NeuralFP: Out-of-distribution Detection using Fingerprints of Neural Networks

Wei-Han Lee, Steve Millman, Nirmit Desai, Mudhakar Srivatsa, Changchang Liu

IBM T. J. Watson Research Center

Background:

Design Details

NN Model Activatio

Cloud

Predict on t

Training Data

Edge

Unlabeled De-

Edge devices use neural network models learnt on cloud to predict labels of its data records

> may lead to incorrect predictions especially for OOD data involved in the training process

However, recent efforts in OOD detection

- > require the retraining of the model
- ➤ assume the existence of a certain amount of OOD records



heen Generative Models

reconstruction errors of OOD records are not always larger than in-distribution data (Layer-6 vs Layer-7)

Fingerprinting on the Cloud

layer based on activations

1) obtain activations for each layer by passing

2) construct deep generative models of each

autoencoders applied on training data

training data through neural network model

compute the reconstruction errors of these

3) different layers have various capability in distinguishing OOD records (Layer-1 vs Layer-3)

Model Fingerprint

Reconstruction Error

Integration Stra



Discussion and Future Directions

NeuralFP can serve as a key technique of detecting outliers in practical edge computing scenarios.

- 1) NeuralFP extracts representative information of the training set by constructing non-linear fingerprints of neural network models.
- 2) NeuralFP successfully distinguishes the differences between in-distribution data and OOD data, through carefully integrate the fingerprints across multiple layers
- 3) We have verified the effectiveness of NeuralFP on multiple real-world datasets, showed its advantages over existing detection methods, and provided useful guidelines for parameter selection in practice

Future Directions

- 1) Investigate the performance of NeuralFP in detecting various types of adversarial outliers
- 2) Integrate other deep generative models such as Generative Adversarial Networks (GAN)



pass a given data through the neural network model to compare with the stored model fingerprints for determining abnormality.

Specifically, a data record \boldsymbol{x}^* would be classified as Outlier if its reconstruction error $e_l^* = \mathcal{L}(a_l(x^*), \hat{a}_l(x^*))$ satisfies $\exists l, e_l^* > \tau_l \text{ or } e_l^* < \mu_l$ One-out integration strategy