



UPPSALA  
UNIVERSITET

UC San Diego

# Uncertainty Estimation with Recursive Feature Machines

Daniel Gedon<sup>1</sup>, Amirhesam Abedsoltan<sup>2</sup>, Thomas B. Schön<sup>1</sup>, Mikhail Belkin<sup>2</sup>  
<sup>1</sup>Uppsala University (Sweden), <sup>2</sup>UC San Diego (USA)



## Motivation

**Setup:** Regression with tabular data

Point estimate

Uncertainty quantification

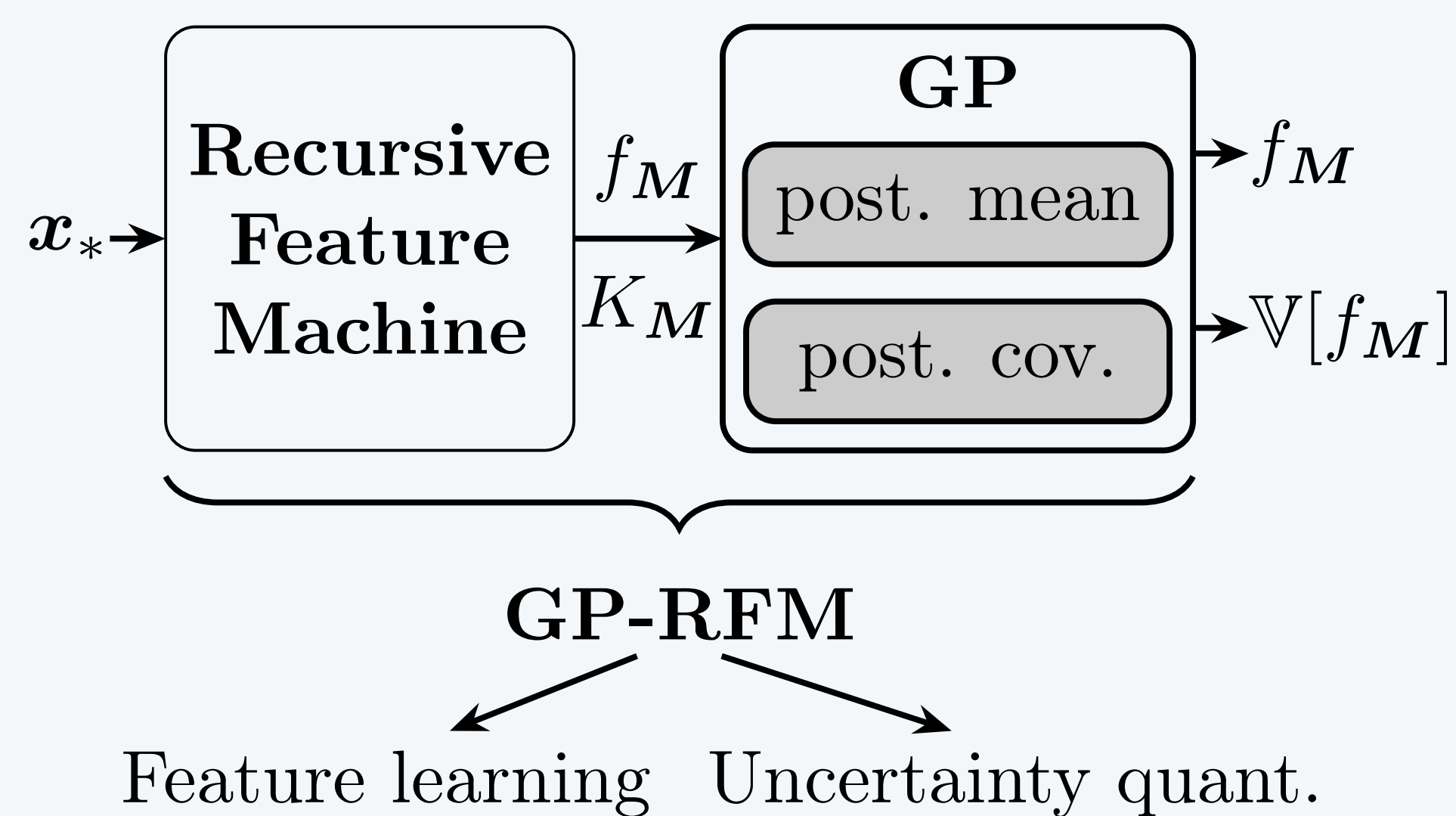
RMSE, ...

NLL, coverage, ...

Boosting vs. Gaussian Process

Can we include feature learning in GPs?

## Method



Restrict to diag. features: **GP-RFM-diag**

Penalize off-diagonal elements:

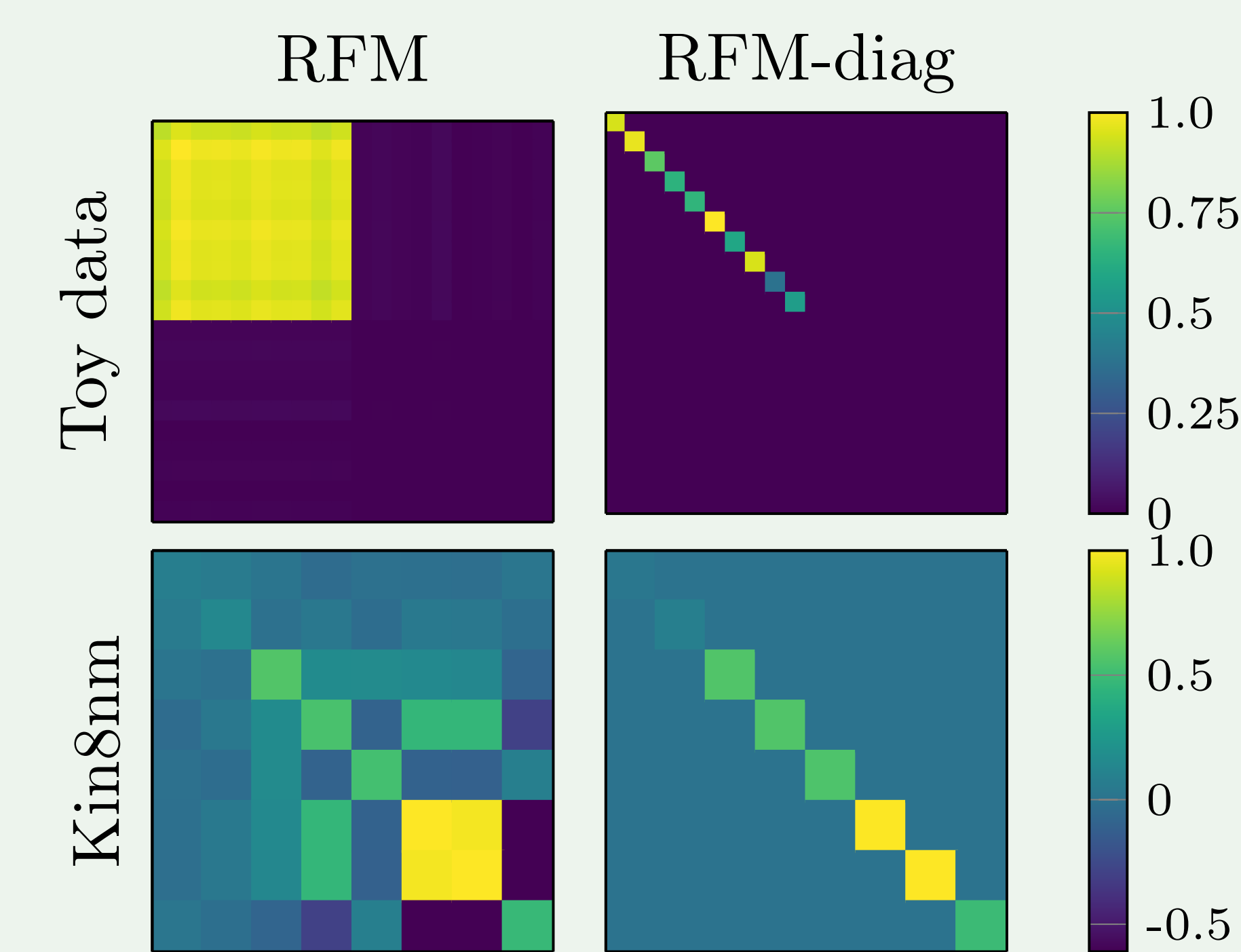
$$M = \frac{1}{n} \sum_{i=1}^n \nabla_x f_M(x) \nabla_x f_M(x)^\top + \lambda_M I_d$$

## Visualising feature matrix $M$

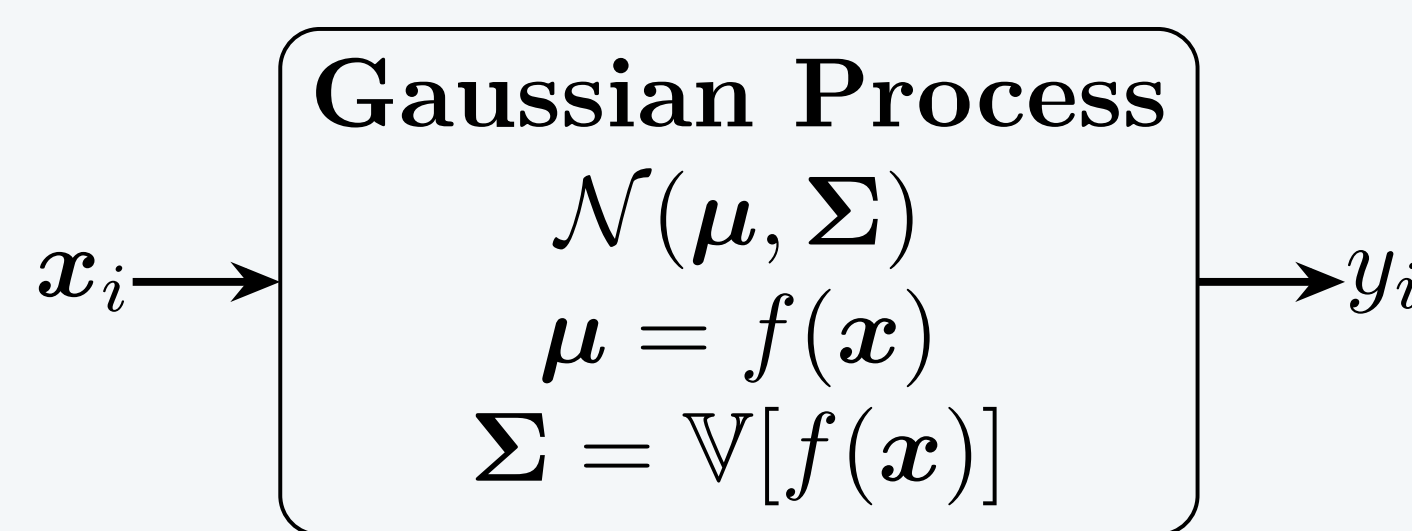
**Data:**  $x \sim \mathcal{U}(0_d, 1_d)$ ,  $y = (\sum_{j=1}^{10} x_{[j]})^2$   
→ introduce correlation.

### Interpretation:

- RFM features capture correlation.
- Correlation in real-world data sets.



## Gaussian Process (GP)



### Model parameterization

pred. function  $f(x) = k(x, X)\alpha$

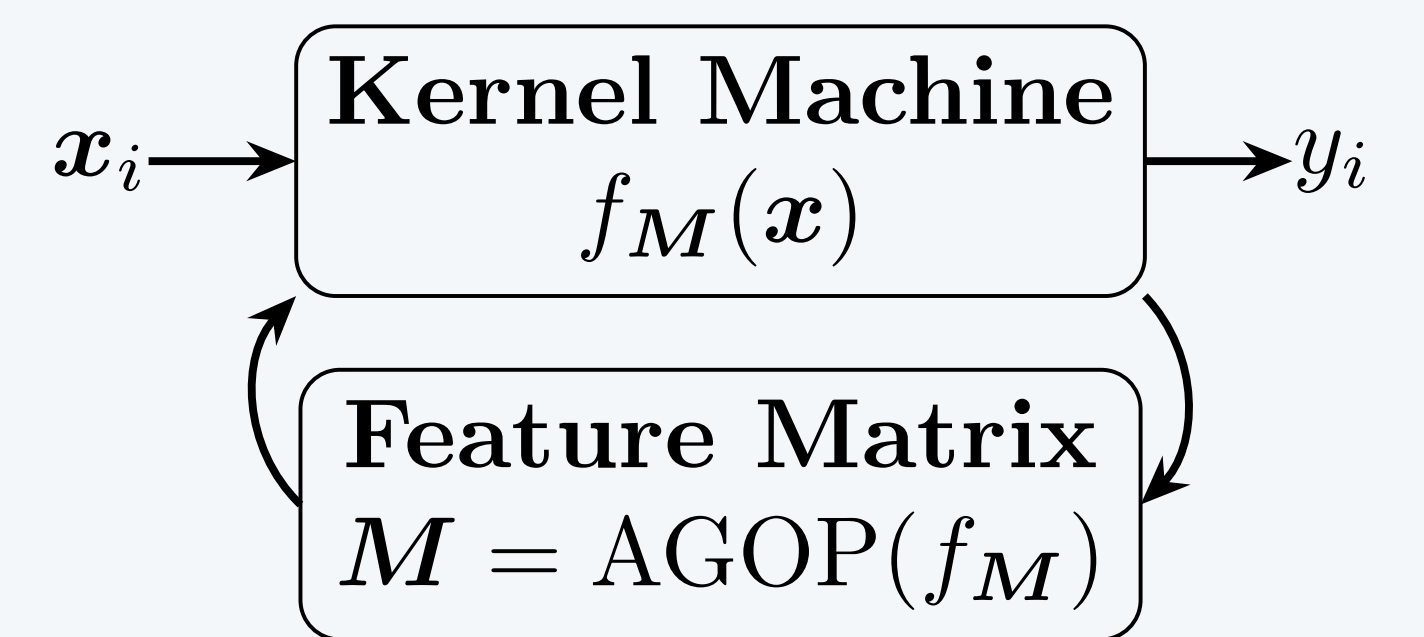
RBf (or Laplace) kernel

$$K_M(x, z) = \exp(-\gamma \|x - z\|_M^2)$$

with Auto. Relevance Det. (ARD)

$$M^{-1} = \text{diag}([\ell_1^2, \dots, \ell_d^2])$$

## Recursive Feature Machine (RFM) [1]



Laplace kernel

$$K_M(x, z) = \exp(-\gamma \|x - z\|_M)$$

with Mahalanobis distance

$$\|x - z\|_M = \sqrt{(x - z)^\top M (x - z)}$$

### Training procedure

Maximum Likelihood Estimation

$$\arg \min_{\theta} -\log p(y | X, \theta)$$

with  $\theta = \{\ell_1, \dots, \ell_d\}$

For  $t$  in  $T$ :

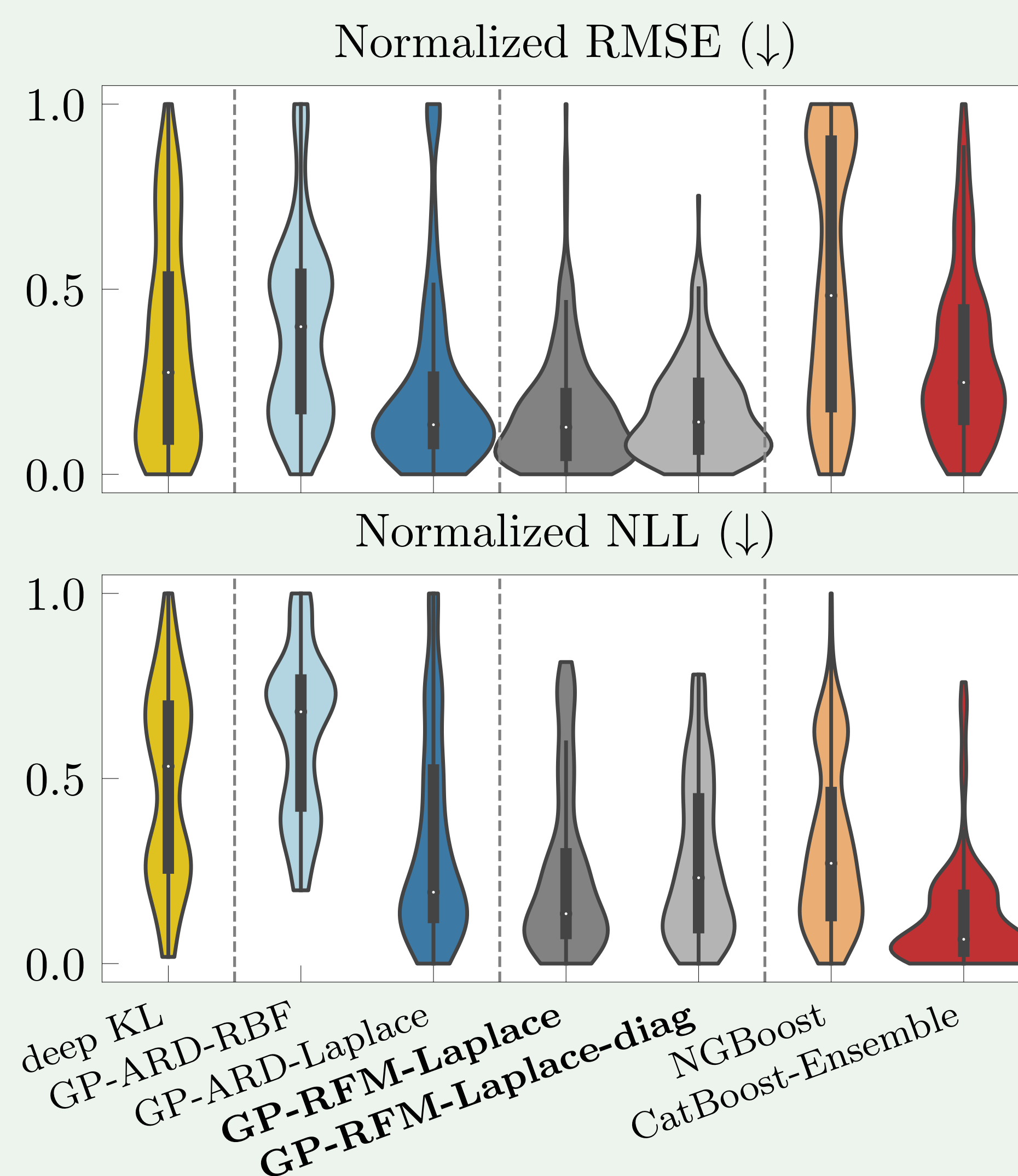
- Solve for kernel weights

$$\alpha = (k_M(x, X) + \lambda_\alpha I_n)^{-1} y$$

- Average gradient outer product (AGOP)

$$M = \frac{1}{n} \sum_{i=1}^n \nabla_x f_M(x_i) \nabla_x f_M(x_i)^\top$$

## Main results



Negative log-likelihood on tabular data.

Dataset	Gauss. Process		Boosting	
	ARD-Lap.	RFM	NG	Cat
cpu-act	2.30	<u>2.21</u>	2.33	<b>2.17</b>
pol	2.84	<u>2.73</u>	3.55	<b>2.09</b>
elevators	-4.75	<b>-4.86</b>	-4.48	-4.73
isolet	3.43	<b>2.34</b>	2.71	<u>2.52</u>
wine	<b>0.95</b>	<b>0.95</b>	1.04	1.03
Ailerons	-7.33	-7.37	<b>-7.42</b>	<u>-7.41</u>
houses	<u>-0.10</u>	-0.07	0.07	<b>-0.12</b>
houses-16H	0.72	0.69	<u>0.57</u>	<b>0.51</b>
Bra-houses	-1.82	-2.11	-2.18	<b>-2.66</b>
bike	6.03	6.04	<u>5.62</u>	<b>5.58</b>
house-sales	-0.30	<b>-0.32</b>	-0.27	<u>-0.31</u>

best, 2nd best

## OpenML datasets [2].

- 16 tabular datasets
- 5 - 613 features
- 6,497 - 22,784 samples

## Setup.

- 20 seeds
- normalised metrics for each dataset

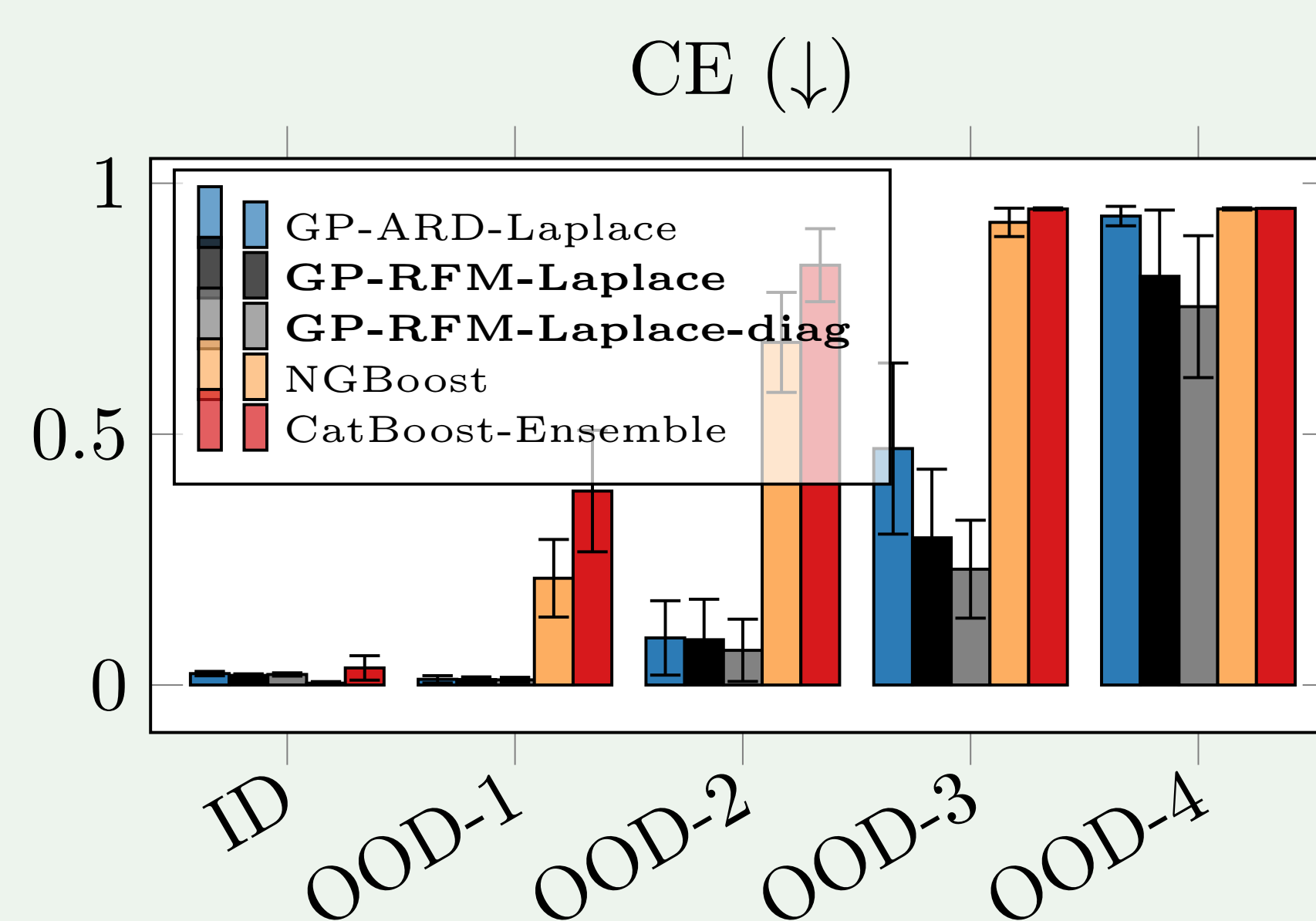
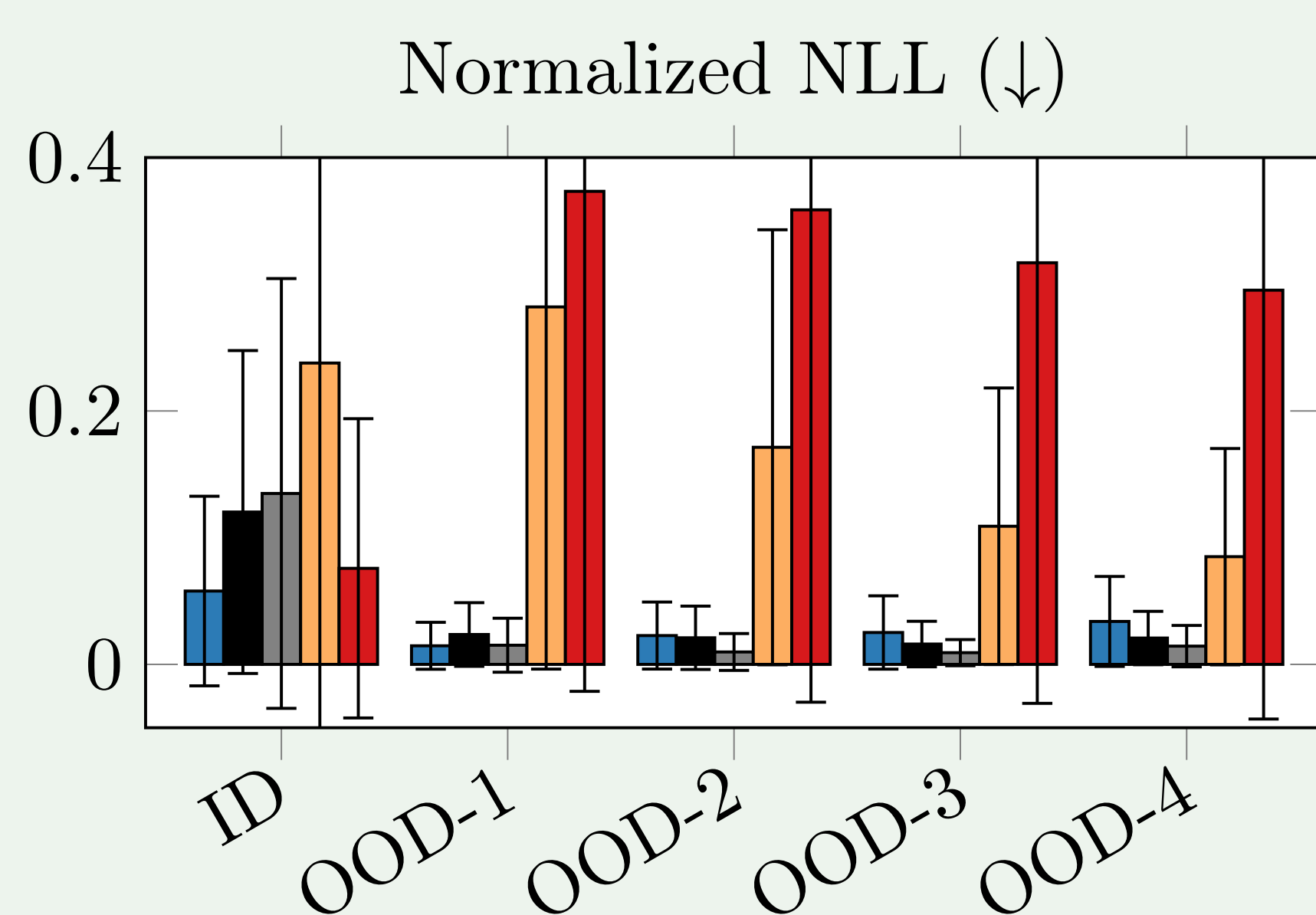
## Interpretation.

- GP-RFM best in RMSE
- GP-RFM overall competitive

## Out-of-distribution data

**Data.** Housing data with increasing target (price) OOD shift.

**Interpretation.** GP-RFM most reliable method under OOD shift.



## Conclusion

**GP + RFM → competitive results**

### Missing in GP literature.

Full feature matrix  $M$  often outperforms diagonal approach

## References

- [1] Radhakrishnan A, Beaglehole D, Pandit P, Belkin M. "Mechanism for feature learning in neural networks and backpropagation-free machine learning models". Science, 2024.
- [2] Vanschoren J, Van Rijn JN, Bischl B, Torgo L. "OpenML: networked science in machine learning." ACM SIGKDD Explorations Newsletter, 2014.