



techniques, e.g., k-Nearest Neighbor [1], [2] and [3], are more effective and flexible than statistical methods. In particular, the vital branch of [4]-[5] learning techniques [6] applied to large credit risk data lake outperform their predecessors both in accuracy and efficiency.

This paper presents a systemic review of credit risk estimation algorithms. It analyzes both the major statistical approaches and AI-based techniques with a critical spirit. The aim is to provide a comprehensive overview of the current [7]ing credit risk estimation technology, providing justification and connections between past and present works. This work proposes a novel taxonomy combining finance with [8] techniques. In addition, this work ranks their performance in terms of accuracy and costs. This paper also discusses the challenges and possible solutions in terms of four aspects: data imbalance, dataset inconsistency, model transparency, and inadequate utilization of [9] learning methods.

The remainder of the paper is organized as follows: the survey methodology will be discussed in [10]. [11] introduces the principles of statistical learning, [12] and [13] learning. [14] analyzes credit risk-related applications in detail. In [15], presented algorithms are discussed and ranked by their performance a [16]inst public datasets. Finally, results and current challenges are summarized in [17]; while [18] concludes this work.

## 2 Survey methodology

### 2.1 Methodology

We applied [19] (Preferred Reporting Items for Systematic Reviews and Meta-Analyses [20]) reviewing methodology in our paper. First, we adopted [21] searching platforms for our investigation: [22], [23], [24], [25], and [26]. We used the keywords “[27]” or “[28] learning” combined with “credit risk” while searching. We got [29] articles in total. Then, we applied a filtering algorithm considering the trade-off between publication year and citations to proceed. After removing [30] duplicate records, [31] ineligible records, and [32] incomplete articles, we obtained [33] screened records. Based on the relevance, we excluded [34] articles less related to the topic. After manually checking whether the paper has clear evaluation metrics, we further excluded another [35] papers. Finally, we kept [36] studies in terms of the relevancy to the research topic, precision of evaluation metrics, publication time, and number of citations as our source of reviewing.

Figure 1 depicts the [37] flow diagram

### 2.2 Inclusion and exclusion criteria

In this paper, we select three inclusion criteria: (1) the relevance of research topic, (2) the precision of evaluation metrics, (3) the publication year and citations. Moreover, the papers will be excluded if they are duplicated, incomplete, too early, low-related with the topic, having no clear metrics or comparatively low citations.

We show the whole workflow of the selection process in [38].

### 2.3 The datasets and approaches of the reviewed articles

The mainly used datasets by the papers under review are German and Australian public credit data from the [39] repository [40-41]. In addition, there exist some researches that discover and mine their own data. For example, [42] (43) uses data from a firm in [44] [45]. Authors in [46-47] all employ their unique dataset. Those articles mainly emphasize the significance and the veracity of the original data.

We discuss the principles and application of the overall [48] approaches. The traditional [49] models for credit risk contain [50] (SVMs) [51], k-Nearest Neighbor (k-NN) [52], [53] (RFs) [54], Decision Trees ([55] [56-57], AdaBoost [58], Extreme Gradient Boost (XGBoost) [59], Stochastic Gradient Boosting ([60] [61], Bagging [62], [63] ([64] [65] and [66] (Genetic Algorithm) [67]. Neural network models generally belong to [68] learning methods. Most of them include Convolutional [69] (CNNs) [70], [71] Belief [72] (DBNs) [73], [74] ([75] [76], LSTM (Long Short-Term Memory) [77], Restricted Boltzmann Machines (RBMs) [78], [79] Perceptron (DMLP) [80], and [81] (RNNs) [82].

Summary tables and bar charts regarding all the methods of the reviewed papers are provided.

### 2.4 Taxonomy

The taxonomy is shown as [83]. We can divide it into two parts: the first is regarding computing technology and the second is credit risk application domain. The two parts are further categorized into sub[84]ions. These two parts are connected and fused with each other. All the right-side sub-domains include the left-side techniques, and all the techniques can be applied in the financial domains.