

Sure, here's a table outlining different combinations of modeling approaches, increasing in complexity up to neural networks (which we'll exclude). Each combination involves data enhancement, feature engineering, and various modeling techniques to improve prediction accuracy.

Complexity Level	Data Enhancement	Feature Engineering	Clustering	Modeling Algorithm
Level 1	Basic Cleaning	None	None	Linear Regression
Level 2	Basic Cleaning	Basic Features	None	Linear Regression
Level 3	Basic Cleaning	Basic Features	None	Ridge/Lasso Regression
Level 4	Basic Cleaning	Polynomial Features	None	Polynomial Regression
Level 5	Scaling/Normalization	Polynomial Features	None	Support Vector Regression (SVR)
Level 6	Scaling/Normalization	Polynomial Features	None	Decision Tree Regression
Level 7	Scaling/Normalization	Advanced Features (e.g., interactions)	None	Random Forest Regression
Level 8	Scaling/Normalization	Advanced Features	None	Gradient Boosting (e.g., XGBoost)
Level 9	Scaling/Normalization	Advanced Features	Yes	Cluster-Based Regression
Level 10	Outlier Treatment	Advanced Features	Yes	Cluster-Based Ensemble
Level 11	Feature Selection/Dimensionality Reduction	Advanced Features	Yes	Stacked Models

Level 12	Advanced Preprocessing (e.g., PCA)	Advanced Features and External Data	Yes	Stacked Ensembles with SVM

Key Components Explained:

- **Data Enhancement/Preprocessing:**
 - **Basic Cleaning:** Handling missing values, correcting data types, basic imputation.
 - **Scaling/Normalization:** Standardizing features for algorithms sensitive to feature scales.
 - **Outlier Treatment:** Removing or capping extreme values to reduce skewness.
 - **Feature Selection/Dimensionality Reduction:** Using techniques like PCA to reduce feature space.
- **Feature Engineering:**
 - **Basic Features:** Creating simple features like property age (2024 - Year Built), total rooms (bedrooms + bathrooms).
 - **Polynomial Features:** Generating interactions and higher-degree terms to capture nonlinearities.
 - **Advanced Features:** Incorporating domain knowledge, ratios (e.g., Lot Size per Room), and external data (e.g., crime rates).
- **Clustering:**
 - Segmenting data using algorithms like K-Means to capture distinct market segments (e.g., luxury vs. affordable homes).
- **Modeling Algorithms:**
 - **Linear Regression:** Establishes a baseline with minimal complexity.
 - **Ridge/Lasso Regression:** Adds regularization to handle multicollinearity.
 - **Polynomial Regression:** Fits nonlinear data by adding polynomial terms.
 - **Support Vector Regression (SVR):** Uses kernel functions to model complex relationships.
 - **Decision Tree Regression:** Splits data based on feature values, capturing nonlinear patterns.
 - **Random Forest Regression:** An ensemble of decision trees to improve generalization.
 - **Gradient Boosting Machines:** Sequentially builds models to correct errors of prior models.
 - **Stacked Models:** Combines various models to leverage their individual strengths.

Progression Explained:

- **Levels 1-3:** Start with basic data cleaning and simple models to establish a baseline.
- **Levels 4-6:** Introduce more complex features and algorithms to capture nonlinear relationships.
- **Levels 7-8:** Apply ensemble methods to improve prediction accuracy and reduce overfitting.
- **Levels 9-10:** Incorporate clustering to handle distinct market segments with specialized models.
- **Levels 11-12:** Utilize advanced preprocessing and combine multiple models for optimal performance.

Implementation Tips:

- **Experimentation:** Test different combinations to see which yield the best validation performance.
- **Cross-Validation:** Use techniques like K-Fold cross-validation to assess model generalization.
- **Hyperparameter Tuning:** Optimize model parameters using Grid Search or Bayesian Optimization.
- **Evaluation Metrics:** Besides RMSE, consider MAE and R^2 to evaluate models comprehensively.

Example Workflow for Level 8:

1. Data Preprocessing:

- Handle missing values.
- Encode categorical variables using One-Hot Encoding.
- Scale features using StandardScaler.

2. Feature Engineering:

- Create `Age = 2024 - Year Built`.
- Compute `Total Rooms = bedrooms + bathrooms`.
- Generate interaction terms or polynomial features if relevant.

3. Modeling:

- Use XGBoost Regressor with tuned hyperparameters.
- Fit the model on training data.

```
import xgboost as xgb
model = xgb.XGBRegressor(n_estimators=500, learning_rate=0.05)
model.fit(X_train, y_train)
```

4. Evaluation:

- Predict on the test set.
- Calculate RMSE and compare with previous models.

By exploring these combinations, you can incrementally improve your model's performance while adhering to industry standards. This approach allows you to understand the impact of each technique and build a robust predictive model for house prices.

Let me know if you need details on implementing any specific combination.