

Learning Modular Structures That Generalize Out-Of-Distribution

Arjun Ashok, Chaitanya Devaguptapu, Vineeth N Balasubramanian

Indian Institute of Technology, Hyderabad
IITH Main Road, Near NH-65, Sangareddy, Kandi, Telangana, India - 502285
arjun.ashok@cse.iith.ac.in, cs19mtech11025@iith.ac.in, vineethnb@cse.iith.ac.in

Abstract

Out-of-distribution generalization remains to be a key challenge that machine learning must overcome to achieve its eventual aims associated with artificial intelligence. We hypothesize that encouraging models to be inherently compositional can enable them to extrapolate to unseen inputs better. We propose two structural regularizers that enforce a network internally to be a composition of specialist features in every layer, and promote the emergence of reusable features. Preliminary evaluation on two benchmark datasets corroborates the promise of our method.

Introduction

Recent work has uncovered that neural networks that are learned on observational data are often prone to spurious correlations, and rely on shortcuts learning for solving the task instead of modelling the underlying causal mechanism (Geirhos et al. 2020). This leads to them failing to transfer to more challenging testing conditions, such as real-world scenarios. Recent works show that modularity is a useful inductive bias that can lead to better systematic generalization (Goyal et al. 2020; Csordás, van Steenkiste, and Schmidhuber 2021). We seek to understand whether networks can be structurally enforced to prefer modular solutions. In this context, (Zhang et al. 2021) show that a fully-trained network contains sub-networks that are less susceptible to spurious correlations, and introduce a method to extract the structure from a trained network. In contrast, we study whether networks can be regularized to avoid fitting spurious correlations in the first place, during training itself. We introduce objectives that explicitly incorporate the structure of the network and induce modular structures to be formed at every layer of the network. Our method enforces the network to be a *compositional hierarchy of expert modules*, promoting the *specialization* of modules and encouraging the emergence of *reusable* modules. Preliminary results show that architectures trained with our objective exhibit more resilience to domain shifts, with boosts in O.O.D. generalization performance across 2 datasets.

Method

In general, a deep neural network contains several layers of neurons, each serving as a feature for every neuron in the next. Every fundamental sub-function (e.g. a single convolutional filter) in the network is associated with separate neurons that arise out of transformations (e.g. dot product) of the function with the input. Our aim is to promote every internal feature to be modular, encouraging the underlying functions to specialize in their own tasks, and to only keep features which are reused by multiple functions in the next layer.

We first regularize such that *every feature in the network should be a different composition of the available sub-features*. That is, every feature should fit as few features as necessary, and should differ as much as possible in the features fit. However, directly encouraging this on the weights would unnecessarily constrain the power of the network.

Hence, we use a differentiable probabilistic binary mask π_i over the weights of the network, relaxed by the Gumbel-Sigmoid estimator (Jang, Gu, and Poole 2017). Each value $\pi_i \in [0, 1]$ represents the probability of sampling the respective weight. During training, the mask is binarized once sampled, as $m_i = \{\text{sigmoid}(\pi_i) > 0.5\} \in \{0, 1\}$. Once trained, we obtain deterministic masks by binarizing the final values.

We impose the following regularization of the continuous masks of the weights:

$$S_1(\pi) = \sum_{l=1}^L \sum_{p=1}^{N_l} \left(\sum_{i=1}^{M_p} \pi_i \right)^2 \quad (1)$$

where L denotes the number of layers, N_l - the number of features in the layer, M_p - the number of outgoing weights from feature p .

Note that we minimize the square of the sum of sampling probabilities of weights outgoing from **each feature** in the current layer, allowing it to be *fit sparsely* by only a few required features from above. Consequently, this promotes a competition among the features in the next layer for the current feature, encouraging them to fit a minimally overlapping set of features from the current layer, leading to each of the former specializing in their underlying function.

Although this objective would encourage specialization, every feature in the current layer may not be necessary, as extra features may correspond to unnecessary functions. The network must automatically be able to decide how many features to keep. However, constructing an objective that can be used to restrict the number of features in a layer is non-trivial, since in the worst case, every feature may be necessary for the task at hand.

Here, we hypothesise that *the necessary features are those that are reused by multiple specialist functions above*. Consequently, we regularize to preserve only those features that have a large number of outgoing weights sampled with high-probability, discarding features not fit by multiple predictors above. We enforce this through the following objective:

$$S_2(\pi) = \sum_{l=1}^L \sum_{p=1}^{N_l} \sqrt{\sum_{i=1}^{M_p} \pi_i^2} \quad (2)$$

This term is inspired from that of group lasso regularization (Kim and Xing 2020); applying this term can effectively zero out the masks of *all* the outgoing weights of some features. Unlike that of group lasso that regularizes the weights and can have overlapping groups, we apply it on the masks and do not have any overlapping groups. Importantly, in our context, the network is regularized to only keep features that are well-used by multiple expert predictors above.

Our final objective is, therefore,

$$L = \ell(\theta) + R(\theta) + \alpha * S_1(\pi) + \beta * S_2(\pi)$$

where ℓ is the loss function used for the task (eg. cross entropy), R being a general regularizer (eg. L_2), and α & β , the weights of each of our regularization terms.

Unlike canonical regularizers like L_2 norm, our regularizers explicitly incorporate the notion of a feature, directly taking the structure of every layer of the network into consideration. As a result, applying the two proposed objectives together would result in a mixture of experts(MoE)-like structure in every layer of the network, with only specialist features present in the network after training.

Preliminary Results

We present preliminary results of our method on two benchmark O.O.D. generalization datasets - Colored MNIST (C-MNIST) and Rotated MNIST (R-MNIST). Each dataset is artificially biased in such a way that in the training dataset, a certain degree of correlation is induced between spurious variables and the class label. In the test dataset, the correlation is reversed. The goal of O.O.D. generalization is to encourage the model to fit the invariant features, and ignore other correlated variables, training and validating only on in-distribution data.

Our method is versatile, and can be used on top of any algorithm. Here, we apply our method on top of empirical

Architecture	Method	C-MNIST	R-MNIST
CNN	ERM	35.23	96.5
	ERM + modReg	38.20	96.7
	IRM	67.69	97.3
	IRM + modReg	71.88	98.1
MLP	ERM	34.27	94.45
	ERM + modReg	36.91	95.43
	IRM	72.58	97.4
	IRM + modReg	75.59	97.9

Table 1: Results of the proposed method on multiple architectures, across datasets.

risk minimization(ERM), the standard approach to machine learning problems, and invariant risk minimization (IRM) (Arjovsky et al. 2020), a method estimates nonlinear, invariant, causal predictors from multiple training environments.

Preliminary results shown in table 1 verify the effectiveness of our method. Our method gives consistent gains across the two datasets and architectures considered. In particular, our gains give considerable boosts in the heavily biased C-MNIST dataset, and also improves performance in the R-MNIST dataset in which existing methods have reached their potential. Further implementation details are deferred to the supplementary material.

Future Work

Although our current formulation does not need multiple environments for training, explicitly incorporating multiple domains when available would be an interesting future direction. Further, we also plan to take our method forward and evaluate on larger architectures such as ResNets, as well as on top of other existing O.O.D. generalization methods.

References

- Arjovsky, M.; Bottou, L.; Gulrajani, I.; and Lopez-Paz, D. 2020. Invariant Risk Minimization. arXiv:1907.02893.
- Csordás, R.; van Steenkiste, S.; and Schmidhuber, J. 2021. Are Neural Nets Modular? Inspecting Functional Modularity Through Differentiable Weight Masks. In *ICLR*.
- Geirhos, R.; Jacobsen, J.-H.; Michaelis, C.; Zemel, R.; Brendel, W.; Bethge, M.; and Wichmann, F. A. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11): 665–673.
- Goyal, A.; Lamb, A.; Sodhani, S.; Hoffmann, J.; Levine, S.; Bengio, Y.; and Scholkopf, B. 2020. Recurrent Independent Mechanisms. In *ICLR*.
- Jang, E.; Gu, S.; and Poole, B. 2017. Categorical Reparameterization with Gumbel-Softmax. arXiv:1611.01144.
- Kim, S.; and Xing, E. P. 2020. Tree-Guided Group Lasso for Multi-Task Regression with Structured Sparsity. In *ICML*.
- Zhang, D.; Ahuja, K.; Xu, Y.; Wang, Y.; and Courville, A. 2021. Can Subnetwork Structure be the Key to Out-of-Distribution Generalization? In *ICML*.