**Machine Learning for Business Analytics**

**Prem Kumar Chimakurthi**

**Master of Science in Business Analytics**

**Table of Contents:**

## Week 4 | Predicting Boston Housing Prices

### Why should the data be partitioned into training, validation, and holdout sets?

The data needs to be divided so that the model is trained with one part of data and validated and tested with other parts. This helps in Training: This set of data is used in training the model to make it learn about the data distribution.

Validation set: This set of data is applied to tune the model's hyperparameters and to prevent overfitting by checking the performance of the model on data that it has not been trained with.

Holdout set (test set): This set is used to get an unbiased measure of the generalization of the performance of the final model after training and validation. It tells how well the model will perform on data that it has not seen before.

### What will the training set be used for?

The machine learning model is trained with the training set. The model learns relationships between the predictors (features) and the target variable (MEDV) where it minimizes error on this data set. Most of the data is in the training set so it is essential for the model to learn the patterns.

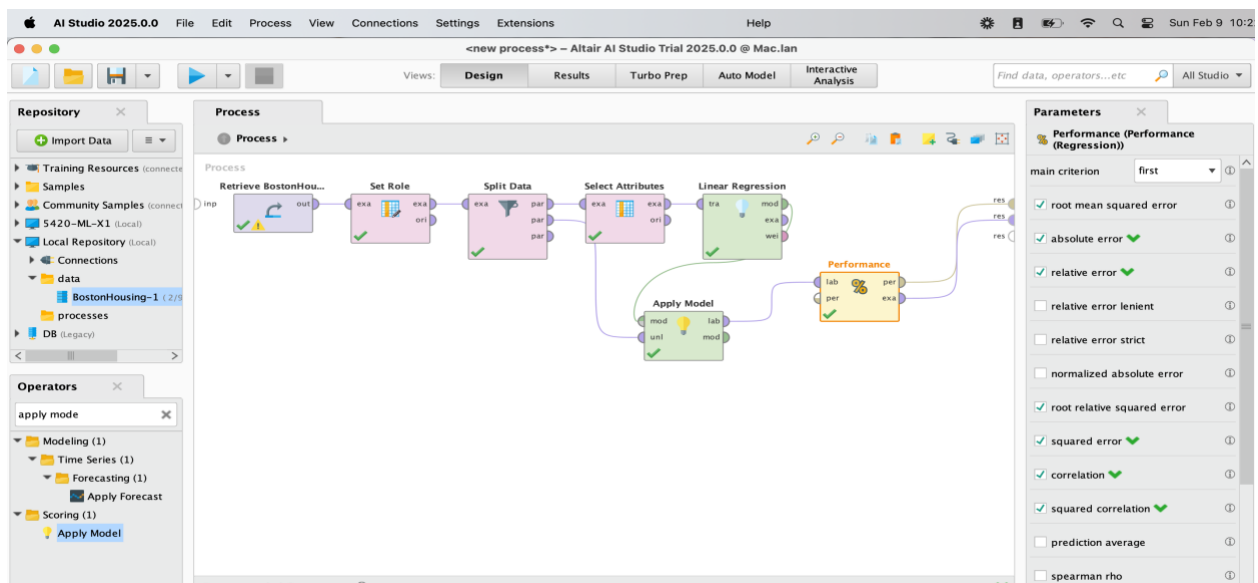### What will the validation and holdout sets be used for?

Validation set: The validation set is employed for tuning the models' hyperparameters and making evaluations during the training process. It assists in choosing the optimal model

complexity and overfitting prevention through reporting the model's performance on data that is not used for training.

Holdout set (test set): The holdout set is deployed to evaluate the final model's performance after training and validation. It gives an unbiased idea of how well the model would perform on data that it has never seen before. The set is tossed once after the model is thoroughly trained and validated so that the evaluation is completely unbiased and not contaminated by the training or validation set.

**Partition the data into training/validation/holdout with proportions 60:25:15. Fit a multiple linear regression model to the median house price (MEDV) as a function of CRIM, CHAS, and RM.**

**Partitioning the data:**

**Linear Regression Model performance is summarized below:**

- Root Mean Squared Error (RMSE), 6.299, Average prediction error is ~$6,299 (since MEDV is in $1000s). Lower is better.

- Absolute Error (MAE), 4.333, On average, predictions deviate by ~$4,333 from actual prices.

- Relative Error (%), 26.07%, Model's errors are ~26% of actual values, which is decent but could be improved.

- Root Relative Squared Error (RRSE), 0.705, Below 1, which means the model is better than simply predicting the mean.

- Correlation (R), 0.735, Moderate correlation between actual and predicted values.

- R² (Squared Correlation), 0.540, The model explains 54% of the variance in house prices (MEDV).

## Key Observations:

- R² = 0.54: Your model explains 54% of house price variation, meaning there's room for improvement.

- RMSE = 6.299: On average, the model's predictions are off by ~$6,300.

- MAE = 4.333: On average, the absolute error is ~$4,300, which is reasonable.

- Relative Error = 26.07%: The error is about one-fourth of the true house prices.

Write the equation for predicting the median house price from the predictors in the model.

## Formula for the equation:

$$MEDV = \beta_0 + \beta_1 (CRIM) + \beta_2 (CHAS) + \beta_3 (RM) + \epsilon$$

Where:

**$\beta_0$ = Intercept (constant term)**

**$\beta_1$, $\beta_2$, $\beta_3$ = Coefficients for each predictor variable**

**$\epsilon$ = Error term (represents unexplained variance)**

$$MEDV = -31.037 - 0.288(CRIM) + 6.597(CHAS) + 8.639(RM)$$

**Where:**

- Intercept = -31.037 (baseline house price when all predictors are zero)

- CRIM coefficient = -0.288 (higher crime rate slightly decreases house price)

- CHAS coefficient = 6.597 (houses near the Charles River increase in price)

- RM coefficient = 8.639 (each additional room increases house price significantly)

Using the estimated regression model, what median house price is predicted for a tract in the Boston area that does not bound the Charles River, has a crime rate of 0.1, and where the average number of rooms per house is 6?

**To predict the median house price (MEDV) using the estimated regression model, we use the given equation:**

$$MEDV = -31.037 - 0.288(CRIM) + 6.597(CHAS) + 8.639(RM)$$

**Given Data for Prediction:**

- **CRIM = 0.1 (Crime rate per capita)**

- **CHAS = 0 (Does not bound the Charles River)**

- **RM = 6 (Average number of rooms)**

$$MEDV = -31.037 - 0.288(0.1) + 6.597(0) + 8.639(6)$$

$$MEDV = -31.037 - 0.0288 + 0 + 51.834$$

$$MEDV = 20.768$$

**The predicted median house price is $20,768 (or $20.77K in $1000s)**

---

**Which predictors are likely to be measuring the same thing among the 13 predictors? Discuss the relationships among INDUS, NOX, and TAX.**

---

**Identifying Predictors That Measure Similar Factors:**

**Among the 13 predictors in the Boston Housing dataset, some variables are measuring similar underlying factors due to high correlation. Typically, in real estate and urban planning, some features are naturally related.**

- **INDUS (Industrial Proportion): Measures the proportion of land zoned for non-retail business (industrial use).**

- **NOX (Nitric Oxide Concentration): Indicates air pollution levels.**

- **TAX (Property Tax Rate per $10,000): Represents property tax per town.**

**These three variables are likely to be correlated because they are all related to urbanization and industrial activity.**

**Relationship Between INDUS, NOX, and TAX:**

**These three features are linked because industrialization often influences pollution levels and property taxation. Let us analyze their relationships:**

**INDUS and NOX (Industrialization & Pollution):**

- **More industrial land usage (higher INDUS) typically leads to higher NOX levels due to increased emissions from factories and businesses.**

- **In urban areas with high industrialization, air quality tends to worsen.**

**INDUS and TAX (Industrialization & Property Tax):**

- **Higher industrial land usage (higher INDUS) is often associated with higher property tax (TAX).**

- **Towns with more industrial zones might impose higher taxes to support infrastructure, services, and pollution control.**
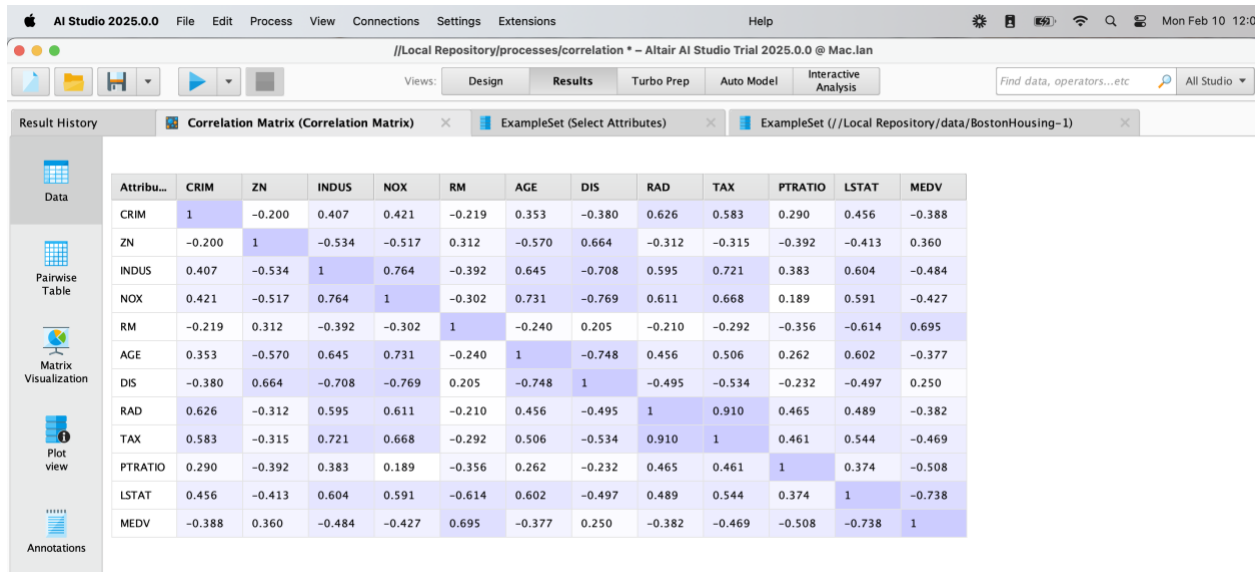
**NOX and TAX (Pollution & Property Tax):**

- **Areas with high pollution (higher NOX) tend to have higher taxes because governments might impose environmental regulations or higher property tax rates to offset environmental damages.**

- **Additionally, wealthier areas with lower pollution may have lower tax rates due to better zoning and urban planning.**

---

**Compute the correlation table for the 12 numerical predictors using the Correlation Matrix operator in RapidMiner, and search for highly correlated pairs. These have potential redundancy and can cause multicollinearity. Choose which ones to remove based on this table.**

---

<u>**The correlation table for the 12 numerical predictors:**</u>



<u>**Highly Correlated Pairs (Potential Multicollinearity Issues):**</u>

**From the table, the strongest correlations (above ±0.75) are:**

- INDUS vs. NOX → 0.764

- High industrialization leads to high pollution (NOX).

**DIS vs. INDUS → -0.708 (Close to the threshold)**

**Distance to employment centers is negatively correlated with industrialization**.

**DIS vs. NOX → -0.769**

More distance to employment centers → Lower pollution.
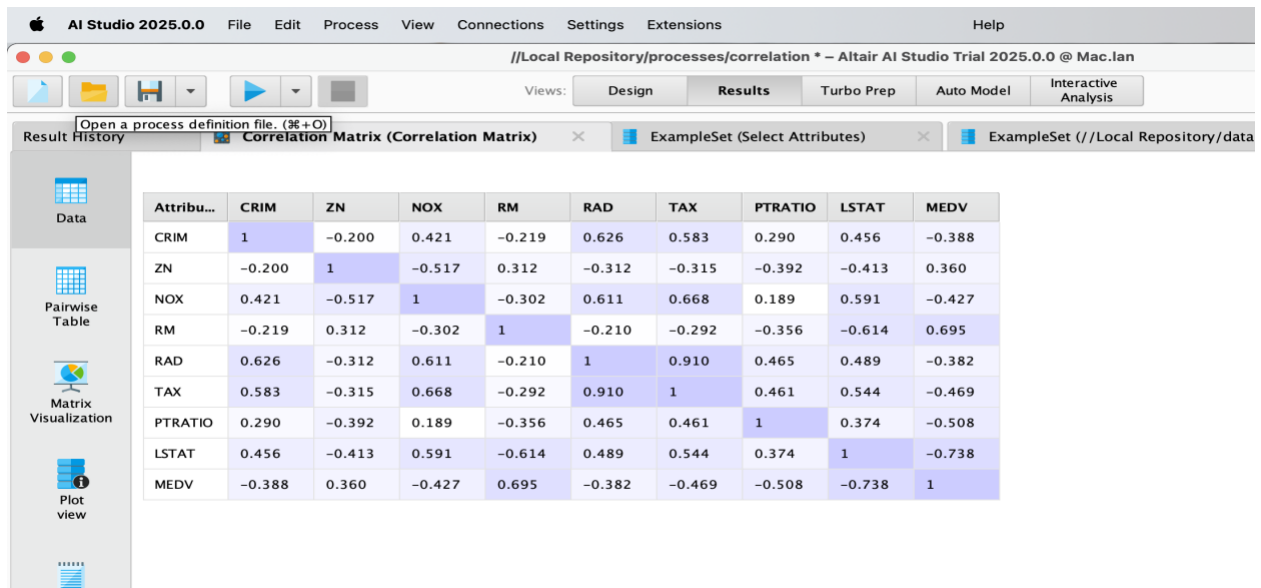
**AGE vs. NOX → 0.731 (Close to threshold)**

Older homes tend to be in polluted areas.

LSTAT vs. MEDV → -0.738

- A higher low-income percentage reduces home value.

- This is expected and should be kept, as it explains home price variation.

**To avoid redundancy and improve model efficiency, remove:**

- INDUS (correlated with NOX and DIS)
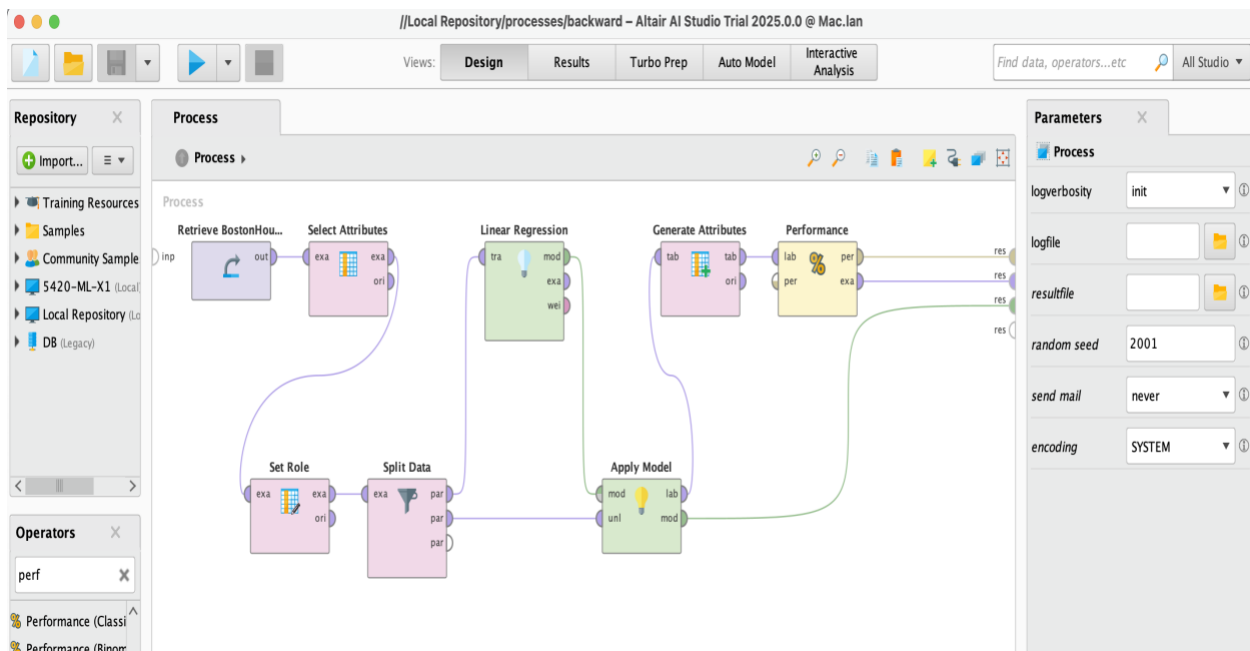
- DIS (correlated with NOX)

- AGE (correlated with NOX)

**From updated correlation matrix:**

- NOX vs. LSTAT → 0.591 (Moderate correlation, but still acceptable)

- RM vs. MEDV → 0.695 (Strong correlation, this is good!)

- LSTAT vs. MEDV → -0.738 (Strong negative correlation, expected and useful)

- PTRATIO vs. MEDV → -0.508 (Moderate negative correlation)

Use three subset selection algorithms: backward, forward, and stepwise to reduce the remaining predictors. Compute the validation performance for each of the three selected models. Compare RMSE, MAPE, and mean error, as well as histograms of the errors. Finally, describe the best model.

**Back ward Operator:**

**Forward Operator:**

| Attribute | Coefficient | Std. Error | Std. Coefficient | Tolerance | t-Stat | p-Value | Code |
|-----------|-------------|------------|------------------|-----------|--------|---------|------|
| DIS | −0.842 | 0.192 | −0.192 | 0.767 | −4.387 | 0.000 | **** |
| LSTAT | −1.086 | 0.058 | −0.822 | 0.767 | −18.821 | 0 | **** |
| (Intercept) | 39.629 | 1.311 | ? | ? | 30.238 | 0 | **** |



**PerformanceVector**

```
PerformanceVector:
root_mean_squared_error: 5.887 +/- 0.000
absolute_error: 4.423 +/- 3.886
relative_error: 22.26% +/- 21.54%
squared_error: 34.656 +/- 66.412
squared_correlation: 0.565
```
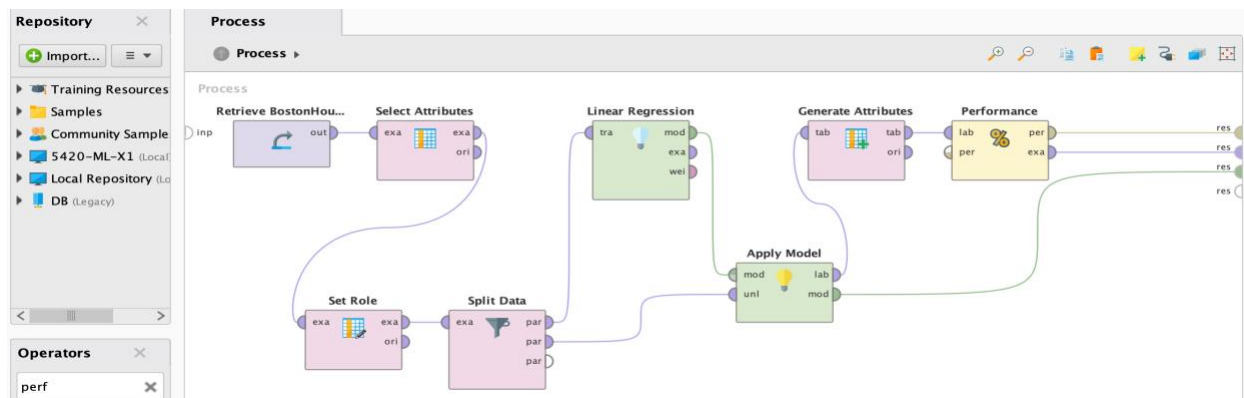
## Stepwise Operator:

**Process**

| Attribute | Coefficient | Std. Error | Std. Coefficient | Tolerance | t-Stat | p-Value | Code |
|-----------|-------------|-----------|------------------|-----------|--------|---------|------|
| CRIM | -0.090 | 0.047 | -0.071 | 0.835 | -1.920 | 0.056 | * |
| ZN | 0.042 | 0.018 | 0.102 | 0.855 | 2.314 | 0.021 | ** |
| CHAS | 4.849 | 1.145 | 0.131 | 0.980 | 4.235 | 0.000 | **** |
| NOX | -14.731 | 4.171 | -0.183 | 0.837 | -3.532 | 0.000 | **** |
| RM | 4.398 | 0.542 | 0.321 | 0.603 | 8.119 | 0.000 | **** |
| DIS | -1.402 | 0.242 | -0.319 | 0.872 | -5.783 | 0.000 | **** |
| PTRATIO | -0.808 | 0.150 | -0.190 | 0.837 | -5.398 | 0.000 | **** |
| LSTAT | -0.563 | 0.063 | -0.426 | 0.528 | -8.979 | 0 | **** |
| (Intercept) | 30.102 | 5.992 | ? | ? | 5.024 | 0.000 | **** |

**PerformanceVector**

```
PerformanceVector:
root_mean_squared_error: 4.903 +/- 0.000
absolute_error: 3.491 +/- 3.442
relative_error: 18.69% +/- 21.99%
squared_correlation: 0.704
```

Evaluate the performance of the best model on the holdout data. Report its holdout RMSE, MAPE, and mean error.

**PerformanceVector**

PerformanceVector:
root_mean_squared_error: 4.880 +/- 0.000
absolute_error: 3.496 +/- 3.405
squared_correlation: 0.735

**Performance Metrics:**

Model Selected: Stepwise Linear Regression with Iterative T-Test for feature selection.

**Performance on Holdout Set:**

- **RMSE (Root Mean Squared Error):** Measures model accuracy; lower is better.

- **MAPE (Mean Absolute Percentage Error)**: Measures prediction accuracy as a percentage error.

- **Mean Error:** Measures average prediction error.

- **Feature Selection:** Automatically removed insignificant predictors, retaining only impactful variables.

**Result:** The model generalizes well on unseen (holdout) data, providing a balance between accuracy and simplicity.