

## **Master's Thesis**

# **Development of a Source Library for Transfer Learning: Leveraging Clustering and Contrastive Learning Techniques**

Sneha Banerjee  
Matriculation Number: 430069

Monday 4<sup>th</sup> March, 2024

---

<b>1<sup>st</sup> Examiner:</b>	<b>Prof. Dr. Christian Bauckhage, University of Bonn</b>
<b>2<sup>nd</sup> Examiner:</b>	<b>Prof. Dr.-Ing. Alexander Verl, ISW, University of Stuttgart</b>
<b>Supervisor:</b>	<b>Daniel Schiller</b>

## Declaration of Originality

Name	Sneha Banerjee
Matriculation number	430069
Address	Allmandring 8A,70569 Stuttgart
Title	<i>Development of a Source Library for Transfer Learning: Leveraging Clustering and Contrastive Learning Techniques</i>

I now declare,

- that I wrote this work independently,
- that no sources other than those stated are used and that all statements taken from other works—directly or figuratively—are marked as such,
- that the work submitted was not the subject of any other examination procedure, either in its entirety or in substantial parts,
- that I have not published the work in whole or in part, and
- that my work does not violate any rights of third parties and that I exempt the University against any claims of third parties.

## Abstract

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Task and Goal . . . . .	1
1.3	Approach . . . . .	1
1.4	Contribution . . . . .	1
<b>2</b>	<b>Related Work</b>	<b>2</b>
<b>3</b>	<b>Fundamentals</b>	<b>3</b>
3.1	Bin Picking . . . . .	3
3.1.1	Definition and Application . . . . .	3
3.1.2	Gripper Types . . . . .	3
3.1.3	Model-Based Approaches . . . . .	3
3.1.4	Model-Free Approaches . . . . .	3
3.1.5	Robot System . . . . .	3
3.2	Deep learning-based 3D Image processing . . . . .	3
3.2.1	Convolution Neural Network . . . . .	3
3.2.2	Clustering Algorithms . . . . .	4
3.2.3	Autoencoders . . . . .	6
3.2.4	PointNet Autoencoder . . . . .	7
3.2.5	Siamese Network . . . . .	7
3.3	Transfer learning . . . . .	9
3.3.1	Definition . . . . .	10
3.3.2	Categories . . . . .	11
3.3.3	Relevant Algorithms . . . . .	11
3.4	Network generation process with PQ-Net++ . . . . .	11
<b>4</b>	<b>Method</b>	<b>12</b>
4.1	Suitable Algorithm and Adjustments . . . . .	12
4.2	Whatever new Concept we propose . . . . .	12
4.3	Architecture . . . . .	12
4.4	Source and Target Objects . . . . .	12
4.5	Data Generation . . . . .	12
4.6	Important Aspects . . . . .	12
<b>5</b>	<b>Application of algorithms</b>	<b>13</b>
<b>6</b>	<b>Experiments</b>	<b>14</b>
<b>7</b>	<b>Implementation in a Bin Picking Cell</b>	<b>15</b>
<b>8</b>	<b>Conclusions</b>	<b>16</b>
8.1	Discussions . . . . .	16
8.2	Limitations . . . . .	16
8.3	Future Scope . . . . .	16
<b>9</b>	<b>Summary</b>	<b>17</b>
	<b>List of Acronyms</b>	<b>i</b>
	<b>List of Figures</b>	<b>ii</b>
	<b>List of Tables</b>	<b>iii</b>

**List of Symbols**

**iv**

**Bibliography**

**v**

# **1 Introduction**

## **1.1 Motivation**

## **1.2 Task and Goal**

## **1.3 Approach**

## **1.4 Contribution**

## **2 Related Work**

## 3 Fundamentals

### 3.1 Bin Picking

#### 3.1.1 Definition and Application

#### 3.1.2 Gripper Types

#### 3.1.3 Model-Based Approaches

#### 3.1.4 Model-Free Approaches

#### 3.1.5 Robot System

### 3.2 Deep learning-based 3D Image processing

#### 3.2.1 Convolution Neural Network

Inspired from the visual cortex of humans, convolutional neural network (Convolutional Neural Network (CNN) or ConvNet) is a type of deep neural network designed for processing data which appear in a grid like manner like images, videos, etc. It was introduced by Yann LeCun et. al in [22] in 1998. It gets its name because of the usage of a special kind of linear mathematical operation called the convolution instead of using matrix multiplication as prevalent in the pre-existing neural networks. The key components of a CNN are - Convolution layer, activation function, pooling layer, loss function, output layer. The principle building block of a CNN is the convolution layer. It consists of a number of learnable filters (kernels) which can be visualised like a cubic block. The success of CNNs can be attributed to three major concepts: sparse interactions, parameter sharing and equivariant representations.

**Sparse interactions** In traditional fully connected network, matrix multiplication is performed which involves a parameter matrix. The interaction between input and output units are captured by a distinct parameter in the parameter matrix. But CNNs have sparse interactions, i.e. only a subset of units or neurons in a layer is connected to a local region in the preceding layer. This is done by using kernels that are significantly smaller than the input. This implies, less number of parameters need to be stored which reduces the memory consumption and also it is computationally efficient because it has to perform fewer operations. For example, if there are  $m$  inputs and  $n$  outputs, the fully connected network would need to store  $m \times n$  parameters and have a runtime of  $O(m \times n)$  time complexity per input. On the other hand, if we restrict the number of connections for each unit to be  $k$ , then we would have  $k \times n$  and have a runtime of  $O(k \times n)$  time complexity per input, where  $k$  is quite some fold lesser in magnitude as compared to  $m$ . Moreover, since convolution layers are stacked upon one another in a deep convolution network, units in the deeper layers have a larger receptive field, because of its indirect interaction with a larger region in the input.

**Parameter sharing** Another important focus behind using the convolution layer was to reduce the number of parameters in a Fully Connected Network (FCN). For example in a  $1024 \times 1024$  image, a FCN would have over 1 million hidden units, which means it would have over 1 trillion trainable parameters. But the pixels in an image are only locally correlated. So, CNNs make use of the kernels to limit the focus on smaller regions on the image at a time known as the receptive field. This significantly reduces the number of trainable parameters, thus reducing the memory consumption. [5]. These layers perform a convolution operation between the input (eg. image)  $I$  and the kernel  $K$  to produce an output  $S$  known as the feature map as in Eq.1.

**Equivariant Representations** A function is said to be equivariant if the input to a function is changed, then the output changes in the similar way. Because of the parameter sharing mechanism, convolutions operations are translation equivariant. When the kernels are applied on an input image, the convolution layers generate a 2D map of where the particular feature occurs in the particular image. Furthermore, a pixel is related to its neighboring pixels to form a meaningful context (create a feature e.g. an edge in an image) but it is not limited to where it can occur throughout the image. Thus it creates multiple filters



each of which look for the same feature throughout the image. But it is also to be kept in mind that convolution is not equivariant to other geometric and affine transformations like rotation, scaling, etc. Since convolution operations are commutative in nature, Eq. 1 can also be written as Eq.2. Typically, Eq. 2 is easier to incorporate into a machine learning library since values for both 'm' and 'n' varies within a small range.[10]

$$S(i,j) = (I * K)(i,j) = \sum_m \sum_n I(m,n)K(i-m,j-n) \quad (1)$$

$$S(i,j) = (I * K)(i,j) = \sum_m \sum_n I(i-m,j-n)K(m,n) \quad (2)$$

**Activation function** The next step in a CNN is to apply an activation function. the purpose of using a non linear activation function is to capture the non-linearity in the data. Moreover if non-linearity is not used in between the multiple layers of a neural network, the network is effectively only one layer deep which is not capable even to capture the non-linearity in real world datasets. Rectified Linear Unit (ReLU) is the most commonly used non-linear activation function used in CNNs because unlike sigmoid function or tanh function it does not penalise "too correct" data points. Another remarkable benefit of using ReLU is it's ability to propagate gradient through deep networks with a constant factor. Also it is more memory efficient to use ReLU as it doesn't require to store the ReLU outputs separately as compared to tanh outputs.

**Pooling** The third step of a CNN is to use the pooling layer. The main purpose of this operation is to make the detection of the features robust to the exact location of the eye (i.e. invariant to small translations). The different types of pooling operations are - max pooling and average pooling. It also helps in reducing the dimensionality of the input without losing too much information. This is done to make the computations faster down the deeper layers of the network. If the pooling operations are performed after every k pixels, then the next layer processes inputs that are k times lesser. Since the number of parameters in a layer are dependent on the size of the input to the layer, it significantly reduces the computational overhead on using pooling operations.

Complete CNN

### 3.2.2 Clustering Algorithms

When the dataset does not have any labelled data, then it is said to be unsupervised learning. In this case, there is no ground truth available to measure the correctness of the outputs generated by the Machine Learning (ML) models. The primary focus of unsupervised learning is to find hidden and interesting pattern in the data. Unsupervised learning is of utmost importance in the Artificial Intelligence (AI) world as several real world datasets do not have available annotations, which requires a lot of human effort. Unsupervised learning algorithms can be broadly categorized into the following domains-clustering, dimensionality reduction, and association analysis. Clustering algorithms aim at grouping unlabelled data into groups or clusters based on how similar or dissimilar the datapoints are to one another. It can reveal underlying hidden pattern in the data and is used in applications like image segmentation, fraud detection, etc. Dimensionality reduction reduces the number of irrelevant features or dimensions of the dataset. The inclusion of more features does help in the better representation of the dataset but it also significantly increases the memory consumption and complexity to work with it. Also it is often difficult to visualise real-life datasets with too many features. Association analysis is a rule-based unsupervised learning method that reveals the relationship between attributes in the dataset. It is used in applications like market analysis, intrusion detection, etc.

Clustering algorithms can be broadly classified into the following categories: density-based, distribution-based, and hierarchical-based. In density-based clustering, the algorithm looks for areas of high concentration of the datapoints and groups them as a cluster. The benefit of this algorithm is that the shape a cluster can be is not limited and hence doesn't necessarily have to be convex in nature as show in Fig. 3.2.1. Moreover, these algorithms are also very robust to outliers as they do not force the outliers to belong to any category but are rather ignored as it is unlikely that outliers can form an area of high concentration of datapoints. I have carried on the experiments on this thesis on two such density-based clustering algorithms called DBSCAN[9] and Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) algorithms[24]. The available scikit-learn implementations were used for this purpose[29]. The DBSCAN[9] algorithm does not require the users to define the number of clusters to be

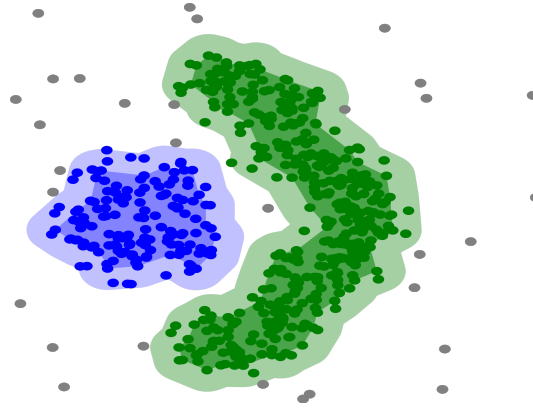


Figure 3.2.1: Density-Based Spatial Clustering of Applications with Noise (DBSCAN) can find non-linearly separable non-convex clusters.[6]

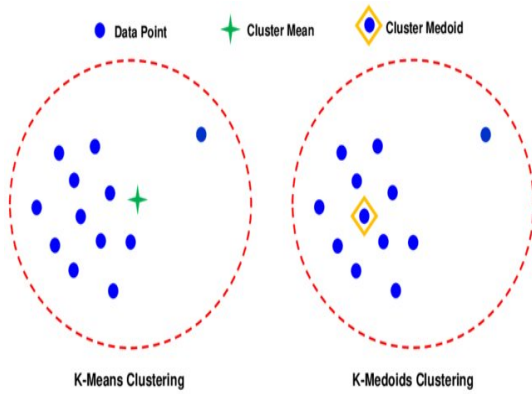


Figure 3.2.2: Difference between K-Means and K-Medoids.[8]

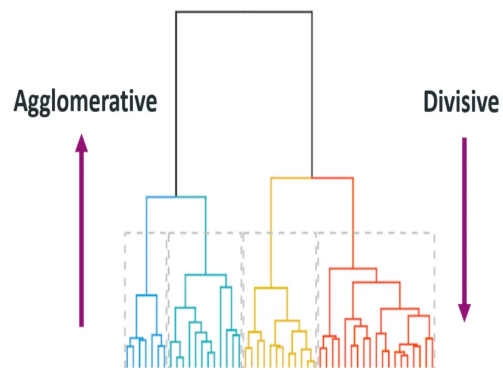


Figure 3.2.3: A dendrogram in hierarchical clustering.[13]

generated which doesn't force dissimilar datapoints to belong to the same cluster. However, this algorithm was still sensitive to two parameters  $\epsilon$ , the maximum distance between the datapoints to be considered in the same cluster and the minimum number of samples in the cluster which the user needs to define [29]. But finding an optimum value for these parameters often require domain expertise and are dependent on the data. The HDBSCAN[24] algorithm mitigates this issue and thus does not require the user to define these two parameters. It is a hierarchical density-based clustering that performs the DBSCAN algorithm over multiple  $\epsilon$  values to find the most stable result. In distribution-based algorithms, a datapoint is said to be a member of the cluster depending on the probability of it's membership to the cluster. The more the distance of a point increases from the centre of the cluster, the less is it's probability of belonging to that cluster. Centroid-based algorithms groups the datapoints based on some initial cluster centres. Once all the datapoints are softly assigned to some cluster membership, the cluster centres are recalculated and this process is iterated until convergence. These algorithms are sensitive to the initial parameters like the cluster centres chosen in the first step. Another major disadvantage of these clustering algorithms are that they always form spherical clusters. The user also needs to define the number of clusters the dataset is to be grouped in, which makes it sensitive to outliers. However, these algorithms can be executed very fast and we have used two such algorithms in our experiments- K-Means[16]and K-Medoids[17]. In K-Means, the mean of the datapoints of the clusters is assigned as the cluster-centroid. It might not be an actual datapoint in the dataset, rather a blurred, noisy average of a datapoints in the cluster. On the contrary, the K-Medoids algorithm assigns an actual datapoint of the dataset, that is most centrally located as the cluster centroid as shown in Fig.3.2.2. Thus K-Medoids is more robust to outliers and noises as compared to K-Means. Hierarchical-based clustering algorithm that form a hierarchy of clusters. Datapoints in a cluster are more

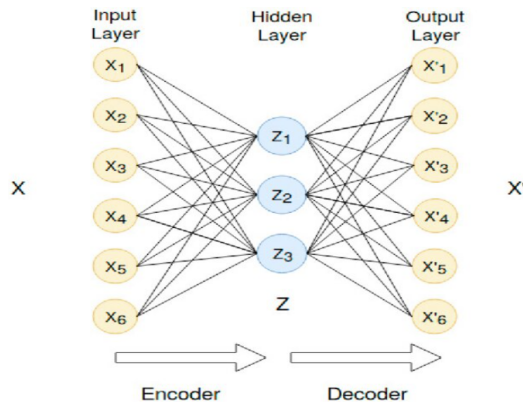


Figure 3.2.4: Architecture of an undercomplete autoencoder.[25]

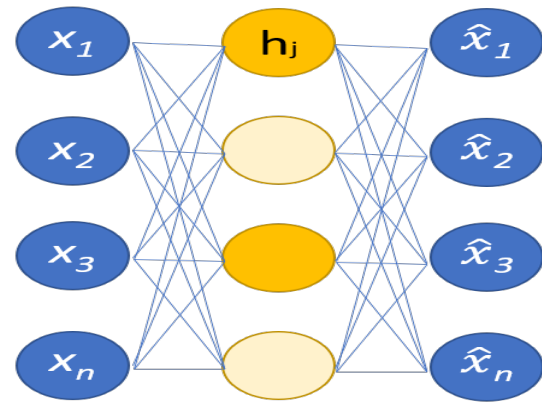


Figure 3.2.5: Architecture of a single layer sparse encoder. The hidden nodes in bright yellow are active, while the ones in light yellow are set to zero, hence inactive[2]

similar to each other as compared to other groups. this hierarchy of clusters is visualised by a hierarchy tree called dendrograms. hierarchical clustering algorithms can be of two types - agglomerative and divisive. In agglomerative clustering, each datapoint is considered as a separate cluster in the first step. Then these clusters are merged into one another until only one cluster remains. Thus at the end, the last level cluster consists of all the datapoints in the dataset. The divisive method is the reverse procedure of the agglomerative method. In the beginning, all the datapoints are considered to be in a single cluster and gradually the cluster is broken into smaller clusters, until each cluster consists of only one datapoint. A visual representation of a dendrogram has been shown in Fig.3.2.3

### 3.2.3 Autoencoders

Autoencoders are a special type of feedforward neural network in which the output tries to reconstruct the input. It is predominantly used in unsupervised learning for the tasks of dimensionality reduction, learning feature representations, and for data compression. It tries to encode the input data into a more compact representation with lower dimensions called the "code"[10] or "latent space". The idea is that this latent space representation should capture the most vital aspects of the input data. This is done by the first part of the network called the encoder. The second part of the network, the decoder, then tries to decode this compressed representation of the input data to reconstruct the original input data as accurately as possible. During the training of an autoencoder, the goal is to minimise this reconstruction loss, i.e. the difference between the original input data and the output of the decoder. There are different types of autoencoders - undercomplete autoencoders, convolutional autoencoders, regularized autoencoders, and variational autoencoders,

Undercomplete autoencoders ensure that the dimension of the latent space representation is less than the dimension of the original input data as shown in Fig.3.2.4. Because, the output of the decoder to be the exact copy of the input is of no use. Rather, if the dimension of the code is less than that of the input, then it ensures that the autoencoder learns those representative features of the input data which are most salient[10]. Convolutional autoencoders are an extension of traditional autoencoders where convolution layers are used as building blocks in both the encoder and the decoder part of the autoencoder. After training the network, the encoder part is used for extracting the features of the data and the decoder part is used for the reconstruction of the input data.

Regularized autoencoders are a special type of autoencoders that employ one or more regularisation techniques to prevent the model from overfitting. The different types of regularisations used are - L1 or L2 Regularisation which adds a penalty to the loss function depending on L1 or L2 norm of the model weights. As a result of this, the model is capable of learning sparse representation of the training data where many weights are set to a very small value or zero as shown in Fig.3.2.5. Dropout is also used as a regulariza-

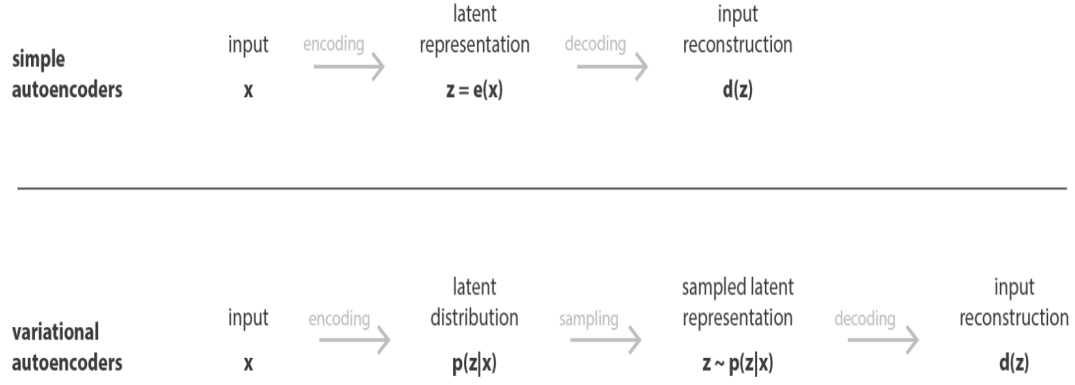


Figure 3.2.6: Difference between a vanilla autoencoder and a VAE.[34]

tion technique where a percentage of the model's nodes are set to zero during training so that the model doesn't rely heavily on any the weight of any particular node, thus preventing overfitting and improving generalisation on unseen data. The previously mentioned regularisation techniques give rise to a variation of regularised autoencoders called the sparse autoencoders[26]. Sometimes noise is also added to the training data or the hidden layers of the model to increase the robustness the model giving rise to denoising autoencoders[35]. One more regularisation technique is to apply a contractive regularization term based on Frobenius norm of the Jacobian matrix of the model's hidden layer activations with respect to the input data[30, 2]. This ensures that a small neighborhood of the input data corresponds to a small neighborhood in the latent space representation, which means small perturbations in the input data leads to small or zero variation in the latent space representation[30, 2]. By doing so, it makes the model more robust to small changes in the input data, thus preventing overfitting.

Variational Autoencoders (VAE)s are a distinct type of autoencoders which produces latent space representations that are continuous, which allows random sampling and interpolation for the generation of new datapoints. The difference between a vanilla autoencoder and a VAE is shown in Fig.3.2.6. Instead of producing a single vector for the latent space representation, it generates two vectors: a vector of means  $\mu$  and a vector of standard deviations  $\gamma$  for all datapoints. An encoding is then sampled from a distribution with mean  $\mu$  and standard deviation  $\gamma$ . Thus even for the same input, i.e. when the mean and the standard deviation are the same, the sampled encoding would vary because of the involvement of the sampling procedure. The mean vector controls the position where the encoding of the input should have its centre and the standard deviation controls the area over which the sampled encoding is allowed to vary from the mean. As a result of this, the decoder learns to decode not just the encoding of a single point but also the points in the neighborhood and hence a continuous latent space representation is obtained. In order to ensure that the latent space representation satisfies that the nearby encoding are similar to each other on a local scale while also facilitating interpolation on a global scale VAEs jointly optimize the reconstruction loss and the Kullback–Leibler divergence (KL)[21] loss.[18, 15]

### 3.2.4 PointNet Autoencoder

### 3.2.5 Siamese Network

Traditional deep learning models have shown ground breaking results in a multitude of applications like image recognition and classification, web scraping, speech recognition and caption generation. But it has a limitation that the models tend to perform well if the training and testing data points are from the same distribution. If the model is to be made capable to predict datapoints belonging to an unknown class, we need to re-train the model with an abundant number of training samples from that class and then predict it with unseen samples of that class. This technique can often be computationally expensive and can increase exponentially if the number of classes increases. Also generating a huge number of samples for a new class of objects can be tedious and quite challenging in multiple scenarios. The Siamese network[19, 4] is one such way to mitigate the above mentioned problem. The idea behind siamese networks is that two

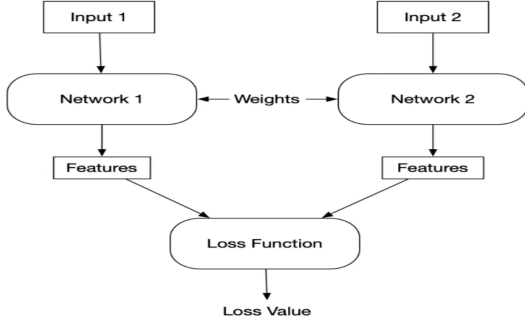


Figure 3.2.7: A generic siamese network.[1]

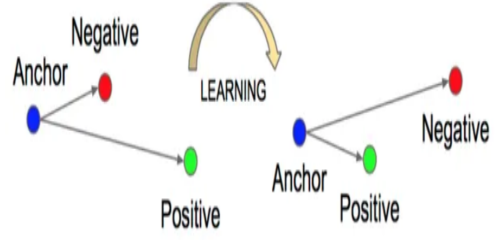


Figure 3.2.8: Before (left) and after (right) minimizing triplet loss function[33]

identical copies of the neural neural network with shared weights process two distinct datapoints as shown in Fig.3.2.7. Because the networks have shared parameters, it ensures that two very similar datapoints cannot have too different locations in the feature maps as both the networks evaluate the same function. In siamese networks, the model learns a similarity function, the loss function, which plays the pivotal role in determining how similar or dissimilar the datapoints are to one another. The two main loss functions used in siamese networks are contrastive loss function[4] and triplet loss function[3]. Since the goal of a siamese network is not to perform classification task on datapoints, rather to compute the similarity between them, contrastive loss functions are more suited for this task as compared to cross-entropy loss functions, as it is for classification tasks.

**Contrastive loss** function evaluates how capable a siamese network is in deciding the similarity between two given datapoints as is given by Eq.3[1]

$$\mathcal{L}_{\text{con}}(Y, D_w) = (1 - Y) \frac{1}{2} (D_w)^2 + (Y) \frac{1}{2} \max(0, m - D_w)^2 \quad (3)$$

where  $\mathcal{L}_{\text{con}}$  is the calculated contrastive loss,  $D_w$  is the Euclidean distance between the outputs of the twin networks[1].  $Y$  takes the value 0 or 1. If the two datapoints belong to the same class then  $Y$  takes the value 0, otherwise it takes the value 1[1].  $\max()$  is a function to choose the higher of the two values, 0 or  $m - D_w$  where  $m$  is the margin and has a value greater than 0[1]. The usage of a margin in this equation ensures that the pair of datapoints which are dissimilar beyond this margin value are excluded while calculating the loss function[1]. The euclidean distance between the embeddings of the two datapoints are given by the Eq.4[1]

$$D_w(X_1, X_2) = \sqrt{\{G_w(X_1) - G_w(X_2)\}^2} \quad (4)$$

where  $G_w(X_i)$  is the output of the twin network such that  $i \in \{1, 2\}$  as shown in Fig.3.2.7.

**Triplet loss** function[36] trains the siamese network in a way that a triplet of datapoints are presented to the siamese network. by applying triplet loss the model tries to group the datapoints similar to the anchor datapoint, close to one another while increasing the distance between the anchor datapoint and the datapoints dissimilar to it as shown in Fig.3.2.8. In order to do so, a positive value called margin is used which increases the distance between dissimilar datapoints and also eliminates the output of any trivial solution. Euclidean distance or cosine distance is used as the distance metric while calculating the triplet loss and shown in Eq.5[12, 1]

$$\mathcal{L}_{\text{tri}}(a, p, n) = \sum_{\substack{a, p, n \\ y_a = y_p \neq y_n}} \max(0, d(a, p) - d(a, n) + \text{margin}) \quad (5)$$

where  $y_x$  is the class label for the datapoint  $x$  for  $x \in \{a, p, n\}$ ,  $a$  is the anchor datapoint,  $p$  is the positive example having the same class label as the anchor datapoint and  $n$  is the negative example having a different class label as compared to the anchor datapoint. but there are some practical problems related to the use of triplet loss function. The number of possible triplets grow cubically with the increase in the size of the dataset. Also most of the triplets are often uninformative. Thus mining "hard triplets" or the ones

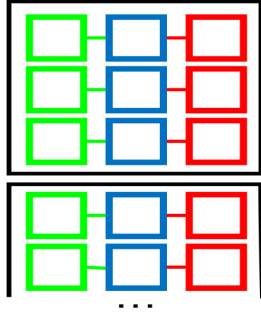


Figure 3.2.9: Minibatches in offline triplet mining. The green boxes are positive datapoints, the blue boxes are anchor datapoints and the red boxes are negative datapoints.

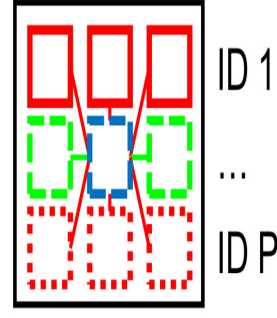


Figure 3.2.10: Triplets constructed only within the minibatch. There are  $P$  object classes in the minibatch. The green boxes are positive datapoints, the blue box is the anchor datapoint and the red boxes are negative datapoints.

crucial for the learning of the model becomes a challenging task. As per the definition of the triplet loss function, the triplets can be of three types: easy triplets, hard triplets and medium hard triplets. The triplets which satisfy the condition  $d(a, p) + margin < d(a, n)$  results to the loss value being 0, hence are called the easy triplets.[33] For the triplet for which the negative datapoint is closer to the anchor as compared to the positive datapoint i.e.  $d(a, n) < d(a, p)$  are the examples of hard triplets.[33] And the medium hard triplets are the ones where the negative datapoints are not closer to the anchor datapoint as compared to the positive datapoint but still results in positive loss i.e.  $d(a, p) < d(a, n) < d(a, p) + margin$ . [33] Thus the medium hard triplets are the most essential for the model training as they constitute for the most useful information. Two triplet mining procedures are used to mitigate the issue of excessive computational overhead with increase in dataset size. They are offline triplet mining and online triplet mining. In offline hard triplet mining, at first the dataset is manually processed to find hard triplets and then they are used for training the model. But manually finding the hard triplets can be quite tedious. Online triplet mining is a better approach as compared to the one mentioned above. The model is trained in minibatches of the dataset as shown in Fig.3.2.9. Using only the triplets that were mined from the dataset could be a wasteful design choice in such scenario. So in a minibatch, it contains way more potential triplets than the ones that were mined. So each member of another triplet becomes an additional negative datapoint example. But both hard positive and hard negative datapoints are required for training. An even better design to mitigate that would be to choose let's say  $K$  datapoints from  $P$  different classes in the dataset, where the  $K$  positive examples serve as hard positives while the rest of the datapoints in the minibatch serve as hard negatives as shown in Fig.3.2.10[12]. Thus siamese networks are very robust to the class imbalance in dataset which is very common in real datasets. Since it depends on very few datapoints of a particular object class, it has good generalisation capabilities for unseen datapoints in the future. They are also very robust to small perturbations in the datapoints due to difference in lighting conditions or orientation or background, as they keep similar objects close to each other even if the datapoints are slightly different. Since the siamese networks places similar datapoints together while increasing the distance between dissimilar datapoints, they are often useful for learning semantic similarity within the datapoints. However despite its benefit, training a siamese network takes longer as compared to traditional neural networks because it involves learning from quadratic pair of datapoints.[1]

### 3.3 Transfer learning

Transfer learning is an optimization in machine learning techniques where a model trained for a particular task let's say task 1 is used for another task, task 2. It uses the knowledge learnt by the model in the previous task to be adapted for the new task. Transfer learning has its benefits in a multitude of applications because real world problems doesn't necessarily always have an abundant amount of labelled data which can be used to train a model from scratch[37]. It can be used for tasks which intend to solve different but related problems[37]. In other words, if the domain of the data for task 2 is similar to that of task 1, transfer



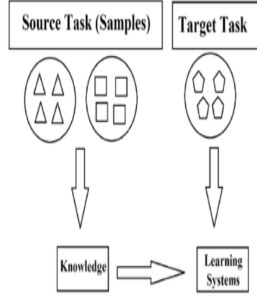


Figure 3.3.1: Model trained for a task in the similar domain is re-used for a new task.[14]

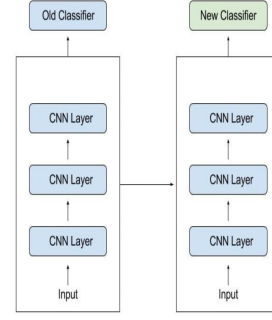


Figure 3.3.2: A model trained on a different task is adapted to be used for a new task.[38]

learning can be used as shown in Fig.3.3.1. Several approaches have been adapted to use transfer learning. Let's say we want to solve task  $T_1$ , but we don't have adequate data to train a deep neural network for it. But we have ample data to solve a similar task  $T_2$ . So we can train a model for task  $T_2$  and then use it for task  $T_1$ . We might need to re-train only the later layers or all the layers of the model, but that heavily depends on the problem at hand. Another approach could be to use an existing pre-trained network. There are a lot of available open-source models which can be re-trained and re-used depending on the problem. Some of the popular models that are used as pre-trained models in transfer learning are Oxford VGG Model[31], Google Inception model[32], Microsoft ResNet model[11]. Another approach could be to use transfer learning for the purpose of feature extraction. Transfer learning has found a lot of application in the domain of computer vision and natural language processing because it requires a lot of data for training such complex models from scratch. In the domain of computer vision, we know that the early layers of a deep neural network focus on learning the low level features like the detection of edges, while the middle layers focus on learning the mid level features like the detection of shapes. In transfer learning, these pre-trained early and middle layers could be used verbatim and only the latter layers which focus on the high level features could be re-trained for the new task as shown in fig.3.3.2. Transfer learning has numerous advantages like significant reduction in the amount of training time[23], improvement in performance for some neural networks[28] and less reliance on a huge amount of data[7, 20, 39].

### 3.3.1 Definition

According to [37, 27], transfer learning can be formally defined as the following. A domain  $\mathcal{D}$  of the input data, with the feature space  $\mathcal{X}$  and a marginal probability distribution  $P(X)$  such that  $X = \{x_1, \dots, x_n\} \in \mathcal{X}$ , where  $n$  is the number of feature vectors in  $X$ . In the given input data domain  $\mathcal{D}$ , the task  $\mathcal{T}$  has a label space  $\mathcal{Y}$  and a predictive function  $f(\cdot)$  which is learned from the feature vector and label pairs  $(x_i, y_i)$  such that  $x_i \in \mathcal{X}$  and  $y_i \in \mathcal{Y}$  and  $f(x)$  is the learner that predicts the label value for  $x$ . As per the above definition, the source domain is defined as  $\mathcal{D}_S$  and the corresponding source task as  $\mathcal{T}_S$ . The target domain is defined as  $\mathcal{D}_T$  and the corresponding target task as  $\mathcal{T}_T$ . Transfer learning aims at improving the target predictive function  $f_T(\cdot)$  by using the related information from  $\mathcal{D}_S$  and  $\mathcal{T}_S$  where  $\mathcal{D}_S \neq \mathcal{D}_T$  or  $\mathcal{T}_S \neq \mathcal{T}_T$ . Since  $\mathcal{D}_S = \{\mathcal{X}_S, P(X_S)\}$  and  $\mathcal{D}_T = \{\mathcal{X}_T, P(X_T)\}$ , here the condition  $\mathcal{D}_S \neq \mathcal{D}_T$  means that  $\mathcal{X}_S \neq \mathcal{X}_T$  and/or  $P(X_S) \neq P(X_T)$ . The scenario where  $\mathcal{X}_S \neq \mathcal{X}_T$  in the context of transfer learning is defined as heterogeneous transfer learning. The scenario where  $\mathcal{X}_S = \mathcal{X}_T$  is defined as homogeneous transfer learning. The scenario where  $P(X_S) \neq P(X_T)$  means the source and the target data are from different domains, i.e., they belong to different marginal distributions. Transfer learning algorithms are not expected to give optimal results in this case. Referring back to the definition of transfer learning, when  $\mathcal{T}_S \neq \mathcal{T}_T$  and it was defined that  $\mathcal{T} = \{(Y), P(Y|X)\}$ , this scenario is possible when  $\mathcal{Y}_S \neq \mathcal{Y}_T$  and/or  $P(Y_S|X_S) \neq P(Y_T|X_T)$ . The scenario where  $P(Y_S|X_S) \neq P(Y_T|X_T)$  implies that the conditional probability distributions between the source and target domains are different. The case  $\mathcal{Y}_S \neq \mathcal{Y}_T$  means the label space of the source and target domain are different.[37]

### **3.3.2 Categories**

### **3.3.3 Relevant Algorithms**

## **3.4 Network generation process with PQ-Net++**



## **4 Method**

### **4.1 Suitable Algorithm and Adjustments**

### **4.2 Whatever new Concept we propose**

### **4.3 Architecture**

### **4.4 Source and Target Objects**

### **4.5 Data Generation**

### **4.6 Important Aspects**

## **5 Application of algorithms**

## 6 Experiments

## **7 Implementation in a Bin Picking Cell**

## **8 Conclusions**

### **8.1 Discussions**

### **8.2 Limitations**

### **8.3 Future Scope**

## 9 Summary

## List of Acronyms

**CNN** Convolutional Neural Network

**FCN** Fully Connected Network

**ReLU** Rectified Linear Unit

**ML** Machine Learning

**AI** Artificial Intelligence

**DBSCAN** Density-Based Spatial Clustering of Applications with Noise

**HDBSCAN** Hierarchical Density-Based Spatial Clustering of Applications with Noise

**VAE** Variational Autoencoders

**KL** Kullback–Leibler divergence

## List of Figures

3.2.1 DBSCAN can find non-linearly separable non-convex clusters.[6]	5
3.2.2 Difference between K-Means and K-Medoids.[8]	5
3.2.3 A dendrogram in hierarchical clustering.[13]	5
3.2.4 Architecture of an undercomplete autoencoder.[25]	6
3.2.5 Architecture of a single layer sparse encoder. The hidden nodes in bright yellow are active, while the ones in light yellow are set to zero, hence inactive[2]	6
3.2.6 Difference between a vanilla autoencoder and a VAE.[34]	7
3.2.7 A generic siamese network.[1]	8
3.2.8 Before (left) and after (right) minimizing triplet loss function[33]	8
3.2.9 Minibatches in offline triplet mining. The green boxes are positive datapoints, the blue boxes are anchor datapoints and the red boxes are negative datapoints.	9
3.2.10 Triplets constructed only within the minibatch. There are P object classes in the minibatch. The green boxes are positive datapoints, the blue box is the anchor datapoint and the red boxes are negative datapoints.	9
3.3.1 Model trained for a task in the similar domain is re-used for a new task.[14]	10
3.3.2 A model trained on a different task is adapted to be used for a new task.[38]	10



## List of Tables

## List of Symbols

## Bibliography

- [1] *A Comprehensive Guide to Siamese Neural Networks*. <https://medium.com/@rinkinag24/a-comprehensive-guide-to-siamese-neural-networks-3358658c0513>.
- [2] *Autoencoder*. <https://en.wikipedia.org/wiki/Autoencoder>.
- [3] Vassileios Balntas et al. "Learning local feature descriptors with triplets and shallow convolutional neural networks." In: *Bmvc*. Vol. 1. 2. 2016, p. 3.
- [4] Jane Bromley et al. "Signature verification using a" siamese" time delay neural network". In: *Advances in neural information processing systems* 6 (1993).
- [5] *Convolutional Neural Networks, Explained*. <https://towardsdatascience.com/convolutional-neural-networks-explained-9cc5188c4939>.
- [6] *DBSCAN Wikipedia*. <https://en.wikipedia.org/wiki/DBSCAN>.
- [7] Lixin Duan, Dong Xu, and Ivor Tsang. "Learning with augmented features for heterogeneous domain adaptation". In: *arXiv preprint arXiv:1206.4660* (2012).
- [8] Alireza Entezami, Hassan Sarmadi, and Behzad Saeedi Razavi. "An innovative hybrid strategy for structural health monitoring by modal flexibility and clustering methods". In: *Journal of Civil Structural Health Monitoring* 10.5 (2020), pp. 845–859.
- [9] Martin Ester et al. "A density-based algorithm for discovering clusters in large spatial databases with noise". In: *kdd*. Vol. 96. 34. 1996, pp. 226–231.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [11] Kaiming He et al. "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [12] Alexander Hermans, Lucas Beyer, and Bastian Leibe. "In defense of the triplet loss for person re-identification". In: *arXiv preprint arXiv:1703.07737* (2017).
- [13] *Hierarchical Clustering*. <https://harshsharma1091996.medium.com/hierarchical-clustering-996745fe656b>.
- [14] Asmaul Hosna et al. "Transfer learning: a friendly introduction". In: *Journal of Big Data* 9.1 (2022), p. 102.
- [15] *Intuitively Understanding Variational Autoencoders*. <https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>.
- [16] Xin Jin and Jiawei Han. "K-Means Clustering". In: *Encyclopedia of Machine Learning*. Ed. by Claude Sammut and Geoffrey I. Webb. Boston, MA: Springer US, 2010, pp. 563–564. ISBN: 978-0-387-30164-8. DOI: 10.1007/978-0-387-30164-8\_425. URL: [https://doi.org/10.1007/978-0-387-30164-8\\_425](https://doi.org/10.1007/978-0-387-30164-8_425).
- [17] Xin Jin and Jiawei Han. "K-Medoids Clustering". In: *Encyclopedia of Machine Learning*. Ed. by Claude Sammut and Geoffrey I. Webb. Boston, MA: Springer US, 2010, pp. 564–565. ISBN: 978-0-387-30164-8. DOI: 10.1007/978-0-387-30164-8\_426. URL: [https://doi.org/10.1007/978-0-387-30164-8\\_426](https://doi.org/10.1007/978-0-387-30164-8_426).
- [18] Diederik P Kingma, Max Welling, et al. "An introduction to variational autoencoders". In: *Foundations and Trends® in Machine Learning* 12.4 (2019), pp. 307–392.
- [19] Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. "Siamese neural networks for one-shot image recognition". In: *ICML deep learning workshop*. Vol. 2. 1. Lille. 2015.
- [20] Brian Kulis, Kate Saenko, and Trevor Darrell. "What you saw is not what you get: Domain adaptation using asymmetric kernel transforms". In: *CVPR 2011*. IEEE. 2011, pp. 1785–1792.
- [21] Solomon Kullback and Richard A Leibler. "On information and sufficiency". In: *The annals of mathematical statistics* 22.1 (1951), pp. 79–86.

- [22] Yann LeCun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
- [23] Tong Liu et al. “Exploring transfer learning to reduce training overhead of hpc data in machine learning”. In: *2019 IEEE International Conference on Networking, Architecture and Storage (NAS)*. IEEE. 2019, pp. 1–7.
- [24] Claudia Malzer and Marcus Baum. “A hybrid approach to hierarchical density-based cluster selection”. In: *2020 IEEE international conference on multisensor fusion and integration for intelligent systems (MFI)*. IEEE. 2020, pp. 223–228.
- [25] Sumit Misra et al. “An autoencoder based model for detecting fraudulent credit card transaction”. In: *Procedia Computer Science* 167 (2020), pp. 254–262.
- [26] Andrew Ng et al. “Sparse autoencoder”. In: *CS294A Lecture notes* 72.2011 (2011), pp. 1–19.
- [27] Sinno Jialin Pan and Qiang Yang. “A survey on transfer learning”. In: *IEEE Transactions on knowledge and data engineering* 22.10 (2009), pp. 1345–1359.
- [28] Weike Pan et al. “Transfer learning to predict missing ratings via heterogeneous user feedbacks”. In: *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain*. 2011, p. 2318.
- [29] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [30] Salah Rifai et al. “Contractive auto-encoders: Explicit invariance during feature extraction”. In: *Proceedings of the 28th international conference on international conference on machine learning*. 2011, pp. 833–840.
- [31] Karen Simonyan and Andrew Zisserman. “Very deep convolutional networks for large-scale image recognition”. In: *arXiv preprint arXiv:1409.1556* (2014).
- [32] Christian Szegedy et al. “Going deeper with convolutions”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9.
- [33] *Triplet Loss and Siamese Neural Networks*. <https://medium.com/@enoshshr/triplet-loss-and-siamese-neural-networks-5d363fdeba9b>.
- [34] *Understanding Variational Autoencoders(VAEs)*. <https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73>.
- [35] Pascal Vincent et al. “Extracting and composing robust features with denoising autoencoders”. In: *Proceedings of the 25th international conference on Machine learning*. 2008, pp. 1096–1103.
- [36] Kilian Q Weinberger and Lawrence K Saul. “Distance metric learning for large margin nearest neighbor classification.” In: *Journal of machine learning research* 10.2 (2009).
- [37] Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. “A survey of transfer learning”. In: *Journal of Big data* 3.1 (2016), pp. 1–40.
- [38] *What Is Transfer Learning? Exploring the Popular Deep Learning Approach*. <https://builtin.com/data-science/transfer-learning>.
- [39] Yin Zhu et al. “Heterogeneous transfer learning for image classification”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 25. 1. 2011, pp. 1304–1309.