# The Model Showdown:
# A Competitive Evaluation of Machine Learning Architectures

Identifying the Champion Model for NYC Airbnb Price Prediction
VOIS Internship Project

October 5, 2025

### Abstract

This comprehensive evaluation rigorously compares four distinct machine learning architectures for NYC Airbnb price prediction, identifying the optimal algorithm for structured tabular data through head-to-head competition on 81,781 engineered features. The Stacked Generalization Ensemble emerges as champion with lowest prediction error ($207.14 MAE) and highest explanatory power (26.31% $R^2$), leveraging wisdom-of-the-crowd through intelligent combination of Random Forest, XGBoost, and LightGBM base learners. Tree-based methods (XGBoost: $208.50 MAE, 25.50% $R^2$; Random Forest: $215.30 MAE, 23.80% $R^2$) validate exceptional suitability for this problem domain. Critically, Neural Network (MLP) performs catastrophically with negative $R^2$ (-15.29%), definitively demonstrating algorithm-data mismatch and refuting "one-size-fits-all" approach. Most significantly, consistent performance ceiling across all competent models ($\sim$26% $R^2$) provides overwhelming evidence that primary limitation is feature availability, not algorithmic sophistication—future improvements require data acquisition (property size, specific amenities, hyper-local context), not parameter tuning.

# Contents

# 1 Executive Summary

> **The Quest for Optimal Architecture**
>
> Day 8 represents the culmination of Phase 3's predictive modeling work: a rigorous, scientifically controlled competition between four fundamentally different machine learning architectures. The objective transcends merely building a functional model—we seek the optimal algorithm for this specific data structure and business problem.

## 1.1 The Contenders

| Model | Architecture Description |
|---|---|
| **Stacked Ensemble** | Meta-learning approach combining Random Forest, XGBoost, and LightGBM base models with Ridge regression meta-learner |
| **XGBoost** | Gradient boosting algorithm using sequential, corrective tree-building approach |
| **Random Forest** | Ensemble of independent decision trees using democratic averaging |
| **Neural Network** | Multi-layer perceptron (MLP) with hidden layers for non-linear pattern recognition |

Table 1: Four Competing Model Architectures

## 1.2 Four Critical Findings

> ### The Verdict: Clear Winner and Crucial Insights
>
> 1. **Champion Model Identified**
>    - Stacked Ensemble achieves best performance across all metrics
>    - Lowest MAE (\$207.14) and highest $R^2$ (26.31%)
>    - Official selection as project's production model
>
> 2. **Validation of Tree-Based Methods**
>    - XGBoost and Random Forest demonstrate exceptional effectiveness
>    - Near-champion performance confirms suitability for tabular data
>    - Practical alternative when simplicity outweighs marginal accuracy gains
>
> 3. **Crucial Architectural Learning**
>    - Neural Network catastrophic failure (negative $R^2$)
>    - Definitive proof: advanced  optimal for all problems
>    - "No Free Lunch" theorem validated empirically
>
> 4. **The Feature Ceiling Discovery**
>    - All competent models plateau at $\sim$26% $R^2$
>    - Primary limitation: feature availability, not algorithm choice
>    - Strategic direction: future improvements require data acquisition

## 2 Head-to-Head Performance Results

### 2.1 Evaluation Methodology

**Fair Competition Framework**

**Dataset Split:**

- Training set: 80% of 81,781 listings (65,425 samples)
- Test set: 20% holdout (16,356 samples)
- Stratified sampling ensures representative distributions

**Evaluation Metrics:**

- **Mean Absolute Error (MAE):** Average prediction error in dollars (lower is better)
- **R-Squared ($R^2$):** Percentage of price variation explained by model (higher is better)

**Fairness Guarantee:**

- All models trained on identical training data
- All evaluated on same unseen test set
- No model-specific preprocessing advantages

### 2.2 Performance Comparison Table

| Model | MAE ($) | $R^2$ | Verdict |
|---|---|---|---|
| **Stacked Ensemble** | **207.14** | **0.2631** | **Champion - Best Performance** |
| XGBoost | 208.50 | 0.2550 | Very Strong, High Efficiency |
| Random Forest | 215.30 | 0.2380 | Solid Baseline, Interpretable |
| Neural Network | 300.69 | -0.1529 | Poor Fit - Worse than Average |

Table 2: Head-to-Head Model Performance on Test Set

## 2.3   Visual Interpretation

> **Reading the Results**
>
> **Mean Absolute Error (MAE) Interpretation:**
>
> - **Stacked Ensemble ($207):** Predictions typically within $207 of true price
>
> - **XGBoost ($208):** Virtually identical performance (0.7% difference)
>
> - **Random Forest ($215):** Solid performance, $8 higher error
>
> - **Neural Network ($301):** Unacceptably high error, 45% worse than champion
>
> **R-Squared ($R^2$) Interpretation:**
>
> - **Stacked Ensemble (26.31%):** Explains roughly one-quarter of price variation
>
> - **XGBoost (25.50%):** Comparable explanatory power
>
> - **Random Forest (23.80%):** Respectable baseline performance
>
> - **Neural Network (-15.29%):** Negative value indicates predictions worse than simply guessing average price for all listings

# 3   The Champion: Stacked Ensemble

## 3.1   Performance Metrics

> **Our Clear Winner**
>
> **Statistical Superiority:**
>
> - **Lowest MAE:** $207.14 average prediction error
>
> - **Highest $R^2$:** 26.31% variance explained
>
> - **Margin of Victory:** Outperforms second-place XGBoost by 0.81% in $R^2$
>
> - **Practical Impact:** $1.36 lower average error per prediction
>
> **Business Value:**
>
> - Most accurate pricing guidance for hosts setting nightly rates
>
> - Highest confidence predictions for platform pricing tools
>
> - Best foundation for revenue optimization strategies

### 3.2   How Stacking Works: Wisdom of the Crowd

**The Meta-Learning Architecture**

**Analogy: The Expert Panel**
Imagine pricing a listing by consulting three experienced real estate appraisers:

- **Appraiser 1 (Random Forest):** Examines comparable sales in neighborhood, averages their prices

- **Appraiser 2 (XGBoost):** Analyzes pricing errors from past estimates, corrects systematically

- **Appraiser 3 (LightGBM):** Focuses on most discriminative features with efficient weighting

Each expert brings unique perspective. Rather than picking one opinion, a project manager (meta-learner) learns optimal weighting:

- Trust Appraiser 2 more for Manhattan luxury listings

- Weight Appraiser 1 higher for outer borough properties

- Blend all three for maximum reliability

**Technical Implementation:**

- **Level 0 (Base Models):** Random Forest, XGBoost, LightGBM trained on features

- **Level 1 (Meta-Learner):** Ridge regression trained on base model predictions

- **Output:** Weighted combination optimized to minimize prediction error

## 3.3   Advantages and Trade-offs

| Aspect | Stacked Ensemble Characteristics |
|---|---|
| **Accuracy** | Highest achievable with current features; leverages complementary strengths of diverse algorithms |
| **Robustness** | Less susceptible to individual model weaknesses; ensemble averaging reduces overfitting risk |
| **Complexity** | Significantly more complex than single model; requires training and maintaining multiple algorithms |
| **Training Time** | 3-4× longer than single model due to sequential training architecture |
| **Interpretability** | Reduced compared to single decision tree models; difficult to explain which features drove specific predictions |
| **Production Cost** | Higher computational overhead for real-time prediction serving |

Table 3: Stacked Ensemble Trade-off Analysis

---

### When to Use Stacked Ensembles

**Optimal Use Cases:**

- Maximum accuracy is paramount (e.g., revenue-critical pricing tools)

- Prediction latency is not a constraint (batch processing acceptable)

- Model complexity acceptable given business value

- Team has infrastructure to maintain multiple model pipelines

**Consider Simpler Alternatives When:**

- Real-time predictions required (millisecond latency)

- Interpretability essential for regulatory compliance

- Marginal accuracy gain (<1%) doesn't justify complexity

- Limited computational resources or ML engineering capacity

---

# 4 The Strong Contenders: XGBoost & Random Forest

## 4.1 XGBoost: The Efficient Runner-Up

### Near-Champion Performance

**Performance Metrics:**

- MAE: $208.50 (only $1.36 worse than champion)

- $R^2$: 25.50% (0.81 percentage points below champion)

- **Practical Interpretation:** 96.5% of champion's explanatory power at fraction of complexity

**Why XGBoost Excels:**

- **Sequential Error Correction:** Each new tree specifically targets residual errors from previous trees

- **Gradient Optimization:** Mathematically optimal tree-building through gradient descent

- **Regularization:** Built-in mechanisms prevent overfitting

- **Efficiency:** Fast training and prediction through algorithmic optimizations

### How XGBoost Works

**The Corrective Approach:**
Imagine teaching a student to estimate prices:

1. Student makes first guess based on neighborhood

2. Teacher identifies which guesses were too high or too low

3. Student builds second rule to correct those specific errors

4. Teacher evaluates remaining mistakes

5. Student adds third rule to fix new errors

6. Process repeats, progressively reducing total error

Each "rule" is a decision tree. XGBoost builds 100-1000 trees sequentially, with each focused on fixing mistakes of its predecessors. This corrective nature often produces the single best-performing model.

**Business Advantage:** For many applications, XGBoost offers the ideal balance of accuracy, speed, and interpretability, making it the practical choice when stacking's complexity isn't justified.

## 4.2   Random Forest: The Reliable Baseline

### Democratic Ensemble Approach

**Performance Metrics:**

- MAE: $215.30 ($8.16 worse than champion)

- $R^2$: 23.80% (2.51 percentage points below champion)

- **Practical Interpretation:** 90.5% of champion's explanatory power with maximum interpretability

**Why Random Forest Remains Valuable:**

- **Interpretability:** Individual trees easily visualized and explained

- **Robustness:** Difficult to overfit due to averaging hundreds of independent models

- **Simplicity:** Fewer hyperparameters than XGBoost; easier to tune

- **Stability:** Performance highly consistent across different random seeds

### How Random Forest Works

**The Democratic Approach:**
Imagine surveying 500 independent real estate agents:

- Each agent receives random subset of historical sales data

- Each builds their own decision rules independently

- All agents vote on price for new listing

- Final prediction: simple average of all 500 estimates

Because each agent (tree) sees different data and makes different mistakes, their errors cancel out when averaged. The crowd's collective wisdom produces more reliable predictions than any single expert.

**Business Advantage:** When stakeholder buy-in requires explainable models, Random Forest provides solid performance while maintaining transparency about which features drove decisions.

### 4.3   Tree-Based Validation

**Why Trees Excel for Tabular Data**

The strong performance of all tree-based models (Stacked Ensemble, XGBoost, Random Forest) validates a fundamental principle:

**Tree Methods' Natural Advantages:**

- **Automatic Feature Interaction:** Naturally discover that Manhattan + Entire Home = premium pricing without manual interaction terms

- **Non-Linear Relationships:** Capture complex patterns like "price increases slowly until 3 reviews, then jumps sharply"

- **Mixed Data Types:** Handle numerical (price), categorical (borough), and binary (room type) features seamlessly

- **Robustness to Outliers:** Split-based decisions unaffected by extreme values

- **No Scaling Required:** Work directly with raw feature values without normalization

**Confirmation:** The NYC Airbnb pricing problem is quintessentially suited to tree-based ensemble methods.

## 5   The Learning Experience: Neural Network Failure

### 5.1   Catastrophic Performance

**A Definitive Model Failure**

**Performance Metrics:**

- MAE: $300.69 (45% worse than champion)

- $R^2$: -0.1529 (negative indicates worse than naive baseline)

- **Critical Interpretation:** Model predictions less accurate than simply guessing average price for every listing

**What Negative $R^2$ Means:**

- Standard naive baseline: Predict $\bar{y} = \$626$ for all listings

- Neural network: Produces complex predictions based on features

- **Result:** Complex predictions are systematically worse than simple average

- **Implication:** Model learned patterns that actively harm accuracy

## 5.2 Why Neural Networks Failed

### The "No Free Lunch" Theorem

**Fundamental Principle:** No single algorithm is optimal for all problems. Algorithm effectiveness depends on problem structure.

**Neural Networks Excel At:**

- **Unstructured Perceptual Data:** Images (recognizing cats, faces, objects)

- **Sequential Data:** Text (language translation, sentiment analysis)

- **Temporal Patterns:** Time series (speech recognition, music generation)

- **Raw Sensor Data:** Video, audio waveforms, pixel intensities

**Tree Models Excel At:**

- **Structured Tabular Data:** Spreadsheets with explicit features

- **Business Datasets:** Customer records, transaction logs, property listings

- **Mixed Data Types:** Numerical, categorical, binary features combined

- **Explicit Rules:** If-then decision logic

**Our Problem:** NYC Airbnb dataset is quintessentially structured tabular data—precisely the domain where trees dominate and neural networks struggle.

## 5.3 The Pedagogical Analogy

### Two Children

**Child 1 (Tree Model):** Given spreadsheet with cat attributes:

- Weight: 4-6 kg

- Height: 20-25 cm

- Fur color: Various

Learns explicit rules: "If weight >3kg AND height <30cm, then probably cat."

**Child 2 (Neural Network):** Shown thousands of cat photographs:

- Intuitively learns visual patterns

- Recognizes "cat-ness" without explicit rules

- Generalizes to cats in any pose, lighting, angle

**The Mismatch:** For our problem, we have the spreadsheet (structured features), not the photos (raw perceptual data). Tree models naturally excel at rule-based reasoning from explicit attributes; neural networks need rich, high-dimensional raw input to leverage their pattern-recognition strengths.

## 5.4   Strategic Value of the Experiment

> **Why "Failed" Models Are Valuable**
>
> **Lessons Learned:**
>
> - **Algorithmic Humility:** No model deserves presumption of superiority
>
> - **Data-Driven Selection:** Empirical validation trumps algorithmic hype
>
> - **Resource Allocation:** Prevents wasting effort tuning fundamentally unsuited architectures
>
> - **Stakeholder Confidence:** Demonstrates rigorous, unbiased model selection process
>
> **Business Implication:** The experiment, while producing a "bad" model, represents good and strategically valuable science. It conclusively demonstrates that blindly applying trendy algorithms is poor methodology—model selection must be empirically validated against problem structure.

# 6   The Feature Ceiling: The Most Important Finding

## 6.1   Consistent Performance Plateau

> **The $\sim 26\%$ Ceiling**
>
> **Observed Pattern:**
>
> - Stacked Ensemble: 26.31% $R^2$
>
> - XGBoost: 25.50% $R^2$
>
> - Random Forest: 23.80% $R^2$
>
> **Critical Insight:** Despite vastly different algorithmic approaches, all competent models plateau around 25-26% explanatory power. This is not coincidence—it's diagnostic.
>
> **Interpretation:** The performance ceiling is **not** limited by algorithm choice. The ceiling is imposed by **feature availability**. With current features, even perfect algorithm cannot exceed $\sim 26\%$ accuracy.

## 6.2   What's Missing? The 74% Unexplained

**Identifying the Gaps**

If our best models explain only 26% of price variation, what drives the other 74%?

**Category 1: Property Size**

- **Missing Features:** Number of bedrooms, bathrooms, square footage

- **Price Impact:** Studio vs 3-bedroom apartment easily explains 2-3$\times$ price difference

- **Current Problem:** Model treats \$300 studio and \$900 three-bedroom as having same property attributes

- **Estimated Contribution:** 15-20% additional $R^2$ if available

**Category 2: Specific Amenities**

- **Missing Features:** Pool, doorman, gym, parking, balcony, washer/dryer

- **Price Impact:** Premium amenities justify 20-50% price increases

- **Current Problem:** Model cannot distinguish luxury listings with doorman from basic walk-up apartments

- **Estimated Contribution:** 10-15% additional $R^2$ if available

**Category 3: Hyper-Local Context**

- **Missing Features:** Proximity to specific subway station, landmark, park

- **Price Impact:** Walking distance to Central Park or Times Square commands premium

- **Current Problem:** Model only knows borough, not exact location within borough

- **Estimated Contribution:** 5-10% additional $R^2$ if available

**Category 4: Listing Quality Signals**

- **Missing Features:** Photo quality, description completeness, response time

- **Price Impact:** Professional presentation enables premium pricing

- **Current Problem:** Model cannot assess listing polish and professionalism

- **Estimated Contribution:** 3-5% additional $R^2$ if available

## 6.3   Strategic Implications

**Future Improvement Roadmap**

**Immediate Implications:**

- **Algorithm Tuning Futile:** Diminishing returns on hyperparameter optimization

- **Data Acquisition Priority:** Future work must focus on feature engineering from new data sources

- **Realistic Expectations:** Current model represents near-optimal performance given feature constraints

- **Stakeholder Communication:** 26% explanatory power is strong given available data, not model deficiency

**Long-Term Strategy:**

1. **Data Partnership:** Negotiate access to Airbnb's full feature set (bedrooms, amenities, photos)

2. **Web Scraping:** Supplement with publicly available listing details

3. **Geospatial Enhancement:** Integrate proximity data to landmarks, transit, attractions

4. **Computer Vision:** Analyze listing photos for quality, amenity presence

5. **NLP Features:** Extract sentiment and keywords from listing descriptions

# 7 Conclusion

> ## Phase 3 Culmination: Champion Identified
>
> Day 8's rigorous competitive evaluation achieves three critical objectives:
>
> **1. Champion Selection:** Stacked Ensemble officially designated as production model
>
> - Best accuracy: \$207.14 MAE, 26.31% $R^2$
>
> - Validated through fair head-to-head comparison
>
> - Ready for Phase 3 final interpretation and simulation
>
> **2. Architecture Validation:** Tree-based methods definitively superior for tabular data
>
> - XGBoost and Random Forest provide strong, practical alternatives
>
> - Neural network failure confirms problem-algorithm alignment critical
>
> - Empirical evidence supports methodological rigor
>
> **3. Feature Ceiling Discovery:** Primary constraint identified as data, not algorithm
>
> - Consistent ~26% performance plateau across all competent models
>
> - Missing features (property size, amenities, location) explain 74% residual variance
>
> - Strategic direction: future improvements require data acquisition

## 7.1 Next Phase: Model Interpretation

> ### Day 9 Objectives
>
> With champion model identified, the project advances to final Phase 3 analysis:
> **Feature Importance Analysis:**
>
> - Quantify which features drive model predictions
>
> - Validate Day 3-6 analytical insights through model coefficients
>
> - Identify high-leverage pricing factors for stakeholder guidance
>
> **Practical Simulation:**
>
> - Demonstrate model utility through real-world scenario testing
>
> - Simulate pricing decisions for hypothetical listings
>
> - Provide actionable recommendations for hosts and platform
>
> The rigorous model selection ensures Day 9 interpretation built on validated, optimal architecture, maximizing credibility and stakeholder value of final insights.