

Yang, M., Lim, M. K., Qu, Y., Li, X. and Ni, D. (2022) Repair missing data to improve corporate credit risk prediction accuracy with multi-layer perceptron. *Soft Computing*, 26(18), pp. 9167-9178. (doi: <u>10.1007/s00500-022-07277-4</u>)

The material cannot be used for any other purpose without further permission of the publisher and is for private use only.

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

https://eprints.gla.ac.uk/270912/

Deposited on 12 March 2022

Enlighten – Research publications by members of the University of Glasgow <u>http://eprints.gla.ac.uk</u>

1 Repair missing data to improve corporate credit risk prediction

2 accuracy with multi-layer perceptron

3

4 Abstract

5 Data loss has become an inevitable phenomenon in corporate credit risk (CCR) prediction. To ensure the integrity of data information for subsequent analysis and 6 7 prediction, it is essential to repair the missing data as accurately as possible. To solve the problem of missing data in credit classification, this study proposes a multi-layer 8 perceptron ensemble (MLP-ESM) model that can perform data interpolation and 9 prediction simultaneously to predict CCR. The model makes full use of non-missing 10 information and interpolates more missing columns with fewer missing values. In this 11 12 way, not only the data features needed for missing data interpolation are extracted, but also the structural relationship features between the predicted target and the existing 13 14 data are extracted, which can achieve the effect of simultaneous interpolation and prediction. The results show that the MLP-ESM model can effectively interpolate and 15 predict the missing dataset of CCR. The prediction accuracy is 83.11%, which is better 16 17 than the traditional machine learning model. This fully shows that the dataset after interpolation can achieve a better prediction effect. 18

19 Keywords

20 Missing data, data interpolation, machine learning, Corporate credit risk, prediction

1 1 Introduction

2 In the credit environment of a developed market economy, credit risk prediction is very important for the development of companies Sivasankar et al. (2016); (Wang et 3 al., 2011; Yuan et al., 2022). Among them, a company's credit rating, especially under 4 the 2014 New Basel Accord, plays a decisive role in the capital market and can be used 5 as an important parameter in the evaluation of the corporate credit risk (CCR) (Yu et 6 al., 2022). However, in the prediction of the CCR, some data will inevitably be lost due 7 to the restructuring of the database, incomplete company records, and the replacement 8 of characteristic indexes (Florez-Lopez, 2010; Moscato et al., 2021; Soldatyuk et al., 9 10 2014). In the CCR prediction, missing data will have a great impact on the prediction results (Garcia et al., 2019; Lan et al., 2021; Twala, 2013). As early as 2006, Mählmann 11 (2006) emphasized the importance of processing missing data for the CCR prediction 12 in their CCR prediction. 13

14 For the missing data in CCR, even a small degree of data loss can lead to a reduction in the number of effective samples, deviation in the analysis results, and the 15 weakening of the credibility of prediction results (Budagaga, 2020; Garcia et al., 2019). 16 If the degree of data loss is serious, it is likely that the original CCR prediction cannot 17 be effectively predicted due to the insufficient sample size (Li et al., 2020; Zhang et al., 18 2016). But for the missing data processing method, relatively simple is the direct 19 deletion method. This method is simple to operate, but it is easy to produce estimation 20 deviation. And simple deletion can lose a lot of information and affect the accuracy of 21 prediction results (Hooke et al., 2021; Nakagawa et al., 2008; Qiu et al., 2008). 22 Therefore, researchers usually do not use such a method when preprocessing missing 23 24 data, except in cases where the missing proportion is very low. However, for the data with low missing degree and great predictive value, data repair through data 25 interpolation is more easily accepted by researchers (Akin Arikan et al., 2018; Catellier 26 27 et al., 2005; De Silva et al., 2016).

The idea of data repair is to use existing data to mine data with predictive value or 28 analytical value from existing samples by advanced technology (Chiang et al., 2016; 29 30 Gu et al., 2017; Wei et al., 2019). Many studies used generative adversarial networks (GANs) and autoencoders to repair the optimal parameters for missing data (Campanile 31 et al., 2020; Chang et al., 2020b; D'Angelo et al., 2021). Chang et al. (2020b) proposed 32 the dual-domain conditional GANs to interpolate seismic data to drive the generative 33 network to learn optimal parameters. Meanwhile, compared with traditional methods 34 35 that only identify outliers, Eduardo et al. (2020) proposed a robust variational autoencoder method that can detect and repair abnormal units in the dataset. In addition, 36 37 the idea of data repair has also been found in practical applications. For example, in the text classification task of natural language processing, Wei and Zou (2019) proposed 38

synonym replacement, random insertion, random exchange, and random deletion to 1 repair data. In addition, in the study of image recognition, Shorten et al. (2019) 2 proposed horizontal flipping, clipping, rotation, translation, adding interference, 3 convolution kernel filter, image mixing, and random deletion, and achieved good data 4 repair effects. Therefore, for the existing CCR prediction dataset, the dataset can 5 contain extremely valuable duration and frequency information. If the data with 6 prediction potential can be mined from the existing samples, the data can be 7 interpolated. It will be possible to improve the accuracy of the CCR forecasts through 8 9 the repair of internal data (Florez-Lopez, 2010; Li et al., 2020).

In the actual prediction of the CCR, scholars mainly discuss manual interpolation, 10 special value interpolation and mean means interpolation based on statistics (Florez-11 Lopez, 2010; Nijman et al., 2021). These methods mainly interpolate the data near the 12 missing dataset so as to realize missing data repair (Chang et al., 2020a; Guo et al., 13 2009; Mählmann, 2006). In recent years, with the rise of artificial intelligence, machine 14 learning algorithms are popular because of their significant advantages in processing 15 complex data and improving model performance (D'Angelo et al., 2021; D'Angelo et 16 al., 2019; Sarker, 2021). These algorithm have also been applied to missing data repair 17 and proved to be better than statistical interpolation methods (Jerez et al., 2010; Moon 18 et al., 2019). For example, Moon et al. (2019) used multi-layer perceptron (MLP) to 19 accurately repair the missing temperature data, and improved data interpolation 20 accuracy. Although machine learning algorithms have been proved to be capable of 21 missing data repair and prediction, some algorithms only perform interpolation in the 22 23 same row or column of data, and the data rows or columns that can be interpolated are also severely limited by data holding amount (Chang et al., 2020a). However, in the 24 25 research field of CCR prediction, due to the limitations of data cost and privacy protection, it is increasingly difficult to find more predictive indicators from the vertical 26 dimension (Shema et al., 2019; Yap et al., 2011; Yue et al., 2016). In addition, the 27 existing interpolation methods of integration models can be limited by data samples, 28 29 which means that the problem of data interpolation under the condition of limited data samples needs to be solved urgently. 30

In order to solve the above problems, this study focuses on the field of CCR 31 prediction. By mining data with predictive potential from existing samples for 32 interpolation, it aims to build a method that can take into account both missing data 33 interpolation and prediction for missing datasets, thus further improve the CCR 34 prediction accuracy. The rest of the study is structured as follows. Section 2 presents 35 the method of this study, Section 3 describes in detail the dataset and data preprocessing 36 37 steps used in this study, and Section 4 analyzes the results of this study method. Finally, the conclusion of this study and the way of future research are presented in Section 5. 38

1 2 Methods

For the topic of missing data in this study, this study intends to propose a model
that can perform data interpolation and prediction simultaneously. The specific process
is as follows.

5 2.1 Theoretical design of the model

In this study, the MLP model with multi-output and strong nonlinear fitting ability
is introduced into the missing data interpolation prediction model (Belue et al., 1997;
Gao et al., 2012; Khoygani et al., 2016). Multi-layer perceptron and interpolation
model-specific process is as follows:

10 (1) The structure of the MLP

MLP, also known as a feedforward neural network, is extended from the perceptron model proposed by Rosenblatt. The perceptron model consists of two layers of neurons. The input layer receives the input signal and transmits it to the output layer to obtain the output result. Figure 1 shows a simple, functional neuron structure model. Each neuron accumulates one or more weighted input values to get the cumulative value, then performs a nonlinear transformation on it using activation function, and then transfers the value to the lower neuron.

18





19

20 21

Figure 1 Functional neuron structure



As the perceptron model has only one layer of functional neurons, its learning ability is very limited, and it can only complete simple linear separable problems. In order to enhance the fitting ability of the model and the ability to deal with nonlinear situations, the multi-layer perceptron model can be obtained by adding hidden layers to the perceptron model. It is mainly characterized by multiple layers, which are generally composed of one input layer, several hidden layers, and one output layer. There is a fully connected structure between the layers. Its model diagram is shown as follows:

Suppose that there is a multi-layer perceptron L layer, in which the first layer is the input layer. $\{x_1, x_2, \dots, x_n\}$ as the input data, and the L layer to output layer. 1 { y_1, y_2, \dots, y_m } for the output data. For the 1 hidden layer, suppose n_l neuron. 2 [$h_1^{(l)}, h_2^{(l)}, h_3^{(l)}, \dots, h_{n_l}^{(l)}$ } as the output data of the 1 layer. Set $w_{ij}^{(l)}$ for the jth 1 - 1 layer 3 of neurons to the ith 1 layer neuron weights. $b_i^{(l)}$ is the bias of the ith neuron in layer 4 L. $f(\cdot)$ is the activation function. Common activation function with sigmoid function 5 tanh function, ReLU function, etc. The 1 layer of the ith neurons $o_i^{(l)}$ and an output layer 6 of the ith neurons y_i output formula respectively as follows:

$$h_i^{(l-1)} = f(o_i^{(l-1)})$$
(1)

$$o_i^{(l)} = \sum_{j=1}^{s_{l-1}} w_{ij}^{(l)} h_j^{(l-1)} + b_i^{(l)}$$
(2)

$$y_i = f(\sum_{j=1}^{s_{L-1}} w_{ij}^{(L)} h_j^{(L-1)} + b_i^{(L)})$$
(3)

Subsequently, the weight and threshold of each neuron at each layer can be learned
by using the error inverse propagation model so as to make the output of the neural
network as close to the real value as possible. The training data set is expressed as
{(x₁, y₁), (x₂, y₃),, (x_n, y_n), y_k ∈ [-1,1]}, (x_k, y_k) is for the one training sample,
multi-layer perceptron output for ŷ_i, error function of the data set is defined as follows:

$$E = \frac{1}{n} \sum_{k=1}^{n} (y_k - \hat{y}_k)^2$$
(4)

12

13 The weight W and bias B in the multi-layer perceptron can be updated iteratively 14 according to the following formula, where α is the learning rate and its value range is 15 (0,1).

$$W^{(l)} = W^{(l)} - \alpha \frac{\partial E}{\partial W^{(l)}}$$
(5)

$$B^{(l)} = B^{(l)} - \alpha \frac{\partial E}{\partial B^{(l)}} \tag{6}$$

16

17 MLP can be trained to complete some complex tasks by modifying network 18 weights and biases.

19 (2) Design of interpolation process

For input dataset containing missing data, according to the data of this research suggests first arranged to organize data, and the columns of the missing parts for standardization, get missing values are arranged trapezoidal datasets X = $(x_1, x_2, ..., x_t, x_{t+1}, ..., x_{t+k}), x_1, x_2, ..., x_t$ is no missing data in the column. For the

- 1 existence of the missing data column $x_{t+1}, x_{t+2}, \dots, x_{t+k}, k$ as the number of columns,
- 2 which explain the index number of variables.
 - For any $i \leq j$, number of x_i and x_j missing values is $n_i \leq n_j$,.
- 4 The dataset X is divided into k + 1 sub-datasets, respectively, $X_1 =$ 5 $(x_1, x_2, ..., x_t), X_2 = x_{t+1}, ..., X_{k+1} = x_{t+k}$
- For the input X_1 , a full connection layer F_1 with 8 neurons was constructed, and a full connection layer G_1 with a single neuron was constructed with the output of the full connection layer F_1 . $G_1(F_1(X_1))$ as the input to interpolate the sub-dataset X_2 . And construct the loss function of interpolation

$$L_1 = \frac{1}{n_2} \sum_{i=1}^m g^*(x_{2i}) \left(G_1(F_1(x_{1i})) - g(x_{2i}) \right)^2$$
(7)1

11

3

12 n_2 is the number of missing values in X_2 , $G_1(F_1(x_{1i}))$ is the value of the dataset 13 X_2 based on X_1 interpolation, and

14 $g(x) = \begin{cases} x, & \text{if } x \text{ is numeral full} \\ 0, & \text{if } x \text{ is umerals empty'} \end{cases}$

15
$$g^*(x) = \begin{cases} 1, & \text{if } x \text{ is numeral full} \\ 0, & \text{if } x \text{ is umerals empty} \end{cases}$$

16 For input X_2 , we construct the interpolation layer P_2 , then

17

$$P_2(X_2) = G_1(F_1(X_1)) * (1 - g^*(X_2)) + X_2 * g^*(X_2)$$
(8)2

18

After nonlinear mapping F_2 , vectors and splicing $F_1(X_1)$ are obtained, denoted as $F_1(X_1) \oplus F_2(P_2(X_2))$, and a fully connected layer G_2 with single neurons is constructed $F_1(X_1) \oplus F_2(P_2(X_2))$ as input so as to $G_2(F_1(X_1) \oplus F_2(P_2(X_2)))$ interpolate sub-datasets X_3 and construct a similar loss function L_2 .

In this way, the interpolation loss function $L_1, L_2, ..., L_k$ constructed $X_2, X_3, ..., X_{k+1}$ by pyramid interpolation can be obtained, at the same time to $F_1(X_1) \oplus F_2(P_2(X_2)) \oplus ... \oplus F_{k+1}(P_{k+1}(X_{k+1}))$ and construct the full connection layer H_1 with 16 hidden layer neurons as the input and construct the full connection layer H_2 with the H_1 output as the input. With the output H_2 as the predicted value, the prediction model F is as follows:

$$F = H_2 \circ H_1 \circ (F_1 \bigoplus F_2 \circ P_2 \bigoplus \dots \bigoplus F_{k+1} \circ P_{k+1})$$
(9)3

1 For the input vector of independent variables x, F(x) is the prediction of the 2 probability of corporate credit downgrade.

With cross-entropy as the loss function of the model, denoted as L_{pred} , then

3 4

 $L_{pred} = \frac{1}{n} \sum_{i=1}^{n} - \left[y^{i} \log F(x_{1}^{i}, x_{2}^{i}, \dots, x_{t}^{i}) + (1 - y^{i}) \log \left(1 - F(x_{1}^{i}, x_{2}^{i}, \dots, x_{t}^{i}) \right) \right]$ (10)4

5 6

7

Then the loss function of this model is:

 $L = L_{pred} + \sum_{i=1}^{k} \lambda_k L_k \tag{11}5$

8

9 λ_k is the weight of interpolation loss of each variable. The above formula indicates 10 that the total loss of the model is the loss predicted L_{pred} by the model to the predicted 11 target y in the sample, plus the loss L_k predicted by the missing independent variable 12 data in the sample multiplied by the penalty weight. Therefore, the overall optimization 13 objective of trapezoidal multi-layer perceptron is:

14

$$z = \min\left(L_{pred} + \sum_{i=1}^{k} \lambda_k L_k\right)$$
(12)6

In the subsequent experiments of this study, as it is difficult for general software packages to directly process the missing values, a large number is used to mark the missing values (after standardization, most of the values in the dataset are between -3 and 3. In this study, -10 is used to mark the missing values).

19 The interpolation process is shown as follows:





Figure 3 The interpolation process of the MLP

1 2.2 Model interpolation idea

The theoretical design ideas of the model proposed in this study are shown in Figure 4, which can be divided into four steps: trapezoidal arrangement of data, model building, model-based training, and model-based interpolation.

Table	1										
Line	Complete data							Missing	g data		
item	Column 1	Column 2	Column 3		Column n	Column 1	Column 2	Column 3	Column 4	Column 5	Column 6
Line 1	Nf ₁₋₁	Nf ₁₋₂	Nf1-3		Nf1-n	Nf ₁₋₁	Nf ₁₋₂	Nf ₁₋₃	Nf ₁₋₄	Nf1-5	Nf1-6
Line 2	Nf ₂₋₁	Nf ₂₋₂	Nf ₂₋₃		Nf _{2-n}	Nf ₂₋₁	Nf ₂₋₂	Nf ₂₋₃	Nf ₂₋₄	Nf2-5	Ne ₂₋₆
Line 3	Nf ₃₋₁	Nf3-2	Nf3-3		Nf3-n	Nf3-1	Nf3-2	Nf3-3	Nf3-4	Ne ₃₋₅	Ne ₃₋₆
Line 4	Nf4-1	Nf4-2	Nf4-3		Nf4-n	Nf4-1	Nf4-2	Nf4-2	Ne ₄₋₄	Ne ₄₋₅	Ne4-6
Line 5	Nf5-1	Nf5-2	Nf5-3		Nf5-n	Nf5-1	Nf5-2	Ne ₅₋₃	Ne ₅₋₄	Ne5-5	Ne ₅₋₆
Line 6	Nf6-1	Nf6-2	Nf6-3		Nf6-n	Nf6-1	Ne ₆₋₂	Ne6-3	Ne ₆₋₄	Ne ₆₋₅	Ne6-6
Line 7	Nf7-1	Nf7-2	Nf7-3		Nf7-n	Ne7-1	Ne7-2	Ne7-3	Ne7-4	Ne7-5	Ne7-6
Table	2										
		Colu	nn 1	Column 2		Colur	nn 3	Colum	ın 4	Column 5	
Lin	e 1	Nf	1-1	Nf ₁₋₂		Nf ₁	-3	Nf ₁₋₄		Nf ₁₋₅	
Line 2		Nf	2-1	Nf ₂₋₂		Nf ₂	3 Nf ₂ .		4	Ne ₂₋₅	
Line 3		Nf	3-1	Nf3	3-2	Nf3	-3	Ne ₃	4	Ne	3-5
Line 4		Nf	4-1	Nf4-2		Ne ₄	-3	Ne ₄₋₄		Ne ₄₋₅	
Line 5		Nf	5-1	Ne	5-2	Ne ₅	-3	Ne ₅ .	4	Ne	5-5
Table	3				Į	ļ					
		Colu	nn 1	Colun	nn 2	Colun	ın 3	Colum	n 4	Colur	nn 5
Lin	e 1	Nf	1-1	Nf ₁	1-2	Nf ₁	-3	Nf ₁ .	4	Nf	1-5
Lin	e 2	Nf	2-1	Nf ₂	2-2	Nf ₂₋₃		Nf ₂ .	4	Nf ₂	2-5
Line 3		Nf	3-1	Nf3	3-2	Nf ₃₋₃		Ne ₃₋₄		Ne ₃₋₅	
Line 4		Nf	4-1	Nf4	4-2	Ne ₄₋₃		Ne ₄₋₄		Ne _{4.5}	
Line 5		Nf	5-1	Neg	5-2	Ne ₅	-3	Ne ₅ .	4	Ne	5-5
Table 4											
		Colu	nn 1	Colur	nn 2	Colun	ın 3	Colum	n 4	Colur	nn 5
Line 1		Nf	1-1	Nf ₁	1-2	Nf ₁₋₃		Nf ₁₋₄		Nf ₁₋₅	
Lin	e 2	Nf	2-1	Nf ₂₋₂		Nf ₂	Nf ₂₋₃ Nf ₂₋₄		4	Nf ₂₋₅	
Lin	e 3	Nf	3-1	Nf3	3-2	Nf3	-3	Nf3.	4	Nf	3-5
Lin	e 4	Nf	4-1	Nf4	4-2	Ne ₄	-3	Ne ₄	4	Ne	4-5
Lin	e 5	Nf	5-1	Ne	5-2	Nes	-3	Ne ₅	4	Ne	5-5

5

6 7

Figure 4 Schematic diagram of interpolation ideas

As shown in the figure above, one is trapezoidal data. As shown in Table 1 of 8 9 Figure 4, complete data and missing data in the dataset are arranged in trapezoidal order 10 according to their completeness. On the left side of the figure is the distribution of complete data. The right side of the figure shows the distribution of missing data after 11 clustering according to the number of missing data. Nf_{1-n} (the light-colored part) 12 represents the attributes of the complete number in the dataset. Nf (N: Numerals full) 13 14 represents the complete number in the dataset. The subscript number "1-n" represents the position of the n column in the first row; Ne₂₋₁ in the figure (in blue) represents the 15

attributes of missing numbers. Where Ne (Ne: Numerals Empty) represents missing numbers; The subscript number "2-1" indicates the location in the data set and is located in row 2 and column 1 of the data set. As shown in Table 1, row 2 is missing to the right (column 6) is missing one data (Ne₂₋₆); The third line is missing the two right-most data (Ne₃₋₅ and Ne₃₋₆); And so on until line 7 is the row with all missing data.

The second is to build models. Based on the trapezoidal data in the first step, as
shown in Table 2, a complete training model was established with the complete data
Nf₁₋₁, Nf₁₋₂, Nf₁₋₃, Nf₁₋₄, and Nf₁₋₅ in the first row of Table 1.

9 The third is model-based training. As shown in Table 3, the complete data in the 10 first left row of Table 2 were used to establish a training model for the complete data in 11 the second row (Nf₂₋₁, Nf₂₋₂, Nf₂₋₃, and Nf₂₋₄).

Fourth, interpolation based on model training. Based on the training in step 3, 12 interpolate Ne2-5 in Table 3 in FIG. 4 and change Ne2-5 into Ne2-5, as shown in Table 4. 13 14 After the interpolation of the second row in Table 3 is completed, the data of the second row is taken as the complete data of the first step, and the analogy continues until all 15 the Spaces of the fifth row are interpolated. The third row of Table 4 also becomes a 16 complete row data after interpolation. Just because the Spaces on the left side of Figure. 17 4 are trapezoidal after sorting, and the interpolation method is one by one from top to 18 19 bottom until all Spaces are interpolated with data, this method is called the "pyramid interpolation method" in this study. 20

In summary, the multi-layer Perceptron integration model (MLP-ESM) based on the pyramid interpolation method makes full use of non-missing information and interpolates the columns with more missing values with fewer missing values in this study. In the process of coding, the model not only extracts the data features needed for interpolation and missing data but also extracts the structural relationship features between the predicted target and the existing data, which can achieve the effect of interpolation and prediction simultaneously.

28

29 **3 Data**

As for the absence of the CCR dataset used in this study, this section will specifically introduce the CCR dataset used in this study from data collection and data preprocessing. The details are as follows.

33 3.1 Data Collection

In this study, Compustat, Bloomberg, Google, and Facebook databases were used to form the missing dataset for the CCR prediction. The dataset Span from January 1, 2009, to December 31, 2019. In addition to corporate credit risk data and its historical value, we also used data variables such as target company financial data, supply chain
data, and non-financial data. We combine financial and non-financial data to include as
much information as possible that affects our prediction results. In this study, the CCR
dataset contained a total of 167,160 pieces of data from 441 S&P rating firms, with a
data loss rate of 50%.

6 3.2 Data preprocessing

To analyze the value of all the data in this study, 118 original variables obtained 7 8 were imported into the database. See Table 1 in the appendix. The table is divided into four columns, which are the serial number, name, number, and rate of missing data, 9 10 respectively. A total of 167,160 databases have been built. It can be seen from Table 1 of the appendix that data loss has certain regularity. It can be summarized as follows:(1) 11 the feature index with the least missing is Google trend (search number), with only 784 12 missing items, with a missing rate of 0.47%; The most missing data included 29 13 characteristic indexes, such as customer Wikipedia (clicks), customer risk and customer 14 money, and 83,580 missing data, with the missing rate up to 50%. (2) Most of the 15 characteristic indexes with low miss rate belong to the financial data reflecting the 16 target company, such as return on assets, adjustment, total assets, capital turnover, and 17 debt-equity ratio; However, the missing rate of customer-related data is high, such as 18 customer-working capital ratio, customer-debt-equity ratio, customer-capital turnover, 19 20 customer-fund debt ratio, and other characteristic indexes, the missing rate is as high as 50%. The missing rate is generally in the middle level of vendor-related indicators, such 21 as vendor-leverage ratio, vendor-total assets, and vendor-cash, and the missing data rate 22 is generally around 20%. 23

24 For the constructed database, this study uses Python 3.0 platform for missing value distribution presentation. The absence degree of the top 30 feature indicators in this 25 study is shown in Figure 5. The more white lines in the figure, the more missing values; 26 27 The darker the color, the fewer missing values. It can also be found from Figure 5 that the degree of data missing from different data sources is different. However, the 28 distribution of missing data among homologous feature indicators is similar. For 29 30 example, most of the features on the right of Figure 5 belong to customer-related feature indicators. 31

- 32
- 33
- 34
- 35
- 36
- 37



7 integrity of feature indexes whose data volume is in the top 30 and 30-60, respectively.

- 8 Figure 6 and 7 are bar charts. The left bar (ranging from 0.0 to 1.0) represents the degree
- 9 of data retention (the highest retention rate is 1=100%). The upper bar represents the
- 10 actual number of data items.
- 11



12 13

Figure 6 Top 30 indicators in terms of data volume

1 As you can see from Figure 6, the biggest misses are the standard deviation of the weekly average daily return of bonds (second bar from right), followed by the working 2 capital ratio (third bar from left) and inverse cash ratio (10th bar from left). In addition, 3 most of the top 30 feature indexes with data volume are concentrated in the target 4 company's own financial data and network data, and most of the feature indexes are 5 relatively complete. It can be seen from the data collection process in Chapter 3 of this 6 study that this is mainly caused by the good degree of information recording and 7 matching. It can also be seen from the distribution in Figure 7 that most of the feature 8 9 indexes with data volume in the top 30-60 belong to the same source. It can be seen that, in the process of data restoration, the missing values in the dataset can be 10 effectively interpolated by the correlation degree between the feature indexes in the 11 12 dataset.



13 14

Figure 7 Top 30-60 indicators in terms of data volume

15

16 In order to explore the correlation between feature indexes of the dataset used in 17 this study and to interpolate other missing data based on their correlation, 30 feature 18 indexes were randomly selected from the dataset for correlation analysis. The drawn

19 thermal map is shown in Figure 8.



Figure 8 Ladder heat map based on 30 indicators

1 2

3

Figure 8 is a trapezoidal ladder diagram. The square "red on the top and bottom of 4 the basket" on the right is the scale of thermal strength. The bluer the color up is, the 5 stronger the positive correlation is. The redder the red down, the stronger the negative 6 7 correlation. As can be seen from the blue depth of the left trapezoidal thermal ladder in Figure 8, the positive correlation intensity of the 30 randomly selected characteristic 8 indexes is obvious. Figure 8 is a macroscopic demonstration of the correlation of feature 9 indicators in the dataset. Therefore, this study also conducted hierarchical clustering to 10 observe whether indicators with low miss rates in the dataset could be substituted with 11 12 indicators with high miss rates to further explore the homology of feature indicators in the dataset of this study. The hierarchical clustering results are shown in Figure 9. 13



1 From the hierarchical clustering in Figure 9, it is possible to cluster data from the 2 same source or different sources. For example, the standard deviation of the weekly average daily return of bonds, the total trading volume of bonds every Sunday, and the 3 quoted standard deviation of the weekly total daily trading volume of bonds are directly 4 5 clustered into one category, belonging to the same characteristic index. On the left, Wikipedia and Google Trends, which belong to different source metrics, also cluster 6 7 into one category after two layers of clustering. Therefore, cross-column prediction is of practical significance when missing data is predicted. 8

9

10 4 Result

Based on the MLP-ESM model constructed in the previous section, this section verifies the model in the actual data set, and compares and analyzes the effects of this research model with other traditional machine learning prediction models to show the prediction results of the model presented in this study. The specific process is as follows.

15 4.1 Prediction effect of the model

This study is based on the MLP-ESM model constructed in Section 2. The missing 16 data collected and processed in Section 3 is used as research data and analyzed on 17 Python 3.0 analysis platform. The study allocates training samples and test samples in 18 19 a ratio of 80% to 20%. The prediction performance of the model was evaluated by receiver operating characteristic curve (ROC) and area under curve (AUC). For the 20 ROC curve, the closer the curve is to the upper left corner, the more accurate the model 21 is. AUC is a performance index to measure the merits of the learner. The higher the 22 23 AUC value, the higher the accuracy of the classifier. Through analysis, the AUC of the model in this study is 83.11%, and ROC is shown in Figure 10. 24





Figure 10 ROC curve of prediction with MLP-ESM model

2 As shown in Figure 10, the prediction accuracy of this research model reaches over 80%, but in the CCR prediction, it is worth noting what level 80% belongs to. 3 Therefore, the prediction results of the model are discussed in this study with reference 4 to previous literature. When evaluating ROC curves, Yang et al. (2017) pointed out that 5 ROC curves and AUC judgment criteria are different in different application fields, 6 such as face recognition, character recognition, and other areas, and AUC standards of 7 90% are acceptable. However, in the field of psychology (Ni et al., 2020), where the 8 prediction requirement is relatively low, the AUC is acceptable at 70%. The prediction 9 of CCR discussed in this study belongs to the prediction research in the field of 10 economics. There are many factors influencing the prediction of CCR, the 11 characteristics of datasets are different, and there is an unavoidable delay in the 12 timeliness of data collection. Therefore, the acceptable accuracy of AUC value in the 13 field of economics is lower than that in the field of psychology. Generally speaking, the 14 classification of prediction grades based on AUC values in the field of economics is 15 consistent despite certain controversies (Yang & Berdine, 2017). On Yang and Berdine 16 (2017) criteria, 80% is an excellent prediction. And in this study, the prediction result 17 of the MLP-ESM model proposed exceeds 80%, which can achieve a superb prediction 18 19 effect.

20

1

21 4.2 Comparison of prediction effects

In order to further verify the prediction ability of the MLP-ESM model constructed in this study, four popular machine learning models, including support vector machine (SVM) model, XGboost (XGB) model, random forest (RF) model, and neural network (NN) model, are adopted in this section to compare their prediction effects. ROC curve and AUC values of the predicted results are shown in Figure 11 and Table5.



Figure 11 ROC plot based on five models

Table 5 AUC based on five models

Prediction method	MLP-ESM	Support vector machine	XGboost	Random forests	Neural network
AUC	83.11%	76.40%	79.01%	76.89%	74.31%

5

1 2

3 4

6 As can be seen from Table 5, the AUC values of the support vector machine, 7 XGboost, random forest, neural network, and support vector machine model are all 8 around 70% and do not exceed 80%. The prediction accuracy of the MLP-ESM model 9 constructed in this study is 83.11%. This proves the validity of the model proposed in 10 this study and shows that data restoration can improve the prediction results of the CCR 11 dataset.

In conclusion, through comparative analysis, it is found that the corporate credit risk prediction model constructed in this study with missing data can achieve an excellent prediction level. At the same time, the interpolation method proposed in this study can repair the missing dataset well, and the interpolated dataset can achieve a better prediction effect compared with the prediction of non-missing data.

17

18 5 Conclusion

In order to advance the research on CCR prediction in the absence of data, this study proposes a MLP-ESM prediction model that can perform data interpolation and prediction simultaneously. The model makes full use of non-missing information and interpolates more missing columns with fewer missing values. In the process of selfcoding, the hidden layer of the model contains variable information related to both the missing value and the predicted target, and the model can output the missing interpolation value of the independent variable and the predicted target at the same time.
In this way, not only the data features needed for missing data interpolation are
extracted, but also the structural relationship features between the predicted target and
the existing data are extracted, which can achieve the effect of simultaneous
interpolation and prediction.

The results show that the MLP-ESM model proposed in this study can realize the 6 simultaneous interpolation and prediction of missing data of CCR, and the AUC of the 7 model is 83.11%. By comparison, the prediction effect of the MLP-ESM model 8 constructed based on missing data in this study is better than that of the traditional 9 machine learning model, which fully indicates that the dataset after interpolation can 10 achieve a better prediction effect. This means that data restoration has application value 11 to improve the accuracy of CCR prediction further and provides a new research idea 12 for CCR prediction under the absence of data. 13

The current work is the basis for many future research directions. First of all, if there are too many values of continuous features in the model operation process, it can lead to sparse problems and low efficiency, which is worth noting. Secondly, how and to what extent data loss affects model performance still needs to be defined and quantified, which can provide more precise guidance for the model selection. These issues can be addressed further in the near future.

1 The appendix

Appended Table 1 List of missing values

The seriel	Appended fusie f List of missing var	The missing	
number	Feature Names	number	Loss rate
1	wiki	20073	12.01%
2	gtrends	784	0.47%
3	Working.Capital.Ratio	33154	19.83%
4	EBIT.Ratio	3509	2.10%
5	Equity.Debt	3665	2.19%
6	Capital.Turnover.Ratio	3556	2.13%
7	Retained.Earnings.Ratio	3672	2.20%
8	Return.on.Assets	3497	2.09%
9	Leverage.Measure	3497	2.09%
10	Inverse.Current.Ratio	33154	19.83%
11	Funds.to.Debt.Ratio	3497	2.09%
12	Change.in.Net.Income	3491	2.09%
13	Adjusted. Size	3497	2.09%
14	AssetsTotal	3497	2.09%
15	Cash	11709	7.00%
16	Long.Term.DebtTotal	3544	2.12%
17	InventoriesTotal	3756	2.25%
18	LiabilitiesTotal	3497	2.09%
19	Net.IncomeLoss.	3491	2.09%
20	Property.Plant.and.EquipmentTotalNet.	10586	6.33%
21	ReceivablesTotal	3923	2.35%
22	RevenueTotal	10997	6.58%
23	Market.ValueTotal	3812	2.28%
24	Price.CloseQuarter	3664	2.19%
25	Price.HighQuarter	3664	2.19%
26	Price.LowQuarter	3664	2.19%
27	weekly_mean_of_daily_sum_entrd_vol_qt	26263	15.71%
28	weekly_mean_of_daily_sd_rptd_pr	26263	15.71%
29	weekly_mean_of_daily_sd_yld_pt	81269	48.62%
30	weekly_mean_of_daily_mean_rptd_pr	26263	15.71%
31	weekly_mean_of_daily_mean_yld_pt	81269	48.62%
32	word.count	35537	21.26%
33	Tone	32451	19.41%
34	posemo	32451	19.41%
35	negemo	32451	19.41%
36	anx	32451	19.41%
37	achieve	32451	19.41%
38	risk	32451	19.41%
39	money	32451	19.41%

The serial number	Feature Names	The missing number	Loss rate
40	Customers_wiki	83580	50.00%
41	Customers_gtrends	83580	50.00%
42	Customers_Working.Capital.Ratio	83580	50.00%
43	Customers_EBIT.Ratio	83580	50.00%
44	Customers_Equity.Debt	83580	50.00%
45	Customers_Capital.Turnover.Ratio	83580	50.00%
46	Customers_Retained.Earnings.Ratio	83580	50.00%
47	Customers_Return.on.Assets	83580	50.00%
48	Customers_Leverage.Measure	83580	50.00%
49	Customers_Inverse.Current.Ratio	83580	50.00%
50	Customers_Funds.to.Debt.Ratio	83580	50.00%
51	Customers_Change.in.Net.Income	83580	50.00%
52	Customers_Adjusted.Size	83580	50.00%
53	Customers_AssetsTotal	83580	50.00%
54	Customers_Cash	83580	50.00%
55	Customers_Long.Term.DebtTotal	83580	50.00%
56	Customers_InventoriesTotal	83580	50.00%
57	Customers_LiabilitiesTotal	83580	50.00%
58	Customers_Net.IncomeLoss.	83580	50.00%
59	Customers_Property.Plant.and.EquipmentTotalNet.	83580	50.00%
60	Customers_ReceivablesTotal	83580	50.00%
61	Customers_RevenueTotal	83580	50.00%
62	Customers_Market.ValueTotal	83580	50.00%
63	Customers_Price.CloseQuarter	83580	50.00%
64	Customers_Price.HighQuarter	83580	50.00%
65	Customers_Price.LowQuarter	83580	50.00%
66	Customers_weekly_mean_of_daily_sum_entrd_vol_qt	83580	50.00%
67	Customers_weekly_mean_of_daily_sd_rptd_pr	83580	50.00%
68	Customers_weekly_mean_of_daily_sd_yld_pt	83580	50.00%
69	Customers_weekly_mean_of_daily_mean_rptd_pr	83580	50.00%
70	Customers_weekly_mean_of_daily_mean_yld_pt	83580	50.00%
71	Customers_word.count	83580	50.00%
72	Customers_Tone	83580	50.00%
73	Customers_posemo	83580	50.00%
74	Customers_negemo	83580	50.00%
75	Customers_anx	83580	50.00%
76	Customers_achieve	83580	50.00%
77	Customers_risk	83580	50.00%
78	Customers_money	83580	50.00%
79	Suppliers_wiki	37191	22.25%
80	Suppliers_gtrends	37191	22.25%
81	Suppliers_Working.Capital.Ratio	37191	22.25%

The serial number	Feature Names	The missing number	Loss rate
82	Suppliers_EBIT.Ratio	37191	22.25%
83	Suppliers_Equity.Debt	37191	22.25%
84	Suppliers_Capital.Turnover.Ratio	37191	22.25%
85	Suppliers_Retained.Earnings.Ratio	37191	22.25%
86	Suppliers_Return.on.Assets	37191	22.25%
87	Suppliers_Leverage.Measure	37191	22.25%
88	Suppliers_Inverse.Current.Ratio	37191	22.25%
89	Suppliers_Funds.to.Debt.Ratio	37191	22.25%
90	Suppliers_Change.in.Net.Income	37191	22.25%
91	Suppliers_Adjusted.Size	37191	22.25%
92	Suppliers_AssetsTotal	37191	22.25%
93	Suppliers_Cash	37191	22.25%
94	Suppliers_Long.Term.DebtTotal	37191	22.25%
95	Suppliers_InventoriesTotal	37191	22.25%
96	Suppliers_LiabilitiesTotal	37191	22.25%
97	Suppliers_Net.IncomeLoss.	37191	22.25%
98	Suppliers_Property.Plant.and.EquipmentTotalNet.	37191	22.25%
99	Suppliers_ReceivablesTotal	37191	22.25%
100	Suppliers_RevenueTotal	37191	22.25%
101	Suppliers_Market.ValueTotal	37191	22.25%
102	Suppliers_Price.CloseQuarter	37191	22.25%
103	Suppliers_Price.HighQuarter	37191	22.25%
104	Suppliers_Price.LowQuarter	37191	22.25%
105	Suppliers_weekly_mean_of_daily_sum_entrd_vol_qt	37191	22.25%
106	Suppliers_weekly_mean_of_daily_sd_rptd_pr	37191	22.25%
107	Suppliers_weekly_mean_of_daily_sd_yld_pt	37191	22.25%
108	Suppliers_weekly_mean_of_daily_mean_rptd_pr	37191	22.25%
109	Suppliers_weekly_mean_of_daily_mean_yld_pt	37191	22.25%
110	Suppliers_word.count	37191	22.25%
111	Suppliers_Tone	37191	22.25%
112	Suppliers_posemo	37191	22.25%
113	Suppliers_negemo	37191	22.25%
114	Suppliers_anx	37191	22.25%
115	Suppliers_achieve	37191	22.25%
116	Suppliers_risk	37191	22.25%
117	Suppliers_money	37191	22.25%
118	pred.weekly_mean_of_daily_mean_yld_pt	81450	48.73%

1 References

2 Akin Arikan, C., & Soysal, S. (2018). Investigation of Reliability Coefficients According to Missing Data 3 Imputation Methods. Hacettepe Universitesi Egitim Fakultesi Dergisi-Hacettepe University 4 Journal of Education, 33(2), 316-336. 5 Belue, L. M., Bauer, K. W., et al. (1997). Selecting optimal experiments for multiple output multilayer 6 perceptrons. Neural Computation, 9(1), 161-183. 7 Budagaga, A. R. (2020). Determinants of banks' dividend payment decisions: evidence from MENA 8 countries. International Journal of Islamic and Middle Eastern Finance and Management, 9 13(5), 847-871. 10 Campanile, L., Iacono, M., et al. (2020). Towards the use of generative adversarial neural networks to 11 attack online resources. Paper presented at the Workshops of the International Conference on 12 Advanced Information Networking and Applications. 13 Catellier, D. J., Hannan, P. J., et al. (2005). Imputation of missing data when measuring physical activity 14 by accelerometry. Medicine and Science in Sports and Exercise, 37(11), S555-S562. 15 Chang, C., Deng, Y., et al. (2020a). Multiple imputation for analysis of incomplete data in distributed 16 health data networks. *Nature communications, 11*(1), 1-11. 17 Chang, D., Yang, W., et al. (2020b). Seismic data interpolation using dual-domain conditional generative adversarial networks. IEEE Geoscience and Remote Sensing Letters, 18(10), 1856-1860. 18 19 Chiang, F., & Sitaramachandran, S. (2016). Unifying Data and Constraint Repairs. Acm Journal of Data 20 and Information Quality, 7(3). D'Angelo, G., & Palmieri, F. (2021). A stacked autoencoder-based convolutional and recurrent deep 21 22 neural network for detecting cyberattacks in interconnected power control systems. 23 International Journal of Intelligent Systems, 36(12), 7080-7102. 24 D'Angelo, G., Ficco, M., et al. (2021). Association rule-based malware classification using common 25 subsequences of API calls. Applied Soft Computing, 105, 107234. 26 D'Angelo, G., Tipaldi, M., et al. (2019). A data-driven approximate dynamic programming approach 27 based on association rule learning: Spacecraft autonomy as a case study. Information Sciences, 28 504, 501-519. 29 De Silva, H., Perera, A. S., et al. (2016, Sep 01-03). Missing Data Imputation using Evolutionary k-30 Nearest Neighbor Algorithm for Gene Expression Data. Paper presented at the 16th International Conference on Advances in ICT for Emerging Regions (ICTer), Negombo, SRI 31 32 LANKA. 33 Eduardo, S., Nazabal, A., et al. (2020). Robust Variational Autoencoders for Outlier Detection and 34 Repair of Mixed-Type Data. Paper presented at the Proceedings of the Twenty Third 35 International Conference on Artificial Intelligence and Statistics. 36 Florez-Lopez, R. (2010). Effects of missing data in credit risk scoring. A comparative analysis of methods 37 to achieve robustness in the absence of sufficient data. Journal of the Operational Research 38 Society, 61(3), 486-501. 39 Gao, D. Q., Yang, Z. P., et al. (2012). Performance evaluation of multilayer perceptrons for discriminating 40 and quantifying multiple kinds of odors with an electronic nose. Neural Networks, 33, 204-215. 41 Garcia, V., Marques, A. I., et al. (2019). Exploring the synergetic effects of sample types on the 42 performance of ensembles for credit risk and corporate bankruptcy prediction. Information 43 Fusion, 47, 88-101.

1	Gu, B., Li, Z., et al. (2017). Web-ADARE: A web-aided data repairing system. Neurocomputing, 253,
2	201-214.
3	Guo, X., Jarrow, R. A., et al. (2009). Credit risk models with incomplete information. Mathematics of
4	Operations Research, 34(2), 320-332.
5	Hooke, M., Mrozinski, J., et al. (2021, Mar 06-13). Salvaging Data Records with Missing Data: Data
6	Imputation using the Multivariate t Distribution. Paper presented at the IEEE Aerospace
7	Conference (AeroConf), Electr Network.
8	Jerez, J. M., Molina, I., et al. (2010). Missing data imputation using statistical and machine learning
9	methods in a real breast cancer problem. Artificial intelligence in medicine, 50(2), 105-115.
10	Khoygani, M. R. R., & Ghasemi, R. (2016). Neural estimation using a stable discrete-time MLP observer
11	for a class of discrete-time uncertain MIMO nonlinear systems. Nonlinear Dynamics, 84(4),
12	2517-2533.
13	Lan, Q. J., & Jiang, S. (2021). A method of credit evaluation modeling based on block-wise missing data.
14	Applied Intelligence, 51(10), 6859-6880.
15	Li, W., Ding, S., et al. (2020). Heterogeneous ensemble learning with feature engineering for default
16	prediction in peer-to-peer lending in China. World Wide Web-Internet and Web Information
17	<i>Systems</i> , 23(1), 23-45.
18	Mählmann, T. (2006). Estimation of rating class transition probabilities with incomplete data. Journal of
19	Banking & Finance, 30(11), 3235-3256.
20	Moon, T., Hong, S., et al. (2019). Interpolation of greenhouse environment data using multilayer
21	perceptron. Computers and Electronics in Agriculture, 166, 105023.
22	Moscato, V., Picariello, A., et al. (2021). A benchmark of machine learning approaches for credit score
23	prediction. Expert Systems with Applications, 165.
24	Nakagawa, S., & Freckleton, R. P. (2008). Missing inaction: the dangers of ignoring missing data. <i>Trends</i>
25	<i>in Ecology & Evolution, 23</i> (11), 592-596.
26	Ni, C., & Jin, X. (2020). Could L2 Lexical Attrition Be Predicted in the Dimension of Valence, Arousal,
27	and Dominance? Frontiers in Psychology, 11, 3464.
28	Nijman, S. W. J., Groenhof, T. K. J., et al. (2021). Real-time imputation of missing predictor values
29	improved the application of prediction models in daily practice. Journal of Clinical
30	Epidemiology, 134 , $22-34$.
31	Qiu, Z., Meng, M. R., et al. (2008, Jul 25-27). Missing Value Treatment of the Data Mining Based on
32	Bayesian Principle. Paper presented at the 3rd International Conference on Computer Science
33	and Education, Kalleng, PEOPLES & CHINA.
34 25	Sarker, I. H. (2021). Machine learning: Algorithms, real-world applications and research directions. 5/
30 26	Computer Science, 2(5), 1-21.
30 27	Data Paper presented at the 10th International Conference on Information and Communication
31 20	Technologies and Development (ICTD). Indian Inst Management, Almedahad, INDIA
30	Shorten C & Khoshaoftaar T M (2019). A survey on image data augmentation for deen learning
<u>40</u>	Journal of Rig Data 6(1) 1-48
41	Sivasankar E. Selvi C. et al (2016 Dec 10-11) A Study of Dimensionality Reduction Techniques with
42	Machine Learning Methods for Credit Risk Prediction Paper presented at the 3rd International
43	Conference on Computational Intelligence in Data Mining (ICCIDM) Rhubaneswar INDIA
44	Soldatvuk, N., & Sopko, S. (2014, Sep 10-12), Methods of solving missing data issues in credit risk

1	scoring and comparison of its effectiveness. Paper presented at the 32nd International
2	Conference on Mathematical Methods in Economics (MME), Olomouc, CZECH REPUBLIC.
3	Twala, B. (2013). Impact of noise on credit risk prediction: Does data quality really matter? Intelligent
4	Data Analysis, 17(6), 1115-1134.
5	Wang, G., & Ma, J. (2011). Study of corporate credit risk prediction based on integrating boosting and
6	random subspace. Expert Systems with Applications, 38(11), 13871-13878.
7	Wei, J., & Zou, K. (2019). Eda: Easy data augmentation techniques for boosting performance on text
8	classification tasks. arXiv preprint arXiv:1901.11196.
9	Yang, S., & Berdine, G. (2017). The receiver operating characteristic (ROC) curve. The Southwest
10	Respiratory and Critical Care Chronicles, 5(19), 34-36.
11	Yap, B. W., Ong, S. H., et al. (2011). Using data mining to improve assessment of credit worthiness via
12	credit scoring models. Expert Systems with Applications, 38(10), 13274-13283.
13	Yu, B. J., Li, C. M., et al. (2022). Forecasting credit ratings of decarbonized firms: Comparative
14	assessment of machine learning models. Technological Forecasting and Social Change, 174.
15	Yuan, K. P., Chi, G. T., et al. (2022). A novel two-stage hybrid default prediction model with k-means
16	clustering and support vector domain description. Research in International Business and
17	Finance, 59.
18	Yue, Y. M., Tian, J. W., et al. (2016, Sep 21-23). Applications of block chain technology in credit rating.
19	Paper presented at the 13th International Conference on Industrial Management (ICIM 2016),
20	Hiroshima, JAPAN.
21	Zhang, S. L., Wang, P., et al. (2016). Missing value data processing based on statistical correlation.
22	Statistics and Decision, (12) : 13-16.
23	