

# Chapter 1

## Blinding with Linear Clustering Removal

Face recognition algorithms work with embedding spaces. They map images of persons into the embedding space such that images of the same person are close to each other in the embedding space. This work investigates discriminatory dimensions in face recognition algorithms. For a given discriminatory dimension, the data can be grouped into clusters. A blinding procedure is proposed to remove the information related to the separation of these clusters. The procedure is a linear operation in the embedding space and uses the following steps:

1. Compute centers of clusters defined by the discriminatory dimension.
2. Use a one-vs-rest (OvR) Ansatz to calculate the directions of discrimination of each cluster relative to the other clusters.
3. Apply singular value decomposition (SVD) on the directions of discrimination to find an orthonormal basis spanning the “discriminatory subspace”.
4. Remove projections onto the “discriminatory subspace” from the embedding vectors. This results in embedding vectors which are orthogonal to the directions of discrimination.

After outlining the method, cluster visualization, awareness and face recognition rates are investigated before and after the blinding procedure.

### 1.1 The math behind

In the following we look at the discriminatory dimension of race. We work with the commonly used racial faces in-the-wild (RFW) data set which groups faces into  $K = 4$  ethnic clusters Caucasian, African, Asian and Indian. I consider a VGG2 model where the embedding space has  $N_e = 128$  dimensions. The procedure outlined above operates on the embedding vectors  $\mathbf{x}_i$  where  $i$  denotes the

sample. Associated to each sample is a cluster label  $k \in \{1, \dots, K\}$ . As stated, the goal is to remove the directions in the embedding space which separate the ethnic clusters. As a first step, we define the centers of each cluster by the average

$$\bar{\mathbf{x}}_k = \frac{1}{n_k} \sum_{i \in C_k} \mathbf{x}_i, \quad (1.1)$$

where  $C_k$  is the set of embedding vectors associated with cluster  $k$  and  $n_k$  is the corresponding size. Following a one-vs-rest (OvR) approach, the normalized direction of discrimination of each cluster  $k$  to the other clusters is given by the vectors

$$\mathbf{u}_k = \frac{\mathbf{v}_k}{\|\mathbf{v}_k\|} \quad \text{with} \quad \mathbf{v}_k = \bar{\mathbf{x}}_k - \frac{1}{K-1} \sum_{k' \neq k} \bar{\mathbf{x}}_{k'}, \quad (1.2)$$

where  $K$  is the number of clusters. Fig. 1.1 shows the projections onto the vectors  $\mathbf{u}_k$ . As it may be expected from the construction, each direction nicely separates the corresponding cluster from the others. These means, that different ethnic groups are literally located in different corners of the embedding space. By construction, the vectors  $\mathbf{u}_k$  are not linearly independent, but span a subspace of rank  $K-1$ . This can be verified by applying a singular value decomposition (SVD) on the matrix  $U = [\mathbf{u}_1 \dots \mathbf{u}_K]$ . SVD also provides an orthonormal basis  $B = [\mathbf{e}_1 \dots \mathbf{e}_{K-1}]$  of the corresponding subspace. The final step is to remove the projections onto this subspace by

$$\mathbf{x}_i^b = \mathbf{x}_i - \sum_{j=1}^{K-1} (\mathbf{x}_i \cdot \mathbf{e}_j) \mathbf{e}_j, \quad (1.3)$$

where  $(\mathbf{x}_i \cdot \mathbf{e}_j)$  is the dot (or scalar) product. Eq. (1.3) yields new embedding vectors  $\mathbf{x}_i^b$  with the same shape as the original ones. The upper index  $b$  stands for *blinded* inspired by the fact that some information with regard to the discriminatory dimension has been removed. Note that the new embeddings depend linearly on the original ones.

## 1.2 Awareness

Awareness is the ability of the model to discriminate between different clusters of the discriminatory dimensions, being the ethnic label in present case. This ability obviously depends on the trained model at hand. Here I benchmark the performance to predict the ethnic labels by the *aware* embeddings  $\mathbf{x}_i$  and the *blinded* ones  $\mathbf{x}_i^b$ . A train/test split of two thirds/one third was used. The accuracy of different classifiers is shown in Tab. 1.1. Not surprisingly, linear classifiers are unable to predict the race for the blinded embeddings. Nearest neighbor approaches still work reasonably. More advanced non-linear classifiers

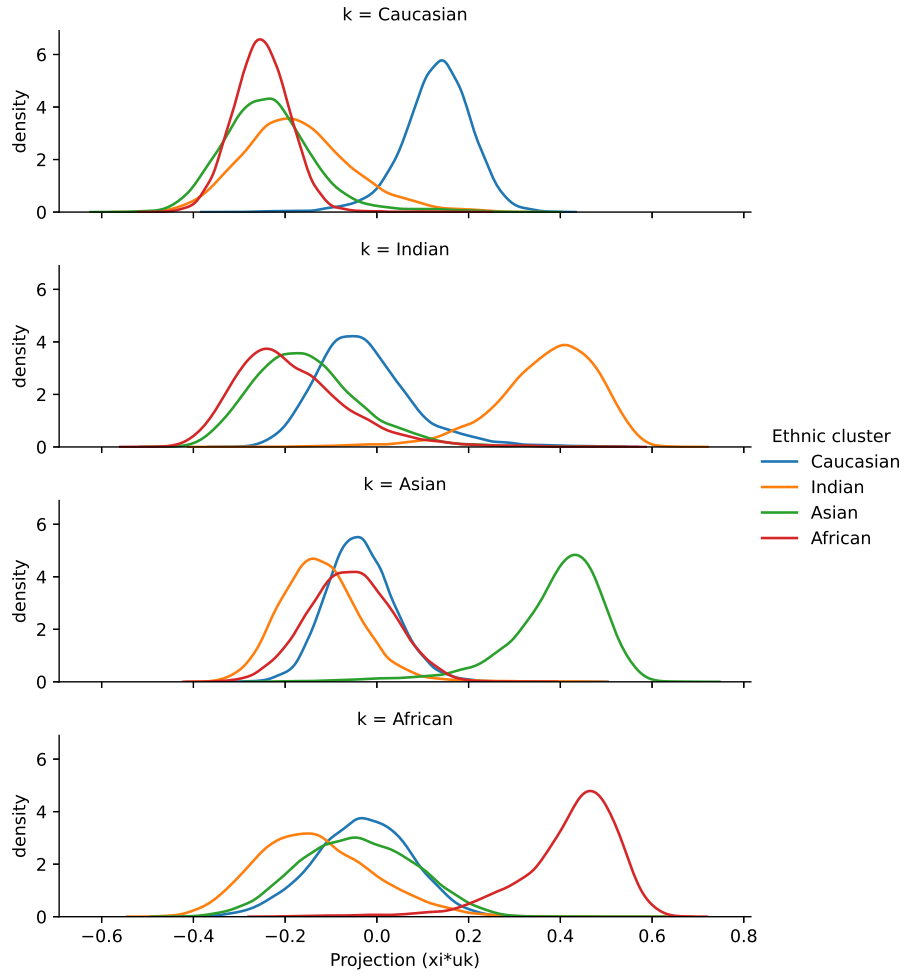


Figure 1.1: Kernel density estimation (KDE) plot of the projections  $(\mathbf{x}_i \cdot \mathbf{u}_k)$ . The normalized directions  $\mathbf{u}_k$  represent the discriminatory directions which separate each cluster from the others.

Model	aware	blinded
Logistic regression	96%	21%
Linear SVM	96%	26%
Nearest neighbor	92%	57%
5 Nearest neighbor	94%	62%
NN with 1 hidden layer (100 nodes), relu	96%	70%
NN with 2 hidden layer (100 nodes each), relu	96%	85%

Table 1.1: Subset accuracy of various classifiers predicting the ethnic labels based on *aware* and *blinded* emdeddings. A train/test split of two thirds/one third was used.

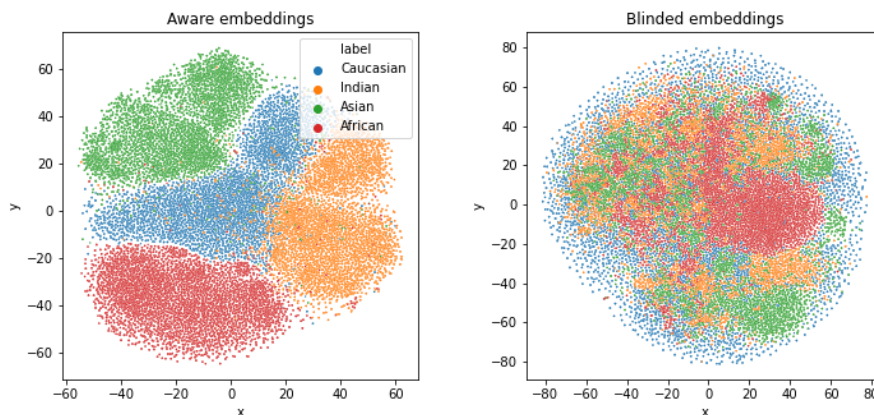


Figure 1.2: tSNE plots of the two embeddings

such as neural networks perform well. The clustering displayed by the corresponding t-SNE plots in Fig. 1.2 is in line with these findings. In the blinded case, clusters can't be separated by a single straight line. However, the data still displays groups defined by ethnic labels. Interestingly, there are clear differences between the groups. Africans are grouped in the center, Caucasians encircle the cloud and Indians/Asians are scattered inbetween.

### 1.3 Cluster scores

Cluster scores give a measure of clustering. They are calculated for both embeddings in Tab. 1.2. The Silhouette cluster score gives a measure between 0 and 1 indicating how well the data is clustered. The Silhouette cluster score of the aware embeddings is only 0.063 - basically indicating the absence of clustering although Fig. 1.1 and Fig. 1.2 show nice clustering for the aware embeddings. This counter-intuitive finding is due to the high dimensionality. Both figures

Cluster score	aware	blinded
Silhouette score	0.063	-0.014
Calinski-Harabasz score	1814	0
Davies-Bouldin score	3.8	2.8x10 <sup>6</sup>

Table 1.2: Cluster scores for *aware* and *blinded* emdeddings.

show projections into lower dimensions and therefore reflect only a marginal part of the information. This is confirmed by the fact that the total variance of the blinded embeddings is still 84% of the original variance.

## 1.4 Face recognition rates and bias

### 1.4.1 Positive / negative pair metric

Face recognition rates and bias are evaluated with the RFW data set. The RFW dataset provides image (i.e. embedding) pairs corresponding to the same or to different persons. The resulting task is a binary classification of the pairs into “same” and “different”. Note that pairs are withing the ethnic group (which turns out to be critical, as shown in the next subsection). The recognition rate is the accuracy of the corresponding classification. The feature used for the classification is the pair distance in the embedding space. Here we use the cosine distance:

$$d_{ij} = 1 - \frac{\mathbf{x}_i \cdot \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \quad (1.4)$$

The face recognition rates calculated in this way are shown in Tab. 1.3. The table includes further Senet models with  $N_e = 256, 2048$ . Surprisingly, the performance increases for the blinded embeddings by about 2% for all clusters. Bias is slightly removed. Fig. 1.3 gives further insights by showing the distribution of the cosine distances for “same” and “different” pairs. The Caucasians are special in that their distributions are not significantly altered. Note that the blinding procedure leads to a better alignment of the thresholds.

### 1.4.2 Nearest neighbor metric

Given a set of images, an alternative way to assess bias is to look at the nearest neighbors and whether the corresponding images belong to the same person or not. This can again be done for the *aware* and *blinded* emdeddings. Again, the cosine distance is used as defined in Eq. (1.4). The result is shown in Tab. 1.4 together with error rates which show whether nearest neighbors belong to the same ethnic group or to a different one (if it is not the same person). With this metric, the blinding reduces performance. Although the confusion with persons of the same ethnic group is reduced upon blinding, confusion with

$N_e$	128	128	256	256	2048	2048
	aware	blinded	aware	blinded	aware	blinded
Total	86%	88%	86%	88%	83%	84%
Caucasian	91%	92%	91%	92%	89%	89%
Indian	86%	88%	86%	88%	85%	85%
Asian	84%	86%	84%	87%	82%	82%
African	84%	86%	84%	85%	76%	79%

Table 1.3: Face recognition rates of the RFW dataset for *aware* and *blinded* embeddings and for different Senet models indicated by the embedding size  $N_e$ . The threshold was optimized for each case (corresponding to a column) with respect to the total dataset.

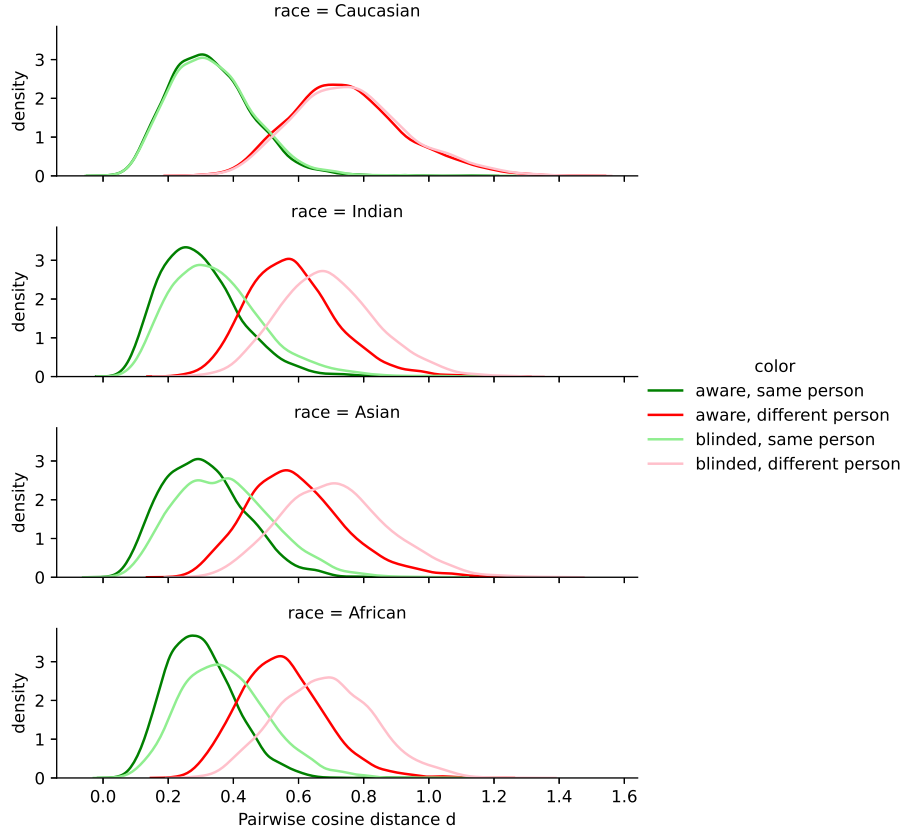


Figure 1.3: Kernel density estimation (KDE) plot of pairwise cosine distances for same persons (green) and different persons (red) and for *aware* (darker color) and *blinded* (lighter color) embeddings.

	aware acc	aware err. rate eq. group	aware err. rate diff. group	blinded acc	blinded err. rate eq. group	blinded err. rate diff. group
Caucasian	90%	8%	2%	86%	6%	8%
Indian	85%	13%	2%	83%	9%	8%
Asian	83%	16%	1%	80%	13%	7%
African	80%	19%	1%	78%	16%	6%

Table 1.4: Face recognition rates by nearest neighbor metric for *aware* and *blinded* emdeddings and the  $N_e = 128$  model. Accuracies are given for each ethnic group. Errors correspond to the remaining accuracy gap. If the nearest neighbor does not belong to the same person, the equal group error rate indicates whether the nearest neighbor belongs to the same group. Similarly, the different group error rate corresponds to nearest neighbors of a different ethnic group.

persons from different ethnic groups overcompensates the reduction leading to an overall performance decrease.

## 1.5 Discussion

The proposed procedure removes linear separability with the effect that linear classifiers can't distinguish between ethnic clusters after blinding. Clustering as measured by cluster validation scores completely vanishes with this blinding procedures. This leads to the first insight

Cluster validation scores are not a measure of bias. They indicate whether linear classifiers can predict clusters of the discriminatory dimension - a property which is related to awareness.

The fact that the ethnic clusters are well separated in the first place shows that the considered model clearly distinguishes the ethnic groups. The blinding removes this separation.

Two measures were considered to assess face recognition rates and bias:

1. Classification of positive and negative pairs
2. Nearest neighbors

To my knowledge, the first method is widely used. Surprisingly, the performance increases and the bias is slightly reduced. With the second measure, performance is significantly reduced and bias again slightly. The reason for the different outcomes is the fact that the pair method only uses pairs within the same ethnic cluster and blinding actually increases the performance within a given cluster. This leads to the second insight

Measuring face recognition rates and bias with a pair approach can be misleading when images of any given pairs are within the same cluster.

Yet another (commonly used) way of how not to measure bias was identified...