

Input-gradient space particle inference for neural network ensembles

Trung Trinh¹, Markus Heinonen¹, Luigi Acerbi², Samuel Kaski^{1,3}
¹Aalto University, ²Helsinki University, ³University of Manchester
 {trung.trinh, markus.o.heinonen, samuel.kaski}@aalto.fi, luigi.acerbi@helsinki.fi



Overview

TL;DR: We learn an ensemble of neural networks that is **diverse** with respect to their **input gradients**.

Repulsive deep ensembles (RDEs) [1]

Description: Train an ensemble $\{\theta_i\}_{i=1}^M$ using Wasserstein gradient descent [2], which employs a **kernelized repulsion term** to diversify the particles to cover the **Bayes posterior** $p(\theta|\mathcal{D})$

$$\theta_i^{(t+1)} = \theta_i^{(t)} + \eta_t \left(\underbrace{\nabla_{\theta_i^{(t)}} \log p(\theta_i^{(t)} | \mathcal{D})}_{\text{Driving force}} - \underbrace{\frac{\sum_{j=1}^N \nabla_{\theta_i^{(t)}} k(\theta_i^{(t)}, \theta_j^{(t)})}{\sum_{j=1}^N k(\theta_i^{(t)}, \theta_j^{(t)})}}_{\text{Repulsion force}} \right)$$

- The **driving force** directs the particles towards high density regions of the posterior.
- The **repulsion force** pushes the particles away from each other to enforce diversity.

Problem: It is unclear how to define the repulsion term for neural networks:

- weight-space repulsion is ineffective due to overparameterization.
- function-space repulsion often results in underfitting.

Defining the input-gradient kernel k

Given a base kernel κ , we define the kernel in the input-gradient space for a minibatch of training samples $\mathcal{B} = \{(\mathbf{x}_b, y_b)\}_{b=1}^B$ as follows:

$$k(\theta_i, \theta_j) = \frac{1}{B} \sum_{b=1}^B \kappa(\nabla_{\mathbf{x}_b} \mathbf{f}(\mathbf{x}_b; \theta_i)_{y_b}, \nabla_{\mathbf{x}_b} \mathbf{f}(\mathbf{x}_b; \theta_j)_{y_b})$$

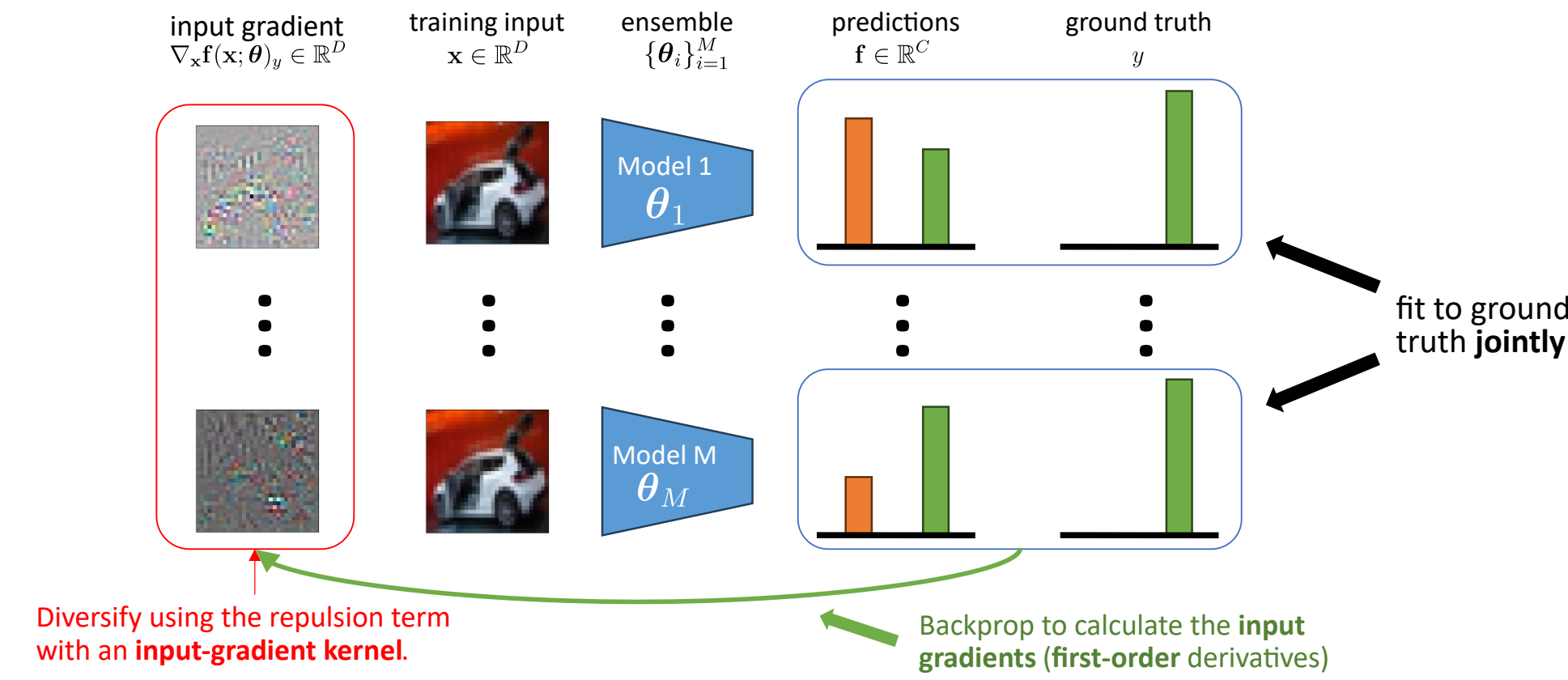
Take the average over all samples in the mini-batch. Compare the **input gradients** of two particles with respect to **the same input**.

We choose the **RBF kernel** on a **unit sphere** as the base kernel κ :

$$\kappa(\mathbf{s}_i, \mathbf{s}_j; \Sigma) = \exp \left(-\frac{1}{2h} (\mathbf{s}_i - \mathbf{s}_j)^\top \Sigma^{-1} (\mathbf{s}_i - \mathbf{s}_j) \right), \quad \mathbf{s}_i = \frac{\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}; \theta_i)_y}{\|\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}; \theta_i)_y\|_2}$$

A scalar adaptively adjusted to prevent kernel vanishing. Diagonal matrix containing the lengthscales. Normalize input gradients to unit vectors.

First-order Repulsive deep ensembles (FoRDEs)



Possible advantages:

- Each member is guaranteed to represent a different function;
- The issues of weight- and function-space repulsion are avoided;
- Each member is encouraged to learn different features, which can improve robustness.

Main takeaways

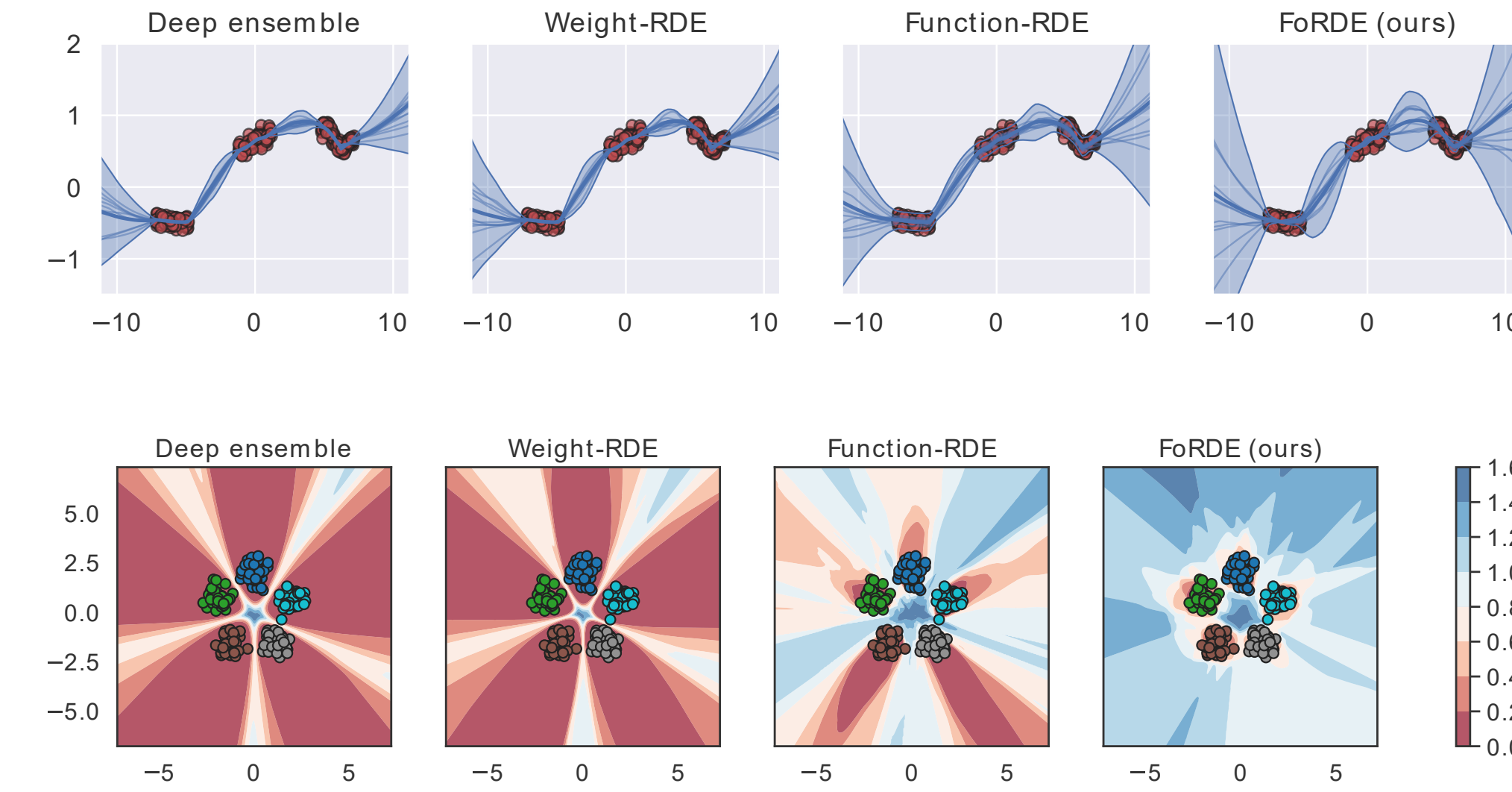
- Input-gradient-space repulsion can perform better than weight- and function-space repulsion.
- Better corruption robustness can be achieved by configuring the repulsion kernel using the eigen-decomposition of the training data.

Benchmark comparison

Table 1: **FoRDE-PCA** achieves the best performance under corruptions while **FoRDE-Identity** outperforms baselines on clean data. **FoRDE-Tuned** outperforms baselines on both clean and corrupted data. Results of RESNET18 / CIFAR-100 averaged over 5 seeds. Each ensemble has 10 members. cA, cNLL and cECE are accuracy, NLL, and ECE on CIFAR-100-C.

METHOD	NLL ↓	ACCURACY (%) ↑	ECE ↓	cA / cNLL / cECE
DEEP ENSEMBLES	0.70±0.00	81.8±0.2	0.041±0.003	54.3 / 1.99 / 0.05
WEIGHT-RDE	0.70±0.01	81.7±0.3	0.043±0.004	54.2 / 2.01 / 0.06
FUNCTION-RDE	0.76±0.02	80.1±0.4	0.042±0.005	51.9 / 2.08 / 0.07
FORDE-PCA (OURS)	0.71±0.00	81.4±0.2	0.039±0.002	56.1 / 1.90 / 0.05
FORDE-IDENTITY (OURS)	0.70±0.00	82.1±0.2	0.043±0.001	54.1 / 2.02 / 0.05
FORDE-TUNED (OURS)	0.70±0.00	82.1±0.2	0.044±0.002	55.3 / 1.94 / 0.05

Illustrative experiments



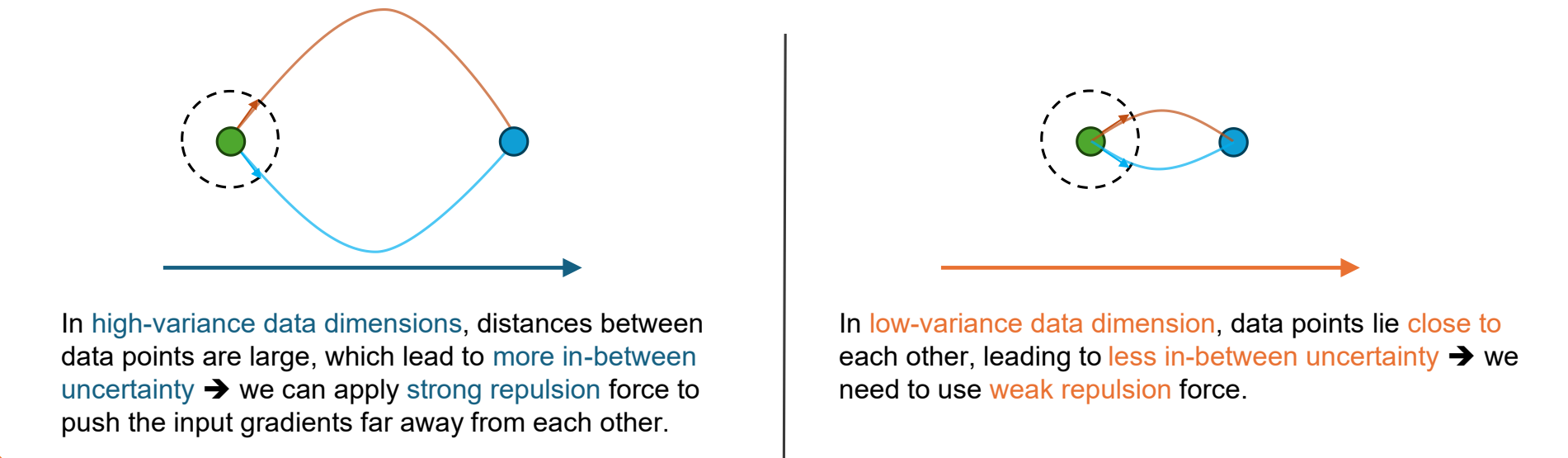
For a 1D regression task (above) and a 2D classification task (below), FoRDEs capture higher uncertainty than baselines in all regions outside of the training data. For the 2D classification task, we visualize the entropy of the predictive posteriors.

Tuning the lengthscales Σ

Each **lengthscale** is inversely proportional to the strength of the repulsion force in the corresponding input dimension.

$$\underbrace{\frac{\partial}{\partial s_d} \kappa(\mathbf{s}, \mathbf{s}'; \Sigma)}_{\text{Repulsion force in the } d\text{-th dimension}} = -\frac{s_d - s'_d}{h \Sigma_{dd}} \kappa(\mathbf{s}, \mathbf{s}'; \Sigma) \propto \frac{1}{\underbrace{\Sigma_{dd}}_{\text{lengthscale in the } d\text{-th dimension}}}$$

Proposition: One should apply **strong forces** in **high-variance dimensions** (more in-between uncertainty) and **weak forces** in **low-variance dimensions** (less in-between uncertainty).

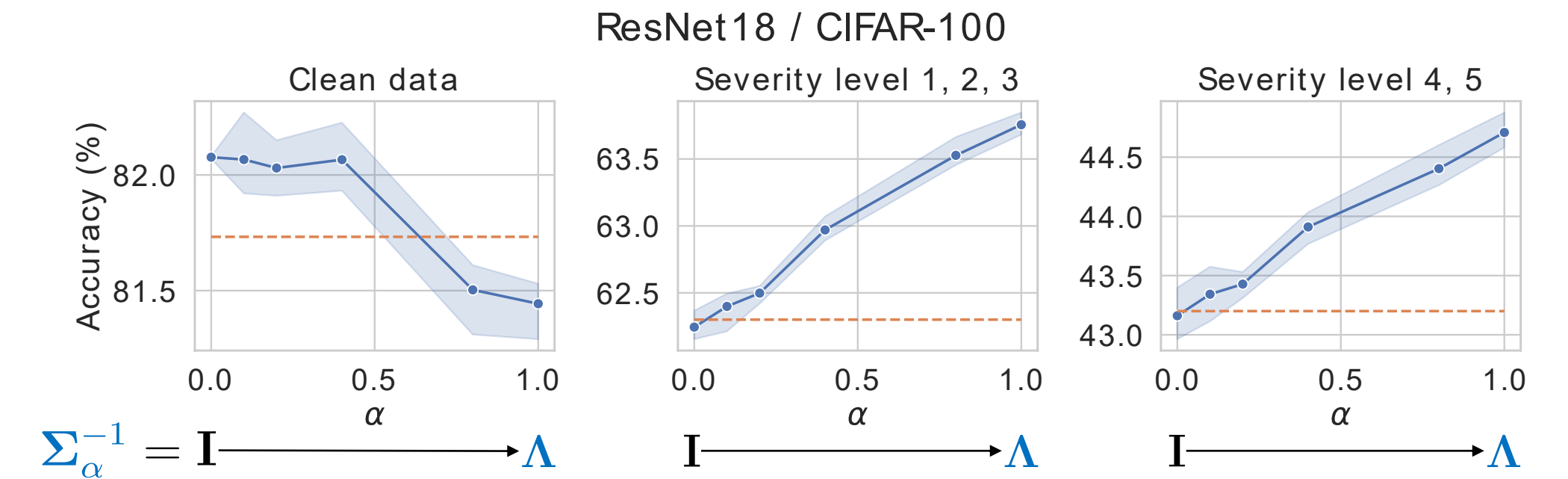


- Use PCA to get the eigenvalues and eigenvectors of the training data: $\{\mathbf{u}_d, \lambda_d\}_{d=1}^D$
- Define the base kernel:

$$\kappa_{\text{PCA}}(\mathbf{s}, \mathbf{s}'; \Sigma_\alpha) = \exp \left(-\frac{1}{2h} (\mathbf{U}^\top \mathbf{s} - \mathbf{U}^\top \mathbf{s}')^\top \Sigma_\alpha^{-1} (\mathbf{U}^\top \mathbf{s} - \mathbf{U}^\top \mathbf{s}') \right)$$

- $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_D]$ is a matrix containing the **eigenvectors** as columns.
- $\Sigma_\alpha^{-1} = (1 - \alpha) \mathbf{I} + \alpha \mathbf{\Lambda}$ where $\mathbf{\Lambda}$ is a diagonal matrix containing the **eigenvalues**.

Lengthscale tuning experiments



- Blue lines show accuracies of FoRDEs, while dotted orange lines show accuracies of Deep ensembles.
- When moving from the identity lengthscales \mathbf{I} to the PCA lengthscales $\mathbf{\Lambda}$:
 - FoRDEs exhibit small performance degradations on clean images of CIFAR-100;
 - while becomes more robust against the natural corruptions of CIFAR-100-C.

References

- [1] F. D'Angelo and V. Fortuin, "Repulsive deep ensembles are Bayesian," Advances in Neural Information Processing Systems, vol. 34, pp. 3451–3465, 2021.
- [2] C. Liu, J. Zhuo, P. Cheng, R. Zhang, and J. Zhu, "Understanding and Accelerating Particle-Based Variational Inference," in International Conference on Machine Learning, 2019.