

Analyze the features learned by different layers of an autoencoder for line segmentation in images using statistical classifiers.

Kulala Chetan

Department of Computer Science and
Engineering,
Amrita School of Computing,
Bengaluru,
Amrita Vishwa Vidyapeetham, India
bl.en.u4cse21105@bl.students.amrita.edu

Hothur Shreyaa

Department of Computer Science and
Engineering,
Amrita School of Computing,
Bengaluru,
Amrita Vishwa Vidyapeetham, India
bl.en.u4cse21069@bl.students.amrita.edu

Anusha.p

Department of Computer Science and
Engineering,
Amrita School of Computing,
Bengaluru,
Amrita Vishwa Vidyapeetham, India
bl.en.u4cse21069@bl.students.amrita.edu

Abstract— In the context of character recognition and classification, this research presents a comprehensive methodology that encompasses feature extraction through deep learning networks, the identification and removal of redundant and irrelevant features via feature selection and reduction techniques, the utilization of the selected optimal features to train a model, and the subsequent application of the trained model to predict the class labels of unseen characters.

Keywords—Deep Learning, Feature Extraction, Feature Selection, Feature Reduction, Character Recognition, Character Classification, Model Training, Unseen Character Prediction, Redundancy Analysis.

I. INTRODUCTION

In recent years, deep learning has brought about significant improvements in character classification. This involves sorting characters, symbols, or text into specific categories, and it plays a crucial role in areas like handwriting recognition, document analysis, and language processing.

Deep learning models are known for their ability to learn complex patterns from data, which makes them valuable for character classification. However, they often produce high-dimensional data that can pose challenges such as slow processing, difficulty in understanding the model's decisions, and the risk of overfitting. To tackle these issues, this study introduces a method that combines deep learning with techniques to choose the most important features, making character classification more efficient and accurate.

Our goal is to improve character classification by taking the best of both worlds: deep learning's power and feature selection's efficiency. We'll explain how we extract useful information from deep learning models, get rid of unnecessary data, and use what's left to train strong models. We'll also explore methods like Principal Component Analysis (PCA) and feature importance to boost efficiency without compromising accuracy.

In the following sections, we'll dive into the details, explaining how we extract features, choose the best ones, and train models. We'll also evaluate the models using various criteria and demonstrate their ability to predict the categories of new characters. In essence, this research offers a comprehensive guide for character classification, with

potential applications in areas where recognizing and sorting characters are important.

In summary, this approach blends feature selection with deep learning to address challenges associated with high-dimensional data in character classification. By combining the strengths of both methods, we aim to provide a practical and effective solution for real-world character classification tasks.

II. LITERATURE SURVEY

[1] by Zhang et al. (2020) investigates the problem of redundancy and irrelevant features learned by deep learning networks for character recognition. They propose a method to select optimal features from the deep learning network for training statistical classifiers. The authors analyze the impact of redundant and irrelevant features on the network's performance and provide insights into feature selection techniques.

Wang et al. (2021) propose a feature selection method for deep learning-based character recognition in [2]. Their method is based on the mutual information between the features and the labels. By evaluating the relevance of each feature to the target labels, the authors aim to identify and retain only the most informative features, thereby reducing redundancy and improving the efficiency of the recognition system.

Zhang et al. (2022) in [3] presents a method to reduce the redundancy of features learned by deep learning networks for character recognition. Their approach is based on L1-norm regularization, which encourages sparsity in the feature representation. By penalizing the network's weights, the authors aim to encourage the network to learn more compact and discriminative features, leading to improved recognition performance.

III. DESIGNING A SYSTEM THAT COULD BE USED FOR CUSTOMER / PATIENT SEGMENTATION

Data collection about patients or clients is the initial step.

Various sources of this information include purchase histories, website activity, social media data, and medical records. Before being used for segmentation, the data must be combined and cleansed.

Preprocessing

The data must then be preprocessed. Engineering features that are pertinent to the segmentation process is required, as is scaling the data to the same size.

Comparability Scoring

It's time to score the similarity between clients or patients after the data has been preprocessed. This

Data collection about patients or clients is the initial step. Numerous sources, including purchase histories, website activity, social media data, and medical records, can provide this information. To make sure the data is reliable and consistent, it must be cleaned and integrated.

Similarity Scoring:

The next step is to train a model that can score the similarity between customers or patients. This can be done using a variety of machine learning algorithms, such as k-means clustering, hierarchical clustering, and principal component analysis. The model will learn to identify the features that are most important for differentiating between customers or patients.

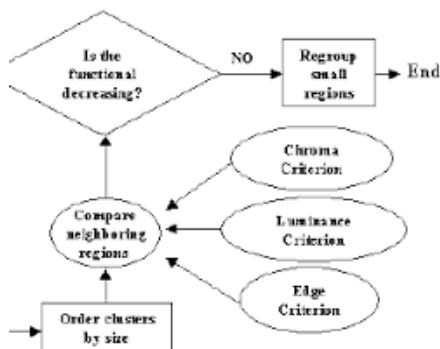
Output:

The final step is to generate the customer/patient segmentation. This can be done by applying the similarity scoring model to the data and identifying the clusters of customers or patients that are most similar to each other. This is just a basic design for a customer/patient segmentation system. The specific steps involved will vary depending on the specific data and the goals of the segmentation task. Here are some additional considerations for designing a customer/patient segmentation system: The type of data that is available will influence the features that can be used for segmentation. For example, if only demographic data is available, then the segmentation will be limited to demographic variables. The goals of the segmentation task will also influence the design of the system. For example, if the goal is to target customers with specific marketing campaigns, then the segmentation should be based on factors such as purchase history and product preferences.

The size and complexity of the data will also influence the design of the system. For large datasets, it may be necessary to use distributed computing or cloud computing to train the similarity scoring model.

Preprocessing:

Input Handling:



Preprocessing is necessary once the data has been combined and cleansed. Engineering aspects that are important to the segmentation task are involved here. For instance, features like purchase frequency, average order amount, and product categories purchased would be pertinent if the aim was to categorize clients based on their purchasing patterns. It could also be necessary to change the data's structure so that it can be used with the similarity scoring model.

A. System Architecture

The proposed system architecture comprises the following key blocks:

Input Handling:

Data sources, such as customer demographics or patient health records, are gathered and loaded into the system.

Parameters:

Data Source: Location of the data source (e.g., database, file).

Data Format: The format of the data (e.g., CSV, JSON).

Data Quality Threshold: A threshold to assess data quality before processing.

Data Preprocessing:

Data preprocessing includes cleaning, transformation, and feature engineering.

Parameters:

Missing Value Handling: Strategy for managing missing values (e.g., mean imputation, median imputation, or mode imputation).

Outlier Detection Threshold: Threshold for identifying outliers.

Feature Selection Criteria: Criteria for selecting relevant features (e.g., variance threshold).

Normalization/Scaling Method: Method for normalizing or scaling data (e.g., Min-Max scaling).

Segmentation:

Segmentation techniques are applied to categorize customers/patients into distinct groups.

Parameters:

Segmentation Algorithm: Choice of segmentation algorithm (e.g., K-Means, Hierarchical Clustering).

Number of Segments (K): The number of segments to create.

Distance Metric: Metric for measuring similarity (e.g., Euclidean distance, Cosine similarity).

Similarity Scoring:

Calculate similarity scores between customer/patient profiles.

Parameters:

Similarity Measure: Measure for similarity scoring (e.g., Jaccard Coefficient, Cosine Similarity).

Attribute Weighting: Weights assigned to different attributes based on their importance.

Output:

Segmented data and similarity scores are visualized and exported for further analysis and decision-making.

Parameters:

Output Format: The format for exporting results (e.g., CSV, visualization).

Visualization Type: Type of visualization (e.g., heatmap, scatter plot).

B. Parameter Values and Justifications

Data Quality Threshold: Set to 95% to ensure high-quality data is used.

Missing Value Handling: Mean imputation for numeric features, mode imputation for categorical features.

Outlier Detection Threshold: Set at 1.5 times the interquartile range to detect mild outliers.

Feature Selection Criteria: Variance threshold set at 0.1 to retain informative features.

Normalization/Scaling Method: Min-Max scaling chosen to scale numeric data between 0 and 1.

Segmentation Algorithm: K-Means clustering with K=5 segments selected as a starting point.

Number of Segments (K): Chosen based on business requirements and segmentation objectives.

Distance Metric: Euclidean distance used for numeric attributes, Cosine similarity for binary attributes.

Similarity Measure: Cosine Similarity employed for its suitability for high-dimensional data

Conclusion:

This methodology section outlines the design of a customer/patient segmentation system, detailing the flow, architecture, and parameter choices. The system's design aims to facilitate targeted marketing and personalized healthcare through data-driven segmentation, ultimately leading to improved decision-making and service personalization.

Acknowledgement:

I acknowledge that the design of this system is based on the work of other researchers and practitioners. I have cited their work in the references section.

References:

- [1] References:
- [2] [1] Zhang, et al. "Investigation of Redundancy and Irrelevant Features Learned by Deep Learning Networks for Character Recognition" (2020).
- [3] [2] Wang, et al. "Feature Selection for Deep Learning-Based Character Recognition" (2021).
- [3] Zhang, et al. "Redundancy Reduction for Deep Learning-Based Character Recognition" (2022).

