# LINEAR EQUATION REPORT

**School or Department:** Wu Han University

**Grade and Specialty:** Excellent engineer 2016

**Name:** Guo Yang, Kong Chuishun, Tang Rui

**Advisor:** Fangling Pu

May 7 2018

【Abstract】In this lab, we consider vector-valued linear and affine functions, and systems of linear equations. And we use a python program to familiar it. In the application, we introduce a method to train Binarized Neural Networks (BNNs) - neural networks with binary weights and activations.

【Key words】 linear equation; BNN, Multidimensional linear function model

# Table of Contents

# Chapter 1 Introduction

## 1.1 Linear and affine functions

1.1.1Vector-valued functions of vectors.

The notation f : $R^n \rightarrow R^m$ means that $f$ is a function that maps real n-vectors to real m-vectors. The value of the function $f$, evaluated at an n-vector x, is an m-vector f(x) = $(f_{1(x)}, f_{2(x)}, ..., f_{m(x)})$. Each of the components $fi$ of $f$ is itself a scalar-valued function of $x$. As with scalar-valued functions, we sometimes write $f(x) = f(x1, x2, ..., xn)$ to emphasize that f is a function of n scalar arguments. We use the same notation for each of the components of f, writing $fi(x) = fi(x_1, x_2, ..., x_n)$ to emphasize that fi is a function mapping the scalar arguments $x_1, ..., x_n$ into a scalar.

## 1.1.2 The matrix-vector product function

Suppose A is an m×n matrix. We can define a function f : $R^n \rightarrow R^m$ by $f(x) = Ax$. The inner product function $f : Rn \rightarrow R$, defined as $f(x) = a^T x$, , is the special case with m = 1.

## 1.1.3 Superposition and linearity.

The function $f : Rn \rightarrow Rm$, defined by $f(x) = Ax$, is linear, i.e., it satisfies the superposition property:
$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y) \tag{1.1}$$
holds for all n-vectors x and y and all scalars $\alpha$ and $\beta$. It is a good exercise to parse this simple looking equation, since it involves overloading of notation. On the left-hand side, the scalar-vector multiplications $\alpha x$ and $\beta y$ involve n-vectors, and the sum αx+βy is the sum of two n-vectors. The function f maps n-vectors to m-vectors, so $f(\alpha x + \beta y)$ is an m-vector. On the right-hand side, the scalar-vector multiplications and the sum are those for m-vectors. Finally, the equality sign is equality between two m-vectors.

We can verify that superposition holds for f using properties of matrix-vector and scalar-vector multiplication:
$$f(\alpha x + \beta y) = A(\alpha x + \beta y)$$
$$= A(\alpha x) + A(\beta y)$$
$$= \alpha(Ax) + \beta(Ay)$$
$$= \alpha f(x) + \beta f(y)$$
Thus we can associate with every matrix $A$ a linear function $f(x) = Ax$.

The converse is also true. Suppose f is a function that maps n-vectors to m vectors, and is linear, i.e., (1.1) holds for all n-vectors x and y and all scalars $\alpha$ and $\beta$. Then there exists an m×n matrix A such that $f(x) = Ax$ for all x. This can be shown in the same way as for scalar-valued functions, by showing that if f is linear, then
$$f(x) = x_1 f_{(e1)} + x_2 f_{(e2)} + \cdots + x_n f_{(en)} \qquad 1.2$$
where $e_k$ is the kth unit vector of size n. The right-hand side can also be written as a matrix-vector product Ax, with
$$A = [\, f_{(e1)}\, f_{(e2)} \cdots f_{(en)}\,]$$

It is easily shown that the matrix-vector representation of a linear function is unique. If f : $R^n \rightarrow R^m$ is a linear function, then there exists exactly one matrix A such that f(x) = Ax for all x.

## 1.1.4 Examples of linear functions

In the examples below we define functions $f$ that map n-vectors $x$ to $n$-vectors $f(x)$. Each function is described in words, in terms of its effect on an arbitrary x. In each case we give the associated matrix multiplication representation.

• Negation. f changes the sign of $x$: $f(x) = -x$.
Negation can be expressed as $f(x) = Ax$ with $A = -I$.

• Reversal. f reverses the order of the elements of $x$: $f(x) = (xn, xn-1, ..., x1)$. The reversal function can be expressed as $f(x) = Ax$ with

$$A = \begin{pmatrix} 0 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 0 \end{pmatrix}$$

(This is the n×n identity matrix with the order of its columns reversed.

## 1.1.5 Affine functions.

A vector-valued function f : $R_n \rightarrow R_m$ is called affine if it can be expressed as f(x) = Ax + b, where A is an m×n matrix and b is an m-vector. It can be shown that a function f : Rn $\rightarrow$ Rm is affine if and only if

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y)$$

holds for all n-vectors x, y, and all scalars α, β that satisfy $\alpha + \beta = 1$. In other words, superposition holds for affine combinations of vectors. (For linear functions, superposition holds for any linear combinations of vectors.) The matrix A and the vector b in the representation of an affine function as f(x) = Ax + b are unique. These parameters can be obtained by evaluating f at the vectors 0, $e1, ..., en$, where $ek$ is the kth unit vector in Rn. We have

$$A = [ f(e1) - f(0) \ f(e2) - f(0) \ \cdots \ f(en) - f(0) ], \quad b = f(0).$$

Just like affine scalar-valued functions, affine vector-valued functions are often called linear, even though they are linear only when the vector b is zero.

## 1.2 Systems of linear equations

Consider a set (also called a system) of $m$ linear equations in $n$ variables or unknowns $x_1, ... ... x_n$,

$$A_{11}x_1 + A_{12}x_2 + \cdots + A_{1n}x_n = b_1$$
$$A_{21}x_1 + A_{22}x_2 + \cdots + A_{2n}x_n = b_2$$
$$.$$
$$.$$
$$.$$
$$A_{m1}x_1 + A_{m2}x_2 + \cdots + A_{mn}x_n = b_m$$

The numbers $A_{ij}$ are called the coefficients in the linear equations, and the numbers $b_i$ are called the right-hand sides (since by tradition, they appear on the right-hand side of the equation). These equations can be written succinctly in matrix notation as

$$Ax = b$$

In this context, the $m \times n$ matrix A is called the *coefficient matrix,* and the *m*-vector b is called the *right-hand* side. An *n*-vector x is called a *solution* of the linear equations if $Ax = b$ holds. A set of linear equations can have no solutions, one solution, or multiple solutions.

The set of linear equations is called over-determined if $m > n$, under-determined if $m < n$, and square if $m = n$; these correspond to the coefficient matrix being tall, wide, and square, respectively. When the system of linear equations is over-determined, there are more equations than variables or unknowns. When the system of linear equations is under-determined, there are more unknowns than equations. When the system of linear equations is square, the numbers of unknowns and equations is the same. A set of equations with zero right-hand side, $Ax = 0$, is called a homogeneous set of equations. Any homogeneous set of equations has $x = 0$ as a solution.

Base on the basic conception, we have several usages. For example, *Leontief input-output model*.

We consider an economy with n different industrial sectors. We let $x_i$ be the economic activity level, or total production output, of sector *i*, for $i = 1, ..., n$, measured in a common unit, such as (billions of) dollars. The output of each sector flows to other sectors, to support their production, and also to consumers. We denote the total consumer demand for sector *i* as $d_i$, for $i = 1, ..., n$. Supporting the output level xj for sector j requires $A_{ij}x_j$ output for sector *i*.

We refer to $A_{ij}x_j$ as the sector i input that flows to sector *i*. (We can have $A_{ii} \neq 0$; for example, it requires some energy to support the production of energy.) Thus, $A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{in}x_n$ is the total sector *i* output required by, or flowing into, the n industrial sectors. The matrix *A* is called the input-output matrix of the economy, since it describes the flows of sector outputs to the inputs of itself and other sectors. The vector $A_x$ gives the sector outputs required to support the production levels given by *x*.

# Chapter 2　Our experiment

## 2.1 Brief Introduction

Based on what we have learned and what we three had done in DSP Laboratory, we want to establish a multidimensional linear function model to train *Cifar10* dataset using binarized neural network.

*Cifar10* is a collection of data sets for universal object recognition. It is an advanced science project research institute led by the Canadian government. *Cifar10* consists of 60, 000 32*32 RGB color images in 10 categories. 50,000 training, 10,000 tests (cross validation). The biggest feature of this dataset is that it migrates to universal objects and applies to multiple categories.

## 2.2 Specific operation

As follow figure, we login DSP laboratory by Xshell5 and have installed our lab environment including annocoda2.7, pylearn2, lasagna, Theano, Cuda and Pycuda. It cost me two days to config all environment well and that tortures me a lot.

Anyway, I successfully solve those problems. Finally we can run our machine learning code using the tesla k80 GPU. Following are our notes about this lab code(train cifar10 dataset using binarized neural network).

```
Last login: Sat May  5 11:40:30 2018 from 127.0.0.1
ctr@:~$ cd ./Desktop/BinaryNet/Train-time/
ctr@:~/Desktop/BinaryNet/Train-time$ ll
total 144220
drwxrwxr-x 2 ctr ctr      4096 5月   5 16:41 ./
drwxrwxr-x 5 ctr ctr      4096 5月   2 17:05 ../
-rw-rw-r-- 1 ctr ctr     10958 5月   2 17:05 binary_net.py
-rw-rw-r-- 1 ctr ctr      9170 5月   2 21:22 binary_net.pyc
-rw-rw-r-- 1 ctr ctr      2157 5月   4 10:49 cifar10.log
-rw-rw-r-- 1 ctr ctr     10814 5月   3 23:30 cifar10.py
-rw-rw-r-- 1 ctr ctr     98614 5月   5 16:35 gpu_mnist.log
-rw-rw-r-- 1 ctr ctr      6973 5月   5 12:05 mnist_gpu.py
-rw-rw-r-- 1 ctr ctr     12288 5月   5 09:11 mnist.log
-rw-rw-r-- 1 ctr ctr 147476810 5月   5 16:41 mnist_parameters.npz
-rw-rw-r-- 1 ctr ctr      6941 5月   2 17:05 mnist.py
-rw-rw-r-- 1 ctr ctr      1531 5月   2 17:05 README.md
-rw-rw-r-- 1 ctr ctr     10854 5月   2 17:05 svhn.py
ctr@:~/Desktop/BinaryNet/Train-time$
```

```
1.     #################################################
2.     #
3.     # 结构:
4.     #
5.     # 输入层
6.     #   |
7.     # 卷积层 1 + BN 层 + 激活层
8.     #   |
9.     # 卷积层 2 + 最大值混合层 2 + BN 层 + 激活层    (binary)
10.    #   |
11.    # 卷积层 3 + BN 层 + 激活层
12.    #   |
13.    # 卷积层 4 + 最大值混合层 4 + BN 层 + 激活层
14.    #   |
15.    # 卷积层 5 + BN 层 + 激活层
16.    #   |
17.    # 卷积层 6 + 最大值混合层 6 + BN 层 + 激活层
18.    #   |
19.    # 全连接层 1 + BN 层 + 激活层   (binary)
20.    #   |
21.    # 全连接层 2 + BN 层 + 激活层
22.    #   |
23.    # 全连接层 3 + BN 层
24.    #   |
25.    # 输出
26.    #
27.    #################################################
```

```python
28.
29.    # theano.function(inputs, outputs=None, mode=None, updates=None, givens=
       None):
30.    #    根据输入 inputs 计算数据 outputs 的函数，其中：
31.    #    inputs:  列表类型，用来保存输入量
32.    #    outputs: 列表或字典类型，用来保存输出量。输入量与输出量的映射关系通常在输
       出量的定义中体现。
33.    #    updates: 一组可迭代更新的量(shared_variable, new_expression),
34.    #          对其中的 shared_variable 输入用 new_expression 表达式更新
35.    #
36.    # lasagne.updates.adam(loss_or_grads, params, learning_rate): 用于参数更
       新，
37.    #    其中: loss_or_grads: 误差或梯度
38.    #    params:              要更新的参数
39.    #    learning_rate:       更新速率(学习速率)
40.    # lasagne.layers.get_output(layer):            对指定网络，计算网络输出;
41.    # lasagne.objectives.categorical_crossentropy(predictions, targets): 计算
       分类结果与目标的交叉熵(误差);
42.    # lasagne.layers.get_all_params(layer):        返回一个列表，包含该层参数的
       theano 共享变量或表达式 ?
43.    # bnn_utils.compute_grads(loss, network):      计算梯度;
44.    # bnn_utilsclipping_scaling(updates, network): 该函数在参数更新后规范化;
45.    # OrderedDict:                                 有序字典类;
46.    # dict.items():                                返回字典的键值;
47.    # tensor.nep():                                相当于"a != b";
48.    # tensor.argmax():                             返回沿指定轴取得最大值的下
       标;
49. from __future__ import print_function
50.
51. import sys
52. import os
53. import time
54. import numpy as np
55. np.random.seed(1234) # for reproducibility?
56.
57. # specifying the gpu to use
58. # import theano.sandbox.cuda
59. # theano.sandbox.cuda.use('gpu1')
60. import theano
61. import theano.tensor as T
62. import lasagne
63. import cPickle as pickle
64. import gzip
65.
```

```python
66. import binary_net
67.
68. from pylearn2.datasets.zca_dataset import ZCA_Dataset
69. from pylearn2.datasets.cifar10 import CIFAR10
70. from pylearn2.utils import serial
71.
72. from collections import OrderedDict
73.
74. if __name__ == "__main__":
75.
76.     # BN parameters (设定超参数)
77.     batch_size = 50
78.     print("batch_size = "+str(batch_size))
79.     # alpha is the exponential moving average factor
80.     alpha = .1
81.     print("alpha = "+str(alpha))
82.     epsilon = 1e-4
83.     print("epsilon = "+str(epsilon))
84.
85.     # BinaryOut
86.     activation = binary_net.binary_tanh_unit
87.     print("activation = binary_net.binary_tanh_unit")
88.     # activation = binary_net.binary_sigmoid_unit
89.     # print("activation = binary_net.binary_sigmoid_unit")
90.
91.     # BinaryConnect     (设定 BNN-神经网络标识)
92.     binary = True
93.     print("binary = "+str(binary))
94.     stochastic = False
95.     print("stochastic = "+str(stochastic))
96.     # (-H,+H) are the two binary values
97.     # H = "Glorot"
98.     H = 1.
99.     print("H = "+str(H))
100.      # W_LR_scale = 1.
101.      W_LR_scale = "Glorot" # "Glorot" means we are using the coefficients fr
    om Glorot's paper
102.     print("W_LR_scale = "+str(W_LR_scale))
103.
104.     # Training parameters
105.     num_epochs = 500
106.     print("num_epochs = "+str(num_epochs))
107.
108.     # Decaying LR    (设定学习速率)
```

```
109.    LR_start = 0.001
110.    print("LR_start = "+str(LR_start))
111.    LR_fin = 0.0000003
112.    print("LR_fin = "+str(LR_fin))
113.    LR_decay = (LR_fin/LR_start)**(1./num_epochs)
114.    print("LR_decay = "+str(LR_decay))
115.    # BTW, LR decay might good for the BN moving average...
116.
117.    train_set_size = 45000
118.    print("train_set_size = "+str(train_set_size))
119.    shuffle_parts = 1
120.    print("shuffle_parts = "+str(shuffle_parts))
121.
122.    print('Loading CIFAR-10 dataset...')    （加载数据）
123.
124.    train_set = CIFAR10(which_set="train",start=0,stop = train_set_size)
125.    valid_set = CIFAR10(which_set="train",start=train_set_size,stop = 50000
    )
126.    test_set = CIFAR10(which_set="test")
127.
128.    # bc01 format
129.    # Inputs in the range [-1,+1]
130.    # print("Inputs in the range [-1,+1]")
131.    train_set.X = np.reshape(np.subtract(np.multiply(2./255.,train_set.X),1
    .),(-1,3,32,32))
132.    valid_set.X = np.reshape(np.subtract(np.multiply(2./255.,valid_set.X),1
    .),(-1,3,32,32))
133.    test_set.X = np.reshape(np.subtract(np.multiply(2./255.,test_set.X),1.)
    ,(-1,3,32,32))
134.
135.    # flatten targets
136.    train_set.y = np.hstack(train_set.y)
137.    valid_set.y = np.hstack(valid_set.y)
138.    test_set.y = np.hstack(test_set.y)
139.
140.    # Onehot the targets
141.    train_set.y = np.float32(np.eye(10)[train_set.y])
142.    valid_set.y = np.float32(np.eye(10)[valid_set.y])
143.    test_set.y = np.float32(np.eye(10)[test_set.y])
144.
145.    # for hinge loss
146.    train_set.y = 2* train_set.y - 1.
147.    valid_set.y = 2* valid_set.y - 1.
148.    test_set.y = 2* test_set.y - 1.
```

```
149.
150.    print('Building the CNN...')      (构建深度网络)
151.
152.    # Prepare Theano variables for inputs and targets
153.    input = T.tensor4('inputs')
154.    target = T.matrix('targets')
155.    LR = T.scalar('LR', dtype=theano.config.floatX)
156.
157.    cnn = lasagne.layers.InputLayer()
158.
159.    # 128C3-128C3-P2
160.    cnn = binary_net.Conv2DLayer()
161.    cnn = lasagne.layers.BatchNormLayer()
162.    cnn = lasagne.layers.NonlinearityLayer()
163.
164.    cnn = binary_net.Conv2DLayer()
165.    cnn = lasagne.layers.MaxPool2DLayer(cnn, pool_size=(2, 2))
166.    cnn = lasagne.layers.BatchNormLayer()
167.    cnn = lasagne.layers.NonlinearityLayer()
168.
169.    # 256C3-256C3-P2
170.    cnn = binary_net.Conv2DLayer()
171.    cnn = lasagne.layers.BatchNormLayer()
172.    cnn = lasagne.layers.NonlinearityLayer()
173.
174.    cnn = binary_net.Conv2DLayer()
175.    cnn = lasagne.layers.MaxPool2DLayer(cnn, pool_size=(2, 2))
176.    cnn = lasagne.layers.BatchNormLayer()
177.    cnn = lasagne.layers.NonlinearityLayer()
178.
179.    # 512C3-512C3-P2
180.    cnn = binary_net.Conv2DLayer()
181.    cnn = lasagne.layers.BatchNormLayer()
182.    cnn = lasagne.layers.NonlinearityLayer()
183.
184.    cnn = binary_net.Conv2DLayer()
185.    cnn = lasagne.layers.MaxPool2DLayer(cnn, pool_size=(2, 2))
186.    cnn = lasagne.layers.BatchNormLayer()
187.    cnn = lasagne.layers.NonlinearityLayer()
188.
189.    # print(cnn.output_shape)
190.    # 1024FP-1024FP-10FP
191.    cnn = binary_net.DenseLayer()
192.    cnn = lasagne.layers.BatchNormLayer()
```

```python
193.    cnn = lasagne.layers.NonlinearityLayer()
194.
195.    cnn = binary_net.DenseLayer()
196.    cnn = lasagne.layers.BatchNormLayer()
197.    cnn = lasagne.layers.NonlinearityLayer()
198.
199.    cnn = binary_net.DenseLayer()
200.    cnn = lasagne.layers.BatchNormLayer()
201.    #(计算网络输出)
202.    train_output = lasagne.layers.get_output(cnn, deterministic=False)
203.
204.    # squared hinge loss    (定义误差计算函数)
205.    loss = T.mean(T.sqr(T.maximum(0.,1.-target*train_output)))
206.
207.    if binary:
208.
209.        # W updates    (权重更新)
210.        W = lasagne.layers.get_all_params(cnn, binary=True)
211.        W_grads = binary_net.compute_grads(loss,cnn)
212.        updates = lasagne.updates.adam(loss_or_grads=W_grads, params=W, lea
    rning_rate=LR)
213.        updates = binary_net.clipping_scaling(updates,cnn)
214.
215.        # other parameters updates    (其它参数更新)
216.        params = lasagne.layers.get_all_params(cnn, trainable=True, binary=
    False)
217.        updates = OrderedDict(updates.items() + lasagne.updates.adam(loss_o
    r_grads=loss, params=params, learning_rate=LR).items())
218.
219.    else:
220.        # (参数更新)
221.        params = lasagne.layers.get_all_params(cnn, trainable=True)
222.        updates = lasagne.updates.adam(loss_or_grads=loss, params=params, l
    earning_rate=LR)
223.
224.    # test prediction and loss expressions (测试数据和误差表示)
225.    test_output = lasagne.layers.get_output(cnn, deterministic=True)
226.    test_loss = T.mean(T.sqr(T.maximum(0.,1.-target*test_output)))
227.    test_err = T.mean(T.neq(T.argmax(test_output, axis=1), T.argmax(target,
    axis=1)),dtype=theano.config.floatX)
228.
229.    (建立训练数据和测试数据的 theano 函数)
230.    # Compile a function performing a training step on a mini-
    batch (by giving the updates dictionary)
```

```
231.      # and returning the corresponding training loss:
232.      train_fn = theano.function([input, target, LR], loss, updates=updates)

233.
234.      # Compile a second function computing the validation loss and accuracy:

235.      val_fn = theano.function([input, target], [test_loss, test_err])
236.
237.      print('Training...')
238.      binary_net.train(
239.              train_fn,val_fn,              # 训练和测试数据的 theano 函数
240.              cnn,                          # 神经网络
241.              batch_size,                   # 小批量数据大小
242.              LR_start,LR_decay,            # 学习速率
243.              num_epochs,                   # 迭代期次数
244.              train_set.X,train_set.y,      # 训练数据
245.              valid_set.X,valid_set.y,      # 验证数据
246.              test_set.X,test_set.y,        # 测试数据
247.              shuffle_parts=shuffle_parts)
```

## 2.3 Analysis about Results

But this code need to run two days. So, we can't write the answer into our report. Fortunately, we also utilize the *mnist* version using binarized neutral network. And it only cost 8 ours to get the parameters. (*mnist* is also a train set which is like *Cifar10*. But it is simpler.)

This is the 954-epoch answer on *mnist* datasets. The k80 GPU memory is too busy。

```
Epoch 954 of 1000 took 110.964931011s
  LR:                          4.62510135885e-07
  training loss:               0.00305322502526
  validation loss:             0.0107206741194
  validation error rate:       1.04999997839%
  best epoch:                  740
  best validation error rate:  0.959999978542%
  test loss:                   0.00897184130102
  test error rate:             0.909999981523%
Traceback (most recent call last):
  File "mnist_gpu.py", line 217, in <module>
    shuffle_parts)
  File "/home/user/ctr/Desktop/BinaryNet/Train-time/binary_net.py", line 290, in train
    train_loss = train_epoch(X_train,y_train,LR)
  File "/home/user/ctr/Desktop/BinaryNet/Train-time/binary_net.py", line 256, in train_epoch
    loss += train_fn(X[i*batch_size:(i+1)*batch_size],y[i*batch_size:(i+1)*batch_size],LR)
  File "/home/user/ctr/anaconda2/lib/python2.7/site-packages/theano/compile/function_module.py",
line 871, in __call__
    storage_map=getattr(self.fn, 'storage_map', None))
  File "/home/user/ctr/anaconda2/lib/python2.7/site-packages/theano/gof/link.py", line 314, in
raise_with_op
    reraise(exc_type, exc_value, exc_trace)
  File "/home/user/ctr/anaconda2/lib/python2.7/site-packages/theano/compile/function_module.py",
line 859, in __call__
    outputs = self.fn()
MemoryError: Error allocating 67108864 bytes of device memory (out of memory).
```

# Bibliography

[1] Binarized Neural Networks : Training Neural Networks with Weights and

Activations Constrained to +1 or−1

[2]Bahdanau,Dzmitry,Cho,Kyunghyun,andBengio,Yoshua. Neural machine translation by jointly learning to align and translate. In ICLR'2015, arXiv:1409.0473, 2015.