

**Auto Encoder Fixed-Target Training Features
Extraction Approach for Binary Classification
Problems**

ABSTRACT

The main issues with machine learning-based feature extraction techniques are the requirement of extensive domain-level knowledge, experience, and the need to be supported by large amounts of data that are sometimes not available. Moreover, it is often difficult to apply domain-level knowledge to extract the necessary features for building a machine-learning classifier. Therefore, it is significantly important to find and develop feature extraction techniques that depend mainly on the training data and don't require or depend on domain-level knowledge and experience. To address these issues for binary classification problems, a novel feature extraction approach, *AE-FT(Fixed Target)* for extracting common features using a Deep Belief Network (DBN)-based Autoencoder (AE) is proposed in this paper. In this approach, common features are extracted by a DBN trained on a dataset sample's binary using the *Fixed Target training approach*.

The proposed common features extraction approach is tested and evaluated on two different data sets. For each dataset, the extracted features are used to train seven of the common machine learning binary classification algorithms and compared their performances. Moreover, the number of extracted features is very small compared to other existing feature extraction methods. Therefore, the proposed common features extraction method improves the performance of the binary classification algorithms by reducing the number of features reducing laborious processes, and increasing the recognition accuracy effectively.

The results show that the proposed common features extraction approach, without any domain-level knowledge or human expertise, provides a very good performance compared to other feature extraction techniques.

11
12
13
14
15

Keywords: Features Extraction; Deep Learning; AE; Fixed-Target Training; Common Features.

16
17

1. INTRODUCTION

18
19
20
21
22

With the rapid development of machine learning technology, as a binary classification problem that helps people to find the law from the massive data to achieve the prediction effectively, data prediction has become an important part of people's daily lives. Feature extraction is a basic and important matter for the classification problem because the original data contain noise and irrelevant information which decreases the classification accuracy.

23 Feature extraction is about finding a good data representation, which is very domain-
24 specific, often requires human expertise, and is related to available measurements. The
25 primary idea behind feature extraction is to compress the data to maintain most of the
26 relevant information.

27 As to feature selection techniques, these techniques are also used for reducing the number
28 of features from the original feature set to reduce model complexity, and model overfitting,
29 enhance model computation efficiency, and reduce generalization error. Therefore, improve
30 the accuracy of the learning algorithm and shorten the training and output time.

31 The feature extraction methods are useful for different applications as mentioned in [1], such
32 as social science, healthcare, environment, agriculture, spam filtering, antivirus technology,
33 economics, medical diagnosis, face recognition, action recognition, speech recognition,
34 gesture recognition, marketing, wireless network, gene expression, software fault detection,
35 internet traffic prediction, etc. Therefore, the research of machine learning algorithms in
36 feature extraction problems is a research hotspot in recent years.

37 The main issue with machine learning-based feature extraction techniques is the
38 requirement of time, extensive domain-level knowledge, and experience as mentioned by
39 Verdonck et al. [2]. Moreover, it is often difficult to apply domain-level knowledge to extract
40 the necessary features for building a machine-learning classifier. Therefore, it is significantly
41 important to find and develop feature extraction techniques that depend mainly on the
42 training data and don't require or depend on domain-level knowledge and experience, and
43 this is our main purpose. Therefore, the focus of this paper is on using machine learning for
44 common features extraction that can be used in binary classification in general, that
45 applicable in many fields such as spam filtering, antivirus technology, and medical diagnosis.

46 This paper presents the use of denoising stacked autoencoders with supervised fixed-target
47 training in order to extract the common features of the training data that can be used in
48 binary classification. The method relies on training a deep belief network (DBN) [3], i.e., a
49 deep unsupervised neural network implemented with a deep stack of denoising
50 autoencoders, in a supervised manner to create an invariant compact representation of the
51 general behavior of the training datasets. In recent years DBNs have proven successful in
52 generating invariant representations for many challenging domains. We used only positive
53 training samples for training the common features extractor. Then, we used the common
54 features extractor for extracting the values of the common features of the training samples
55 and used it for training the binary classifiers.

56 In contrast to most existing approaches that normally have a separate stage for data
57 preprocessing followed by domain-dependent feature extraction. We developed a domain-
58 independent deep neural network framework for common feature extraction which enables
59 us to easily, without the need for domain-level knowledge or expertise, extract features that
60 can be effectively used in binary classification problems.

61 We trained the proposed feature extractor, using the binary representation of the datasets,
62 to extract the common features and then used it for training binary classification models. In
63 the experiments, we used the extracted common features to build binary classifiers using
64 seven binary classification methods Naive Bayes, Logistic Regression, K-Nearest
65 Neighbors, Support Vector Machine, Decision Tree, Random Forest, and Voting
66 Classification for binary classification.

67 The proposed method, fixed-target training of a deep stacked autoencoder, enabled
68 a good recognition accuracy, better generalization, and more stability than that which could
69 be achieved with the other methods. The proposed approach achieves 73.50% accuracy,
70 which is so far a good result that does not need any domain expertise.

71 The remainder of this paper is organized as follows: The next section presents the study
72 background, Section 3 describes our proposed approach, and Section 4 presents the
73 experimental results. In section 5, we discuss and evaluate the results of the study and
74 present our conclusions in Section 6, while presenting some directions for future studies.

75

76 **2. BACKGROUND**

77 The success of machine learning often depends strongly on the success of feature
78 extraction, the features used influence the result more than everything else. No algorithm
79 alone can supplement the information gain given by correct feature extraction. So, feature
80 extraction is a basic and important matter for the classification problem because the original
81 data contain noise and irrelevant information which decreases the classification accuracy.

82 There are two broad categories for feature extraction algorithms: linear and nonlinear. Linear
83 feature extraction assumes that the data lies in a linear subspace. Use matrix factorization to
84 protect them. On the other hand, in nonlinear feature extraction ordimensionality reduction, a
85 low-dimensional surface can be mapped into a high-dimensional space so that a nonlinear
86 relationship among the features can be found and easily detected. Theoretically, a
87 transformation function $f(x)$ can be used to map the features into a higher-dimensional space
88 and then mapped back into the lower-dimensional space, so that the relationship can be
89 viewed as nonlinear. We focus on the nonlinear feature extraction algorithms.

90 **2.1. Kernel Principal Component Analysis (KPCA)**

91 KPCA introduced by Scholkopf et al. [4], is an extension of Principal Component
92 Analysis (PCA) that allows for the separability of nonlinear data by making use of kernels.
93 The basic idea behind it is to project the linearly inseparable data onto a higher dimensional
94 space where it becomes linearly separable. Unfortunately, it has a serious limitation in terms
95 of space complexity since it stores all dot products of the training set and therefore the size
96 of the matrix increases quadratically with the number of data points as presented in [5].

97 Another drawback of the KPCA, however, is the cost of computation could be
98 extremely high, which could lead to the attendant numerical problems of diagonalizing large
99 matrices, which limits its applicability in many large dataset problems. But, an Expectation-
100 Maximization (EM) algorithm for KPCA to overcome these drawbacks was proposed in [6],
101 which is an expectation-maximization approach for performing kernel principal component
102 analysis. Experimental results showed that EM is an efficient method computationally,
103 especially for a large number of data points.

104 **2.2. Locally Linear Embedding (LLE)**

105 Locally Linear Embedding, proposed by Saul et al. [7], is a dimensionality reduction
106 technique based on Manifold Learning that involves the computation of low-dimensional
107 neighborhood preserving embeddings of inputs that are of high dimension in nature.
108 Manifold Learning aims to make a manifold object, an object of D dimensions that is
109 embedded in a higher-dimensional space, representable in its original D dimension instead
110 of being represented in an unnecessarily greater space.

111 LLE has the ability to learn the global structure of nonlinear manifolds like those from
112 images of faces or documents of text by exploiting the local symmetries of linear
113 reconstructions. LLE has been applied successfully in a wide range of applications which
114 includes face recognition and remote sensing, MRI, shape analysis of the hippocampus in
115 AD, diffusion tensor imaging, breast lesion segmentation, feature fusion, and image
116 classification according to [8].

117 LLE is popular among researchers because of its ability to deal with large data sets
118 of high-dimensional data and its non-iterative way of finding embeddings. However, it has
119 some drawbacks which include sensitivity to noise, the inability to deal with novel data, and
120 the inevitable ill-conditioned Eigen problems. Some efforts have recently been made to
121 develop extensions of the classical LLE.

122 Supervised and semi-supervised versions of LLE were proposed by [9] and [10],
123 respectively, for plant classification based on images of leaves.

124 **2.3. Linear Discriminant Analysis (LDA)**

125 LDA is a supervised learning dimensionality reduction, feature extraction technique,
126 and Machine Learning classifier that was invented by Fisher et al. [11]. LDA uses within-
127 classes and between-classes measures by maximizing the distance between the mean of
128 each class and minimizing the spreading within the class itself. This is a good choice
129 because maximizing the distance between the means of each class when projecting the data
130 in a lower-dimensional space can lead to better classification results.

131 An advantage of LDA is that it is able to use information from both features to create
132 a new axis which in turn minimizes the variance and maximizes the class distance of the
133 variables. Although the LDA is one of the most well-used data reduction techniques, it has
134 some limitations. The small sample problem (SSS), is one of the main problems of LDA,
135 which happens when the dimensions are much higher than the number of samples in the
136 data matrix, LDA is unable to find the lower dimensional space resulting in the within-class
137 matrix becoming singular. Different approaches have been proposed to solve this problem,
138 such as what was proposed in [12] and [13]. In addition to the assumption that the input data
139 follows a Gaussian Distribution, therefore applying LDA to not Gaussian data can lead to
140 poor classification results.

141 A semisupervised variant of LDA, which performed better than the classical LDA,
142 was proposed by Zhang et al. [14] that mainly combines both labeled and unlabeled data for
143 training LDA and allows using LDA for the situation where the labeled data are few.

144 Application of LDA includes facial recognition, text recognition, automatic diagnosis
145 of machine operations, early detection of diseases, person reidentification, hand movement
146 classification, motor imagery EEG, and groundwater redox conditions.

147 **2.4. t-distributed Stochastic Neighbor Embedding (t-SNE)**

148 t-Stochastic Neighbor Embedding (t-SNE) is an unsupervised Non-linear Dimension
149 Reduction Technique (NLDRT) that was introduced by Maaten et al. [15]. The technique is a
150 variation of the Stochastic Neighbor Embedding introduced by Hinton et al.[10], whose main
151 objective is the construction of probability distributions from pairwise distances such that
152 larger distances correspond to smaller probabilities and vice versa. t-SNE is typically used to
153 visualize high-dimensional datasets, it works by minimizing the divergence between a
154 distribution constituted by the pairwise probability similarities of the input features in the
155 original high-dimensional space, which is modeled using a Gaussian Distribution and its
156 equivalent in the reduced low-dimensional space, modeled using a Student's t-distribution.

157 t-SNE makes use of the Kullback-Leiber (KL) divergence in order to measure the
158 dissimilarity of the two different distributions, as mentioned in [16]. The KL divergence is
159 then minimized using gradient descent. t-SNE is the most commonly used in single-cell
160 analysis. However, it has some limitations as mentioned in [17]. The limitations include slow
161 computation time, the inability to meaningfully represent very large datasets, and the loss of
162 large-scale information.

163 **2.5. Deep Learning Approach**

164 The major difference between deep learning and traditional pattern recognition
165 methods is that deep learning automatically learns features from big data, instead of
166 adopting handcrafted features, as stated in [18]. Deep learning is able to quickly acquire new
167 effective feature representations from training data.

168 In recent years DBNs [3], deep unsupervised neural networks, have proven successful in
169 generating invariant representations for many challenging domains. Autoencoders are feed-
170 forward DBNs that were first introduced by Rumelhart et al. [19]. They can learn a
171 compressed and distributed representation of data, which can be used as a dimensionality
172 reduction or feature extraction technique. They use nonlinear transformations to project data
173 from a high dimension to a lower one. An autoencoder usually has at least one hidden layer
174 between the input and output layers. The number of neurons in the hidden layer is usually
175 set to less than those in the input and output layers, thus creating a bottleneck, with the
176 intention of forcing the network to learn a higher-level representation of the input as
177 presented in [20].

178 Autoencoders are typically trained in an unsupervised manner, using
179 backpropagation with stochastic gradient descent, to approximate a function by which data
180 can be classified, as mentioned in [21]. For every training input, the difference between the
181 input and the output is measured (using squared error) and it is back-propagated through the
182 neural network to perform weight updates on the different layers.

183 Compared with other machine learning methods, deep learning is able to detect
 184 complicated interactions in features, learning lower-level features from nearly unprocessed
 185 original determine characteristics that are not easy to be detected. Furthermore, they hand
 186 class members with high cardinal numbers and process untapped data. Unfortunately, if all
 187 input features are independent of each other, then the autoencoder will find it particularly
 188 difficult to encode the input data into a lower-dimensional space.

189 The advantages are higher discriminating power and control overfitting when it is
 190 unsupervised. On the other hand, there are some bottlenecks in deep learning based feature
 191 extraction methods, time-consuming data pre-processing, domain expertise, the need for
 192 large amounts of data, loss of data interpretability, and transformation may be expensive.

193 There are many other deep learning-based feature extraction approaches. A
 194 skeleton-based abnormal gait recognition approach was proposed by Jun et al. [22]. They
 195 proposed a feature extraction method using the RNN AEs to minimize the irrelevant
 196 information of the original skeleton data. They used two-step training of a hybrid RNN and
 197 AE-DM model and approved that it is more effective than the single-step training of the End-
 198 to-End model that has the same data flow. Ma and Yuan [23] proposed a method for
 199 extracting features from images based on deep CNN and PCA. They used a neural network
 200 to extract features and a PCA algorithm for feature dimension reduction. Then they
 201 compared the performance of the PCA before and after the improvement claim achieving
 202 memory, and time optimization. Moreover, the SVM classifier accuracy was enhanced.
 203 Dahouda et al. [24] proposed a deep learning-based feature extraction approach with a
 204 modular neural network, they employed a pre-trained neural architecture search net
 205 (NASNet) as a feature extractor on a custom dataset of raw copper and cobalt image. Then,
 206 they used the extracted features to build a deep neural network and machine learning
 207 algorithms for the image classification of copper and cobalt raw minerals. However, it is an
 208 empirical not an exhaustive study, and the data preprocessing was ignored. Petrovska et al.
 209 [25] used pre-trained neural networks to extract features, then applied PCA to reduce the
 210 dimensionality of the extracted features. However, pre-trained neural networks were used in
 211 both [24] and [25]. Moreover, these are domain-dependent approaches that work only for the
 212 specific domain not for binary classification problems in general.

213

214 **Table 1. Comparison of Feature Extraction Methods.**

Feature Extraction Technique	Domain level Knowledge	Data preprocessing	Limitations
Kernel PCA	yes	yes	space complexity
LLE	yes	yes	sensitivity to noise, the inability to deal with novel data and the inevitable ill-conditioned Eigen problems
LDA	yes	yes	small sample problem (SSS)
t-SNE	yes	yes	slow computation time
AE	often	yes	Loss of data interpretability Transformation may be expensive

215 The main issue with machine learning-based feature extraction techniques is the
216 requirement for extensive domain-level knowledge and experience. Moreover, it is often
217 difficult to apply domain-level knowledge to extract the necessary features for building a
218 machine-learning classifier. Table 1 shows that all the current approaches have the same
219 problems. Therefore, it is significantly important to find and develop feature extraction
220 techniques that depend mainly on the training data and don't require or depend on domain-
221 level knowledge and experience, and this is our main purpose.

222 3. METHODOLOGY

223

224 In this section, we describe and discuss our proposed novel deep learning-based approach
225 for common feature extraction, our datasets, and training methods in detail. The main
226 question we are trying to answer is the following:

227 *Is it possible to extract the common features from the raw binary representations without any*
228 *domain expertise of a given dataset that could be used in binary classification?*

229 In recent years, deep learning methods have proven very successful in
230 accomplishing dimensionality reduction and feature extraction tasks in many domains,
231 especially computer vision, and cybersecurity according to [27, 28, 29, 30]. The proposed
232 methodology works as follows:

233 Firstly, training the proposed common features extraction model to extract features
234 to use in binary classification using various common machine learning classification
235 algorithms. Secondly, select seven binary classification algorithms and use the extracted
236 common features to build binary classifiers using common algorithms.

237 Our method uses stacked denoising autoencoders for extracting the common
238 features of the training dataset that can be used effectively in binary classification. The input
239 to the DBN has a fixed length, but the dataset sample length is variable. In order to
240 represent the dataset binary as a fixed-sized vector, which would be the input to the neural
241 network, we repeatedly pad the sample binary until the specified size is met. This process
242 gives better results than 0's or 1's padding that the DBN learns as a common feature. Then
243 we use two different datasets, the IMDB dataset [31] and the Enron-Spam dataset [32], to
244 train and test the proposed feature extractor for extracting the common features and then
245 used for training binary classification models using the extracted features of the binary
246 representation of the datasets. We focus on the top seven most common binary
247 classification algorithms Naive Bayes, Logistic Regression, K-Nearest Neighbours, Support
248 Vector Machine, Decision Tree, Random Forest, and Voting Classification.

249 3.1. Deep Belief Network Fixed-Target Training

250 In stacked denoising autoencoders, first introduced by Vincent et al. [26], the data at
251 the input layer is replaced by noised data while the data at the output layer stays the same;
252 therefore, the autoencoder can be trained with much more generalization power, according
253 to [18].

254 Usually, denoising autoencoder training is unsupervised according to Vincent et al.
255 [26]. The input sample is corrupted by adding noise (or more often by zeroing the values).
256 That is, given an input \mathbf{S} , first it is corrupted to \mathbf{S}^{\wedge} and then fed to the input layer of the
257 network. The objective function of the network in the output layer remains to generate the

258 uncorrupted version of the input (see Figure 1-(a)). But in our approach, we use a novel
259 supervised training strategy, which we call a *fixed-target training* strategy.

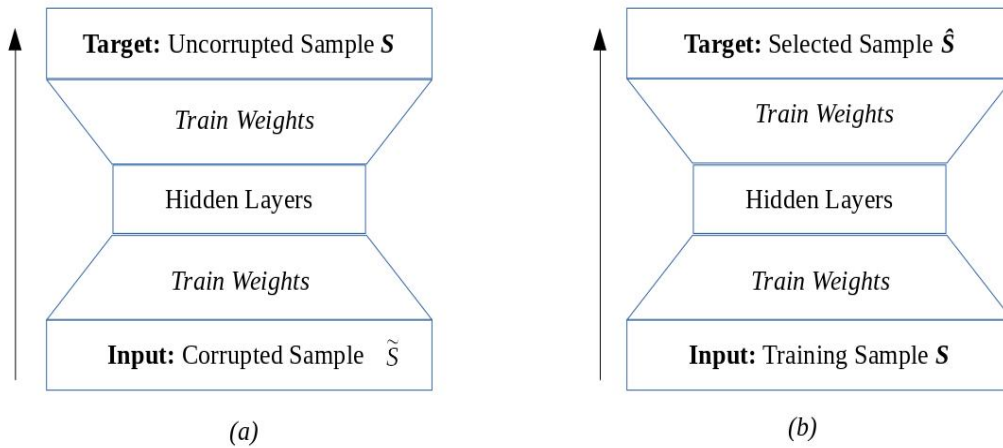


Fig. 1. Comparison between Denoising autoencoder training. (a)-traditional unsupervised training, (b)-fixed-target supervised training.

260 In the *fixed-target training* strategy, we randomly select one of the training samples
261 \hat{S} , fix it in the output layer, and for every input S of the training samples, the objective
262 function of the network is to generate \hat{S} , (i.e. we consider all training samples as corrupted
263 versions of the selected sample \hat{S} (see Figure 1-(b)). So, the “hidden units” of DBN
264 compute internal representations analogous to the extracted common features.

265 This training approach works better than traditional training. By fixing the target
266 output, the network is forced to generalize better and determine more high-level common
267 patterns. Moreover, the network learns better even when few training samples are available.
268 When a DBN’s training is complete, we discard the decoder layer, fix the values of the
269 encoder layer, and use the encoder as the common features extractor. In a typical
270 implementation, the extracted features may then be used for supervised binary classification.

271 In order to achieve our goal, we create a DBN by training a deep stack of denoising
272 autoencoders. We use fixed-target training to train a deep denoising autoencoder consisting
273 of five layers: 8,192–2,048–512–128–32. At the end of this training phase, we have a deep
274 network that is capable of converting 8,192 length input vectors into 32 floating point feature
275 values. Note that the network is trained only using the samples in the training set, and for all
276 future samples it will be run in prediction mode;(i.e., receiving the 8,192-sized vector it will
277 produce 32 output values, without modifying the weights). See Figure 2.

278

279

280

281

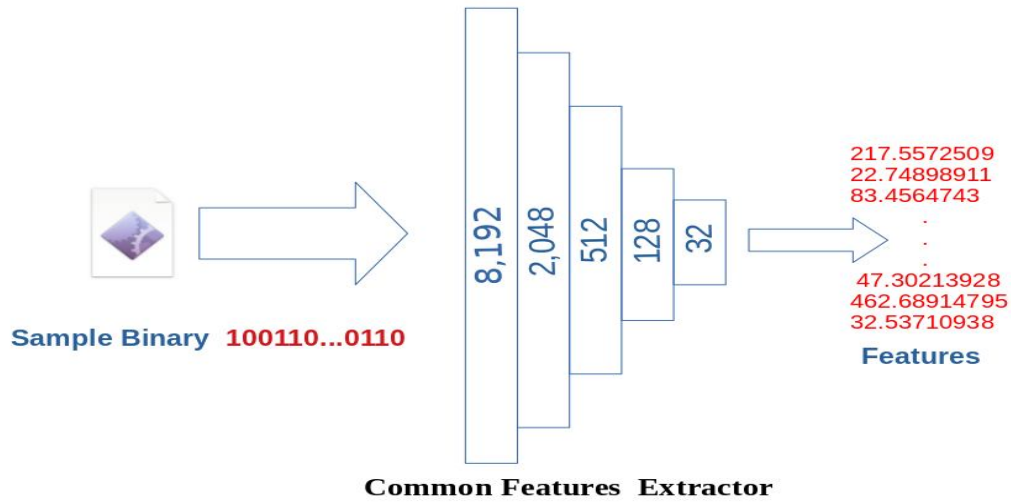


Fig. 2. Illustration of the common features extraction stages from feeding the sample binary to features extraction using DBN.

282 Our goal is to train the proposed feature extractor for extracting the common
 283 features and then used it for training binary classification models using the extracted
 284 features of the binary representation of the datasets. We focus on the top seven most
 285 common binary classification algorithms Naive Bayes, Logistic Regression, K-Nearest
 286 Neighbours, Support Vector Machine, Decision Tree, Random Forest, and Voting
 287 Classification. The sample's binary bit string is fed to the neural network, and the deep
 288 neural network generates a 32-sized vector at its output layer, which we treat as the
 289 common feature values of the sample.

290 3.2. Datasets

291 We are using two different datasets, the Internet Movie Database (IMDB) dataset
 292 [31] and the Enron-Spam dataset [32], for testing our approach. The IMDB dataset, Large
 293 Movie Review Dataset v1.0, of highly polar movie reviews in the form of text comments on
 294 different movies and a positive or negative score. This dataset contains movie reviews along
 295 with their associated binary sentiment polarity labels. It is intended to serve as a benchmark
 296 for sentiment classification. The core dataset contains 50,000 reviews split evenly into 25k
 297 train and 25k test sets. The overall distribution of labels is balanced (25k pos and 25k neg).
 298 Furthermore, the train and test sets contain a disjoint set of movies, so no significant
 299 performance is obtained by memorizing the movie's unique terms and they are associated
 300 with the observed labels. In the labeled train/test sets, a negative review has a score < 4 out
 301 of 10, and a positive review has a score > 7 out of 10. Thus, reviews with more neutral
 302 ratings are not included in the train/test sets.

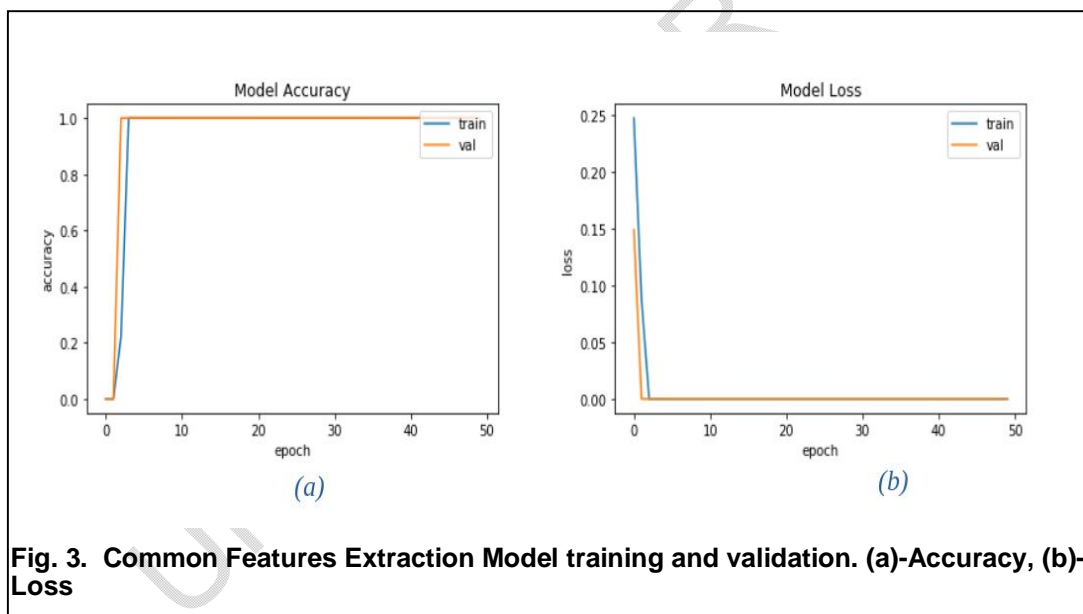
303 These datasets have been used as a benchmark in several research papers on
 304 spam filtering, text classification, and natural language processing as mentioned in [33, 34].
 305 Therefore, the results of this work are hopefully comparable to other similar works within the
 306 area, without having to account for unique datasets.

307 3.3. Common Features Extraction Method

308 As described in the previous section, we use a fixed-target training strategy to train a
309 deep denoising autoencoder consisting of five layers (8,192–1,024–512–128–32). We
310 randomly pick one of the training samples as a pivot sample, and each time a new input is
311 given to the network we put the pivot sample in the output as the target. This way, we
312 enforce the autoencoder to learn the common features of the training samples.

313 *Fixed-target training* is essential for the network to learn common features in a very
314 small time and even when few samples are available. We use rectified linear units (ReLU)
315 for the nonlinearity function when training deep neural networks to diminish the gradient
316 vanishing problem and result in faster convergence, as approved in [1, 20].

317 We build the model using machine learning algorithms in *the Keras library*. Other
318 parameters we use are 50 training epochs, a learning rate of 0.003, a batch size of 10, and
319 no momentum. Figure 3 shows the learning curves, which approve the feasibility of our
320 proposed common feature extraction approach.



321 EXPERIMENTAL RESULTS

322 In this section, we conducted various experiments to evaluate the features of the
323 DBN AEs for binary classification. First, we compared the number of features used in [31,
324 32] and the extracted common features. Then, we conducted experiments to show the
325 effectiveness of the proposed DBN AEs. Finally, we showed the performance of the features
326 of the DBN AEs used in binary classification problems, we used seven binary classification
327 methods Naive Bayes, Logistic Regression, K-Nearest Neighbours, Support Vector Machine,
328 Decision Tree, Random Forest, and Voting Classification- for binary classification.

329 In the experiments, we trained a deep denoising autoencoder consisting of five
330 layers, 8192–2048–512–128–32, using a Fixed-Target training strategy. We used only
331 positive training samples for training the common features extractor. Then, we used the
332 common features extractor for extracting the values of the common features of the training
333 samples and used it for training the binary classifiers.

334

In the testing phase, we first extract the values of the common features of the

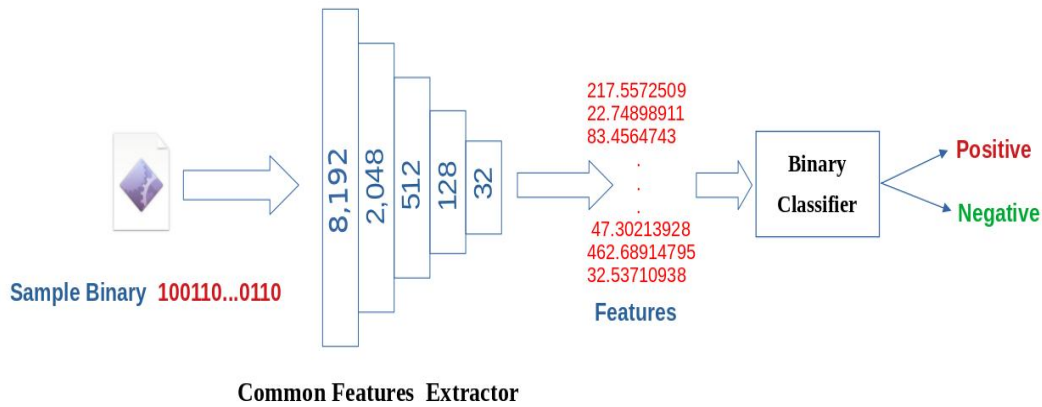


Fig. 4. Overview of our deep learning approach for common features extraction, illustration of all stages from feeding the sample binary to the Common Features Extractor ending with the decision made by the binary classifier.

335

sample and pass them to the binary classifier for classification. See Figure 4.

336

We used the training score and test score to evaluate our model, the test score measures how well our model learn from our training data, while the test score measures the accuracy. Higher the test score better the model is generalized. The results demonstrate that, with proper structure and parameters, the performance of the proposed deep learning method on common feature extraction is useful even in the lack of domain expertise in binary classification.

342

Table 2. Performance of Binary Classification Methods, trained by the common features on IMDB and Enron-Spam datasets.

343

Binary Classification Method	IMDB Dataset		Enron-Spam Dataset	
	Training Score	Test Score	Training Score	Test Score
Naive Bayes	63.0%	60.7%	63.0%	66%
Logistic Regression	88.0%	73.5%	87.0%	73%
K-Nearest Neighbours	87.3%	73.2%	88.7%	73.2%
Support Vector Machine	67.0%	66.3%	67.3%	67%
Decision Tree	100%	70.5%	100%	72%
Random Forest	93.5%	73.4%	93.5%	73.5%
Voting Classification	85.0%	72.7%	84.7%	73%

344

The accuracy results are given in Table 2, test samples are 73.5% correctly classified by conducting a Logistic Regression classifier for the proposed approach with the extracted common features of the IMDB dataset. For the Enron-Spam dataset, the best results are obtained by Random Forest 73.5%. The binary classification performances of the approaches are given in Table 2. According to the results, the common features extracted by

349 the proposed approach give a very good performance for all seven binary classification
350 methods. Logistic Regression and Random Forest are the most effective classifiers for this
351 approach.

352 Compared with other methods, the proposed approach fixed-target training of a
353 denoising deep stacked autoencoder achieves 73.50% accuracy, which is very good for
354 binary classification without any need for domain-level expertise or data preprocessing.

355 **5. DISCUSSION AND EVALUATION**

356 We use the proposed fixed-target training strategy to train a deep denoising
357 autoencoder consisting of five layers (8,192–1,024–512–128–32). We randomly pick one of
358 the training samples as a pivot sample, and each time a new input is given to the network we
359 put the pivot sample in the output as the target. This way, we enforce the autoencoder to
360 learn the common features of the training samples. Unfortunately, it is not easy to find the
361 optimal autoencoder structure.

362 We trained and tested, using two different datasets, the proposed feature extractor
363 for extracting the common features. Then, we used the extracted features for training seven
364 of the most common binary classification algorithms: Naive Bayes, Logistic Regression, K-
365 Nearest Neighbours, Support Vector Machine, Decision Tree, Random Forest, and Voting
366 Classification.

367 The learning curves in Figure 3 approve the feasibility of our proposed common
368 feature extraction approach. According to the results reported in Table 2, the best accuracy
369 values are obtained by conducting Logistic Regression and Random Forest classifiers.
370 When comparing the number of features used, 32 common features were extracted using
371 the proposed feature extraction algorithm, on the other hand, 3000 and 10000 features were
372 used in [31] and [32], respectively, in addition to the preprocessing carried to extract the
373 features used in [31] and [32]. However, by applying the proposed common features
374 extraction algorithm, it is observed that high success rates are achieved with very few
375 features and increase the overall process performance.

376 **6. CONCLUSION**

377 In this paper, we reviewed past approaches for feature extraction and proposed a
378 novel method based on deep belief networks for common features extraction. Current
379 approaches for feature extraction are time-consuming and require extensive domain-level
380 knowledge and experience. Therefore, it is significantly important to find and develop
381 feature extraction techniques that depend mainly on the training data and don't require or
382 depend on domain level knowledge and experience.

383 Our proposed feature extraction approach, AE fixed-target supervised training,
384 extracts the common characteristics of the original data and minimizes the irrelevant
385 information without the need for any domain-level expert or expensive data preprocessing.
386 Furthermore, it needs less training data and produces fewer features compared to other
387 feature extraction techniques. Therefore, the extracted features improve the binary
388 classification performance. The most interesting result of our evaluation was the very good

389 performance of our common features extension approach using all binary classification
390 algorithms tested on both datasets that have been used, but further experiments are needed
391 to be more confident.

392 Future work could include examining how to improve the accuracy of the proposed
393 common feature extraction method by finding the optimal autoencoder structure and
394 activation function. Furthermore, examining different ways for the pivot sample selection. Of
395 course, we could also look at different types of classifiers and use different datasets in
396 different domains. One could also see, whether it is useful to use the proposed method in
397 data collection; when only a few samples are available. Finally, it would be interesting to see
398 if the proposed method could mitigate the problem of the need for domain expertise and data
399 preprocessing for feature extraction and the need for a big dataset to train binary classifiers.

400

401

402 Data Availability: Two Datasets are used, The IMDB and Enron-Spam data supporting this
403 research are from previously reported studies and datasets, which have been cited. The
404 processed data are available at <http://www.cs.cornell.edu/people/pabo/movie-review-data/>
405 and <http://www2.aueb.gr/users/ion/data/enron-spam/>, respectively.

406

407 ACKNOWLEDGEMENTS

408

409 Funding: The study attracted no funding.

410 Conflict of Interest: All authors declare that they have no conflict of interest.

411

412 REFERENCES

- 413 1. Ghojogh, B., Samad, M. N., Mashhadi, S. A., Kapoor, T., Ali, W., Karray, F., & Crowley, M.
414 (2019). Feature selection and feature extraction in pattern analysis: A literature review. *arXiv*
415 *preprint arXiv:1905.02845*.
- 416 2. Verdonck, T., Baesens, B., Óskarsdóttir, M., & vanden Broucke, S. (2021). Special issue on
417 feature engineering editorial. *Machine Learning*, 1-12.
- 418 3. Hinton, G. E., Osindero, S., & Teh, Y. W. (2006). A fast learning algorithm for deep belief
419 nets. *Neural computation*, 18(7), 1527-1554.
- 420 4. Schölkopf, B., Smola, A., & Müller, K. R. (1998). Nonlinear component analysis as a kernel
421 eigenvalue problem. *Neural computation*, 10(5), 1299-1319.
- 422 5. Liu, X., & Yang, C. (2009, October). Greedy kernel PCA for training data reduction and nonlinear
423 feature extraction in classification. In *MIPPR 2009: Automatic Target Recognition and Image*
424 *Analysis*(Vol. 7495, pp. 768-775). SPIE.
- 425 6. Rosipal, R., & Girolami, M. (2001). An expectation-maximization approach to nonlinear component
426 analysis. *Neural Computation*, 13(3), 505-510.
- 427 7. Saul, L. K., & Roweis, S. T. (2003). Think globally, fit locally: unsupervised learning of low
428 dimensional manifolds. *Journal of machine learning research*, 4(Jun), 119-155.

- 429 8. Nanga, S., Bawah, A. T., Acquaye, B. A., Billa, M. I., Baeta, F. D., Odai, N. A., ... & Nsiah, A. D.
430 (2021). Review of Dimension Reduction Methods. *Journal of Data Analysis and Information*
431 *Processing*, 9(3), 189-231.
- 432 9. Kouropteva, O., Okun, O., & Pietikäinen, M. (2005, June). Incremental locally linear embedding
433 algorithm. In *Scandinavian Conference on Image Analysis*(pp. 521-530). Springer, Berlin,
434 Heidelberg.
- 435 10. Pan, Y., Ge, S. S., & Al Mamun, A. (2009). Weighted locally linear embedding for dimension
436 reduction. *Pattern Recognition*, 42(5), 798-811.
- 437 11. Xanthopoulos, P., Pardalos, P. M., & Trafalis, T. B. (2013). Linear discriminant analysis. In *Robust*
438 *data mining*(pp. 27-33). Springer, New York, NY.
- 439 12. Sharma, A., & Paliwal, K. K. (2012). A new perspective to null linear discriminant analysis method
440 and its fast implementation using random matrix multiplication with scatter matrices. *Pattern*
441 *Recognition*, 45(6), 2205-2213.
- 442 13. Belhumeur, P. N., Hespanha, J. P., & Kriegman, D. J. (1997). Eigenfaces vs. fisherfaces:
443 Recognition using class specific linear projection. *IEEE Transactions on pattern analysis and*
444 *machine intelligence*, 19(7), 711-720.
- 445 14. Zhang, Y., & Yeung, D. Y. (2011). Semisupervised generalized discriminant analysis. *IEEE*
446 *Transactions on Neural Networks*, 22(8), 1207-1217.
- 447 15. Van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of machine*
448 *learning research*, 9(11).
- 449 16. Hinton, G. E., & Roweis, S. (2002). Stochastic neighbor embedding. *Advances in neural*
450 *information processing systems*, 15.
- 451 17. Xie, B., Mu, Y., Tao, D., & Huang, K. (2011). m-SNE: Multiview stochastic neighbor embedding.
452 *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 41(4), 1088-1096.
- 453 18. Wang, H., & Raj, B. (2017). On the origin of deep learning. *arXiv preprint arXiv:1702.07800*.
- 454 19. Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1985). *Learning internal representations by error*
455 *propagation*. California Univ San Diego La Jolla Inst for Cognitive Science.
- 456 20. David, O. E., & Netanyahu, N. S. (2015, July). Deepsign: Deep learning for automatic malware
457 signature generation and classification. In *2015 International Joint Conference on Neural Networks*
458 *(IJCNN)*(pp. 1-8). IEEE.
- 459 21. Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural
460 networks. *Science*, 313(5786), 504-507.
- 461 22. Jun, K., Lee, D. W., Lee, K., Lee, S., & Kim, M. S. (2020). Feature extraction using an RNN
462 autoencoder for skeleton-based abnormal gait recognition. *IEEE Access*, 8, 19196-19207.
- 463 23. Ma, J., & Yuan, Y. (2019). Dimension reduction of image deep feature using PCA. *Journal of*
464 *Visual Communication and Image Representation*, 63, 102578.
- 465 24. Dahouda, M.K. & Joe, I. (2022). Neural Architecture Search Net-Based Feature Extraction With
466 Modular Neural Network for Image Classification of Copper/ Cobalt Raw Minerals. *IEEE*
467 *Access*, 10, pp. 72253-72262. DOI: 10.1109/ACCESS.2022.3187420.
- 468 25. Petrovska, B., Zdravevski, E., Lameski, P., Corizzo, R., Štajduhar, I., & Lerga, J. (2020). Deep
469 learning for feature extraction in remote sensing: A case-study of aerial scene
470 classification. *Sensors*, 20(14), 3906.

- 471 26. Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., Manzagol, P. A., & Bottou, L. (2010). Stacked
472 denoising autoencoders: Learning useful representations in a deep network with a local denoising
473 criterion. *Journal of machine learning research*, 11(12).
- 474 27. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2017). Imagenet classification with deep
475 convolutional neural networks. *Communications of the ACM*, 60(6), 84-90.
- 476 28. Lee, H., Ekanadham, C., & Ng, A. (2007). Sparse deep belief net model for visual area V2.
477 *Advances in neural information processing systems*, 20.
- 478 29. Le, Q. V. (2013, May). Building high-level features using large scale unsupervised learning. In
479 *2013 IEEE international conference on acoustics, speech and signal processing*(pp. 8595-8598).
480 IEEE.
- 481 30. Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image
482 recognition. *arXiv preprint arXiv:1409.1556*.
- 483 31. Maas, A., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., & Potts, C. (2011, June). Learning word
484 vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for*
485 *computational linguistics: Human language technologies*(pp. 142-150).
- 486 32. Metsis, V., Androutsopoulos, I., & Paliouras, G. (2006, July). Spam filtering with naive bayes-which
487 naive bayes?. In *CEAS*(Vol. 17, pp. 28-69).
- 488 33. Eyheramendy, S., Lewis, D. D., & Madigan, D. (2003, January). On the naive bayes model for text
489 categorization. In *International workshop on artificial intelligence and statistics*(pp. 93-100). PMLR.
- 490 34. Gómez-Hidalgo, J. M., López, M. J. M., & Sanz, E. P. (2000). Combining text and heuristics for
491 cost-sensitive spam filtering. In *Fourth Conference on Computational Natural Language Learning*
492 *and the Second Learning Language in Logic Workshop*.
- 493