

2
3 **AE-FT: Auto Encoder Fixed-Target Training**
4 **Features Extraction Approach for Binary**
5 **Classification Problems**
6
7
8
9

10
11 **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. 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 Deep Belief Network (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.

12
13 *Keywords: Features Extraction; Deep Learning; AE; Fixed-Target Training; Common*
14 *Features.*

15
16
17 **1. INTRODUCTION**

18
19 With the rapid development of machine learning technology, as a binary classification
20 problem that helps people to find the law from the massive data to achieve the prediction
21 effectively, data prediction has become an important part of people's daily lives.

22 Feature extraction is about finding a good data representation, which is very domain-
23 specific, often requires human expertise, and is related to available measurements. The
24 primary idea behind feature extraction is to compress the data with the goal of maintaining
25 most of the relevant information. As to feature selection techniques, these techniques are

26 also used for reducing the number of features from the original feature set to reduce model
27 complexity, and model overfitting, enhance model computation efficiency, and reduce
28 generalization error. Therefore, it is better to better improve the accuracy of the learning
29 algorithm and shorten the training and output time.

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

37 Linear feature extraction assumes that the data lies in a linear subspace. Use matrix
38 factorization to protect them. On the other hand, in nonlinear feature extraction or
39 dimensionality reduction, a low-dimensional surface can be mapped into a high-dimensional
40 space so that a nonlinear relationship among the features can be found and easily detected.
41 Theoretically, a transformation function $f(x)$ can be used to map the features into a higher-
42 dimensional space and then mapped back into the lower-dimensional space, so that the
43 relationship can be viewed as nonlinear.

44 **1.1. Kernel PCA (KPCA)**

45 Kernel PCA introduced by Scholkopf et al. [2], is an extension of PCA that allows for
46 the separability of nonlinear data by making use of kernels. The basic idea behind it is to
47 project the linearly inseparable data onto a higher dimensional space where it becomes
48 linearly separable. Unfortunately, it has a serious limitation in terms of space complexity
49 since it stores all dot products of the training set and therefore the size of the matrix
50 increases quadratically with the number of data points as presented in [3].

51 Another drawback of the KPCA, however, is the cost of computation could be
52 extremely high, which could lead to the attendant numerical problems of diagonalizing large
53 matrices, which limits its applicability in many large dataset problems. But, an Expectation-
54 Maximization (EM) algorithm for KPCA to overcome these drawbacks was proposed in [4],
55 which is an expectation-maximization approach for performing kernel principal component
56 analysis. Experimental results showed that EM is an efficient method computationally,
57 especially for a large number of data points..

58 **1.2. Locally Linear Embedding (LLE)**

59 Locally Linear Embedding, proposed by Saul et al. [5], is a dimensionality reduction
60 technique based on Manifold Learning that involves the computation of low-dimensional
61 neighborhood preserving embeddings of inputs that are of high dimension in nature.
62 Manifold Learning aims to make a manifold object, an object of D dimensions that is
63 embedded in a higher-dimensional space, representable in its original D dimension instead
64 of being represented in an unnecessarily greater space.

65 LLE has the ability to learn the global structure of nonlinear manifolds like those from
66 images of faces or documents of text by exploiting the local symmetries of linear
67 reconstructions. LLE has been applied successfully in a wide range of applications which
68 includes face recognition and remote sensing, MRI, shape analysis of the hippocampus in
69 AD, diffusion tensor imaging, breast lesion segmentation, feature fusion, and image
70 classification according to [6].

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

76 Supervised and semi-supervised versions of LLE were proposed by [7] and [8],
77 respectively, for plant classification based on images of leaves.

78 **1.3. Linear Discriminant Analysis (LDA)**

79 LDA is a supervised learning dimensionality reduction, feature extraction technique,
80 and Machine Learning classifier that was invented by Fisher et al. [9]. LDA uses within-
81 classes and between-classes measures by maximizing the distance between the mean of
82 each class and minimizing the spreading within the class itself. This is a good choice
83 because maximizing the distance between the means of each class when projecting the data
84 in a lower-dimensional space can lead to better classification results.

85 An advantage of LDA is that it is able to use information from both features to create
86 a new axis which in turn minimizes the variance and maximizes the class distance of the
87 variables.

88 Although the LDA is one of the most well-used data reduction techniques, it has a
89 number of limitations. The small sample problem (SSS), is one of the main problems of LDA,
90 which happens when the dimensions are much higher than the number of samples in the
91 data matrix, LDA is unable to find the lower dimensional space resulting in the within-class
92 matrix becoming singular. Different approaches have been proposed to solve this problem,
93 such as what was proposed in [10] and [11]. In addition to the assumption that the input data
94 follows a Gaussian Distribution, therefore applying LDA to not Gaussian data can lead to
95 poor classification results.

96 A semisupervised variant of LDA, which performed better than the classical LDA,
97 was proposed by [12] that mainly combines both labeled and unlabeled data for training LDA
98 and allows using LDA for the situation where the labeled data are few.

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

102 **1.4. t-distributed Stochastic Neighbor Embedding (t-SNE)**

103 t-Stochastic Neighbor Embedding (t-SNE) is an unsupervised Non-linear Dimension
104 Reduction Technique (NLDRT) that was introduced by Maaten et al. [13]. The technique is a
105 variation of the Stochastic Neighbor Embedding introduced by Hinton et al.[8], whose main
106 objective is the construction of probability distributions from pairwise distances such that
107 larger distances correspond to smaller probabilities and vice versa. t-SNE is typically used to
108 visualize high-dimensional datasets, it works by minimizing the divergence between a
109 distribution constituted by the pairwise probability similarities of the input features in the
110 original high-dimensional space, which is modeled using a Gaussian Distribution and its
111 equivalent in the reduced low-dimensional space, modeled using a Student's t-distribution.

112 t-SNE makes use of the Kullback-Leiber (KL) divergence in order to measure the
113 dissimilarity of the two different distributions, as mentioned in [14]. The KL divergence is
114 then minimized using gradient descent.

115 t-SNE is the most commonly used in single-cell analysis. However, it has some
116 limitations as mentioned in [15]. The limitations include slow computation time, the inability to
117 meaningfully represent very large datasets, and the loss of large-scale information.

118 **1.5. Deep Learning Approach**

119 The major difference between deep learning and traditional pattern recognition
120 methods is that deep learning automatically learns features from big data, instead of
121 adopting handcrafted features, as stated in [16]. Deep learning is able to quickly acquire new
122 effective feature representations from training data.

123 Autoencoders

124 Autoencoders are feed-forward neural networks that were first introduced by
125 Rumelhart et al. [17]. They can learn a compressed and distributed representation of data,
126 which can be used as a dimensionality reduction or feature extraction technique. They use
127 nonlinear transformations to project data from a high dimension to a lower one. An
128 autoencoder usually has at least one hidden layer between the input and output layers. The
129 number of neurons in the hidden layer is usually set to less than those in the input and
130 output layers, thus creating a bottleneck, with the intention of forcing the network to learn a
131 higher-level representation of the input as presented in [18].

132 Autoencoders are typically trained in an unsupervised manner, using
133 backpropagation with stochastic gradient descent, to approximate a function by which data
134 can be classified, as mentioned in [19] and [20]. For every training input, the difference
135 between the input and the output is measured (using squared error) and it is back-
136 propagated through the neural network to perform weight updates on the different layers.

137 Compared with other machine learning methods, deep learning is able to detect
138 complicated interactions in features, learning lower-level features from nearly unprocessed
139 original determine characteristics that are not easy to be detected. Furthermore, they hand
140 class members with high cardinal numbers and process untapped data. Unfortunately, if all

141 input features are independent of each other, then the autoencoder will find it particularly
142 difficult to encode the input data into a lower-dimensional space.

143 The advantages are higher discriminating power and control overfitting when it is
144 unsupervised, and the disadvantages are loss of data interpretability and transformation may
145 be expensive.

146 There are many other deep learning-based feature extraction approaches. A
147 skeleton-based abnormal gait recognition approach was proposed by Jun et al. [21]. They
148 proposed a feature extraction method using the RNN AEs to minimize the irrelevant
149 information of the original skeleton data. They used two-step training of a hybrid RNN and
150 AE-DM model and approved that it is more effective than the single-step training of the End-
151 to-End model that has the same data flow. Ma and Yuan [22] proposed a method for
152 extracting features from images based on deep CNN and PCA. They used a neural network
153 to extract features and a PCA algorithm for feature dimension reduction. Then they
154 compared the performance of the PCA before and after the improvement claim achieving
155 memory, and time optimization. Moreover, the SVM classifier accuracy was enhanced.
156 However, these are domain-dependent approaches that work only for the specific domain
157 not for binary classification problems in general.

158

159

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

160 The main issue with machine learning-based feature extraction techniques is the
161 requirement of extensive domain-level knowledge and experience, see Table 1. Moreover, it
162 is often difficult to apply domain-level knowledge to extract the necessary features for
163 building a machine-learning classifier. Therefore, it is significantly important to find and
164 develop feature extraction techniques that depend mainly on the training data and don't
165 require or depend on domain-level knowledge and experience, and this is our main purpose.

166 This paper presents the use of denoising stacked autoencoders with supervised
167 fixed-target training in order to extract the common features of the training data that can be
168 used in binary classification. The proposed approach achieves 73.50% accuracy, which is so
169 far a good result that does not need any domain expertise.

170 In the experiments, we used the extracted common features to build binary
171 classifiers using seven binary classification methods Naive Bayes, Logistic Regression, K-
172 Nearest Neighbors, Support Vector Machine, Decision Tree, Random Forest, and Voting
173 Classification- for binary classification.

174 The proposed method, fixed-target training of a deep stacked autoencoder, enabled
175 a good recognition accuracy, better generalization, and more stability than that which could
176 be achieved with the other methods. In contrast to most existing approaches that normally
177 have a separate stage for data preprocessing, followed by domain-dependent feature
178 extraction. We developed a domain-independent deep neural network framework for
179 common feature extraction which enables us to easily, without the need for domain-level
180 knowledge or expertise, extract features that can be effectively used in binary classification
181 problems.

182

183 2. METHODOLOGY

184

185 In this section, we describe our proposed novel deep learning-based approach for common
186 feature extraction in detail. The main question we are trying to answer is the following:

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

189 In recent years, deep learning methods have proven very successful in
190 accomplishing dimensionality reduction and feature extraction tasks in many domains,
191 especially computer vision and cybersecurity according to [24, 25, 26, 27]. The proposed
192 methodology works as follows:

193 Firstly, training the proposed common features extraction model to extract features
194 to use in binary classification using various common machine learning classification
195 algorithms.

196 Secondly, select seven binary classification algorithms and use the extracted
197 common features to build binary classifiers using common algorithms.

198 Our method uses stacked denoising autoencoders for extracting the common
199 features of the training dataset that can be used effectively in binary classification. The input
200 to the Deep Belief Network (DBN) has a fixed length, but the dataset sample length is
201 variable. In order to represent the dataset binary as a fixed-sized vector, which would be the
202 input to the neural network, we repeatedly pad the sample binary until the specified size is
203 met. This process gives better results than 0's or 1's padding that the DBN learns as a
204 common feature. Then we use two different datasets, the IMDB dataset [28] and the Enron-
205 Spam dataset [29], to train and test the proposed feature extractor for extracting the
206 common features and then used for training binary classification models using the extracted
207 features of the binary representation of the datasets. We focus on the top seven most
208 common binary classification algorithms Naive Bayes, Logistic Regression, K-Nearest
209 Neighbours, Support Vector Machine, Decision Tree, Random Forest, and Voting
210 Classification.

211 2.1. Deep Belief Network Fixed-Target Training

212 In stacked denoising autoencoders, first introduced by Vincent et al. [23], the data at
 213 the input layer is replaced by noised data while the data at the output layer stays the same;
 214 therefore, the autoencoder can be trained with much more generalization power, according
 215 to [16].

216 Usually, denoising autoencoder training is unsupervised, the input sample is
 217 corrupted by adding noise (or more often by zeroing the values). That is, given an input \mathbf{S} ,
 218 first it is corrupted to $\tilde{\mathbf{S}}$ and then fed to the input layer of the network. The objective function
 219 of the network in the output layer remains to generate the uncorrupted version of the input
 220 (see Figure 1-(a)). But in our approach, we use a novel supervised training strategy, which
 221 we call a *fixed-target training* strategy.

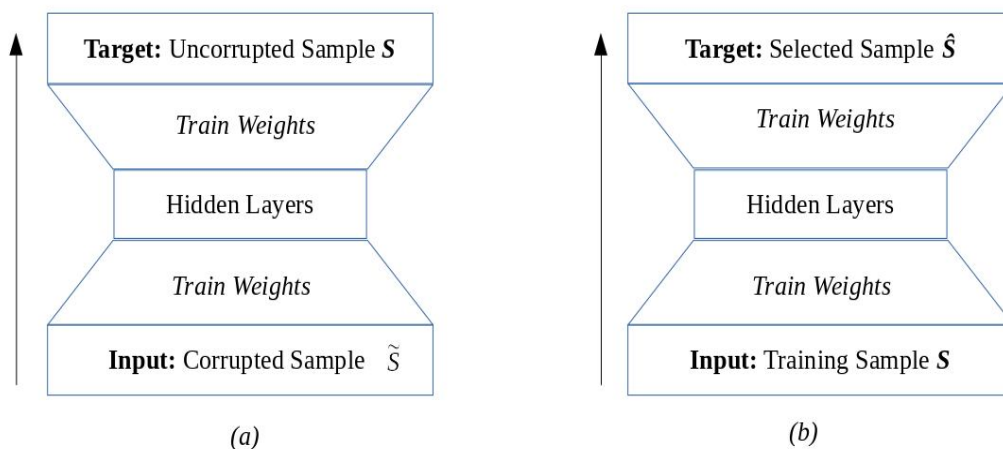


Fig. 1. Denoising autoencoder training. (a)-traditional unsupervised training, (b)-fixed-target supervised training.

222 In the *fixed-target training* strategy, we randomly select one of the training samples
 223 \mathbf{S}^{\wedge} , fix it in the output layer, and for every input $\tilde{\mathbf{S}}$ of the training samples, the objective
 224 function of the network is to generate \mathbf{S}^{\wedge} , (i.e. we consider all training samples as corrupted
 225 versions of the selected sample \mathbf{S}^{\wedge} (see Figure 1-(b)). So, the "hidden units" of DBN
 226 compute internal representations analogous to the extracted common features.

227 This training approach works better than traditional training. By fixing the target
 228 output, the network is forced to generalize better and determine more high-level common
 229 patterns. Moreover, the network learns better even when few training samples are available.

230 When a DBN's training is complete, we can discard the decoder layer, fix the values
 231 of the encoder layer, and use the encoder as the common features extractor. In a typical
 232 implementation, the extracted features may then be used for supervised binary classification.

233 In order to achieve our goal, we create a deep belief network (DBN) by training a
 234 deep stack of denoising autoencoders. We use fixed-target training to train a deep denoising
 235 autoencoder consisting of five layers: 8,192–2,048–512–128–32. At the end of this training
 236 phase, we have a deep network that is capable of converting 8,192 length input vectors into
 237 32 floating point feature values. Note that the network is trained only using the samples in
 238 the training set, and for all future samples it will be run in prediction mode;(i.e., receiving the

239 8,192-sized vector it will produce 32 output values, without modifying the weights). See
240 Figure 2.

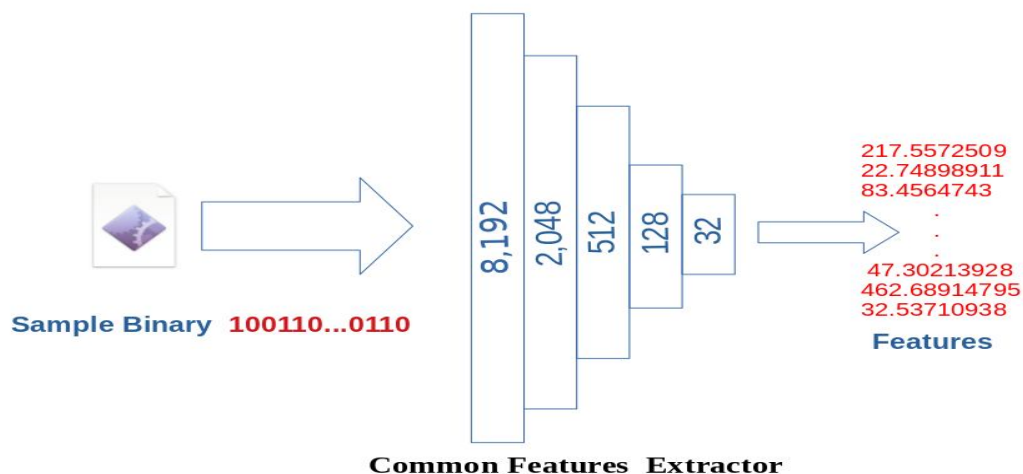


Fig. 2. Illustration of the common features extraction stages from feeding the sample binary to features extraction using DBN.

241 The sample's binary bit string is fed to the neural network, and the deep neural
242 network generates a 32-sized vector at its output layer, which we treat as the common
243 feature values of the sample.

244 The next section provides the implementation details and experimental results and
245 demonstrates that the resulting 32-sized vector (i.e., the extracted common features) indeed
246 provides a good representation of the data common feature.

247 2.2. Datasets

248 We are using two different datasets, the Internet Movie Database (IMDB) dataset
249 [28] and the Enron-Spam dataset [29], for testing our approach. The IMDB dataset, Large
250 Movie Review Dataset v1.0, of highly polar movie reviews in the form of text comments on
251 different movies and a positive or negative score. This dataset contains movie reviews along
252 with their associated binary sentiment polarity labels. It is intended to serve as a benchmark
253 for sentiment classification. The core dataset contains 50,000 reviews split evenly into 25k
254 train and 25k test sets. The overall distribution of labels is balanced (25k pos and 25k neg).
255 Furthermore, the train and test sets contain a disjoint set of movies, so no significant
256 performance is obtained by memorizing the movie's unique terms and they are associated
257 with the observed labels. In the labeled train/test sets, a negative review has a score < 4 out
258 of 10, and a positive review has a score > 7 out of 10. Thus, reviews with more neutral
259 ratings are not included in the train/test sets.

260 These datasets have been used as a benchmark in several research papers on
261 spam filtering, text classification, and natural language processing [30, 31]. Therefore, the
262 results of this work are hopefully comparable to other similar works within the area, without
263 having to account for unique datasets.

264 The goal is to train the proposed feature extractor for extracting the common
265 features and then used it for training binary classification models using the extracted
266 features of the binary representation of the datasets. We focus on the top seven most
267 common binary classification algorithms Naive Bayes, Logistic Regression, K-Nearest
268 Neighbours, Support Vector Machine, Decision Tree, Random Forest, and Voting
269 Classification.

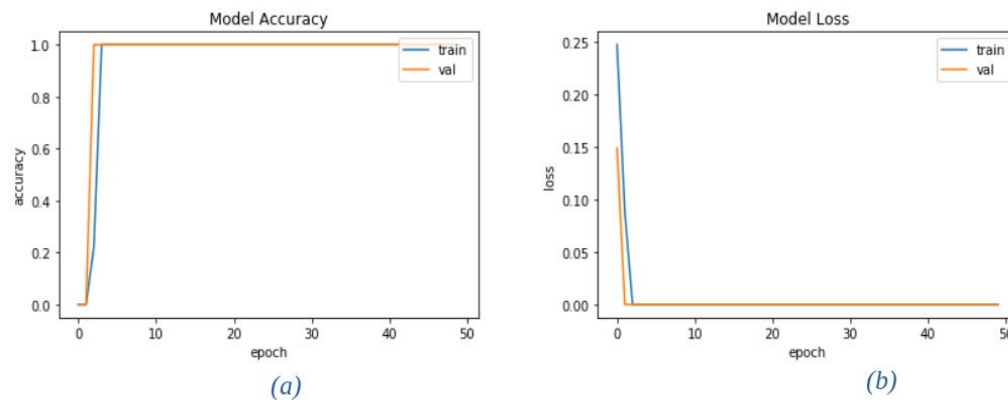
270 2.3. Common Features Extraction Method

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

276 *Fixed-target training* is essential for the network to learn common features in a very
277 small time and even when few samples are available. We use rectified linear units (ReLU)
278 for the nonlinearity function when training deep neural networks to diminish the gradient
279 vanishing problem and result in faster convergence, as approved in [12, 18].

280 We build the model using machine learning algorithms in *the Keras library*. Other
281 parameters we use are 50 training epochs, a learning rate of 0.003, a batch size of 10, and
282 no momentum.

283 Figure 3 shows the learning curves, which approve the feasibility of our proposed
284 common feature extraction approach.



285 **Fig. 3. Common Features Extraction Model training and validation. (a)-Accuracy, (b)-Loss**

286

286 3. EXPERIMENTAL RESULTS

287 In this section, we conducted various experiments to evaluate the features of the
288 DBN AEs for binary classification. First, we compared the number of features used in [28,

289 29] and the extracted common features. Then, we conducted experiments to show the
 290 effectiveness of the proposed DBN AEs. Finally, we showed the performance of the features
 291 of the DBN AEs used in binary classification problems, we used seven binary classification
 292 methods Naive Bayes, Logistic Regression, K-Nearest Neighbours, Support Vector Machine,
 293 Decision Tree, Random Forest, and Voting Classification- for binary classification.

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

299 In the testing phase, we first extract the values of the common features of the
 300 sample and pass them to the binary classifier for classification. See Figure 4.

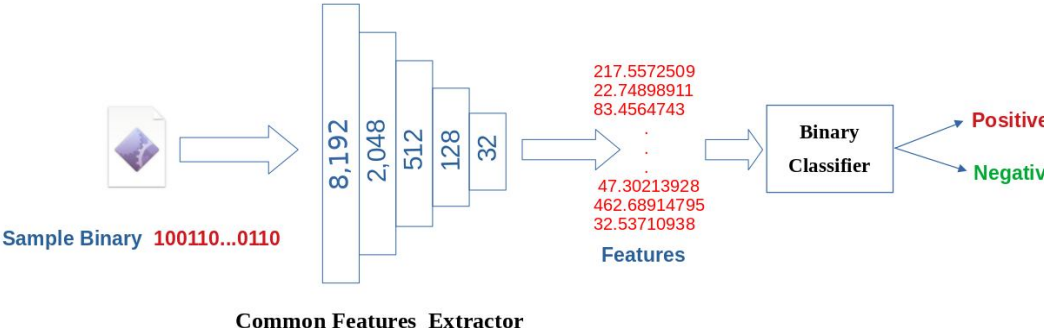


Fig. 4. Overview of our deep learning approach for common features extraction, illustration of all stages from feeding the sample binary to the Common Features Extractor ending with the decision made by the binary classifier.

301
 302 The results demonstrate that, with proper structure and parameters, the
 303 performance of the proposed deep learning method on common feature extraction is useful
 304 even in the lack of domain expertise.

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

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%
<u>Voting Classification</u>	<u>85.0%</u>	<u>72.7%</u>	<u>84.7%</u>	<u>73%</u>

307 The accuracy results are given in Table 2, test samples are 73.5% correctly
308 classified by conducting a Logistic Regression classifier for the proposed approach with the
309 extracted common features of the IMDB dataset. For the Enron-Spam dataset, the best
310 results are obtained by Random Forest 73.5%. The binary classification performances of the
311 approaches are given in Table 2. According to the results, the common features extracted by
312 the proposed approach give a very good performance for all seven binary classification
313 methods. Logistic Regression and Random Forest are the most effective classifiers for this
314 approach.

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

318 4. DISCUSSION AND EVALUATION

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

324 According to the results reported in Table 2, the best accuracy values are obtained
325 by conducting Logistic Regression and Random Forest classifiers. When comparing the
326 number of features used, 32 common features were extracted using the proposed feature
327 extraction algorithm, on the other hand, 3000 and 10000 features were used in [28] and
328 [29], respectively, in addition to the preprocessing carried to extract the features used in [28]
329 and [29]. However, by applying the proposed common features extraction algorithm, it is
330 observed that high success rates are achieved with very few features and increase the
331 overall process performance.

332 5. FUTURE WORK

333 Future work could include examining how to improve the accuracy of the proposed
334 common feature extraction method by finding the optimal autoencoder structure and
335 activation function. Furthermore, examining different ways for the pivot sample selection. Of
336 course, we could also look at different types of classifiers and use different datasets in
337 different domains. One could also see is it useful to use the proposed method in data
338 collection; when only a few samples are available.

339 Finally, it would be interesting to see if the proposed method could mitigate the
340 problem of the need for domain expertise and data preprocessing for feature extraction and
341 the need for a big dataset to train binary classifiers.

342 6. CONCLUSION

343 In binary classification, feature extraction is a basic and important matter because
344 the original data contain noise and irrelevant information which decreases the classification
345 accuracy. to reduce the dimensionality of the feature vector space.

346 In this paper, we propose a feature extraction method using the DBN. The proposed
347 approach, AE fixed-target supervised training, extracts the common characteristics of the
348 original data and minimizes the irrelevant information without the need for any domain-level
349 expert or expensive data preprocessing. Furthermore, it needs less training data and
350 produces fewer features compared to other feature extraction techniques. Therefore, the
351 extracted features improve the binary classification performance.

352 The most interesting result of our evaluation was the very good performance of our
353 common features extension approach using all binary classification algorithms tested on
354 both datasets that have been used, but further experiments are needed to be more
355 confident.

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

360

361 ACKNOWLEDGEMENTS

362

363 Funding: The study attracted no funding.

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

365

366 REFERENCES

- 367 1. Ghojogh, Benjamin, Maria N. Samad, Sayema Asif Mashhadi, Tania Kapoor, Wahab Ali, Fakhri
368 Karray, and Mark Crowley. Feature selection and feature extraction in pattern analysis: A literature
369 review. *arXiv preprint arXiv:1905.02845* (2019).
- 370 2. Schölkopf, Bernhard, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis
371 as a kernel eigenvalue problem. *Neural computation* 10, no. 5 (1998): 1299-1319.
- 372 3. Liu, Xiaofang, and Chun Yang. Greedy kernel PCA for training data reduction and nonlinear
373 feature extraction in classification. In *MIPPR 2009: Automatic Target Recognition and Image*
374 *Analysis*, vol. 7495, pp. 768-775. SPIE, 2009.
- 375 4. Rosipal, Roman, and Mark Girolami. An expectation-maximization approach to nonlinear
376 component analysis. *Neural Computation* 13, no. 3 (2001): 505-510.
- 377 5. Saul, Lawrence K., and Sam T. Roweis. Think globally, fit locally: unsupervised learning of low
378 dimensional manifolds. *Journal of machine learning research* 4, no. Jun (2003): 119-155.
- 379 6. Nanga, Salifu, Ahmed Tijani Bawah, Benjamin Ansah Acquaye, Mac-Issaka Billa, Francis Delali
380 Baeta, Nii Afotey Odai, Samuel Kwaku Obeng, and Ampem Darko Nsiah. Review of Dimension
381 Reduction Methods. *Journal of Data Analysis and Information Processing* 9, no. 3 (2021): 189-
382 231.

- 383 7. Kouropteva, Olga, Oleg Okun, and Matti Pietikäinen. Incremental locally linear embedding
384 algorithm. In *Scandinavian Conference on Image Analysis*, pp. 521-530. Springer, Berlin,
385 Heidelberg, 2005.
- 386 8. Pan, Yaozhang, Shuzhi Sam Ge, and Abdullah Al Mamun. Weighted locally linear embedding for
387 dimension reduction. *Pattern Recognition* 42, no. 5 (2009): 798-811.
- 388 9. Xanthopoulos, Petros, Panos M. Pardalos, and Theodore B. Trafalis. Linear discriminant analysis.
389 In *Robust data mining*, pp. 27-33. Springer, New York, NY, 2013.
- 390 10. Sharma, Alok, and Kuldip K. Paliwal. A new perspective to null linear discriminant analysis method
391 and its fast implementation using random matrix multiplication with scatter matrices. *Pattern*
392 *Recognition* 45, no. 6 (2012): 2205-2213.
- 393 11. Belhumeur, Peter N., Joao P. Hespanha, and David J. Kriegman. Eigenfaces vs. fisherfaces:
394 Recognition using class specific linear projection. *IEEE Transactions on pattern analysis and*
395 *machine intelligence* 19, no. 7 (1997): 711-720.
- 396 12. Zhang, Yu, and Dit-Yan Yeung. Semisupervised generalized discriminant analysis. *IEEE*
397 *Transactions on Neural Networks* 22, no. 8 (2011): 1207-1217.
- 398 13. Van der Maaten, Laurens, and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine*
399 *learning research* 9, no. 11 (2008).
- 400 14. Hinton, Geoffrey E., and Sam Roweis. Stochastic neighbor embedding. *Advances in neural*
401 *information processing systems* 15 (2002).
- 402 15. Xie, Bo, Yang Mu, Dacheng Tao, and Kaiqi Huang. m-SNE: Multiview stochastic neighbor
403 embedding. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 41, no. 4
404 (2011): 1088-1096.
- 405 16. Wang, Haohan, and Bhiksha Raj. On the origin of deep learning. *arXiv preprint*
406 *arXiv:1702.07800* (2017).
- 407 17. Rumelhart, David E., Geoffrey E. Hinton, and Ronald J. Williams. *Learning internal*
408 *representations by error propagation*. California Univ San Diego La Jolla Inst for Cognitive
409 Science, 1985.
- 410 18. David, Omid E., and Nathan S. Netanyahu. Deepsign: Deep learning for automatic malware
411 signature generation and classification. In *2015 International Joint Conference on Neural Networks*
412 *(IJCNN)*, pp. 1-8. IEEE, 2015.
- 413 19. Werbos. *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Science*,
414 *Appl. Math.* Harvard University (1974).
- 415 20. Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. Reducing the dimensionality of data with neural
416 networks. *science* 313, no. 5786 (2006): 504-507.
- 417 21. Jun, Kooksung, Deok-Won Lee, Kyoobin Lee, Sanghyub Lee, and Mun Sang Kim. Feature
418 extraction using an RNN autoencoder for skeleton-based abnormal gait recognition. *IEEE*
419 *Access* 8 (2020): 19196-19207.
- 420 22. Ma, Ji, and Yuyu Yuan. Dimension reduction of image deep feature using PCA. *Journal of Visual*
421 *Communication and Image Representation* 63 (2019): 102578.
- 422 23. Vincent, Pascal, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and
423 Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network
424 with a local denoising criterion. *Journal of machine learning research* 11, no. 12 (2010).
- 425 24. Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep
426 convolutional neural networks. *Communications of the ACM* 60, no. 6 (2017): 84-90.
- 427 25. Lee, Honglak, Chaitanya Ekanadham, and Andrew Ng. Sparse deep belief net model for visual
428 area V2. *Advances in neural information processing systems* 20 (2007).
- 429 26. Le, Quoc V. Building high-level features using large scale unsupervised learning. In *2013 IEEE*
430 *international conference on acoustics, speech and signal processing*, pp. 8595-8598. IEEE, 2013.
- 431 27. Simonyan, Karen, and Andrew Zisserman. Very deep convolutional networks for large-scale
432 image recognition. *arXiv preprint arXiv:1409.1556* (2014).

- 433 28. Maas, Andrew, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher
434 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of*
435 *the association for computational linguistics: Human language technologies*, pp. 142-150. 2011.
- 436 29. Metsis, Vangelis, Ion Androutsopoulos, and Georgios Paliouras. Spam filtering with naive bayes-
437 which naive bayes?. In *CEAS*, vol. 17, pp. 28-69. 2006.
- 438 30. Eyheramendy, Susana, David D. Lewis, and David Madigan. On the naive bayes model for text
439 categorization. In *International workshop on artificial intelligence and statistics*, pp. 93-100. PMLR,
440 2003.
- 441 31. Gómez-Hidalgo, José María, Manuel J. Maña López, and Enrique Puertas Sanz. Combining text
442 and heuristics for cost-sensitive spam filtering. In *Fourth Conference on Computational Natural*
443 *Language Learning and the Second Learning Language in Logic Workshop*. 2000.

UNDER PEER REVIEW