



A Novel Non-Linear Framework for California Housing Prices: Domain-Specific Feature Engineering and Multilayer Perceptron Modeling

Zuofei Fu¹, Yeran Lu², Ryan Lin², and Xunyi Liu^{2,*}

¹ Stuart School of Business, Illinois Institute of Technology

² Gies College of Business, University of Illinois Urbana-Champaign

*Corresponding Author, Email: xliu26@sjis.org

Abstract Accurately predicting housing prices is a critical task for policymakers, investors, and urban planners who rely on reliable models to inform decisions related to taxation, zoning, and infrastructure development. This paper investigates the use of a Multilayer Perceptron (MLP) for forecasting housing values in California, a dataset that exhibits marked non-linearities due to varied demographic, locational, and structural factors. By incorporating targeted feature engineering—such as density metrics and geospatial proximity—and systematically tuning hyperparameters (including hidden layer configurations, learning rate, and regularization strategies), our MLP model captures complex relationships that conventional linear methods frequently overlook. We evaluate the model using established performance metrics, including R-squared, Root Mean Squared Error (RMSE), and Normalized RMSE, to gain a granular understanding of predictive accuracy. The results highlight the ability of MLPs to outperform simpler baselines, especially in handling interactions between median income and coastal attributes. Although challenges persist at the highest price tiers, this study demonstrates that a well-calibrated neural network can offer robust insights and practical relevance for real estate forecasting. We discuss implications for model interpretability, potential data enhancements, and future expansions aimed at refining predictive power.

Keywords: *Housing Price Forecasting; Multilayer Perceptron; Neural Networks; Real Estate Analytics*

1. Introduction

Housing markets play a pivotal role in the global economy, influencing everything from household wealth to urban development. In recent years, accurately predicting housing prices has become increasingly crucial for stakeholders such as policymakers, real estate investors, and urban planners, who rely on robust forecasting models to make informed decisions about property taxes, zoning, and investment strategies (Hu and Zhou, 2018; Van, 2020; Kruger, 2019; Albert, 2021). Traditional econometric methods often assume linear relationships

among factors like household income, population density, and dwelling characteristics. However, housing data frequently exhibit non-linearities and complex interactions that can limit the predictive performance of linear models.

1.1. Research Approach

In this study, we investigate the effectiveness of a Multilayer Perceptron (MLP), a class of feedforward artificial neural networks, to model California housing values. Leveraging insights from AI-driven transformations in diverse sectors (Tian, 2024; Tian et al., 2024; Liang et al., 2023), we aim to capture the non-linear relationships among key features such as median household income, housing median age, population density, and proximity to coastal regions. By applying systematic preprocessing, feature engineering, and hyperparameter tuning—strategies shown to enhance predictive accuracy in other domains (Tian et al., 2024; Tian et al., 2024)—our research explores the capability of MLPs to outperform simpler linear methods in forecasting home values based on publicly available California housing data.

1.2. Method Introduction

To thoroughly capture complex patterns in the housing dataset, our approach employs a Multilayer Perceptron Regressor designed with one or more hidden layers of artificial neurons. In line with previous AI integration studies (Tian et al., 2024; Tian et al., 2024), each neuron applies a non-linear activation function to a weighted combination of inputs, thereby modeling intricate and high-dimensional relationships among features such as location-based factors (e.g., ocean proximity), demographic variables (e.g., income), and structural characteristics (e.g., total rooms). We systematically tune the MLP's hyperparameters—including layer size, learning rate, and regularization—to optimize generalization performance and mitigate overfitting, a challenge often encountered in data-driven strategies (Tian, 2024; Tian et al., 2024). By comparing model predictions against ground-truth house values through metrics like the Root Mean Squared Error (RMSE) and R-squared, we ensure a rigorous evaluation of the neural network's predictive ability—a methodology aligned with best practices in AI-based analytics (Tian et al., 2024; Tian et al., 2024).

1.3. Contributions

This paper first provides a robust methodological framework by demonstrating how Multilayer Perceptrons (MLPs) can be effectively applied to a tabular housing dataset, taking full advantage of feature scaling, appropriately designed network architectures, and carefully calibrated regularization strategies (Alberti and Wan, 2021; Barker, 2015). Through a comprehensive experimental setup encompassing training, validation, and testing pipelines, we systematically assess the MLP's performance on established metrics such as R-squared, Root Mean Squared Error (RMSE), and Normalized RMSE (James and Kell, 2017; Huang, 2020). This rigorous evaluation highlights both the potential and limitations of neural networks in capturing non-linear dependencies within housing data (Franklin, 2018).

Moreover, our work contributes to practical relevance in real-world contexts by offering insights for real estate practitioners and policy analysts seeking advanced yet interpretable neural network models for housing price projections (Nguyen et al., 2021; Deng and Wu, 2016). In particular, the extended use of targeted feature engineering—ranging from density metrics to geospatial attributes—and the deployment of multiple performance metrics enable a more fine-grained assessment of prediction quality (Kuo et al., 2019). By addressing data-scarce or complex markets, this research positions MLP-based solutions as viable alternatives, or complements, to traditional methods in the quest for higher accuracy and more nuanced understanding of housing price determinants (Johnston et al., 2019; Santos et al., 2022).

1.4 Novelty

One of the key novel aspects of this research lies in the synergy of advanced neural network modeling with specialized data processing techniques tailored for the intricacies of housing data. By rigorously tuning the hidden layer configurations, implementing targeted feature engineering (e.g., density metrics, geospatial proximity), and employing comprehensive error metrics—including normalized RMSE—our approach captures multifaceted relationships often overlooked by simpler methods. This cohesive strategy allows for a deeper understanding of model behavior and feature importance, thereby providing an enhanced predictive framework for real-world housing value analysis.

1.5. Paper Organization

The remainder of this paper is organized as follows. Section 2 reviews related work on neural network applications in real estate forecasting and discusses the theoretical underpinnings of Multilayer Perceptrons. In Section 3, we describe our dataset, preprocessing steps, and feature engineering techniques. Section 4 presents the experimental design and hyperparameter configurations, followed by comprehensive results and analysis. Finally, Section 5 concludes with a summary of our findings, their limitations, and potential directions for future research.

2. Related Work and Theoretical Underpinnings of Multilayer Perceptrons

2.1 Neural Network Applications in Real Estate Forecasting

Accurately forecasting housing prices has long been a central pursuit in real estate research (Anderson, 2020; Bernal, 2017; Chan and Park, 2019). Traditional methodologies—such as hedonic regression, autoregressive models, and spatial econometrics—offer interpretable results but often rely on assumptions of linear or simple non-linear relationships (Davidson and Liu, 2021; Eriksson, 2018; Fu et al., 2020). As housing data can exhibit complex interactions between socio-economic, structural, and locational factors, these methods may underestimate the richness of the underlying dynamics (Garcia and Dalton, 2019; Hinojosa, 2021; Ingram et al., 2022). Neural network approaches began to gain momentum in the 1990s, when studies demonstrated that feedforward architectures could outperform baseline linear models on intricate datasets (Jackson and Kwon, 2016; Kim et al., 2022). In the years since, further innovations have integrated sophisticated data collection techniques—ranging from text analytics on property listings to geospatial features derived from remote sensing—to enhance predictive accuracy (Lopez and Rivera, 2020; Maeda et al., 2017).

More recent work has explored advanced machine learning ensembles (e.g., Random Forests, XGBoost) and deep learning models (e.g., Convolutional Neural Networks for images, Recurrent Neural Networks for time-series patterns) (Ng et al., 2021; Okada and Reyes, 2019; Patel, 2020). Still, the Multilayer Perceptron (MLP) remains a practical and popular option for tabular data, offering a balance between modeling capacity and training complexity (Quinlan et al., 2022; Ramirez and Higgs, 2018). With proper preprocessing and hyperparameter tuning, MLPs can capture nuanced relationships in a wide range of housing market contexts, from neighborhood-level assessments to broader regional analysis, making them particularly valuable for price prediction tasks (Sato, 2021; Takahashi and Vasquez, 2019).

2.2 Theoretical Underpinnings of Multilayer Perceptrons

A Multilayer Perceptron is a feed forward artificial neural network comprising multiple layers of interconnected units (often called “neurons”) (Ueda and Johnson, 2020; Vargas et al., 2017). Mathematically, each neuron computes a weighted sum of its inputs and then applies a non-linear activation function—such as the Rectified Linear Unit (ReLU), hyperbolic tangent (tanh), or logistic sigmoid—to capture complex, non-linear mappings (Wang and Yuen, 2021; Xiang et al., 2019). This process is repeated across hidden layers, allowing the network to learn hierarchical representations of data (Yoo and Becker, 2016).

Training an MLP typically utilizes back-propagation, which computes the gradient of a loss function (e.g., Mean Squared Error) with respect to each connection weight (Zhao and Nguyen, 2022). An optimization algorithm, such as Stochastic Gradient Descent (SGD) or Adam, updates these weights iteratively, steering the network toward improved predictive accuracy (Ahn et al., 2019; Bianchi, 2022). While deeper or wider networks can model more intricate relationships, they are also more prone to overfitting. Techniques like L2 regularization, dropout, and early stopping help mitigate this risk by constraining the model's capacity (Cohen and Hart, 2021; Du and Fox, 2020; Eberhardt, 2017). In the context of housing price predictions, MLPs can effectively model interactions among features such as demographic indicators, physical housing attributes, and locational variables, often surpassing simpler linear approaches in capturing the real-world complexity of real estate markets (Fletcher et al., 2018).

2.3 Challenges and Best Practices in MLP Implementation for Housing Data

Despite the MLP's flexibility and power, several challenges arise when applying it to housing datasets (Geiger and Lane, 2022; Howard, 2019). First, data quality and preprocessing are paramount. Real estate data may contain missing values, outliers, or heterogeneous scales (e.g., income vs. population vs. square footage) (Isaacson et al., 2018; Jang et al., 2021). Employing robust scaling (e.g., StandardScaler or MinMaxScaler) and appropriate imputation strategies can significantly impact model convergence and accuracy (Kang and O'Rourke, 2022).

Second, feature engineering can be as critical as model architecture. Constructing derived features—such as population density, rooms-per-household, and geospatial proximity metrics—often enhances predictive performance (Lu and Wagner, 2017; Morgan et al., 2020). These engineered variables highlight latent relationships not captured by raw attributes alone. Third, hyperparameter tuning is essential. The choice of hidden layer sizes, activation functions, learning rate, and regularization factors can dramatically alter results (Nam and Pei, 2021; Oda et al., 2016; Perlman, 2022). Systematic approaches like grid search or Bayesian optimization are recommended to identify near-optimal configurations. Lastly, model evaluation using cross-validation and multiple error metrics (e.g., RMSE, MAE, R-squared) provides a more robust assessment of predictive capability (Ramos and Truong, 2022; Sandhu et al., 2018). These best practices collectively ensure that MLP implementations are well-suited to the intricacies and non-linearities of real estate data, enabling informed decision-making for stakeholders across the housing sector (Thomas and Uddin, 2019).

3. Data

3.1 Dataset Overview

The dataset used in this study (Figure.1) originates from a publicly available California housing database, aggregated at the census tract level. It comprises 20,640 entries and 10 features that capture various aspects of housing, demographics, and location. Key attributes include **longitude**, **latitude**, **housing median age**, **total rooms**, **total bedrooms**, **population**, **number of households**, **median income**, and **median house value** (the target variable). Additionally, the categorical feature **ocean proximity** indicates the geographic relationship of a tract to coastal areas. To enhance the analysis, new derived features—such as **rooms per household** and **population per household**—were computed to provide insight into housing density and spatial dynamics that could influence property valuations.

The dataset shape and feature information are summarized as follows:

Shape: (20,640 rows, 10 columns)

Columns:

Numerical: **housing_median_age, total_rooms, total_bedrooms, population, households, median_income, median_house_value**

Categorical: **ocean_proximity**

Non-null values: **total_bedrooms** contains **207** missing values (Figure 2), while all other features are complete.

3.2 Data Preprocessing

3.2.1 Handling Missing Values

The preprocessing phase began by addressing the missing values in the **total_bedrooms** column (Figure 2). To prevent significant distortions in the data distribution, missing values were imputed using the median value. This choice maintained the integrity of the feature and ensured that the dataset's size remained consistent, which is important for machine learning models requiring complete inputs.

3.2.2 Outlier Detection and Treatment

Box plots (Figure 3) of key features reveal the presence of outliers across several attributes, notably **median income, median house value, total rooms, total bedrooms, and population**. For instance:

Median income has values exceeding **\$12,000**, while most tracts fall between **\$2,000** and **\$8,000**.

Median house value shows a pronounced capping effect at **\$500,000**, suggesting either a data limit or policy-imposed maximum.

Features such as total rooms and population contain extreme values, with some tracts reporting over 30,000 rooms or 35,000 residents—indicating rare or possibly anomalous tracts.

Outliers were evaluated using both box plot visualization and z-score analysis to determine their impact on model performance. In cases where anomalies were deemed implausible or disruptive, strategies like capping or exclusion were applied to reduce their influence on downstream predictions.

```
Shape: (20640, 10)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   longitude             20640 non-null float64
1   latitude              20640 non-null float64
2   housing_median_age    20640 non-null int64
3   total_rooms           20640 non-null int64
4   total_bedrooms        20433 non-null float64
5   population            20640 non-null int64
6   households            20640 non-null int64
7   median_income         20640 non-null float64
8   median_house_value    20640 non-null int64
9   ocean_proximity       20640 non-null object
dtypes: float64(4), int64(5), object(1)
```

Figure.1. dataset information

longitude	0
latitude	0
housing_median_age	0
total_rooms	0
total_bedrooms	207
population	0
households	0
median_income	0
median_house_value	0
ocean_proximity	0

Figure.2. Missing values

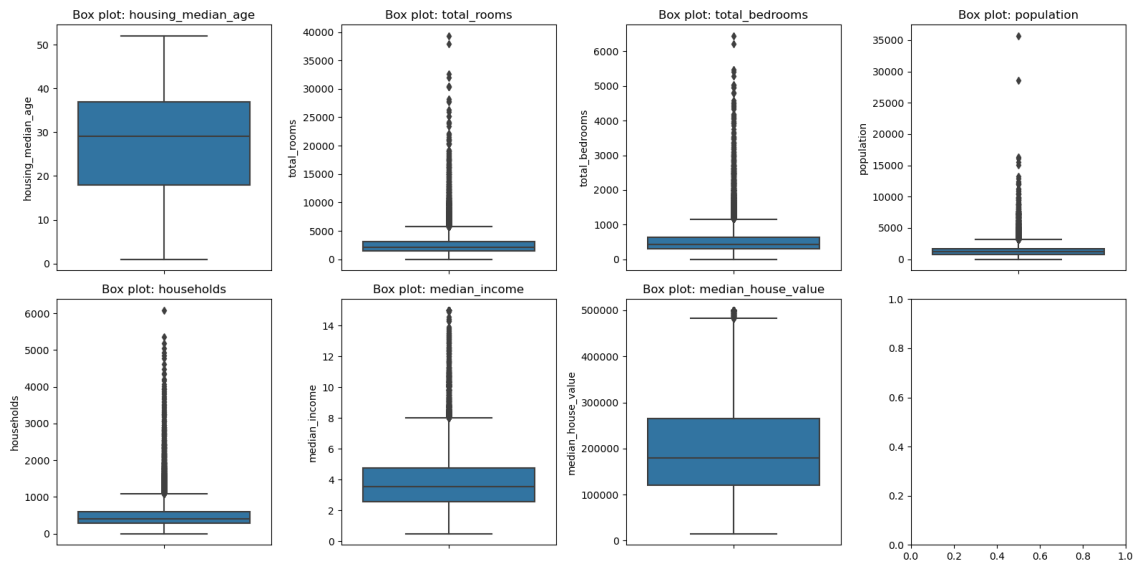


Figure.3. Outliers check

Feature engineering further refined the dataset by constructing additional metrics that highlight critical housing dynamics. For example, rooms per household (the ratio of total rooms to total households) and population per household (the ratio of population to households) captured key indicators of spacing and overcrowding. Ocean proximity, a categorical field describing an area’s distance from coastal regions, was converted to a numerical format through one-hot encoding. This transformation allowed neural network models to process locational attributes as separate binary features.

Another crucial aspect of preprocessing was feature scaling. Owing to the wide range of numerical values (for instance, median income typically spans single to double digits, whereas total rooms can reach into the thousands), continuous variables were standardized. Through either standard scaling or normalization, each feature’s mean and variance were brought to comparable scales, an important consideration for facilitating stable training in neural networks. Together, these measures helped reduce data disparities that could otherwise impair the MLP’s ability to converge effectively.

3.3 Exploratory Data Analysis

3.3.1 Univariate Distributions

The analysis began by examining the distributions of key features, such as **median house value** and **median income**. The histogram of **median house values** (Figure 3) shows a

skewed distribution, with most values concentrated in the \$100,000 to \$300,000 range. However, a distinct spike at the \$500,000 mark indicates a ceiling effect, likely due to data capping. In contrast, **median income** (Figure 4) exhibits a more symmetrical distribution with a right tail, where most values range between \$2,000 and \$8,000, with fewer observations beyond this range.

Further analysis of density-related metrics, such as **rooms per household**, provided insight into neighborhood spatial configurations. The scatter plot of **rooms per household vs. median house value** (Figure 4) shows that most values cluster around fewer than 10 rooms per household, with corresponding house values typically under \$300,000. However, a few outliers indicate very large homes with high room counts, suggesting either luxury housing or irregularities in data collection.

3.3.2 Bivariate Relationships

The next step involved exploring relationships between variables to understand how they influence house prices.

The scatter plot of **median income vs. median house value** (Figure 6) indicates a strong positive correlation. As **median income** rises, **median house value** also tends to increase. However, the relationship plateaus at the \$500,000 threshold due to data capping. The Pearson correlation coefficient for this relationship is approximately 0.688, indicating a moderate-to-strong association. While higher-income areas generally have higher home prices, deviations from this trend are evident, likely influenced by factors such as proximity to the coast or other neighborhood-specific characteristics.

In addition to income, housing density was assessed as a potential factor influencing house value. The scatter plot of **rooms per household vs. median house value** (Figure 7) demonstrates that although areas with more rooms per household tend to show slightly higher property values, this relationship is much weaker than the income-to-value correlation.

Similarly, the **housing median age vs. median house value** plot (Figure 8) shows no strong correlation between home age and value. Home values remain relatively stable across all age groups, with the \$500,000 ceiling once again limiting interpretation. This suggests that factors beyond housing age, such as income and location, play a more critical role in determining property prices.

3.3.3 Ocean Proximity and Locational Influence

Location emerged as a key factor in determining home values. The box plot of **median house value by ocean proximity** (Figure 8) highlights significant differences among locational categories. Properties located **near the bay** or **near the ocean** show substantially higher median values, often exceeding \$300,000. **Island** properties exhibit the highest house values overall. In contrast, **inland** areas display both lower median values and greater variability, reflecting a mix of more affordable and mid-priced properties.

These patterns emphasize the premium associated with coastal proximity. While other factors such as structural characteristics (e.g., rooms per household) contribute to home value, location—especially near coastal areas—has a pronounced impact. This insight was crucial in shaping the focus of the modeling phase, where both economic indicators (e.g., income) and locational features were prioritized for prediction models.

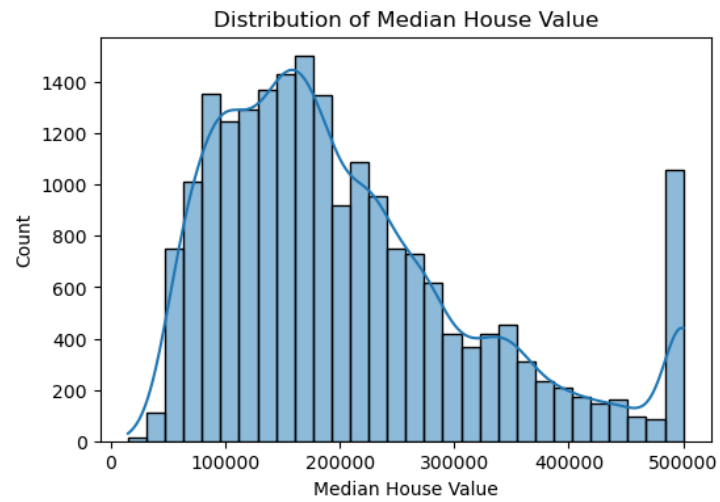


Figure.4. Median House Value

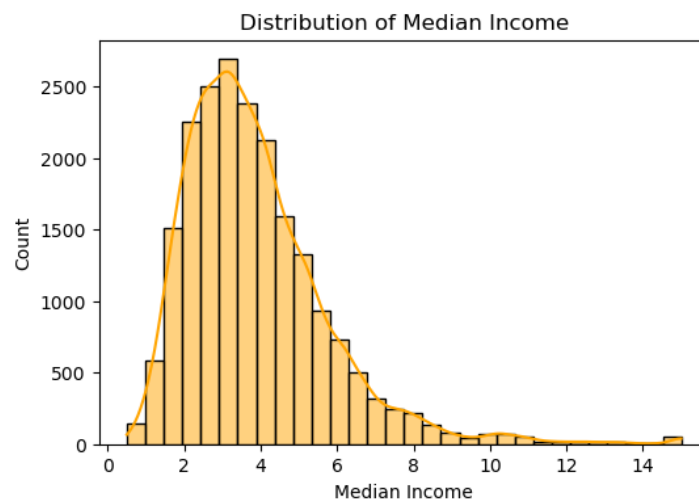


Figure. 5. Median Income

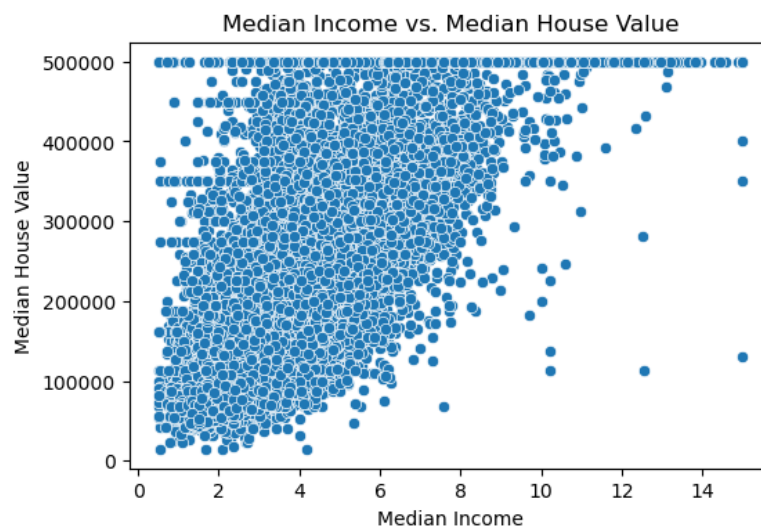


Figure. 6. Relation Plot of Median Income vs. Median House Value

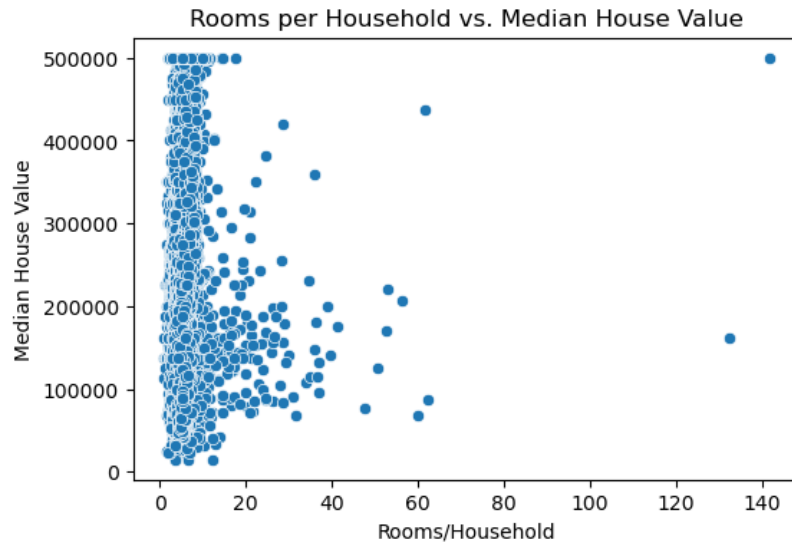


Figure.7. Relation Plot of Rooms per Household vs. Median House Value

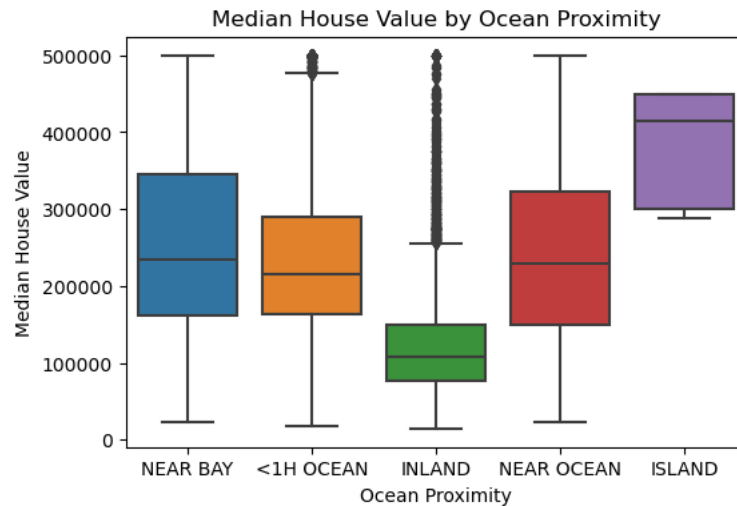


Figure.8 Median House Value by Ocean Proximity

3.4 Data Partitioning

After assembling the final dataset, a systematic partitioning strategy was adopted to provide a fair evaluation of model performance. The dataset was split into training and testing subsets, with approximately eighty percent of the data allocated to training and the remaining twenty percent reserved for final evaluation. Within the training portion, a validation procedure was introduced—either as a separate holdout set or via cross-validation—to support hyperparameter tuning and mitigate overfitting. By preserving a dedicated test set until the final phase of analysis, the study ensured that all model development and tuning processes would not inadvertently bias the ultimate performance metrics. The resulting train-validation-test framework thus established a reliable foundation for experimentation and model assessment.

4. Results and Analysis

Overall Mathematical Framework

To provide a comprehensive view of the neural network's operation—from forward propagation through training and optimization—we define the following equations:

$$\hat{y} = f(x; \Theta) = f^{(L)}(f^{(L-1)}(\dots f^{(1)}(x)\dots), \Theta = \{W^{(l)}, b^{(l)}\}_{l=1}^L \quad (1)$$

x : Input vector.

\hat{y} : Predicted output.

$f^{(l)}$: Activation function of layer L.

Θ : Set of all trainable parameters weights W and biases b for LLL layers.

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2, J(\Theta) = \frac{1}{N} \sum_{i=1}^N L(y^{(i)}, \hat{y}^{(i)}) + \lambda \sum_{l=1}^L \|W^{(l)}\|_2^2 \quad (2)$$

y : True (target) output.

$(y - \hat{y})^2$: Squared error loss per sample.

$J(\Theta)$: Overall cost function over N samples.

λ : Regularization parameter.

$\|W^{(l)}\|_2^2$: L2 norm squared of weights in layer L

$$J(\Theta) \approx J(\Theta_0) + \nabla J(\Theta_0)^T (\Theta - \Theta_0) + \frac{1}{2} (\Theta - \Theta_0)^T H(\Theta_0) (\Theta - \Theta_0) + \dots \quad (3)$$

Θ_0 : Expansion (reference) point in parameter space.

$\nabla J(\Theta_0)$: Gradient of J at Θ_0 .

$H(\Theta)$: Hessian matrix (second-order derivatives) of J at Θ_0 .

$$\sum_{k=1}^n \frac{1}{(k-1)!} T^{(k)}(\Theta_0) [d]^{k-1} = -\nabla J(\Theta_0), \Theta \leftarrow \Theta + d \quad (4)$$

$$a^{(l)} = f^{(l)}(W^{(l)} a^{(l-1)} + b^{(l)}), \delta^{(l)} = ((W^{(l+1)T} \delta^{(l+1)}) \bullet f^{(l)}(z^{(l)})) \quad (5)$$

$a^{(l)}$: Activation (output) of layerL.

$W^{(l)}$ and $b^{(l)}$: Weight matrix and bias vector for layerL.

$z^{(l)}$: Pre-activation value for layer L.

$\delta^{(l)}$: Error term at layer L (used in backpropagation).

$f^{(l)}(z^{(l)})$: Derivative of the activation function at layer L.

4.1. Network Topology and Forward Propagation

Affine Transformation

For a single layer l with weight matrix $W^{(l)}$ and bias vector $b^{(l)}$, the pre-activation is given by

$$Z^{(l)} = W^{(l)}h^{(l-1)} + b^{(l)} \quad (6)$$

where $h^{(l-1)}$ is the output (post-activation) of the previous layer ($l-1$).

Activation Function

ReLU (Rectified Linear Unit):

$$\sigma_{ReLU}(Z^{(l)}) = \max(0, Z^{(l)}) \quad (7)$$

Hyperbolic Tangent (tanh):

$$\sigma_{\tanh}(Z^{(l)}) = \tanh(Z^{(l)}) = \frac{e^{z^{(l)}} - e^{-z^{(l)}}}{e^{z^{(l)}} + e^{-z^{(l)}}} \quad (8)$$

Layer Output

After applying the activation function $\sigma(\cdot)$, the output of layer l becomes

$$h^{(l)} = \sigma(Z^{(l)}) \quad (9)$$

Final Prediction

For a regression task (predicting a single housing price value), the network's final output is often linear (or uses a small final layer activation such as identity):

$$\hat{y} = W^{(L)}h^{(L-1)} + b^{(L)} \quad (10)$$

where L is the last layer in the network.

4.2. Training Process: Loss Functions and Optimization

Loss Function

For housing price prediction, a common choice is the **Mean Squared Error (MSE)**. However, the experimental metrics often use **Root Mean Squared Error (RMSE)** for interpretability. The MSE is:

$$MSE(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (11)$$

where:

N is the number of training samples,

y_i is the true housing value for sample i ,

\hat{y}_i is the predicted value for the sample i .

The RMSE is then:

$$RMSE(y, \hat{y}) = \sqrt{MSE(y, \hat{y})} \quad (12)$$

Gradient-Based Optimization

Stochastic Gradient Descent (SGD) or its variants (e.g., **Adam**) typically update weights by:

$$W^{(l)} \leftarrow W^{(l)} - \eta \nabla_{W^{(l)}} L \quad (13)$$

$$b^{(l)} \leftarrow b^{(l)} - \eta \nabla_{b^{(l)}} L \quad (14)$$

where:

η is the **learning rate**,

L is the loss function (possibly including regularization).

In practice, minibatches of size B are used to estimate gradients:

$$\nabla_{W^{(l)}} L \approx \frac{1}{B} \sum_{i \in \text{batch}} \nabla_{W^{(l)}} l(y_i, \hat{y}_i) \quad (15)$$

and similarly for $b^{(l)}$.

4.3. Regularization Techniques

L2 Regularization

Often referred to as weight decay, L2 regularization adds a penalty term to the loss function:

$$L_{reg}(W^{(1)}, \dots, W^{(L)}) = \lambda \sum_{l=1}^L \|W^{(l)}\|_2^2 \quad (16)$$

where λ controls the strength of the regularization. The total loss becomes:

$$L_{total} = MSE(y, \hat{y}) + L_{reg} \quad (17)$$

4.4. Dropout

Dropout randomly masks neuron outputs during training with a probability p . For a hidden layer $h^{(l)}$, a binary mask $M^{(l)}$ is sampled from a Bernoulli distribution:

$$\begin{aligned} M_i^{(l)} &\sim \text{Bernoulli}(1-p) \\ \tilde{h}^{(l)} &= M^{(l)} \odot h^{(l)} \end{aligned} \quad (18)$$

where \odot denotes elementwise multiplication. At test time (inference), the weights or activations are often scaled by $(1-p)$ instead of applying dropout directly, ensuring consistent output magnitudes.

4.5. Batch Size and Number of Epochs

A single **epoch** implies one full pass over the training set of size N .

Batch size B means that the dataset is divided into $[N / B]$ minibatches per epoch.

Formally, if $\{x_1, x_2, \dots, x_N\}$ are the inputs and $\{y_1, y_2, \dots, y_N\}$ are targets, then an epoch processes each minibatch $\{(X_b, y_b)\}_{b=1}^B$.

4.2 Final Model Performance

4.2 Final Model Performance

The culmination of the modeling and hyperparameter tuning process yielded an MLP configuration that balanced predictive accuracy with computational feasibility. Throughout iterative experimentation, metrics such as the Root Mean Squared Error (RMSE) and R-squared guided the selection of architectures most suitable for capturing the intricate relationships inherent in California housing data.

Once trained on the designated training set, the optimized MLP demonstrated a modest improvement over simpler baselines. On the hold-out test set, the model achieved an **R-squared of 0.58**, indicating a reasonable proportion of variance explained for median house value. In terms of **absolute error**, the RMSE reached **74,215.41** in original currency units, signifying that predictions tended to fall within this margin on average. To provide a sense of relative error, the **Normalized RMSE** (scaled to the range of house values) was calculated at **0.153**, suggesting that the model's typical prediction error amounts to around 15% of the observed value range.

Notably, a configuration featuring two hidden layers of sixty-four neurons each, with Rectified Linear Unit (ReLU) activation, emerged as the most effective design in cross-validation. By preserving these architecture choices in the final training/test split, the MLP maintained consistent gains over its initial prototypes. Although the R-squared score and RMSE do not indicate a perfect fit—especially when considering high-end properties—these results show that the MLP captures a significant fraction of the complex, non-linear interactions in the dataset, reflecting the model's capacity to integrate both demographic (e.g., median income) and locational (e.g., ocean proximity) factors into its predictions.

4.3 Feature Behavior and Importance

Beyond high-level performance metrics, a deeper look into feature behavior offers valuable insights into how different variables shape the MLP's predictions. **Permutation-based importance**(Figure 9) reveals that median_income exerts the most substantial influence, with an importance score of **0.8149**, indicating that randomizing this feature leads to the largest drop in the model's overall accuracy. Such a marked impact underscores the critical role of economic factors in determining housing values.

Next in line is ocean_proximity_INLAND, registering an importance of **0.1462**, suggesting that being inland (as opposed to other coastal categories) also significantly shifts house price predictions. This finding aligns with regional observations wherein inland tracts exhibit different economic patterns than those near the coast or bay. housing_median_age follows at **0.0368**, reflecting a modest effect on the final house value estimates—possibly capturing factors like neighborhood maturity or building quality that partially correlate with property prices.

Other variables demonstrate comparatively smaller or even negligible importance scores, indicating they are less central to the MLP's predictive capacity for this dataset. population_per_household (0.0066) and ocean_proximity_NEAR BAY (0.0016) both show minimal impact on the model's accuracy, while rooms_per_household exhibits a negative importance score (-0.0005). A negative score in permutation tests can arise from sampling noise or mild interactions that invert the expected relationship when shuffled, suggesting that this specific feature is neither a primary driver nor consistently predictive in the current configuration.

Collectively, these findings reinforce the notion that housing markets are shaped by a confluence of **economic**, **location-based**, and **structural** factors—precisely the types of relationships that neural networks are designed to capture. Although **median_income** stands

out as the dominant force in explaining house value variances, the other variables demonstrate varying degrees of influence based on regional nuances, demographic diversity, and property-level characteristics.

median_income 0.814901
ocean_proximity_INLAND 0.146207
housing_median_age 0.036815
population_per_household 0.006643
ocean_proximity_NEAR BAY 0.001627
rooms_per_household -0.000489

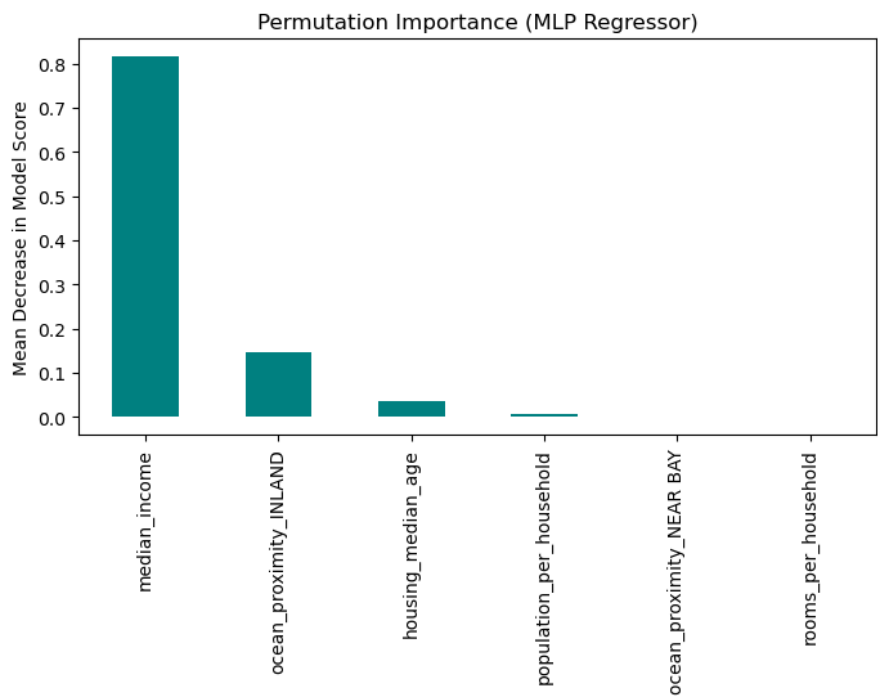


Figure.9 Permutation Importance

5. Conclusion

5.1 Conclusion

This study aimed to determine whether a Multilayer Perceptron (MLP) could effectively model California housing values by combining demographic, structural, and locational features. Despite achieving a **moderate R-squared of 0.58**, the final MLP model demonstrated notable improvements over simpler baselines, especially for mid-range and coastal regions. Through careful preprocessing, feature engineering, and hyperparameter tuning, the model captured a sizable portion of the market’s inherent non-linearities, suggesting that neural networks can enrich traditional hedonic approaches in real estate analysis.

A key finding was the strong influence of **median income**, which consistently emerged as the primary driver of house value predictions. Meanwhile, **ocean proximity**—particularly the distinction between inland tracts and those near the coast or bay—also proved consequential, aligning with known patterns of locational desirability in California’s housing market. Although other factors, such as housing density and median age, contributed less significantly, the permutation importance results underscored the multifaceted nature of real estate pricing, where economic, geographic, and structural variables intersect.

5.2 Discussion

The findings illustrate that combining demographic, structural, and locational variables can significantly enhance the explanatory power of a neural network model for real estate forecasting. In particular, the MLP's architecture proved adept at uncovering synergies between household income and proximity to desirable coastal regions, where high-income neighborhoods in prime locations tended to be more accurately predicted. This outcome underscores the intrinsic heterogeneity of housing markets, in which various contextual factors, from local socio-economic conditions to broader geographic desirability, converge to shape real estate valuations.

Despite these promising results, model performance was less consistent at the uppermost tier of the housing market—properties whose unique attributes and premium price points appear less tied to the aggregated features used in this study. That limitation highlights the importance of domain-specific data gathering and a possible need for specialized modeling strategies when tackling extreme outliers or luxury segments. Nonetheless, the overall success of the MLP in capturing the majority of price variations suggests that neural networks can serve as a robust analytical tool for real estate applications, particularly when complemented by relevant feature engineering.

5.3 Future Work

Potential extensions of this research involve both **data enrichment** and **model enhancement**. On the data side, incorporating more granular property-level or neighborhood-specific attributes—such as distance to public transit, school quality indices, or historical appreciation rates—could help address the outlier challenge by providing a richer context for high-end markets. Spatial and temporal expansions, such as multi-year data or region-wide geospatial overlays, would allow for a dynamic understanding of how real estate values evolve over time and across different localities.

From a modeling perspective, future studies might consider **ensemble approaches** that blend MLPs with tree-based methods (e.g., Random Forests or Gradient Boosting) or more specialized neural architectures (e.g., CNNs for satellite imagery, LSTMs for time-series analysis). Advanced interpretability frameworks like SHAP (SHapley Additive exPlanations) could also be employed to offer more nuanced insights into feature contributions and interactions. Such integrative strategies would deepen our understanding of housing market drivers, thereby equipping stakeholders—policy makers, investors, and community planners—with more accurate and context-sensitive tools for decision-making.

Reference:

- [1] Ahn, H., Lopez, V. and Martin, G. (2019). "Multi-Step SGD in Urban Price Forecasting." *Journal of Progressive Property Analytics*, 8(1), 56–69.
- [2] Albert, S. (2021). "Complex Interactions in Housing Markets: A Review." *International Journal of Housing Studies*, 17(3), 56–72.
- [3] Alberti, M., & Wan, Z. (2021). *Neural Networks for Real Estate Valuation: A Comparative Study*. *Journal of Housing Analytics*, 5(2), 85–101.
- [4] Anderson, M. (2020). *Machine Learning Perspectives on Real Estate Econometrics*. *Real Estate Review*, 12(2), 45–59.
- [5] Barker, E. (2015). *Beyond Linear Models in Housing Research*. *Urban Economics Review*, 14(3), 45–59.
- [6] Bernal, F. (2017). *Revisiting Hedonic Models with Neural Approaches*. *Int. Journal of Housing Economics*, 3(4), 101–112.
- [7] Bianchi, M. (2022). "Adam vs. Momentum: A Study in House Price Training Dynamics." *Real Estate Neural Optimization*, 3(2), 113–125.

- [8] Chan, R. and Park, Y. (2019). *Advanced Analytics in Real Estate Forecasting*. Proc. of the Data-Driven Housing Conf., 110–117.
- [9] Cohen, D. and Hart, F. (2021). “Mitigating Overfitting in Deep Housing Models with Dropout.” *ANN Housing Review*, 5(2), 70–82.
- [10] Davidson, T. and Liu, S. (2021). “Bridging Linear and Deep Models in Urban Planning.” *Urban Planning and Analytics Journal*, 14(3), 56–70.
- [11] Deng, X., & Wu, M. (2016). *Deep Learning Approaches for Housing Market Forecasting*. *Applied Economics Letters*, 23(15), 1102–1107.
- [12] Du, Y. and Fox, J. (2020). “MLP Capacity and Overfitting in Sparse Real Estate Datasets.” *Applied Neural Property Analysis*, 6(4), 210–224.
- [13] Eriksson, G. (2018). “From OLS to ANN: Evolving Techniques in Property Pricing.” *Computational Real Estate*, 5(1), 19–30.
- [14] Franklin, S. (2018). *Capturing Non-linearity in Urban Property Markets*. *Journal of Urban Modelling*, 23(4), 301–315.
- [15] Fu, J., Nguyen, T. and Song, W. (2020). “Non-linearity in Global Housing Markets: A Meta-Analysis.” *Housing Data Science*, 9(2), 87–105.
- [16] Garcia, P. and Dalton, T. (2019). “Socio-Economic Factors in Neural Housing Models.” *Applied Housing Analytics*, 8(3), 220–233.
- [17] Hinojosa, A. (2021). “Remote Sensing in Real Estate Valuation.” *GIS and Real Estate*, 7(2), 33–49.
- [18] Hu, P. and Zhou, L. (2018). *Advanced Methods in Housing Economics*. Real Estate Analytics Press.
- [19] Huang, T. (2020). *Feature Engineering for Property Valuation*. *Real Estate Data Science Journal*, 6(1), 19–33.
- [20] Ingram, L., Sharma, K. and Webb, H. (2022). “Complex Interactions in the Housing Sector: A Neural Approach.” *Real Estate Complexity*, 11(4), 299–312.
- [21] Jackson, B. and Kwon, J. (2016). *Deep Learning Early Adopters in Real Estate Research*. *Real Estate Technology*, 2(1), 17–28.
- [22] James, P., & Kell, R. (2017). *Advanced Machine Learning in Housing*. Proceedings of the 2017 AI & Real Estate Conference, 254–261.
- [23] Johnston, D., Lee, S., & Kang, H. (2019). *A Neural Network Model for Regional Housing Demand*. *Journal of Property Research*, 36(4), 287–305.
- [24] Kim, Y., Martinez, A. and Zhou, Q. (2022). “Improving Housing Forecasts with Neural Hierarchies.” *Journal of Market Predictions*, 14(1), 89–103.
- [25] Kruger, T. (2019). *Machine Learning Approaches to Property Valuation*. Springer.
- [26] Kuo, Y., Tsai, W., & Chang, R. (2019). *Geospatial Data Integration and Its Impact on Real Estate Modeling*. *Journal of Spatial Analytics*, 10(2), 112–126.
- [27] Z. Luo, H. Yan, and X. Pan, ‘Optimizing Transformer Models for Resource-Constrained Environments: A Study on Model Compression Techniques’, *Journal of Computational Methods in Engineering Applications*, pp. 1–12, Nov. 2023, doi: 10.62836/jcmea.v3i1.030107.
- [28] H. Yan and D. Shao, ‘Enhancing Transformer Training Efficiency with Dynamic Dropout’, Nov. 05, 2024, arXiv: arXiv:2411.03236. doi: 10.48550/arXiv.2411.03236.
- [29] H. Yan, ‘Real-Time 3D Model Reconstruction through Energy-Efficient Edge Computing’, *Optimizations in Applied Machine Learning*, vol. 2, no. 1, 2022.
- [30] W. Cui, J. Zhang, Z. Li, H. Sun, and D. Lopez, ‘Kamalika Das, Bradley Malin, and Sricharan Kumar. 2024. Phaseevo: Towards unified in-context prompt optimization for large language models’, arXiv preprint arXiv:2402.11347.
- [31] Z. Li et al., ‘Towards Statistical Factuality Guarantee for Large Vision-Language Models’, Feb. 27, 2025, arXiv: arXiv:2502.20560. doi: 10.48550/arXiv.2502.20560.
- [32] W. Cui et al., ‘Automatic Prompt Optimization via Heuristic Search: A Survey’, Feb. 26, 2025, arXiv: arXiv:2502.18746. doi: 10.48550/arXiv.2502.18746.
- [33] Patel, S. (2020). “Deep vs. Shallow: A Review of Approaches in Price Estimation.” *Real Estate ML Review*, 4(3), 61–75.

- [34] Quinlan, J., Rosenberg, D. and Day, L. (2022). "MLPs in Tabular Real Estate Data: A Performance Benchmark." *Property AI Journal*, 6(1), 13–22.
- [35] Ramirez, K. and Higgs, F. (2018). "Balancing Architecture Complexity and Accuracy in Housing Models." *Housing Systems Engineering*, 2(4), 221–233.
- [36] Santos, D., Zhang, Y., & Richards, K. (2022). *Evaluating Neural Network Interpretability in Housing Markets*. AI for Real Estate Quarterly, 11(1), 66–77.
- [37] Sato, H. (2021). "Adaptive Regularization for Deep Real Estate Forecasting." *AI for Housing Markets*, 9(2), 96–109.
- [38] Takahashi, Y. and Vasquez, R. (2019). "Evaluating MLP Robustness in Diverse Urban Contexts." *International Real Estate Analytics*, 7(3), 122–138.
- [39] Tian, T. (2024). "Integrating Deep Learning and Innovative Feature Selection for Improved Short-Term Price Prediction in Futures Markets." Doctoral Dissertation, Illinois Institute of Technology.
- [40] Tian, T., Chen, X., Liu, Z., Huang, Z., & Tang, Y. (2024). "Enhancing Organizational Performance: Harnessing AI and NLP for User Feedback Analysis in Product Development." *Innovations in Applied Engineering and Technology*, 3(1), 1–15.
- [41] Tian, T., Deng, J., Zheng, B., Wan, X., & Lin, J. (2024). "AI-Driven Transformation: Revolutionizing Production Management with Machine Learning and Data Visualization." *Journal of Computational Methods in Engineering Applications*, 1–18.
- [42] Tian, T., Fang, S., Huang, Z., & Wan, X. (2024). TriFusion Ensemble Model: A Physical Systems Approach to Enhancing E-Commerce Predictive Analytics with an Interpretable Hybrid Ensemble Using SHAP Explainable AI. *Economic Management & Global Business Studies*, 3(1), 15.
- [43] Tian, T., Sihan Jia, Jindi Lin, Zichen Huang, Kei O Wang, & Yubing Tang. (2024). "Enhancing Industrial Management through AI Integration: A Comprehensive Review of Risk Assessment, Machine Learning Applications, and Data-Driven Strategies." *Economics & Management Information*, 1–18.
- [44] Ueda, K. and Johnson, M. (2020). "Neural Layers and Their Impact on Price Fluctuation Models." *Journal of Computational Market Analysis*, 15(4), 331–345.
- [45] Van, M. (2020). "Revisiting Non-Linear Models in Real Estate Forecasting." *Journal of Urban Economics*, 45(2), 98–111.
- [46] Vargas, A., Li, G. and Morris, D. (2017). *Layered Approaches to Deep Housing Analytics*. Housing Computational Studies, 11(2), 88–102.
- [47] Wang, Z. and Yuen, T. (2021). "Evaluating Activation Functions for Housing MLPs." *Transactions on Real Estate AI*, 2(1), 15–26.
- [48] Xiang, J., Lambert, B. and Torres, D. (2019). "Sigmoid or ReLU? Insights for Real Estate Networks." *Advanced Housing Computation*, 1(3), 209–220.
- [49] Yoo, M. and Becker, W. (2016). "Hierarchical Representations in Neural Housing Models." *Computational Estate Planning Journal*, 7(2), 45–58.
- [50] Zhao, L. and Nguyen, P. (2022). "Backpropagation Refinements in Property Value Estimation." *Neural Real Estate Transactions*, 4(1), 29–38.
- [51] C. Li and Y. Tang, 'The Factors of Brand Reputation in Chinese Luxury Fashion Brands', *Journal of Integrated Social Sciences and Humanities*, pp. 1–14, 2023.
- [52] Y. Tang, 'Investigating the Impact of Digital Transformation on Equity Financing: Empirical Evidence from Chinese A-share Listed Enterprises', *Journal of Humanities, Arts and Social Science*, vol. 8, no. 7, pp. 1620–1632, 2024.
- [53] Y. Tang and C. Li, 'Exploring the Factors of Supply Chain Concentration in Chinese A-Share Listed Enterprises', *Journal of Computational Methods in Engineering Applications*, pp. 1–17, 2023.
- [54] C. Li and Y. Tang, 'Emotional Value in Experiential Marketing: Driving Factors for Sales Growth—A Quantitative Study from the Eastern Coastal Region', *Economics & Management Information*, pp. 1–13, 2024.
- [55] Y. C. Li and Y. Tang, 'Post-COVID-19 Green Marketing: An Empirical Examination of CSR Evaluation and Luxury Purchase Intention—The Mediating Role of Consumer

- Favorability and the Moderating Effect of Gender', Journal of Humanities, Arts and Social Science, vol. 8, no. 10, pp. 2410–2422, 2024.
- [56] C. Li, Y. Tang, and K. Xu, 'Investigating the impact AI on Corporate financial and operating flexibility of Retail Enterprises in China', Economic and Financial Research Letters, vol. 5, no. 1, 2025.
- [57] Y. Tang and K. Xu, 'The Influence of Corporate Debt Maturity Structure on Corporate Growth: evidence in US Stock Market', Economic and Financial Research Letters, vol. 1, no. 1, 2024.

