

FEW-SHOT LEARNING DRIVEN BY CONTRASTIVE TECHNIQUES

ANVAYA RAI



BHARTI SCHOOL OF TELECOMMUNICATION
TECHNOLOGY AND MANAGEMENT

INDIAN INSTITUTE OF TECHNOLOGY DELHI

FEBRUARY 2024

© Indian Institute of Technology Delhi (IITD), New Delhi, 2024

FEW-SHOT LEARNING DRIVEN BY CONTRASTIVE TECHNIQUES

by

ANVAYA RAI

BHARTI SCHOOL OF TELECOMMUNICATION
TECHNOLOGY AND MANAGEMENT

submitted

in fulfillment of the requirements of the degree of
Doctor of Philosophy
to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI
FEBRUARY 2024

*Dedicated to
My Teachers
Family
&
Friends*

THESIS CERTIFICATE

This is to certify that the thesis titled **Few-shot Learning driven by Contrastive Techniques**, submitted by **Anvaya Rai**, to the Indian Institute of Technology, Delhi, for the award of the degree of **Doctor of Philosophy**, is a bonafide record of the research work done by him under my supervision. He has fulfilled the requirements for the submission of the thesis, which to the best of our knowledge has reached the required standard.

The material contained in the thesis has not been submitted in part or full to any other University or Institute for the award of any degree or diploma.

Prof. Brejesh Lall

Department of Electrical Engineering
Indian Institute of Technology Delhi
New Delhi -110016

Date: FEBRUARY 2024

ACKNOWLEDGEMENTS

This thesis is the outcome of a wonderful seven year working experience under the guidance of **Dr. Brejesh Lall**, to whom I wish to express my appreciation and gratitude. His vision, creativeness, and enthusiasm have not just fueled this research but also have a lasting influence in my career as a role model. I would also like to express my deep sense of gratitude towards **Dr. R. K. P. Bhatt** and **Dr. S. D. Joshi** for his support and guidance whenever required during the entire course work.

I take this opportunity to thank my friends and colleagues at IIT Delhi and C-DOT for their generous help and fruitful collaborations, especially Ms. Astha Zalani Jain and Mr. Raghwendra Prakash Singh. I would specially like to thank my friends, Mr. Dharmesh Kumar, Mr. Adarsh Tripathi and Mr. Alok Sharma, for constantly boosting my moral and motivating me in my research.

My deepest gratitude goes to my wife **Ms. Pragya**. She had given me the much needed comfort and inspiration for completing this thesis.

Finally, I would like to dedicate this work to my parents, whose constant support and encouragement have really brought me here. Lastly, I am eternally grateful to the entire universe for the love, support, and opportunities it has bestowed upon me to make a positive impact on the world.

Anvaya Rai

ABSTRACT

KEYWORDS: Hyperspherical Manifold; Extreme Learning Machines; Principal Component Analysis Extreme Learning Machine; Linear Discriminant Analysis Extreme Learning Machine; Few-shot; Contrastive Learning; Visual SLAM; Vision based navigation; Transfer Learning; Manifold Generation.

The few-shot learning concept using contrastive techniques has come into existence because the Deep neural networks (DNN) demand huge amounts of manually annotated training data. In many scenarios, this type of data acquisition is highly tedious and at times impossible. Thus, there was a need to learn from few examples, using few iterations and augmentations. In this work, we have proposed a novel contrastive learning based few-shot technique and applied it in the area of Face Recognition. We call it as ***Latent Feature Transformed Contrastive Learning (LFT-CLR)***. The suggested framework can be extended to other domains as well. Additionally, we also propose a novel contrastive loss, called ***Weighted Normalized Temperature-scaled Cross Entropy Loss (wNT-Xent)***, to refine a pre-trained DNN. The suggested approach has been validated using finely curated datasets based upon LFW [20].

In the proposed inference pipeline, the fine-tuned DNN is followed by PCA/ LDA based latent feature space transformation, resulting in extracting a manifold from the feature space to project the test samples. This feature space transformation helps in dimensionality reduction of the latent feature space and projection of query images onto an optimally dense subspace of the original latent feature space, leading to accurate and faster inference.

In addition to Face Recognition, we also extend this work in the domain of Visual SLAM and introduce a novel approach to accurately localize a subject in indoor environments by using the scene images captured from the subject’s mobile phone camera. We present a novel deep neural network (DNN), called ***InPosNet***, that generates a concise representation of an indoor scene while being able to distinguish between their inherent symmetry. It also enables the user in real time distinction between the images of the same location but captured from different orientations, there by enabling the user to detect the orientation along with position. The novel DNN presented in the work is motivated by MobileNetv3-Small [54]. A localization accuracy of less than 1 meter from ground truth is achieved and enumerated through the experimental results. The goal is to present a vision based system that will have

the ability to be used for indoor positioning, without any need for additional infrastructure.

अमूर्त

प्रमुख शब्द: हाइपरस्फेरिकल मैनिफोल्ड; चरम शिक्षण मशीनें; प्रमुख घटक विश्लेषण चरम शिक्षण मशीन; रैखिक विभेदक विश्लेषण चरम शिक्षण मशीन; कुछ-शॉट; विरोधाभासी शिक्षा; विजुअल स्लैम; दृष्टि आधारित नेविगेशन; स्थानांतरण सीखना; अनेक गुना पीढ़ी.

तकनीकों का उपयोग करके कुछ-शॉट सीखने की अवधारणा अस्तित्व में आई है क्योंकि डीप न्यूरल नेटवर्क (डीएनएन) बड़ी मात्रा में मैनुअल रूप से एनोटेड प्रशिक्षण डेटा की मांग करते हैं। कई परिदृश्यों में, इस प्रकार का डेटा अधिग्रहण अत्यधिक कठिन और कभी-कभी असंभव होता है। इस प्रकार, कुछ पुनरावृत्तियों और संवर्द्धनों का उपयोग करके कुछ उदाहरणों से सीखने की आवश्यकता थी। इस कार्य में, हमने एक नवीन विरोधाभासी शिक्षण आधारित कुछ-शॉट तकनीक का प्रस्ताव रखा है और इसे चेहरा पहचान के क्षेत्र में लागू किया है। इसे हम **वेडेड सिमसीएलआर** कहते हैं। सुझाए गए ढांचे को अन्य डोमेन तक भी बढ़ाया जा सकता है। इसके अतिरिक्त, हम पूर्व-प्रशिक्षित डीएनएन को परिष्कृत करने के लिए एक नवीन विपरीत हानि का भी प्रस्ताव करते हैं, जिसे **भारित सामान्यीकृत तापमान-स्केल्ड क्रॉस एन्ट्रॉपी लॉस** कहा जाता है। सुझाए गए दृष्टिकोण को एलएफडब्ल्यू पर आधारित बारीक क्यूरेटेड डेटासेट का उपयोग करके मान्य किया गया है।

सुझाई गई अनुमान पाइपलाइन में, फाइन-ट्यून किए गए डीएनएन के बाद पीसीए/एलडीए आधारित अव्यक्त फीचर स्पेस परिवर्तन होता है। यह फीचर स्पेस परिवर्तन अव्यक्त फीचर स्पेस की आयामीता में कमी और मूल अव्यक्त फीचर स्पेस के इष्टतम घने उप-स्थान पर क्वेरी छवियों के प्रक्षेपण में मदद करता है, जिससे सटीक और तेज़ अनुमान लगाया जा सकता है।

फेस रिकग्निशन के अलावा, हम विजुअल स्लैम के क्षेत्र में भी इस काम का विस्तार करते हैं और विषय के मोबाइल फोन कैमरे से कैप्चर की गई दृश्य छवियों का उपयोग करके इनडोर वातावरण में किसी विषय को सटीक रूप से स्थानीयकृत करने के लिए एक नया दृष्टिकोण पेश करते हैं। हम **इनपोसनेट** नामक एक नया डीप न्यूरल नेटवर्क डीएनएन प्रस्तुत करते हैं, जो उनकी अंतर्निहित समरूपता के बीच अंतर करने में सक्षम होने के साथ-साथ एक इनडोर दृश्य का संक्षिप्त प्रतिनिधित्व उत्पन्न करता है। यह उपयोगकर्ता को एक ही स्थान की लेकिन अलग-अलग ओरिएंटेशन से कैप्चर की गई छवियों के बीच वास्तविक समय में अंतर करने में सक्षम बनाता है, जिससे उपयोगकर्ता स्थिति के साथ-साथ ओरिएंटेशन का पता लगाने में सक्षम होता है। कार्य में प्रस्तुत उपन्यास डीएनएन मोबाइलनेट बनाम 3-छोटा से प्रेरित है। प्रायोगिक परिणामों के माध्यम से जमीनी सच्चाई से 1 मीटर से कम की स्थानीयकरण सटीकता प्राप्त की जाती है और उसकी गणना की जाती है। लक्ष्य एक दृष्टि आधारित प्रणाली प्रस्तुत करना है जिसमें अतिरिक्त बुनियादी ढांचे या बाहरी हार्डवेयर की आवश्यकता के बिना, इनडोर पोजिशनिंग के लिए उपयोग करने की क्षमता होगी।

ABBREVIATIONS

ELM	Extreme Learning Machines
PCA-ELM	Principal Component Analysis Extreme Learning Machine
LDA-ELM	Linear Discriminant Analysis Extreme Learning Machine
FSL	Few-shot Learning
CL	Contrastive Learning
SLAM	Simultaneous Localisation and Motion
NT-Xent	Normalized Temperature-scaled Cross Entropy Loss
MoCo	Momentum Contrast
SimCLR	Simple framework for Contrastive Learning
SimSiam	Simple Siamese Representation
NNCLR	Nearest Neighbour Contrastive Learning
NNSiam	Nearest Neighbour Siamese Representation
BYOL	Bootstrap Your Own Latent
SLFN	Single Hidden Layer Feed forward Network
CNN	Convolutional Neural Networks
DNN	Deep Neural Networks
LARS	Layer-wise Adaptive Rate Scaling
DWT	Discrete Wavelet Transform
SOTA	State of the Art
LFT-CLR	Latent Feature Transformed Contrastive Learning
wNT-Xent	Weighted NT-Xent Loss
InPosNet	Indoor Positioning Network

NOTATION

S	Dimension of Feature Transformed Latent Space
D	Dimension of Original Latent Feature Space
J	Number of Novel Dataset classes
U	Feature Transformed Latent Space
x	Image sample from Novel Dataset
z	Image sample representation in Projected Space
τ	Normalising Temperature for Loss Function
θ	Online Network Parameter
ζ	Target Network Parameter
Φ	Activation Function

Contents

ACKNOWLEDGEMENTS	ii
ABSTRACT	iii
ABBREVIATIONS	vi
NOTATION	vii
LIST OF TABLES	x
LIST OF FIGURES	xiii
1 INTRODUCTION	1
1.1 Motivation, Objective and Scope	3
1.1.1 Motivation	3
1.1.2 Objective	4
1.1.3 Scope	5
1.2 LITERATURE REVIEW	6
1.3 PRELIMINARIES	11
1.3.1 Extreme Learning Machines	12
1.3.2 Contrastive Learning Frameworks	15
2 CONTRASTIVE LEARNING FOR FACE	21
2.1 SimCLR adaption of Novel Contrastive Learning	21
2.1.1 Choice of Projection Head for Online & Target Network	24
2.2 SimSiam adaption of Novel Contrastive Learning	25

2.3	Results and Discussions	26
3	OPTIMAL FEATURE REPRESENTATION FOR FACE	36
3.1	Estimation of the Optimal Latent Feature Space	36
3.2	Novel Feature Space Transformation Framework	38
3.3	Results and Discussions	44
4	OPTIMAL FEATURE REPRESENTATION FOR Visual SLAM	54
4.1	InPosNet : Novel Deep Neural Network for Indoor Positioning	55
4.2	Statistically Matched DWT assisted Visual SLAM	59
4.2.1	Estimation of Multidimensional Statistically Matched Wavelet Filter Bank	60
4.3	Results and Discussions	63
5	CONCLUSION AND FUTURE WORK	72
	Bibliography	75
5.1	Journal Paper	86
5.2	Conference Paper	86

List of Tables

3.1	Results on various Pre-Trained Networks & Dataset	48
3.2	Results on Pre-Trained & Fine-Tuned ResNet-34 Network followed by PCA-ELM based Feature Space Transformer Network, over various Dataset . . .	52
3.3	Results of experiments on Un-normalised and Diverse Test and Train Dataset using FaceNet/ ResNet-50 as Pre-trained network	52
3.4	Results of experiments on Un-normalised and Noisy Test Dataset with Normalised Train Dataset with FaceNet/ ResNet-50 as Pre-trained network . .	53
4.1	Specifications for MobileNetv3-Small	56
4.2	Specifications for InPosNet	56
4.3	Top 1 match results on different Evaluation Dataset	66
4.4	Top 3 match results on different Evaluation Dataset	66
4.5	Top 5 match results on different Evaluation Dataset	66
4.6	Timing details executing experiments on single Testing images using Super-Point and SuperGlue algorithms	67
4.7	Timing details executing experiments on single dataset with 1957 Testing images and 4923 Training images	68
4.8	Top 3 Match Accuracy vs Feature Space Dimension	68

List of Figures

1.1	Face Recognition Pipeline	3
1.2	Correlation between Human Facial Features: Various Facial Attributes modelled using GLM and Face Embedding from Facenet.	5
1.3	Proposed approach with various components that are used for addressing the Low-shot Recognition/ Verification problem.	7
1.4	Extreme Learning Machines implementation using Single hidden Layer Feed-forward Neural network.	14
1.5	MoCo Contrastive Learning Framework.	16
1.6	SimCLR Contrastive Learning Framework.	16
1.7	BYOL Contrastive Learning Framework.	17
1.8	SimSiam Contrastive Learning Framework.	18
1.9	Barlow Twins Contrastive Learning Framework.	19
1.10	Contrastive Learning Framework for: (a) NNCLR (b) NNSiam	20
2.1	Novel Contrastive Learning Framework.	22
2.2	SimCLR adaption of Novel Contrastive Learning Framework. (a) SimCLR architecture (b)LFT-CLR architecture.	23
2.3	SimSiam adaption of Novel Contrastive Learning Framework. (a) SimSiam architecture (b)wSimSiam architecture. The encoder network may be ResNet-34 or ResNet-50.	26
2.4	Test Accuracy over LFW of SimCLR with (a) ResNet-34 (b) ResNet-50 as backbone. Training Loss of SimCLR with (c) ResNet-34 (d) ResNet-50 using NT-Xent loss. Backbone DNN, Resnet-34 and Resnet-50, were pre-trained over MsCeleb-1M. In the first setting, weights of the Projection Head were randomly initialised, while in the second setting, the weights were initialised through the eigen basis vectors obtained by projecting the Novel Training Dataset in the Latent Feature Space.	29

2.5	Accuracy achieved using SimCLR adaption of Novel Contrastive Learning Framework when trained (a) under a regularizer having decay parameter set to 0.001 and (b) without a regularizer. Backbone DNN was Resnet-34 and Resnet-50, pre-trained over MsCeleb-1M.	30
2.6	(a) Accuracy and (b) Loss achieved using SimCLR adaption of Novel Contrastive Learning Framework for varied number of features selected in the Transformed Feature Space. Backbone DNN was Resnet-34 pre-trained over MsCeleb-1M. In index, pca_n stands for Top n pca components selected for the Transformed Latent Space.	31
2.7	(a) Accuracy and (b) Loss achieved using SimCLR adaption of Novel Contrastive Learning Framework for varied number of features selected in the Transformed Feature Space. Backbone DNN was Resnet-50 pre-trained over MsCeleb-1M. In index, pca_n stands for Top n pca components selected for the Transformed Latent Space.	32
2.8	Accuracy achieved using SimCLR adaption of Novel Contrastive Learning Framework when trained without a regularizer. Backbone DNN was Resnet-34 and Resnet-50, pre-trained over Imagenet.	33
2.9	(a)Accuracy and (b) Loss achieved using SimSiam adaption of Novel Contrastive Learning Framework for varied number of features selected in the Transformed Feature Space. Backbone DNN was Resnet-34 pre-trained over MsCeleb-1M. In index, pca_n stands for Top n PCA components selected for the Transformed Latent Space.	34
2.10	(a)Accuracy and (b) Loss achieved using SimSiam adaption of Novel Contrastive Learning Framework for varied number of features selected in the Transformed Feature Space. Backbone DNN was Resnet-50 pre-trained over MsCeleb-1M. In index, pca_n stands for Top n PCA components selected for the Transformed Latent Space.	35
3.1	Steps to transform Original Latent Space to Feature Transformed Latent Space using PCA and Feature Space Transformer (PCA-ELM).	39
3.2	Proposed pipelines using the PCA-ELM, LDA-ELM and LDA-PCA-ELM initialised Feature Space Transformer Network following a Pre-trained Deep Neural Network.	40
3.3	Proposed fine-tuning process using the PCA-ELM, LDA-ELM and LDA-PCA-ELM initialised Feature Space Transformer Network following a Pre-trained Deep Neural Network.	43
3.4	Benchmark comparison over Ultra-Lowshot Dataset after appending (a) PCA and (b) LDA based Feature Space Transformer Network with Pre-Trained Facenet.	46

3.5	Benchmark comparison over Ultra-Lowshot Dataset after appending (a) PCA and (b) LDA based Feature Space Transformer Network with Fine-Tuned ResNet-50.	47
3.6	Benchmark comparison over Lowshot Dataset after appending (a) PCA and (b) LDA based Feature Space Transformer Network with Pre-Trained Facenet.	49
3.7	Benchmark comparison over Lowshot Dataset after appending (a) PCA and (b) LDA based Feature Space Transformer Network with Fine-Tuned Resnet-50.	51
4.1	Suggested pipeline for inference and position the user based upon the scene image.	55
4.2	Comparison of the last layers of MobileNetv3-Small and InPosNet.	55
4.3	Sample indoor map depicting similar looking structures and inherent structural symmetry.	57
4.4	With reference to floor plan shown in Figure 18, (a) and (b) show Training image samples captured from markers A and B . (c) shows Test image sample captured from location T near marker A	57
4.5	Feature Point matching between the 2 Training images and the Test image sample. (a) show the case when image taken from marker B is matched with Test image taken from marker T leading to incorrect matching. (b) show the case when image taken from marker A is matched with Test image taken from marker T leading to correct matching.	58
4.6	Optimal subspace selection using PCA-ELM.	58
4.7	(a) 2 Chanel 2D-multirate filter bank. (b) Suggested pipeline for inference and position the user based upon the scene image.	59
4.8	Benchmark comparison for Top 1 match over different datasets.	64
4.9	Benchmark comparison for Top 3 match over different datasets.	65
4.10	Benchmark comparison for Top 5 match over different datasets.	65
4.11	Transfer learning of MobileNetv3-Small network with triplet loss. (a) Training vs Validation Loss till 15 epochs (b) Train, Validation and Test Accuracy over Original Evaluation Dataset.	69
4.12	Top 3 Accuracy for various (α, β) experimental setting over different dataset.	70
4.13	Shapes of various Statistically Matched filters of the filter bank for (a)L (b)H (c)LL (d)LH (e)HL (f)HH subbands.	71