

Manu Varma

Kin Recognition Using Weighted Graph Embeddings

Computer Science Tripos – Part II

St John's College

April 23, 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 April 23, 2021

Proforma

Name: **Manu Varma**
College: **St John's College**
Project Title: **Kin Recognition Using Weighted Graph Embeddings**
Examination: **Computer Science Tripos – Part II, June 2021**
Word Count: **TBA**
Line Count: **2907**
Project Originator: **The Dissertation Author**
Supervisor: **Daniel Bates**

Original Aims of the Project

Work Completed

Special Difficulties

Contents

1	Introduction	7
1.1	Problem Overview	7
1.2	Motivation	8
1.3	Related Work	8
1.4	Project Overview	8
2	Preparation	10
2.1	Convolutional Neural Networks	10
2.1.1	Convolutional Layers	11
2.1.2	Pooling Layers	11
2.1.3	Activation Functions	12
2.1.3.1	ReLU	12
2.1.3.2	Softmax	12
2.2	Face Detection	13
2.3	Face Descriptors	13
2.3.1	Local Binary Patterns	13
2.3.2	Histogram of Gradients	15
2.3.2.1	Calculating the Gradients	15
2.3.2.2	Weighted Vote into Histogram	16
2.3.3	Scale-Invariant Feature Transform	17
2.3.4	VGG	19
2.4	Metric Learning	20
2.5	K-Nearest Neighbors	21
2.6	Principal Component Analysis	21
2.7	Requirements Analysis	23
2.8	Software Engineering Practices	24
2.8.1	Starting Point	24
2.8.2	Tools Used	24
2.8.3	Datasets	24
2.8.4	Testing	25
2.8.5	Licensing	25
3	Implementation	26
3.1	Repository Overview	26
3.2	Data Preparation	26
3.2.1	Cross-Validation	26
3.2.2	Positive and Negative Pairs	27

3.2.3	PCA	28
3.2.4	Saving Results to Disk	28
3.3	Face Descriptors	29
3.3.1	SIFT Implementation	29
3.3.2	VGG Implementation	29
3.3.3	CifarNet Extension	30
3.4	WGEML	30
3.4.1	Problem	31
3.4.2	Approach	32
3.5	Prediction	33
3.5.1	Tri-kin Relationship Prediction	34
3.6	Overview of Workflow	34
3.6.1	Preprocessing	34
3.6.2	Training	34
3.6.3	Testing	35
4	Evaluation	37
4.1	Original Results	37
4.2	Success Criterion	38
4.3	Potential Biases in Datasets	38
4.4	Ablation Studies	39
4.4.1	Blocking Face Descriptors	39
4.5	CifarNet Extension	41
4.6	Unit Tests	42
5	Conclusion	44
5.1	Lessons Learnt	44
5.2	Future Work	45
	Bibliography	45
A	Cascade Classifier	48
B	Project Proposal	49

List of Figures

1.1	An example of an input into a classical kin recognition problem	7
2.1	An example of an artificial neural network with two hidden layers	10
2.2	A convolution done on an input. Image sourced from Yakura et al. (2018) [22] .	11
2.3	An example of the max pooling operation done on a single slice of the input . .	12
2.4	Potential neighborhoods of the pixel	14
2.5	Example of LBP operator on a pixel	14
2.6	The Difference of Gaussians being created. Image reproduced from Lowe (2004) [11]	17
2.7	The different configurations that VGG can have. Image reproduced from Simonyan et al. (2015) [20]	19
2.8	A high-level view of metric learning	21
2.9	Principal components of an arbitrary 2 dimensional dataset	22
3.1	Folder Structure of the project	27
3.2	The Architecture of the CifarNet model	30
3.3	The preprocessing pipeline	34
3.4	The training pipeline	35
3.5	A sample output of the testing stage on KinFaceW-I unrestricted	35
3.6	The testing pipeline	36
4.1	The average accuracy differences of the two test datasets for each training dataset and relationship	39
4.2	The average accuracy differences of the two test datasets and relationship for each training dataset	39
4.3	The plot of accuracy versus number of face descriptors used for KinFaceW-II Unrestricted	40
4.4	The plot of accuracy versus number of face descriptors used for KinFaceW-I Unrestricted	41
4.5	Differences in accuracy when VGG is replaced with CFN grouped by relationship	42
4.6	Differences in accuracy when VGG is replaced with CFN grouped by dataset . .	42
4.7	Coverage of the unit tests of the project	43

Chapter 1

Introduction

Kinship recognition is the way of recognizing whether two given people are related to each other or not based on how they look. This is an evolutionary trait in humans as it is advantageous to inclusive fitness for an organism to be able to recognize which of their neighbors are close relatives [6]. 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 as we have seen that strangers are able to match photographs of mothers to their infants without any prior contact with any of the family [14].

Computational kinship recognition is the field of computationally figuring out kinship relationships between people without any prior knowledge of the family.

1.1 Problem Overview

There are multiple problems in the field of computational kin recognition, and I will focus on solving one of them. The first of which 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.

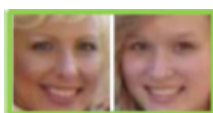


Figure 1.1: An example of an input into a classical kin recognition problem

These relationships tend to be more commonly parent-child relationships as opposed to sibling-sibling relationships, for example. A further extension of the main kin recognition problem is Tri-Subject Kinship which takes 3 images, two parents and a child, and determines if they are related or not. These problems are the ones that are being looked at in this dissertation.

Some other major problems in the field include search-and-retrieval and family classification. Search-and-retrieval is the problem which takes an image of a person as input and searches through a database to find people with whom they are most likely to be related to. This outputs a list of these people that 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

Accurate kinship recognition software has multiple applications in humanitarian issues. With regards to humanitarian issues, by using kin recognition, we can better recognize missing children and match them with their parents [19]. It can also be used for stopping human traffickers from claiming they are family members of the victim or to reunite families across refugee camps.

Furthermore, there are potential use-cases in social media and it can also pose privacy protections.

1.3 Related Work

One of the first papers in the field of computational kinship recognition used color, facial parts, facial distances and a Histogram of Gradients vector as the features for each image. Using these features, K-Nearest-Neighbors and a Support Vector Machine (SVM) were used in order to classify the image pairs into true and false parent-child pairs [5]. A classification accuracy of 70.67% was obtained overall and an SVM with a radial basis kernel obtained an accuracy of 68.6%.

Other approaches include using deep learning to solve the problem. One paper used a Siamese Convolutional Neural Network (CNN) approach by putting both of the face images in the pair through a SqueezeNet network which was trained on VGGFace2 which creates a feature vector for each image [15]. Using a similarity criterion, a new feature vector is created using the two feature vectors and a fully connected layer and a sigmoid activation function 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. Another paper that used a Siamese approach aimed to solve both the standard kinship verification problem and the Tri-Subject problem [23]. Both networks for the problems had three stages, a feature extraction stage, which used the ResNet50 or SENet50 models pretrained on VGGFace2 to extract the features from each face into a vector, a feature fusion stage which combines the feature vectors together, and a similarity quantization stage to get the similarity score. They then use a jury system to fuse together models which allows them to achieve an accuracy of 75.9%.

However, a prominent approach that is taken includes using metric learning to solve the problem. For example, the paper on Neighborhood Repulsed Metric Learning [13] uses a supervised metric learning approach. Using the approach, they aim to find a metric which 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. Furthermore, the paper that will be implemented in this dissertation, using Weighted Graph Embedding-Based Metric Learning, involves a metric learning approach to the problem [10].

1.4 Project Overview

The project deals with a 2019 method using Weighted Graph Embedding-Based Metric Learning [10], or WGEML for short, and aims to implement the paper and verify the accuracies that were obtained. The following is shown in the dissertation:

1. The algorithm, WGEML, is implemented and accuracies that are within an acceptable range are obtained, as shown in section 4.1.

2. Ablation studies on the face descriptors are performed and the results are in section 4.4.1.
3. We use the models that were created to discover biases in the datasets that were used and find out the impact that has overall in section 4.3.
4. We further test the face descriptors by replacing VGG with a smaller CNN in which the implementation is discussed in section 3.3.3.

Chapter 2 explains much of the needed technical background for the implementation, from what a Neural Network is and what metric learning is to each of the face descriptors used and what face descriptors are. Chapter 3 will then discuss 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

In this chapter, I discuss the necessary technical background information needed for the project in sections 2.1 through 2.6, then the requirements of the project are elaborated on in section 2.7 before the starting point and Software Engineering practices are discussed in section 2.8.

2.1 Convolutional Neural Networks

To talk about Convolutional Neural Networks, first we must discuss what an artificial neural network (ANN) is. This is just a set of layers like we see in figure 2.1 which have neurons in each of the layers which each output one value.

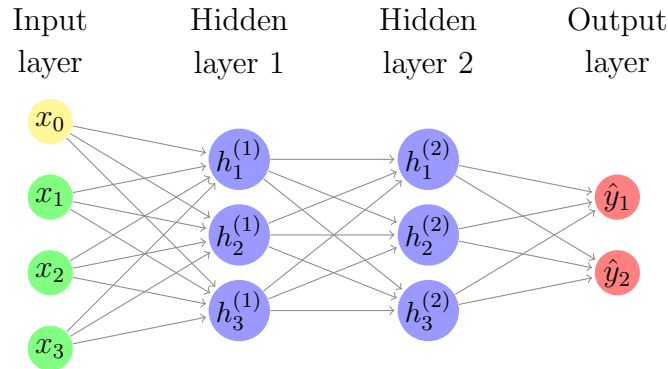


Figure 2.1: An example of an artificial neural network with two hidden layers

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 ANN 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 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 [16]. 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 regular neural networks.

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 in order to create the output. A kernel is a matrix which is usually small and has sizes less than that of the input image but expands 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 started in the top left of the input and the dot product of the kernel and the input sub-array is taken to be the output at that cell. Then the kernel is moved along as much as specified in the stride. So, if the stride was 2×2 , then the kernel would move 2 to the right and do the same thing. An illustration of this is seen in figure 2.2.

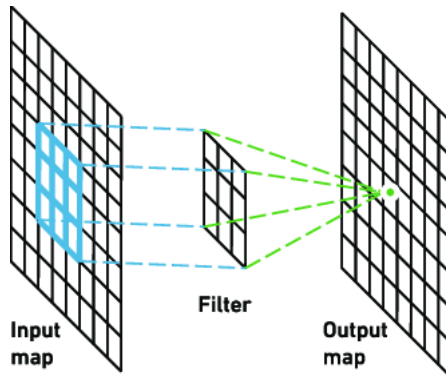


Figure 2.2: A convolution done on an input. Image sourced from Yakura et al. (2018) [22]

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, which is usually ReLU, is applied to the output.

2.1.2 Pooling Layers

Pooling layers in a CNN tend to reduce the dimension of the representation. This then reduces the number of parameters in the representation. There are multiple pooling layers, 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 $F \times F$ subarray 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 2×2 in figure 2.3.

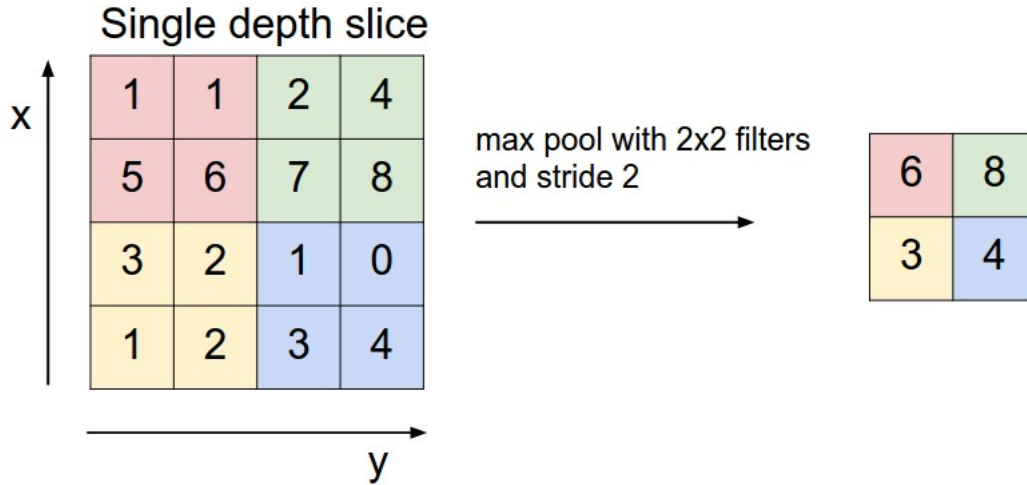


Figure 2.3: An example of the max pooling operation done on a single slice of the input

2.1.3 Activation Functions

There are a multitude of functions which are applied to the output of a layer which have different effects. These are called activation functions and we'll talk about two of them here, ReLU and Softmax.

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 if it is used as the activation layer for some layers as opposed to other ones.

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 it makes it 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 cascade classifier and is discussed in appendix A.

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 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. There are multiple face descriptors that can be made of 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 value¹:

$$0.299 \times R + 0.587 \times G + 0.114 \times B$$

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.

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

¹Information taken from https://docs.opencv.org/3.4/de/d25/imgproc_color_conversions.html using the RGB to Gray color conversion

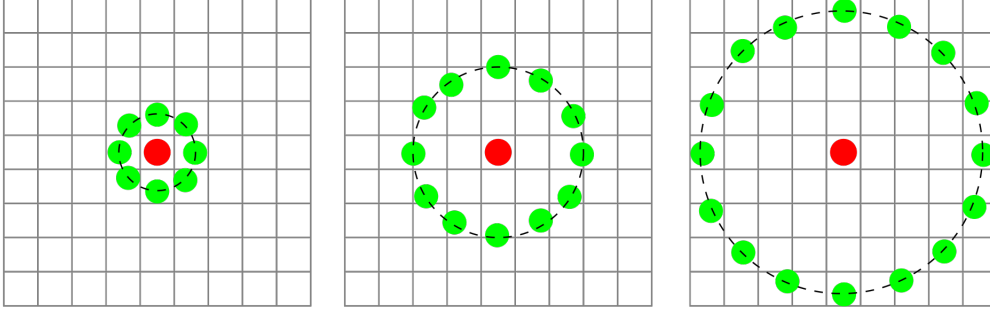


Figure 2.4: Potential neighborhoods of the pixel

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.

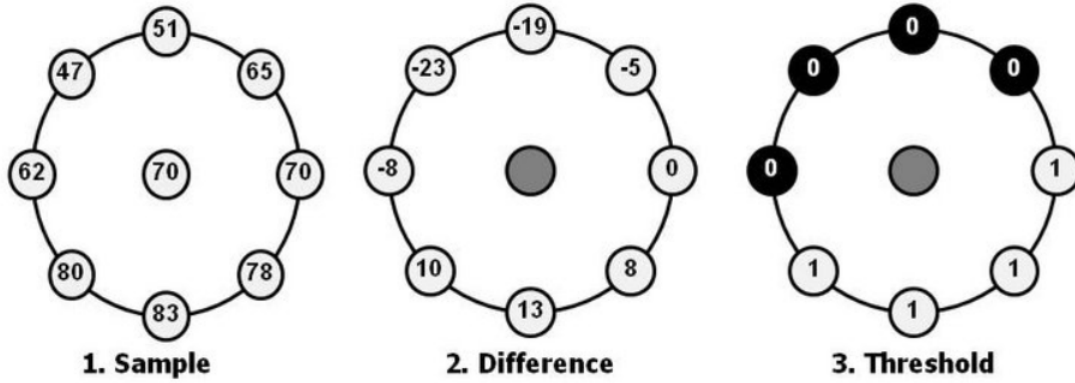


Figure 2.5: Example of LBP operator on a pixel

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:

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

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 64×64 , we split it up into non-overlapping blocks of size 8×8 , of which there are 8×8 of. 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

Another face descriptor is Histogram of Gradients (HOG) [3]. 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 due to the fact that 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

First, we need to calculate the gradients of the image. As the image function is discrete, in other words 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 3 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 which has the maximum magnitude among the 3 channels and the corresponding angle is used as well.

Furthermore, the angles returned range between 0° and 360° . 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 however, 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 ranges of $[0, 20)$, for example. 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 θ , if $0 \equiv \theta \bmod 20$ then:

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

Otherwise, we have it that θ is between two angles, ϕ_1 and ϕ_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 θ is to the label. So, for the label ϕ_1 , we have that we add to the label associated with ϕ_1 the value:

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

And, similarly for ϕ_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 ϕ_2 as 180° for the sake of this calculation but we add the value to the label 0 since $0 \equiv 180 \bmod 180$.

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 16×16 blocks of size 4×4 first and then 8×8 blocks of size 8×8 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 [11] is another way of obtaining face descriptors, although the original SIFT algorithm differs from how it is used in the project, though it is still important to briefly go over.

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 64×64 , then the images in the second octave will be 32×32 and in the third it would be 16×16 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, σ 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 2.6.

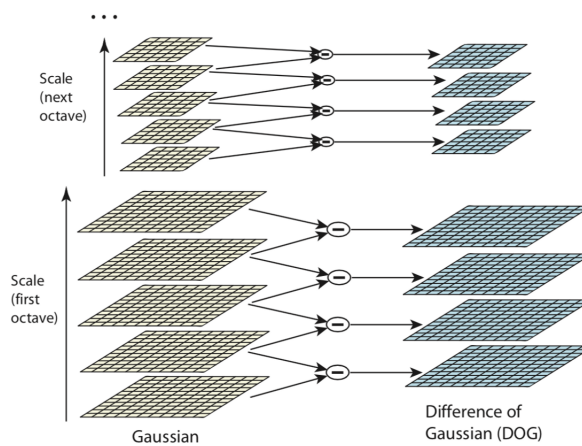


Figure 2.6: The Difference of Gaussians being created. Image reproduced from Lowe (2004) [11]

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

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

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:

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

And if the value of the extrema is less than 0.03×255 , 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 peak above 80% of it is also considered to calculate the orientation.

Now, for each keypoint, we have the location, scale, and orientation so we can finally create the descriptor. A 16×16 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 talked about in section 3.3.1.

2.3.4 VGG

Another way of getting face descriptors for an image is to use the VGG network. The original VGG network [20] is a CNN which takes in a 224×224 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 with the probability that the i^{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, more parameters. The architecture of the second deepest configuration is configuration D in Figure 2.7.

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 (224×224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	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 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Figure 2.7: The different configurations that VGG can have. Image reproduced from Simonyan et al. (2015) [20]

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 [17]. 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 dataset² was created by the authors of [17].

The VGG model for faces is then the model described above with the softmax layer having 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 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 [21]. 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

²The dataset is available at https://www.robots.ox.ac.uk/~vgg/data/vgg_face/

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 hold, 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 2.8.

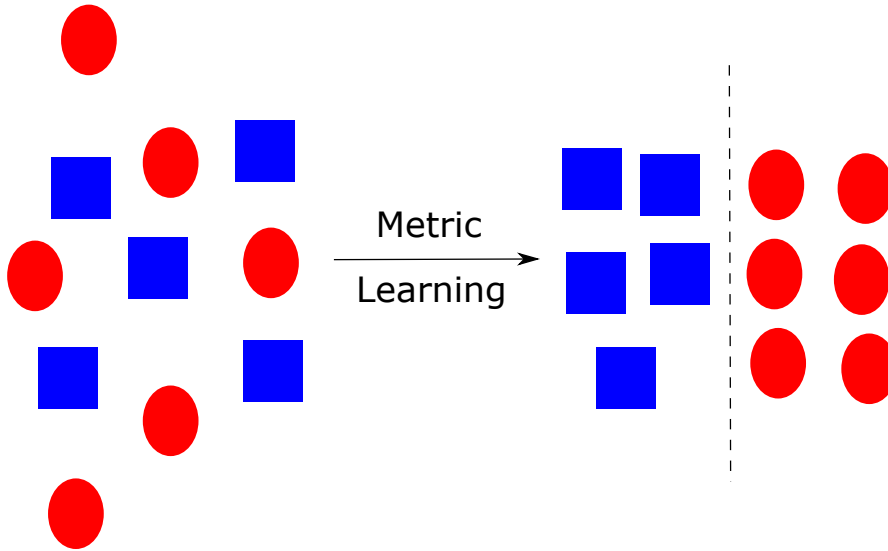


Figure 2.8: A high-level view of metric learning

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.

2.6 Principal Component Analysis

Principal Component Analysis (PCA) [8] 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

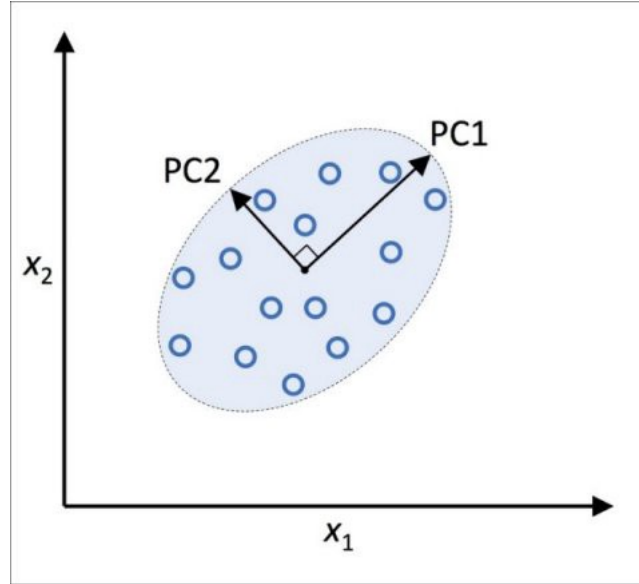


Figure 2.9: Principal components of an arbitrary 2 dimensional dataset

algorithm identifies *principal components* which are the directions of maximum variance in the dataset, which can be visualized in figure 2.9.

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:

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

Where:

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

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 Σ , 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 $k \leq p$.

2.7 Requirements Analysis

The main requirement of the project is stated in appendix B which is to replicate the results from Liang et al. (2019) [10] 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 for Kin Recognition. 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 64×64 image on disk.
- **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 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 its 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 of a medley of scripts which were used to integrate the functions that were created for each of the requirements with the datasets, which required a data preparation module.

2.8 Software Engineering Practices

2.8.1 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.

2.8.2 Tools Used

The project was written in Python 3.6. In order to separate the environment that the project needs with 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 performance of many parallel computations, such as matrix multiplication, and provided a simple interface to do these with.
2. **Keras and Tensorflow**: 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 and PCA.
5. **Scipy**: This library provided a function that solved the general eigenvalue problem, given two numpy arrays which was 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.

2.8.3 Datasets

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

The KinFaceW datasets [12] [13] 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 64×64 . 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 [18] 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.

2.8.4 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.

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.

2.8.5 Licensing

The SIFT algorithm had been patented in the US. However, the patent expired last year, and the patent was only to protect 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.

Chapter 3

Implementation

In this chapter, the repository is examined in detail in section 3.1, each of the modules are discussed in sections 3.2, 3.3, 3.4, and 3.5, in which the extension to replace VGG with another face descriptor is discussed in section 3.3.3, and how the project is 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 section 2.2. The rest of the modules are discussed in sections 3.2, 3.3, 3.4, and 3.5. The `test` folder contains all of the unit tests for each of the files in the aforementioned modules. The `scripts` folder contains any scripts that would be run directly with Python which integrate the modules together in order 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.

The `data` folder contains the three datasets which contain the images and the meta data 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.

3.2 Data Preparation

3.2.1 Cross-Validation

We split the datasets for each relationship into 5 folds each in which each fold within a relationship has around the same number of pairs. Each training set composed of 4 of the folds and the last fold was designated as the test set, which composed of 5 training/testing set combinations. 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.

```

KinRecognitionWGEML
├── data/
├── out
├── ablation_studies/
├── pairwise_accs/
├── src
│   ├── data_preparation
│   │   ├── PCA.py
│   │   ├── prep_cross_valid.py
│   │   ├── properly_formatted_inputs.py
│   │   └── save_and_load.py
│   ├── face_descriptors
│   │   ├── CifarNet.py
│   │   ├── HOG.py
│   │   ├── LBP.py
│   │   ├── SIFT.py
│   │   └── VGG.py
│   ├── face_detection/face_detection.py
│   ├── prediction/predictor.py
│   ├── scripts
│   │   ├── ablation_study.py
│   │   ├── preprocessing_fds.py
│   │   ├── preprocessing_TSK.py
│   │   ├── testing.py
│   │   ├── training.py
│   │   └── get_pairwise_accuracies.py
│   ├── test/
│   └── WGEML
│       ├── constants.py
│       └── WGEML_training.py
├── venv/
├── .coveragerc
├── .gitignore
├── Makefile
├── README.md
└── requirements.txt

```

Figure 3.1: Folder Structure of the project

3.2.2 Positive and Negative Pairs

Within each dataset, there are pairs of images which 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. 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.

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 [13]. 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 PCA

In order to save on computation time for the training process as well as on disk space, PCA is used in order to reduce the dimensions of each of the face descriptors. 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 dimensions to 100 through the method defined in Section 2.6.

3.2.4 Saving Results to Disk

Running the project end-to-end each time would take too long and, in order to better explore some of the results that were taken, the pipeline needed to be split into modules. This ended up being the preprocessing stage, the training stage and then 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`.

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, there are some extra details that 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 [11] 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 instead obtained from dividing the image into 7×7 overlapping patches on a 16×16 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 [20] 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 and the pretrained weights were uploaded to their site¹. Therefore, it was easier to download the weights from the website and reformat it as needed.

Firstly, due to the fact that the weights were used 566MB, this wasn't something that could just be uploaded to GitHub. Thus, when predicting, I made sure to check if there was a weights file with a specific name in the `src/face_descriptors` folder path. If there wasn't, then the weights would be downloaded from the site. Due to the fact that I was using the Torch weights, but I was implementing this in Keras, I needed to reshape a few of the fully-connected layers. For example, one layer had the shape $(7, 7, 512, 4096)$ for the weights but for Keras, it needed to be in the shape $(7 \times 7 \times 512, 4096)$. This was easily solved by using the `np.reshape` function for each of the fully-connected layers and these slightly modified weights were saved to an `h5` file on disk.

From there, the weights would be on disk so, to get a face descriptor for a set of images, I would create the VGG face model with the softmax layer, then load the weights into the model before returning the model without the softmax layer. This model then was used to get the face descriptors for the set of images given. Due to the fact that VGG takes images of size 224×224 , the images needed to be resized to 224×224 before it 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 which had size $(n, 4096)$.

¹https://www.robots.ox.ac.uk/~vgg/software/vgg_face/

3.3.3 CifarNet Extension

In the CifarNet extension, I replaced the VGG network with a smaller network that was trained on CIFAR-10 [9], which is called CifarNet², 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.

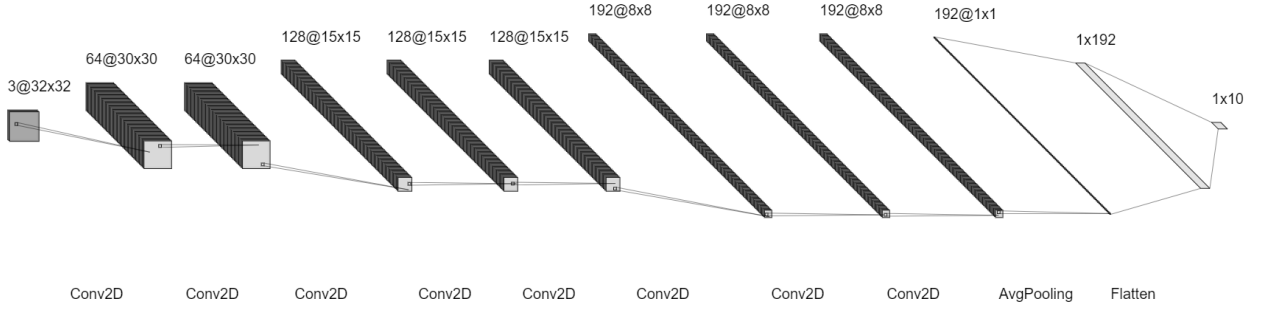


Figure 3.2: The Architecture of the CifarNet model

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 in a similar manner 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) [10]. 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 [10].

²Architecture taken from <https://github.com/deep-fry/mayo/blob/master/models/cifarnet.yaml>

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} .

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:

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

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. In order 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 [7]:

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

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 β is a regularization parameter which 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 sim ranges from 0 to 1 as needed of a similarity function. This differs from the version mentioned in [10] 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, θ . If $\text{sim}(\mathbf{x}, \mathbf{y}) \geq \theta$, then the pair is labelled

as having that relationship and, otherwise, it isn't. As this value θ wasn't mentioned in [10], 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

Due to the fact that tri-kin relationships have 3 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 3 main stages to running the project, which are, for each dataset, the preprocessing stage, the training stage and then 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.

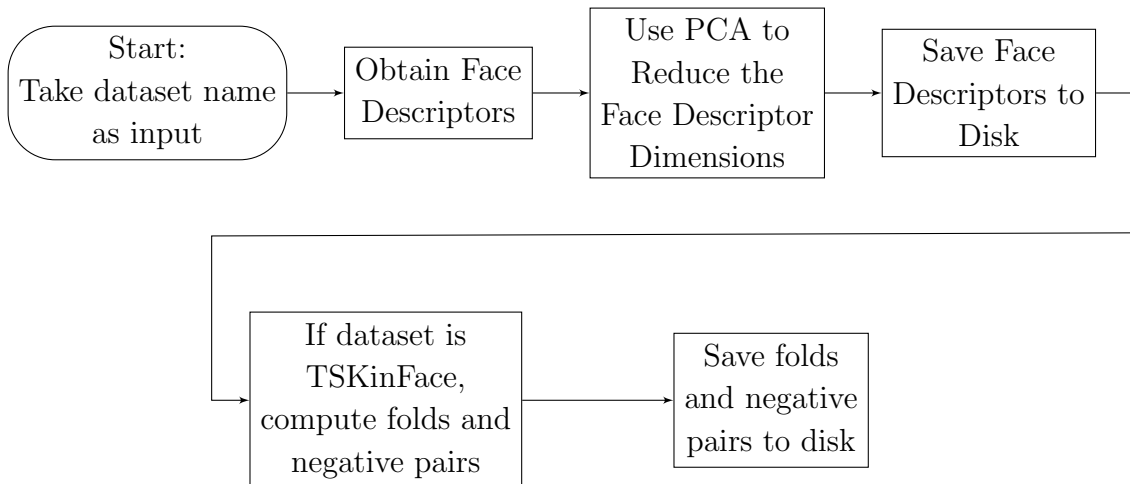


Figure 3.3: The preprocessing pipeline

3.6.2 Training

Once all of the face descriptors for each dataset are computed and, for each relationship in the dataset and for 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.

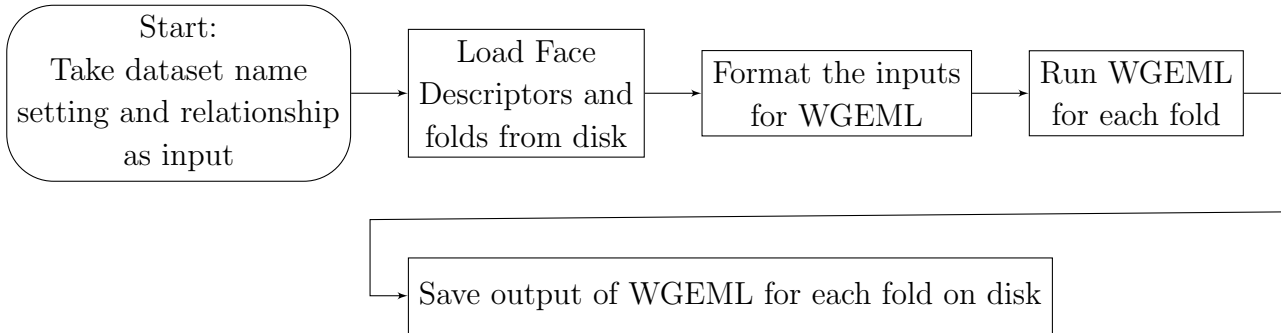


Figure 3.4: The training pipeline

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 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

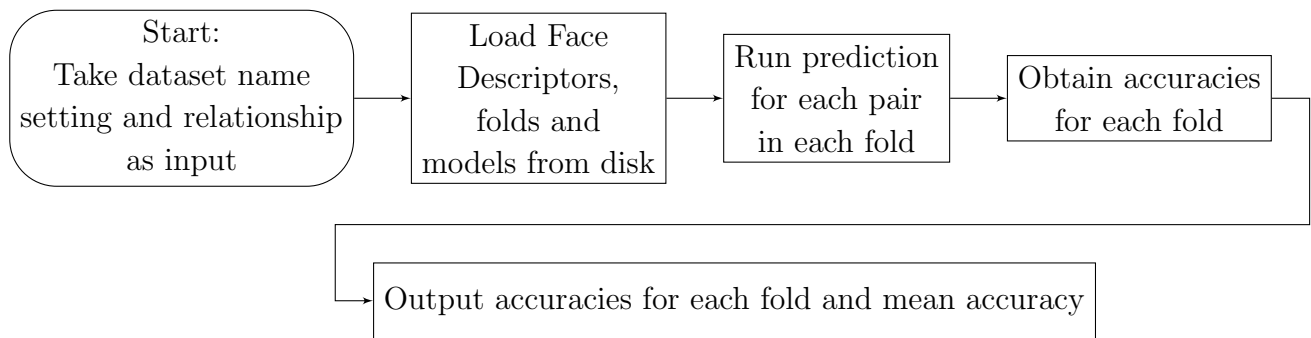


Figure 3.6: The testing pipeline

Chapter 4

Evaluation

In this chapter, I will discuss the accuracies that were obtained from the project in section 4.1 whether the success criterion was achieved in section 4.2. Then, the potential biases in these datasets are discussed in section 4.3, the ablation studies are discussed in section 4.4, the CifarNet extension results will be discussed in section 4.5 and finally the unit tests and their coverage will be briefly discussed in section 4.6.

4.1 Original Results

The project is run end-to-end for each dataset, setting and relationship and the raw accuracies are displayed in table 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 displayed in table 4.2 where a positive number means my accuracies were better.

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 4.1: Accuracies of WGEML applied to each dataset for each relationship

We can see that none of the absolute values of the differences go higher than 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

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	2.21	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	-8.98	-6.78	-6.86	-8.69

Table 4.2: Differences in accuracies between my implementation and [10] in %

pairs for the testing phase of TSKinFace as TSKinFace didn't come with any premade negative pairs. 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 [10] within a $\pm 15\%$ accuracy by using the WGEML methodology which means that the project is a success according to appendix B. Furthermore, the extensions were also completed and the results of which are discussed in sections 4.3, 4.4, and 4.5.

4.3 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 datasets that were available. 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 3 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 3 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 4.1 and the average accuracy differences for each dataset across all relationships and test datasets is shown in figure 4.2.

The main reason for this is due to the fact that 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 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 in order 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

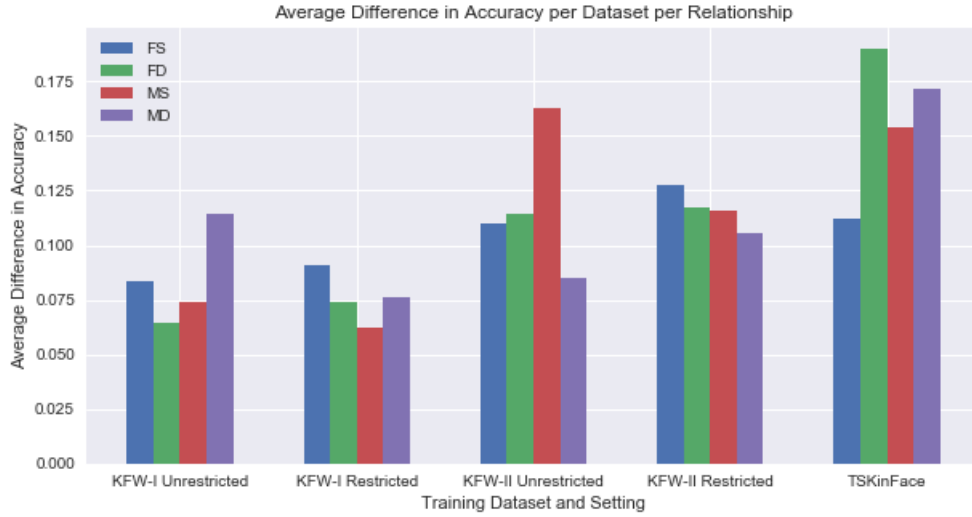


Figure 4.1: The average accuracy differences of the two test datasets for each training dataset and relationship

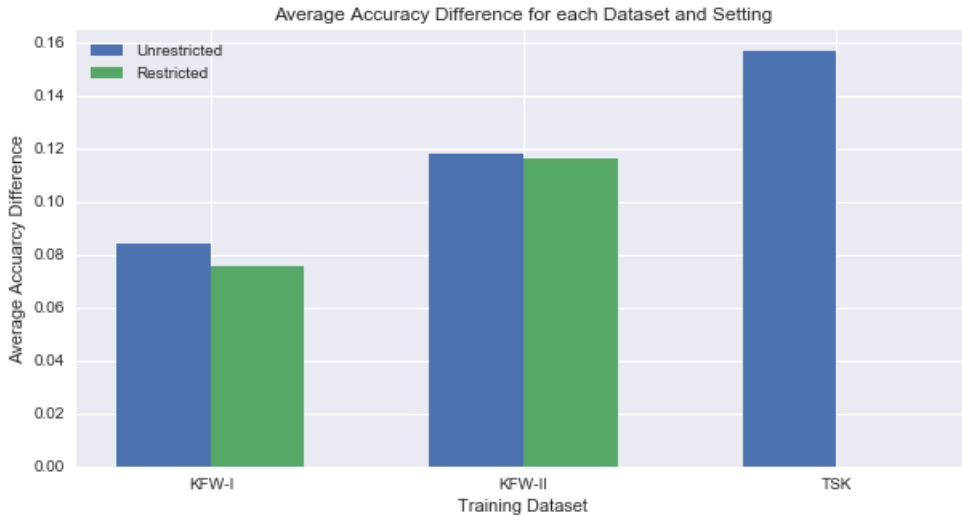


Figure 4.2: The average accuracy differences of the two test datasets and relationship for each training dataset

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 is able to generalize better which explains why the differences are smaller than that of KinFaceW-II and TSKinFace.

4.4 Ablation Studies

4.4.1 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 amount 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 4.3.

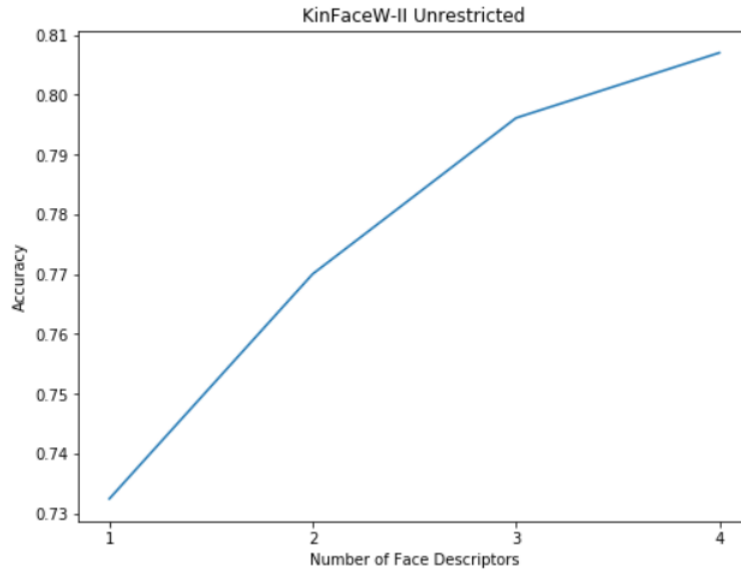


Figure 4.3: The plot of accuracy versus number of face descriptors used for KinFaceW-II Unrestricted

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 4.4.

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, from what I could see. This makes sense since VGG is a deep feature detector whereas the other face descriptors use the texture only.

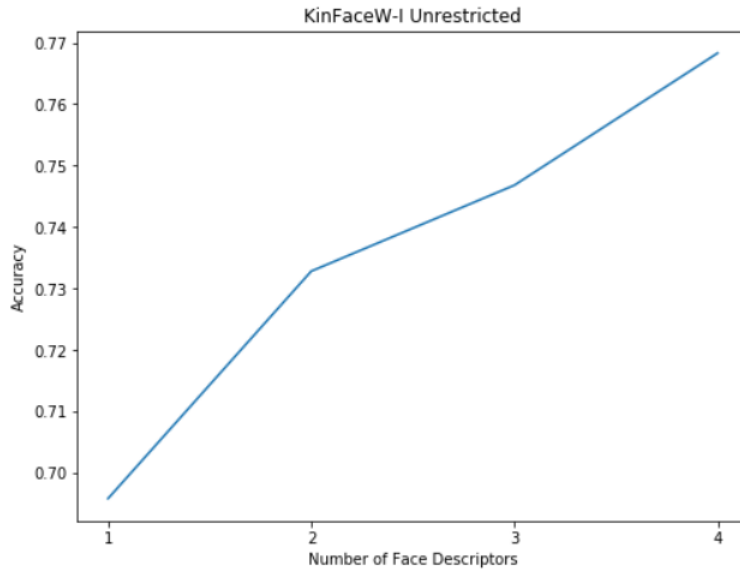


Figure 4.4: The plot of accuracy versus number of face descriptors used for KinFaceW-I Unrestricted

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 3-4 instances where this occurred.

4.5 CifarNet Extension

When I replaced VGG with CifarNet (CFN), I found that 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 biggest difference being a 1% difference in accuracy). When grouped by relationship, I found that the father-son relationship had the least change in accuracy with an average 0.092% accuracy difference. This is shown in figure 4.5 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 is shown in figure 4.6.

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 mother-daughter 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 mother-son relationship. The KinFaceW-II Unrestricted setting with the mother-son 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 less 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 some cases.

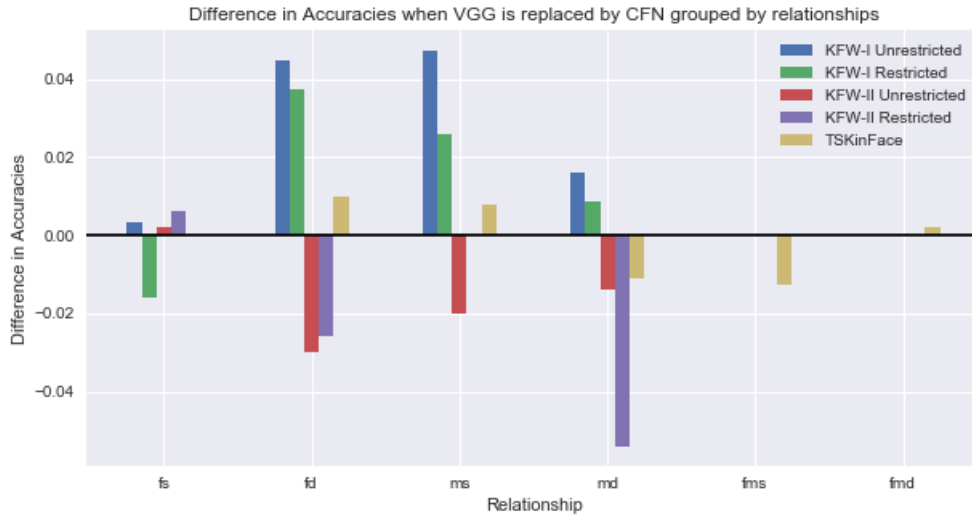


Figure 4.5: Differences in accuracy when VGG is replaced with CFN grouped by relationship

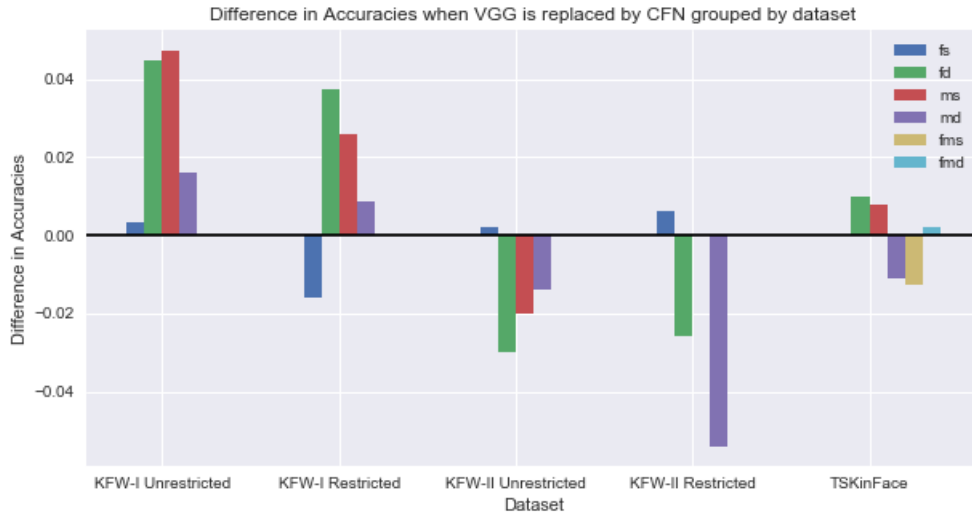


Figure 4.6: Differences in accuracy when VGG is replaced with CFN grouped by dataset

4.6 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 due to the fact that 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 4.7. 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 which trains the model and another function which 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 the most of the untested code.

Name	Stmts	Miss	Cover
src/WGEML/constants.py	3	0	100%
src/WGEML/create_values.py	57	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	27	0	100%
TOTAL	592	38	94%

Figure 4.7: Coverage of the unit tests of the project

Chapter 5

Conclusion

Overall, the project was successful since I was able to replicate implement each of the core components, as specified in section 2.7 and, this allowed me to replicate the findings of 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 4.4.1. Furthermore, sources of bias in the datasets were found in section 4.3.

5.1 Lessons Learnt

One of the main lessons I learnt 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 individual function I wrote and used is correct. For example, solving the eigenvalue problem and getting the top d vectors with respect to 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.

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, due to the fact that 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, I think using more sophisticated software engineering tools such as some continuous integration tool such as Jenkins together with a containerization tool such as 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

In the future, the model can be extended to work on videos. This was an extension at the beginning of the project which ended up not coming to fruition due to the fact that I couldn't obtain a dataset for it. 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.

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] W.D. Hamilton. The genetical evolution of social behaviour. ii. *Journal of Theoretical Biology*, 7(1):17 – 52, 1964.
- [7] Y. Jia, F. Nie, and C. Zhang. Trace ratio problem revisited. *IEEE Transactions on Neural Networks*, 20(4):729–735, 2009.
- [8] I.T. Jolliffe. *Principal component analysis*. Springer-Verlag New York, 2002.
- [9] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [10] 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.
- [11] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91110, November 2004.
- [12] 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.
- [13] 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.

- [14] Jill M. Mateo. Perspectives: Hamilton’s legacy: Mechanisms of kin recognition in humans. *Ethology*, 121(5):419–427, 2015.
- [15] 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.
- [16] Keiron O’Shea and Ryan Nash. An introduction to convolutional neural networks. *ArXiv e-prints*, 11 2015.
- [17] 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.
- [18] Xiaoqian Qin, Xiaoyang Tan, and Songcan Chen. Tri-subject kinship verification: Understanding the core of a family, 2015.
- [19] Joseph P Robinson, Ming Shao, and Yun Fu. Visual kinship recognition: A decade in the making, 2020.
- [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [21] 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.
- [22] 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.
- [23] Jun Yu, Mengyan Li, Xinlong Hao, and Guochen Xie. Deep fusion siamese network for automatic kinship verification. *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*, pages 892–899, 2020.

Appendix A

Cascade Classifier

Appendix B

Project Proposal