

Notes on Math and Architectures of Deep Learning

Benjamin Santana

April 2023

1 Introduction

Before jumping into deep learning, let us talk about machine learning a bit. What is machine learning? Living creatures are taking decisions at every moment. Say for a example a dog sees an object and has to take a decision, does it get away from it, does it runs to towards it, or does it simply ignore it.

This will depend on certain factors about the object, is it sharp? does it looks cozy? Believe it or not the dog has a machine inside its head, taking this decisions. Its brain.

Machine learning is trying to emulate this machine, taking different parameters and output a decision.

Computers already do this, but the new thing is that now we do not give the computer step by step instructions on how to take the decision. Instead we develop a mathematical model for the problem.

Let us continue with the dog example I was taking about (which I completely stole from the book).

We have two *features*, hardness and sharpness.

We could define a *parameterized function* with a *weighted sum of inputs*:

$$y(\text{hardness}, \text{sharpness}) = w_0 \cdot \text{hardness} + w_1 \cdot \text{sharpness} + b \quad (1)$$

The weights w_0 , w_1 and the bias b are the parame-

ters of the function. The output y can be interpreted as a threat score.

The weights are not know at first, we need to estimate them, this is what we call *model training*.

The steps to solve any machine learning problem are the following:

1. We first design a parameterized model function (e.g. weighted sum) with unknown parameters (weights). This makes the *model architecture*.
2. Estimate the weights via model training.
3. Now we have a complete model, we can send inputs not seen before and it should work, this is called *inference*.

In supervised learning, us humans will create the training data, this means creating example inputs, each corresponding with the desired output.

The training the model part is iterating through different weight values and checking with the training data to see if we are close to the desire output.

We optimize the way we look for those weights, we will see that in detail in the future.

At first, our model will be dumb, it will not be able to perform well, like a baby. But once we have run some iterations and we correct it, it will start performing better.

2 Simple Model

The book used a simple example on how to apply machine learning to the dog problem we stated above, in the book it was a cat but I am more of a dog person so...

We have a hypothetical dog which only needs to take one decision in life, run away from an object, ignore it, or approach it. It takes this decision based only on two quantitative inputs.

2.1 Input Features

$x_0 = \text{hardness}$, $x_1 = \text{sharpness}$

we can normalize this input

$$v_{norm} = \frac{(v - v_{min})}{(v_{max} - v_{min})} \quad (2)$$

where v_{min} is the min value possible and v_{max} is the max value possible.

This will generate values between 1 and 0, thus we can define the input as

$$v \in [v_{min}, v_{max}] \rightarrow v_{norm} \in [0, 1]$$

A single input becomes

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}, \in [0, 1]^2 \quad (3)$$

2.2 Output Decisions

We will have our model estimate a threat score (y) and based on a threshold (δ) we will assign a class to it. If $y > \delta$ high threat, run away. If y is around 0, do nothing. If $y < \delta$ negative threat, approach.

2.3 Model Estimation

We need to estimate the function which transforms the input vector to the output. This consists on two

parts

- model architecture selection
- model training

2.3.1 Model Architecture Selection

Our model has three parameters w_0, w_1, b they can be represented as

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}, \quad (4)$$

and b as a scalar. All three being real numbers.

We will define a super simple function such as

$$y(x_0, x_1) = w_0 x_0 + w_1 x_1 + b = \vec{w}^T \vec{x} + b \quad (5)$$

Where b is the bias.

$$y(\vec{x}) = \vec{w}^T \vec{x} + b \quad (6)$$

2.3.2 model training

We defined $y(\vec{x})$, now we have to find the best weights for the function. We will choose the parameters so that the outputs on the training data inputs, match the corresponding outputs as close as possible.

How do we calculate the error?

On the i^{th} training data pair $(\vec{x}(i), y_{gt}(i))$ (gt ground truth) we can say $y_{predicted}(i) = \vec{w}^T \vec{x}(i) + b$, then we compare $y_{predicted}(i)$ with y_{gt}

$$e(i)^2 = (y_{predicted}(i) - y_{gt})^2 \quad (7)$$

It is squared to take into account negative values.

This way we calculate the error of one iteration, the overall error (aka loss) is the sum of all $e(i)^2$

$$E^2 = \sum_{i=0}^{i=N} e(i)^2 \quad (8)$$

our goal is to find the \vec{w} that minimizes E^2 .

The ideal goal would be to actually get E^2 to 0, but in practice, one often stops when the results are accurate enough to solve the problem we are presented.

2.4 Geometrical View of Machine Learning

In the model specified above, each input is a vector with 2 values x_0 and x_1 (hardness and sharpness). If we map this into a 2D space, each input is a point on the plane. Usually it would have more dimensions, but this being a super simple problem only have two. This space is called *feature space*.

The model will work as a translator, that takes points in the feature space and map them to another space with less dimensions called, the *output space* (they really break their heads thinking for a name for this last one).

It is supposed to be easier for the computer to take decisions in the *output space*

2.5 Regression vs Classification in ML

We have two types of machine learning models: *regressors* and *classifiers*.

A regressor attempts to provide us with a value based on a specific input, while classifiers aim to assign a class to which the input belongs.

It is possible to create a regressor and then, based on a threshold, select a class. Alternatively, a model can directly output one of these predefined classes.

2.6 Linear vs Nonlinear Models

Usually when we try and classify points, they are not in a clear straight line. Meaning that Function 9 not going to cut it. This is also shown in Figure 1.

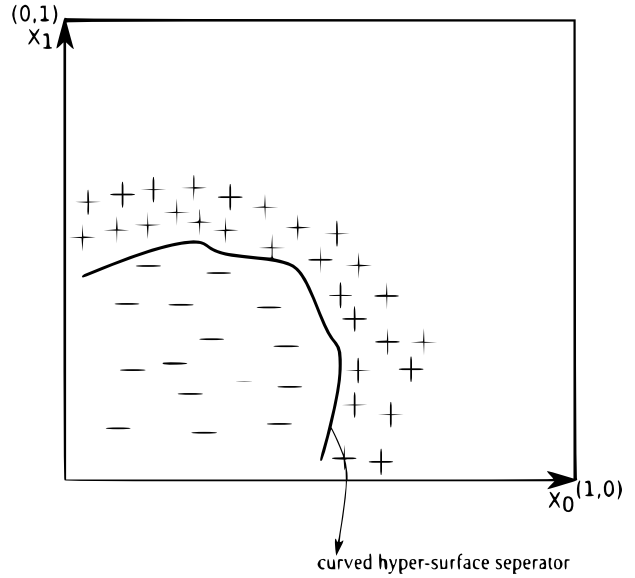


Figure 1: Feature map where we cannot use a linear function to classify the inputs

$$y(\vec{x}) = \vec{w}^T \vec{x} + b \quad (9)$$

We will need a non-linear function to classify these points. A popular function in machine learning is the Sigmoid function, (named because of its S shape). Shown in Equation 10

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (10)$$

Thus we can use the model architecture shown in Equation 11. It is the building blocks of a neural network.

$$y = \sigma(\vec{w}^T \vec{x} + b) \quad (11)$$