

# MACHINE LEARNING

ERNEST YEUNG [ERNESTYALUMNI@GMAIL.COM](mailto:ERNESTYALUMNI@GMAIL.COM)

From the beginning of 2016, I decided to cease all explicit crowdfunding for any of my materials on physics, math. I failed to raise *any* funds from previous crowdfunding efforts. I decided that if I was going to live in *abundance*, I must lose a scarcity attitude. I am committed to keeping all of my material **open-sourced**. I give all my stuff *for free*.

In the beginning of 2017, I received a very generous donation from a reader from Norway who found these notes useful, through *PayPal*. If you find these notes useful, feel free to donate directly and easily through [PayPal](#), which won't go through a 3rd. party such as indiegogo, kickstarter, patreon. Otherwise, under the *open-source MIT license*, feel free to copy, edit, paste, make your own versions, share, use as you wish.

gmail : ernestyalumni  
linkedin : ernestyalumni  
tumblr : ernestyalumni  
twitter : ernestyalumni  
youtube : ernestyalumni

## CONTENTS

<b>Part 1. Data; Data Wrangling, Data cleaning, Web crawling, Data input</b>	2
1. Sample, example data; input data	2
1.1. sklearn, from sci-kit learn, sample data, datasets	2
<b>Part 2. Introduction</b>	2
1.2. Supervised Learning	2
2. Deep Learning	3
3. Parallel Computing	3
3.1. Udacity Intro to Parallel Programming : Lesson 1 - The GPU Programming Model	3
4. Pointers in C; Pointers in C categorified (interpreted in Category Theory)	7
<b>Part 3. Machine Learning with Deep Learning</b>	8
5. Feedforward; Feedforward Propagation and Prediction	10
6. Backpropagation; Backpropagation algorithm	10
6.1. Cost functional	12
7. Universal approximation theorem	12
8. LSTM; Long Short Term Memory	12
8.1. How to choose the number of hidden layers and nodes in a neural net	13
<b>Part 4. Support Vector Machines (SVM)</b>	13
9. From linear classifier as a hyperplane, (big) margin, to linear support vector machine (SVM), and Lagrangian dual (i.e. conjugate variables, conjugate momenta)	13
10. So-called “Kernel trick”; feature space is a Hilbert space	14
10.1. Dealing with Errors, (non-negative) slack variables, dealing with not-necessarily perfectly separable data	14
11. Dual Formulation	14
11.1. Implementation	14
12. Support Vector Machines (SVM) natively implemented in theano, entirely on the GPU with the CUDA backend, with constrained gradient descent	15

*Date:* 24 avril 2016.

*Key words and phrases.* Machine Learning, statistical inference, statistical inference learning.

12.1.	Executive Summary	15
12.2.	Motivations and Introductions	15
12.3.	Concise Mathematical Review/Summary of the theory for SVM	16
12.4.	So-called “Kernel trick”; feature space is a Hilbert space	17
12.5.	Dual Formulation	17
12.6.	Constrained Optimization	18
12.7.	Constrained Gradient Descent (Implementation)	19
12.8.	Immediate Results from training on sample datasets	20
12.9.	Conclusions/Summary/Dictionary between Math and Code	20
<b>Part 5.</b>	<b>Notes</b>	20
	References	21

ABSTRACT. Everything about Machine Learning.

**Part 1. Data; Data Wrangling, Data cleaning, Web crawling, Data input**

1. SAMPLE, EXAMPLE DATA; INPUT DATA

1.1. **sklearn, from sci-kit learn, sample data, datasets.** cf. `sampleinputdataX_sklearn.ipynb`

For  $j = 0, 1, \dots, d - 1$ ,  $d =$  number of “features”,

$$x_i^{(j)} \in (\mathbb{R}^N)^d = \underbrace{\mathbb{R}^N \times \mathbb{R}^N \times \dots \times \mathbb{R}^N}_d$$

e.g.  $N = 442$  (number of given observations/data)

$y_i \in \mathbb{R}^N$  (represents target or result)

Given data  $(x_i^{(j)}, y_i) \in (\mathbb{R}^N)^d \times \mathbb{R}^N$ ,

we can restrict data  $(x_i^{(j)}, y_i)$  to subsets to train and test, for training and testing.

So let  $I_{\text{train}}, I_{\text{test}} \subset \{0, 1, \dots, N - 1\}$  s.t.  $I_{\text{train}} \cap I_{\text{test}} = \emptyset$ .

*Want:*

$$\begin{aligned} (x_i^{(j)}, y_i)_{i \in I_{\text{train}}} &\mapsto \theta_\alpha \\ (\mathbb{R}^{|I_{\text{train}}|})^d \times \mathbb{R}^{|I_{\text{train}}|} &\rightarrow \mathbb{R}^{|d|} \end{aligned}$$

and so further, I think the idea is

$$\begin{aligned} (x_i^{(j)}, y_i)_{i \in I_{\text{test}}} &\xrightarrow{L_{\theta_\alpha}} L_{\theta_\alpha}(\theta_\alpha(x_i^{(j)}, y_i)) \\ (\mathbb{R}^{|I_{\text{test}}|})^d \times \mathbb{R}^{|I_{\text{test}}|} &\rightarrow \mathbb{R} \end{aligned}$$

**Part 2. Introduction**

1.1.1. *Terminology.*

inputs  $\equiv$  independent variables  $\equiv$  predictors (cf. statistics)  $\equiv$  features (cf. pattern recognition)

outputs  $\equiv$  dependent variables  $\equiv$  responses

cf. Chapter 2 Overview of Supervised Learning, Section 2.1 Introduction of Hastie, Tibshirani, and Friedman (2009) [1]

cf. Chapter 2 Overview of Supervised Learning, Section 2.2 Variable Types and Terminology of Hastie, Tibshirani, and Friedman (2009) [1]

<sup>1</sup><sub>nlab</sub> FinSet <https://ncatlab.org/nlab/show/FinSet>

1.1.2. *FinSet.*

The category  $\text{FinSet} \in \text{Cat}$  is the category of all finite sets (i.e.  $\text{Obj}(\text{FinSet}) \equiv$  all finite sets) and all functions in between them; note that  $\text{FinSet} \subset \text{Set}$  <sup>1</sup>

Recall that the  $\text{FinSet}$  *skeletal* is

1.2. **Supervised Learning.** cf. <http://cs229.stanford.edu/notes/cs229-notes1.pdf>

Consider data to belong to the category of all possible data:

$$\text{Data} \equiv \text{Dat} = (\text{Obj}(\text{Dat}), \text{MorDat}, 1, \circ), \quad \text{Dat} \in \text{Cat}$$

Consider the **training set**:

$$\text{training set} := \{(x^{(i)}, y^{(i)}) | i = 1 \dots m, x^{(i)} \in \mathcal{X}, y^{(i)} \in \mathcal{Y}\}$$

where  $\mathcal{X}$  is a manifold (it can be topological or smooth, EY:20160502 I don’t know exactly because I need to check the topological and/or differential structure);  $\mathcal{Y} \in \text{Obj}(\text{FinSet})$ , or  $(\mathcal{Y} \in \text{Obj}(\text{Top}))$ (or  $\mathcal{Y} \in \text{Obj}(\text{Man}))$ .

So training set  $\subset \mathcal{X} \times \mathcal{Y} \in \text{Obj}(\text{Dat})$ .

I propose that there should be a functor  $H$  that represents the “learning algorithm”:

$$\text{Dat} \xrightarrow{H} \text{ML}$$

s.t.

$$H : \mathcal{X} \times \mathcal{Y} \rightarrow \text{Hom}(\mathcal{X}, \mathcal{Y})$$

$$H(\text{training set}) = H(\{(x^{(i)}, y^{(i)}) | i = 1 \dots m\}) = h$$

When  $\mathcal{Y} \in \text{Obj}(\text{FinSet})$ , *classification*.

When  $\mathcal{Y} \in \text{Obj}(\text{Top})$  (or  $\text{Obj}(\text{Man})$ ), *regression*.

1.2.1. *Linear Regression.* Keeping in mind

$$\text{Dat} \xrightarrow{H} \text{ML}$$

Consider

$$h : \mathbb{R}^p \rightarrow \text{Hom}(\mathcal{X}, \mathcal{Y})$$

$$h : \theta \mapsto h_\theta$$

s.t.

$$h_\theta : \mathcal{X} \rightarrow \mathcal{Y}$$

so (possibly)  $h \in \text{ObjML}$  (or is  $h$  part of the functor  $H$ ?)

Consider the cost function  $J$

$$J : \mathbb{R}^p \rightarrow \text{Hom}(\mathfrak{X} \times \mathfrak{Y}, \mathbb{R}) = C^\infty(\mathcal{X} \times \mathcal{Y})$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

1.2.2. *LMS algorithm (least mean square (or Widrow-Hoff learning rule))*. Define **gradient descent** algorithm:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

with  $:=$  being assignment (I'll use  $:=$  for “define”, in mathematical terms, use context to distinguish the 2), where  $\alpha$  is the *learning rate*.

Rewriting the above,

$$\theta := \theta - \alpha \text{grad} J(\theta)$$

where  $\text{grad} : C^\infty(M) \rightarrow \mathfrak{X}(M)$ , with  $M$  being a smooth manifold.

This is *batch gradient descent*:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) = \theta_j - \alpha \frac{\partial}{\partial \theta_j} \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \left( \frac{\partial h_\theta(x^{(i)})}{\partial \theta} \right)$$

Simply notice how the entire training set of  $m$  rows is used.

I will expound on the so-called distinguished object  $1 \xrightarrow{P} X$  on pp. 8, in Section 2 The Category of Conditional Probabilities of Culbertson and Sturtz (2013) [2] because it wasn't clear to me in the first place (the fault is mine; the authors wrote a very lucid and very fathomable, pedagogically-friendly exposition).

$\forall Y$  with indiscrete  $\sigma$ -algebra  $\Sigma_Y = \{Y, \emptyset\}$   
(remember,  $((Y, \Sigma_Y), \mu_Y)$ ,  $\mu_Y(\phi) = 0$ ,  $\mu_Y(Y) = 1$ ),

$\exists!$  unique morphism in  $\text{Mor}\mathcal{P}$ ,  $X \rightarrow Y$ , since

$\forall P : X \rightarrow Y$ ,  $P \in \text{Mor}\mathcal{P}$ ,  $P_x$  must be a probability measure on  $Y$ , because

$$(X, \Sigma_X) \xrightarrow{P} (Y, \Sigma_Y)$$

$$P : \Sigma_Y \times X \rightarrow [0, 1]$$

$$P(\cdot|x) : \Sigma_Y \rightarrow [0, 1] \equiv \begin{array}{l} P_x : \Sigma_Y \rightarrow [0, 1] \text{ s.t.} \\ P_x(\emptyset) = 0, P_x(Y) = 1 \end{array}$$

i.e. EY: 20160503, Given  $x \in X$  occurs,  $Y$  must occur.

By def. of terminal object ( $\forall (X, \Sigma_X) \in \text{Obj}\mathcal{P}$ ,  $\exists!$  morphism  $P$  s.t.  $(X, \Sigma_X) \xrightarrow{P} (Y, \Sigma_Y)$ ,  $Y$  *terminal* object, and denote unique morphism  $!_X : X \rightarrow Y$ ,  $!_X \in \text{Mor}\mathcal{P}$ .

Up to isomorphism, canonical terminal object is 1-element set denoted by  $1 = \{*\}$ , with the only possible  $\sigma$ -algebra ( $\mu(*) = 1$ ,  $\mu(\emptyset) = 0$ ),

$$\forall P : 1 \rightarrow X, P \in \text{Mor}\mathcal{P}, P \in \text{Hom}_{\mathcal{P}}(1, X), \forall X \in \text{Mor}\mathcal{P}$$

$P$  is an “absolute” probability measure on  $X$  because “there’s no variability (conditioning) possible within singleton set  $1 = \{*\}$ .” [2]

Now

$$P : \Sigma_X \times 1 \rightarrow [0, 1]$$

$$P(\cdot|*) : \Sigma_X \rightarrow [0, 1]$$

where  $P(\cdot|*) : \Sigma_X \rightarrow [0, 1]$  perfect probability measure on  $X$ ,  $P(\cdot|*) : \Sigma_X \rightarrow [0, 1] \equiv P_*$ , i.e.  $P(\cdot|*) = p(\cdot)$  (usual probability on  $X$ ).

$\forall A \in \Sigma_X$ ,  $P(A|\cdot) : 1 \rightarrow [0, 1]$ , but  $P(A|*) = P(A)$ ,  $P(A|\emptyset) = 0$ .

Refer to

$$1 \xrightarrow{P} X$$

morphism  $P : 1 \rightarrow X \in \text{Mor}\mathcal{P}$  as probability measure or distribution on  $X$ .

## 2. DEEP LEARNING

Deep Learning Tutorial [6]

## 3. PARALLEL COMPUTING

3.1. **Udacity Intro to Parallel Programming : Lesson 1 - The GPU Programming Model**. Owens and Luebki pound fists at the end of this video. =)))) **Intro to the class**.

3.1.1. *Running CUDA locally*. Also, **Intro to the class**, in Lesson 1 - The GPU Programming Model, has links to documentation for running CUDA locally; in particular, for Linux: **<http://docs.nvidia.com/cuda/cuda-getting-started-guide-for-linux/index.html>**. That guide told me to go download the NVIDIA CUDA Toolkit, which is the **<https://developer.nvidia.com/cuda-downloads>**.

For *Fedora*, I chose Installer Type **runfile (local)**.

Afterwards, installation of CUDA on Fedora 23 workstation had been nontrivial. Go see either my github repository **ML-grabbag** (which will be updated) or my **wordpress blog** (which may not be upgraded frequently).

$P = VI = I^2 R$  heating.

3.1.2. *Definitions of Latency and throughput (or bandwidth)*. cf. **Building a Power Efficient Processor**

**Latency vs Bandwidth**

latency [sec]. From the title “Latency vs. bandwidth”, I’m thinking that throughput = bandwidth (???). throughput = job/time (of job).

Given total task, velocity  $v$ ,

total task /  $v$  = latency. throughput = latency/(jobs per total task).

Also, in **Building a Power Efficient Processor**. Owens recommends the article David Patterson, “Latency...”

cf. **GPU from the Point of View of the Developer**

$n_{\text{core}} \equiv$  number of cores

$n_{\text{vecop}} \equiv (n_{\text{vecop}} - \text{wide axial vector operations} / \text{core core})$

$n_{\text{thread}} \equiv$  threads/core (hyperthreading)

$n_{\text{core}} \cdot n_{\text{vecop}} \cdot n_{\text{thread}}$  parallelism

There were various websites that I looked up to try to find out the capabilities of my video card, but so far, I’ve only found these commands (and I’ll print out the resulting output):

```
$ lspci -vnn | grep VGA -A 12
03:00.0 VGA compatible controller [0300]: NVIDIA Corporation GM200 [GeForce GTX 980 Ti] [10de:17c8] (rev a1) (prog-if 00 [VGA])
Subsystem: eVga.com. Corp. Device [3842:3994]
Physical Slot: 4
Flags: bus master, fast devsel, latency 0, IRQ 50
Memory at fa000000 (32-bit, non-prefetchable) [size=16M]
Memory at e0000000 (64-bit, prefetchable) [size=256M]
Memory at f0000000 (64-bit, prefetchable) [size=32M]
I/O ports at e000 [size=128]
[virtual] Expansion ROM at fb000000 [disabled] [size=512K]
Capabilities: <access denied>
Kernel driver in use: nvidia
Kernel modules: nouveau, nvidia
```

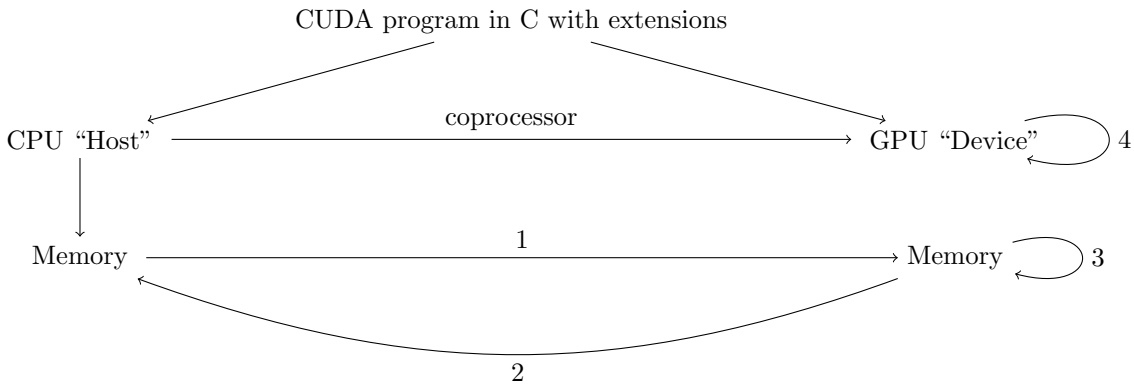
```
$ lspci | grep VGA -E
03:00.0 VGA compatible controller: NVIDIA Corporation GM200 [GeForce GTX 980 Ti] (rev a1)
```

```
$ grep driver /var/log/Xorg.0.log
[ 18.074] Kernel command line: BOOT_IMAGE=/vmlinuz-4.2.3-300.fc23.x86_64 root=/dev/mapper/fedora-root ro rd.lvm.lv=fedora
[ 18.087] (WW) Hotplugging is on, devices using drivers 'kbd', 'mouse' or 'vmmouse' will be disabled.
[ 18.087] X.Org XInput driver : 22.1
```

```
[ 18.192] (II) Loading /usr/lib64/xorg/modules/drivers/nvidia_drv.so
[ 19.088] (II) NVIDIA(GPU-0): Found DRM driver nvidia-drm (20150116)
[ 19.102] (II) NVIDIA(0):      ACPI event daemon is available, the NVIDIA X driver will
[ 19.174] (II) NVIDIA(0): [DRI2]      VDPAU driver: nvidia
[ 19.284]      ABI class: X.Org XInput driver, version 22.1
...

$ lspci -k | grep -A 8 VGA
03:00.0 VGA compatible controller: NVIDIA Corporation GM200 [GeForce GTX 980 Ti] (rev a1)
      Subsystem: eVga.com. Corp. Device 3994
      Kernel driver in use: nvidia
      Kernel modules: nouveau, nvidia
03:00.1 Audio device: NVIDIA Corporation GM200 High Definition Audio (rev a1)
      Subsystem: eVga.com. Corp. Device 3994
      Kernel driver in use: snd_hda_intel
      Kernel modules: snd_hda_intel
05:00.0 USB controller: VIA Technologies, Inc. VL805 USB 3.0 Host Controller (rev 01)
```

CUDA Program Diagram



CPU “host” is the boss (and issues commands) -Owen.  
Coprocessor : CPU “host” → GPU “device”  
Coprocessor : CPU process ↦ (co)-process out to GPU

With

- 1 data cpu → gpu
- 2 data gpu → cpu (initiated by cpu host)

- 1., 2., uses `cudaMemcpy`
- 3 allocate GPU memory: `cudaMalloc`
- 4 launch kernel on GPU

Remember that for 4., this launching of the kernel, while it’s acting on GPU “device” onto itself, it’s initiated by the boss, the CPU “host”.

Hence, cf. **Quiz: What Can GPU Do in CUDA**, GPUs can respond to CPU request to receive and send Data CPU → GPU and Data GPU → CPU, respectively (1,2, respectively), and compute a kernel launched by the CPU (3).

A CUDA Program A typical GPU program

- `cudaMalloc` - CPU allocates storage on GPU
- `cudaMemcpy` - CPU copies input data from CPU → GPU
- *kernel launch* - CPU launches kernel(s) on GPU to process the data
- `cudaMemcpy` - CPU copies results back to CPU from GPU

Owens advises minimizing “communication” as much as possible (e.g. the `cudaMemcpy` between CPU and GPU), and do a lot of computation in the CPU and GPU, each separately.

Defining the GPU Computation

Owens circled this

BIG IDEA 

This is Important

Kernels look like serial programs  
Write your program as if it will run on **one** thread  
The GPU will run that program on **many** threads

Squaring A Number on the CPU

Note

- (1) Only 1 thread of execution: (“thread” := one independent path of execution through the code) e.g. the `for` loop
- (2) no explicit parallelism; it’s serial code e.g. the `for` loop through 64 elements in an array

GPU Code A High Level View

CPU:

- Allocate Memory
- Copy Data to/from GPU
- Launch Kernel - species degree of parallelism

GPU:

- Express Out = In · In - says *nothing* about the degree of parallelism

Owens reiterates that in the GPU, everything looks serial, but it’s only in the CPU that anything parallel is specified.

pseudocode: CPU code: square kernel <<< 64 >>> (outArray,inArray)

Squaring Numbers Using CUDA Part 3

From the example

```
// launch the kernel
square<<<1, ARRAY_SIZE>>>(d_out , d_in)
```

we’re introduced to the “CUDA launch operator”, initiating a kernel of 1 block of 64 elements (`ARRAY_SIZE` is 64) on the GPU. Remember that `d_` prefix (this is naming convention) tells us it’s on the device, the GPU, solely.

With CUDA launch operator  $\equiv \langle \langle \langle \rangle \rangle \rangle$ , then also looking at this explanation on **stackexchange** (so surely others are confused as well, of those who are learning this (cf. **CUDA kernel launch parameters explained right?**). From **Eric**’s answer,

threads are grouped into blocks. all the threads will execute the invoked kernel function.

Certainly,

$$\langle \langle \langle \rangle \rangle \rangle: (n_{\text{block}}, n_{\text{threads}}) \times \text{kernelfunctions} \mapsto \text{kernelfunction} \langle \langle \langle n_{\text{block}}, n_{\text{threads}} \rangle \rangle \rangle \in \text{End} : \text{Dat}_{\text{GPU}}$$
$$\langle \langle \langle \rangle \rangle \rangle: \mathbb{N}^+ \times \mathbb{N}^+ \times \text{Mor}_{\text{GPU}} \rightarrow \text{EndDat}_{\text{GPU}}$$

where I propose that GPU can be modeled as a category containing objects  $\text{Dat}_{\text{GPU}}$ , the collection of all possible data inputs and outputs into the GPU, and  $\text{Mor}_{\text{GPU}}$ , the collection of all kernel functions that run (exclusively, and this *must* be the class, as reiterated by Prof. Owen) on the GPU.

Next,

$$\text{kernelfunction} \langle \langle \langle n_{\text{block}}, n_{\text{threads}} \rangle \rangle \rangle: \text{din} \mapsto \text{dout} \quad (\text{as given in the “square” example, and so I propose})$$

$$\text{kernelfunction} \langle \langle \langle n_{\text{block}}, n_{\text{threads}} \rangle \rangle \rangle: (\mathbb{N}^+)^{n_{\text{threads}}} \rightarrow (\mathbb{N}^+)^{n_{\text{threads}}}$$

But keep in mind that `dout`, `din` are pointers in the C program, pointers to the place in the memory.

`cudaMemcpy` is a functor category, s.t. e.g.  $\text{Obj}_{\text{CudaMemcpy}} \ni \text{cudaMemcpyDeviceToHost}$  where

$$\text{cudaMemcpy}(-, -, n_{\text{thread}}, \text{cudaMemcpyDeviceToHost}) : \text{Memory}_{\text{GPU}} \rightarrow \text{Memory}_{\text{CPU}} \in \text{Hom}(\text{Memory}_{\text{GPU}}, \text{Memory}_{\text{CPU}})$$

### Squaring Numbers Using CUDA 4

Note the C language construct *declaration specifier* - denotes that this is a kernel (for the GPU) and not CPU code. Pointers need to be allocated on the GPU (otherwise your program will crash spectacularly -Prof. Owen).

3.1.3. *What are C pointers?* Is  $\langle \text{type} \rangle *$ , a pointer, then a mapping from the category, namely the objects of types, to a mapping from the specified value type to a memory address?

e.g.

$$\langle \rangle * : \text{float} \mapsto \text{float} *$$

$$\text{float} * : \text{din} \mapsto \text{some memory address}$$

and then we pass in mappings, not values, and so we're actually declaring a square *functor*.

What is `threadIdx`? What is it mathematically? Consider that  $\exists 3$  “modules”:

$$\text{threadIdx}.x$$

$$\text{threadIdx}.y$$

$$\text{threadIdx}.z$$

And then the line

```
int idx = threadIdx.x;
```

says that `idx` is an integer, “declares” it to be so, and then assigns `idx` to `threadIdx.x` which surely has to also have the same type, integer. So (perhaps)

$$idx \equiv \text{threadIdx}.x \in \mathbb{Z}$$

is the same thing.

Then suppose  $\text{threadIdx} \subset \text{FinSet}$ , a subcategory of the category of all (possible) finite sets, s.t. `threadIdx` has 3 particular morphisms,  $x, y, z \in \text{MorthreadIdx}$ ,

$$x : \text{threadIdx} \mapsto \text{threadIdx}.x \in \text{Obj}_{\text{FinSet}}$$

$$y : \text{threadIdx} \mapsto \text{threadIdx}.x \in \text{Obj}_{\text{FinSet}}$$

$$z : \text{threadIdx} \mapsto \text{threadIdx}.x \in \text{Obj}_{\text{FinSet}}$$

### Configuring the Kernel Launch Parameters Part 1

$n_{\text{blocks}}, n_{\text{threads}}$  with  $n_{\text{threads}} \geq 1024$  (this maximum constant is GPU dependent). You should pick the  $(n_{\text{blocks}}, n_{\text{threads}})$  that makes sense for your problem, says Prof. Owen.

3.1.4. *Memory layout of blocks and threads.*  $\forall (n_{\text{blocks}}, n_{\text{threads}}) \in \mathbb{Z} \times \{1 \dots 1024\}$ ,  $\{1 \dots n_{\text{block}} \times \{1 \dots n_{\text{threads}}\}$  is now an ordered index (with lexicographical ordering). This is just 1-dimensional (so possibly there's a 1-to-1 mapping to a finite subset of  $\mathbb{Z}$ ).

I propose that “adding another dimension” or the 2-dimension, that Prof. Owen mentions is being able to do the Cartesian product, up to 3 Cartesian products, of the block-thread index.

### Quiz: Configuring the Kernel Launch Parameters 2

Most general syntax:

Configuring the kernel launch

```
kernel<<<grid of blocks , block of threads >>>(...)
```

```
// for example
```

```
square<<<dim3(bx,by,bz) , dim3(tx,ty,tz) , shmem>>>(...)
```

where `dim3(tx,ty,tz)` is the grid of blocks  $bx \cdot by \cdot bz$

`{dim3}(tx,ty,tz)` is the block of threads  $tx \cdot ty \cdot tz$

`shmem` is the shared memory per block in bytes

**Problem Set 1** “Also, the image is represented as an 1D array in the kernel, not a 2D array like I mentioned in the video.”

Here's part of that code for squaring numbers:

```
--global__ void square(float *d_out , float *d_in) {
    int idx = threadIdx.x;
    float f = d_in[idx];
    d_out[idx] = f*f;
}
```

3.1.5. *Grid of blocks, block of threads, thread that's indexed; (mathematical) structure of it all.* Let

$$\text{grid} = \prod_{I=1}^N (\text{block})^{n_I^{\text{block}}}$$

$$\begin{aligned} \text{where } N = 1, 2, 3 \text{ (for CUDA) and by naming convention} \\ I = 1 \equiv x \\ I = 2 \equiv y \\ I = 3 \equiv z \end{aligned}$$

Let's try to make it explicit (as others had difficulty understanding the grid, block, thread model, cf. [colored image to greyscale image using CUDA parallel processing, Cuda gridDim and blockDim](#)) through commutative diagrams and categories (from math):

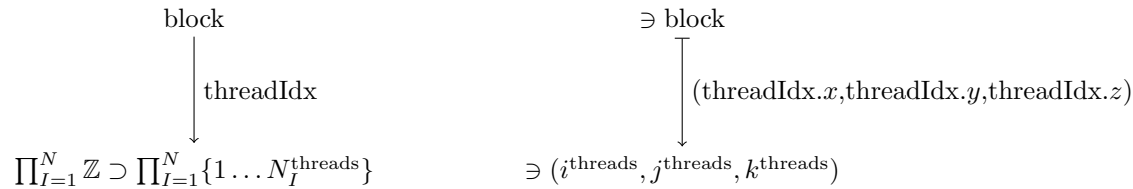
$$\begin{array}{ccc} \prod_{I=1}^N \mathbb{Z}^+ & \ni (N_x^{\text{blocks}}, N_y^{\text{blocks}}, N_z^{\text{blocks}}) & \\ \text{gridDim} \left( \begin{array}{c} \uparrow \\ \text{dim3} \\ \downarrow \end{array} \right) & & \left( \begin{array}{c} \uparrow \\ \text{dim3} \\ \downarrow \end{array} \right) (\text{gridDim}.x, \text{gridDim}.y, \text{gridDim}.z) \\ \text{grid} & & \ni \text{gridSize}(N_x^{\text{blocks}}, N_y^{\text{blocks}}, N_z^{\text{blocks}}) \end{array}$$

$$\begin{array}{ccc} \text{grid} & \ni \text{d\_rgbaImage} & \\ \downarrow \text{blockIdx} & & \downarrow (\text{blockIdx}.x, \text{blockIdx}.y, \text{blockIdx}.z) \\ \prod_{I=1}^N \mathbb{Z} \supset \prod_{I=1}^N \{1 \dots N_I^{\text{blocks}}\} & & \ni (i^{\text{blocks}}, j^{\text{blocks}}, k^{\text{blocks}}) \end{array}$$

and then similar relations (i.e. arrows, i.e. relations) go for a block of threads:

$$\begin{array}{ccc} \prod_{I=1}^N \mathbb{Z}^+ & \ni (N_x^{\text{threads}}, N_y^{\text{threads}}, N_z^{\text{threads}}) & \\ \text{blockDim} \left( \begin{array}{c} \uparrow \\ \text{dim3} \\ \downarrow \end{array} \right) & & \left( \begin{array}{c} \uparrow \\ \text{dim3} \\ \downarrow \end{array} \right) (\text{blockDim}.x, \text{blockDim}.y, \text{blockDim}.z) \\ \text{block} & & \ni \text{blockSize}(N_x^{\text{threads}}, N_y^{\text{threads}}, N_z^{\text{threads}}) \end{array}$$





**gridsize help assignment 1 Pp** explains how threads per block is variable, and remember how Owens said Luebki says that a GPU doesn't get up for more than a 1000 threads per block.

3.1.6. *Generalizing the model of an image.* Consider vector space  $V$ , e.g.  $\dim V = 4$ , vector space  $V$  over field  $\mathbb{K}$ , so  $V = \mathbb{K}^{\dim V}$ . Each pixel represented by  $\forall v \in V$ .

Consider an image, or space,  $M$ .  $\dim M = 2$  (image),  $\dim M = 3$ . Consider a local chart (that happens to be global in our case):

$$\begin{array}{ccc}
 \varphi: M \rightarrow \mathbb{Z}^{\dim M} \supset \{1 \dots N_1\} \times \{1 \dots N_2\} \times \dots \times \{1 \dots N_{\dim M}\} \\
 \varphi: x \mapsto (x^1(x), x^2(x), \dots, x^{\dim M}(x)) \\
 \begin{array}{ccc}
 E & \xrightarrow{\varphi} & M \times V \\
 \pi \downarrow & \swarrow & \\
 M & & 
 \end{array}
 \end{array}$$

$$\begin{array}{ccc}
 E & \xrightarrow{\varphi} & \text{grid} \times \text{block of threads} \\
 \pi \downarrow & \swarrow & \\
 \text{grid} & & 
 \end{array}$$

Consider a “coarsing” of underlying  $M$ :

$$\begin{array}{ccc}
 M \times V & \xrightarrow{\text{proj}} & \text{proj}(M) \times \text{proj}(V) \\
 \pi \downarrow & & \downarrow \text{proj}(\pi) \\
 M = \{1 \dots N_1\} \times \{1 \dots N_2\} \times \dots \times \{1 \dots N_{\dim M}\} & \xrightarrow{\text{proj}} & \text{proj}(M) = \{1 \dots \frac{N_1}{N_1^{\text{threads}}}\} \times \{1 \dots \frac{N_2}{N_2^{\text{threads}}}\} \times \dots \times \{1 \dots \frac{N_{\dim M}}{N_{\dim M}^{\text{threads}}}\}
 \end{array}$$

e.g.  $N_1^{\text{thread}} = 12$

$N_2^{\text{thread}} = 12$

Just note that in terms of syntax, you have the “block” model, in which you allocate blocks along each dimension. So in

*const dim3 blockSize( $n_x^b, n_y^b, n_z^b$ )*

*const dim3 gridSize( $n_x^{\text{gr}}, n_y^{\text{gr}}, n_z^{\text{gr}}$ )*

Then the condition is  $n_x^b/\dim V, n_y^b/\dim V, n_z^b/\dim V \in \mathbb{Z}$  (condition),  $(n_x^{\text{gr}} - 1)/\dim V, n_y^{\text{gr}}/\dim V, n_z^{\text{gr}}/\dim V \in \mathbb{Z}$

**Transpose Part 1**

Now

$$\text{Mat}_{\mathbb{F}}(n, n) \xrightarrow{T} \text{Mat}_{\mathbb{F}}(n, n)$$

$$A \mapsto A^T \text{ s.t. } (A^T)_{ij} = A_{ji}$$

$$\text{Mat}_{\mathbb{F}} \xrightarrow{T} \mathbb{F}^{n^2}$$

$$A_{ij} \mapsto A_{ij} = A_{in+j}$$

$$\begin{array}{ccc}
 \text{Mat}_{\mathbb{F}}(n, n) & \longrightarrow & \mathbb{F}^{n^2} \\
 T \downarrow & & \downarrow T \\
 \text{Mat}_{\mathbb{F}}(n, n) & \longrightarrow & \mathbb{F}^{n^2}
 \end{array}
 \quad
 \begin{array}{ccc}
 A_{ij} & \longmapsto & A_{in+j} \\
 T \downarrow & & \downarrow T \\
 (A^T)_{ij} = A_{ji} & \longmapsto & A_{jn+i}
 \end{array}$$

### Transpose Part 2

Possibly, transpose is a functor.

Consider struct as a category. In this special case,  $\text{Objstruct} = \{\text{arrays}\}$  (a struct of arrays). Now this struct already has a hash table for indexing upon declaration (i.e. “creation”): so this category struct will need to be equipped with a “diagram” from the category of indices  $J$  to struct:  $J \rightarrow \text{struct}$ .

So possibly

$$\begin{array}{ccc}
 \text{struct} & \xrightarrow{T} & \text{array} \\
 \text{ObjStruct} = \{ \text{arrays} \} & \xrightarrow{T} & \text{Objarray} = \{ \text{struct} \} \\
 J \rightarrow \text{struct} & \xrightarrow{T} & J \rightarrow \text{array}
 \end{array}$$

**Quiz: What Kind Of Communication Pattern** This quiz made a few points that clarified the characteristics of these so-called communication patterns (amongst the memory?)

- map is bijective, and  $\text{map} : \text{Idx} \rightarrow \text{Idx}$
- gather - not necessarily surjective
- scatter - not necessarily surjective
- stencil - surjective
- transpose (see before)

### Parallel Communication Patterns Recap

- map - bijective
- transpose - bijective
- gather - not necessarily surjective, and is many-to-one (by def.)
- scatter - one-to-many (by def.) and is not necessarily surjective
- stencil - several-to-one (not injective, by definition), and is surjective
- reduce - all-to-one
- scan/sort - all-to-all

### Programmer View of the GPU

thread blocks: group of threads that cooperate to solve a (sub)problem

### Thread Blocks And GPU Hardware

CUDA GPU is a bunch of SMs:

Streaming Multiprocessors (SM)s

SMs have a bunch of simple processors and memory.

Dr. Luebki:

Let me say that again because it's really important  
GPU is responsible for allocating blocks to SMs

Programmer only gives GPU a pile of blocks.

### Quiz: What Can The Programmer Specify

I myself thought this was a revelation and was not intuitive at first:

Given a single kernel that's launched on many thread blocks include  $X$ ,  $Y$ , the programmer cannot specify the sequence the blocks, e.g. block  $X$ , block  $Y$ , run (same time, or run one after the other), and which SM the block will run on (GPU does all this).

### Quiz: A Thread Block Programming Example

Open up `hello blockIdx.cu` in Lesson 2 Code Snippets (I got the repository from github, repo name is cs344).

At first, I thought you can do a single file compile and run in Eclipse without creating a new project. No. cf. [Eclipse creating projects every time to run a single file?](#).

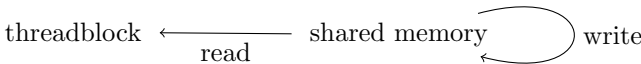
I ended up creating a new CUDA C/C++ project from File -> New project, and then chose project type Executable, Empty Project, making sure to include Toolchain CUDA Toolkit (my version is 7.5), and chose an arbitrary project name (I chose cs344single). Then, as suggested by [Kenny Nguyen](#), I dragged and dropped files into the folder, from my file directory program.

I ran the program with the “Play” triangle button, clicking on the green triangle button, and it ran as expected. I also turned off Build Automatically by deselecting the option (no checkmark).

GPU Memory Model



Then consider threadblock  $\equiv$  thread block  
Objthreadblock  $\supset \{ \text{ threads } \}$   
FinSet  $\xrightarrow{\text{threadIdx}}$  thread  $\in$  Morthreadblock



$\forall$  thread,



Synchronization - Barrier  
Quiz: The Need For Barriers

3 barriers were needed (wasn’t obvious to me at first). All threads need to finish the write, or initialization, so it’ll need a barrier.

While

```
array[idx] = array[idx+1];
```

is 1 line, it’ll actually need 2 barriers; first read. Then write.

So *actually* we’ll need to *rewrite* this code:

```
int temp = array[idx+1];
__syncthreads();
array[idx] = temp;
__syncthreads();
```

kernels have implicit barrier for each.

Writing Efficient Programs

- (1) Maximize *arithmetic intensity* arithmetic intensity  $:= \frac{\text{math}}{\text{memory}}$

video: Minimize Time Spent On Memory

local memory is fastest; global memory is slower

local > shared >> global >> CPU

kernel we know (in the code) is tagged with `__global__`

quiz: A Quiz on Coalescing Memory Access

Work it out as Dr. Luebki did to figure out if it’s coalesced memory access or not.

Atomic Memory Operations

atomicadd atomicmin atomicXOR atomicCAS Compare And Swap

4. POINTERS IN C; POINTERS IN C CATEGORIFIED (INTERPRETED IN CATEGORY THEORY)

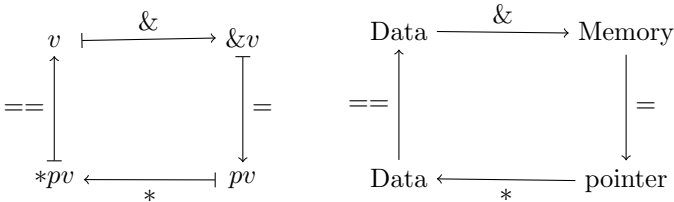
Suppose  $v \in \text{ObjData}$ , category of data **Data**,  
e.g.  $v \in \text{Int} \in \text{ObjType}$ , category of types Type.

$$\begin{aligned} \text{Data} &\xrightarrow{\&} \text{Memory} \\ v &\mapsto \&v \end{aligned}$$

with address  $\&v \in \text{Memory}$ .  
With  
assignment  $pv = \&v$ ,

$pv \in \text{Objpointer}$ , category of pointers, pointer  
 $pv \in \text{Memory}$  (i.e. not  $pv \in \text{Dat}$ , i.e.  $pv \notin \text{Dat}$ )

$$\text{pointer} \ni pv \mapsto^* *pv \in \text{Dat}$$



Examples. Consider `passfunction.c` in Fitzpatrick [5].

Consider the type `double`, `double`  $\in \text{ObjTypes}$ .

`fun1`, `fun2`  $\in \text{MorTypes}$  namely  
`fun1`, `fun2`  $\in \text{Hom}(\text{double}, \text{double}) \equiv \text{Hom}_{\text{Types}}(\text{double}, \text{double})$

Recall that

$$\begin{aligned} \text{pointer} &\xrightarrow{*} \text{Dat} \\ \text{pointer} &\xrightarrow{\&} \text{Memory} \end{aligned}$$

$*$ ,  $\&$  are functors with domain on the category pointer.  
Pointers to functions is the “extension” of functor  $*$  to the codomain of  $\text{MorTypes}$ :

$$\begin{aligned} \text{pointer} &\xrightarrow{*} \text{MorTypes} \\ \text{fun1} &\mapsto^* *fun1 \in \text{Hom}_{\text{Types}}(\text{double}, \text{double}) \end{aligned}$$

Part 3. Machine Learning with Deep Learning

cf. Machine Learning - Introduction, from Coursera. Dr. Andrew Ng.

(1) Week 1

- Linear Regression with One Variable
  - Model and Cost Function
    - \* Model Representation
    - \* Cost Function
    - \* Cost Function - Intuition I
    - \* Cost Function - Intuition II
  - Parameter Learning
    - \* Gradient Descent
    - \* Gradient Descent Intuition
    - \* Gradient Descent For Linear Regression

cf. Linear Regression with One Variable

cf. [Model Representation; Week 1 Linear Regression with 1 Variable, Coursera Machine Learning, Ng](#)

For hypothesis  $h$ ,

$$\begin{aligned} h_\theta &: \mathbb{R}^d \rightarrow \mathbb{R} \\ h_\theta &: x \mapsto h_\theta(x) \quad (\text{prediction of } y \text{ for } x) \end{aligned}$$

$$h_\theta \in L(\mathbb{R}^d, \mathbb{R})$$

$$\begin{aligned} h_\theta &: \mathbb{R}^{|\theta|} \rightarrow L(\mathbb{R}^d, \mathbb{R}) \\ \theta &\mapsto h_\theta \end{aligned}$$

[Cost Function; Week 1, Coursera, Machine Learning, Ng](#)

So for parameters

$$\theta \in \mathbb{R}^{|\theta|}$$

define a *cost function*

$$(1) \quad J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$

Find

$$\min_{\theta} J(\theta) = ?(???)$$

for

$$J : \mathbb{R}^{|\theta|} \rightarrow \mathbb{R}$$

Actually,

$$(2) \quad \begin{aligned} &J(\theta, (x_i, y_i)_{i \in I_{\text{train}}}) \\ &J : \mathbb{R}^{|\theta|} \times (\mathbb{R}^d)^m \times \mathbb{R}^m \rightarrow \mathbb{R} \end{aligned}$$

$m$  = number of training examples =  $|I_{\text{train}}|$ .

Considering

$$H(\theta + \Delta\theta) \approx J(\theta) + \text{grad}J(\theta) \cdot \Delta\theta + \frac{1}{2t} \|\Delta\theta\|^2$$

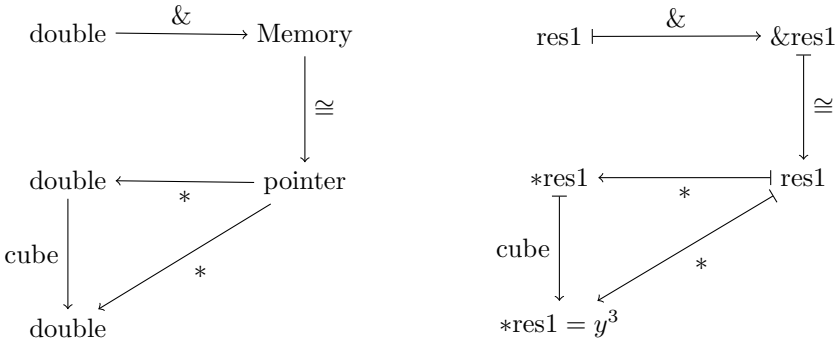
Suppose  $\Delta\theta \equiv \Delta\theta(t) = t\Delta\theta$

$\Delta\theta \approx -\gamma \text{grad}J(\theta)$  is an ansatz,  $\gamma$  small enough.

Then assume  $J$  convex, use this ansatz by plugging in, with Lipshitz condition

$$\|\text{grad}J(\theta + \Delta\theta) - \text{grad}J(\theta)\| \leq L\|\Delta\theta\|$$

some constant  $L > 0$ ,



It’s unclear to me how `void cube` can be represented in terms of category theory, as surely it cannot be represented as a mapping (it acts upon a functor, namely the `*` functor for pointers). It doesn’t return a value, and so one cannot be confident to say there’s explicitly a domain and codomain, or range for that matter.

But what is going on is that

$$\begin{aligned} \text{pointer}, \text{double}, \text{pointer} &\xrightarrow{\text{cube}} \text{pointer}, \text{pointer} \\ \text{fun1}, x, \text{res1} &\xrightarrow{\text{cube}} \text{fun1}, \text{res1} \end{aligned}$$

s.t.  $*\text{res1} = y^3 = (*\text{fun1}(x))^3$

So I’ll speculate that in this case, `cube` is a functor, and in particular, is acting on `*`, the so-called deferencing operator:

$$\begin{aligned} \text{pointer} &\xrightarrow{*} \text{float} \in \text{Data} \xrightarrow{\text{cube}} \text{pointer} \xrightarrow{\text{cube}(*)} \text{float} \in \text{Data} \\ \text{res1} &\xrightarrow{*} *\text{res1} \quad \text{res1} \xrightarrow{\text{cube}(*)} \text{cube}(*\text{res1}) = y^3 \end{aligned}$$

cf. Arrays, from Fitzpatrick [5]

$$\text{Types} \xrightarrow{\text{declaration}} \text{arrays}$$

If  $x \in \text{Objarrays}$ ,

$$\&x[0] \in \text{Memory} \xrightarrow{=} x \in \text{pointer (to 1st element of array)}$$

cf. Section 2.13 Character Strings from Fitzpatrick [5]

```
char word[20] = ‘‘four’’;
char *word = ‘‘four’’;
```

cf. C++ extensions for C according to Fitzpatrick [5]

- simplified syntax to pass by reference pointers into functions
- inline functions
- variable size arrays

```
int n;
double x[n];
```

- complex number class

4.0.7. *Need a CUDA, C, C++, IDE? Try Eclipse!* This website has a clear, lucid, and pedagogical tutorial for using Eclipse: [Creating Your First C++ Program in Eclipse](#). But it looks like I had to pay. Other than the well-written tips on the webpage, I looked up stackexchange for my Eclipse questions (I had difficulty with the Eclipse documentation).



$$\begin{aligned} \theta_{n+1}^i &= \theta_n^i - \gamma_n (\text{grad} J(\theta))^i \\ \gamma_n &= \frac{(\theta_n^i - \theta_{n-1}^i)(\text{grad}_\theta J(x_n) - \text{grad}_\theta J(x_{n-1}))^i}{\|\text{grad}_\theta J(x_n) - \text{grad}_\theta J(x_{n-1})\|^2} = \frac{(\theta_n - \theta_{n-1}) \cdot (\text{grad}_\theta J(x_n) - \text{grad}_\theta J(x_{n-1}))}{\|\text{grad}_\theta J(x_n) - \text{grad}_\theta J(x_{n-1})\|^2} \end{aligned}$$

(3)

or as Ng points out in the [Gradient Descent lesson recap](#), the correct way is to store in temporary variables first:

$$\begin{aligned} \text{temp} &= \theta_n^i - \gamma_n (\text{grad} J(\theta))^i \\ \theta_{n+1}^i &= \text{temp} \end{aligned}$$

(4)

where  $\text{temp} \in \mathbb{R}^{|\theta|}$

In the lesson recap for [Gradient Descent Intuition](#), Ng denotes the learning rate  $\alpha \in \mathbb{R}$  with  $\alpha$ , but note that it’s denoted as  $\gamma$  or **gamma** for **sci-kit learn**. So be aware of different notations. Nevertheless, the learning rate can be a constant, but even then, choosing it is nontrivial.

4.0.8. *Testing many hypotheses at the same time, via refactoring the matrix.* In [Linear Algebra Review of Week 1, Matrix Matrix Multiplication](#), Ng provided a useful tip in refactoring the matrix of hypotheses  $h_\theta$  so to test multiple number of hypotheses at the same time on the same input data,  $X$ .

Mathematically, beginning with

$$h : \mathbb{R}^{|\theta|} \longrightarrow L(\mathbb{R}^d, \mathbb{R})$$

$$\theta \quad \longmapsto \quad h_\theta$$

Consider testing  $H$  different hypotheses,  $\underbrace{\mathbb{R}^{|\theta|} \times \cdots \times \mathbb{R}^{|\theta|}}_H \equiv \otimes_{i=1}^H \mathbb{R}^{|\theta|}$ ,

so treat

$$\otimes_{i=1}^H \mathbb{R}^{|\theta|} = \text{Mat}_{\mathbb{R}}(|\theta|, H)$$

and so

$$h : \otimes_{i=1}^H \mathbb{R}^{|\theta|} = \text{Mat}_{\mathbb{R}}(|\theta|, H) \longrightarrow \otimes_{i=1}^H L(\mathbb{R}^d, \mathbb{R})$$

$$\theta^{(i)} \longmapsto h_{\theta^{(i)}}$$

cf. [Week 4, Non-linear Hypotheses video of Motivations for Coursera’s Machine Learning by Ng](#)

For a sigmoid function  $g$ , consider

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

If  $n$  large (Ng’s notation),  $d = \dim \mathbb{R}^d$ , number of features for training (data) set, for including quadratic features,

$$\begin{aligned} &x_1^2, x_1 x_2, x_1 x_3, x_1 x_4 \dots x_1 x_{100} \\ &x_2^2, x_1 x_3, \dots \end{aligned}$$

$$\approx \mathcal{O}(n^2) \approx \frac{n^2}{2} \qquad (\mathcal{O}(d^2) \approx \frac{d^2}{2})$$

e.g. computer vision,  
e.g.  $50 \times 50$  pixel images,  
 $n = 2500$

pixel intensity  $\in [0, 255]$   
rgb  $\in [0, 255]^3$

$$\begin{aligned} g : \mathbb{R}^{|\theta|} &\rightarrow L(\mathbb{R}^d, \mathbb{R}) \\ \theta &\mapsto g(\theta) \equiv g_\theta \end{aligned}$$

$n \equiv d = 2$ .

Consider

$$\sum_{\substack{a_1, a_2 = 0 \\ i = a_1 + 2a_2}} \theta^{(i)} x_1^{a_1} x_2^{a_2}$$

and so for this example

$$g(\theta)(x_1, x_2) = g\left(\sum_{\substack{a_1, a_2 = 0 \\ i = a_1 + 2a_2}} \theta^{(i)} x_1^{a_1} x_2^{a_2}\right)$$

For computer vision, consider

$$x \in \mathbb{R}^d \text{ with } d = n^x \times n^y$$

and in particular, given pixel intensity or rgb range,

$$\begin{aligned} x &\in [0, 255]^d \\ x &\in [0, 255]^{3d} \end{aligned}$$

cf. [Model Representation I of Week 4, Coursera’s Machine Learning Introduction with Ng](#)

The notes at the end of each video segment **help very much**.

For input

$$\mathbf{x} \in \mathbb{R}^d$$

e.g.  $d = 1, 2, 3$ , or  $4, \dots$

$x_0$  = “bias unit”, input node 0,  $x_0 = 1$  always (Ng).

Sigmoid (logistic) activation function  $\equiv a$ .

$$a_i^{(j)} \equiv \text{“activation” of unit } i \text{ in layer } j$$

$j \in \{2, \dots, N - 1\}$ ,  $j = 1$  is input layer,  $j = N$  is output layer.

$$a_i^{(j)} = g(\Theta_{ik}^{(j-1)} x_k)$$

$$j \xrightarrow{\Theta^{(j)}} j + 1$$

$\Theta^{(j)}$  matrix of weights controlling function mapping from layer  $j$  to layer  $j + 1$ .

$$h_\Theta(x) = a_1^{(N)} = g(\Theta_{1k}^{(N-1)} a_k^{(N-1)})$$

$\forall$  layer  $j$ ,  $\exists$  matrix of weights  $\Theta^{(j)}$ .

If  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j + 1$ ,  $\dim \Theta^{(j)} = s_{j+1} \times (s_j + 1)$

If  $N = 2$ , (1 neuron or only 1 hidden layer)

$$x = (x_i)_{i=1\dots d} \in \mathbb{R}^d, \qquad y \in \mathbb{R}, x_0 = 1$$

$$y = h(\Theta_{1k}^{(1)} x_k^{(1)}) = h(\Theta_{1k}^{(1)} x_k) = h(\Theta^{(1)})(x)$$

e.g.  $h(z) = \frac{1}{1+e^z}$  logistic function.

Neural Network, input layer, output layer, and hidden layers.

$$(5) \quad \Theta_{ik}^{(j)} x_k \mapsto g a_i^{(j+1)} \quad \begin{array}{l} k = 0, 1, \dots s_j \\ i = 1, 2, \dots s_{j+1} \end{array}$$

Note that  $y$  can be  $y \in \mathbb{R}^M$ , not just  $M = 1$ .

### Model Representation II

$z_i^{(j)}$ ,  $i = 1, \dots s_j$ , layer  $j = 1, \dots N$ .

$$(6) \quad g : z_i^{(j)} \mapsto a_i^{(j)}$$

e.g.  $z_i^{(j)} = \Theta_{ik}^{(j-1)} x_k$ ,  $k = 0, 1, \dots d$ .

Set  $x = a^{(1)}$  for input layer.

$$(7) \quad \Theta^{(j-1)} \in \text{Mat}_{\mathbb{R}}((d+1), s_j)$$

$$\Theta^{(j-1)} : a^{(j-1)} \in \mathbb{R}^{d+1} \mapsto z^{(j)} \in \mathbb{R}^{s_j} \xrightarrow{g} a^{(j)} \in \mathbb{R}^{s_j} \xrightarrow{a_0^{(j)}=1} a^{(j)} \in \mathbb{R}^{s_j+1}$$

For the  $j = N$  case, “output” layer,

$$(8) \quad \Theta^{(N-1)} : a^{(N-1)} \mapsto z^N \in \mathbb{R} \xrightarrow{g} g(z^N) = a^N = h_{\Theta}(x) \in \mathbb{R} \quad \Theta^{(N-1)} \in \text{Mat}_{\mathbb{R}}(s_{N-1} + 1, 1)$$

In general,

$$\Theta^{(N-1)} : a^{(N-1)} \mapsto z^N \in \mathbb{R} \xrightarrow{g} g(z^N) = a^N = h_{\Theta}(x) \in \mathbb{R}^M \quad \Theta^{(N-1)} \in \text{Mat}_{\mathbb{R}}(s_{N-1} + 1, M)$$

cf. [Learning With Large Datasets](#), Quiz of Week 10, Gradient Descent with Large Datasets; Learning with Large Datasets.

Suppose you are facing a supervised learning problem and have a very large dataset ( $m = 100,000,000$ ). How can you tell if using all of the data is likely to perform much better than using a small subset of the data (say  $m = 1,000$ )?

Plot a learning curve ( $J_{\text{train}}(\theta)$  and  $J_{CV}(\theta)$ , plotted as a function of  $m$ ) for a range of values of  $m$  and verify that the algorithm has high variance when  $m$  is small.

cf. 1.4 Regularized cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[ -y_k^{(i)} \log((h_{\theta}(x^{(i)}))_k) - (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right] +$$

$$+ \frac{\lambda}{2m} \left[ \sum_{j=1}^{s_2} \sum_{k=1}^d (\Theta_{j,k}^{(1)})^2 + \sum_{j=1}^K \sum_{k=1}^{s_2} (\Theta_{j,k}^{(2)})^2 \right]$$

### 5. FEEDFORWARD; FEEDFORWARD PROPAGATION AND PREDICTION

Given ordered sequence of linear transformations  $L$ ,  $L \geq 2$ ,

$$\Theta^{(l)} \in \text{Mat}_{\mathbb{R}}(s_l + 1, s_{l+1}) \text{ i.e. } s_{l+1} \times (s_l + 1) \text{ matrix size, } \forall l = 1, 2, \dots L - 1$$

$$(9) \quad \Theta^{(l)} : \mathbb{R}^{s_l+1} \rightarrow \mathbb{R}^{s_{l+1}}$$

$$\Theta^{(l)} : a^{(l)} \mapsto z^{(l+1)} = \Theta^{(l)} a^{(l)} = \Theta_{ij}^{(l)} a_j^{(l)} = z_i^{(l+1)}$$

$a^{(l)} \equiv$  “activation” of layer  $l$ .

$s_l \equiv$  “layer size” of layer  $l$ , number of units or nodes in layer  $l$

$$(10) \quad \begin{array}{l} g : \mathbb{R}^{s_l} \rightarrow \mathbb{R}^{s_l} \\ g : z^{(l)} \mapsto g(z^{(l)}) \end{array}$$

e.g.  $g$  sigmoid function.

Remember to add  $a_0^{(l)} = 1$ ,  $\forall l = 1, \dots L - 1$ , i.e.  $\forall$  input layer and hidden layers.

For  $l = 1$ , the so-called *input layer*, is such that

$$(11) \quad (a_0^{(1)} = 1, x) = a^{(1)}$$

For  $l = 1, 2, \dots L - 1$ ,

$$\mathbb{R}^{s_l} \xrightarrow{a_0^{(l)} = 1} \mathbb{R}^{s_l+1} \xrightarrow{\Theta^{(l)}} \mathbb{R}^{s_{l+1}} \xrightarrow{g} \mathbb{R}^{s_{l+1}}$$

$$(12) \quad a^{(l)} \mapsto \xrightarrow{a_0^{(l)} = 1} (a_0^{(l)} = 1, a^{(l)}) \mapsto \xrightarrow{\Theta^{(l)}} z^{(l+1)} \mapsto \xrightarrow{g} g(z^{(l+1)}) = a^{(l+1)}$$

### 6. BACKPROPAGATION; BACKPROPAGATION ALGORITHM

First, do feedforward on *each* training example  $t$ , i.e.

$$\forall t = 1, 2, \dots m$$

$$(13) \quad \mathbb{R}^d \xrightarrow{(g \circ \Theta^{(l)} \circ (a_0^{(l)} = 1) \times \cdot)^{L-1}} \mathbb{R}^K$$

$$x^{(t)} \mapsto \xrightarrow{(g \circ \Theta^{(l)} \circ (a_0^{(l)} = 1) \times \cdot)^{L-1}} a^{(L)}$$

For  $K = 1$  or  $K > 1$  e.g.  $K = 10$  for multi-class logistic regression.

In fact, we obtain an ordered sequence of “activation” vectors:

$$\forall t = 1, 2, \dots m$$

$$(14) \quad \mathbb{R}^d \xrightarrow{(g \circ \Theta^{(l)} \circ (a_0^{(l)} = 1) \times \cdot)^{L-1}} \mathbb{R}^{s_2} \times \mathbb{R}^{s_2} \times \mathbb{R}^{s_3} \times \mathbb{R}^{s_3} \times \dots \times \mathbb{R}^K$$

$$x^{(t)} \mapsto \xrightarrow{(g \circ \Theta^{(l)} \circ (a_0^{(l)} = 1) \times \cdot)^{L-1}} z^{(2)}, a^{(2)}, z^{(3)}, a^{(3)}, \dots, a^{(L)}$$

From [Backpropagation algorithm, Cost Function and Backpropagation, Week 5 of Coursera’s Machine Learning Introduction by Ng](#), [Backpropagation algorithm](#) notes, and [ex4.pdf](#)

Calculate

$$(15) \quad \begin{array}{l} \delta^{(L)} := a^{(L)} - y \\ \delta_k^{(L)} := a_k^{(L)} - y_k \quad \forall k = 1, 2, \dots K \end{array}$$

For the term  $((\Theta^{(L-1)})^T \delta^{(L)})$ , the matrix size dimensions of the  $(\Theta^{(L-1)})^T$  are  $\dim(\Theta^{(L-1)})^T = (s_{L-1} + 1) \times s_L$ .

It seems that the element-wise or component-wise multiplication that seems obvious in Matlab/Octave or numpy is called the *Hadamard product*, denoted  $\circ$  or  $\odot$ . There ought to be a homomorphism that maps this operation onto “vectorized” forms of these vectors that allows for, or is equipped with the operation, Hadamard product.

For  $m = 1$ ,

$$\delta^{(L-1)} := \left( (\Theta^{(L-1)})^T \delta^{(L)} \right) \odot g'(z^{(L-1)}) \in \mathbb{R}^{s_{L-1}+1} \quad \forall k = 0, 1, \dots s_{L-1}$$

i.e.

$$\text{vec}(\delta^{(L-1)}) = \text{vec}((\Theta^{(L-1)})^T \delta^{(L)}) \odot \text{vec}(g'(z^{(L-1)})) \mapsto \delta^{(L-1)} \in \mathbb{R}^{s_{L-1}+1}$$

$$(s^{(L-1)})_K := ((\Theta^{(L-1)})^T \delta^{(L)})_K (g'(z^{(L-1)}))_K$$

Then add this term to the so-called “accumulator matrix”  $\Delta^{(l)}$ :

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

Note that prior, skip or *remove*  $\delta_0^{(l+1)}$  entry:

$$\delta^{(l+1)} \in \mathbb{R}^{s_{l+1}+1} \xrightarrow{r} \delta^{(l+1)} \in \mathbb{R}^{s_{l+1}}$$

The whole purpose is to obtain

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)} = a_j^{(l)} ((\Theta^{(l+1)})^T \delta^{(l+2)} \odot g'(z^{(l+1)}))_i$$

which can be shown.

So first we had set

$$\Delta_{ij}^{(l)} = 0$$

for

$$\Delta^{(l)} \in \text{Mat}_{\mathbb{R}}(s_l, s_{l+1}) \in \mathbb{R}^{s_l} \otimes \mathbb{R}^{s_{l+1}}$$

Again, it can be shown that

$$(16) \quad \Delta_{ij}^{(l)} = \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$$

and so

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \delta_j^{(l+1)} (a^{(l)})_i$$

and so for  $\forall t$ ,

$$\begin{cases} D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} & \text{if } j = 0 \end{cases}$$

$$D^{(l)} = \frac{1}{m} \sum_{t=1}^m (\Delta^{(l)})^{(t)} + \lambda \Theta^{(l)} \in \text{Mat}_{\mathbb{R}}(s_l, s_{l+1})$$

In summary, we have, for the first step,

$$(17) \quad \delta^{(L)} := a^{(L)} - y \in \mathbb{R}^K$$

$$(18) \quad \delta^{(l)} = (\Theta^{(l)})^T \delta^{(l+1)} \odot g'(z^{(l)}) \in \mathbb{R}^{s_l+1}, \quad l = L-1, L-2, \dots, 2, \quad (L-2) \text{ steps}$$

and so for

$$(19) \quad (\Delta^{(l)})^{(t)} := (\delta^{(l+1)} (a^{(l)})^T)^{(t)}$$

$$(20) \quad D^{(l)} = \frac{1}{m} \sum_{t=1}^m (\Delta^{(l)})^{(t)} + \lambda \Theta^{(l)} \in \text{Mat}_{\mathbb{R}}^l(s_l, s_{l+1})$$

with  $D^{(l)} \sim \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ .

And so

$$(\mathbb{R}^{s_2})^2 \times (\mathbb{R}^{s_3})^2 \times \dots \times (\mathbb{R}^K)^2 \longrightarrow \text{Mat}_{\mathbb{R}}(s_1, s_s) \times \text{Mat}_{\mathbb{R}}(s_2, s_3) \times \dots \times \text{Mat}_{\mathbb{R}}(s_{L-1}, s_L)$$

$$(21) \quad z^{(2)}, a^{(2)}, z^{(3)}, a^{(3)}, \dots, z^{(L)}, a^{(L)} \longmapsto (\Delta^{(1)})^{(t)}, (\Delta^{(2)})^{(t)}, \dots, (\Delta^{(L-1)})^{(t)}$$

$\forall t = 1, \dots, m$ , then obtaining

$$(22) \quad D^{(l)} \sim \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) \in \text{Mat}_{\mathbb{R}}(s_l, s_{l+1}) \quad \forall l = 1, 2, \dots, L-1$$

To collect our facts, consider that we’re given  $x \in (\mathbb{R}^d)^m$ , with  $x_i^{(t)}$ ,  $i = 1 \dots d$ , with  $y \in (\mathbb{R}^K)^m$ .  
 $t = 1 \dots m \quad y \in \{1, 2, \dots, K\}^m$  (classifier)

“layer”  $l = 1, 2, \dots, L-1$  For input layer  $\Theta^{(1)} : \mathbb{R}^{d+1} \rightarrow \mathbb{R}^{s_2}$   
 $\Theta^{(1)} : a^{(1)} \mapsto \Theta^{(1)} a^{(1)} = z^{(1)}$ , with  $a^{(1)} = (1, x^{(t)})$ .

Instead of thinking of separate “layers”, one should really think of encapsulating the relation, or arrows, or mappings between “layers”:

$$\mathbb{R}^d \xrightarrow{a_0^{(1)} = 1} \mathbb{R}^{d+1} \xrightarrow{\Theta^{(1)}} \mathbb{R}^{s_2} \xrightarrow{g} \mathbb{R}^{s_2}$$

$$x \xrightarrow{a_0^{(1)} = 1} (a_0^{(l)} = 1, x) \xrightarrow{\Theta^{(1)}} z^{(2)} \xrightarrow{g} g(z^{(2)}) = a^{(2)}$$

$$(23)$$

$$\mathbb{R}^{s_l} \xrightarrow{a_0^{(l)} = 1} \mathbb{R}^{s_{l+1}} \xrightarrow{\Theta^{(l)}} \mathbb{R}^{s_{l+1}} \xrightarrow{g} \mathbb{R}^{s_{l+1}}$$

$$(24) \quad a^{(l)} \xrightarrow{a_0^{(l)} = 1} (a_0^{(l)} = 1, a^{(l)}) \xrightarrow{\Theta^{(l)}} z^{(l+1)} \xrightarrow{g} g(z^{(l+1)}) = a^{(l+1)}$$

I found that Theano wasn’t like the ‘.stack’ method, the “addition” of adding the  $a_0 = 1$  component to a vector or matrix, as a shared variable, very much on the GPU (it indeed is a bug, [Merge fails on GPU but passes on CPU #152](#)), and so I rewrote the mathematical formulation to fit in with separating the intercepts from the “weights” or  $\Theta$ .

For

$$(25) \quad \Theta^{(l)}, b^{(l)} : \mathbb{R}^{s_l} \rightarrow \mathbb{R}^{s_{l+1}}$$

where

$$(26) \quad \Theta^{(l)} \in \text{Mat}_{\mathbb{R}}(s_l, s_{l+1}) = \mathbb{R}^{s_{l+1}} \otimes (\mathbb{R}^{s_l})^*$$

$$b^{(l)} \in \mathbb{R}^{s_{l+1}}$$

$$\mathbb{R}^{s_l} \xrightarrow{\Theta^{(l)}, b^{(l)}} \mathbb{R}^{s_{l+1}} \xrightarrow{g} \mathbb{R}^{s_{l+1}}$$

$$a^{(l)} \xrightarrow{\Theta^{(l)}, b^{(l)}} z^{(l+1)} \xrightarrow{g} g(z^{(l+1)}) = a^{(l+1)}$$

(27)

6.1. **Cost functional.** The cost function  $J$  is really a cost *functional*, to first input in the output values  $y$ . So

$$(28) \quad \begin{aligned} J : (\mathbb{R}^K)^m &\rightarrow L((\Theta, \mathbf{b}), \mathbb{R}) \\ J : y &\mapsto J_y \equiv J \end{aligned}$$

for a “vector-valued” regression, with the usual linear regression being the case of  $K = 1$ .

For  $y$  taking on discrete values,

$$(29) \quad \begin{aligned} J : \{1, 2, \dots, K\}^m &\rightarrow L((\Theta, \mathbf{b}), \mathbb{R}) \\ J : y &\mapsto J_y \equiv J \end{aligned}$$

Then, we can find the cost  $J((\Theta, b))$ , for a particular choice of the parameters,  $(\Theta, b) \in (\Theta, \mathbf{b})$ :

$$(30) \quad \begin{aligned} J_y &\equiv J : (\Theta, \mathbf{b}) \rightarrow \mathbb{R} \\ J & : (\Theta, b) \rightarrow J(\Theta, b) \end{aligned}$$

i.e.  $J \in C^\infty((\Theta, b))$  (hopefully  $J$  is smooth or at least  $C^2$  differentiable, so that a Hessian can be obtained).

For the above  $(\Theta, \mathbf{b})$  was notation or shorthand as follows:

$$(\Theta, \mathbf{b}) \equiv (\text{Mat}_{\mathbb{R}}(s_1, s_2) \times \mathbb{R}^{s_2}) \times (\text{Mat}_{\mathbb{R}}(s_2, s_3) \times \mathbb{R}^{s_3}) \times \cdots \times (\text{Mat}_{\mathbb{R}}(s_{L-1}, s_L) \times \mathbb{R}^{s_L})$$

## 7. UNIVERSAL APPROXIMATION THEOREM

**Wikipedia: Universal Approximation Theorem**

From Hornik (1991) [7], pp. 252, Section 2. Results,

$$(31) \quad \mathcal{N}_k^{(n)}(\psi) = \{h : \mathbb{R}^k \rightarrow \mathbb{R} \mid h(x) = \sum_{j=1}^n \beta_j \psi(a'_j x - \theta_j)\}$$

with  $a = (\alpha_1, \dots, \alpha_k)$  with  $a' \equiv a^T \equiv$  transpose of  $a$ .  
 $x = (\xi_1, \dots, \xi_k)$

For arbitrary number of hidden layers,

$$\mathcal{N}_k(\psi) = \bigcup_{n=1}^{\infty} \mathcal{N}_k^{(n)}(\psi)$$

The 2 very important theorems from Hornik (1991) are the following:

**Theorem 1.** *If  $\psi$  unbounded and nonconstant, then  $\mathcal{N}_k(\psi)$  dense in  $L^p(\mu)$ ,  $\forall$  finite measure  $\mu$  on  $\mathbb{R}^k$*

**Theorem 2.** *If  $\psi$  cont., bounded, nonconstant, then  $\mathcal{N}_k(\psi)$  dense in  $C(X)$ ,  $\forall$  compact subsets  $X$  of  $\mathbb{R}^k$ , i.e.  $\forall f \in C(X)$ ,  $\exists$  sequence,  $(h_n)$  s.t.  $h_n \xrightarrow{n} f$  uniformly i.e.  $\forall$  given  $\epsilon > 0$ ,  $\exists N = N(\epsilon)$  (independent of  $x \in X \subset \mathbb{R}^k$ ), s.t.  $|h_n(x) - f(x)| < \epsilon \quad \forall x \in X, \quad \forall n \geq N(\epsilon)$*

I will write now a dictionary between Hornik’s notation and my notation (take note, Hornik’s notation  $\equiv$  my notation).

$f \in C(X)$ ,  $f : \mathbb{R}^k \rightarrow \mathbb{R}$ ,  $k \equiv d$ , so  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$\psi \equiv g$ , e.g.  $g(z) = \frac{1}{1+\exp(-z)}$  or  $g(z) = \tanh(z)$ , but equip  $g$  with element-wise (component-wise) action, i.e.  $g$  as a functor,

$g : \mathbb{R}^k \rightarrow \mathbb{R}^k$  , i.e.  $g : \mathbf{Vec} \rightarrow \mathbf{Vec}$ .

$g : x_j \mapsto g(x_j)$

Now  $a \equiv \Theta \in \text{Mat}_{\mathbb{R}}(d, n)$ ,

$$g(\Theta x + b) = g(z)$$

i.e.  $z \in \mathbb{R}^n$ ,

$$z := \Theta x + b \text{ i.e. } z_j = \Theta_{jk} x_k + b_j \quad g(z) \in \mathbb{R}^n$$

and so, notation-wise,

$$\sum_{j=1}^n \beta_j \psi(a'_j x - \theta_j) \equiv \sum_{j=1}^n \beta_j g(\Theta_{jk} x_k + b_j)$$

Consider

$$\Theta^{(1)} \in \text{Mat}_{\mathbb{R}}(d, s_2), b^{(1)} \in \mathbb{R}^{s_2}$$

$$\Theta^{(2)} \in \text{Mat}_{\mathbb{R}}(s_2, 1), b^{(2)} \in \mathbb{R}$$

$$h(x) = \Theta^{(2)} g(\Theta^{(1)} x + b^{(1)}) + b^{(2)} \in \mathcal{N}_d^{(s_2)}(g)$$

so the neural net of  $L$  total layers  $d = 1$  “input layer”,  $l = L$  is “output layer” is a tuple  $((\Theta, b), g) \in \mathcal{N}_d^{(L)}(g)$

Hornik, Stinchcombe, and White (1989) [8] deals with multi-(hidden) layer networks on pp. 363, on and after Corollary 2.6.

Given training data,

$$(32) \quad \begin{aligned} (X, y) : \mathbb{R} &\rightarrow (\mathbb{R}^d \times \mathbb{R}^k)^m \\ (X, y)(t) &\mapsto (X(t), y(t)) \end{aligned}$$

discretize time  $t \in \mathbb{R}$ ,

$$(33) \quad \begin{aligned} \mathbb{R} &\xrightarrow{\text{discretize}} \mathbb{Z} \\ [0, T] &\text{ where } T \in \mathbb{R}^+ \rightarrow \{0, 1, \dots, T-1\} \text{ where } T \in \mathbb{Z}^+ \end{aligned}$$

Consider 4 different feedforwards. Note  $y(-1) = 0$ .

## 8. LSTM; LONG SHORT TERM MEMORY

**LSTM (Long Short Term Memory), according to Christian Herta**

Rewriting Herta’s formulation of LSTM, which actually puts in the “cell” memory into some of the input, forget gates, that’s different from a “traditional” LSTM (see Wikipedia),

$$(34) \quad \begin{aligned} \text{input gates } i_t &= \psi_{(i)}(\Theta^{(i)} X_t + b^{(i)} + \theta^{(i)} h_{t-1} + W^{(i)} c_{t-1}) \\ f_t &= \psi_{(f)}(\Theta^{(f)} X_t + b^{(f)} + \theta^{(f)} h_{t-1} + W^{(f)} c_{t-1}) \\ c_t &:= f_t \odot c_{t-1} + i_t \odot g_t \\ \text{output gates } o_t &= \psi_{(o)}(\Theta^{(o)} X_t + b^{(o)} + \Theta^{(o)} g_{t-1} + W^{(o)} c_t) \end{aligned}$$

and then finally, not predict yet (I was mistaken) but  $h$  here denotes some other “hidden” variable,

$$(35) \quad h_t = o_t \odot \psi_h(c_t)$$

$o_t, c_t \in \mathbb{R}^H$ , and so  $W^{(i)}, W^{(f)}, W^{(o)} \in \text{Mat}_{\mathbb{R}}(H, s_2)$ .

$$(36) \quad \begin{aligned} y_t &= \psi_{(y)}(\Theta^{(y)} h_t + b^{(y)}) \\ \Theta^{(y)} : \mathbb{R}^{s_L} &\rightarrow \mathbb{R}^K \end{aligned}$$

$$(37) \quad (\mathbb{R}^d \times \mathbb{R}^H)^m \times (\mathbb{R}^H)^m \xrightarrow{\{\psi_\alpha \circ ((\Theta^{(\alpha)}, b^{(\alpha)}), \theta^{(\alpha)}, W^{(\alpha)}))\}_{\alpha=1, \dots, m}} (\mathbb{R}^H \times \mathbb{R}^H)^m \xrightarrow{\quad} (\mathbb{R}^H \times \mathbb{R}^H)^m \xrightarrow{\quad} (\mathbb{R}^H)^m$$

$$X_t, h_{t-1}, c_{t-1} \mapsto (i_t, f_t, g_t) \mapsto c_t, o_t \mapsto h_t$$

Consider what we’re essentially doing at time step  $t$ :

$$((\mathbb{R}^d \times \mathbb{R}^H) \times (\mathbb{R}^H)^m) \xrightarrow{\{\psi_\alpha \circ ((\Theta^{(\alpha)}, b^{(\alpha)}), \theta^{(\alpha)}, W^{(\alpha)}))\}_{\alpha=1, \dots, m}} (\mathbb{R}^H \times \mathbb{R}^H \times \mathbb{R}^H)^m \xrightarrow{(\cdot, \cdot, (\Theta^{(y)}, b^{(y)}))} (\mathbb{R}^H \times \mathbb{R}^H \times \mathbb{R}^H)^m$$

$$(38) \quad X_t, h_{t-1}, c_{t-1} \mapsto c_t, h_t \mapsto c_t, h_t, y_t$$

The recurrence relation is essentially this:

$$(39) \quad X(t), h(t-1), c(t-1) \mapsto c(t), h(t), y(t) \quad \forall t = 0, 1, \dots, T-1$$

This is the recurrence relation that changes with time  $t$  and in the language of **theano**, for **theano.scan** it is the argument value for argument **sequences**.

Notice how  $h, c$  change over time. These are the sequences we want to take in as input and output.  $X(t)$  is the sequences we want to “iterate over.”  $X(t)$  doesn’t get modified by our operations over time  $t$  (than what is given).  $y(t)$  is an output we desire. So in the language of **theano**, for **theano.scan**,  $X(t)$  goes to the argument value for argument **sequences**, as it’s part of the “list of Theano variables or dictionaries describing the sequences **scan** has to iterate over” and since  $X = X(t)$  for time  $t$ , the “**taps**” is  $[0]$ .  $h, c, y$  is expected to be the return value of the Python function describing a single time step, and “the order of the outputs is the same as the order of **outputs\_info**.”

Look at the parameters:

$$(40) \quad \begin{aligned} (\Theta^{(g)}, b^{(g)}), \theta^{(g)} &\in (\text{Mat}_{\mathbb{R}}(d, H) \times \mathbb{R}^H \times \text{Mat}_{\mathbb{R}}(H, H)) \\ (\Theta^{(\alpha)}, b^{(\alpha)}), \theta^{(\alpha)}, W^{(\alpha)} &\in (\text{Mat}_{\mathbb{R}}(d, H) \times \mathbb{R}^H \times \text{Mat}_{\mathbb{R}}(H, H) \times \text{Mat}_{\mathbb{R}}(H, H)) \\ (\Theta^{(y)}, b^{(y)}) &\in (\text{Mat}_{\mathbb{R}}(H, K)) \end{aligned}$$

These parameters are what you put into, in the language of **theano**, for the argument value of **non\_sequences** of **theano.scan**.

### 8.1. How to choose the number of hidden layers and nodes in a neural net.

## Part 4. Support Vector Machines (SVM)

The clearest and most mathematically rigorous (and satisfying) introductory exposition on support vector machines (SVM) comes out of a Bachelor’s thesis from Nowak (2008) [9]. There is a lot of material that tries to talk about SVM, but the implementation either boils down to showing how to turn the crank on a black-box solution, or is too verbose without saying anything substantial. I’ll include references and links of the material I looked at and didn’t find as helpful as Nowak (2008) [9].

[Lecture12 pdf slides for Ng’s Machine Learning Intro. for coursera](#)

[Support Vector Machine \(and Statistical Learning Theory\) Tutorial by Jason Weston, NEC Labs America](#)

[Wikipedia page for Support Vector Machine](#)

[Support Vector Machines and Generalisation in HEP](#) Not much real generalization going on here other than a recap of literally what’s exactly in Shawe-Taylor and Cristianini (2000) [10].

[https://www.cs.cornell.edu/people/tj/publications/joachims\\_99a.pdf](https://www.cs.cornell.edu/people/tj/publications/joachims_99a.pdf)

## 9. FROM LINEAR CLASSIFIER AS A HYPERPLANE, (BIG) MARGIN, TO LINEAR SUPPORT VECTOR MACHINE (SVM), AND LAGRANGIAN DUAL (I.E. CONJUGATE VARIABLES, CONJUGATE MOMENTA)

Intuitively, we seek to find a boundary line that’ll draw a line that separates the data points into distinct  $K$  (usually  $K = 2$ ) classes to classify the data points. Then, this boundary line will help to predict what class a new data point would fall into, be classified to be. For a linear model, i.e. “linear discriminator”, what we’re trying to do is

find

$$\theta \in \mathbb{R}^d \setminus \{0\}, b \in \mathbb{R}$$

s.t.

$$(41) \quad y^{(i)}(\langle \theta, x^{(i)} \rangle + b) - 1 \geq 0 \quad \forall i = 1, \dots, m$$

where  $\|\theta\|$  is minimal. It is minimal because, since the distance between 2 hyperplanes,

$$\langle \theta, x \rangle - b = \pm 1 \quad (\text{defining equations for hyperplanes})$$

is

$$\frac{2}{\|\theta\|} \quad (\text{distance between 2 hyperplanes})$$

Thus, we want the “margins”, that distance between hyperplanes separating the input data points, to be as big as possible, and so we want  $\|\theta\|$  small.

Consider this cost functional, called “Lagrangian”, that we want to minimize:

$$(42) \quad \mathcal{L}((\theta, b), \lambda) = \frac{1}{2} \|\theta\|^2 - \sum_{j=1}^m \lambda_j y^{(j)} (\langle \theta, x^{(j)} \rangle - b - 1) = \frac{1}{2} \|\theta\|^2 - \sum_{j=1}^m \lambda_j y^{(j)} (\langle \theta, x^{(j)} \rangle - b) + \sum_{i=1}^m \lambda_i$$

Note that

$$(43) \quad f_0((\theta, b)) := \frac{1}{2} \|\theta\|^2 (\text{objective function (slightly modified)})$$

is the objective function, what we want to minimize.

The KKT condition tells us that  $(\theta, b)$  makes  $\mathcal{L}$  a minimum for a certain  $\lambda$ :

$$(44) \quad \begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta_j} &= 0 = \theta_j - \sum_{i=1}^m \lambda_i y^{(i)} x_j^{(i)} \\ \frac{\partial \mathcal{L}}{\partial b} &= 0 = - \sum_{i=1}^m \lambda_i y^{(i)} \end{aligned}$$

Note that this step in taking the partial derivatives of  $\mathcal{L}$  in Eq. 69 is analogous to the construction/computation of dual “conjugate” variables, conjugate momentum, in physics.

Notice then that

$$(45) \quad \begin{aligned} \frac{1}{2} \|\theta\|^2 &= \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle \text{ and} \\ \sum_{i=1}^m \lambda_i y^{(i)} (\langle \theta, x^{(i)} \rangle - b) &= \sum_{i=1}^m \lambda_i y^{(i)} \left( \sum_{j=1}^m \lambda_j y^{(j)} \langle x^{(j)}, x^{(i)} \rangle \right) \\ \implies \mathcal{L}((\theta, b), \lambda) &= -\frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle + \sum_{i=1}^m \lambda_i \end{aligned} \quad (46)$$

## 10. SO-CALLED “KERNEL TRICK”; FEATURE SPACE IS A HILBERT SPACE

The so-called “feature space”  $F$  is a Hilbert space  $H$ ,  $\Phi : \mathbb{R}^d \rightarrow H$ , equipped with inner product

$$(47) \quad \langle \Phi(x), \Phi(y) \rangle = K(x, y)$$

with  $K : \mathbb{K}^d \times \mathbb{K}^d \rightarrow \mathbb{K}^K$  being called the kernel function. Recall that the feature space  $F$  had been introduced to represent the process of preprocessing input data  $X$ . For example, given a single input data example,  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$ , maybe we’d want to consider polynomial features, linear combinations of various orders of monomials  $X_i X_j$  or  $X_i^2 X_j$ , and so on. Then  $\Phi$  represents the map from  $X$  to all these features.

The essence of the kernel trick is this: the explicit form of  $\Phi$  need not be known, nor even the space  $H$ . Only the kernel function  $K$  form needs to be guessed at.

And so even if we now have to modify our Eq. 42 to account for this preprocessing map  $\Phi$ , applied first to our training data  $x^{(i)} \equiv X^{(i)}$  (Novak’s notation vs. Andrew Ng’s notation), we essentially still have the same form, formally.

Keep in mind the whole point of this nonlinear preprocessing map  $\Phi$  - we want to keep the linear discrimination procedure with the weight, or parameter  $\theta$ , and intercept  $b$ , being this linear model on the feature space (Hilbert space)  $F$ . We’re linear in  $F$ . But we’re nonlinear in the input data  $X = \{X^{(1)}, \dots, X^{(m)}\}$ .

So,

$$(48) \quad \begin{aligned} \mathcal{L}((\theta, b), \lambda) &= \frac{1}{2} \|\theta\|^2 - \sum_{j=1}^m \lambda_j y^{(j)} (\langle \theta, \Phi(x^{(j)}) \rangle - b - 1) = \frac{1}{2} \|\theta\|^2 - \sum_{j=1}^m \lambda_j y^{(j)} (\langle \theta, x^{(j)} \rangle - b) + \sum_{i=1}^m \lambda_i \text{ and so} \\ \frac{\partial \mathcal{L}}{\partial \theta_j} &= 0 = \theta_j - \sum_{i=1}^m \lambda_i y^{(i)} \Phi(x)_j^{(i)} \\ \implies \mathcal{L}((\theta, b), \lambda) &= -\frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \langle \Phi(x^{(i)}), \Phi(x^{(j)}) \rangle + \sum_{i=1}^m \lambda_i \end{aligned}$$

### 10.1. Dealing with Errors, (non-negative) slack variables, dealing with not-necessarily perfectly separable data.

First, loosen the strict constraint  $y^{(i)} (\langle \theta, x^{(i)} \rangle - b) \geq 1$  by introducing *non-negative* slack variables  $\xi_i$ ,  $i = 1 \dots m$ ,

$$(49) \quad y^{(i)} (\langle \theta, x^{(i)} \rangle - b) \geq 1 - \xi_i, \quad \forall i = 1, 2, \dots, m$$

Simply add  $\xi$  to the objective function to implement penalty (for “too much slack”):

$$(50) \quad f_0(\theta, b, \xi) = \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \xi_i$$

So then the total Lagrangian becomes

$$(51) \quad \mathcal{L}(\theta, b, \xi, \lambda, \mu) = \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \lambda_i (y^{(i)} (\langle \theta, x^{(i)} \rangle - b) - 1 + \xi_i) - \sum_{i=1}^m \mu_i \xi_i$$

where the constraint is turned into a Lagrange-multiplier type relation:

$$(52) \quad \xi_i \geq 0 \implies \mu_i (\xi_i - 0) \quad \forall i = 1, \dots, m$$

$-\mu_i \xi_i$  is indeed a valid cost (penalty) functional (if  $\xi_i < 0$ ,  $-\mu_i \xi_i > 0$ , and there’s more penalty as  $\xi_i$  gets more negative. Note that I understood this cost or penalty accounting, given an *inequality constraint*, from reading notes from here, .

## 11. DUAL FORMULATION

$$\min. \quad W(\lambda) = - \sum_{i=1}^m \lambda_i + \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} K(x^{(i)}, x^{(j)})$$

(53)

$$\begin{aligned} \text{s.t.} \quad & \sum_{i=1}^m \lambda_i y^{(i)} = 0 \\ & 0 \leq \lambda_i \leq C \end{aligned}$$

At this point, Eq. 82 is what I could consider the “theoretical gold” version. Further modification of this formulation are really to efficiently implement this on the computer (or microprocessor!). But the schemes should respect this “gold” version and compute what this is and say.

**11.1. Implementation.** Wotao Yin’s notes had a terse, but to-the-point, survey/summary of optimization, in particular non-linear optimization with inequality constraints, for his courses 273a and Math 164, Algorithms for constrained optimization. In both course notes, the material is “taken from the textbook Chong-Zak, 4th. Ed.” So we’ll refer to Chong and Zak (2013) [11].

From Ch. 22 “Algorithms for Constrained Optimization”, 2nd. Ed., from pp. 439, Sec. 22.2 Projections, consider  $\Omega \subset \mathbb{R}^d$ ,

$$\Omega = \{\mathbf{x} | l_i \leq x_i \leq u_i, i = 1 \dots d\}$$

Let  $\Pi \equiv$  projection operator. Define the above case as such:

$$\forall \mathbf{x} \in \mathbb{R}^d, y := \Pi[x] \in \mathbb{R}^d$$

$$y_i \equiv \begin{cases} u_i & \text{if } x_i > u_i \\ x_i & \text{if } l_i \leq x_i \leq u_i \\ l_i & \text{if } x_i < l_i \end{cases}$$

### 11.1.1. Projected Gradient descent. .

Implement  $\sum_{i=0}^m \lambda_i y^{(i)} = 0$ , consider the orthogonal projector matrix (operator)

$$(54) \quad \mathbf{P} := \mathbf{1}_{\mathbb{R}^d} - A^T (A A^T)^{-1} A$$

If  $m = 1$ , then

$$\text{Proj}_{\Omega}(\mathbf{y}) = \mathbf{y} - \frac{\mathbf{a}_1^T \mathbf{y} - b}{\|\mathbf{a}_1\|^2} \mathbf{a}_1$$

If  $m > 1$ , then

$$\text{Proj}_{\Omega}(\mathbf{y}) = (\mathbf{1}_{\mathbb{R}^d} - A^T (A A^T)^{-1} A) \mathbf{y} + A^T (A A^T)^{-1} \mathbf{b}$$

For the linear (but it’s an equality) constraint

$$\sum_{i=1}^m \lambda_i y^{(i)} = 0$$

so

$$(55) \quad \mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)} = 0}(\mathbf{y}) = \left( \mathbf{y} - \frac{\sum_{i=1}^m y^{(i)}(\mathbf{y})_i}{\sum_{i=1}^m (y^{(i)})^2} (y^{(i)} \mathbf{e}_i) \right)$$

Narasimhan’s Optimization Tutorial 3, Projected Gradient Descent, Duality had some concrete pseudocode for the projected gradient descent [12].



In summary,

we seek to minimize

$$(56) \quad W(\lambda) = -\sum_{i=1}^m \lambda_i + \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} K(X^{(i)}, X^{(j)}) \quad \forall i = 1, 2, \dots, m$$

by iterating  $t = 0, 1, \dots$ , as such:

$$\lambda'_i(t+1) := \lambda_i(t) - \alpha \text{grad} W(\lambda)$$

$$\lambda''_i(t+1) := \mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)} = 0}(\lambda'_i(t+1))$$

$$\lambda_i(t+1) := \Pi_{0 \leq \lambda_i \leq C}(\lambda''_i(t+1))$$

where

$$(57) \quad \mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)} = 0}(\lambda'_i(t+1)) = \lambda'_i(t+1) - \frac{\sum_{i=1}^m y^{(i)} \lambda'_i(t+1)}{\sum_{i=1}^m (y^{(i)})^2} y^{(i)}$$

$$\Pi_{0 \leq \lambda_i \leq C}(\lambda''_i(t+1)) = \begin{cases} C & \text{if } \lambda''_i(t+1) > C \\ \lambda''_i(t+1) & \text{if } 0 \leq \lambda''_i(t+1) \leq C \\ 0 & \text{if } \lambda''_i(t+1) < 0 \end{cases}$$

11.1.2. *Computing  $b$ , the intercept, with a good algebra tip: multiply both sides by the denominator.* Bishop (2007) [13], on pp. 330 of Ch. 7, Sparse Kernel Machines, gave a very good (it resolved possible numerical instabilities) prescription on how to compute the intercept  $b$ , given  $\lambda$ , which would then give us the function that can make predictions  $\hat{y}$  on input data example  $X^{(i)} \in X$ . It's worth expounding upon here.

For any support vector (Bishop called it a support vector; what I think it's equivalent to is that we've trained on our training set  $(X, y)^{\text{train}}$ , and this is 1 of the training examples)  $X^{(i)}$ ,  $i = 1 \dots m$ ,

$$(58) \quad y^{(i)} f(X^{(i)}) = 1$$

. Then using

$$(59) \quad f(x) := \sum_{i=1}^m y^{(i)} \lambda_i^* K(X^{(i)}, x) + b$$

$$\implies y^{(i)} \left( \sum_{j=1}^m y^{(j)} \lambda_j^* K(X^{(j)}, X^{(i)}) + b \right) = 1$$

Although we can solve this equation for  $b$  with algebra/arithmetic for our arbitrarily chosen support vector, it's numerically more stable to 1st. multiply through by  $y^{(i)}$ , using  $(y^{(i)})^2 = 1$ , and then averaging over all support vectors.

$$(60) \quad \sum_{j=1}^m y^{(j)} \lambda_j^* K(X^{(j)}, X^{(i)}) + b = y^{(i)}$$

$$\implies b = \frac{1}{m} \left( \sum_{i=1}^m y^{(i)} - \sum_{i,j=1}^m y^{(j)} \lambda_j^* K(X^{(j)}, X^{(i)}) \right)$$

11.1.3. *Prediction (with SVM).*

$$(61) \quad \hat{y}(X) = \sum_{i=1}^m y^{(i)} \lambda_i^* K(X^{(i)}, X) + b^*$$

$$\hat{y} : \mathbb{R}^d \rightarrow \{0, 1, \dots, K-1\}$$

Clarke, Fokoue, and Zhang (2009) [14]

## 12. SUPPORT VECTOR MACHINES (SVM) NATIVELY IMPLEMENTED IN THEANO, ENTIRELY ON THE GPU WITH THE CUDA BACKEND, WITH CONSTRAINED GRADIENT DESCENT

**12.1. Executive Summary.** I implemented SVM natively in theano, and can run entirely on the GPU(s) (through the CUDA C/C++ backend). Solving the *constrained optimization* problem to train a SVM here used *parallel reduce* algorithms; in fact a parallel reduce nested in side a parallel reduce. The work complexity achieved in this case should be of  $O(2 \log m)$  where  $m$  is the total number of training examples, as opposed to  $O(m^2)$  for Quadratic Programming (QP), such as Sequential Minimal Optimization (SMO). Training on the same data set for vehicles as previous work for C/C++ library `libsvm` that uses SMO, `SVM_parallel` (what I call the described implementation here) achieves an accuracy of 95.1% on test data, as opposed to 87.8% for `libsvm`. What I'd like to do in the future is to train and test on larger datasets ( $m > 10000$ ) to test `SVM_parallel` for its promising scalability, and to implement it as the outer layer of a deep neural network (DNN) for “Deep SVM”.

**12.2. Motivations and Introductions.** Support Vector Machines (SVM) can be used for (binary) classification in supervised learning on labeled data, being able to learn non-linear, higher-dimensional features and to predict a boundary line or discriminator between classes, through the so-called *kernel trick*, which is really to presume a higher-dimensional Hilbert space to represent feature space  $\mathcal{F}$  (i.e.  $\mathcal{F}$  is a Hilbert space).

I considered the proposition of having, as a final, outer “layer,” of a deep neural network (DNN) to be a support vector machine. Could it outperform the same DNN with a sigmoid or softmax function at the final layer?

To my knowledge, there was not a native implementation of SVM on theano, a Python framework for deep learning/DNNs. Being that I sought to speed up learning/computation on the GPU(s) through the theano CUDA backend, it would seem to defeat the GPU speedup advantages if from the very last layer, a large global memory transfer had to occur from the GPU (with the last DNN layer), to a host CPU, *serial*, implementation of SVM. Global memory transfers are prohibitive expensive, for latency, i.e. time-wise, between host CPU and GPUs.[3]

The outline/plan/highlights for this (short) paper is as follows: in

- **12.3**, I review or summarize points in the theory for SVM, give derivations, etc., in which (this has all been done before; I sought to give a concise review)
  - **12.3.1**, I recap basic, elementary concepts motivating the linear discriminator concept and what hyperplanes are, and how they're defined by a *linear* function
  - **12.3.2** - I continue to review and present derivations for the *Lagrangian*, a Lagrange multiplier problem, we want to minimize, and apply the usual Karush-Kuhn-Tucker (KKT) condition to make progress in deriving the constrained optimization problem we seek to solve,
  - **12.4** kernel trick, with the feature space  $\mathcal{F}$  as a Hilbert space, **12.4.1** slack variables to deal with non-perfectly-separable data, and more derivation that was done before, but made explicitly here
- **12.5**, the *constrained optimization* problem we wish to solve for SVM
- **12.6**, I translate how our constrained optimization problem is to be solved with *projected gradient descent* or “constrained gradient descent” (as the projection operators enforce constraint equalities and inequalities). As noted, this method/algorithm was chosen to utilize the (very useful) `grad` method in the `theano` software package.
  - The Eqns. **86**, **87** at the end of **12.6.1** is the crux of the training method considered in this paper and code for SVM and was directly referred to when implemented in code.
    - I also show how to compute, in a numerically stable manner, the intercept  $b$ , after  $\lambda_i$  Lagrange multipliers are found, and how to compute predictions  $\hat{y}$ , in **12.6.2**, **12.6.3**
- **12.7** the *constrained gradient descent* to solve our constrained optimization problem is implemented and I detail its implementation using software package `theano`, and especially its CUDA backend, so to run solely on the *GPU*. I give the rationale in deploying this constrained gradient descent as opposed to the Sequential Minimal Optimization (SMO) usually used in the predominant SVM software package (in C/C++) `libsvm`. Noteworthy, I also show that **work complexity goes from  $O(m^2)$  to  $O(2 \log m)$** .
- **12.7.1** briefly tells where the code is made available and the 1-to-1 correspondence between the code and mathematical formulation. I also note that *novel use of theano's reduce within reduce*.

- **12.8** has *results* that I try to compare with sample datasets used previously, that I’ve trained as quickly as possible. Look here for the results; better yet, feel free to try the code and jupyter notebook and share benchmarks.<sup>2</sup>

### 12.3. Concise Mathematical Review/Summary of the theory for SVM.

12.3.1. *Hyperplanes and distances to motivate the linear discriminator concept; Support Vector Machine name.* I’ll recap basic, elementary concepts, from Clarke, Fokoue, and Zhang (2009) [14], that motivate the concept of a linear discriminator classifying input data  $X$ .

Consider  $\theta \in \mathbb{R}^d$ , and a linear function  $y$ ,

$$(62) \quad \begin{aligned} y &: \mathbb{R}^d \rightarrow \mathbb{R} \\ y(x) &:= \langle \theta, x \rangle + b \end{aligned}$$

Consider a “level set” at real number value  $c \in \mathbb{R}$ ,  $H_c(\theta, b)$ :

$$(63) \quad H_c(\theta, b) := \{x | y(x) = \langle \theta, x \rangle + b = c\}$$

where  $\dim H_c(\theta, b) = d - 1$  is a *hyperplane*.

$\theta \in \mathbb{R}^d$  is the normal vector to this hyperplane, since,

$$\begin{aligned} \forall x^{(i)}, x^{(j)} \in H_c(\theta, b), \text{ then} \\ \langle \theta, x^{(i)} \rangle + b = c = \langle \theta, x^{(j)} \rangle + b \implies \langle \theta, x^{(i)} - x^{(j)} \rangle = 0 \end{aligned}$$

and since  $x^{(i)} - x^{(j)} \in TH_c(\theta, b)$ , i.e.  $x^{(i)} - x^{(j)}$  belongs in the tangent space to  $H_c(\theta, b)$ ,  $TH_c(\theta, b)$ , then  $\theta$ , in general, is normal to the hyperplane ( $\langle \theta, x^{(i)} - x^{(j)} \rangle = 0$ ).

Given  $z \in \mathbb{R}^d$ , what is the distance from  $z$  to this hyperplane  $H_c(\theta, b)$ ,  $d(z, H_c(\theta, b))$ ? Consider  $z^* = z + t\theta \in H_c(\theta, b)$ . Then

$$\langle \theta, z^* \rangle + b = c = \langle \theta, z \rangle + t\langle \theta, \theta \rangle + b = c \implies t = \frac{c - b - \langle \theta, z \rangle}{\|\theta\|^2}$$

$$\text{and so } d(z, H_c(\theta, b)) = \|t\theta\| = \frac{|\langle \theta, z \rangle + b - c|}{\|\theta\|}$$

Thus, for the perpendicular distance between 2 parallel hyperplanes,  $H_c(\theta, b)$ ,  $H_{c'}(\theta, b)$ , can be found: choose a pt. from  $H_c(\theta, b)$ , without loss of generality, s.t.  $z = \left(\frac{c-b}{\theta_1}, 0, \dots, 0\right)$ , so that

$$\langle \theta, z \rangle + b = c \implies \theta_1 z^1 = c - b$$

Then

$$(64) \quad d(H_c(\theta, b), H_{c'}(\theta, b)) = \frac{|\langle \theta, z \rangle + b - c'|}{\|\theta\|} = \frac{|c - b + b - c'|}{\|\theta\|} = \frac{|c - c'|}{\|\theta\|}$$

Given an input (data) domain  $\mathcal{X} \subseteq \mathbb{R}^d$ , for the case of binary classification, with total number of classes  $K = 2$ , we can consider representing the outcomes  $y$  for each input data example,  $X \in \mathbb{R}^d$ , in 2 ways:

$$(65) \quad y \in \{-1, 1\} \text{ or } y \in \{0, 1\} \text{ for } y \in \{0, 1, \dots, K-1\} (K=2)$$

What ends up happening is that the distance between 2 hyperplanes,  $c = -1, c' = 1$  vs.  $c = 0, c' = 1$ , respectively, changes, as  $d(H_c(\theta, b), H_{c'}(\theta, b)) = \frac{|c-c'|}{\|\theta\|}$ , but its absolute value doesn’t matter. What matters is the form of  $y : \mathbb{R}^d \rightarrow \mathbb{R}$ , of Eq. 62 which defines the hyperplane in Eq. 63, notably in  $\theta, b$ . The lesson is to *be consistent with what the value of  $y$  is to define what class  $X$  belongs to*. For instance, Bishop (2007) [13] and Clarke, Fokoue, and Zhang (2009) [14] chooses to consider  $y \in \{-1, 1\}$ ,  $\forall X$  and I’ll do the same here.

The name “support vectors” seems to come from this intuitive notion:  $\theta \in \mathbb{R}^d, b$  are determined from  $m$  input data examples  $X^{(i)} \in \mathbb{R}^d$ ,  $\forall i = 1, 2, \dots, d$ , and  $\forall i$ , the corresponding class label  $y \in \mathbb{Z}$ .  $\forall X^{(i)} \in \mathbb{R}^d$ , imagine attaching normal vectors of the form  $t\theta$ ,  $t \in \mathbb{R}$  that extend out to the respective hyperplane, determined by  $y^{(i)}$ . These imagined vectors “support” the respective hyperplane.

An important takeaway is that the equation defining the hyperplane  $H_c(\theta, b)$  in Eq. 63 is *linear*.

12.3.2. *Margins, cost functional or “Lagrangian”, dual formulation.* With output, outcome  $y \in \{-1, 1\}$ ,  $\forall X$  input data example, the distance between the 2 hyperplanes, which are level sets of  $c = -1, c' = 1$ , is

$$\frac{2}{\|\theta\|}$$

The method of SVM seeks to maximize this distance, also known as “margin”, to make margins as big as possible.

Clearly, this is equivalent to minimizing  $\frac{1}{2}\|\theta\|^2$ , with  $\frac{1}{2}$  multiplication factor chosen, without loss of generality, to make taking derivatives of  $\theta$  easier.

But we also have the following constraints. We want to have a “margin” of  $\frac{2}{\|\theta\|}$  between the hyperplanes that’ll separate the input data examples  $X^{(i)}$ ,  $\forall i = 1, 2, \dots, m$ , for different classes, in this binary classification class, of those with  $y^{(i)} \in \{-1, 1\}$ , and so those  $X^{(i)}$ ’s will “fall far away” from this “margin” and remain within its corresponding hyperplane  $H_{c'=1}(\theta, b)$  or  $H_{c=1}(\theta, b)$ , thus defining these inequalities:

$$(66) \quad y^{(i)}(\langle \theta, x^{(i)} \rangle + b) - 1 \geq 0 \quad \forall i = 1, \dots, m$$

So we want to find

$$\theta \in \mathbb{R}^d \setminus \{0\}, b \in \mathbb{R}$$

s.t.

$$y^{(i)}(\langle \theta, x^{(i)} \rangle + b) - 1 \geq 0 \quad \forall i = 1, \dots, m$$

where  $\frac{2}{\|\theta\|}$  is maximized, or equivalently, defining the so-called *objective function*  $f_\theta(\theta, b)$ , minimize  $f_\theta(\theta, b)$ :

$$(67) \quad f_\theta(\theta, b) := \frac{1}{2}\|\theta\|^2 \quad (\text{objective function})$$

Consider then this cost functional, also known as the “Lagrangian”, which we want to *minimize*.

$$(68) \quad \mathcal{L}((\theta, b), \lambda) = \frac{1}{2}\|\theta\|^2 - \sum_{j=1}^m \lambda_i y^{(i)}(\langle \theta, x^{(i)} \rangle + b - 1) = \frac{1}{2}\|\theta\|^2 - \sum_{j=1}^m \lambda_i y^{(i)}(\langle \theta, x^{(i)} \rangle + b) + \sum_{i=1}^m \lambda_i$$

Of note, we introduced *Lagrangian multipliers*  $\lambda_i$ ,  $\forall i \in 1, 2, \dots, m$ ,  $\lambda_i \in \mathbb{R}$ , to account for each of the constraints given in Eq. 66.

The Karush-Kuhn-Tucker (KKT) condition tells us that  $(\theta, b)$  makes  $\mathcal{L}$  a minimum for a certain  $\lambda$  (and that these  $\lambda_i$ ’s exist), and that these relations hold:[9], [11]:

$$(69) \quad \begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta_j} &= 0 = \theta_j - \sum_{i=1}^m \lambda_i y^{(i)} x_j^{(i)} \quad j = 1, 2, \dots, d \\ \frac{\partial \mathcal{L}}{\partial b} &= 0 = - \sum_{i=1}^m \lambda_i y^{(i)} \end{aligned}$$

and

$$(70) \quad \lambda_i \geq 0 \quad \forall i = 1, 2, \dots, m$$

,

$$(71) \quad \sum_{j=1}^m \lambda_i y^{(i)}(\langle \theta, x^{(i)} \rangle + b - 1) = 0$$

$\forall i = 1, 2, \dots, m$ , we want input data example  $X^{(i)}$  to be “far away” from the boundary line, or, i.e. to give enough “margin” from the other class’s hyperplane, and so in general,  $(\langle \theta, x^{(i)} \rangle + b - 1)$  will be non-zero in Eq. 71. So this condition is equivalently

$$(72) \quad \sum_{j=1}^m \lambda_i y^{(i)} = 0$$

<sup>2</sup>[github:ernestyalumni/MLgrabbag/ML](https://github.com/ernestyalumni/MLgrabbag/ML), [github:ernestyalumni/MLgrabbag SVM](https://github.com/ernestyalumni/MLgrabbag/SVM) [theano.ipynb](https://github.com/ernestyalumni/MLgrabbag/SVM/blob/master/theano.ipynb)

It’s interesting to see that the step in taking the partial derivatives of  $\mathcal{L}$  in Eq. 69 is analogous to the construction/computation of dual “conjugate” variables, conjugate momentum, in physics.

Notice then that

$$(73) \quad \frac{1}{2}\|\theta\|^2 = \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle \text{ and}$$

$$\sum_{i=1}^m \lambda_i y^{(i)} (\langle \theta, x^{(i)} \rangle + b) = \sum_{i=1}^m \lambda_i y^{(i)} \left( \sum_{j=1}^m \lambda_j y^{(j)} \langle x^{(j)}, x^{(i)} \rangle \right)$$

$$(74) \quad \implies \mathcal{L}((\theta, b), \lambda) = -\frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle + \sum_{i=1}^m \lambda_i$$

**12.4. So-called “Kernel trick”; feature space is a Hilbert space.** The so-called “feature space”  $\mathcal{F}$  is a Hilbert space  $\mathcal{H}$ ,  $\Phi : \mathbb{R}^d \rightarrow \mathcal{H}$ , equipped with inner product

$$(75) \quad \langle \Phi(x), \Phi(y) \rangle = K(x, y)$$

with  $K : \mathbb{K}^d \times \mathbb{K}^d \rightarrow \mathbb{K}^K$  being called the kernel function. Recall that the feature space  $\mathcal{F}$  had been introduced to represent the process of preprocessing input data  $X$ . For example, given a single input data example,  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$ , maybe we’d want to consider polynomial features, linear combinations of various orders of monomials  $X_i X_j$  or  $X_i^2 X_j$ , and so on. Then  $\Phi$  represents the map from  $X$  to all these features.

As both a pedantic remark and academic question, I had denoted  $\mathbb{K}$  to be, in general, a *field* - (very familiar) examples of fields are  $\mathbb{K} = \mathbb{R}, \mathbb{C}, \mathbb{Z}$ , the real numbers, complex numbers, integers, respectively. Many times, input data that we receive could take on discrete values, meaning  $X^{(i)} \in \mathbb{Z}$ . Would there be issues, in proving existence, continuity, and differentiability, throughout derivations for these SVM algorithms, if the underlying field  $\mathbb{K}$  is not the real number line  $\mathbb{R}$ ?

Nevertheless, the essence of the kernel trick is this: the explicit form of  $\Phi$  need *not be known*, nor even the space  $\mathcal{H}$ . Only the kernel function  $K$  form needs to be guessed at.

And so even if we now have to modify our Eq. 68 to account for this preprocessing map  $\Phi$ , applied first to our training data  $X^{(i)}$ , we essentially still have the same form, formally.

Keep in mind the whole point of this nonlinear preprocessing map  $\Phi$  - we want to keep the linear discrimination procedure with the weight, or parameter  $\theta$ , and intercept  $b$ , being this linear model on the feature space (Hilbert space)  $F$ . We’re *linear* in  $\mathcal{F}$ . But we’re *nonlinear* in the input data  $X = \{X^{(1)}, \dots, X^{(m)}\}$ ’s domain.

So,

$$\mathcal{L}((\theta, b), \lambda) = \frac{1}{2}\|\theta\|^2 - \sum_{j=1}^m \lambda_j y^{(j)} (\langle \theta, \Phi(x^{(j)}) \rangle + b - 1) = \frac{1}{2}\|\theta\|^2 - \sum_{j=1}^m \lambda_j y^{(j)} (\langle \theta, x^{(j)} \rangle + b) + \sum_{j=1}^m \lambda_j \text{ and so}$$

$$(76) \quad \frac{\partial \mathcal{L}}{\partial \theta_j} = 0 = \theta_j - \sum_{i=1}^m \lambda_i y^{(i)} \Phi(x)_j^{(i)} \\ \implies \mathcal{L}((\theta, b), \lambda) = -\frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} \langle \Phi(x^{(i)}), \Phi(x^{(j)}) \rangle + \sum_{i=1}^m \lambda_i = \mathcal{L}(X, y, \lambda)$$

Note that we’ll now want to *maximize* this dual formulation  $\mathcal{L}(X, y, \lambda)$ .

**12.4.1. Dealing with Errors, (non-negative) slack variables, dealing with not-necessarily perfectly separable data.** First, “loosen the strict constraint”  $y^{(i)} (\langle \theta, x^{(i)} \rangle + b) \geq 1$  by introducing *non-negative* slack variables  $\xi_i$ ,  $i = 1 \dots m$ ,

$$(77) \quad y^{(i)} (\langle \theta, x^{(i)} \rangle + b) \geq 1 - \xi_i, \quad \forall i = 1, 2, \dots, m$$

Simply add  $\xi$  to the objective function to implement penalty (for “too much slack”), with a “regularization” constant  $C$  (in analogy to regularization in the linear regression or logistic regression classifier methods):

$$(78) \quad f_0(\theta, b, \xi) = \frac{1}{2}\|\theta\|^2 + C \sum_{i=1}^m \xi_i$$

So then the total Lagrangian becomes

$$(79) \quad \mathcal{L}(\theta, b, \xi, \lambda, \mu) = \frac{1}{2}\|\theta\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \lambda_i (y^{(i)} (\langle \theta, x^{(i)} \rangle + b) - 1 + \xi_i) - \sum_{i=1}^m \mu_i \xi_i$$

where the constraint is turned into a Lagrange-multiplier type relation:

$$(80) \quad \xi_i \geq 0 \implies \mu_i (\xi_i - 0) \quad \forall i = 1, \dots, m$$

$-\mu_i \xi_i$  is indeed a valid cost (penalty) functional (if  $\xi_i < 0$ ,  $-\mu_i \xi_i > 0$ , and there’s more penalty as  $\xi_i$  gets more negative. I understood this cost or penalty accounting, given an *inequality constraint*, from reading notes from here, <http://www.pitt.edu/~jrclass/opt/notes4.pdf>).

If we “turn the crank” and take partial derivatives of  $\mathcal{L}$ , with respect to  $\xi_i$ , finding its “conjugate momentum dual”, we’ll actually see that  $\mathcal{L}$  has *no* dependence on  $\xi_i$ :

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \lambda_i - \mu_i = 0$$

$$(81) \quad \begin{array}{l} \text{since} \quad C - \lambda_i = \mu_i \quad \implies C \geq \lambda_i \\ \mu_i \geq 0 \text{ is given} \end{array}$$

$$\mathcal{L}(\theta, b, \xi, \lambda, \mu) = \frac{1}{2}\|\theta\|^2 - \sum_{i=1}^m \lambda_i (y^{(i)} (\langle \theta, \Phi(x^{(i)}) \rangle)) + \sum_{i=1}^m \lambda_i = \mathcal{L}((\theta, b), \lambda)$$

$\xi, \mu$  no longer appear in the dual Lagrangian,  $\mathcal{L}(X, y, \lambda)$ , which we want to *maximize*, nor in the so-called “primal” Lagrangian,  $\mathcal{L}((\theta, b), \lambda)$ .

**12.5. Dual Formulation.** Denoting  $W(\lambda) := -\mathcal{L}(X, y, \lambda)$ ,

$$(82) \quad \boxed{\begin{array}{ll} \text{minimize.} & W(\lambda) = -\sum_{i=1}^m \lambda_i + \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} K(x^{(i)}, x^{(j)}) \\ & \text{s.t.} \quad \sum_{i=1}^m \lambda_i y^{(i)} = 0 \\ & 0 \leq \lambda_i \leq C \quad \forall i = 1, 2, \dots, m \end{array}}$$

At this point, Eq. 82 is what I could consider the “theoretical gold” version. Further modification of this formulation are really to efficiently implement this on the computer (or microprocessor!). But the schemes should respect this “gold” version and compute what this is and say.

**12.6. Constrained Optimization.** Wotao Yin’s notes had a terse, but to-the-point, survey/summary of optimization, in particular nonlinear optimization with inequality constraints, for his courses [273a](#) and [Math 164, Algorithms for constrained optimization](#). In both course notes, the material is “taken from the textbook Chong-Zak, 4th. Ed.” So we’ll refer to Chong and Zak (2013) [\[11\]](#).

From Ch. 22 “Algorithms for Constrained Optimization”, 2nd. Ed., pp. 439, Sec. 22.2 “Projections”, consider  $\Omega \subset \mathbb{R}^d$ , with

$$\Omega = \{\mathbf{x} | l_i \leq x_i \leq u_i, i = 1 \dots d\} \quad (86)$$

Let us denote  $\Pi \equiv$  projection operator. Let us mathematically formulate how projection operator  $\Pi$  maps a point  $\mathbf{x} \in \mathbb{R}^d$  onto the subset  $\Omega \subset \mathbb{R}^d$  defined above:

$$(83) \quad \begin{aligned} &\forall \mathbf{x} \in \mathbb{R}^d, y := \Pi[x] \in \mathbb{R}^d \\ &y_i \equiv \begin{cases} u_i & \text{if } x_i > u_i \\ x_i & \text{if } l_i \leq x_i \leq u_i \\ l_i & \text{if } x_i < l_i \end{cases} \end{aligned}$$

**12.6.1. Projected Gradient descent.** We want to minimize  $W(\lambda)$  in Eq. [82](#). The software package **theano** provides the graph-generating method **grad**, which automatically computes the symbolic gradient of a scalar-valued function of symbolic (theano) variables. This **grad** has been very useful for automating the computation of the so-called “back-propagation” step of machine learning/deep learning.

We would like to reuse this useful theano method for SVM. Therefore I sought out a solution to our constrained optimization problem that’ll involve computing gradients at each iteration, but subject to our constraint equality and inequalities.

We already know how to deal with constraint *inequalities* via the projection operator in Eq. [83](#). And note that this can be simply implemented in Python/theano with a **if/else** statement(s) and **theano.tensor.switch**, respectively.

To implement the constraint *equality*,  $\sum_{i=0}^m \lambda_i y^{(i)} = 0$ , consider the orthogonal projector matrix (operator)

$$(84) \quad \mathbf{P} := \mathbf{1}_{\mathbb{R}^d} - A^T(AA^T)^{-1}A$$

with  $A$  being a transformation from  $\mathbb{R}^d$  to  $\mathbb{R}^m$ , i.e.  $A : \mathbb{R}^d \rightarrow \mathbb{R}^m$ , and where  $Ax = b$  is the constraint equality (written in its most general form) [\[11\]](#).

So for where  $\Omega = \{X | AX = b\}$ , if  $m = 1$ , then

$$\text{Proj}_{\Omega}(\mathbf{y}) = \mathbf{y} - \frac{\mathbf{a}_1^T \mathbf{y} - b}{\|\mathbf{a}_1\|^2} \mathbf{a}_1$$

If  $m > 1$ , then

$$\text{Proj}_{\Omega}(\mathbf{y}) = (\mathbf{1}_{\mathbb{R}^d} - A^T(AA^T)^{-1}A)\mathbf{y} + A^T(AA^T)^{-1}\mathbf{b}$$

For the linear (equality) constraint

$$\sum_{i=1}^m \lambda_i y^{(i)} = 0$$

we have

$$(85) \quad \mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)}=0}(\mathbf{y}) = \left( \mathbf{y} - \frac{\sum_{i=1}^m y^{(i)}(\mathbf{y})_i}{\sum_{i=1}^m (y^{(i)})^2} (y^{(i)} \mathbf{e}_i) \right)$$

In summary,

we seek to minimize

$$W(\lambda) = - \sum_{i=1}^m \lambda_i + \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} K(X^{(i)}, X^{(j)})$$

by iterating  $t = 0, 1, \dots$ , as such:

$$\lambda'_i(t+1) := \lambda_i(t) - \alpha \text{grad} W(\lambda)$$

$$\lambda''_i(t+1) := \mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)}=0}(\lambda'_i(t+1))$$

$$\lambda_i(t+1) := \Pi_{0 \leq \lambda_i \leq C}(\lambda''_i(t+1))$$

where

$$(87)$$

$$\mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)}=0}(\lambda'_i(t+1)) = \lambda'_i(t+1) - \frac{\sum_{i=1}^m y^{(i)} \lambda'_i(t+1)}{\sum_{i=1}^m (y^{(i)})^2} y^{(i)}$$

$$\Pi_{0 \leq \lambda_i \leq C}(\lambda''_i(t+1)) = \begin{cases} C & \text{if } \lambda''_i(t+1) > C \\ \lambda''_i(t+1) & \text{if } 0 \leq \lambda''_i(t+1) \leq C \\ 0 & \text{if } \lambda''_i(t+1) < 0 \end{cases}$$

The  $\alpha$  parameter is the analogue to the *learning rate* of gradient descent and will need to be tuned.

**12.6.2. Computing  $b$ , the intercept, with a good algebra tip: multiply both sides by the denominator.** Bishop (2007) [\[13\]](#), on pp. 330 of Ch. 7, Sparse Kernel Machines, gave a very good (it resolved possible numerical instabilities) prescription on how to compute the intercept  $b$ , given  $\lambda$ , which would then give us the function that can make predictions  $\hat{y}$  on input data example  $X^{(i)} \in X$ . It’s worth expounding upon here.

For any support vector (Bishop called it a support vector; what I think it’s equivalent to is that we’ve trained on our training set  $(X, y)^{\text{train}}$ , and this is 1 of the training examples)  $X^{(i)}$ ,  $i = 1 \dots m$ ,

$$(88) \quad y^{(i)} f(X^{(i)}) = 1$$

. Then using

$$(89) \quad \begin{aligned} f(x) &:= \sum_{i=1}^m y^{(i)} \lambda_i^* K(X^{(i)}, x) + b \\ \implies y^{(i)} \left( \sum_{j=1}^m y^{(j)} \lambda_j^* K(X^{(j)}, X^{(i)}) + b \right) &= 1 \end{aligned}$$

Although we can solve this equation for  $b$  with algebra/arithmetic for our arbitrarily chosen support vector, it’s numerically more stable to 1st. multiply through by  $y^{(i)}$ , using  $(y^{(i)})^2 = 1$ , and then averaging over all support vectors.

$$(90) \quad \begin{aligned} \sum_{j=1}^m y^{(j)} \lambda_j^* K(X^{(j)}, X^{(i)}) + b &= y^{(i)} \\ \implies b &= \frac{1}{m} \left( \sum_{i=1}^m y^{(i)} - \sum_{i,j=1}^m y^{(j)} \lambda_j^* K(X^{(j)}, X^{(i)}) \right) \end{aligned}$$



12.6.3. *Prediction (with SVM)*. Compute predictions with this formula: [14]

$$(91) \quad \begin{aligned} \hat{y}(X) &= \sum_{i=1}^m y^{(i)} \lambda_i^* K(X^{(i)}, X) + b^* \\ \hat{y} : \mathbb{R}^d &\rightarrow \{0, 1, \dots, K-1\} \quad (\text{with } K=2 \text{ for binary classification}) \end{aligned}$$

12.7. **Constrained Gradient Descent (Implementation)**. From Eqns. 86, 87, with the algorithm or iterative, computational steps that we should take mathematically formulated (clearly), I had sought out to implement these steps using **theano** and on the GPU, in the hopes of speeding up computation and developing a method that can scale with  $m$  input data examples.

Take a look at this double summation term in Eqn. 86 for  $W(\lambda)$ :

$$(92) \quad f_1(\lambda) := \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} K(X^{(i)}, X^{(j)}) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i y^{(i)} K(X^{(i)}, X^{(j)}) \lambda_j y^{(j)} = (\mathbf{q}^{(i)})^T K(X^{(i)}, X^{(j)}) \mathbf{q}^{(j)}$$

Quadratic programming (denoted “QP” in computer science literature) is essentially trying to put the calculation of double sums, such as the above, into quadratic form, as in the very last equality. Techniques for efficient calculation, on the CPU, of this quadratic form, after reformulation of the original problem, make up QP.

The prevailing software package used for SVM, written in C/C++, that also underlies SVM module for sci-kit learn (**sklearn**) [16] is **libsvm** [17]. **libsvm** employs the method of Sequential Minimal Optimization (SMO) [18]. The main advantage of SMO is that only 2  $\lambda_i$ ’s, Lagrange multipliers, are considered in the working set at each stage in time and the optimal solution is computed analytically at this point.

For instance, suppose we are considering 2 Lagrange multipliers  $\lambda_1$  and  $\lambda_2$ . We first compute the optimal value changing  $\lambda_2$  only. Then, using the inequality constraints  $0 \leq \lambda_1, \lambda_2 \leq C$ , and equality constraint  $\sum_{i=1}^m \lambda_i y^{(i)} = 0$  (but adapted to the fact that we’re only changing 2  $\lambda_i$ ’s), we can analytically compute the other  $\lambda_1$ .

The (serial) computation in C/C++ of this analytical problem at this single step for SMO is fast. However, for the fitting for large data sets (large  $m$ ), the fit time complexity is more than quadratic with the number of examples  $m$ , which makes it difficult to scale to datasets of more than a couple of 10000 examples ( $m > 10000$ )<sup>3</sup> [16].

Instead, I considered the idea behind *All-pairs N-body* algorithm of Nyland, Harris, and Prins (2007) in Ch. 31 of **GPU Gems 3** [15]. It was also explained in Udacity’s CS344 with Owens and Luebke [3]<sup>4</sup>.

Look again at Eq. 92,  $f_1$ , which clearly requires  $m^2$  fetches, or reads, for  $\lambda_i \lambda_j y^{(i)} y^{(j)} K(X^{(i)}, X^{(j)})$  term and  $m^2$  computations, for each  $(i, j)$  pairs (and there are  $m^2$  total pairs). It could also help to imagine a  $m \times m$  matrix:

$$\begin{array}{cccc} i=1, j=1 & i=1, j=2 & \dots & i=1, j=m \\ i=2, j=1 & i=2, j=2 & \dots & i=2, j=m \\ \vdots & \vdots & \ddots & \vdots \\ i=m, j=1 & i=m, j=2 & \dots & i=m, j=m \end{array}$$

and observing that  $\forall i=1, 2, \dots, m$ , we’re doing  $m$  computations for  $j$  and needing to fetch  $m$  values for each  $\lambda_j, y_j, X^{(j)}$ , and so on.

Consider this computation: for a given, single  $i \in \{1, 2, \dots, m\}$ , define

$$(93) \quad f_{1i}(\lambda) := \frac{1}{2} \sum_{j=1}^m \lambda_j y^{(j)} K(X^{(i)}, X^{(j)})$$

For this step, we’ll only need to do  $m$  fetches for the  $\lambda_j, y^{(j)}, X^{(j)}$  values, and  $X^{(i)}$  value will be fetched once. As this is a summation over a potentially large vector ( $m$  can be big), this looks like a good case/candidate for the usage of parallel *reduce* algorithm. The work complexity of parallel reduce is  $O(\log m)$  [3]<sup>5</sup>. Theano has an implementation of reduce in **theano.reduce**.

Once all  $m$   $f_{1i}$ ’s are obtained, for  $i=1, 2, \dots, m$ , then parallel reduce can be used again (especially if  $m$  is large!). Also, empirically, I found that using **theano.reduce** again helped to circumvent the problem of the maximum recursion limit for Python<sup>6</sup>, which is inherent with Python (cf. `import sys sys.getrecursionlimit()`). In practice, above about 10000 recursions, the Python script fails with run-time errors.

Nevertheless, in this second (parallel) reduce step, we are doing

$$f_1 = \sum_{i=1}^m \lambda_i y^{(i)} f_{1i}(\lambda)$$

with  $m$  fetches of values for  $\lambda_i, y^{(i)}$ . The work complexity here for this reduce step is again  $O(\log(m))$

Thus, we are doing, for 2 (parallel) reduces,  $2m$  fetches (or reads), for each  $\lambda_i$  or  $y^{(i)}$  or  $X^{(i)}$ .

The total work complexity is  $O(2 \log(m))$ .

Likewise, for the computation of the intercept  $b$  in Eqn. 90, after minimizing  $W(\lambda)$  by varying  $\lambda$ , I also employed parallel reduce via **theano.reduce** (but only once for the single sum) and for the prediction step for  $\hat{y}$  in Eq. 91

12.7.1. *Code (theano/Python script), jupyter notebook accompanying code*. SVM is implemented as described above, in particular Eqns. 86, 87, in the Python class **SVM\_parallel**. The default kernel function  $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is the radial basis function, which takes the form of a gaussian, is first implemented and I can implement other kernel functions easily, as a Python function object and Python class member, in the future.

Take note that for the (currently) implemented radial basis function, Python function (object) **rbf** in **SVM.py** of [github:ernestyalumni/MLgrabbag/ML](https://github.com/ernestyalumni/MLgrabbag/ML), what’s formulated is this:

$$(94) \quad K(X^{(i)}, X^{(j)}) = \exp \left( -\frac{\|X^{(i)} - X^{(j)}\|^2}{2\sigma^2} \right)$$

with  $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

Take a look at the  $\sigma \in \mathbb{R}$  parameter in Eq. 94.  $\sigma$  is analogous to the variance of a Gaussian (normal) distribution. For other implementations, notably **libsvm** and sci-kit learn, they use this form of the radial basis function:

$$K(X^{(i)}, X^{(j)}) = \exp \left( -\gamma \|X^{(i)} - X^{(j)}\|^2 \right)$$

So  $\gamma$  parameter here is equivalent to  $\sigma$ :

$$\gamma = \frac{1}{2\sigma^2}$$

While this redefinition makes no change to the formulation above, this is something to note when when using **libsvm**, sci-kit learn, or **SVM.py** here when manually inputting the parameters to train models.

The theano/Python code follows directly from Eqns. 86, 87 and is in the /ML subfolder of the github repository **MLgrabbag**<sup>7</sup>, in **SVM.py**. Wherever a summation is seen in the mathematical formulation, **theano.reduce** is used.

In the **SVM\_parallel** Python class method **build\_W**, I code a **theano.reduce** inside a **theano.reduce** and show it’s possible to be done. This represents, both formally and the parallel reduction on the GPU, the double summation that we sought to compute in 86 for  $W(\lambda)$ .

The jupyter notebook **SVM\_theano.ipynb** in the same github repository steps through how I developed and used **SVM\_parallel**, training it on a number of sample datasets. Because of the interactivity of jupyter notebook, I invite others to explore and play with the notebook if further clarification on **SVM\_parallel**, or how to use it, is needed<sup>8</sup>

<sup>3</sup>[sklearn.svm.SVC](https://sklearn.org/stable/modules/generated/sklearn.svm.SVC.html)

<sup>4</sup>Quiz: All Pairs N-Body

<sup>5</sup>Step Complexity of Parallel Reduce - Intro to Parallel Programming, Udacity

<sup>6</sup>max recursion limit #689

<sup>7</sup>[github:ernestyalumni/MLgrabbag](https://github.com/ernestyalumni/MLgrabbag)

<sup>8</sup>[github:ernestyalumni/MLgrabbag SVM theano.ipynb](https://github.com/ernestyalumni/MLgrabbag/blob/master/SVM_theano.ipynb)

12.8. Immediate Results from training on sample datasets.

12.8.1. *Real-World Examples.* I trained a SVM on 2 of the real-world data sets provided by Hsu, Chang, and Lin [19], one for astroparticles and another for vehicles, using, for hardware, a NVIDIA GeForce GTX 980 Ti. Checking the computational graph generated by theano (using `theano.function.maker.fgraph.toposort()`), `nvidia-smi -l 2` (monitoring real-time GPU usage), and the (usual, in Utilities) CPU resources System Monitor.

I will copy the results from Hsu, Chang, and Lin [19] for comparison. The accuracy measure is determined from the given *test* data, *not* on the training data (which is part of good machine learning and scientific practice).

Applications	# training data	# testing data	# features	# classes	$C =$	$\gamma =$	Accuracy by <code>libsvm</code>
Astroparticle <sup>9</sup>	3089	4000	4	2	2.0	2.0	96.9%
Vehicle <sup>10</sup>	1243	41	21	2	128.0	0.125	87.8%

Table 1: Sample Dataset Problem characteristics and accuracy performance [19].

Applications	$C =$	$\sigma =$	$\alpha =$	# iterations	Time to train (on GTX 980Ti)	Accuracy by <code>SVM_parallel</code>
Astroparticle	2.0	0.30	0.001	15	1h 7min 18s	96.1%
Vehicle	128.0	2.0	0.001	20	14min 54s	95.1%

Table 2: Results of training on Sample Datasets with `SVM_parallel`

The very last result testing on the test data for vehicles is promising for `SVM_parallel`. At this point, I would invite others to suggest sample and real-world datasets to train and test on, using `SVM_parallel`, as I also try to find other datasets, and add onto the jupyter notebook [SVM theano.ipynb on github](#). It’d be interesting to vary the *number of training examples*, to find a dataset with more than 10000 ( $m > 10000$ ) examples and see how `SVM_parallel` can scale with large data sets (indeed, for  $m > 10000$ , the SVM would have  $m > 10000$  support vectors in the model), and vary the *number of features* (whether SVM does better with large or small number of features, relative to  $m$ ).

12.9. Conclusions/Summary/Dictionary between Math and Code. I had reviewed the motivation and derivations for SVM.

What’s novel is that, given the GPU(s), I implemented *constrained gradient descent* or *projected gradient descent*, for training models, instead of Quadratic Programming, that computes a quadratic form (to tackle the double summation in the dual formulation), through SMO, as used before (e.g. `libsvm`, sci-kit learn). Its (*constrained gradient descent* or its implementation here `SVM_parallel`) work complexity is  $O(2 \log m)$ , as opposed to  $O(m^2)$ . This was achieved by using theano’s `reduce`, inside a `reduce`.

Its (i.e. `SVM_parallel`) promising to be scalable to large datasets ( $m > 10000$ ). I seek to find large datasets to train and test on and are appropriate for binary classification, and invite others to make suggestions or play with the code and jupyter notebook itself.

I’ll provide a 1-to-1 dictionary here between the mathematical formulation and the Python code. As a note on software engineering, object-oriented programming (OOP) and how to code classes, I had sought to identify (make isomorphisms) and design Python classes and function objects with 1-to-1 correspondence to the mathematical formulation. The hope is that it would allow other developers to rapidly make progress in improving upon the code or to rapidly understand its usage and apply it as they’d like to see fit.

we seek to minimize

$$W(\lambda) = - \sum_{i=1}^m \lambda_i + \frac{1}{2} \sum_{i,j=1}^m \lambda_i \lambda_j y^{(i)} y^{(j)} K(X^{(i)}, X^{(j)})$$

SVM\_parallel.build\_W

by iterating  $t = 0, 1, \dots$ , as such:

SVM\_parallel.train\_mode\_full(max\_iters=250)

$$\lambda'_i(t+1) := \lambda_i(t) - \alpha \text{grad} W(\lambda)$$

$$\lambda''_i(t+1) := \mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)} = 0}(\lambda'_i(t+1))$$

SVM\_parallel.build\_update

$$\lambda_i(t+1) := \Pi_{0 \leq \lambda_i \leq C}(\lambda''_i(t+1))$$

where

$$\mathbf{P}_{\sum_{i=1}^m \lambda_i y^{(i)} = 0}(\lambda'_i(t+1)) = \lambda'_i(t+1) - \frac{\sum_{i=1}^m y^{(i)} \lambda'_i(t+1)}{\sum_{i=1}^m (y^{(i)})^2} y^{(i)}$$

$\Longleftrightarrow$

updatelambda\_mult=updatelambda\_mult-T.dot(y,updatelambda\_mult)/T.dot(y,y)\*y in SVM.build\_update

$$\Pi_{0 \leq \lambda_i \leq C}(\lambda''_i(t+1)) = \begin{cases} C & \text{if } \lambda''_i(t+1) > C \\ \lambda''_i(t+1) & \text{if } 0 \leq \lambda''_i(t+1) \leq C \\ 0 & \text{if } \lambda'_i(t+1) < 0 \end{cases}$$

updatelambda\_mult=T.switch(T.lt(C,updatelambda\_mult),C,updatelambda\_mult) in SVM.build\_update

$\Longleftrightarrow$

updatelambda\_mult=T.switch(T.lt(updatelambda\_mult,lower\_bound),lower\_bound,updatelambda\_mult) in SVM.build\_update

Finally, to tie it back into my original motivation, now that SVM is natively implemented in theano, it would be interesting to try to develop (and of course find appropriate datasets to train and test on) a DNN that will have as its “outer” or last layer to be a SVM. Since SVM is now part of the theano computational graph, optimization (the so-called “backpropagation” step) will be done automatically and simply with theano’s `grad`, on all the parameters or “weights” of the entire model.

Part 5. Notes

Restricted Boltzmann machine - estimate a probability distribution

Recurrent neural network - creates an internal state of the network which allows it to exhibit dynamic temporal behavior

[How to choose the number of hidden layers and nodes in a feedforward neural network?](#)

“In sum, for most problems, one could probably get decent performance (even without a second optimization step) by setting the hidden layer configuration using just two rules: (i) number of hidden layers equals one; and (ii) the number of neurons in that layer is the mean of the neurons in the input and output layers.”



## REFERENCES

- [1] Trevor Hastie, Robert Tibshirani, Jerome Friedman. **The Elements of Statistical Learning: Data Mining, Inference, and Prediction**, Second Edition (Springer Series in Statistics) 2nd ed. 2009. Corr. 7th printing 2013 Edition. ISBN-13: 978-0387848570. [https://web.stanford.edu/~hastie/local.ftp/Springer/OLD/ESLII\\_print4.pdf](https://web.stanford.edu/~hastie/local.ftp/Springer/OLD/ESLII_print4.pdf)
- [2] Jared Culbertson, Kirk Sturtz. *Bayesian machine learning via category theory*. [arXiv:1312.1445](https://arxiv.org/abs/1312.1445) [math.CT]
- [3] John Owens. David Luebki. *Intro to Parallel Programming. CS344*. **Udacity** <http://arxiv.org/abs/1312.1445> Also, <https://github.com/udacity/cs344>
- [4] CS229 Stanford University. <http://cs229.stanford.edu/materials.html>
- [5] Richard Fitzpatrick. “Computational Physics.” <http://farside.ph.utexas.edu/teaching/329/329.pdf>
- [6] LISA lab, University of Montreal. Deep Learning Tutorial. <http://deeplearning.net/tutorial/deeplearning.pdf> September 2015.
- [7] Kurt Hornik. “Approximation Capabilities of Muiltlayer Feedforward Networks.” **Neural Networks**, Vol. 4, pp. 251-257. 1991
- [8] Kurt Hornik. Maxwell Stinchcombe and Halbert White. “Multilayer Feedforward Networks are Universal Approximators.” **Neural Networks**, Vol. 2, pp. 359-366, 1989.
- [9] Thomas Nowak. “Implementation and Evaluation of a Support Vector Machine on an 8-bit Microcontroller.” Univ.Ass. Dipl.-Ing. Dr.techn. Wilfried Elmenreich Institut für Technische Informatik Fakultät für Informatik Technische Universität Wien. Juli 2008. <https://www.lri.fr/~nowak/misc/bakk.pdf>
- [10] J. Shawe-Taylor and N. Cristianini, Support Vector Machines and other kernel-based learning methods, Cambridge University Press (2000).
- [11] Edwin K. P. Chong and Stanislaw H. Zak. **An Introduction to Optimization**. 4th Edition. Wiley. (January 14, 2013). ISBN-13: 978-1118279014
- [12] Lecture by Harikrishna Narasimhan. *Optimization Tutorial 3: Projected Gradient Descent, Duality*. **E0 270 Machine Learning**. Jan 23, 2015. <http://drona.csa.iisc.ernet.in/~e0270/Jan-2015/Tutorials/lecture-notes-3.pdf>
- [13] Christopher M. Bishop. **Pattern Recognition and Machine Learning** (Information Science and Statistics). Springer (October 1, 2007). ISBN-13: 978-0387310732
- [14] Bertrand Clarke, Ernest Fokoue, Hao Helen Zhang. **Principles and Theory for Data Mining and Machine Learning** (Springer Series in Statistics) Springer; 2009 edition (July 30, 2009). ISBN-13: 978-0387981345
- [15] Hubert Nguyen. **GPU Gems 3**. Addison-Wesley Professional (August 12, 2007). ISBN-13: 978-0321515261. Also made available in its entirety online at [https://developer.nvidia.com/gpugems/GPUGems3/gpugems3\\_pref01.html](https://developer.nvidia.com/gpugems/GPUGems3/gpugems3_pref01.html)
- [16] Scikit-learn: Machine Learning in Python, Pedregosa *et al.*, **JMLR** **12**, pp. 2825-2830, 2011.
- [17] C.-C. Chang and C.-J. Lin. *LIBSVM : a library for support vector machines*. **ACM Transactions on Intelligent Systems and Technology**, 2:27:1–27:27, 2011.
- [18] J. Platt. *Fast training of support vector machines using sequential minimal optimization*. In A. Smola B. Schölkopf, C. Burges, editor, **Advances in Kernel Methods: Support Vector Learning**. MIT Press, Cambridge, MA, 1998.
- [19] Chih-Wei Hsu, Chih-Chung Chang, and Chih-Jen Lin. *A Practical Guide to Support Vector Classification*. <http://www.ee.columbia.edu/~sfchang/course/spr/papers/svm-practical-guide.pdf>