

The First Step-by-Step Guide for Implementing Neural Architecture Search with Reinforcement Learning Using TensorFlow



[Wallarm](#)

[Dec 12, 2017](#) · 7 min read

Our team is no stranger to various flavors of AI including deep learning (DL). That's why we've immediately noticed when Google came out with [AutoML](#) project, designed to make AI build other AIs.

Neural networks have recently gained popularity and wide practical applications. However, to get good results with neural networks, it is critical to pick the right network topology, which has always been a difficult manual task.

Google's recent project promises to help solve this task automatically with a meta-AI which will design the topology for neural network architecture. Google, however, did not offer documentation or examples of how to use this new wonderful technology. We liked the idea and, among the first, came up with a practical implementation that other people can follow, using it as an example. This is similar in concept to AlphaGo, for instance.

Google's approach is based on the AI concept called [Reinforcement Learning](#), meaning that the parent AI reviews the efficiency of the child AI and makes adjustments to the neural network architecture, such as adjusting the number of layers, weights, regularization methods, etc. to improve efficiency.

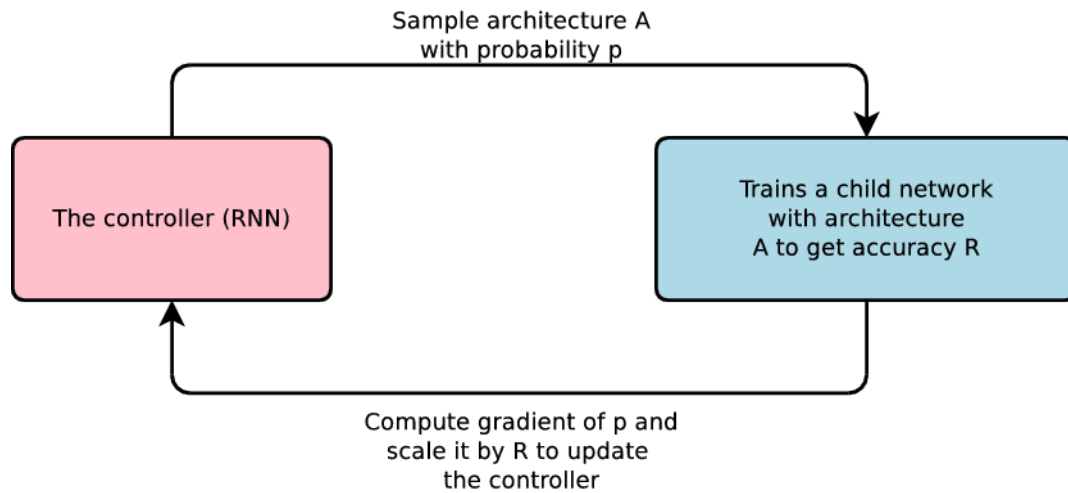


Image from [Google Blog](#)

The advantage of automation is the ability to eliminate guesswork from the manual neural network model design as well as significantly reducing the time required for each problem, since designing the neural network model is the most labor-intensive part of the task.

Although Google has recently [open sourced](#) an example of NASnet, how they found the architecture of NASnet is still unclear to most folks.

In addition, in our opinion, the name itself adds to the confusion with these technologies.

In this post, we will take a detailed look (with a step by step explanation) at implementing a simple model for neural architecture search with AutoML and reinforcement learning.

Note: To understand this post, you will need to have sufficient background understanding of the convolutional neural networks, recurrent neural networks, and reinforcement learning.

Links below will provide you with good background information:

- https://en.wikipedia.org/wiki/Reinforcement_learning
- <https://www.cs.cmu.edu/afs/cs/project/jair/pub/volume4/kaelbling96a-html/node37.html>
- <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- https://en.wikipedia.org/wiki/Convolutional_neural_network

Neural Architecture Search (NAS) with Reinforcement Learning is a method for finding good neural networks architecture. For this post, we will try to find optimal architecture for Convolutional Neural Network (CNN) which recognizes handwritten digits.

For this implementation, we use TensorFlow 1.4, but if you want to try this at home, you can use any version after 1.1, since NASCell first became available in TensorFlow 1.1. It is important not to confuse AutoML and NAS.

Data and preprocessing.

To train the model we will use the [MNIST](#) database of handwritten digits, which has a training set of 55,000 examples and a test set of 10,000 examples.

The Model.

The network we are building in this exercise consists of a controller and the actual neural network that we are trying to optimize. The Controller is an rnn tensorflow with [NAS cells](#) and special reinforcement learning methods for training and getting rewards. We will define “rewards” as maximizing the accuracy of the desired neural network and train the Controller to improve this outcome. The controller should generate Actions to modify the architecture of CNN. **Specifically, Actions can modify filters: the dimensionality of the output space, kernel_size (integer, specifying the length of the 1D convolution window), pool_size (integer, representing the size of the pooling window) and dropout_rate per layer.**

Details of architecture search space

All convolutions employ Rectified Linear Units (ReLU) nonlinearity. Weights were initialized by the [Xavier initialization](#) algorithm.

Implementation.

For the Controller, we built a method for policy network based on NASCell. This network takes, as inputs, the current state (in this task, state and action are the same things) and maximum number of searching layers and outputs new Action to update the desired neural network. If for some reason, NASCell is not available, you can use any RNNCell.

<https://gist.github.com/d0znpp/db1caaefd6c9ed5c72e65721fbd859bb>

To allow [hyperparameter tuning](#) we put our code into a Reinforce class.

```
1 class Reinforce():
2     def __init__(self, sess, optimizer, policy_network, max_layers,
3         global_step,
4         division_rate=100.0, reg_param=0.001,
5         discount_factor=0.99, exploration=0.3):
6         self.sess = sess
7         self.optimizer = optimizer
8         self.policy_network = policy_network
9         self.division_rate = division_rate
```

```

8         self.reg_param = reg_param
9         self.discount_factor=discount_factor
10        self.max_layers = max_layers
11        self.global_step = global_step
12
13        self.reward_buffer = []
14        self.state_buffer = []

```

sess and optimizer — TensorFlow session and optimizer, will be initialized separately.

- policy_network — Method initialized above.
- max_layers — The maximum number of layers
- division_rate — Normal distribution values of each neuron from -1.0 to 1.0.
- reg_param — Parameter for regularization.
- exploration — The probability of generating random action.

Of course, we also must create variables and placeholders, consisting of logits and gradients. To do this, let's write a method create_variables:

```

1  def create_variables(self):
2      with tf.name_scope("model_inputs"):
3          # raw state representation
4          self.states = tf.placeholder(tf.float32, [None,
self.max_layers*4], name="states")
5
6          with tf.name_scope("predict_actions"):
7              # initialize policy network
8              with tf.variable_scope("policy_network"):
9                  self.policy_outputs =
self.policy_network(self.states, self.max_layers)
10                 self.action_scores =
tf.identity(self.policy_outputs, name="action_scores")
11                 self.predicted_action = tf.cast(tf.scalar_mul(
12                     self.division_rate, self.action_scores),
tf.int32, name="predicted_action")
13
14                 # regularization loss
15                 policy_network_variables =
tf.get_collection(tf.GraphKeys.TRAINABLE_VARIABLES,
16                 scope="policy_network")
17
18                 # compute loss and gradients
19                 with tf.name_scope("compute_gradients"):

```

```

20             # gradients for selecting action from policy
network
21             self.discounted_rewards =
tf.placeholder(tf.float32, (None,),
22
name="discounted_rewards")
23
24             with tf.variable_scope("policy_network",
reuse=True):
25                 self.logprobs =
self.policy_network(self.states, self.max_layers)
26                 # compute policy loss and regularization
loss
27                 self.cross_entropy_loss =
tf.nn.softmax_cross_entropy_with_logits(
28                     logits=self.logprobs,
labels=self.states)
29                 self.pg_loss =
tf.reduce_mean(self.cross_entropy_loss)
30                 self.reg_loss =
tf.reduce_sum([tf.reduce_sum(tf.square(x))
31                                     for x in
policy_network_variables])
32                 self.loss = self.pg_loss + self.reg_param *
self.reg_loss
33
34                 #compute gradients
35                 self.gradients =
self.optimizer.compute_gradients(self.loss)
36
37                 # compute policy gradients
38                 for i, (grad, var) in
enumerate(self.gradients):
39                     if grad is not None:
40                         self.gradients[i] = (grad *
self.discounted_rewards, var)
41                     # training update
42                     with
tf.name_scope("train_policy_network"):
43                         # apply gradients to update
policy network
44                         self.train_op =
self.optimizer.apply_gradients(
45                             self.gradients)

```

After computing the initial gradients, we launch the [gradient descent](#) method. Now let's take a look at how reinforcement learning is implemented.

First, we can multiply gradient value to the discounted reward.

After defining the variables, we should initialize it in a TensorFlow graph in end of `__init__`:

```
1 self.create_variables()
2 var_lists = tf.get_collection(tf.GraphKeys.GLOBAL_VARIABLES)
3 self.sess.run(tf.variables_initializer(var_lists))
```

Every Action depends on the previous state, but sometimes, for more effective training, we can generate random actions to avoid local minimums.

In each cycle, our network will generate an Action, get rewards and after that, take a training step.

The implementation of the training step includes `store_rollout` and `train_step` methods below:

```
1 def store_rollout(self, state, reward):
2     self.reward_buffer.append(reward)
3     self.state_buffer.append(state[0])
4
5 def train_step(self, steps_count):
6     states = np.array(self.state_buffer[-
7 steps_count:])/self.division_rate
8     rewards = self.reward_buffer[-steps_count:]
9     _, ls = self.sess.run([self.train_op, self.loss],
10                           {self.states: states,
11 self.discounted_rewards: rewards})
12     return ls
```

As mentioned above, we need to define rewards for each Action\State.

This is accomplished by generating a new CNN network with new architecture per Action, training it and assessing its accuracy. Since this process generates a lot of CNN networks, let's write a manager for it:

```

1 class NetManager():
2     def __init__(self, num_input, num_classes, learning_rate, mnist,
3                 max_step_per_action=5500, bathc_size=100,
4                 dropout_rate=0.85):
5         self.num_input = num_input
6         self.num_classes = num_classes
7         self.learning_rate = learning_rate
8         self.mnist = mnist
9
10        self.max_step_per_action = max_step_per_action
11        self.bathc_size = bathc_size
12        self.dropout_rate = dropout_rate #Dropout after dense layer
in CNN

```

Then we formed bathc with hyperparameters for every layer in “action” and we created *cnn_drop_rate* — list of dropout rates for every layer.

```

1 def get_reward(self, action, step, pre_acc):
2     action = [action[0][0][x:x+4] for x in range(0, len(action[0]
3     [0]), 4)]
4     cnn_drop_rate = [c[3] for c in action]
5
6     with tf.Graph().as_default() as g:
7         with g.container('experiment'+str(step)):
8             model = CNN(self.num_input, self.num_classes, action)
9             loss_op = tf.reduce_mean(model.loss)
10            optimizer =
tf.train.AdamOptimizer(learning_rate=self.learning_rate)
            train_op = optimizer.minimize(loss_op)

```

Here we define a convolution neural model with CNN class. It can be any class that is able to generate the neural model by some action.

We created a separate container to avoid confusion in TF graph.

After creating a new CNN model, we can train it and get a reward.

```

1 with tf.Session() as train_sess:
2     init = tf.global_variables_initializer()
3     train_sess.run(init)
4
5     for step in range(self.max_step_per_action):

```

```

6         batch_x, batch_y =
self.mnist.train.next_batch(self.batch_size)
7         feed = {model.X: batch_x,
8                 model.Y: batch_y,
9                 model.dropout_keep_prob: self.dropout_rate,
10                model.cnn_dropout_rates: cnn_drop_rate}
11         _ = train_sess.run(train_op, feed_dict=feed)
12
13
14         batch_x, batch_y = self.mnist.test.next_batch(10000)
15         loss, acc = train_sess.run(
16             [loss_op, model.accuracy],
17             feed_dict={model.X: batch_x,
18                       model.Y: batch_y,
19                       model.dropout_keep_prob: 1.0,
20                       model.cnn_dropout_rates:
[1.0]*len(cnn_drop_rate)})
21         if acc - pre_acc <= 0.01:
22             return acc, acc
23         else:
24             return 0.1, acc

```

As defined, the reward improves accuracy on all test datasets; for MNIST it is 10000 examples.

Now that we have everything in place, let's find the best architecture for MNIST. First, we will optimize the architecture for the number of layers. Let's set the maximum number of layers to 2. Of course, you can set this value to be higher, but every layer needs a lot of computing power.

```

1  def train(mnist, max_layers):
2      sess = tf.Session()
3      global_step = tf.Variable(0, trainable=False)
4      starter_learning_rate = 0.1
5      learning_rate = tf.train.exponential_decay(0.99, global_step,
500, 0.96, staircase=True)
6
7      optimizer =
tf.train.RMSPropOptimizer(learning_rate=learning_rate)
8      reinforce = Reinforce(sess, optimizer, policy_network,
args.max_layers, global_step)
9      net_manager = NetManager(num_input=784, num_classes=10,
learning_rate=0.001, mnist=mnist)
10
11     MAX_EPISODES = 250

```



```

12     step = 0
13     state = np.array([[10.0, 128.0, 1.0, 1.0]*max_layers],
dtype=np.float32)
14     pre_acc = 0.0
15     for i_episode in range(MAX_EPISODES):
16         action = reinforce.get_action(state)
17         print("current action:", action)
18         if all(ai > 0 for ai in action[0][0]):
19             reward, pre_acc = net_manager.get_reward(action, step,
pre_acc)
20         else:
21             reward = -1.0
22             # In our sample action is equal state
23             state = action[0]
24             reinforce.store_rollout(state, reward)
25
26             step += 1
27     ls = reinforce.train_step(MAX_STEPS)
28     log_str = "current time: " +
str(datetime.datetime.now().time()) +
29         " episode: " + str(i_episode) + " loss: " + str(ls) + "
last_state: " +
30         str(state) + " last_reward: " + str(reward)
31     print(log_str)
32
33 def main():
34     max_layers = 3
35     mnist = input_data.read_data_sets("MNIST_data/", one_hot=True)
36     train(mnist, max_layers)
37
38 if __name__ == '__main__':
39     main()

```

We couldn't be sure about what we should feed to our policy network. First, we tried to always feed in the array of 1.0 to our RNN per episode, but it yielded no results. Then we tried feeding every new state per episode and it resulted in a good architecture. We concluded that the first state can be any non-zero array, to expedite finding a suitable architecture we set the first state: `[[10.0, 128.0, 1.0, 1.0]*args.max_layers]`

We have updated the weights after every episode. Otherwise, our calculations would have been useless. That's why our "batch size" for reinforce = 1.

After 100 cycles, we get the following architecture:

- input layer : 784 nodes (MNIST images size)

- first convolution layer : 61x24
- first max-pooling layer: 60
- second convolution layer : 57x55
- second max-pooling layer: 59
- output layer : 10 nodes (number of class for MNIST)

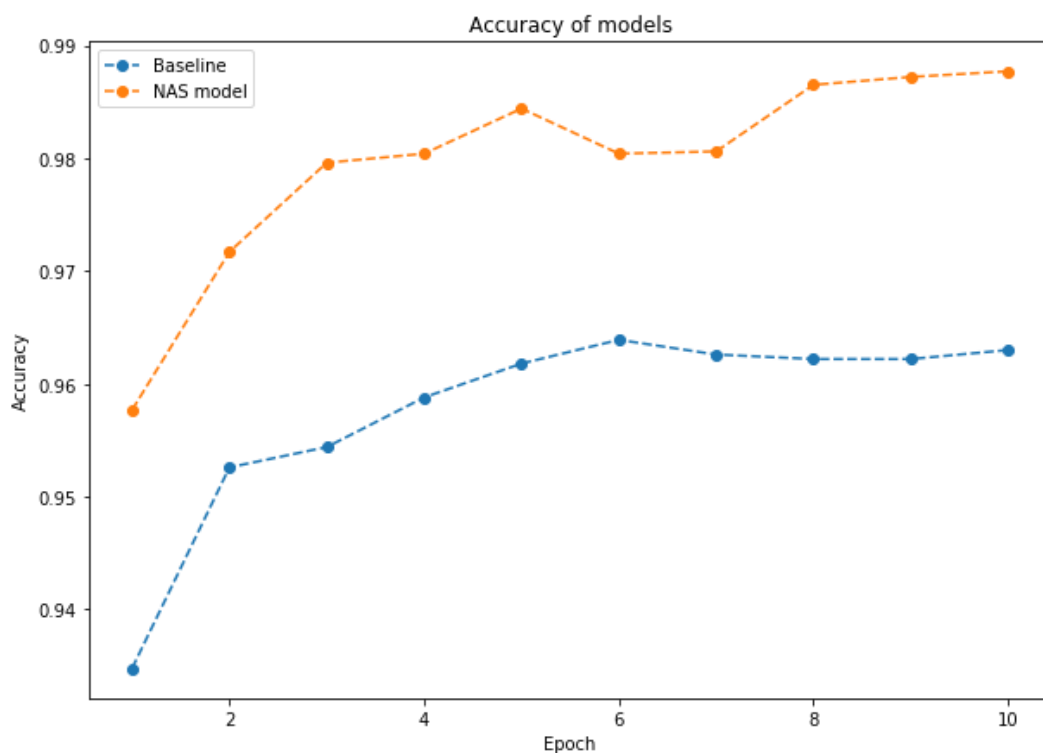
Measuring Results

Now that we've trained our "NAS model" on MNIST dataset, we should be able to compare the architecture our AI has created with the other architectures created manually. For comparable results we will use popular [Convolutional Neural Network \(CNN\) architecture](#) for MNIST [It's not the state-of-the-art architecture, but it's good for comparing]:

- input layer : 784 nodes (MNIST images size)
- first convolution layer : 5x32
- first max-pooling layer: 2
- second convolution layer : 5x64
- second max-pooling layer: 2
- output layer : 10 nodes (number of class for MNIST)

All weights were initialized by the Xavier algorithm.

We trained our models on 10 epochs and got of the accuracy of 0.9987 for our "NAS model", compared to 0.963 for the popular manually defined neural network architecture.



Conclusion

We have presented a code example of a simple implementation that automates the design of machine learning models and:

- doesn't require any human time to design
- actually delivered a better performance than a manually designed network

Going forward, we will continue working on careful analysis and testing of these machine-generated architectures to help refine our understanding of them. Naturally, if we search for more parameters using our model, we'll achieve better results for MNIST, but more importantly, this simple example illustrates how this approach can be applied to the problems that are much more complicated. We built this model using some assumptions which are quite difficult to justify if you notice any mistakes, please write in issues on GitHub.

[The full code of the project is available on Github](#)