

LinkedIntern Class of 2019 Summer

Scalable Automated Machine Learning

in GLMix 2.0

Carnegie Mellom University Yuwei Qiu Al Algorithms Foundation Team

Agenda

Motivation

Algorithm

Implementation

Results



Motivation Algorithm Implementation Results

Deep Learning in the enterprise domain

Deep learning and neural network are not only limited to researches.
Oltimate goal: Learn more effectively, less trial-and-error.



Problem of Learning

Performance is very sensitive to parameters
Architectural parameters & Others



Units per layer

Algorithm

Motivation

-

Results

Motivation

Algorithm

Implementation

Results

Current: Solution = Expertise + Data + Hours and Hours

But can we turn this into: Solution = Data + Less Hours (+ Improvement of Performance)



Algorithm

How to determine your model?

\circ Some possible answers...



Motivation

Algorithm

Implementation

Results

Grid Search and Random Search

- \circ Shoot all of combinations.
- \circ Then find the best set of parameters





Grid Layout

Random Layout

Algorithm

Grid Search and Random Search

If you want to grid search on N sets of parameters...
 Typically much more time or much more resources.



N times more time

N times more resources

One AI Joke...



Motivation

Algorithm

Implementation

Results

Vertically Horizontally

Algorithm

AdaNet

- Instead of stacking heavyweight layers, AdaNet stacks with lightweight subnetworks.
- \circ To automatically find the best model for your data input.

• Presented in "AdaNet: Adaptive Structural Learning of Artificial Neural Networks" at ICML 2017



Algorithm

Bird View



Select next subnetworks



 \mathbf{O}

00

Iteration t+1

 $\mathbf{O}\mathbf{O}$

Iteration t







A Closer Look into AdaNet Structure











Yuwei Qiu



A Closer Look into AdaNet Structure

o Iteration 1

- Train new subnetworks separately
- Ensemble with newly added weights
 - Each unit in layer k of this subnetwork may have **connections** to the existing units in layer k+1 of current adanet in addition to units in layer k+1 of the subnetwork.





Motivation

Algorithm

Implementation

Results

Select new candidates for the next iteration

Candidates for the next iteration may include a subnetwork with:
 (Just an example)

- $\circ\,$ the same as
- **o 1** layer deeper subnetwork than
- $\circ\,$ more units than
- $\circ\,$ that of the subnetworks of current candidates pool.

\circ Only keep best candidates for the next iteration.







A Closer Look into AdaNet Structure



Implementation

Motivation

Algorithm

Implementation

Results

Challenging Distributed Training of AdaNet

Training neural network models is computationally demanding
 Even the most lightweight network can cost an impractically long time on a single machine.

Motivation Algorithm Implementation Results

LinkedIntern Yuwei Qiu

Challenging Distributed Training of AdaNet

 $\,\circ\,$ Suppose each time we try 2 different subnetworks.

 $\circ~$ Subnetwork 2 is slightly larger than Subnetwork 1 ~





Motivation Algorithm Implementation Results

LinkedIntern Yuwei Qiu

Challenging Distributed Training of AdaNet

Suppose each time we try 2 different subnetworks.
 Subnetwork 2 is slightly larger than Subnetwork 1

Synchronization process of workers



Worker is idle



Subnetwork 1's training process

Su tr

Subnetwork 2's training process



ideally

• Suppose each time we try 2 different subnetworks. • Subnetwork 2 is slightly larger than Subnetwork 1 Select Ensemble Train subnetworks subnetworks worker 1 Implementation (Chief) worker 2 worker 3 worker 4 worker 5 (Eval) ideally

Challenging Distributed Training of AdaNet

Synchronization process of workers





Subnetwork 2's training process



Subnetwork 1&2's ensemble



Structure Adjusting On Chief

Evaluation On Eval

• Suppose each time we try 2 different subnetworks • Subnetwork 2 is slightly larger than Subnetwork 1 **Fix Structure** Train subnetworks Ensemble Select subnetworks worker 1 Implementation (Chief) worker 2 worker 3 worker 4 worker 5 (Eval)

Challenging Distributed Training of AdaNet

LinkedIntern Yuwei Qiu Practically

worker 1 (Chief) worker 2 worker 3 worker 4 worker 5 (Eval) worker 1 (Chief) worker 2 worker 3 worker 4 worker 5 (Eval)

Ideally vs Practically

Algorithm Implementation

Results



Observations

Many workers involved in the training of one subnetwork
 Heavier straggler effect and longer synchronization time



Motivation Algorithm Implementation Results

LinkedIntern Yuwei Qiu

Reduce Straggler Effect and Sync Time

Subnetwork 1 occupies 4 workers from t_o to t_1 Subnetwork 2 occupies 4 workers from t_2 to t_3



Reduce Straggler Effect and Sync Time

Subnetwork 1 occupies 2 workers from t'_o to t'_1 Subnetwork 2 occupies 2 workers from t'_o to t'_2



Observations

Many workers involved in the training of one subnetwork

 Heavier straggler effect and longer synchronization time

 Resource hogging in the gap between the training of subnetworks.



LinkedIntern Yuwei Qiu

Implementation

Reduce Resource Hogging

 Workers involved in one subnetwork-training task never wait for other workers to finish before they are scheduled to a new subnetwork-training task.



new subnetworktraining task

Reduce Resource Hogging

 \circ When all subnetworks of current iteration are finished, start building the ensemble.



Reduce Resource Hogging

 ${\rm \circ}$ After evaluation, Chief decides subnetworks for next iteration.

- $\circ\,$ Keep luck guess
- $\circ\,$ Throw away bad guess





Resources

All of the following experiments are conducted on *mlearn-alpha* by:
 Training with 20 workers, each with 32 GiB RAM
 Validation with 20 workers, each with 8 GiB RAM

LinkedIntern Yuwei Qiu

Results

jymbii: Logistic Regression

Model	Steps	Training Time	Eval Time	AUC
Logistic Regression	120,232	105 mins	20 mins	66.3%

Motivation

Algorithm

Implementation

Results

Optimizer: 1e-1 Decay LR SGD Reg: L2, 1e-4

jymbii: Best NN with Grid Search

~ 5000

Model	Steps	Training Time	Eval Time	AUC
Logistic Regression	120,232	105 mins	20 mins	66.3%
Grid Search_BST	120,232	155 mins	31 mins	67.2%

Motivatio

Algorithm

Implementation

Results

 Optimizer: 1e-1 Decay LR SGD
 R D
 R D
 R D

 Reg: L2, 1e-4
 B e
 B e
 B e
 B e

 Act: ReLU
 B e
 N L O
 B e
 N L O
 U

 Dropout: 0.1
 B e
 N L O
 U
 U
 U
 U
 U

 BatchNorm
 U.
 U.
 U.
 U.
 U.
 U.
 U.
 U.

 # Layer_2 / # Layer_1 = 2/3
 @Input
 @FC_1 (ReLU)
 @FC_2 (ReLU)
 @Output (SoftMax)

~ 5000

~ 3000

jymbii: Best NN with Grid Search



LinkedIntern Yuwei Qiu

Results

jymbii: AdaNet

Model	Steps	Training Time	Eval Time	AUC
Logistic Regression	4,000	6 mins	19 mins	63.3%
Logistic Regression	120,232	105 mins	20 mins	66.3%
Grid Search_BST	4,000	8 mins	28 mins	64.5%
Grid Search_BST	120,232	155 mins (x12)	31 mins	67.2%
AdaNet	2,000 x2	36 mins	31 mins	67.6%

*Randomly down-sampled 4000 steps, 2 iterations for structure tuning.

LinkedIntern Yuwei Qiu

Results

pymk: Scalable Automated ML eases your way

n	Model	Steps	Training Time	Eval Time	AUC	
n tion	Logistic Regression (Default Setting)	20,000	10 mins	4 mins	67.5%	Default setting for every parameter

Results

pymk: Scalable Automated ML eases your way

otivation	Model	Steps	Training Time	Eval Time	AUC	
gorithm ementation	Logistic Regression (Default Setting)	20,000	10 mins	4 mins	67.5%	Default setting for every parameter
Results	Logistic Regression (Manual Tuned)	20,000	9 mins	3 mins	73.1%	l spent hours to manual tune it

pymk: Scalable Automated ML eases your way

otivation	Model	Steps	Training Time	Eval Time	AUC	
lgorithm ementation	Logistic Regression (Default Setting)	20,000	10 mins	4 mins	67.5%	Default setting for every parameter
Results	Logistic Regression (Manual Tuned)	20,000	9 mins	3 mins	73.1%	l spent hours to manual tune it
	Grid Search_BST	20,000	18 mins (x4)	6 mins	74.8%	Grid search spend around 1 hour

pymk: Scalable Automated ML eases your way

otivation	Model	Steps	Training Time	Eval Time	AUC	
lgorithm ementation	Logistic Regression (Default Setting)	20,000	10 mins	4 mins	67.5%	Default setting for every parameter
Results	Logistic Regression (Manual Tuned)	20,000	9 mins	3 mins	73.1%	l spent hours to manual tune it
	Grid Search_BST	20,000	18 mins (x4)	6 mins	74.8%	Grid search spend around 1 hour
	AdaNet	10,000 x2	17 mins	4 mins	76.1%	

LinkedIntern Yuwei Qiu As for AdaNet, I just run the code. And after **20 minutes**...

Summary

Motivation

Algorithm

Implementation

Results

Tuning a Neural Network is exhausting.

Select optimal subnetworks and implement an adaptive learning process to attain an ensemble with AdaNet

Implement a network-level asynchronous distributed learning strategy to maintain the feasibility of scalable AdaNet

Achieve promising results on *Job You May Be Interested In(jymbii)* dataset and *People You May Know(pymk)* dataset

Thank you soooooo much!



AdaNet - Theory

 \circ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$



Motivation Algorithm Design

Result

AdaNet - Theory

 \circ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$



Motivation Algorithm Design

Result

AdaNet - Theory

 \circ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$



Motivation Algorithm Design

Result

AdaNet - Theory

 \circ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$

Motivation Algorithm Design

Result



$$\varphi_s \circ \mathbf{h}_s =$$

 $(\varphi_s \circ h_{s,1}, \dots, \varphi_s \circ h_{s,n_s})$

Motivation

Algorithm

Design

Result

AdaNet - Theory

 \circ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$



Motivation Algorithm Design

Result

LinkedIntern Yuwei Qiu

AdaNet - Theory

 \circ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$ (Almost)

Every network could be presented as a member of Hidden Layer Family.



Motivation

Algorithm

Design

Result

AdaNet - Theory

○ Hidden Layer Family $\mathcal{H} \longrightarrow \mathbf{h}_s \longrightarrow h_{s,n_s}$ ○ Objective Function with Hidden Layer Family





AdaNet - Algorithm

 $\ensuremath{\circ}$ Iteratively "ensemble" new subnetworks to the current network

 \circ During the iteration t

Motivation

Algorithm

Design

Result

\circ "Ensemble" a new candidate subnetwork ${\boldsymbol u}$

 $\circ\,$ Minimize the objective function to train it



Ensemble mixture weights

AdaNet - Algorithm

- $\ensuremath{\circ}$ Iteratively "ensemble" new subnetworks to the current network
- \circ During the iteration t

Motivation

Algorithm

Design

Result

- \circ "Ensemble" a new candidate subnetwork **u** *"under specific rules"*
- $\,\circ\,$ Minimize the objective function to train it
- $\circ\,$ Choose the one that leading to the best reduction of objective function, i.e.

if $\min_{\mathbf{w}} F_t(\mathbf{w}, \mathbf{h}) \leq \min_{\mathbf{w}} F_t(\mathbf{w}, \mathbf{h}')$, then

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^B} F_t(\mathbf{w}, \mathbf{h}), \quad \mathbf{h}_t = \mathbf{h}$$

And otherwise

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^B} F_t(\mathbf{w}, \mathbf{h}'), \quad \mathbf{h}_t = \mathbf{h}'$$

