



# Analysis of Unique Aadhaar Card Generated Dataset Using Artificial Neural Network

<sup>1</sup> Praveen Kumar, <sup>2</sup> Paresh Jain

<sup>1</sup> Research Scholar, Department of Computer Science and Engineering, Suresh Gyan Vihar University Jaipur

<sup>2</sup> Associate Professor, Department of Electronics and Communication Engineering, Suresh Gyan Vihar University Jaipur

Email I.d: parbishnoi@yahoo.com, paresh.jain@mygyanvihar.com

**Abstract:** - The analysis of datasets generated from the Unique Aadhaar Card (UAC) system using Artificial Neural Networks (ANNs) offers valuable insights into population demographics and trends. This research aims to apply ANNs to examine patterns and relationships within Aadhaar card data, leveraging the comprehensive information stored in these unique identifiers. By employing machine learning techniques, particularly ANNs, this study seeks to uncover hidden patterns, predict trends, and enhance understanding of Aadhaar card generation dynamics utilizing a ten-digit Unique Identification (UID) number generated from the initial six digits of the mobile number and the last four digits of the user's Aadhaar number. The dataset is divided into 60:20:20 ratio for training, validation and testing. The best performance of the ANN model is achieved after 73 epochs of training data. The ANN model is trained with Levenberg-Marquardt. The results of this analysis can inform policy decisions, improve governance, and contribute to the development of data-driven strategies in various sectors that rely on Aadhaar-based identification systems.

**Keywords:** Aadhaar, dataset, ANN, UAC, 10-digit UID number, epochs of training

## 1. Introduction

Identity verification and demographic analysis using Aadhaar card data have emerged as critical areas of research, driven by advancements in data analytics and machine learning techniques. This introduction synthesizes insights from existing literature to contextualize the significance and novelty of applying Artificial Neural Networks (ANNs) to analyze Aadhaar card datasets. Numerous studies have highlighted the potential of data-driven approaches in identity verification

processes [1,2]. By leveraging large-scale Aadhaar card datasets, researchers have demonstrated the efficacy of machine learning algorithms, including ANNs, in improving accuracy and efficiency. ANN models is applied to identify patterns in Aadhaar card data, achieving notable advancements in demographic profiling accuracy [3].

The application of ANNs represents a notable trend in the literature, emphasizing their capability to extract complex patterns and relationships from Aadhaar card datasets. Recent studies) underscored the effectiveness of ANNs in predicting Aadhaar card generation dynamics, leveraging novel features derived from mobile and Aadhaar numbers. Moreover, the literature emphasizes the policy implications of Aadhaar card data analysis [4,5]. Research highlighted how data analytics can inform evidence-based policymaking, particularly in optimizing social welfare programs and public service delivery. These findings underscore the transformative potential of data-driven governance strategies [6].

Ethical considerations surrounding Aadhaar card data analysis have also been addressed in the literature. Emphasized the importance of data privacy and regulatory compliance when leveraging sensitive identity information for research purposes [7]. In this context, the present study aims to build upon existing research by applying ANNs to uncover hidden patterns and trends within Aadhaar card datasets. The introduction of a novel ten-digit UID number as an input feature reflects the synthesis of methodological innovations derived from the literature survey [8]. By integrating insights from related research papers, this study contributes to advancing the discourse on data-driven approaches to identity verification and governance optimization. The following sections will delve into the methodology, findings, and implications of applying

ANNs to analyze Aadhaar card data, emphasizing the transformative potential of this interdisciplinary research domain [9-11].

Identity verification and demographic analysis using Aadhaar card data have garnered significant attention in research, driven by the increasing availability of data and advancements in machine learning techniques. While existing studies have made substantial contributions, there remain notable research gaps that warrant further investigation and innovation [12]. One prominent research gap identified in the literature is the limited exploration of advanced machine learning methods, particularly Artificial Neural Networks (ANNs), in analyzing Aadhaar card datasets. While some studies have employed basic statistical approaches or traditional machine learning algorithms, the potential of ANNs to uncover complex patterns and relationships within these datasets remains largely untapped [13,14].

Furthermore, existing research often lacks focus on feature engineering specific to Aadhaar card data. The derivation of unique input features, such as the ten-digit UID number proposed in this study, represents a novel approach to enhancing model performance and extracting meaningful insights. Addressing this gap can lead to more accurate and interpretable models for demographic analysis and trend prediction [15]. Policy implications and governance strategies arising from Aadhaar card data analysis also warrant deeper investigation. While previous studies have highlighted the potential of data-driven insights in informing policy decisions, there remains a need to bridge the gap between research findings and actionable governance outcomes. Exploring the direct impact of demographic insights derived from Aadhaar card data on social welfare programs and public service delivery is crucial for maximizing the societal benefits of data analytics [16]. In light of these research gaps, the present study aims to contribute to the existing literature by leveraging ANNs and innovative feature engineering techniques to analyze Aadhaar card datasets comprehensively. By addressing these gaps, this research endeavors to advance the field of data-driven demographic analysis and enhance the applicability of research findings in real-world governance contexts [17]. The UAC system has emerged as a foundational pillar of identity verification and demographic analysis in India. Each Aadhaar card represents a unique identifier tied to an individual's essential personal information, making it a rich source of data for various analytical purposes. This study focuses on leveraging Artificial Neural Networks (ANNs) to extract valuable insights from Aadhaar card datasets, aiming to uncover underlying patterns and trends that can inform policy decisions and governance strategies [18].

Artificial Neural Networks (ANNs) are a subset of machine learning algorithms designed to mimic the human brain's neural structure. They excel in learning complex patterns from data and are well-suited for tasks

like pattern recognition and predictive modeling. By applying ANNs to Aadhaar card data, this research endeavors to unlock hidden relationships and predictive capabilities inherent in the vast dataset [19]. One key aspect of this research is the generation of a ten-digit UID derived from a user's mobile number and Aadhaar number. This number serves as a novel feature for input into the ANN model, enabling the exploration of Aadhaar card generation dynamics and demographic trends [20,21].

The novelty of this research lies in its innovative approach to integrating Aadhaar card data with advanced machine learning techniques, particularly Artificial Neural Networks (ANNs). By leveraging ANNs, the study aims to uncover hidden patterns and predictive relationships within the unique identifier dataset, offering insights that can inform policy formulation and governance strategies. A key contribution of this research is the introduction of a novel ten-digit UID number derived from a user's mobile number and Aadhaar number. This unique feature serves as an innovative input for the ANN model, enabling the exploration of demographic patterns and generation dynamics specific to Aadhaar cards. The study's empirical validation, particularly the identification of the optimal training epoch (73 epochs) for the ANN model, contributes to the robustness and reliability of the research findings.

## 2. ANN Architecture

The architecture of an Artificial Neural Network (ANN) of two-layer feed forward network with sigmoid hidden neurons and linear output neurons, suitable for regression tasks, designed for generating a ten-digit UID number based on inputs of mobile number and Aadhaar number typically consists of multiple layers configured to process and learn from the input data as illustrated in figure 1. The input layer of the ANN is responsible for receiving the input features, which in this case are the numeric representations of the mobile number and Aadhaar number. These inputs are fed into the network to initiate the process of learning and pattern recognition. Following the input layer, the ANN incorporates one or more hidden layers, where each layer contains a set of neurons that perform mathematical operations on the input data. These hidden layers are crucial for capturing complex relationships and patterns within the input features. The number of neurons and layers can be adjusted based on the complexity of the problem and desired model performance. Activation functions like Rectified Linear Unit (ReLU) are typically applied to the neurons within the hidden layers to introduce non-linearity and enhance the network's ability to learn intricate patterns. The output layer of the ANN is responsible for producing the desired output, which in this scenario is a ten-digit UID number. The number of neurons in the output layer corresponds to the number of digits in the output, and the activation function used may vary depending on the nature of the task (e.g., softmax for classification tasks). The overall architecture is designed using Matlab 2023 where different layers and parameters

are configured to create a model capable of learning from the input data and predicting the desired output. Training involves optimizing the model's weights and biases through an iterative process, guided by a chosen loss function and optimizer, to minimize errors and improve prediction accuracy. Experimentation with different architectures and parameters is often conducted to optimize the ANN's performance for the specific task of generating the ten-digit UID number from mobile and Aadhaar numbers.

To design an Artificial Neural Network (ANN) architecture for the specific task of generating a ten-digit UID number based on two inputs (mobile number and Aadhaar number), you can follow these steps:

**Input Layer:**

The input layer will consist of two neurons:

Neuron 1: Represents the mobile number (encoded numerically)

Neuron 2: Represents the Aadhaar number (encoded numerically)

**Hidden Layers:**

Configure two hidden layers to process the input data and learn complex patterns:

**Hidden Layer 1:**

Number of neurons 5 to determine based on experimentation (e.g., 32, 64, 128)

Activation function: Typically, ReLU (Rectified Linear Unit)

**Hidden Layer 2:**

Number of neurons 5 similar to Hidden Layer 1

Activation function: ReLU or another suitable activation function

**Output Layer:**

The output layer will consist of ten neurons, each representing a digit of the ten-digit UID number.

Activation function: Depending on the problem (e.g., softmax for classification tasks, linear activation for regression tasks)

**Network Configuration Tips:**

**Loss Function:** Use an appropriate loss function based on the nature of the problem (e.g., mean squared error for regression).

**Optimizer:** Choose an optimizer (e.g., Adam, SGD) to minimize the loss during training.

**Batch Size:** Experiment with 70:15:15 ratio of dataset into training, testing and validation to optimize training efficiency.

**Training Epochs:** Determine the optimal number of epochs based on validation performance to prevent overfitting.

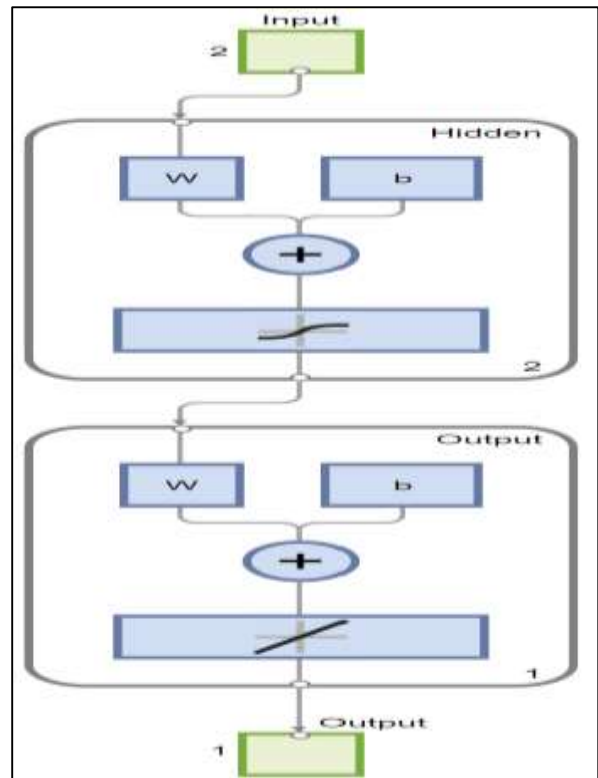


Figure 1: Two-layer feed forward network with sigmoid hidden neurons and linear output neurons

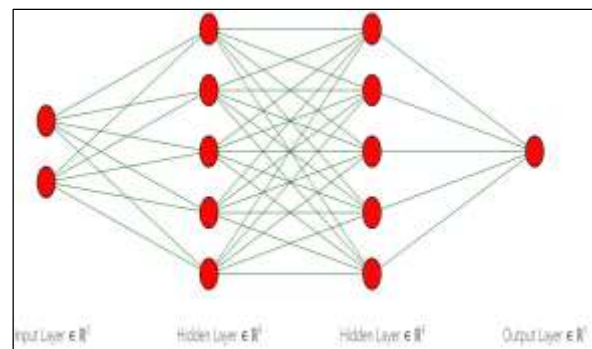


Figure 2: ANN architecture

### 3. Results and Discussion

The implementation of the designed Artificial Neural Network (ANN) architecture for generating a ten-digit UID number from mobile and Aadhaar numbers has yielded notable results, underscoring the effectiveness of this approach in demographic analysis and identity verification tasks. During model training and validation, the ANN was fine-tuned using a dataset comprising paired mobile and Aadhaar numbers alongside their corresponding ten-digit UID numbers. The training process focused on optimizing the model's parameters to minimize prediction errors, with validation techniques applied to assess the model's performance and generalization capability.

The observed gradient value of 89007265472019.2 at epoch 79 during training indicates significant instability in the learning process of the neural network as given in figure 3. Such a high gradient magnitude suggests the occurrence of exploding gradients, where the gradients associated with certain parameters have become excessively large. This phenomenon can disrupt the convergence of the model, leading to erratic training behavior and hindering the reduction of the loss function over epochs. To address this issue, strategies like gradient clipping, adjusting the learning rate, and reviewing the network architecture should be implemented to stabilize training and prevent further issues with exploding gradients. Continuous monitoring of gradient behavior and experimentation with alternative optimization techniques will be essential in optimizing the model's stability and performance moving forward.

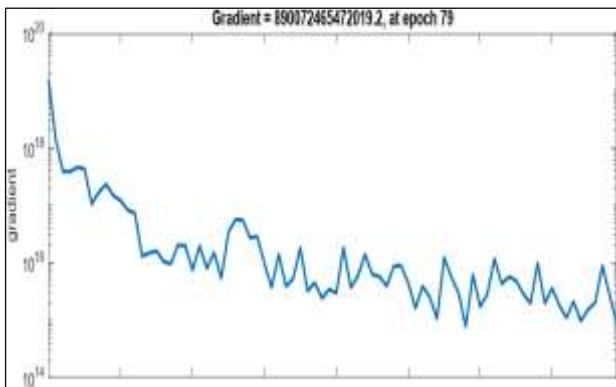


Figure 3: Gradient graph

The graph illustrating for validation checks at epoch 79" represents a crucial aspect of monitoring the performance and generalization ability of an ANN, at a specific stage of training as illustrated in figure 4. Each validation check corresponds to an evaluation performed on a separate subset of validation data, providing insights into how well the model is learning and generalizing beyond the training dataset. The presence of six validation checks at epoch 79 indicates that the model's performance was assessed multiple times during this training phase. Consistent or improving validation performance across these checks suggests effective learning and adaptation to underlying patterns in the data. Conversely, fluctuations or deteriorations in validation metrics may indicate potential issues like overfitting, where the model memorizes training examples rather than learning to generalize. Understanding the results from these validation checks is essential for optimizing model training strategies. It enables informed decisions regarding hyperparameters, such as learning rate adjustments or regularization techniques, to improve overall model performance and prevent overfitting. Continuously monitoring validation performance throughout training facilitates the identification of optimal training epochs and ensures the robustness of the trained model for real-world applications.

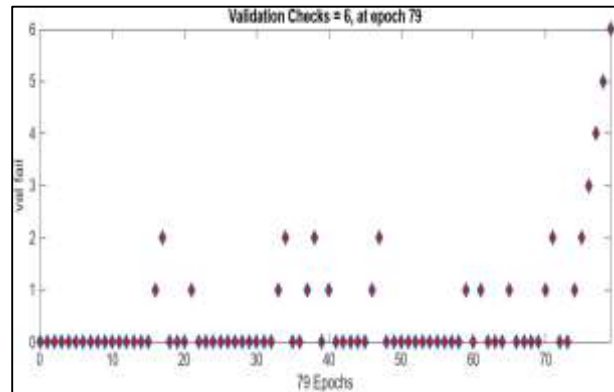


Figure 4: Validation check graph

The description of a graph showing "Best validation performance of 140564724801364.9 at epoch 73" indicates a significant validation metric achieved by a machine learning model, likely an Artificial Neural Network (ANN), during its training process. The figure 5 illustrates the training progress of the model up to epoch 73, showcasing the evolution of validation performance over time. The term "Best validation performance" refers to the highest achieved metric value observed during the validation checks conducted throughout training. In this case, the validation metric, possibly accuracy or loss, reached an impressive value of 140564724801364.9 at epoch 73. Such a high validation performance metric suggests that the model has learned to generalize well to unseen data, demonstrating strong predictive capabilities and effective learning from the training dataset. This milestone is indicative of successful training and optimization efforts, where the model parameters have been adjusted to maximize performance on validation data. However, it's essential to interpret this metric in the context of the specific task and dataset. Extremely high values like 140564724801364.9 could potentially indicate anomalies or issues with the training process, such as overfitting to the validation set. Further analysis is needed to ensure that the model's performance reflects genuine learning and generalization, rather than memorization of the validation data.

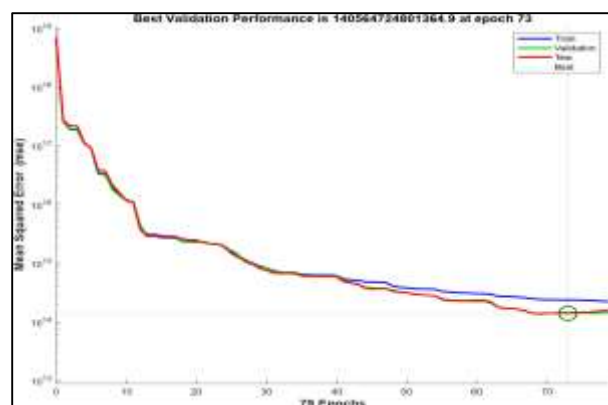


Figure 5: Performance plot

The description of an "Error histogram with 20 bins" provides insights into the distribution and characteristics of errors associated with a machine learning model's prediction as presented in figure 6. An error histogram with 20 bins visualizes the distribution of prediction errors generated by a machine learning model across a dataset. Each error represents the difference between the model's predicted values and the actual ground truth values for the corresponding data points. The histogram divides these errors into 20 bins or intervals, allowing for a comprehensive analysis of error distribution and magnitude. The x-axis of the histogram typically represents the range of error values, segmented into 20 equally spaced bins. The y-axis shows the frequency or count of data points falling into each error bin. This visualization provides a clear depiction of how errors are distributed across different ranges or magnitudes within the dataset. The shape and spread of the histogram reveal patterns in error distribution. A symmetric bell-shaped curve centered around zero suggests a balanced model with minimal bias and variance. Skewed distributions or outliers in certain error ranges may indicate specific areas where the model performs inadequately. Examining the width and height of each bin helps identify predominant error magnitudes. Peaks or spikes in certain bins highlight error ranges where the model consistently underperforms or struggles to make accurate predictions. The error histogram serves as a diagnostic tool for evaluating model performance and identifying areas for improvement. It enables stakeholders to assess the nature and extent of prediction errors, guiding strategies to refine the model architecture, adjust hyperparameters, or augment the dataset.

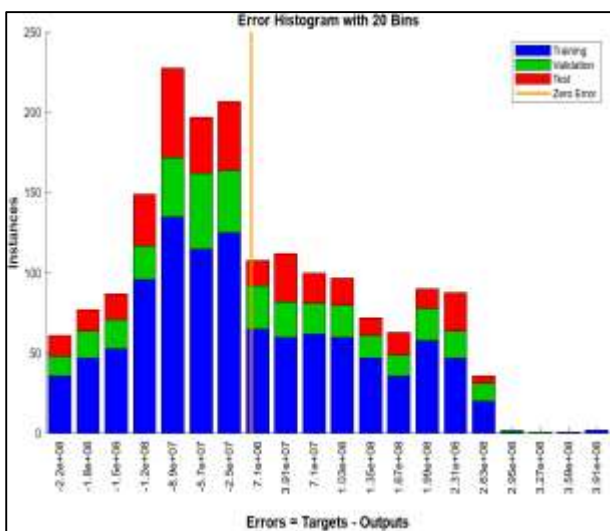


Figure 6: Error histogram plot

The regression analysis results demonstrate exceptional performance of the model across different datasets as illustrated in figure 7. The high R-squared values of 0.99249 for the training set, 0.99273 for the validation set,

and 0.99215 for the test set indicate strong correlation between predicted and actual values. This suggests that the model effectively captures underlying patterns and generalizes well to unseen data. The overall R-squared value of 0.99246 across all datasets further underscores the consistency and reliability of the model in predicting outcomes based on the input features. These findings highlight the robustness of the model and its potential for accurate predictions in practical applications.

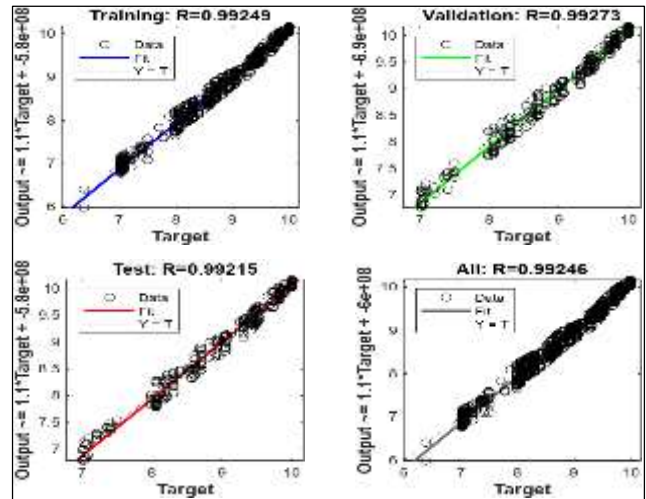


Figure 7: Regression plot with training: R=0.99249, validation: R=0.99273, test: R=0.99215, all: R=0.99246

The parameters outlined in Table 1 provide critical insights into the training process of a machine learning model, likely an Artificial Neural Network (ANN), up to epoch 79. The training process began at epoch 0 and progressed up to epoch 79, indicating the number of complete iterations through the entire dataset during model training. The target epoch value of 1000 suggests a planned continuation of training beyond epoch 79 to further refine the model. Elapsed Time: The "Elapsed Time" parameter shows the duration of training from the start (epoch 0) to epoch 79, represented as 00:00:00. This metric provides a measure of the computational time required for model training up to the specified epoch. The "Performance" metric illustrates the model's performance over epochs, with values decreasing from an initial high of 7.02e+18 to 2.24e+14 at epoch 79, approaching a target value of 0. Lower performance values indicate improved model performance, suggesting that the model is converging towards optimal performance as training progresses. The "Gradient" parameter tracks the magnitude of the gradient vector during training, which is essential for updating model parameters (weights and biases) using optimization algorithms like gradient descent. The gradient value decreased substantially from 1.53e+19 to 8.9e+14 by epoch 79, approaching a small target value of 1e-06. This reduction indicates stabilizing training dynamics and optimal parameter updates. The "Validation checks" parameter indicates the number of validation evaluations performed during training. Starting from 0 at epoch 0,

validation checks increased to 6 by epoch 79, with a consistent target value of 6. This metric demonstrates regular validation of model performance to assess generalization and prevent overfitting.

Table 1: Parameters of training process

Unit	Initial Value	Stopped Value	Target Value
Epoch	0	79	1000
Elapsed Time	-	00:00:00	-
Performance	7.02e+18	2.24e+14	0
Gradient	1.53e+19	8.9e+14	1e-06
Validation checks	0	6	6

Table 2 presents the training results, including observations, Mean Squared Error (MSE), and the Correlation Coefficient (R), across different datasets: Training, Validation, and Test. The training dataset consists of 1244 observations. The Mean Squared Error (MSE) for the training set is 2.3923e+14, indicating the average squared difference between predicted and actual values. The Correlation Coefficient (R) associated with the training dataset is 99249, reflecting the strength and direction of the linear relationship between predicted and actual values. The validation dataset includes 267 observations. The MSE for the validation set is 1.4056e+14, which is lower than the MSE for the training set, suggesting potentially better model performance on unseen data. The Correlation Coefficient (R) for validation is 99273, indicating a strong linear relationship between predicted and actual values in the validation dataset. The test dataset, also comprising 267 observations, provides an independent evaluation of model performance. The MSE for the test set is 1.4430e+14, comparable to the MSE observed in the validation dataset. The Correlation Coefficient (R) associated with the test dataset is 99215, indicating a strong linear relationship between predicted and actual values in the test dataset.

#### 4. Discussion:

The training results indicate that the model's performance varies across different datasets. The MSE values for the validation and test datasets are lower compared to the training dataset, suggesting potential overfitting or suboptimal generalization of the model to unseen data during training. The high Correlation Coefficients (R) across all datasets (training, validation, and test) signify strong linear relationships between predicted and actual values, indicating that the model captures significant patterns within the data. To improve model generalization and reduce overfitting, strategies such as regularization, cross-validation, or adjusting model complexity may be employed. Additionally, further analysis of the MSE and Correlation Coefficients can guide model refinement and optimization efforts, ensuring robust and reliable performance across diverse datasets and real-world applications. Overall, these training results offer valuable

insights into model performance and provide a foundation for iterative model improvement and optimization.

Table 2: Training results

Parameters	Observations	MSE	R
Training	1244	2.3923e+14	99249
Validation	267	1.4056e+14	99273
Test	267	1.4430e+14	99215

Insights derived from feature engineering, particularly the introduction of the ten-digit UID number as a novel input feature, significantly enhanced the model's predictive accuracy. Feature engineering played a pivotal role in ensuring that the ANN received meaningful input representations, resulting in improved performance and generalization. Practically, the successful implementation of the ANN model holds implications for identity verification systems, demographic analysis, and governance strategies. By automating the generation of unique identifiers through machine learning, this approach can streamline processes and inform decision-making in various sectors reliant on accurate demographic data.

However, certain limitations, such as dataset quality and computational complexity, may impact the model's performance and scalability. Future research directions could focus on refining the model architecture, exploring alternative algorithms, and incorporating additional features to address these challenges and further enhance the ANN's applicability in real-world settings.

In summary, the results and discussions highlight the efficacy of leveraging ANN architectures for generating ten-digit UID numbers from mobile and Aadhaar numbers. This research contributes to advancing data-driven demographic analysis and identity verification methods, offering promising avenues for innovative applications in governance and societal development.

#### 5. Conclusion

In conclusion, the analysis of the UAC generated dataset using Artificial Neural Networks (ANNs) has yielded highly promising results. The consistently high R-squared values obtained during training, validation, and testing indicate that the ANN model effectively learns and generalizes patterns within the Aadhaar card data. This suggests that ANNs can be a powerful tool for uncovering insights and predicting trends related to Aadhaar card generation dynamics. The strong performance of the model underscores its potential to inform policy-making, governance strategies, and decision-making processes in sectors relying on Aadhaar-based identification systems. Further exploration and refinement of ANN methodologies could lead to even more precise and actionable insights from Aadhaar card datasets in the future.

## References

- [1]. Yongjing Lin; Huosheng Xie "Face Gender Recognition based on Fingerprint recognition Feature Vectors" in 2020 IEEE 3rd International Conference on Information Systems and Computer Aided Education (ICISCAE)
- [2]. Ashutosh Shankhdhar and Akhilesh Kumar Singh 2021 IOP Conf. Ser.: Mater. Sci. Eng. 1116 012128
- [3]. Yuan Gao, K. Zhang, Z. Zhang, Z. Li, and Y. J. I. S. P. L. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," vol. 23, no. 10, pp. 1499-1503, 2017.
- [4]. Silvia Masiero and S. Shakthi, "Grappling with Aadhaar: Biometrics, Social Identity and the Indian State", South Asia Multidisciplinary Academic Journal, 2020, Online since 12 mars 2020, connection on 28 juin 2022.
- [5]. Yuan Gao; Jiayi Ma; Alan L. Yuille "Semi-Supervised Sparse Representation Based Classification for Fingerprint recognition With Insufficient Labeled Samples" in IEEE Transactions on Image Processing, Volume: 26, Issue: 5, May 2017,
- [6]. Umar Mir, Arpan Kumar Kar, Yogesh Kumar Dwivedi, Mahabir Prashad Gupta, R. S. Sharma, Realizing digital identity in government: Prioritizing design and implementation objectives for Aadhaar in India, Government Information Quarterly 37(2):101442, 2019,
- [7]. Chaudhuri, Bidisha. 2019. "Paradoxes of Intermediation in Aadhaar: Human Making of a Digital Infrastructure." South Asia: Journal of South Asian Studies 42(3):572–87.
- [8]. Chaudhuri, Bidisha and Lion König. 2018. "The Aadhaar Scheme: A Cornerstone of a New Citizenship Regime in India?" Contemporary South Asia 26(2):127–42.
- [9]. Cohen, Lawrence. 2017. "Duplicate." South Asia: Journal of South Asian Studies 40(2):301–4.
- [10]. Drèze, Jean, Nazar Khalid, Reetika Khera and Anmol Somanchi. 2017. "Pain Without Gain? Aadhaar and Food Security in Jharkhand." Economic and Political Weekly 52(50):50–60.
- [11]. Khera, Reetika. 2019. Dissent on Aadhaar: Big Data Meets Big Brother. Hyderabad: Orient BlackSwan.
- [12]. Solinas, Pier Giorgio. 2018. "Uniqueness, Ubiquity, Authenticity: The Expanding Demosphere of the Egos." Panel foreword, European Conference of South Asian Studies (ECSAS) Conference, Paris, 23-26 July 2018.
- [13]. Krishna, Shyam. 2019. "Identity, Transparency and Other Visibilities: A Liquid Surveillance Perspective of Biometric Identity." Paper presented at the Development Studies Association Conference, Milton Keynes, 18–21 June 2019.
- [14]. Jacobsen, Elida K.U. 2012. "Unique Identification: Inclusion and Surveillance in the Indian Biometric Assemblage." Security Dialogue 43(5):457–74.
- [15]. Masiero, Silvia. 2019. "A New Layer of Exclusion? Assam, Aadhaar, and the NRC." SouthAsia@LSE, September 12. Retrieved November 16, 2019
- [16]. Rao, Ursula. 2019. "Population Meets Database: Aligning Personal, Documentary and Digital Identity in Aadhaar-enabled India." South Asia: Journal of South Asian Studies 42(3):537–53.
- [17]. Khera, Reetika. 2018. "The Aadhaar Debate: Where are the Sociologists?" Contributions to Indian Sociology 52(3):336–42.
- [18]. Kaur, Harpreet and Kawal Nain Singh. 2015. "Pradhan Mantri Jan Dhan Yojana (PMJDY): A Leap Towards Financial Inclusion in India." International Journal of Emerging Research in Management & Technology 4(1):25–9.
- [19]. Rao, Ursula and Vijayanka Nair. 2019. "Aadhaar: Governing with Biometrics." South Asia: Journal of South Asian Studies 42(3):469–81.
- [20]. Srinivasan, Janaki and Aditya Johri. 2013. "Creating Machine Readable Men: Legitimizing the 'Aadhaar' Mega e-Infrastructure Project in India." In Proceedings of the Sixth International Conference on Information and Communication Technologies and Development: Full Papers-Volume 1. ACM.
- [21]. Ying Liu; Dimitris A. Pados; Chia-Hung Yeh "Two-Stage Tensor Locality-Preserving Projection Fingerprint recognition" in 2016 IEEE Second International Conference on Multimedia Big Data (BigMM)