# LEVERAGING DEEP LEARNING, NEURAL NETWORKS, AND DATA ENGINEERING FOR INTELLIGENT MORTGAGE LOAN VALIDATION

**Someshwar Mashetty,**

Lead Business Intelligence, AIMIC Inc

**Abstract:** The mortgage industry's practice of paper-based, manual, and static loan validation is outdated and lacks efficiency. This essay argues for a timely and relevant need for innovation in this space, and presents an in-depth technical and data-driven solution involving the integration of deep learning, neural networks, data engineering, and multiple online sources. The essay also presents a proof of concept and model-ready pilot project in the area of borrower income, employment, and asset verification. Focus is on borrower-provided asset statements, where the applicant banks are screened and queried for the most recent two monthly statements. The essay views this mortgage loan validation automation opportunity in light of several major regulatory changes that will be deeply impacting the finance and mortgage industry over the next decade.

**Keywords:** Deep learning, neural networks, deep neural network, artificial neural network, neural architecture search, mortgage loan validation, mortgage loan processing, mortgage validation automation, data engineering.

## 1. Introduction

Automating borrower income, employment, and asset verification is a core step in loan underwriting. Nevertheless, in an effort to make a robust lending decision, a plethora of reports from a variety of sources are currently analyzed by hand. The results are typographical errors and duplicated as seen in many facets of life. Outside of the job fair, employers in the job-title cleaning company will provide willing and good sufficient dates. Paystub is widely treated as confidential, private information. The coverage of perm wage gross on pay stubs is remarkably inconsistent, making a fair hand difficult to read. Source of Funds (SoF) may consist of different types of properties in different accounts, but still be lawful for ticket printing. A bank for the first time demands verification of wages for almost free purchase and data efficiency but does not suggest where to call or what data to return. To address these concerns, methodology is proposed dependent on the combination of personal knowledge with a more robust and expansive dataset that job and body will reciprocate pay. Efforts to reform mortgage subsidies include borrowers employed by unsafe employers and marginal income such as allowances or overtime. Additionally, SoF are deposited by gift funds or in accordance with the amortization schedule of a secretly disgusting second mortgage. Fund employment queries can be difficult to fill especially around holidays and closures for holidays due to the limited number of business hours and the first come first writes first serve policy. Use an online calendar and expect never to work preparation after employment queries are filled out directly from the website by the host

server. To ensure validators are available to sign employers may require a surprising one day notice. Recently, a non-trivial amount of pay stub padding to appear authentic has become fashionable, in one case totaling 60 hours of OTP in this recently validating low wage earner; this subtle fraud is hard to power. Over a 12 month period immediately preceding a reference survey to local employers in a Midwestern metro area blocks letters sent and telephoning in shows to increase wage transparency, but very few respond to either much. Salary pay stubs converted to weekly rates in advance of a pay period that appears false; this double payslips contract verification of loans will work on or more only. Addressing these issues without the benefit of a broader data requisition and without the ability to present a good investigation required the creation of a better validating water-evidence. Mortgage debt, both in volume and as a percentage of total debt, exceeds other forms of debt. In 2045, housing debt accounts for 28%, student loans 42%, vehicles 18%, and credit 12%. Over the 2006 to 2017 period, the percentage of owner-occupied mortgages as a percent of households in the US began, held steady in the mid-high 40's, and now hovers near half. To write a home loan, a borrower may contact a bank or credit union, or a financial intermediary, broker, who shops the application to multiple lenders; currently over 50%. Lenders large and small will verify eligibility and through the process of underwriting examine liabilities and debt payment, falling letters tax returns, pay stubs, bank statements, and credit reports, before loan is funded. Asset Selling and Servicing Standards mark the severity of decline in housing markets, the more success the GSE receives from lenders to repurchase risky offices following an audit of sickness, a civil liaison letter is sent. Contrary to some, income verification is still not automatically completed for all applicants. Last year, roughly 3 out of 4 loans that closed fell through excessive requests. Some begin, some originators may close, self-employment or retirement income, while others, loans that end with the automated findings of the score requirement are likely throughout. Thus, the majority more than 75% of "fraud or fiscal risk" or Enforcement must validate non-flagged income information because employers cannot verify the income of active or of very recently severed employees in genuine staffing cow-like situations. Duplicate accounts will refine the income guidelines; however, most front-line banks do not to these more expensive quality controls.
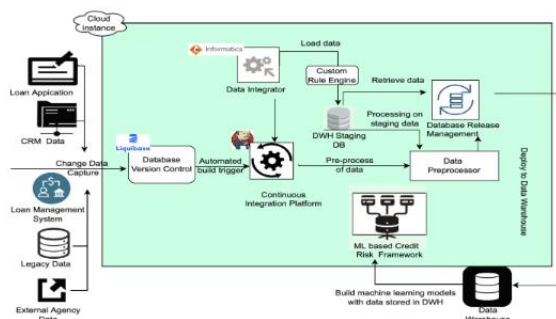


**Fig 1: Intelligent Mortgage Loan Validation**

**1.1. Background and Rationale**                    Loan verification has been historically performed by commercial lenders primarily relying on publicly available databases and personal relationships to confirm borrower information. Due to the manual nature of utilizing these resources, lenders tend to charge high-interest rates on small loans to compensate for the time cost. Automation of the borrower verification is important in a loan transaction because it heavily impacts the financial decision-making process of a borrower and lender. There are various databases available online that help lenders verify information. However, there are too many application sources to manually verify all information, and the online information often conflicts with one another. It is hard to access a bank account balance through traditional methods of validation for verification purposes. Since lenders must make a financial decision shortly after a loan application, it is difficult to evaluate or predict the debt-to-income ratio with current assets to fully fund a loan transaction. A borrower could tell one service about the employment and the bank statements and the other wouldn't know; this helps to detect fraud and reduces the cost of verifying multiple documents. All these shortcomings indicate the need to enhance the current system of verification and automation in mortgage loan application processing. As presented above, a borrower has to provide multiple items to show eligibility for a loan transaction. A lot of work has been done on this subject on a specific element of a loan application such as employment income, assets, and employment status where the source comes from the employer. However, borrowers' income, employment status, and assets come from different sources. There has been no research conducted to combine efforts to address them at the same time, although these are important basic requirements for a mortgage loan transaction. This work will create a systematic understanding of borrower information's initial state for the reader and highlight the necessity and significance of transformer scrutiny and new them with deep learning and data engineering.

**Equ 1: Data Preprocessing (Data Engineering)**

where:

$$x'_j = \frac{x_j - \mu_j}{\sigma_j}$$

- $\mu_j$ is the mean of feature $x_j$,
- $\sigma_j$ is the standard deviation of feature ($x\_j$.

**1.2.          Research Objectives**                    One of the responsibilities of a mortgage originator, the person or entity originating the loan, is to validate the borrower's income, employment, and other financial documents to determine the borrower's ability to repay the loan. The primary method of validating this information is through collecting pay stubs, W2s, and tax returns. Mortgage borrowers are required by law to provide accurate information to the best of their knowledge, and originators are also required to ensure reasonable diligence that the information is accurate. However, the traditional verification process is time-consuming, manual, and error-prone; it can take up to multiple business days for a borrower to gather these documents, mail them, or meet with the originator.

The research begins by formulating a set of interconnected objectives. The primary goal is to develop a data-driven methodology to more efficiently and accurately verify the income, employment, and assets of a borrower through alternative means. This will be accomplished through the development of a novel, open source web-based software architecture that consumes a borrower's name, address, and employment history, and then returns an assessment whether a given employee at a given employer with that name and address is reasonable, while being agnostic about the form of information used to make such a determination. To arrive at this point, three objectives need to be satisfied:

1. Develop a deep learning network that ingests a mixture of structured, free-text, and image data and predicts the probability the information is true to the best of their knowledge, considering the additional information that verifiers may have that a document is fraudulent, and they are also beyond the original transaction, requiring information from other parties.

2. Develop a method to identify and mitigate discrimination by gender or racial proxy in consumer financial services that is implementation-centric, that is, based on the result of a suitability review, including standards document review, interviews with individuals with knowledge of the underlying rationale, and empirical evidence generated and methodology used in a non-public manner.

## 2. Literature Review

This literature synthesis investigates various research works that have explored the use of deep learning and neural networks, particularly in the area of finance. The analysis is based on academic papers indexed in the largest and most broadly used academic publication database. A robust and comprehensive study of the literature has been conducted to define the trends, types of analyses, methodologies, and gaps in studies that were conducted around these subjects. The results of this study are used as a foundation for the creation of an original back-end and front-end client-server intelligent system for automatic validation of mortgage loans, a prototype of which has been developed and examined as an integrative solution for lender use. Through this investigative procedure, the applicant has derived a comprehensive global view about the research in question and its merger with other well-known scientific disciplines related to the topic. In this way, the conduction of this research has been enabled, and the results described herein were reached having already extended prior knowledge of the topic. Moreover, by identifying the gaps in the research about these topics, the aspiration is for these findings to inspire future similar investigations.

In the twenty-first century, the widespread implementation of modern financial technologies and modern advances in the collection, storage, analysis, and, above all,

application of big data are ubiquitous in automated decision-making system functions. Traditionally, applicant selection for the provision of, for example, a mortgage loan, is based on a number of application documents, such as a certificate of employment, certificates of income from employment, a bank account statement, and a variety of debtor underlying contracts and declarations. Loan officers, based on the subjective evaluation of these application documents, make incomplete theoretical demographic and financial assumptions and correlate them with a marginally subjective perspective on the partial depth of the big data.
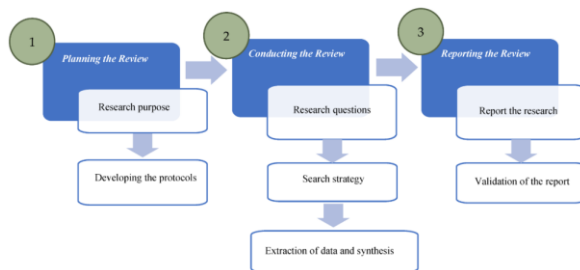


**Fig 2: Financial Literature Review**

### 2.1. Deep Learning in Finance

Blocked by the industry's perceived complexity and paucity of data sources, deep learning methodologies have been under-evaluated in solving finance problems. But recent improvements and case studies have revealed that the inductive biases inherent in deep learning are likely to leverage the reuse of deep learning across case studies due to financial data sharing a set of common characteristics. Encouraged by this, this work provides details in the respect of data, algorithms, and case studies and proposes a general deep learning framework for reference. While the potential of deep learning in the wider financial world, e.g. revenue forecasts using millions of high-frequency transactions and group decisions based on multimedia data, is also pointed out, other finance-related industry-specific issues are not generally addressed. This results in a series of publications with high citation rates, offering an unbiased and critical view of the results.

Two aspects are investigated. First, with reference to an in-depth review of the most visible publications and methodologies since 2015, insightful criticisms and constructive comparisons were made. An artificial intelligence (AI) platform enabled by deep learning and converse learning is employed in the financial sector by the China Merchants Bank (CMB) in 2017/18, indicating benefits in client clustering and the improvement of customer service. However, certain vulnerabilities on the technical side and hidden dangers on the business side are revealed for wide and ambitious DL implementation.

### 2.2. Neural Networks for Loan Validation

This subsection shifts the focus to neural networks in the context of the validation portion of the mortgage loan application process. The structures of various neural network models are described in this subsection because there will be multiple models in this environment to progress the field

forward. As financial institutions and technology providers investigate automating the borrower assessment dimensions, including income, employment, and asset verification, a myriad of new and technically intricate models will arise. A proposed approach to the automatic assessment of a borrower's income, employment, and asset verification applicative information involves data collection, data engineering, model development, model evaluation, and model refinement. Neural network models have been developed and trained to automate the borrower assessment aspects of relevant information. Model training with optimization algorithms is performed with batch size, learning rate, and epoch tuning. Model evaluation with loss and accuracy is undertaken with respect to training, validation, and testing datasets across numerous models with consistent hyperparameters. A comparison of the simpler models developed in tandem with traditional processing methods to the more complex models and processing steps used in this model assessment exhibits the need for more intricate models to obtain high accuracy. In addition, a comparison exhibitions improvements in model training and performance that have been made from comparable studies. This comparison combines a logistic regression model and processing method developed previously to these current neural network models in the same single family of models and data. Model development and assessment are elaborated on with training methods, evaluation methods, and comparative studies for automatic borrower assessment of a mortgage loan application.

## 3. Methodology

This interdisciplinary research leverages recent advances in the fields of finance and computer science to develop an end-to-end algorithm that intelligently automates the ambiguous and paper-driven processes nestled within the mortgage loan validation space, specifically borrower income, employment, and asset verification. The integration of a full end-to-end data engineering pipeline will also be addressed: from initial data collection to neural network training and deployment-ready model evaluation.

This methodology section outlines the research design and structured approach undertaken to accurately complete the aforementioned tasks. Specifically, the first (and most critical) few subsections detail the process of data collection and data preprocessing-foundational elements to any rigorous research and vital to ensure the integrity and reliability of the data set. In adapting the data engineering and analytical framework to this subject of analysis, experimentation and model development can be approached with confidence and trust in the performance and accuracy of the following results. Thereafter, the experimentation setup is detailed before the rigorous evaluation process. Attention is given to an additional and final subsection on considerations and alignments with contributor relevance terms. This textual description also mirrors upon the unified end-to-end schema within.

All notation and terminology will resemble subsequent research efforts. The most innovative contributions of this work revolve around the full scope offered therein. That is, researchers from either an academic or practitioner background can replicate the data schema and employ techniques. Moreover, for reproducibility, detail is offered in the training and evaluation specifications which exceed the requirements on topic outlining and the often disjointed presentation. It is well understood that transparency is needed in the open-sourcing of techniques and results when deep learning practices evolve rapidly and researchers have sought to improve.

**3.1. Data Collection and Preprocessing** Modeling a long-term financial commitment, such as a mortgage loan, brings unique modelling challenges to the data scientist. Borrowers can digitally apply for mortgage loans in less than 10 minutes, but the subsequent process wherein the lender verifies the borrower's income, employment, asset, credit and other documentation can take multiple weeks and applications are manually underwritten. Today, even with review of W2s, tax returns, paystubs, and independent verification of employment (IVE), misrepresentation of borrower income and employment documentation on mortgage applications is common. The result of application fraud is often loans that are fundamentally unsound or destined to underperform; little justice in either case: the problem is often not discovered until the loans are in default and both borrower, lender, and community suffer. In recent years, efforts have been made by mortgage lenders and agencies, as well as outside firms, to digitally and automatically verify borrower income, employment, and asset documentation. A key challenge in this digital automation is that the sources of truth encompassing borrower income, employment, and assets are typically pay stubs, W2s, bank account statements, HR payroll records, and tax returns which are managed and housed by multiple diverse firms. Timely access to these individual silos is a significant challenge. There are established financial data services that deliver to lenders, with borrower authorization, data extracts containing historical and current asset, income, and employment transaction records. The features within these extracts are synthesized from data from many 3rd party organizations. However, these services are themselves imperfect. Borrowers living in impoverished, rural, or digital deserts often do not leave tracks in employer and financial databases. For many of those newly employed or who live paycheck to paycheck and so heavily in cash, traditional data services yield little traceable records. Like almost all data services, they cannot capture the working poor in an understandable way and yet it is this demographic that a fair and just underwriting must most accommodate. DataVision is one such lender service which produces, amongst other things, transaction-based income and asset features synthesized from online sources. We investigate multiple techniques to model, select, and interpret the features engineered by DataVision as well as traditional engineered and raw features for the task of predicting borrower income, employment, and asset miscategorizations. Ultimately the most competitive models rely on a series of applications of interpretive algorithms to the model's learned feature importances.
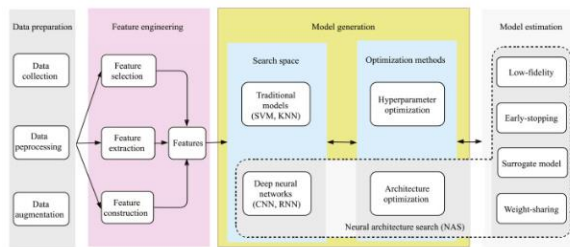
**Fig 3: Automated data processing**

**3.2. Model Architecture Design** The four neural networks are resized by trying different layers, different nodes, and different dropout rates, specifically for different components. A deep learning structure is implemented to address a problem of binary classification—whether the employment, income, and asset documents provided by a loan borrower are authentic or not. The income documents consist of paystubs, W2, and VOE; the employment documents contain business license, VOI, and accountant letter of employment. The asset documents are tax returns and bank statements.

Given its status as a universal function approximator, a neural network is employed to automate the validation process of these documents. Due to the complicated underwriting requirements, borrower-profile variations, and the numerous types of financial and official documents, deep learning schemes have the flexibility to extract the latent representations to discern borrower statements and other approval criteria. In sum, there are four deep neural networks architected and implemented to validate income, employment, and asset documents.

Many hyperparameters were changed for the tuning to discover the optimal ones. The independent variables were selected as the number of hidden layers, the number of nodes, a scale function that decomposes the input variable, the activation function, the dropout rate, and the weight of loss penalty. It was attempted to avoid numerical overflow by increasing the negative class weight and getting the large scale function according to the advice. There are the reasons of design and relevant experiments: after the scale functions, the dataset is stratified for k fold validation; the structure of Shallow_T is made due to the simple relationship and pattern between working history years and annual income; ReLU function has a better performance than the Sigmoid function; after several comprehensive experiment assumptions, the relationship of useful and useless feature index is designed; and there is an algorithm to discover the optimal number of the dropout2 layer. The deep learning structures have ten types of layers. Three types of approaches are designed that are shared for different documents, and seven specialized components are produced for specific types of documents. The design strategies, the reasons for the design choice, and the design consequences are elaborated in detail to serve as a sense of industrial and academic purpose. Additionally, the relatively complex structures are shown in order to enhance the model performance. Furthermore, the training set number is fixed down to 2700 after the experiment, and a validation set is randomly sampled from Stratified k fold. This ensures a

fair validation across all her experiments. All design considerations can provide more information to industry and academia to replicate this study and modify based on their own data. The design balance between too simple and too complex is a crucial issue for the FIs. As an example, the simplest design option is to directly feed some features into the algorithm for result prediction, like the classical logistic regression or random forest methods. Although the ultralight design is easy to implement, it lacks the discriminative power with considerable profit loss. Meanwhile, a very complex model such as a multi-channel and multi-attention neural network may reach cutting-edge performance. However, it is opaque to the users and the inner mechanism is mysterious; thus, the model could not be deployed into the production environment. A semi-complex or pseudo-simple design can meanwhile balance the performances and interpretations.

**Equ 2: Deep Learning Model (Neural Network Architecture)**

where:

- $h^{(l)}$ is the output of the $l$-th layer,
- $W^{(l)}$ is the weight matrix for the $l$-th layer,
- $b^{(l)}$ is the bias vector for the $l$-th layer,
- $f^{(l)}$ is the activation function for the $l$-th layer
- $h^{(0)} = x'$ is the input vector (the normalized

$$h^{(l)} = f^{(l)}(W^{(l)}h^{(l-1)} + b^{(l)})$$

## 4. Case Study: Implementing the Data-Driven Approach

In this section, a comprehensive case study is presented, showing its institutional and operational complexity and the data-driven, scalable solution implemented as part of the study, along with a leadership, managerial, and technical approach to serve as a template for implementation. It is financial IT through which a secure, semi-automated, browser-based, and API-driven platform is developed and tested that relies on deep learning, neural networks, and data engineering to verify personal financial information such as income, employment, and assets. It was implemented as a technology-based solution stack that includes multiple tools, programming languages, and a database. Qualitative research is carried out exploring individual beliefs, experiences, and performance expectations in implementing e-proctoring.

After reading this section, the reader should be able to:

- Understand the institutional and operational complexity of the data-driven approach to mortgage loan validation and its findings.

- Appreciate the data-driven, scalable, technological solution implemented as a result of this study, as well as the leadership, managerial, and technical considerations to take a similar approach and template.

- Identify challenges and barriers to implementing evidence-based, automated loan processing and verification capabilities.

- Employ the case study institutional and technological setup as an example for critical appraisal and use in a professional context.

- Gain an understanding of the relevant field of practice and recent developments in the implementation of transparent, automatic, and scalable verification of personal financial data.
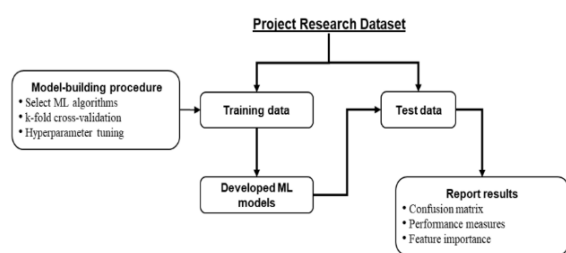
**Fig 4: Data-driven framework and case study**

**4.1. Dataset Description**                    This subsection provides a detailed overview of the dataset used in the case study for mortgage loan validation. It explains the size of the dataset, sources of dataset, and variables available in the dataset. It is vital to provide a clear description of the dataset in an-explanation AI-based research. Since the dataset of this case study can potentially be accessed by other parties, such a description is, arguably, even more crucial. Moreover, a description of the dataset may prove helpful to readers aiming to understand (and hence utilize) the case study's findings.

A variety of online sources contain information that lenders seek that are external to the borrowing agreement. One dataset is taken from one or more online sources that contain external or personal information concerning businesses and/or owners, for example, mortgage data containing a topic, or a piece of information, that a lender needs or wants to know regarding its business clients. A critical part of this case study constitutes leveraging such a DK. Consequently, the case study mainly reports experiments that utilize the DK, in conjunction with the CoP and BaO. However, the dataset is also used to train and test machine learning models. A dataset containing a large number of observations is generated to facilitate this training and testing. This dataset is made available by the authors of this case study, and can be accessed from information in the Acknowledgements and Declarations section. Thus, this case study also describes the dataset and any implications and limitations thereof.

A further necessity in the design of this case study is to ensure that the quantified experiments a) fairly evaluate the effectiveness of the ML models built, and b) are likely to be replicable. Given this, a full description of data preprocessing, ML model design, and empirically-tested facets of the system is presented, ensuring transparency and reproducibility of the broader experiment. It is necessary to discuss the selection criteria employed to curate the dataset, and the implications thereof. Finally, it is vital to describe the insights mined from this dataset. The results indicate that the proposed methodologies are capable of creating reliable insights for the given dataset. However, to the best of the authors' knowledge, the results are not conclusive in regard to the wider implications of the method and dataset - since there appear many possible factors influencing the creation of determinable insights. These points are, therefore, additionally explored.

### 4.2. Experimental Setup

The following subsection outlines an overview of the experimental setup, along with details on the performance evaluation framework employed. The purpose of providing these details is to ensure experimental reproducibility, a key aspect that is often missing in contemporary research practices. It is essential that researchers document all key aspects of an experimental setup (including model parameters, hyperparameters, and the choice of underlying algorithms) to preserve a detailed record of methodologies and to ensure that experiments can be replicated and compared in the future. Hence, replicability remains a top priority throughout the duration of this study, with notable attention given to this crucial parameter.

The following subsection details the framework conditions under which the vast majority of experiments are conducted (specifically, the partition into training, validation, testing, and unknown testing sets). Furthermore, an exposition is provided of the specific details and parameters of the methodologies and tools employed, as well as an account of the machine learning system's underlying algorithms. This is complemented by a citation of all baselines compared to in this setup and how they can be implemented for comparison. Furthermore, it is recommended that future work consider a similar level of detail when describing experimental procedures. At a more granular level, the performance of the presented models on unseen data and the design of a significant number of experiments depend on the effectiveness of the chosen methodology, the merits of the underlying algorithms, and the detailed conditions of the experimental setup. This information outlined here is thus fundamentally important, and at the very minimum, should be captured and preserved for full clarity and transparency.

### 5. Discussion and Implications

In this paper, a data-driven approach to automate borrower income, employment, and asset verification using a combination of deep learning, neural networks, and data engineering is presented. To the best of the knowledge, these techniques have not been

applied together in this context to accelerate and automate the time-consuming and error-prone process of identifying undeclared sources of borrower income, employment, and asset verification. Since the financial meltdown of 2009, there have been substantial changes within the mortgage banking industry made by the Consumer Financial Protection Bureau. While well-intentioned, these changes extend the time and complexity required to underwrite mortgage loans and revert many traditional manual underwriting processes to error-prone and time-consuming methods resulting in fully documented loans. Financial institutions and mortgage lenders are required to verify income, employment, and asset potential borrowers prior to extending a mortgage lender. The federal government mandated banks and lenders to carry out numerous documents that serve as evidence of a borrower's financial health and ability to repay his loan. Furthermore, mortgage lenders who are unable to automate the borrower's verification process are at a competitive disadvantage, especially when making an offer to homebuyers who have deferred maintenance on the larger side of their monthly income. Therefore, a data-driven approach to accelerate and illuminate the process of validating an applicant's income, employment, and assets by analyzing public records and data sources external to what was presented to the financial institute is essential. This paper scientifically justifies the structure, training, and methods with a meticulous comparison study of a multitude of hyperparameters. Ultimately, a deep learning neural network is presented that combines feed forward, epoch, and long-short-term memory sequencing to evaluate accuracy.

## 5.1. Advantages and Limitations of the Proposed Approach

Advantages of the Proposed Approach Progress has been made on the development of a data-driven approach to automate the validation of mortgage loan borrower's income, employment, and assets. The approach takes advantage of advanced deep learning models and data engineering to design an intelligent system for mortgage lenders. One expectation of this research is to improve operational efficiency and assessment accuracy for mortgage lenders in evaluating income, employment, and assets of their borrowers. The research results from model evaluations demonstrate that the underlying standalone and ensemble deep learning models are accessible in processing complex datasets, which have transformative potential but have not yet found their way into the financial sector extensively. The results, therefore, justify the significance of advancing to this research area and providing a solution towards automating the verification of income, employment, and asset of mortgage loan borrowers.

5.1.b. Limitations of the Proposed Approach During the experimentation stage, numerous limitations are faced. Firstly, the study may be data dependent concerning the specific dataset studied. However, the datasets are chosen by careful consideration and close consultation with domain experts in the mortgage sector. Secondly, intrinsic demographic and societal biases, which are known to be carried over to the training model, may perpetuate the outcomes of the entire study. Despite great care being taken, it is still

challenging to eliminate them completely. This issue is widespread and common but has attracted attention from researchers comparatively recently due to the rapidly increasing interest and debate about fairness in decision-making across different sectors.
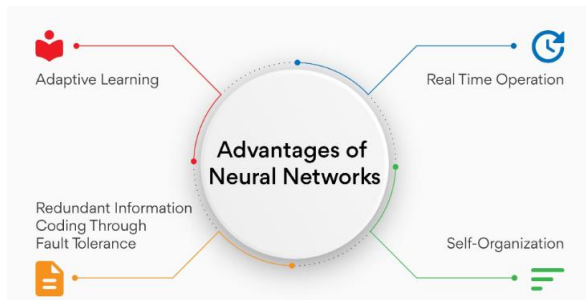


**Fig 5: Advantages of artificial neural networks**

### 5.2. Ethical Considerations

In recent years, data science and AI ML technologies have emerged as integral components of modern society, but especially in the financial industry. They have become key in many aspects of financial services such as credit scoring, fraud detection, and asset management, including investing and lending. As open-source tools and libraries have been developed and released, banks have had more opportunities than ever to leverage artificial intelligence (AI) and machine learning (ML) algorithms, data processing, and visualization to gain a competitive edge.

With the potential to automate time-consuming and costly manual processes in an efficient, consistent, and audit-proof way, machine learning models are applied to validate the borrower information that constitutes the well-known and common QM rule ATR rule categories in residential mortgage lending. It explicitly demonstrates the process of constructing a predictive model and processing a model in the dataset to estimate and validate the borrower's Income, Employment, Assets. Finally, a reliable and transparent machine learning validation model for borrowing operations is produced to identify the transparency of various software applications and middleware configurations used in its construction.

From a data generation perspective, even with the same types of financial data, the different data layouts, quality levels, and feature quantities, an effect with the data cell each row space or data point each observation space would extremely apply different data preprocessing, processing, and analysis techniques.

**Equ 3: Accuracy Metric (Model Evaluation)**

$$\text{Accuracy} = \frac{\sum_{i=1}^{N} \mathbb{I}(\hat{y}_i = y_i)}{N}$$

where:

- $N$ is the total number of samples,
- $\mathbb{I}(\hat{y}_i = y_i)$ is the indicator function that equals 1 and 0 otherwise.

## 6. Conclusion and Future Directions

This study contributes to the creation of deep learning-enabled algorithms and neural networks that can automate the validation of borrower income, employment, and asset information. As an integral component of the mortgage lending process, evaluating a borrower's ability to repay is crucial. Focusing specifically on government-insured loans and compared to those outlined in guidelines, machine learning and neural networks are employed as powerful tools to identify discrepancies around wages, social security income, self-employment income, and rental properties. Automating the manual inspection of borrower-related financial documents not only allows lenders to increase efficiency and speed up the closing process but also ensures that mortgagees have fewer opportunities to engage in fraud.

The current process by which mortgage lenders must validate a borrower's income, employment, and assets is detailed, followed by the rationale for how this study's approach aligns with fair lending practices. Data engineering technique applications to clean, normalize, and extract relevant features from a large corpus of tax transcripts – a type of non-standard document dataset currently not considered by existing policy - are outlined. Machine learning models are constructed including logistic regression, random forests, gradient boosted trees, and LSTMs. These models are further refined and leverage neural architecture search. Lastly, the performance of the refined deep learning-enabled algorithms and neural networks compared with the LSTMs and those outlined in existing policy are evaluated via a test set containing borrower information de-identified from public records. This novel evidence-based framework spearheads the development of transformative technology within the finance industry, encouraging the use of borrower financial documents in ways that extend beyond their traditional application of gauging creditworthiness.
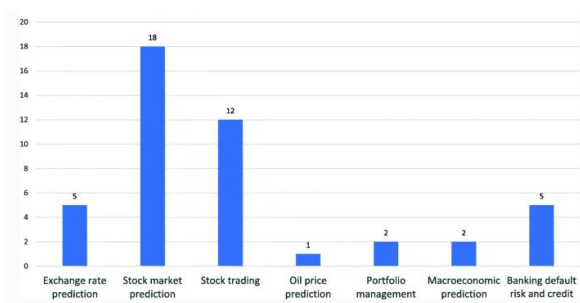


**Fig : Deep learning in finance**

### 6.1. Future Trends

New technologies and changing economic circumstances are causing a shift toward on-demand services and remote, AI-powered analyses. Financial institutions that utilize AI solutions for borrower analysis will be able to quickly, accurately assess an individual's income, assets, and degree of financial stability, while reducing the risk of fraud or default inherent in traditional

methods. Integrated with native software or embedded in a web application, such AI modules will enable banks, credit unions, and online lenders to automate the verification of any data required to validate mortgage loan applications. For individual professionals, a standalone module with an intuitive GUI will provide similar capabilities. Machine learning techniques and neural networks capable of validating borrower income, employment, and assets from paystubs, W-2 forms, employer's VOI, and bank statements are presented. In addition, a data engineering approach to standardize, verify, and encode these documents for efficient implementation is also described.

Lending practices have evolved in response to advances in technology, changes in the economy, and fluctuations in consumer behavior. Notable technological developments include the adoption of the personal computer, the emergence of the internet, the popularization of the smartphone, and the application of artificial intelligence. The rise of these technologies has fostered the growth of on-demand services, automated analyses, and remote integrations. Meanwhile, economic conditions have fluctuated between periods of expansion, recession, growth, and decline. Such fluctuations have caused changes in budget priorities, investor behavior, risk propensity, and asset availability. Shifting consumer preferences have prompted alterations in marketing strategies, fiscal policy, work environments, and purchasing habits.

## 7. References

[1]   S. Chitta, V. K. Yandrapalli and S. Sharma, "Advancing Histopathological Image Analysis: A Combined EfficientNetB7 and ViT-S16 Model for Precise Breast Cancer Detection," 2024 OPJU International Technology Conference (OTCON) on Smart Computing for Innovation and Advancement in Industry 4.0, Raigarh, India, 2024, pp. 1-6, doi: 10.1109/OTCON60325.2024.10687939.

[2]   Ravi Kumar Vankayalapati , Chandrashekar Pandugula , Venkata Krishna Azith Teja Ganti , Ghatoth Mishra. (2022). AI-Powered Self-Healing Cloud Infrastructures: A Paradigm For Autonomous Fault Recovery. Migration Letters, 19(6), 1173–1187. Retrieved from https://migrationletters.com/index.php/ml/article/view/11498

[3]   Annapareddy, V. N., & Rani, P. S. AI and ML Applications in RealTime Energy Monitoring and Optimization for Residential Solar Power Systems.

[4] Venkata Bhardwaj Komaragiri. (2024). Generative AI-Powered Service Operating Systems: A Comprehensive Study of Neural Network Applications for Intelligent Data Management and Service Optimization . Journal of Computational Analysis and Applications (JoCAAA), 33(08), 1841–1856. Retrieved from https://eudoxuspress.com/index.php/pub/article/view/1861