

Software Packages for Deep Learning

Chensong Zhang

with Zheng Li and Ronghong Fan

DL Seminar — April 27, 2017

Outline

Introduction

Numerical Comparisons

Python

Torch

Caffe

TensorFlow

MXNet

§1. Introduction

- 1 Machine Learning
- 2 General Tasks of ML
- 3 Packages for General Machine Learning
- 4 Deep Learning: Pros and Cons
- 5 Popular Packages: Basic Information
- 6 Popular Packages: Performance

Machine Learning



- Unlike traditional numerical simulation, “ML gives computers the ability to learn without being explicitly programmed” [Samuel 1959]
- As a research field, ML explores the study and construction of algorithms that can **learn** from and **make predictions** on **data**
- Related fields: data mining, computational statistics, optimization, ...
- Fourth paradigm, big data, artificial intelligence, Internet of things, deep learning, ...

General Tasks of ML

- Classification: Inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more (multi-label classification) of these classes
- Clustering: Inputs are divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task
- Regression: Similar to classification, but the outputs are continuous rather than discrete
- Density estimation
- Dimensionality reduction
- ...

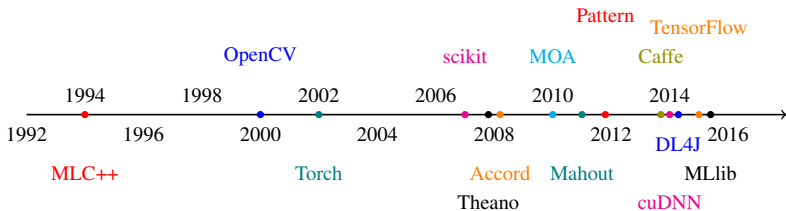
Packages for General Machine Learning

What is the purpose?

- Solving problems from practical applications (user interface)
- Developing algorithms and optimizing implementation (development)
- Theoretical analysis for machine learning

What do we want for a ML package?

- Easy for new tasks and new network structures (less steep learning curve)
- Easy for debugging (with good support and large community)
- Performance and scalability



Deep Learning: Pros and Cons

Deep Learning has been introduced with the objective of moving ML closer to one of its original goals—AI. The main motivations includes:

- Insufficient depth can hurt
- The brain has a deep architecture
- Cognitive processes seem deep

Pros:

- conceptually simple
- nonlinear
- highly flexible and configurable
- learned features can be extracted
- can be fine-tuned with more data
- efficient for multi-class problems
- world-class at pattern recognition

Cons:

- hard to interpret
- theory not well understood
- slow to train and score
- overfits, needs regularization
- more parameters
- inefficient for categorical variables
- data hungry, learns slowly

Popular Packages: Basic Information

Viewpoint	Torch	Caffe	TensorFlow	MXNet
Started	2002	2013	2015	2015
Main Developers	Facebook, Twitter, Google, ...	BVLC (Berkeley)	Google	DMLC
License	BSD	BSD	Apache	Apache
Core Languages	C/Lua	C++	C++ Python	C++
Supported Interface	Lua	C++/Python Matlab	C++/Python R/Java/Go	C++/Python R/Julia/Scala

- Other worth noting packages: CNTK/DMTK (Microsoft), PyTorch (beta)
- BVLC, Berkeley Vision and Learning Center
- DMLC, Distributed (Deep) Machine Learning Community, supported by Amazon, Intel, Microsoft, nVidia, Baidu, ...

Popular Packages: Performance

Viewpoint	Torch	Caffe	TensorFlow	MXNet
Pretrained Models	Yes	Yes	No	Yes
High-level Support	Good	Good	Good	Good
Low-level Operators	Good	Good	Fairly good	Very few
Speed One-GPU	Great	Great	Not so good	Excellent
Memory Management	Great	Great	Not so good	Excellent
Parallel Support	Multi-GPU	Multi-GPU	Multi-GPU	Distributed

<http://blog.revolutionanalytics.com/2016/08/deep-learning-part-1.html>

§2. Numerical Comparisons

- 7 Hardware Platforms
- 8 Neural Networks and Data Sets
- 9 CPU Scalability: FCN Synthetic
- 10 CPU Scalability: FCN Real
- 11 CPU Scalability: CNN Synthetic
- 12 CPU Scalability: CNN Real
- 13 CPU Scalability: RNN
- 14 GPU Scalability: FCN Synthetic
- 15 GPU Scalability: CNN Synthetic
- 16 GPU Scalability: CNN Real
- 17 GPU Scalability: RNN Real
- 18 GPU Scalability: LSTM

Hardware Platforms

- Use one quad-core desktop CPU (i.e., Intel i7-3820 CPU @ 3.60GHz) and two 8-core server-grade CPUs (i.e., Intel Xeon CPU E5-2630 v3 @ 2.40GHz)
- Use two generations of GPU cards, GTX 1080 @ 1607MHz with Pascal architecture, and Telsa K80 @ 562MHz with Kepler architecture

Computational Unit	Cores	Memory	OS	CUDA
Intel CPU i7-3820	4	64 GB	Ubuntu 14.04	-
Intel CPU E5-2630x2	16	128 GB	CentOS 7.2	-
GTX 980	2048	4 GB	Ubuntu 14.04	8.0
GTX 1080	2560	8 GB	Ubuntu 14.04	8.0
Telsa K80 GK210	2496	12 GB	CentOS 7.2	8.0

Figure: The experimental hardware setting for data parallelization

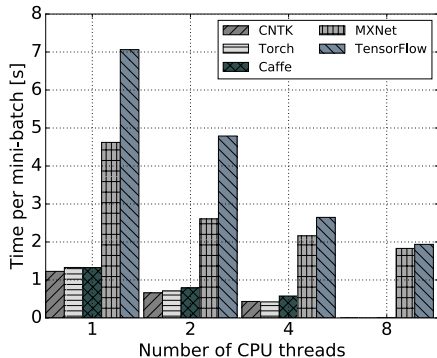
Neural Networks and Data Sets

- A large fully-connected neural network (**FCN-S**) with around 55 million parameters is used to evaluate the performance of FCN
- The classical AlexNet (**AlexNet-S**) is used as an representative of CNN
- A smaller FCN (**FCN-R**) is constructed for MNIST data set
- An AlexNet (**AlexNet-R**) architecture is used for Cifar10 data set
- For RNNs, considering that the main computation complexity is related to the length of input sequence, 2 LSTM layers are selected for testing, with input length of 32.

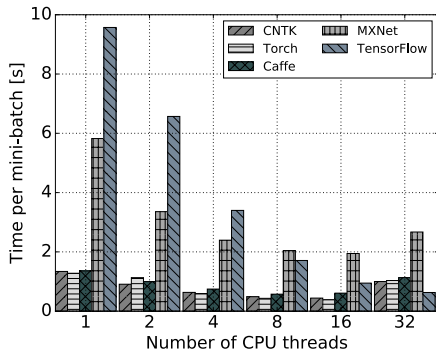
Networks		Input	Output	Layers	Parameters
FCN	FCN-S	26752	26752	5	~55 millions
CNN	AlexNet-S	150528	1000	4	~61 millions
FCN	FCN-R	784	10	5	~31 millions
CNN	AlexNet-R	3072	10	4	~81 thousands
RNN	LSTM	10000	10000	2	~13 millions

Figure: The experimental setup of neural networks for synthetic data and real data

CPU Scalability: FCN Synthetic



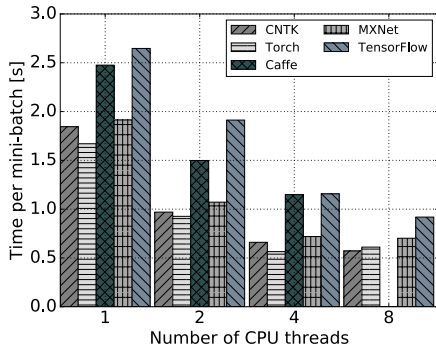
(a) Results on i7-3820.



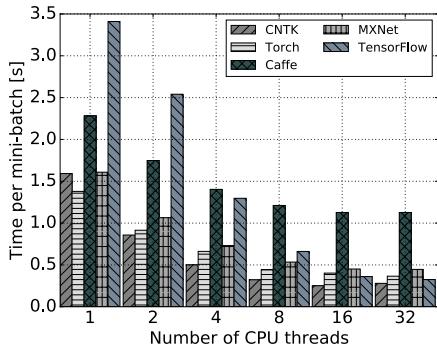
(b) Results on E5-2630.

Figure: FCN-S performance comparison on CPU platform with a mini-batch size of 64.(The lower the better.)

CPU Scalability: FCN Real



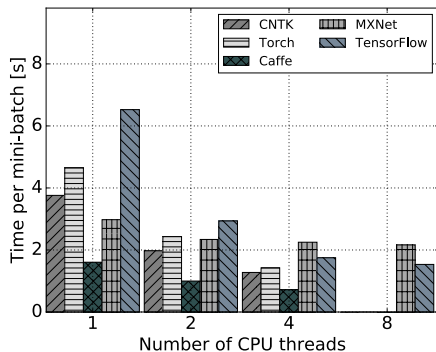
(a) Results on i7-3820.



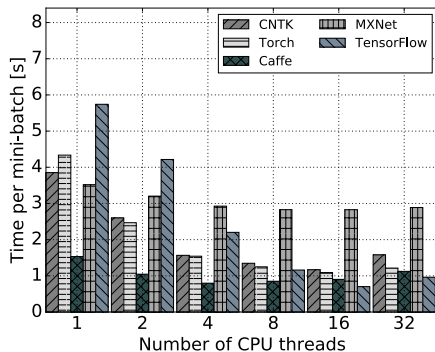
(b) Results on E5-2630.

Figure: The FCN-R performance comparison on CPU platform with a mini-batch size of 1024.(The lower the better.)

CPU Scalability: CNN Synthetic



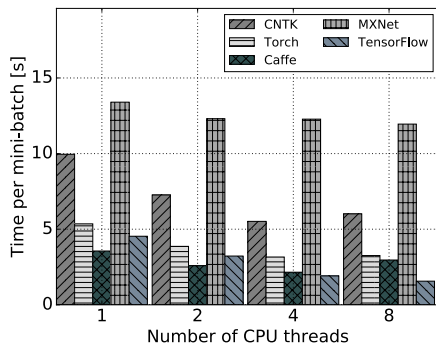
(a) Results on i7-3820.



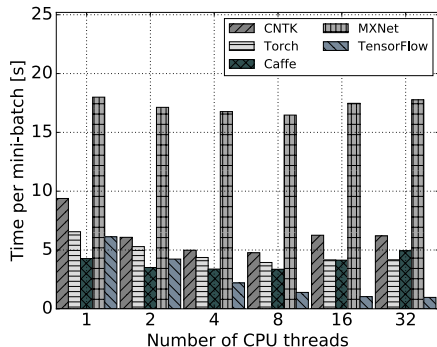
(b) Results on E5-2630.

Figure: AlexNet-S performance comparison on CPU platform with a mini-batch size of 16 (The lower the better.)

CPU Scalability: CNN Real



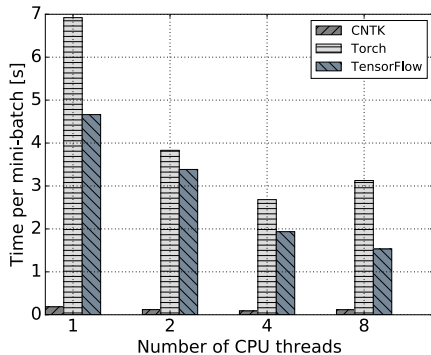
(a) Results on i7-3820.



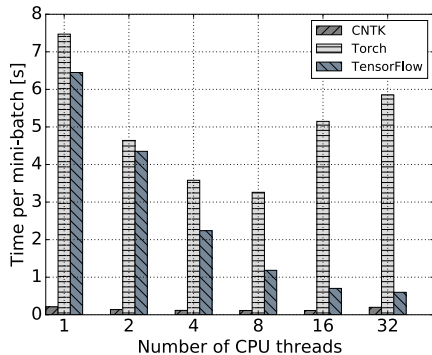
(b) Results on E5-2630.

Figure: AlexNet-R performance comparison on CPU platform with a mini-batch size of 1024.(The lower the better.)

CPU Scalability: RNN



(a) Results on i7-3820.



(b) Results on E5-2630.

Figure: LSTM performance comparison on CPU platform with a mini-batch size of 256.(The lower the better.)

GPU Scalability: FCN Synthetic

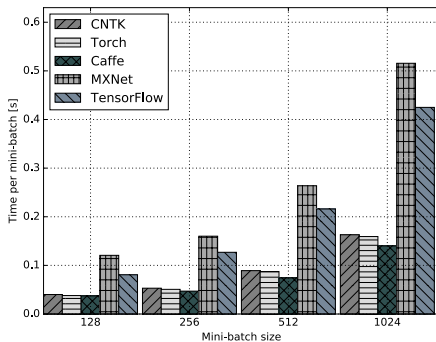
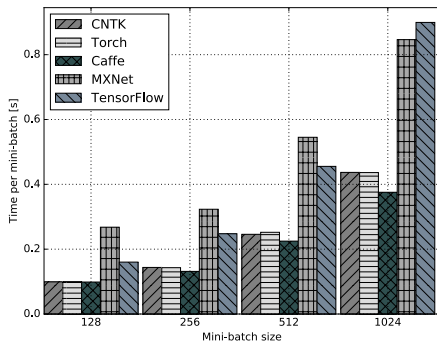


Figure: The performance comparison of FCN-S on GPU platforms.

GPU Scalability: CNN Synthetic

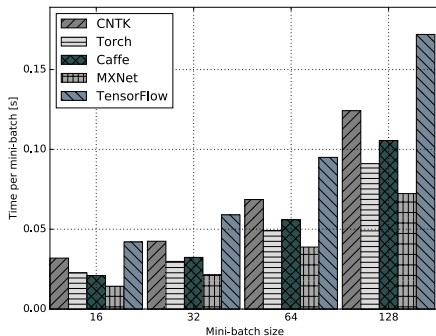
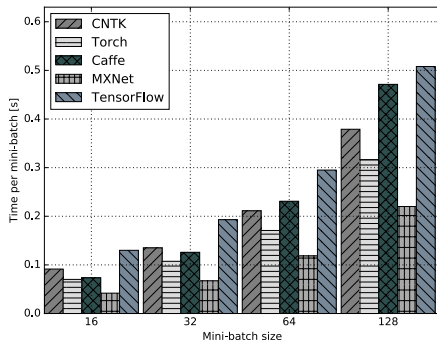


Figure: The performance comparison of AlexNet-S on GPU platforms.

GPU Scalability: CNN Real

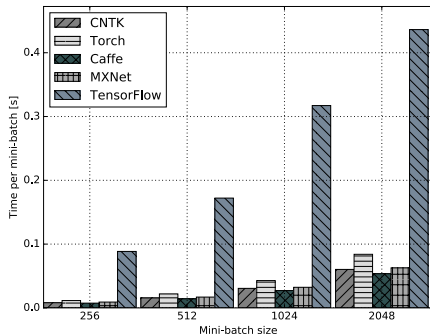
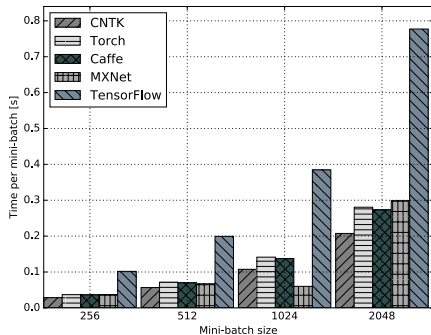


Figure: The performance comparison of AlexNet-R on GPU platforms.

GPU Scalability: RNN Real

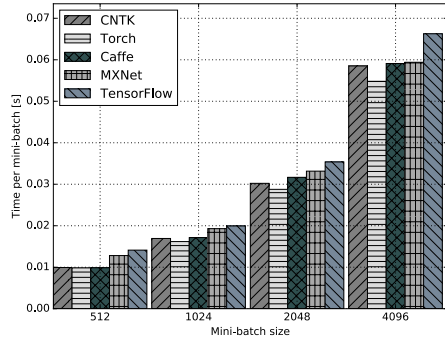
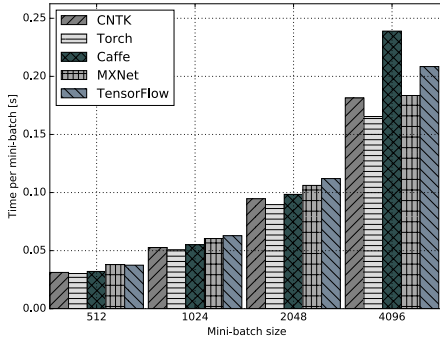
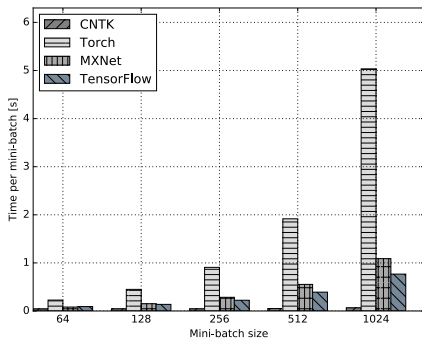
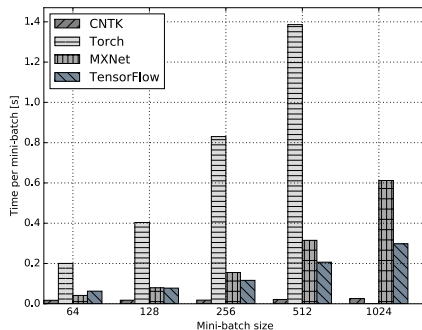


Figure: The performance comparison of FCN-R on GPU platforms.

GPU Scalability: LSTM



(a) Results on Tesla K80.



(b) Results on GTX1080.

Figure: The performance comparison of LSTM on GPU platforms.

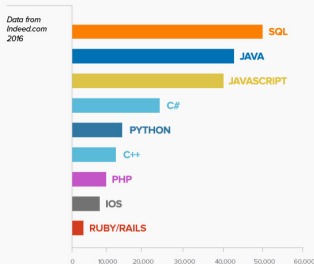
§3. Python

- 19 Python: A general-purpose programming language
- 20 Python for Scientific Computing
- 21 Modules
- 22 Built-in Data Structures
- 23 Control Flow
- 24 Additional Comments

Python: A general-purpose programming language

- Created by Guido van Rossum in 1989 and first released in 1991
- Named after “the Monty Python” (British comedy group)
- An interpreted language—simple, clear, and readable
- Python has many excellent packages for machine learning
- The language of choice in introductory programming courses

Languages ranked by number of programming jobs



Feb 2017	Change	Programming language	Share	Trends
1		Java	22.6 %	-1.3 %
2		Python	14.7 %	+2.8 %
3		PHP	9.4 %	-1.2 %
4		C#	8.3 %	-0.3 %
5	↑↑	Javascript	7.7 %	+0.4 %
6		C	7.0 %	-0.2 %
7	↓↓	C++	6.9 %	-0.6 %
8		Objective-C	4.2 %	-0.6 %
9	↑	R	3.4 %	+0.4 %
10	↓	Swift	2.9 %	+0.1 %

Python for Scientific Computing

Why Python for scientific computing?

- Dynamic data types and automatic memory management
- Full modularity, supporting hierarchical packages
- Strong introspection capabilities¹
- Exception-based error handling

Why consider such a slow language for simulation?

- Good for proof-of-concept prototyping
- Implementation time versus execution time
- Code readability and maintenance — short code, fewer bugs
- Well-written Python code is “fast enough” for most computational tasks
- Time critical parts executed through compiled language or **available packages**

¹Code introspection is the ability to examine classes, functions and keywords to know what they are, what they do and what they know. Python provides several functions and utilities for code introspection, like `dir()`, `help()`, `type()`.

Modules

Using modules

- 1 This way will only introduce the name 'math' into the name space in which the import command was issued. The names within the math module will not appear in the enclosing namespace: they must be accessed through the name math. For example: `math.sin(3.14)`.

```
1 import math
```

- 2 This way does not introduce the name math into the current namespace. It does however introduce all public names of the math module into the current namespace, directly using: `sin(3.14)`

```
1 from math import *
```

- 3 This will only import the sin function from math module and introduce the name sin into the current namespace, but it will not introduce the name math into the current namespace, directly using: `sin(3.14)`

```
1 from math import sin
```

Built-in Data Structures

Numeric types: int, float, complex

```
1 b=1L      # long int
2 c=0xf     # int (hex format)
3 d=010     # int (octal format)
4 e=1.0     # float
5 f=1+2j    # complex
```

Sequence types: list, tuple, str, dict

```
1 t=(3.14, True, 'Yes', [1], (0xf,))      # tuple example
2 l=[3.14, True, 'Yes', [1], (1L, 0xf)] + [None]*3  # list example
3 s='Hello' + ", " + 'world!'            # str example 1
4 s=("Hello, " "world!")                  # str example 2
5 d={1: 'int', 'pi': 3.14}                 # dict example
6 s="Python"; s.find('thon')              # find substring
```

Formatted output

```
1 print('%(lang)s has %(num)02d quote types.' %{"lang":"Python", "num":3})
```

User defined functions

```
1 def square(x):
2     return x*x
```

Control Flow

If-then-else

```
1 a = 1
2 if a > 0:
3     print "a is positive"
4 elif a=0:
5     print "a is zero"
6 else:
7     print "a is negative"
```

For loop

```
1 for i in range(10):
2     print i
```

While loop

```
1 sum = 0; i = 0
2 while i < 10:
3     sum += i
4     i += 1
```

Additional Comments

- ❶ In Python, everything (including functions, modules, and files) are objects. A variable is created through assignment:

```
1 x = y = z = 0.1
```
- ❷ When objects are passed to a function, Python always passes (the value of) the reference to the object to the function.
- ❸ `help()` is a function which gives information about the object. For example, `help('modules')` will generate a list of all modules which can be imported into the current interpreter.
- ❹ Some useful and important packages
 - NumPy: for scientific computing
 - Matplotlib/Pylab: for visualising data
 - SciPy: providing lots of numerical algorithms
 - SymPy: for symbolic mathematics

§4. Torch

- 25 Programming interface
- 26 Example: Linear Regression
- 27 Example: Linear Regression
- 28 Example: Linear Regression
- 29 Example: Two-Layer Network
- 30 Example: Two-Layer Network

Programming interface

- ① Wide range of applications
 - Speech, image and video applications
 - Large-scale machine-learning applications
- ② Fastest scripting language Lua is used
- ③ Easily ported to any platform
 - Torch can run on iPhone with no modification to scripts
- ④ Easy extensibility
 - Easy to integrate any library into Torch

Example: Linear Regression

```
1  require 'torch'
2  require 'optim'
3  require 'nn'
4
5  # write the loss to a text file and read from there to plot it as training proceeds
6  logger = optim.Logger('loss_log.txt')
7
8  # input data
9  data = torch.Tensor({{40, 6, 4},{44, 10, 4},{46, 12, 5},
10 {48, 14, 7},{52, 16, 9},{58, 18, 12},{60, 22, 14},
11 {68, 24, 20},{74, 26, 21},{80, 32, 24}}
12
13 # define the container
14 model = nn.Sequential()
15 ninputs = 2; noutputs = 1
16
17 # define the only module
18 model:add(nn.Linear(ninputs , noutputs))
19
20 # Define a loss function
21 criterion = nn.MSECriterion()
22
23 # retrieve its trainable parameters
24 x, dl_dx = model:getParameters()
25
26 # compute loss function and its gradient
27 feval = function(x_new)
28   # set x to x_new, if different
29   if x ~= x_new then
30     x:copy(x_new)
31   end
```


Example: Linear Regression

```
32 # select a new training sample
33 _nidx_ = (_nidx_ or 0) + 1
34 if _nidx_ > (#data)[1] then _nidx_ = 1 end
35
36 local sample = data[_nidx_]
37 local target = sample[{ {1} }]
38 local inputs = sample[{ {2,3} }]
39
40 # reset gradients
41 dl_dx:zero()
42
43 # evaluate the loss function and its derivative wrt x
44 local loss_x = criterion:forward(model:forward(inputs), target)
45 model:backward(inputs, criterion:backward(model.output, target))
46
47 # return loss(x) and dloss/dx
48 return loss_x, dl_dx
49 end
50
51 # define SGD
52 sgd_params = {
53   learningRate = 1e-3,
54   learningRateDecay = 1e-4,
55   weightDecay = 0,
56   momentum = 0
57 }
58 # we cycle 1e4 times over our training data
59 for i = 1,1e4 do
60   #this variable is used to estimate the average loss
61   current_loss = 0
```

Example: Linear Regression

```

62  #an epoch is a full loop over our training data
63  for i = 1,(#data)[1] do
64      # return new x and value of the loss functions
65      _, fs = optim.sgd(feval, x, sgd_params)
66      # update loss
67      current_loss = current_loss + fs[1]
68  end
69
70  # report average error on epoch
71  current_loss = current_loss / (#data)[1]
72  print('current loss = ' .. current_loss)
73
74  logger:add(['training error'] = current_loss)
75  logger:style(['training error'] = '-')
76  logger:plot()
77 end
78
79 # Test the trained model
80 text = {40.32, 42.92, 45.33, 48.85, 52.37, 57, 61.82, 69.78,
81         72.19, 79.42}
82
83 for i = 1,(#data)[1] do
84     local myPrediction = model:forward(data[i][{2,3}]))
85     print(string.format("%2d %6.2f %6.2f", i, myPrediction[1], text[i]))
86 end

```

Example: Two-Layer Network

```
1 import torch
2 from torch.autograd import Variable
3
4 # N is batch size; D_in is input dimension;
5 # H is hidden dimension; D_out is output dimension.
6 N, D_in, H, D_out = 64, 1000, 100, 10
7
8 # Create random Tensors to hold inputs and outputs, and wrap them in Variables.
9 x = Variable(torch.randn(N, D_in))
10 y = Variable(torch.randn(N, D_out), requires_grad=False)
11
12 # Use the nn package to define our model as a sequence of layers.
13 model = torch.nn.Sequential( torch.nn.Linear(D_in, H),
14                               torch.nn.ReLU(),
15                               torch.nn.Linear(H, D_out), )
16
17 # The nn package also contains definitions of popular loss functions;
18 loss_fn = torch.nn.MSELoss(size_average=False)
19
20 learning_rate = 1e-4
21
22 for t in range(500):
23     # Forward pass: compute predicted y by passing x to the model.
24     y_pred = model(x)
25     # Compute and print loss.
26     loss = loss_fn(y_pred, y)
27     print(t, loss.data[0])
28     # Zero the gradients before running the backward pass
29     model.zero_grad()
```

Example: Two-Layer Network

```
30     # Backward pass: compute gradient of the loss
31     loss.backward()
32
33     # Update the weights using gradient descent
34     for param in model.parameters():
35         param.data -= learning_rate * param.grad.data
36     end
37 end
```

<https://github.com/jcjohnson/pytorch-examples>

§5. Caffe

- 31 Programming interface
- 32 Example: Image Classification
- 33 Example: Extend Layers
- 34 Example: Extend Layers

Programming interface

- ① Expressive architecture
 - Define models and optimization by configuration without hard-coding
 - With protocol tool to define parameters for nets and solvers . . .
- ② Support GPUs
- ③ Mainly focus CNN for images
- ④ Not well documented

Example: Image Classification

```
1 import caffe
2 import matplotlib.pyplot as plt
3
4 # paste your image URL here
5 my_image_url = "https://wikipedia/Orang_Utan/2 C_Malaysia.JPG"
6 !wget -O image.jpg $my_image_url
7
8 # transform it and copy it into the net
9 image = caffe.io.load_image('image.jpg')
10 caffe.net.blobs['data'].data[...] = transformer.preprocess('data', image)
11
12 # perform classification
13 caffe.net.forward()
14
15 # obtain the output probabilities
16 output_prob = net.blobs['prob'].data[0]
17
18 # sort top five predictions from softmax output
19 top_inds = output_prob.argsort()[::-1][:5]
20
21 plt.imshow(image)
22 print 'probabilities and labels:'
23 zip(output_prob[top_inds], labels[top_inds])
```

Example: Extend Layers

```
1 import caffe
2 import numpy as np
3
4 class EuclideanLoss(caffe.Layer):
5     def setup(self, bottom, top):
6         #check input pair
7         if len(bottom) != 2:
8             raise Exception("Need two inputs to compute distance")
9
10    def reshape(self, bottom, top):
11        #check input dimensions match
12        if bottom[0].count != bottom[1].count:
13            raise Exception("Inputs must have the same dimension")
14        #difference in shape of inputs
15        self.diff = np.zeros_like(bottom[0].data, dtype=np.float32)
16        # loss output is scalar
17        top[0].reshape(1)
18
19    def forward(self, bottom, top):
20        self.diff[...] = bottom[0].data - bottom[1].data
21        top[0].data[...] = np.sum(self.diff**2)/bottom[0].num/2.
22
23    def backward(self, top, propagate_down, bottom):
24        for i in range(2):
25            if not propagate_down[i]:
26                continue
27            if i == 0:
28                sign = 1
29            else:
30                sign = -1
```


Example: Extend Layers

```
31 bottom[i].diff[...] = sign.self.diff / bottom[1].num
```

Define a class in Python to extend Layer

```
32 layer{  
33     type: "Python"  
34     python_param {  
35         module: "layers"  
36         layer: "EuclideanLoss"  
37     }  
38 }
```

https://docs.google.com/presentation/d/1UeKXVgRvvxg9OUdh_UiC5G71UMscNPlvArsWER41PsU/edit#slide=id.gc2fcdcce7_216_0

§6. TensorFlow

- 35 Computational graph
- 36 Programming interface
- 37 Visualization: TensorBoard
- 38 Example: SoftMax in TF
- 39 Example: SoftMax in TF

Computational graph

```
1 import tensorflow as tf
2
3 graph = tf.Graph()
4
5 with graph.as_default():
6     a = tf.constant(1.0)
7     tf.summary.scalar('aa', a)
8     b = tf.constant(2.0)
9     tf.summary.scalar('bb', b)
10    c = tf.multiply(a, b)
11    tf.summary.scalar('c', c)
12
13    merged = tf.summary.merge_all()
14    writer = tf.summary.FileWriter('./board', graph)
15
16 with tf.Session(graph=graph):
17     tf.global_variables_initializer().run()
18     writer.add_summary(merged.eval())
```



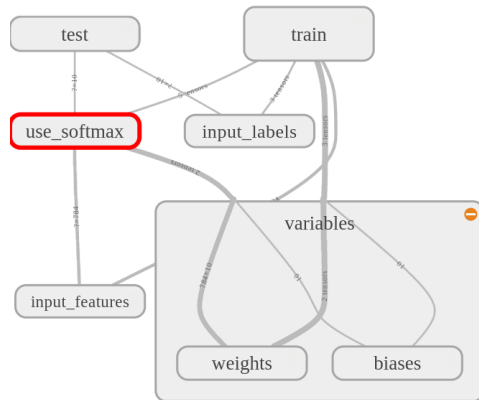
Programming interface

- **Graph:** In TensorFlow, machine learning algorithms are represented as computational graph. A computational or dataflow graph is a form of directed graph where vertices or nodes describe operations, while edges represent data flowing between these operations.
- **Operation:** An operation may represent a mathematical equation, a variable or constant, a control flow directive, a file I/O operation or even a network communication port.
- **Tensor:** A tensor is a multi-dimensional collection of homogeneous values with a fixed, static type.
- **Variable:** Variables can be described as persistent, mutable handles to in-memory buffers storing tensors.
- **Session:** In TensorFlow the execution of operations and evaluation of tensors may only be performed in a special environment called session.

Visualization: TensorBoard

Computation graphs are powerful but complicated

- thousands of nodes or more
- network is deep
- graph visualization tool TensorBoard is helpful



Example: SoftMax in TF

```
1 import tensorflow as tf
2
3 # Import training and test data
4 import tensorflow.examples.tutorials.mnist.input_data as input_data
5 mnist = input_data.read_data_sets("MNIST_data/", one_hot=True)
6
7 # Create a new graph
8 graph = tf.Graph()
9 with graph.as_default():
10     with tf.name_scope('input_features'):
11         # Placeholder for input variables (None = variable dimension)
12         x = tf.placeholder(tf.float32, shape=[None, 784],
13                             name='input_x')
14
15     with tf.name_scope('input_labels'):
16         # Placeholder for labels
17         y_ = tf.placeholder(tf.float32, shape=[None, 10],
18                             name='labels')
19
20     with tf.name_scope('variables'):
21         W = tf.Variable(tf.zeros([784, 10]), name='weights')
22         tf.summary.histogram('WEIGHTS', W)
23         b = tf.Variable(tf.zeros([10]), name='biases')
24         tf.summary.histogram('BIASES', b)
25
26     with tf.name_scope('use_softmax'):
27         # Apply softmax regression model to the input data and get prediction labels y
28         y = tf.nn.softmax(tf.matmul(x, W) + b)
```

Example: SoftMax in TF

```

29 with tf.name_scope('train'):
30     # Compute the cross entropy of real label y_ and prediction label y
31     cross_entropy = -tf.reduce_sum(y_*tf.log(y))
32     # Create a gradient-descent optimizer, minimizing the entropy, learning rate 0.01
33     train_step = tf.train.GradientDescentOptimizer(0.01).minimize(cross_entropy)
34
35 with tf.name_scope('test'):
36     correct_prediction = tf.equal(tf.argmax(y,1), tf.argmax(y_,1))
37     accuracy = tf.reduce_mean(tf.cast(correct_prediction, "float"))
38     tf.summary.scalar('Accuracy', accuracy)
39 merged = tf.summary.merge_all()
40 writer = tf.summary.FileWriter('./board', graph)
41
42 with tf.Session(graph=graph) as sess:
43     # Initialize all variables
44     tf.global_variables_initializer().run()
45
46 for step in range(1000):
47     if (step%10) == 0:
48         feed = {x: mnist.test.images, y_: mnist.test.labels}
49         _, acc = sess.run([merged, accuracy], feed_dict=feed)
50         print('Accuracy at %s step: %s' % (step, acc))
51     else:
52         batch_x, batch_y = mnist.train.next_batch(100)
53         sess.run(train_step, feed_dict={x: batch_x, y_: batch_y})
54         writer.add_summary(merged.eval(feed_dict={x: batch_x,
55                                                     y_: batch_y}), global_step=step)

```

§7. MXNet

40 Programming interface

41 Example: SoftMax in MXNet

42 Example: SoftMax in MXNet



Programming interface

- ① Mxnet.ndarray
 - Similar to numpy.ndarray
 - Supports both CPU and GPU
- ② Support building neural network graphs
 - Call `mx.viz.plot_network()`
- ③ Mixed programing
 - Support both imperative and declarative programming
- ④ Provide intermediate-level and high-level interface modules
- ⑤ Provide data parallelism with multi-devices
- ⑥ Provide abundant IO functions
- ⑦ Support many scope applications(e.g. computer vision, natural language processing, speech recognition , unsupervised machine learning, support embedded APIs, visualization)

Example: SoftMax in MXNet

```
1 import mxnet
2 import mxnet.symbol as sym
3 import numpy as np
4 import numpy.random as random
5 import time
6 from minpy.core import function
7 from minpy.core import grad_and_loss
8
9 # define softmax symbol
10 x_shape = (num_samples, num_classes)
11 label_shape = (num_samplesm,)
12 softmax_symbol = sym.SoftmaxOutput(data=sym.Variable('x'),
13                                   name='softmax', grad_scale=1.0/num_samples)
14
15 # convert MXNet symbol into a callable function
16 # corresponding gradient function
17 softmax = function(softmax_symbol, [('x', x_shape),
18                                   ('softmax_label', label_shape)])
19
20 # make softmax_label;
21 # MXNet's softmax operator does not use one-of-many label format
22 softmax_label = np.argmax(label, axis=1)
23
24 # Redefine loss function using softmax as one operator
25 def train_loss(w, x):
26     y = np.dot(x, w)
27     prob = softmax(x=y, softmax_label=softmax_label)
28     loss = -np.sum(label * np.log(prob)) / num_samples
29     return loss
```

Example: SoftMax in MXNet

```
30 # Initialize weight matrix (again)
31 weight = random.randn(num_features, num_classes)
32
33 # Calculate gradient function automatically
34 grad_function = grad_and_loss(train_loss)
35
36 # Now training it for 100 iterations
37 start_time = time.time()
38 for i in range(100):
39     dw, loss = grad_function(weight, data)
40     if i % 10 == 0:
41         print 'Iter {}, training loss {}'.format(i, loss)
42     weight -= 0.1 * dw
43 print 'Training time: {}'.format(time.time() - start_time)
```

Thank You!

