

DL – Exercise 1 – Report

Antigen Discovery for SARS-CoV-2 (“Corona”) Virus Vaccine

Eldan Chodorov

ID: 201335965

Amit Keinan

ID: 208296426

November 2021

Table of Contents

2.....	Table of Contents
3.....	Data Handling
3.....	Models Architectures
3.....	Experiments
3.....	First Training
4.....	Hyper-Parameters Tuning
5.....	"Basic" Hyper-Parameters Tuning
5.....	Batch size
5.....	Learning rate
9.....	Epochs number
9.....	Activation function
9.....	Architectures Comparison
10.....	Adjusting learning rate and epochs number
13.....	Dealing With Imbalanced Data
13.....	Evaluation
13.....	Oversampling
15.....	Final training
16.....	Inference on Spike Protein
17.....	Optimize Input
19.....	Appendix – Problematic Experiments
19.....	First Training
20.....	Hyper-Parameters Tuning
20.....	"Basic" Hyper-Parameters Tuning
20.....	Batch size
21.....	Learning rate
22.....	Epochs number
23.....	Activation function
23.....	Architectures Comparison

Data Handling

Each data item is a 9-lengthed sequence of characters, and all of the characters are in a vocabulary of 20 characters. We decided to map each sequence to 9×20 vector – concatenation of nine (one for character) one-hot encoding vectors, each one is of size 20. This seems like a natural choice because we do not have prior assumptions about the data, we want to keep the order and we do not know if there is any connection between some of the characters.

Models Architectures

We will try multi-layered perceptron networks with different parameters. The parameters are numbers of layers, number of neurons in each hidden layer and activation function. In all cases the first layer will get 9×20 -sized input vector and the last layer will output 1-sized output vector after going through sigmoid activation. In inference, we round the output, so effectively the threshold is 0.5.

Experiments

We actually wrote this report twice. In the first time we trained our network and saw good results – accuracy, recall and precision were all above 0.9. We tried different hyper-parameters and optimized the results. We knew we would have to deal with the imbalanced data, but were misled by the high scores and thought this specific issue is a matter of some precents improvement. However, when we measured a relevant metric directly (false negative ratio) we noticed that we trained models that only output “True”. We easily fixed it with oversampling the negative samples and repeated all our experiments. We left here (in the end under “Appendix – Problematic Experiments”) also the first set of experiments to show the full research process.

First Training

For the first training, we tried feed forward network with one hidden layer with 256 neurons with relu activation. We trained the network for 5 epochs. We used batch size of 64.

We used binary cross entropy loss and learning rate of 0.001.

Results on the (oversampled) test set:

(tn:2170, fp:171, fn:381, tp:2120)

true positive rate (AKA recall): 0.848, true negative rate: 0.923

accuracy: 0.886

recall: 0.848

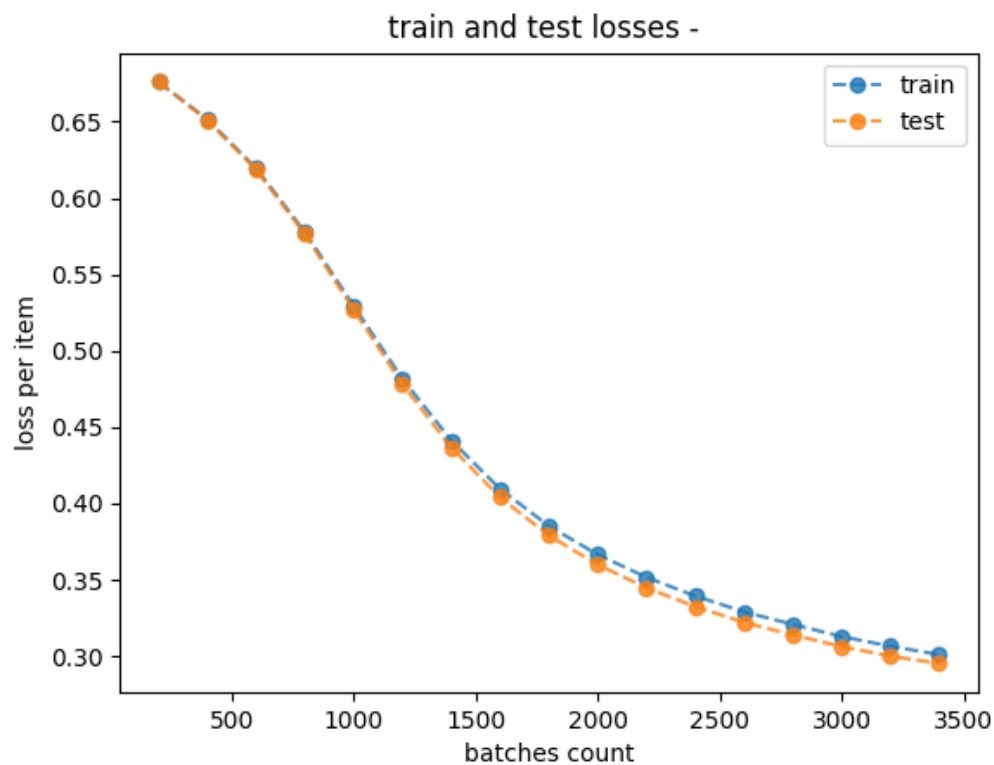
precision: 0.925

f1: 0.885

The results were measured on the oversampled data. It is not the best practice and we will improve it in the section that deals directly with imbalanced data, but it is fine for now.

We can see that the results are fine – all of the metrics are above 84%, and the model succeeds both on positive and negative results.

Those are the learning curves:



This was a sanity check and now we can go forward to hyper-parameters tuning.

Hyper-Parameters Tuning

There are many architectural choices and parameters to tune when training a neural net. We decided to split them to two groups: "basic" hyper-parameters and architectural parameters.

Basic hyper-parameters

- Learning rate
- Batch size
- Epochs number
- Activation function

Architectural parameters

- Number of layers
- Number of neurons in each layer

We will start by choosing the basic hyper-parameters with an arbitrary architecture, and once find a good configuration test the different architectural parameters with this configuration.

"Basic" Hyper-Parameters Tuning

The basic architecture I will try the parameters with is multi-layered perceptron with one hidden layer of 256 neurons. The loss function is binary cross-entropy.

Batch size

Usually choosing batch size has less to do with model quality and more with performance issues so we just tried some batch sizes and chose 64, which gave fine results in short training time. There is a connection between batch size and learning rate but once we set a fixed batch size, we can just choose the proper learning rate.

Batch size=64

Learning rate

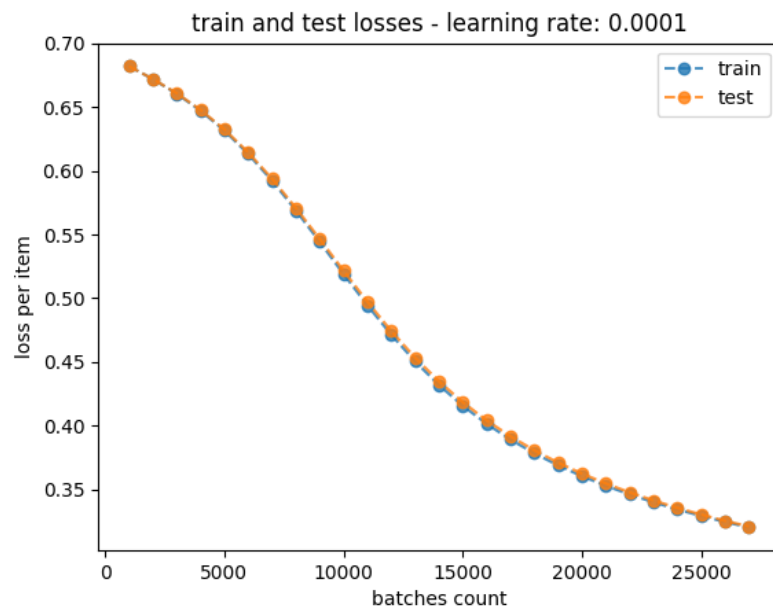
We trained the network with different learning rates, from 0.0001 to 0.1, for 40 epochs.

The full list:

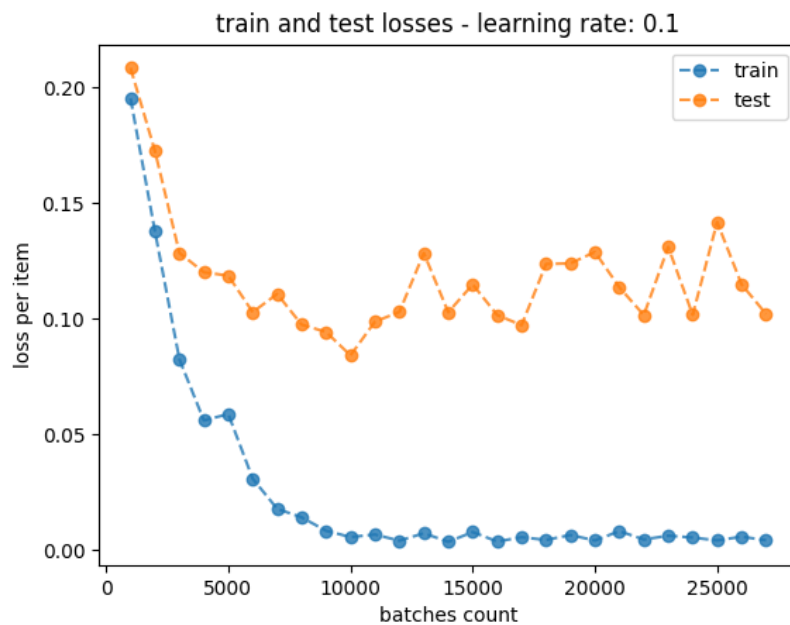
```
learning_rates = [10 ** (-4), 5 * 10 ** (-4), 10 ** (-3), 5 * 10 ** (-3), 10 ** (-2), 5 * 10 ** (-2), 10 ** (-1)]
```

Let's take a look in some learning curves to choose a good value:

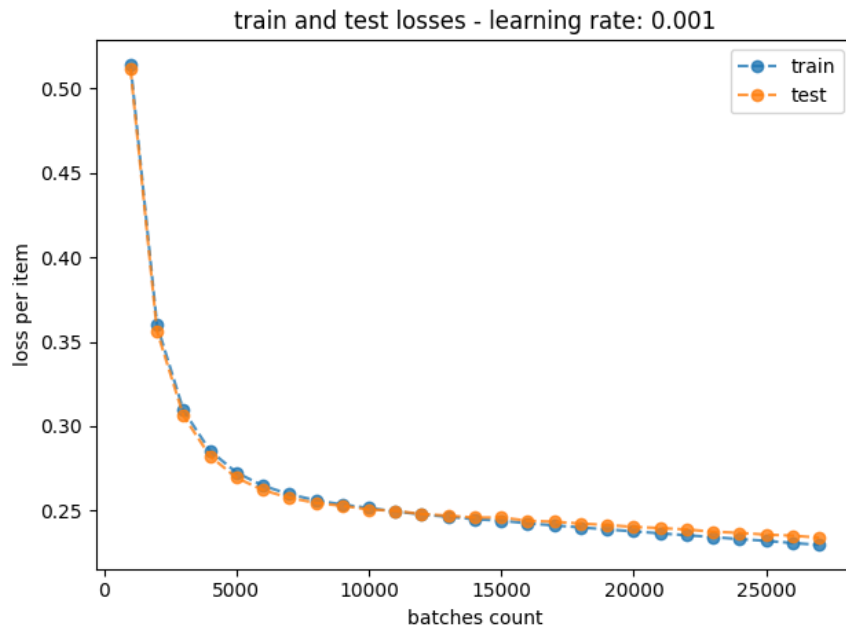
When we have a very small learning rate (0.0001), we have stable but slow learning:



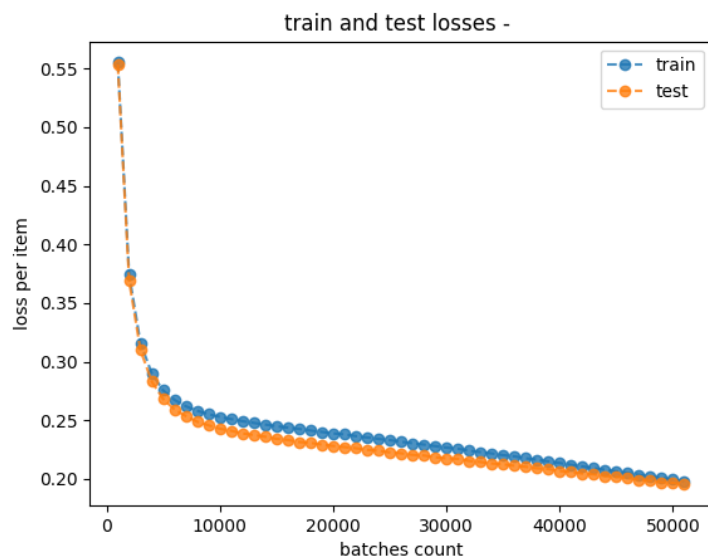
When we have a very big learning rate, (0.1) the train loss goes down but the test loss does not – overfit.



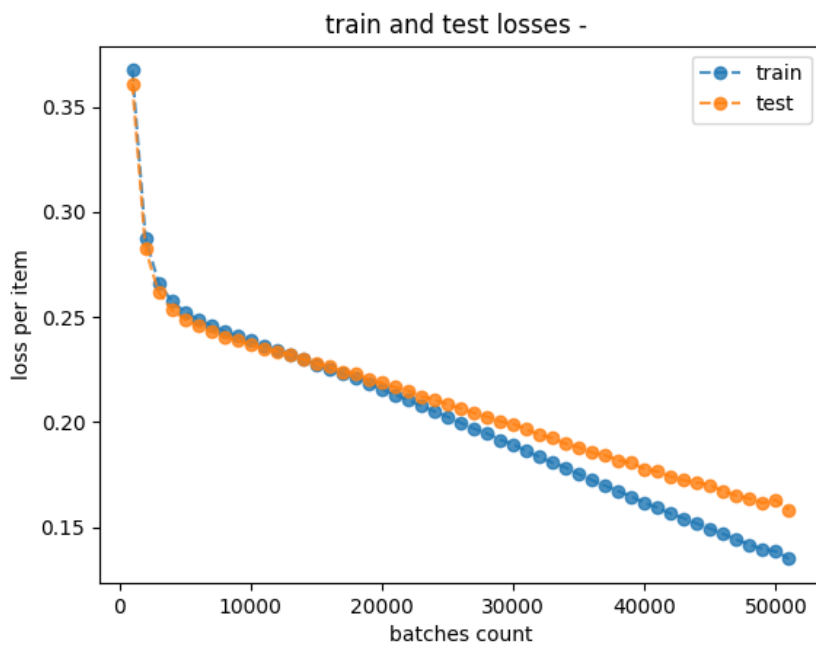
With a medium learning rate (0.001), we get fast convergence as wanted:



We can see that after 40 epochs there is still a downward trend of the test loss, so we will try to train a little more – for 75 epochs.

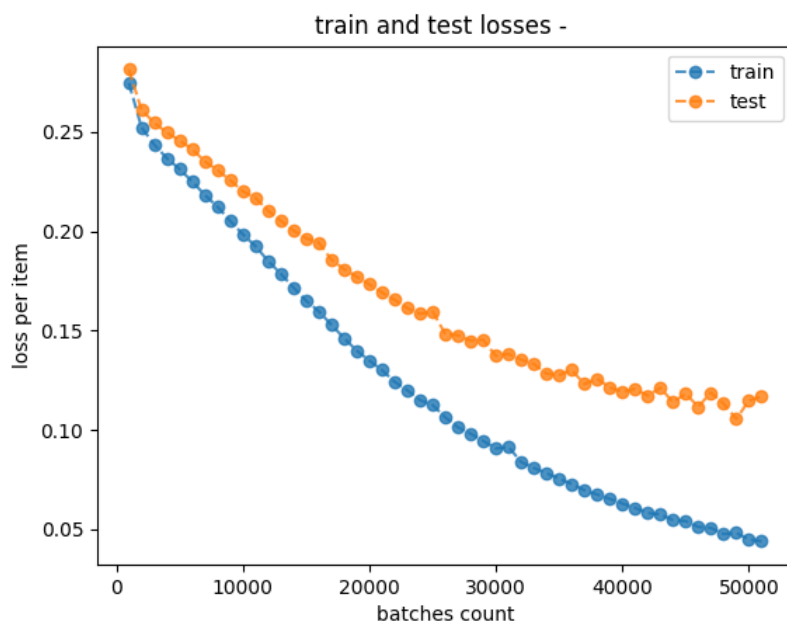


We also tried a slightly higher learning rate – 0.002:



Here there is a stronger fit – the train loss is better than the test loss, but also the test loss and the other metrics on the test set are better so we will go with this learning rate.

That's how the learning looks with learning rate of 0.005:



Here the test loss converges faster, so we will go with this learning rate.

Learning rate = 0.005

Epochs number

According to last learning curve 70 epochs seems like a good choice, about where the test loss converges.

Epochs number = 70

Activation function

We will try four different activation functions for the hidden layer:

- Sigmoid
- Relu
- LeakyRelu
- Tanh

We will compare them by test accuracy as instructed in the exercise definition.

The results:

activation function	accuracy
Sigmoid	0.903
Relu	0.967
LeakyReLU	0.970
Tanh	0.915

We can see that LeakyRelu gives us the best accuracy so we chose it.

Activation function = LeakyReLU

Architectures Comparison

After choosing the basic hyper-parameters, we searched the best network architecture. The parameters are number of layers and number of neurons in each layer, and of course, they are dependent so we searched for both of them together.

We tried:

Linear layer with sigmoid activation (kind of logistic regression) – **Zero** hidden layers

Shallow network – **One** hidden layer of sizes: 32, 64, 128, 256, 512, 1024, 2048

Deep network – **Two** hidden layers from same size each: 16, 32, 64, 128, 256, 512, 1024

Very deep network – **Five** hidden layers with 128 neurons each.

For now, I follow the instructions and look at the accuracy (and not recall or precision).

We looked at the learning curves as sanity check and they were just fine, similar to first part the curve I showed before, just as expected.

The results are:

# hidden layers	# neurons in each hidden layer	accuracy
0	-	90.0%
1	32	92.4%
1	64	93.2%
1	128	92.4%
1	256	91.6%
1	512	91.4%
1	1024	91.1%
1	2048	90.6%
2	16	95.8%
2	32	96.4%
2	64	97.2%
2	128	97.8%
2	256	97.4%
2	512	97.0%
2	1024	97.6%
5	128	97.4%

We can see that two hidden layers models are significantly better than one-layer models which are better than zero hidden layers models. We can also see that the number of neurons in each layer (width) does not make a dramatic difference in the accuracy.

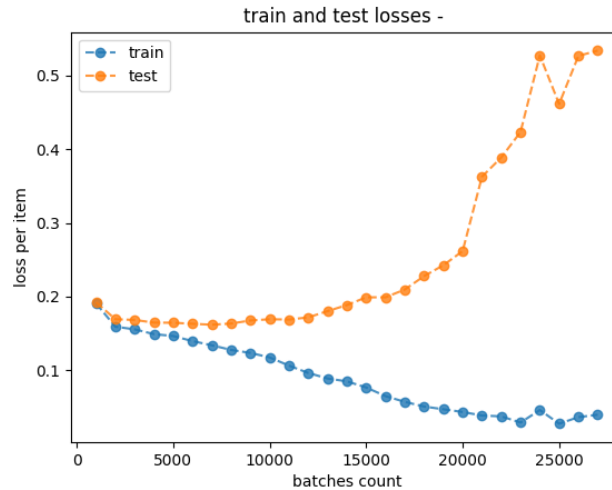
Therefore, the chosen architecture is:

Number of hidden layers = 2

Neurons in each layer = 128

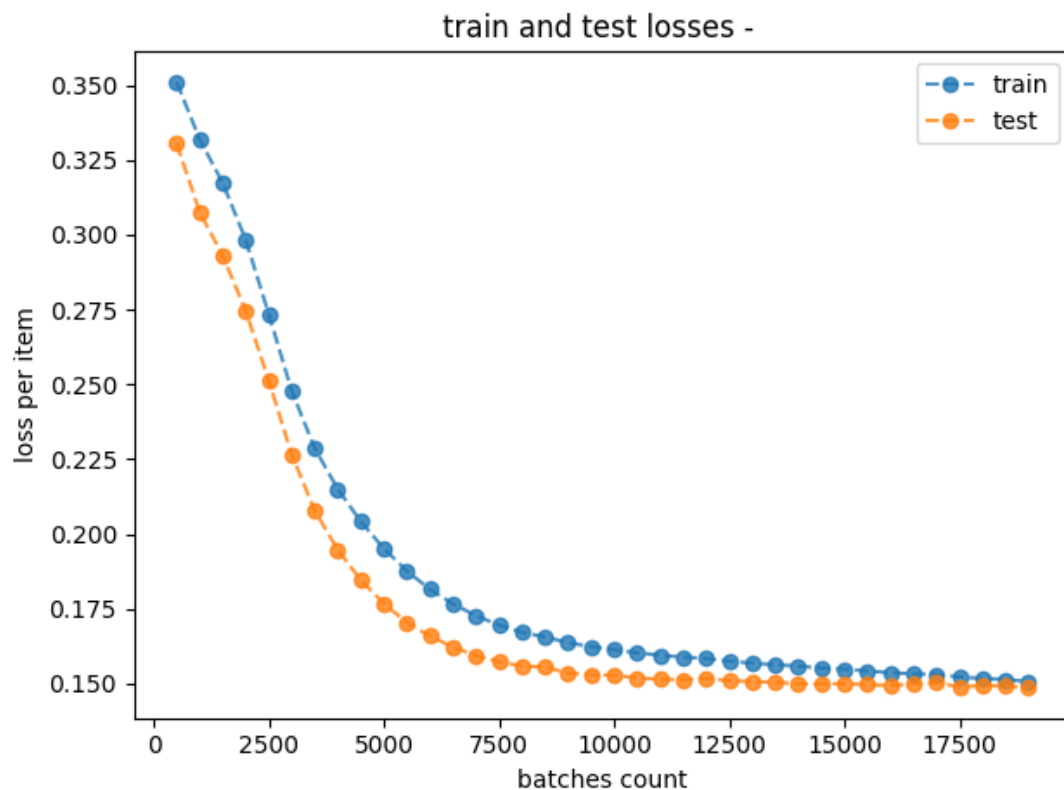
Adjusting learning rate and epochs number

When we look in the learning curves, we see that we get overfit:



That's because we tuned each parameter as it is independent of the others, and this assumption is not true. The best way is to try some kind of grid search, but this is really expensive in computation time. For now, we will just re-tune the learning rate and epochs number to prevent overfit.

We found out that training the network with smaller learning (0.001) rate and less epochs (50) gives good learning:



We even can train for less epochs and get close results so we will go with 40 epochs.

To sum up, the chosen params are:

- Batch size=64
- Learning rate = 0.001
- Epochs number = 40
- Activation function = LeakyReLU
- Number of hidden layers = 2
- Neurons in each layer = 128

Dealing With Imbalanced Data

In the previous sections we dealt with the imbalanced data with oversampling. In this section we will improve our pipeline and measurement and focus specifically on this problem.

Evaluation

The first step in dealing with imbalanced data situation, just like most machine learning situations, is defining a proper metric. At first, we did not do it (you can see it in the appendix) and did not understand out model is doing a very bad job.

First of all, we will oversample our data better – we will oversample the negative samples only in the train set and leave the test set as is. Next, we will separate the single accuracy metric to two – one only for the positive samples, and one only for the negative ones.

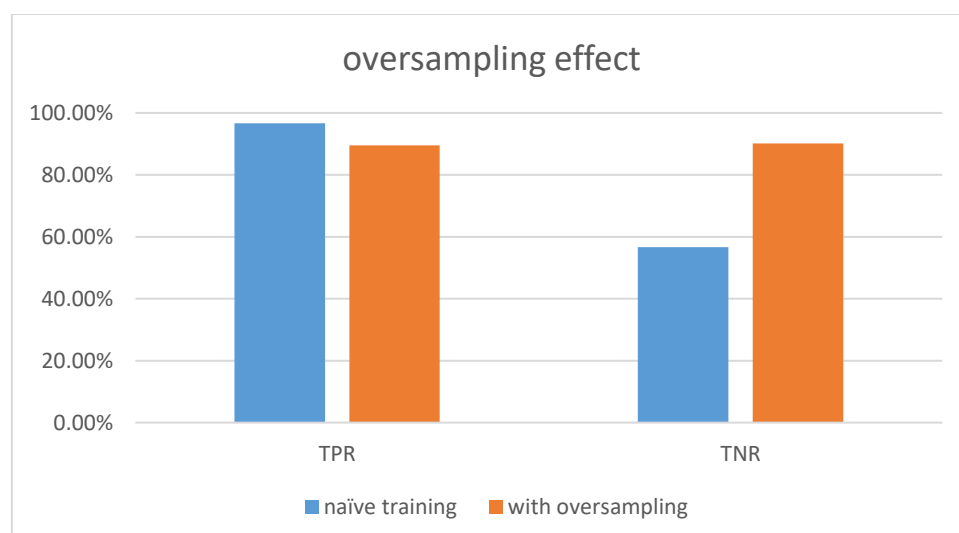
$$\text{True Positive Ratio} = \text{Recall} = \frac{TP}{P}$$

$$\text{True Negative Ratio} = \frac{TN}{N}$$

Those metrics are agnostic to the positive-negative ratio in the data.

Oversampling

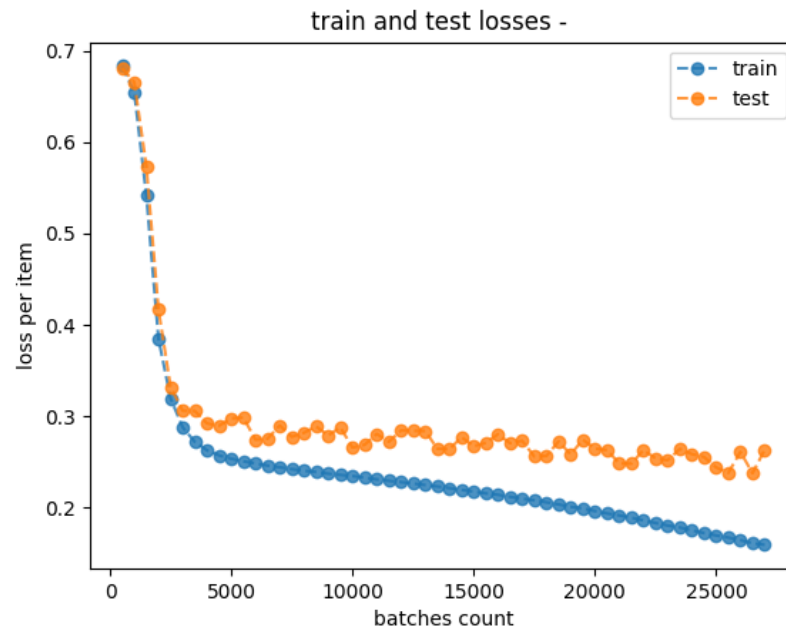
We used a simple oversample technique. We calculate the ratio between the positive samples and the negative samples in the train set (it is about 8.23, the train-test split is random so it changes a bit each run), and then duplicate each negative sample for $\text{int}(\text{ratio})$ times (usually 8). That way the loss function will be affected also by the negative samples, and the positive samples won't take over so easily.



In the plot above we can see comparison of two trainings – one with oversampling and one without. We can see that the oversampling improves the prediction of the negative samples dramatically (about 4 times less mistakes), and hurt the TPR a little bit. Overall, it's worth it.

Final training

We chose all of the hyper-params and decided how to deal with imbalanced data. Now it is time to show the final learning curve and metrics.



The naïve reader may think that the test loss does not go down fast enough, but the wise reader will remember that the test loss does not say a lot because it is not oversampled, so it gives far much more weight to performance on positive samples.

	test	Train – oversampled
TPR	0.89	0.90
TNR	0.86	0.99
accuracy	0.89	0.94

The metrics are fine – pretty high results on the test set. An interesting observation is that we got overfit on the negative samples because they were aggressively over sampled. Maybe it was better to try a lower over sampling factor or try other methods (under sampling, oversampling + dropout, ect.).

Inference on Spike Protein

We used our trained model to identify sequences in the Spike protein. The protein is from length 1273, and therefore we have 1265 peptides ($1273 - 9 + 1$).

Most of the predictions were 1171 positive predictions out of 1265 predictions, which is about 92.6%.

The most detectable peptides are:

peptide	score
VYYPDKVFR	1.0
VYADSFVIR	1.0
RVDFCGKGY	1.0
ERDISTEY	~1.0 (9.999...)
KTSVDCTMY	~1.0 (9.999...)

Of course, we first notified the CDC 😊

Optimize Input

In the last section we optimized the input with respect to the pretrained model. It turns out that our model gives high scores to random values, so the loss was usually pretty low from the beginning.

To convert the arbitrary tensor to peptide we took the argmax of each 20-sized segment of the tensor and took the matching letter of the peptide. (kind of inverse one-hot encoding)

We got different results each run. For example:

ADEVEQSWV

Because the model tends to give high scores to random values, we tried to change the mean and the std of the normal distribution, but it didn't help much. So, in order to make sure the process works properly we tried to do the opposite – find the least detectable peptide. That way we started with high loss and watched it goes down while the input changes. A very undetectable peptide is:

NVEAIHYSA

Appendix – Problematic Experiments

In this appendix we show the first time we wrote this report – with bad models that only output 1...

First Training

For the first training, we tried feed forward network with one hidden layer with 128 neurons with relu activation. We trained the network for 5 epochs. We used batch size of one for simplicity.

I used binary cross entropy loss and learning rate of 0.001.

Results on the test set:

Accuracy: 0.917

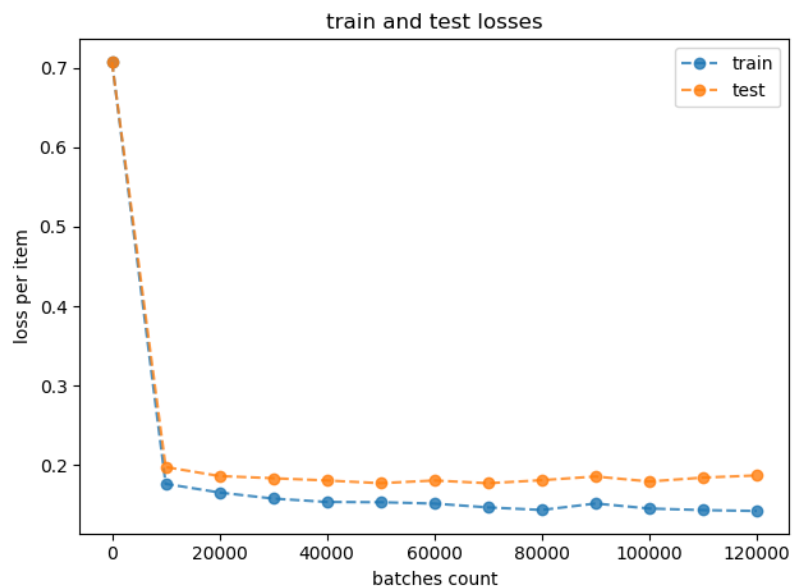
Recall: 0.968

Precision: 0.941

f1: 0.954

We can see that the results are fine – all of the metrics I checked are above 90%. Especially it is good to see that both the recall and the precision are pretty high even though the data set is imbalanced.

Those are the learning curves:



This was sanity check and now we can go forward to hyper-parameters tuning.

Hyper-Parameters Tuning

There are many architectural choices and parameters to tune when training a neural net. We decided to split them to two groups: "basic" hyper-parameters and architectural parameters.

Basic hyper-parameters

- Learning rate
- Batch size
- Epochs number
- Activation function

Architectural parameters

- Number of layers
- Number of neurons in each layer

We will start by choosing the basic hyper-parameters with an arbitrary architecture, and once find a good configuration test the different architectural parameters with this configuration.

"Basic" Hyper-Parameters Tuning

The basic architecture I will try the parameters with is multi-layered perceptron with one hidden layer of 256 neurons. The loss function is binary cross-entropy.

Important note: The results showed here are not always statistically significant, but because this is an exercise, we focused on exploration and trying a lot of parameters and not on being sure our parameters are the best.

Batch size

Usually choosing batch size has less to do with model quality and more with performance issues so we just tried some batch sizes and chose 64, which gave fine results in short training time. There is a connection between batch size and learning rate but once we set a fixed batch size, we can just choose the proper learning rate.

Batch size=64

Learning rate

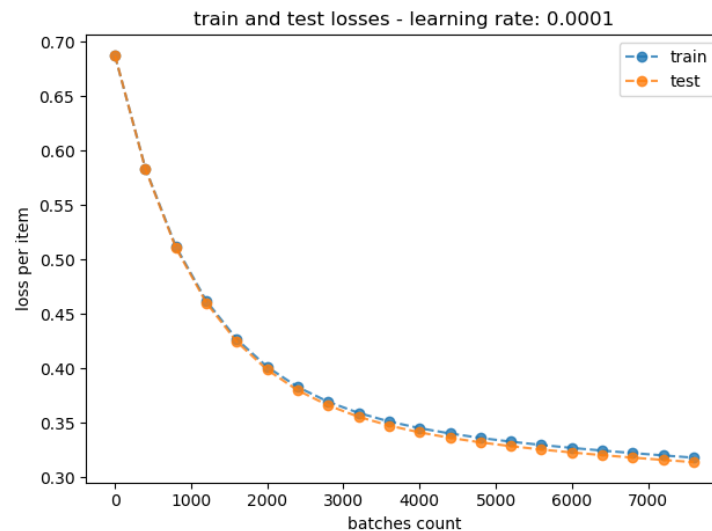
We trained the network with different learning rates, from 0.0001 to 0.1.

The full list:

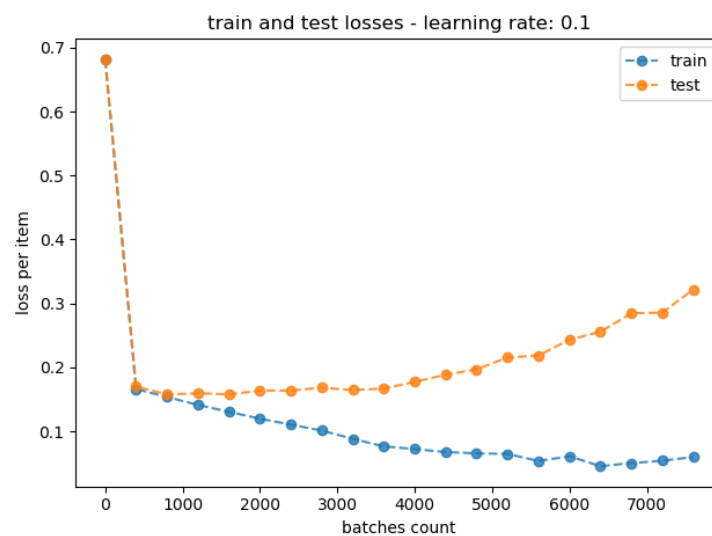
`learning_rates = [10 ** (-4), 5 * 10 ** (-4), 10 ** (-3), 5 * 10 ** (-3), 10 ** (-2), 5 * 10 ** (-2), 10 ** (-1)]`

Let's take a look in some learning curves to choose a good value:

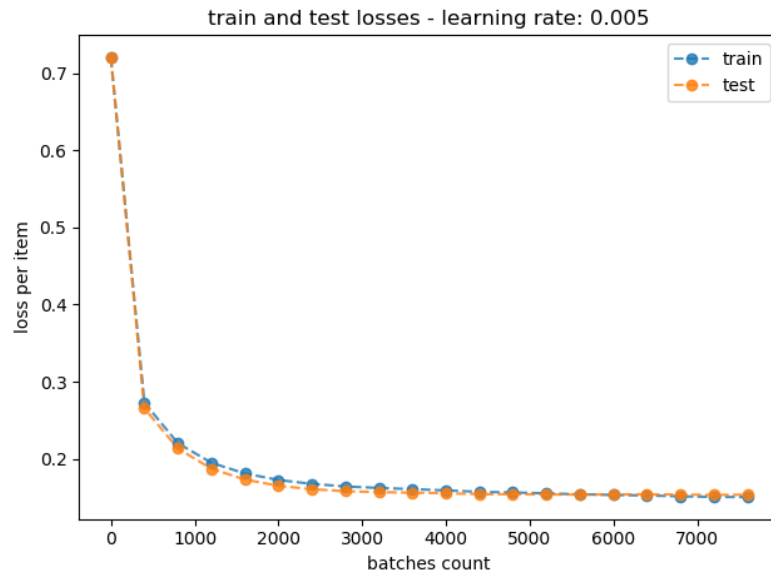
When we have a very small learning rate (0.0001), we have stable but slow learning:



When we have a very big learning rate, (0.1) the train loss goes down but the test loss goes up – overfit.



With a medium learning rate (0.005), we get fast convergence as wanted:

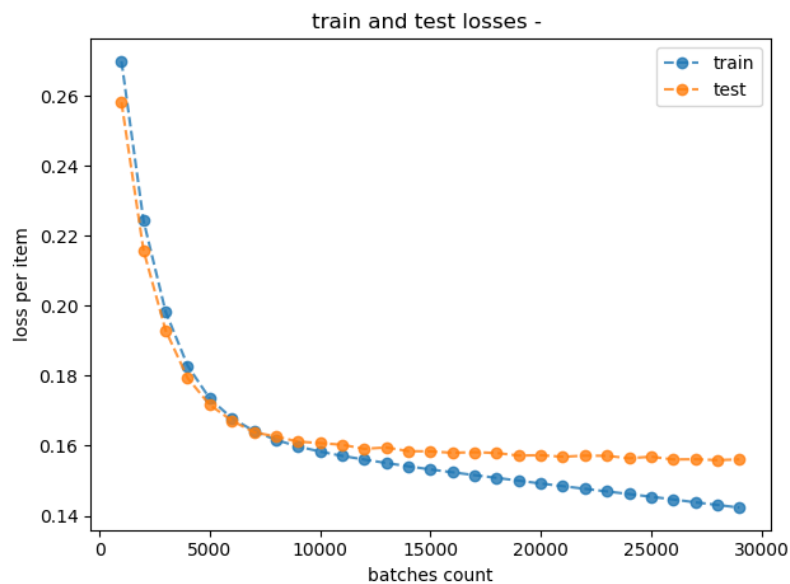


We will go on the safe side and choose a slightly smaller learning rate of 0.002 to avoid problems later.

Learning rate = 0.002

Epochs number

We trained the network for 75 epochs with the chosen params and got this learning curve:



The test loss does not get much better from the 10,000 batch, which corresponds to about 25 epochs. From this point we get to a kind of overfit – the test loss does not get worse but the train loss goes to zero without real improvement of the model.

Epochs number = 25

Activation function

We will try four different activation functions for the hidden layer:

- Sigmoid
- Relu
- LeakyRelu
- Tanh

We will compare them by test accuracy as instructed in the exercise definition.

The results:

activation function	accuracy
Sigmoid	0.907
Relu	0.930
LeakyRelu	0.884
Tanh	0.907

We can see that Relu gives us the best accuracy so we chose it.

Activation function = Relu

Architectures Comparison

After choosing the basic hyper-parameters, we searched the best network architecture. The parameters are number of layers and number of neurons in each layer, and of course, they are dependent so we searched for both of them together.

We tried:

Linear layer with sigmoid activation (kind of logistic regression) – **Zero** hidden layers

Shallow network – **One** hidden layer of sizes: 32, 64, 128, 256, 512, 1024, 2048

Deep network – **Two** hidden layers from same size each: 16, 32, 64, 128, 256, 512, 1024

Very deep network – **Five** hidden layers with 128 neurons each.

For now, I follow the instructions and look at the accuracy (and not recall or precision).

We looked at the learning curves as sanity check and they were just fine, similar to first part the curve I showed before, just as expected.

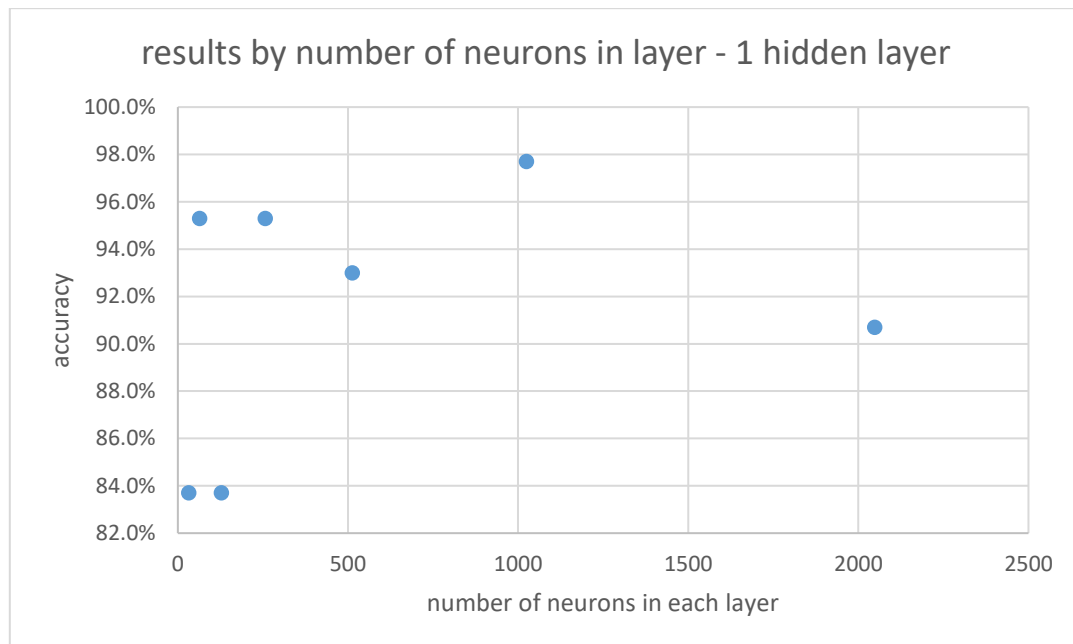
The results are:

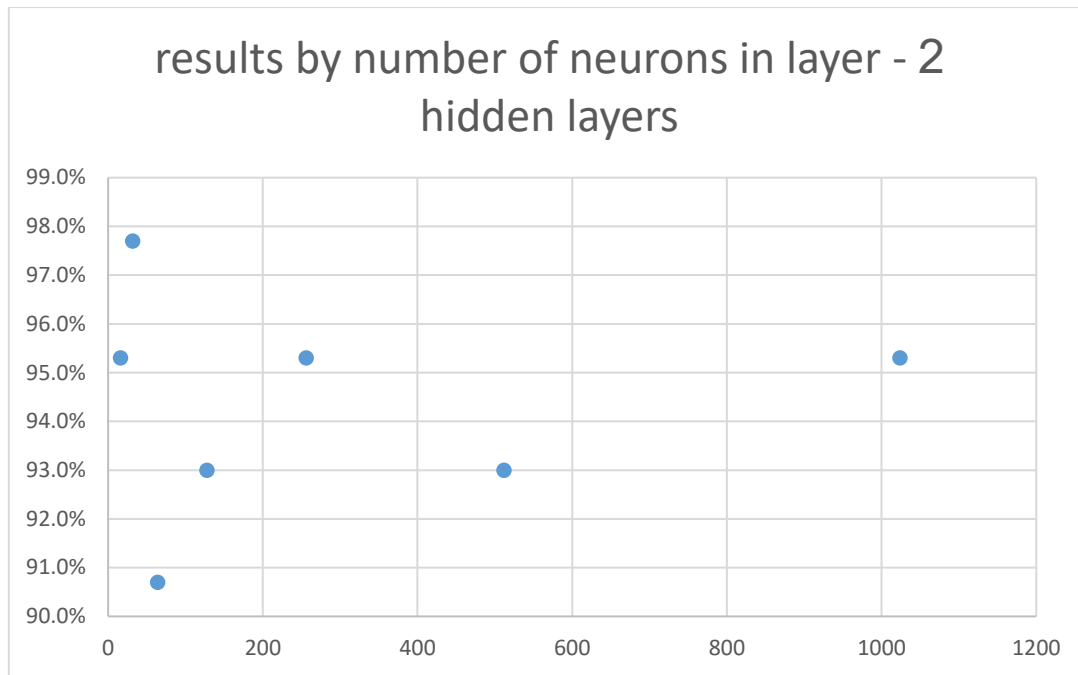
# hidden layers	# neurons in each hidden layer	accuracy
-----------------	--------------------------------	----------

0	-	93.0%
1	32	83.7%
1	64	95.3%
1	128	83.7%
1	256	95.3%
1	512	93.0%
1	1024	97.7%
1	2048	90.7%
2	16	95.3%
2	32	97.7%
2	64	90.7%
2	128	93.0%
2	256	95.3%
2	512	93.0%
2	1024	95.3%
5	128	97.7%

We can see big differences – the accuracies goes from ~84% to ~98% (8 times less errors). In general we see that the deeper the net – the better the accuracy.

For each positive number of hidden layers I plotted the results as function of the number of neurons in each layer.





We did not connect the dots with a line because we think that what we see here is mainly statistical noise. We think that because there is no clear trend and because when we notices that when we run the same training several times we sometimes get different results. However, when we look at the average accuracy for each network depth we can see a trend – the deeper, the better.

Therefore, we decided to check this trend. We trained 3 different models – three MLPs with 256 neurons in each layer, but with different depths: 2, 5, 10. Each model was trained three times to make the results more reliable.

The results:

Network depth	Mean accuarcy
2	93.8%
5	96.9%
10	85.3%

Therefore the chosen architecture is:

Number of layers = 5

Neurons in each layer = 256

מבוא ללמידה עמוקה תרגיל 1

עמית קינן 2088296426 אלדן חודורוב 201335965

3 בנובמבר 2021

1 בהרכבה של פונקציות

1.1 הרכבה של פונקציות ליניאריות

הראו כי הרכבה של פונקציות ליניאריות הינה לינארית, נסתכל במקרה הוקטורי כאשר פונקציה לינארית מוגדרת בתור

$$\begin{aligned}f(x) &= Ax \wedge g(y) = By \\x &\in \mathbb{R}^n, y \in \mathbb{R}^m, A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{k \times m} \\g(f(x)) &= g(Ax) = B(Ax) = (BA)x\end{aligned}$$

על כן אם נסמן את $BA = C \in \mathbb{R}^{k \times n}$ נקבל כי $g(f(x)) \in \mathbb{R}^{k \times n}$ הינה לינארית

1.2 הרכבה של פונקציות אפיניות הינה אפינית

באופן דיי ישיר מהסעיף הקודם נגדיר את $a \in \mathbb{R}^m$ וכן את $b \in \mathbb{R}^k$ ונגדיר מחדש את $f \wedge b$ באופן הבא

$$f(x) = Ax + a, g(y) = By + b$$

ועל ידי הרכבה נקבל

$$\begin{aligned}g(f(x)) &= g(Ax + a) = B(Ax + a) + b = (BA)x + \underbrace{B \cdot a}_{\in \mathbb{R}^k} + b \\&= \underbrace{(BA)}_{=C}x + \underbrace{B \cdot a + b}_{=c} \\&= Cx + c\end{aligned}$$

אזי פונקציה אפינית גם כן

2 calculus behind the Gradient Descent method

2.1 מהו תנאי העצירה לסכימה האיטרטיבית

$$\theta^{n+1} = \theta^n - \alpha \nabla f_{\theta^n}(x)$$

באופן דיי ישיר תנאי העצירה הינו ש $\nabla f_{\theta^n}(x) = 0$ ואז נקבל $\theta^{n+1} = \theta^n$

2.2 הראו עם התור טיילור מסדר שני של הפונקציה מה התנאי הנדרש להסיק את הסיווג של נקודה סטציונרית כמינימום\קסימום

$$f(x + dx) = f(x) + \nabla f(x) \cdot dx + dx^T \cdot H(x) \cdot dx + O(\|dx\|^3)$$

$$H_{i,j}(x) = \frac{d^2 f}{dx_i dx_j}(x)$$

נראה כי התנאי הנדרש לסיווג הנקודה כמינימום\מקסימום לוקאלי הוא ש $H(x)$ היא פונקציה חיובית או שלילית בהחלט בהתאמה וכן כי $\nabla f(x) = 0$ כלומר נניח כעת כי $H(x)$ חיובית בהחלט והנגזרת מתאפסת ו- נקבל כי

$$f(x+dx) = f(x) + dx^T \cdot H(x) \cdot dx + O(\|dx\|^3)$$

כעת נרצה להראות ש $f(x)$ היא נקודת מקסימום לוקאלי כלומר לכל dx קטן מספיק מתקיים:

$$f(x+dx) \leq f(x)$$

מהנחה כי $H(x) > 0$ נקבל כי

$$dx^T \cdot H(x) \cdot dx \geq 0$$

על כן נרצה להראות כי

$$dx^T \cdot H(x) \cdot dx \geq O(\|dx\|^3)$$

וכן ידוע כי עבור dx קטן מספיק

$$O(\|dx\|^3) \leq O(\|dx\|^2) \approx dx^T \cdot H(x) \cdot dx$$

על כן

$$dx^T \cdot H(x) \cdot dx + O(\|dx\|^3) \geq 0$$

ומכך נסיק

$$f(x+dx) = f(x) + \underbrace{\nabla f(x)}_{=0} \cdot dx + \underbrace{dx^T \cdot H(x) \cdot dx}_{\geq 0} + O(\|dx\|^3) \geq f(x)$$

כלומר $f(x)$ הינה נקודת מינימום לוקאלי, ובאופן סימטרי ניתן להוכיח את הדבר על נקודת מקסימום לוקאלי.

3 תיצור פונקציית Loss של רשת שמנבאת זווית בין 0-360 ושיודעת לקחת בחשבון מרחק זוויתי - כלומר המרחק בין 2 ל-360 הינו 2 ולא 358

```
def degrees_loss(y_pred,y_true):
    abs_diff = abs(y_true-y_pred)
    return min(abs_diff, 360-abs_diff)
```

ניתן לראות כי המערכת מחשבת את המרחק בין הפרדיקציה ל- GT ומחזירה מרחק הזוויתי בין התוצאה לפרדיקציה. אציין כי יש עוד אפשרויות, למשל העברה לרדיאנים וחשוב המרחק קוסינוס מבין התוצאות, משהו קצת יותר דומה ל-cosine similarity

4 גזירת פונקציות עם כלל השרשרת

$$\frac{d}{dx} f(x + y, 2x, z) \quad 4.1$$

בשביל לגזור את הביטוי נסמן את התתי פונקציות של f באופן הבא

$$g(x) = x + y$$

$$h(x) = 2x$$

$$s(x) = z$$

על כן אם נפעיל את כלל השרשרת נקבל

$$\begin{aligned} \frac{df(x + y, 2x, z)}{dx} &= \frac{df(g(x), h(x), s(x))}{dx} = \\ &= \frac{dg(x)}{dx} \cdot \frac{df(g(x), h(x), s(x))}{dg(x)} + \frac{dh(x)}{dx} \cdot \frac{df(g(x), h(x), s(x))}{dh(x)} + \frac{ds(x)}{dx} \cdot \frac{df(g(x), h(x), s(x))}{ds(x)} \\ &= 1 \cdot \frac{df(g(x), h(x), s(x))}{dg(x)} + 2 \cdot \frac{df(g(x), h(x), s(x))}{dh(x)} + 0 \cdot \frac{df(g(x), h(x), s(x))}{ds(x)} \\ &= 1 \cdot \frac{df(g(x), h(x), s(x))}{dg(x)} + 2 \cdot \frac{df(g(x), h(x), s(x))}{dh(x)} \end{aligned}$$

$$f(x) = f_1(f_2(\dots f_n(x))) \quad 4.2$$

$$\frac{df}{dx} = \frac{df_1}{df_2} \cdot \frac{df_2}{df_3} \dots \frac{df_n}{dx}$$

כלומר נקודת הגזירה של $\frac{df_i}{df_{i+1}}$ תגזר לפי הנקודה $f_{i+1}(f_{i+2}(\dots f_n(x)))$ לכל i למעט f_n שתגזר לפי הנקודה x

$$f_1(x, f_2(x, f_3(\dots f_{n-1}(x, f_n)))) \quad 4.3$$

בשביל לגזור את הביטוי הנ"ל נתחיל קודם להסתכל על ההנגזרת של $f_{n-1}(x, f_n(x))$

$$\frac{df_{n-1}(x, f_n(x))}{dx} = \frac{df_{n-1}(x)}{dx(1, 0)} + \frac{df_{n-1}(f_n(x))}{df_n(x)(0, 1)} \cdot \frac{df_n(x)}{dx}$$

אם נסתכל עוד צעד אחורה על $f_{n-2}(x, f_{n-1}(x, f_n))$ נקבל כי

$$\frac{df_{n-2}(x, f_{n-1}(x, f_n))}{dx} = \frac{df_{n-2}(x)}{dx(1, 0)} + \frac{df_{n-2}(f_{n-1}(x))}{df_{n-1}(x)(0, 1)} \cdot \frac{df_{n-1}(x, f_n(x))}{dx}$$

כאשר הביטוי האחרון היא הנגזרת של f_{n-1} שהראינו. כלומר אם נמשיך בצורה רקורסיבית נגיע חזרה ל f_1 כלומר

$$\frac{df_1(x, f_2(x, f_3(\dots f_{n-1}(x, f_n))))}{dx} = \frac{df_1(x)}{dx(1, 0)} + \frac{df_1(f_2(x, f_3(\dots)))}{df_2(x)(0, 1)} \cdot \frac{df_2(x, f_3(x, f_4(\dots f_n(x))))}{dx}$$

$$f(x + g(x + h(x))) \quad 4.4$$

נגדיר את

נגזור את הפונקציה לפי הנגזרות הפנימיות שלה ונקבל

$$\frac{df^*}{dx} = \frac{df(x + g(x + h(x)))}{d(x + g(x + h(x)))} \cdot (1 + \frac{dg(x + h(x))}{d(x + h(x))}) \cdot (1 + \frac{dh(x)}{dx})$$

5 הוכח כי $D_{kl}(P||Q) \geq 0$ היא לא שלילית כלומר

להוכיח את הטענה נשתמש בכך שאנו יודעים כי $\sum_i^n p_i \log(\frac{p_i}{q_i}) \geq (\sum_i p_i) \cdot \log \frac{\sum_i p_i}{\sum_i q_i}$ אם כן, ניתן לראות כי,

$$D_{kl}(P||Q) = \sum_i^n p_i \log(\frac{p_i}{q_i}) \geq (\sum_i p_i) \cdot \log \frac{\sum_i p_i}{\sum_i q_i}$$

$$1 \cdot \log \frac{1}{1} = 1 \cdot 0 = 0$$

כאשר $p = q$ יש שיוויון במעבר כלומר $D_{kl}(P||Q) = 0$

6 הוכח כי $D_{kl}(P||Q)$ היא קמורה

נדרש להוכיח כי

$$D_{kl}(\lambda p_1 + (1 - \lambda)p_2 || \lambda q_1 + (1 - \lambda)q_2) \leq \lambda D_{kl}(p_1 || q_1) + (1 - \lambda)D_{kl}(p_2 || q_2)$$

כאשר

$$(p_1, p_2) \wedge (q_1, q_2)$$

הם זוגות של הסתברויות לא שליליות. וכן $\lambda \in [0, 1]$ נשתמש באותה טענה מהסעיף הקודם

$$\sum_i^n p_i \log(\frac{p_i}{q_i}) \geq (\sum_i p_i) \cdot \log \frac{\sum_i p_i}{\sum_i q_i}$$

אם כן,

$$D_{kl}(\lambda p_1 + (1 - \lambda)p_2 || \lambda q_1 + (1 - \lambda)q_2)$$

$$= \sum (\lambda p_1 + (1 - \lambda)p_2) \cdot \log \left(\frac{\lambda p_1 + (1 - \lambda)p_2}{\lambda q_1 + (1 - \lambda)q_2} \right)$$

$$\leq \sum \left(\lambda p_1 \log \left(\frac{\lambda p_1}{\lambda q_1} \right) + (1 - \lambda)p_2 \log \left(\frac{(1 - \lambda)p_2}{(1 - \lambda)q_2} \right) \right)$$

$$= \sum \lambda p_1 \log \left(\frac{\lambda p_1}{\lambda q_1} \right) + \sum (1 - \lambda)p_2 \log \left(\frac{(1 - \lambda)p_2}{(1 - \lambda)q_2} \right)$$

$$\lambda D_{kl}(p_1 || q_1) + (1 - \lambda)D_{kl}(p_2 || q_2)$$

כנדרש

7 הראה כי הטענות של Cybenko and Hornik נתנות להחלה על פונקציה Relu

נראה כי ניתן לייצר פונקציה σ שהינה סכום של פונקציות $Relu$ עם הזזה וניפוח שינה רציפה ומונוטונית המקיימת כי $\lim_{x \rightarrow -\infty} \sigma(x) = 0$ וכן $\lim_{x \rightarrow \infty} \sigma(x) = 1$ כלומר הפונקציה תקיים את התנאים הנדרשים ועל כן חלים עליה הטענות. נגדיר את הפונקציה

$$\sigma(x) = Relu(x + 0.5) - Relu(x - 0.5)$$

ניתן לראות כי הפונקציה מקיימת

$$\sigma(x) = \begin{cases} 0 & x < -0.5 \\ x + 0.5 & x \in [-0.5, 0.5] \\ 1 & x > 0.5 \end{cases}$$

על כן ניתן לבטא באמצעות סכום של פונקציות σ

$$f = \sum_i \alpha_i \sigma'(\omega_i x + b_i)$$

כלומר היא הינה צפופה ב- $C[0, 1]$ תחת הנורמת סופרימום

8 הרחב את הבניה שהראינו בכיתה לביטוי רשת רדודה עם רשת עמוקה עם $O(n)$ נוירונים

8.1 תבטא רשת עם שכבה 1 על ידי רשת עמוקה עם $O(n)$ נוירונים לפונקציה ה- Sigend

כלומר נקבל כי

$$f(x) = \sum_i \alpha_i \sigma(\omega_i x + b_i)$$

כך ש α_i יכול להיות שלילי או חיובי. ולא רק חיובי כמו שהראינו בשיעור. בשיעור השתמשנו בסימון של h_i לסמן את הזרמי חיבור של הנוירונים ברשת. נמשיך את הסימון הזה. אז במטרה להשיג את הרשת המתבקשת, נשנה את הזרוע של h_1 שעקבה וסכמה את הערכים, כך שזרוע היא רק תסכום את הנוירונים עם $\alpha_i > 0$ ונוסיף לרשת זרוע נוספת שתקרא h_4 והיא תסכום את כל הנוירונים עם $\alpha_i < 0$ כאשר h_3, h_2 ישארו ללא שינוי. h_4 תפעל על ידי הכפלה ב-1 של ה- σ המתאימה ושל האינפוט מהשכבה הקודמת $h_4 = -\sigma_i(-1 \cdot h_2)$. כאשר אינפוט בסימן הפוך יתאפסו. נזכור שבסוף הרשת נכפול שוב נסכום את התוצאות שקיבלנו ב- h_1 ו- h_4 ונקבל את הנדרש. כי בפועל עבור h_1 נקבל מ $\alpha \geq 0$ מ h_2 את הערך 0 בעוד שב h_4 נקבל את ההפך. כלומר קיבלנו יכולת לבטא את הרשת השטוחה על ידי רשת עם n שכבות כאשר בכל שכבה יש $O(1)$ נוירונים - (בפועל עד 4 נוירונים).