

Deep learning: Neural Network from scratch
creating a neural network library to solve the mnist dataset

Adrien PELFRESNE - Alexis VAPAILLE

May 30, 2024

1 Abstract

We have created a deep neural network "library" from scratch, in rust, without the use of any machine learning or deep learning library.

The goal of this project was to solve the mnist dataset of handwritten digits, and to provide performances comparison between two multi-layer-perceptron, with and without convolution.

Our library let you create your own sequential neural network, with a simple and declarative api and is widely open to extension. This project also include a drawing mode, when you can draw your own digit and see what the model guess.

2 Motivation

Our motivation with this *from scratch* approach was to deepen our understanding of neural networks. Because high level library like [keras](#) and the incredible level of abstraction that come with it, let us sometimes forgot the inner machinery. nevertheless, this underlying work that those library are doing is passionating, and we through it was a shame to be a simple user of those interface. The only (methemathical) helper library that we have used is [ndarray](#) which provides an n-dimensional container for general elements and for numerics, and efficients matrix operations on the gpu.

You can find our project on [github](#), we widely encourage reader to check the source code, as in this report, we won't go into code specific details and focus more on the abstract explanation.

The project is split into : the actual neural network lib, the mnist exemple, which is using the library, as well a our report sources and the gui.

3 Sequential neural network

The multi-layer perceptron is known as a *Sequential* neural network, meaning that the **output** Y_N of the layer L_N is the **input** X_{N+1} of the next layer L_{N+1} .

$$Y_N = X_{N+1}$$

and conversely

$$X_N = Y_{N-1}$$

with subscript N indentifying layers in sequential order. from the input layer, when we pass our data, through the output layer, data flows through the network and through different kind of layers.

3.1 General layer interface

We have defined a general *Layer* interface, to be implemented by all the differents layer we gonna use.

This interface has two main method : **feed_forward** and **propagate_backward**, The *feed forward* method take the input vector X as parameter and return the output vector Y of the layer. and the *propagate backward* method take as parameters the layer output gradient $\frac{\partial C}{\partial Y}$

and return the layer input gradient $\frac{\partial C}{\partial X}$. We will cover the use of these function later, alongside the theoretical explanation of these concepts.

3.2 Dense layer

In a dense layer, also knowed as *Linear layer* each input neuron is connected to every output neurone, each connexion is called a *weight*, w_{ij} is the weight connecting neurone x_i (from the input) to neurone y_j (to the output).

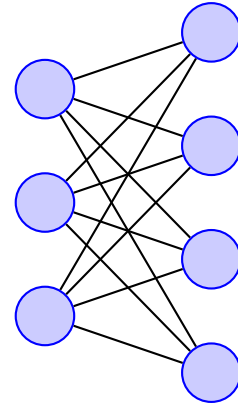


Figure 1: a Dense Layer

The output value y_j of a dense layer neuron is given by

$$y_j = \sum_{i=1}^n x_i w_{i,j} + b_j$$

with b_j the bias term for the output neuron j . This calculation can be expressed using matrices operation on whole vectors.

$$Y = XW + B$$

3.3 Activation layer

In a lot of deep neural network content, an activation function is applied inside the dense layer to produce the y output value from the weighted sum z . We choose, in our implementation, to enhance the seperation of concerns, with a decoupled version.

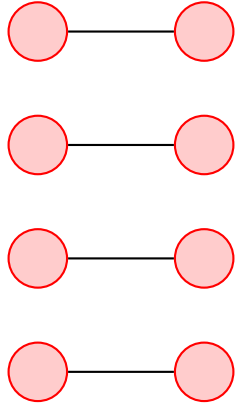


Figure 2: Element-wise Connections Between Two Layers

The Activation layer take a vector of size n as input and produce a same length vector as a output, with the activation function applied to each neuron. we must note that some activations function transform a vector by mapping individual component through the function, like the ReLU

$$ReLU(z) = \max(0, z)$$

while other activation functions such as *softmax* take vector z of K real number, and is used to convert a vector of K real number into a probability distribution of K possible outcomes.

$$\sigma : \mathbb{R}^K \rightarrow (0, 1)^K, K \geq 1$$

is computed by taking a vector $z = (z_1, \dots, z_k) \in \mathbb{R}^K$ and compute each component of a vector $\sigma(\mathbf{z}) \in (0, 1)^K$

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

we have implemented a variety of activations function, to try different options and see what works the best with the mnist dataset.

4 Cost function

To mesure how well our neural network perform, we use a cost function, which output a real number, by comparing the neural network *output*, and the effective *observed* value. A simple cost function is the mean squared error, which is defined as :

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

this cost function is great for regression neural network when we want to predict a continuous value, but for a classification task, like the mnist data set (when we need to classifie each image) we can use the cross entropy cost function, defined as :

$$\text{Cross Entropy} = - \sum_{c=1}^C y_c \log(\hat{y}_c)$$

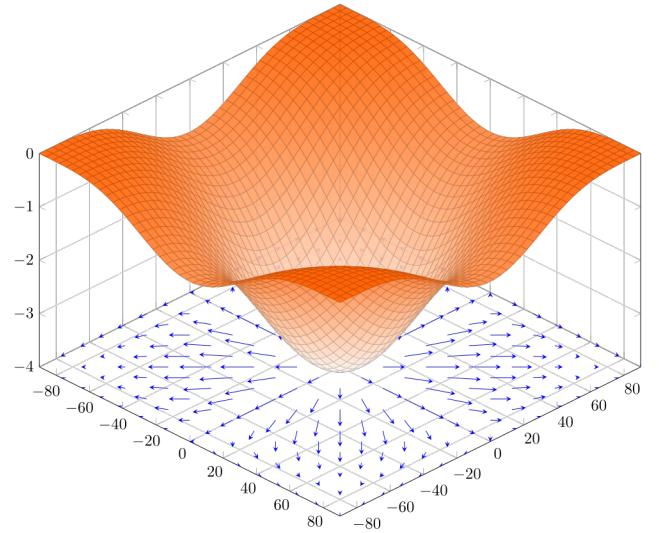


Figure 3: Function representation and gradient field of the function $-(\cos^2 x + \cos^2 y)^2$

with C the number of class (eg: 10 for the mnist dataset), y_c the observed probability, and \hat{y}_c the neural network output.

and even

$$\text{Cross Entropy} = -\log(\hat{y}_c)$$

if y_c is a one hot encoded, vector of observed values, (ie: 1 for the observed class and 0 for others).

5 The learning process

5.1 Gradient

The gradient of a function ∇f is the vector field that represents the direction and rate of the quickest increase of the function from a point p . We can visualize the gradient of the function $-(\cos^2 x + \cos^2 y)^2$

at the bottom of the image, we can visualize the vector field for each point (x, y) . with the arrow representing the quickest increase.

5.2 Minimization

To make our neural network lean, we want to adjust his parameters (ie: weights and biases) to **minimize the cost function** C . The algorithm used to minimize the cost function is known as the *Gradient descent*, and the idea, is to get down our function by going through the stepest descent in the function, repeating that way we are sure to reach, at least a local minimum.

Here is the process for a gradient descent step :

1. feed **all** your input data points into the network, calculating the cost for every output.
2. calculate the averaged cost function.
3. compute the gradient of the averaged cost function **with respect to the network parameters**.

4. update the parameters in the opposite direction of the gradient

5.3 backpropagation

The algorithm used to calculate the gradient of the cost function with respect to the neural network parameters is called backpropagation. It is a way to precisely calculate (and not only estimate like the finite difference method) this gradient.

The gradient will be calculated for each layer, starting from the last layer and going all the way back to the first layer. The key mathematical tool used in backpropagation is the chain rule, which allow the gradient of the cost function to be split into gradient of simpler functions at each neuron within the network.

For each layer, to calculate the parameters gradient $\frac{\partial C}{\partial P}$, we need to have the output gradient $\frac{\partial C}{\partial Y}$, this gradient will be given to us by the next layer.

we also want to calculate $\frac{\partial C}{\partial X}$, the input gradient X , to give it to the previous layer, which will use it as her own **output gradient**, because, as we said, the output gradient of one layer is the input gradient of the next layer

5.4 Dense layer backpropagation

The dense layer have two trainable parameters, the weights, W and the biases B . Starting with the gradient of the cost function with respect to the weights, we are given the output gradient :

$$\frac{\partial C}{\partial Y} = \begin{bmatrix} \frac{\partial C}{\partial y_1} \\ \frac{\partial C}{\partial y_2} \\ \vdots \\ \frac{\partial C}{\partial y_j} \end{bmatrix} \quad (1)$$

The cost function C depend on Y which depend on W thus C depend on W

$$\frac{\partial C}{\partial W} = \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial W}$$

We calculate each weight gradient component $\frac{\partial C}{\partial W_{ij}}$

$$\frac{\partial C}{\partial W_{ij}} = \sum_{l=1}^N \frac{\partial L}{\partial Y_{lj}} \frac{\partial Y_{lj}}{\partial W_{ij}}$$

because W_{ij} is used in every example for calculating the j^{th} column of the output matrix. Let's look at the formula for Y_{lj} to get the derivative

$$Y_{lj} = X_{l1}W_{1j} + \dots + X_{li}W_{ij} + \dots + X_{lp}W_{pj} + b_j$$

$$\begin{aligned} \frac{\partial Y_{lj}}{\partial W_{ij}} &= X_{li} \\ \Rightarrow \frac{\partial C}{\partial W_{ij}} &= \sum_{l=1}^N \frac{\partial C}{\partial Y_{lj}} X_{li} = \sum_{l=1}^N X_{li}^T \frac{\partial C}{\partial Y_{lj}} \end{aligned}$$

So finally

$$\frac{\partial C}{\partial W} = X^T \frac{\partial C}{\partial Y}$$

for the biases B .

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial b} \quad (2)$$

$$\frac{\partial C}{\partial b_i} = \sum_{l=1}^N \frac{\partial C}{\partial Y_{li}} \frac{\partial Y_{li}}{\partial b} \quad (3)$$

since the bias term b_i is used in the evaluation of the entire column of Y .

$$\frac{\partial Y_{li}}{\partial b_i} = 1$$

$$\frac{\partial C}{\partial b_i} = \sum_{l=1}^N \frac{\partial C}{\partial Y_{li}}$$

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial Y}$$

and lastly, we calculate the gradient with respect to X

$$\frac{\partial C}{\partial X} = \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial X} \quad (4)$$

$$\frac{\partial C}{\partial X_{ij}} = \sum_{l=1}^k \frac{\partial C}{\partial Y_{il}} \frac{\partial Y_{il}}{\partial X_{ij}} \quad (5)$$

since the data point i will only influence the data point i in Y . Other data points will not be affected. Further, X_{ij} is used in calculation of every dimension of $Y_{i,:}$. To calculate the gradient,

$$Y_{il} = \sum_{t=1}^p X_{it}W_{tl}$$

$$\frac{\partial Y_{il}}{\partial X_{ij}} = W_{jl}$$

$$\frac{\partial C}{\partial X_{ij}} = \sum_{l=1}^k \frac{\partial C}{\partial Y_{il}} W_{jl} = \sum_{l=1}^k \frac{\partial C}{\partial Y_{il}} W_{lj}^T$$

$$\Rightarrow \frac{\partial C}{\partial X} = \frac{\partial C}{\partial Y} W^T$$

5.5 Activation layer backpropagation

An activation layer doesn't have any trainable parameters, so we just gonna focus on the input gradient calculation.

we are given the output gradient $\frac{\partial C}{\partial Y}$: we know that

$$Y = \begin{bmatrix} f(X_1) \\ f(X_2) \\ \vdots \\ f(X_i) \end{bmatrix}$$

$$\frac{\partial C}{\partial X} = \frac{\partial C}{\partial Y} \frac{\partial Y}{\partial X}$$

$$\frac{\partial C}{\partial X_{ij}} = \sum_{l=1}^N \frac{\partial C}{\partial Y_{lj}} \frac{\partial Y_{lj}}{\partial X_{ij}}$$

using the same procedure as for the previous calculations

$$\Rightarrow \frac{\partial C}{\partial X} = \frac{\partial C}{\partial Y} \odot Y'$$

with

$$Y' = \begin{bmatrix} f'(X_1) \\ f'(X_2) \\ \vdots \\ f'(X_i) \end{bmatrix}$$

6 Optimizers

At this point, to update a neural network parameter θ , we use the gradient descent formula, with a fixed learning rate η

$$\theta := \theta - \eta \frac{\partial C}{\partial \theta}$$

This is known, in the Keras API for example as the *Stochastic gradient descent optimizer*, but it exists other optimizers such as *Adam*.

6.1 Adam

Adam accelerates the convergence towards the optimal set of weights by taking into account past gradients, and Adam will adapt the learning rate for each individual weight.

ADAM initializes two vectors, m and v , which store the exponential moving averages of past gradients and past squared gradients, respectively. These vectors are used to scale the gradient updates adaptatively.

1. **Initialization:** Initialize two vectors, m and v , to store the moving averages of the gradients and their squares, respectively.
2. **Computing Moving Averages of the Gradients:**

$$\begin{aligned} m_t &= \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t, \\ v_t &= \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2, \end{aligned}$$

where g_t is the gradient at time step t , and β_1 and β_2 are factors that control the decay rates.

3. **Bias Correction:** Correct initial bias in m and v :

$$\begin{aligned} \hat{m}_t &= \frac{m_t}{1 - \beta_1^t}, \\ \hat{v}_t &= \frac{v_t}{1 - \beta_2^t}. \end{aligned}$$

4. **Update Weights:** Adjust weights based on the corrected moments:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \cdot \hat{m}_t,$$

where θ are the parameters, η is the learning rate, and ϵ is a small constant for numerical stability.

7 All about batch

We went through several feed forward implementations. The first one was to pass data point one by one in the network $x \in \mathbb{R}^i \dots$. This was very slow, for example, for the MNIST dataset, which has 60k training examples, we were going through the whole neural network: feed forward, cost evaluation, back propagation, 60k times. But we then learned that we can instead pass *batch* of data in the neural network, in one time, and this batch can be of any length. The input then becomes $x \in M_{n \times i}(\mathbb{R})$ with each line of the input matrix representing a single data point.

The calculation became much faster, because we can now feed n elements at a time in the neural network.

8 Initializer

Initializers determine the starting weights of a neural network, it is important to choose the right initializer, depending on our network architecture, to avoid problems like vanishing and exploding gradients, and to also optimize convergence.

At the moment we are writing this report, we have implemented two initializers inside our library,

1. **He** The He initializer, designed for ReLU activation functions, initializes weights from a normal distribution centered at zero with a standard deviation of $\sqrt{\frac{2}{n}}$ (where n is the number of input nodes). This setup helps prevent the vanishing gradient problem in deep networks with ReLU activations, ensuring that gradients remain large enough to sustain effective learning through many layers.
2. **Glorot uniform** Conversely, the Glorot Uniform initializer, also known as Xavier Uniform, is ideal for networks using sigmoid or tanh activations. It selects weights from a uniform distribution within $[-c, c]$, where $c = \sqrt{\frac{6}{n_{in} + n_{out}}}$. This approach maintains a balance in the variance of neurons' outputs across the network, aiding in stable gradient flow in deep networks with saturating activations.

9 Convolutional neural network

ALEXIS a toi de travailler

10 Metrics

With the mnist dataset, which is a classification problem, we can implement multiple metrics to evaluate the network performances, we will explain the one we implemented in our library.

1. Accuracy measures the proportion of correct prediction to the total observations in the dataset. for example if we have correctly classified 950 image out of 1000, the accuracy will be 95percent.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

Precision and Recall are two metrics acting on a specified class. **Retrieved elements** are all the network prediction for a specific class. **Relevant elements** are all the correct instances of the class [2].

2. Precision answer the question : *How many retrieved items are relevant*

$$\text{Precision} = \frac{\text{Relevant retrieved instance}}{\text{All retrieved instance}}$$

3. Recall answers the question: *How many relevant item are retrieved*

$$\text{Recall} = \frac{\text{Relevant instances}}{\text{All relevants instance}}$$

11 Training, Validation, and Test Sets

To monitor the performances of a neural network effectively, it is important to break down the dataset into three distinct subsets: training, validation, and test sets.

- **Training set** : The training set is used during the learning phase, with training data directly influencing the parameters updates
- **Validation set** : The validation set is used to provide an unbiased evaluation of a model fit on the training dataset while tuning the model's hyperparameters. This set is used during training, to prevent overfitting when the model performs well on the training data but poorly on unseen data.
- **Test set** : The test set is used only after the model's training and validation phases are complete. It provides the final, unbiased performance metric of the neural network.

12 Results

We include a small code snippet, showing our custom library api to declare a sequential neural network.

```
let net = NeuralNetworkBuilder::new()
  .push(DenseLayer::new(
    28 * 28,
    32,
    InitializerType::He,
  ))
  .push(
    ActivationLayer::from(Activation::ReLU)
  )
  .push(
    DenseLayer::new(32, 10,
      ↪ InitializerType::He)
  )
  .push(
    ActivationLayer::from(
      Activation::Softmax
    )
  );
  .watch(MetricsType::Accuracy)
  .watch(MetricsType::Precision)
  .watch(MetricsType::Recall)
Ok(net.compile(GradientDescent::new(0.01),
  ↪ CostFunction::CrossEntropy?))
```

the layer are added in sequential order with the push method, we watched Accuracy, Precision and Recall. the watch method will add a *MetricsType* to the list of metrics that the program need to calculate. those metrics can be given over epochs for the training and validation dataset. and for the test dataset.

To create a great network, we will try to reach the *VC trade-off principle* [6] [8], this principle state that:

1. A more complex model might fit the training data very well (low training error) but risks overfitting, meaning it could perform poorly on new data because it is too tailored to the training set (high variance).
2. A less complex model might not fit the training data as closely (higher training error) but can generalize better on new data because it does not capture the noise and specific details of the training set (low variance).

Our challenge will thus be, to find the right level of complexity that minimizes both the error on the training data.

12.1 Multi-layer perceptron

before trying to get a great fitting model, we will analyse the result of an overfitting model, to see how this behavior translate into metrics. To make a model overfit, we have multiple options

- increase the model complexity (number of hidden layers, number of neuron per layer)
- reduce the batch size to make the learning process very granular and increase the chance of the network learning to much specific patterns

In our case we have choosen to lower the batch size, here is our network hyperparameters

hyperparameters	epochs	batch size
Value	30	5

Table 1: hyperparameters overfitting model

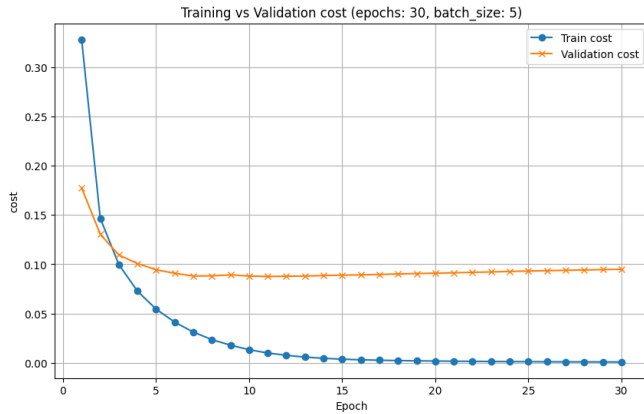


Figure 4: Training vs Validation cost overfitting model

In the cost comparaisn between our training data and our validation data, the overfitting clearly stand out, our train cost (or loss) get very low, while the validation cost remain stable this mean that our neural network has learn to much specific pattern in the training data, making the error going very low on this set, but on the validation set, the network doesn't achieve a great generalization.

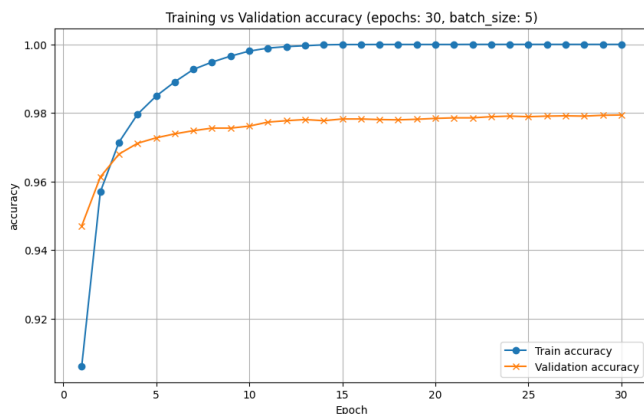


Figure 5: Training vs Validation cost overfitting model

The same observation applies to the evolution of accuracy over epochs. the accuracy reach a stunning 100percent on our training data, meaning that with the point used to train the network, it dont make any mistakes, but the validation accuracy struggle to get up.

We have achieved our best accuracy with the following hyperparameters and network definition.

hyperparameters	epochs	batch size
Value	15	100

Table 2: hyperparameters greatfitting model

```
let net = SequentialBuilder::new()
.push(
  DenseLayer::new(784, 256,
    InitializerType::He))
.push(
  ActivationLayer::from(
    Activation::ReLU)
)
.push(
  DenseLayer::new(256, 256,
    InitializerType::He))
.push(
  ActivationLayer::from(
    Activation::ReLU)
)
.push(
  DenseLayer::new(256, 10,
    InitializerType::GlorotUniform))
.push(
  ActivationLayer::from
    (
      Activation::Softmax)
)
.watch(MetricsType::Accuracy);
Ok(net.compile(GradientDescent::new(0.04),
  ↪ CostFunction::CrossEntropy?))
```

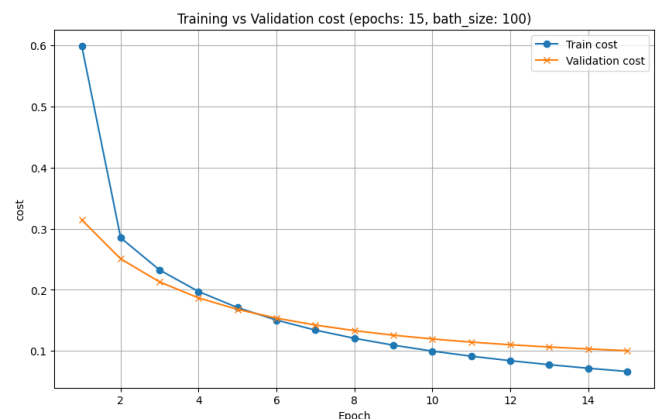


Figure 6: Training vs Validation cost greatfitting model

Here, our model seems to have generalized patterns, because the training cost and the validation cost follow the same path, slowly increasing, with the gap between curves never blowing off the roof.

we reach a validation cost of 0.097, the cost for the testing data set is 0.091.

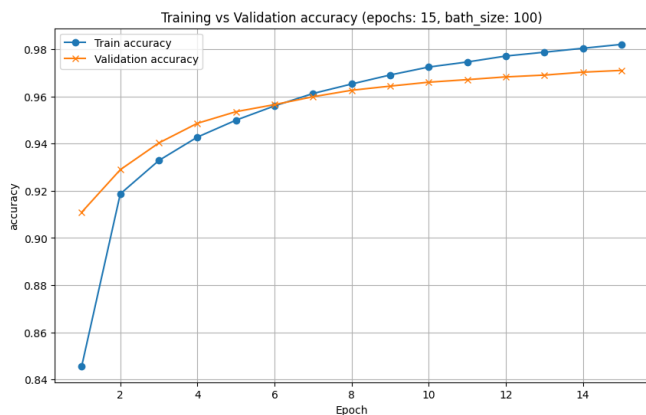


Figure 7: Training vs Validation cost greatfitting model

we reach a validation accuracy of 97.10percent for the final eppochs, the accuracy for the testing data set is 97.31percent. which is not a *state of the art* accuracy but is a good result.

12.2 Convolutional multi-layer perceptron

13 Conclusion

This project of creating our own *from scratch* library has brought us a lot of knowledge, It as made us greatly understand the whole neural network process. We have planned to add more functionality to our library and to optimize the code base even after the end of this project. we want to add more optimizers, the dropout layer, and a learning rate scheduler. and over the longer term, we'd like to add other neural network paradigm, such as Re-current neural network.

References

- [1] Papers with Code. *He Initialization*. n.d. URL: <https://paperswithcode.com/method/he-initialization>.
- [2] Wikipedia contributors. *Precision and recall* — *Wikipedia, The Free Encyclopedia*. [Online; accessed 28-May-2024]. 2024. URL: https://en.wikipedia.org/wiki/Precision_and_recall.
- [3] Stack Exchange. *Backpropagation with Softmax Cross Entropy*. 2016. URL: <https://stats.stackexchange.com/questions/235528/backpropagation-with-softmax-cross-entropy>.
- [4] Google. *Cross Entropy*. 2024. URL: <https://www.google.com/search?q=cross+entropy&sourceid=chrome&ie=UTF-8>.
- [5] V7 Labs. *A Guide to Cross-Entropy Loss*. n.d. URL: <https://www.v7labs.com/blog/cross-entropy-loss-guide>.
- [6] Yann LeCun. *Quand la machine apprend: La révolution des neurones artificiels et de l'apprentissage profond*. Odile Jacob, 2019.
- [7] Towards Data Science. *Derivative of the Softmax Function and the Categorical Cross Entropy Loss*. 2019. URL: <https://towardsdatascience.com/derivative-of-the-softmax-function-and-the-categorical-cross-entropy-loss-ffceefc081d1>.
- [8] Vladimir Vapnik and Alexey Chervonenkis. “Theory of pattern recognition”. In: (1974).
- [9] Wikipédia. *Matrice Jacobienne*. 2024. URL: https://fr.wikipedia.org/wiki/Matrice_jacobienne.
- [10] Wikipedia. *Hadamard product (matrices)*. 2024. URL: [https://en.wikipedia.org/wiki/Hadamard_product_\(matrices\)](https://en.wikipedia.org/wiki/Hadamard_product_(matrices)).
- [11] Wikipedia. *One-hot*. 2024. URL: <https://en.wikipedia.org/wiki/One-hot>.
- [12] Wikipedia. *Precision and recall*. 2024. URL: https://en.wikipedia.org/wiki/Precision_and_recall.
- [13] Wikipedia. *Softmax function*. 2024. URL: https://en.wikipedia.org/wiki/Softmax_function.