

Toward Face Biometric De-identification using Adversarial Examples

Mahdi Ghafourian, Julian Fierrez, Luis F. Gomez, Ruben Vera-Rodriguez, Aythami Morales

Universidad Autonoma de Madrid,

Email: {mahdi.ghafourian, julian.fierrez, luisf.gomez, ruben.vera, aythami.morales}@uam.es

Zohra Rezgui, Raymond Veldhuis

Universtity of Twente, Email: z.rezgui@utwente.n, r.n.j.veldhuis@utwente.nl

Abstract—The remarkable success of face recognition (FR) has endangered the privacy of internet users particularly in social media. Recently, researchers turned to use adversarial examples as a countermeasure to privacy attacks. In this paper, we assess the effectiveness of using two widely known adversarial methods (BIM and ILLC) for de-identifying personal images. We discovered, unlike previous claims in the literature, that it is not easy to get a high protection success rate (suppressing identification rate) with imperceptible adversarial perturbation to the human visual system. Finally, we found out that the transferability of adversarial examples is highly affected by the training parameters of the network with which they are generated.

Index Terms—Artificial Intelligence, Face Biometrics, De-identification, Adversarial Attacks

I. INTRODUCTION

Deep learning has evolved as a strong and efficient tool to be applied to a broad spectrum of complex learning problems that were difficult to solve using traditional machine learning. The development of deep convolutional neural networks (CNNs) has been so revolutionary that today it can exceed human-level performance. As a consequence, deep networks are being extensively used in many recent applications including face recognition. Now, face recognition (FR) systems have become an exceptionally accurate technology in identifying people from images [1], [2]. While being useful, face recognition may invade the privacy of individuals [3] when used to exploit and process illicitly their face images [4], [5] and videos [6], [7] found on the internet, particularly social media.

In recent years, several reports revealed unauthorized collections of large datasets of identified face data from social media. Reports on Cambridge Analytica [8] in 2018, and Clearview AI in 2020 [9] are glaring examples of privacy leakage related to face biometrics. So far, the most common defense against this threat has been to set all social media profiles to ‘private’, allowing only chosen friends access to your images [10].

To mitigate these privacy threats, some studies [11], [12] turned to generate adversarial perturbations called cloaks to de-identify face biometrics in personal images before uploading them to social media. These perturbations are being generated by applying a very slight (imperceptible to human eyes) modification to the input and optimizing it to maximize the probability of misclassification by a machine learning classifier [13]. Using attacks to preserve privacy in biometrics has

attracted attention [14], [15] which also includes adversarial examples.

In another line of work, instead of introducing imperceptible artifacts at the raw image level to harden automatic identification, one can operate at the feature level by disentangling there the identification information and reducing it while preserving other information of interest (e.g., facial emotions [16], soft biometrics [17], etc.) See the work by Morales et al. [18] and the references therein for further information in this line.

In the present paper, we conduct an experimental evaluation of the effectiveness of two popular adversarial methods, i.e. Basic Iterative Method (BIM) and Iterative Least Likely Class (ILLC) [19], for de-identifying face biometrics in personal photos at the raw image level. In particular, we focussed on the transferability of the de-identified face biometrics across different classifiers. To this end, we used three popular pre-trained face recognition models (*FaceNet*, *ResNet-50*, and *SENet-50*) interchangeably to create an adversarial example by one model and defend against it using all three models.

By analyzing the quantitative results of BIM and ILLC methods, we obtained some important findings. First, it is not likely to obtain a high protection success rate together with quite imperceptible adversarial perturbation. In particular, when it comes to black-box scenarios and any preprocessing (e.g. image compression, resizing) that affects the adversarial trigger, this goal would be ambitious. Second, we discuss that the definition of feature embeddings of the adversarial class are highly dependent on the other training classes in the attacker network. Therefore, the transferability of generated adversarial examples (i.e. de-identified personal images) conforms with the similarity of the attacker network to that of the defender in terms of training parameters. Third, unlike our expectation, although the BIM method is an untargeted method (i.e. adversarial method without an specific target), it is more protective than the targeted ILLC method.

II. PROTECTION MODEL

In this section, we introduce the protector’s goal, capabilities, and knowledge under which the de-identified samples are generated. Since the goal of our study is to preserve the privacy using adversarial examples, we call the party who generates the examples the protector and the party whose network is used for classifying the examples, the invader. For a better

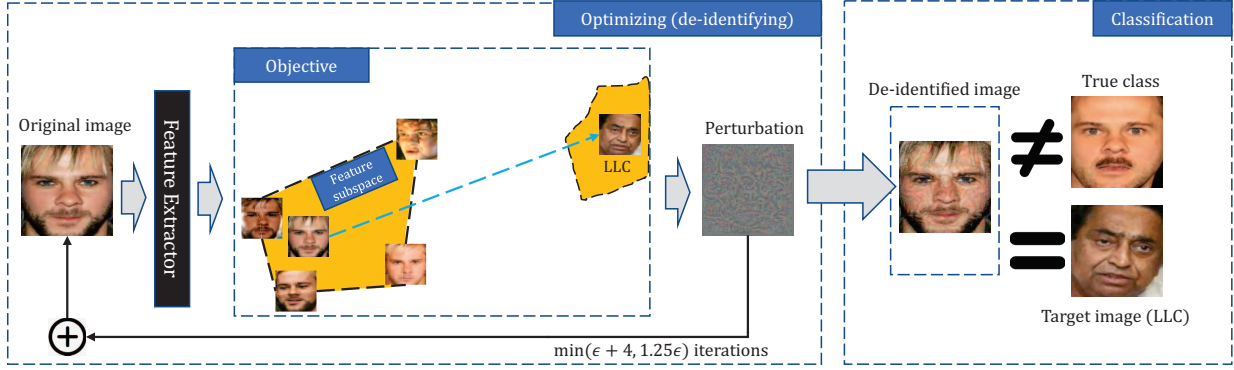


Fig. 1: Overview of the targeted adversarial examples to de-identify face images.

understanding of the paper, we provide definitions from their original sources with which we conducted our experiments. Therefore, in the remaining of the paper, we use the following notations:

- x : the input face biometric of the identity who wants to be de-identified. It is an RGB image in the shape of a 3D tensor ($width \times height \times depth$) whose values range is in $[0, 255]$.
- x_{adv} : the adversarial example (i.e. de-identified image) for x .
- y_{true} : the true class label for the image x .
- y_{target} : the target class label that the attacker is trying to optimize the input image to fool the defender classifier with, in our case the least likely class (y_{LLC}).
- ϵ the noise budget to add to one pixel of X .
- $C(x)$: it denotes the classifier $C(x) : X \rightarrow Y$ where $x \in X \subset \mathbb{R}^d$, and $y = \{1, 2, \dots, N\}$ with N being the total number of classes.
- $J(x, y_{target})$: the cross-entropy cost function for computing the loss of x given the target class label y_{target} .
- $Clip_{\epsilon}\{x_{adv}\}$: clipping function to confine the alteration of each pixel in the de-identified image x_{adv} to the noise budget ϵ to keep the result in the L_P ϵ -neighbourhood of the input image x .

In general, adversarial methods are divided into two categories: *Untargeted*, where $C(x_{adv}) \neq y_{true}$ and *Targeted*, where $C(x_{adv}) = y_{target}$ (see Fig. 1).

A. Protector's goal

The goal of the protector is to craft a constrained adversarial perturbation to de-identify face biometrics in their personal image. To this end, the protector adds a small perturbation measured by L_P norm to the original face biometric in a specific number of iterations. For the adversarial method we used, the upper bound of this number of iterations is determined by $\min(\epsilon + 4, 1.25\epsilon)$.

B. Protector's capability

To achieve the goal, the generated adversarial examples must satisfy $\|x_{adv} - x\|_p \leq \epsilon$ to mislead the model of the

privacy invader. Therefore, the protector is able to conduct the following optimization problems in the aforementioned number of iterations according to the method he adopts. Regarding the untargeted methods, the protector generates the de-identified face by maximizing the cost function $J(x_{adv}, y_{true})$:

$$x_{adv} = \operatorname{argmax}_{x_{adv}: \|x_{adv} - x\|_p \leq \epsilon} J(x_{adv}, y_{true}) \quad (1)$$

while for the targeted method, de-identified face images are crafted by minimizing the cost function $J(x_{adv}, y_{target})$:

$$x_{adv} = \operatorname{argmin}_{x_{adv}: \|x_{adv} - x\|_p \leq \epsilon} J(x_{adv}, y_{target}) \quad (2)$$

C. Protector's knowledge

Similar to the real-world scenarios, we conducted our assessment in a black-box setting. In black-box attacks, it is assumed that the protector has no prior knowledge of the invader's network or its parameters. With this assumption, the protector can only acquire the classification output of the invader model. Therefore, in an oracle attack, the protector evaluates the protection success rate by providing crafted inputs with various perturbation budgets. However, the protector can use the same dataset for generating adversarial examples with which the invader's model has been trained.

III. GENERATING DE-IDENTIFIED FACES

The aim of de-identification on face biometrics is to preserve the privacy of the subjects by protecting their true identity against unwanted face identifications. To this end, the use of adversarial perturbations through a technique called Image Cloaking has been proposed recently. Shan et al. [11] proposed a method called Fawkes, improving image cloaking by reducing the effectiveness of face recognition software while preserving the quality of the image to human eyes. This method is a targeted approach choosing k random target classes, picking the centroid G_k of them, and selecting the most dissimilar class T to that of the user's face x by computing:

$$T = \operatorname{argmax}_{k=1..K} \min_{x \in X} \operatorname{Dist}(C(x), G_k) \quad (3)$$

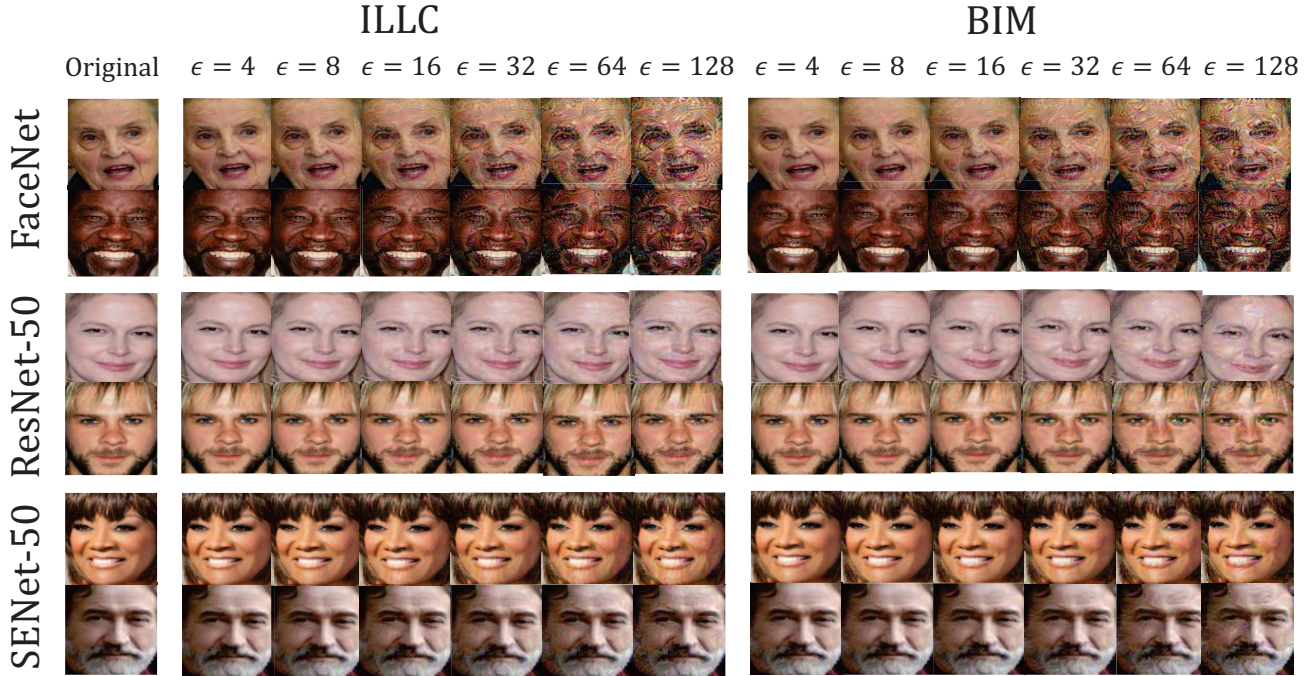


Fig. 2: Example of de-identified face images for all the models with various perturbation budgets (ϵ)

Another similar work called LowKey [12] did the image cloaking by updating x_{adv} iteratively adding the gradient toward the maximization objective. They applied Gaussian smoothing to maintain the quality of the image reducing the accuracy of Amazon Rekognition to 32.5% (i.e. 67.5% protection rate). In the current paper, we generate de-identified face images with various perturbation budgets using BIM and ILLC adversarial methods as it is shown in Fig. 2.

A. Basic Iterative Method (BIM)

According to [20], the simplest method to generate an adversarial image is to find the perturbation that maximizes the cost function with respect to a L_∞ constraint with just one back-propagation (FGSM method). Later, [19] extended that method by doing back-propagation iteratively while clipping values changes in pixels after each iteration to keep the alteration to the ϵ -neighbourhood of the original image. This kind of iterative hill-climbing attacks have for long being studied in the biometrics security literature [21], [22]. This method is called BIM and the adversarial image in each iteration is crafted as below:

$$x_{adv}^{(i+1)} = \text{Clip}_\epsilon \{ x_{adv}^{(i)} + \alpha \cdot \text{sign}(\nabla_{x_{adv}^{(i)}} J(x_{adv}^{(i)}, y_{true})) \} \quad (4)$$

where α is the step size and $x_{adv}^{(0)} = x$ at the initialization of BIM method. Therefore, by maximizing the cost, the classification result of the de-identified face image x_{adv} would lie far from the original image x .

B. Iterative Least Likely Class (ILLC)

Unlike BIM, the only difference of this method is to reduce the cost but toward a specific target. In this case, the target is the least likely class when the original image is classified. As a result, the crafted de-identified face will avoid the original image mistakenly identified as another person in the classification database. However, the effectiveness of this method for de-identification relies on the dissimilarity rate of all the subjects in the training dataset. This method is also an iterative method initiated with $x_{adv}^{(0)} = x$ and the adversarial image in each iteration is crafted as below:

$$x_{adv}^{(i+1)} = \text{Clip}_\epsilon \{ x_{adv}^{(i)} - \alpha \cdot \text{sign}(\nabla_{x_{adv}^{(i)}} J(x_{adv}^{(i)}, y_{LLC})) \} \quad (5)$$

IV. EVALUATION

A. Evaluation metric

So far, the most common metric that has been used to evaluate the performance of adversarial examples is transferability. This metric denotes that the examples produced to deceive a particular model can be used to deceive other models regardless of the underlying architecture. To estimate the transferability of the generated adversarial examples we use the protection success rate also called the suppressing identification rate. In our case, it would be the misclassification rate of the de-identified faces by the target classifier. Thus, given the adversarial method Adv_ϵ to generate the de-identified face image as $x_{adv} = \text{Adv}_\epsilon(x)$ for the input face x under the constraints of perturbation budget ϵ and l_p -norm,

and target classifier $C(x)$, the Protection Success Rate (PSR) is defined as:

$$\text{PSR}(\text{Adv}_\epsilon, C) = 100 - \left(\frac{100}{N} \sum_{i=1}^N 1(C(\text{Adv}_\epsilon(x_i)) = y_{\text{true}}) \right) \quad (6)$$

where N is the number of test samples and $1(\cdot)$ is the indicator function. The higher the PSR, the more resilient the example is to be identified in the target classifier.

B. Evaluation settings

Our experiments are divided into two phases: *Generating the de-identified image* of the input face in the source network by the protector, and *Classifying the example* in target networks to evaluate the Protection Success Rate (PSR). To this end, we used three widely used pre-trained face recognition models (all trained on the VGGFace2 dataset [23]): FaceNet [1], ResNet-50 [2], SENet-50 [24].

We start the process of generating de-identified faces in the source network as follows:

- First, we select N random subjects from the VGGFace2 dataset to protect their identity.
- Second, the perturbation budget ϵ is picked from the $set_\epsilon = \{4, 8, 16, 32, 64, 128\}$ [19]. In terms of transferability, we will assess the proportion of Protection Success Rate with respect to the image quality degradation. The ideal output is to achieve the largest PSR using the smallest possible perturbation budget.
- Third, the number of iterations for optimizing the input face toward the adversarial goal is calculated as $n_{\text{iter}} = \min(\epsilon + 4, 1.25 \times \epsilon)$.
- Finally, for every $Model \in \{\text{FaceNet}, \text{ResNet-50}, \text{SENet-50}\}$ as source network, for each random input face $x \in \{x_i\}_{i=1}^N$, and for every $\epsilon \in set_\epsilon$, we iterate the input image x by n_{iter} doing backpropagation toward y_{target} for ILLC method and y_{true} for BIM method.

Some examples of de-identified face images regarding both adversarial methods for each $\epsilon \in set_\epsilon$ are depicted in Fig. 2.

Once the de-identified face is crafted, for each $Model \in \{\text{FaceNet}, \text{ResNet-50}, \text{SENet-50}\}$ as target networks, we assess the PSR of the crafted examples via the following steps:

- First, the face is extracted using MTCNN [25] to check if the perturbation makes the face undetectable.
- Second, the detected face is fed to the classifier of the selected $Model$.
- Third, based on the classification maximum probability, we compute Top1, Top5, Top10, Top25, Top50 where $C(\text{Adv}_\epsilon(x_i)) = y_{\text{true}}$.
- Finally, for each top, we calculate PSR according to Eq. 6.

The resulting PSR for the n_{iter} corresponding to each $\epsilon \in set_\epsilon$ for *FaceNet*, *ResNet-50*, and *SENet-50* is depicted in Figs. 3, 4, 5 respectively.

C. Evaluation results

To obtain our results, we crafted examples on one model per experiment then we evaluated them against all networks

independently. To assess the effect of compression to the adversarial trigger, all the input faces are fed into networks uncompressed, and crafted adversarial examples are stored with JPEG compression. Another important aspect that we included in our investigation is the effect of resizing crafted examples. FaceNet is different from the other two networks in terms of input image size. While FaceNet accepts images with size 160×160 , ResNet and SENet accept 224×224 . This means that de-identified faces experience image resizing when they are crafted in FaceNet as source network and classified in ResNet and SENet as target network and vice versa. Looking at Figs. 3, 4, 5, the first apparent understanding that spring to mind is that all adversarial examples crafted using a specific source model (FaceNet, ResNet, or SENet) transfer particularly well when considering identification based on the same recognition model. In addition, it is clear that the examples generated by FaceNet are more transferable compared to those crafted by ResNet and SENet. Comparing Fig. 3 with Figs. 4, 5, it can be seen that examples crafted by FaceNet using BIM method at $\epsilon = 32$ reported high transferability as they are highly protective when they were classified by the other two networks.

It is also obvious that, in all figures, when the perturbation budget increases (i.e. as the quality of the image is decreasing due to adding more noise), the protection success rate increases as well, but at the cost of sacrificing image quality. Considering these charts, we see that BIM (an untargeted approach) outperforms ILLC (a targeted method). Taking into account Fig. 3, Top-25 charts, it can be noticed that while in BIM chart at $\epsilon = 32$, $\text{PSR} \geq 95\%$ for ResNet and SENet while the corresponding ones for ILLC are $\text{PSR} \leq 65\%$.

These results show that using BIM and ILLC adversarial methods to preserve privacy for face images can only be achievable with $\epsilon > 32$ at the cost of degrading the quality of the image. It also indicates that the protection success rate of the crafted examples is highly affected by resizing the examples and the difference of training parameters between source and target networks. Finally, these results point out that untargeted methods need further attention as in our experiments BIM performed better than the ILLC.

V. CONCLUSION

Using adversarial methods to de-identify face biometrics, it is likely that untargeted method are more protective than targeted ones. Yet, further studies are needed to prove this hypothesis. Besides, using these two methods, it's not possible to get a high de-identification with completely imperceptible perturbation. That's why most of the current literature suggests keeping the balance between the suppressing identification rate and the image quality. To this end, in our future study, we will focus on the effectiveness and transferability of less destructive adversarial methods to preserve the quality of the image including one-pixel attack, jacobian-based saliency map attack (JSMA), and deepfool.

Future work will also explore image quality aspects [26] of perturbed images related to their biometric content [27],

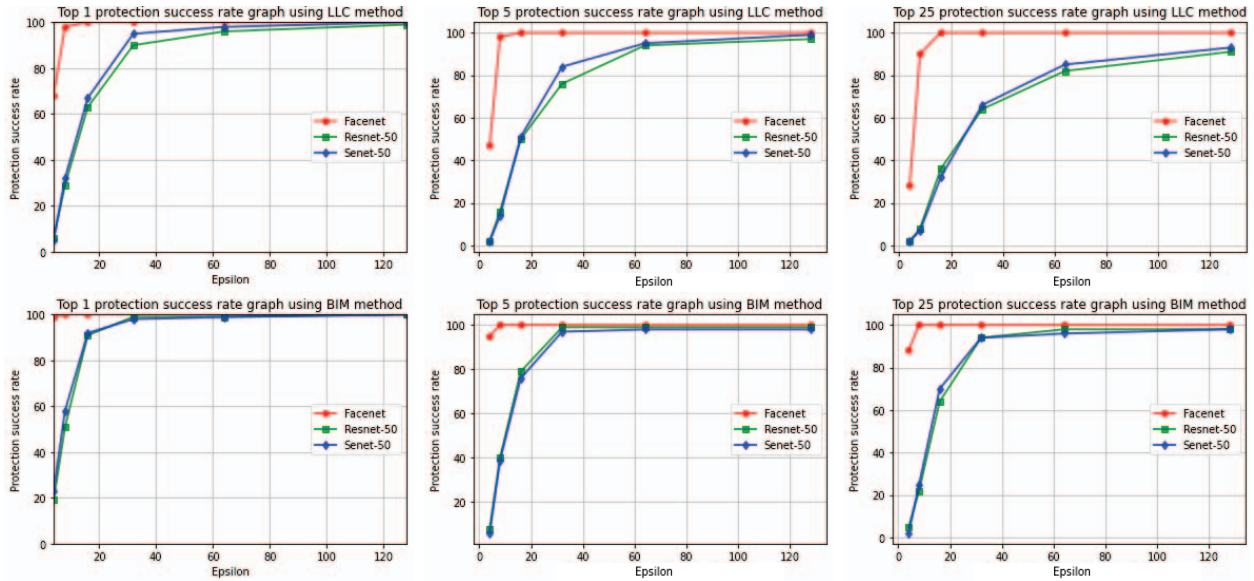


Fig. 3: Protection Success Rate (PSR) as perturbation budget increases for adversarial examples crafted using *FaceNet*. First row: LLC method (left to right: Top1, Top5, Top25). Second row: BIM method.

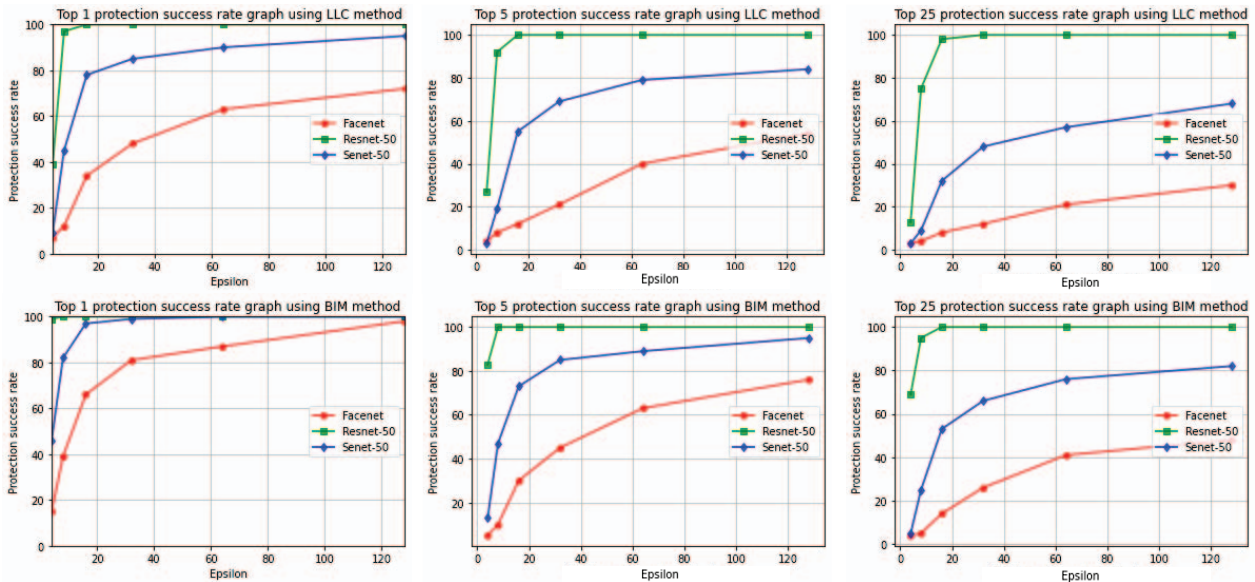


Fig. 4: Protection Success Rate (PSR) as perturbation budget increases for adversarial examples crafted using *ResNet*. First row: LLC method (left to right: Top1, Top5, Top25). Second row: BIM method.

[28] and human perception, and based on that we will try to optimize in a multimodal machine learning setup [29] the protection success rate as target function conditioned to human perception and biometric content restrictions.

ACKNOWLEDGMENT

This work has been supported by projects: PRIMA (H2020-MSCA-ITN-2019-860315), TRESPASS-ETN (H2020-MSCA-ITN-2019-860813), and BBforTAI (PID2021-127641OB-I00

MICINN/FEDER). Mahdi and Zohra are PRIMA Early-Stage Researchers. Luis is a TRESPASS-ETN Early-Stage Researcher.

REFERENCES

- [1] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 815–823, 2015. 1, 4
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE CVPR*, pp. 770–778, 2016. 1, 4

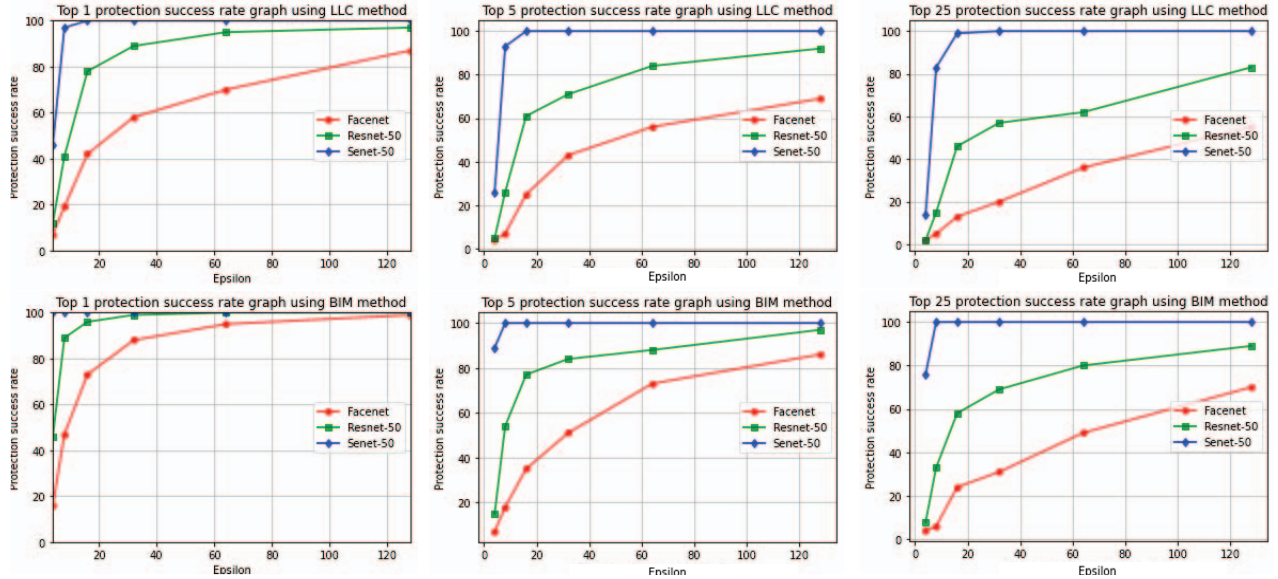


Fig. 5: Protection Success Rate (PSR) as perturbation budget increases for adversarial examples crafted using *SENet*. First row: LLC method (left to right: Top1, Top5, Top25). Second row: BIM method.

- [3] A. Hassanpour, M. Moradikia, B. Yang, A. Abdelhadi, C. Busch, and J. Fierrez, "Differential privacy preservation in robust continual learning," *IEEE Access*, vol. 10, pp. 24273–2428, February 2022. 1
- [4] A. Hadid, N. Evans, S. Marcel, and J. Fierrez, "Biometrics systems under spoofing attack: an evaluation methodology and lessons learned," *IEEE Signal Processing Magazine*, vol. 32, pp. 20–30, Sept. 2015. 1
- [5] J. Hernandez-Ortega, J. Fierrez, et al., *Handbook of Biometric Anti-Spoofing*, ch. Introduction to Presentation Attack Detection in Face Biometrics and Recent Advances. Springer, 2023. 3rd Ed. 1
- [6] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, *Handbook of Digital Face Manipulation and Detection*, ch. An Introduction to Digital Face Manipulation, pp. 3–26. Springer, 2022. 1
- [7] R. Tolosana, S. Romero-Tapiador, R. Vera-Rodriguez, E. Gonzalez-Sosa, and J. Fierrez, "Deepfakes detection across generations: Analysis of facial regions, fusion, and performance evaluation," *Engineering Applications of Artificial Intelligence*, vol. 110, April 2022. 1
- [8] A. Samuel, "The shady data-gathering tactics used by Cambridge Analytica were an open secret to online marketers. I know, because I was one." <https://www.theverge.com/2018/3/25/17161726/facebook-cambridge-analytica-data-online-marketers>, 2018. 1
- [9] K. Hill, "The secretive company that might end privacy as we know it," in *Ethics of Data and Analytics*, pp. 170–177, Auerbach, 2020. 1
- [10] B. S. Ledford, *An Assessment of Image-Cloaking Techniques Against Automated Face recognition for Biometric Privacy*. PhD thesis, Florida Institute of Technology, Melbourne, Florida, 2021. 1
- [11] S. Shan, E. Wenger, J. Zhang, H. Li, H. Zheng, and B. Y. Zhao, "Fawkes: Protecting privacy against unauthorized deep learning models," in *USENIX Security Symposium*, pp. 1589–1604, 2020. 1, 2
- [12] V. Cherepanova, M. Goldblum, H. Foley, S. Duan, J. Dickerson, G. Taylor, and T. Goldstein, "Lowkey: Leveraging adversarial attacks to protect social media users from facial recognition," *arXiv preprint arXiv:2101.07922*, 2021. 1, 3
- [13] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrnđić, P. Laskov, G. Giacinto, and F. Roli, "Evasion attacks against machine learning at test time," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 387–402, Springer, 2013. 1
- [14] M. Gomez-Barrero, J. Galbally, and J. Fierrez, "Efficient software attack to multimodal biometric systems and its application to face and iris fusion," *Pattern Recognition Letters*, vol. 36, pp. 243–253, Jan. 2014. 1
- [15] M. Ghafourian, J. Fierrez, et al., "OTB-morph: one-time biometrics via morphing applied to face templates," in *IEEE/CVF Winter Conf. on Applications of Computer Vision*, pp. 321–329, 2022. 1
- [16] A. Peña, J. Fierrez, A. Lapedriza, and A. Morales, "Learning emotional-blinded face representations," in *ICPR*, January 2021. 1
- [17] E. Gonzalez-Sosa, J. Fierrez, R. Vera-Rodriguez, and F. Alonso-Fernandez, "Facial soft biometrics for recognition in the wild: Recent works, annotation and cots evaluation," *IEEE Trans. on Information Forensics and Security*, vol. 13, pp. 2001–2014, August 2018. 1
- [18] A. Morales, J. Fierrez, R. Vera-Rodriguez, and R. Tolosana, "SensitiveNets: learning agnostic representations with application to face recognition," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 43, pp. 2158–2164, June 2021. 1
- [19] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *Artificial Intelligence Safety and Security*, pp. 99–112, Chapman and Hall/CRC, 2018. 1, 3, 4
- [20] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014. 3
- [21] J. Galbally, C. McCool, J. Fierrez, S. Marcel, and J. Ortega-Garcia, "On the vulnerability of face verification systems to hill-climbing attacks," *Pattern Recognition*, vol. 43, pp. 1027–1038, March 2010. 3
- [22] M. Gomez-Barrero, J. Galbally, J. Fierrez, and J. Ortega-Garcia, "Face verification put to test: a hill-climbing attack based on the uphill-simplex algorithm," in *Proc. Intl. Conf. on Biometrics*, March 2012. 3
- [23] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "VGGFace2: A dataset for recognising faces across pose and age," in *IEEE Intl. Conf. on Automatic Face & Gesture Recognition (FG)*, pp. 67–74, 2018. 4
- [24] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *IEEE Conf. on Computer Vision and Pattern Recognition*, 2018. 4
- [25] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016. 4
- [26] T. Schlett, C. Rathgeb, O. Henniger, J. Galbally, J. Fierrez, and C. Busch, "Face image quality assessment: A literature survey," *ACM Computing Surveys*, vol. 54, no. 10, pp. 1–49, 2022. 4
- [27] F. Alonso-Fernandez, J. Fierrez, and J. Ortega-Garcia, "Quality measures in biometric systems," *IEEE Security & Privacy*, vol. 10, pp. 52–62, December 2012. 4
- [28] J. Hernandez-Ortega, J. Galbally, J. Fierrez, and L. Beslay, "Biometric quality: Review and application to face recognition with FaceQnet," *arXiv 2006.03298*, 2021. 4
- [29] A. Peña, I. Serna, A. Morales, J. Fierrez, et al., "Human-centric multimodal machine learning: Recent advances and testbed on AI-based recruitment," *SN Computer Science*, 2023. 5