# Probabilistic Neural Architecture Search

Chen Zhao

## 1 Motivations and Contributions

- The trade-off between high computational cost and high memory cost has resulted in the usage of different types of surrogates (or proxies) during architecture search, architecture surrogates and dataset surrogates.

    - Architecture surrogates are small versions of the final network that are cheaper to store in memory and faster to train.
    - Dataset surrogates are datasets that can be used as a proxy to perform a quicker search. For example, one could perform architecture search on CIFAR-10 and "transfer" the resulting architecture (with some modifications, e.g. changing the number of filters) to ImageNet.

- While good performance on a surrogate task does not guarantee good performance on the final task (i.e. fully-sized network on the target dataset), the vast majority of architecture search methods use at least one, if not both, of these surrogates.

- The proposed procedure PARSEC (Probabilistic neural ARchitecture SEarCh) is memory-efficient.

- PARSEC transfers probability distributions over architectures learnt on small surrogates to larger networks and datasets, enabling us to further reduce the computational cost.

- PARSEC is computationally efficient and can be run in less than a day on a single GPU on CIFAR-10.

- PARSEC outperforms other methods that consider the same architecture search space, while drastically reducing the search time.

## 2 Architecture Search Space

In this work, we consider the same architecture search space as in DARTS. Neural network architectures are obtained by stacking a recurrent convolutional unit, denoted as cell. A cell is defined as a directed acyclic graph of $N$ ordered

nodes, $\{\mathbf{z}_1, ..., \mathbf{z}_N\}$. An intermediate node, $\mathbf{z}_k$, is defined as the sum of two of the previous nodes, $\mathbf{z}_i$ and $\mathbf{z}_j$ with $i, j < k$, after applying primitive operations $o_{i,k}$ and $o_{j,k}$, respectively:

$$\mathbf{z}_k = o_{i,k}(\mathbf{z}_i) + o_{j,k}(\mathbf{z}_j) \tag{1}$$

## 3  Deterministic Approach to Architecture Search

It consider an over-parametrized parent network containing all possible paths between nodes:

$$\mathbf{z}_k = \sum_{i,p} \alpha_{i,k}^p \cdot o_{i,k}^p(\mathbf{z}_i) \quad \text{where} \quad \sum_p \alpha_{i,k}^p = 1 \tag{2}$$

where indices $i$ and $p$ run over all possible inputs and operations for node $k$ respectively. $o_{i,k}^p$ denotes primitive operation $p$ on input node $i$ to contribute to output node $k$, and $\alpha_{i,k}^p$ is the weight of the path.

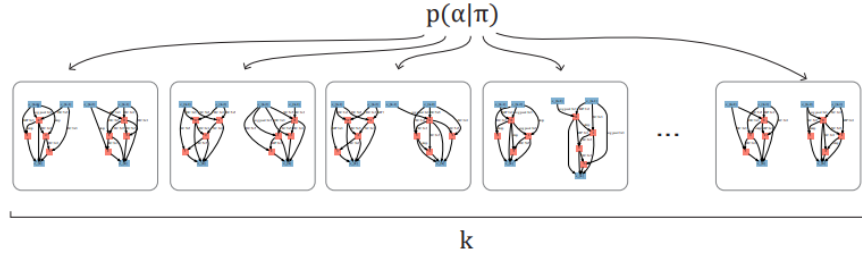## 4  Probabilistic Neural Architecture Search



Figure 1:  Pictorial representation of architecture sampling. Each sample consists of a normal and a reduction cell, which are shown side-to-side.

In PARSEC, child networks from the search space are directly obtained as samples from the specified architecture distribution, and thus neural architecture is reduced to inferring the probability distribution $p(\alpha|\pi)$ of high-performing architectures on the task at hand. Given a supervised task and denoting with $\mathbf{y}$ the targets and with $\mathbf{X}$ the input features, this is achieved by optimizing the continuous prior hyper-parameters $\pi$ through an empirical Bayes Monte Carlo procedure. We optimize

$$p(\mathbf{y}|\mathbf{X}, \mathbf{v}, \pi) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{v}, \alpha) p(\alpha|\pi) d\alpha \tag{3}$$

2

with respect to the network weights $\mathbf{v}$, which are shared across all child architectures (i.e., samples), and the prior hyper-parameters $\pi$, i.e., the architecture parameters.

We assume independent categorical distributions over each input/operation pair. Specifically, for normal and reduction cells we set

$$p(\boldsymbol{\alpha}^{(\mathrm{n})}|\boldsymbol{\pi}^{(\mathrm{n})}) = \prod_{n=1}^{N}\prod_{i=1}^{2}\mathrm{Cat}(\boldsymbol{\alpha}_{n,i}^{(\mathrm{n})}|\boldsymbol{\pi}_{n,i}^{(\mathrm{n})}),$$

$$p(\boldsymbol{\alpha}^{(\mathrm{r})}|\boldsymbol{\pi}^{(\mathrm{r})}) = \prod_{n=1}^{N}\prod_{i=1}^{2}\mathrm{Cat}(\boldsymbol{\alpha}_{n,i}^{(\mathrm{r})}|\boldsymbol{\pi}_{n,i}^{(\mathrm{r})}),$$

where $n$ runs over the nodes, $i$ runs over the inputs of each node (in our search space each node has two inputs), $\pi_{n,i}^{(n)}$ and $\pi_{n,i}^{(r)}$ are the vectors of probabilities for all possible incoming node/operation pairs and finally, we introduced $\alpha^{(\cdot)} = \{\alpha_{n,i}^{(\cdot)}\}$ and $\pi^{(\cdot)} = \{\pi_{n,i}^{(\cdot)}\}$.

# 5 Importance-Weighted Monte Carlo empirical Bayes

We develop an importance weighted EB procedure for jointly optimizing $\pi$ and $v$. We begin by introducing the following estimator:

$$p(\boldsymbol{y}|\mathbf{X},\boldsymbol{v},\boldsymbol{\pi}) = \int p(\boldsymbol{y}|\mathbf{X},\boldsymbol{v},\boldsymbol{\alpha})p(\boldsymbol{\alpha}|\boldsymbol{\pi})\mathrm{d}\boldsymbol{\alpha}$$

$$\approx \frac{1}{K}\sum_{k}p(\boldsymbol{y}|\mathbf{X},\boldsymbol{v},\boldsymbol{\alpha}_k),$$

with $\alpha_k \sim p(\alpha|\pi)$. Note that the gradients can be written as

$$\nabla_{\boldsymbol{v},\boldsymbol{\pi}}\log p(\boldsymbol{y}|\boldsymbol{X},\boldsymbol{v},\boldsymbol{\pi}) = \frac{1}{p(\boldsymbol{y}|\boldsymbol{X},\boldsymbol{v},\boldsymbol{\pi})}\int \nabla_{\boldsymbol{v},\boldsymbol{\pi}}p(\boldsymbol{y},\boldsymbol{\alpha}|\boldsymbol{X},\boldsymbol{v},\boldsymbol{\pi})\mathrm{d}\boldsymbol{\alpha}$$

$$= \frac{1}{p(\boldsymbol{y}|\boldsymbol{X},\boldsymbol{v},\boldsymbol{\pi})}\int \nabla_{\boldsymbol{v},\boldsymbol{\pi}}\log p(\boldsymbol{y},\boldsymbol{\alpha}|\boldsymbol{X},\boldsymbol{v},\boldsymbol{\pi})p(\boldsymbol{y},\boldsymbol{\alpha}|\boldsymbol{X},\boldsymbol{v},\boldsymbol{\pi})\mathrm{d}\boldsymbol{\alpha}$$

Finally, taking an importance sampling estimator to the expectation, we obtain tractable gradient estimators for $\theta$ and $\pi$:

$$\nabla_{\boldsymbol{v}} \log p(\boldsymbol{y}|\mathbf{X}, \boldsymbol{v}, \boldsymbol{\pi}) \quad \approx \quad \sum_{k=1}^{K} \tilde{\omega}_k \nabla_{\boldsymbol{v}} \log p(\boldsymbol{y}|\mathbf{X}, \boldsymbol{v}, \boldsymbol{\alpha}_k) \triangleq \tilde{\nabla}_{\boldsymbol{v}},$$

$$\nabla_{\boldsymbol{\pi}} \log p(\boldsymbol{y}|\mathbf{X}, \boldsymbol{v}, \boldsymbol{\pi}) \quad \approx \quad \sum_{k=1}^{K} \tilde{\omega}_k \nabla_{\boldsymbol{\pi}} \log p(\boldsymbol{\alpha}_k|\boldsymbol{\pi}) \triangleq \tilde{\nabla}_{\boldsymbol{\pi}},$$

where $\alpha_k \sim p(\alpha|\pi)$ and $\tilde{\omega}_k = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{v}, \alpha_\mathbf{k})}{\sum_j p(\mathbf{y}|\mathbf{X}, \mathbf{v}, \alpha_\mathbf{j})}$.

---

**Algorithm 1** Importance weighted Monte-Carlo EB algorithm used for joint training of the network and architecture parameters.

---

1: Define train and search splits of the data
2: $\boldsymbol{\theta}, \boldsymbol{\pi} \leftarrow$ Initial values
3: Initialize $\boldsymbol{\theta}$ and $\boldsymbol{\pi}$
4: **repeat**
5:     Sample batch $(\mathbf{X}^{(s)}, \boldsymbol{y}^{(s)})$ from search set
6:     **for** $k \in \{1, \ldots, K\}$ **do**
7:         Sample architecture $\boldsymbol{\alpha}_k$ from $p(\boldsymbol{\alpha}|\boldsymbol{\pi})$
8:         $\omega_k \leftarrow p(\boldsymbol{y}^{(s)}|\mathbf{X}^{(s)}, \boldsymbol{v}, \boldsymbol{\alpha})$
9:     **end for**
10:     Compute normalized weights $\tilde{\omega}_k = \frac{\omega_k}{\sum_{j=1}^{K} \omega_j}$
11:     Sample batch $(\mathbf{X}^{(s)}, \boldsymbol{y}^{(t)})$ from train set
12:     **for** $k \in \{1, \ldots, K\}$ **do**
13:         $g_{\boldsymbol{v},k} \leftarrow \nabla_{\boldsymbol{v}} \log p(\boldsymbol{y}^{(t)}|\mathbf{X}^{(t)}, \boldsymbol{v}, \boldsymbol{\alpha}_k)$
14:         $g_{\boldsymbol{\pi},k} \leftarrow \nabla_{\boldsymbol{\pi}} \log p(\boldsymbol{\alpha}_k|\boldsymbol{\pi})$
15:     **end for**
16:     Update $\boldsymbol{v}$ by $\tilde{\nabla}_{\boldsymbol{v}} = \sum_k \tilde{\omega}_k g_{\boldsymbol{w},k}$
17:     Update $\boldsymbol{\pi}$ by $\tilde{\nabla}_{\boldsymbol{\pi}} = \sum_k \tilde{\omega}_k g_{\boldsymbol{\pi},k}$
18: **until** Number of epochs is reached or convergence is achieved
19: **Return:** $(\boldsymbol{\theta}, \boldsymbol{\pi})$

---