

1 Machine Learning Theory

Throughout this section the theory of machine learning that will be used in this thesis is described and explained.

1.1 Motivation

It can be difficult for humans to recognize certain patterns and trends in data. This becomes more difficult the greater the quantity of the data is, which is becoming more and more common with the rapidly growing topic of *Big Data*. For this reason, computers are often used instead of humans to recognize patterns and trends in the data by analyzing the data, which is what is called *Machine Learning*. In this thesis, we will use machine learning in section **MANGLER REFERENCE** to develop a model to estimate the 2D pose of a single human in an image. Later, in section **MANGLER REFERENCE**, we will use machine learning to improve our understanding of the model.

1.2 Machine Learning Paradigms

Machine learning is usually split into the following three paradigms

1. *Supervised learning* where the data consists of features and labels. By analyzing the data the algorithm learns to predict the labels given the features [7]. Supervised learning is further split into *classification* and *regression*. If the value of each label is limited, then the task is a classification task. If the value of each label is not limited, then the task is a regression task.
2. *Unsupervised learning* where the data only consists of features. The algorithm then learns properties of the data, without any provided labels [7].
3. *Reinforcement learning* where the algorithm learns to perform the action in a given environment that yields the highest reward [1].

In this thesis we will make use of supervised learning when developing our model for pose estimation. Later, unsupervised learning is used when we explore our developed model.

1.3 Evaluation of Machine Learning Models

When developing a machine learning model it is important to know how trustworthy the developed model is. This is usually done by testing how good the model is at generalizing unseen data, which is done by making use of *evaluation metrics*.

1.3.1 Splitting the dataset

When developing a machine learning model, the data needs to both create the model, but also to evaluate the model. For the evaluation of the model, one of the two following techniques is usually used

1. *Cross validation* where the data is split into K random non-overlapping chunks of equal size. The model is then trained for K rounds on $K - 1$ of the chunks, where the last chunk is used for evaluating the model. After each round the parameters of the model is reset to ensure one round does not affect another round. After the K rounds the average loss of the K rounds is the loss of the model [6].

2. *Train-validation-test* where the data is split into 3 random non-overlapping chunks. The training dataset is then used for training the model and the validation dataset is used for evaluating the model as it is being developed - this often means, that the *hyperparameters*, the parameters that are not possible to fit from the data, are being tweaked to yield the best validation loss. Lastly, the testing dataset is used as a final evaluation of the model to yield an unbiased evaluation of the model. Once the testing dataset has been used it can no longer be used for evaluating the data, as this ensure an unbiased evaluation [4].

Throughout this thesis the train-validation-test technique will be used over cross validation for evaluating the developed models. This is done, since cross validation is better suited for smaller datasets, as the runtime is much greater than the runtime of the train-validation-test technique.

1.3.2 Evaluation Metrics for Supervised Machine Learning (Loss Functions)

When we have trained a model, we need to somehow evaluate how well the model performs on unseen data. This is usually done by making use of evaluation metrics or *loss functions*. There are many different loss functions, each with their own advantages and disadvantages. One of the most common loss functions for regression is the *Mean Squared Error (MSE)*, defined as

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

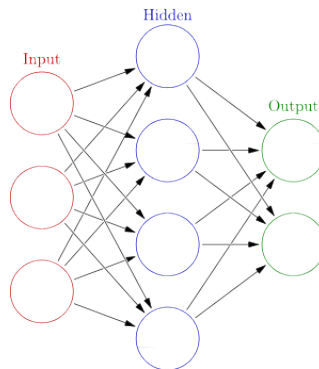
where y_i is the true value of the i th observation and \hat{y}_i is the estimated value of the i th observation. Thus, *MSE* measures the average squared difference between the true observation and the estimated observation. The aim of a model is thus to make the *MSE* as small as possible [2].

1.4 Neural Networks

In recent years *deep learning* and *neural networks* have revolutionized the use of machine learning. In this thesis a neural network will be used for performing the human pose estimation.

1.4.1 The Mathematics Behind Neural Networks

Figure 1: Visualization of a feedforward neural network with a single hidden layer [8]



The architecture

One of the most common types of neural networks are *feedforward neural networks*, where the data flows unidirectionally through the network. Such a network is visualized in Figure 1. The network is built up of three types of components: the *input layer*, the *hidden layers* and the

output layer. Each layer is built up of *units*, also called *neurons* (which are visualized as circles in Figure 1), where each neuron has a *bias* assigned to it, and is connected to one or two other layers through *edges* (which are visualized as arrows in Figure 1), where each edge has a *weight* assigned to it. Hidden layers are connected to two other layers - one before the hidden layer and one after the hidden layer - where the input layer is only connected to the next layer in the network and the output layer is only connected to the previous layer in the network.

We can define the network mathematically by letting $h_n^{(i)}$ denote the value of the n th node in the i th layer, $w_{m,n}$ denote the value of the weight of the edge connecting the n th node in the i th layer to the m th node in layer $i + 1$ and $b_n^{(i)}$ denote the bias corresponding to the n th node in the i th layer.

When data flows through the model it follows the following formula

$$\mathbf{h}^{(i+1)} = g^{(i+1)} \left(\mathbf{W}^{(i+1)} \mathbf{h}^{(i)} + \mathbf{b}^{(i+1)} \right)$$

where $\mathbf{W}^{(i+1)}$ is the weights between layer i and layer $i + 1$ defined by

$$\mathbf{W}^{(i+1)} = \begin{pmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,0} & w_{m,1} & \cdots & w_{m,n} \end{pmatrix},$$

$\mathbf{h}^{(i)}$ is the values of the nodes in the i th layer defined by

$$\mathbf{h}^{(i)} = \begin{pmatrix} h_0^{(i)} \\ h_1^{(i)} \\ \vdots \\ h_n^{(i)} \end{pmatrix},$$

$\mathbf{b}^{(i+1)}$ is the values of the biases of layer $i + 1$ defined by

$$\mathbf{b}^{(i+1)} = \begin{pmatrix} b_0^{(i+1)} \\ b_1^{(i+1)} \\ \vdots \\ b_m^{(i+1)} \end{pmatrix}$$

and g is an *activation function*, that is typically applied element-wise [3] [5]. One often used activation function is the *rectified linear activation function* (or *ReLU* for short) defined by

$$g(x) = \max\{0, x\}.$$

The ReLU-function is very close to being linear, making the function keep many of the properties of linear functions that make them easy to optimize, as well as keep many of the properties that make linear models well at generalizing, which are two great advantages of using the ReLU-function. Another great advantage of using the ReLU-function is stated by the *universal approximation theorem* that states, that a feedforward network with a linear output layer and at least one hidden layer with the ReLU-function (or another activation function from a wide class of activation functions) can approximate any continuous function on a closed and bounded subset of \mathbb{R}^n , as long as the network has enough hidden neurons [3].

Backpropagation

1.4.2 Termonology

Overfitting and Regularization

The main goal of a machine learning model is to generalize well on unseen data. This can often be difficult, as the model simply "remembers" the training data instead of learning the patterns in the training data. In other words, the gap between the training error and the test error is too large, which is a concept called *overfitting*. Certain techniques are designed to reduce the test error - these techniques are collectively called *regurlization* [3].

Mini-batch

Batch Normalization

Epoch

Optimizer

1.4.3 Convolutional Neural Networks

NN Upsampling

Maxpooling

1.4.4 Stacked Hourglass

Reasoning behind using the Stached Hourglass

The Residual Modules

The Hourglass

The Stacked Hourglass

1.5 Principal Components Analysis and K-means Clustering

1.5.1 Principal Components Analysis (PCA)

1.5.2 K-means Clustering