

# REVISITING ACTIVE SETS FOR GAUSSIAN PROCESS DECODERS

Pablo Moreno-Muñoz\*, Cilie W. Feldager\* and Søren Hauberg

## HIGHLIGHTS

- ▷ We revisit active sets for scaling GP decoders, a *sparse* approximation predominantly used before the seminal work of Snelson and Ghahramani (2006).
- ▷ New links between active sets and cross-validation based on the recent result from Fong and Holmes (2020).
- ▷ Formulation of the stochastic active sets (SAS) approach for both deterministic and Bayesian versions of GP decoders.

## HISTORICAL REMARKS

“Traditionally, sparse models have very often been built upon a carefully chosen subset of the training inputs. [...] In sparse Gaussian processes it has also been suggested to select the inducing inputs  $\mathbf{X}_u$  from among the training inputs. Since this involves a prohibitive combinatorial optimization, greedy optimization approaches have been suggested [...]. Recently, Snelson and Ghahramani (2006) have proposed to relax the constraint that the inducing variables must be a subset of training/test cases, turning the discrete selection problem into one of continuous optimization.”

Quiñero-Candela and Rasmussen (2005) on the main difficulties behind the optimal selection of subsets.

## GAUSSIAN PROCESS DECODERS

The Gaussian process latent variable model (GP-LVM) (Lawrence, 2005) defines a decoder which is a non-linear mapping  $\mathbf{x} = f(\mathbf{z})$  from the latent space  $\mathcal{Z} \in \mathbb{R}^Q$  to observation space  $\mathcal{X} \in \mathbb{R}^D$ . The prior on this map is a Gaussian process (GP) so it is drawn like  $f \sim \mathcal{GP}(0, k_\theta(\cdot, \cdot))$ , where  $k_\theta$  is the covariance function or *kernel*.  $\mathbf{K}_{NN}$  denotes the evaluated kernel function so the  $i, j$ th element of  $\mathbf{K}_{NN}$  equals  $k_\theta(\mathbf{z}_i, \mathbf{z}_j)$ .

$$p(\mathbf{x}|f, \mathbf{z}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | f(\mathbf{z}_n), \sigma^2) \quad p(f|\mathbf{z}) = \mathcal{N}(f(\mathbf{z})|0, \mathbf{K}_{NN})$$

$$p(\mathbf{x}|\mathbf{z}) = \int p(\mathbf{x}|f, \mathbf{z})p(f|\mathbf{z})df = \mathcal{N}(\mathbf{x}|0, \mathbf{K}_{NN} + \sigma^2\mathbb{I}).$$

$$\mathcal{L} = -\frac{DN}{2} \log 2\pi - \frac{D}{2} \log |\mathbf{K}_{NN} + \sigma^2\mathbb{I}| - \frac{1}{2} \text{tr}((\mathbf{K}_{NN} + \sigma^2\mathbb{I})^{-1} \mathbf{x}\mathbf{x}^\top)$$

## WHITEBOARD

## ACKNOWLEDGMENTS

This work was supported by research grants (15334, 42062) from VILLUM FONDEN. This project has also received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement 757360). This work was funded in part by the Novo Nordisk Foundation through the Center for Basic Machine Learning Research in Life Science (NNF20OC0062606). This work was further supported by the Pioneer Centre for AI, DNRF grant number P1.

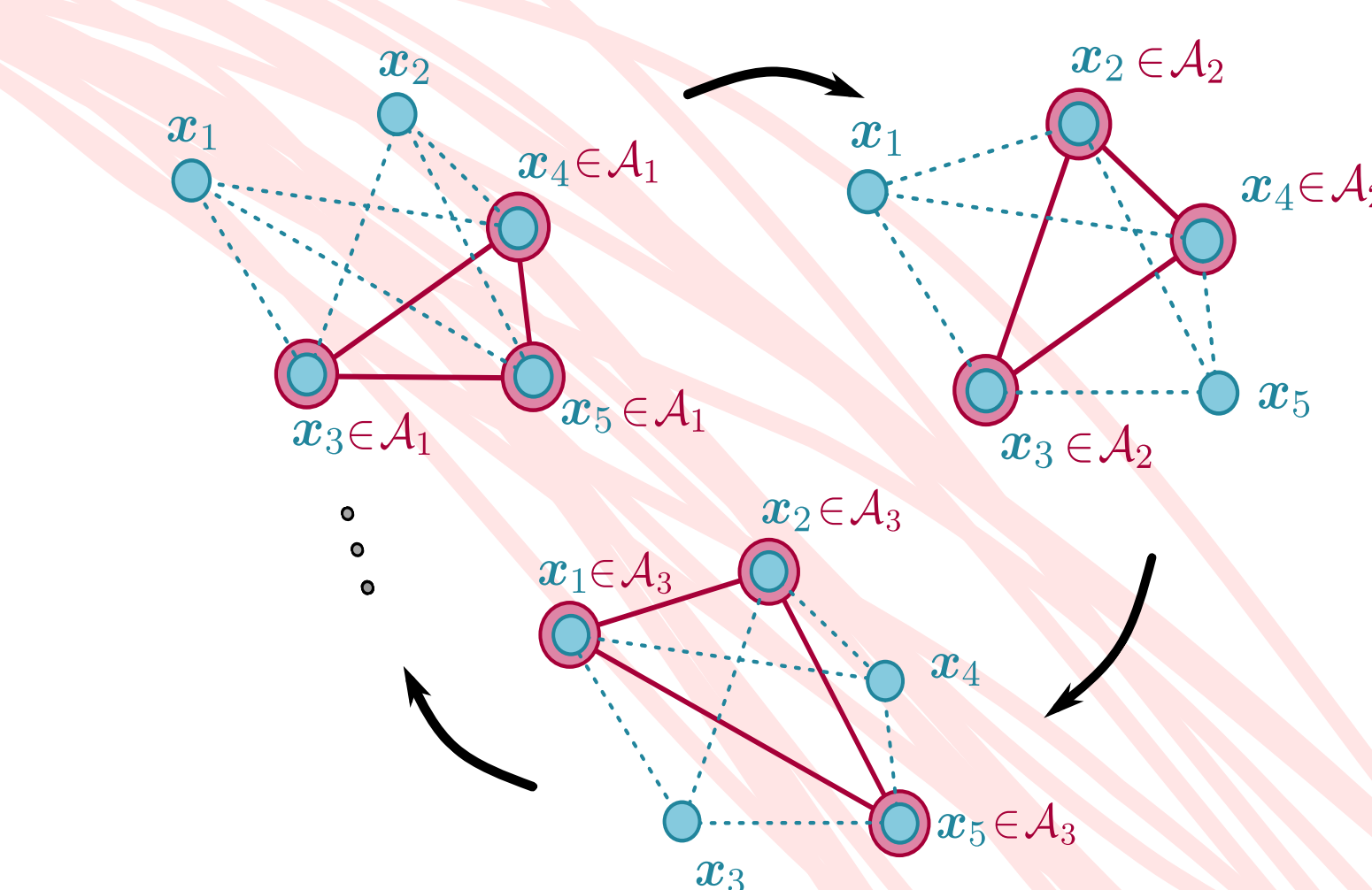
## STOCHASTIC ACTIVE SETS

Recently, Fong and Holmes (2020) linked *cross validation* (CV) with the log-marginal likelihood, effectively showing that the latter is equivalent to the average over exhaustive leave- $R$ -out CV scores. In particular, the average is w.r.t. the size of the hold-out set.

$$\mathcal{S}_{CV}(\mathbf{x}|R) = \frac{1}{C} \sum_{p=1}^C \frac{1}{R} \sum_{n \in \mathcal{R}_p} \log p(\mathbf{x}_n | \mathbf{x}_{\mathcal{A}_p}, \mathbf{z}) = \frac{1}{R} \mathbb{E}_{\mathcal{A}} \left[ \sum_{n \in \mathcal{R}_p} \log p(\mathbf{x}_n | \mathbf{x}_{\mathcal{A}}, \mathbf{z}) \right],$$

$$\log p(\mathbf{x}|\mathbf{z}) = \sum_{r=1}^R \mathcal{S}_{CV}(\mathbf{x}|r) = \mathcal{S}_{CCV}(\mathbf{x}|R) + \mathcal{S}_{PCV}(\mathbf{x}|R). \quad \mathcal{S}_{CCV}(\mathbf{x}|R) = \sum_{r=1}^R \mathcal{S}_{CV}(\mathbf{x}|r)$$

$$\mathcal{S}_{PCV}(\mathbf{x}|R) = \mathbb{E}_{\mathcal{A}}[\log p(\mathbf{x}_{\mathcal{A}}|\mathbf{z}_{\mathcal{A}})]$$



## STOCHASTIC APPROXIMATION

$$p(\mathbf{x}_{\mathcal{A}}|\mathbf{z}_{\mathcal{A}}) = \mathcal{N}(\mathbf{x}_{\mathcal{A}}|0, \mathbf{K}_{\mathcal{A}\mathcal{A}} + \sigma^2\mathbb{I}), \quad p(\mathbf{x}_n|\mathbf{x}_{\mathcal{A}}, \mathbf{z}) = \mathcal{N}(\mathbf{x}_n|\mathbf{m}_{n|\mathcal{A}}, \mathbf{c}_{n|\mathcal{A}}),$$

$$\log p(\mathbf{x}|\mathbf{z}) \approx \sum_{n \in \mathcal{R}} \log p(\mathbf{x}_n|\mathbf{x}_{\mathcal{A}}, \mathbf{z}) + \log p(\mathbf{x}_{\mathcal{A}}|\mathbf{z}_{\mathcal{A}})$$

## TRAINING ALGORITHMS

### Algorithm 1 SAS for GP decoders

- 1: **Input:** Observed data  $\mathbf{x}$
- 2: **Parameters:** Initialize  $\theta, \phi$  //  $\theta, \mathbf{z}$  if NA
- 3: **for**  $e$  **in** epochs **do**
- 4:   **for**  $b$  **in** batches **do**
- 5:     Sample  $\mathbf{x}_{\text{batch}} \sim \mathbf{x}$
- 6:      $\mathbf{x}_{\mathcal{R}}, \mathbf{x}_{\mathcal{A}} \leftarrow \text{random\_split}(\mathbf{x}_{\text{batch}})$
- 7:     **if** amortized **then**
- 8:       Get  $\{\mathbf{z}_{\mathcal{R}}, \mathbf{z}_{\mathcal{A}}\} \leftarrow g(\mathbf{x}_{\mathcal{R}}, \mathbf{x}_{\mathcal{A}}|\phi)$
- 9:     **end if**
- 10:     Compute  $\mathbf{K}_{\mathcal{A}\mathcal{A}}^{-1}$  // via Cholesky
- 11:     Evaluate  $\log p(\mathbf{x}_{\mathcal{A}}|\mathbf{z}_{\mathcal{A}})$
- 12:     Evaluate  $\log p(\mathbf{x}_n|\mathbf{x}_{\mathcal{A}}, \mathbf{z})$ ,  $\forall \mathbf{x}_n \in \mathbf{x}_{\mathcal{R}}$
- 13:     Evaluate Eq. 6
- 14:     **do** Adam( $\theta, \phi$ ) **step for**  $\mathcal{L}$
- 15:   **end for**
- 16: **end for**

NA: Non-amortized.

### Algorithm 2 SAS for Bayesian GP decoders

- 1: **Input:** Observed data  $\mathbf{x}$
- 2: **Parameters:** Initialize  $\theta, \phi$  //  $\theta, \mu, \sigma$  if NA
- 3: **for**  $e$  **in** epochs **do**
- 4:   **for**  $b$  **in** batches **do**
- 5:     Sample  $\mathbf{x}_{\text{batch}} \sim \mathbf{x}$
- 6:      $\mathbf{x}_{\mathcal{R}}, \mathbf{x}_{\mathcal{A}} \leftarrow \text{random\_split}(\mathbf{x}_{\text{batch}})$
- 7:     **if** amortized **then**
- 8:       Get  $\mu_{\mathbf{z}} \leftarrow g_{\mu}(\mathbf{x}_{\mathcal{R}}, \mathbf{x}_{\mathcal{A}}|\phi)$
- 9:       Get  $\sigma_{\mathbf{z}} \leftarrow g_{\sigma}(\mathbf{x}_{\mathcal{R}}, \mathbf{x}_{\mathcal{A}}|\phi)$
- 10:     **end if**
- 11:     Sample  $\{\mathbf{z}_{\mathcal{R}}, \mathbf{z}_{\mathcal{A}}\} \sim q(\mu_{\mathbf{z}}, \sigma_{\mathbf{z}})$  // RT
- 12:     Compute  $\mathbf{K}_{\mathcal{A}\mathcal{A}}^{-1}$  // via Cholesky
- 13:     Evaluate  $\mathcal{L}$  in Eq. 8
- 14:     **do** Adam( $\theta, \phi$ ) **step for**  $\mathcal{L}_{\text{ELBO}}$
- 15:   **end for**
- 16: **end for**

NA: Non-amortized, RT: Reparametrization trick.

## EXPERIMENTS & RESULTS

Table 1: Comparative metrics for SAS and Bayesian SAS on MNIST, FMNIST and CIFAR-10.

MODEL	SAS			BAYESIAN SAS		
	A = 100	A = 200	A = 400	A = 100	A = 200	A = 400
MNIST / RMSE ↓	2.55 ± 0.98	2.47 ± 0.98	2.41 ± 0.93	2.16 ± 0.02	2.08 ± 0.02	1.99 ± 0.02
MNIST / MAE ↓	1.61 ± 0.97	1.55 ± 0.99	1.51 ± 0.96	1.11 ± 0.02	1.04 ± 0.02	0.96 ± 0.01
MNIST / NLPD ↓	2.99 ± 1.41	2.92 ± 1.38	2.84 ± 1.31	2.33 ± 0.03	2.26 ± 0.02	2.17 ± 0.02
FMNIST / RMSE ↓	2.37 ± 0.95	2.31 ± 0.94	2.25 ± 0.90	1.99 ± 0.17	1.88 ± 0.20	1.85 ± 0.13
FMNIST / MAE ↓	1.48 ± 0.91	1.42 ± 0.91	1.39 ± 0.89	1.11 ± 0.02	1.02 ± 0.03	0.98 ± 0.02
FMNIST / NLPD ↓	2.76 ± 1.33	2.71 ± 1.31	2.65 ± 1.23	2.16 ± 0.18	2.07 ± 0.19	2.04 ± 0.12
CIFAR10 / RMSE ↓	2.66 ± 1.08	2.55 ± 1.06	2.55 ± 1.03	2.74 ± 1.07	2.64 ± 1.08	2.57 ± 1.02
CIFAR10 / MAE ↓	1.77 ± 1.06	1.69 ± 1.06	1.69 ± 1.02	1.84 ± 1.03	1.76 ± 1.05	1.71 ± 1.03
CIFAR10 / NLPD ↓	3.20 ± 1.55	3.07 ± 1.44	3.32 ± 1.89	3.24 ± 1.53	3.14 ± 1.53	3.06 ± 1.45

All metrics are ( $\times 10^{-1}$ ).

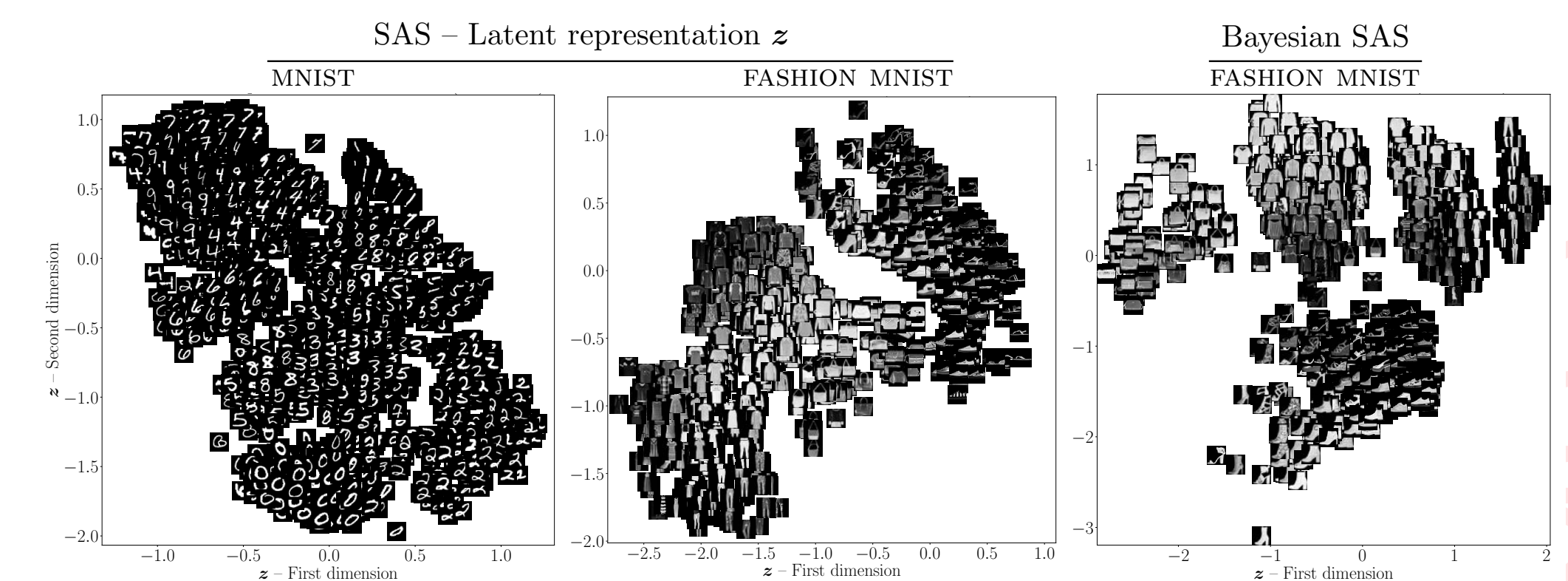
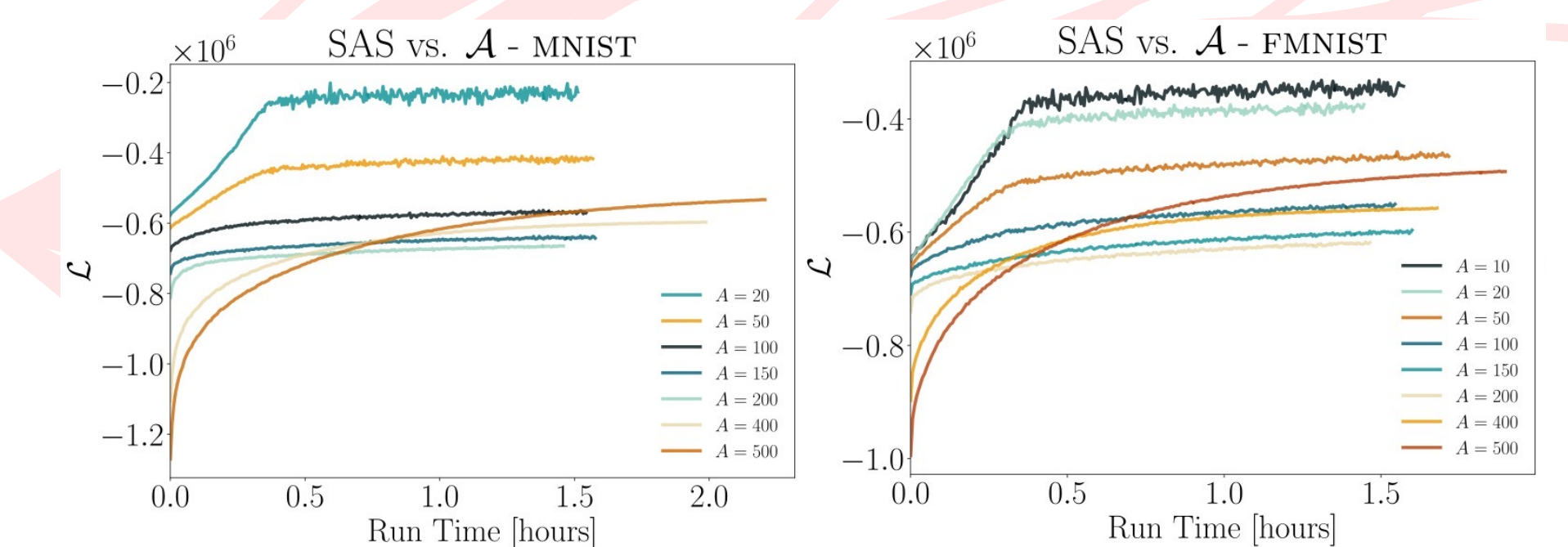


Table 2: Classification accuracy ( $\uparrow$ ) on 2-dim. latent space  $\mathcal{Z}$ .

MODEL	MNIST	FMNIST
BAYESIAN SAS-GP DEC. (ours)	0.63 ± 0.022	0.63 ± 0.020
BAYESIAN GP-LVM	0.18 ± 0.033	0.24 ± 0.043
VAE	0.54 ± 0.026	0.58 ± 0.008



## REFERENCES

- J. Quiñero-Candela and C. E. Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research (JMLR)*, 6(Dec):1939–1959, 2005.
- E. Fong and C. C. Holmes. On the marginal likelihood and cross-validation. *Biometrika*, 107(2):489–496, 2020.
- N. Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *Journal of Machine Learning Research (JMLR)*, 6(11), 2005.
- E. Snelson and Z. Ghahramani. Sparse Gaussian processes using pseudo-inputs. *In Advances in Neural Information Processing Systems (NIPS)*, 2006.

