

SWIS: SELF-SUPERVISED REPRESENTATION LEARNING FOR WRITER INDEPENDENT OFFLINE SIGNATURE VERIFICATION

Siladitya Manna[†]

Soumitri Chattopadhyay^{*}

Saumik Bhattacharya[◊]

Umapada Pal[†]

[†] Indian Statistical Institute, Kolkata; [◊] Indian Institute of Technology, Kharagpur; ^{*} Jadavpur University

ABSTRACT

Writer independent offline signature verification is one of the most challenging tasks in pattern recognition as there is often a scarcity of training data. To handle such data scarcity problem, in this paper, we propose a novel self-supervised learning (SSL) framework for writer independent offline signature verification. To our knowledge, this is the first attempt to utilize self-supervised setting for the signature verification task. The objective of self-supervised representation learning from the signature images is achieved by minimizing the cross-covariance between two random variables belonging to different feature directions and ensuring a positive cross-covariance between the random variables denoting the same feature direction. This ensures that the features are decorrelated linearly and the redundant information is discarded. Through experimental results on different data sets, we obtained encouraging results.

Keywords - Self-supervised, Cross-covariance, Decorrelation, Writer-independent, SVM

1. INTRODUCTION

Signature verification has been used as one of the most essential steps for identity verification of person-specific documents like forms, bank cheques, or even the individual themselves. This makes signature verification an important task in domain of computer vision and pattern recognition. There are mainly two types of signature verification processes: (1) offline and (2) online. In offline signature verification, the input is basically a 2D image which is scanned from the original signature or captured into an image by some electronic device. Whereas, in online signature verification, the writer usually pens down his signature on an electronic tablet using a stylus and the information is recorded at some regular timestep along with the position of the stylus.

Offline signature verification can again be divided into two types: (1) Writer dependent and (2) writer independent. In writer dependent scenario, the system needs to be updated and retrained for every new user signature that gets added to the system. This makes the process cumbersome and less feasible. However, in writer independent scenario, a generalized system needs to be built which can differentiate between genuine and forged signatures without repeated retraining.

Most researchers have leveraged supervised learning methods [1–6] for offline signature verification. While hand-crafted feature analyses have comprised the bulk of studies in this domain [6–9], various deep learning-based methods have also been proposed, particularly dwelling on metric learning approaches [1–4]. Nevertheless, all the aforementioned works are fully supervised methods and therefore, share the common bottleneck of data scarcity. To this end, we demonstrate the first use of self-supervision for offline signature verification.

Self-supervised learning aims at developing a pre-training paradigm to learn a robust representation from an unlabelled corpus for generalization to any given downstream task. Widely studied in recent years, several pretext tasks have been proposed, such as solving jigsaw puzzles [10], image colorization [11] to name a few. Contrastive learning based self-supervised algorithms, like SimCLR [12], MoCo [13] has also gained popularity, which aim at learning similarity between augmented views of the same image while distancing views from different images. [14] aimed at simultaneously maximizing similarity and minimizing redundancy between embeddings of two distorted views of an image.

In this work, we propose a self-supervised learning algorithm for offline writer-independent signature verification. Self-supervised learning is a sub-domain of unsupervised learning that aims at learning representations from the data without any ground truth or human annotations. As a skilled forgery is supposed to be very close to the genuine signature, it is necessary to distinguish between each constituting element of the signatures for correct classification. However, since it is not possible to obtain a large number of annotated genuine signatures from the individuals for training a large model, we use self-supervised learning for training the model to learn representations which are generalized for signatures over a large number of individuals. This work is the first of its kind to apply self-supervised learning framework for learning representations from signature images. Also, in the downstream stage, we do not use any siamese type architecture in the downstream task for the offline signature verification, and show the capability of the pretrained encoder to effectively cluster the genuine signatures of the different unknown writers.

The main contributions of this work are as follows:

- A novel self-supervised approach is introduced here for offline writer independent signature verification purpose.
- To the best of our knowledge, this is the first work of the use of self-supervised learning in signature verification.
- We have shown that the proposed SSL is better than the state-of-the art self-supervised contrastive learning approaches used in Computer vision and Medical image analysis areas.

The rest of the paper is organized as follows. Sec. 2 describes the self-supervised learning methodology that is used in this work. Sec. 3 presents the details about the datasets we use. In Sec. 4, we present the experimental results and the comparison with the base models. Finally, we conclude the paper in Sec. 5.

2. METHODOLOGY

In this section, we discuss the pre-processing and the algorithm steps that are used to train the proposed encoder.

2.1. Pretraining Methodology

In signature images, it is essential to capture the stroke information from the different authors as well as to learn the variations in the signatures of the same individual. To feed the stroke information without any human supervision, we divided the signature images into patches of dimensions 32×32 with an overlap of 16 pixels from a signature image reshaped to 224×224 . This gives 169 patches from a single image of dimensions 32×32 . As the base encoder we choose ResNet-18 [15]. When the patches are passed through the encoder, we obtain an output of $1 \times 1 \times 512$ from each patch. We rearrange the patches into a grid of 13×13 to obtain an output of shape $13 \times 13 \times 512$. After applying global average pooling (GAP), we obtain an output feature vector of dimension 1×512 . This feature vector is then passed through a non-linear projector with 1 hidden layer and output dimension 512 to obtain the final output.

For forming positive pairs, we augment a single signature image in two randomly chosen augmentations. The augmentation details are mentioned in Sec. 3.2. The images are then divided into patches as mentioned before and then passed through the encoder and the projector.

Thus, the proposed loss function has the form:

$$\mathcal{L}_C = \frac{1}{N} \sum_{i=1}^D \left(\sum_{\substack{j=1 \\ j \neq i}}^D \left(\sum_{k=1}^N z_k^i \cdot z_k^j \right)^2 + \left(\sum_{k=1}^N z_k^i \cdot z_k^i - 1 \right)^2 \right) \quad (1)$$

where z_k^i is a scalar value at i -th dimension of the k -th centered and normalized feature vector z_k . Thus, the pre-processing steps before feeding the feature vector z_k^i to the loss function are as follows

$$\bar{z}_k^i = \frac{\tilde{z}_k^i}{\sqrt{\sum_{k=1}^N (z_k^i)^2}} \quad \forall i \in [1, D] \quad (2)$$

$$z_k^i = \bar{z}_k^i - \mu_{z_k}, \quad \text{where } \mu_{z_k} = \frac{1}{N} \sum_{k=1}^N \bar{z}_k^i \quad \forall i \in [1, D]$$

It is to be noted that z_k^i and z_k^j are obtained from the each element of a positive pair. Thus, the proposed loss function does not optimize the terms of a cross-covariance matrix in the true meaning of the term. We can refer to this matrix as a Pseudo cross-covariance matrix.

From eq. 1, we can see that optimizing the proposed loss function allows us to decorrelate the dimensions of the output. We treat each dimension as a random variable Z_i . As Z_i is the output feature vector from the last Batch Normalization layer in the projecto, $Z_i \sim \mathcal{N}(0, 1)$. Normalizing Z_i and subtracting mean along each dimension in Eqn. 2, bring the feature vectors inside a unit hyper-sphere S^D , where D is the dimension of the feature vector, and centers each dimension at 0, i.e., $Z_i \sim \mathcal{N}(0, \sigma_i^2)$. Since, we are making the cross-covariance matrix to an Identity matrix,

$$Cov(Z_i, Z_j) = 0 \Rightarrow \rho = 0 \quad (3)$$

For Normal Random Variables Z_i ,

$$\mathbb{E}[Z_i, Z_j] = \mathbb{E}[Z_i] \cdot \mathbb{E}[Z_j] \quad \forall i, j \in [1, D] \wedge i \neq j \quad (4)$$

The diagonal terms of the cross-covariance matrix are optimised such that it equates to 1. Hence, the PDF of the feature vectors $f_{Z_1, \dots, Z_D} \sim \mathcal{N}(0, \mathcal{I}_{D \times D})$. Consequently, each output dimension becomes independent.

2.2. Pretraining Model Architecture

The model architecture used in the pretraining phase is given in Figure 1. The diagram shows the input that is fed to the ResNet18 [15] encoder. The input is reshaped to $169 \times 32 \times 32 \times 3$ before passing it through the encoder. Figure 1 also shows an example of the input used in the pretraining phase.

2.3. Downstream Evaluation

For predicting whether a signature is forged or genuine, we take 8 reference signature for each user and use them to train a Support Vector Machine (SVM) classifier with radial basis function kernel. We assume that the user for which the signature is being verified is known. We also assume that the forged signature will be mapped outside the decision boundary of that particular user. If the user is predicted correctly and

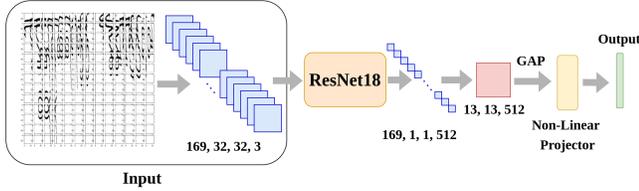


Fig. 1: Model architecture used in the pretraining phase of the proposed method.

the signature is genuine, we count it as a correct prediction. Similarly, if the predicted user is not correct and the signature is actually forged, then also it is counted as a correct prediction. In all the other cases, the prediction is considered as wrong.

By using a SVM classifier, we depend on the feature extraction capability of the pretrained encoder to express the input in terms of its linearly decorrelated factors. Whereas all the contemporary state-of-the-art supervised algorithms use siamese type architecture or supervised contrastive learning framework for the offline signature verification task.

3. EXPERIMENTAL DETAILS

In this section, we are going to discuss the details of the datasets that were used in our experiments, and the configurations used for training our encoder in the pretext (or pretraining) task.

3.1. Datasets

In this work, we used two datasets, namely, BHSig260 [17] and ICDAR 2011 [16]. BHSig260 dataset contains signatures from 100 writers for Bengali and 160 writers for Hindi signatures. For each writer of both the languages, there are 24 genuine and 30 forged signatures. Among the 100 writers in the Bengali subset, we randomly select 50 writers for the training set and the rest 50 are used for testing. For the Hindi subset, we randomly selected 50 writers for self-supervised pretraining and the rest 110 writers were left for testing. Similarly, for ICDAR 2011 Signature Verification dataset, there are signatures for Dutch and Chinese languages. The subset of the Dutch signatures contains signatures from 10 writers for training and 54 writers for testing. In the test set, however, there are 8 reference genuine signatures for each writer. To adhere to this structure, we randomly selected 8 genuine signatures from the test set of BHSig260 dataset for each writer and used it as the reference set, for both Bengali and Hindi languages.

3.2. Pretraining Experiments Configuration

For the pretraining phase, we used different number of epochs for different datasets. The models were trained by optimizing the loss function given by 1 using LARS [19] optimizer. We

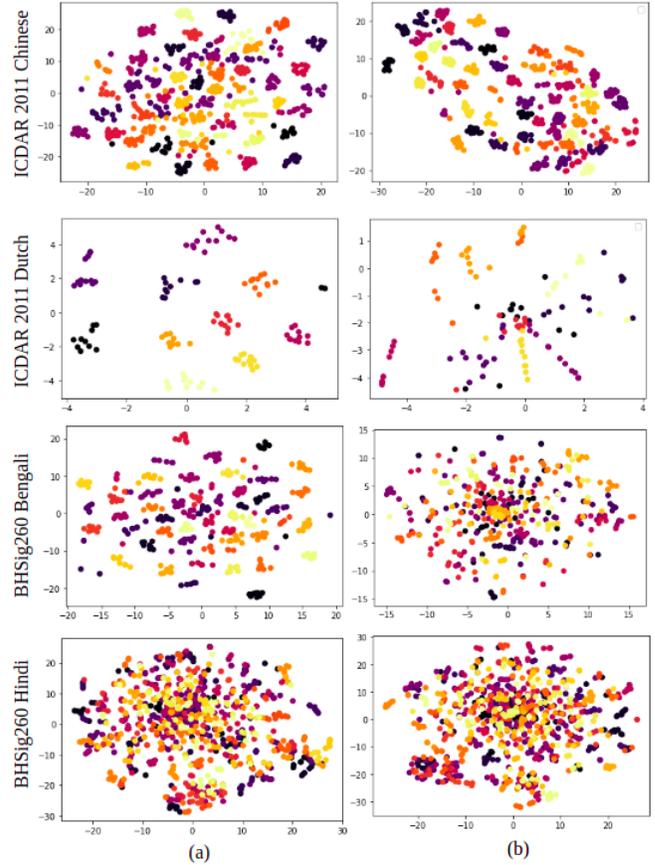


Fig. 2: t-SNE visualisations obtained by (a) the proposed method compared with those obtained by (b) SimCLR [12] on different datasets. The color coding scheme denotes each writer cluster. .

used a learning rate of 0.1 and a momentum value of 0.9. The batch-normalization and bias parameters were excluded from weight normalization. We decayed the learning rate following a cosine decay schedule with a linear warmup period of 10 epochs at the start. The decay was scheduler for 1000 epochs irrespective of the number of training epochs.

For the ICDAR datasets, we pretrained the model for 500 epochs. Whereas for the BHSig260 dataset, the pretraining was carried out for 200 epochs only. For both the datasets, the batch size used was 32.

To ensure that the pretrained models learn generalized and robust features, we applied several augmentations, such as, color jittering, affine transformation and random cropping to 224×224 . The images obtained after augmentation were normalized to the range $[-1.0, +1.0]$.

As not all images in the datasets contain perfectly cropped signature images, we cropped the images such that the input to the encoder contained is a tightly bounded signature image. To achieve this objective, we performed Otsu's thresholding [20] followed by finding the bounding box with least

Table 1: Comparison of the proposed method with state-of-the-art self-supervised learning baselines.

Method	ICDAR 2011 Dutch [16]			ICDAR 2011 Chinese [16]			BHSig260 Bengali [17]			BHSig260 Hindi [17]		
	Accuracy (%)	FAR	FRR	Accuracy (%)	FAR	FRR	Accuracy (%)	FAR	FRR	Accuracy (%)	FAR	FRR
SimCLR [12]	69.46	0.554	0.060	59.76	0.431	0.317	73.45	0.117	0.543	72.45	0.103	0.599
Proposed	77.62	0.316	0.133	64.68	0.278	0.583	72.04	0.367	0.116	72.43	0.104	0.598

Table 2: Comparison of the proposed method with supervised learning methods in literature.

Method	BHSig260 Bengali [17]			BHSig260 Hindi [17]		
	Accuracy (%)	FAR	FRR	Accuracy (%)	FAR	FRR
Pal et al. [17]	66.18	0.3382	0.3382	75.53	0.2447	0.2447
Dutta et al. [18]	84.90	0.1578	0.1443	85.90	0.1310	0.1509
Dey et al. [2]	86.11	0.1389	0.1389	84.64	0.1536	0.1536
Alaei et al. [7]	–	0.1618	0.3012	–	0.1618	0.3012
Proposed	72.04	0.367	0.116	72.43	0.104	0.598

area containing all non-zero pixels around the centre of mass of the image. After this preprocessing step, the images were divided into patches of dimension 32×32 with an overlap of 16 pixels and fed to the encoder for training.

4. EXPERIMENTAL RESULTS

4.1. Downstream Results

The downstream task we considered in our work is the writer-independent classification of signatures into two classes: genuine or forged. The predictions were obtained using the procedure described in Section 2.3. The results obtained by the proposed model in the downstream task on the datasets ICDAR 2011 and BHSig260 signature verification datasets are given in Table 1. We also pre-trained and validated our proposed method on GPDS300 [21] and CEDAR [22] dataset, and we achieved accuracies of 69.28% and 83.8%, respectively.

4.2. Ablation on Hyperparameters

We tested the robustness of the representations learnt by our proposed model using Gaussian noise(AWGN) with $\mu = 0.0$, $\sigma^2 = 0.01$ and obtained accuracy(ACC), FAR and FRR of 76.84% ($\sigma = 0.26533$), 0.3242 ($\sigma = 0.005$) and 0.17 ($\sigma = 0.003$), respectively for the CEDAR dataset. Using Random cropping, we obtained ACC, FAR and FRR of 79.3% ($\sigma = 0.94$), 0.344 ($\sigma = 0.0124$) and 0.1157 ($\sigma = 0.0128$), respectively. We also consider ablation on projector depth, augmentation and patch overlap on the CEDAR dataset. Increasing the overlap of patches from 0 to 8 pixels shows accuracy(ACC), FAR and FRR of 83.8%, 0.118 and 0.187, respectively. Increasing the number of layers in the projector did not improve the performance. Removing color jitter as augmentation from the above model yielded ACC, FAR and FRR of 83.1%, 0.11 and 0.19, respectively.

4.3. Comparison with SOTA Self-supervised Algorithms

In this section, we show how the proposed loss function fares at training the encoder to learn representations from the data. As shown in Table 1, in spite of trained in a self-supervised manner, the proposed framework performs satisfactorily on both the multilingual datasets. Table 1 also presents the comparative results of one of the state-of-the-art self-supervised algorithm (SimCLR) on the same data. From Fig. 2, we can see that the proposed algorithm performs better at producing distinct clusters for ICDAR 2011 Chinese and BHSig260 Bengali dataset, whereas the plots for ICDAR 2011 Dutch and BHSig260 Hindi datasets look equally well-clustered for both the proposed model and SimCLR. It should be mentioned here that the SimCLR algorithm was trained for 1000 epochs on the ICDAR 2011 dataset (both, Dutch and Chinese).

4.4. Comparison with Supervised Methods

To further validate our proposed self-supervised pipeline, we compare its performance with some fully supervised methods in literature. The results have been tabulated in Table 2. We observe that the proposed framework performs competitively against the fully supervised works on the BHSig260 datasets, outperforming [17] by a large margin on the Bengali signature dataset. Moreover, the low FAR and FRR values obtained by the proposed method on the signature datasets affirm its potential in separating forged signatures from the genuine ones.

5. CONCLUSION

In this work, we proposed a self-supervised representation learning framework where a novel loss function is used that aims at decorrelating the dimensions from each other to discard redundant features and encourage learning of linearly uncorrelated generative features of the input. Through t-SNE plots we show that the proposed algorithm extracts better uncorrelated information from the input than the SOTA SSL methods on the same datasets. From the comparative results, it is evident that the proposed method performs better than or at par with the state-of-the-art algorithm SimCLR. This work shows the extensive scope and applicability of the proposed method in the field of signature verification and paves a way for further research in this direction.

6. REFERENCES

- [1] H. Rantzsch, H. Yang, and C. Meinel, "Signature embedding: Writer independent offline signature verification with deep metric learning," in *International symposium on visual computing*. Springer, 2016, pp. 616–625. 1
- [2] S. Dey, A. Dutta, J. I. Toledo, S. K. Ghosh, J. Lladós, and U. Pal, "Signet: Convolutional siamese network for writer independent offline signature verification," *arXiv preprint arXiv:1707.02131*, 2017. 1, 4
- [3] V. Ruiz, I. Linares, A. Sanchez, and J. F. Velez, "Off-line handwritten signature verification using compositional synthetic generation of signatures and siamese neural networks," *Neurocomputing*, vol. 374, pp. 30–41, 2020. 1
- [4] Q. Wan and Q. Zou, "Learning metric features for writer-independent signature verification using dual triplet loss," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 3853–3859. 1
- [5] E. Parcham, M. Ilbeygi, and M. Amini, "Cbcapsnet: A novel writer-independent offline signature verification model using a cnn-based architecture and capsule neural networks," *Expert Systems with Applications*, vol. 185, p. 115649, 2021. 1
- [6] A. K. Bhunia, A. Alaei, and P. P. Roy, "Signature verification approach using fusion of hybrid texture features," *Neural Computing and Applications*, vol. 31, no. 12, pp. 8737–8748, 2019. 1
- [7] A. Alaei, S. Pal, U. Pal, and M. Blumenstein, "An efficient signature verification method based on an interval symbolic representation and a fuzzy similarity measure," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 10, pp. 2360–2372, 2017. 1, 4
- [8] L. G. Hafemann, R. Sabourin, and L. S. Oliveira, "Offline handwritten signature verification—literature review," in *2017 Seventh International Conference on Image Processing Theory, Tools and Applications (IPTA)*. IEEE, 2017, pp. 1–8. 1
- [9] D. Banerjee, B. Chatterjee, P. Bhowal, T. Bhattacharyya, S. Malakar, and R. Sarkar, "A new wrapper feature selection method for language-invariant offline signature verification," *Expert Systems with Applications*, vol. 186, p. 115756, 2021. 1
- [10] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *European conference on computer vision*. Springer, 2016, pp. 69–84. 1
- [11] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *European conference on computer vision*. Springer, 2016, pp. 649–666. 1
- [12] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607. 1, 3, 4
- [13] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738. 1
- [14] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy reduction," in *ICML*, 2021. 1
- [15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. 2
- [16] G. Alvarez, B. Sheffer, and M. Bryant, "Offline signature verification with convolutional neural networks," *Technical report, Stanford University*, 2016. 3, 4
- [17] S. Pal, A. Alaei, U. Pal, and M. Blumenstein, "Performance of an off-line signature verification method based on texture features on a large indic-script signature dataset," in *2016 12th IAPR workshop on document analysis systems (DAS)*. IEEE, 2016, pp. 72–77. 3, 4
- [18] A. Dutta, U. Pal, and J. Lladós, "Compact correlated features for writer independent signature verification," in *2016 23rd international conference on pattern recognition (ICPR)*. IEEE, 2016, pp. 3422–3427. 4
- [19] Y. You, I. Gitman, and B. Ginsburg, "Large batch training of convolutional networks," *arXiv preprint arXiv:1708.03888*, 2017. 3
- [20] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Transactions on Systems, Man and Cybernetics*, pp. 62–66, 1979. 3
- [21] J. F. Vargas-Bonilla, M. A. Ferrer, C. M. Travieso, and J. B. Alonso, "Off-line handwritten signature GPDS-960 corpus," in *9th International Conference on Document Analysis and Recognition*. IEEE Computer Society, 2007, pp. 764–768. 4
- [22] M. K. Kalera, S. N. Srihari, and A. Xu, "Offline signature verification and identification using distance statistics," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 18, no. 7, pp. 1339–1360, 2004. 4