



Bachelor Thesis

2D Articulated Human Pose Estimation

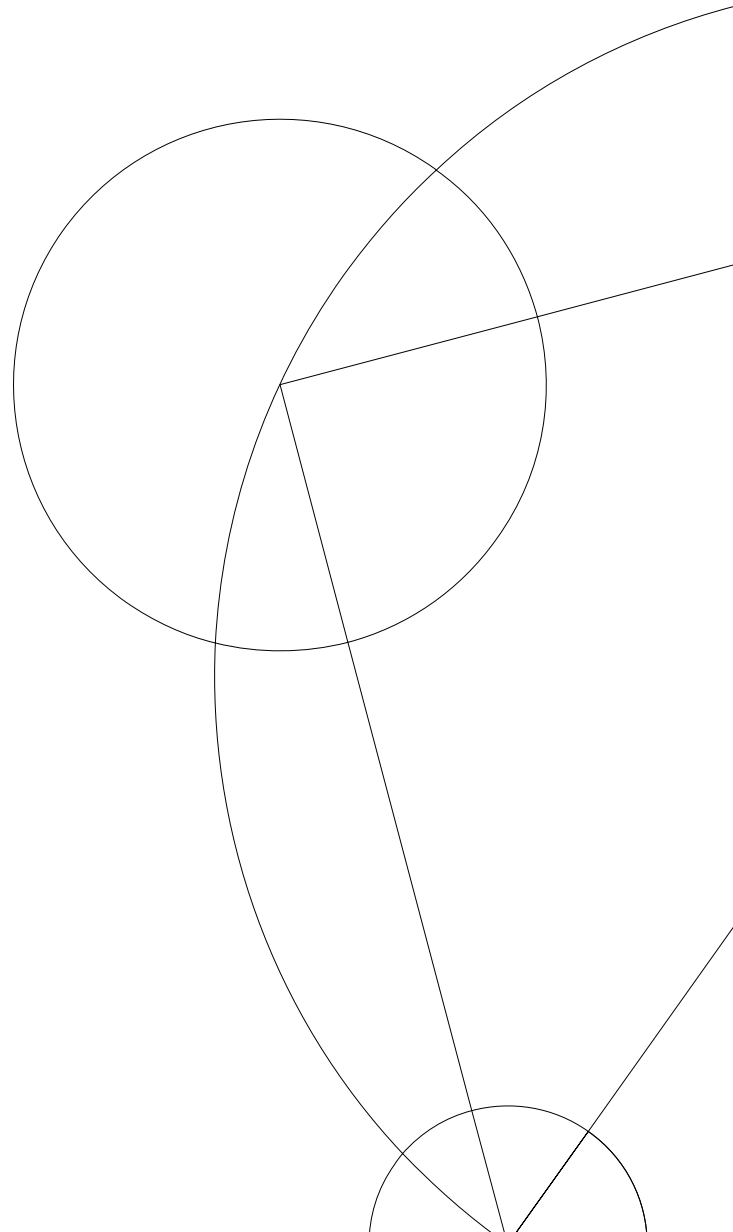
Using Explainable Artificial Intelligence

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

March 18, 2021

Supervisor

Kim Steenstrup Pedersen kimstp@di.ku.dk



Contents

1	Machine Learning Theory	3
1.1	Motivation	3
1.2	Machine Learning Paradigms	3
1.3	Evaluation of Machine Learning Models	3
1.3.1	Splitting the dataset	3
1.3.2	Evaluation Metrics for Supervised Machine Learning (Loss Functions)	4
1.4	Neural Networks	4
1.4.1	The Mathematics Behind Neural Networks	4
1.4.2	Termonology	6
1.4.3	Convolutional Neural Networks	6
1.4.4	Stacked Hourglass	6
1.5	Principal Components Analysis and K-means Clustering	6
1.5.1	Principal Components Analysis (PCA)	6
1.5.2	K-means Clustering	6
2	The Dataset	7
2.1	The COCO Dataset	7
2.2	Data Preprocessing	8
2.2.1	Creating the test dataset	8
2.2.2	Preprocessing the images	8
2.2.3	Handling the labels	9

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 [8]. 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 [8].
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 [7].

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

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

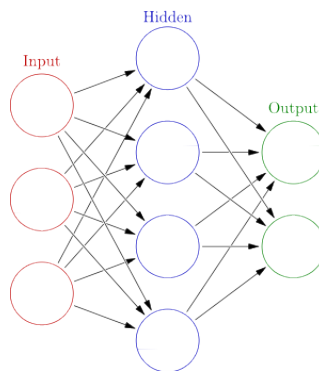
where y_i is the true value of the i th observation and \hat{y}_i is the estimated value of the i th observation. Thus, *MSE* measures the average squared difference between the true observation and the estimated observation. The aim of a model is thus to make the *MSE* as small as possible [2].

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 [10]



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] [6]. 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

2 The Dataset

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

2.1 The COCO Dataset

Figure 2: Example of an image from the COCO dataset with the keypoints drawn on [9]



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

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 [9], 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 2.

The annotation of each person consists of an array with a length of 51. Each joint corresponds to three sequential elements in the array, where the first index tells the x -location of the joint in the image, the second index tells the 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. Thus, up to 17 keypoints of a person can be annotated. v has three outcomes: if $v = 0$, then the joint is not labeled, if $v = 1$, then the joint is labeled but not visible, and if $v = 2$, then 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 unlabeled, hence, why 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. 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 for two reasons

1. 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.

2. Each joint could correspond to multiple pixels in the image, hence why it is not correct to use only a single pixel as the location of the joint in the image, which is the current case.

2.2 Data Preprocessing

2.2.1 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 set and the 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.

2.2.2 Preprocessing the images

Figure 3: Data distribution

	Amount of images	Percentage
Training set	118.304	92.684
Validation set	5.064	3.658
Testing set	5.064	3.658
Total	128.432	100

Figure 4: The results of processing the image from Figure 2 with the corresponding labels [9]



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. 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. Each of these squared images are

then resized to a 256×256 image, centralized by subtracting the mean rgb of all of the images, and saved. Doing all of these steps results in the distribution of images displayed in Figure 3. In Figure 4 the results of processing the image from Figure 2 are shown with the corresponding labels. Lastly, the data is shuffled to help the developed model generalize the data better.

2.2.3 Handling the labels

Figure 5: An example of the heatmaps of a single image fused together and put over the original image [5]



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

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 5.

The 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.

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Ed. by B. Schölkopf M. Jordan J. Kleinberg. (accessed: 10.3.2021).
- [2] Trevor Hastie Robert Tibshirani Gareth James Daniela Witten. *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).
- [3] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [4] 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).
- [5] Bin Wang Wenqing Zheng Qi Dang Jianqin Yin. “Deep Learning Based 2D Human Pose Estimation: A Survey”. In: 6.24 (December 2019).
- [6] Grant Sanderson. *But what is a Neural Network? | Deep learning, chapter 1*. 3Blue1Brown. URL: <https://www.youtube.com/watch?v=aircAruvnKk>.
- [7] Mark Girolami Simon Rogers. *A First Course in Machine Learning*. Chapman and Hall/CRC, 2017.
- [8] Jerome Friedman Trevor Hastie Robert Tibshirani. *Elements of Statistical Learning*. URL: https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII_print12_toc.pdf. (accessed: 10.3.2021).
- [9] Serge Belongie James Hays Pietro Perona Deva Ramanan C. Lawrence Zitnick Piotr Dollar Tsung-Yi Lin Michael Maire. “Microsoft COCO: Common Objects in Context”. In: (2014).
- [10] *What is a Neural Network?* URL: <https://deeptai.org/machine-learning-glossary-and-terms/neural-network>. (accessed: 17.3.2021).