

IN GOD WE TRUST
DEEP LEARNING COURSE
2021-2022 FALL SEMESTER

HW01

...

Omid Sharafi
(400201518)

Instructor:
Dr.Emad Fatemizadeh

November 19, 2021



**Sharif
University
of
Technology**



Contents

1	Theoretical Problems	2
1.1	Vapnik-Chervonenkis	2
1.2	MSE	5
1.3	ML	6
1.4	Regression	7
2	Finger spelling	8
2.1	Visualize data set	8
2.2	Train model and draw loss and weights histogram	9
2.2.1	SGD optimizer result	9
2.2.2	ADAM optimizer result	12
2.2.3	Dropout	15
2.2.4	Batch Normalization	18
2.2.5	Unbalanced data-set	18
2.2.6	Webcam	18
3	Inverted Pendulum	20



1 Theoretical Problems

1.1 Vapnik-Chervonenkis

$VC_dimension$ برای کلاس فرضیات H ، مجموعه نقاط C با سیزده تریه H shatter می شود. اگر به ازای هر انتخاب H برای هر کلاس از نقاط C ، کلاس H می تواند جدا سازی با دقت 100 درصد این مجموعه را داشته باشد.

حالت $VC_dimension$ کلاس فرضیات H برابر با n خواهد بود اگر و تنها اگر داشته باشیم:

- ① مجموعه C با سیزده رجه داشته باشد که تریه H shatter شود.
- ② به ازای هر مجموعه نقاط با سیزده n ، رجه برای n تین H وجود دارد که نمی گذارد مجموعه تریه H shatter شود.

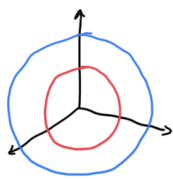
این معرکه عمداً توانایی کلاس فرضیات H را در فضای ورودی نشان می دهد.

Figure 1: VC dimension



A: VC $\text{sgn}(x^T x + \theta)$? (in 3D space)

answer: VC = 1



$$f(\vec{x}) = \begin{cases} \|\vec{x}\|^2 + \theta < 0 \rightarrow - \\ \|\vec{x}\|^2 + \theta > 0 \rightarrow + \end{cases}$$

It's obvious that for single node, VC can shatter it.

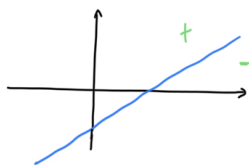
(just adjust θ)

For 2 points x, y :

if $(\|x\| = \|y\|) \rightarrow f(\vec{x}) = f(\vec{y})$: just label \vec{x}, \vec{y} different. ~~X~~
 if $(\|x\| > \|y\|) \rightarrow \|x\|^2 + \theta > \|y\|^2 + \theta$: $\begin{cases} \text{Label}(\vec{x}) = - \\ \text{Label}(\vec{y}) = + \end{cases}$ ~~X~~

B: VC $\text{sgn}(x_1\theta_1 + x_2\theta_2 + \theta_3)$?

answer: VC = 3



for 3 points:



at worst

one label different from two others, so we can draw a line to separate that point ✓

for 4 points: if at least 3 points on a line



else, we have a convex hull

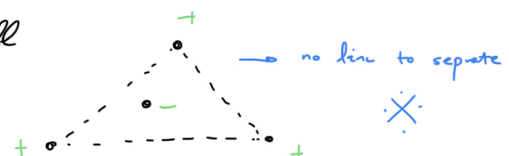


Figure 2: Part A, B



C. VC superplane in N 'D space ?

answer : $VC = N+1$

$$f(x) = \text{sign}(w^T x + b), w \in \mathbb{R}^N, b \in \mathbb{R}$$

consider $N+1$ point : $\vec{p}_1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \vec{p}_2 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \dots, \vec{p}_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \vec{p}_{n+1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

so for each point i ($1 \leq i \leq n$) $f(x) = \text{sign}(w_i + b)$ \rightarrow for each labely
for point $i = n+1$ $f(x) = \text{sign}(b)$ \rightarrow we can assign perfect w_i and b ✓

For upper bound : ① proof that set of $f(x) = \text{sign}(w^T x)$ can not shatter $N+1$

: assume that $S = \{x_1, \dots, x_{N+1}\}$ shattered by $f(x)$. so for each 2^{N+1} possible labeling, we have a proper w_k so that $f(w_k x)$ produce true labeling for all points in S .

$$H \triangleq \begin{bmatrix} w_1^T x_1 & w_2^T x_1 & \dots & w_{2^{N+1}}^T x_1 \\ \vdots & \vdots & \dots & \vdots \\ w_1^T x_{N+1} & w_2^T x_{N+1} & \dots & w_{2^{N+1}}^T x_{N+1} \end{bmatrix}$$

\hookrightarrow each of the columns make our labeling

Linear independent of rows

$$\Rightarrow \text{no } v \text{ that } v^T H = 0$$

output labeling $_k = k$ th entry $v^T H = v^T X w_k$

$$1 \leq k \leq 2^{N+1}$$

$$\Rightarrow \text{there is a } k \text{ that } \text{sign}(X w_k) = \text{sign}(v) \Rightarrow v^T X w_k > 0$$

$$\Rightarrow \text{rank}(H) = \text{number of rows} = N+1$$

$$H = XW \Rightarrow \text{rank}(H) \leq \min\{\text{rank}(X), \text{rank}(W)\} \leq N$$

$\leq N$ we are inside \mathbb{R}^N space

$$\leq N \cdot X$$

② now we have $x = \begin{bmatrix} x' \\ 1 \end{bmatrix}, w = \begin{bmatrix} w' \\ b \end{bmatrix}$ and maximum VC-dimension

$$\text{of a } f(x') = \text{sign}(w'^T x') \leq N, \text{ so VC-dimension } f(x) \leq N+1$$

Figure 3: Part C



1.2 MSE

Problem 2. MSE as loss function in regression problem $\xrightarrow{\text{means}}$ Assuming normal distribution of data in ML problem

• MSE as loss function :
$$w_{\text{opt}} = \underset{w}{\operatorname{argmin}} \sum (y_i - w^T x_i)^2$$

• normal distribution assumption : $y_i = w^T x_i + \varepsilon_i$, $\varepsilon_i : \mathcal{N}(\cdot, \sigma^2)$

$\xrightarrow{\text{solve ML}}$
$$w_{\text{opt}} = \underset{w}{\operatorname{argmax}} \prod f(y_i | w, x_i, \sigma) =$$
$$\underset{w}{\operatorname{argmax}} \prod \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(y_i - w^T x_i)^2}{2\sigma^2}} \stackrel{\ln}{=} \underset{w}{\operatorname{argmax}} \sum \left(\underbrace{-\ln(\sqrt{2\pi} \sigma)}_{\text{cte}} - \underbrace{\frac{(y_i - w^T x_i)^2}{2\sigma^2}}_{\text{cte}} \right)$$
$$= \underset{w}{\operatorname{argmax}} \sum -(y_i - w^T x_i)^2 = \underset{w}{\operatorname{argmin}} \sum (y_i - w^T x_i)^2 \quad \checkmark$$



1.3 ML

Problem 3. calculate normal distribution parameters of $\{x_i\}_{i=1}^n$ with MSE

$$(\mu, \sigma)_{\text{ans}} = \underset{\mu, \sigma^2}{\operatorname{argmax}} \prod f(x_i | \mu, \sigma^2) = \underset{\mu, \sigma^2}{\operatorname{argmax}} \prod \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\stackrel{\ln}{=} \underset{\mu, \sigma^2}{\operatorname{argmax}} \sum \left(-\ln(\sqrt{2\pi}\sigma) - \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$\frac{\partial}{\partial \mu} = 0 \Rightarrow \frac{1}{2\sigma^2} \sum 2(x_i - \mu) = 0 \Rightarrow n\mu = \sum x_i \Rightarrow \boxed{\mu = \frac{1}{n} \sum x_i}$$

$$\frac{\partial}{\partial \sigma^2} = 0 \Rightarrow \sum \left(-\frac{1}{2\sigma} \frac{1}{\sigma} + \frac{1}{2\sigma} \frac{(x_i - \mu)^2}{\sigma^3} \right) = 0$$

$$\Rightarrow \frac{n}{2\sigma^2} = \frac{1}{2\sigma^4} \sum (x_i - \mu)^2 \Rightarrow \boxed{\sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2}$$



1.4 Regression

Problem 4.

A) Regression: $w_{\text{ans}} = \underset{w}{\operatorname{argmin}} (Y - W \cdot X)^T (Y - W \cdot X)$

$$= \underset{w}{\operatorname{argmin}} Y^T Y - Y^T (W \cdot X) - (W^T X^T) Y + (W^T X^T) (W \cdot X)$$

$$\frac{\partial}{\partial w} = -Y^T X - X^T Y + X^T (X \cdot W) + (W^T X^T) X \stackrel{\text{scalar}}{=} X^T (X \cdot W) = X^T Y \Rightarrow \boxed{w = (X^T X)^{-1} X^T Y}$$

B) if $X^T X$ has no inverse

$$w_{\text{ans}} = \underset{w}{\operatorname{argmin}} \|Y - W \cdot X\|_2^2 + \lambda \|w\|_2^2 \quad (\text{Regularization } L_2)$$

$$\frac{\partial}{\partial w} = -2 X^T Y + 2 X^T X w + 2 \lambda w = 0 \Rightarrow \boxed{w = (X^T X + \lambda I)^{-1} X^T Y}$$

2 Finger spelling

2.1 Visualize data set

In this part, we separated train, validation data. Made the output in one hot format. Shuffled the data and finally normalized the data of each part of the data set.

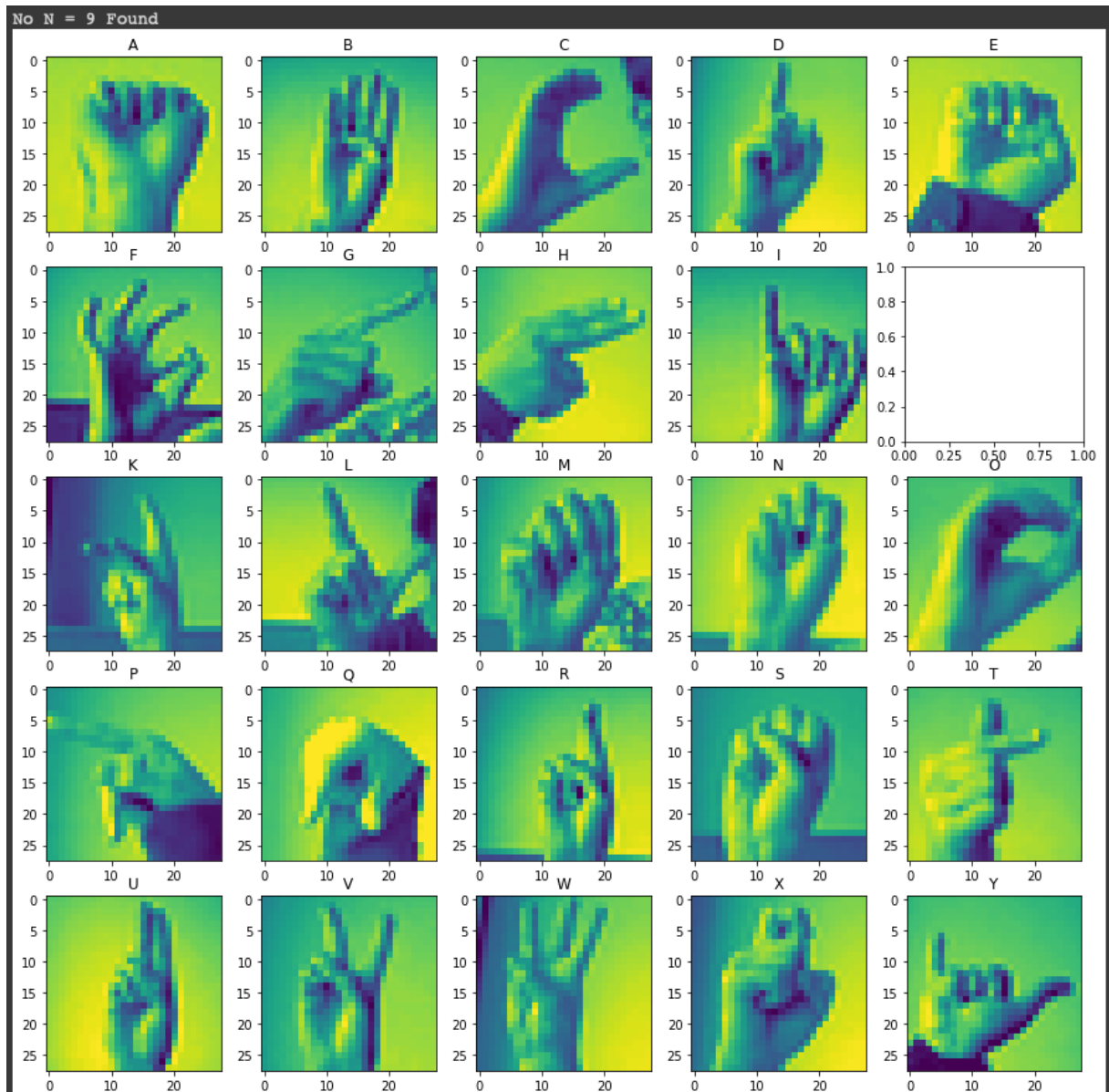


Figure 4: Finger Spelling Data set



2.2 Train model and draw loss and weights histogram

2.2.1 SGD optimizer result

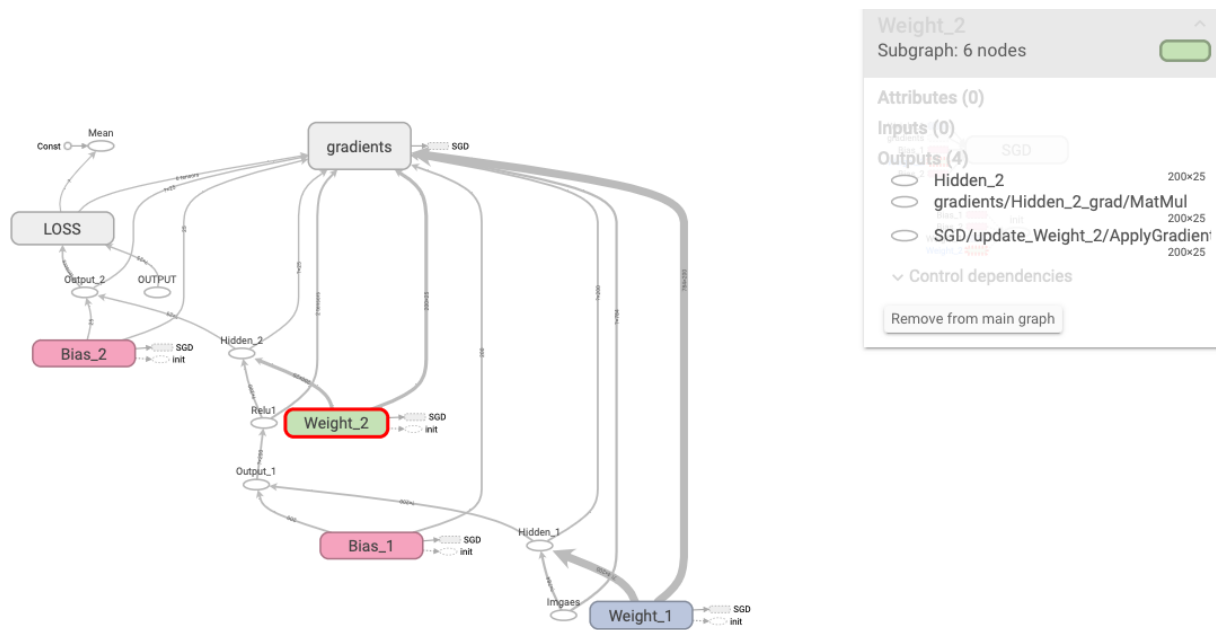


Figure 5: SGD network design(one 200 node hidden layer)

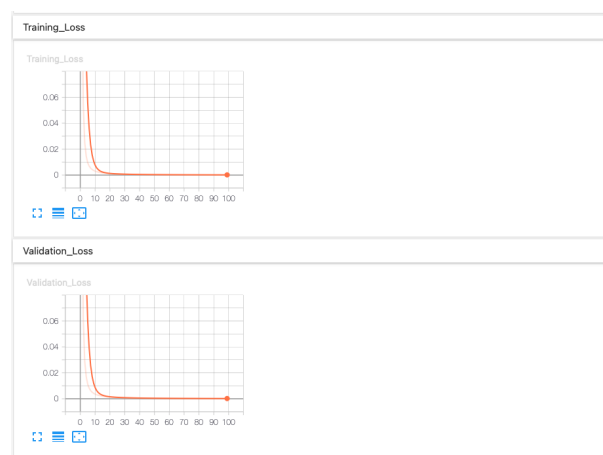


Figure 6: Loss

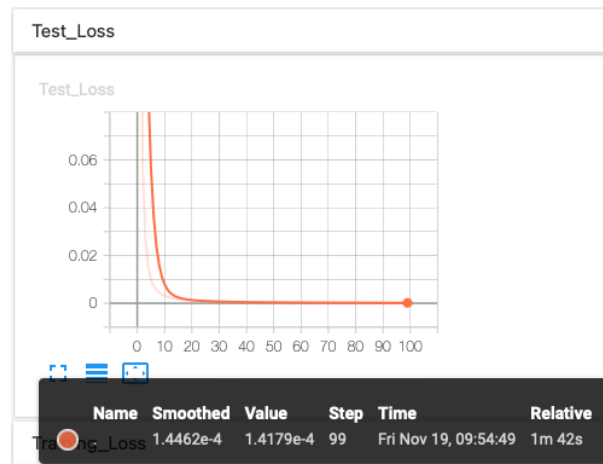


Figure 7: Loss

```
INFO:tensorflow:Restoring parameters from ./drive/MyDrive/graphs/sgd/sgd_model.ckpt  
SGD Validation Accuracy = 1.0  
SGD Test Accuracy = 1.0
```

Figure 8: Accuracy

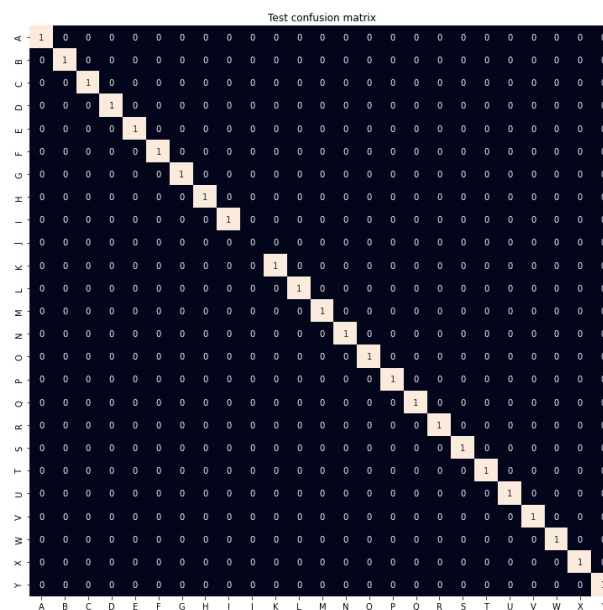


Figure 9: Confusion matrix



As you can see, the network learn data perfect and it is not over-fitting because the accuracy on the test is also perfect. Also to check if there is a any problem in training method or not I just reduced the epochs number and increased batch size and you can see the result at results folder. More general and deeper network is available at the end of the report and saved in the best folder.

And about the weighs distribution, you can see change its distribution at first epochs.

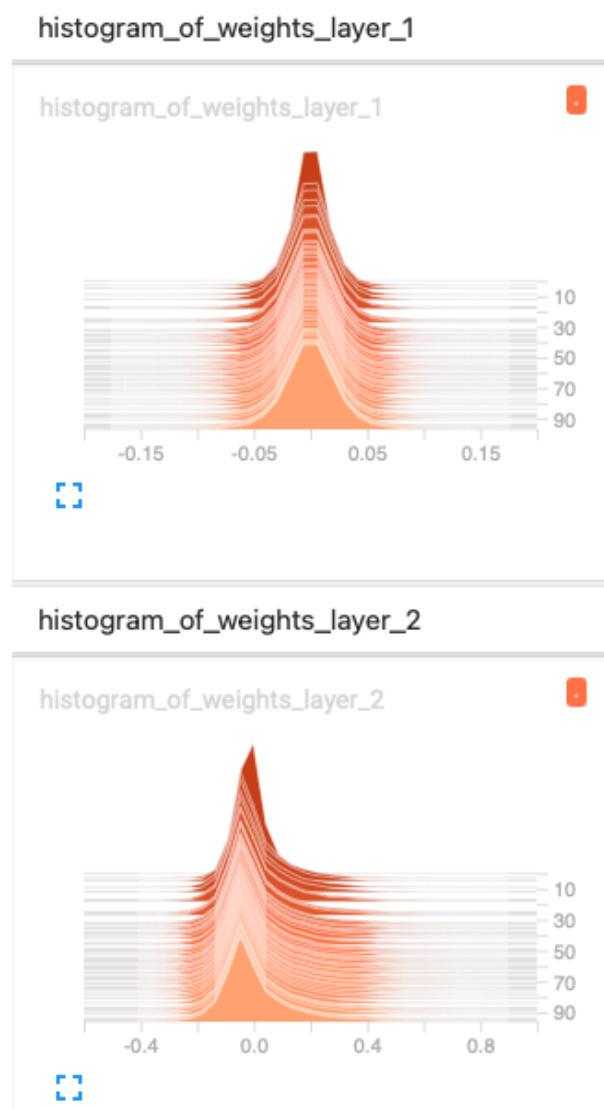


Figure 10: Weighs distribution



2.2.2 ADAM optimizer result

As you can see, both ADAM and SGD results are good. Although ADAM is more faster and better but it has lots of oscillation in its accuracy and convergence. For more approaches for example we can use NADAM, AdaGrad and other training approaches.

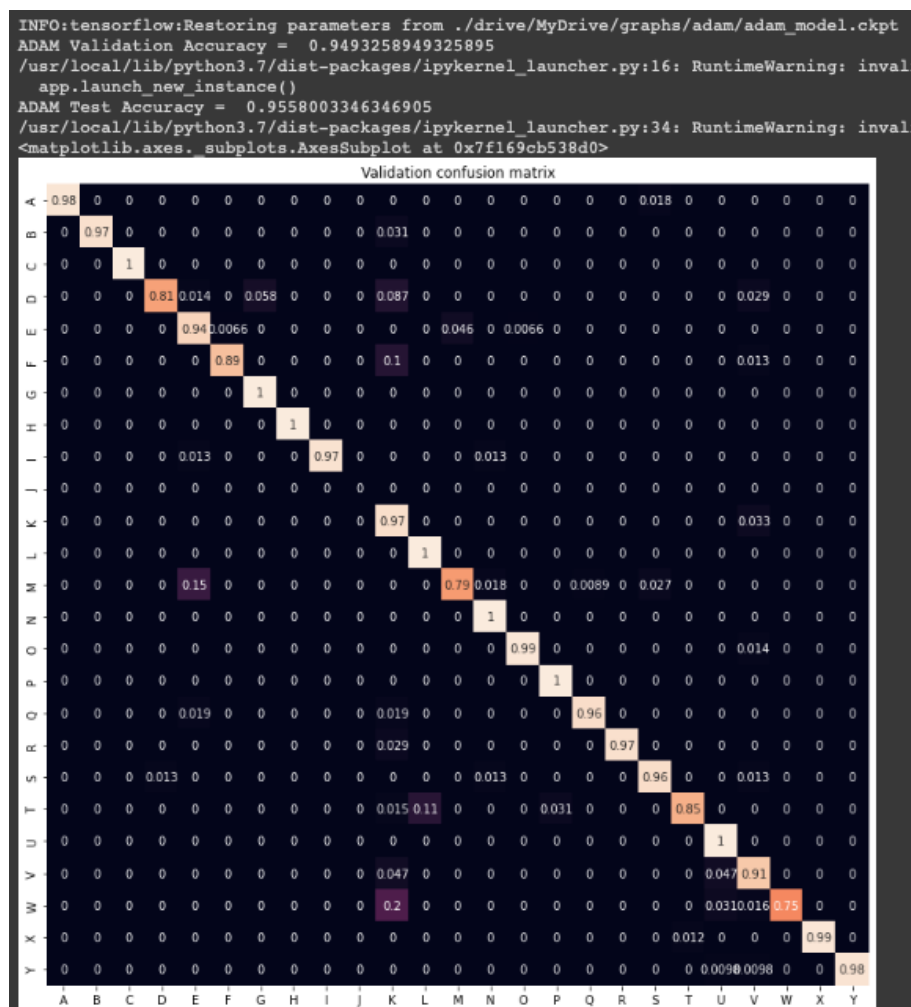


Figure 11: Confusion matrix (ADAM)



Department of Electrical Engineering

And about the weights distribution, you can see change in distribution and increment of variance and separation during epochs in ADAM method.



Figure 12: Loss (ADAM)

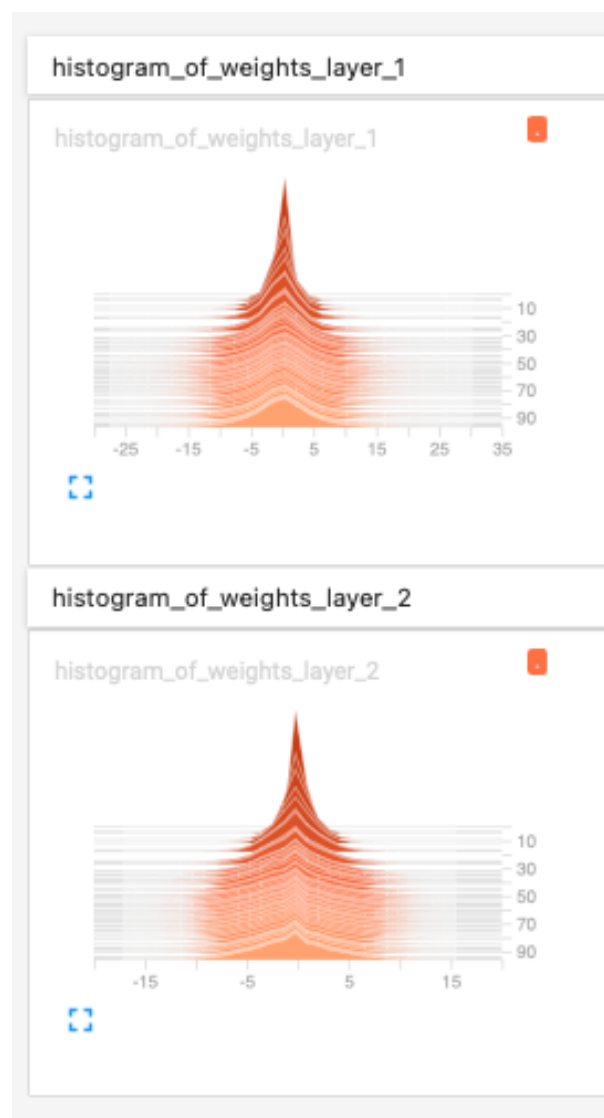
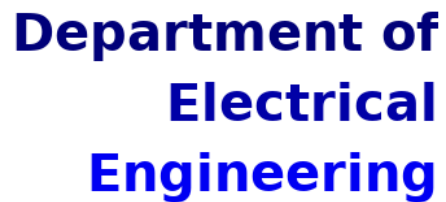


Figure 13: Weighs distribution (ADAM)



From this part, to see the change result first we try to reduce our normal SGD network accuracy by increasing batch-size and run network just 5 epochs. Then I tried different combination of dropout rate and epochs numbers and its really fun that when we have:

then we can see some improvement in result. I mean that it seems that number of training on data is constant but one approach without dropout and the other one with dropout. Also the drop out rate should not be high.





Department of Electrical Engineering

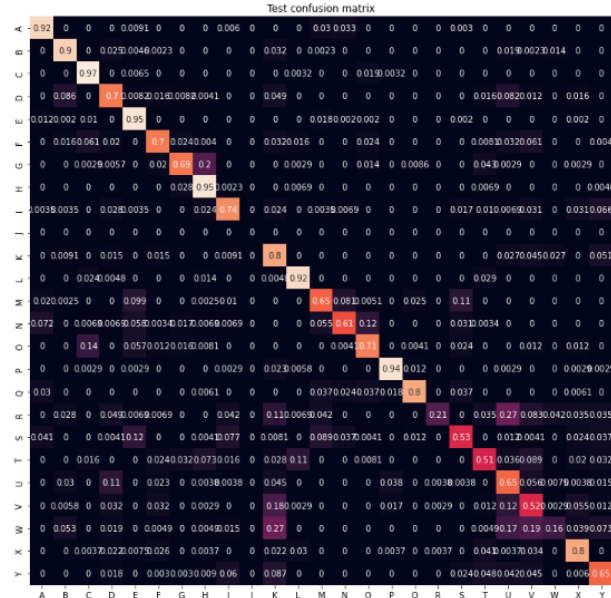


Figure 15: Dropout (5 epochs with rate = 0.5) (Acc = 73)

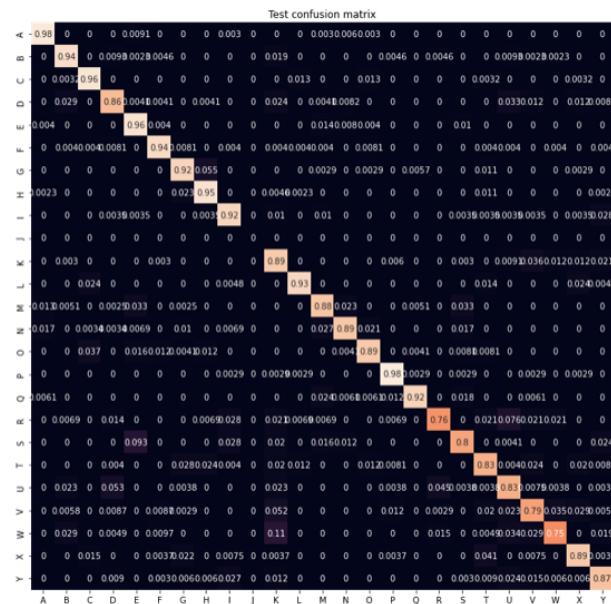
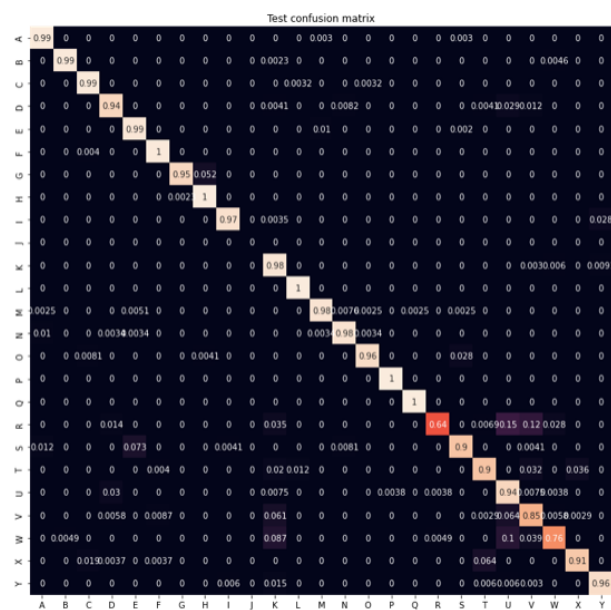
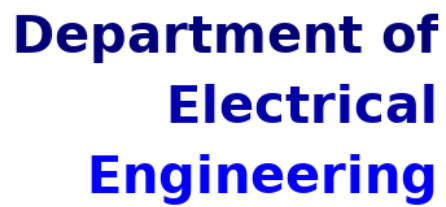


Figure 16: Dropout (10 epochs with dropout rate = 0.5) (Acc = 89.8)





2.2.4 Batch Normalization

For this part I used 5 layers with batch normalization and there was a small improvement in accuracy.

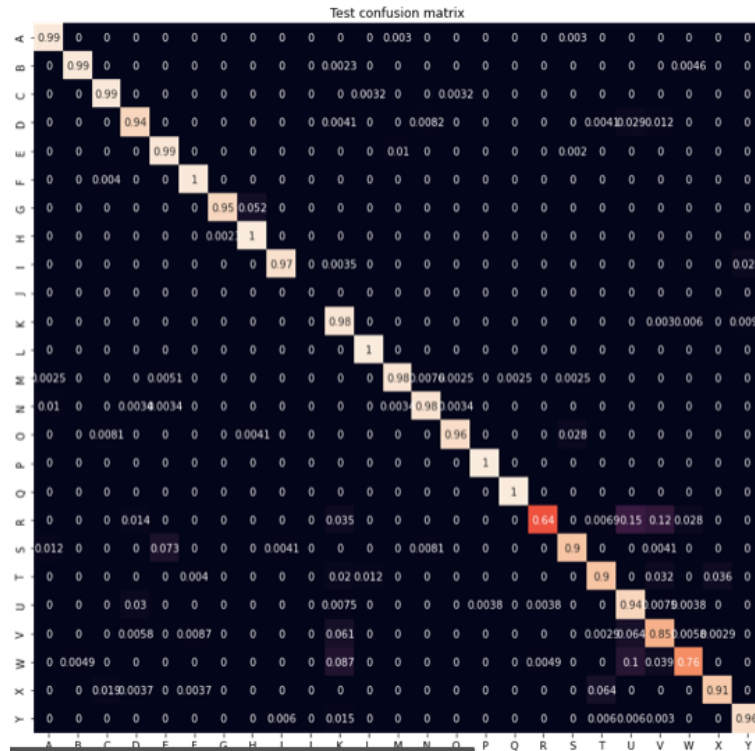


Figure 18: (Acc = 95.2)

2.2.5 Unbalanced data-set

For Label D, I had 176 samples in training data-set and I just kept 26 of them randomly. As you can see, error increased in label D. For solving the unbalanced data-set problem, we can repeat the available data or we can estimate its distribution and make more random samples from that. We can also use rotation or scaling approaches for that category images to increase their number.

2.2.6 Webcam

I also did this part in corporation with my friend Ali Arasteh. Although I trained and tested different networks and approaches, the result was not

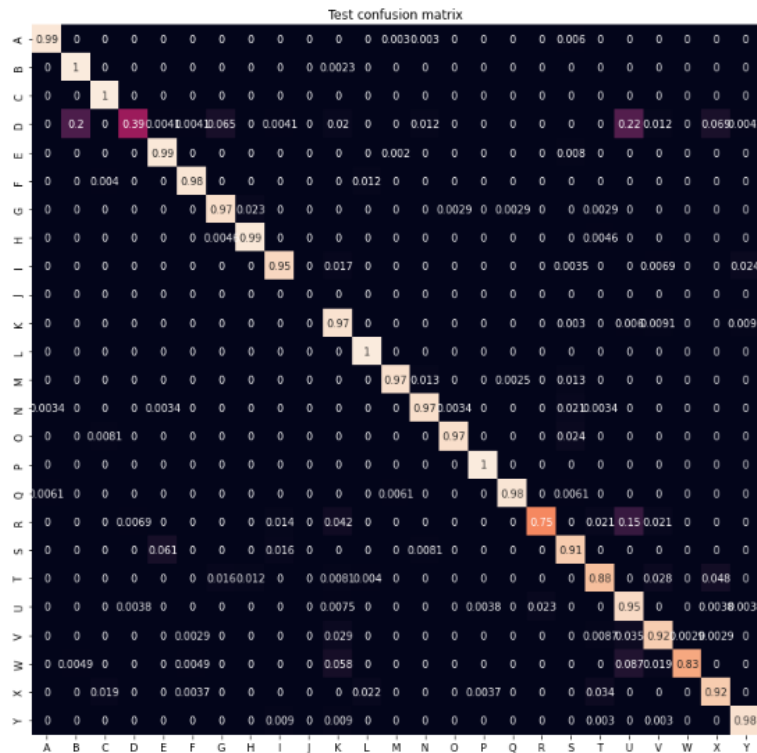


Figure 19: Unbalanced data-set for label 'D'

really perfect but for example I got the answer for L label. The problem here is that we can not normalize just one picture and we should have more images from that environment. Also as I overlooked the data-set, the images was really simple usually with solid background. We should have more powerful data-set for training.



3 Inverted Pendulum

I completed this part. All of the codes and results are available in this github repository.