

Manu Varma

# **Kin Recognition Using Weighted Graph Embeddings**

Computer Science Tripos – Part II

St John's College

May 13, 2021

# Declaration

I, Manu Varma of St. John's College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

Signed Manu Varma

Date May 13, 2021

# Proforma

Candidate Number: **2356E**  
Project Title: **Kin Recognition Using Weighted  
Graph Embeddings**  
Examination: **Computer Science Tripos – Part II, June 2021**  
Word Count: **11998<sup>1</sup>**  
Line Count: **3027<sup>2</sup>**  
Project Originator: **The Dissertation Author**  
Supervisor: **Daniel Bates**

## Original Aims of the Project

Kin recognition is the ability to recognize whether two people are related to each other just by looking at them. The goal of the project was to implement a state-of-the-art model, Weighted Graph Embedding-Based Metric Learning (WGEML), to solve the Kin Verification problem and to verify the accuracies obtained. I also wanted to further evaluate the models obtained by looking at potential biases in the datasets used in the paper.

## Work Completed

I successfully implemented WGEML and was able to recreate the original accuracies within an acceptable range, thus fulfilling my success criterion. Furthermore, I performed ablation studies in order to determine if there is a relationship between the number of face descriptors used and the accuracy of the model. I also replaced one face descriptor with a less computationally expensive one. I explored biases in the datasets used using the model that was created and discussed the implications of such biases on the results.

## Special Difficulties

None.

---

<sup>1</sup>Calculated using `texcount`

<sup>2</sup>Calculated using `find . -name '*.py' | xargs wc -l`

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem Overview . . . . .	1
1.2	Motivation . . . . .	2
1.3	Related Work . . . . .	2
1.4	Project Overview . . . . .	3
<b>2</b>	<b>Preparation</b>	<b>4</b>
2.1	Convolutional Neural Networks . . . . .	4
2.1.1	Convolutional Layers . . . . .	5
2.1.2	Pooling Layers . . . . .	6
2.1.3	Activation Functions . . . . .	6
2.1.3.1	ReLU . . . . .	7
2.1.3.2	Softmax . . . . .	7
2.2	Face Detection . . . . .	7
2.3	Face Descriptors . . . . .	7
2.3.1	Local Binary Patterns . . . . .	8
2.3.2	Histogram of Gradients . . . . .	9
2.3.2.1	Calculating the Gradients . . . . .	9
2.3.2.2	Weighted Vote into Histogram . . . . .	10
2.3.3	Scale-Invariant Feature Transform . . . . .	11
2.3.4	VGG . . . . .	11
2.4	Metric Learning . . . . .	12
2.5	K-Nearest Neighbors . . . . .	14
2.6	Requirements Analysis . . . . .	14
2.7	Software Engineering Practices . . . . .	15
2.7.1	Starting Point . . . . .	15
2.7.2	Tools Used . . . . .	15
2.7.3	Datasets . . . . .	16
2.7.4	Testing . . . . .	17
2.7.5	Licensing . . . . .	17
<b>3</b>	<b>Implementation</b>	<b>18</b>
3.1	Repository Overview . . . . .	18
3.2	Data Preparation . . . . .	19
3.2.1	Cross-Validation . . . . .	19
3.2.2	Positive and Negative Pairs . . . . .	20
3.2.3	Dimensionality Reduction . . . . .	20

3.2.4	Saving Results to Disk . . . . .	21
3.3	Face Descriptors . . . . .	21
3.3.1	SIFT Implementation . . . . .	22
3.3.2	VGG Implementation . . . . .	22
3.3.3	CifarNet Extension . . . . .	22
3.4	WGEML . . . . .	23
3.4.1	Problem . . . . .	23
3.4.2	Approach . . . . .	24
3.5	Prediction . . . . .	26
3.5.1	Tri-kin Relationship Prediction . . . . .	26
3.6	Overview of Workflow . . . . .	27
3.6.1	Preprocessing . . . . .	27
3.6.2	Training . . . . .	27
3.6.3	Testing . . . . .	27
<b>4</b>	<b>Evaluation</b>	<b>29</b>
4.1	Overall Accuracies of the Model . . . . .	29
4.2	Success Criterion . . . . .	30
4.3	Receiver Operating Characteristic (ROC) Curves . . . . .	30
4.4	Potential Biases in Datasets . . . . .	32
4.5	Ablation Studies . . . . .	33
4.5.1	Blocking Face Descriptors . . . . .	33
4.6	CifarNet Extension . . . . .	35
4.7	Unit Tests . . . . .	36
<b>5</b>	<b>Conclusion</b>	<b>37</b>
5.1	Lessons Learned . . . . .	37
5.2	Future Work . . . . .	38
	<b>Bibliography</b>	<b>39</b>
<b>A</b>	<b>Algorithms Implemented in Libraries</b>	<b>42</b>
A.1	Haar-Based Cascade Classifier . . . . .	42
A.2	SIFT Keypoint Extraction . . . . .	42
A.3	Principal Component Analysis . . . . .	44
<b>B</b>	<b>Raw Table Data</b>	<b>46</b>
<b>C</b>	<b>Project Proposal</b>	<b>47</b>

# List of Figures

1.1	An example of an input into the kin verification problem . . . . .	1
2.1	An example of an artificial neural network with two hidden layers . . . . .	4
2.2	A visualization of a convolution . . . . .	5
2.3	Max-pooling operation . . . . .	6
2.4	Local Binary Patterns neighborhoods . . . . .	8
2.5	Example of the LBP operator on a pixel . . . . .	9
2.6	VGG network configurations . . . . .	12
2.7	A high-level view of metric learning . . . . .	14
3.1	Folder structure of the project . . . . .	19
3.2	The architecture of the CifarNet model . . . . .	23
3.3	The preprocessing pipeline . . . . .	27
3.4	The training pipeline . . . . .	28
3.5	A sample output of the testing stage on KinFaceW-I unrestricted . . . . .	28
3.6	The testing pipeline . . . . .	28
4.1	WGEML implementation accuracies . . . . .	29
4.2	Difference in WGEML accuracies . . . . .	30
4.3	ROC curves for each dataset . . . . .	31
4.4	ROC curves for each relationship . . . . .	31
4.5	Differences in accuracy when a the test dataset is different from the training set	32
4.6	Average differences in accuracy when a the test dataset is different from the training set . . . . .	33
4.7	Accuracy versus the number of face descriptors used . . . . .	34
4.8	Differences in accuracy when VGG is replaced with CFN grouped by relationship	35
4.9	Differences in accuracy when VGG is replaced with CFN grouped by dataset . .	35
4.10	Coverage of the unit tests of the project . . . . .	36
A.1	Haar-based features . . . . .	42
A.2	Difference of Gaussians . . . . .	43

# Chapter 1

## Introduction

Kinship recognition is the ability to recognize how two people are related to each other without prior information about who they are. In organisms, it is advantageous to inclusive fitness to be able to recognize which of their neighbors are close relatives [7]. Thus, it stands to reason that the ability to recognize kinship relationships has evolved in humans. In humans, specifically, facial resemblance is expected to serve as an indicator of kinship and it has been demonstrated that strangers are able to match photographs of mothers to their infants without any prior contact with any of the family [16].

This idea can be modeled which gives rise to computational kinship recognition. This is the field of studying how kinship relationships between people can be identified without any prior knowledge of the family.

### 1.1 Problem Overview

In the field of computational kinship recognition, there are a variety of problems that can be explored. The most common problem, kinship verification, takes as input a pair of images and a proposed kinship relationship, for example, mother-daughter, and recognizes whether the relationship exists. We can see an example set of images in figure 1.1.

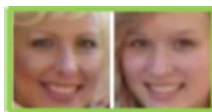


Figure 1.1: An example of an input into the kin verification problem

These relationships are conventionally parent-child relationships as opposed to sibling-sibling relationships. A further extension of the main kin verification problem is Tri-Subject Kinship Verification which takes three images, two parents and a child, and determines if they are related or not. These are the problems that I consider in this dissertation.

Some other major problems in the field include search and retrieval and family classification. Search and retrieval takes an image of a person as input and searches through a database to find who they are most likely related to. This outputs a list of people they could be related to. Family classification deals with a similar task of taking an image of a person as input and figuring out which family they may belong to.

## 1.2 Motivation

Applications of accurate kinship recognition software include being able to better recognize missing children and match them with their parents [22]. As the models tend to focus on facial features, even as the missing child ages, as long as the right face descriptors are used, it is still possible to find missing children. Furthermore, it can also be used for stopping human traffickers from claiming they are family members of a victim or to reunite families across refugee camps. These humanitarian issues can be solved by using an application of one of the problems discussed in section 1.1. For example, to reunite families across refugee camps, we can use a method that solves the search and retrieval problem.

In the field of computer vision, there is a goal to investigate the relationships among multiple images [21] rather than just looking at one entity. Advancements in kinship verification and tri-subject kinship verification help further this goal as we are directly investigating the relationships among multiple people.

## 1.3 Related Work

The first paper in this field used handcrafted features such as color, facial features, and distances as well as a Histogram of Gradients vector [5]. A classification accuracy of 70.67% was obtained overall using K-Nearest Neighbors and 68.6% when a Support Vector Machine was used. It was also found that the average accuracy of human performance for the dataset used in the paper was 67.19% and the accuracies ranged from 50% to 90% [5]. As such, the classifiers were performing better than the humans that were tested in the paper.

Following this paper, approaches to the problems can be generally grouped into methods that use deep learning and those that use metric learning.

Deep learning approaches tend to use Siamese models in their approach. For example, a Siamese Convolutional Neural Network (CNN) was utilized by inputting the two input face images into a SqueezeNet network trained on VGGFace2 which creates a feature vector for each image [17]. Using a similarity criterion, a new feature vector is created using the two feature vectors and a fully connected layer. A sigmoid activation function then creates the predicted similarity score. This approach yielded an average test accuracy of 67.66% over all of the relationships that were in the dataset. These networks have three stages, a feature extraction stage, which uses a pre-trained CNN on each input image to extract the features onto a vector, a feature fusion stage which combines the feature vectors, and a similarity quantization stage to get the similarity score of the inputs.

A widely-used approach for the problem is to use metric learning. One of the first papers that used metric learning was the paper on Neighborhood Repulsed Metric Learning [15]. It aims to find a metric that minimizes the distance between the vectors of images that have a kinship relationship and maximize the distances between the vectors of pairs of images that don't have a kinship relationship. The Discriminative Multimetric Learning method [30] extends this model to use multiple feature vectors instead of just one for each image in the training and testing set. The paper that will be implemented in this dissertation, Weighted Graph Embedding-Based Metric Learning for Kinship Verification, is a state-of-the-art metric learning approach to the problem [12].

## 1.4 Project Overview

The project aims to reproduce the results of the paper, Weighted Graph Embedding-Based Metric Learning for Kinship Verification [12], or WGEML for short, and to verify the accuracies that were obtained. I also aim to explore the environments the model could work under and infer properties of the model and the datasets used from these environments. I was able to achieve the following:

1. I implemented WGEML and was able to obtain accuracies that are within an acceptable range of the original, as shown in section 4.1. Furthermore, I also used ROC curves in section 4.3 to further evaluate the models that were created. The mean ROC curves are created for each dataset as well as for each relationship.
2. I performed ablation studies on the face descriptors and the results show diminishing returns as more face descriptors are used for the model, which is in section 4.5.1.
3. I used the models that were created to discover biases in the datasets that were used. These biases came from the fact that the pairs of images in the KinFaceW-II and TSKinFace datasets came from the same photograph. I then explore the consequences of these biases that were found to be present in KinFaceW-II and TSKinFace in section 4.4.
4. I further tested the face descriptors by replacing VGG with a smaller CNN in which the implementation is discussed in section 3.3.3 and found that the model does worse but only by 5% at most in section 4.6.

Chapter 2 explains much of the needed technical background for the implementation, from what metric learning is to each of the face descriptors used and what face descriptors are. Chapter 3 then discusses the specifics of how WGEML and the extension that used a different network than VGG was implemented. Chapter 4 then discusses all of the results that were obtained from experimentation with WGEML and the implications of the results.

# Chapter 2

## Preparation

The approach that I'm using requires the use of metric learning, which is discussed in section 2.4, as well as the use of multiple face descriptors, which is discussed in section 2.3. I then discuss the core requirements and the extension in 2.6 and the prerequisites for the project.

### 2.1 Convolutional Neural Networks

To talk about Convolutional Neural Networks, first, we must discuss what an artificial neural network (ANN) is. An ANN is a collection of connected nodes. These are organized into layers such that there is an input and output layer as well as multiple layers in between which do some computation. The layers are made up of the connected nodes.

A basic type of this is a Multilayer Perceptron (MLP) which is a set of layers like we see in figure 2.1 which have neurons in each of the layers which each output one value.

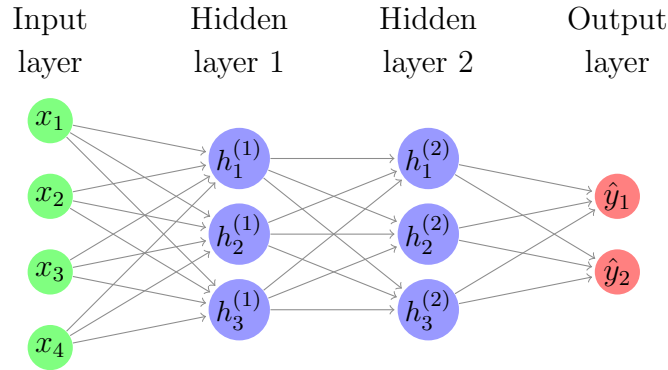


Figure 2.1: An example of an artificial neural network with two hidden layers

If we have that  $\mathbf{x}$  is the input layer,  $\mathbf{h}_i$  is the  $i$ th hidden layer and  $\mathbf{y}$  is the output layer, then an MLP can be written as:

$$\begin{aligned}
\mathbf{h}_1 &= f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \\
\mathbf{h}_2 &= f_2(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \\
&\dots = \dots \\
\mathbf{h}_i &= f_i(\mathbf{W}_i \mathbf{h}_{i-1} + \mathbf{b}_i) \\
&\dots = \dots \\
\mathbf{y} &= f_{n+1}(\mathbf{W}_{n+1} \mathbf{h}_n + \mathbf{b}_{n+1})
\end{aligned}$$

Where  $f_i$  represents the activation function for the corresponding layer and  $\mathbf{W}_i$  and  $\mathbf{b}_i$  are the weights and biases for the  $i$ th hidden layer which are the trainable parameters and  $n$  is the number of hidden layers.

A Convolutional Neural Network (CNN) is a type of neural network which is used to process data with a grid pattern which have some spatial locality, such as images. As opposed to an artificial neural network which is just composed of fully-connected layers, a CNN also uses convolutional layers and pooling layers [18]. Fully-connected layers have it so that each neuron in the layer have all of the connections to each of the input values, which is the same as in MLPs.

### 2.1.1 Convolutional Layers

A convolutional layer takes an input image and, as hyperparameters, takes the number of output filters, kernel size, and stride size in order to create a set of kernels and biases to create the output. A kernel is a matrix that is, usually, small and has sizes less than that of the input image but covers the depth of the entire input image. So, for example, if the input to a convolutional layer was  $32 \times 32 \times 64$ , then a kernel could have size  $3 \times 3 \times 64$  where the first two sizes are specified. The kernel is then applied to the image and, for each sub-array, the dot product of the kernel and the input sub-array is taken to be the output at that cell. An illustration of this is seen in figure 2.2.

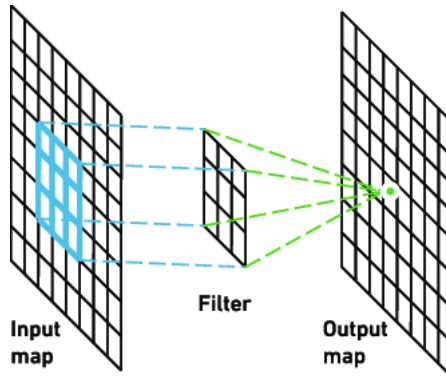


Figure 2.2: A convolution done on an input. Image sourced from Yakura et al. (2018) [28]

The sub-arrays the kernel must work on are separated by however much the stride is. So, if the stride was  $2 \times 2$ , then the sub-arrays would be separated from each other by 2 cells in each direction.

However, this produces a 2-dimensional array whereas we want a 3-dimensional volume as the output with a certain specified depth, which is the number of output filters. Thus, that

many kernels are created and used for the layer to create the output. These 2-dimensional arrays are then stacked on top of each other to create a 3-dimensional output. A bias is then applied and the activation function is applied to the output, which is usually ReLU which is described in section 2.1.3.1.

### 2.1.2 Pooling Layers

Pooling layers in a CNN tend to reduce the dimension of the representation which, in turn, reduces the size of the representation. There are multiple pooling layer types, such as max-pooling and average-pooling, which are the two types that are used in this project in VGG and CifarNet which are discussed in sections 2.3.4 and 3.3.3. These pooling layers take an input of size  $W_1 \times H_1 \times D_1$ , have hyperparameters which are the stride,  $S$ , and filter size,  $F$ , and output a tensor of size:

$$\frac{W_1 - F}{S + 1} \times \frac{H_1 - F}{S + 1} \times D_1$$

This output is calculated by taking each  $F \times F$  sub-array for each slice of the input and doing the corresponding operation on it and returning that as the output for that cell. So, max-pooling would take the max of each of the  $F \times F$  cells and average-pooling would take the average. The stride is the same as the stride in the convolutional layer. We can see an example of this with a max-pooling layer of stride 2 and has filter size  $2 \times 2$  in figure 2.3.

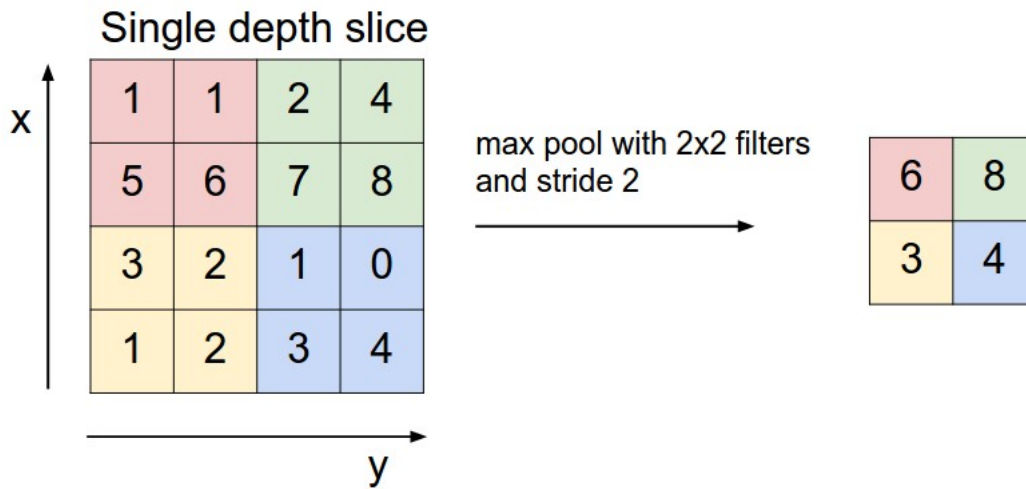


Figure 2.3: An example of the max-pooling operation done on a single slice of the input. Image sourced from Karpathy [10]

### 2.1.3 Activation Functions

There are a multitude of functions that can be applied to the output of a layer which each have different effects. These are called activation functions, two of them being ReLU and Softmax. Only nonlinear activation functions are considered since we want to be able to have nonlinear decision boundaries. If we had linear activation functions, the neural network would still be linear in nature and, thus, we wouldn't be able to deal with non-linear problems.

### 2.1.3.1 ReLU

ReLU applies the following on each output of a layer:

$$f(x) = \max(0, x)$$

In other words, it makes sure that all negative values become 0. This is simple to calculate and has the property that the derivative is either 1 or 0 which makes the gradient computation simpler. Due to the simplicity of the activation function, this helps speed up training a neural network.

### 2.1.3.2 Softmax

Another one is the softmax function which takes a vector,  $\mathbf{z}$  and the component-wise output is:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$$

Where  $N$  is the dimension of the input vector. In other words, it takes a vector and outputs a vector of the same size where the sum of the values of the vector adds up to 1, so you can interpret the individual values as probabilities. This is a useful activation function for the last layer of a neural network that solves a classification problem. For example, in the CifarNet model in section 3.3.3, the last layer has the softmax function applied to it in order for each of the 10 values to be interpreted as a probability that the input image was the corresponding class.

## 2.2 Face Detection

Face detection is the task of finding faces within a given image and returning the set of faces found. This differs from face recognition since the task is only to find the faces and not to recognize who the faces belong to. The project uses OpenCV's version of face detection which is a Haar-based cascade classifier and the specifics of which are discussed in Appendix A.1.

## 2.3 Face Descriptors

We wish to be able to create a mapping from a colored image into a vector to make computations easier and to be able to determine the similarity between images. These are called image descriptors for general images. When we try and map face images to vectors, we call these face descriptors. We want these face descriptor mappings to be able to match the same person in different poses and illuminant geometries to face descriptors that are close to each other in distance. Thus, these mappings try and capture color, texture, and shapes, for example. Multiple face descriptors can be made for a face image, each of which extracts different features from the face. We use the Local Binary Patterns, Histogram of Gradients, Scale-Invariant Feature Transform, and VGG face descriptors to extract features from each face.

### 2.3.1 Local Binary Patterns

The Local Binary Patterns (LBP) visual descriptor which is adapted for faces is a texture descriptor [1] which means that the algorithm attempts to describe the texture of the image, rather than anything to do with the color. As such, given an image, we must first convert it to grayscale. This can be done in various ways but in the project, the OpenCV method is used which maps each RGB pixel to the grayscale value<sup>1</sup>:

$$(R, G, B) \mapsto 0.299 \times R + 0.587 \times G + 0.114 \times B$$

Then, for each pixel in the image, a neighborhood of pixels is obtained from it. There are multiple ways to define this neighborhood, as shown in figure 2.4.

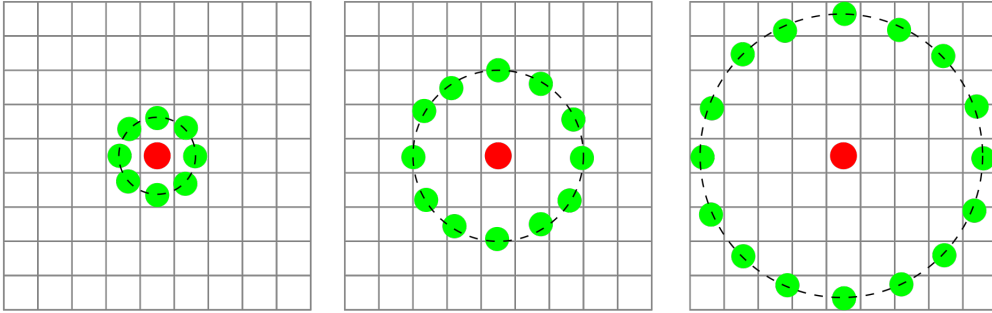


Figure 2.4: Potential neighborhoods of the pixel. Image sourced from Xiawi [27]

The simplest neighborhood, which is the neighborhood used in this project, is the direct neighbors of the pixel, which is the first neighborhood in figure 2.4. However, on the edges, some of the neighbors won't exist. To mitigate this, in the implementation, I pad the grayscale image with 0s on the outside of the image such that each pixel in the image has the same number of neighbors. Once we have our neighborhood of the pixel, we compare the grayscale value of the main pixel with each of the grayscale values in the neighborhood. For each neighboring value, if it is greater than the main pixel, we make it a 1, otherwise, we make it a 0. We are then able to create a binary number from the neighboring values. In practice, where the number is started from doesn't matter but in my implementation, the number starts from the left cell and goes counterclockwise. Applying this to figure 2.5, we get that our LBP value for this pixel would be  $01111000_2$  which, in decimal, is 120.

To summarize the LBP operator mathematically, given a coordinate in the image,  $(x, y)$  which has grayscale value  $g$ , a neighborhood of  $P$  points with radius  $R$  enumerated by  $g_p$  where  $p \in \{0, P-1\}$ , we have that [2]:

$$LBP_{P,R}(x, y) = \sum_{p=0}^{P-1} s(g_p - g) 2^p$$

Where:

$$s(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{else} \end{cases}$$

<sup>1</sup>Information taken from [https://docs.opencv.org/3.4/de/d25/imgproc\\_color\\_conversions.html](https://docs.opencv.org/3.4/de/d25/imgproc_color_conversions.html) using the RGB to Gray color conversion

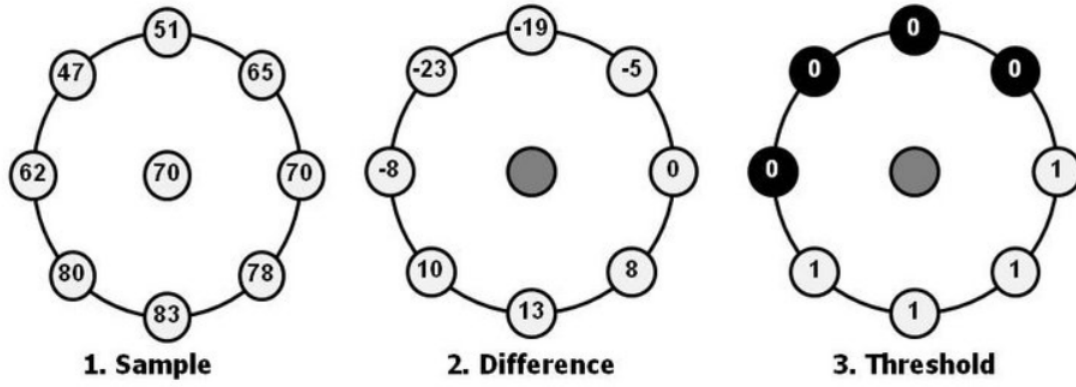


Figure 2.5: Example of the LBP operator on a pixel. Image sourced from Matti Pietikinen [20]

Once we get the values of the pixel, an added extension, which we do for face description, is to check whether it is a *uniform value* which we define as a value such that, in binary, there are only, at most, 2 bitwise transitions when traversed circularly. For example, 11000111 is a uniform value since it only transitions from 1 to 0 and 0 to 1 but 11001000 isn't since it has 4 bitwise transitions. With this extension, we have 58 possible uniform values that a pixel's LBP value can be and an extra value for the LBP value not being uniform. In other words, there are 59 LBP values that a pixel in an image can take.

Once each pixel in the image has an associated LBP value, we can split up the image into blocks. For our implementation, since our input images have size  $64 \times 64$ , we split it up into non-overlapping blocks of size  $8 \times 8$ , of which there are  $8 \times 8$ . For each block, we obtain a histogram of the LBP values that were in the block, which gives us a 59-dimensional vector. To obtain the vector for the entire image, each of these vectors are appended together and a 3776-dimensional LBP face descriptor is obtained for our case.

## 2.3.2 Histogram of Gradients

Histogram of Gradients (HOG) [3] is a face descriptor that uses the gradients of an image at each point to form a face descriptor. Given a colored image, we can think of the image as a function,  $I : \mathbb{N}_m \times \mathbb{N}_n \rightarrow \mathbb{R}^3$ , which takes a coordinate in the image and returns a vector of values, which correspond to the red, green, and blue values of the pixel. This then means that we are able to calculate the gradient of the image. The gradient of an image can characterize local object appearance and shape since the gradient of the image can be used to help find edges in the image. As such, it is useful to obtain such a histogram of gradients.

### 2.3.2.1 Calculating the Gradients

In order to compute the HOG vector for an image, we must compute the gradients of the image. As the image function is discrete, so we cannot analytically find the gradients, we must find approximations to do so. This is done by convolving specific kernels on the image. What this means is that the kernel,  $K$  which is a matrix, is applied to each pixel in the image and its neighbors and the values are multiplied with the corresponding value in the kernel and all of the values are then summed up to create the new value for the convolved image. For example, if we had the kernel:

$$K = \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$$

And the image:

$$I = \begin{bmatrix} 0 & 5 & 4 \\ 3 & 3 & 3 \\ 2 & 1 & 2 \end{bmatrix}$$

We get that the kernel convolved on the image is:

$$K * I = \begin{bmatrix} 0 \times 1 + 5 \times 2 + 4 \times 1 \\ 3 \times 1 + 3 \times 2 + 3 \times 1 \\ 2 \times 1 + 1 \times 2 + 2 \times 1 \end{bmatrix} = \begin{bmatrix} 14 \\ 12 \\ 6 \end{bmatrix}$$

We can then approximate the derivative of an image using a kernel being convolved on the image. To do so, we use the kernels:

$$G_x = \begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$$

$$G_y = \begin{bmatrix} -1 & 0 & 1 \end{bmatrix}^T$$

By convolving these kernels on the image, we are able to get an estimate for the derivative in the  $x$ -direction and in the  $y$ -direction, respectively. Abusing notation slightly where here the square of the matrix just means an element-wise square, and division just means element-wise division, we get that the magnitude of the gradient is:

$$\sqrt{(G_x * I)^2 + (G_y * I)^2}$$

And that the angles at each point are:

$$\tan^{-1}((G_y * I)/(G_x * I))$$

However, there are two things that need to be addressed before we move on to the rest of the algorithm. Firstly, there are still three channels for the image, so we have obtained the gradient in the  $x$  and  $y$  direction for each pixel and each channel. As such, we define the gradient of each pixel to be the gradient vector that has the maximum magnitude among the three channels.

Furthermore, the angles returned range between  $0^\circ$  and  $360^\circ$ . However, we require that the angles be “unsigned”, so the angle at each pixel,  $\theta_{(x,y)}$ , becomes:

$$\theta_{(x,y)} := \theta_{(x,y)} \bmod 180$$

### 2.3.2.2 Weighted Vote into Histogram

At this point, we have the magnitude of the gradient at each point in the image as well as the angle of the gradient. Now, similarly to LBP, we break up the image into blocks. We create a histogram for each of these blocks and append them together to make the face descriptor for the entire image. Unlike LBP, we don’t have a finite set of labels that each pixel can neatly fall into, however, since neither the magnitude nor the angles are discrete. Thus, first, the labels of the histogram are going to be the angles of the gradients. The labels will then be:

$$[0, 20, 40, 60, 80, 100, 120, 140, 160]$$

However, these aren't blocks with ranges of  $[0, 20)$ . Instead, for each pixel in the block, we distribute the magnitude of the gradient of the pixel between the angles that the angle falls between. Given a pixel with magnitude  $m$  and angle  $\theta$ , if  $0 \equiv \theta \pmod{20}$  then:

$$\text{histogram}[\theta/20] += m$$

Otherwise,  $\theta$  is between two angles,  $\phi_1$  and  $\phi_2$ , both of which are divisible by 20, where  $\phi_1 < \theta < \phi_2$ . In this case, we weight the amount that we add to each label based on how far away  $\theta$  is to the label. So, for the label  $\phi_1$ , we add to the label associated with  $\phi_1$  the value:

$$\frac{\phi_2 - \theta}{20} \times m$$

And, similarly for  $\phi_2$ :

$$\frac{\theta - \phi_1}{20} \times m$$

In other words, the closer the angle is to the label, the more of the magnitude is contributed to the label's histogram value. As a caveat, if  $\phi_2 = 180^\circ$ , we treat  $\phi_2$  as  $180^\circ$  for the sake of this calculation but we add the value to the label 0 since  $0 \equiv 180 \pmod{180}$ .

We can do this for each pixel in the block and, thus, we are able to get a 9-dimensional vector for each block. In the project, the image is split up into  $16 \times 16$  blocks of size  $4 \times 4$  first and then  $8 \times 8$  blocks of size  $8 \times 8$  next and each of these blocks contributes a histogram to the overall vector which leads us with a face descriptor with dimension:

$$16 \times 16 \times 9 + 8 \times 8 \times 9 = 2880$$

### 2.3.3 Scale-Invariant Feature Transform

SIFT [13] is another way of obtaining face descriptors, although the original SIFT algorithm differs from how it is used in the project. The algorithm, generally, is split into finding the keypoints, fine-tuning the keypoints, assigning an orientation to each keypoint, and then getting a descriptor from each keypoint. The method of obtaining each of the keypoints in an image, which are just points of interest which can be near important features of the image, is explained in Appendix A.2.

Once we have the keypoints in the image, a  $16 \times 16$  window is obtained around the keypoint and divided into 16 blocks. For each block, a histogram of the gradients is taken with 8 bins which are then appended together which gives a 128-dimensional vector for the keypoint. The SIFT descriptor of the image is then the vectors for each keypoint appended together to create a  $128n$ -dimensional vector where  $n$  is the number of keypoints in the image.

The actual implementation of SIFT that is used differs from the original version which is discussed in section 3.3.1.

### 2.3.4 VGG

Another way of getting face descriptors for an image is to use a CNN. One CNN we can use is VGG. The original VGG network [23] is a CNN which takes in a  $224 \times 224$  colored image and outputs a probability vector of size 1000 for the ImageNet classes. In other words, if `out` is the output of the model, `out[i]` corresponds to the probability that the  $i^{\text{th}}$  ImageNet class is the

class of the input image. There are 5 configurations of the network, which have varying depth and, as a result, varying parameters. The architecture of the second deepest configuration is configuration D in Figure 2.6 which is used in this implementation.

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Figure 2.6: The different configurations that VGG can have. Image reproduced from Simonyan et al. (2015) [23]

By training this model on the ImageNet dataset, the model is able to get a top-5 classification error of 7.5%, which means that for 7.5% of the entries in the test set, none of the top 5 classes that the image could be were the actual image.

However, in order to use the model for face description, we must change the output layer and the training dataset [19]. Instead of training on general objects in an image, we train the model on faces, specifically celebrity faces. The celebrities are curated to a list of 2622 people in which 2000 images are obtained for each celebrity and the images are then curated. The curated set of images constitutes the training set for the VGG model for faces. This dataset<sup>2</sup> was created by the authors of the VGG model for faces paper [19].

The VGG model for faces is then the model described above with the softmax layer having a dimension of 2622. The model is then trained with the dataset we described which allows us to get rid of the softmax layer and use the output of the last fully-connected layer as our 4096-dimensional face descriptor for the image.

## 2.4 Metric Learning

We wish to measure the similarity between a pair of faces in order to determine whether they are related or not. One way to approach this is to use similarity learning in which the goal is

<sup>2</sup>The dataset is available at [https://www.robots.ox.ac.uk/~vgg/data/vgg\\_face/](https://www.robots.ox.ac.uk/~vgg/data/vgg_face/)

to learn a similarity function in order to measure how similar two objects are. However, we often use distance to help measure the similarity between data points. By learning what this distance function ought to be, we would be able to find a better similarity function so finding this distance function is the goal of metric learning [24]. To go more in depth, we must first define a few concepts.

Given a non-empty set  $A$ , we define a *distance* over  $A$  as a function  $d : A \times A \rightarrow \mathbb{R}$  such that the following holds:

1. **Non-Negativity:**  $\forall x, y \in A, d(x, y) \geq 0$
2. **Coincidence:**  $\forall x, y \in A, d(x, y) = 0$  if and only if  $x = y$ .
3. **Symmetry:**  $\forall x, y \in A, d(x, y) = d(y, x)$
4. **Triangle Inequality:** The triangle inequality must hold. That is,  $\forall x, y, z \in A$ :

$$d(x, y) + d(y, z) \geq d(x, z)$$

Given the definition of a distance, we define a *Mahalanobis distance* that corresponds to the positive semidefinite matrix  $M$  to be a distance function,  $d_M : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , where  $d$  is the number of dimensions the input vector has which is defined as:

$$d_M(x, y) = \sqrt{(x - y)^T M (x - y)}$$

If  $M$  is positive semidefinite, then  $M$  must be decomposable into  $B^T B$  where  $B$  is a real matrix. Thus, we can also write the Mahalanobis distance as:

$$d_M(x, y) = \sqrt{(x - y)^T M (x - y)} = \sqrt{(x - y)^T B^T B (x - y)} = \sqrt{(Bx - By)^T (Bx - By)}$$

So the Mahalanobis distance is essentially a Euclidean distance after applying a linear transformation to each of the input vectors. This makes it so that we only have to learn what this linear transformation is which is simpler to find than searching through the set of all possible distance functions to see which is the best overall.

The Mahalanobis distance is a *pseudometric*, which is a distance function in which the coincidence property doesn't necessarily hold, which occurs when the matrix  $M$  isn't full-rank. A matrix is *full-rank* when the rank of the matrix is either equal to the number of rows or columns in the matrix.

Thus, given a dataset  $\mathcal{X} = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$  which has the sets:

$$S = \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} \mid x_i \text{ and } x_j \text{ are similar}\}$$

$$D = \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} \mid x_i \text{ and } x_j \text{ are not similar}\}$$

We wish to find the distance metric  $M$  such that we can minimize some loss function,  $l$ :

$$\min_M l(d_M, S, D)$$

In other words, in metric learning, we wish to find out what the matrix  $M$  should be defined as in order to make data points which are similar closer together and data points which aren't similar further away, which is encapsulated by the loss function. We can visualize this in figure 2.7.

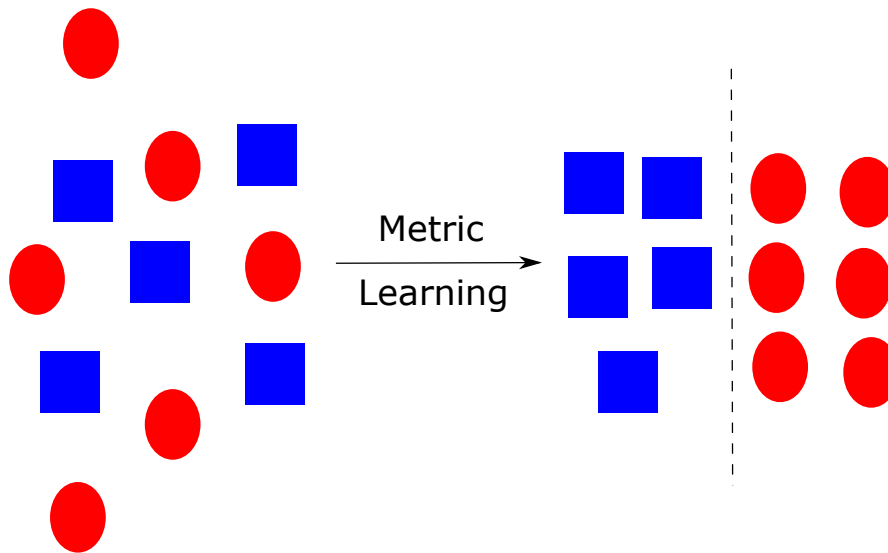


Figure 2.7: A high-level view of metric learning

## 2.5 K-Nearest Neighbors

K-Nearest Neighbors is an algorithm that finds the  $K$  nearest neighbors for each datum in a given set of data. Given a list of datapoints in Euclidean space, we wish to find, for each datapoint, the  $K$  closest datapoints to it. The `sklearn` library is used to implement it which uses a brute force algorithm, a ball tree or a *kd*-tree depending on the input.

## 2.6 Requirements Analysis

The main requirement of the project is stated in Appendix C which is to replicate the results from Liang et al. (2019) [12] within a 15% error range or to reject the results. In other words, the project should implement the Weighted Graph Embedding-Based Metric Learning (WGEML) algorithm. The core project can be broken down into the following requirements:

- **Face Detection:** Given an image, the faces in the image must be identified and saved as a  $64 \times 64$  image on disk. This is a low risk and low priority requirement as this is mainly just using OpenCV methods and isn't necessary for testing the model, as the images are already faces of size  $64 \times 64$ .
- **Face Descriptors:** The face descriptors, LBP, HOG, SIFT, and VGG, must be implemented such that, given an image, each of these face descriptors should be able to be obtained from them.
- **WGEML:** The WGEML algorithm must be implemented such that, given a relationship, the positive pairs of faces, the negative pairs of faces, if any, it returns the distance metric matrix for each face descriptor used and how much each face descriptor should be weighted in the prediction step.
- **Prediction:** Given the model obtained from the WGEML algorithm, a pair of images, and a relationship, it must return how similar the images are and whether the people in the images have the kin relationship given.

- **End-to-end Replication:** Running the project end-to-end on the given datasets either gives similar results to the original paper or rejects the original results.

The extensions can be summarized as follows:

- Replace the VGG face descriptor with a smaller network to see how it affects accuracies.
- Look at how biases in the datasets affect the results and identify these biases.
- Use different combinations of face descriptors for WGEML to see how it affects the accuracy of the model.

**Analysis** There were two main phases for the project, the phase to finish the core requirements and one for the extensions. Each of the first 4 core requirements were split up into their own modules and thus constituted their own sub-phase in the main part of the project. They were also linear in that each requirement depended on the output of the last modules. The last requirement constituted a medley of scripts that were used to integrate the functions that were created for each of the requirements with the datasets, which required a data preparation module.

## 2.7 Software Engineering Practices

### 2.7.1 Starting Point

Before I started my project, I had knowledge in Python and surface-level knowledge of Keras and Tensorflow. I had also worked with face recognition before. Furthermore, I had used Numpy many times before, so I had enough knowledge about it to use it comfortably in my project. However, I hadn't ventured into computer vision before the project aside from knowing very generally what a face descriptor is.

In terms of what was already available, face detection had been done by OpenCV already, thus my code used the classifier was pre-trained by OpenCV. Furthermore, VGG was also already pre-trained, though I had to modify it slightly to work with my implementation. However, the rest of the project was created from scratch in terms of the code and how the project was structured.

### 2.7.2 Tools Used

The project was written in Python 3.6. To separate the environment that the project needs from my local environment, I used a virtual environment to contain the installations of the required libraries. I also had a requirements.txt file to contain the names and versions of the packages that were used. I also used the following libraries:

1. **Numpy:** Using `numpy` helped increase the performance of many parallel computations, such as matrix multiplication, and provided a simple interface to do these with.
2. **Keras:** I used Keras to create the VGG model and the CifarNet model. We then use it to train the CifarNet model and then make predictions for both models. The weights of VGG were already pre-computed so no training was necessary for VGG.

3. **OpenCV**: I used OpenCV for general image processing, such as converting an image into grayscale or loading in images, and for face detection.
4. **Sklearn**: This library provided an implementation for K-Nearest Neighbors, PCA, and obtaining folds for cross-validation.
5. **Scipy**: This library provided a function that solved the general eigenvalue problem, given two `numpy` arrays which were used in the main WGEML algorithm.

Finally, Git and GitHub were used for version control and backups. Branches were created for each feature that was to be added to the project and I merged them into the master branch once enough testing was done, whether it was unit testing or integration testing. On GitHub, Dependabot was used to notify me of vulnerabilities in the packages I used and to tell me to upgrade the version in the `requirements.txt` file.

### 2.7.3 Datasets

I used the KinFaceW-I, KinFaceW-II, and TSKinFace datasets to train the model.

The KinFaceW datasets<sup>3</sup> [14] [15] are composed of face images from the internet which include public figures as well as their parents or children. Both datasets contain no restriction on pose, lighting background, expression, age, ethnicity, or partial occlusion. The main difference between the datasets is that the pairs of face images that have a kin relationship in KinFaceW-I are from different images whereas, in KinFaceW-II, they are from the same image. In terms of the specifics of the datasets, they both support four kin relationships, Father-Son (FS), Father-Daughter (FD), Mother-Son (MS), and Mother-Daughter (MD). Each face image has size  $64 \times 64$ . The datasets also contain pre-computed 5 folds for cross-validation for each relationship which contained the names of the pairs of images that had the relationship and those that didn't, henceforth positive and negative pairs respectively. Finally, the dataset has 2 main settings which are used: the restricted setting in which only positive pairs of images are used and the unrestricted setting in which negative pairs of images are used.

The TSKinFace dataset<sup>4</sup> [21] contains images for the relationships, Father-Mother-Son (FMS), Father-Mother-Daughter (FMD), and Father-Mother-Son-Daughter (FMSD). The dataset contained folders for each relationship and a positive set comprised of the images in which the names of the images had the form “[relationship]-N-[member].jpg” in which “relationship” was the relationship,  $N$  was a consistent number and “member” refers to which member of the relationship they were. For example, “FMS-10-F.jpg”, “FMS-10-M.jpg”, and “FMS-10-S.jpg” would form a positive triplet. The dataset only contained the images in this form so negative pairs of images and the folds for cross-validation had to be computed in the project.

I examine the effects of the WGEML algorithm on the FS, FD, MS, MD, FMS, and FMD relationships, as defined above.

---

<sup>3</sup>Datasets obtained from <https://www.kinfacew.com/download.html>

<sup>4</sup>Dataset obtained from [http://parnec.nuaa.edu.cn/\\_upload/tp1/02/db/731/template731/pages/xtan/TSKinFace.html](http://parnec.nuaa.edu.cn/_upload/tp1/02/db/731/template731/pages/xtan/TSKinFace.html)

### 2.7.4 Testing

Throughout the project, I created unit tests for any modules and for each function in the module. I would only end up merging a feature branch into the master branch if all of the unit tests passed for that feature. There were minor problems with unit testing the modules that created the VGG model and CifarNet model which is talked about in section 4.7. However, though there were problems with testing the neural networks, I was able to unit test WGEML fully due to the nature of the algorithm being deterministic given the same training dataset.

Along with this, I did integration tests in the form of feeding the scripts small, controlled inputs to see if the scripts would output what was expected. These scripts each used functions from different modules and did a part of the workflow.

Finally, an end-to-end test was done once each script was tested individually which came in the form of trying the small, controlled input for the first script and then running it through each of the scripts successively.

### 2.7.5 Licensing

The SIFT algorithm had been patented in the US. However, the patent expired last year, and the patent was only to stop commercial use of the algorithm, whereas this is an academic use of the algorithm.

Furthermore, I am able to use VGG for non-commercial purposes under the Creative Commons Attribution License.

The `numpy`, `sklearn`, and `scipy` libraries as well as the version of OpenCV I used are all able to be used since they are BSD licensed. Furthermore, `keras` is able to be used since it is MIT licensed.

# Chapter 3

## Implementation

In this chapter, the structure of the project is examined in detail in section 3.1, and each of the modules are discussed in sections 3.2, 3.3, 3.4, and 3.5. I also aimed to replace VGG with another face descriptor and the implementation of this is discussed in section 3.3.3. Finally, how the project is broken down and run from end-to-end is discussed in section 3.6.

### 3.1 Repository Overview

The repository is shown in Figure 3.1

The `src` folder is split up into the different overarching modules which contain the functions needed, `data_preparation`, `face_descriptors`, `face_detection`, `prediction`, and `WGEML`. The face detection module implements face detection by using OpenCV which uses the method described in Appendix A.1. A pre-trained OpenCV Haar Cascade Classifier was used<sup>1</sup> and the `CascadeClassifier` class was utilized to take advantage of the pre-trained model.

The rest of the modules are discussed in sections 3.2, 3.3, 3.4, and 3.5 in which each of these modules were implemented mainly from scratch, aside from the VGG implementation in which a pre-trained model was used and modified for my purposes. The `test` folder contains all of the unit tests I wrote for each of the files in the aforementioned modules.

The `scripts` folder contains the scripts I wrote that would be run directly with Python which integrates the modules together to run the overall project. The majority of the computation is done in the modules and the scripts mainly run the functions from the modules with the proper data.

Each folder in the `src` folder includes an `__init__.py` file for it to be able to be referenced by other folders in the project. The `.coveragerc` file was used to make sure the test files, `venv` and `__init__.py` weren't included in the coverage report. The `Makefile` was mainly used as a way to create shortcuts to run certain scripts multiple times, such as training on a given dataset for each relationship, or running all of my unit tests and creating a coverage report.

The `data` folder contains the three datasets which contain the images and the metadata about the images, such as which are positive and negative pairs. Certain extra information is saved to these folders as well which is discussed in section 3.2.4.

Finally, the `out` folder contains CSVs that contain the accuracies of certain experiments which are also discussed in section 3.2.4 as well as the ROC curves created from the `ROC_curves.py` script.

---

<sup>1</sup>[https://github.com/opencv/opencv/blob/master/data/haarcascades/haarcascade\\_frontalface\\_default.xml](https://github.com/opencv/opencv/blob/master/data/haarcascades/haarcascade_frontalface_default.xml)

```

KinRecognitionWGEML
├── data/
├── out
│   ├── ablation_studies/
│   ├── pairwise_accs/
│   ├── ROC_dataset.png
│   └── ROC_rel.png
├── src
│   ├── data_preparation (The folder the data preparation functions as explained in section 3.2)
│   │   ├── PCA.py
│   │   ├── prep_cross_valid.py
│   │   ├── properly_formatted_inputs.py (Helps create the inputs for training WGEML)
│   │   └── save_and_load.py
│   ├── face_descriptors (The folder has the implementations of each of the face descriptors)
│   │   ├── CifarNet.py
│   │   ├── HOG.py
│   │   ├── LBP.py
│   │   ├── SIFT.py
│   │   └── VGG.py
│   ├── face_detection
│   │   ├── face_detection.py
│   │   └── haarcascade_frontalface_default.xml (The classifier obtained from OpenCV)
│   ├── prediction/predictor.py
│   ├── scripts
│   │   ├── ablation_study.py
│   │   ├── get_pairwise_accuracies.py
│   │   ├── preprocessing_fds.py (Obtains face descriptors for each image in a dataset)
│   │   ├── preprocessing_TSK.py (Creates TSKinFace negative pairs and the folds)
│   │   ├── ROC_curves.py
│   │   ├── testing.py (Used to test a set of models with the corresponding test images)
│   │   └── training.py (Used to train the models for each fold)
│   ├── test/
│   └── WGEML
│       ├── constants.py
│       └── WGEML_training.py (Contains the functions to create and train WGEML)
├── venv/
├── .coveragerc
├── .gitignore
├── Makefile
├── README.md
└── requirements.txt

```

Figure 3.1: Folder structure of the project

## 3.2 Data Preparation

### 3.2.1 Cross-Validation

I split the datasets for each relationship into 5 folds in which each fold within a relationship has around the same number of pairs. Once I have the 5 separate folds, I obtain 5 different

train/test splits since I use each fold as a test set and the remaining 4 folds as the training set. For example, if we had the folds  $\{A, B, C, D, E\}$ , our splits are then shown in table 3.1.

Training Set	Testing Set
$B, C, D, E$	$A$
$A, C, D, E$	$B$
$A, B, D, E$	$C$
$A, B, C, E$	$D$
$A, B, C, D$	$E$

Table 3.1: Training/Testing Splits

I then run the algorithm end-to-end for each train/test split and average the accuracies in the end. For the KinFaceW datasets, as the images were already split into folds by the dataset, that was used in the project, whereas with TSKinFace, it had to be generated randomly, which was done using `sklearn` which had a `KFold` function.

### 3.2.2 Positive and Negative Pairs

Within each dataset, there are pairs of images that either have the kin relationship or don't, which we call positive and negative pairs respectively.

KinFaceW had these negative and positive pairs prepared in the `mat` files that came with the dataset. However, as TSKinFace only came with the images, only the positive images could be found from them. Thus, negative pairs for the dataset had to be created randomly and the negative pairs generated were saved on disk. For each relationship, the number of negative pairs was set to be equal to the number of positive pairs in the corresponding set, which was 404 for the training set and 101 for the test set. Then, two images were picked out randomly such that they didn't belong to the same positive pair, which meant that they didn't have the same photo ID number, as explained in section 2.7.3. As we wanted to create pairs for each relationship, if we had the relationship FS, then we'd want to pick a certain number of negative pairs from the FMS and FMSD image sets such that the ratio of negative pairs picked from each set is similar to the ratio of positive pairs in each set. To create a pair from a given set, a number was randomly picked from 1 to  $n$  where  $n$  was the number of pairs in that set and another number was picked from 1 to  $n - 1$ . If the two numbers were equal, the second number became  $n$ , thus creating 2 random numbers that aren't equal.

Finally, each of the KinFaceW datasets had a setting, restricted or unrestricted, in which the restricted setting meant that no negative pairs are used in the training splits, and unrestricted means that they are used [15]. TSKinFace didn't have this type of setting, however, but the differences in whether negative pairs were or weren't used in the training process was still explored.

### 3.2.3 Dimensionality Reduction

In order to save on computation time for the training process as well as on disk space, Principal Component Analysis (PCA) is used to reduce the dimensions of each of the face descriptors, which is described in detail in Appendix A.3. In the implementation, after each image had its face descriptors obtained for a dataset, for each type of face descriptor, PCA is used to

reduce the dimension of each face descriptor to 100. This is done using the `sklearn` function which implements the probabilistic PCA model [25]. Although the dimensionality of the face descriptor decreases, we are still able to keep the most important features of the original dataset which means that the face descriptors are still usable, even though the dimensionality has been vastly decreased.

### 3.2.4 Saving Results to Disk

Running the project end-to-end each time would take too long and, to better explore some of the results that were taken, the pipeline needed to be split into modules. I split it up into the preprocessing stage, the training stage, and the testing stage. These stages are explained in further detail in section 3.6. What was saved to disk is as follows, split by stage:

1. *Preprocessing:*

- **Face Descriptors:** Under the corresponding dataset folder in the data folder, each face descriptor for each image is stored as a pickle file which contains a map from the image name to the PCA reduced face descriptor. The files are split up by face descriptor type so `VGG_face_descriptors.pkl` would correspond to the VGG face descriptors of each image.
- **TSKinFace Splits and Negative Pairs:** Since the TSKinFace splits and negative pairs are created randomly and both are used in the training and testing step, it is imperative that these are saved on disk since trying to recompute them would result in different splits and negative pairs.

2. *Training:*

- **WGEML Output:** The output of the training is saved on disk for each fold of the relationship. It is saved to the `data` folder under the corresponding dataset and setting. Each file is called `[relationship]_out.pkl` which is under the folder path `data/[dataset]/WGEML_out/[setting]/[relationship]_out.pkl`.

3. *Testing:*

- **Ablation Studies Results:** The accuracies obtained from the ablation studies are stored in a CSV for each dataset and setting configuration which is stored in the path `out/ablation_studies/[dataset]_[setting].csv`.
- **Pairwise Accuracies:** Similarly, the accuracies obtained from changing which dataset the test data comes from for each model is saved as a CSV to the path `out/pairwise_accs/[dataset]_[setting].csv`.
- **ROC Curves:** The ROC curves created were saved as images in order to be analyzed in detail and shown in this dissertation.

## 3.3 Face Descriptors

Each of the LBP and HOG face descriptors were implemented exactly as mentioned in section 2.3 without the use of any libraries other than Numpy and OpenCV. However, some extra details need to be explained for the SIFT and VGG face descriptors.

### 3.3.1 SIFT Implementation

The major difference between the original version of SIFT described in the original SIFT paper [13] and what was needed was that keypoints of each image weren't calculated. Instead, the descriptor that was obtained from each keypoint, which is described in section 2.3.3, is obtained from dividing the image into  $7 \times 7$  overlapping patches on a  $16 \times 16$  grid.

This leads to there being a 128-dimensional vector for each patch which leads to a  $128 \times 7 \times 7 = 6272$ -dimensional vector for each image. This also helps to standardize the dimension of the vector for each image as different images can have a different number of keypoints which would lead to different dimensionalities of the vectors. This decision was made by the WGEML paper in their description of what face descriptors were used.

Due to the differences, this version of SIFT had to be implemented directly rather than using a library function like OpenCV's SIFT implementation.

### 3.3.2 VGG Implementation

The original VGG network that was trained on ImageNet had 138 million parameters [23] which would then increase if we were to increase the number of outputs of the softmax layer for the face network. This would take too long to train in practice. However, the pre-trained weights for the VGG face network were uploaded to the authors' site<sup>2</sup>. Therefore, it was easier to download the weights from the website and reformat it as needed than to train the model myself.

The VGG face model was implemented using `keras` to create the base model without the weights. To create the pre-trained model, it had to be checked whether the weights were already on disk or not. If they weren't, they would be downloaded and modified in order to fit in the model and then saved onto disk in the proper file. The modification that was required was to reshape the weights to the dimension that `keras` wanted rather than what they originally were, since the weights were originally to be used with Torch.

From there, the weights would be on disk so, to get a face descriptor for a set of images, the VGG face model with the softmax layer was instantiated, then the weights were loaded into the model before the model without the softmax layer is returned. This model is then used to get the face descriptors for a given set of images. Since VGG takes images of size  $224 \times 224$ , the images need to be resized to  $224 \times 224$  before they could be used as an input into the model. Thus, the model took a `numpy` array of size  $(n, 224, 224, 3)$  and outputted a list of face descriptors that had size  $(n, 4096)$ .

### 3.3.3 CifarNet Extension

In the CifarNet extension, I replaced the VGG network with a smaller network that was trained on CIFAR-10 [11], which is called CifarNet<sup>3</sup>, until it was able to achieve a 91% accuracy with it. The architecture is shown in figure 3.2. After each of the convolutional layers, a batch normalization layer was applied with momentum 0.9997 and after the fourth and seventh batch normalization layer, a dropout layer was applied. The last layer after the flatten layer is a fully-connected layer with the softmax function as its activation function.

<sup>2</sup>[https://www.robots.ox.ac.uk/~vgg/software/vgg\\_face/](https://www.robots.ox.ac.uk/~vgg/software/vgg_face/)

<sup>3</sup>Architecture taken from <https://github.com/deep-fry/mayo/blob/master/models/cifarnet.yaml>

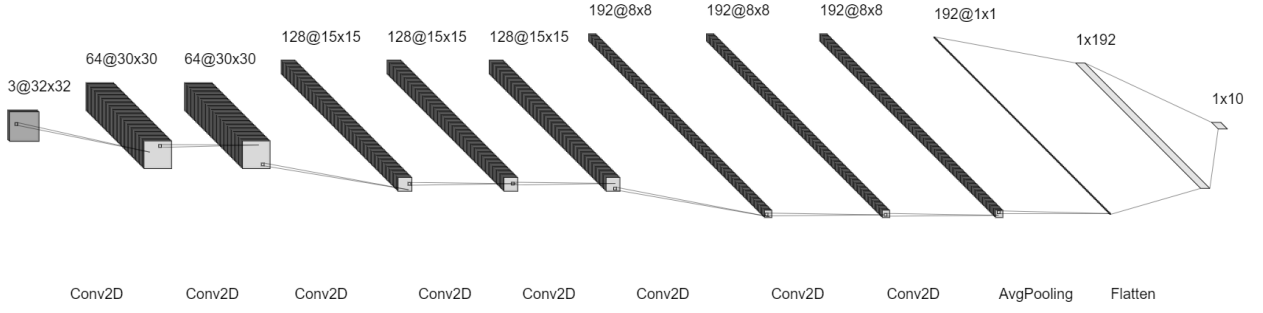


Figure 3.2: The architecture of the CifarNet model

This was achieved by using Stochastic Gradient Descent as the optimizer and an adaptive learning rate and using a learning rate of 0.05 for the first 50 epochs before shrinking it to 0.005 for the next 25 epochs and finally 0.0005 for the last 25. The face descriptor is then obtained similarly to VGG where the softmax layer is taken out and the layer before is used as the face descriptor.

### 3.4 WGEML

The main algorithm used in the project is Weighted Graph Embedding-Based Metric Learning (WGEML) [12]. The algorithm is a form of metric learning which takes in, as input, a training positive set for each face descriptor,  $\mathcal{S}^p = \{(\mathbf{x}_i^p, \mathbf{y}_i^p) \mid 1 \leq i \leq N\}$  where  $N$  is the number of image pairs in the set and a negative pair set  $\mathcal{D}^p = \{(\mathbf{x}_i^p, \mathbf{y}_j^p) \mid 1 \leq i \leq N, j \neq i\}$ . Let there be  $M$  face descriptors for each image in the set. It also takes in a tuning parameter  $r$  and a neighborhood size  $K$  as input but those have been experimentally found to be 5 for each, each of which will be discussed later in this section. This section first gives an overview of the algorithm in Algorithm 1 before discussing the problem statement in section 3.4.1 and then the specifics of the approach in section 3.4.2.

The overall algorithm is written in pseudocode in Algorithm 1 from [12].

#### 3.4.1 Problem

Given these inputs, the goal is to find the distance metrics and weights for:

$$\begin{aligned} d^2(\mathbf{x}_i, \mathbf{y}_i) &= \sum_{p=1}^M w_p (\mathbf{x}_i^p - \mathbf{y}_i^p)^T \mathbf{A}_p (\mathbf{x}_i^p - \mathbf{y}_i^p) \\ &= \sum_{p=1}^M w_p d_{\mathbf{A}_p}^2(\mathbf{x}_i^p, \mathbf{y}_j^p) \end{aligned}$$

Where  $\mathbf{x}_i^p$  represents the  $p$ th face descriptor of the image  $\mathbf{x}_i$ ,  $w_p$  is a weight and  $\mathbf{A}_p$  is a  $D \times D$  semidefinite positive matrix, where  $D$  is the dimensionality of the corresponding face descriptor. This distance function finds the distance between each pair of face descriptors and weights them accordingly. We wish to find  $\mathbf{A}_p$  such that the between-class variance is maximized, and the within-class variance is minimized. This can be formalized as the optimization problem:

**Algorithm 1** WGEML

**Inputs:** The positive and negative pair sets  $\mathcal{S}^p$  and  $\mathcal{D}^p$  for each face descriptor, the tuning parameter  $r$  and the number of neighbors to be considered  $K$

**Outputs:** The matrices  $\mathbf{U}_p$  and the weights  $w_p$  for each face descriptor.

- 1: Initialize  $\mathbf{w}$
- 2: Initialize  $\mathbf{U}$
- 3: **for**  $p \in \{1, \dots, M\}$  **do**
- 4:   Use KNN to find the nearest neighbors of each  $\mathbf{x}_i^p$  and  $\mathbf{y}_i^p$ .
- 5:   Create the matrices  $\mathbf{S}_p, \mathbf{D}_p, \mathbf{D}_{1p}, \mathbf{D}_{2p}$  using equations 3.1, 3.2, 3.3, and 3.4
- 6:   Regularize the matrix  $\mathbf{S}_p$  using equation 3.5
- 7:   Solve the eigenvalue problem in equation 3.6 to get  $\mathbf{U}_p$  and append it to  $\mathbf{U}$
- 8:   Compute  $w'_p$  using equation 3.7 and append it to  $\mathbf{w}$
- 9: Divide each  $w'_p$  in  $\mathbf{w}$  by  $\sum_{i=1}^M w'_p$ , as in equation 3.8, to get the vector  $\mathbf{w}$  such that the values sum to 1
- 10: Return  $\mathbf{w}$  and  $\mathbf{U}$ .

$$\begin{aligned}
 \max_{\mathbf{A}, \mathbf{w}} \mathcal{F} = & \sum_{p=1}^M w_p^r \left[ \frac{1}{2} \left( \frac{1}{NK} \sum_{i=1}^N \sum_{n_1=1}^K d_{\mathbf{A}_p}^2(\mathbf{x}_i^p, \mathbf{y}_{i,n_1}^p) + \frac{1}{NK} \sum_{i=1}^N \sum_{n_2=1}^K d_{\mathbf{A}_p}^2(\mathbf{x}_{i,n_2}^p, \mathbf{y}_i^p) \right) \right. \\
 & \left. + \frac{1}{N} \sum_{i=1, j \neq i}^N d_{\mathbf{A}_p}^2(\mathbf{x}_i^p, \mathbf{y}_j^p) \right] / \frac{1}{N} \sum_{i=1}^N d_{\mathbf{A}_p}^2(\mathbf{x}_i^p, \mathbf{y}_i^p) \\
 \text{s.t. } & \sum_{p=1}^M w_p = 1 \\
 & \forall p \in \{1, \dots, M\}, w_p \geq 0
 \end{aligned}$$

Where  $\mathbf{y}_{i,n_1}^p$  is the  $n_1$ th nearest neighbor of  $\mathbf{y}_i$  and similarly for  $\mathbf{x}_{i,n_2}$ . We have  $w_p^r$  in order to avoid over-fitting and use information from each face descriptor. This optimization function tries to minimize the denominator which, in turn, pulls the samples that have the kin relationship together and maximizes the numerator which means that the pairs which don't have the kin relationship are pushed further away from each other as well as those of their neighbors.

### 3.4.2 Approach

Now that the problem is defined, we break down  $\mathbf{A}_p = \mathbf{U}_p \mathbf{U}_p^T$  since  $\mathbf{A}_p$  is symmetric and positive semidefinite.  $\mathbf{U}_p$  has size  $D \times d$  such that  $d \ll D$ . In my project,  $d$  was set to be 10 by varying the value and seeing how it affects the accuracies, and picking the best values. Thus, the optimization problem can be rewritten as:

$$\max_{\mathbf{U}, \mathbf{w}} \sum_{p=1}^M w_p^r \frac{\text{tr}[\mathbf{U}_p^T (\frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p) \mathbf{U}_p]}{\text{tr}[\mathbf{U}_p^T \mathbf{S}_p \mathbf{U}_p]}$$

Such that  $\mathbf{U}_p^T \mathbf{U}_p = \mathbf{I}$ ,  $\sum_{p=1}^M w_p = 1$  and  $\forall p \in \{1, \dots, M\}, w_p \geq 0$ , where:

$$\mathbf{S}_p = \frac{1}{N} \sum_{(\mathbf{x}_i^p, \mathbf{y}_i^p) \in \mathcal{S}^p} (\mathbf{x}_i^p - \mathbf{y}_i^p)(\mathbf{x}_i^p - \mathbf{y}_i^p)^T \quad (3.1)$$

$$\mathbf{D}_p = \frac{1}{N} \sum_{(\mathbf{x}_i^p, \mathbf{y}_j^p) \in \mathcal{D}^p} (\mathbf{x}_i^p - \mathbf{y}_j^p)(\mathbf{x}_i^p - \mathbf{y}_j^p)^T \quad (3.2)$$

$$\mathbf{D}_{1p} = \frac{1}{NK} \sum_{\substack{(\mathbf{x}_i^p, \mathbf{y}_i^p) \in \mathcal{S}^p \\ \mathbf{y}_k^p \in \mathcal{N}_K(\mathbf{y}_i^p)}} (\mathbf{x}_i^p - \mathbf{y}_k^p)(\mathbf{x}_i^p - \mathbf{y}_k^p)^T \quad (3.3)$$

$$\mathbf{D}_{2p} = \frac{1}{NK} \sum_{\substack{(\mathbf{x}_i^p, \mathbf{y}_i^p) \in \mathcal{S}^p \\ \mathbf{x}_k^p \in \mathcal{N}_K(\mathbf{x}_i^p)}} (\mathbf{x}_k^p - \mathbf{y}_i^p)(\mathbf{x}_k^p - \mathbf{y}_i^p)^T \quad (3.4)$$

Where  $\mathcal{N}_K(\mathbf{x}_i^p)$  represents the  $K$  nearest neighbors of  $\mathbf{x}_i^p$ . This is when K-nearest neighbors is used, and  $K$  is set to 5 in the experiments. To solve this, we first let  $\mathbf{w}$  be constant in order to solve  $\mathbf{U}$  and then use that to find  $\mathbf{w}$ . By letting  $\mathbf{w}$  be constant, the problem is reduced to:

$$\max_{\mathbf{U}_p^T \mathbf{U}_p = \mathbf{I}} \frac{\text{tr}[\mathbf{U}_p^T (\frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p) \mathbf{U}_p]}{\text{tr}[\mathbf{U}_p^T \mathbf{S}_p \mathbf{U}_p]}$$

For each  $p \in \{1, \dots, M\}$ . This can then be converted to an alternative problem [8]:

$$\max_{\mathbf{U}_p} \text{tr} \left[ (\mathbf{U}_p^T \mathbf{S}_p \mathbf{U}_p)^{-1} \mathbf{U}_p^T \left( \frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p \right) \mathbf{U}_p \right]$$

Which can be solved by solving the following generalized eigenvalue problem:

$$\left( \frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p \right) \mathbf{u} = \lambda \mathbf{S}_p \mathbf{u}$$

Thus,  $\mathbf{U}_p = [u_1, u_2, \dots, u_d]$  such that the corresponding eigenvalues for  $u_1, \dots, u_d$  are the top  $d$  largest eigenvalues (ie.  $\lambda_1 \geq \dots \geq \lambda_d$ ). However, this can be problematic in practice since  $\mathbf{S}_p$  can be near singular when  $D > N$ . Thus, the identity matrix is added as a regularizer as follows:

$$\mathbf{S}_p = (1 - \beta) \mathbf{S}_p + \beta \frac{\text{tr}(\mathbf{S}_p)}{N} \mathbf{I} \quad (3.5)$$

Where  $\beta$  is a regularization parameter that is set to 0.5 experimentally. This allows us to then not need to solve the generalized eigenvalue problem, as that takes more computation time, and we can solve:

$$\mathbf{S}_p^{-1} \left( \frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p \right) \mathbf{u} = \lambda \mathbf{u} \quad (3.6)$$

Thus, we have found  $\mathbf{U}_p, \forall p \in \{1, \dots, M\}$ .

To find the weights,  $\mathbf{w}$ , we can construct the Lagrangian:

$$\mathcal{L}(\mathbf{w}, \lambda) = \sum_{p=1}^M w_p^r \frac{\text{tr}[\mathbf{U}_p^T (\frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p) \mathbf{U}_p]}{\text{tr}[\mathbf{U}_p^T \mathbf{S}_p \mathbf{U}_p]} - \lambda \left( \sum_{p=1}^M w_p - 1 \right)$$

And thus, solve it:

$$\frac{\partial \mathcal{L}}{\partial w_p} = r w_p^{r-1} \frac{\text{tr}[\mathbf{U}_p^T (\frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p) \mathbf{U}_p]}{\text{tr}[\mathbf{U}_p^T \mathbf{S}_p \mathbf{U}_p]} - \lambda = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{p=1}^M w_p - 1 = 0$$

Which means:

$$w'_p = \left( \frac{\text{tr}[\mathbf{U}_p^T \mathbf{S}_p \mathbf{U}_p]}{\text{tr}[\mathbf{U}_p^T (\frac{1}{2}(\mathbf{D}_{1p} + \mathbf{D}_{2p}) + \mathbf{D}_p) \mathbf{U}_p]} \right)^{\frac{1}{r-1}} \quad (3.7)$$

$$w_p = \frac{w'_p}{\sum_{i=1}^M w'_i} \quad (3.8)$$

Thus, we now have the matrices to calculate the distance and the weights for each face descriptor.

## 3.5 Prediction

Given the matrices  $U_p$  and weights  $w_p$  from WGEML for all face descriptors where  $U_p$  corresponds to the  $p$ 'th face descriptor, we can use this to predict whether a given pair of images are of the kin relationship specified. First, we transform each of the faces into their  $M$  face descriptors and use PCA to reduce the dimensions to the proper dimensionality. We then have  $\mathbf{x} = \{x_p \mid 1 \leq p \leq M\}$  and  $\mathbf{y} = \{y_p \mid 1 \leq p \leq M\}$ . Then, the similarity of these images can be calculated as follows, letting  $A_p = U_p U_p^T$ :

$$\text{sim}(\mathbf{x}, \mathbf{y}) = \sum_{p=1}^M \frac{w_p}{2} \left( \frac{x_p^T A_p y_p}{\sqrt{x_p^T A_p x_p} \sqrt{y_p^T A_p y_p}} + 1 \right)$$

This ends up finding a weighted cosine similarity for a non-Euclidean space, since our metric here is  $A_p$ , for each face descriptor that is used. The function  $\text{sim}$  ranges from 0 to 1 as needed for a similarity function. This differs from the version mentioned in [12] which was:

$$\text{sim}(\mathbf{x}, \mathbf{y}) = \sum_{p=1}^M \frac{w_p}{2} \left( \frac{x_p^T U_p^T U_p y_p}{\sqrt{x_p^T U_p^T U_p x_p} \sqrt{y_p^T U_p^T U_p y_p}} + 1 \right)$$

However, as mentioned in section 3.4, the model was trained such that  $U_p^T U_p = I$  which would make this boil down to a weighted cosine similarity which wouldn't make much use of the metric learning. As such, this was changed up to use  $A_p$  in the implementation rather than  $U_p^T U_p$ .

In order to then determine whether two images are of the kin relationship specified, the similarity must then be compared to a threshold,  $\theta$ . If  $\text{sim}(\mathbf{x}, \mathbf{y}) \geq \theta$ , then the pair is labeled as having that relationship and, otherwise, it isn't. As this value  $\theta$  wasn't mentioned in [12], it had to be found experimentally. By ranging over samples from 0 to 1, it was found that  $\theta = 0.6$  yielded the best results for the accuracies and seemed the most similar to what the original paper had.

### 3.5.1 Tri-kin Relationship Prediction

Since tri-kin relationships have three images as the input, two parents and one child, the similarity is defined as the mean similarity between each of the parents and the child. In other words:

$$\text{sim}(\mathbf{p}_1, \mathbf{p}_2, \mathbf{c}) = \frac{\text{sim}(\mathbf{p}_1, \mathbf{c}) + \text{sim}(\mathbf{p}_2, \mathbf{c})}{2}$$

Where  $\mathbf{p}_1, \mathbf{p}_2$  represent the face descriptors for each parent and  $\mathbf{c}$  represents the face descriptors of the child.

## 3.6 Overview of Workflow

There are three main stages to running the project, which are, for each dataset, the preprocessing stage, the training stage, and the testing stage. The important outputs of each of these stages are saved onto disk so that each stage doesn't have to be run right after another.

### 3.6.1 Preprocessing

Given a dataset, we wish to precompute the face descriptors for each image and set up the 5-fold cross-validation for TSKinFace. The pipeline is shown in figure 3.3.

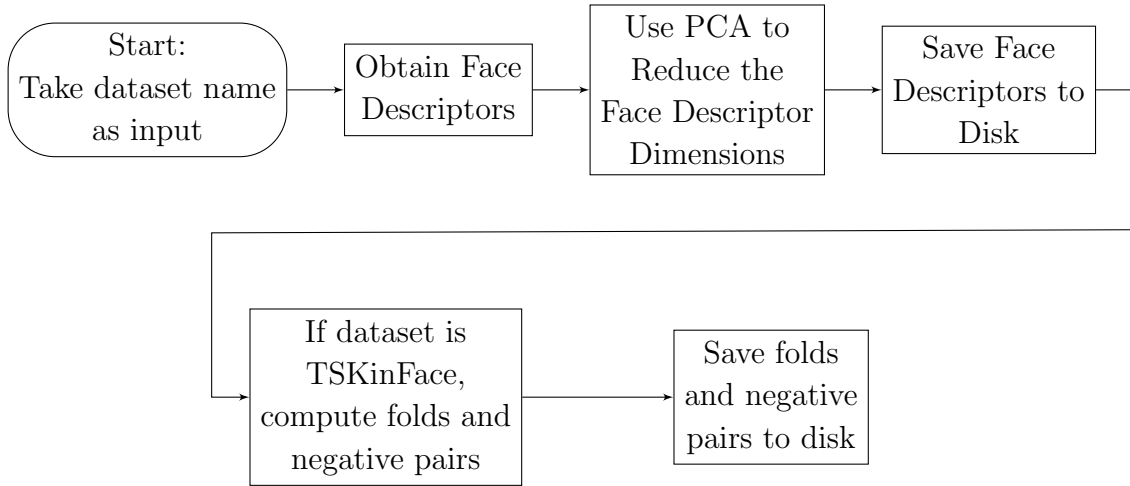


Figure 3.3: The preprocessing pipeline

### 3.6.2 Training

Once all of the face descriptors for each dataset are computed and, for each relationship in the dataset and each setting (restricted or unrestricted), the cross-validation folds are created, the training set is properly made from the data and WGEML, which is described in section 3.4, is run. The output of WGEML is then saved onto disk, which are the metrics and weights for each face descriptor for each fold. In the end, this is done for each relationship, setting, and fold, so for KinFaceW-I, for example, there would be  $4 \times 2 \times 5 = 40$  models saved. The pipeline is shown in figure 3.4.

### 3.6.3 Testing

Given a dataset, setting, and relationship, the folds and the corresponding models are loaded from on disk, and, for each test set, the corresponding model is used to predict whether the images in the test set are of the given relationship or not. The accuracy is then recorded, and

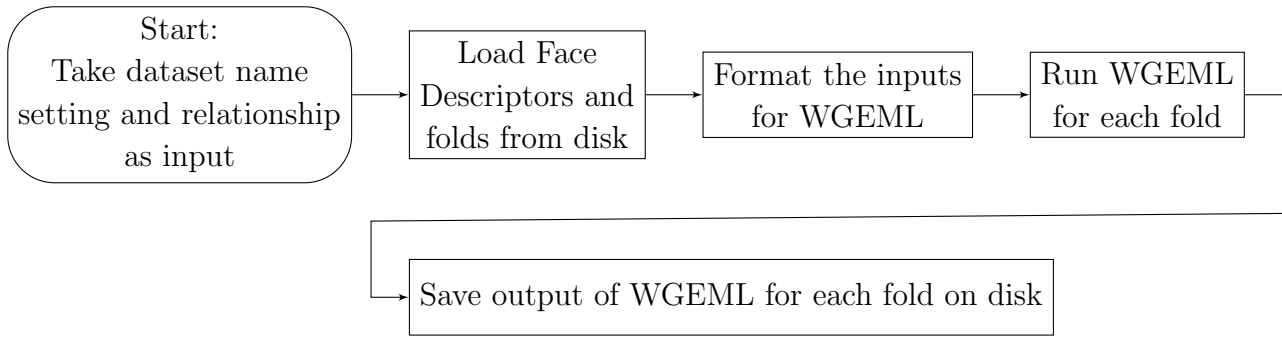


Figure 3.4: The training pipeline

the average accuracy is outputted for the given configuration along with each of the individual accuracies, as shown in Figure 3.5. The pipeline is shown in figure 3.6.

```

(venv) mv465@idun:~/KinRecognition$ CUDA_VISIBLE_DEVICES=0 make runPredictionKFW1Unrestricted
python3 -m src.scripts.testing "KinFaceW-I" "fs" "unrestricted"
KinFaceW-I-fs-unrestricted: [0.7581 0.8387 0.8387 0.8065 0.7656]
KinFaceW-I-fs-unrestricted: 0.8015
python3 -m src.scripts.testing "KinFaceW-I" "fd" "unrestricted"
KinFaceW-I-fd-unrestricted: [0.7037 0.7407 0.7593 0.6296 0.75 ]
KinFaceW-I-fd-unrestricted: 0.7167
python3 -m src.scripts.testing "KinFaceW-I" "ms" "unrestricted"
KinFaceW-I-ms-unrestricted: [0.8913 0.8043 0.6522 0.6957 0.7708]
KinFaceW-I-ms-unrestricted: 0.7629
python3 -m src.scripts.testing "KinFaceW-I" "md" "unrestricted"
KinFaceW-I-md-unrestricted: [0.84 0.72 0.82 0.84 0.7407]
KinFaceW-I-md-unrestricted: 0.7921
  
```

Figure 3.5: A sample output of the testing stage on KinFaceW-I unrestricted

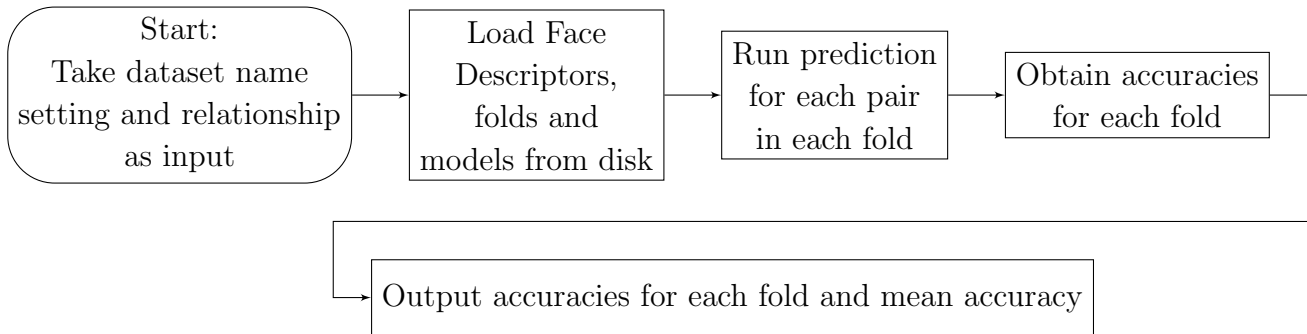


Figure 3.6: The testing pipeline

# Chapter 4

## Evaluation

In this chapter, we see that the accuracies obtained were close to that of the original paper in section 4.1 which achieves the success criterion. We further evaluate the model using ROC curves looking at the models for each dataset and each relationship. The biases in KinFaceW-II and TSKinFace are found and discussed in section 4.4. I also performed ablation studies on the face descriptors used in the model and we end up seeing a diminishing returns in accuracy in section 4.5. Finally, the results of using CifarNet instead of VGG are found to be, generally, slightly worse but not by much, which we see in section 4.6 and, finally, the unit tests and their coverage will be briefly discussed in section 4.7.

### 4.1 Overall Accuracies of the Model

We abbreviate “father” as “F”, “mother” as “M”, “son” as “S”, and “daughter” as “D” so the relationships are noted as a combination of these acronyms, so “FS” stands for the Father-Son relationship, etc., as discussed in section 2.7.3. The project is run end-to-end for each dataset, setting, and relationship, and the accuracies are shown in figure 4.1. For the success criterion, we look at the difference between the accuracies my implementation had obtained, and the original paper had obtained, which is shown in figure 4.2 where a positive number means my accuracies were better.



Figure 4.1: Accuracies of WGEML applied to each dataset for each relationship grouped by dataset

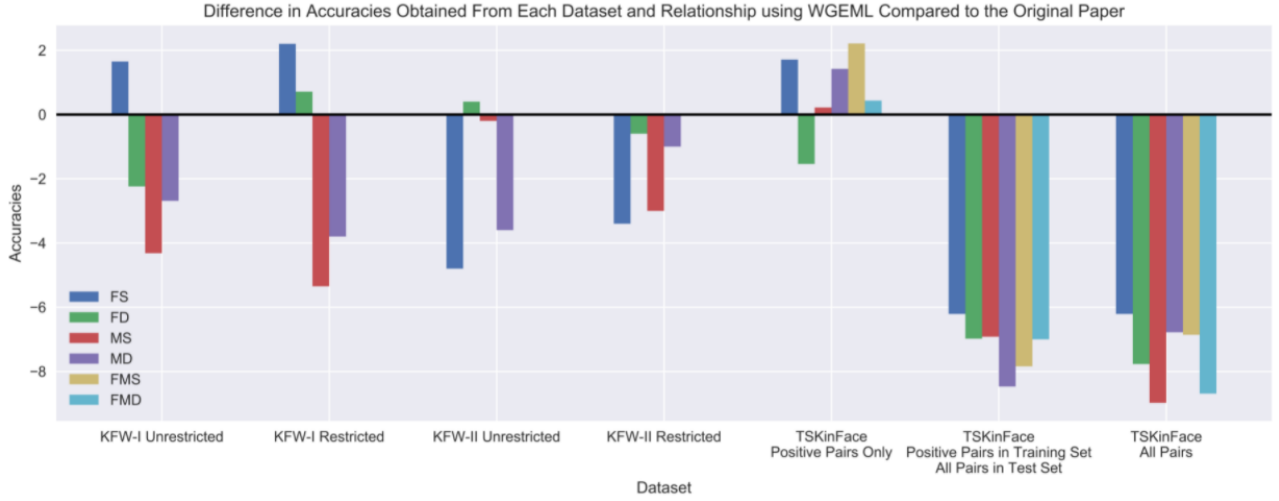


Figure 4.2: Differences in accuracies between my implementation and [12] in %

We can see that none of the absolute values of the differences are above 8.89% which, as will be discussed in 4.2, means that the success criterion is fulfilled.

In general, we see that the accuracies of the Father-Son relationship for each dataset are higher than those of the other relationships. However, there is one exception to this which is in the TSKinFace dataset when we use only positive pairs in the training set and both positive and negative pairs in the testing set.

When only positive pairs were used for the TSKinFace dataset for both training and testing, we get accuracies that are fairly close to that of the original paper with an average difference of 0.742%. However, when negative pairs were added to the test set, the accuracies drop by, on average, 8.29%. This led me to believe that the original paper had only used the positive pairs for the testing phase of TSKinFace as TSKinFace didn't come with any premade negative pairs. I attempted to confirm this but as of writing, I have yet to receive a response. As a note, in the later sections of this chapter, TSKinFace refers to having all pairs in both the training and test set.

## 4.2 Success Criterion

As seen in section 4.1, I was able to replicate the accuracies from the paper by Liang et al [12] within a  $\pm 15\%$  accuracy by using the WGEML methodology which means that the project is a success according to the goals defined in my proposal. Furthermore, the extensions were also completed and the results of which are discussed in sections 4.4, 4.5, and 4.6.

## 4.3 Receiver Operating Characteristic (ROC) Curves

We can also examine the ROC curves for each dataset in figure 4.3. These are curves that graph the false positive rate against the true positive rate for a model. By looking at the area under the curve (AUC), we can obtain the probability that our model ranks a random positive pair of images higher than that of a random negative pair, where rank refers to the similarity of the two images as defined in section 3.5. This allows us to compare the models for each dataset

by looking at the area under the ROC curves for each dataset which can give us an idea of how good our model is.

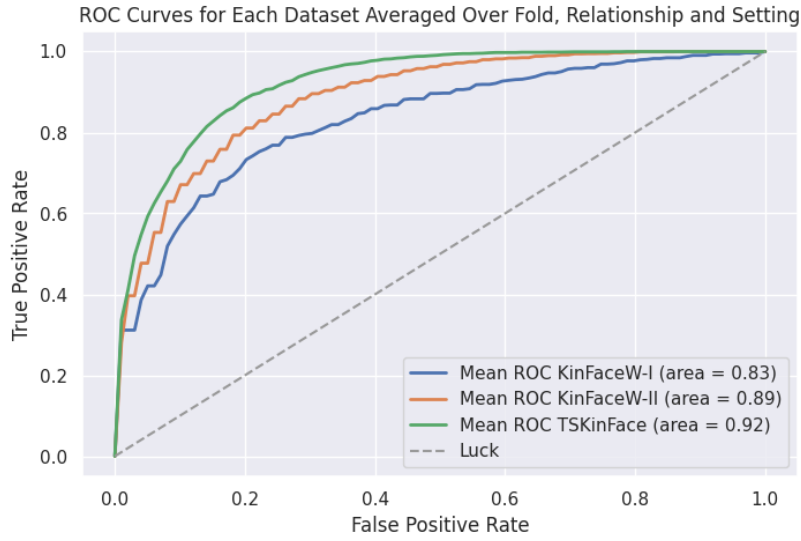


Figure 4.3: The mean ROC curve for each dataset averaged over each fold, relationship and setting

Here, we see that the AUC is the most for the TSKinFace dataset and the least for the KinFaceW-I dataset, though both curves are far from that of a classifier that is luck-based. We see that the model performs the best on TSKinFace, then KinFaceW-II, and then KinFaceW-I. However, as we will see in section 4.4, this could be due to the implicit biases in the TSKinFace and KinFaceW-II datasets.

The ROC curves can also be created for each relationship averaged over each fold, dataset, and setting, which is shown in figure 4.4.

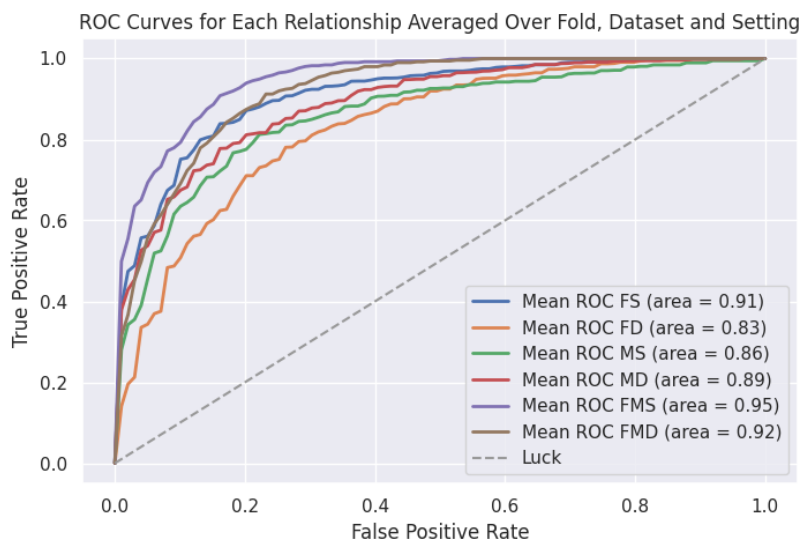


Figure 4.4: The mean ROC curve for each relationship averaged over each fold, dataset and setting

We can see that the FMS and FMD relationships do the best, but this is since these relationships are only represented in the TSKinFace dataset. Thus, the biases that are present in TSKinFace are also present in these relationships, which artificially boosts the area under their respective ROC curves. Other than that, we can also see that the models for the FS relationship tend to do the best on average whereas the FD relationship does the worst, which we can also see in figure 4.1. The two relationships with the highest AUC are the ones where the parent and child are of the same sex, FS and MD.

## 4.4 Potential Biases in Datasets

In order to look at potential biases that could occur in each of the datasets, each model was tested on all of the available datasets. So, for example, the model that was returned from training on the KinFaceW-I unrestricted dataset for the Father-Son relationship will have been tested on all of the three datasets that were used.

In every case, the accuracies went down when the test set that was used came from a different dataset than the training set. This is to be expected as the model might be overfitting to small biases that are inherent in every dataset. The models trained with the KinFaceW-I dataset generalizes fairly well, compared to KinFaceW-II and TSKinFace, since it seems to have the lowest differences of the three datasets with an average of an 8% decrease in accuracy compared to the accuracies when it is tested with itself. Then, KinFaceW-II was worse at a difference of 11.7% and then TSKinFace was by far the worst with an average difference of 15.7%. The average accuracy difference across the two test datasets for each training dataset and relationship is shown in figure 4.5 and the average accuracy differences for each dataset across all relationships and test datasets is shown in figure 4.6.

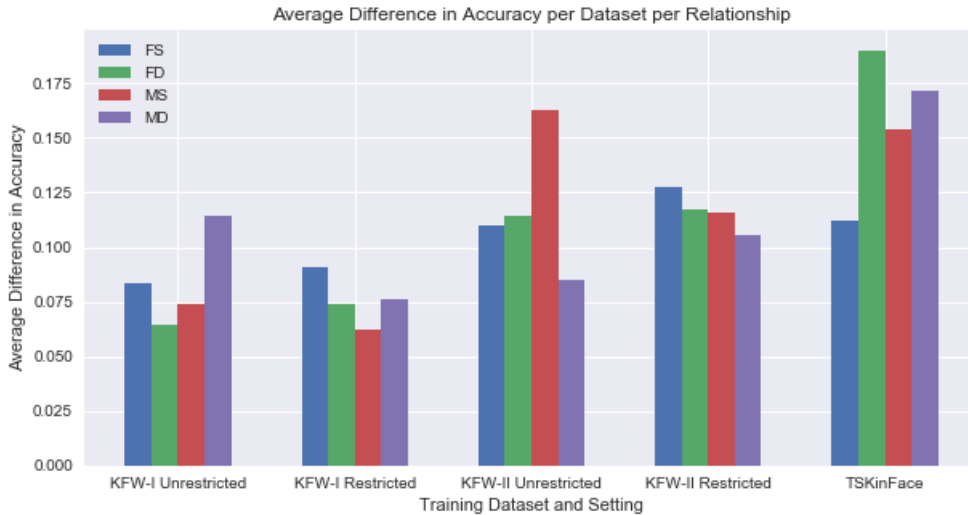


Figure 4.5: The average accuracy differences of the two test datasets for each training dataset and relationship

The main reason for this is since both KinFaceW-II and TSKinFace get the positive pairs from the same image whereas KinFaceW-I doesn't. At a high-level view, this means that WGEML might be taking into account the lighting of the image, the distortion of the face

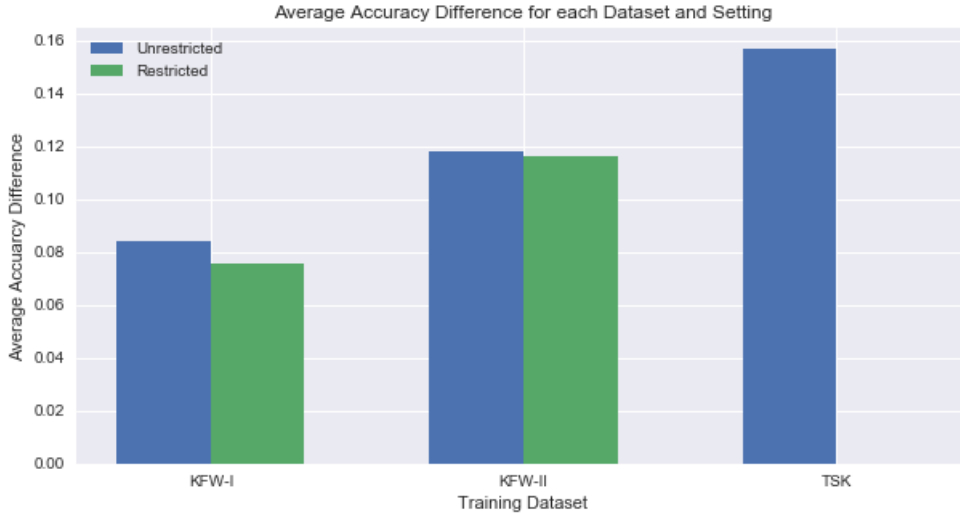


Figure 4.6: The average accuracy differences of the two test datasets and relationship for each training dataset

due to the camera settings, or other non-face related features of the image and WGEML is using the similarity of those features more in the KinFaceW-II and TSKinFace datasets than actual facial features. This means that WGEML might be picking up on these cues to increase performance [4]. For example, if a pair of images that have a similar lighting comes up, at a high-level, WGEML might be using the fact that it has similar lighting to infer that the faces come from the same photo and thus are more likely to be related. Since pairs in KinFaceW-I aren't from the same image, necessarily, WGEML can generalize better which explains why the differences are smaller than that of KinFaceW-II and TSKinFace. While it might be possible to normalize the lighting and do other techniques to try and get rid of the biases, we see that there are several cues that could be used to tell if two faces come from the same image [4]. It isn't feasible to normalize for all of these cues and it might not be possible for some of them.

## 4.5 Ablation Studies

### 4.5.1 Blocking Face Descriptors

Ablation studies were done by using each possible subset of the 4 face descriptors that were originally used for WGEML for training and testing. For example, only the VGG face descriptor was used for training and testing and the accuracies were obtained for each dataset, setting, and relationship.

After obtaining each accuracy for each dataset, setting, relationship, and subset of face descriptors, I decided to group up the accuracies by the number of face descriptors used by averaging the accuracies for each subset of that number.

As expected, in general, as the number of face descriptors are increased, the accuracy increases as well. However, as you increase the number of face descriptors, there are diminishing returns on the accuracy (ie. the second derivative is negative). An example of this is shown with KinFaceW-II unrestricted which is shown in figure 4.7.

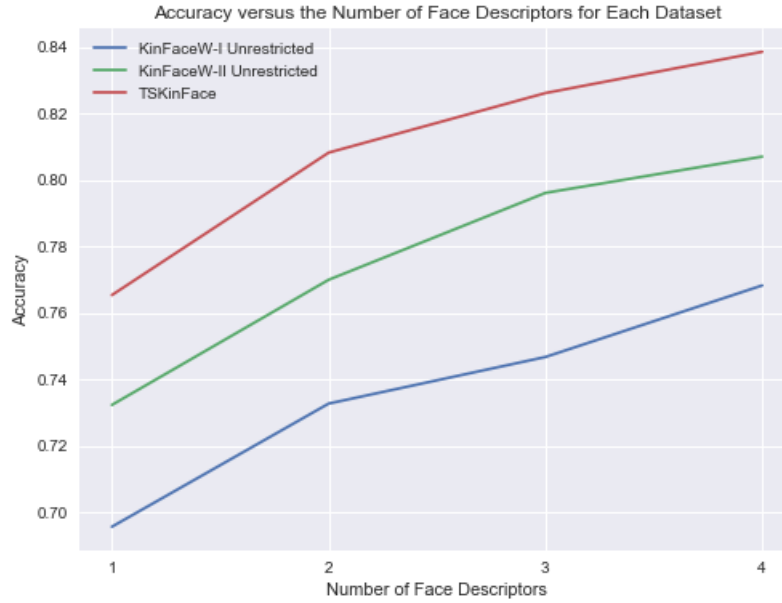


Figure 4.7: The plot of accuracy versus number of face descriptors used for each unrestricted dataset

The only dataset this doesn't seem to apply to is KinFaceW-I unrestricted in which adding a second face descriptor increases the accuracy quite a bit but the third doesn't increase it too much and then the fourth increases it more, as seen in figure 4.7.

However, this is mainly due to two outliers, for the FD relationship, increasing from 2 to 3 face descriptors only increased the accuracy by 0.2%, and for the MS relationship, going from 3 to 4 increased accuracy by 3%.

At a high level, as more face descriptors are added, the features that each face descriptor capture are more likely to overlap with each other, compared to when there's only one face descriptor, in which case another one could add a lot more information. However, if there are already 3 face descriptors, another face descriptor might share a lot of the same information the others already have captured and only be able to add, relatively, a small amount of information.

I then looked at which individual configurations of face descriptors used, and relationship resulted in a better accuracy than if all of the face descriptors were used. One thing I saw was that, out of the 22 instances this happened among all of the datasets, all of them used the VGG face descriptor. This led me to believe that the VGG descriptor is quite useful when used in conjunction with at least one other face descriptor since we also had that the configurations that used VGG tended to do better in accuracy than the ones that didn't. This makes sense since VGG is a deep feature detector whereas the other face descriptors use the texture only.

Furthermore, the KinFaceW-I restricted configuration had by far the highest number of configurations which did better than using all 4 face descriptors for that relationship, with 11 of the 22 configurations. The other datasets had about 3 or 4 instances where this occurred.

## 4.6 CifarNet Extension

When I replaced VGG with CifarNet (CFN), it does worse in the KinFaceW-I dataset, it does better in the KinFaceW-II dataset, and it doesn't have a huge effect on the TSKinFace dataset (with the largest difference being a 1% difference in accuracy). When grouped by relationship, the FS relationship had the least change in accuracy with an average 0.092% difference. This is shown in figure 4.8 where a positive number implies that the original face descriptors did better. There is more variation in the other relationships than there is in the father-son relationship, and this is reflected in the standard deviations for each relationship, ignoring the tri-kin relationships as there is only one dataset for this.

The differences grouped by dataset are shown in figure 4.9.

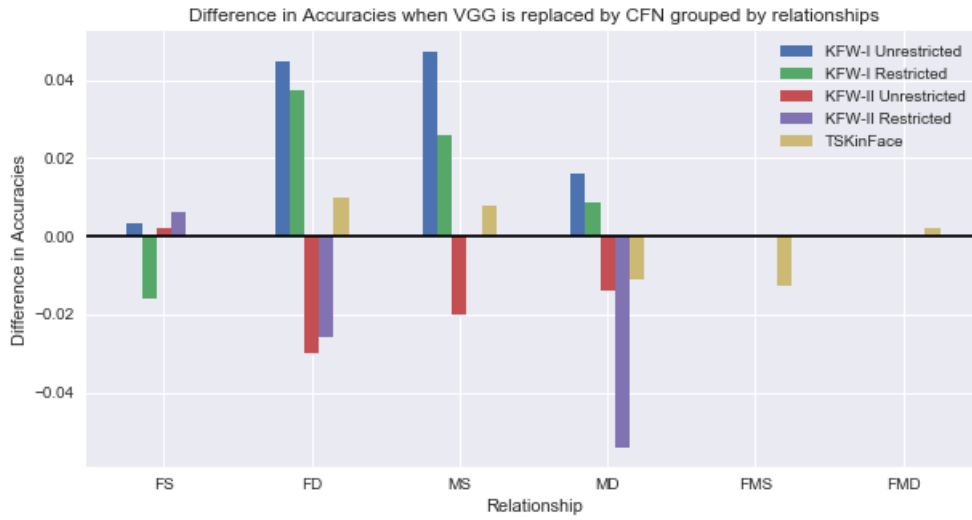


Figure 4.8: Differences in accuracy when VGG is replaced with CFN grouped by relationship

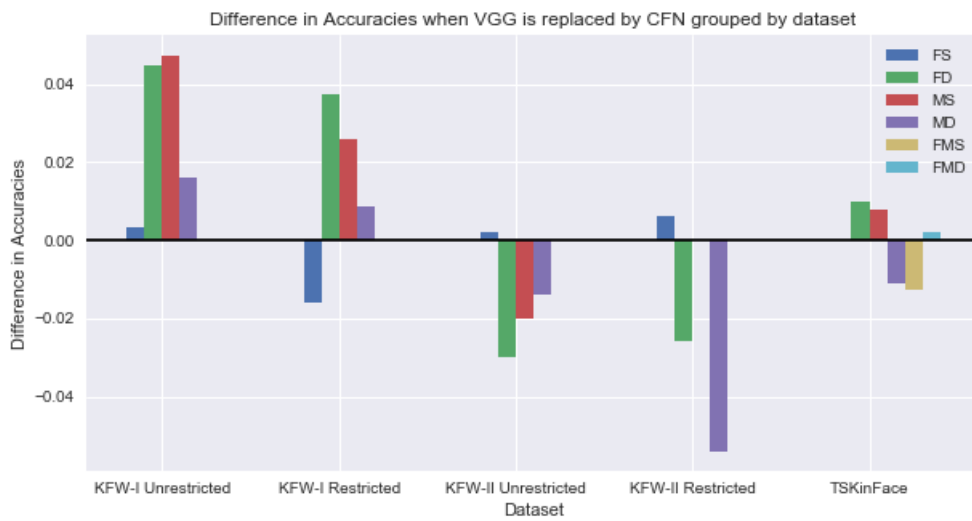


Figure 4.9: Differences in accuracy when VGG is replaced with CFN grouped by dataset

In general, VGG did better than CFN for KinFaceW-I and did worse for KinFaceW-II. The maximum difference in accuracies was 5.4% which occurred in the KinFaceW-II unrestricted

setting with the MD relationship in which CFN made the accuracy increase by 5.4%. The maximum decrease in accuracy was 4.7% with the KinFaceW-I Restricted dataset for the MS relationship. The KinFaceW-II Unrestricted setting with the MS relationship had no difference whatsoever in accuracies between the face descriptors.

Although VGG generally did better accuracy-wise, it might be useful to use the CFN descriptor instead due to the smaller model which means that there is less computation when it comes to prediction. The fact that it is a smaller model also means that there are fewer weights so the file the weights are saved to is smaller, which can be seen since the CFN weights are 5MB whereas the VGG weights are 566MB. This means that using the CFN model can be useful in mobile and embedded systems whereas having the extra computation and size cost might not be worth it for the 5% gain in the best case, and, on average, VGG has a 0.1% gain in accuracy, when the CFN weights take up around 1% of the disk space.

## 4.7 Unit Tests

To make sure the project was doing everything that it was supposed to be doing, unit tests were created for each file except the scripts. This is since the scripts were for integrating the modules together and thus didn't require unit testing. The statement coverage for each file and overall is shown in figure 4.10. In general, most of the files had 100% statement coverage and the main reason that CifarNet doesn't have a higher coverage is due to a function that trains the model and another function that packages up the creation of the model and training the model into one function. This is fairly difficult to unit test from my knowledge, which is what makes up most of the untested code.

Name	Stmts	Miss	Cover
src/WGEML/WGEML_training.py	57	0	100%
src/WGEML/constants.py	3	0	100%
src/data_preparation/PCA.py	5	0	100%
src/data_preparation/prep_cross_valid.py	41	0	100%
src/data_preparation/properly_formatted_inputs.py	17	0	100%
src/data_preparation/save_and_load.py	31	2	94%
src/face_descriptors/CifarNet.py	73	20	73%
src/face_descriptors/HOG.py	47	0	100%
src/face_descriptors/LBP.py	45	0	100%
src/face_descriptors/SIFT.py	151	10	93%
src/face_descriptors/VGG.py	79	6	92%
src/face_detection/face_detection.py	16	0	100%
src/prediction/predictor.py	29	0	100%
TOTAL	594	38	94%

Figure 4.10: Coverage of the unit tests of the project

# Chapter 5

## Conclusion

Overall, the project successfully reproduced the implementation and results presented in the WGEML paper within a 15% error range. Furthermore, I was able to further explore the intricacies of the algorithm and how it reacts to variations in the input. By doing ablation studies, I was able to discover that adding face descriptors to the algorithm has diminishing returns in section 4.5.1. Furthermore, sources of bias in the datasets were found in section 4.4.

### 5.1 Lessons Learned

One of the main lessons I learned was with respect to how important testing is. By writing these unit tests, even for seemingly simple functions, I saved myself a lot of potential future problems as I was able to catch many of the bugs while I was writing the code instead of at the integration stage. Of course, there were still bugs when I tried running the project end-to-end, but these were much simpler to solve knowing that each function I wrote and used is correct. For example, solving the eigenvalue problem and getting the top  $d$  vectors for WGEML is a fairly simple function to write. However, had I not created some non-trivial matrices to test that function with, I wouldn't have realized that the function was only returning  $d$  elements of each eigenvector instead of  $d$  eigenvectors. This is especially important in machine learning since these bugs could have gone unnoticed without proper testing.

As I hadn't had much experience with Computer Vision prior to this project, a lot of research had to be done into the intricacies of the original WGEML paper, such as what each face descriptor is and how they work. I was able to learn how to read multiple research papers and consolidate their information into my project.

Dependencies posed minor problems at points. Although I had a `requirements.txt` file and a virtual environment, running the testing pipeline on my laptop caused different results to the external GPU I was using for running the project end-to-end. The testing pipeline was the only pipeline I could reasonably run on my laptop due to its simplicity. However, since I had changed some of the dependencies in the `requirements.txt` file on the external GPU, I had obtained different results. I realized how important maintaining the same dependencies was to a project.

If I were to do the project over again, using more sophisticated software engineering tools such as Jenkins, which offers continuous integration, and Docker would help immensely. Sometimes I would forget to run my unit tests before pushing my changes to the feature branch and only afterwards I would realize that the unit tests fail and one time I only noticed after the

broken changes had already reached the master branch which continuous integration could help prevent. Docker could have also helped to prevent the dependency issue I had.

## 5.2 Future Work

It would be interesting to extend the model to work on videos since videos allow for more information, such as the various angles of a person's face, to be taken into account [29]. Furthermore, the project can be extended such that the potential relationship doesn't have to be specified in the input, which would answer whether two people are related at all rather than if they are related in a specific way.

Aside from extensions regarding the implementation, we can also further explore whether, in the ablation studies, the graph of accuracy versus the number of face descriptors continues in the way I'd expect with it leveling off after enough face descriptors or does it go down after a bit since the abundance of information could potentially interfere with each other? Furthermore, how would the different sets of face descriptors that were used affect this? Furthermore, in the original results, there wasn't much difference in accuracies between the unrestricted and restricted settings in many of the cases. In this case, the effect the negative pairs have on the algorithm can be further explored and whether different sets of negative pairs can give statistically significant differences in the accuracies.

# Bibliography

- [1] T. Ahonen, A. Hadid, and M. Pietikainen. Face description with local binary patterns: Application to face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(12):2037–2041, 2006.
- [2] Laleh Armi and Shervan Fekri-Ershad. Texture image analysis and texture classification methods - a review, 2019.
- [3] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 886–893 vol. 1, 2005.
- [4] Mitchell Dawson, Andrew Zisserman, and Christoffer Nellåker. From same photo: Cheating on visual kinship challenges. *CoRR*, abs/1809.06200, 2018.
- [5] R. Fang, K. D. Tang, N. Snavely, and T. Chen. Towards computational models of kinship verification. In *2010 IEEE International Conference on Image Processing*, pages 1577–1580, 2010.
- [6] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [7] W.D. Hamilton. The genetical evolution of social behaviour. ii. *Journal of Theoretical Biology*, 7(1):17 – 52, 1964.
- [8] Y. Jia, F. Nie, and C. Zhang. Trace ratio problem revisited. *IEEE Transactions on Neural Networks*, 20(4):729–735, 2009.
- [9] I.T. Jolliffe. *Principal component analysis*. Springer-Verlag New York, 2002.
- [10] Andrej Karpathy. Cs231n convolutional neural networks for visual recognition. <https://cs231n.github.io/convolutional-networks/>. Accessed: 2021-04-30.
- [11] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [12] J. Liang, Q. Hu, C. Dang, and W. Zuo. Weighted graph embedding-based metric learning for kinship verification. *IEEE Transactions on Image Processing*, 28(3):1149–1162, 2019.
- [13] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91110, November 2004.

- [14] J. Lu, J. Hu, X. Zhou, Y. Shang, Y. Tan, and G. Wang. Neighborhood repulsed metric learning for kinship verification. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2594–2601, 2012.
- [15] J. Lu, X. Zhou, Y. Tan, Y. Shang, and J. Zhou. Neighborhood repulsed metric learning for kinship verification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(2):331–345, 2014.
- [16] Jill M. Mateo. Perspectives: Hamilton’s legacy: Mechanisms of kin recognition in humans. *Ethology*, 121(5):419–427, 2015.
- [17] A. Nandy and S. S. Mondal. Kinship verification using deep siamese convolutional neural network. In *2019 14th IEEE International Conference on Automatic Face Gesture Recognition (FG 2019)*, pages 1–5, 2019.
- [18] Keiron O’Shea and Ryan Nash. An introduction to convolutional neural networks. *ArXiv e-prints*, 11 2015.
- [19] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In Mark W. Jones Xianghua Xie and Gary K. L. Tam, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 41.1–41.12. BMVA Press, 2015.
- [20] Matti Pietikinen. Scholarpedia - local binary patterns example. <http://www.scholarpedia.org/article/File:LBP.jpg>. Accessed: 2021-04-30.
- [21] Xiaoqian Qin, Xiaoyang Tan, and Songcan Chen. Tri-subject kinship verification: Understanding the core of a family, 2015.
- [22] Joseph P Robinson, Ming Shao, and Yun Fu. Visual kinship recognition: A decade in the making, 2020.
- [23] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [24] Juan Luis Surez-Daz, Salvador Garca, and Francisco Herrera. A tutorial on distance metric learning: Mathematical foundations, algorithms, experimental analysis, prospects and challenges (with appendices on mathematical background and detailed algorithms explanation), 2020.
- [25] Michael E. Tipping and Christopher M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 61(3):611–622, 1999.
- [26] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 1, pages I–I, 2001.
- [27] Xiawi. Wikipedia - local binary patterns - neighborhood examples. [https://en.wikipedia.org/wiki/Local\\_binary\\_patterns#/media/File:Lbp\\_neighbors.svg](https://en.wikipedia.org/wiki/Local_binary_patterns#/media/File:Lbp_neighbors.svg). Accessed: 2021-04-30.

- [28] Hiromu Yakura, Shinnosuke Shinozaki, Reon Nishimura, Yoshihiro Oyama, and Jun Sakuma. Malware analysis of imaged binary samples by convolutional neural network with attention mechanism. pages 127–134, 03 2018.
- [29] Haibin Yan and Junlin Hu. Video-based kinship verification using distance metric learning. *Pattern Recogn.*, 75(C):1524, March 2018.
- [30] Haibin Yan, Jiwen Lu, Weihong Deng, and Xiuzhuang Zhou. Discriminative multimetric learning for kinship verification. *IEEE Transactions on Information Forensics and Security*, 9(7):1169–1178, 2014.

# Appendix A

## Algorithms Implemented in Libraries

### A.1 Haar-Based Cascade Classifier

For face detection using Haar-based cascade classifiers [26], we first need a training set which contains images with a face in it and images which don't. From here, we use a set of Haar features like those shown in figure A.1 to obtain the features from each image.

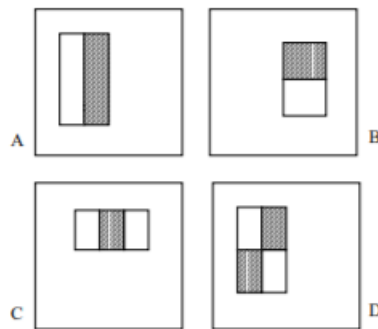


Figure A.1: Haar-based features to be applied on an image. Sourced from Viola et al. [26]

These Haar features work similar to convolution where, when we apply it to a region in the image, we add up the values in the black part and subtract the values in the white part of the kernel. The rectangles in A and B would find edge features, the rectangle in C would find line features and the last one finds four-rectangle features. We have to use each of these features for every part of the image which would require a large amount of computation. However, an optimization called integral images can be used to simplify this down [26].

Now, we have a set of features and a training set of positive and negative images. By using a variant of AdaBoost [6], we can select a small set of features for each image and train the classifier. This essentially trains a cascade of classifiers in stages such that only a subset of the features are computed in each stage and if one stage fails, we discard the feature.

### A.2 SIFT Keypoint Extraction

The first thing to do is to create octaves for the given image. An octave is a set of the given image being blurred multiple times. For example, in the first octave, the original image is the first image, and then the next image is blurred slightly which is then blurred further for the

next image in the octave, etc. In the second octave, the image is halved in size and the same blurring effect happens. So, if the original image was of size  $64 \times 64$ , then the images in the second octave will be  $32 \times 32$  and in the third it would be  $16 \times 16$  and so on. A specified number of octaves and blurred images are used for the SIFT algorithm. The way that the image is blurred is as follows:

$$L(x, y, \sigma) = G(x, y, \sigma) * I(x, y)$$

Where  $x, y$  is the coordinate in the image,  $I$  is the function mapping coordinates to the value of the image at that coordinate,  $\sigma$  is the amount of blurring,  $*$  represents convolving  $G$  on the image and:

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

Which is the Gaussian blur.

From here, for each octave, a difference of Gaussians is created to help find the keypoints of the image. The difference of Gaussians is just the difference between the consecutive Gaussians. This can be visualized in figure A.2.

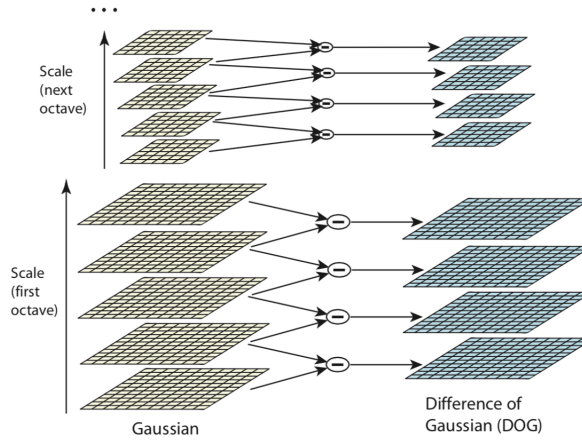


Figure A.2: The Difference of Gaussians being created. Image reproduced from Lowe (2004) [13]

This is used to approximate the Laplacian of Gaussians as the LoG helps find the edges of the image by blurring the image a bit and then finding the second order derivatives. It is first blurred as taking the Laplacian straight away would be sensitive to noise. However, this is computationally expensive. The Difference of Gaussians is a good approximation of the scale invariant Laplacian,  $\sigma^2 \nabla^2 G$ .

From here we look for the keypoints in the image using the Difference of Gaussians. This is done by finding the local maxima and minima of the DoG. Once the extrema of the Difference of Gaussians are obtained, we need to refine the approximation of the keypoint because the actual keypoint is more likely to be between pixels. Thus, we can use a Taylor expansion around the proposed keypoint of the scale-space function  $D(\sigma, x, y)$  where  $\sigma$  is the blur level in the DoG. This Taylor expansion looks like:

$$D(\mathbf{x}) = D + \frac{\partial D}{\partial \mathbf{x}} \mathbf{x} + \frac{1}{2} \mathbf{x}^T \frac{\partial^2 D}{\partial \mathbf{x}^2} \mathbf{x}$$

Where  $D$  is the value of  $D$  at the proposed keypoint,  $\mathbf{x} = (\sigma, x, y)^T$ ,  $\frac{\partial D^T}{\partial \mathbf{x}} = (\frac{\partial D}{\partial \sigma}, \frac{\partial D}{\partial x}, \frac{\partial D}{\partial y})$ . Letting this equal 0, we get that the offset from our keypoint is:

$$\hat{\mathbf{x}} = -\frac{\partial^2 D^{-1}}{\partial \mathbf{x}^2} \frac{\partial D}{\partial \mathbf{x}}$$

If our offset is greater than 0.5 in any dimension, then we want to try again since that means it's closer to another sample point. We keep trying again until we get an offset which is close to the sample point. We then find the value at the subpixel extrema:

$$D(\hat{\mathbf{x}}) = D + \frac{1}{2} \frac{\partial D^T}{\partial \mathbf{x}} \hat{\mathbf{x}}$$

And if the value of the extrema is less than  $0.03 \times 255$ , then we throw it out since it is an unstable extrema and has low contrast.

Furthermore, we also eliminate any keypoints that are potentially on an edge. We can do this by looking at the Hessian of the keypoint. We have that the Hessian is:

$$\mathbf{H} = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{yx} & D_{yy} \end{bmatrix}$$

To determine whether something is a corner, we care about the eigenvalues of  $\mathbf{H}$ , or more specifically, the ratio between the eigenvalues. Using the trace and determinant of  $\mathbf{H}$ , we can find that:

$$\frac{\text{Tr}(\mathbf{H})^2}{\text{Det}(\mathbf{H})} = \frac{(r+1)^2}{r}$$

Where  $r$  is the ratio of the eigenvalues. Thus, letting the maximum ratio that the eigenvalues are allowed to be at be  $r_0 = 10$ , we just need to find if:

$$\frac{\text{Tr}(\mathbf{H})^2}{\text{Det}(\mathbf{H})} < \frac{(r_0+1)^2}{r_0}$$

If it is, then this is a proper keypoint. If it isn't then this is an edge so we can discard it.

We now have each of the keypoints which are scale-invariant, so we need to make it rotation-invariant as well. This is done by assigning an orientation to the keypoints. A neighborhood is taken around the keypoint in which the gradient is obtained for each pixel in the area and a histogram of these gradients is created based on the angle, this time in bins of 0 to 10, 10 to 20, etc. until 350 to 360. The peak in this histogram is calculated and taken and any other peaks which have a value above 80% of the original peak is also considered to calculate the orientation. Now, for each keypoint, we have the location, scale, and orientation.

### A.3 Principal Component Analysis

Principal Component Analysis (PCA) [9] is a statistical tool that we can use to reduce the dimensionality of a dataset while maintaining as much of the variability as possible. The algorithm identifies *principal components* which are the directions of maximum variance in the dataset.

Let  $p$  be the number of variables that are being measured in the dataset, which in our case is the dimension of the face descriptor. Let  $\mathbf{X}$  be the matrix  $n \times p$  data matrix that is composed

of the  $p$   $n$ -dimensional vectors  $\mathbf{x}_1, \dots, \mathbf{x}_p$  in which  $x_i$  represents the  $n$  observations of the  $i$ th variable. The goal of PCA is then to create a  $p \times k$  matrix  $\mathbf{W}$ , where  $k \leq p$ , which maps each row-vector of  $\mathbf{X}$  from the original  $p$ -dimensional feature space to a new  $k$ -dimensional feature subspace such that the variability of the dataset is kept intact, as much as possible.

First, the dataset must be standardized which means that the mean of each feature must be 0, which we denote as the matrix  $\mathbf{X}'$  and the column vectors as  $\mathbf{x}'_1, \dots, \mathbf{x}'_p$ . From there, the covariance matrix is obtained. This matrix is the following:

$$\Sigma = \begin{bmatrix} Cov(\mathbf{x}'_1, \mathbf{x}'_1) & \dots & Cov(\mathbf{x}'_1, \mathbf{x}'_p) \\ \vdots & \ddots & \vdots \\ Cov(\mathbf{x}'_p, \mathbf{x}'_1) & \dots & Cov(\mathbf{x}'_p, \mathbf{x}'_p) \end{bmatrix}$$

Where:

$$\Sigma_{ij} = Cov(\mathbf{x}'_i, \mathbf{x}'_j) = \frac{1}{n}(\mathbf{x}'_i \cdot \mathbf{x}'_j)$$

Since we have that the mean of both are 0 now, after standardization. Thus:

$$\Sigma = \frac{1}{N} \mathbf{X}' \mathbf{X}'^T$$

We then find the eigenvalues and eigenvectors of  $\Sigma$ , of which we get  $p$  of them. The eigenvalues get sorted from largest to smallest and, thus, the first principal component is the eigenvector that corresponds to the largest eigenvalue, the second principal component is the eigenvector with the second largest eigenvalue, etc. Denoting  $\mathbf{w}_i$  as the  $i$ th principal component, we can make the matrix  $\mathbf{W}$  as:

$$\mathbf{W} = [\mathbf{w}_1 \quad \dots \quad \mathbf{w}_k]$$

Which is a  $p \times k$  matrix. Thus, to reduce the dimensions of the original dataset, we can create the matrix:

$$\mathbf{T} = \mathbf{XW}$$

In which the  $i$ th row of  $\mathbf{T}$  corresponds with the  $i$ th observation of  $\mathbf{X}$  and this vector has is  $k$ -dimensional where  $k \leq p$ .

# Appendix B

## Raw Table Data

The tables for figures 4.1 and 4.2 is as follows:

Dataset	Setting	FS	FD	MS	MD	FMS	FMD
KFWI	Unrestricted	0.8015	0.7166	0.7628	0.7921		
KFWI	Restricted	0.811	0.7391	0.7405	0.769		
KFWII	Unrestricted	0.838	0.778	0.832	0.78		
KFWII	Restricted	0.848	0.768	0.802	0.814		
TSK	Only Positive Pairs	0.9201	0.8826	0.9162	0.9182	0.9571	0.9343
TSK	All Pairs In Test/Pos Pairs in Train	0.8409	0.8282	0.8448	0.8193	0.8566	0.86
TSK	All Pairs In Test/Train	0.8409	0.8203	0.8242	0.8362	0.8664	0.8431

Table B.1: Accuracies of WGEML applied to each dataset for each relationship

Dataset	Setting	FS	FD	MS	MD	FMS	FMD
KFWI	Unrestricted	1.65	-2.24	-4.32	-2.69		
KFWI	Restricted	2.2	0.71	-5.35	-3.8		
KFWII	Unrestricted	-4.8	0.4	-0.2	-3.6		
KFWII	Restricted	-3.4	-0.6	-3	-1		
TSK	Only Positive Pairs	1.71	-1.54	0.22	1.42	<b>2.21</b>	0.43
TSK	All Pairs In Test/Pos Pairs in Train	-6.21	-6.98	-6.92	-8.47	-7.84	-7
TSK	All Pairs In Test/Train	-6.21	-7.77	<b>-8.98</b>	-6.78	-6.86	-8.69

Table B.2: Differences in accuracies between my implementation and [12] in %

# Appendix C

## Project Proposal

The proposal is available on the next page.

# Kin Recognition using Weighted Graph Embeddings

Computer Science - Part II Project Proposal

Manu Varma

## 1 Introduction and Background

Kin recognition is the way of recognizing whether two people are related or not based on their faces and, if they are related, how are they related. It is advantageous to inclusive fitness for an organism to be able to recognize which of their neighbors were close relatives [4]. Thus, it stands to reason that the ability to recognize kin relationships has evolved in humans. In humans, facial resemblance is expected to serve as an indicator of kinship. Strangers are able to match photographs of mothers to their infants without any prior contact with the family [7].

In the field of creating computational models of kin recognition, there are multiple types of problems that are being solved. The first and most common kind is kinship verification which is what will be explored in this project. The task is to determine whether a pair of faces are blood relatives or not and what type of relatives they are. The next type of task is family classification which is the task of determining which family a single member belongs to. The next type of task is tri-subject verification where it is determined whether a child is related to a pair of parents. Finally, there is the search and retrieval task where a child is checked against all of the faces in a database to see if there is anyone they are related to. This has applications in helping find missing children specifically [8].

## 2 Starting Point

### 2.1 Personal Starting Point

Currently, I have minimal experience with Tensorflow and Keras. I have also worked with facial recognition a few years ago. I am very familiar with Python, overall, however. The main courses that could be helpful would be Machine Learning and Bayesian Inference and Deep Neural Networks in which I will be attending both in Lent Term and I will try to read ahead in before I start my implementation.

## 2.2 Current Existing Literature

One of the earliest approaches to kinship recognition [3] obtained the main facial features from each person using a pictorial structure model which then got the features and then a feature vector was created from that. Using the feature vectors, the differences between the corresponding feature vectors were calculated before applying K-nearest neighbors and Support Vector Machine methods to train the model.

Another method used in 2012 was to first partition the face into regions in 5 layers by slowly breaking down the face in each layer [9]. This analyzes the facial features and compares the features between two people. However, due to aging, these features can be deformed a bit so this is mitigated using transfer subspace learning. They do this transfer learning on child-old parent and child-young parent pairs in order to create a new problem. This way, the old and young parents can be operating on the same distribution of features rather than different distributions. This allows the feature difference that young and old parents might have to decrease and establishes a new standard for getting the features of the face from each person. After this, they use the gender relation, age difference and the distance between the two people in the photo along with the results of the transfer learning to get a kinship score and pick the relationship that has the highest score.

Lastly, a more recent method from 2019 uses Weighted Graph Embedding-Based Metric Learning [5] to obtain the classification. Instead of obtaining a single feature vector, four types of feature descriptors are obtained. These descriptors are feature vectors obtained from Local Binary Patterns, a Histogram of Gradients, a scale-invariant feature transform vector and a 4096-dimensional VGG-Face CNN descriptor using the VGG-Very-Deep-16 CNN architecture. Once these descriptors are obtained, for each type of relationship (father-son, mother-daughter, etc.), an intrinsic graph and two penalty graphs are created. K-nearest neighbors is used on these graphs to construct a set of matrices which are subsequently used to figure out whether the inputted people have the relationship the graph was testing. When the method was evaluated on KinFaceW-I and KinFaceW-II, the mean verification accuracy was about 78% and 83% respectively.

## 3 Substance and Structure of the Project

### 3.1 Aims of the Project

I aim to replicate the kin recognition results using the Weighted Graph Embedding Based method [5]. Once that is done, I aim to extend the model with the following extensions:

- Many current implementations of kin recognition use the fact that they faces came from the same image which makes kin recognition easier [1]. We can try to extend

the current implementation and create our own model for this problem to better account for this.

- The functionality can be extended to also include video rather than just images.
- Extend the functionality beyond kin recognition to relationship recognition to also identify if people are friends as well or if they are strangers
- Add on a system where you can search for all of the people in a database that you might be related to in general or to find pictures of a specific family member.

### 3.2 Dataset

I will be using a combination of the KinFaceW-I [6] and Family101 [2] datasets. The KinFaceW-I dataset contains pairs of faces collected from the internet which are labeled with one of four kin relations: Father-Son (F-S), Father-Daughter (F-D), Mother-Son (M-S), and Mother-Daughter (M-D). For each pair of faces, each face comes from a different photo than the other face in the pair. The photos are taken from uncontrolled environments where there was no restriction in the pose, lighting, background, expression, age, ethnicity, and partial occlusion which would mean that the dataset is not very biased.

The Family101 dataset contains 101 different families with 607 individuals and 14,816 images composed of renowned public families.

For the potential extension of including videos, a variant of the KinFaceW dataset can be used called Kinship Face Videos in the Wild [10]. This is similar to the KinFaceW-I dataset except the resolution of each frame is  $900 \times 500$  pixels and each video contains 100-500 frames each.

All of these datasets are public so long as they are cited as done above. In addition to being cited, the video dataset requires the creators to be emailed as well and the Family101 dataset requires them to be informed about the accuracies obtained.

### 3.3 Structure of the Project

To achieve this goal, a similar method to what is used in the weighted graph embedding approach [5]. The main components are as follows:

- I will require face recognition first in order to get the person's face from a given image into the format that will be used by the rest of the system.
- From the face that is a standardized size, the necessary feature descriptors need to be taken from the face image.

- The next component is the intrinsic graph and penalty graphs that are created for each kin relationship.
- Finally, the needed matrices that are obtained from said graphs are created to figure out whether the pair of faces has the relationship that is being tested.

### 3.4 Evaluation

One method of evaluation is to use the training/testing data split that the datasets already provide and use that as a benchmark for accuracy. We can use the KinFaceW-I and II datasets to benchmark the implementation with the paper and the Family101 dataset to see how it extends to other data. We can use currently existing models as a benchmark and compare our accuracy that was obtained with existing accuracies. We can also create a confusion matrix and calculate the  $F_1$ -score of the model based on the results of using the test data. Furthermore, the validation and learning curves can be analyzed in order to figure out if the model is underfitting or overfitting.

## 4 Success Criterion

The project can be considered a success if it can successfully replicate the accuracies that the Weighted Graph Embedding-Based methodology towards kin recognition [5] obtained in their evaluation with an error range of  $\pm 15\%$  due to factors such as the training/testing split being different from that used, using a different random number seed and using a different framework than was used in the paper or if I can invalidate the results of the paper.

## 5 Timeline and Milestones

### 1. Oct 23 - Nov 06:

This time will be mainly used for research and getting familiar with everything I need to know about the subject such as:

- Reading up more in depth on existing literature
- Getting a more in depth idea of what architecture I should be using
- Figuring out which libraries are necessary to be used
- Getting an understanding of the libraries to be used by doing small quick projects in either Tensorflow or PyTorch.

2. **Nov 06 - Nov 20:**

Basic face recognition will be implemented so that generic images can be given and not just images of a specific size of just their face.

- Using Dlib or OpenCV, create basic facial recognition software that outputs the image of the person's face in the required format.

**Deliverable:** A basic, usable facial recognition model that can be used for the rest of the project.

3. **Nov 20 - Dec 04:**

Create the ability to get the necessary feature vectors like local binary patterns and the histogram of gradients.

- Given the input of the pictures of the faces, get the local binary patterns and histogram of gradients out from them.

4. **Dec 04 - Dec 18:**

The VGG-Face CNN descriptors and SIFT face descriptors will be obtained in this time which would then be fed into the graphs. If this is done before the sprint is up, work will be started on implementing the graphs.

**Deliverable:** The face descriptor methods are all created.

5. **Dec 18 - Jan 01:**

The implementation of the intrinsic graph would be created in this sprint. This should create the graph based on the class information.

6. **Jan 01 - Jan 15:**

The penalty graphs should be created in this sprint and thus get the calculations necessary to figure out the kin relationship. This overall model that is created should be able to fulfill the success criterion.

**Deliverable:** A model that can predict kin relationships between pairs of images.

7. **Jan 15 - Jan 29:**

At this point, ablation studies will be done to help to evaluate the network by seeing which inputs are necessary for the network to work with more focus being given to the progress report and presentation:

- Start work on an ablation study.
- Write up the Progress Report and create a presentation for it to be handed in on February 5th.

**Deliverable:** Progress Report

8. **Jan 29 - Feb 12:**

Finish work on the ablation studies and start work on the first extension:

- Finish up ablation studies for evaluation and finish model
- Start working on extending the model to work for videos as well.

**Deliverable:** A finalized model that reports the kin relation between a people in a pair of images.

9. **Feb 12 - Feb 26:**

Finish work on the first extension and start doing work on a second extension to try and improve on the model for images.

**Deliverable:** A new model which is the model implemented which is extended to work on videos

10. **Feb 26 - Mar 12:**

At this point, the first draft of the dissertation starts to be written alongside some work on improving the model, with priority given to the draft:

- Write out a draft of the Introduction and Preparation Chapters
- Work on second extension for improving the existing model.

11. **Mar 12 - Mar 26:**

Continue working on writing the dissertation and extension with the same priority:

- Write out a draft of the Implementation chapter
- Work on second extension for improving the existing model. If it is finished, then this will be included in the implementation chapter. If not, it will be abandoned in favor of finishing up the draft of the dissertation.

**Deliverable:** A draft of the Introduction, Preparation and Implementation chapters to be given to supervisor for feedback

12. **Mar 26 - Apr 09:**

Continue working on writing the dissertation:

- Write out a draft of the Evaluation chapter
- Revise the Introduction, Preparation and Implementation chapter based on supervisor feedback

13. **Apr 09 - Apr 23:**

Continue working on writing the dissertation:

- Write out a draft of the Conclusion chapter

**Deliverable:** A draft dissertation to be given to supervisor for feedback

14. **Apr 23 - May 14:**

If everything has gone well, this should be where final touches are applied.

**Deliverable:** The completed dissertation

## 6 Resources Declaration

I plan to use my current laptop which is a Windows laptop with an Intel i7 at 2.8GHz, 16GB RAM, an NVIDIA GTX 1060 graphics card, 1TB of HDD space and 256GB of SSD space. I accept full responsibility for this machine and I have made contingency plans to protect myself against hardware and/or software failure. Such contingencies include backing up my files to Github, Google Drive and an external 1TB hard drive. The Github repository will be updated fairly often as I work on my project and the Google Drive and hard drive will be updated every 1-2 weeks. Should my machine encounter software failures, I would buy a new laptop and work on the MCS machines until the new laptop arrives. Furthermore, I will utilize GPUs provided by the Computer Laboratory when I need to train any models.

## References

- [1] Mitchell Dawson, Andrew Zisserman, and Christoffer Nellåker. From same photo: Cheating on visual kinship challenges. *CoRR*, abs/1809.06200, 2018.
- [2] R. Fang, A. C. Gallagher, T. Chen, and A. Loui. Kinship classification by modeling facial feature heredity. In *2013 IEEE International Conference on Image Processing*, pages 2983–2987, 2013.
- [3] R. Fang, K. D. Tang, N. Snavely, and T. Chen. Towards computational models of kinship verification. In *2010 IEEE International Conference on Image Processing*, pages 1577–1580, 2010.
- [4] W.D. Hamilton. The genetical evolution of social behaviour. ii. *Journal of Theoretical Biology*, 7(1):17 – 52, 1964.
- [5] J. Liang, Q. Hu, C. Dang, and W. Zuo. Weighted graph embedding-based metric learning for kinship verification. *IEEE Transactions on Image Processing*, 28(3):1149–1162, 2019.
- [6] J. Lu, X. Zhou, Y. Tan, Y. Shang, and J. Zhou. Neighborhood repulsed metric learning for kinship verification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(2):331–345, 2014.
- [7] Jill M. Mateo. Perspectives: Hamilton’s legacy: Mechanisms of kin recognition in humans. *Ethology*, 121(5):419–427, 2015.
- [8] Joseph P Robinson, Ming Shao, and Yun Fu. Visual kinship recognition: A decade in the making, 2020.

- [9] S. Xia, M. Shao, J. Luo, and Y. Fu. Understanding kin relationships in a photo. *IEEE Transactions on Multimedia*, 14(4):1046–1056, 2012.
- [10] Haibin Yan and Junlin Hu. Video-based kinship verification using distance metric learning. *Pattern Recogn.*, 75(C):1524, March 2018.