

Theory-Inspired Path-Regularized Differential Network Architecture Search

Pan Zhou, Caiming Xiong, Richard Socher, Steven C.H. Hoi

Salesforce Research pzhou@salesforce.com

Oct 2020

Outline



Background: what is network architecture search (NAS)

Theoretical analysis: why DARTS often select so many skip connections

Solution: group-structured sparse gate and path-depth-wise regularization

Experiments: higher efficiency and classification accuracy

Conclusion

Outline



Background: what is network architecture search (NAS)

Theoretical analysis: why DARTS often select so many skip connections

Solution: group-structured sparse gate and path-depth-wise regularization

Experiments: higher efficiency and classification accuracy

Conclusion

Background: What Is NAS?



NAS (network architecture search) aims to automatically select a proper operation from an

operation set for each edge in a dense graph



Background: What Is NAS?



Solution: reinforcement learning (RL) and evolutionary algorithms (EA) are used to solve this discrete operation selection problem.

Issue: huge search space

E.g. a graph of 10 nodes has $7C_{10}^2 = 7^{45}$ possible operation selections if the operation set is of size 7 high computational cost (more than 3000 GPU days)



discrete search space: $\left\{ oldsymbol{lpha} | oldsymbol{lpha}_i \in \{0,1\}, \sum_i oldsymbol{lpha}_i = 1
ight\}$

To reduce cost, one often search a small network and then stack several cells to build a large one.

Differential Architecture Search (DARTS)



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



Differential Architecture Search (DARTS)



Since the weights of most operations are not exactly zero, one often needs to prune the operations with small weight.



This posteriors pruning often leads to information loss, as it destroys the learnt architecture.





Observations: dominated skip connections in the architectures selected by DARTS



Problems:

1) why DARTS prefers to select so many skip connections?

2) how to avoid the dominated skip connections?

Outline



Background: what is network architecture search

Theoretical analysis: why DARTS often select so many skip connections

Solution: group-structured sparse gate and path-depth-wise regularization

Experiments: higher efficiency and classification accuracy

Conclusion

Formulations of DARTS



• **Dense directed graph** via connecting current node with all previous nodes

$$\boldsymbol{X}^{(l)} = \sum_{s=0}^{l-1} \left[\boldsymbol{\alpha}_{s,1}^{(l)} \operatorname{zero}(\boldsymbol{X}^s) + \boldsymbol{\alpha}_{s,2}^{(l)} \operatorname{skip}(\boldsymbol{X}^s) + \boldsymbol{\alpha}_{s,3}^{(l)} \operatorname{conv}(\boldsymbol{W}_s^{(l)}; \boldsymbol{X}^{(s)}) \right] \in \mathbb{R}^{m \times p}$$

$$(l = 1, \dots, h-1)$$

where $\alpha_{s,1}^{(l)}, \alpha_{s,2}^{(l)}, \alpha_{s,3}^{(l)}$ respectively denote weights of zero, skip and convolution operations.

• **Prediction** by feeding the features in all layers into a linear classifier

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



Formulations of DARTS



• **Dense directed graph** via connecting current node with all previous nodes

$$\begin{split} \boldsymbol{X}^{(l)} = & \sum_{s=0}^{l-1} \left[\boldsymbol{\alpha}_{s,1}^{(l)} \mathsf{zero}(\boldsymbol{X}^s) + \boldsymbol{\alpha}_{s,2}^{(l)} \mathsf{skip}(\boldsymbol{X}^s) + \boldsymbol{\alpha}_{s,3}^{(l)} \mathsf{conv}(\boldsymbol{W}_s^{(l)}; \boldsymbol{X}^{(s)}) \right] \in \mathbb{R}^{m \times p} \\ & (l = 1, \cdots, h-1) \end{split}$$

where $\alpha_{s,1}^{(l)}, \alpha_{s,2}^{(l)}, \alpha_{s,3}^{(l)}$ respectively denote weights of zero, skip and convolution operations.

• **Prediction** by feeding the features in all layers 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 $oldsymbol{W}$

 $\min_{\alpha} F_{val}(\boldsymbol{W}^*(\alpha), \alpha), \quad \text{s.t. } \boldsymbol{W}^*(\alpha) = \operatorname{argmin}_{\boldsymbol{W}} F_{train}(\boldsymbol{W}, \boldsymbol{\alpha}),$

optimize architecture parameter α

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

• **Gradient descent** for optimization:

 $\begin{array}{ll} \text{inner optimization:} & \boldsymbol{W}_{s}^{(l)}(k+1) = \boldsymbol{W}_{s}^{(l)}(k) - \eta \nabla_{\boldsymbol{W}_{s}^{(l)}(k)} F_{\text{train}}(\boldsymbol{W}, \boldsymbol{\alpha}) \ (\forall l, s) \\ \text{outer optimization:} & \boldsymbol{\alpha}_{s}^{(l)}(k+1) = \boldsymbol{\alpha}_{s}^{(l)}(k) - \eta \nabla_{\boldsymbol{\alpha}_{s}^{(l)}(k)} F_{\text{val}}(\boldsymbol{W}_{s}^{(l)}(k+1), \boldsymbol{\alpha}) \ (\forall l, s) \\ \end{array}$

• **Gradient descent** for optimization:

inner optimization: $\boldsymbol{W}_{s}^{(l)}(k+1) = \boldsymbol{W}_{s}^{(l)}(k) - \eta \nabla_{\boldsymbol{W}_{s}^{(l)}(k)} F_{\mathrm{train}}(\boldsymbol{W}, \boldsymbol{\alpha}) \ (\forall l, s)$

Theorem 1 (Convergence for inner problem, informal).

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

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

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 convolution and skip connections, C is a constant and η is learning rate.

 $zero \times \alpha_{i}^{(s)}$

Theorem 1 (Convergence for inner problem, informal). For inner problem, the gradient descent algorithm can enjoy linear convergence rate:

$$\begin{split} F_{\mathrm{train}}(\boldsymbol{W}(k+1),\boldsymbol{\alpha}) &\leq (1-\lambda) \, F_{\mathrm{train}}(\boldsymbol{W}(k),\boldsymbol{\alpha}) \quad (\forall k \geq 1), \\ \text{where } \lambda &= c\eta \sum_{s=0}^{h-2} (\boldsymbol{\alpha}_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\boldsymbol{\alpha}_{t,2}^{(s)})^2 \text{ with constant } C \text{ and learning rate } \eta \text{ .} \end{split}$$

• Convergence rate $(1 - \lambda)$ depends on the weight $\alpha_{t,2}^{(s)}$ of skip connects more heavily: $\lambda = c\eta \sum_{s=0}^{h-2} \lambda_s$ with $\lambda_s = (\alpha_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$

Theorem 1 (Convergence for inner problem, informal). For inner problem, the gradient descent algorithm can enjoy linear convergence rate:

$$\begin{split} F_{\mathrm{train}}(\boldsymbol{W}(k+1),\boldsymbol{\alpha}) &\leq (1-\lambda) \, F_{\mathrm{train}}(\boldsymbol{W}(k),\boldsymbol{\alpha}) \quad (\forall k \geq 1), \\ \text{where } \lambda &= c\eta \sum_{s=0}^{h-2} (\boldsymbol{\alpha}_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\boldsymbol{\alpha}_{t,2}^{(s)})^2 \text{ with constant } C \text{ and learning rate } \eta \text{ .} \end{split}$$

• Convergence rate $(1 - \lambda)$ depends on the weight $\alpha_{t,2}^{(s)}$ of skip connects more heavily: $\lambda = c\eta \sum_{s=0}^{h-2} \lambda_s$ with $\lambda_s = (\alpha_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$

weights of convolutions which connect the last node $X^{(h-1)}$



Theorem 1 (Convergence for inner problem, informal). For inner problem, the gradient descent algorithm can enjoy linear convergence rate:

$$\begin{split} F_{\mathrm{train}}(\boldsymbol{W}(k+1),\boldsymbol{\alpha}) &\leq (1-\lambda) \, F_{\mathrm{train}}(\boldsymbol{W}(k),\boldsymbol{\alpha}) \quad (\forall k \geq 1), \\ \text{where } \lambda &= c\eta \sum_{s=0}^{h-2} (\boldsymbol{\alpha}_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\boldsymbol{\alpha}_{t,2}^{(s)})^2 \text{ with constant } C \text{ and learning rate } \eta \text{ .} \end{split}$$

• Convergence rate $(1 - \lambda)$ depends on the weight $\alpha_{t,2}^{(s)}$ of skip connects more heavily: $\lambda = c\eta \sum_{s=0}^{h-2} \lambda_s$ 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 $X^{(h-1)}$



Theorem 1 (Convergence for inner problem, informal). For inner problem, the gradient descent algorithm can enjoy linear convergence rate:

$$F_{ ext{train}}(oldsymbol{W}(k+1),oldsymbol{lpha}) \leq (1-\lambda) \, F_{ ext{train}}(oldsymbol{W}(k),oldsymbol{lpha}) \quad (orall k \geq 1),$$
where $\lambda = c\eta \sum_{s=0}^{h-2} (oldsymbol{lpha}_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (oldsymbol{lpha}_{t,2}^{(s)})^2$ with constant c and learning rate η .

• Convergence rate $(1 - \lambda)$ depends on the weight $\boldsymbol{\alpha}_{t,2}^{(s)}$ of skip connects more heavily: $\lambda = c\eta \sum_{s=0}^{h-2} \lambda_s$ with $\lambda_s = (\boldsymbol{\alpha}_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\boldsymbol{\alpha}_{t,2}^{(s)})^2$

weight product gives heavier dependence



• Since training and validation data are drawn from a same distribution, we have

$$\mathbb{E}[F_{\text{train}}(\boldsymbol{W}), \boldsymbol{\alpha})] = \mathbb{E}[F_{\text{val}}(\boldsymbol{W}), \boldsymbol{\alpha})]$$

- When skip connections have larger weights, the validation loss can decrease faster
- Since all types of operations between two nodes share a softmax distribution $zero \times \alpha'$

$$\boldsymbol{\alpha}_{t,i}^{(s)} = \frac{\exp(\boldsymbol{\beta}_{t,i}^{(s)})}{\sum_{i=1}^{3} \exp(\boldsymbol{\beta}_{t,i}^{(s)})} \longrightarrow \sum_{i} \boldsymbol{\alpha}_{t,i}^{(s)} = 1$$

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

• In the outer level, DARTS will increase the weights of skip connections and reduce the weights of other operations.

outer optimization: $\boldsymbol{\alpha}_{s}^{(l)}(k+1) = \boldsymbol{\alpha}_{s}^{(l)}(k) - \eta \nabla_{\boldsymbol{\alpha}_{s}^{(l)}(k)} F_{\mathrm{val}}(\boldsymbol{W}_{s}^{(l)}(k+1), \boldsymbol{\alpha}) \ (\forall l, s)$

• After searching, the **posterior pruning will preserve** most of **skip connections** and **prune** most of **other operations**.



• Our theoretical result can answer the first question:

why DARTS prefers to select so many skip connections?

Outline



Background: what is network architecture search

Theoretical analysis: why DARTS often select so many skip connections

Solution: group-structured sparse gate and path-depth-wise regularization

Experiments: higher efficiency and classification accuracy

Conclusion

One solution: independent gate implemented by Bernoulli distribution for each operation





Theorem 2 (Convergence for inner problem, informal).

When we replace the weights from a softmax distribution with the independent gate, then increasing $g_{t,i}^{(s)}$ of any operations can reduce or maintain the validation loss.

Issue:

independent gate leads to a dense network and thus performance degradation,

since posterior pruning for a compact network prunes operations with non-zero weights.

ales*f*orc



Solution: group-structured sparsity regularization on the stochastic gates

- Step 1. use Gumbel trick to produce an approximate Bernoulli variable \boldsymbol{u}

$$u = \frac{\exp\left(\frac{v_1 + \ln v_2}{\tau}\right)}{1 + \exp\left(\frac{v_1 + \ln v_2}{\tau}\right)} \approx \text{Bernoulli}(v_2), \ v_1 \sim \text{Uniform}(0, 1), \ v_2 = \frac{\exp(\boldsymbol{\beta}_{t,i}^{(s)})}{1 + \exp(\boldsymbol{\beta}_{s,i}^{(t)})}$$

• Step 2. rescale u from [0,1] to [a,b] (a<0, b>1), and feed $\boldsymbol{g}_{t,i}^{(s)}$ into a hard threshold gate

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

 $\mathbf{g}_{t,i}^{(s)} = \begin{cases} 0, & \text{if } u \in (0, -\frac{a}{b-a}], \text{ sparse} \\ \mathbf{g}_{t,i}^{(s)} = \begin{cases} 0, & \text{if } u \in (0, -\frac{a}{b-a}], \text{ sparse} \\ \mathbf{g}_{t,i}^{(s)}, & \text{if } u \in (-\frac{a}{b-a}, \frac{1-a}{b-a}], \\ 1, & \text{if } u \in (\frac{1-a}{b-a}, 1], \end{cases} & \mathbb{P}(\mathbf{g}_{t,i}^{(s)} \neq 0) = \Theta\left(\beta_{t,i}^{(s)} - \tau \ln \frac{-a}{b}\right), \\ \text{gate activation probability} \end{cases}$

where $\boldsymbol{\Theta}$ denotes the sigmoid function.



Solution: group-structured sparsity regularization on the stochastic gates

• Step 3. divide the operations in the cell into two groups, skip connection group and non-skip connection group, and compute their average gate activation probabilities:

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

• Step 4. we penalize these two terms independently to avoid competition between skip connection and other type operations.



Solution: group-structured sparsity regularization on the stochastic gates

Advantages: this solution greatly reduce the skip connections in the selected network.





Issues of independent gates: searching algorithm prefers to select shallow networks due to their faster convergence rate over deep ones.

Theorem 2 (Convergence Comparison between shallow and deep networks, informal) With proper assumptions, shallow network B can converge faster than the deep network A.





Solution: path-depth-wise regularization

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





Solution: path-depth-wise regularization

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

$$\mathcal{L}_{\text{path}}(\boldsymbol{\beta}) = \prod_{l=1}^{h-1} \mathbb{P}_{l,l+1}(\boldsymbol{\beta}) = \prod_{l=1}^{h-1} \sum_{O_t \in \mathcal{O}_p} \Theta\left(\boldsymbol{\beta}_{l,t}^{(l+1)} - \tau \ln \frac{-a}{b}\right)$$

- Step 2. we encourage the selected model to be deep via penalizing small $\mathcal{L}_{ ext{path}}(m{eta})$

Solution: path-depth-wise regularization

Advantages: this solution can search much deeper networks than DARTS.



salesforce

Network Architecture Search Model



• Architecture search model (PR-DARTS):

optimize network parameter ${oldsymbol W}$

$$\min_{\boldsymbol{\beta}} F_{\mathrm{val}}(\boldsymbol{W}^{*}(\boldsymbol{\beta}),\boldsymbol{\beta}) + \lambda_{1}\mathcal{L}_{\mathrm{skip}}(\boldsymbol{\beta}) + \lambda_{2}\mathcal{L}_{\mathrm{non-skip}}(\boldsymbol{\beta}) - \lambda_{3}\mathcal{L}_{\mathrm{path}}(\boldsymbol{\beta}), \ \mathbf{s.t.} \boldsymbol{W}^{*}(\boldsymbol{\beta}) = \operatorname{argmin}_{\boldsymbol{W}} F_{\mathrm{train}}(\boldsymbol{W},\boldsymbol{\beta})$$

optimize architecture parameter lpha

• Advantages:

(1) it avoids unfair competition between skip and non-skip connections $\min_{\boldsymbol{\beta}} F_{\text{val}}(\boldsymbol{W}^{*}(\boldsymbol{\beta}), \boldsymbol{\beta}) + \lambda_{1} \mathcal{L}_{\text{skip}}(\boldsymbol{\beta}) + \lambda_{2} \mathcal{L}_{\text{non-skip}}(\boldsymbol{\beta}) - \lambda_{3} \mathcal{L}_{\text{path}}(\boldsymbol{\beta}), \text{ s.t. } \boldsymbol{W}^{*}(\boldsymbol{\beta}) = \operatorname{argmin}_{\boldsymbol{W}} F_{\text{train}}(\boldsymbol{W}, \boldsymbol{\beta})$

(2) it avoids unfair competition between shallow and deep networks

$$\min_{\boldsymbol{\beta}} F_{\text{val}}(\boldsymbol{W}^{*}(\boldsymbol{\beta}),\boldsymbol{\beta}) + \lambda_{1} \mathcal{L}_{\text{skip}}(\boldsymbol{\beta}) + \lambda_{2} \mathcal{L}_{\text{non-skip}}(\boldsymbol{\beta}) - \lambda_{3} \mathcal{L}_{\text{path}}(\boldsymbol{\beta}), \text{ s.t.} \boldsymbol{W}^{*}(\boldsymbol{\beta}) = \operatorname{argmin}_{\boldsymbol{W}} F_{\text{train}}(\boldsymbol{W},\boldsymbol{\beta})$$

Outline



Background: what is network architecture search

Theoretical analysis: why DARTS often select so many skip connections

Solution: group-structured sparse gate and path-depth-wise regularization

Experiments: higher efficiency and classification accuracy

Conclusion

Experimental Results

Search Time on CIFAR10: much higher search efficiency



Accuracy on CIFAR10 and ImageNet: much smaller classification error





Outline



Background: what is network architecture search

Theoretical analysis: why DARTS often select so many skip connections

Solution: group-structured sparse gate and path-depth-wise regularization

Experiments: higher efficiency and classification accuracy

Conclusion

Conclusion



• Problems:

(1) why DARTS prefers to select so many skip connections?

more skip connections lead to faster convergence speed and thus are selected.

(2) how to avoid the dominated skip connections?

we propose a group-structured sparsity regularization and a path-depth-wise regularization



Thanks !