

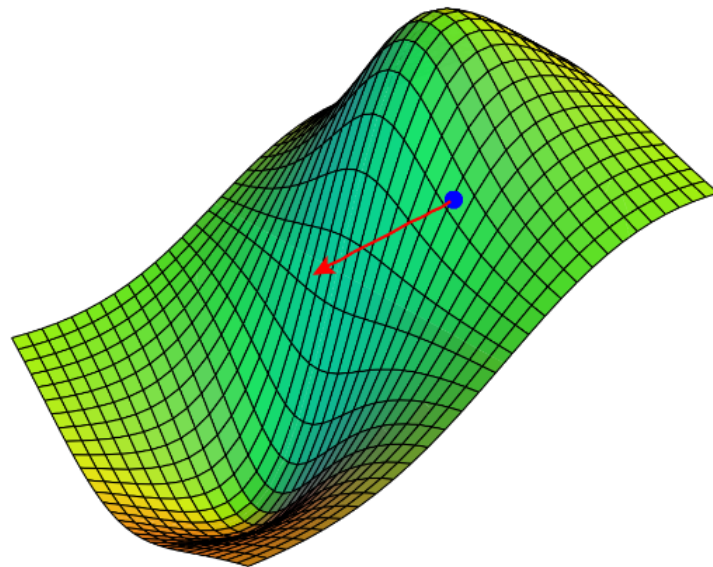


UNIVERSITÉ DE NANTES



IAE NANTES
ÉCONOMIE & MANAGEMENT

Quantitative Research on Market Predictability: Machine Learning and the Efficient Market Hypothesis



Aurouet Lucas

Master EKAP
IAE Nantes

Under the supervision of Prof. O.DARNE 01/03/2021

Contents

1	Introduction	2
2	Theory	3
2.1	Statistics & Machine Learning	3
2.1.1	The Perceptron: a simple example to understand how machines learn	3
2.1.2	Networks & Backpropagation	7
2.1.3	Activation functions and enhanced gradient descent	11
2.1.4	Support Vector Machine & the kernel trick	11
2.1.5	Soft Margin SVM	14
2.2	Finance	14
3	Application	14
4	Conclusion	14

1 Introduction

Trading is futile. This is an oversimplification of the Efficient Market Hypothesis. According to Eugene Fama who formulated this hypothesis in 1970¹, prices at a certain date incorporate all the available and relevant information at that date. There are never overvalued nor undervalued assets, fundamental and/or technical analysis are therefore pointless as all the signals are already reflected by prices on the market. One cannot beat the market, is another way of expressing the EMH, which brings us to our initial statement; trading is futile. As trading basically consists in predicting an asset performance, under the EMH, trading is a wasted effort. It is quite hard to understand how, if the EMH holds, financial institutions have poured significant amounts of money into trying to outperform the market. Even less understandable is how some have succeeded. I think it is reasonable to define what 'outperforming the market' means. 'Outperforming the market' encompasses the achievement of greater returns than a portfolio containing every asset available on the market, with the same level of risk offered by such portfolio. Alternatively, it could mean achieving the same return, with less risk. Eugene Fama states that the only way to expect returns greater than the market is to take additional risks, or that you cannot achieve that same return on investment with a risk inferior to the market risk. We are now able to ask the following question: does the EMH holds with real life data ? Is trading futile ? Can you achieve greater returns than the market for a corresponding market risk ? Or is one constantly better off holding a large index like the S&P 500 for instance ? In my opinion, questioning the EMH is the most riveting topic in quantitative finance as it answers questions from a very theoretical point of view, typically what is observed in academics but also from a very practical point of view, the kind we observe on a trading floor. Econometrics is the science of combining these two aspects, using theoretical tools to shine lights on very down to earth problems. Therefore this work will try to reflect both dimensions. There are many ways to prove or disprove the EMH, but consistently beating the market is one of them (consistency is essential here). The objective here is to use data science to try and outperform the market. The results should definitely not be taken as an attempt to generate profits but rather as an exercise where the objective is to make use of quantitative techniques to elucidate a question. Machine Learning is a very powerful tool when a problem comes to identifying patterns and regularities in the data. As we will discuss later on, it is particularly adapted in our case. One fact often played down is that Machine Learning boils down to a list of algorithms applied in a particular order, and it is very easy to make mistakes when using automated techniques. We hope to have sufficient knowledge to question our own results, especially when working in quantitative finance. Although we are not, as mentioned above, attempting to concretely apply our findings, We think it is only natural to add some perspective when working with financial data. Finance is connected to the economy by many pipelines and misleading Artificial Intelligence can induce disastrous consequences. Artificial Intelligence is a very broad term and Machine Learning is one of its component, the later relies on the idea of being able to find patterns without explicitly defining these patterns, whereas Artificial Intelligence relies on the idea of writing a program capable of the same thought process as a human. They can sometimes be interchangeable and the frontier between the two is often blurry, here we are focusing on learning patterns from data, hence we will almost exclusively refer to Machine Learning as opposed to Artificial intelligence. We will also refer to Statistical Learning, which could be seen as the gap between traditional statistical

¹Efficient Capital Markets: A Review of Theory and Empirical Work - E.Fama

models and Machine Learning models. Statistical Learning solves data related problems by formulating them in a formal statistical problem which will then be solved using parts of Machine Learning algorithms. This gives us a large panel of tools to extract useful information from the data, coerce this information into a functional form that can then be exploited to infer the behavior of unseen data. We will first gaze over the Statistical and Machine Learning methods we will be using ranging from penalized regression to neural networks, we will then detail the financial side of the problem, going over technical analysis, fundamental analysis and the necessary financial mathematics. The second part is a step-by-step journal of the progress of the overall project and what modifications have been made to improve the data processing².

This project is coded in Python 3.8 and is associated with a Jupyter Notebook, available on [GitHub](#). The integrity of the code is open source.

I often compare open source to science. To where science took this whole notion of developing ideas in the open and improving on other peoples' ideas and making it into what science is today and the incredible advances that we have had. - Linus Torvalds, developer of Linux and Git

2 Theory

2.1 Statistics & Machine Learning

2.1.1 The Perceptron: a simple example to understand how machines learn

Artificial Neural Networks (ANN or NN for short) are Machine Learning models designed to replicate the way the human brain processes information and is heavily inspired by biology. In a human brain, the information travels through a network of neurons which can activate to pass the signal to the next neurons. Neurons are connected by branches in which the information transits, starting from a particular point, these branches, along with the neurons that activated create a path to an end point which outputs a decision. The choice for each neuron to carry on the information is what determines the final output. ANN replicate this logic by creating a stream of branches and neurons inside which the information is processed to eventually reach a particular output.

In this section, we will focus on the Perceptron, a particular ANN which consists of only one neuron. It is easier to understand the logic behind ANN using a Perceptron, the results can then be generalized to more complex networks with more neurons. We will try to understand how a Perceptron is able to find patterns in the data without being explicitly programmed. This is what we refer to as the learning algorithm, or how the machine learns from the data.

We begin by giving the model inputs (x_p ³) that account for the exogenous variables. Each input is multiplied by a weight w_j and all weighted inputs are summed. The

²Although data processing often refers to the pre-screening of the data before feeding it into a model, here we refer to the entire process from the data gathering to the diagnosis of the results

³ $p = 1, \dots, P$ where P is the total number of inputs

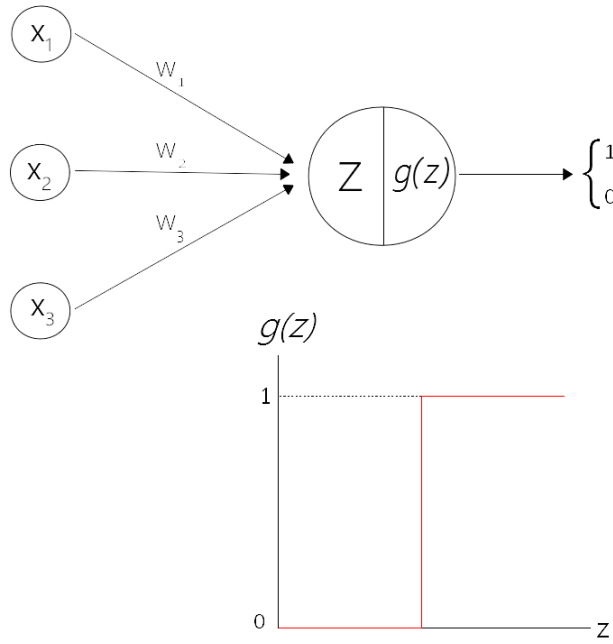


Figure 1: Diagram of a Perceptron

weighted sum, z goes through an activation function $g(z)$, that outputs 1 if the sum exceeds a certain threshold, 0 otherwise. For instance, we can think of inputs as values for technical indicators such as the RSI, the MACD and the stochastic oscillator and the output as 1 if we should "buy" and 0 if we should "sell". We take any data point and extract the values for the inputs and for the output, we compute the weighted sum and if that number exceeds the threshold, the Perceptron will output 1, or "buy".

To learn how to accurately classify data, the Perceptron will update each weight such that we eventually reach a point where every data point is associated to the correct class. If the Perceptron outputs "buy", and the correct position was indeed buying, we can repeat the procedure for another data point. If the correct position was to sell, we update the weights and then we move on to the next data point. Weights are updated whenever the Perceptron outputs the wrong value, and kept the same whenever the Perceptron outputs the correct value. As weights are the only elements that can be manipulated, by constantly updating them we will eventually come to the right combination of weights, which outputs the correct buy or sell output for every data point. The weights are then no longer updated and the Perceptron has finished training. For each new set of values for the technical indicators, the weights will be such that the weighted sum of inputs activate the neuron only when necessary, or when buying is the right call, and stay deactivated when the right position is to sell. The weights update is not random, we desire to increment weights in order to increase the accuracy at each update. In order to do so, we can calculate a loss function, the loss will be 0 if the output corresponds to the target value and different from 0 if the Perceptron fails. This loss function represents the error, or how accurate the Perceptron is, with respect to every weights combination. We can minimize this loss function by differentiating it with respect to each weights, which solves for the weights that outputs the minimal error⁴. This closely relates to Ordinary Least Squares where we minimize the Mean Squared Error in order to find coefficients that minimize the

⁴0 ideally, but in reality the minimal error lies somewhere between 0 and 1

squared error. Although with OLS, one can find the solution to the differential equation analytically, it is often impossible with Perceptrons, and ANN, as they do not have closed-form solutions. To find the correct weights we are required to use numerical methods, and update them gradually towards the minimum instead of directly deriving weights that minimize the loss function.

Gradient descent is a numerical algorithm to help find the local minimum of a function (denoted as C), assuming it is differentiable. Starting with random initial values for the weights we compute the vector of partial derivatives of the loss function with respect to every weights, also called the gradient vector of the loss function ($\nabla(C)$). This vector of partial derivatives will point in the direction of steepest ascent where any increment in the weights results in the largest possible step on the curve. If for a particular set of weights we land at a point a on the loss function, the gradient tells us in which direction the loss function increases the most or which new combination of weights will land on a' as far up as possible from a . As our objective is to find the minimum of the loss function, we take the inverse of the gradient vector, which then leads to the the direction of steepest descent. The inverse gradient tells us how to update weights so that the loss function decreases the most rapidly towards a local minimum. Once we updated the weights, we repeat the procedure and compute the gradient vector of the loss function evaluated at a the new weights combination. We take a step in the direction of steepest descent (inverse of the gradient vector at that point) and this procedure is iterated until we reach a minimum. At a local minimum all the partial derivatives are 0, hence the gradient vector is a null vector. Less formally, once we reach a local minimum, the gradient descent tells us that we should move by 0 in every direction, which is equivalent to not moving at all.

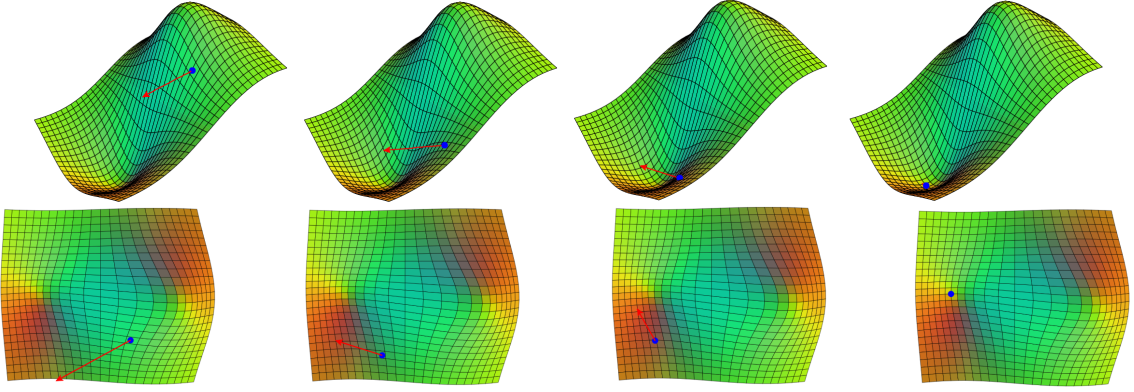


Figure 2: Gradient Descent

For each data point we compute the error and the gradient at that point, take the inverse and increment the weights in the direction pointed by the inverse gradient. By always choosing the direction of steepest descent we will eventually converge to a local minimum, given that the step size is adapted.

$$\overrightarrow{w^{new}} = \overrightarrow{w^{old}} + \eta(-\nabla(C)) \quad (1)$$

Where $\overrightarrow{w^{new}}$ is the vector of updated weights, $\overrightarrow{w^{old}}$ is the vector of old weights,

η is the learning rate which determines the size of the step and $\nabla(C) = \begin{bmatrix} \frac{\partial C}{\partial w_j} \\ \dots \\ \frac{\partial C}{\partial w_J} \end{bmatrix}$ is the

vector of partial derivatives of the loss function (C) with respect to each weights in $\overrightarrow{w^{old}}$. We often determine stopping rules that stop the gradient descent algorithm instead of reaching a point where all partial derivatives are 0. The algorithm is efficient nonetheless, the gradient vector does not point towards the minimum but the direction in which the loss function decreases. If we take a step in that direction once we are really close from the minimum, we might overshoot and land beyond the minimum, therefore if no stopping rules are applied the algorithm keeps iterating indefinitely. Once we reach a certain number of iterations, or that we are satisfied with the model's accuracy, we can voluntarily stop the procedure and get values for the weights that minimize the loss function. It is a possibility that the minimum identified by the gradient descent algorithm is a local minimum, while we ideally want to find the global minimum of the loss function. By construction, gradient descent is not able to tell whether we attained a local or global minimum which can lead to significant inaccuracies in the model. This is one of the crucial drawbacks in gradient descent.

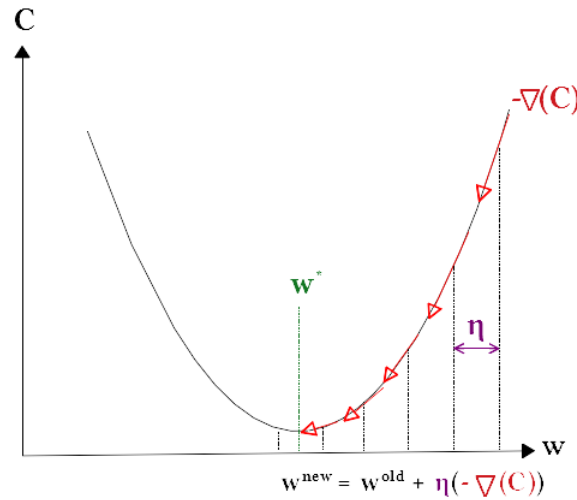


Figure 3: Overshooting the local minimum

In the simple example above, we see how inappropriate step size η can lead the algorithm to oscillate around the local minimum of the loss function which is found at weight w^* . The algorithm will never converge as weight will "bounce" from right to left, when the weight is too large the gradient points towards the left and because the step size is too large, we land beyond the minimum. The weight is then too small, the gradient points towards the right and we take a step in that direction landing beyond w^* . The algorithm will iterate indefinitely as we never reach w^* where $\nabla(C)$ is null. This shows the importance of determining stopping rules such as number of iterations or minimum acceptable accuracy to prevent infinite iterations.

The Perceptron is limited to linearly separable problems, where it is possible to find a linear decision frontier in the input space that separates classes. In two dimensions, the frontier is a straight line, in three dimensions it would be a plane and for all further dimensions the frontier is a hyperplane. Once the neuron has been trained and we

the final weights are obtained we can calculate $b(x) = W^T x + w_0$ to trace a decision boundary. Every point falling on a particular side of the boundary will be associated with the corresponding class and any point falling on the other side of the frontier will be associated to the other class. The diagram below illustrates two different problems, the first one which is linearly separable and the second one which is not. Because Perceptron is a linear classifier, we need to upgrade it in order to solve non linear problems such as the one presented on the diagram above.

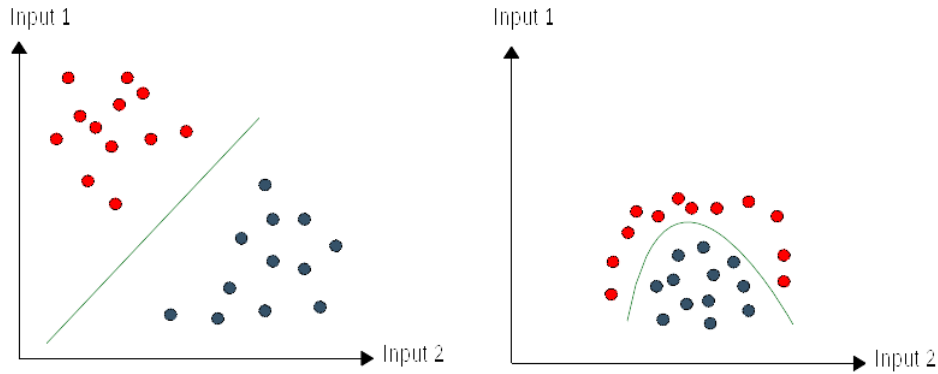


Figure 4: Linear separability

2.1.2 Networks & Backpropagation

The use of Neural Networks allows us to find non-linear decision boundaries and separate classes that could not be identified by a Perceptron. Neural Networks are a series of Perceptrons interconnected with each other. They consist in one entry layer, one output layer and intermediate "hidden layers" which size may vary. The size of a layer is determined by the number of neurons and the ultimate size of the network is determined by the size of each layer multiplied by the number of layers. The following example is a Neural Network with one entry layer of 3 neurons, two hidden layers of 3 neurons each and one output layer of 1 neuron.

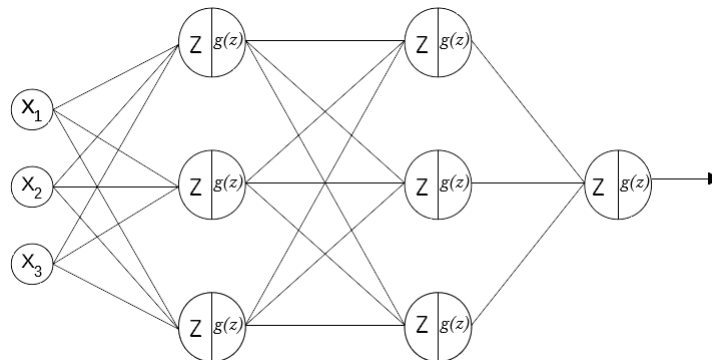


Figure 5: Neural Network with two hidden layer

The core principles are identical to the Perceptron, although one major modifi-

cation has to be made. The gradient descent algorithm requires the ability to compute all partial derivatives with respect to each weights. Because the first weights are passed through an activation function before reaching the next layers, we need this activation function to be differentiable⁵. In that case, we can apply the chain rule to compute partial derivatives from one end of the network to the other. Let us first define the notation.

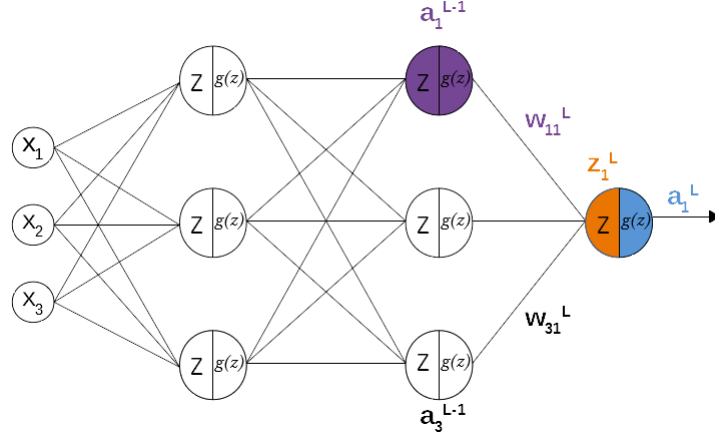


Figure 6: simple network

We associate each weight w to the two nodes it is connecting such that w_{jk}^L connects node j in the previous layer to node k in the current layer. The current layer is denoted as L and previous layers are denoted as $L-1, L-2, \dots, L-M$. We also define the loss function as $C_n, n = 1, 2, \dots, N$, N being the sample size and the activation function as g . The weighted sum of inputs arriving to node j is denoted as z_j^L and the output from node j is denoted as a_j^L . To use gradient descent we need to compute the gradient of the loss function to minimize. The gradient is therefore expressed as:

$$\nabla(C) = \begin{bmatrix} \frac{\partial C_n}{\partial w_{j,k}^L} \\ \dots \\ \frac{\partial C_n}{\partial w_{j,K}^{L-m}} \end{bmatrix} \quad (2)$$

We will show how to compute the gradient for one particular data point $n = 0$. We need to find the partial derivative of the loss function with respect to each weight in the network. For this example we will use the Mean Squared Error but other functions are appropriate to classification problems. This will be a two step procedure; first we compute partial derivatives with respect to weights on the last layer (output layer $L-1$). Then we generalize the results to every weight in the network. For simplicity, we assume the penultimate layer ($L-1$) contains only one neuron (a_1^{L-1}). The last two layers can then be represented as:

⁵We give further details on why that is when we compute partial derivatives

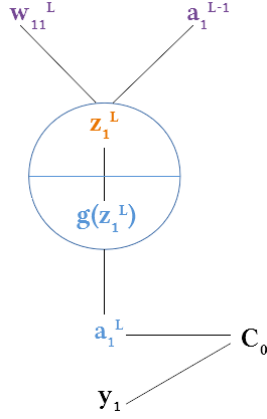


Figure 7: output layer

This diagram is inspired by the work of Grant Sanderson ([3Blue1Brown](#)) on [backpropagation](#).

The diagram shows how the output from node 1 in layer $L - 1$ is multiplied by a particular weight $w_{1,1}^L$ to obtain z_1^L , $z_1^L = \sum w_{i,j}^L * a_j^{L-1}$. z_1^L is passed through the activation function $g(z_1^L)$ and produces the output for node 1 in layer L , a_1^L . Finding partial derivatives with respect to weights on the last hidden layer $L - 1$ to the output layer L is straightforward because we can decompose the effect of weight $w_{1,1}^L$ with the chain rule:

$$\frac{\partial C_0}{\partial w_{1,1}^L} = \frac{\partial C_0}{\partial a_1^L} * \frac{\partial a_1^L}{\partial z_1^L} * \frac{\partial z_1^L}{\partial w_{1,1}^L} \quad (3)$$

The partial derivative of the loss function with respect to the weight $w_{1,1}^L$ is equal to the product of the partial derivative of the loss function with respect to the output of the last layer, the partial derivative of the output of the last layer with respect to the weighted sum z_1^L and the partial derivative of the weighted sum with respect to the weight $w_{1,1}^L$. because the loss function is a function of the output given by the network, which is a function of the activation of the previous node which in turn, is a function of the weights, the partial derivative of the loss function is equal to the product of the partial derivatives of the functions that compose it. We yet have to compute these derivatives.

- $\frac{\partial C_0}{\partial a_1^L}$, the loss function is equal to: $(a_j^L - y_j)^2$. Therefore its partial derivative is: $2(a_j^L - y_j)$.
- $\frac{\partial a_1^L}{\partial z_1^L}$, the output a_1^L is the result of the activation function evaluated at z_1^L . We can also write; $a_1^L = g(z_1^L)$ Therefore its partial derivative is equal to the partial derivative of the activation function, $g'(z_1^L)$. This is the reason why we need differentiable activation functions such as the sigmoid.
- $\frac{\partial z_1^L}{\partial w_{1,1}^L}$, z_1^L is the weighted sum of the previous output multiplied by the weight $w_{1,1}^L$, a slight change in the weight will cause z_1^L to change by the value of the previous

output, a_1^{L-1} . The partial derivative of z_j^L with respect to the weight $w_{j,k}^L$ is equal to the previous output⁶.

We were able to compute one element of the gradient vector, for one particular training example:

$$\frac{\partial C_0}{\partial w_{1,1}^L} = 2(a_j^L - y_j) * g'(z_i^L) * a_1^{L-1} \quad (4)$$

Recall that we are using MSE to measure loss, MSE is the average error on all data points hence we need to compute a value for each data point $n = 0, \dots, N$. The gradient of the average loss is equal to the average gradient of the loss for each example therefore we repeat the process for each data point in the training set to obtain one element of the gradient vector which will be the average gradient for all points:

$$\frac{\partial \bar{C}}{\partial w_{j,k}^L} = \frac{1}{N} \sum_{n=1}^N \frac{\partial C_n}{\partial w_{j,k}^L} \quad (5)$$

This is only true for the partial derivatives with respect to weights on the last layer $L - 1$, because the loss function is a direct function of the output from the last layer $C = C(a_j^L, y_j)$. However, it is not a direct function of the outputs from previous layers. This is where backpropagation is used, to rewrite $\frac{\partial C_0}{\partial a_j^{L-m}}$, $m \neq 0$ such that C_0 becomes a function of a_j^{L-m} using the chain rule. We want to solve:

$$\frac{\partial C_0}{\partial w_{j,k}^{L-1}} = \frac{\partial C_0}{\partial a_j^{L-1}} * \frac{\partial a_j^{L-1}}{\partial z_j^{L-1}} * \frac{\partial z_j^{L-1}}{\partial w_{j,k}^{L-1}} \quad (6)$$

$\frac{\partial a_j^{L-1}}{\partial z_j^{L-1}}$ and $\frac{\partial z_j^{L-1}}{\partial w_{j,k}^{L-1}}$ are derived exactly like if they were on the last layers as a_j^{L-1} , z_j^{L-1} and $w_{j,k}^{L-1}$ are direct functions of each other ($a_j^{L-1} = g(z_j^{L-1})$ and $z_j^{L-1} = \sum w_{j,k}^{L-1} * a_j^{L-2}$). We need to rewrite $\frac{\partial C_0}{\partial a_j^{L-1}}$ to be able to compute the gradient and identify the effect of weight $w_{j,k}^{L-1}$ on the loss function. Although the loss function is not a direct function of nodes that are not on the last layer, every node plays a role in the activation of all further neurons. Eventually, each neurons impacts the last output neuron, which in turn is a variable on which the loss function directly depends. The idea of backpropagation is to take advantage of the structure of the network to reverse engineer the propagation of the activation of neurons to be able to compute partial derivative of the loss functions with respect to all weights in the network. We would like to remind the reader that with all the partial derivatives we can calculate the gradient and optimize the weights in order to minimize the loss function. The first step to reverse engineer the network is to find $\frac{\partial C_0}{\partial a_j^{L-1}}$ by applying the chain rule on that expression:

⁶We use indices j and k as this is true for every weight

$$\frac{\partial C_0}{\partial a_j^{L-1}} = \sum \left[\frac{\partial C_0}{\partial a_j^L} * \frac{\partial a_j^L}{\partial z_j^L} * \frac{\partial z_j^L}{\partial a_j^{L-1}} \right] \quad (7)$$

We already demonstrated how to calculate $\frac{\partial C_0}{\partial a_j^L} = 2(a_j^L - y_i)$ and $\frac{\partial a_j^L}{\partial z_j^L} = g'(z)$. We will now focus on the new term:

$\frac{\partial z_j^L}{\partial a_j^{L-1}}$, the partial derivative of the weighted sum with respect to previous output is equal to the weight $w_{j,k}^L$ and $\frac{\partial C_0}{\partial a_j^{L-1}} = 2(a_j^L - y_i)g'(z)w_{j,k}^L$. Overall we were able to calculate the impact of the weights of the previous layer by applying a chain rule inside another chain rule:

$$\frac{\partial C_0}{\partial w_{j,k}^{L-1}} = \underbrace{\frac{\partial C_0}{\partial a_j^{L-1}}}_{2(a_j^L - y_i)g'(z)w_{j,k}^L} * \frac{\partial a_j^{L-1}}{\partial z_j^{L-1}} * \frac{\partial z_j^{L-1}}{\partial w_{j,k}^{L-1}} \quad (8)$$

This is where the idea of backpropagation stems from, to compute partial derivatives with respect to previous weights, one needs to compute the partial derivatives with respect to the last weights. The gradient must be calculated on the last layers first, and then the process can be expanded backward until the beginning of the network using the chain rule. The computation of partial derivatives is repeated for each data point and the average derivative is calculated $\frac{\partial \bar{C}}{\partial w_{j,k}^L} = \frac{1}{N} \sum_{n=1}^N \frac{\partial C_n}{\partial w_{j,k}^L}$ before adding it to the gradient vector. After iterating through the entire data set, the gradient vector is used to update weights such that we take a step towards a local minimum of the loss function.

$$\overrightarrow{w^{new}} = \overrightarrow{w^{old}} + \eta(-\nabla(C))$$

The same stopping rules we discussed before are applied here to stop the gradient descent. We have seen how Neural Networks treat the information using a simplified example involving a Perceptron. We also explained how the model finds the optimal parameters by gradually minimizing the loss function by taking iterative steps following the gradient (inverse gradient). We then generalized the results to Neural Networks and showed how backpropagation can be used to solve the same problem with more complex networks. Until now we have considered vanilla examples, the following section deepens our understanding of various adjustments that can reinforce the network.

2.1.3 Activation functions and enhanced gradient descent

2.1.4 Support Vector Machine & the kernel trick

We are trying to find the maximum margin, in other words we want to maximize the width of the margin that is calculated as:

$$width = (X_+ - X_-) \frac{w}{||w||} \quad (9)$$

This is equivalent to maximizing $\frac{w}{||w||}$ or equivalently to minimize $||w||$. For mathematical we further modify the expression above to arrive at the following minimization problem: the largest margin is attained when $\frac{1}{2}||w||^2$ is minimum. This minimizing problem is constrained by: $y_i(X_i w + b) = 1$ as discussed earlier. Hence we are looking to solve:

$$\begin{aligned} & \min \frac{1}{2} ||w||^2 \\ \text{s.t. } & y_i(X_i w + b) - 1 = 0 \end{aligned} \quad (10)$$

This is a standard constrained optimization problem. We will find the maximum margin using a Lagrangian defined as:

$$\mathcal{L} = \frac{1}{2} ||w||^2 - \sum \alpha_i (y_i (w \cdot X_i + b) - 1) \quad (11)$$

To solve a Lagrangian we start by solving the partial derivatives with respect to each variable (α_i , b and w).

- $\frac{\partial \mathcal{L}}{\partial w} = w - \sum \alpha_i y_i X_i = 0 \longrightarrow w = \sum \alpha_i y_i X_i$
- $\frac{\partial \mathcal{L}}{\partial b} = - \sum \alpha_i y_i = 0 \longrightarrow \sum \alpha_i y_i = 0$

We can plug these two equalities in \mathcal{L} :

$$\mathcal{L} = \frac{1}{2} \underbrace{\left(\sum \alpha_i y_i X_i \right) \left(\sum \alpha_j y_j X_j \right) - \sum \alpha_i y_i X_i \cdot \left(\sum \alpha_j y_j X_j \right)}_{=\frac{1}{2}a - a = -\frac{1}{2}a, \text{ with } a = \sum \alpha_i y_i X_i} - \underbrace{\sum \alpha_i y_i}_{=0} + \sum \alpha_i$$

We arrive at the final result:

$$\mathcal{L} = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j X_i \cdot X_j \quad (12)$$

This result tells us that the optimization problem depends on the dot product of pairs of data points $X_i \cdot X_j$. To find the optimum, we now have to solve for α_i . Which means every pair of data points i and j will be weighted by some α_i . As we will discuss, this indicates that finding which $\alpha_i \neq 0$ is equivalent to finding support vectors. Once we find the support vectors (points on the marginal hyperplanes), we can easily find the intercept b and the slope of the decision boundary w .

At this point we would like to take a step back to illustrate the importance of the α_i term. Before we express w and b as sums of other variables we had:

$$\mathcal{L} = \frac{1}{2}||w||^2 - \sum \alpha_i \underbrace{f_i(w, b, X_i)}_{y_i(w \cdot X_i + b) - 1}$$

In order to maximize this expression with respect to α_i we would like $\sum \alpha_i f_i(w, b, X_i)$ to be as small as possible. We know from previous constraints that $\sum \alpha_i f_i(w, b, X_i) \geq 0$ hence, to maximize \mathcal{L} we need $\sum \alpha_i f_i(w, b, X_i) = 0$. This is possible in two different scenarios⁷:

1. $f_i(w, b, X_i) > 0$, which implies $\alpha_i = 0$
2. $f_i(w, b, X_i) = 0$, which implies $\alpha_i > 0$

These two scenarios tell us that for points X_i on the margin ($f_i(w, b, X_i) = 0$), α_i is strictly positive, for points not on the margin ($f_i(w, b, X_i) > 0$), $\alpha_i = 0$. Because w is defined as $\sum \alpha_i y_i X_i$, w only exists (w is not null) for points on the margin. For every other points, α_i is equal to 0 then, $\alpha_i y_i X_i = 0$ and w at that particular point i is null.

To maximize the margin, SVM only has to look at points on the margin, the other points will be irrelevant in the optimization problem. In other terms, once we solve for α_i , whenever $\alpha_i \neq 0$; the point is relevant and is on the margin. Solving for α_i is equivalent to finding points on the margin (support vectors), hence finding margins.

We then only need to differentiate \mathcal{L} with respect to α_i by setting the partial derivatives equal to 0:

$$\max_{\alpha_i} \mathcal{L} = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j X_i \cdot X_j \quad (13)$$

We saw that optimizing this expression required that α_i is not null only for points on the margin. Therefore, once we set partial derivatives equal to 0, we will find optimal α_i which, plugged into $w = \sum \alpha_i y_i X_i$ will give us the w vector, orthogonal to the decision boundary with the largest margins.

We also know that $w^T X_i + b = y$, or $b = y - X_i w^T$. This strict equality is only true when X_i is on the margin. Hence for a subset s of i , $s \subset i$, for which X_s is on the margin, we have $b = y - X_s w^T$. Since we solved for X_s by solving for α_i we can compute b . The classification problem is solved, for any new point u we have:

$$\begin{aligned} &+ \text{ if } w \cdot u + b > 0 \\ &- \text{ if } w \cdot u + b < 0 \end{aligned} \quad (14)$$

w and b are known so we can classify every new point according on the hyperplane we found by solving for w and b .

⁷Both terms cannot be equal to 0

2.1.5 Soft Margin SVM

2.2 Finance

3 Application

4 Conclusion