# EXPERIMENT REPORT

| | |
|---|---|
| **Student Name** | Vishal Raj (14277627) |
| **Project Name** | AT3 - Data Product with Machine Learning |
| **Date** | 9th November, 2023 |
| **Deliverables** | main_vishal_raj.ipynb<br>train_model_vishal_raj.py<br>best_model-vishal_raj<br>Wide and Deep Neural Network<br>https://github.com/vishalraj247/Flight_Fare_Prediction.git<br>https://github.com/vishalraj247/Flight_Fare_App.git |

## 1. EXPERIMENT BACKGROUND

Provide information about the problem/project such as the scope, the overall objective, expectations. Lay down the goal of this experiment and what are the insights, answers you want to gain or level of performance you are expecting to reach.

| | |
|---|---|
| **1.a. Business Objective** | The business aims to establish a robust flight fare prediction system to inform pricing strategies. Such a system is designed to provide airlines with a tool for better revenue management through dynamic pricing, help travel agencies craft competitive package deals, and assist travellers in budgeting by predicting fare trends. An accurate prediction model could significantly impact revenue streams and customer satisfaction by minimising flight costs' uncertainty. |
| **1.b. Hypothesis** | My hypothesis posits that the wide and deep neural network, with its hybrid structure, is superior for the task of predicting flight fares compared to traditional models. This theory is premised on the belief that this architecture will handle both the feature interaction intricacies and the unique hierarchical patterns within the fare data effectively, thus achieving a lower RMSE and MAE on the test set. |

| | |
|---|---|
| **1.c. Experiment Objective** | The anticipated outcome is achieving precise fare predictions with an RMSE of less than 180 and an MAE close to 120 for the test dataset. These figures were chosen based on industry standards and previous model performances. The experiment seeks to assess the model's predictive strength, identify any performance gaps, and determine the feasibility of implementing more complex modelling techniques if necessary. |

| **2. EXPERIMENT DETAILS** |
|---|
| Elaborate on the approach taken for this experiment. List the different steps/techniques used and explain the rationale for choosing them. |

| | |
|---|---|
| **2.a. Data Preparation** | The data preparation for our experiment involved aggregating flight fare data from various sources into a unified dataset. Using make_dataset.py, we loaded and combined zipped files into one CSV file, creating a consistent data structure for analysis. The data_preprocessor_dl.py script then enriched the dataset with date-related features and transformed categorical variables like airport codes and cabin classes into numerical form using a robust encoding strategy. We handled anomalies by normalising extreme fare values and optimising data types to enhance model training efficiency. After rigorous cleaning and augmentation, the processed data was archived ready for model ingestion and future updates. |
| **2.b. Feature Engineering** | For feature engineering, we considered temporal features such as day of the week and time of day, which previous studies suggested could impact fares. We employed embedding layers for categorical features to reduce dimensionality and improve model interpretability. Some features, such as airline-specific codes and names, were excluded to prevent overfitting and ensure the model's applicability across various input features. |

| 2.c. Modelling | I employed a wide and deep neural network model, leveraging a combination of linear and neural network layers to capture different data abstractions. The model's hyperparameters, including the number of neurons in dense layers (256 and 128) and dropout rates (0.5), were derived from empirical evidence suggesting their effectiveness in balancing learning capacity and generalization. I decided against a full grid search for hyperparameter tuning due to time and resource constraints but noted this as an area for potential future exploration. |
|---|---|

## 3. EXPERIMENT RESULTS

Analyse in detail the results achieved from this experiment from a technical and business perspective. Not only report performance metrics results but also any interpretation on model features, incorrect results, risks identified.

| 3.a. Technical Performance | The model's technical performance, with an RMSE of 177.048 and MAE of 120.444 on the test dataset, suggests that it is capable but still has room for improvement, especially in handling outliers and rare events that cause fare spikes. |
|---|---|
| 3.b. Business Impact | From a business standpoint, the model's current performance can provide value in enhancing pricing strategies, although the relatively higher RMSE may lead to less competitive pricing in some scenarios. The model's predictions can be integrated into current systems with caution, using a phased approach to monitor and improve prediction accuracy with ongoing data collection. |
| 3.c. Encountered Issues | The main issues encountered revolved around computational efficiency and managing the extensive requirements of deep learning models. One strategy adopted was to prioritise the depth of feature learning over extensive hyperparameter tuning. We also acknowledged the need to enhance our data pre-processing in future experiments to handle anomalies more effectively. |

## 4. FUTURE EXPERIMENT

| | |
|---|---|
| Reflect on the experiment and highlight the key information/insights you gained from it that are valuable for the overall project objectives from a technical and business perspective. | |
| **4.a. Key Learning** | The key learning from this experiment is the potential viability of the wide and deep model for flight fare prediction, with certain limitations in its current form. Continuous model refinement and incorporation of new data are critical for improving the model's accuracy and reliability. |
| **4.b. Suggestions / Recommendations** | Future recommendations include increasing the diversity and volume of data to train more generalized models, implementing more sophisticated hyperparameter tuning when resources allow, and conducting ablation studies to understand feature contributions. If the model stability is confirmed, we suggest deploying it in a controlled production environment with real-time monitoring and performance tracking systems. |