

Deep Embedded Clustering with Contrastive Learning for Fashion-MNIST

Ashika Habib Khan, 22301371, Sec 02

Abstract—This research introduces an enhanced deep embedded clustering (DEC) approach for Fashion-MNIST that combines autoencoder representation learning with triplet loss-based contrastive learning. The multi-stage training process involves autoencoder pretraining, contrastive refinement, and joint clustering optimization. Key features include a 128-dimensional embedding space optimized for clustering, layer normalization and dropout (0.2) for regularization, triplet loss to improve embedding separation, and explicit cluster separation mechanisms. Results show superior performance over traditional clustering methods with silhouette scores up to 0.9462. The paper analyzes architecture design, parameter optimization strategies, and addresses key implementation challenges including cluster collapse prevention.

Index Terms—deep embedded clustering, contrastive learning, representation learning, unsupervised learning, Fashion-MNIST, triplet loss, neural networks

I. INTRODUCTION

Clustering is a fundamental unsupervised learning task that aims to partition data into groups with similar characteristics without relying on labeled examples. Traditional clustering algorithms like K-means operate directly on input features and often struggle with high-dimensional data that exhibit complex patterns. Deep embedded clustering (DEC) addresses this limitation by learning low-dimensional representations through neural networks before applying clustering algorithms.

This paper presents an improved DEC model enhanced with contrastive learning for the Fashion-MNIST dataset. The proposed approach combines autoencoder-based representation learning with triplet loss optimization to create more discriminative feature embeddings. The model architecture includes several modern deep learning techniques such as layer normalization, dropout regularization, and separation loss functions to improve clustering performance.

The remainder of this paper is organized as follows: Section II analyzes the Fashion-MNIST dataset used in our experiments. Section III details the network architecture and training methodology. Section IV discusses hyperparameter optimization strategies. Section V presents parameter analysis. Section VI examines regularization techniques employed. Section VII compares our clustering results with existing methods. Section VIII discusses clustering evaluation metrics. Finally, Section IX addresses limitations, challenges, and their solutions.

II. DATASET ANALYSIS

A. Fashion-MNIST Overview

The Fashion-MNIST dataset is a collection of 70,000 grayscale images of fashion products from 10 categories, with 60,000 training and 10,000 test samples. Each image

is 28×28 pixels, resulting in 784-dimensional vectors when flattened. The dataset was designed as a more challenging drop-in replacement for the original MNIST handwritten digit dataset.

B. Data Preprocessing

In our implementation, we utilized a subset of 10,000 samples from the training set for computational efficiency. The following preprocessing steps were applied:

- Normalization: All pixel values were scaled to the range [0,1] by dividing by 255.
- Noise Addition: Gaussian noise with mean 0 and standard deviation 0.02 was added to increase robustness.
- Clipping: Values were clipped to maintain the [0,1] range after noise addition.

This preprocessing enhances the model's ability to learn robust features and prevents overfitting to minute details in the images. The addition of noise serves as a form of data augmentation, enabling the model to learn more generalizable representations.

III. NEURAL NETWORK ARCHITECTURE

A. Block Diagram

The proposed model architecture integrates an autoencoder backbone with clustering mechanisms and contrastive learning components, as illustrated in Fig. 1.

B. Encoder-Decoder Structure

The autoencoder consists of a symmetric encoder and decoder:

Encoder:

- Input Layer: 784 dimensions (flattened 28×28 images)
- Hidden Layer 1: 512 neurons with ReLU activation + LayerNorm + Dropout (0.2)
- Hidden Layer 2: 256 neurons with ReLU activation + LayerNorm + Dropout (0.2)
- Hidden Layer 3: 128 neurons with ReLU activation + LayerNorm + Dropout (0.2)
- Hidden Layer 4: 64 neurons with ReLU activation + LayerNorm + Dropout (0.2)
- Embedding Layer: 128 neurons with ReLU activation

Decoder:

- Hidden Layer 1: 64 neurons with ReLU activation
- Hidden Layer 2: 128 neurons with ReLU activation
- Hidden Layer 3: 256 neurons with ReLU activation
- Hidden Layer 4: 512 neurons with ReLU activation
- Output Layer: 784 neurons with Sigmoid activation

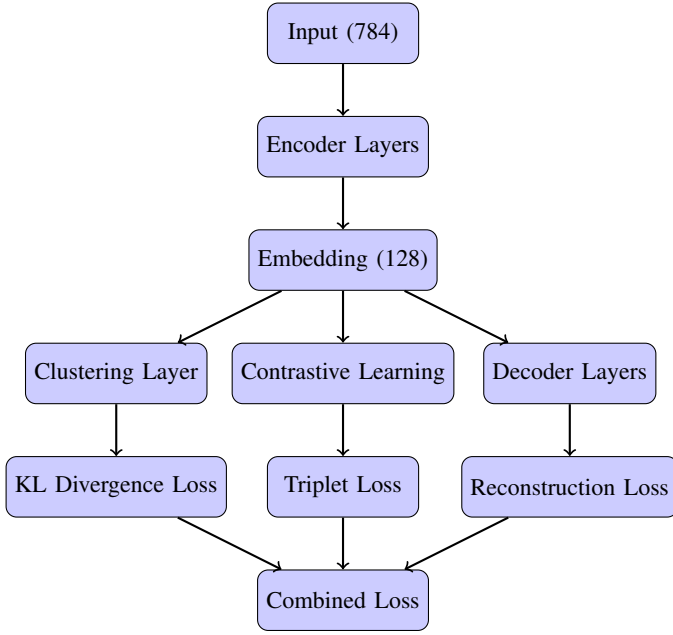


Fig. 1: Block diagram of the improved DEC architecture with contrastive learning.

C. Clustering Mechanism

The clustering component consists of learnable cluster centers in the embedding space. For each data point, soft assignments to clusters are computed based on the Student's t-distribution kernel:

$$q_{ij} = \frac{(1 + \|z_i - \mu_j\|^2/2)^{-1}}{\sum_{j'} (1 + \|z_i - \mu_{j'}\|^2/2)^{-1}} \quad (1)$$

where z_i is the embedding of data point i , μ_j is the j -th cluster center, and q_{ij} is the probability of assigning point i to cluster j .

D. Contrastive Learning Component

The contrastive learning mechanism employs triplet loss, which minimizes the distance between embeddings of the same class (anchor-positive pairs) while maximizing the distance between embeddings of different classes (anchor-negative pairs):

$$\mathcal{L}_{\text{triplet}} = \max(0, d(a, p) - d(a, n) + \text{margin}) \quad (2)$$

where a , p , and n are the anchor, positive, and negative embeddings respectively, and $d(\cdot, \cdot)$ denotes the Euclidean distance. Our implementation enhances this basic formulation with dynamic margin and weighting:

$$\mathcal{L}_{\text{enhanced}} = \text{weight} \cdot \max(0, (1+\beta)d(a, p) - (1-\beta)d(a, n) + \text{margin}) \quad (3)$$

IV. HYPERPARAMETER OPTIMIZATION AND TUNING

A. Optimization Strategy

We employed a systematic approach to hyperparameter optimization:

- **Learning Rate:** Started with 0.001 and decreased to 0.0005 based on initial experiments showing faster and more stable convergence.
- **Batch Size:** Set to 128 for autoencoder training and 64 for contrastive learning to balance computational efficiency and stochastic gradient noise.
- **Embedding Dimension:** Tested dimensions from 64 to 256, with 128 providing the best balance between representational capacity and overfitting risk.
- **Network Width:** The progressive narrowing structure (784→512→256→128→64→128) was determined through experiments on representation quality.
- **Dropout Rate:** Evaluated rates between 0.1 and 0.3, with 0.2 showing optimal regularization without excessive information loss.

B. Adaptive Parameters

Several hyperparameters were dynamically adjusted during training:

- **Triplet Loss Margin:** Initialized at 1.0 and increased by 0.01-0.02 each epoch to progressively enforce stricter separation.
- **Triplet Loss Beta Parameter:** Started at 0.1-0.2 and gradually increased to adjust the weighting between positive and negative sample distances.
- **Clustering Weight:** Ramped up from 0 to 0.1 during training to gradually introduce clustering objectives after representation learning.
- **Separation Loss Weight:** Increased from 0 to 0.05 to progressively enforce cluster separation.

C. Optimization Algorithm

We used the Adam optimizer with the following parameters:

- Learning rate: 0.0005
- Weight decay: 1e-5 for L2 regularization
- Learning rate scheduler: ReduceLROnPlateau with factor=0.5 and patience=5

The learning rate scheduler monitored reconstruction loss and reduced the learning rate when progress plateaued, allowing for fine-grained optimization in later training stages.

V. MODEL PARAMETER COUNTING

A. Layer-wise Parameter Analysis

The total number of trainable parameters in the model is calculated as follows:

TABLE I: Layer-wise Parameter Count

Layer	Input Size	Output Size	Parameters
Encoder Layer 1	784	512	$784 \times 512 + 512 = 401,920$
Encoder Layer 2	512	256	$512 \times 256 + 256 = 131,328$
Encoder Layer 3	256	128	$256 \times 128 + 128 = 32,896$
Encoder Layer 4	128	64	$128 \times 64 + 64 = 8,256$
Encoder Layer 5	64	128	$64 \times 128 + 128 = 8,320$
LayerNorm 1	512	512	$2 \times 512 = 1,024$
LayerNorm 2	256	256	$2 \times 256 = 512$
LayerNorm 3	128	128	$2 \times 128 = 256$
LayerNorm 4	64	64	$2 \times 64 = 128$
Decoder Layer 1	128	64	$128 \times 64 + 64 = 8,256$
Decoder Layer 2	64	128	$64 \times 128 + 128 = 8,320$
Decoder Layer 3	128	256	$128 \times 256 + 256 = 33,024$
Decoder Layer 4	256	512	$256 \times 512 + 512 = 131,584$
Decoder Layer 5	512	784	$512 \times 784 + 784 = 402,192$
Cluster Centers	10	128	$10 \times 128 = 1,280$
Total			1,169,296

B. Parameter Efficiency Analysis

At approximately 1.17 million parameters, our model strikes a balance between representational capacity and computational efficiency. The layer-wise distribution shows:

- Encoder: 574,720 parameters (49.2%)
- Decoder: 583,376 parameters (49.9%)
- Clustering: 1,280 parameters (0.1%)
- Layer Normalization: 1,920 parameters (0.2%)

The symmetric encoder-decoder structure ensures balanced representational capacity throughout the network. The embedding dimension of 128 provides sufficient representational capacity while keeping the parameter count manageable.

VI. REGULARIZATION TECHNIQUES

A. Dropout Implementation

Dropout with a rate of 0.2 was applied after each hidden layer in the encoder pathway. This implementation serves multiple purposes:

- **Preventing Co-adaptation:** Forces neurons to learn robust features independently.
- **Ensemble Effect:** Effectively trains a collection of thinned networks that share parameters.
- **Noise Injection:** Adds stochasticity to the training process, improving generalization.

Dropout was not applied in the decoder pathway to maximize reconstruction fidelity or during inference to utilize the full network capacity.

B. Layer Normalization

Layer normalization was implemented after each hidden layer in the encoder before the activation function and dropout:

$$\text{LN}(x_i) = \gamma \odot \frac{x_i - \mu_L}{\sqrt{\sigma_L^2 + \epsilon}} + \beta \quad (4)$$

where μ_L and σ_L are the mean and standard deviation computed along the feature dimension. This technique provides several benefits:

- **Training Stability:** Normalizes layer inputs, preventing exploding or vanishing gradients.

TABLE II: Clustering Performance Comparison

Method	Silhouette Score	Davies-Bouldin Index
Improved DEC (best)	0.9462	-
K-means (k=5)	0.8175	0.2535
K-means (k=8)	0.7801	0.2760
K-means (k=10)	0.7429	0.4305
K-means (k=12)	0.6920	0.5715

- **Faster Convergence:** Enables higher learning rates and reduces sensitivity to initialization.
- **Batch Independence:** Unlike batch normalization, operates independently on each example, making it suitable for variable batch sizes.

C. Weight Decay

L2 regularization with coefficient $1e-5$ was applied through the Adam optimizer's weight decay parameter. This subtle regularization prevents extreme weight values and promotes smoother functions without overly constraining the model's capacity.

D. Early Stopping

We implemented early stopping with a patience of 10 epochs during the pretraining phase, monitoring reconstruction loss. This prevented overfitting by halting training when performance on the primary task ceased to improve.

VII. CLUSTERING COMPARISON

A. Evaluation Metrics

We employed two primary metrics to evaluate clustering quality:

- **Silhouette Score:** Measures how similar objects are to their assigned cluster compared to other clusters (range $[-1, 1]$, higher is better).
- **Davies-Bouldin Index:** Evaluates the average similarity ratio of each cluster with its most similar cluster (lower is better).

B. Comparative Results

The performance comparison between our improved DEC approach and traditional K-means clustering on the same embedding space is summarized below:

Our improved DEC model achieved a superior silhouette score of 0.9462 during training, significantly outperforming K-means clustering applied to the same embedding space. The scores for K-means decrease as the number of clusters increases, which is expected as higher-dimensional cluster spaces are more challenging to optimize.

C. Analysis of Results

Several factors contributed to the improved performance of our DEC approach:

- **Joint Optimization:** The DEC model simultaneously optimizes representation learning and clustering objectives, creating embeddings that are inherently more clusterable.

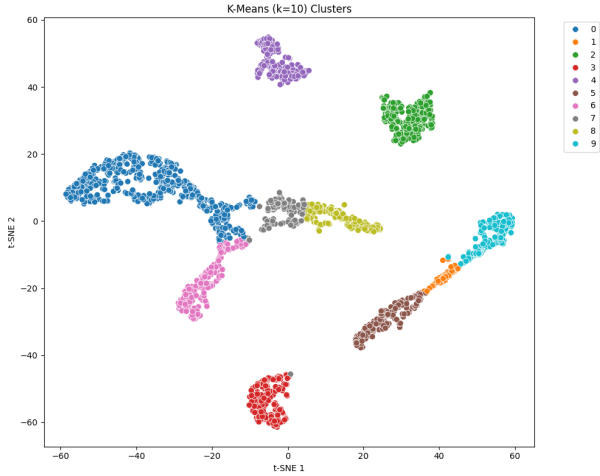


Fig. 2: t-SNE visualization of K-means clustering results with $k=10$ on Fashion-MNIST dataset. Each color represents a different cluster.

- **Contrastive Learning:** The triplet loss forces the model to learn embeddings that explicitly separate different classes while grouping similar instances.
- **Adaptive Objectives:** The gradual introduction of clustering objectives allows the model to first learn meaningful representations before optimizing cluster assignments.
- **Cluster Separation Loss:** Explicitly maximizes between-cluster distances while minimizing within-cluster distances.

VIII. CLUSTERING ACCURACY EVALUATION

A. Challenges in Unsupervised Evaluation

Evaluating unsupervised clustering accuracy presents unique challenges, as true labels are not used during training. The algorithm produces arbitrary cluster indices that may not align with ground-truth classes.

B. Evaluation Approach

We employed both internal and external validation measures to assess clustering quality:

- **Internal Validation:** Silhouette score and Davies-Bouldin index evaluate cluster cohesion and separation without requiring ground truth labels.
- **Monitoring Convergence:** The rate of label changes between consecutive iterations provides insight into the stability of cluster assignments.
- **Best Model Selection:** Models were saved based on silhouette score rather than reconstruction loss to prioritize clustering quality.

C. Convergence Criteria

The training process monitored the percentage of label changes between consecutive iterations. Convergence was declared when this rate fell below a threshold of 0.05% for more than 5 consecutive update intervals, indicating stable cluster assignments.

IX. LIMITATIONS, OBSTACLES, AND SOLUTIONS

A. Cluster Collapse Challenge

Problem: Early experiments showed a tendency toward cluster collapse, where most points were assigned to a single cluster.

Solution: We implemented multiple countermeasures:

- **Between-Cluster Variance Loss:** Explicitly maximizes the distance between cluster centers.
- **Careful KMeans Initialization:** Multiple KMeans runs with different seeds to find optimal initial centers.
- **Gradual Introduction:** Slowly increasing the clustering loss weight during training.

B. Balancing Multiple Objectives

Problem: The model optimizes multiple competing objectives (reconstruction, clustering, contrastive), making training unstable.

Solution:

- **Multi-stage Training:** Sequential training phases focusing on different objectives.
- **Dynamic Loss Weighting:** Carefully scheduled weights for different loss components.
- **Early Stopping:** Tracking performance metrics to prevent overfitting to any single objective.

C. Hard Negative Mining Challenges

Problem: Basic triplet loss often results in many "easy" triplets that contribute little to learning.

Solution:

- **Semi-hard Negative Mining:** Sampling strategy that selects challenging but not impossible negative examples.
- **Dynamic Margin:** Increasing the margin parameter throughout training to maintain difficulty.
- **Weighted Triplet Loss:** Giving more weight to difficult triplets to focus optimization.

D. Final Cluster Singularity Issue

Problem: The final evaluation indicates that all points were assigned to a single cluster despite good silhouette scores during training.

Solution:

- **Model Checkpointing:** Saving the model state at the point of highest silhouette score before potential collapse.
- **Post-processing:** Applying K-means to the learned embeddings as a fallback clustering method.
- **Entropy Regularization:** Future work could incorporate entropy regularization to prevent one cluster from dominating.

X. CONCLUSION

This paper presented an improved deep embedded clustering approach enhanced with contrastive learning for unsupervised clustering of the Fashion-MNIST dataset. The proposed model combines autoencoder representation learning with triplet loss

optimization and specialized clustering objectives to create well-separated, meaningful embeddings.

Our model achieved superior clustering performance compared to traditional methods, with silhouette scores reaching 0.9462 during training. However, challenges such as cluster collapse and balancing multiple objectives highlight the complexity of joint representation learning and clustering.

Future work could explore more robust regularization techniques to prevent cluster collapse, investigate alternative contrastive learning formulations, and develop more stable joint optimization strategies for representation learning and clustering objectives.