

Data Uncertainty Learning in Face Recognition

Jie Chang¹, Zhonghao Lan², Changmao Cheng¹, Yichen Wei¹
¹Megvii Inc. ²University of Science and Technology of China

{changjie, lanzhonghao, chengchangmao, weiyicheng}@megvii.com

Abstract

Modeling data uncertainty is important for noisy images, but seldom explored for face recognition. The pioneer work [35] considers uncertainty by modeling each face image embedding as a Gaussian distribution. It is quite effective. However, it uses fixed feature (mean of the Gaussian) from an existing model. It only estimates the variance and relies on an ad-hoc and costly metric. Thus, it is not easy to use. It is unclear how uncertainty affects feature learning.

This work applies data uncertainty learning to face recognition, such that **the feature (mean) and uncertainty (variance) are learnt simultaneously**, for the first time. Two learning methods are proposed. They are easy to use and outperform existing deterministic methods as well as [35] on challenging unconstrained scenarios. We also provide insightful analysis on how incorporating uncertainty estimation helps reducing the adverse effects of noisy samples and affects the feature learning.

1. Introduction

Data uncertainty¹ captures the “noise” inherent in the data. Modeling such uncertainty is important for computer vision application [22], e.g., face recognition, because noise widely exists in images.

Most face recognition methods represent each face image as a deterministic point embedding in the latent space [7, 27, 41, 42, 33]. Usually, high-quality images of the same ID are clustered. However, it is difficult to estimate an accurate point embedding for noisy face images, which are usually out of the cluster and have larger uncertainty in the embedding space. This is exemplified in Fig 1 (a). The positive example is far from its class and close to a noisy negative example, causing a mismatch.

Probabilistic face embeddings (PFE) [35] is the first work to consider *data uncertainty* in face recognition. For each sample, it estimates a Gaussian distribution, instead of

¹Uncertainty could be characterised into two main categories. Another type is *model uncertainty*.

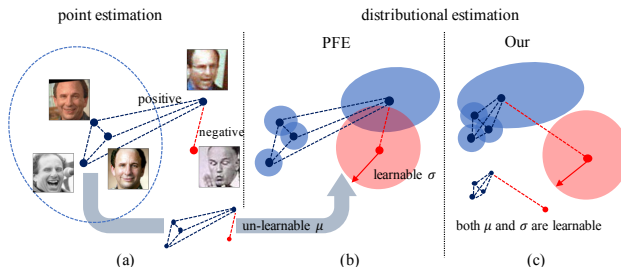


Figure 1: (a) Deterministic model gives point embedding without considering the data uncertainty; (b) probabilistic model gives distributional estimation parameterised with estimated mean and estimated variance. PFE leverages the pre-trained point embedding as the mean μ , **only** learn the uncertainty σ for each sample; (c) our method simultaneously learn σ as well as μ , leading to better intra-class compactness and inter-class separability for μ in the latent space. Different classes are marked as blue or red. **Best viewed in color.**

a fixed point, in the latent space. Specifically, given a pre-trained FR model, *the mean of the Gaussian for each sample is fixed as the embedding produced by the FR model*. An extra branch is appended to the FR model and trained to estimate the variance. The training is driven by a new similarity metric, mutual likelihood score or *MLS*, which measures the “likelihood” between two Gaussian distributions. It is shown that PFE estimates small variance for high-quality samples but large variance for noisy ones. Together with the *MLS* metric, PFE can reduce the mismatches on noisy samples. This is illustrated in Fig 1, (b). While being effective, PFE is limited in that *it does not learn the embedded feature (mean) but only the uncertainty*. As a result, it is unclear how uncertainty affects feature learning. Also, the conventional similarity metric such as cosine distance cannot be used. The more complex *MLS* metric is in demand, which takes more runtime and memory.

For the first time, this work applies *data uncertainty learning* (DUL) to face recognition such that *feature (mean) and uncertainty (variance) are learnt simultaneously*. As illustrated in Fig 1 (c), this improves the features such that the instances in the same class are more compact and the instances in different classes are more separated. In this case, the learned feature is directly usable for conventional similarity metric. *MLS* metric is no longer necessary.

Specifically, we propose two learning methods. The first

is classification based. It learns a model from scratch. The second is regression based. It improves an existing model, similar as PFE. We discuss how the learned uncertainty affects the model training in two methods, from the perspective of *image noise*. We provide insightful analysis that the learned uncertainty will improve the learning of identity embeddings by *adaptively reducing the adverse effects of noisy training samples*.

Comprehensive experiments demonstrate that our proposed methods improve face recognition performance over existing deterministic models and PFE on most public benchmarks. The improvement is more remarkable on benchmarks with low quality face images, indicating that model with *data uncertainty learning* is more suitable to unconstrained face recognition scenario, thus important for practical tasks.

2. Related Work

Uncertainty in Deep Learning The nature of uncertainties as well as the manner to deal with them have been extensively studied to help solve the reliability assessment and risk-based decision making problems for a long time [9, 31, 8]. In recent years, uncertainty is getting more attention in deep learning. Many techniques have been proposed to investigate how uncertainty specifically behaves in deep neural networks [3, 10, 11, 22]. Specific to deep uncertainty learning, *uncertainties* can be categorized into *model uncertainty* capturing the noise of the parameters in deep neural networks, and *data uncertainty* measuring the noise inherent in given training data. Recently, many computer vision tasks, i.e., semantic segmentation [19, 21], object detection [6, 25] and person Re-ID [50], have introduced deep uncertainty learning to CNNs for the improvement of model robustness and interpretability. In face recognition task, several works have been proposed to leverage *model uncertainty* for analysis and learning of face representations [13, 51, 23]. Thereinto PFE [35], is the first work to consider *data uncertainty* in face recognition task.

Noisy Data Training Large-scale datasets, i.e., CASIA-WebFace [47], Vggface2 [5] and MS-Celeb-1M [14], play the important role in training deep CNNs for face recognition. It is inevitable these face datasets collected online have lots of *label noise* — examples have erroneously been given the labels of other classes within the dataset. Some works explore the influence of *label noise* [39] and how to train robust FR models in this case [17, 44, 29]. Yu *et al.* [50] claims in person Re-ID that another *image noise* brought by poor quality images also has detrimental effect on the trained model. Our methods are not specifically proposed for noisy data training, however, we provide insightful analysis about how the learned data uncertainty affect the model training from the perspective of *image noise*. Additionally,

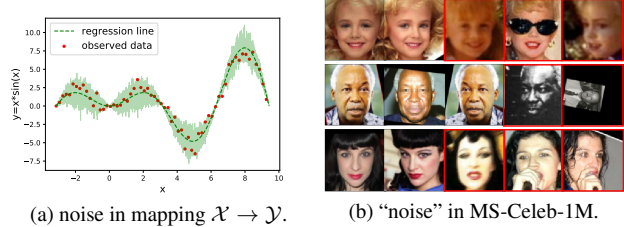


Figure 2: (a): Target y in observed data pair (red-dot) is corrupted by x dependent noise. *Data uncertainty regression* will give us the “noise level” (green-shaded) beyond the particular predicted value (green line); (b): Samples labeled with the same ID are presented in each row. Samples with red box are regarded as noisy data compared with other intra-class samples. **Best viewed in color.**

we experimentally demonstrate the proposed methods perform more robustly on noisy dataset.

3. Methodology

In Section 3.1, we first reveals the data uncertainty inherently existed in continuous mapping space and our specific face datasets. In Section 3.2, we propose DUL_{cls} to consider *data uncertainty learning* in a standard face classification model. We next propose another regression-based method, DUL_{rgs} to improve existing deterministic models in Section 3.3. Last in Section 3.4, we clarify some differences between proposed methods and existing works.

3.1. Preliminaries

Uncertainty in Continuous Mapping Space Supposing a continuous mapping space $\mathcal{X} \rightarrow \mathcal{Y}$ where each $y_i \in \mathcal{Y}$ is corrupted by some input-dependent noise, $n(\mathbf{x}_i), \mathbf{x}_i \in \mathcal{X}$, then we say this mapping space carries *data uncertainty* in itself. Considering a simple case, the noise is additive and drawn from Gaussian distribution with mean of zero and x -dependent variance. Then each observation target $y_i = f(\mathbf{x}_i) + \epsilon\sigma(\mathbf{x}_i)$, where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and $f(\cdot)$ is the embedding function we want to find. Conventional regression model only trained to approximate $f(\mathbf{x}_i)$ given the input \mathbf{x}_i . However, regression model with *data uncertainty learning* also estimates $\sigma(\mathbf{x}_i)$, representing the uncertainty of the predicted value $f(\mathbf{x}_i)$ (see Fig 2, (a)). This technique has been used by many tasks [22, 4, 30, 12, 2].

Uncertainty in Face Datasets Similar to the above continuous mapping space, face datasets composed with $\mathcal{X} \rightarrow Y$ also carries *data uncertainty*. Here \mathcal{X} is the continuous image space while Y is the discrete identity labels. Typically, large amount of face images collected online are visually ambiguous (poorly aligned, severely blurred or occluded). It is difficult to filter out these poor quality samples from training set (see Fig 2, (b)). During deep learning era, each sample is represented as an embedding \mathbf{z}_i in the latent space. If we

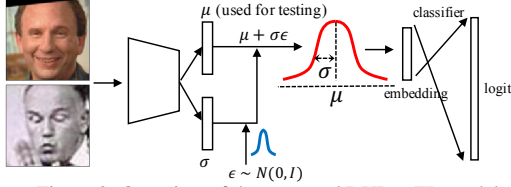


Figure 3: Overview of the proposed DUL_{cls} FR model.

hypothesize that each $\mathbf{x}_i \in \mathcal{X}$ has an ideal embedding $f(\mathbf{x}_i)$ mostly representing its identity and less unaffected by any identity irrelevant information in \mathbf{x}_i , then the embedding predicted by DNNs can be reformulated as $\mathbf{z}_i = f(\mathbf{x}_i) + n(\mathbf{x}_i)$ where $n(\mathbf{x}_i)$ is the uncertainty information of \mathbf{x}_i in the embedding space.

3.2. Classification-based DUL for FR

We propose DUL_{cls} to firstly introduce *data uncertainty learning* to the face classification model which can be trained end-to-end.

Distributional Representation Specifically, we define the representation \mathbf{z}_i in latent space of each sample \mathbf{x}_i as a Gaussian distribution,

$$p(\mathbf{z}_i|\mathbf{x}_i) = \mathcal{N}(\mathbf{z}_i; \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2 \mathbf{I}) \quad (1)$$

where both the parameters (mean as well as variance) of the Gaussian distribution are input-dependent predicted by CNNs: $\boldsymbol{\mu}_i = f_{\theta_1}(\mathbf{x}_i)$, $\boldsymbol{\sigma}_i = f_{\theta_2}(\mathbf{x}_i)$, where θ_1 and θ_2 refer to the model parameters respectively w.r.t output $\boldsymbol{\mu}_i$ and $\boldsymbol{\sigma}_i$. Here we recall that the predicted Gaussian distribution is diagonal multivariate normal. $\boldsymbol{\mu}_i$ can be regarded as the identity feature of the face and the $\boldsymbol{\sigma}_i$ refers to the uncertainty of the predicted $\boldsymbol{\mu}_i$. Now, the representation of each sample is not a deterministic point embedding any more, but a stochastic embedding *sampled* from $\mathcal{N}(\mathbf{z}_i; \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2 \mathbf{I})$, in the latent space. However, sampling operation is not differentiable preventing the backpropagation of the gradients flow during the model training. We use re-parameterization trick [24] to let the model still take gradients as usual. Specifically, we first sample a random noise ϵ from a normal distribution, which is independent of the model parameters, and then generate \mathbf{s}_i as the equivalent sampling representation (see Fig 3 for an overview pipeline),

$$\mathbf{s}_i = \boldsymbol{\mu}_i + \epsilon \boldsymbol{\sigma}_i, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (2)$$

Classification Loss Since \mathbf{s}_i is the final representation of each image \mathbf{x}_i , we feed it to a classifier to minimize the following softmax loss,

$$\mathcal{L}_{softmax} = \frac{1}{N} \sum_i \log \frac{e^{\mathbf{w}_{y_i} \mathbf{s}_i}}{\sum_c e^{\mathbf{w}_c \mathbf{s}_i}}, \quad (3)$$

In practice, we use different variants of $\mathcal{L}_{softmax}$ such as additive margin [40], feature ℓ_2 normalization [32] and arc-face [7], to train our face classification model.

KL-Divergence Regularization Eq. 2 indicates that all identity embeddings $\boldsymbol{\mu}_i$ are corrupted by $\boldsymbol{\sigma}_i$ during the training period, this will prompt the model to predict small $\boldsymbol{\sigma}$ for all samples in order to suppress the unstable ingredients in \mathbf{s}_i such that Eq. 3 can still converge at last. In this case, the stochastic representation can be reformulated as $\mathbf{s}_i = \boldsymbol{\mu}_i + c$ which is actually degraded to the original deterministic representation². Inspired by the variational information bottleneck [1], we introduce a regularization term during the optimization by explicitly constraining $\mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i)$ to be close to a normal distribution, $\mathcal{N}(\mathbf{0}, \mathbf{I})$, measured by Kullback-Leibler divergence (KLD) between these two distributions. This KLD term is,

$$\begin{aligned} \mathcal{L}_{kl} &= KL[\mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2)|\mathcal{N}(\epsilon|\mathbf{0}, \mathbf{I})] \\ &= -\frac{1}{2}(1 + \log \boldsymbol{\sigma}^2 - \boldsymbol{\mu}^2 - \boldsymbol{\sigma}^2) \end{aligned} \quad (4)$$

Noted that \mathcal{L}_{kl} is monotonely decreasing w.r.t $\boldsymbol{\sigma}$ under the restriction that $\boldsymbol{\sigma}_i^{(l)} \in (0, 1)$ (l refers to the l^{th} dimension of the embedding). \mathcal{L}_{kl} works as a good ‘‘balancer’’ with Eq. 3. Specifically, DUL_{cls} is discouraged from predicting large variance for all samples, which may lead to extremely corruption on $\boldsymbol{\mu}_i$, thus making $\mathcal{L}_{softmax}$ hard to converge. Simultaneously, DUL_{cls} is also discouraged from predicting lower variance for all samples, which may lead to larger \mathcal{L}_{kl} to punish the model in turn.

Last, we use $\mathcal{L}_{cls} = \mathcal{L}_{softmax} + \lambda \mathcal{L}_{kl}$ as the total cost function, and λ is a trade-off hyper-parameter, and it is further analysed in Section 4.6.

3.3. Regression-based DUL for FR

DUL_{cls} is a general classification model with data uncertainty learning. Next we propose another regression based method, DUL_{rgs} , improving existing FR models by data uncertainty learning.

Difficulty of Introducing Data Uncertainty Regression to FR

DUL_{rgs} is inspired from data uncertainty regression [26, 22] for continuous mapping space $\mathcal{X} \rightarrow \mathcal{Y}$ as described in Section 3.1. However, mapping space in face datasets is constructed by continuous image space \mathcal{X} and discrete identity label \mathcal{Y} , which cannot be directly fitted via data uncertainty regression. The key point lies in that the identity labels $y_c \in Y$ cannot serve as continuous target vector to be approximated. This difficulty is also mentioned in PFE [35] but is not resolved.

²Here c refers to the estimated $\boldsymbol{\sigma}$ which nearly constant and small.

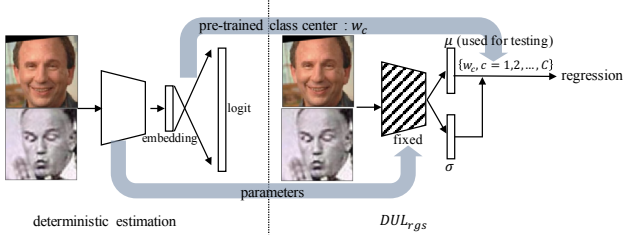


Figure 4: Overview of the proposed DUL_{rgs} model. All parameters in the convolution layers are pre-trained by a deterministic FR model and are fixed during the training of DUL_{rgs} .

Constructing New Mapping Space for FR We construct a new target space, which is continuous, for face data. Most importantly, it is nearly equivalent to the original discrete target space Y , which encourages the correct mapping relationship. Specifically, we pre-train a classification-based deterministic FR model, and then utilize the weights in its classifier layers, $\mathcal{W} \in \mathbb{R}^{D \times C}$ as the expected target vector³. Since each $\mathbf{w}_i \in \mathcal{W}$ can be treated as the typical center of the embeddings with the same class, $\{\mathcal{X}, \mathcal{W}\}$ thus can be regarded as the new equivalent mapping space. Similar to the uncertainty in continuous mapping space as described in Section 3.1, $\{\mathcal{X}, \mathcal{W}\}$ has inherent noise. We can formulate the mapping from $\mathbf{x}_i \in \mathcal{X}$ to $\mathbf{w}_i \in \mathcal{W}$ as $\mathbf{w}_i = f(\mathbf{x}_i) + n(\mathbf{x}_i)$, where $f(\mathbf{x}_i)$ is the “ideal” identity feature and each observed \mathbf{w}_i is corrupted by input dependent noise.

Distributional Representation Next we can estimate above $f(\mathbf{x}_i)$ and $n(\mathbf{x}_i)$ by data uncertainty regression. Specifically, a Gaussian distribution is assumed for the likelihood: $p(\mathbf{z}_i | \mathbf{x}_i) = \mathcal{N}(\mathbf{z}_i; \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2 \mathbf{I})$ where $\boldsymbol{\mu}_i$ as well as $\boldsymbol{\sigma}_i$ are also parameterised by the weights in neural networks⁴ (see Fig. 4). If we take each \mathbf{w}_c as the target, we should maximize the following likelihood for each \mathbf{x}_i ,

$$p(\mathbf{w}_c | \mathbf{x}_i \in c, \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi\boldsymbol{\sigma}_i^2}} \exp\left(-\frac{(\mathbf{w}_c - \boldsymbol{\mu}_i)^2}{2\boldsymbol{\sigma}_i^2}\right). \quad (5)$$

Actually, we take the log likelihood as follows,

$$\ln p(\mathbf{w}_c | \mathbf{x}_i \in c, \boldsymbol{\theta}) = -\frac{(\mathbf{w}_c - \boldsymbol{\mu}_i)^2}{2\boldsymbol{\sigma}_i^2} - \frac{1}{2} \ln \boldsymbol{\sigma}_i^2 - \frac{1}{2} \ln 2\pi. \quad (6)$$

Assumed that $\mathbf{x}_i, i \in 1, 2, \dots$ are independently and identically distributed (*iid.*), the likelihood over all data-points is $\prod_c \prod_i \ln p(\mathbf{w}_c | \mathbf{x}_i \in c, \boldsymbol{\theta})$. Practically, we train the network to predict the log variance, $\mathbf{r}_i := \ln \boldsymbol{\sigma}_i^2$, to stabilize the numerical during the stochastic optimization. Last, the likelihood maximization is reformulated as the minimization of

³Here D refers to the dimensions of the embedding and C refers to the numbers of classes in training set.

⁴Noted here $\boldsymbol{\mu}_i \approx f(\mathbf{x}_i)$ and $\boldsymbol{\sigma}_i \approx n(\mathbf{x}_i)$.

cost function,

$$\mathcal{L}_{rgs} = \frac{1}{2N} \sum_c \sum_{i \in c} \left[\frac{1}{D} \sum_l (\exp(-r_i^{(l)}) (w_c^{(l)} - \mu_i^{(l)})^2 + r_i^{(l)} \right], \quad (7)$$

where D , N and l refers to the size of embedding dimension, the size of data-points and the l^{th} dimension of each feature vector, respectively. We omit the constant term, $\frac{D}{2} \ln 2\pi$ during the optimization.

Loss Attenuation Mechanism By qualitatively analyzing Eq. 6, our learned variance $\boldsymbol{\sigma}_i$ could actually be regarded as the uncertainty score measuring the confidence of the learned identity embedding, $\boldsymbol{\mu}_i$, belonging to c^{th} class. Specifically, for those ambiguous $\boldsymbol{\mu}_i$ located far away from its class center \mathbf{w}_c , DUL_{rgs} will estimate large variance to temper the error term, $\frac{(\mathbf{w}_c - \boldsymbol{\mu}_i)^2}{2\boldsymbol{\sigma}_i^2}$, instead of overfitting on these noisy samples. DUL_{rgs} is discouraged from predicting large variance for all samples, which may lead to underfitting of $(\mathbf{w}_c - \boldsymbol{\mu})^2$ and larger $\log \boldsymbol{\sigma}$ term will punish the model in turn. Simultaneously, DUL_{rgs} is also discouraged from predicting very small variance for all samples, which may lead to exponentially increases of error term. Thus, Eq. 6 allows DUL_{rgs} to adapt the weighting of error term. This makes the model learn to attenuate the effect from those ambiguous $\boldsymbol{\mu}_i$ caused by poor quality samples.

3.4. Discussion of Related Works

We first discuss the connection between DUL_{cls} and variational information bottleneck (VIB) [1]. VIB [1] is a variational approximation to information bottleneck (IB) principle [38] under the framework of deep learning. VIB seeks a stochastic mapping from input data X to latent representation Z , in terms of the fundamental trade-off between making Z as concise as possible but still have enough ability to predict label Y [38]. It is noted that \mathcal{L}_{cls} is similar to the objective function in VIB. However, we analyze this classification method from data uncertainty perspective while VIB derived this objective function from the view of information bottleneck.

We next clarify some differences between DUL_{rgs} and PFE [35]. Although both PFE and DUL_{rgs} formally encode the input uncertainty as variance representation. However, PFE essentially measures the likelihood of each positive pair of $\{\mathbf{x}_i, \mathbf{x}_j\}$ sharing the same latent embedding: $p(\mathbf{z}_i = \mathbf{z}_j)$. While DUL_{rgs} interprets a conventional *Least-Square Regression* technique as a *Maximum likelihood Estimation* with a data uncertainty regression model.

Last, both DUL_{cls} and DUL_{rgs} learn identity representation $\boldsymbol{\mu}$ as well as uncertainty representation $\boldsymbol{\sigma}$, which ensure our predicted $\boldsymbol{\mu}$ can be directly evaluated by common-used matching metric. However, PFE has to use mutual likelihood score (*MLS*) as the matching metric to improve the

performance of deterministic model because identity representation is not learnt in PFE.

4. Experiments

In this section, we first evaluate the proposed methods on standard face recognition benchmarks. Then we provide qualitative and quantitative analysis to explore what is the meaning of the learned data uncertainty and how data uncertainty learning affects the learning of FR models. Last, we conduct experiments on the noisy MS-Celeb-1M dataset to demonstrate that our methods perform more robustly than deterministic methods.

4.1. Datasets and Implementation Details

We describe the public datasets that are used, and our implementation details.

Datasets We use MS-Celeb-1M datasets with 3,648,176 images of 79,891 subjects as training set. 2 benchmarks including LFW [18] and MegaFace [20]⁵, and 3 unconstrained benchmarks: CFP [34]⁶, YTF [43] and IJB-C [28], are used to evaluate the performance of $DUL_{cls/rgs}$ following the standard evaluation protocols.

Architecture We train baseline models on ResNet [15] backbone with SE-blocks [16]. The head of the baseline model is: `BackBone-Flatten-FC-BN` with embedding dimensions of 512 and dropout probability of 0.4 to output the embedding feature. Compared with baseline model, DUL_{cls} has an additional head branch sharing the same architecture to output the variance. DUL_{rgs} also has an additional head branch whilst its architecture is: `BackBone-Flatten-FC-BN-ReLU-FC-BN-exp`, to output the variance.

Training All baseline models and DUL_{cls} models are trained for 210,000 steps using a SGD optimizer with a momentum of 0.9, weight decay of 0.0001, batch size of 512. We use triangular learning rate policy [36] with the *max_lr* of 0.1 and *base_lr* of 0. For most DUL_{cls} models, we set trade-off hyper-parameter λ as 0.01. For the proposed DUL_{rgs} , we first train baseline model for 210,000 steps and then fix parameters in all convolution layers (*step 1*). Then we train the mean branch as well as the variance branch in head from scratch for additional 140,000 steps with batch size of 256 (*step 2*). During *step 2*, we set learning rate starting at 0.01, and then decreased to 0.001 and 0.0001 at 56,000 and 84,000 steps.

⁵Noted that we use rank1 protocol of MegaFace

⁶Noted that we only use “frontal-profile” protocol of CFP

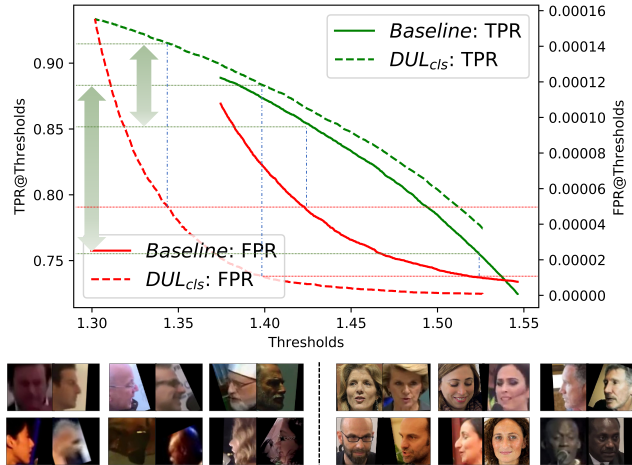


Figure 5: Top: TPR&FPR vs. threshold on IJB-C; Bottom: false acceptance cases mostly happened in the baseline model (left); false acceptance cases mostly happened in DUL_{cls} (right). Both baseline model and DUL_{cls} are trained by ResNet18 with AM-Softmax on MS-Celeb-1M dataset. **Best viewed in color.**

4.2. Comparing DUL with Deterministic Baselines

In this part, all baseline models are trained with ResNet18 backbone [15], equipped with different variants of softmax loss, i.e., AM-Softmax [40], ArcFace [7] and L2-Softmax [32]. Both the embedding features and the weights in classifier are ℓ_2 -normalized during the training. Our proposed DUL_{cls} models are trained with the same backbone and loss functions. Our proposed DUL_{rgs} models are trained based on the different pre-trained baseline models, as described in Section 4.1.

Table 1 reports the testing results obtained by the baseline models (“Original”) and the proposed DUL models. Cosine similarity is used for evaluation. Our proposed methods outperform the baseline deterministic models on most benchmarks⁷. This demonstrates that the proposed methods are effective on different state-of-the-art loss functions. These results indicate that the identity embeddings (μ in our methods) trained with data uncertainty (σ in our method) present better intra-class compactness and inter-class separability than the point embeddings estimated by baseline models, especially on those unconstrained benchmarks: CFP with frontal/profile photos and YTF/IJB-C with most blur photos collected from YouTube videos, compared with benchmarks with most clear and frontal photos (LFW and MegaFace).

The proposed DUL achieves most remarkable improvement on verification protocols of IJB-C benchmark, which is also the most challenging one. We thus plot how true acceptance rate (TPR) and false acceptance rate (FPR) perform along with the change of thresholds. As illustrated in Fig 5, DUL_{cls} achieves higher TPR and lower FPR than

⁷Noted that DUL_{rgs} combined with L2-Softmax deteriorates on IJB-C, which should be further explored in the future.

| Base Model | Representation | LFW | CFP-FP | MegaFace(R1) | YTF | IJB-C (TPR@FPR) | | | |
|-----------------|--------------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------|--------------|
| | | | | | | 0.001% | 0.01% | 0.1% | AUC |
| AM-Softmax [40] | Original | 99.63 | 96.85 | 97.11 | 96.09 | 75.43 | 88.65 | 94.73 | 87.51 |
| | PFE [35] | 99.68 | 94.57 | 97.18 | 96.12 | 86.24 | 92.11 | 95.71 | 91.71 |
| | DUL _{cls} | 99.71 | 97.28 | 97.30 | 96.46 | 88.25 | 92.78 | 95.57 | 92.40 |
| | DUL _{rgs} | 99.66 | 97.61 | 96.85 | 96.28 | 87.02 | 91.84 | 95.02 | 91.44 |
| ArcFace [7] | Original | 99.64 | 96.77 | 97.08 | 96.06 | 73.80 | 88.78 | 94.30 | 86.94 |
| | PFE [35] | 99.68 | 95.34 | 96.55 | 96.32 | 86.69 | 92.28 | 95.66 | 91.89 |
| | DUL _{cls} | 99.76 | 97.01 | 97.22 | 96.20 | 87.22 | 92.43 | 95.38 | 92.10 |
| | DUL _{rgs} | 99.66 | 97.11 | 96.83 | 96.38 | 86.21 | 91.03 | 94.53 | 90.79 |
| L2-Softmax [32] | Original | 99.60 | 95.87 | 90.34 | 95.89 | 77.60 | 86.19 | 92.55 | 85.83 |
| | PFE [35] | 99.66 | 86.45 | 90.64 | 95.98 | 79.33 | 87.28 | 93.41 | 87.01 |
| | DUL _{cls} | 99.63 | 97.24 | 93.19 | 96.56 | 79.90 | 87.80 | 93.44 | 87.38 |
| | DUL _{rgs} | 99.66 | 96.35 | 89.66 | 96.08 | 74.46 | 83.23 | 91.09 | 83.10 |

Table 1: Results of models (ResNet18) trained on MS-Celeb-1M. ‘‘Original’’ refers to the deterministic embeddings. The better performance among each base model are shown in bold numbers. We use σ both for fusion and matching (with mutual likelihood scores) in PFE. AUC is calculated when FPR spans on the interval [0.001%, 0.1%] and we rescale it.

baseline model at different settings of matching threshold. Additionally, the lower FPR is set, the better DUL_{cls} performs on TPR. Fig 5 also shows the vast majority cases of false acceptance respectively happened in baseline model and DUL_{cls}. We can see that DUL_{cls} resolves more FP cases with extreme noises, which are typically occurring in the baseline model. This indicates that model with data uncertainty learning is more applicable to the unconstrained face recognition scenario than deterministic model.

We have the similar conclusion for DUL_{rgs}.

4.3. Comparing DUL with PFE

For comparison, we re-implemented PFE on all baseline models according to the recommended settings of implementation details in [35]⁸. We note that our re-implementation has achieved similar or slightly better results than those in [35]. Our DUL_{cls/rgs} use averaged pooling aggregation for features in template and are evaluated by cosine similarity. Compared with PFE, our proposed DUL_{cls} achieves better performances in all cases, and the proposed DUL_{rgs} also shows competitive performances. Results are reported in Table 1.

PFE interprets the point embedding learned by deterministic FR models as the mean of its output distributional estimation and only learn the uncertainty (variance) for each sample. Thus, PFE has to use *MLS* metric, which takes the predicted variance into account. Although PFE achieves better results with the help of the matching measurement with more precision, it still suffers more computational complexity for matching. Specifically, for verification of 6000 face pairs (LFW), standard cosine metric takes less than 1 second via matrix multiplication, while *MLS* takes 1min28s, on two GTX-1080.

⁸which means we use mutual likelihood score (*MLS*) for matching and its proposed fusion strategy for feature aggregation in template/video benchmarks, i.e., YTF and IJB-C.

| Method | Training Data | LFW | YTF | MegaFace | CFP-FP |
|------------------------|---------------|-------|--------------|----------|--------|
| FaceNet [33] | 200M | 99.63 | 95.1 | - | - |
| DeepID2+ [37] | 300K | 99.47 | 93.2 | - | - |
| CenterFace [42] | 0.7M | 99.28 | 94.9 | 65.23 | 76.52 |
| SphereFace [27] | 0.5M | 99.42 | 95.0 | 75.77 | 89.14 |
| ArcFace [7] | 5.8M | 99.83 | 98.02 | 81.03 | 96.98 |
| CosFace [41] | 5M | 99.73 | 97.6 | 77.11 | 89.88 |
| L2-Face [32] | 3.7M | 99.78 | 96.08 | - | - |
| Yin <i>et al.</i> [49] | 1M | 98.27 | - | - | 94.39 |
| PFE [35] | 4.4M | 99.82 | 97.36 | 78.95 | 93.34 |
| Baseline | 3.6M | 99.83 | 96.50 | 98.30 | 98.75 |
| PFE _{rep} | 3.6M | 99.82 | 96.50 | 98.48 | 97.28 |
| DUL _{cls} | 3.6M | 99.78 | 96.78 | 98.60 | 98.67 |
| DUL _{rgs} | 3.6M | 99.83 | 96.84 | 98.12 | 98.78 |

Table 2: Comparison with the state-of-the-art methods on LFW, YTF, MegaFace (MF) and CFP-FP. ‘‘-’’ indicates that the author did report the performance on the corresponding protocol. ‘‘PFE_{rep}’’ means we reproduce PFE by ourself. Backbone: ResNet64.

4.4. Comparison with State-Of-The-Art

To compare with state-of-the-art, we use a deeper and stronger backbone, ResNet64, trained with AM-Softmax loss on MS-Celeb-1M dataset, as our baseline model. Then we train the proposed DUL models following the setting described in section 4.1.

The results are illustrated in Table 2. Noted that performances of baseline model have been saturated on LFW and CFP-FP, where the merit of data uncertainty learning is not obvious. However, DUL_{cls/rgs} still slightly improve the accuracy on YTF and MegaFace⁹. Table 3 reports the results of different methods on IJB-C. Both PFE and DUL achieve much better performances over baseline models.

4.5. Understand Uncertainty Learning

In this part, we qualitatively and quantitatively analyze the proposed DUL to gain more insights about data uncertainty learning.

⁹Noted that our used MegaFace datasets is refined, while previous reported SOTA results in Table 2 usually use non-refined MegaFace.

| Method | Training Data | IJB-C (TPR@FPR) | | | |
|------------------------|---------------|-----------------|--------------|--------------|--------------|
| | | 0.001% | 0.01% | 0.1% | AUC |
| Yin <i>et al.</i> [48] | 0.5M | - | - | 69.3 | - |
| Cao <i>et al.</i> [5] | 3.3M | 74.7 | 84.0 | 91.0 | - |
| Multicolumn [46] | 3.3M | 77.1 | 86.2 | 92.7 | - |
| DCN [45] | 3.3M | - | 88.5 | 94.7 | - |
| PFE [35] | 4.4M | 89.64 | 93.25 | 95.49 | - |
| Baseline | 3.6M | 83.06 | 92.16 | 95.83 | 91.97 |
| PFE _{rep} | 3.6M | 89.77 | 94.14 | 96.37 | 93.74 |
| DUL _{cls} | 3.6M | 88.18 | 94.61 | 96.70 | 93.97 |
| DUL _{rgs} | 3.6M | 90.23 | 94.21 | 96.32 | 93.88 |

Table 3: Comparison with the state-of-the-art methods on IJB-C. Backbone: ResNet64.

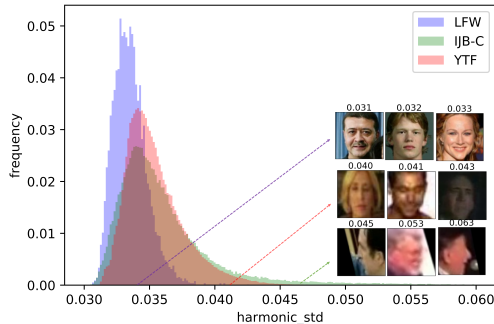


Figure 6: Uncertainty distribution on different dataset for DUL_{rgs}. Similar uncertainty distribution has also been observed in DUL_{cls}. **Best viewed in color.**

What is the meaning of the learned uncertainty? The estimated uncertainty is closely related to the quality of face images, for both DUL_{cls} and DUL_{rgs}. This is also observed in PFE [35]. For visualization, we show the learned uncertainty¹⁰ of different dataset in Figure 6. It illustrates that the learned uncertainty increases along with the image quality degradation. This learned uncertainty could be regarded as the quality of the corresponding identity embedding estimated by the model, measuring the proximity of the predicted face representation to its genuine (or true) point location in the latent space.

Therefore, two advantages are obtained for face recognition with data uncertainty learning. First, the learned variance can be utilized as a “risk indicator” to alert FR systems that the output decision is unreliable when the estimated variance is very high. Second, the learned variance also can be used as the measurement of image quality assessment. In this case, we note that it is unnecessary to train a separate quality assessment model which requires explicit quality labels as before.

How the learned uncertainty affect the FR model? In this part, we attempt to shed some light on the mechanism of how the learned data uncertainty affects the model training

¹⁰Specifically, we use harmonic mean of the predicted variance $\sigma \in \mathbb{R}^{512}$ as the approximated measurement of the estimated uncertainty. The same below.

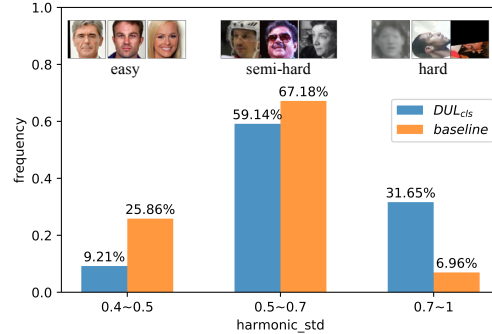


Figure 7: Bad case analysis between baseline model and DUL_{cls}. **Best viewed in color.**

and helps to obtain better feature embeddings.

We classify the training samples in MS-Celeb-1M dataset into three categories according to the degree of estimated uncertainty by DUL_{cls}: *easy* samples with low variance, *semi-hard* samples with medium variance and *hard* samples with large variance. We calculated the proportion of mis-classified samples in each of the three categories to all mis-classified samples respectively produced by baseline model and our DUL_{cls}. Fig 7 illustrates that our DUL_{cls} causes relatively less bad cases on easy samples as well as semi-hard samples, compared with the baseline model. However, for those hard samples with extreme noises, baseline model produces less bad cases, when compared with DUL_{cls}. This demonstrates that FR networks with data uncertainty learning focus more on those training samples which **should** be correctly classified and simultaneously “give up” those detrimental samples, instead of over-fitting them. This supports our previous discussion in Section 3.2.

We also conduct similar experiment for DUL_{rgs}. We calculate the averaged euclidean distances¹¹ between the class center w_c and its intra-class estimated identity embedding, $\mu_{i \in c}$, respectively for baseline model and DUL_{rgs}. As illustrated in Fig 8, DUL_{rgs} pulls the easy and semi-hard samples closer to their class center whilst pushes those hard samples further away. This also supports our discussion in Section 3.3 that Eq. 6 effectively prevents model over-fitting on extremely noisy samples by the adaptive weighting mechanism w.r.t σ .

Last, we manually construct imposter/genuine test pair with different blurriness to compare the cosine similarity respectively obtained by baseline model and our methods. As illustrated in Fig 9, along with the increase of blurriness, both baseline model and DUL deteriorate rapidly. However, our proposed DUL achieves higher similarity score for genuine pair and lower similarity score for imposter pair than baseline model, indicating that it is more robust.

¹¹Noted this averaging distances are further averaged over all classes in MS-Celeb-1M.

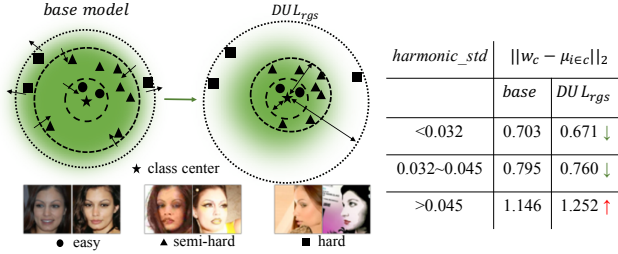


Figure 8: Averaged intra-class distances $\|\mu_{i \in c} - \mathbf{w}_c\|_2$ between base model and DUL_{rgs} .

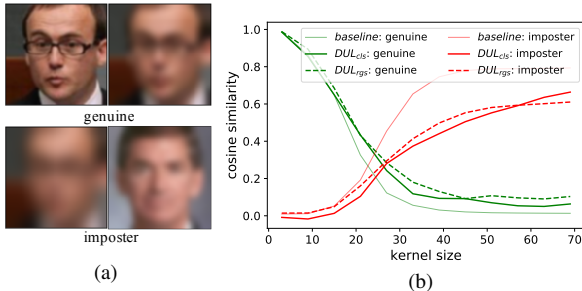


Figure 9: (a) Blur genuine and imposter pair; (b) Cosine similarity score obtained by baseline and proposed DUL for each pair.

4.6. Other Experiments

Impact of hyper-parameter of DUL_{cls} In this part, we qualitatively analyze what the trade-off hyper-parameter λ controls in DUL_{cls} . As mentioned in VIB [1], KLD term works as a regularization to trade off the conciseness and the richness of the information preserved in bottleneck embeddings. We experimentally find the KL divergence in our method affects the representational capacity of σ . As illustrated in Table 4, DUL_{cls} without the optimization of KLD term ($\lambda = 0$) performs close to baseline model. In this case, DUL_{cls} estimates relatively small σ_i for all samples, which makes the sampled representation $\mu_i + \epsilon\sigma_i$ nearly deterministic. With the enhancement of the optimization strength of KLD term ($\lambda \uparrow$), DUL_{cls} is prone to “assign” larger variance for noisy samples and small variance for high quality ones (as illustrated in Fig 7). However, overly minimizing KLD ($\lambda = 1$) term will prompt the model to predict large variance for all samples, which makes \mathcal{L}_{cls} in Eq. 3 hard to converge, thus the performances deteriorate rapidly (see Table 4).

DUL performs more robustly on noisy training data. Based on the analysis of Section 3.4 about how the learned variance affect the model training. We further conduct experiments on noisy MS-Celeb-1M to prove it. We randomly select different proportions of samples from MS-Celeb-1M to pollute them with Gaussian blur noise. Table 5 demonstrates that our proposed $DUL_{cls/rgs}$ perform more robustly on noisy training data.

| λ | $\bar{\sigma}$ | YTF | MegaFace | IJB-C (TPR@FPR) | | |
|-----------|----------------|-------|----------|-----------------|-------|-------|
| | | | | 0.001% | 0.01% | 0.1% |
| baseline | - | 96.09 | 97.11 | 75.32 | 88.65 | 94.73 |
| 0.0 | 0.2562 | 96.14 | 97.13 | 64.92 | 88.55 | 94.64 |
| 0.00010 | 0.3074 | 96.36 | 97.25 | 65.44 | 85.22 | 94.44 |
| 0.001 | 0.3567 | 96.26 | 97.38 | 62.88 | 86.65 | 94.46 |
| 0.01 | 0.5171 | 96.46 | 97.30 | 88.25 | 92.78 | 95.57 |
| 0.1 | 0.8505 | 96.42 | 95.07 | 87.19 | 91.78 | 95.13 |
| 0.5 | 0.9012 | 87.40 | 85.73 | 40.23 | 52.70 | 58.52 |
| 1.0 | 0.9520 | 75.14 | 63.90 | 1.770 | 4.530 | 13.02 |

Table 4: Results of DUL_{cls} trained with different trade-off λ . $\bar{\sigma}$ represents we average the harmonic mean of the estimated variance over all training samples in MS-Celeb-1M. The backbone is ResNet18 with AM-Softmax loss.

| percent | Model | MegaFace | LFW | YTF | IJB-C (TPR@FPR) | | |
|---------|-------------|----------|-------|-------|-----------------|-------|-------|
| | | | | | 0.001% | 0.01% | 0.1% |
| 0% | baseline | 97.11 | 99.63 | 96.09 | 75.32 | 88.65 | 94.73 |
| | baseline | 96.64 | 99.63 | 96.16 | 64.96 | 86.00 | 94.82 |
| | PFE [35] | 97.02 | 99.63 | 96.1 | 83.39 | 91.33 | 95.54 |
| 10% | DUL_{cls} | 96.88 | 99.75 | 96.44 | 88.04 | 93.21 | 95.96 |
| | DUL_{rgs} | 96.05 | 99.71 | 96.46 | 84.74 | 91.56 | 95.30 |
| | baseline | 96.20 | 99.61 | 96.00 | 43.52 | 80.48 | 94.22 |
| 20% | PFE [35] | 96.90 | 99.61 | 95.86 | 82.03 | 90.89 | 95.38 |
| | DUL_{cls} | 96.37 | 99.71 | 96.68 | 89.01 | 93.24 | 95.97 |
| | DUL_{rgs} | 95.51 | 99.66 | 96.64 | 81.10 | 90.91 | 95.27 |
| 30% | baseline | 95.72 | 99.60 | 95.45 | 31.51 | 76.09 | 93.11 |
| | PFE [35] | 96.82 | 99.61 | 96.12 | 80.92 | 90.31 | 95.29 |
| | DUL_{cls} | 95.86 | 99.73 | 96.38 | 86.05 | 91.80 | 95.02 |
| 40% | DUL_{rgs} | 94.96 | 99.66 | 96.66 | 81.54 | 91.20 | 95.32 |
| | baseline | 95.14 | 99.56 | 95.51 | 39.69 | 77.12 | 93.73 |
| | PFE [35] | 96.59 | 99.59 | 95.94 | 77.72 | 89.46 | 94.82 |
| | DUL_{cls} | 95.33 | 99.66 | 96.54 | 84.15 | 92.60 | 95.85 |
| | DUL_{rgs} | 94.28 | 99.58 | 96.68 | 78.13 | 87.64 | 94.67 |

Table 5: Comparison of baseline model and proposed $DUL_{cls/rgs}$ trained on noisy MS-Celeb-1M. Backbone model is ResNet18 with AM-Softmax loss.

5. Conclusion

In this work, we propose two general learning methods to further develop and perfect the *data uncertainty learning* (DUL) for face recognition: DUL_{cls} and DUL_{rgs} . Both methods give a Gaussian distributional estimation for each face image in the latent space and simultaneously learn identity feature (mean) and uncertainty (variance) of the estimated mean. Comprehensive experiments demonstrate that our proposed methods perform better than deterministic models on most benchmarks. Additionally, we discuss how the learned *uncertainty* affects the training of model from the perspective of *image noise* by both qualitative analysis and quantitative results.

6. Acknowledgement

This paper is supported by the National key R&D plan of the Ministry of science and technology (Project Name: “Grid function expansion technology and equipment for community risk prevention”, Project No.2018YFC0809704)

References

- [1] Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational information bottleneck. In *Proceedings of the International Conference on Learning Representations*, 2017.
- [2] Christopher M Bishop and Cazaow S Quazaz. Regression with input-dependent noise: A bayesian treatment. In *Advances in Neural Information Processing Systems*, pages 347–353, 1997.
- [3] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- [4] Axel Brando, Jose A Rodríguez-Serrano, Mauricio Ciprian, Roberto Maestre, and Jordi Vitrià. Uncertainty modelling in deep networks: Forecasting short and noisy series. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 325–340. Springer, 2018.
- [5] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 67–74. IEEE, 2018.
- [6] Jiwoong Choi, Dayoung Chun, Hyun Kim, and Hyuk-Jae Lee. Gaussian yolov3: An accurate and fast object detector using localization uncertainty for autonomous driving. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [7] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019.
- [8] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural Safety*, 31(2):105–112, 2009.
- [9] Michael Havbro Faber. On the treatment of uncertainties and probabilities in engineering decision analysis. *Journal of Offshore Mechanics and Arctic Engineering*, 127(3):243–248, 2005.
- [10] Yarin Gal. *Uncertainty in deep learning*. PhD thesis, PhD thesis, University of Cambridge, 2016.
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.
- [12] Paul W Goldberg, Christopher KI Williams, and Christopher M Bishop. Regression with input-dependent noise: A gaussian process treatment. In *Advances in neural information processing systems*, pages 493–499, 1998.
- [13] Sixue Gong, Vishnu Naresh Boddeti, and Anil K Jain. On the capacity of face representation. *arXiv preprint arXiv:1709.10433*, 2017.
- [14] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *European Conference on Computer Vision*, pages 87–102. Springer, 2016.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [17] Wei Hu, Yangyu Huang, Fan Zhang, and Ruirui Li. Noise-tolerant paradigm for training face recognition cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11887–11896, 2019.
- [18] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. 2008.
- [19] Shuya Isobe and Shuichi Arai. Deep convolutional encoder-decoder network with model uncertainty for semantic segmentation. In *2017 IEEE International Conference on Innovations in Intelligent Systems and Applications (INISTA)*, pages 365–370. IEEE, 2017.
- [20] Ira Kemelmacher-Shlizerman, Steven M Seitz, Daniel Miller, and Evan Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4873–4882, 2016.
- [21] Alex Kendall, Vijay Badrinarayanan, and Roberto Cipolla. Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *BMVC*, 2015.
- [22] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017.
- [23] Salman Khan, Munawar Hayat, Syed Waqas Zamir, Jianbing Shen, and Ling Shao. Striking the right balance with uncertainty. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 103–112, 2019.
- [24] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *ICLR*, 2014.
- [25] Florian Kraus and Klaus Dietmayer. Uncertainty estimation in one-stage object detection. *arXiv preprint arXiv:1905.10296*, 2019.
- [26] Quoc V Le, Alex J Smola, and Stéphane Canu. Heteroscedastic gaussian process regression. In *Proceedings of the 22nd international conference on Machine learning*, pages 489–496. ACM, 2005.
- [27] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphreface: Deep hypersphere embedding for face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 212–220, 2017.
- [28] Brianna Maze, Jocelyn Adams, James A Duncan, Nathan Kalka, Tim Miller, Charles Otto, Anil K Jain, W Tyler Niggel, Janet Anderson, Jordan Cheney, et al. Iarpa janus benchmark-c: Face dataset and protocol. In *2018 International Conference on Biometrics (ICB)*, pages 158–165. IEEE, 2018.

- [29] Hong-Wei Ng and Stefan Winkler. A data-driven approach to cleaning large face datasets. In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 343–347. IEEE, 2014.
- [30] David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*, volume 1, pages 55–60. IEEE, 1994.
- [31] M Elisabeth Paté-Cornell. Uncertainties in risk analysis: Six levels of treatment. *Reliability Engineering & System Safety*, 54(2-3):95–111, 1996.
- [32] Rajeep Ranjan, Carlos D Castillo, and Rama Chellappa. L2-constrained softmax loss for discriminative face verification. *arXiv preprint arXiv:1703.09507*, 2017.
- [33] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [34] Soumyadip Sengupta, Jun-Cheng Chen, Carlos Castillo, Vishal M Patel, Rama Chellappa, and David W Jacobs. Frontal to profile face verification in the wild. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–9. IEEE, 2016.
- [35] Yichun Shi, Anil K Jain, and Nathan D Kalka. Probabilistic face embeddings. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [36] Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 464–472. IEEE, 2017.
- [37] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deeply learned face representations are sparse, selective, and robust. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2892–2900, 2015.
- [38] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- [39] Fei Wang, Liren Chen, Cheng Li, Shiyao Huang, Yanjie Chen, Chen Qian, and Chen Change Loy. The devil of face recognition is in the noise. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 765–780, 2018.
- [40] Feng Wang, Jian Cheng, Weiyang Liu, and Haijun Liu. Additive margin softmax for face verification. *IEEE Signal Processing Letters*, 25(7):926–930, 2018.
- [41] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5265–5274, 2018.
- [42] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515. Springer, 2016.
- [43] Lior Wolf, Tal Hassner, and Itay Maoz. Face recognition in unconstrained videos with matched background similarity. In *CVPR 2011*.
- [44] Xiang Wu, Ran He, Zhenan Sun, and Tieniu Tan. A light cnn for deep face representation with noisy labels. *IEEE Transactions on Information Forensics and Security*, 13(11):2884–2896, 2018.
- [45] Weidi Xie, Li Shen, and Andrew Zisserman. Comparator networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 782–797, 2018.
- [46] Weidi Xie and Andrew Zisserman. Multicolumn networks for face recognition. In *BMVC*, 2018.
- [47] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014.
- [48] Bangjie Yin, Luan Tran, Haoxiang Li, Xiaohui Shen, and Xiaoming Liu. Towards interpretable face recognition. *arXiv preprint arXiv:1805.00611*, 2018.
- [49] Xi Yin and Xiaoming Liu. Multi-task convolutional neural network for pose-invariant face recognition. *IEEE Transactions on Image Processing*, 27(2):964–975, 2017.
- [50] Tianyuan Yu, Da Li, Yongxin Yang, Timothy M Hospedales, and Tao Xiang. Robust person re-identification by modelling feature uncertainty. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 552–561, 2019.
- [51] Umara Zafar, Mubeen Ghafoor, Tehseen Zia, Ghufuran Ahmed, Ahsan Latif, Kaleem Razzaq Malik, and Abdullahi Mohamud Sharif. Face recognition with bayesian convolutional networks for robust surveillance systems. *EURASIP Journal on Image and Video Processing*, 2019(1):10, 2019.