TRAINING PRIVACY-PRESERVING VIDEO ANALYTICS PIPELINES BY SUPPRESSING FEATURES THAT REVEAL INFORMATION ABOUT PRIVATE ATTRIBUTES

Chau Yi Li, Andrea Cavallaro

Centre for Intelligent Sensing, Queen Mary University of London, UK

ABSTRACT

Deep neural networks are increasingly deployed for scene analytics, including to evaluate the attention and reaction of people exposed to out-of-home advertisements. However, the features extracted by a deep neural network that was trained to predict a specific, consensual attribute (e.g. emotion) may also encode and thus reveal information about private, protected attributes (e.g. age or gender). In this work, we focus on such leakage of private information at inference time. We consider an adversary with access to the features extracted by the layers of a deployed neural network and use these features to predict private attributes. To prevent the success of such an attack, we modify the training of the network using a confusion loss that encourages the extraction of features that make it difficult for the adversary to accurately predict private attributes. We validate this training approach on image-based tasks using a publicly available dataset. Results show that, compared to the original network, the proposed PrivateNet can reduce the leakage of private information of a state-ofthe-art emotion recognition classifier by 2.88% for gender and by 13.06% for age group, with a minimal effect on task accuracy.

Index Terms— feature learning, video analytics, privacy, private attributes.

1. INTRODUCTION

Deep neural network classifiers are used to understand the reactions of audiences to advertisements in public spaces and during interactions with wayfinding kiosks. Such reactions are typically measured using computer vision by first detecting faces and then analysing attributes associated with facial expressions. However, neural networks trained to extract features that predict a task-specific, consensual attribute (the target task), may also encode information about private attributes [1], such as gender and age [2]. An adversary with access to the information extracted by the layers of a deep neural network can therefore use the pipeline for purposes different from the consensual ones.

An adversary may perpetrate a *membership* or *attribute* inference attack. Privacy-preserving methods to protect from membership inference attacks focus on preventing the identi-

fication of individuals whose data were included in a training dataset [3-6] or are ingested by a deployed system. Methods that protect from attribute inference attacks aim to prevent the prediction of (private) attributes from a training dataset [6] or from a classification pipeline deployed on a system, such as a smart kiosk. The latter scenario is the scope of our work: we consider a system whose untrusted rich execution environment is compromised by malware, but is equipped with security measures that ensure confidentiality and integrity of the visual data (e.g. a limited-memory, isolated Trusted Execution Environment, or TEE). The memory limitation of a TEE can only host a few of the last layers of a feed-forward network, typically the layers after the backbone structure or only the fully-connected layer [7]. We consider the adversary to have access to the untrusted execution environment of the system and thus in a position to exploit the features extracted by these intermediate layers to infer private information [8]. Existing works that aim to prevent attribute inference attacks modify the input data with adversarial perturbations [9] or train the network to cause mis-classification of private attributes, for instance using the cross-entropy loss [10].

In this work, we aim to prevent the estimation of private information at inference time by training the network to extract features that are useful for the target task only and that cause a known adversary classifier to perform similarly to a *random classifier* on a protected attribute. To conceal the private information from the adversary classifier, we train the pipeline with a confusion loss. The proposed approach does not modify the network architecture or the number of its parameters, and hence has no latency impact on the inference pipeline. In the specific implementation reported in this paper, we consider visual emotion recognition as the target task, and age and gender as protected attributes¹. We also investigate the robustness of the proposed approach by training an adversary to classify the private information from the privacypreserving network.

¹The source code is available at the following URL: https://github.com/smartcameras/PrivateNet.

2. METHOD

Let $C(\cdot)$ be a *D*-class deep neural network classifier of *N* layers that, given an image *x*, predicts its *consensual* class as one in the set:

$$\mathcal{Y} = \{y_1, y_2, \dots, y_D\}.$$
 (1)

Let $\mathcal{F}_i(x)$ be the output of layer *i* of $C(\cdot)$, with $i \in \{1, ..., N\}$. The output of the last layer N is a vector that represents the confidence of C that x belongs to any of the classes:

$$\mathcal{F}_N(x) = (p_1^y, p_2^y, \dots, p_D^y). \tag{2}$$

Let $\mathcal{A}(\cdot)$ be an adversary K-class classifier that aims to predict a *protected*, *private* class from the output, $\mathcal{F}_i(x)$, of an intermediate layer as one in the set:

$$\mathcal{S} = \{s_1, s_2, \dots, s_K\}.\tag{3}$$

Let M be the number of layers in $\mathcal{A}(\cdot)$. The output of the last layer M is a vector representing the confidence of $\mathcal{A}(\cdot)$ that $\mathcal{F}_i(x)$ belongs to any of the (private) classes:

$$\mathcal{Q}_M(\mathcal{F}_i(x)) = (p_1^s, p_2^s, \dots, p_K^s).$$
(4)

Let \mathcal{X} be a set of face images, each annotated with an emotion (target) attribute, \hat{y} , and a private attribute, \hat{s} . We measure the leakage of private information as the accuracy, $T_{\hat{s}}$, of the private attribute inference from \mathcal{X} :

$$T_{\hat{s}} = \frac{|\{x \in \mathcal{X} : \mathcal{A}(\mathcal{F}_i(x)) = \hat{s}\}|}{|\mathcal{X}|},$$
(5)

where $|\cdot|$ is the cardinality of a set. The higher $T_{\hat{s}}$, the higher the leakage of private information through $\mathcal{C}(\cdot)$.

A classifier $C(\cdot)$ that is privacy-preserving should maintain a high accuracy in the prediction of the consensual task, while concealing the private attributes from the adversary classifier $A(\cdot)$, i.e. $T_{\hat{s}}$ should be close to a random classifier.

To predict the consensual class, the classifier $C(\cdot)$ is typically trained with a cross-entropy loss:

$$\mathcal{L}_{CE}\Big(\hat{y}, \mathcal{F}_N(x)\Big) = -\log\left(\frac{exp(p_{\hat{y}})}{\sum_i exp(p_i^y)}\right), \qquad (6)$$

where $exp(\cdot)$ is the exponential function.

To prevent the leakage of private information associated with a protected attribute, we propose to obfuscate with a confusion loss, \mathcal{L}_{con} , the features in the intermediate layers that are useful to the adversary:

$$\mathcal{L}_{con}\Big(\mathcal{F}_{i}(x)\Big) = \|\mathcal{Q}_{M}(\mathcal{F}_{i}(x)) - \mathcal{U}_{D}\|^{2}, \tag{7}$$

where $\mathcal{U}_D = (\frac{1}{D}, ..., \frac{1}{D})$ denotes equal probability for each private attribute class, hence causing the adversary to perform similarly to a random classifier. As the features extracted by an intermediate layer embeds the features extracted by



Fig. 1. Architecture of the ARM network [11] used as the emotion recognition classifier. Given an image x, the network outputs the probability p_i^y , $i \in \{1, ..., 7\}$, of each emotion class (Eq. 2). An adversary could infer private, protected attributes from the features extracted by the intermediate layers of the network (red dashed arrows). Coloured blocks represent layers with trainable parameters.

the layers preceding it, backpropagating the confusion loss through the entire network optimizes the layers preceding the targeted layer to extract generic features that cause the adversary to perform similarly to a random classifier on the protected attribute. Thus, in training we encourage the deep neural network to extract privacy-preserving features by combining Eq. 6 and 7 as the overall loss function, \mathcal{L} :

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_{CE} + \lambda\mathcal{L}_{con}, \tag{8}$$

where λ determines the relative importance between the losses: $\lambda < 0.5$ gives more importance to maintaining the utility and $\lambda > 0.5$ gives more importance to confusing the adversary.

We use as dataset the Real-world Affective Faces Database (RAF-DB) [12], which contains 15,338 images. Each image was annotated, on average, by 40 annotators into D = 7emotion attributes (surprise, fear, disgust, happiness, sadness, anger, and neutral) and demographics including race, gender and age group. Specifically, we focus on gender $(K = 2, \text{ male and female})^2$ and age group $(K = 5, \text{ namely})^2$ 0-3, 4-19, 20-39, 40-69, and over 70). As for the emotion classifier $\mathcal{C}(\cdot)$, we use the Amend Representation Module network proposed by Shi and Zhou [11]. On RAF-DB, the ARM network attains a state-of-the-art emotion recognition accuracy of 91.10%. The ARM network aims to remove the distortion from the edges of the image on the features, referred to as albino erosion. The ARM network consists of a ResNet-18 [13] backbone, an Amend Representation Module (ARM), and a fully connected layer. The output of the ResNet-18 backbone is a feature map of size 7×7 and 512 channels. An ARM consists of 3 blocks, namely a feature

 $^{^2 \}rm We$ excluded the images with *unsure* gender class, which contribute to 6.3% of the dataset.

Table 1. The leakage of private information is measured as the accuracy of an adversary classifier that infers private attributes from the features extracted from the ARM network [11], trained for emotion recognition. Note that the accuracy of the adversary is lower than that of a network classifying the original image but higher than that of a random classifier.

Attributo	Random	Baseline on	Adversary on feature		
Attribute	classifier	original image	Backbone	ARM	
Gender	50	88.22	74.97	63.30	
Age	20	77.35	68.20	57.20	

rearrangement block, a convolution layer that is followed by batch normalisation, and a sharing affinity block. The rearrangement block distributes the backbone feature map into 2 channels, each of size 112×112 , while maintaining the relative positions between the features in the same channel [14], whereas the sharing affinity block ensures the global average representation of faces in the mini-batch is propagated. For this paper, we choose to study the leakage of private information from the features extracted by the ResNet-18 backbone and that by the ARM, which represent layers that cannot be executed in the TEE. Fig. 1 shows the network architecture and the features studied in this paper.

3. RESULTS

In this section, we discuss the utility and the privacy level of the proposed privacy-preserving network, PrivateNet. The *utility* of the classifier network is its accuracy on emotion recognition. The *privacy level* provided for a private attribute by the network is how close the accuracy of an adversary classifier is to a random guess on that attribute. Furthermore, the *robustness* of the network is the difference in accuracy between the adversary accessing the features of the original network (Tab. 1) and the privacy-preserving networks. The more negative the percentage difference, the more robust the privacy-preserving network.

To establish a fair comparison with the private information that can be inferred from the original image and to ensure that any difference is due to the input and not to the change in the architecture, we use the same ARM network architecture in our experiments. The adversary for the backbone features consists of an ARM and a fully connected layer, whereas that for ARM features consists of a fully connected layer only. All networks are trained on an Nvidia Tesla V100-SXM2 GPU, using a learning rate of 0.001 with ADAM optimiser, with 200 epochs. The networks were tested on the same GPU.

Tab. 1 reports the accuracy of the adversaries on features extracted through different layers of the networks to infer two private attributes, namely gender and age. A random classifier would achieve an accuracy of 50% for gender (K = 2) and 20% for age (K = 5). The classification accuracy of

Table 2. Utility, privacy and robustness of privacy-preserving networks trained with an adversarial loss [10] and with the proposed confusion loss (Eq. 7). The level of privacy guaranteed for an attribute is measured as the accuracy of an adversary classifier on the features of the network. The lower the accuracy, the more difficult to infer private information from the features. Robustness is the difference in the accuracy of the adversary classifier on the features of the original network (Tab. 1) and the privacy-preserving network, reported in relative percentage difference. The more negative the difference, the less susceptible (i.e. more robust) the privacy-preserving network to the adversary classifier. The results of the most robust network are shown in bold. KEY – Att.: attribute; Adv.: adversarial; Prop.: proposed; K: known adversary; U: unknown adversary.

Att.	Feature	Loss	Utility	Privacy		Robustness	
				K	U	K	U
Gender	Backbone	Adv.	88.43	56.47	88.50	-24.68%	+18.04%
		Prop.	89.40	51.27	72.81	-31.61%	-2.88%
	ARM	Adv.	89.51	43.53	87.84	-36.17%	+37.99%
		Prop.	89.47	49.91	62.43	-26.82%	+1.75%
Age	Backbone	Adv.	89.63	10.72	62.58	-83.06%	-8.24%
		Prop.	89.83	20.18	59.29	-68.12%	-13.06%
	ARM	Adv.	89.47	10.72	54.17	-81.26%	-5.30%
		Prop.	89.86	21.28	56.88	-62.80%	-0.56%

the ARM network architecture, using the original images as baselines, is 88.22% for gender and 73.53% for age. The leakage of private information decreases along the network, as the network extracts features that are more relevant to the target task. In particular, the accuracy of the adversaries on the gender attribute drops from 74.97%, with features extracted by the backbone, to 63.30%, with features extracted by the ARM. Nonetheless, the adversary classifier can still predict the private attributes with high accuracy from the extracted features. For example, the adversary classifier who has access to the backbone feature can predict the gender attribute with 73.75% accuracy (83.59% of the baseline accuracy) and the age attribute with 68.20% accuracy (88.17% of the baseline).

Tab. 2 compares the utility, privacy and robustness of 8 networks, trained with 2 loss functions, namely the proposed confusion loss and the adversarial loss by Edwards and Storkey [10], against 2 known adversaries on the backbone and ARM features, for 2 private attributes, namely gender and age. In the rest of the paper, we refer to the networks in the format *loss-feature-attribute*, hence *confusion-feature-attribute* is an example of the proposed PrivateNet. The utility of the 8 networks ranges from 88.43% to 89.86%, corresponding to a slight drop (1.36% to 2.93%) from the original ARM network (accuracy of 91.10%). Therefore training the network in a privacy-preserving way only slightly reduces the network performance for the consensual task.

For the binary (K = 2) gender attribute, the accuracy of



Fig. 2. Impact of the relative weighting of the cross-entropy loss and the proposed confusion loss (λ in Eq. 8) on the accuracy of the consensual task (utility) and of adversaries that use the backbone features to predict the age group (K = 5). Note that $\lambda = 1$ means that the network is not trained for the target task and therefore we do not report the accuracy of the unknown adversary for this value. Dashed lines show the ideal behaviours, i.e. maximum utility for the consensual attribute and random results for the private attribute.

the known adversary on PrivateNet is 51.27% on confusionbackbone-gender and 49.91% on confusion-ARM-gender. Thus the confusion loss has achieved the goal of obfuscating the known adversary to perform similarly to a random classifier (50%). However, the high accuracy of the unknown adversary (72.81% and 62.43% on backbone and ARM features, respectively) indicates that the networks are not robust against the attack of an adversary that is unknown at the time of training. Training the network with adversarial loss forces the known adversary to mis-classify the binary attribute (from male to female and from female to male), hence the accuracy of the known adversary on adversarial-backbonegender is lower than that on confusion-backbone-gender (43.53%). However, the accuracy of an unknown adversary on adversarial-backbone-gender (87.84%) is higher than that of confusion-backbone-gender, and also higher than that of the original network (68.20%). Moreover, while training with the confusion loss has increased the robustness by 2.88% on confusion-backbone-gender and decreased it by 1.75% on confusion-ARM-gender, training with the adversarial loss decreases the robustness against an unknown adversary by 18.04% and 37.99% on adversarial-backbone-gender and adversarial-ARM-gender, respectively. Therefore, training with an adversarial loss in fact encourages the network to extract features that are more useful for inferring the binary attribute, hence reducing the ability of the network to protect the attribute.

For the age group attribute (K = 5), training with the

Table 3. Inference time of the privacy-preserving network, reported as average \pm standard deviation (in milliseconds), on a RAF-DB image. These times are similar to those of the original network (10.76 \pm 0.34 milliseconds).

Protected	Against adversary on			
attribute	backbone	ARM		
Gender	10.86 ± 0.34	10.82 ± 0.35		
Age	10.77 ± 0.32	10.68 ± 0.34		

adversarial loss protects the private attribute against a known adversary better than training it with the confusion loss: the accuracy of the adversary is 10.72% on *adversarial-backbone-age* and 20.18% on *confusion-backbone-age*. However, the accuracy of an unknown adversary is 62.58% on *adversarial-backbone-age* and 59.29% on *confusion-backbone-age*. Overall, the proposed PrivateNet are more robust than networks trained with the adversarial loss.

Fig. 2 reports the results obtained when varying the relative weight of the cross-entropy and confusion losses (λ in Eq. 8) on the consensual task and the robustness of the proposed privacy-preserving network. Note that $\lambda = 0$ is equivalent to the original ARM network that is optimised for emotion recognition only; increasing λ gives more importance to the confusion loss; and $\lambda = 1$ discards the cross-entropy loss and hence the network is optimised for misleading the known adversary only. As λ increases, the utility (emotion recognition accuracy) decreases, as expected. While training with most values of λ can cause the accuracy of the known adversary to be close to that of a random classifier, the robustness of the PrivateNet increases with λ . Tab. 3 reports the inference time of the proposed networks on RAF-DB. As the numbers of parameters are the same, the inference times of the networks are similar to that of the original ARM network $(10.76 \pm 0.34 \text{ milliseconds}).$

To summarise, PrivateNet maintains comparable utility in emotion recognition, while protecting the private attributes from known and unknown adversaries. Also, this approach to privacy preservation has no impact on the inference time of the network.

4. CONCLUSION

We addressed the problem of private information leakage through the features extracted by the intermediate layers of a deep learning classifier. Unlike works that use an adversarial loss to cause the mis-classification of a private attribute, we obfuscate its associated features using a confusion loss. The proposed approach was validated in a scenario where the goal is to conceal age group and gender attributes from a known adversary with access to the output of the layers of an emotion recognition network. The proposed PrivateNet reduces the accuracy of the adversary to close to that of a random classifier, with negligible effects on the accuracy of the target task. Moreover, the proposed confusion loss is preferable to the adversarial loss in reducing the leakage of private information with an unknown adversary classifier.

Future work will consider more granular private attributes and the protection of multiple private attributes in a single network.

5. REFERENCES

- Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov, "Machine learning models that remember too much," in *Proceedings of ACM SIGSAC Conference* on Computer and Communications Security, November 2017, p. 587–601.
- [2] Hwansoo Lee, Siew Fan Wong, Jungjoo Oh, and Younghoon Chang, "Information privacy concerns and demographic characteristics: Data from a Korean media panel survey," *Government Information Quarterly*, vol. 36, no. 2, pp. 294–303, 2019.
- [3] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov, "Membership inference attacks against machine learning models," in *Proceedings of IEEE Symposium on Security and Privacy*, May 2017, pp. 3–18.
- [4] Milad Nasr, Reza Shokri, and Amir Houmansadr, "Machine learning with membership privacy using adversarial regularization," in *Proceedings of ACM SIGSAC Conference on Computer and Communications Security*, October 2018, p. 634–646.
- [5] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes, "ML-leaks: Model and data independent membership inference attacks and defenses on machine learning models," *Annual Network and Distributed System Security Symposium*, February 2019.
- [6] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov, "Exploiting unintended feature leakage in collaborative learning," *Proceedings of IEEE Symposium on Security and Privacy*, pp. 691–706, May 2018.
- [7] Fan Mo, Ali Shahin Shamsabadi, Kleomenis Katevas, Soteris Demetriou, Ilias Leontiadis, Andrea Cavallaro, and Hamed Haddadi, "Darknetz: towards model privacy at the edge using trusted execution environments," *Proceedings of the 18th International Conference on Mobile Systems, Applications, and Services*, p. 161–174, June 2020.
- [8] Mohammad Malekzadeh, Anastasia Borovykh, and Deniz Gündüz, "Honest-but-curious nets: Sensitive attributes of private inputs can be secretly coded into the

classifiers' outputs," in *Proceedings of ACM SIGSAC Conference on Computer and Communications Security*, November 2021.

- [9] Chau Yi Li, Ali Shahin Shamsabadi, Ricardo Sanchez-Matilla, Riccardo Mazzon, and Andrea Cavallaro, "Scene privacy protection," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, May 2019, pp. 2502–2506.
- [10] Harrison Edwards and Amos Storkey, "Censoring representations with an adversary," in *Proceedings of International Conference in Learning Representations*, May 2016, pp. 1–14.
- [11] Jiawei Shi and Songhao Zhu, "Learning to amend facial expression representation via de-albino and affinity," *arXiv:2103.10189*, September 2021.
- [12] Shan Li and Weihong Deng, "Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 356–370, September 2019.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of IEEE Conference on Computer Vision* and Pattern Recognition, June 2016, pp. 770–778.
- [14] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2016, pp. 1874–1883.