# Assessing the Influencing Factors on the Accuracy of Underage Facial Age Estimation

Felix Anda, Brett A. Becker, David Lillis, Nhien-An Le-Khac and Mark Scanlon

*Forensics and Security Research Group*
*University College Dublin*
Dublin, Ireland
felix.anda@ucdconnect.ie, {brett.becker, david.lillis, an.lekhac, mark.scanlon}@ucd.ie

*Abstract*—Swift response to the detection of endangered minors is an ongoing concern for law enforcement. Many child-focused investigations hinge on digital evidence discovery and analysis. Automated age estimation techniques are needed to aid in these investigations to expedite this evidence discovery process, and decrease investigator exposure to traumatic material. Automated techniques also show promise in decreasing the overflowing backlog of evidence obtained from increasing numbers of devices and online services. A lack of sufficient training data combined with natural human variance has been long hindering accurate automated age estimation – especially for underage subjects. This paper presented a comprehensive evaluation of the performance of two cloud age estimation services (Amazon Web Service's Rekognition service and Microsoft Azure's Face API) against a dataset of over 21,800 underage subjects. The objective of this work is to evaluate the influence that certain human biometric factors, facial expressions, and image quality (i.e. blur, noise, exposure and resolution) have on the outcome of automated age estimation services. A thorough evaluation allows us to identify the most influential factors to be overcome in future age estimation systems.

*Index Terms*—Machine Learning, Digital Forensics, Facial Age Estimation, Human Biometrics

## I. INTRODUCTION

The number of internet users is constantly rising and each year increasing numbers of young people are online. The most vulnerable groups in cyberspace are subject to possible exposure to cybercrimes such as phishing attacks, hacking, sextortion, child sexual exploitation material (CSEM), and child grooming.

Digital Forensic (DF) laboratories are frequently handling evidence involving minors. These cases involve the identification of victims of human trafficking and the detection of CSEM, which has been regarded by many as one of the most damaging crimes [1]. Exposure to the analysis of illicit content affects law enforcement officers by causing psychological distress such as secondary traumatic stress disorder [2]. Incorporating technologies such as Artificial Intelligence into DF has potential to avert the impact on investigators.

Today, digital information is widely shared through social media, IoT devices, surveillance, cloud services, etc. Each source compounds evidence acquisition and processing, and contributes to the extensive backlog of cases requiring digital forensic analysis [3]. This variety of sources is a hindrance frequently encountered in modern policing [4]. Automated facial age estimation is a critical service that can potentially elevate the overflow through automatically classifying data on behest of investigators and focusing their analysis efforts.

As part of this work, the VisAGe dataset[1] is assessed against two of the best performing cloud age estimation services, i.e., Microsoft Azure's Face API and Amazon Web Service's (AWS's) Rekognition service [5]. VisAGe is fully human annotated with the values of the ground-truth age per single-faced image – this facilities the performance evaluation of each aforementioned cloud services in terms of Mean Absolute Error (MAE), which is a measure between the actual and predicted age. A variation of this measurement was evaluated against each feature to compute the Pearson Correlation Coefficient (PCC) between the two variables. Whilst weak correlations throughout the entire age range were recurrent, important results consisting of mild and strong correlations were obtained and the major trends between them were evaluated.

A summary of the contribution of this work includes:

- Identification of the influencing factors for accurate facial age estimation for underage subjects and their weighting on the accuracy obtained.
- Analysis of trends within both strong positive and negative linear correlations and how they affect the underage facial age estimations for different ages.
- Comprehensive evaluation of Microsoft Azure Face API and AWS Rekognition's facial image attributes and their association to facial age estimation.
- Analysis of the VisAGe underage dataset facial attribute distribution.

## II. LITERATURE REVIEW

### A. Digital Forensic Backlog

The requirement for DF investigation has exploded due to the rapid increase of both the number of cases requiring DF analysis and the volume of information to be processed per case (due to increases in the number of relevant devices and their capacities) [4], [6]. This puts prosecutions at risk and can lead to cases dismissals. The use of data mining, triage processes and data reduction has been suggested to alleviate this backlog [3].

---

[1] https://visage.forensicsandsecurity.com/

### B. Influencing Factors

The factors affecting facial ageing have been categorised into intrinsic and extrinsic components [7]. For the former, there are internal factors such as size of the bone, genetics or facial changes due to the development of a child. For the latter, any presence of external factors including the environment, habits, diet, makeup and cosmetics, etc.

*1) Facial Expressions:* One example of influencing factors in age estimation is facial expressions. Voelkle et al. [8] found that happy facial expressions are mostly underestimated whereas, smiling, frowning, surprise and laughing may introduce facial lines that are confused for wrinkles and thus impact on the age estimation performance.

*2) Noise:* Noise introduces more error onto the estimation depending on its magnitude. It is a randomness that affects an image due to either brightness, colour or digital encoding, and often occurs during image capture, digital sharing, etc. [9]. The presence of noise in an image is expected to be linearly correlated with performance.

*3) Makeup:* Facial cosmetics have been found to influence perceived facial age estimation; a simple cosmetic alteration is capable of compromising the outcome of a biometric system [10]. Lip makeup was found to be the most prominent of the cosmetic range with a mild correlation to the decay in age estimation accuracy for specific ages. Moreover, Chen et al. [11] found that the presence of cosmetics can hide facial imperfections caused by age, e.g., wrinkles and dark spots, resulting in underestimation.

### C. Data Bias

Wang et al. [12] states that biased databases are more commonplace; therefore, trained models are unable to handle race/ethnicity and gender without bias and thus cause the performance to decline. The influence of race and gender seems to be the most common as both of these attributes play an important role in age estimation. Anda et al. [13] evaluated the influence of gender in automated age estimation and determined that for four age prediction services, the accuracy for female subjects is lower than for males. In previous studies, the effect of ageing has also been found to vary within gender, with male faces tending to age slower compared to female faces [14]. Models trained with unbalanced datasets will produce biased results thus leading to compromised accuracy.

## III. METHODOLOGY

The VisAGe dataset was processed by Azure's Face API and AWS Rekognition and the age estimations obtained from the two cloud services were measured against the ground-truth age in the dataset. The difference between the two values has been denoted as the error difference ($Er_d$). This has been used as the principle measurement in assessing the accuracy of the underage facial age estimation. Additional features of both cloud facial analysis services were utilised to classify and annotate the data as per Tables I and II. To process the correlations between variables, the object attributes have been broken down into categorical values.

Having determined the attributes of each image and their associated $Er_d$, the correlation between the two variables of data was then calculated to identify which attributes were the larger influencing factors of $Er_d$ and by what gravity, e.g., weak, mild, or strong.

Attributes with mild to strong correlations had influence in the accuracy of the underage facial age estimation. Through analysing the distribution of errors, as discussed in Section IV-A2, the error bin of 0 to 5 contains the largest amount of occurrences in comparison to succeeding error margins. Henceforth, the investigation has been split into the gravity of errors in order to identify traits that most of the data adhere to, versus the traits of the minorities, i.e., data that lies within $Er_d > 5$.

### A. VisAGe Dataset

The VisAGe dataset was created to address the shortage of adequate underage databases available to investigators [15]. It is composed of a three-stage validation process comprising of both automatic age and gender classifications provided by Microsoft Azure Cognitive Face API, and a manual Quality and Control system through the VisAGe web voting application.

### B. Cloud Services

Two cloud services were used in this study to provide the underage facial age estimations of each image within the VisAGe single-faced dataset; Amazon AWS Rekognition Service and the Microsoft Azure Face API service.

*1) Microsoft Azure: Face API:* This service assisted the annotation of each record according to the detected facial attributes such as perceived emotion, presence of facial hair and makeup, facial expressions like happiness, contempt, neutrality, and fear, etc. A comprehensive list is presented in Table I.

TABLE I
MICROSOFT AZURE COGNITIVE SERVICES FACE API ATTRIBUTES [16].

| Field | Description |
|---|---|
| emotion | Neutral, anger, contempt, disgust, fear, happiness, sadness, and surprise. |
| noise | Noise level of face pixels. |
| age | "Visual age" number in years. |
| gender | Estimated gender with male or female values. |
| makeup | Presence of lip and eye makeup. |
| accessories | Accessories around face, including 'headwear', 'glasses' and 'mask'. |
| facialHair | Moustache, beard and sideburns. |
| hair | Group of hair values indicating whether the hair is visible, bald, and hair colour if hair is visible. |
| headPose | 3-D roll/yaw/pitch angles for face direction. |
| blur | Face is blurry or not. 'Low', 'Medium' or 'High'. |
| smile | Smile intensity, a number between [0,1]. |
| exposure | Face exposure level. Level returns 'GoodExposure', 'OverExposure' or 'UnderExposure'. |
| occlusion | Values are Booleans and include 'foreheadOccluded', 'mouthOccluded' and 'eyeOccluded'. |
| glasses | Glasses type. Values include 'NoGlasses', 'ReadingGlasses', 'Sunglasses', 'SwimmingGoggles'. |

*2) Amazon AWS: Rekognition Service:* Amazon Rekognition is a pre-trained image analysis service. Its face detection and analysis service was used to perform several visual analyses on VisAGe; extracting facial attributes such as facial hair, expressions, etc., detected on each single-faced image. The attributes, as outlined in Table II, were then correlated against Amazon's facial age estimator to provide a comprehensive evaluation on the accuracy of underage facial age estimation against the influencing factors.

TABLE II
AMAZON AWS REKOGNITION ATTRIBUTES [17]

| Field | Description |
|---|---|
| Age.Range | Estimated age range. |
| Smile.Value | Smile value detected true or false. |
| Eyeglasses.Value | Eyeglasses detected true or false. |
| Sunglasses.Value | Sunglasses detected true or false. |
| Gender.Value | detected gender on subject. |
| Beard.Value | Beard detected true or false. |
| Moustache.Value | Moustache detected true or false. |
| EyesOpen.Value | Open eyes detected true or false. |
| MouthOpen.Value | Open mouth detected true or false. |
| Emotions | Detection true or false for each array. |
| Landmarks[0] | X-axis and Y-axis positions. |
| Roll (Degree) | Face titled to the side. |
| Yaw (Degree) | Face turned to the side. |
| Pitch (Degree) | Face titled up or down. |
| Brightness | Brightness of the image. |
| Sharpness | Sharpness of the image. |
| Confidence | Certainty of the estimation. |

*C. Skin Tone Classifiers: Simple Skin Detection and Face Colour Extraction*

Automated detection of skin tone has received considerable attention from researchers – specifically for biometrics and computer vision applications [18], [19]. For this study, the impact of two approaches has been evaluated: Simple Skin Detection (SSD) and Face Colour Extraction (FCE). Both approaches are based on k-means clustering[2] in order to determine and classify a subject's skin tone.

SSD refers to unsupervised skin tone estimation/segmentation; the approach predicts skin tone from an image of a subject, while doing a rough segmentation of the skin based on a pixel-wise classifier [20]. The algorithm consists of two main components: foreground/background separation using Otsu's Binarisation and pixel-wise skin classifier based on HSV and YCbCr colour spaces [21].

The FCE approach initially detects the facial landmarks using the *Dlib* library [22]. Subsequently, noise is removed by applying the convex hull algorithm[3] on the facial land-marked point. Finally, the RGB values of the skin are computed using a histogram-based clustering algorithm. These values can be seen in Table III and have contributed in a mild inverse fashion to the error difference, i.e., the more "red" the values, the less the error.

[2]k-means clustering is a method for vector quantization – mainly used for cluster analysis in the data mining field.
[3]Convex hull is a fundamental structure for both mathematics and computational geometry [23]

*D. Pearson Correlation Coefficient*

The Pearson Correlation Coefficient (PCC) measures the linear correlation between two variables. In this work, these are the attribute and $\text{Er}_\text{d}$. The value of the coefficient lies between +1 and -1; where ±1 indicates a perfect correlation and 0 represents no correlation at all. A negative coefficient signifies an inverse relationship between the variables. For a sample of data, such as that examined here, the PCC is often represented as $r_{xy}$ and is defined in Equation 1:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}} \quad (1)$$

where $n$ is the size of the sample, $x_i$ and $y_i$ are individual sample pairs and $\bar{x}$ and $\bar{y}$ are the mean of $x$ and $y$. The correlation value obtained for each sample, i.e., the facial attribute and $\text{Er}_\text{d}$ pair, was matched inline with a scale of weak, mild, or high. It is important to note that for the purpose of this work, weak, mild and strong correlations are characterised with $0.1 - 0.29$, $0.30 - 0.49$, $0.50 - 1$ correlation values respectively (whereby the negatives of these values represent inverse correlations). These definitions have been defined in a computer forensic related study regarding analysis of correlations of Internet usage [24]. Conversely, correlation close to zero, specifically within the $-0.1 - -0.1$ range has been referenced as minuscule correlation.

## IV. EXPERIMENTS AND RESULTS

Due to the different rates of performance, the two cloud services have been assessed independently. Overall, Microsoft Azure achieves a MAE of 2.082 for the VisAGe dataset, whilst AWS has a MAE of 4.075. Furthermore, the distribution of $\text{Er}_\text{d}$ for each class service has been analysed.

It must be noted that for all succeeding correlation Figures, the attribute error is shown to have a positive perfect degree of correlation to $\text{Er}_\text{d}$. This is expected as any attributed examined with itself produces this behaviour.

*A. Microsoft Azure*

Influencing factors affecting Azure's facial age estimation have been evaluated. Section IV-A1 looks into the distribution of correlations between the $\text{Er}_\text{d}$ and other attributes in order to identify the influencing factors and their gravity towards the $\text{Er}_\text{d}$. The distribution of significant correlations of greater than or equal to 5 between attributes are outlined in Table I and the $\text{Er}_\text{d}$ for different ages are represented in Figure 2.

*1) Strong PCC Distribution per Age with $Er_d \geqslant 0$:* The distribution of strong correlation values have been evaluated per age between the variables: $\text{Er}_\text{d} \geqslant 0$ and the attributes detected. It was observed that one-year-olds were the only age that demonstrated any linear correlations. These positive strong correlations were produced by the facial hair attributes: moustache, beard and sideburns. It was anticipated that the presence of facial hair will hinder accurate estimation of facial age. However the cause of facial hair being detected for 1-year-olds was produced by incorrect detection of moustaches and beards (typically from food around the subject's mouth).

Furthermore, no attribute was identified to be of strong influencing factor towards the accuracy of the age estimator for all succeeding ages, when the $Er_d \geqslant 0$ is considered.

*2) Error Distribution:* Figure 1 is the univariate distribution of observations of the $Er_d$ value. It can be concluded that the general consensus of Azure's underage facial age estimation is reasonably accurate, i.e., the majority of scores obtained were relatively low with the bulk of the result being less than or equal to 5. It can be observed that there is a great difference on the amount of results achieving accuracy of $Er_d < 5$ versus larger $Er_d$ values of greater than or equal to 5. The distribution of strong correlations achieved in Section IV-A1 was further filtered by $Er_d < 5$ and $Er_d \geqslant 5$, as discussed in Sections IV-A3 and IV-A4 respectively.

*3) Strong PCC Distribution per Age with $Er_d$ lsess than 5:* Whenever Azure's facial age estimation demonstrates a high level of accuracy, achieving error margins $\leqslant 5$, the distribution of $|PCC| \geqslant 0.5$ presented no correlating data attributes across all ages. These results were similar to that obtained in Section IV-A1. It can be concluded that no influencing factors have been identified to be associated with the estimator achieving good results.

*4) Strong PCC Distribution per Age with $Er_d$ greater than 5:* Conversely, when the accuracy of the estimator declines beyond the error margins of 5, the distribution of strong correlations have been identified between attributes and the estimator's $Er_d$ occurring on ages 1, 2, 4 to 7 and again on 9 to 10 years old, as shown on Figure 2. No distribution of strong correlation was detected for ages 3, 8 and 11 to 18 where $Er_d$ is set to greater than 5. To delve further into identifying the attributes triggering these results and by what magnitude, Figures 3 and 4 outline the distribution of the aforementioned PCC values according to specific attributes for each age.

For one year-old subjects, an interquartile range (IQR) of strong correlation values around 0.5 to 0.75 were detected, as shown in the Figure 2, with outliers lying in the negative region. Figure 3 confirms that this outcome was the result of the age displaying strong correlations of 0.81 for exposureLevel_underExposure, noiseLevel_medium, sideburns, moustache and beard. These noted attributes, as shown in
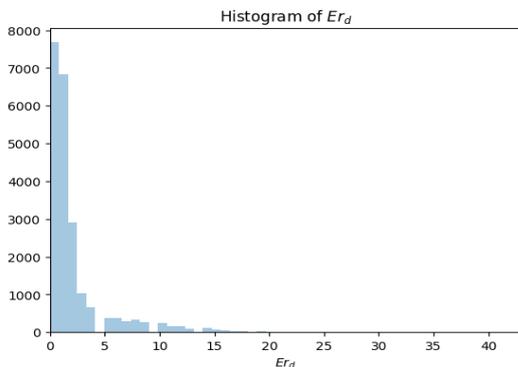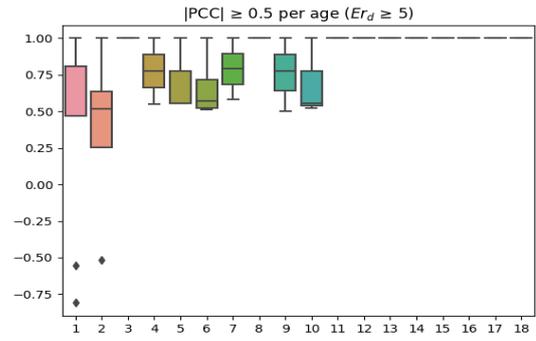


Fig. 2. Azure: Box-plot of PCC Distribution per Age, where $Er_d > 5$.

Table III, were found to have a strong linear influence to the decline in the accuracy of the Azure's age estimator for one year olds. Equally, attributes that displayed strong negative correlations of -0.55 and -0.81 for noiseLevel_low and exposureLevel_goodExposure respectively, were found to have linear influence in the improvement of estimator's performance accuracy. Furthermore, these two negative attributes contributed to the outliers in the data for age one. Such strong PCC values obtained in the age were not expected as a more diverse set of PCC figures were thought to be more probable.

In Figure 2, a similar IQR has been found for 2 year olds. This IQR lies just above the 0.25 to 0.75 range denoting that attributes with PCC values $\geqslant 0.5$ were close to the 0.5 benchmark. By referring to Figure 3, it is confirmed that strong correlating attributes had a magnitude of 0.52 and 0.51. In comparison to the preceding age, two year olds presented with more diverse assortment of attributes, only 3 of the attributes managed to achieve strong correlation values of over 0.5 in magnitude; these strong influencing attributes ($|PCC| > 0.5$) are outlined in Table III. Gender was the key prominent attribute that influenced the increase and decrease of $Er_d$ for two year olds; female subjects caused a decline in accuracy for 2 year olds, whilst male subjects were found to linearly influence the incline of the accuracy. Additionally, emotion of contempt was also found to be a strong influencing factor affecting the accuracy for two year olds. All succeeding ages, as shown on both Figures 2 and 5, present no strong negative correlation above the -0.5 threshold. Therefore no influencing factors have been identified that elevate the gravity of the $Er_d$ for ages 3 and above. Moreover, for 4 year olds, Figure 3 illustrates only one attribute to have strong association with the accuracy of the age estimator; emotion of anger with value 0.55.

Ages 5, 6, 9 and 10 all exhibit forms of facial hair correlations with the performance of Azure's facial age estimation on the underage dataset (again, miscategorisaitons at these ages). In particular, age five has both beard and sideburns attributes with strong correlation PCC values of 0.55. Similarly, both attributes have also been connected to age 6 with correlation PCC values of 0.52 and 0.62 respectively and again on



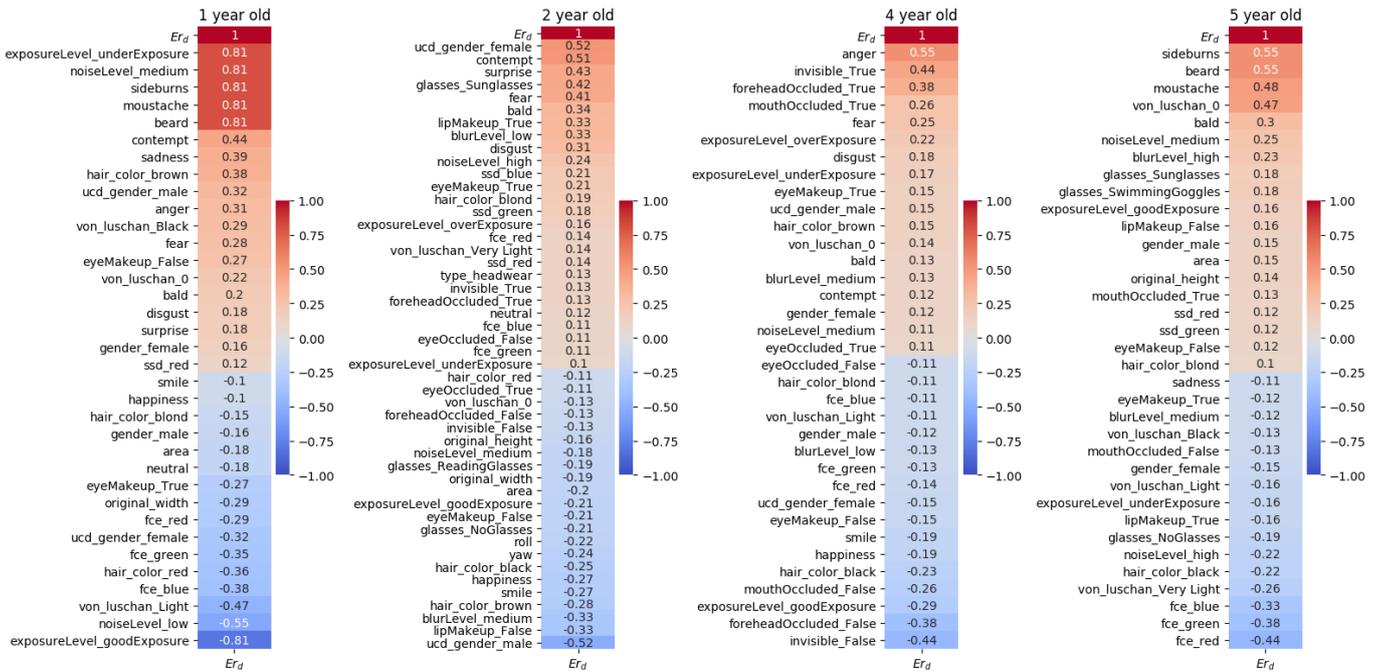Fig. 1. Azure: Distribution of $Er_d$.

Fig. 3. Azure: Strong correlations between attributes and $Er_d > 5$ for ages 1, 2, 4 and 5.
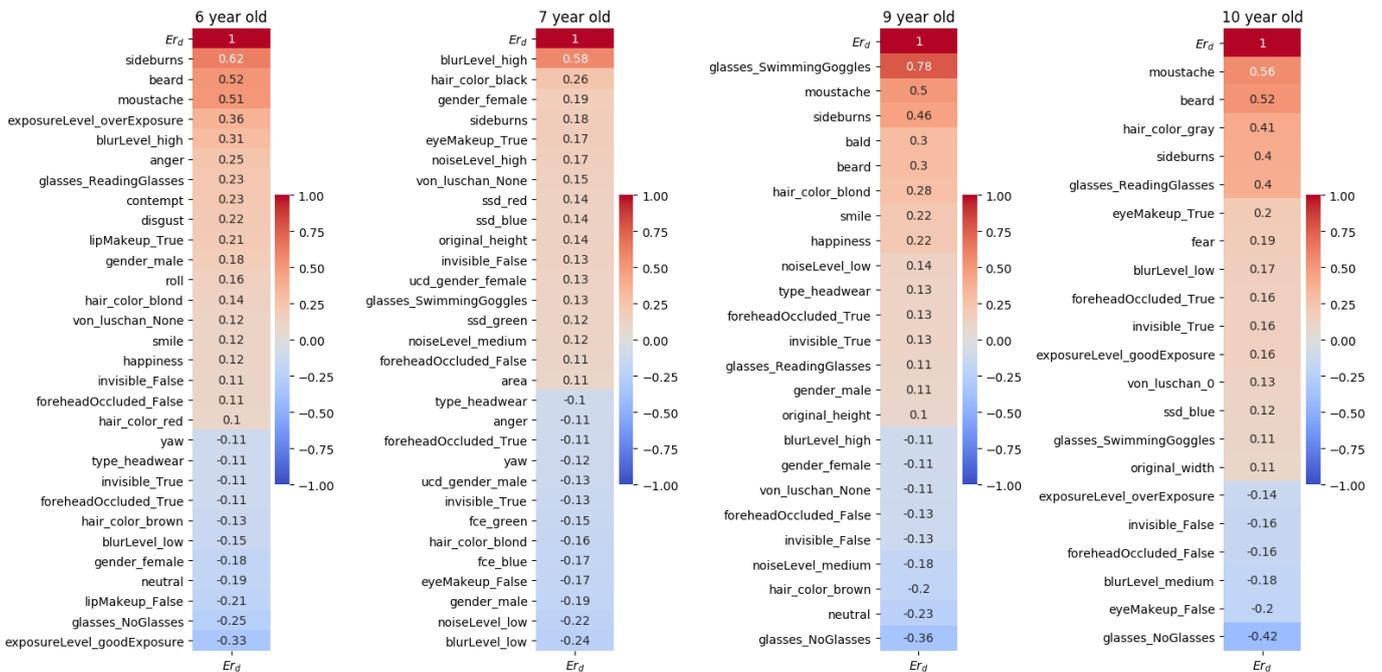


Fig. 4. AZURE: Strong correlations between attributes and $Er_d > 5$ for ages 6, 7, 9 and 10.

age 10 for beard. Another facial hair attribute, moustache, has also been consistently detected across the 6 to 10 age range as shown on Table III. Overall, it can be deduced that misidentification of facial hair has shown prominence in influencing the decline in the facial age estimator's accuracy for the underage age group. Further research is required to identify the underlying cause of these attributes being detected for the underage age group, particularly under 10s. Conversely, age seven did not present with any correlation towards facial hair. Instead, as shown on Figure 4, blurLevel_high was the only strong correlating attribute detected. Moreover, for age nine, along with the correlation to the moustache facial hair, glasses_SwimmingGoggles were also found to have strong correlation to the $Er_d$ with PCC value of 0.78.

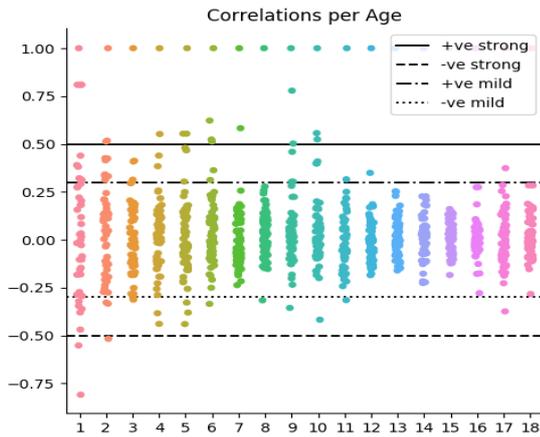| Degree of Correlation | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Attribute Name** | Age | | | | | | | | | | |
| | 1 | 2 | 4 | 5 | 6 | 7 | 9 | 10 | 11 | 12 | 17 |
| exposureLevel_underExposure | 0.81 | | | | | | | | | | |
| exposureLevel_goodExposure | -0.81 | | | | -0.33 | | | | | | |
| exposureLevel_overExposure | | | | | 0.36 | | | | | | |
| noiseLevel_medium | 0.81 | | | | | | | | | | |
| noiseLevel_low | -0.55 | | | | | | | | | | |
| blurLevel_low | | 0.33 | | | | | | | | | |
| blurLevel_medium | | 0.33 | | | | | | | | | |
| blurLevel_high | | | | | 0.31 | 0.58 | | | | | |
| sideburns | 0.81 | | | 0.55 | 0.62 | 0.46 | | 0.4 | | | |
| moustache | 0.81 | | | 0.48 | 0.51 | | 0.5 | 0.56 | | 0.35 | |
| beard | 0.81 | | | 0.55 | 0.52 | 0.3 | | 0.52 | | | |
| bald | | 0.34 | | 0.3 | | 0.3 | | | | | |
| hair_color_brown | 0.38 | | | | | | | | | | |
| hair_color_gray | | | | | | | | 0.41 | | | |
| hair_color_red | -0.36 | | | | | | | | | | |
| fce_red | | | | -0.44 | | | | | | | |
| fce_blue | -0.38 | | | -0.33 | | | | | | | |
| fce_green | -0.35 | | | -0.38 | | | | | | | |
| ucd_gender_female | -0.32 | 0.52 | | | | | | | | | |
| ucd_gender_male | 0.32 | -0.52 | | | | | | | | | |
| contempt | 0.44 | 0.51 | | | | | | | | | |
| anger | 0.31 | | 0.55 | | | | | | | | |
| sadness | 0.39 | | | | | | | | | | |
| fear | | 0.41 | | | | | | | | | |
| disgust | | 0.31 | | | | | | | | | |
| surprise | | 0.43 | | | | | | | | | |
| foreheadOccluded_True | | | 0.38 | | | | | | | | |
| foreheadOccluded_False | | | -0.38 | | | | | | | | |
| invisible_True | | | 0.44 | | | | | | | | |
| glasses_NoGlasses | | | | | | | -0.36 | -0.42 | -0.32 | | -0.37 |
| glasses_SwimmingGoggles | | | | | | | 0.78 | | | | |
| glasses_ReadingGlasses | | | | | | | | 0.4 | 0.32 | | 0.37 |
| glasses_Sunglasses | | 0.42 | | | | | | | | | |
| lipMakeup_True | | 0.33 | | | | | | | | | |
| lipMakeup_False | | -0.33 | | | | | | | | | |



Fig. 5. Azure: Correlations per Age with $Er_d > 5$.

*5) Mild Correlations:* Mild correlations have been defined as PCC values between 0.30 to 0.49. Human biometric factors have been playing both strong and mild roles in influencing the accuracy of the age estimations. In addition to the aforementioned biometric attributes, hair colour and skin tone have been found to have mild correlation with $Er_d > 5$, as shown in Table III. The presence of bald, and brown and grey hair colours on subjects contribute to a higher $Er_d$. The correlation of hair_color_gray with $Er_d$ was expected as the hair colour is often associated with older adult age ranges. Red hair colour, however, was found to have negative correlation value of -0.36 for one year olds. Furthermore, skin tone (as measured by the FCE attribute) have been detected to have mild correlation to the accuracy of the facial age estimation. In particular, it can be observed that presence of any detected level of FCE on a subject linearly correlates to a more accurate age estimation; fce_red, fce_blue and fce_green all demonstrate a negative correlation for ages one and five.

Other biometrics that showed strong correlations have also displayed mild correlation values. Mild correlation results for facial hair, as shown on Table III, are inline with the findings in Section IV-A4. Conversely, a bias towards male subjects
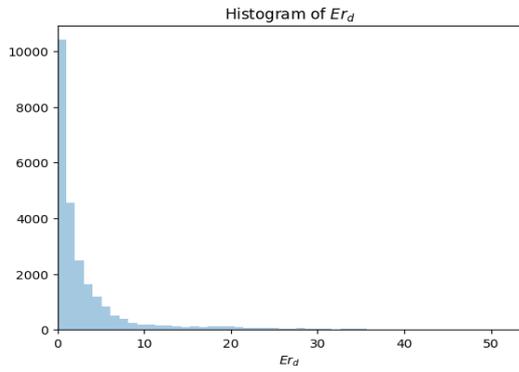
Fig. 6. AWS: Distribution of $Er_d$.



Fig. 7. AWS: Correlations per Age with $Er_d > 5$.

was highlighted in Section IV-A4. Upon analysing the mild correlations, the female gender attribute has a mild negative PCC of -0.32 verses 0.32 for males. Contempt and anger were the two emotional attributes detected to strongly influence the accuracy of age estimation. In addition to these, emotion of sadness, fear, disgust and surprise were also detected to influence the accuracy of age estimation – however, only in a mild manner. In general, the detection of emotion whether with strong or mild correlation, has linear influence in the decrease of the age estimation performance.

Similarly, the same can be said for the quality of image; the higher level of noise and exposure, presence of blur and occlusion all have linear correlations to higher values of $Er_d$. Glasses were predominantly found in the older age range – particularly on ages 9, 10, 11 and 17. Its correlation values imply mild to strong correlation with $Er_d$. Therefore, this is identified as an influencing factor towards Azure's facial age estimation. Moreover, the detection of mild negative correlations of the attribute glasses_NoGlasses substantiates this finding. This was expected as presence of glasses can distort and provide occlusion to a subjects' face. Similar to glasses, the detection of lip makeup has been found to be mildly associated with $Er_d$ with attribute lipMakeup_False substantiating the result through an opposite correlation with equal gravity.

### B. Amazon AWS

In this section we explore the functionality of Amazon's Rekognition service, analyse its age estimation accuracy and identify factors that contributed to our results.

*1) Error Distribution:* Figure 6 shows the error tolerance distribution. The majority of errors had low $Er_d$ (between 0 and 5) signifying that AWS Rekognition's accuracy was within a degree of approximately ±5 for most underage single-faced images processed. A significantly smaller portion of the age estimations had $Er_d \geq 10$.

*2) Strong PCC Distribution per Age with $Er_d$ greater than 5:* Figure 7 illustrates the correlations between the attributes and the AWS $Er_d > 5$. This figure verifies that there are no strong or mild linear correlations between attributes, as shown
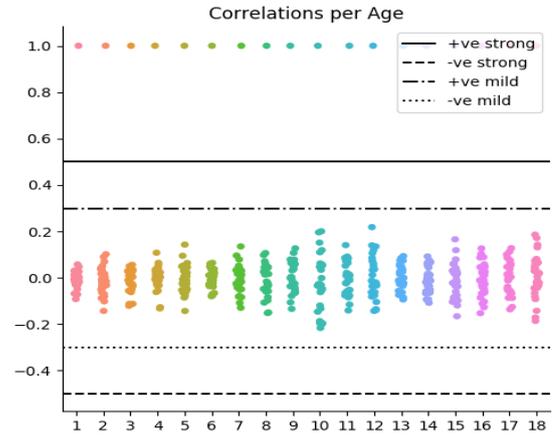
in Table II. While there are a variety of attributes found to have weak associations with $Er_d > 5$, there are no strong influencing factors that affect the AWS accuracy when the error margin is greater than 5, as shown in Table II. This investigation was replicated for $Er_d \geqslant 0$ and $Er_d \leqslant 5$ inline with the investigation process used for Azure. A similar result with $Er_d > 5$ obtained for all other values of $Er_d$. There were no strong correlations identified. Therefore, from the conclusive results obtained for AWS Rekognition, it can be concluded that there are no influencing factors that contribute to the magnitude of its facial age estimation accuracy. Baring in mind that these correlation results are based on PCC, mild to strong nonlinear correlation may still exist. Further study is required to investigate potential nonlinear correlations.

### V. CONCLUDING REMARKS AND DISCUSSION

For Microsoft Azure's Face API, it can be concluded that when the predicted age is close to the ground-truth age, no single attribute was found having prominent association with high level of errors (error difference of 5) or high accuracy (error difference of 1). The majority of strong correlations of 0.5 and greater was only found between the ages 1 to 10. A small number of factors were found to influence the $Er_d$, such as the quality of the image, i.e., a good exposure level and a low level of noise. Additionally, Azure was also found to have higher accuracy when processing male subjects in comparison with females. A total of 0.52 linear correlation was found between the attribute ucd_gender_female and $Er_d$, whereas a negative correlation of equal magnitude was found for ucd_gender_female for age 2. These factors were only noted within the lower limit of the ages evaluated.

Attributes that were found to have strong linear correlation to $Er_d$ can be encapsulated into three main types: quality of image, emotions, and human biometric factors (gender and facial hair). These categories have been identified as the key influencing factors to the accuracy of the tested facial age estimators. While the quality of image impacts the accuracy in the age estimation, emotions are believed to be linked with

facial lines on a subject and therefore can be misinterpreted as wrinkles by the estimator [8]. Furthermore, detection of facial hair and makeup were frequent and often associated with having mild correlation to $Er_d$. It was further found that subjects detected with facial hair were due to them wearing fake moustaches, beards or having food on their face. Eye and lip makeup was also misclassified as present in one year old's. Other biometric factors including hair colour and skin tone (measured by FCE and ssd values) were not identified to have strong influence towards $Er_d$.

Regarding Amazon AWS Rekognition, there were no strong or mild influencing factors that displayed linear correlation with the accuracy of the cloud service.

The distribution of error rates for both AWS and Azure are illustrated in Figures 6 and 1 respectively. The majority of difference between the predicted age and the ground-truth age are relatively low with the majority laying on the $0 \leqslant Er_d \leqslant 5$ for both cloud services. Hence, it can be concluded that their accuracy in underage estimation is relatively high, but that such a MAE may not be accurate enough for some specific law enforcement use cases.

*A. Future Work*

Amazon AWS's and Azure's image classification for facial hair and makeup attributes can be improved. Further investigation can be conducted on the identification and segregation of negative influencing factors, as highlighted in this paper. Exploring the effects of isolating negative influencing factors and the inclusion of only positive influencing factors has an impact on the accuracy of underage facial age estimation. Next, linear correlations between the Amazon AWS Rekognition facial feature detector and the $Er_d$ were predominantly poor. Acknowledging that the coefficient values obtained were based on Pearson's linear approach, it must be considered that a potential strong correlation may exist between the two variables non-linearly. As a result, future work is to explore with nonlinear correlation [25]. Finally, the distributions should be evaluated with different datasets and address the question of how to tackle biased datasets.

## REFERENCES

[1] R. Moore, *Cybercrime: Investigating high-technology computer crime.* Routledge, 2010.

[2] M. L. Bourke and S. W. Craun, "Secondary traumatic stress among internet crimes against children task force personnel: Impact, risk factors, and coping strategies," *Sexual Abuse*, vol. 26, no. 6, pp. 586–609, 2014, pMID: 24259539. [Online]. Available: https://doi.org/10.1177/1079063213509411

[3] D. Lillis, B. Becker, T. O'Sullivan, and M. Scanlon, "Current Challenges and Future Research Areas for Digital Forensic Investigation," in *The 11th ADFSL Conference on Digital Forensics, Security and Law (CDFSL 2016).* Daytona Beach, FL, USA: ADFSL, 05 2016, pp. 9–20.

[4] M. Scanlon, "Battling the digital forensic backlog through data deduplication," in *2016 Sixth International Conference on Innovative Computing Technology (INTECH)*, Aug 2016, pp. 10–14.

[5] F. Anda, D. Lillis, A. Kanta, B. A. Becker, E. Bou-Harb, N.-A. Le-Khac, and M. Scanlon, "Improving borderline adulthood facial age estimation through ensemble learning," in *Proceedings of the 14th International Conference on Availability, Reliability and Security.* ACM, 2019. [Online]. Available: https://doi.org/10.1145/3339252.3341491

[6] D. Quick and K.-K. R. Choo, "Impacts of increasing volume of digital forensic data: A survey and future research challenges," *Digital Investigation*, vol. 11, no. 4, pp. 273 – 294, 2014. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1742287614001066

[7] R. Angulu, J. R. Tapamo, and A. O. Adewumi, "Age estimation via face images: a survey," *EURASIP Journal on Image and Video Processing*, vol. 2018, no. 1, p. 42, 2018.

[8] M. C. Voelkle, N. C. Ebner, U. Lindenberger, and M. Riediger, "Let me guess how old you are: Effects of age, gender, and facial expression on perceptions of age." *Psychology and aging*, vol. 27, no. 2, p. 265, 2012.

[9] M. A. Farooque and J. S. Rohankar, "Survey on various noises and techniques for denoising the color image," *International Journal of Application or Innovation in Engineering & Management (IJAIEM)*, vol. 2, no. 11, pp. 217–221, 2013.

[10] A. Dantcheva, C. Chen, and A. Ross, "Can facial cosmetics affect the matching accuracy of face recognition systems?" in *2012 IEEE Fifth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*, Sep. 2012, pp. 391–398.

[11] C. Chen, A. Dantcheva, and A. Ross, "Impact of facial cosmetics on automatic gender and age estimation algorithms," in *2014 International Conference on Computer Vision Theory and Applications (VISAPP)*, vol. 2. IEEE, 2014, pp. 182–190.

[12] X. Wang, V. Ly, G. Lu, and C. Kambhamettu, "Can we minimize the influence due to gender and race in age estimation?" in *2013 12th International Conference on Machine Learning and Applications*, vol. 2, Dec 2013, pp. 309–314.

[13] F. Anda, D. Lillis, N.-A. Le-Khac, and M. Scanlon, "Evaluating automated facial age estimation techniques for digital forensics," in *2018 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2018, pp. 129–139.

[14] A. M. Albert, K. Ricanek Jr, and E. Patterson, "A review of the literature on the aging adult skull and face: Implications for forensic science research and applications," *Forensic science international*, vol. 172, no. 1, pp. 1–9, 2007.

[15] F. Anda, N.-A. Le-Khac, and M. Scanlon, "DeepUAge : Improving Underage Age Estimation Accuracy to Aid CSEM Investigation," *Forensic Science International: Digital Investigation*, vol. 32, 2020. [Online]. Available: https://doi.org/10.1016/j.fsidi.2020.300921

[16] Microsoft, "Azure Cognitive Services Face API Documentation," 2019, [Online; accessed 02-12-2019]. [Online]. Available: https://docs.microsoft.com/en-us/azure/cognitive-services/face/

[17] AWSChris, "Detecting faces in an image," https://github.com/awsdocs/amazon-rekognition-developer-guide/blob/master/doc_source/faces-detect-images.md, 2020.

[18] U. Khan, M. Cheema, and N. Sheikh, "Adaptive video encoding based on skin tone region detection," in *IEEE Students Conference, ISCON'02. Proceedings.*, vol. 1. IEEE, 2002, pp. 129–134.

[19] C. Manders, F. Farbiz, and C. J. Herng, "The effect of linearization of range in skin detection," in *2007 6th International Conference on Information, Communications & Signal Processing.* IEEE, 2007, pp. 1–5.

[20] C. Yao, "Skin tone estimation and segmentation in matlab and opencv," https://github.com/colin-yao/simple-skin-detection, 2018.

[21] E. Buza, A. Akagic, and S. Omanovic, "Skin detection based on image color segmentation with histogram and k-means clustering," in *2017 10th International Conference on Electrical and Electronics Engineering (ELECO)*, Nov 2017, pp. 1181–1186.

[22] D. E. King, "Dlib-ml: A machine learning toolkit," *Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.

[23] C. B. Barber, D. P. Dobkin, and H. Huhdanpaa, "The quickhull algorithm for convex hulls," *ACM Transactions on Mathematical Software (TOMS)*, vol. 22, no. 4, pp. 469–483, 1996.

[24] S. Satpathy, S. K. Pradhan, and S. Mohapatra, "Internet Usage Analysis Using Karl Pearson's Coefficient of Correlation-A Computer Forensic Investigation," *International Journal of Science and Research*, vol. 3, no. 11, pp. 2791–2794, 2014.

[25] J. Hauke and T. Kossowski, "Comparison of values of pearson's and spearman's correlation coefficients on the same sets of data," *Quaestiones geographicae*, vol. 30, no. 2, pp. 87–93, 2011.