

	PREDICTIONS	OBJECTIVES	DATA
3RO	<div><div>End-user</div><div>Who will use the predictive system / who will be affected by it?</div><div>Potential Homebuyers, Property Sellers, Real Estate Agents and Brokers, Investors, HDB and Policy Makers, Marketing and Sales Teams</div></div>	<div><div>Value proposition</div><div>What are we trying to do for the system's users? (e.g. spend less time on X, increase Y...)</div><div>Provide Accurate Price Estimates Save Time and Effort Enhance Market Understanding Improve Market Efficiency</div></div>	<div><div>Data sources</div><div>Where do/can we get data from? (internal database, 3rd party API, etc.)</div><div>Data.gov.sg Google Maps API / OneMap API HDB Official Websites and Publications Real Estate Portals</div></div>
E SF	<div><div>Problem</div><div>Question to predict answers to (on behalf of user)</div><div>What will be the resale price of a specific HDB flat given its characteristics and market conditions?</div><div>Input (i.e. question "parameter")</div><div>Flat Characteristics, Location Data, Temporal Data</div><div>Possible outputs (i.e. "answers")</div><div>Predicted Resale Price</div><div>Type of problem (e.g. classification, regression, recommendation...)</div><div>Regression</div><div>Baseline: simple, alternative way of making predictions (e.g. manual rules)</div><div>Comparative Market Analysis (CMA) Manual Rules</div></div>	<div><div>Performance evaluation</div><div>Domain-specific / bottom-line metrics for monitoring performance in production</div><div>User Satisfaction Transaction Volume Revenue Impact</div><div>Prediction accuracy metrics (e.g. MSE if regression; % accuracy, #FP for classification)</div><div>Root Mean Squared Error (RMSE) Mean Absolute Error (MAE) R² (Coefficient of Determination)</div><div>Offline performance evaluation method (e.g. cross-validation or simple training/test split)</div><div>Cross-Validation Training/Test Split</div></div>	<div><div>Data preparation</div><div>How do we get training data (inputs, and outputs if supervised learning)? How many data points?</div><div>Data Source: The training data is obtained from the comprehensive HDB resale transaction records available on data.gov.sg. This dataset includes historical resale prices and detailed flat characteristics from January 2017 to June 2024.</div><div>Data Points: The dataset contains approximately 181,000 rows, each representing a unique resale transaction.</div><div>Supervised Learning: This project is based on supervised learning, where the inputs are the features extracted from the dataset (such as flat type, floor area, location, and lease details), and the output is the resale price of the HDB flats.</div><div>Input features (extracted from data sources). If too many, list types of features and mention key ones.</div><div>Feature Types:<ul style="list-style-type: none"><li>Numerical Features: Quantitative values such as floor area in square meters, lease commencement date, remaining lease duration, and resale price.</li><li>Categorical Features: Qualitative values that describe the flat and its location, such as town, flat type, flat model, storey range, block number, and street name.</li><li>Temporal Features: Time-related data such as the transaction month.</li></ul></div><div>Key Input Features:<ul style="list-style-type: none"><li>Flat Type: Classification of the flat based on the number of rooms (e.g., 3-room, 4-room).</li><li>Floor Area (sqm): Total interior space of the flat in square meters.</li><li>Flat Model: Design type of the flat (e.g., Improved, Model A).</li><li>Town: The locality or district where the flat is located.</li><li>Lease Commencement Date: The year the flat's lease started.</li><li>Remaining Lease: Duration left on the lease, usually expressed in years and months.</li><li>Storey Range: The range of floors where the flat is situated (e.g., 04 TO 06).</li><li>Transaction Month: The month and year of the transaction, providing temporal context.</li></ul></div></div>
DAT	<div><div>Using predictions</div><div>When do we make predictions and how many?</div><div><ul style="list-style-type: none"><li>Prediction Timing: Predictions are made whenever a user inputs the characteristics of an HDB flat into the system. This could occur in real-time scenarios, such as during a property search by a potential buyer or a valuation request by a seller.</li><li>Batch predictions might also be run periodically to provide market overviews or updates to stakeholders.</li><li>Number of Predictions: The system could handle thousands of predictions daily, depending on user demand and batch processing needs.</li><li>Real-time predictions are typically single or few per request, while batch predictions can cover all available listings or historical data points.</li></ul></div><div>What is the time constraint for making those predictions?</div><div>Real-Time Predictions: These should be made within a few seconds to provide a seamless user experience, particularly when users are actively exploring or evaluating properties. Batch Predictions:</div></div>		<div><div>Learning models</div><div>When do we create/update models? With which data / how much?</div><div>Model Creation and Update Frequency:<ul style="list-style-type: none"><li>Initial Model Creation: A baseline model is created using the entire historical dataset (from January 2017 to June 2024).</li><li>Periodic Updates: The model should be updated periodically, for instance, monthly or quarterly, to incorporate new data and adapt to changes in the market.</li><li>Event-Driven Updates: Models may also be updated in response to significant market changes or policy updates affecting HDB resale prices.</li></ul></div><div>Data for Updates:<ul style="list-style-type: none"><li>Each update should use the most recent data available, ensuring that the model reflects the latest trends and conditions.</li><li>Typically, all available data (including new and historical) should be used to maintain a comprehensive understanding of the market.</li></ul></div><div>What is the time constraint for creating a model?</div><div>The initial training process may take several hours to days, depending on the data size and complexity of the model. Subse</div></div>

How do we use predictions and confidence values?

Using Predictions:

- Predictions provide estimated resale prices for HDB flats, which can be used by buyers, sellers, and agents to guide decision-making.
- These estimates help in setting competitive prices, evaluating investment opportunities, and understanding market trends.

Using Confidence Values:

- Confidence values, or prediction intervals, give users a range within which the actual resale price is likely to fall.
- These values enhance user trust by providing transparency about the uncertainty or reliability of the predictions.
- They can be used to assess the risk associated with a given price estimate and to make more informed decisions.

Criteria for deploying model (e.g. minimum performance value — absolute, relative to baseline or to previous model)

- Minimum Performance Threshold: The model should meet or exceed a predefined performance threshold, such as an RMSE less than SGD 30,000 or an  $R^2$  greater than 0.80 on the test data.
- Improvement Over Baseline: The new model should perform better than the baseline or previous models, as measured by key metrics like RMSE, MAE, and  $R^2$ .
- Stability and Generalization: The model must demonstrate stability and the ability to generalize well across different subsets of the data, as verified through cross-validation or other evaluation methods.

Reset Form