Why Normalizing Flows Fail to Detect Out-of-Distribution Data

Polina Kirichenko

Pavel Izmailov
Oracle Pavel Izmailov

Andrew Gordon Wilson

- Normalizing flows often surprisingly fail to distinguish between in- and out-of-distribution data: *why*?
- Flows learn local pixel correlations and generic image-to-latent-space transformations not specific to the target image dataset
- Modifying the architecture of flow coupling layers, we can bias the flow towards learning the semantic structure of the target data, improving OOD detection



Out-of-distribution detection

- Goal: detect inputs that are different from the training data
- Important in many safety-critical applications
- Likelihood-based generative models are attractive for OOD detection: we can reject inputs that get low likelihood

Normalizing flows

- Deep generative models based on invertible neural networks
- We can compute density of the data exactly via change of variables

$$z \sim p_Z$$
 $x = f_{\theta}^{-1}(z), \quad p(x) = p_Z(f_{\theta}(x)) \cdot \left| \det\left(\frac{\partial f_{\theta}}{\partial x}\right) \right|$

Inductive biases

- Maximum likelihood training encourages solutions which concentrate all mass on training data, i.e. overfit
- The likelihood assignment outside training data will be determined by inductive biases of the model



Latent space

• There is a direct correspondence between image and latent representation coordinates, the input shape is often visible in latent representation



 $y_{\rm id} = x_{\rm id}$

Inductive biases of flows

• Coupling layers:

$$f_{\text{aff}}^{-1}(x_{\text{id}}, x_{\text{change}}) = (y_{\text{id}}, y_{\text{change}}),$$

• x is split into x_{id} and x_{change} using masking:



What are coupling layers doing?

- Affine coupling layer predicts scale and shift which directly model masked pixels
- The objective pushes values *s* to be large while keeping the norm ||z|| small: as a result (-t) is approximating masked input and *s* represents the confidence of approximation



(a) Checkerboard

(b) Checkerboard, OOD







neural network

e.g. ResNet

 $y_{\text{change}} = (x_{\text{change}} + t(x_{\text{id}})) \odot \exp(s(x_{\text{id}}))$

(c) Horizontal

(d) Horizontal, OOD

- Horizontal mask should make it more challenging to predict masked pixels
- However, the layers learn to *co-adapt* and encode information to be used in subsequent layers



Bigger datasets

• We observe similar results for flow trained on ImageNet: coupling layers can predict inputs for both in-distribution and OOD data





(b) ImageNet input, in-distribution



(c) CelebA input, OOD

Changing inductive biases

• We change the masking strategy to avoid coupling layer coadaptation and overfitting to local color correlations



- We introduce low-dimensional bottlenecks to the *st*-networks
- The lower the dimension of the bottleneck, the better the ranking of



Image embeddings

 Flows applied to embeddings extracted with a pre-trained CNN instead of raw pixels can detect OOD inputs

