

Neural Optimizer Search with Reinforcement Learning

Irwan Bello¹ Barret Zoph¹ Vijay Vasudevan¹ Quoc V. Le¹

¹Google Brain

ICLR, 2017/ Presenter: Anant Kharkar

1 Introduction

- Motivation
- Approach

2 Methods

- Domain-Specific Language
- Controller RNN

3 Experiments

- Optimizer Discovery
- Transfer Experiment

4 Summary

Outline

1 Introduction

- Motivation
- Approach

2 Methods

- Domain-Specific Language
- Controller RNN

3 Experiments

- Optimizer Discovery
- Transfer Experiment

4 Summary

Motivation

Classical optimizers:

- SGD
- SGD w/Momentum
- Adam
- RMSProp

Combination of stochastic methods and heuristic approximations

Want to automate process of generating update rules

Produce equation, not just numerical updates

1 Introduction

- Motivation
- Approach

2 Methods

- Domain-Specific Language
- Controller RNN

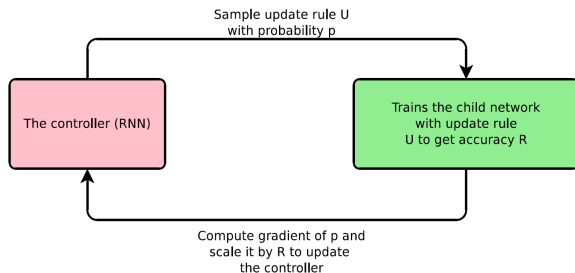
3 Experiments

- Optimizer Discovery
- Transfer Experiment

4 Summary

Approach

- RNN controller produces update rule string
- Controller updated based on performance of optimizer
- RL approach to training



- How to generate update rules? First define space of update rules

- LSTM for numerical updates (Andrychowicz et al., 2016)
 - Equations are more transferrable
- Genetic programming for update equations (Orchard & Wang, 2016)
 - Slow and needs heuristics
- Neural Architecture Search (Zoph & Le, 2017) - seen earlier
 - RNN produces network architecture

Outline

- 1 Introduction
 - Motivation
 - Approach
- 2 **Methods**
 - **Domain-Specific Language**
 - Controller RNN
- 3 Experiments
 - Optimizer Discovery
 - Transfer Experiment
- 4 Summary

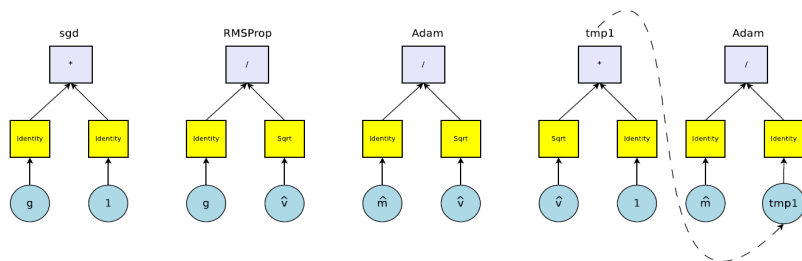
Domain-Specific Language

- Each optimizer has computational graph - binary expression tree

Components:

- 2 operands
- Unary function for each operand
- Binary function to combine

$$\Delta w = \lambda * b(u_1(op_1), u_2(op_2))$$

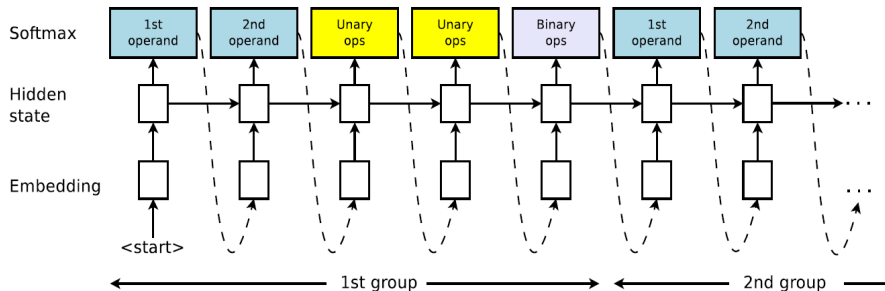


Outline

- 1 Introduction
 - Motivation
 - Approach
- 2 **Methods**
 - Domain-Specific Language
 - **Controller RNN**
- 3 Experiments
 - Optimizer Discovery
 - Transfer Experiment
- 4 Summary

Controller RNN

Trained with Adam



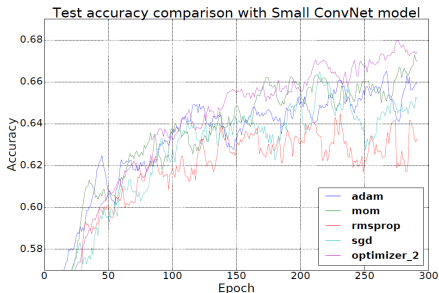
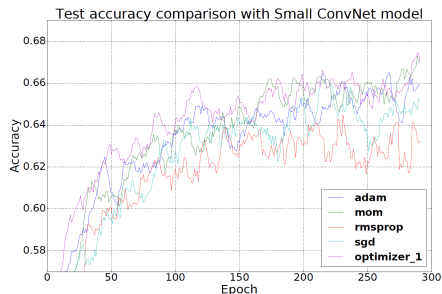
Objective function: $J(\theta) = \mathbb{E}_{\Delta \sim p_{\theta}(\cdot)}[R(\Delta)]$
Optimize reward (accuracy of target model)

- Operands
 - Gradients: $g, g^2, g^3, \text{sign}(g)$
 - Moving averages: $\hat{m}, \hat{v}, \hat{y}, \text{sign}(\hat{m})$
 - Weights: $10^{-4}w, 10^{-3}w, 10^{-2}w, 10^{-1}w$
 - ADAM, RMSProp, 1, small noise
- Unary Functions
 - $x, -x, e^x, \log|x|, \text{clip}, \text{drop}, \text{sign}$
- Binary Functions
 - $x + y, x - y, x * y, \frac{x}{y+\epsilon}$
- Optimizers tested on 3x3 ConvNet (32 filters) for 5 epochs
- Favors early progress

Outline

- 1 Introduction
 - Motivation
 - Approach
- 2 Methods
 - Domain-Specific Language
 - Controller RNN
- 3 Experiments
 - **Optimizer Discovery**
 - Transfer Experiment
- 4 Summary

Optimizer Discovery



Recurring element:

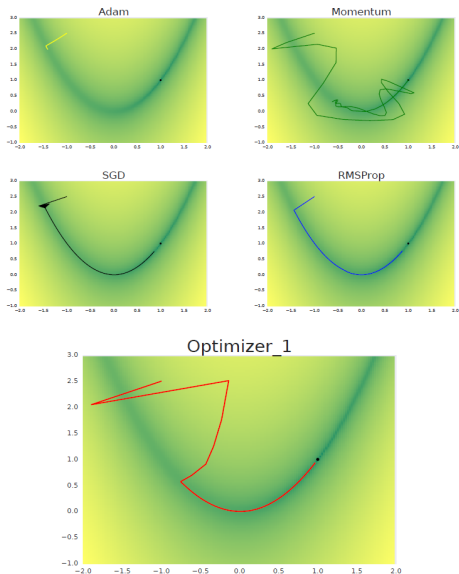
$$e^{\text{sign}(g) * \text{sign}(m)} * g$$

- If $\text{sign}(g)$ agrees with running average, scale e - g keeps decreasing
- Else scale $\frac{1}{e}$ - gradient direction changed

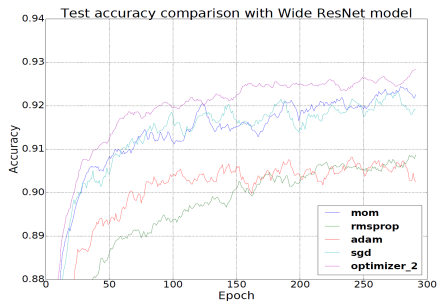
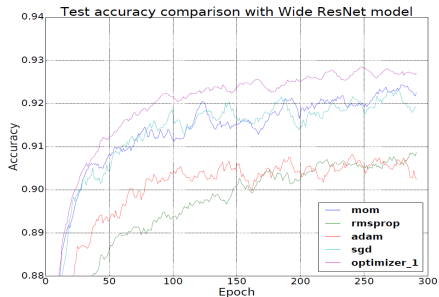
Outline

- 1 Introduction
 - Motivation
 - Approach
- 2 Methods
 - Domain-Specific Language
 - Controller RNN
- 3 Experiments
 - Optimizer Discovery
 - Transfer Experiment
- 4 Summary

Rosenbrock Function



Wide ResNet (Zagoruyko & Komodakis, 2016) - 300 epochs



Optimizer	Final Val	Final Test	Best Val	Best Test
SGD	92.0	91.8	92.9	91.9
Momentum	92.7	92.1	93.1	92.3
ADAM	90.4	90.1	91.8	90.7
RMSProp	90.7	90.3	91.4	90.3
$[e^{\text{sign}(g)*\text{sign}(m)} + \text{clip}(g, 10^{-4})] * g$	92.5	92.4	93.8	93.1
$\text{clip}(\hat{m}, 10^{-4}) * e^{\hat{v}}$	93.5	92.5	93.8	92.7
$\hat{m} * e^{\hat{v}}$	93.1	92.4	93.8	92.6
$g * e^{\text{sign}(g)*\text{sign}(m)}$	93.1	92.8	93.8	92.8
$\text{drop}(g, 0.3) * e^{\text{sign}(g)*\text{sign}(m)}$	92.7	92.2	93.6	92.7
$\hat{m} * e^{g^2}$	93.1	92.5	93.6	92.4
$\text{drop}(\hat{m}, 0.1)/(e^{g^2} + \epsilon)$	92.6	92.4	93.5	93.0
$\text{drop}(g, 0.1) * e^{\text{sign}(g)*\text{sign}(m)}$	92.8	92.4	93.5	92.2
$\text{clip}(\text{RMSProp}, 10^{-5}) + \text{drop}(\hat{m}, 0.3)$	90.8	90.8	91.4	90.9
$\text{ADAM} * e^{\text{sign}(g)*\text{sign}(m)}$	92.6	92.0	93.4	92.0
$\text{ADAM} * e^{\hat{m}}$	92.9	92.8	93.3	92.7
$g + \text{drop}(\hat{m}, 0.3)$	93.4	92.9	93.7	92.9
$\text{drop}(\hat{m}, 0.1) * e^{g^3}$	92.8	92.7	93.7	92.8
$g - \text{clip}(g^2, 10^{-4})$	93.4	92.8	93.7	92.8
$e^g - e^{\hat{m}}$	93.2	92.5	93.5	93.1
$\text{drop}(\hat{m}, 0.3) * e^{10^{-3}w}$	93.2	93.0	93.5	93.2

Neural Machine Translation

Completely different model & task: WMT 2014 English → German task
GNMT model - 8 LSTM layers

Optimizer	Train perplexity	Test BLEU
Adam	1.49	24.5
$g * e^{\text{sign}(g) * \text{sign}(m)}$	1.39	25.0

Summary

- RNN generates optimizer equations
- Train RNN via RL setup
- Optimizers tested on small ConvNet
- New optimizers on par with state of the art