

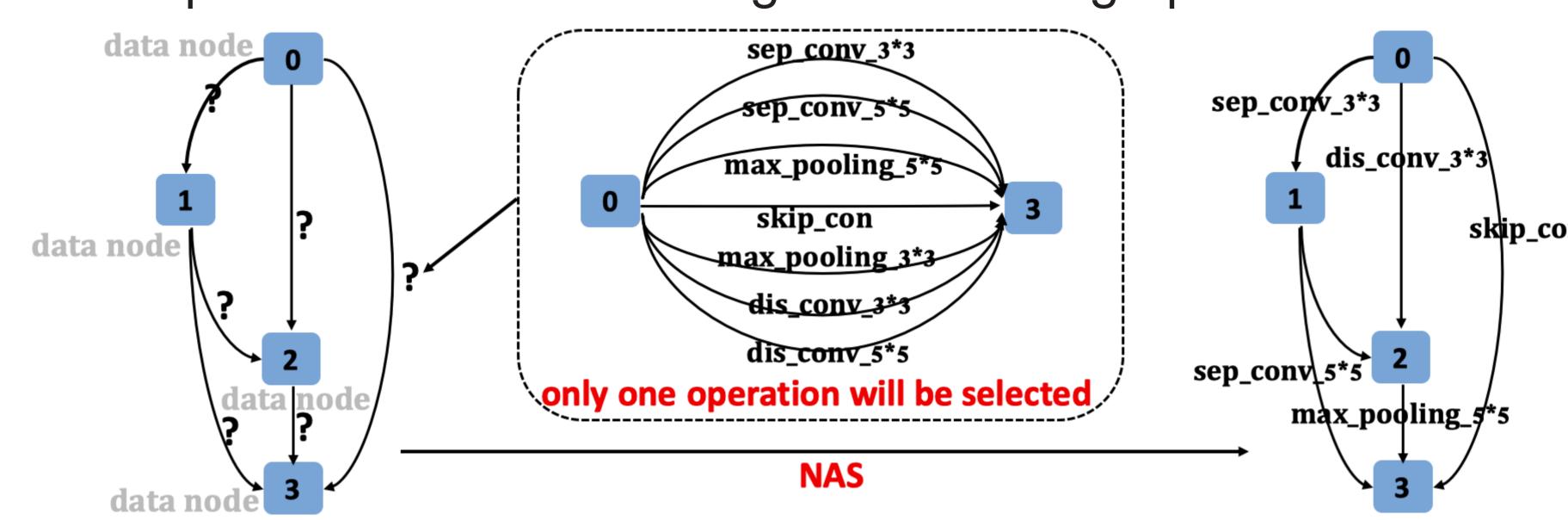
Theory-Inspired Path-Regularized Differential Network Architecture Search

nomepage nomepage

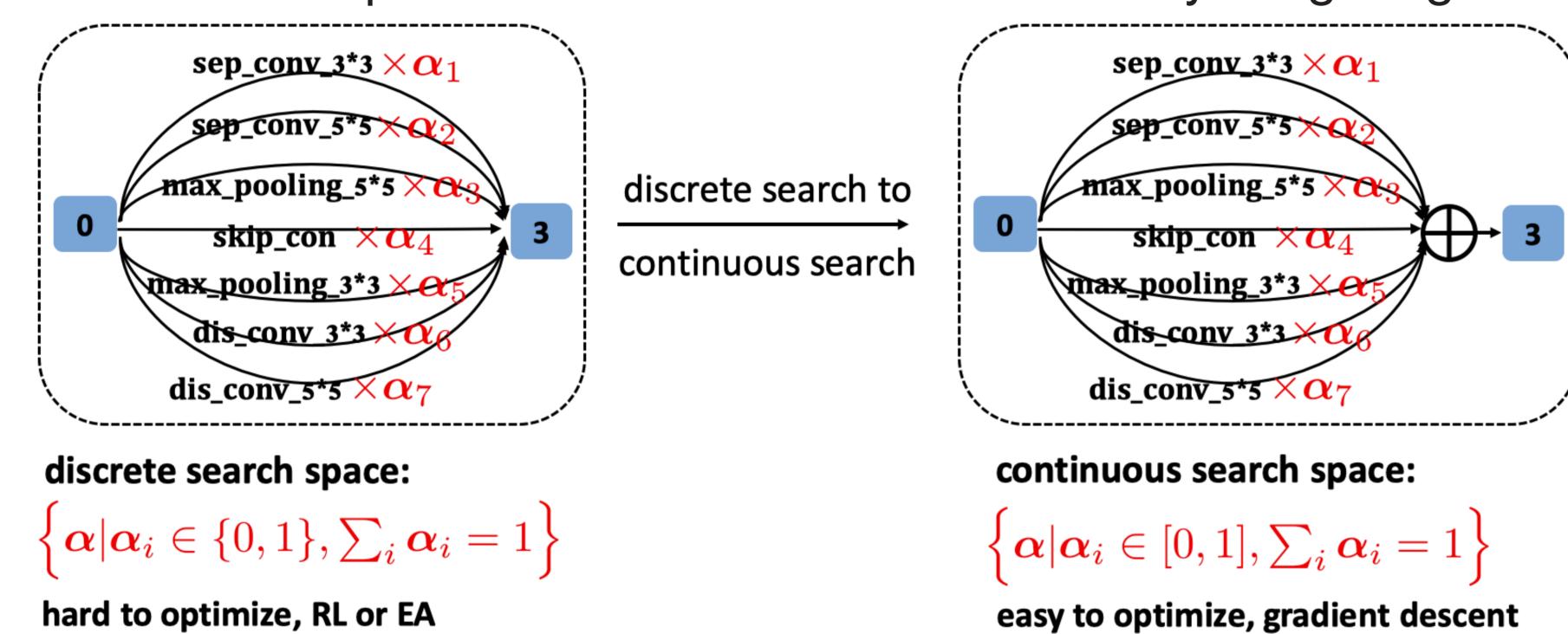
Pan Zhou, Caiming Xiong, Richard Socher, Steven C.H. Hoi Salesforce Research {pzhou, cxiong, rsocher, shoi}@salesforce.com

Problem Setup

What is NAS (network architecture search)? it aims to automatically select a proper operation from an operation set for each edge in a dense graph.



DARTS converts discrete operation selection into continuously weighting a set of operations.

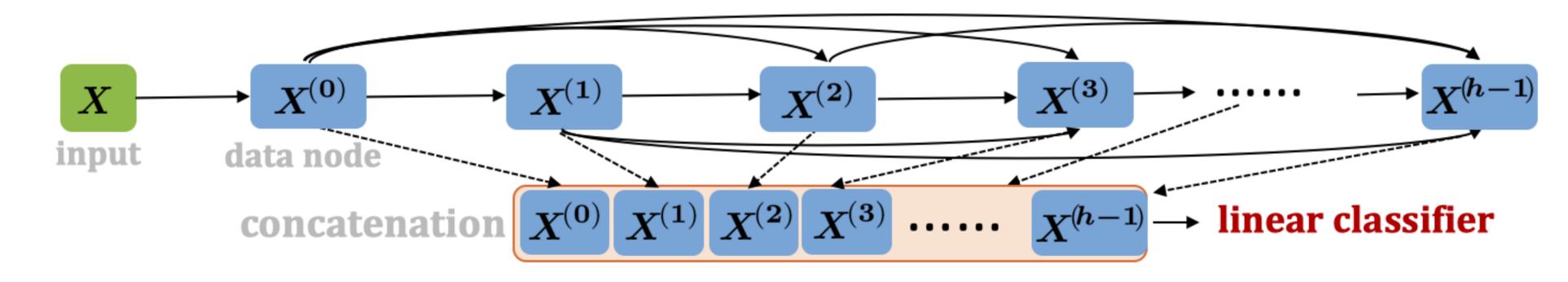


Observations: dominated skip connections in the architectures selected by DARTS

Problem:

- ▶ 1) why DARTS prefers to select so many skip connections?
- ▶ 2) how to avoid the dominated skip connections?

Formulations of DARTS



Dense directed graph via connecting current node with all previous nodes

$$m{X}^{(l)} = \sum_{s=0}^{l-1} \left[m{lpha}_{s,1}^{(l)} \mathsf{zero}(m{X}^s) + m{lpha}_{s,2}^{(l)} \mathsf{skip}(m{X}^s) + m{lpha}_{s,3}^{(l)} \mathsf{conv}(m{W}_s^{(l)}; m{X}^{(s)})
ight] \in \mathbb{R}^{m imes p}$$

where $\alpha_{s,1}^{(l)}$, $\alpha_{s,2}^{(l)}$ and $\alpha_{s,3}^{(l)}$ respectively denote the weights of zero operation zero(\mathbf{X}^s) = 0, skip connection skip(\mathbf{X}^s) = \mathbf{X}^s and convolution conv($\mathbf{W}_s^{(l)}$; $\mathbf{X}^{(s)}$).

Prediction by feeding the features of all nodes into a linear classifier

$$u_i = \sum_{s=0}^{h-1} \langle \boldsymbol{W}_s, \boldsymbol{X}_i^{(s)} \rangle \in \mathbb{R}$$

DARTS Model:

optimize network parameter \boldsymbol{W} on training data $\min_{\alpha} F_{\text{val}}(\boldsymbol{W}^*(\alpha), \alpha)$, optimize architecture parameter $\boldsymbol{\alpha}$ on validation data

where the squared loss $F(\mathbf{W}, \boldsymbol{\beta}) = \frac{1}{2n} \sum_{i=1}^{n} (u_i - y_i)^2$.

Main Results

Gradient descent for optimization:

inner optimization: $\mathbf{W}_{s}^{(l)}(k+1) = \mathbf{W}_{s}^{(l)}(k) - \eta \nabla_{\mathbf{W}_{s}^{(l)}(k)} F_{\text{train}}(\mathbf{W}, \alpha)$ outer optimization: $\alpha_{s}^{(l)}(k+1) = \alpha_{s}^{(l)}(k) - \eta \nabla_{\alpha_{s}^{(l)}(k)} F_{\text{val}}(\mathbf{W}_{s}^{(l)}(k+1), \alpha)$

Convergence of inner optimization. Under proper assumptions, for inner problem, the gradient descent algorithm can enjoy linear convergence rate:

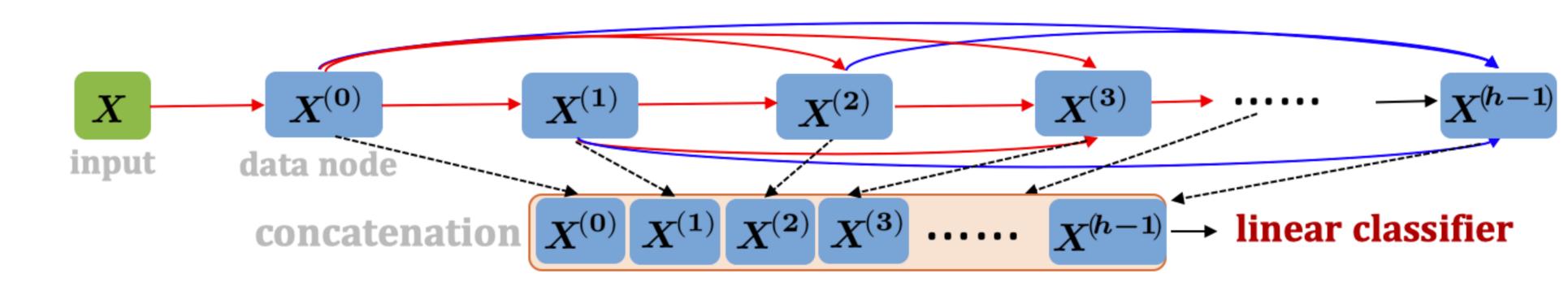
$$F_{\text{train}}(\boldsymbol{W}(k+1), \boldsymbol{lpha}) \leq (1-\lambda) F_{\text{train}}(\boldsymbol{W}(k), \boldsymbol{lpha}) \quad (\forall k \geq 1),$$

where $\lambda = c\eta \sum_{s=0}^{h-2} (\alpha_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$ in which $\alpha_{t,2}^{(s)}$ and $\alpha_{t,3}^{(s)}$ respectively denote weights of skip and convolution connections, is a constant c and η is learning rate.

Remark: convergence rate $1 - \lambda$ relies on weight $\alpha_{t,2}^{(S)}$ of skip connects more heavily weights of convolutions which connect the last node $X^{(h-1)}$

$$\lambda = c\eta \sum_{s=0}^{h-2} \lambda_s$$
 with $\lambda_s = (\alpha_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$

weights of skip connections which do not connect the last node $X^{(h-1)}$



Analysis to outer optimization:

- when skip connections have larger weights, the validation loss can decrease faster, since $\mathbb{E}[F_{\text{train}}(\mathbf{W}), \alpha)] = \mathbb{E}[F_{\text{val}}(\mathbf{W}), \alpha)].$
- as all types of operations between two nodes share a softmax distribution

$$oldsymbol{lpha_{t,i}^{(\mathcal{S})}} = rac{\exp(oldsymbol{eta_{t,i}^{(\mathcal{S})}})}{\sum_{i=1}^{3} \exp(oldsymbol{eta_{t,i}^{(\mathcal{S})}})} \longrightarrow \sum_{i} oldsymbol{lpha_{t,i}^{(\mathcal{S})}} = 1\,,$$

if weight of skip connection becomes larger, other weights become smaller.

Conclusion:

- outer level optimization increases weights of skip connections and reduce other weights
- posterior pruning preserves most skip connections, prunes most other operations

Solution to Alleviate Unfair Competition among Operations

Natural solution: independent gate $g_{ti}^{(S)}$ (Bernoulli distribution) for each operation

Issue 1 of independent gate. By using independent gate, increasing $g_{t,i}^{(s)}$ of any operations can provably reduce or maintain the validation loss.

Solution to Issue 1: group-structured sparsity regularization on gates

► Step 1. use Gumbel trick to produce an approximate Bernoulli variable u

$$u = rac{\exp\left(rac{V_1 + \ln V_2}{ au}
ight)}{1 + \exp\left(rac{V_1 + \ln V_2}{ au}
ight)} pprox ext{Bernoulli}ig(v_2ig), v_1 \sim ext{Uniform}(0,1), v_2 = rac{\exp(oldsymbol{eta}_{t,i}^{(oldsymbol{S})}ig)}{1 + \exp(oldsymbol{eta}_{s,i}^{(oldsymbol{S})}ig)}$$

 $g_{t,i}^{(s)} = a + (b - a)u, \quad g_{t,i}^{(s)} = \min(1, \max(0, g_{t,i}^{(s)}))$

Step 2. rescale to [a,b] (a < 0,b > 1), feed into hard threshold gate

$$\begin{array}{l} \text{if } u \in (0, -\frac{a}{b-a}] \text{ (sparse)}, \\ \text{if } u \in (-\frac{a}{b-a}, \frac{1-a}{b-a}], & \mathbb{P}(\boldsymbol{g}_{t,i}^{(s)} \neq 0) &= \Theta \\ \text{sigmoid} \Big(\boldsymbol{\beta}_{t,i}^{(s)} - \tau \ln \frac{-a}{b} \Big) \Big) \end{array}$$

Solutions to Alleviate Unfair Competition among Operations

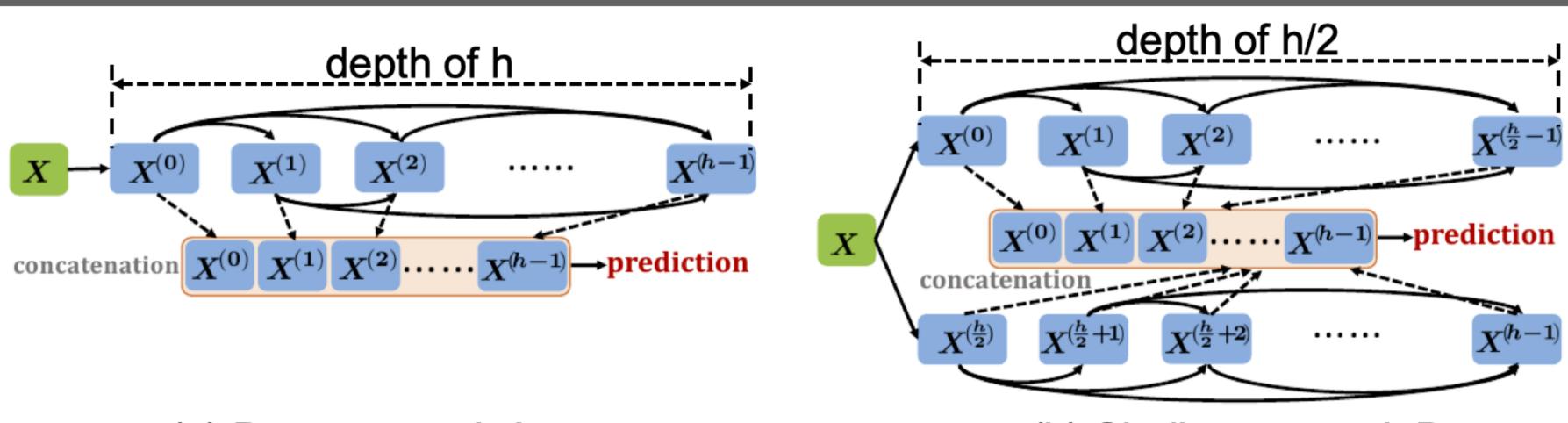
► Step 3. divide operations in network into skip connection and non-skip connection groups, compute their average gate activation probabilities:

$$\mathcal{L}_{\mathsf{skip}}(\boldsymbol{\beta}) = \zeta \sum_{l=1}^{h-1} \sum_{s=0}^{l-1} \Theta\left(\boldsymbol{\beta}_{s,t_{\mathsf{skip}}}^{(l)} - \tau \ln \frac{-a}{b}\right), \ \mathcal{L}_{\mathsf{non-skip}}(\boldsymbol{\beta}) = \frac{\zeta}{r-1} \sum_{l=1}^{h-1} \sum_{s=0}^{l-1} \sum_{1 \leq t \leq r, t \neq t_{\mathsf{skip}}}^{h-1} \Theta\left(\boldsymbol{\beta}_{s,t}^{(l)} - \tau \ln \frac{-a}{b}\right),$$

Step 4. penalize $\mathcal{L}_{\text{skip}}(\beta)$ and $\mathcal{L}_{\text{non-skip}}(\beta)$ independently to avoid competition between skip connection and other type operations.

Issue 2 of independent gate: searching algorithm prefers to select shallow networks due to their faster convergence rate over deep ones.

Convergence Comparison. With proper assumptions, shallow network B can converge faster than the deep network A



(a) Deep network A

(b) Shallow network B

Solution to Issue 2: path-depth-wise regularization

► Step 1. compute probability that all neighboring nodes are connected via parameterized operations

$$\mathcal{L}_{\mathrm{path}}(\boldsymbol{\beta}) = \prod_{l=1}^{h-1} \mathbb{P}_{l,l+1}(\boldsymbol{\beta}) = \prod_{l=1}^{h-1} \sum_{O_t \in \mathcal{O}_p} \Theta \left(\boldsymbol{\beta}_{l,t}^{(l+1)} - \tau \ln \frac{-a}{b}\right).$$
 corresponding activation probability of parameterized operations connected by learnable parameterized operations, e.g. various types of convolutions put data node concatenation $X^{(0)} \times X^{(1)} \times X^{(2)} \times X^{(3)} \times X^{(h-1)}$ linear classifier

ightharpoonup Step 2. encourage deeper model via penalizing small $\mathcal{L}_{path}(\beta)$

Final Model: $\min_{\beta} F_{\text{val}}(W^*(\beta), \beta) + \lambda_1 \mathcal{L}_{\text{skip}}(\beta) + \lambda_2 \mathcal{L}_{\text{non-skip}}(\beta) - \lambda_3 \mathcal{L}_{\text{path}}(\beta),$ $\text{s.t. } W^*(\beta) = \operatorname{argmin}_{W} F_{\text{train}}(W, \beta)$

Experiments

