



Bachelor Thesis

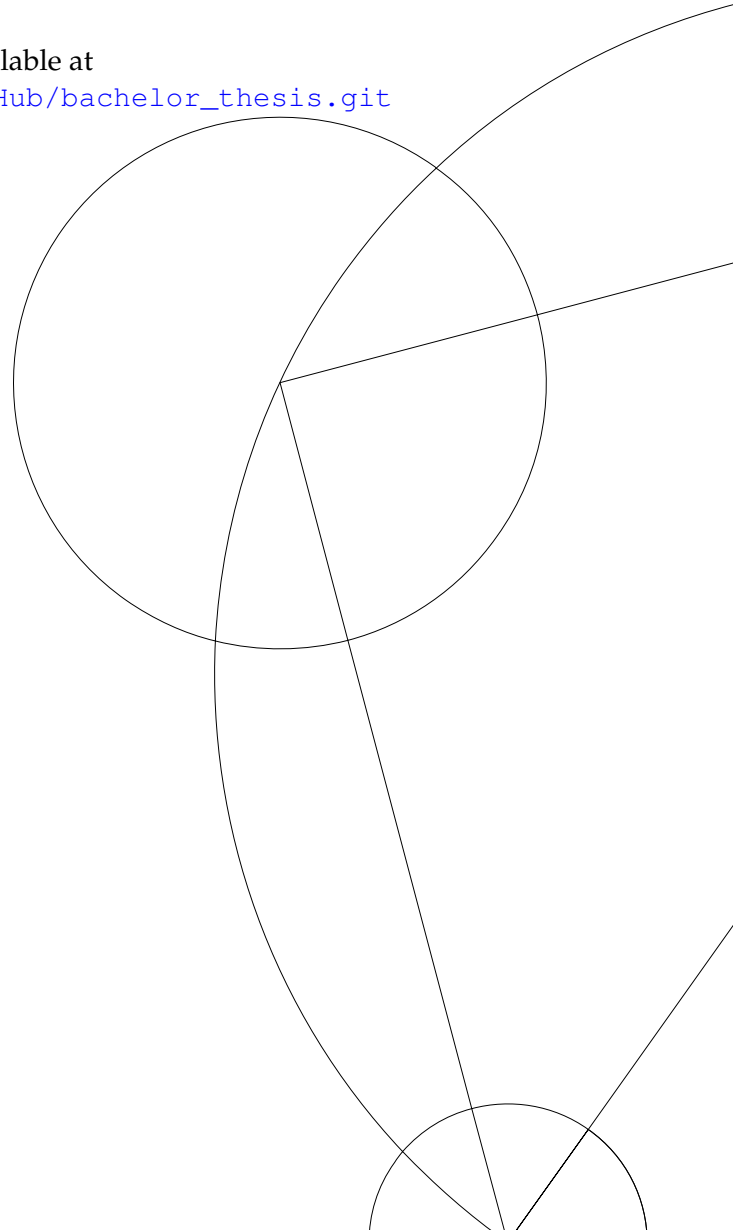
2D Articulated Human Pose Estimation Using Explainable Artificial Intelligence

André Oskar Andersen (wpr684)
wpr684@alumni.ku.dk

April 22, 2021

Source code available at
https://github.com/KeenbitGitHub/bachelor_thesis.git

Supervisor
Kim Steenstrup Pedersen kimstp@di.ku.dk



Contents

| | | |
|----------|---|-----------|
| 1 | Notation | 3 |
| 2 | Machine Learning Theory | 4 |
| 2.1 | Motivation | 4 |
| 2.2 | Machine Learning Paradigms | 4 |
| 2.3 | Evaluation of Machine Learning Models | 4 |
| 2.3.1 | Splitting the dataset | 4 |
| 2.3.2 | Evaluation Metrics for Supervised Machine Learning (Loss Functions) | 5 |
| 2.4 | Neural Networks | 5 |
| 2.4.1 | Feedforward Neural Networks | 5 |
| 2.4.2 | Convolutional Neural Networks | 10 |
| 3 | Stacked Hourglass, t-SNE and K-Means Clustering | 12 |
| 3.1 | Stacked Hourglass | 12 |
| 3.1.1 | Reasoning behind using the Stacked Hourglass | 12 |
| 3.1.2 | The Residual Module | 12 |
| 3.1.3 | The Hourglass | 12 |
| 3.1.4 | The Stacked Hourglass | 13 |
| 3.2 | t-SNE | 13 |
| 3.3 | K-means Clustering | 13 |
| 4 | The Dataset | 14 |
| 4.1 | The COCO Dataset | 14 |
| 4.2 | Data Preprocessing | 15 |
| 4.2.1 | Creating the test dataset | 15 |
| 4.2.2 | Preprocessing the images | 15 |
| 4.2.3 | Handling the labels | 16 |
| 5 | Experiment | 17 |
| 5.1 | Configuration Details | 17 |
| 5.2 | Results | 19 |
| 5.3 | Training Details | 19 |
| 6 | Discussion | 20 |

1 Notation

The notation used throughout this thesis is summarized below

- x : A scalar
- \mathbf{x} : A vector
- \mathbf{X} : A matrix
- \mathbb{R} : The set of real numbers
- $\nabla_x f$: Gradient of f with respect to x
- $\frac{\partial y}{\partial x}$: Partial derivative of y with respect to x
- \mathcal{O} : Big O-notation
- \odot : Element-wise multiplication
- $\mathcal{N}(\mu, \sigma^2)$: Normal/Gaussian distribution with mean μ and standard deviation σ^2

2 Machine Learning Theory

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

2.1 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.

2.2 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 [5]. 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 [5].
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.

2.3 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*.

2.3.1 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 [15].

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 [6].

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.

2.3.2 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 [7].

2.4 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. Throughout subsection 2.4 the theory and mathematics behind neural networks is described and explained.

2.4.1 Feedforward Neural Networks

Algorithm 1 GradientDescent [19]

Require: Learning rate η

Require: Starting position θ

- 1: **while** stopping criterion not met **do**
 - 2: Apply update: $\theta \leftarrow \theta - \eta \nabla f(\theta)$
 - 3: **return** θ
-

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* [4].

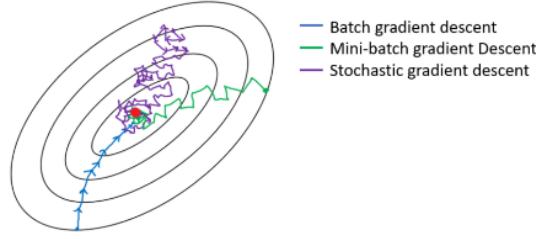


Figure 1: Comparison of batch, mini-batch and stochastic gradient descent [12]

Gradient Descent

The goal of a machine learning model when training is to minimize its loss. There are different methods to do so, however, the most common algorithms are variants of *gradient descent*, whose algorithm is described in Algorithm 1. The algorithm works by taking a learning rate η , a starting position θ and a function f as input, where f is the function to minimize. It then computes the gradient of f with respect to θ , and subtracts the gradient times η from θ . This is done until a stopping condition is met, such as when the magnitude of the gradient $|\nabla f(\theta)|$ is small or until a maximum amount of iterations has been reached [19].

Online, Mini-batch and Batch methods

When gradient descent is used in machine learning, computing $\nabla f(x)$ is usually done by averaging the gradient of each of the n observations of the trainingset, which is called a *batch gradient method* and is computationally inefficient, as the cost is $\mathcal{O}(n)$. It is therefore common to use variants of gradient descents, that reduces the cost of computing the gradient. In *online gradient methods* (also known as *stochastic gradient descent*) a single observation from the dataset is used to compute the gradient, which brings the cost down to $\mathcal{O}(1)$. In *mini-batch gradient methods* a subset of the dataset is used to compute the gradient, making the cost $\mathcal{O}(|\mathcal{B}|)$, where $|\mathcal{B}|$ is the mini-batch size [19].

Choosing the right batch size can be difficult, however, there are a few guidelines which one can follow [12] [19]

1. Looking at Figure 1 we see, that batch gradient descent uses the fewest iterations, however, each iteration takes the longest to compute. On the other hand, in online gradient descent each iteration is the fastest to compute, however, it is also the method that uses the most iterations. Lastly, mini-batch gradient descent combines the two: it uses less iterations than online gradient descent, but more than batch gradient descent, and each iteration takes less time than in the case with batch gradient descent, but longer than in the case with online gradient descent.
2. A batch size that is of power of 2 can offer in better runtime for some hardware. A batch size that is often used for larger models is 16, however, they typically range between 32 and 256.
3. Smaller batch sizes can offer a regularizing effect, as it is difficult for the model to "remember" the complete dataset from batches that does not represent the whole dataset.

Optimization Algorithms

Online, mini-batch and batch gradient descent are all optimization algorithms used for estimating the minimum of a function. One problem of these algorithms is, that the learning rate can be difficult to choose. Therefore, there have been developed a range of various optimization algorithms that uses a separate learning rate for each parameter and automatically adapt these learning rates. One of which is `RMSPprop`, which has been visualized in Algorithm 2. The algorithm works by using an decaying average that discards knowledge from the past, so

Algorithm 2 RMSProp [4]

Require: Learning rate η

Require: Decay rate ρ

Require: Starting position θ

Require: Small constant δ , usually 10^{-6}

- 1: Initialize accumulation variables: $\mathbf{r} \leftarrow \mathbf{0}$
 - 2: **while** stopping criterion not met **do**
 - 3: Sample a minibatch of m observations from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$
 - 4: Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 - 5: Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$
 - 6: Compute parameter update: $\Delta \theta = -\frac{\eta}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$
 - 7: Apply update: $\theta \leftarrow \theta + \Delta \theta$
 - 8: **return** θ
-

that it can converge after finding a convex bowl. The algorithm uses a hyperparameter ρ , that controls the length scale of the moving average [4].

Batch Normalization

Batch normalization is a reparametrization method, which is applied to individual layers. If $\mathbf{x} \in \mathcal{B}$ is an input to the batch normalization, BN, then batch normalization is done by the following

$$\text{BN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \hat{\boldsymbol{\mu}}_{\mathcal{B}}}{\hat{\boldsymbol{\sigma}}_{\mathcal{B}}} + \boldsymbol{\beta}$$

where

$$\hat{\boldsymbol{\mu}} = \frac{1}{|\mathcal{B}|} \sum_i \mathcal{B}_i$$

and

$$\hat{\boldsymbol{\sigma}} = \sqrt{\epsilon + \frac{1}{|\mathcal{B}|} \sum_i (\mathcal{B}_i - \boldsymbol{\mu})^2},$$

which makes the minibatch have 0 mean and unit variance. γ and $\boldsymbol{\beta}$ are then used to make the mini-batch have an arbitrary mean and standard deviation and are two parameters that needs to be learned when the network is being fitted. This helps the network to converge, as the batch normalization keeps centering the mean and standard deviation of the mini-batches [19].

Epoch

An *epoch* is an iteration through the whole dataset during fitting of the network. Multiple epochs are often needed to reach the minimum of the loss function [19].

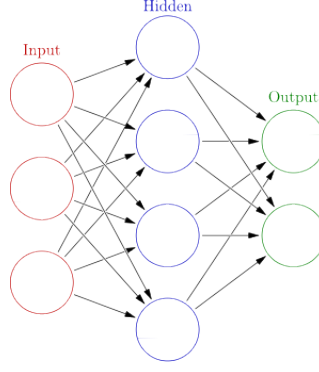


Figure 2: Visualization of a feedforward neural network with a single hidden layer [18]

The Architecture and Forwardpropagation

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 2. 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 2), 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 2), 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 $a_n^{(i)}$ 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

$$\mathbf{a}^{(i+1)} = g^{(i+1)} \left(\mathbf{z}^{(i+1)} \right)$$

where

$$\mathbf{z}^{(i+1)} = \mathbf{W}^{(i+1)} \mathbf{a}^{(i)} + \mathbf{b}^{(i+1)},$$

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

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

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

$$\mathbf{a}^{(i)} = \begin{pmatrix} a_0^{(i)} \\ a_1^{(i)} \\ \vdots \\ a_n^{(i)} \end{pmatrix},$$

$\mathbf{b}^{(i+1)}$ 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 [4] [17]. 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 and 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* which 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 (and actually some functions outside of this class), as long as the network has enough hidden neurons [4].

Backpropagation

Backpropagation is an algorithm used to compute the gradient of the network. It is used together with an optimization algorithm, such as `RMSPROP`, to train the model by minimizing the training loss of the model. Backpropagation happens after data has flowed through the model from the input to the output, and works by computing the gradient of each parameter sequentially from the output to the input of the model. The procedure makes heavily use of the *chain rule* from calculus, which states, that if we let $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$, g be a function that maps from \mathbb{R}^m to \mathbb{R}^n and f be a function that maps from \mathbb{R}^n to \mathbb{R} , then, if we let $\mathbf{y} = g(\mathbf{x})$ and $z = f(\mathbf{y})$, we can then compute $\frac{\partial z}{\partial x_i}$ by

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

[4]. If we use this to find the gradient of each parameter, we will find, that the partial derivative for each weight is

$$\frac{\partial L}{\partial w_{jk}^{(i)}} = \frac{\partial z_j^{(i)}}{\partial w_{jk}^{(i)}} \frac{\partial a_j^{(i)}}{\partial z_j^{(i)}} \frac{\partial L}{\partial a_j^{(i)}} = a_k^{(i-1)} g'^{(i)}(z_j^{(i)}) \frac{\partial L}{\partial a_j^{(i)}}$$

and the partial derivative of each bias is

$$\frac{\partial L}{\partial b_j^{(i)}} = \frac{\partial z_j^{(i)}}{\partial b_j^{(i)}} \frac{\partial a_j^{(i)}}{\partial z_j^{(i)}} \frac{\partial L}{\partial a_j^{(i)}} = g'^{(i)}(z_j^{(i)}) \frac{\partial L}{\partial a_j^{(i)}}$$

where for both cases

$$\frac{\partial L}{\partial a_j^{(i)}} = \sum_{k=0}^{n_i-1} w_{jk}^{(i+1)} g'^{(i+1)}(z_j^{(i+1)}) \frac{\partial L}{\partial a_j^{(i+1)}}.$$

if $\mathbf{a}^{(i)}$ is not the output-layer. Once the partial derivative of all weights and biases has been found, the gradient vector can be formed and an optimization method can be used to optimize the parameters of the model [16].

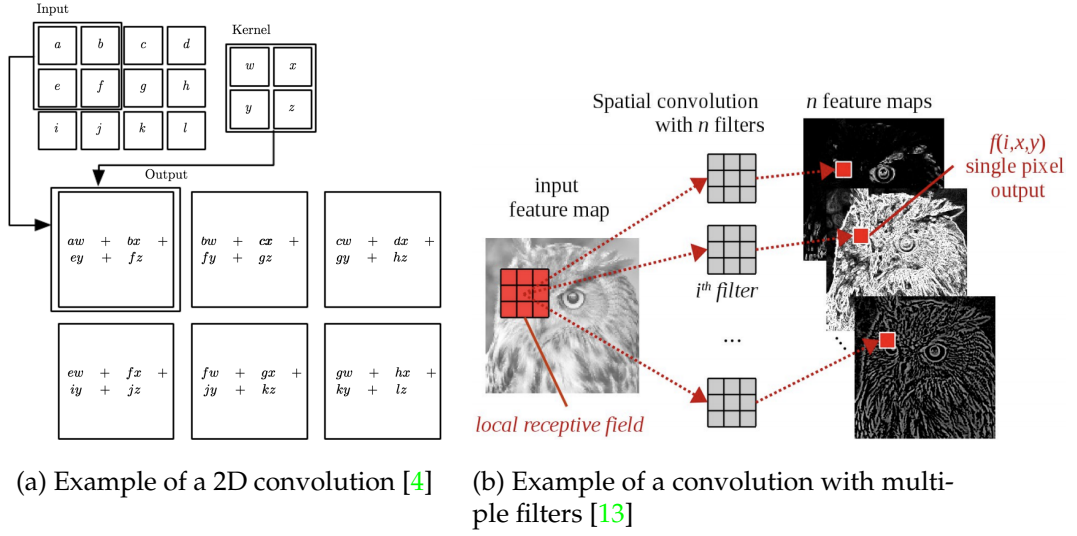


Figure 3: Convolutions visualized

2.4.2 Convolutional Neural Networks

Feedforward neural networks introduced in 2.4.1 can be used for pattern recognition within images, however, they are usually not used for this task. Consider a colored input image of dimension 64×64 . If we were to use a feedforward neural network on this image, each neuron in the first hidden layer would be connected to the input layer through 12.288 weights. Not only would this use a lot of computational power and time to train, however, a network of this size would also be prone to overfitting. Instead, *convolutional neural networks* (also known as CNN's) are usually used. A CNN usually consists of *convolutional layers*, *pooling layers* and *fully-connected layers*, where the fully-connected layers are analogous to the layers in a feedforward neural network [10].

Convolutional layers

A convolutional layer is composed of a set of *kernels* (also known as *filters*), which are matrices of weights of dimension $k \times k$, where k usually is 5, 3 or 1, and each weight is a parameter for the model to be learned [13]. Each kernel is used on the input to produce a *feature map*. The kernels are applied to the input by "sliding" over the input (where the step size is called *stride* and is usually by default equal to 1). Each $k \times k$ grid of the input (called the *local receptive field*) is then used to compute the dot-product between the grid and each kernel, which is then placed in the corresponding feature map of each kernel, as visualized in Figure 3. When all of the feature maps have been computed, the feature maps are stacked together, forming a tensor, which is then returned by the layer and a activation function can be applied.

As described previously, by using a convolutional layer we can dramatically decrease the amount of weights used by the layer. If we were to use a 3×3 kernel on a colored image, we would reduce the amount of weights on each neuron in the following layer from 12.288 down to just 27, reducing both the training time and making the network less prone to overfitting [10].

Pooling layers

Pooling layers are layers used to reduce the dimension of the input. The most common pooling layer is the *maxpooling*-operation. The operation works by considering each $k \times k$ "grid", like in the case with the convolutional layer, in which the maximum entry in that grid is being inserted into the output [10].

Algorithm 3 Nearest Neighbour Upsampling [14]

Require: Input image, X , of size $m \times n$

Require: Wanted output size $xm \times yn$

- 1: Create empty image, O , of size $xm \times yn$
 - 2: **for each** pixel, p , in X **do**
 - 3: $i, j \leftarrow$ index of p in X
 - 4: Insert p at index (xi, yj) in O
 - 5: **for each** empty pixel, p , in O **do** $p \leftarrow$ value of nearest neighbour
 - return O
-

Nearest Neighbour Upsampling

Sometimes we want to increase the size of an image. This is done by making use of *upsampling* (also known as *interpolation*) techniques. One of the most common upsampling techniques is *nearest neighbour upsampling*, whose pseudocode has been written in 3. The algorithm starts off by taking an image, X , of size $m \times n$, as input, which we wish to upsample to size $xm \times yn$. The algorithm then loops over each pixel, p , in X , finds the corresponding index, i, j , of p in X , and places p at index (xi, yj) in the output image O of size $xm \times yn$. When this is done, it assigns each of the empty pixel in O the value of their nearest neighbour, making each pixel in O have a value, and then returns O [14].

3 Stacked Hourglass, t-SNE and K-Means Clustering

In this section the various algorithms and architectures used throughout this thesis is described and explained in details.

3.1 Stacked Hourglass

When performing the pose estimation in section **REFERENCE MANGLER**, we will be implementing and using the *stacked hourglass* described by Newell *et al.* [9]. The following description and explanation of the architecture is based on an interpretation of Newell *et al.* [9] and Camilla Olsen [11].

3.1.1 Reasoning behind using the Stacked Hourglass

We have decided to make use of the stacked hourglass described by Newell *et al.*, as it is a intermediate architecture that has shown state-of-the-art results. At the same time the architecture of the network is similar to the architecture of *autoencoders*, making the architecture useful for encoding the data into a lower dimension, which can be useful in section **REFERENCE MANGLER**, when we will be doing the interpretation of the model.

3.1.2 The Residual Module

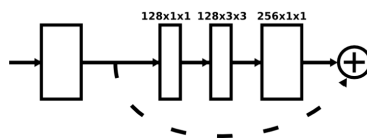


Figure 4: Visualization of the residual module [9]

The stacked hourglass makes heavily use of so-called *residual modules*, which is visualized in Figure 4. The module works by taking an input, which is sent through a 1×1 and a 3×3 convolution, which both return 128 featuremaps. Then, the 128 featuremaps are sent through a 1×1 convolution, which returns 256 featuremaps. Lastly, element-wise addition is then used to add the 256 featuremaps to the input of the module, which the model then returns. All convolutions are followed by an activation function and are *same convolutions*, meaning the outputted featuremaps are of the same dimensions as the inputted featuremaps.

3.1.3 The Hourglass

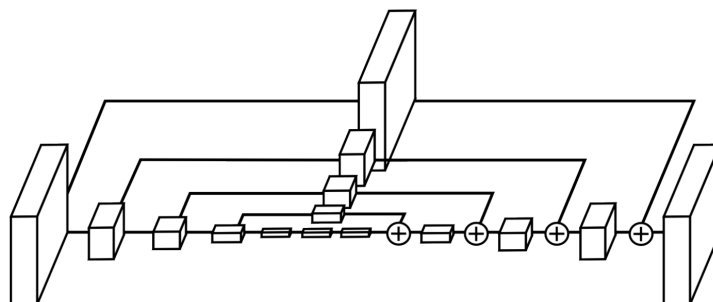


Figure 5: Visualization of a single hourglass [9]

The stacked hourglass consists of hourglasses, where each hourglass is split into an encoder, where the featuremaps are downsampled, and a decoder, where the featuremaps are upsampled. The hourglass is symmetric, in the sense, that it has an equal amount of downsampling layers in the encoder as there are upsampling layers in the decoder. In Figure 5 a single hourglass has been visualized, where each box is a residual module. The hourglass works by using residuals and max poolings to process features down to a low resolution. Then, nearest neighbor upsampling is used to upsample the featuremaps until the featuremaps have the same dimensions as the input of the hourglass. Before each max pool in the encoder, the network branches off and applies a residual to the feature maps, which are then added element-wise before the corresponding residual in the decoder, which helps to ensure that lost information in the encoder is kept. Between the encoder and decoder the network has a bottleneck, where no downsampling or upsampling happens, instead only residuals are processing the featuremaps. After the decoder two 1×1 convolution layers are applied to produce the final predictions of the network.

3.1.4 The Stacked Hourglass

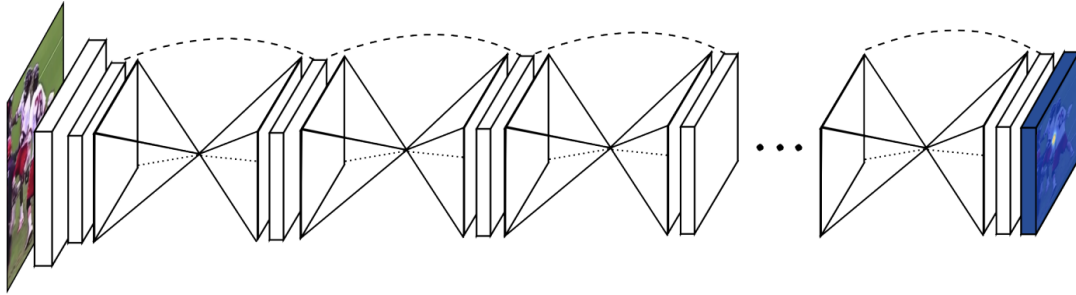


Figure 6: Visualization of the stacked hourglass [9]

The full network is built by stacking multiple hourglasses end-to-end, making the output of one hourglass be the input of the next hourglass, as shown in Figure 6, which makes each hourglass reestimate features of the image. The input of the network is a 256×256 RGB-image. To lower the memory usage, the network starts off with a 7×7 convolution layer with stride 2, followed by a residual module and max pooling to bring the resolution down to 64×64 , which is then input to the first hourglass. By the end the whole network outputs n heatmaps corresponding to the n joints it should predict for a single person. The prediction of a joint is thus the maximum activation of the corresponding heatmap.

3.2 t-SNE

3.3 K-means Clustering

4 The Dataset

To perform the pose estimation, we need some data on which to train, validate and test our model. Throughout this section the used data and preprocessing are described.

4.1 The COCO Dataset



Notice how the image contains multiple people, each with their own keypoints and amount of joints labeled

Figure 7: Example of an image from the COCO dataset with the keypoints drawn on [8]

The data needed for our model has to fit to our problem and has to be annotated, as our model will perform supervised learning. There are multiple datasets that fits these requirements. One of these datasets is the Common Objects in Context (COCO) dataset [8], which we will be using. The dataset contains annotations for different purposes, however, for our pose-estimation-task, only the keypoint annotations of human bodies are needed. An example of such a picture with the keypoints labeled can be seen in Figure 7.

The annotation of each person consists of an array with a length of 51, which annotates 17 keypoints of a person. Thus, each joint corresponds to three sequential elements in the array, where the first and second indices corresponds to the x and y -location of the joint in the image, and the third index is a flag, v , telling the visibility of the joint in the image. v has three outcomes: if $v = 0$, the joint is not labeled, if $v = 1$, the joint is labeled but not visible, and if $v = 2$, the joint is visible and labeled.

The creators of the dataset has already split the data into three parts: a part used for training the model, a part used for validating the model and a part used for testing the model. However, the part used for testing the model is unlabel, hence, it is unusable for our purpose, as our model will be doing supervised learning. As both the training dataset and the validation dataset will be used for training and tuning the model, we will need to create our own hold-out dataset for testing to provide an unbiased evaluation of the final model.

The training and validation sets contains a total of about 123.000 various images. As we only need the images that contain humans, we will be discarding the images without any humans, leaving us with a total of about 66.808 images of humans doing various tasks, with a total of 149.813 humans annotated with keypoints. Each image can contain multiple people, which we need to handle before training our model, as we will be focusing on single-human pose estimation. Besides this, each image also has different resolution and aspect ratio, which we also need to handle, as our model requires the images to have a fixed resolution. Lastly, we should also do some handling of the labels before training the model, as there could have been some inaccuracies, when the joints were labeled. This especially applies when $v = 1$, that is, when the joint is labeled but not visible, as there are more inaccuracies or uncertainty when labeling a non-visible joint than when labeling a visible joint.

4.2 Data Preprocessing

4.2.1 Creating the test dataset

To create the dataset which will be used for testing we take the training set, since it is the larger of the training and validation set, and sample 5.064 images randomly without replacement, to create a test set. This ensures that the test-set and validation-set are of the same size. This new test set will not be used when training the model nor used when tuning the parameters. Instead, it will only be used to evaluate the very final model.

4.2.2 Preprocessing the images

| | Amount of images | Percentage |
|----------------|------------------|------------|
| Training set | 124.040 | 92,45 |
| Validation set | 5.064 | 3,77 |
| Testing set | 5.064 | 3,77 |
| Total | 134.168 | 100 |

Table 1: Data distribution

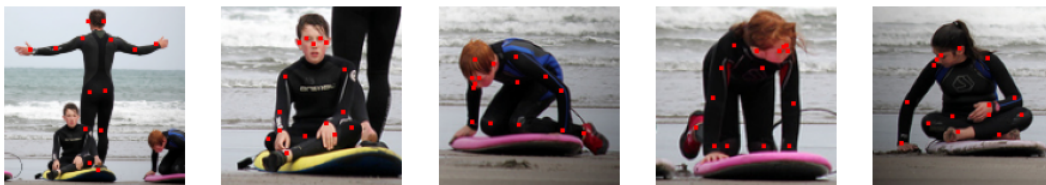


Figure 8: The results of processing the image from Figure 7 with the corresponding labels [8]

We start the preprocessing of the images by creating multiple bounding boxes, where each bounding box surrounds a single person, which is done by making use of the bounding box annotations provided by COCO. Then each bounding box is transformed into a square by making the shorter sides have the same length as the longer sides - this is done to ensure that the aspect ratio of the image is kept, when it is later resized.

An issue can happen, where the bounding box still contains multiple people, which will confuse our model, since it does not know which person it should annotate. An example of this can be seen in the first image of Figure 8. To fix this we center the bounding box around the person it should annotate, making the model annotate the person in the center of the input image. This is done by centering the bounding box with respect to the outermost keypoints of the person.

Since each keypoint does not necessarily lie on the edge of the person, the current bounding boxes would result in not all of the pixels of the corresponding person being in the bounding box. For this reason, each bounding box is expanded with 10% in the height and width. If, however, the image cannot contain the expanded bounding box, the bounding box is then expanded as much as possible, while still being a square. If it is the case, that one of the corners of the bounding box lies outside of the image, then the bounding box is moved either up or down, making the corner lie inside the image and keeping the annotated person centered along the x -axis.

When all of the above is done, the image is finally cropped to each bounding box, resulting in multiple squared images, each containing an unique person at the center. Each of these squared images are then resized to a 256×256 image, has their rgb -values centered by subtracting the mean rgb of all of the images from the training set, and saved. Doing all of these

steps results in the distribution of images displayed in Table 1. In Figure 8, the results of processing the image from Figure 7 are shown with the corresponding labels. Lastly, the data is shuffled to help the developed model generalize the data better.

4.2.3 Handling the labels



Left: The original image. Right: The heatmaps of all the keypoints, fused together to a single image.

Figure 9: An example of the heatmaps of a single image fused together and put over the original image [2]

For each image of a single person, our model outputs 17 heatmaps, one for each possible joint in the image, which tells the probability of the joint being in each pixel. An example of the heatmaps fused together can be seen in Figure 9.

The ground truth heatmap of a single joint is created firstly by initializing an all-zero 2D array with size 256×256 for each of the 17 heatmaps. Next, in the i th 2D array at position (x_i, y_i) , corresponding to the position of the i th joint, a 1 is placed - this 1 now corresponds to where the i th joint is placed in the image according to the keypoint annotation of the image. Next, a Gaussian filter is used to smear out the image, where the standard deviation depends on the visibility of the joint: if the joint is visible, then the standard deviation is 0.5, whereas the standard deviation is 1 if the joint is not visible, as there are more uncertainty with the labeling of such keypoints. Lastly, as the model outputs 17 64×64 heatmaps, our heatmaps are resized from a dimension of 256×256 to a dimension of 64×64 .

We do all of this for all of the 17 joints for each image, resulting in the keypoints which will be used for developing our model.

5 Experiment

Throughout the following section the results and configuration details of our trained model are described and explained.

5.1 Configuration Details

- PCK og plateau
- Skriv hvorfor jeg ikke laver flere eksperimenter men kun ét

To perform the pose estimation, we make use of the stacked hourglass architecture, however, our model only consists of a single hourglass of 4 downsamples and upsamples, 1 residual module between each down- and upsample, 1 residual module in each skip-connection, and 3 residual modules in the bottleneck. Newell *et al.* [9] and Olsen [11] experiment with different amount of hourglasses and with different amount of residual modules in each hourglass. They both come to the conclusion, that stacking multiple hourglasses or using hourglasses with multiple residual modules between each down- and upsample increases the performance of the model. However, as the main purpose of this thesis is not to create a model with state-of-the-art results, but instead to create a model that can be interpreted and explored, we have chosen to reduce the size of the model.

To prevent the model from overfitting we follow Newell *et al.* [9] and use batch normalization. Newell *et al.* does not describe where to perform the batch normalization, so we follow Olsen [11] and perform the batch normalization before each convolutional layer in each residual module, after the first convolutional layer of the entire network and before the last convolutional layer of the entire network. For the choice of activation function we follow Olsen [11] and use the *ReLU*-function after each batch normalization. Each max pooling and nearest neighbor upsampling uses a kernel size of 2, which halves and doubles the size of the input, respectively. The full network has been visualized in Figure 10.

For the initial values of the weights we follow Olsen [11] and initialize each weight by sampling from a *Glorot normal distribution* (also known as a *Xavier normal distribution*), described as

$$\mathcal{N}\left(0, \frac{2}{fan_{in} + fan_{out}}\right)$$

where fan_{in} is the amount of input connections and fan_{out} is the amount of output connections to the layer of the weight [3]. By doing so we make all layers have the same activation variance and gradient variance, essentially helping the model to converge [4].

We make use of a mini-batch size of 16 and no data augmentation, since the dataset is already rather large and captures a lot of the variances.

To optimize the network we follow Newell *et al.* and make use of *MSE* as our loss function, *RMSPROP* as our optimizer, as well as use a initial learning rate of $2.5e - 4$ [9]. After each epoch we find the validation accuracy of the model by computing the *PCK* of the model.

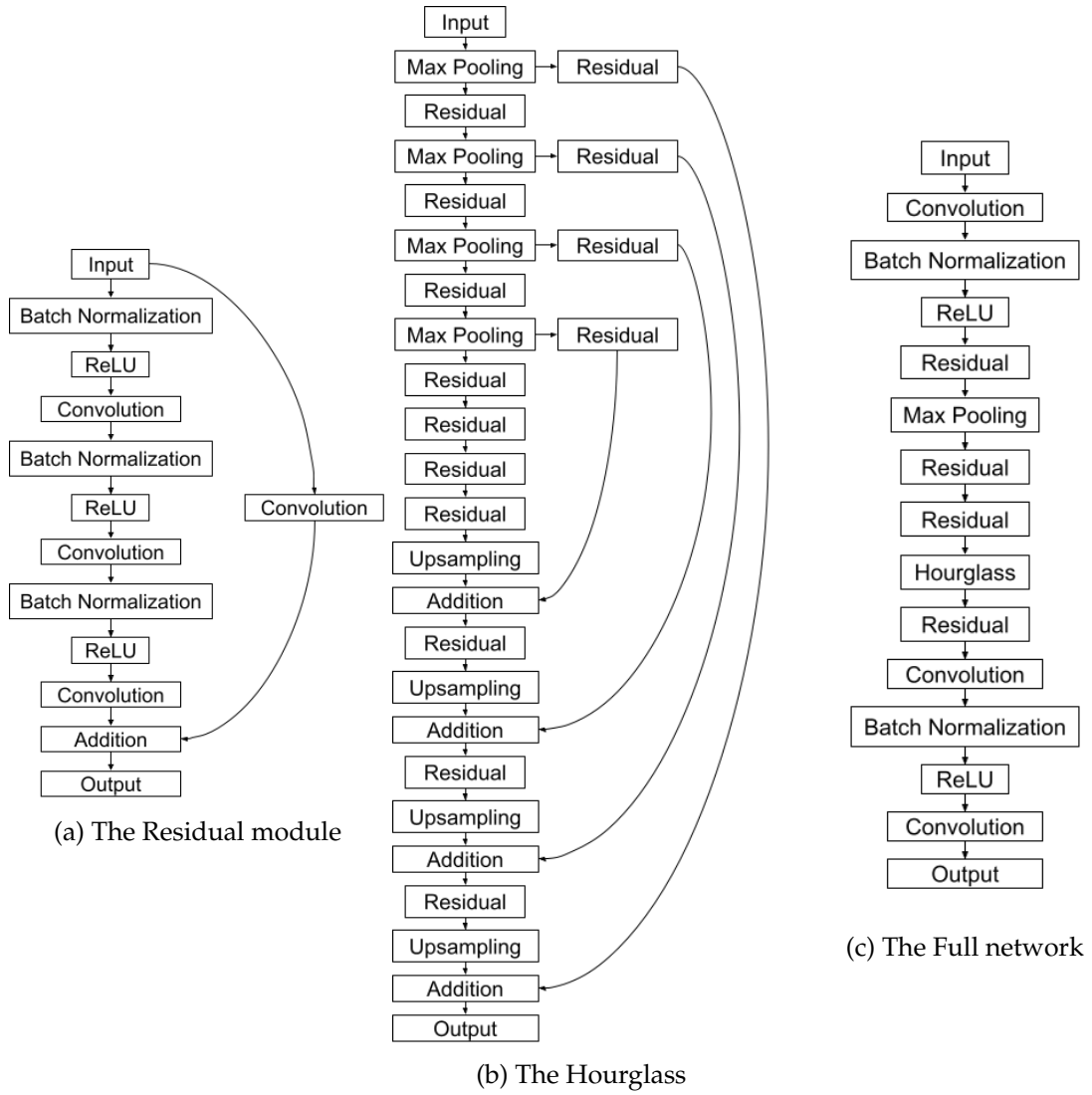


Figure 10: Overview of the used architecture

5.2 Results

- vis grafer
- Overfitting

5.3 Training Details

The stacked hourglass was implemented in Python 3.8.2 using PyTorch version 1.7.1 and Cuda version 10.2 on a machine using Windows 10 version 20H2, build 19042. The network was trained on an 8 GB NVIDIA GeForce GTX 1070 GPU using a Samsung 840 EVO SSD for data storage. Training the network takes about 70 minutes per epoch, totalling to about 58 hours for 50 epochs.

6 Discussion

- PCK: Scaler ikke med input størrelse
- Glorot Initialisering - kan være dårligt, se EML forelæsning om optimization.

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Ed. by B. Schölkopf M. Jordan J. Kleinberg. URL: <https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>. (accessed: 10.3.2021).
- [2] Qi Dang, Jianqin Yin, Bin Wang, and Wenqing Zheng. "Deep Learning Based 2D Human Pose Estimation: A Survey". In: 6.24 (December 2019).
- [3] Daniel Falbel, JJ Allaire, and François Chollet. *Glorot normal initializer, also called Xavier normal initializer*. Keras. URL: https://keras.rstudio.com/reference/initializer_glorot_normal.html. (accessed: 22.4.2021).
- [4] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016. (accessed: 18.3.2021).
- [5] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *Elements of Statistical Learning*. URL: https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII_print12_toc.pdf. (accessed: 10.3.2021).
- [6] Bulat Ibragimov. *Modelling and Analysis of Data, Lecture 3 - Nonlinear Regression*. 25.11.2020. URL: https://absalon.ku.dk/courses/42639/files/4289570?module_item_id=1145250. (accessed: 17.3.2021).
- [7] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning with Applications in R*. URL: <https://static1.squarespace.com/static/5ff2adbe3fe4fe33db902812/t/601cc86d7f828c4792e0bcae/1612499080032/ISLR+Seventh+Printing.pdf>. (accessed: 17.3.2021).
- [8] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollar. "Microsoft COCO: Common Objects in Context". In: (2014).
- [9] Alejandro Newell, Kaiyu Yang, and Jia Deng. "Stacked Hourglass Networks for Human Pose Estimation". In: (2016).
- [10] Keiron O'Shea and Ryan Nash. *An Introduction to Convolutional Neural Networks*. Tech. rep. 1511.08458. Department of Computer Science, Aberystwyth University, Ceredigion, School of Computing, and Communications, Lancaster University, 2015.
- [11] Camilla Maach Brønnum Olsen. "Articulated Pose Estimation of Humans". MA thesis. University of Copenhagen, Department of Computer Science, 2019.
- [12] Jens Petersen. *Elements of Machine learning - Optimization in Deep Learning*. 2021. URL: https://absalon.ku.dk/courses/46845/files/4490846?module_item_id=1206787. (accessed: 25.3.2021).
- [13] Moacir A. Ponti, Leonardo S. F. Ribeiro, Tiago S. Nazare, Tu Bui, and John Collomosse. *Everything you wanted to know about Deep Learning for Computer Vision but were afraid to ask*. Tech. rep. ICMC – University of Sao Paulo and CVSSP – University of Surrey, 2017.
- [14] Mike Pound. *Resizing Images - Computerphile*. Computerphile. URL: https://www.youtube.com/watch?v=AqscP7rc8_M. (accessed: 19.4.2021).
- [15] Simon Rogers and Mark Girolami. *A First Course in Machine Learning*. Chapman and Hall/CRC, 2017.
- [16] Grant Sanderson. *Backpropagation calculus | Deep learning, chapter 4*. 3Blue1Brown. URL: https://www.youtube.com/watch?v=tIeHLnjs5U8&ab_channel=3Blue1Brown. (accessed: 7.4.2021).

- [17] Grant Sanderson. *But what is a Neural Network? | Deep learning, chapter 1*. 3Blue1Brown. URL: <https://www.youtube.com/watch?v=aircAruvnKk>. (accessed: 18.3.2021).
- [18] *What is a Neural Network?* URL: <https://deeptai.org/machine-learning-glossary-and-terms/neural-network>. (accessed: 17.3.2021).
- [19] Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. *Dive into Deep Learning*. <https://d2l.ai>. 2020.