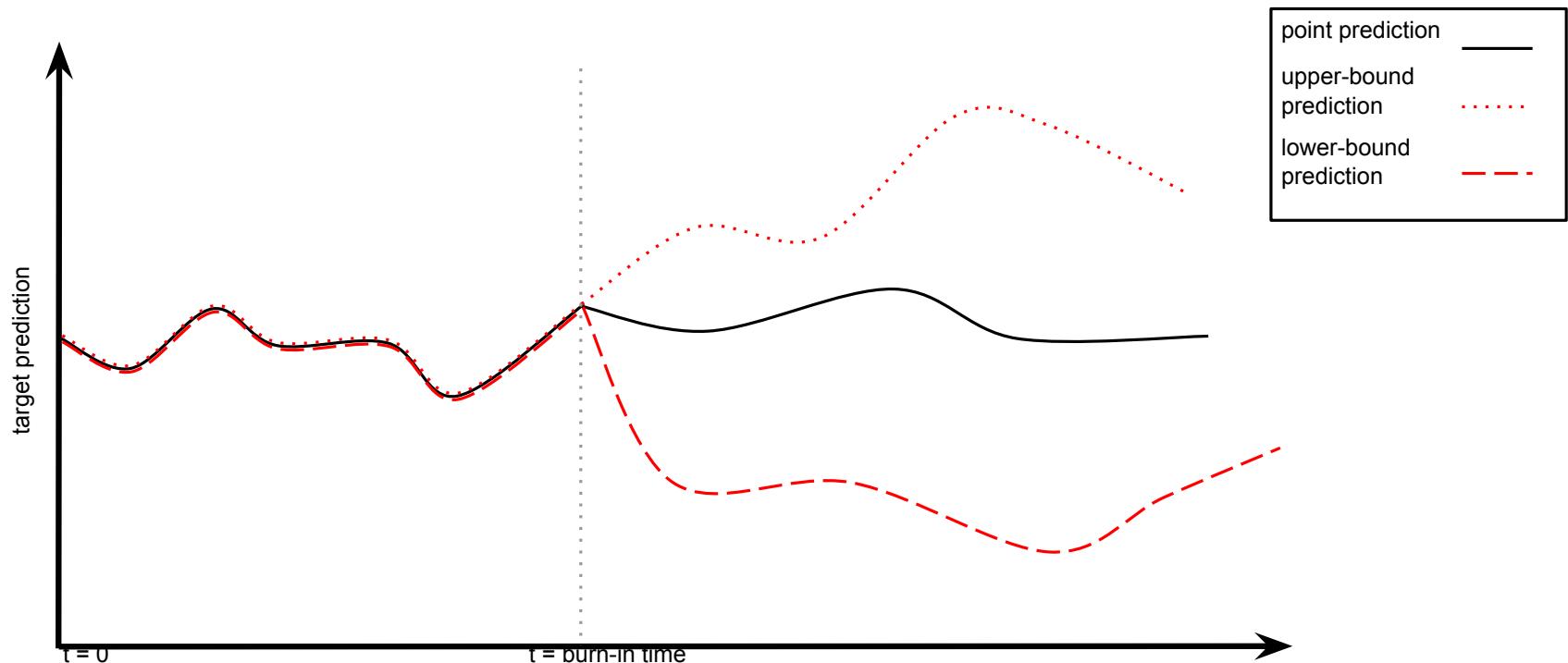


$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i$$

$$y_i \in [\mathbf{w}^\top \mathbf{x}_i \pm z_{1-\frac{\alpha}{2}} \sqrt{s^2 \mathbf{x}_i^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{x}_i + s^2}]$$

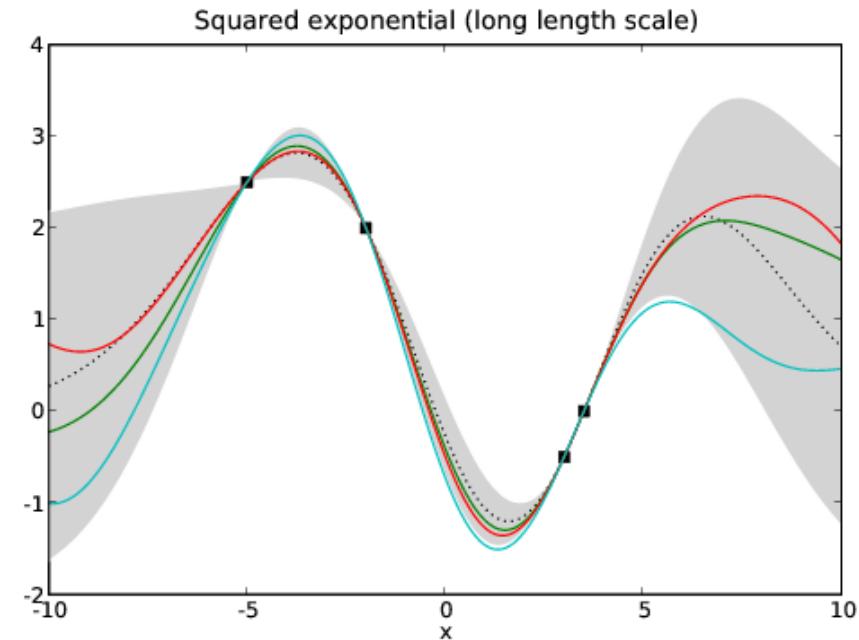
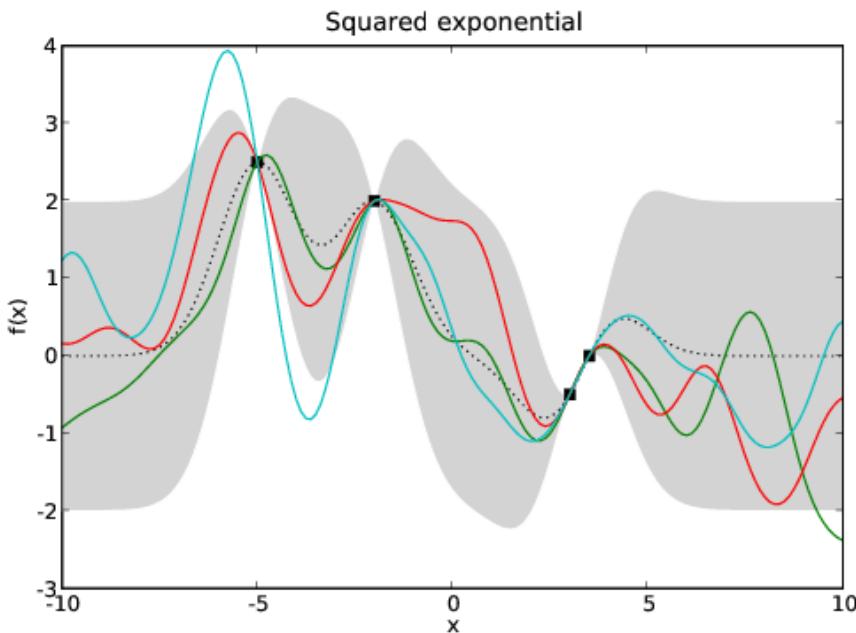

instantaneous error

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i$$



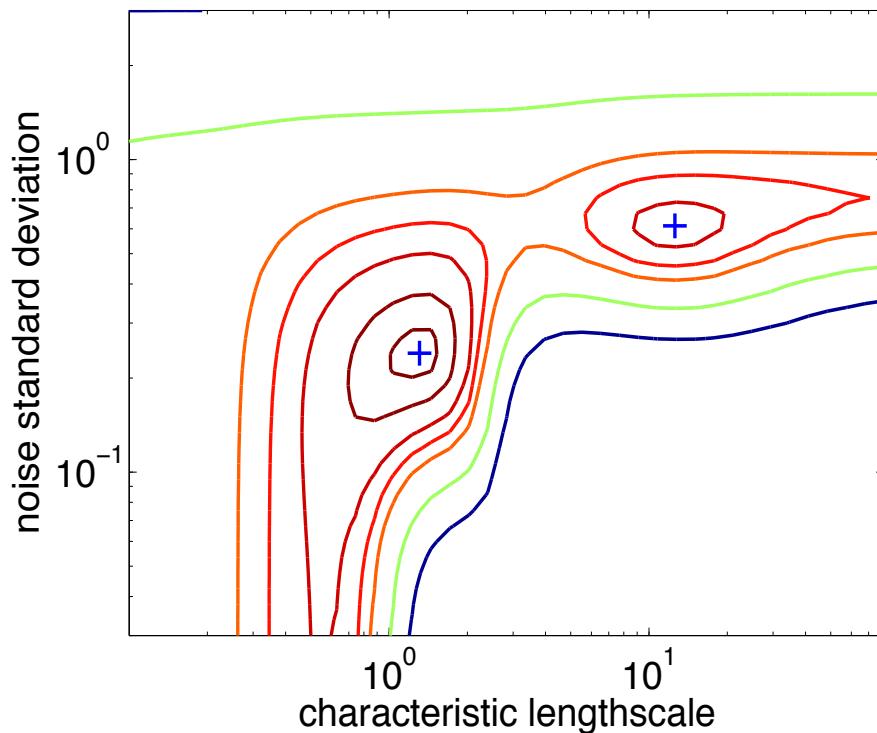
$$k(\mathbf{x}_p, \mathbf{x}_q) = \sigma_w^2 \exp\left(-\frac{1}{2}(\mathbf{x}_p - \mathbf{x}_q)^\top D^{-2}(\mathbf{x}_p - \mathbf{x}_q)\right) + \sigma_y^2 \delta_{pq}$$

signal variance
characteristic lengthscale matrix
noise variance



Images taken from <http://pythonhosted.org/infpy/gps.html>

$$\log(p(\mathbf{y}|X, \theta)) = -\frac{n}{2} \underbrace{\log(2\pi)}_{\text{constant}} - \frac{1}{2} \underbrace{\log(|K(X, X)|)}_{\text{complexity penalty}} - \frac{1}{2} \underbrace{(\mathbf{y} - \mathbf{M})^\top K(X, X) (\mathbf{y} - \mathbf{M})}_{\text{data-fit}}$$



Graph taken from Williams, Christopher KI, and Carl Edward Rasmussen. "Gaussian processes for machine learning." *the MIT Press* 2.3 (2006): 4.

$$\hat{f}(\mathbf{x}) = \frac{1}{n|H|} \sum_{i=1}^n k((\mathbf{x}_i - \mathbf{x})H^{-1})$$

bandwidth matrix

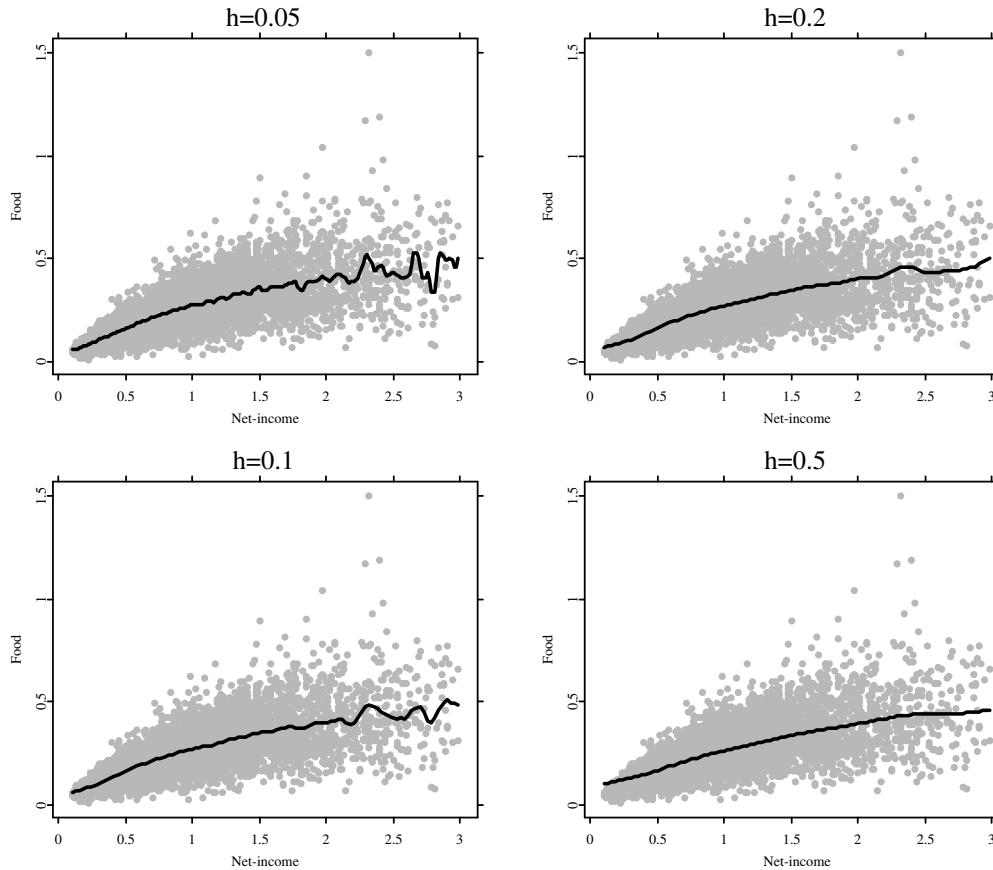


Image taken from <http://pvermees.andropov.org/noble/disc/discPaper/node19.html>

- 52 Online learners & 12 Batch learners
- Rigorous testing on synthetic streams

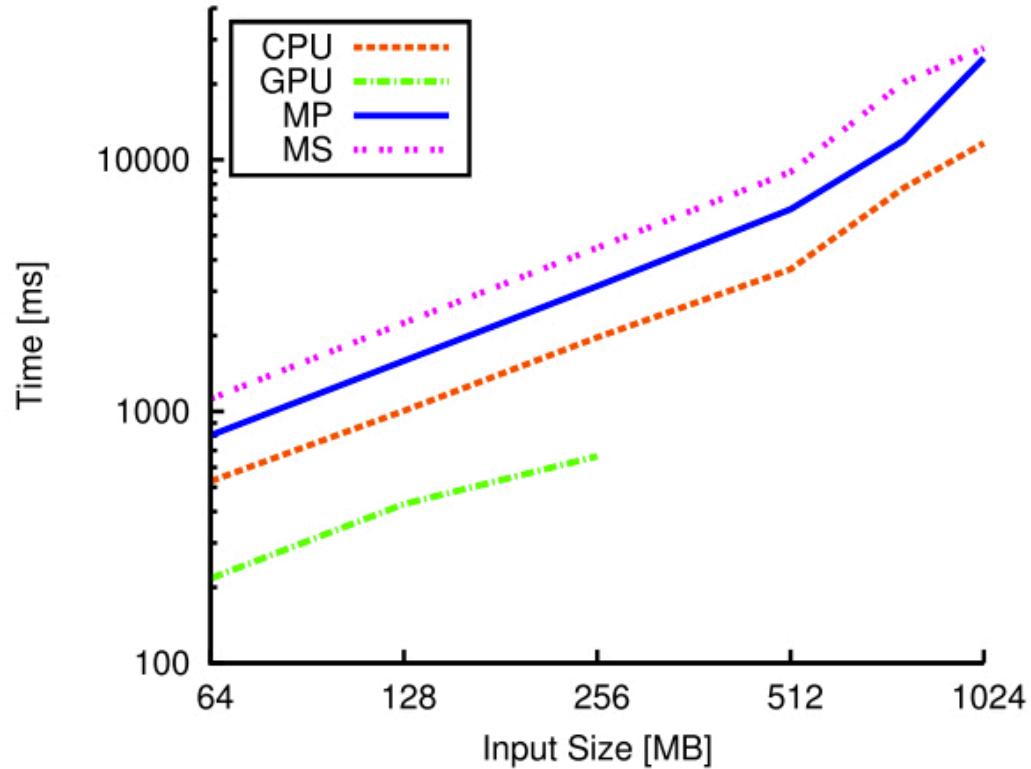


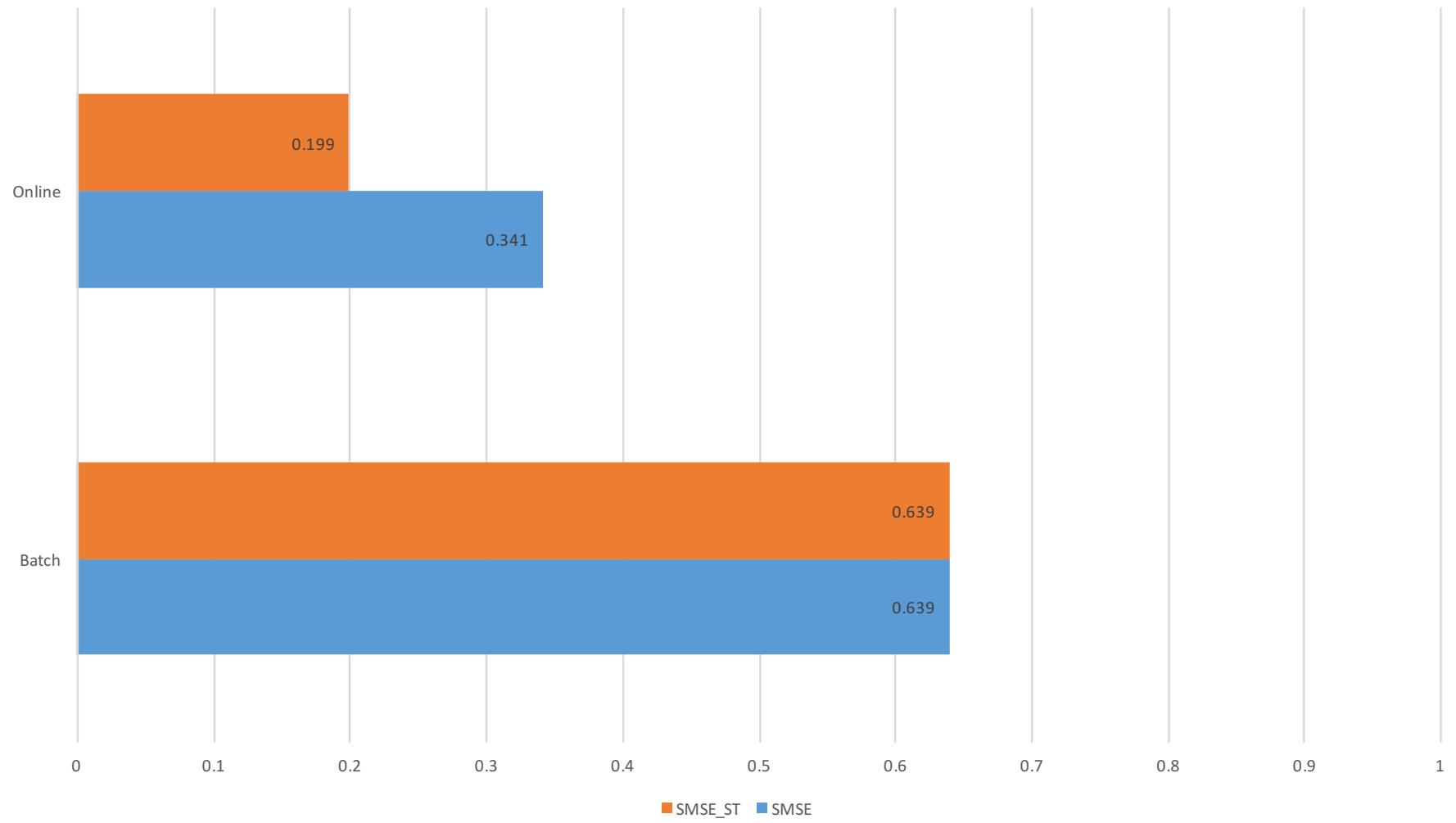
Image taken from Heimel, Max, et al. "Hardware-oblivious parallelism for in-memory column-stores." *Proceedings of the VLDB Endowment* 6.9 (2013): 709-720.

- 52 Online learners & 12 Batch learners
- Rigorous testing on synthetic streams
- Evaluation dimensions: accuracy, prediction bounds quality, time efficiency

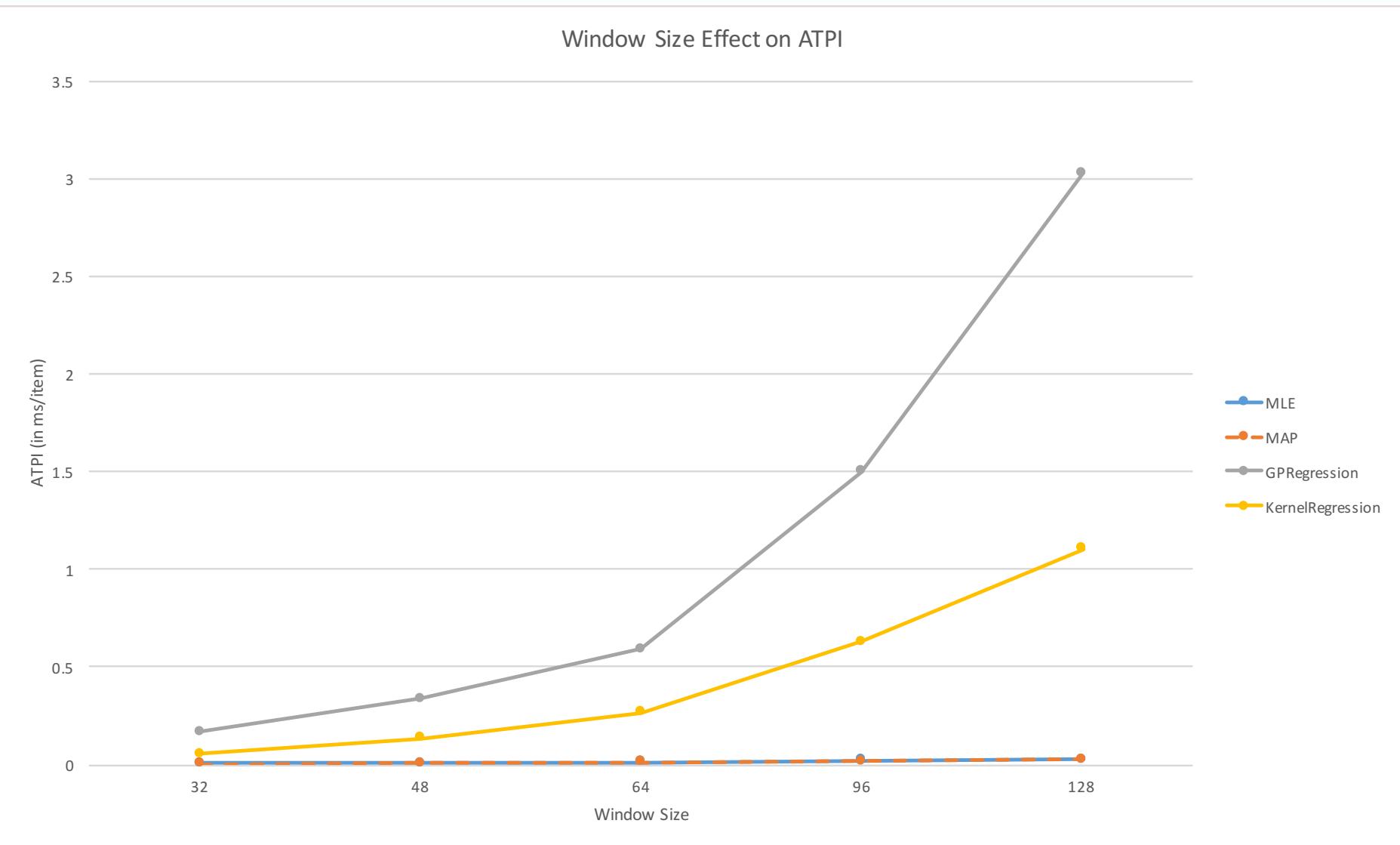
Accuracy	Prediction Bounds	Time Efficiency
RMSE	ICR	ATPI (ms/item)
RMSE_ST	AIW	Max. Data Rate (item/ms)
SMSE	SAIW	
SMSE_ST		

- 52 Online learners & 12 Batch learners
- Rigorous testing on synthetic streams
- Evaluation dimensions: accuracy, prediction bounds quality, time efficiency
- Prequential method and windowed prequential method

Online Learners vs Batch Learners (In)accuracy Comparison on Drifting Data

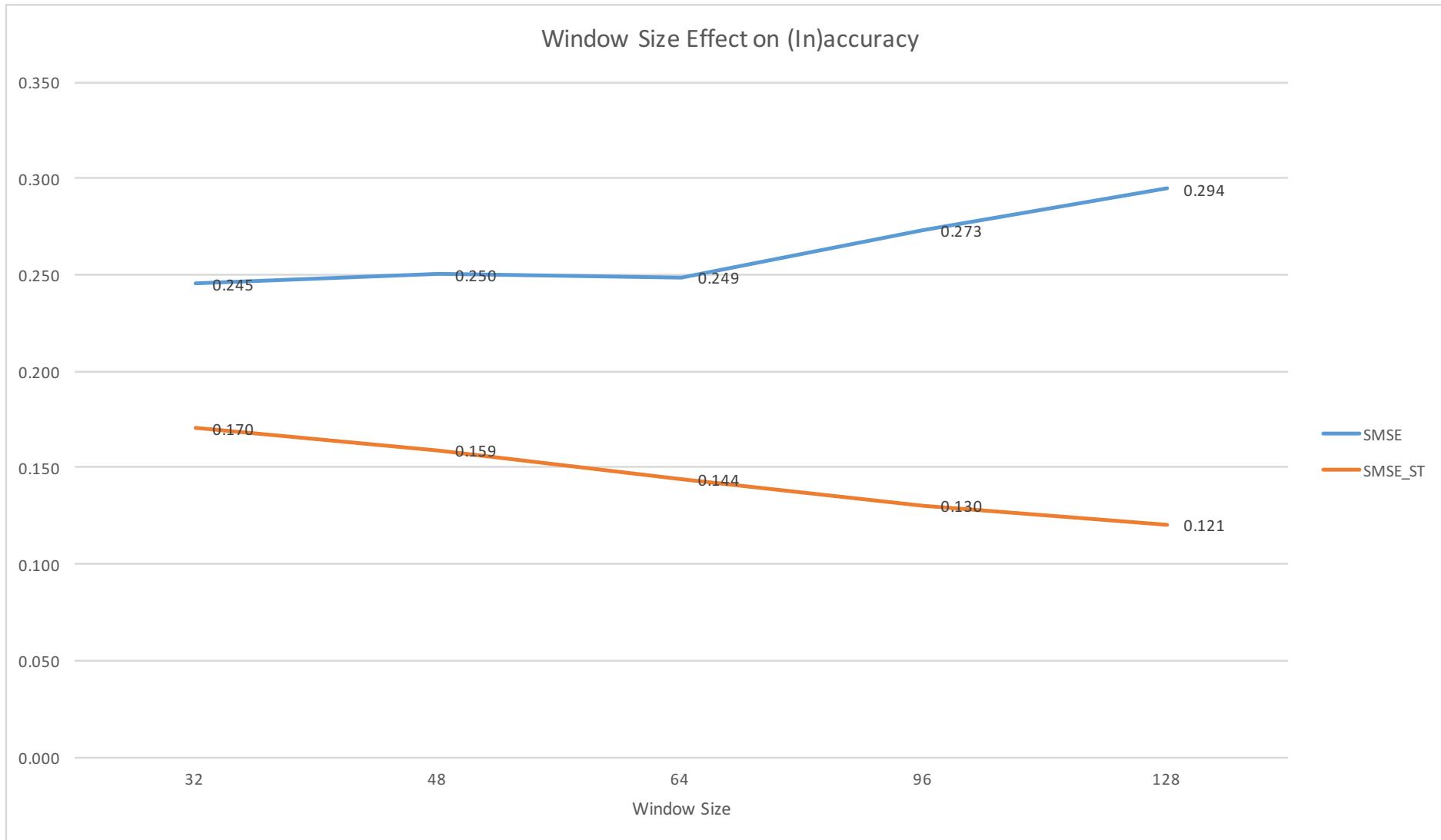


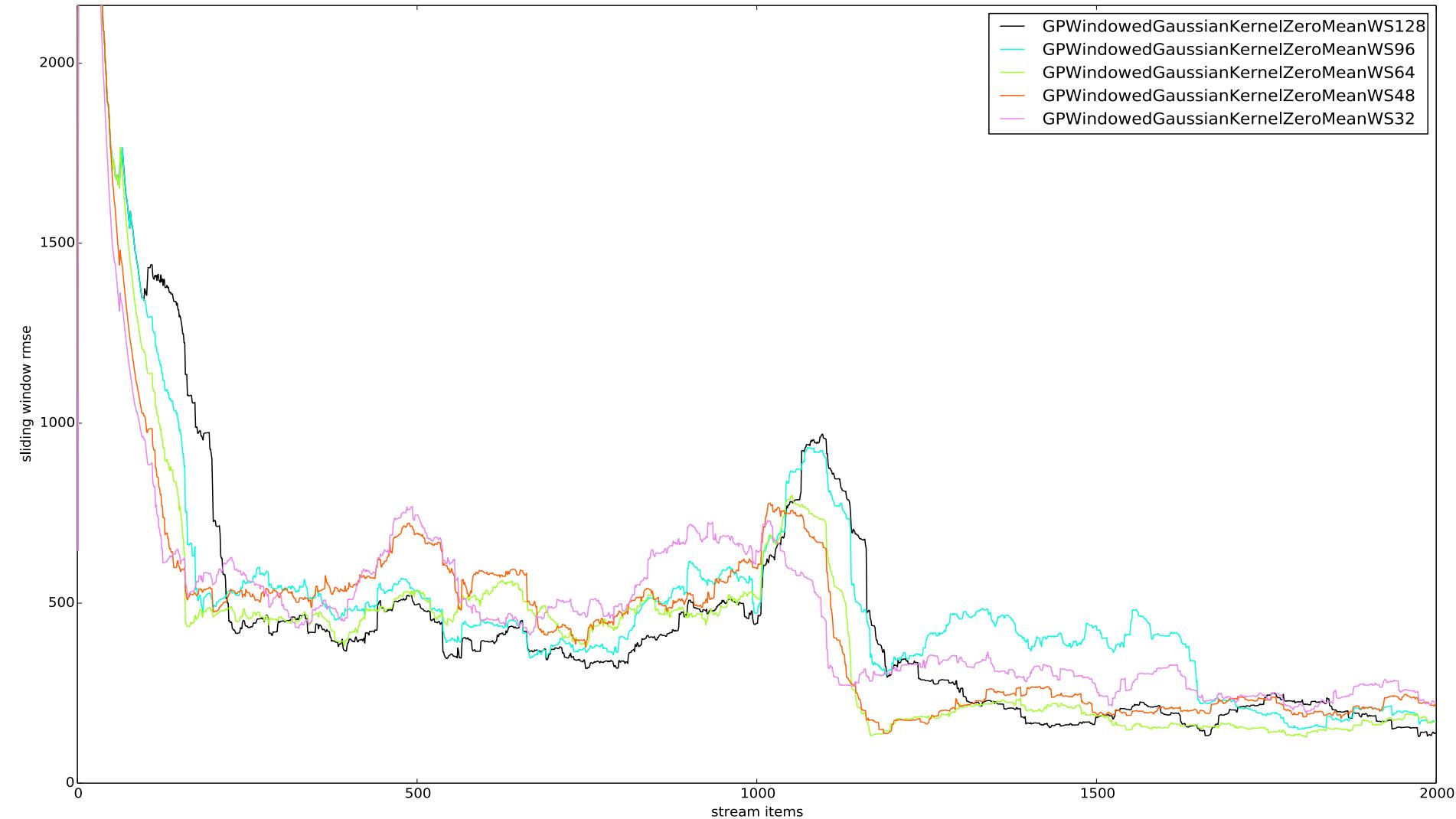
Window size on efficiency



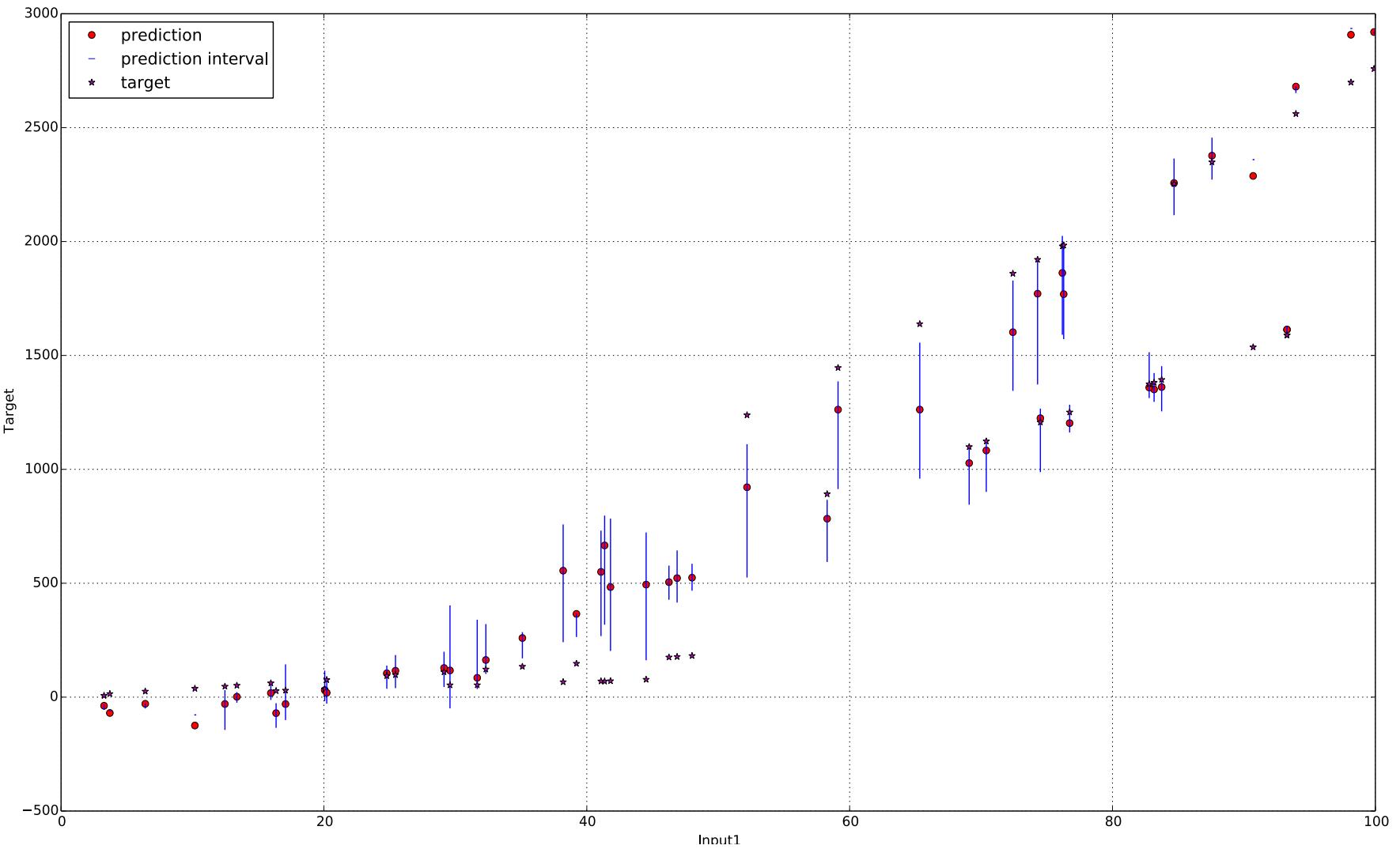
Window size on efficiency



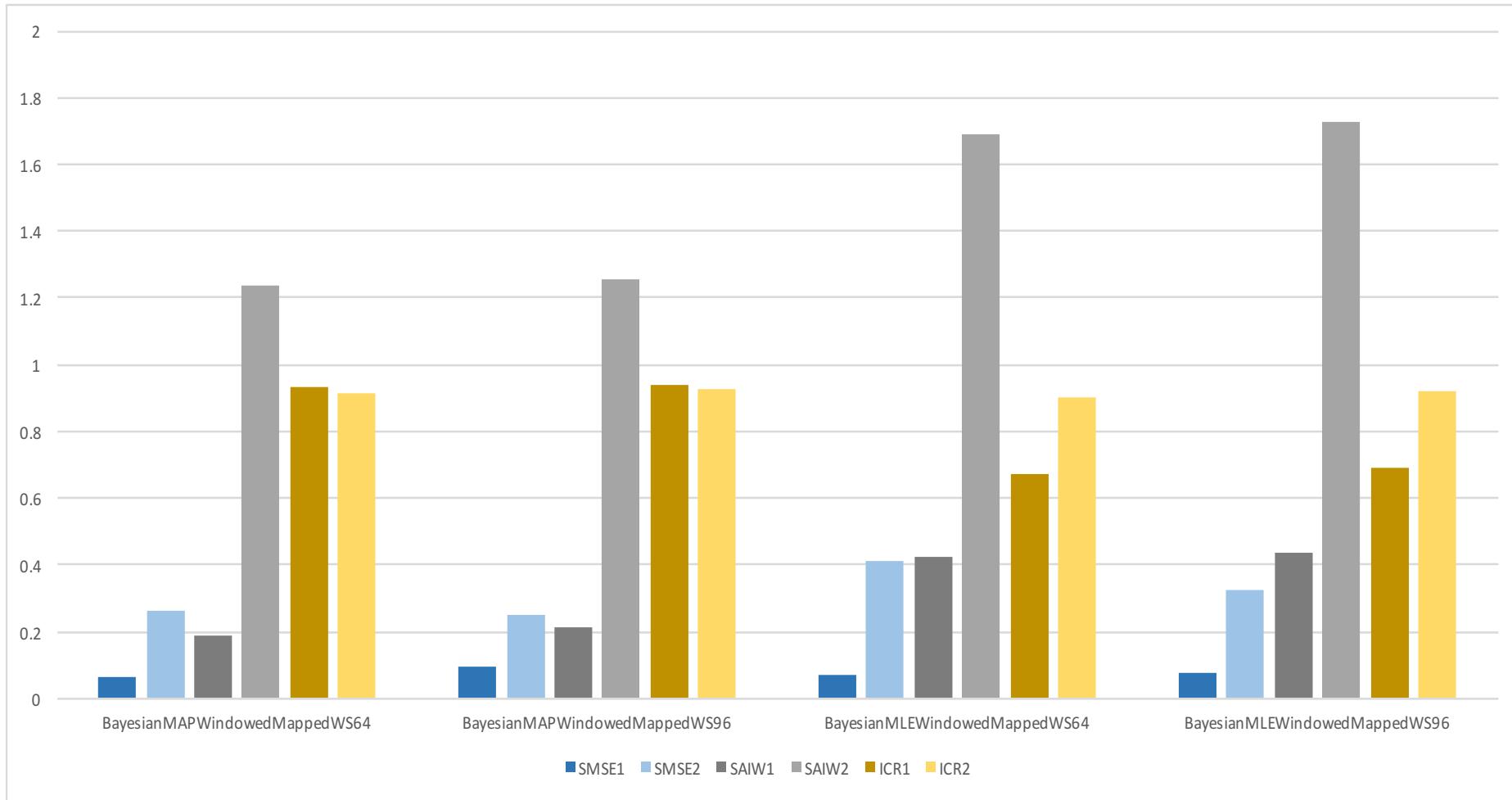


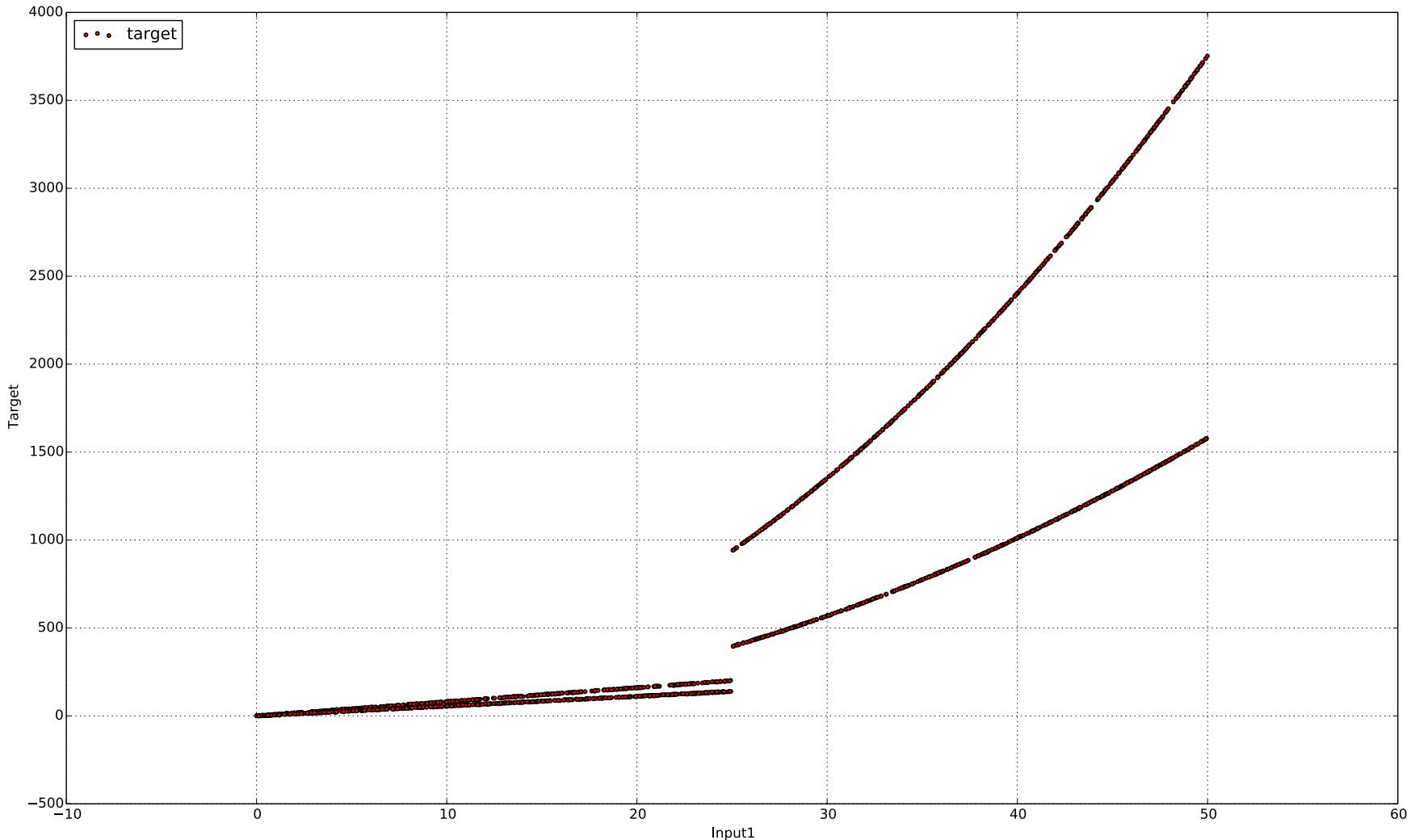


GPWindowedGaussianKernelZeroMean variants with different window sizes on **SYNTH_D_CD_2000_4_10_1_24**

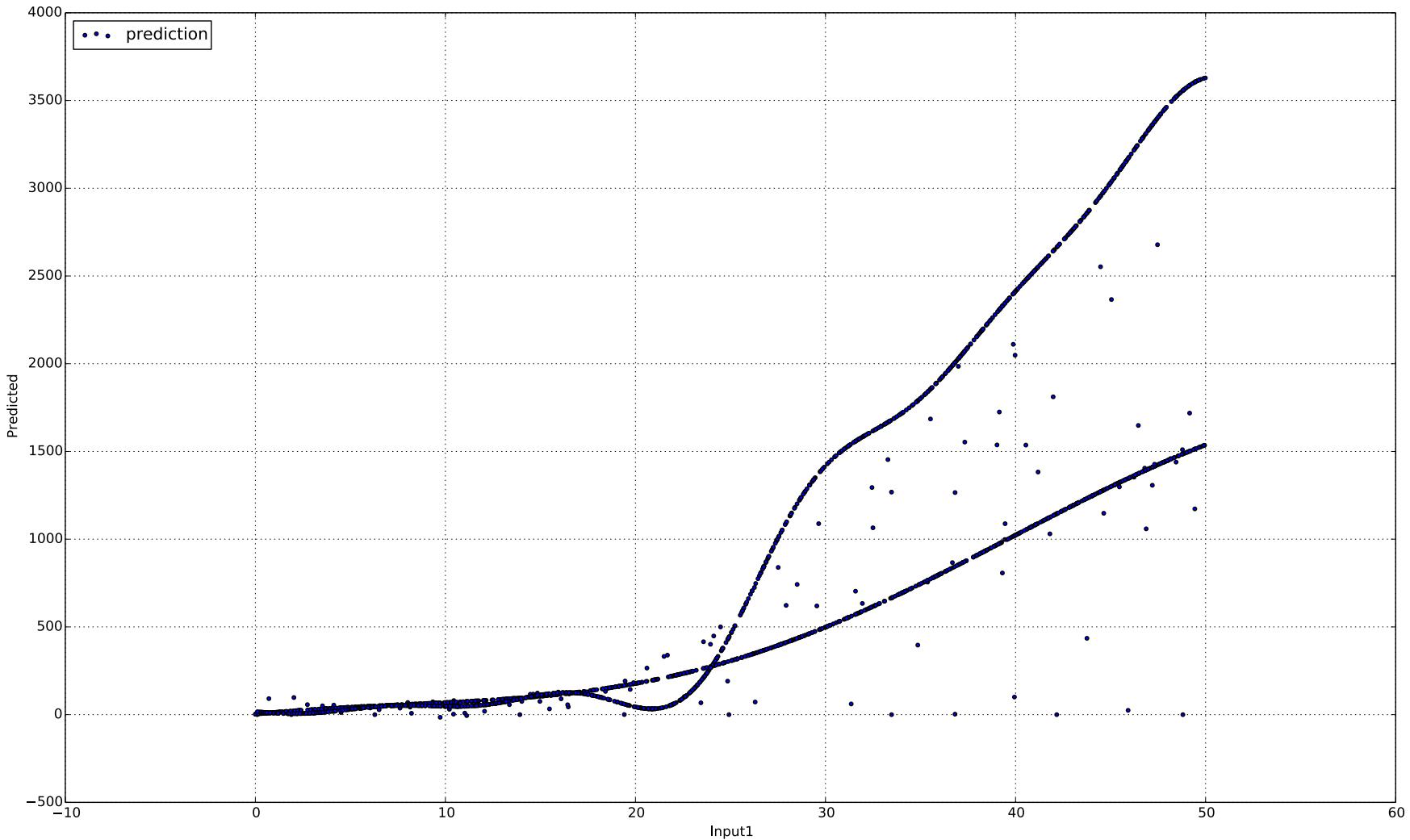


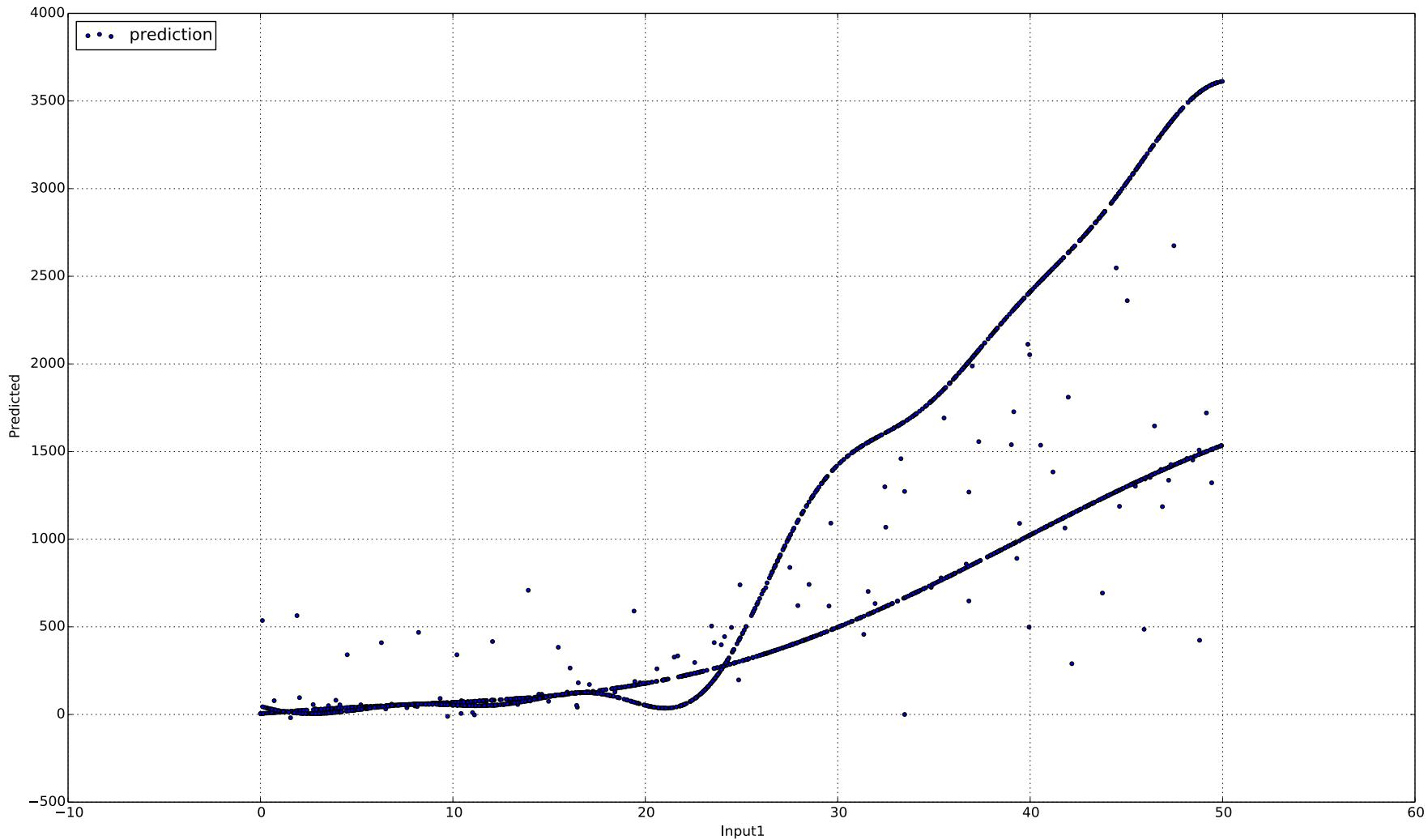
50 sample predictions of **BayesianMAPForgettingMapped_FF0.05** on **SYNTH_D_CD_2000_1_100_1_12**



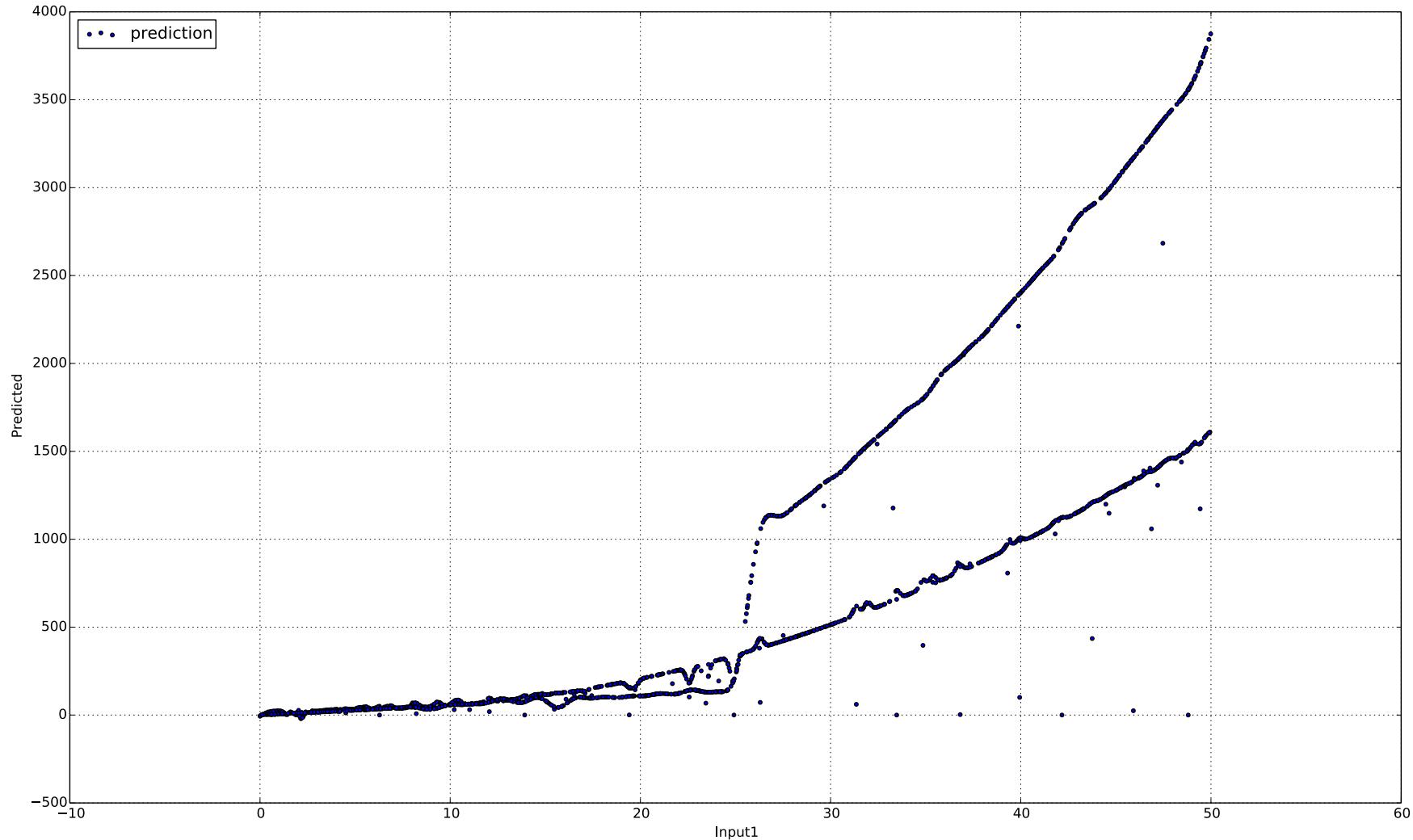


SYNTH_D_CD_2000_1_50_1_13

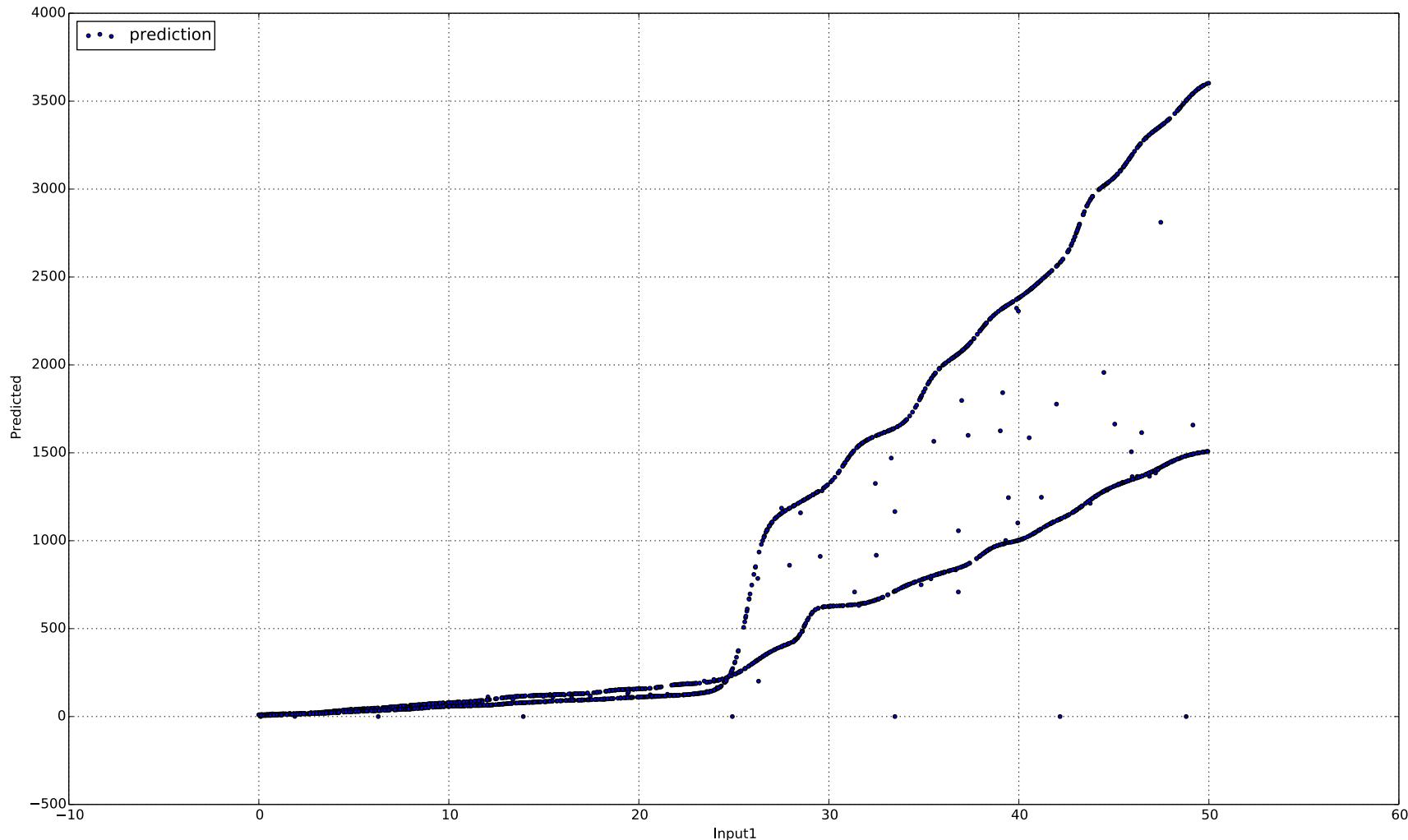




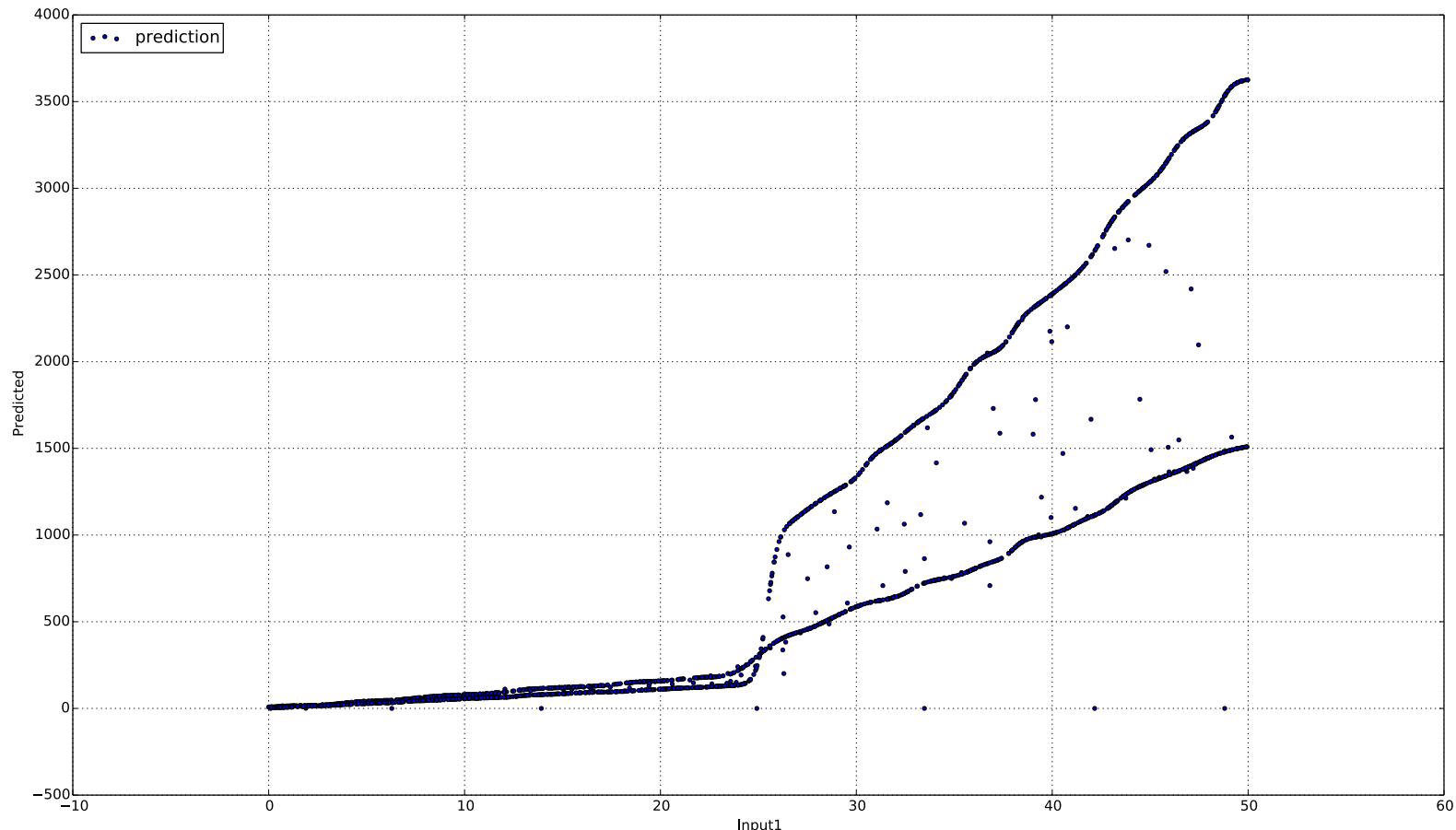
GPRegressionAvgMean_WS64 on SYNTH_D_CD_2000_1_50_1_13



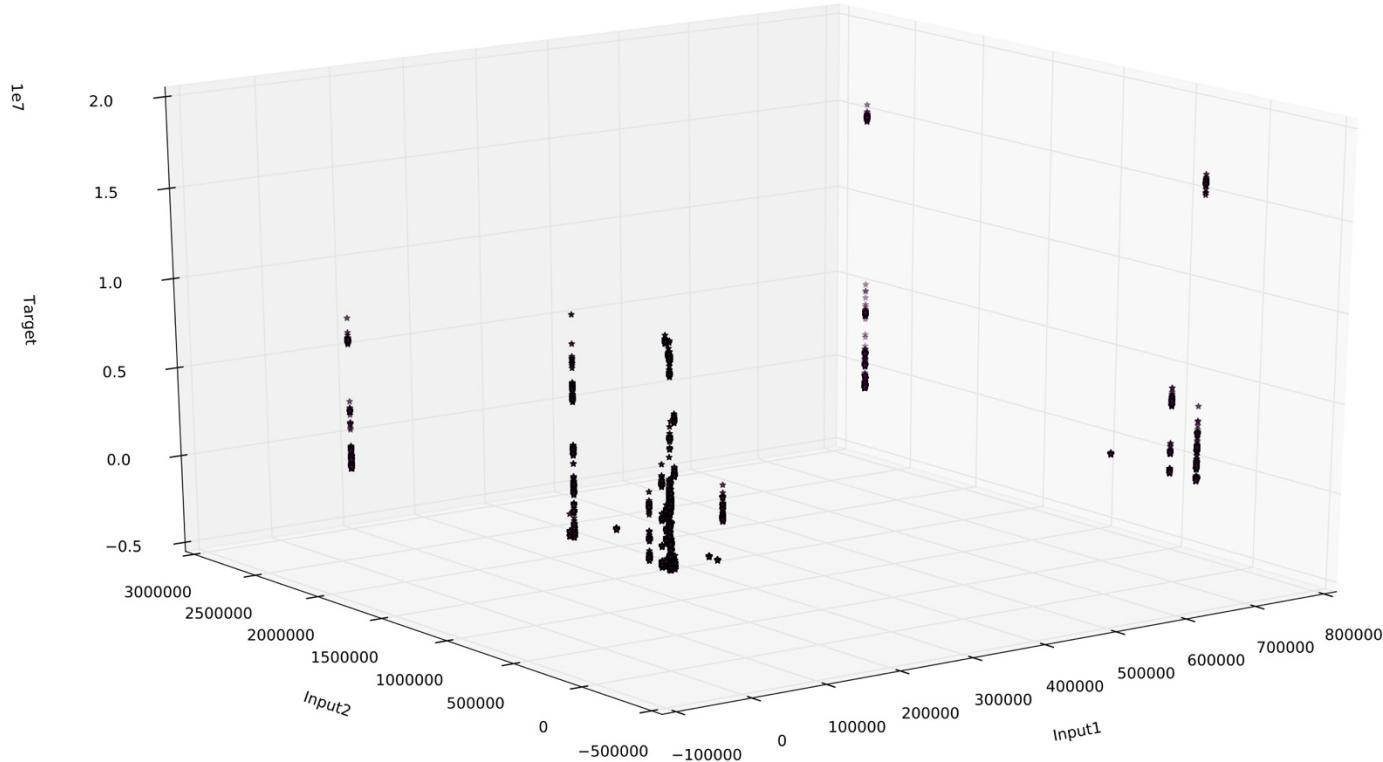
GPRegressionOLSMean_WS64 on SYNTH_D_CD_2000_1_50_1_13



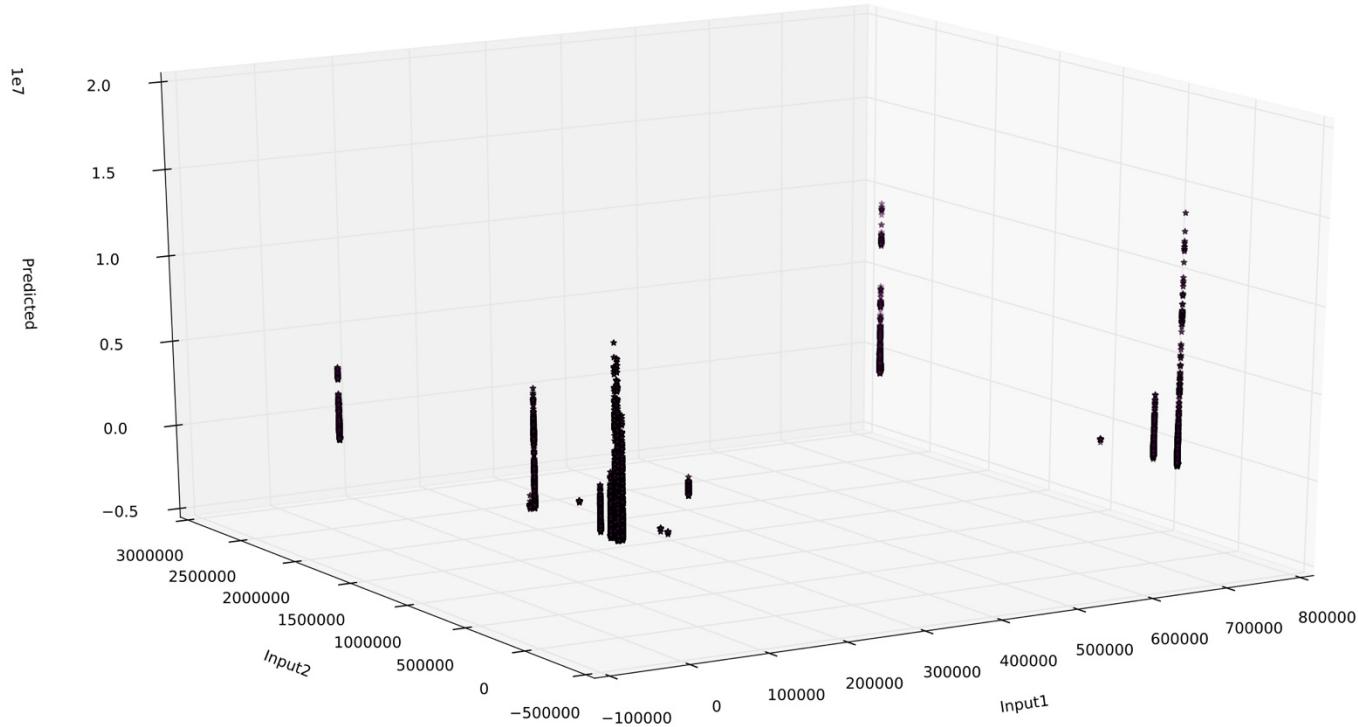
KernelRegression_WS64 on SYNTH_D_CD_2000_1_50_1_13



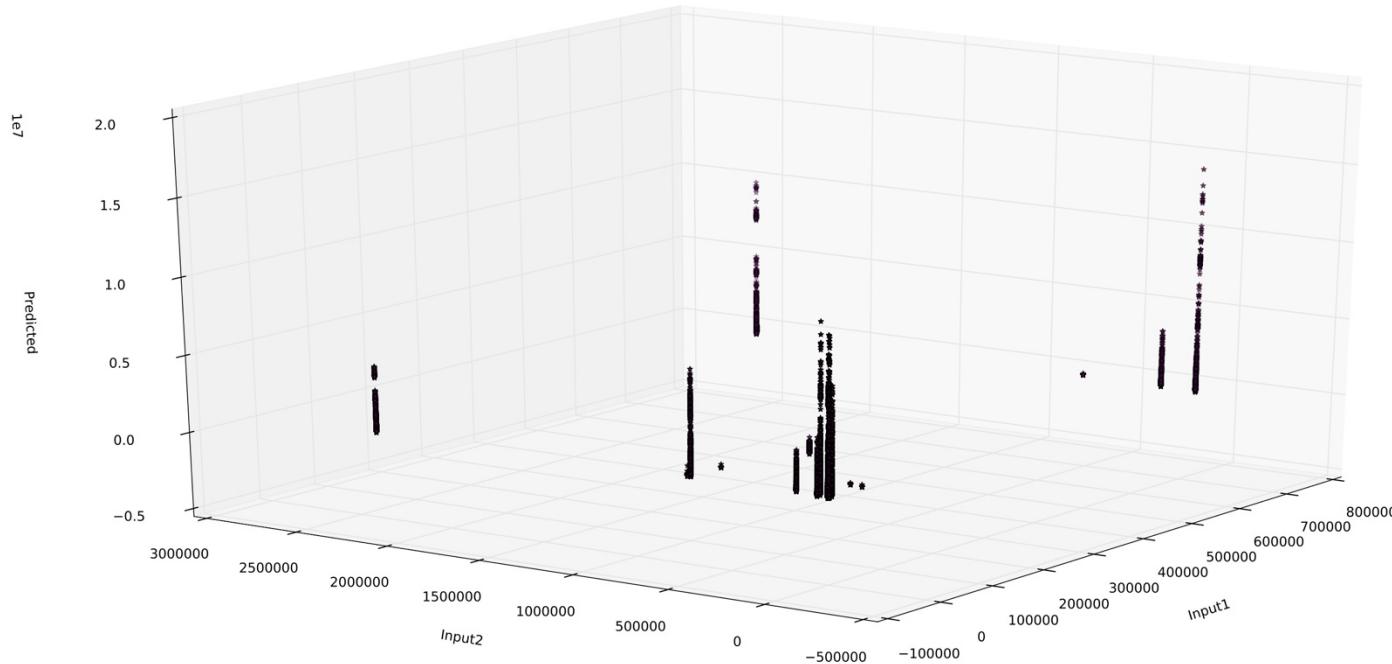
KernelRegression_WS96 on SYNTH_D_CD_2000_1_50_1_13



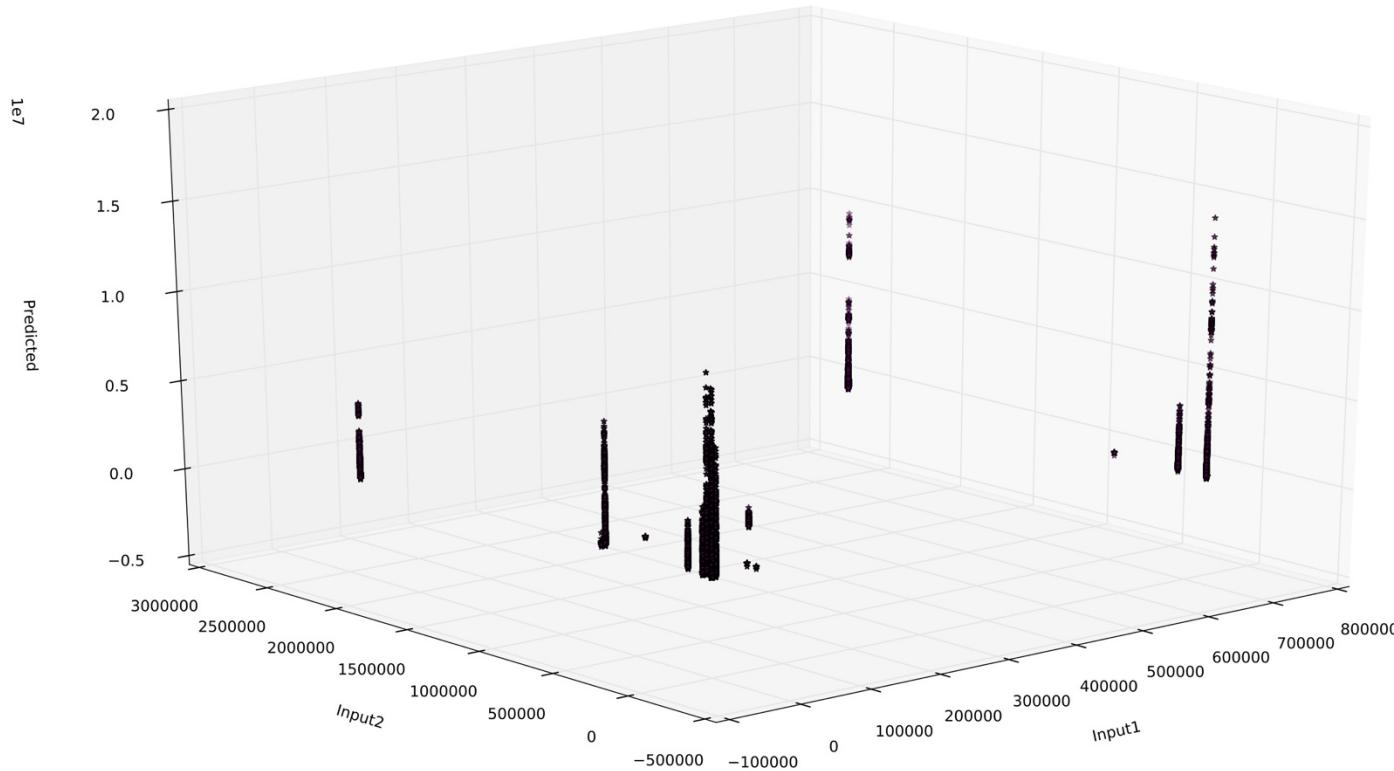
OCL_LEFTFETCHJOIN_GPU



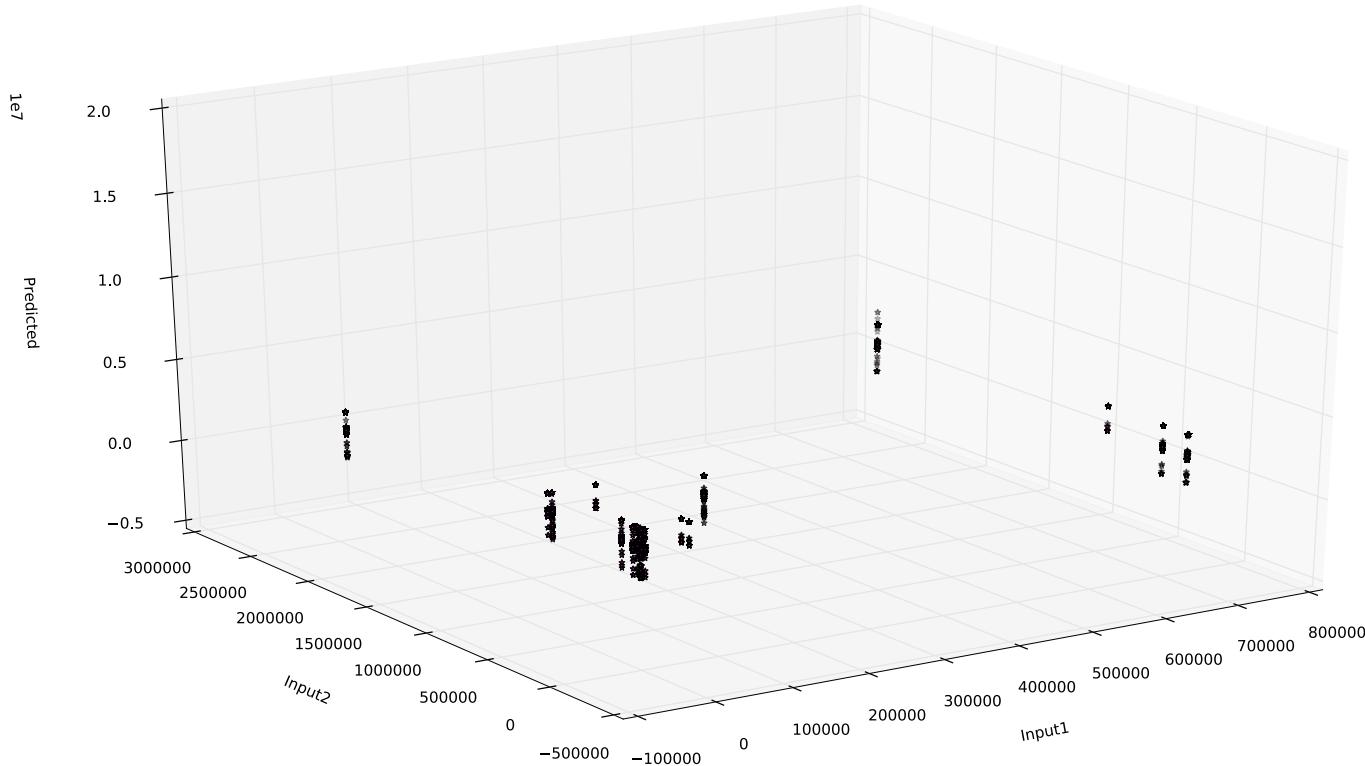
GPRegressionZeroMean_WS64 on OCL_LEFTFETCHJOIN_GPU



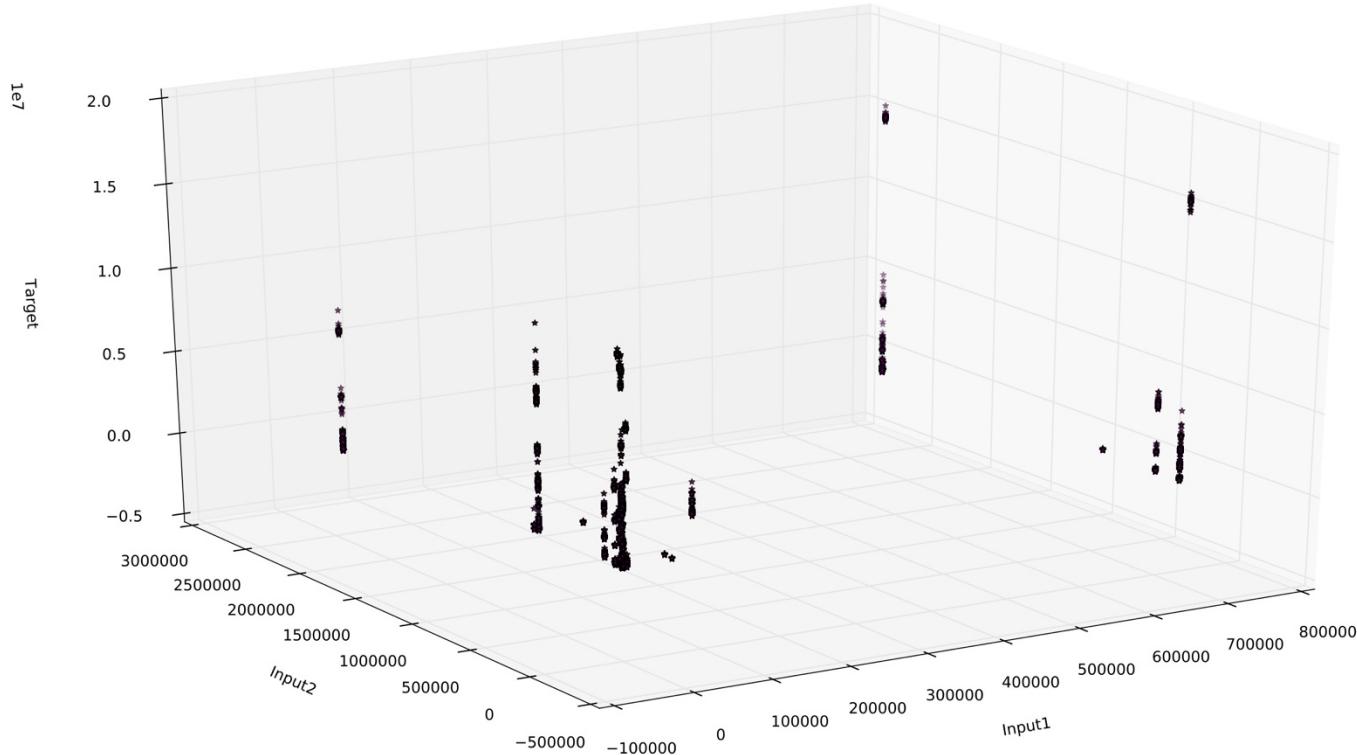
GPRegressionAvgMean_WS64 on OCL_LEFTFETCHJOIN_GPU



GPRegressionOLSMean_WS64 on **OCL_LEFTFETCHJOIN_GPU**



KernelRegression_WS64 on OCL_LEFTFETCHJOIN_GPU



KernelRegression_WS96 on OCL_LEFTFETCHJOIN_GPU

- Sliding-windowed non-parametric online learners are effective in learning from non-stationary data streams
- Sliding-window size has major implications on accuracy and time efficiency
- Forgetting-based parametric learners do not perform well on the drifting data
- Finding good hyperparameter configurations is a challenging task

Algorithm Comparison Table

	Accuracy	Prediction Bounds	Time Efficiency	Space Efficiency
MLE	+	+	+++	+++
MAP	+	+	+++	+++
GP Regression	+++	+++	+	++
Kernel Regression	+++	++	++	++

- Online extension for Quantile Regression and SV-r
- Prediction Bounds Estimation mechanism for SV-r

Carl Edward Rasmussen and Christopher K. I. Williams.
Gaussian Processes for Machine Learning (Adaptive
Computation and Machine Learning). The MIT Press, 2005.
ISBN 0-262-18253-X.

Bernard W Silverman. Density estimation for statistics
and data analysis, volume 26. CRC press, 1986.

Heimel, Max, et al. "Hardware-oblivious parallelism for
in-memory column-stores." *Proceedings of the VLDB
Endowment* 6.9 (2013): 709-720.

Joao Gama, Raquel Sebastiao, and Pedro Pereira
Rodrigues. Issues in evaluation of stream learning
algorithms. *Proceedings of the 15th ACM SIGKDD
international conference on Knowledge discovery and data
mining – KDD '09*, pages 329–337, 2009. doi:
[10.1145/1557019.1557060](https://doi.org/10.1145/1557019.1557060).

Questions?

Thanks for listening! ☺