
A webcam-based human-computer interface for defining smooth 3D trajectories by tracking 2D hand landmarks

Master's Thesis submitted to the
Faculty of Informatics of the *Università della Svizzera Italiana*
in partial fulfillment of the requirements for the degree of
Master of Science in Informatics

presented by
Umberto Cocca

under the supervision of
Dr. Alessandro Giusti
co-supervised by
Dr. Gianluigi Ciocca, Dr. Loris Roveda

February 2022

I certify that except where due acknowledgement has been given, the work presented in this thesis is that of the author alone; the work has not been submitted previously, in whole or in part, to qualify for any other academic award; and the content of the thesis is the result of work which has been carried out since the official commencement date of the approved research program.

Umberto Cocca
Lugano, 23 February 2022

-dedica-

“Gutta cavat lapidem.”
lat. «la goccia scava la pietra»

Abstract

Robotic systems are increasingly being adopted for filming and photography purposes. By exploiting robots, in fact, it is possible to program complex motions for the camera, achieving high-quality videos/photographs. Our contribution is to propose an alternative to the joystick: being able to control a drone using just a hand and giving space to new ways of expressing video content, even to those who are novices. To achieve this, we capture 3D movements (using a hand tracking system) after building a hand gesture recognition and setting up a solid pipeline to detect 3D trajectories obtained from 2D pixel landmarks, hence estimating the orientation and distance from the camera. Trajectories are interpolated and smoothed with ridge regression. As proof, the captured trajectory was launched on a drone, in simulation, implemented in the ROS framework. The whole pipeline can be easily translated into another kind of task, e.g. interaction and communication in AR / VR.

Acknowledgements

Work in progress

Contents

List of Figures	viii
List of Tables	ix
List of Equations	x
List of Listings	xi
Introduction	1
1 Literature review	2
2 Background	3
2.1 Regression	3
2.1.1 Linear Models for Regression	3
2.1.2 Normal equations	4
2.2 Regularization	6
2.2.1 The Problem Of Overfitting	7
2.2.2 Cost Function	7
2.2.3 (Batch) Gradient Descent	8
2.3 Artificial Neural Networks	11
2.3.1 Feedforward Fully-Connected Neural Networks	11
2.3.2 Neural Network (NN) Setup for Classification	15
2.3.3 Optimisation algorithms	17
3 Tools	18
4 Methodologies	19
5 Evaluation	20
Conclusion and perspectives	21
A List of Acronyms	22

Bibliography**23**

Figures

2.1	Gradient descent.	9
2.2	Image of a human neuron.	11
2.3	Feed forward neural network.	12
2.4	Image of a human neuron.	14
2.5	Leaky Rectified linear activation.	15

Tables

Equations

2.1	Sum of squares regression.	5
2.2	Sum of squares regression in matrix notation.	5
2.3	Design matrix.	5
2.4	Weights of $M - 1$ features and target of N examples.	5
2.5	Solve the optimization problem for ridge regression.	6
2.6	Normal equations.	6
2.7	Cost function Regularisation term.	7
2.8	Mean Squared Error.	9
2.9	Regularisation term.	9
2.10	Gradient Descent.	10
2.11	Gradient Descent.	10
2.12	Cost function for Gradient Descent.	10
2.13	Forward propagation.	13
2.14	Leaky Rectified linear activation.	14
2.15	Forward propagation.	15
2.16	Forward propagation.	16
2.17	Forward propagation.	16
2.18	Forward propagation.	16
2.19	Softmax cross entropy.	17

Listings

Introduction

Chapter 1

Literature review

This chapter discusses previous research about the topic, providing a brief introduction to the approach we adopted.

Chapter 2

Background

This chapter provides some background concepts to understand the material presented in this thesis.

2.1 Regression

The goal of regression is to predict the value of one or more continuous target variables t given the value of a D -dimensional vector x of input variables. [Bishop, 2006]

Given a training data set comprising N observations x_n , where $n = 1, \dots, N$, together with corresponding target values t_n , the goal is to predict the value of t for a new value of x .

From a probabilistic perspective, we aim to model the predictive distribution $p(t|x)$ because this expresses our uncertainty about the value of t for each value of x .

2.1.1 Linear Models for Regression

The simplest linear model for regression is one that involved a linear combination of the input variables:

$$y(x, w) = w_0 + w_1x_1 + \dots w_Dx_D \quad (2.1)$$

where $x = (x_1, \dots, x_D)^T$. This is often simply known as linear regression. The key property of this model is that it is a linear function of the parameters w_i , but also of x_i and establishes significant limitations on the model. It is possible to extend the class of models by considering linear combinations of fixed nonlinear functions of the input variables:

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x) \quad (2.2)$$

where $\phi_j(x)$ are known as basis functions. The total number of parameters in this model will be M .

The parameter w_0 is called bias parameter. It is often convenient to define an additional dummy "basis function" $\phi_0(x) = 1$, so that:

$$y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \phi(x) \quad (2.3)$$

where $w = (w_0, \dots, w_{M-1})$ and $\phi = (\phi_0, \dots, \phi_{M-1})^T$

By using non linear basis functions, we allow the function $y(x, w)$ to be a non linear function of the input vector x . A particle example of this model where there is a single input variable x is the polynomial regression. The basis functions take the form of powers of x so that $\phi_j(x) = x^j$.

There are other possible choices for the basis functions as:

$$\phi_j(x) = e^{\frac{-(x-u_j)^2}{2s^2}} \quad (2.4)$$

where u_j regulates the locations of the basis functions in input space, while the parameter s is their spatial scale. These are usually referred to as "Gaussian" basis functions.

We can use simply the identity basis functions in which the vector $\phi(x) = x$.

2.1.2 Normal equations

The values of the coefficients will be determined by fitting the polynomial to the training data. This can be done by minimizing an error function that measures the misfit between the function $y(x, w)$, for any given value of w , and the training set data points. One simple choice of error function, which is widely used, is given by the Sum of Squared estimate of Errors (SSE) between the predictions $y(x_n, w)$ for each data point x_n and the corresponding target values t_n (called also sum of squares regression [Valchanov, 2018]), so that we minimize:

$$E(w) = \frac{1}{2} \sum_{n=1}^N [y(x_n, w) - t_n]^2 \quad (2.5)$$

Equation 2.1. Is a statistical technique used in regression analysis to determine the dispersion of data points and the function that best fits (varies least) from the data.

where the factor of $1/2$ is included for mathematical convenience. It is a nonnegative quantity that would be zero if, and only if, the function $y(x, w)$ were to pass exactly through each training data point.

We can solve the curve fitting problem by choosing the value of w for which $E(w)$ is as small as possible. Because the error function is a quadratic function of the coefficients w , its derivatives with respect to the coefficients will be linear in the elements of w , and so the minimization of the error function has a unique solution w^* . The resulting polynomial is given by the function $y(x, w^*)$.

We can write (2.5) in matrix notation as:

$$\frac{1}{2}(\phi w - t)^T(\phi w - t) \quad (2.6)$$

Equation 2.2. This is the sum of squares regression in matrix notation.

where ϕ is an $N \times M$ matrix, so that:

$$\phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix} \quad (2.7)$$

Equation 2.3. This is called the design matrix whose elements are given by $\phi_{nj} = \phi_j(x_n)$.

$$w = \begin{pmatrix} w_0 \\ \vdots \\ w_{M-1} \end{pmatrix} \quad t = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix} \quad (2.8)$$

Equation 2.4. Weights of $M - 1$ features and target of N examples.

Therefore, we have to solve the optimization problem finding the minimum of the cost function $E(w)$:

$$w^* = \arg \min_w \frac{1}{2}(\phi w - t)^T(\phi w - t) \quad (2.9)$$

Equation 2.5. The goal is to find the weight w that minimize the cost function $E(w)$.

So we compute the gradient and solving for w we obtain:

$$\begin{aligned} \nabla E(w) &= \phi^T(\phi w - t) = 0 \\ \phi^T \phi w - \phi^T t &= 0 \\ \phi^T \phi w &= \phi^T t \\ w^* &= (\phi^T \phi)^{-1} \phi^T t \end{aligned} \quad (2.10)$$

Equation 2.6. They are known as the normal equations for the least squares problem

This is the real minimum because if we take the second derivative $\nabla^2 E(w) = \phi^T \phi$ and this is a symmetric matrix, so it's also positive definitive matrix, which means that this objective function that we try to minimize is convex. So if we find a stationary point, such that derivative is zero, we also find a global minimum. Furthermore, to find a solution we need to invert this matrix $\phi^T \phi$, so we need some condition that assure this is invertible and this is the case when the columns of the matrix are linearly independent.

Once we find the solution w^* and we receive a new data point that we never seen during training, we just predict the new target t^* as:

$$t^* = w^{*T} \phi(x) \quad (2.11)$$

2.2 Regularization

If we use a set of features that is too expressive, for example a ten polynomial grade (x^{10}), then this interpolates very close our training data, because we have a model with 10 parameters where we can perfectly represents the data points. The risk is that the model became something like this:

[picture]

We want a trade-off between fits data and able to generalize. Infact, on the other hand if our model is not to expressive and not to complex our data will be linearly rapresentable in the feature space and this means that the performance will be very poor.

2.2.1 The Problem Of Overfitting

We want to penalize for some features with parameters with high values, if our model is overfitting is very likely that the parameters of our model will have a big magnitude, this means that also the features supply to this parameters will be higher and very low and this cause a lot of problems.

In order to control over-fitting we introduce the idea of adding a regularization term to an error fuction, so that the total error function to minimized takes the form:

$$E(w) = E_D(w) + \lambda E_W(w) \quad (2.12)$$

Where λ is the regularization coefficient and it is the trade-off between how we well fit training set and how to establish the parameters w with low values, therefore having simple hypotesys avoiding over-fitting. E_D is the error based on dataset, while E_W is based on weights.

2.2.2 Cost Function

One of the simplest forms of regularizer is given by the sum-of-squares of the weight vector elements:

$$E_W(w) = \frac{1}{2} w^T w = \frac{\|w\|_2^2}{2} \quad (2.13)$$

Equation 2.7. Cost function | Regularisation term.

where $\|w\|_2$ is the euclidean norm $\sqrt{\sum_{i=1}^n x_i^2}$.

This is also called ridge regression, a method of estimating the coefficitents of multiple-regression models in scenarios where indipendent variables are highly correlated.

If we also consider the sum-of-squares error function given by:

$$E_D(w) = \frac{1}{2} [t_n - w^T \phi(x_n)]^2 \quad (2.14)$$

then the total error function becomes:

$$E(w) = \frac{1}{2} \sum_{n=1}^N [t_n - w^T \phi(x_n)]^2 + \frac{\lambda}{2} w^T w \quad (2.15)$$

This particular choice of regularizer is known in the machine learning literature as weight decay because in sequential learning algorithms, it encourages weight values to decay towards zero.

Setting the gradient of $E(w)$ wrt w to zero, and solving for w , we obtain:

$$w^* = (\lambda I + \phi^T \phi)^{-1} \phi^T t \quad (2.16)$$

and this is an extension of the least-squares solution (2.10). This is a better version than before because it is also possible to prove that $(\lambda I + \phi^T \phi)$ is always invertible if $\lambda > 0$, therefore w^* always exists.

2.2.3 (Batch) Gradient Descent

Gradient descent is a first-order iterative optimization algorithm for finding the minimum w of a function $E(w)$. To achieve this goal, it performs two steps iteratively, until convergence:

- Compute the slope (gradient) that is the first-order derivative of the function at the current point
- Move in the opposite direction of the slope increase from the current point by the computed amount

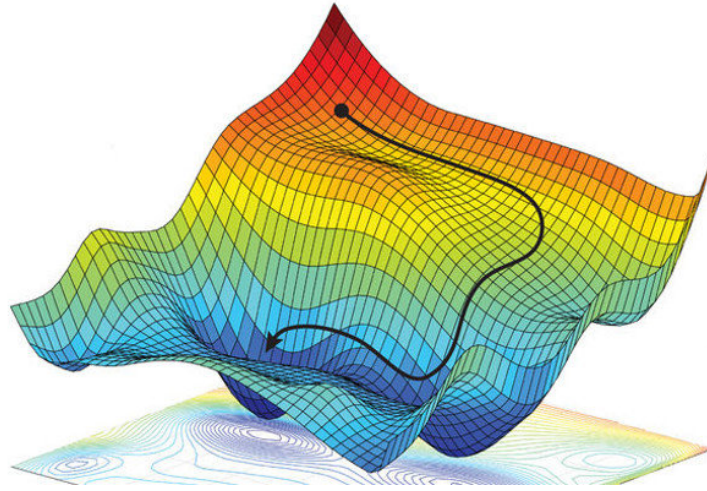


Figure 2.1. Gradient descent is based on the observation that if the multi-variable function $E(w)$ is defined and differentiable in a neighborhood of a point w_0 , then w_0 decreases fastest if one goes from w_0 in the direction of the negative gradient of $E(w)$ at w_0 , $-\nabla E(w_0)$.

Let's now using Mean Squared Error (MSE), instead of using (2.5) due to the fact that we are able to reach obvious benefits: first of all to keep the value in a expressible range usable by computers, then to make the results comparable across samples regardless of the size of the sample. Infact, the SSE depends on how many terms are added up (note the case of millions/billions of data points). In addition, using SSE or MSE it still leads to find an equivalent solution.

$$E_D(w) = \frac{1}{2N} \sum_{n=1}^N [y(x_n, w) - t_n]^2 \quad (2.17)$$

Equation 2.8. Mean Squared Error. $1/2$ is added, as in SSE, so the derivative doesn't need a constant out front. We get away with it because the minima of $E_D(w)$ and $E_D(w)/2$ are achieved at the same value(s) of w .

Concerning the regularisation term $E_W(w)$ (2.13), we could also write:

$$E_W(w) = \frac{\sum_{m=1}^M w_m^2}{2} \quad (2.18)$$

Equation 2.9. Regularisation term.

Note that, conventionally, m starts from 1, and not from 0, even if it exists (2.8). Nev-

ertheless, it's important to know that nothing change consistently even if we consider the 0th weight.

Combining together following (2.12)

$$E(w) = \frac{1}{2N} \left\{ \sum_{n=1}^N [y(x_n, w) - t_n]^2 + \lambda \sum_{m=1}^M w_m^2 \right\} \quad (2.19)$$

Equation 2.10. Gradient Descent.

Have some function $E(w)$

Want $\arg \min_w E(w)$

Keep changing w to reduce $E(w)$ until we hopefully end up at a minimum:

$$\begin{aligned} \text{Repeat} = \{ \\ & w_0 := w_0 - \alpha \frac{1}{N} \sum_{n=1}^N [y(x_n, w) - t_n]^2 x_0 \\ & w_j := w_j - \alpha \frac{1}{N} \sum_{n=1}^N [y(x_n, w) - t_n]^2 x_j + \frac{\lambda}{M} w_j \\ & \} \end{aligned} \quad (2.20)$$

Equation 2.11. Cost function for Gradient Descent.

where in compact form is:

$$\begin{aligned} \text{Repeat} = \{ \\ & w_j := w_j (1 - \alpha \frac{\lambda}{M}) - \alpha \frac{1}{N} \sum_{n=1}^N [y(x_n, w) - t_n]^2 x_j \\ & \} \end{aligned} \quad (2.21)$$

Equation 2.12. Cost function for Gradient Descent.

2.3 Artificial Neural Networks

The term Neural Network (NN) has its origins in attempts to find mathematical representations of information processing in biological systems. In fact, Artificial Neural Networks (ANNs), subset of Machine Learning (ML) field, are models inspired from the biological performance of human brain [Andina et al., 2007].

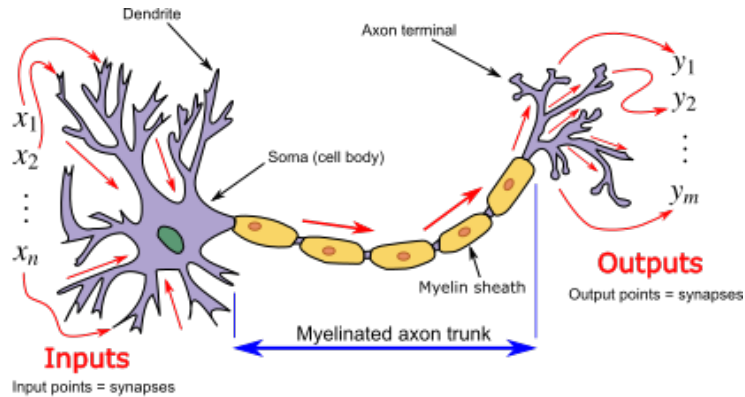


Figure 2.2. Neuron and myelinated axon, with signal flow from inputs at dendrites to outputs at axon terminals. The signal represents a short electrical pulse called 'spike'.

Neurons have cell body (fig. 2.2) and a number of input wires called dendrites. Neurons also have an output wire called axon, used to send signals to other neurons. At a simplistic level the neuron is a computational unit that gets a number of inputs through its input wires, then does some computation and finally it sends outputs to other neurons connected to it in the brain.

2.3.1 Feedforward Fully-Connected Neural Networks

In (fig. 2.3) is visible a NN, seen as mathematical model. It's just a group of this different neurons strung together. Trying to underline an analogy with the biological systems: circles identify the cell body where they are fed with some inputs that pass through the input wires, similar to the dendrites. The neuron does some computation and outputs some value on an output wire, where in the biological neuron it identifies the axon.

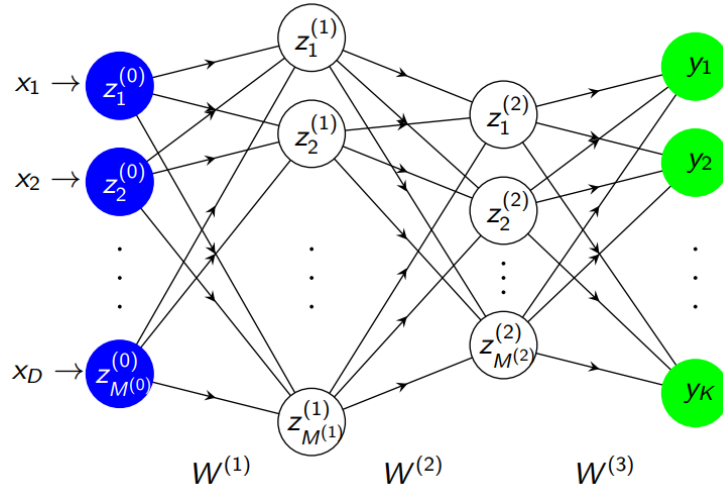


Figure 2.3. The image graphically shows a NN with three layers (or two hidden layers). The input layer has $M^{(0)}$ neurons and the output layer has K outputs. Neurons are organized in layers to allow parallel computation to avoid cyclic dependencies. The process of computing NN outputs from inputs is called forward propagation. This kind of network is also called Multi-layer perceptron.

The $z_m^{(l)}$ are neurons, each takes its input values and computes a single output value from them. Neurons are organized in layers $1, \dots, L$ and usually the starting input is considered the 0th layer. Inputs x_1, \dots, x_D are occasionally called input layer/neurons (even though they do not compute anything). The output of the entire network is then $y = z^{(L)}$, called output layer. Instead, the internal layers are called hidden layers. Each layer $l \in \{1, \dots, L\}$ has $M^{(l)}$ neurons.

$M^{(1)}$ neurons perform a perceptron-like computation:

$$u_m^{(1)} = (w_m^{(1)})^T x + b_m^{(1)}, \quad z_m^{(1)} = f(u_m^{(1)}), \quad m = 1, \dots, M^{(1)} \quad (2.22)$$

with a differentiable activation function f for gradient descent. This step is iterated multiple times taking the outputs of the previous step:

$$z^{(l-1)} = (z_m^{(l-1)})_{m=1, \dots, M^{(l-1)}} \quad (2.23)$$

as input of:

$$u_m^{(l)} = (w_m^{(l)})^T z^{(l-1)} + b_m^{(l)}, \quad z_m^{(l)} = f(u_m^{(l)}) \quad (2.24)$$

where $m = 1, \dots, M^{(l)}$ and $l = 2, \dots, L$. Weights w are usually independent for each step. Additionally define $z^{(0)}$ to be the input, i.e.

$$z^{(0)} = x \quad (2.25)$$

For each layer $l \in 1, \dots, L$ the computation is

$$z_m^{(l)} = f((w_m^{(l)})^T z^{(l-1)} + b_m^{(l)}) \quad (2.26)$$

which can be written as a matrix multiplication:

$$z^{(l)} = f(W^{(l)} z^{(l-1)} + b^{(l)}) \quad (2.27)$$

Equation 2.13. Function that identifies input transformation at each step l of the net.

The weights w are directed connections between the neurons, e.g. the neurons of layer 2 are connected to the ones of layer 1 by the weights $w_{mn}^{(2)}, m = 1, \dots, M^{(1)}, n = 1, \dots, M^{(2)}$. The bias b varies according to the propensity of the neuron to activate, influencing its output.

$f()$ is an activation function and needs to be differentiable, because we wish to apply gradient descent training. For the hidden layers of the network, the activation function must be nonlinear, because multiple linear computations can be collapsed to a single one, therefore in order to gain power from iterative computation we thus need nonlinear steps.

Many possible activation functions for the hidden layers of a NN exist:

- Sigmoid, Hyperbolic Tangent: Monotonic, squeeze output to a fixed range
- ReLU: "almost linear" (a clipped identity function)

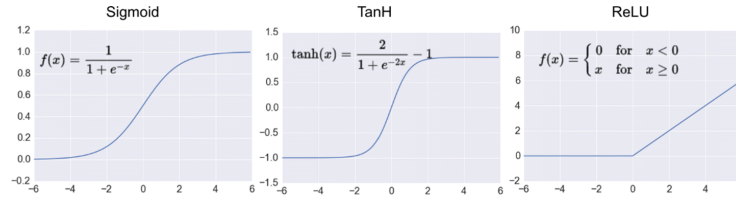


Figure 2.4. Neuron and myelinated axon, with signal flow from inputs at dendrites to outputs at axon terminals. The signal represents a short electrical pulse called 'spike'.

One of the key characteristics of modern deep learning system is to use non-saturated activation function (e.g. ReLU) to replace its saturated counterpart (e.g. sigmoid, tanh). The advantage of using non-saturated activation function lies in two aspects: the first is to solve the so called "exploding/vanishing gradient" [Bengio et al., 1994], in particular on the difficulty of training Recurrent Neural Network (RNN) [Pascanu et al., 2013], while the second is to accelerate the convergence speed. More sophisticated activation function as the "leaky ReLU" trying to solve the dying ReLU problem [Lea]. In contrast to ReLU, in which the negative part is totally dropped, leaky ReLU assigns a non-zero slope to it. Leaky ReLU and its variants are consistently better than ReLU in Convolutional Neural Network (CNN) [Xu et al., 2015].

Leaky Rectified linear activation is introduced in acoustic model [Maas et al.]. Mathematically, it is defined as follows:

$$f(x) = \begin{cases} \frac{x}{a} & \text{if } x < 0 \\ x & \text{otherwise} \end{cases} \quad (2.28)$$

Equation 2.14. Function that identifies input transformation at each step l of the net.

where a is a fixed parameter in range $(1; +\infty)$. In original paper, the authors suggest to set a to a large number like 100.

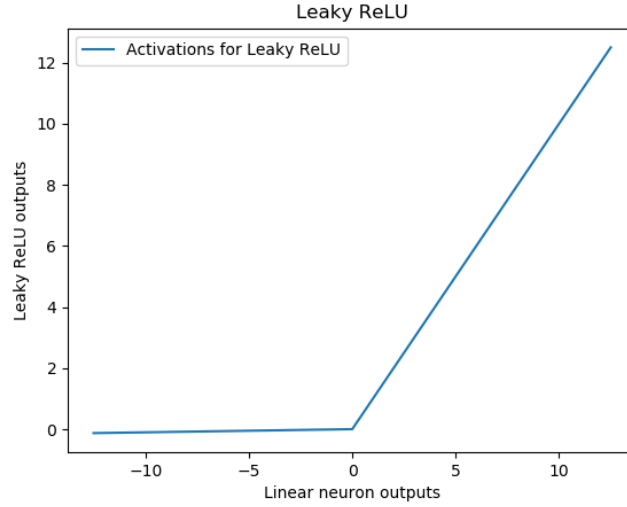


Figure 2.5. Neuron and myelinated axon, with signal flow from inputs at dendrites to outputs at axon terminals. The signal represents a short electrical pulse called 'spike'.

2.3.2 NN Setup for Classification

For a classification task with K classes, we use a K -dimensional output layer. A sample $x \in R^D$ is classified as belonging to class k if the output neuron y_k has the maximal value:

$$c^* = \arg \max_k y_k \quad (2.29)$$

The problem is that the *argmax* function is not differentiable. Therefore, this is solved by letting the NN output a probability distribution over classes:

$$y = (y_k)_{k=1,\dots,K} \quad y_k \geq 0, \quad \sum_k y_k = 1 \quad (2.30)$$

Equation 2.15. Function that identifies input transformation at each step l of the net.

The advantage is that we can derive a differentiable measure of the quality of the output on theoretical grounds, using probability theory. In order to make the network output a probability distribution, we take exponentials and normalize. This is the softmax nonlinearity:

$$S(y) = \left(\frac{e^{y_1}}{\sum_k e^{y_k}}, \dots, \frac{e^{y_k}}{\sum_k e^{y_k}} \right) \quad (2.31)$$

Equation 2.16. Function that identifies input transformation at each step l of the net.

In contrast to other activation functions, it is applied to the full last layer of the NN, not to each independent component. The hidden layers can have any nonlinear activation function.

The learning process is structured as a non-convex optimisation problem in which the aim is to minimise a cost function, which measures the distance between a particular solution and an optimal one.

If we assume a NN with softmax output, we can compute the loss by measuring the cross-entropy between the output distribution and the target distribution.

Encoding the target in one-hot style, e.g. if a sample belongs to class k , the target is:

$$t = (0, \dots, 0, 1, 0, \dots, 0) \quad (2.32)$$

Equation 2.17. Function that identifies input transformation at each step l of the net.

We treat this as a probability distribution: in an ideal world, a perfect hypothesis y would exactly match this t , assigning probability 1 to the correct class, and probability 0 otherwise. The cross-entropy loss is defined as:

$$E_{CE} = - \sum_k (t_k \log y_k) \quad (2.33)$$

Equation 2.18. Function that identifies input transformation at each step l of the net.

The intuition about the cross-entropy corresponds to the number of additional bits needed to encode the correct output, given that we have access to the (possibly wrong) prediction of the network. One property of the cross-entropy loss is that it is always non-negative. For efficiency and numerical stability, one should merge softmax loss and cross-entropy criterion into one function:

$$E_{CE+SM} = - \sum_k (t_k \log S_k(y)) \quad (2.34)$$

Equation 2.19. To train the network with backpropagation, you need to calculate the derivative of the loss. In the general case, that derivative can get complicated, but using the softmax and the cross entropy loss, that complexity fades away.

2.3.3 Optimisation algorithms

The choice of optimization algorithms strongly influence the effectiveness of the learning process as they update and calculate the appropriate and optimal values of that model. Specifically if we consider the gradient descent, the most popular optimization strategy used in machine learning, the extent of the update is determined by the learning rate λ , which guarantees convergence to the global minimum for convex error surfaces and to a local minimum for non-convex surfaces. By the way, there are better optimization as the non linear conjugate gradient [Hager and Zhang, 2006], BFGS, the improved version to decrease memory usage L-BFGS [Saputro and Widyaningsih, 2017], etc... the main advantages are that no need to manually pick λ and often are faster than gradient descent, although more complex algorithms.

The optimiser chosen for this thesis project is Adam, an algorithm for first-order gradient-based optimisation of stochastic objective functions, based on adaptive estimates of lower-order moments [Kingma and Ba, 2017]. It is an extension to stochastic gradient descent that has recently seen broader adoption for Deep Learning (DL) applications in computer vision and Natural Language Processing (NLP).

Chapter 3

Tools

Chapter 4

Methodologies

Chapter 5

Evaluation

Conclusion and perspectives

Appendix A

List of Acronyms

ANN	Artificial Neural Network.....
CNN	Convolutional Neural Network
DL	Deep Learning
ML	Machine Learning.....
MSE	Mean Squared Error.....
NLP	Natural Language Processing.....
NN	Neural Network
RNN	Recurrent Neural Network.....
SSE	Sum of Squared estimate of Errors

Bibliography

- Leaky relu: improving traditional relu – machinecurve. <https://www.machinecurve.com/index.php/2019/10/15/leaky-relu-improving-traditional-relu/>. (Accessed on 01/11/2022).
- Diego Andina, D. Pham, D. Andina, Antonio Vega-Corona, Juan Seijas, and J. Torres-García. *Neural Networks Historical Review*. 02 2007. ISBN 978-0-387-37450-5. doi: 10.1007/0-387-37452-3_2.
- Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.
- Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- W. W. Hager and H. Zhang. A survey of nonlinear conjugate gradient methods, 2006. *Pacific journal of Optimization*, vol. 2, no. 1, pp. 35–58.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. Citeseer.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR, 2013.
- Dewi Retno Sari Saputro and Purnami Widyaningsih. Limited memory broyden-fletcher-goldfarb-shanno (l-bfgs) method for the parameter estimation on geographically weighted ordinal logistic regression model (gwolr). In *AIP Conference Proceedings*, volume 1868, page 040009. AIP Publishing LLC, 2017.
- Iliya Valchanov. Sum of squares total, sum of squares regression and sum of squares error, 2018. URL <https://365datascience.com/tutorials/statistics-tutorials/sum-squares/>. [sum of squares total, denoted SST | Published on 5 Nov 2018].

Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.

Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. Mediapipe hands: On-device real-time hand tracking, 2020.