# Revisiting Skin Tone Fairness in Dermatological Lesion Classification

Thorsten Kalb[1], Kaisar Kushibar[1], Celia Cintas[2], Karim Lekadir[1], Oliver Diaz[1], and Richard Osuala[1]

[1] Departament de Matemàtiques i Informàtica, Universitat de Barcelona, Spain
richard.osuala@ub.edu
[2] IBM Research Africa, Nairobi, Kenya

**Abstract.** Addressing fairness in lesion classification from dermatological images is crucial due to variations in how skin diseases manifest across skin tones. However, the absence of skin tone labels in public datasets hinders building a fair classifier. To date, such skin tone labels have been estimated prior to fairness analysis in independent studies using the Individual Typology Angle (ITA). Briefly, ITA calculates an angle based on pixels extracted from skin images taking into account the lightness and yellow-blue tints. These angles are then categorised into skin tones that are subsequently used to analyse fairness in skin cancer classification. In this work, we review and compare four ITA-based approaches of skin tone classification on the ISIC18 dataset, a common benchmark for assessing skin cancer classification fairness in the literature. Our analyses reveal a high disagreement among previously published studies demonstrating the risks of ITA-based skin tone estimation methods. Moreover, we investigate the causes of such large discrepancy among these approaches and find that the lack of diversity in the ISIC18 dataset limits its use as a testbed for fairness analysis. Finally, we recommend further research on robust ITA estimation and diverse dataset acquisition with skin tone annotation to facilitate conclusive fairness assessments of artificial intelligence tools in dermatology. Our code is available at https://github.com/tkalbl/RevisitingSkinToneFairness.

**Keywords:** Dermatology · Fairness · Deep Learning · Skin Cancer

## 1 Introduction

Skin cancer is one of the most prevalent cancer types worldwide [23]. According to [2], early detection increases the survival rate to 99% compared to 32% in late stage detection. The 5-year survival rates after surgical removal of melanoma have been shown to be lower for black patients (73%) compared to white (88%), although melanoma is 23 times more prevalent in white patients [8]. Dick et al. [10] found that black patients were significantly more likely to present with advanced-stage disease, even after adjusting for tumour characteristics and demographic factors. Deep learning models have demonstrated a remarkable performance in skin lesion analysis and classification [11]. Therefore, they are promising

tools to detect skin cancer earlier and, thus, in theory, bear the potential to reduce the aforementioned disparities. In practice, however, deep learning models have been shown to be prone to and exacerbate existing societal biases [16,24,6]. Although bias and fairness assessment in skin lesion classification has been an active research area [17,24,4], there is a substantial limitation on developing an unbiased classifier. That is, many publicly available datasets lack information about ethnicity or skin types [1]. Hence, such labels are usually obtained using different automated methods. One of the common approaches is based on Individual Topology Angle (ITA) (see Section 2.3) that have been used in several studies [5,9,14,15,18,17]. For instance, Kinyanjui et al. [15] estimated skin tones to assess their effect on lesion classification performance, while Bevan et al. [5] labelled skin tones for model debiasing. These previous works [15,18,5,17] form the basis of our analysis and are further described in detail in Section 2.3.

These existing studies utilise ITA-based estimation of skin tones as a proxy to ground truth. However, in this work, we show that there is a large disagreement in the assigned skin tones on the same dataset, which suggests that the reported results may be inconclusive. Therefore, we investigate the causes and extent of the discrepancies between the estimated skin tones in previous studies. We uncover the complexities, pitfalls, and differences across ITA-based skin tone estimation techniques that question the conclusions derived in previous studies. In summary, our contributions are as follows:

- We compare different ITA-based automatic skin tone estimation algorithms and highlight common pitfalls.
- We demonstrate the impact of different skin tone estimations on the outcome of fairness analyses of the same model.
- We show the limitations of the commonly used ISIC18 dataset [22] to assess fairness of skin lesion classifiers.

## 2    Methods and Materials

### 2.1    Dataset

The ISIC18 dataset [22] includes 10015 dermatoscopic images from Austria and Australia. The present work focuses on classification fairness of seven skin lesion types distributed as: 1113 melanoma (MEL), 6705 melanocytic nevus (NV), 514 basal cell carcinoma (BCC), 327 actinic keratosis (AKIEC), 1099 benign keratosis (BKL), 115 dermatofibroma (DF), and 142 vascular lesions (VASC).

### 2.2    Evaluation of skin lesion classification

To build a skin lesion classification model, we used a MobileNetV2 [21] initialised with ImageNet weights [12]. Only the last layer was replaced to correspond with the seven skin lesion classes. The choice of this model was driven by its universal applicability [16] due to its lightweight architecture that facilitates deployment across a wide range of devices including portable devices in low-resource settings.
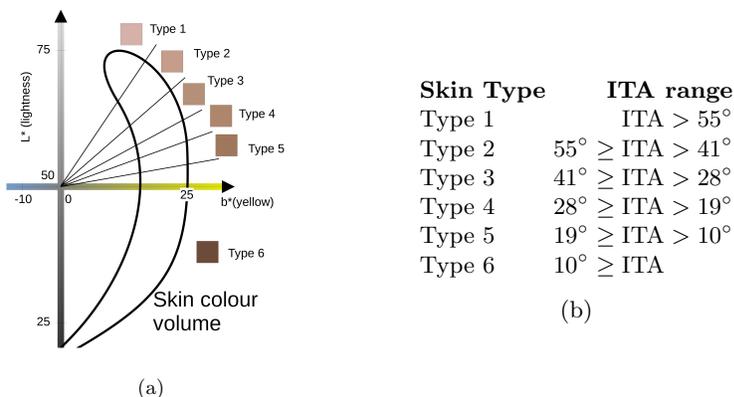
| Skin Type | ITA range |
|-----------|-----------|
| Type 1 | ITA $> 55°$ |
| Type 2 | $55° \geq$ ITA $> 41°$ |
| Type 3 | $41° \geq$ ITA $> 28°$ |
| Type 4 | $28° \geq$ ITA $> 19°$ |
| Type 5 | $19° \geq$ ITA $> 10°$ |
| Type 6 | $10° \geq$ ITA |

(b)

(a)

Fig. 1: (a) Skin colour volume in the $L*$ - $b*$ plane of CIELab colour space with ITA thresholds [5]. (b) Skin types from ITA thresholds, as in [5].

A grid search was performed for a combination of (a) a constant learning rate ranging from 1e-3 to 1e-6 in steps of 10 on a logarithmic scale and (b) the batch size, where 12, 24, 32, 48, and 64 were tested. Based on validation loss, the grid search yielded a learning rate of 1e-5 and a batch size of 16. The network was optimised using Adam to minimise the sparse categorical cross entropy loss for a maximum of 60 epochs with an early stopping policy to terminate training if the validation loss did not improve for 30 epochs. All experiments were performed on an NVIDIA GeForce RTX 2080 SUPER with 8GiB memory. All evaluations were conducted three times with different seeds for splitting into training (57%), validation (14%) and test (29%) sets, stratified by lesion, with a slightly different split for the RP data shift experiment (see Section 3.3).

### 2.3  Skin tone estimation

**Individual Topology Angle (ITA)**  In the absence of dermatologist-confirmed skin type labels, researchers proposed automatic skin tone and skin type estimations [15,18,14,9,17,5], which are commonly based on the ITA [7]. The ITA is defined within the L*-b*-plane of the CIELab colour space, according to Equation 1, and illustrated in Figure 1a.

$$\text{ITA} = \arctan\left(\frac{L - 50}{b}\right) \cdot \frac{180°}{\pi} \qquad (1)$$

ITA values from different studies only become comparable when the lighting conditions and measurement devices are known and corrected for [19]. Categorical skin types are obtained by binning ITA values [15,14,19] as shown in Figure 1b. We note that there is no consensus for ITA binning thresholds and high uncertainty for any ITA (colour) to Fitzpatrick [13] (sun reactivity) mapping. Given our analysis of four existing ITA-based automatic skin tone estimation methods [15,18,5,17], we note the following key issues that such methods need to address.

**I.1** *Lighting conditions*: The ITA is sensitive to illumination, especially to brightness or lack of yellow chroma.

**I.2** *Non-skin imaging artefacts*: The ITA is defined for any colour, but only meaningful for skin. Estimating ITA, in part, for hair, artefacts and dark borders can create misleading results.

**I.3** *Lesion to skin contrast*: The pigmented lesion is not representative for the skin colour of the patient and needs to be excluded.

**I.4** *From pixel to image-level*: The ITA is defined for each pixel. As one representative ITA value needs to be assigned to an image, there is no consensus how to address variance and outliers of the ITA distribution.

In the following, we describe the aforementioned four existing skin tone estimation methods [15,18,5,17], which are empirically analysed in Section 3.

**Method 1: Deep learning-based skin segmentation** We adopt a skin segmentation model kindly provided by the authors of [15], which is a Mask R-CNN trained with manually segmented ISIC18 images to address I.2 and I.3. The segmented pixels of healthy skin are converted to CIELab colour space, for which the median within one standard deviation of the mean of ITA's $L*$ and $b*$ components is calculated. Based on these $L*$ and $b*$ median values, an image's ITA is calculated according to Equation 1. Selecting $L*$ and $b*$ median values separately can help to avoid outliers, but may also lead to less precise ITAs, as $L*$ and $b*$ medians likely do not correspond to the same pixels.

**Method 2: Colour-based skin segmentation** This method follows the skin segmentation algorithm proposed in [18]. Input images are converted to grayscale before applying Otsu binarisation and thresholding [20] to detect pixels that are non-lesion. For original values of healthy skin pixels, different thresholds in HSV and YCrCb colour spaces are applied to define potential skin colours. For the pixels within these thresholds, the mean values for red, green and blue are computed and define a representative skin colour. This skin colour is then converted into CIELab space before applying Equation 1 to calculate the ITA.

**Method 3: Random patch algorithms** Next method is based on semi-random patches proposed in [5]. It is assumed that a lesion is in the centre of the image and healthy skin is presumably found in at least one of eight patches of $20 \times 20$ pixels in the periphery. Before patch extraction, input images are centre-cropped, resized and small dark artefacts such as hair are removed via black-head morphology to address I.2. For each patch, the mean ITA is calculated. The ITA of the patch is selected that corresponds to the brightest skin type, which arguably has the effect of having excluded the pigmented lesion. In contrast to previous methods, this method used *arctan* in Equation 1 for the ITA instead of *arctan2*. Unlike *arctan*, *arctan2* takes into account the signs of the catheti as described in Equation 2.

$$\text{arctan2}(x, y) = \begin{cases} \arctan(y/x), & \text{if } x > 0 \\ \arctan(y/x) - \text{sgn}(y) * \pi, & \text{if } x < 0 \\ \text{sgn}(y) * \pi & \text{if } x = 0 \end{cases} \qquad (2)$$

Using *arctan* in Equation 1 assumes $b*$ cannot be negative. However, a negative $b*$ value encodes blue or absence of yellow chroma. Although the blue colour in skin is not intuitive, its appearance may depend on variations in illumination, dermatoscope, and camera. Therefore, in our analysis, we include both versions of ITA estimation using *arctan* and *arctan2* referred to as RP and RP2, respectively.

**Method 4: Histogram thresholding with grey-world white balancing**
The next method is adopted from [17]. After centre-cropping and resizing, a grey-world white balance algorithm is applied to correct for light influence addressing I.1. Next, non-lesion skin is segmented using the Generalized Histogram Threshold [3] algorithm. The segmented area is transformed from RGB to CIELab-space to calculate the ITA. However, in [17], reproducibility was limited regarding how one representative ITA is obtained from the multiple segmented pixels. A further limitation is the assumption that skin images are grey on average.

## 3   Experiments and Results

### 3.1   Comparison of ITA estimation methods

A comparison of the estimated skin tone distributions is shown in Figure 2. Although the same thresholds (see Figure 1b) are applied to all four skin tone estimation methods for ITA binning, their estimates on the ISIC18 dataset clearly differ. A detailed comparison of the agreement among the ITA estimation methods is shown in Figure 3, where a diagonal matrix would represent perfect agreement. There is little agreement between RP and DLHSS, especially for dark skin. Most type 6 images in RP appear as type 2 in DLHSS and most type 2 images in DLHSS are classified as type 1 in RP. All type 5 and 6 images in DLHSS are labelled as type 1 in RP. Comparing the RP to RP2, the entries of the matrix lay on the diagonal and in the first row. This implies that in most cases, RP and RP2 agree on the same skin type. However, approximately 21.6% of the samples are classified as type 1 by RP2 and darker by RP. Thus, using the *arctan* can account for the over-estimation of dark skin images and affect all except for type 1. Comparing DLHSS and RP2, both agree that approximately 99.0% of the samples show skin at least as light as type 3 (ITA $> 28°$), while only eight out of 10015 samples are classified at least as dark as type 4 (ITA $\leq 28°$).

These eight samples are shown in row 1 of Figure 4 and indicate that images labelled as dark skin are actually dark images of skin. We qualitatively evaluated these images with a dermatologist who confirmed that they likely correspond to Fitzpatrick (FST) [13] skin type III or lower. When RP2 labels images as
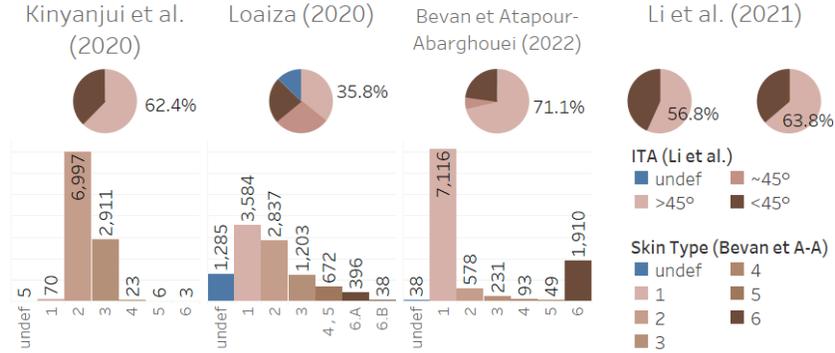
Fig. 2: Comparison of the four different ITA-based skin type estimation methods based on values reported in the literature [15,18,5,17]. Note that Li et al. [17] only reported darker (ITA $< 45°$, 43.2%) and lighter (ITA $> 45°$, 63.8%) skin and that their sum equals 107%; it is assumed that one of these values is correct and both possible distributions are shown.

dark skin and DLHSS does not, this can be explained by RP2 segmenting the lesion in all eight patches. Since RP2 reports the brightest patch and DLHSS a median, DLHSS is likely more reliable in these cases. On visual inspection of the five darkest images (DLHSS), this algorithm appears to be susceptible to hair and rare lesion sites, such as the tongue or ears; and dark skin labels can be explained by segmentation failures of DLHSS. These results suggest that the ISIC18 dataset is not sufficiently diverse for a fairness analysis, as it presumably does not contain any images of dark skin (i.e. FST IV-VI [13]).

### 3.2   Fairness analysis

The question arises as to whether and how the different ITA estimation methods impact the result of a downstream fairness analysis. To this end, we analyse the balanced accuracy per skin type for the baseline experiment without data shift for DLHSS, RP, and RP2. As shown in Figure 5, we obtain different fairness results per method despite using the exact same classification model and test set. For DLHSS ITA values, we note a decline in average balanced accuracy for types 1-4 (according to Figure 1b). In contrast, the balanced accuracy of very light skin (type 1) and very dark skin (type 6) according to RP is lower than for intermediate skin tones (types 2-5). According to RP2, the light skin samples (type 1) show the worst performance, and the balanced accuracy appears to increase overall for darker skin types. Thus, analysing the fairness of the same lesion predictions, the outcome changes depending on the applied ITA estimation method.

|  |  | DLHSS |  |  |  |  |  |  | RP (reproduced) |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | Type 1 | Type 2 | Type 3 | Type 4 | Type 5 | Type 6 | undef | Type 1 | Type 2 | Type 3 | Type 4 | Type 5 | Type 6 |
| RP (Original) | Type 1 | 32 | 5,225 | 1,836 | 10 | 5 | 5 | 3 | 7,080 | 9 | 2 | 1 | 2 | 22 |
|  | Type 2 | 1 | 86 | 489 | 2 |  |  |  | 18 | 559 | 1 |  |  |  |
|  | Type 3 | 2 | 28 | 200 | 1 |  |  |  |  | 9 | 219 | 1 | 1 | 1 |
|  | Type 4 | 1 | 26 | 66 |  |  |  |  |  |  | 4 | 83 | 6 |  |
|  | Type 5 |  | 14 | 32 | 3 |  |  |  |  |  |  | 3 | 42 | 4 |
|  | Type 6 | 36 | 1,606 | 258 | 8 |  |  | 2 | 4 |  |  |  | 2 | 1,904 |
|  | undef |  | 26 | 12 |  |  |  |  | 35 |  |  |  |  | 3 |
| RP2 | Type 1 | 72 | 6,963 | 2,209 | 14 | 5 | 5 | 5 | 7,136 | 80 | 52 | 42 | 37 | 1,926 |
|  | Type 2 |  | 45 | 451 | 1 |  |  |  |  | 497 |  |  |  |  |
|  | Type 3 |  | 1 | 172 | 1 |  |  |  |  |  | 174 |  |  |  |
|  | Type 4 |  |  | 46 |  |  |  |  |  |  |  | 46 |  |  |
|  | Type 5 |  | 1 | 13 | 2 |  |  |  |  |  |  |  | 16 |  |
|  | Type 6 |  | 1 | 2 | 6 |  |  |  | 1 |  |  |  |  | 8 |

Fig. 3: Comparison of ITA estimates of Deep Learning-based Healthy Skin Segmentation (DLHSS) [15], Random Patch [5] (RP) and Random Patch with *arctan2* (RP2).
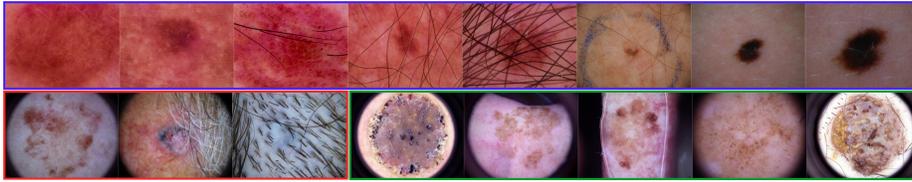


Fig. 4: Blue and red: All samples where both DLHSS and RP estimate dark skin tones (types 4-6, ITA $\leq 28°$). Blue: Dark skin according to DLHSS and RP2. Green: Samples with darkest skin according to DLHSS (type 6, ITA $\leq 10°$).

### 3.3　Simulated data shifts

To simulate data shift, light skin images are used during training (ITA $> 41°$) while dark skin images are used in testing (ITA $\leq 41°$). To assess whether the ITA estimation method impacts classification fairness in the presence of data shift, reproduced DLHSS and original RP ITA values are used. For comparison, the size of the baseline test set was defined equal to the DLHSS test set size and the train-validation ratio was 80:20 for all experiments. The results of the data shift experiments are visualised in Figure 6. Despite a higher baseline accuracy, the weighted precision, recall, and f1-score, remain similar between the baseline and the data shift experiment with DLHSS labels. The balanced accuracy is the *unweighted* macro-averaged recall, hence the difference in balanced accuracy and *weighted* (macro-averaged) recall suggests that the classification performance differs per lesion type. This difference appears to depend on the train-test split induced by the ITA estimation method, yielding different lesion distributions. Thus, not only the choice of evaluation metric, and the lesion type distribution, but also the ITA estimations can alter the conclusions of the fairness analysis.
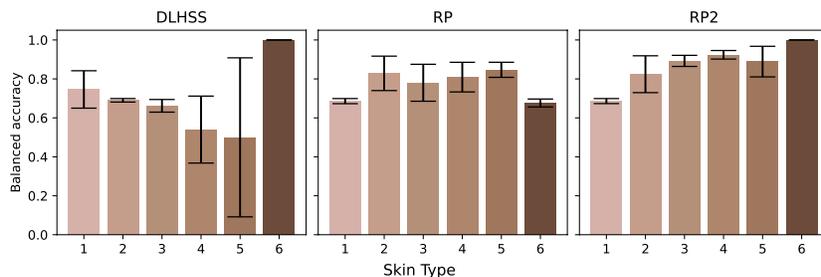
Fig. 5: Average balanced accuracy per skin type in the baseline experiment according to different automatic ITA estimation algorithms: Deep Learning-based Healthy Skin Segmentation (DLHSS), Random Patch algorithm with *arctan* (RP) and with *arctan2* (RP2). All bar charts stem from the same predictions.
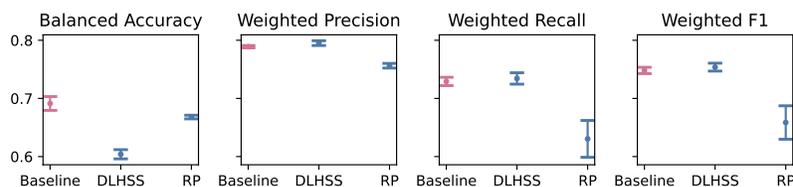


Fig. 6: Metrics for baseline and data shift experiments with the test set of the latter containing exclusively skin types 3-6 (ITA $\leq 41$ according to Figure 1b). ITAs are obtained from DLHSS and RP estimation methods. Error bars show the standard deviation between three experiment repetitions with different random train-validation (DLHSS, RP) or train-validation-test (Baseline) splits.

## 4   Conclusions

We compared ITA-based skin tone estimation methods that revealed common pitfalls. Namely, susceptibility to lighting conditions, colour space calibration, presence of hair or dark edges. Moreover, we showed the effects of the differences in extracting the healthy skin, and differences in mapping ITA values from pixel to image-level. Further, we observed disagreements between ITA estimation methods, which did not necessarily refer to the same images as dark, and overestimation of dark samples in the ISIC18 dataset. A qualitative analysis with a dermatologist revealed that the images, where different ITA estimation methods agree on their dark skin tones (ITA $\leq 28°$), do not represent moderate brown to black skin types (FST IV-VI). Furthermore, our skin lesion classification experiments demonstrated that the choice of ITA estimation method substantially impacts the results of classification fairness analyses. Our data shift experiments further confirmed that ITA estimation altered conclusions drawn from the fairness analysis. We illustrated the need for more diverse dermatological datasets with diligently annotated skin tones, lighting conditions, dermatoscope

and camera information. Apart from improving ITA estimation, further avenues of research include the measurement of the differences in model calibration and in epistemic and aleatoric uncertainty per skin tone type. In the presence of limited datasets, synthetic samples with controllable skin tones and lesion types can be explored to quantify their effect on fairness. Overall, our work shows that current skin tone fairness assessments are inconclusive and further research is needed into unified and robust algorithms for automatic skin tone estimation to avoid carrying unwanted biases into dermatology practice when deploying deep learning models.

# References

1. Alipour, N., Burke, T., Courtney, J.: Skin Type Diversity: a Case Study in Skin Lesion Datasets (Jul 2023). https://doi.org/10.21203/rs.3.rs-3160120/v1
2. American Cancer Society: Cancer Facts & Figures 2023 (2023), https://www.cancer.org/content/dam/cancer-org/research/cancer-facts-and-statistics/annual-cancer-facts-and-figures/2023/2023-cancer-facts-and-figures.pdf
3. Barron, J.T.: A generalization of otsu's method and minimum error thresholding. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (eds.) Computer Vision – ECCV 2020. pp. 455–470. Springer International Publishing, Cham (2020)
4. Bevan, P.J., Atapour-Abarghouei, A.: Skin Deep Unlearning: Artefact and Instrument Debiasing in the Context of Melanoma Classification (2021). https://doi.org/10.48550/ARXIV.2109.09818
5. Bevan, P.J., Atapour-Abarghouei, A.: Detecting Melanoma Fairly: Skin Tone Detection and Debiasing for Skin Lesion Classification. In: Kamnitsas, K., Koch, L., Islam, M., Xu, Z., Cardoso, J., Dou, Q., Rieke, N., Tsaftaris, S. (eds.) Domain Adaptation and Representation Transfer. pp. 1–11. Springer Nature Switzerland, Cham (2022)
6. Birhane, A., Prabhu, V., Han, S., Boddeti, V.N.: On hate scaling laws for data-swamps. arXiv preprint arXiv:2306.13141 (2023)
7. Chardon, A., Cretois, I., Horseau, C.: Skin colour typology and suntanning pathways. International Journal of Cosmetic Science **13**(4), 191–208 (1991). https://doi.org/10.1111/j.1467-249 4.1991.tb00561.x
8. Collins, K.K., Fields, R.C., Baptiste, D., Liu, Y., Moley, J., Jeffe, D.B.: Racial differences in survival after surgical treatment for melanoma. Annals of Surgical Oncology **18**(10), 2925–2936 (Apr 2011). https://doi.org/10.1245/s10434-011-1706-3

9. Corbin, A., Marques, O.: Exploring strategies to generate fitzpatrick skin type metadata for dermoscopic images using individual typology angle techniques. Multimedia Tools and Applications (Nov 2022). https://doi.org/10.1007/s11042-022-14211-1

10. Dick, M., Aurit, S., Silberstein, P.: The odds of stage iv melanoma diagnoses based on socioeconomic factors. Journal of cutaneous medicine and surgery **23**(4), 421–427 (2019)

11. Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., Thrun, S.: Dermatologist-level classification of skin cancer with deep neural networks. nature **542**(7639), 115–118 (2017)

12. Fei-Fei, L., Deng, J., Li, K.: ImageNet: Constructing a large-scale image database. Journal of Vision **9**(8), 1037–1037 (Mar 2010). https://doi.org/10.1167/9.8.1037

13. Fitzpatrick, T.B.: The Validity and Practicality of Sun-Reactive Skin Types I Through VI. Archives of Dermatology **124**(6), 869–871 (06 1988). https://doi.org/10.1001/archderm.1988.01670060015008

14. Groh, M., Harris, C., Soenksen, L., Lau, F., Han, R., Kim, A., Koochek, A., Badri, O.: Evaluating Deep Neural Networks Trained on Clinical Images in Dermatology with the Fitzpatrick 17k Dataset (2021)

15. Kinyanjui, N.M., Odonga, T., Cintas, C., Codella, N.C., Panda, R., Sattigeri, P., Varshney, K.R.: Fairness of classifiers across skin tones in dermatology. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 320–329. Springer, Springer International Publishing (2020)

16. Lekadir, K., Osuala, R., Gallin, C., Lazrak, N., Kushibar, K., Tsakou, G., Aussó, S., Alberich, L.C., Marias, K., Tsiknakis, M., et al.: FUTURE-AI: guiding principles and consensus recommendations for trustworthy artificial intelligence in medical imaging. arXiv preprint arXiv:2109.09658 (2021). https://doi.org/10.48550/arXiv.2109.09658

17. Li, X., Cui, Z., Wu, Y., Gu, L., Harada, T.: Estimating and Improving Fairness with Adversarial Learning. arXiv (2021). https://doi.org/10.48550/arXiv.2103.04243

18. Loaiza, K.: The skin tone problem in artificial intelligence. In: 1st Congress of Women in Bioinformatics and Data Science Latin America (09 2020). https://doi.org/10.13140/RG.2.2.20564.63361/1

19. Ly, B.C.K., Dyer, E.B., Feig, J.L., Chien, A.L., Bino, S.D.: Research techniques made simple: Cutaneous colorimetry: A reliable technique for objective skin color measurement. Journal of Investigative Dermatology **140**(1), 3–12.e1 (Jan 2020). https://doi.org/10.1016/j. jid.2019.11.003

20. Otsu, N.: A threshold selection method from gray-level histograms. IEEE transactions on systems, man, and cybernetics **9**(1), 62–66 (1979)

21. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: MobileNetV2: Inverted residuals and linear bottlenecks. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. IEEE (Jun 2018). https://doi.org/10.1109/cvpr.2018.00474

22. Tschandl, P., Rosendahl, C., Kittler, H.: The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. Scientific Data **5**(1) (08 2018). https://doi.org/10.1038/sdata.2018.161

23. World Health Organization: Skin cancer - IARC (2023), `https://www.iarc.who.int/cancer-type/skin-cancer/`, last accessed 28 July 2023

24. Wu, Y., Zeng, D., Xu, X., Shi, Y., Hu, J.: FairPrune: Achieving Fairness Through Pruning for Dermatological Disease Diagnosis. In: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (eds.) Medical Image Computing and Computer Assisted

Intervention – MICCAI 2022. pp. 743–753. Springer, Springer Nature Switzerland, Cham (2022)