



Learnable Embedding Space for Efficient Neural Architecture Compression

Shengcao Cao*, Xiaofang Wang* and Kris M. Kitani

Peking University, Carnegie Mellon University



Compressed Architecture Search

- ❖ **Goal:** automatically search for a small and accurate network architecture based on a given large network architecture
- ❖ **Bottleneck:** the need to repeatedly evaluate different architectures
- ❖ **Proposal:** a learnable embedding space over the domain of network architectures that can be used to generate a priority ordering of architectures for evaluation

Formulation

Reward Function

$$x^* = \arg \max_{x \in \mathcal{X}} f(x)$$

Domain of Architectures

$$f(x) = C(x) (2 - C(x)) \cdot \frac{A(x)}{A(x_{\text{teacher}})}$$

Validation Accuracy

$$C(x) = 1 - \frac{\# \text{params}(x)}{\# \text{params}(x_{\text{teacher}})}$$

Compression Ratio

Search with Bayesian Optimization (BO)

- ❖ Gaussian Process (GP) Prior: $\mu(\cdot) : \mathcal{X} \mapsto \mathcal{R}$ Mean Function, $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \mapsto \mathcal{R}$ Kernel Function

- ❖ At architecture search step t :

Posterior distribution $p(f(x) | f(x_{1:t})) \sim \mathcal{N}(\mu_t(x), \sigma_t^2(x))$

- ❖ Select the next architecture for evaluation

Maximize the acquisition function $\text{EI}_t(x) = \mathbb{E}_t[\max(0, f(x) - f_t^*)]$

How to define the kernel function?

Learnable Embedding Space

$$h(\cdot; \theta) : \mathcal{X} \mapsto \mathcal{R}^d$$

Architecture Embedding

$$k(x, x'; \theta) = \exp\left(-\frac{\|h(x; \theta) - h(x'; \theta)\|^2}{2\sigma^2}\right)$$

RBF Kernel

1. Layer Type: a one-hot vector
2. Layer Attributes: filter size, stride ...
3. Layer Connectivity

Objective Function

$$L(\theta) = -\frac{1}{|D|} \sum_{i: x_i \in D} \log p(f(x_i) | f(D \setminus x_i); \theta)$$

Posterior Probability

Learn a θ such that the function f is consistent with the GP prior.

$$D = \{x_1, \dots, x_t\}$$

At architecture step t , we have t evaluated architectures.

Algorithm Sketch

Algorithm 1 Neural Architecture Search with Bayesian Optimization

Input: Number of steps T . Number of kernels K . Teacher network x_{teacher} .
 Randomly sample K architectures x_1^k, \dots, x_K^k from the search space defined based on x_{teacher} .
 Initialize the set of evaluated architectures $D = \emptyset$.
for $k = 1, \dots, K$ **do**
 Evaluate the K architectures x_1^k, \dots, x_K^k .
 $D = D \cup \{x_1^k, \dots, x_K^k\}$.
for $t = 1, \dots, T$ **do**
 Randomly initialize the weights of kernel k , denoted as θ^k .
 Randomly sample a subset of D , denoted as D^k .
 Learn θ^k on D^k by minimizing the negative log posterior probability.
 Compute the posterior distribution conditioned on the architectures in D^k with kernel k .
 Maximize the acquisition function and denote the solution as $x_{k,t+1}^*$.
end for
end for
 Return the best architecture in D as the solution.

Search Space

- ❖ Layer Removal
- ❖ Layer Shrinkage
- ❖ Add Skip Connections

Code



Compression Results

Model	Teacher	Accuracy	#Params	Ratio	Times	$f(x)$
CIFAR-100	VGG-19	73.71%	20.09M	-	-	-
	Random Search	68.17%	2.83M	0.8593	8.04x	0.9046
	Ours	71.41%	2.61M	0.8699	7.99x	0.9518
ResNet-18	Teacher	78.68%	11.23M	-	-	-
	Random Search	69.86%	1.26M	0.8878	10.10x	0.8752
	Ours	73.85%	1.87M	0.8335	6.01x	0.9123
ResNet-34	Teacher	78.71%	21.33M	-	-	-
	Random Search	72.33%	3.61M	0.8308	5.95x	0.8924
	Ours	74.05%	3.18M	0.8508	6.88x	0.9192

Model	Teacher	Accuracy	#Params	Teacher	Accuracy	#Params
CIFAR-100	VGG-19	71.41%	2.61M	ResNet-18	73.68%	1.87M
	ShuffleNet	68.45%	0.23M	ResNet-34	73.68%	2.36M

Comparison to ShuffleNet [Zhang et al.]

Compression Results on CIFAR-100 Choice of the Objective Function

Architecture Visualization

