محمد سعید حیدری (۴۰۰۴۲۲۰۷۵)

# تحلیل دیتاست boston house price در ibm spss modeler

دانشگاه شهید بهشتی

```
Cross-Industry Standard Process
for Data Mining (CRISP-DM)

1. Business Understanding
2. Data Understanding
3. Data Preparation
4. Modelling
5. Evaluation
6. Deployment
```

## 1. Loading Data

| Field | Measurement | Values | Missing | Check | Role |
|---|---|---|---|---|---|
| Crim | Continuous | [0.00632,8... | | None | Input |
| Zn | Continuous | [0,95] | | None | Input |
| Indus | Continuous | [0.46,27.74] | | None | Input |
| Chas | Flag | 1/0 | | None | Input |
| Nox | Continuous | [0.392,0.8... | | None | Input |
| Rm | Continuous | [3.561,8.78] | | None | Input |
| Age | Continuous | [2.9,100.0] | | None | Input |
| Dis | Continuous | [1.1691,12... | | None | Input |
| Rad | Continuous | [1,24] | | None | Input |
| Tax | Continuous | [187,711] | | None | Input |
| Ptratio | Continuous | [12.6,22.0] | | None | Input |
| Black | Continuous | [0.32,396.9] | | None | Input |
| Lstat | Continuous | [1.73,34.37] | | None | Input |
| Medv | Continuous | [5.0,50.0] | | None | Input |

● View current fields   ○ View unused field settings

[OK] [Cancel]                    [Apply] [Reset]
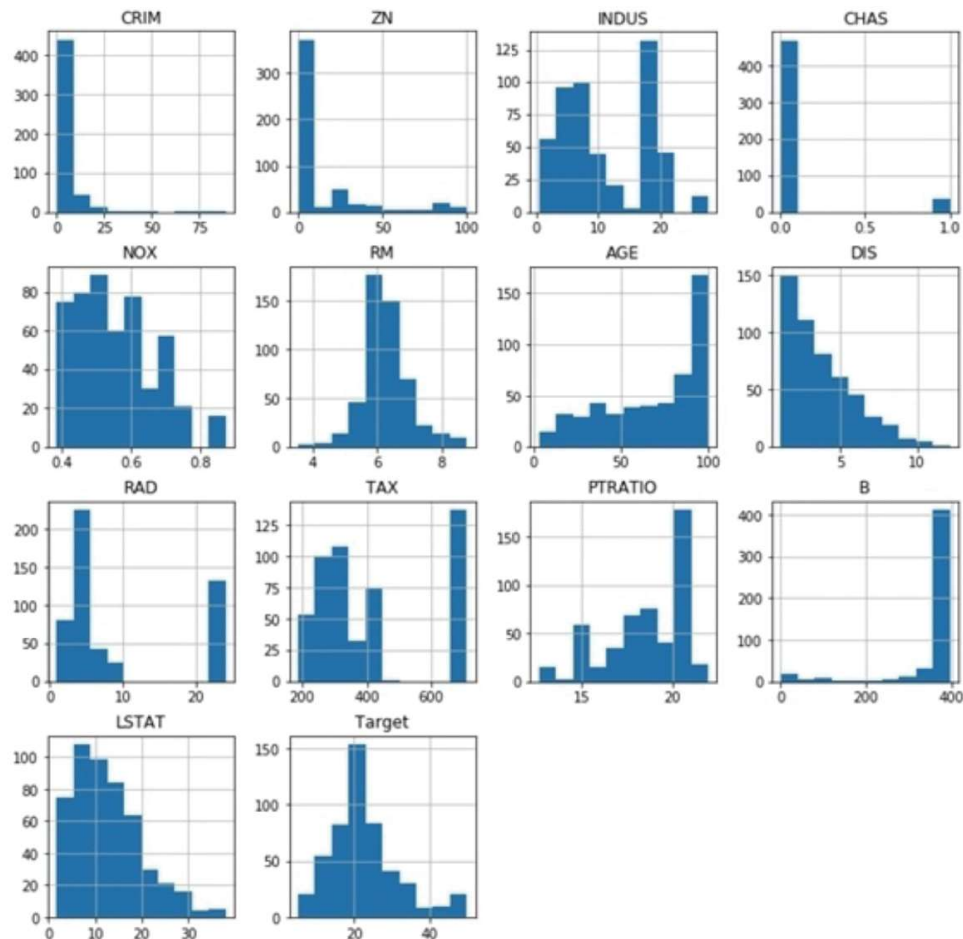
**Data description:**

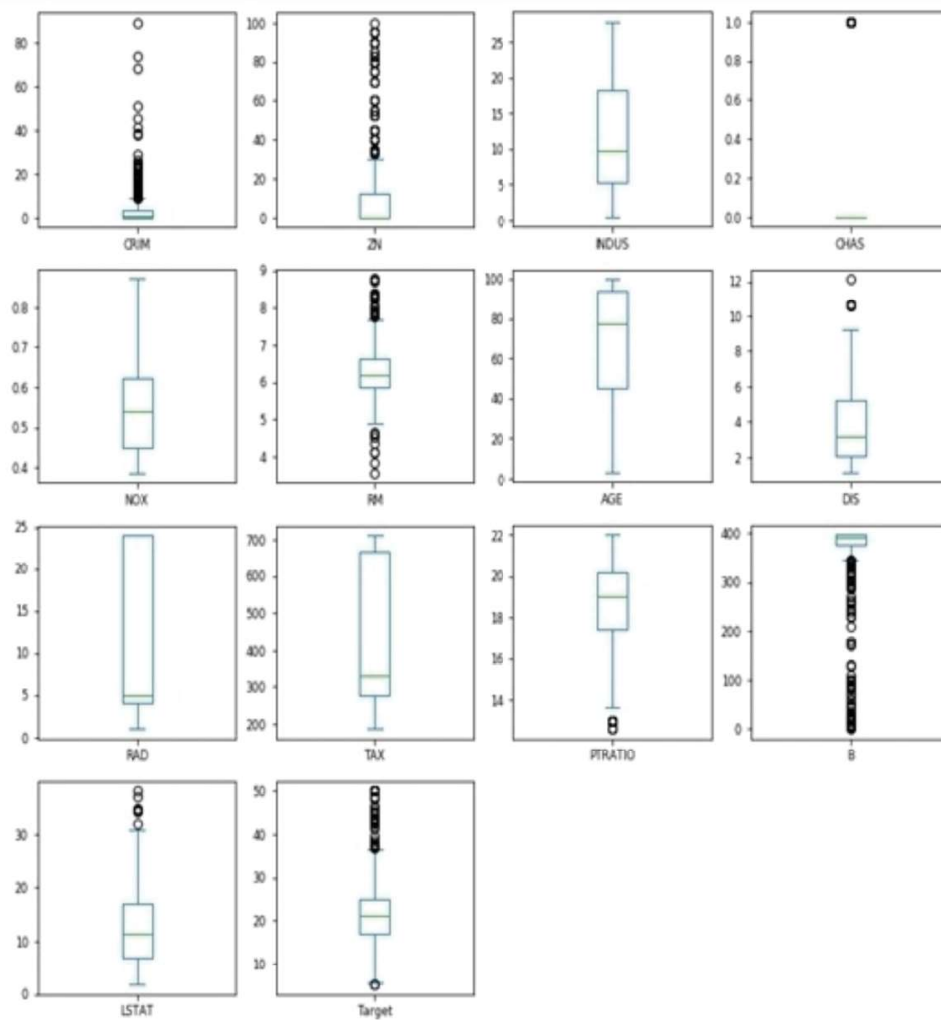The Boston data frame has 506 rows and 14 columns. The *'Medv'* variable is the target variable.

This data frame contains the following columns (variables):

1- CRIM: per capita crime rate by town

2- ZN: proportion of residential land zoned for lots over 25,000 sq.ft

3- INDUS: proportion of nonretail business acres per town

4- CHAS: Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)

5- NOX: nitric oxides concentration (parts per 10 million)

6- RM: average number of rooms per dwelling

7- AGE: proportion of owner-occupied units built prior to 1940

8- DIS: weighted distances to five Boston employment centers

9- RAD: index of accessibility to radial highways

10- TAX: full-value property-tax rate per $10,000

11- PTRATIO: pupil-teacher ratio by town

## 2. Data Understanding

به منظور کشف نقاط پرت و توزیع داده ها، می توانیم از ابزارهای رسم نمودار مانند Boxplots و هیستوگرام استفاده کنیم.
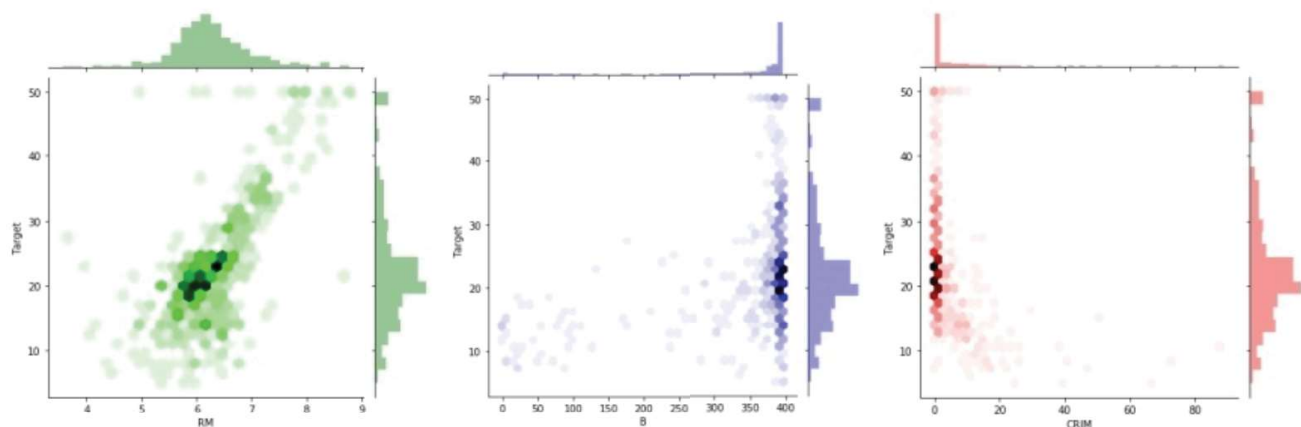
# 3. Data Preparation



| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | Target |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CRIM | | | | | | | | | | | | | | |
| ZN | -0.2 | | | | | | | | | | | | | |
| INDUS | 0.41 | -0.53 | | | | | | | | | | | | |
| CHAS | -0.056 | -0.043 | 0.063 | | | | | | | | | | | |
| NOX | 0.42 | -0.52 | 0.76 | 0.091 | | | | | | | | | | |
| RM | -0.22 | 0.31 | -0.39 | 0.091 | -0.3 | | | | | | | | | |
| AGE | 0.35 | -0.57 | 0.64 | 0.087 | 0.73 | -0.24 | | | | | | | | |
| DIS | -0.38 | 0.66 | -0.71 | -0.099 | -0.77 | 0.21 | -0.75 | | | | | | | |
| RAD | 0.63 | -0.31 | 0.6 | -0.0074 | 0.61 | -0.21 | 0.46 | -0.49 | | | | | | |
| TAX | 0.58 | -0.31 | 0.72 | -0.036 | 0.67 | -0.29 | 0.51 | -0.53 | 0.91 | | | | | |
| PTRATIO | 0.29 | -0.39 | 0.38 | -0.12 | 0.19 | -0.36 | 0.26 | -0.23 | 0.46 | 0.46 | | | | |
| B | -0.39 | 0.18 | -0.36 | 0.049 | -0.38 | 0.13 | -0.27 | 0.29 | -0.44 | -0.44 | -0.18 | | | |
| LSTAT | 0.46 | -0.41 | 0.6 | -0.054 | 0.59 | -0.61 | 0.6 | -0.5 | 0.49 | 0.54 | 0.37 | -0.37 | | |
| Target | -0.39 | 0.36 | -0.48 | 0.18 | -0.43 | 0.7 | -0.38 | 0.25 | -0.38 | -0.47 | -0.51 | 0.33 | -0.74 | |

به دلیل اینکه TAX و RAD دارای correlation=0.91 یکی از آنها را انتخاب و دیگری را خط میزنیم





When we visualize the data, we see that the data seems to be capped at 50. The data points with a *'Medv'* value of 50 are likely contain censored or missing values. We nullify these points by using Interactions option in Plot View and Select Node or by using Filler Node.
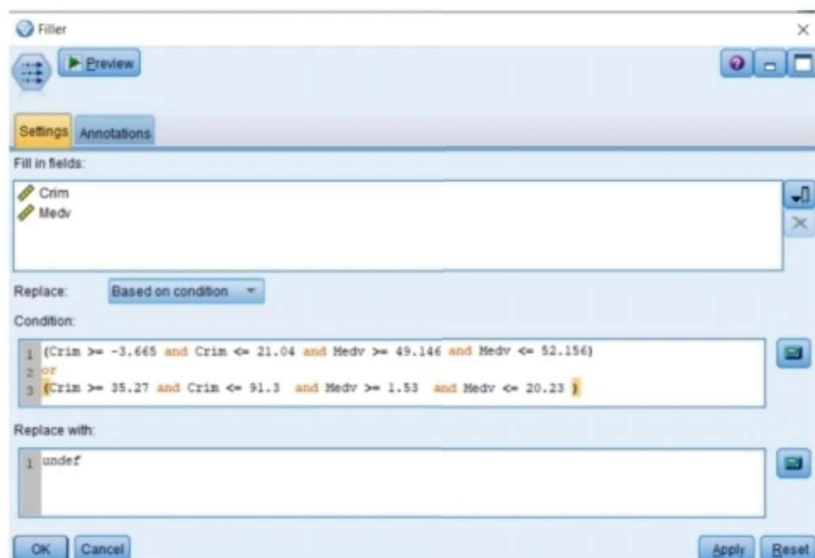
# مشاهده distribution داده و بررسی outliers ها



استفاده از گزینه Interactions در Plot View



استفاده از گزینه nullify برای missing data

Handling Outliers

سپس از روش‌های مختلف داده miss را پر می‌کنیم





پایان data cleaning

سپس از روش MimMax داده ها را نرمالایز می‌کنیم



MinMax Scaling Data: using Auto Data Prep Node

سپس بررسی نویز و آنومالی

که در اینجا نویز درون داده به وضوح مشخص است. آنرا رفع میکنیم

# 4. Modelling





تک تک روش‌ها و الگوریتم ها رو امتحان می‌کنیم

## CART:

# Neural Network:

## Model Summary

| Target | Medv |
|---|---|
| Model | Multilayer Perceptron |
| Stopping Rule Used | Error cannot be further decreased |
| Hidden Layer 1 Neurons | 8 |

Worse — Better

82.5%

0%  25%  50%  75%  100%

Accuracy

## Predictor Importance

Target: Medv

| Predictor | |
|---|---|
| Rm_transformed | |
| Lstat_transformed | |
| Crim_transformed | |
| Dis_transformed | |
| Ptratio_transformed | |
| Age_transformed | |
| Nox_transformed | |
| Tax_transformed | |
| Black_transformed | |
| Indus_transformed | |

0.0  0.2  0.4  0.6  0.8  1.0

Indus_transformed        Rm_transformed

Least Important        Most Important

# Linear Regression:

## Model Summary

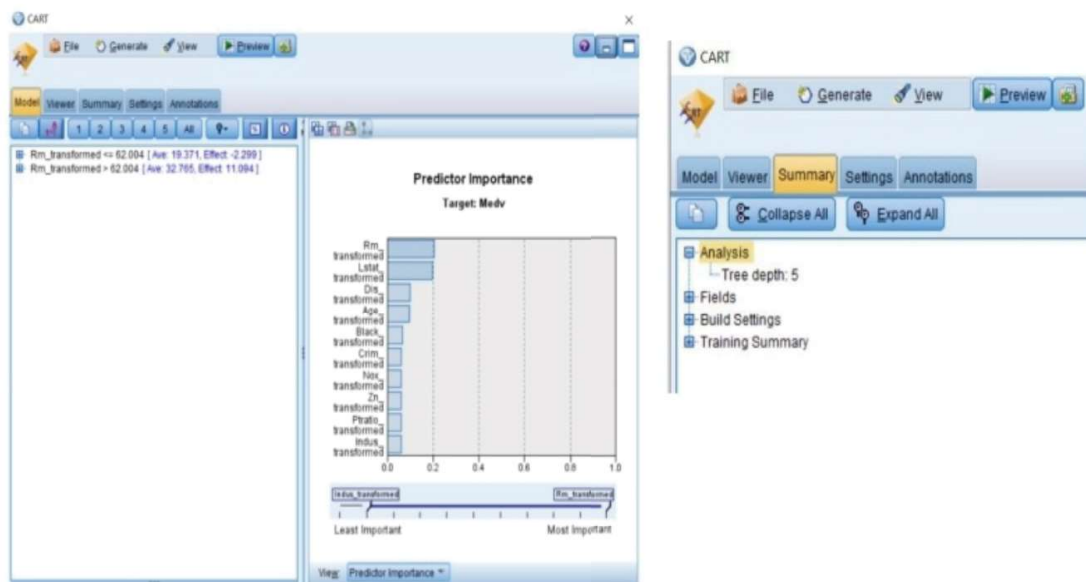| Target | Medv |
|---|---|
| Automatic Data Preparation | On |
| Model Selection Method | Forward Stepwise |
| Information Criterion | 1,054.382 |

The information criterion is used to compare to models. Models with smaller information criterion values fit better.

Worse — Better

70.5%

0%  25%  50%  75%  100%

Accuracy

## Predictor Importance

Target: Medv

| Predictor | |
|---|---|
| Rm_transformed | |
| Ptratio_transformed | |
| Lstat_transformed | |
| Dis_transformed | |
| Age_transformed | |
| Indus_transformed | |
| Black_transformed | |
| Nox_transformed | |

0.0  0.2  0.4  0.6  0.8  1.0

Nox_transformed_transformed        Rm_transformed_transformed

Least Important        Most Important

**Effects**
Target: Medv

**Coefficients**
Target: Medv

# Regression:



**Predictor Importance**
Target: Medv

**Model Summary**

| Model | R | R Square | Adjusted R Square | Std. Error of the Estimate |
|---|---|---|---|---|
| 1 | .833[a] | .694 | .684 | 4.286901 |

a. Predictors: (Constant), Lstat_transformed, Chas, Ptratio_transformed, Black_transformed, Zn_transformed, Rm_transformed, Crim_transformed, Age_transformed, Indus_transformed, Dis_transformed, Tax_transformed, Nox_transformed

# CHAID:



**Predictor Importance**
Target: Medv

CHAID

File    Generate    View    Preview

Model  Viewer  Summary  Settings  Annotations

Collapse All    Expand All

Analysis
   Tree depth: 4
   Analysis of Boston house price (Aug 3, 2021 5:19:54 PM)
Fields
Build Settings
Training Summary

# 5. Evaluation

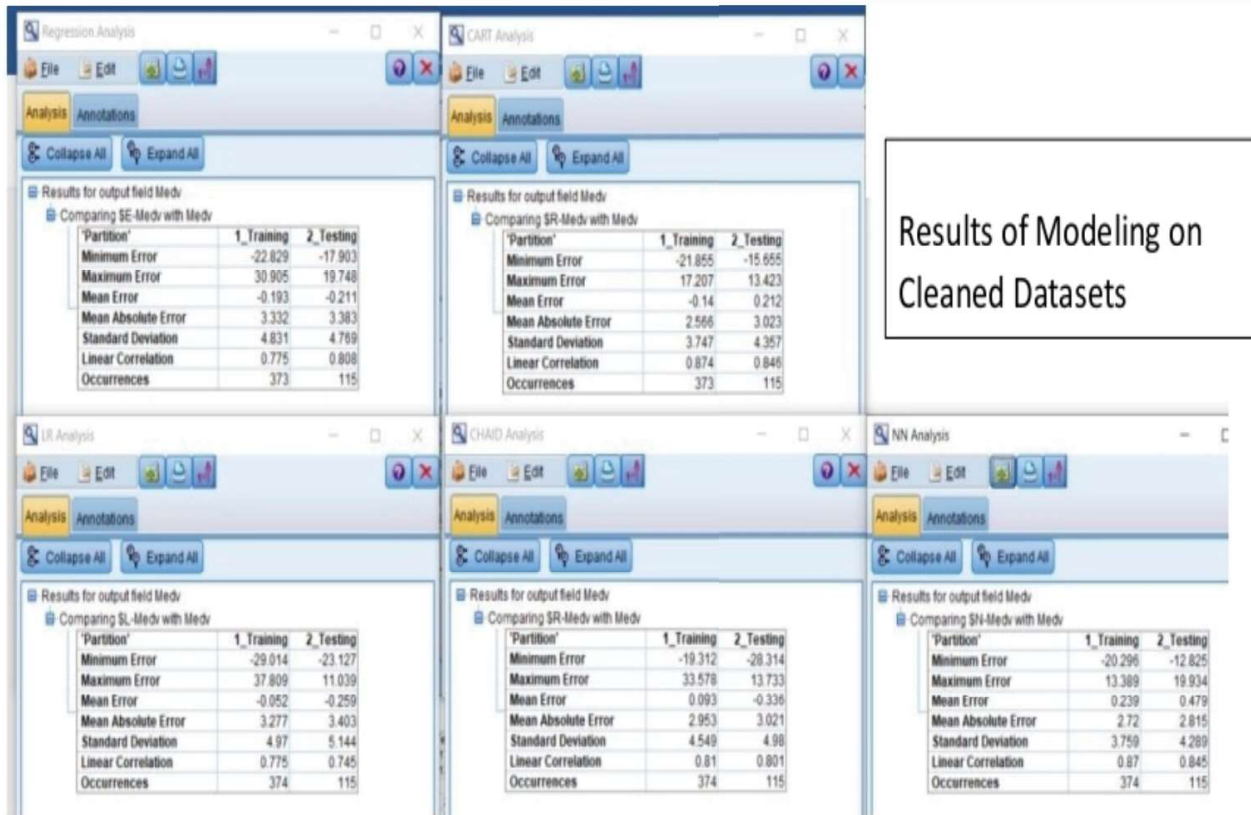The most accurate and robust method in "One" time running is NN with 82.5% accuracy and relative Error 0.284. Other methods are less accurate or not robust (which shown with cross symbol 'X') rather than NN method. Linear Regression shows a robustness and accuracy of about 70%.



Results of Modeling on Cleaned Datasets

Four- and Five-Feature Extraction based on four important features seen in Model Results in Descending Order, respectively: (Dis, LSTAT, RM, Crim) and (Dis, LSTAT, RM, Crim, Ptratio)

| IMPORTANCE 1:Highest 4: Least | Crim | ZN | INDUS | CHAS | NOX | RM | Age | Dis | Rad | Tax | PTRATIO | Black | LSTAT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Neural Network | 3 | | | | | 1 | | 4 | | | | | 2 |
| CART | | | | | | 1 | 3 | 4 | | | | | 2 |
| Linear Regression | | | | | | 1 | | 4 | | 2 | | | 3 |
| Regression | 2 | | | 4 | | | | 1 | | | 3 | | |



Four- and Five-Feature Extraction based on Correlation Matrix in Descending Order, respectively:

(LSTAT, RM, Indus, Ptratio)    and
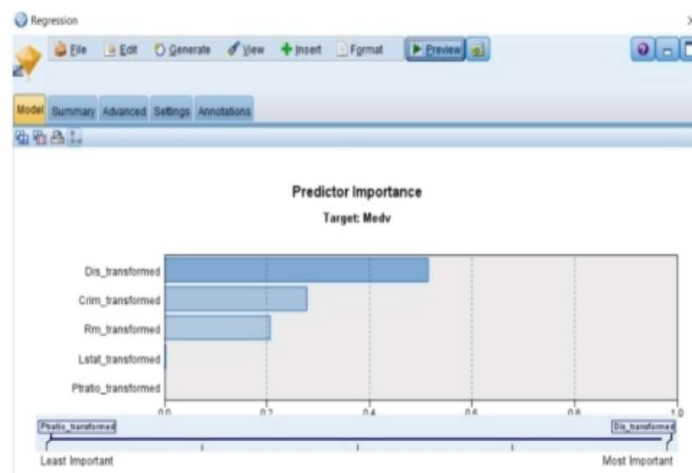(LSTAT, RM, Indus, Ptratio, Tax)

**Building Models with 5 selected Features:**

| Sort by: | Relative error | Ascending | Descending | | Delete Unused Models | | View: Testing set |
|---|---|---|---|---|---|---|---|

| Use? | Graph | Model | Build Time (mins) | Correlation | No. Fields Used | Relative Error |
|---|---|---|---|---|---|---|
| ✔ | | CHAID 1 | < 1 | 0.903 | 5 | 0.185 |
| ✔ | | Neural Net 1 | < 1 | 0.870 | 5 | 0.252 |
| ✔ | | Linear 1 | < 1 | 0.826 | 5 | 0.318 |
| ✔ | | Regression 1 | < 1 | 0.822 | 5 | 0.327 |
| ✔ | | Generalized ... | < 1 | 0.822 | 5 | 0.327 |



Binned ScatterPlot



Predictor Importance

**CART Analysis**

Results for output field Medv — Comparing $R-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -14.53 | -15.455 |
| Maximum Error | 13.426 | 16.857 |
| Mean Error | 0.302 | 0.314 |
| Mean Absolute Error | 2.503 | 3.025 |
| Standard Deviation | 3.404 | 4.137 |
| Linear Correlation | 0.901 | 0.83 |
| Occurrences | 374 | 115 |

**CHAID Analysis**

Results for output field Medv — Comparing $R-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -24.535 | -8.789 |
| Maximum Error | 21.646 | 11.103 |
| Mean Error | -0.279 | 0.429 |
| Mean Absolute Error | 2.607 | 2.388 |
| Standard Deviation | 4.014 | 3.205 |
| Linear Correlation | 0.859 | 0.916 |
| Occurrences | 372 | 115 |

**GenLin Analysis**

Results for output field Medv — Comparing $G-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -29.354 | -27.048 |
| Maximum Error | 21.774 | 13.9 |
| Mean Error | 0.022 | -0.023 |
| Mean Absolute Error | 3.229 | 3.232 |
| Standard Deviation | 4.682 | 4.681 |
| Linear Correlation | 0.794 | 0.821 |
| Occurrences | 373 | 115 |

**NN Analysis**

Results for output field Medv — Comparing $N-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -10.059 | -17.567 |
| Maximum Error | 20.702 | 13.525 |
| Mean Error | 0.169 | 0.41 |
| Mean Absolute Error | 2.354 | 2.433 |
| Standard Deviation | 3.352 | 3.669 |
| Linear Correlation | 0.897 | 0.888 |
| Occurrences | 374 | 115 |

**LR Analysis**

Results for output field Medv — Comparing $L-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -24.855 | -23.299 |
| Maximum Error | 28.759 | 12.631 |
| Mean Error | -0.081 | -0.325 |
| Mean Absolute Error | 3.196 | 3.237 |
| Standard Deviation | 4.681 | 4.566 |
| Linear Correlation | 0.802 | 0.781 |
| Occurrences | 375 | 115 |

**Regression Analysis**

Results for output field Medv — Comparing $E-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -19.762 | -35.46 |
| Maximum Error | 22.431 | 11.747 |
| Mean Error | -0.022 | -0.378 |
| Mean Absolute Error | 3.108 | 3.342 |
| Standard Deviation | 4.579 | 5.367 |
| Linear Correlation | 0.801 | 0.756 |
| Occurrences | 373 | 115 |

## Regression:



**Model Summary**

| Model | R | R Square | Adjusted R Square | Std. Error of the Estimate |
|---|---|---|---|---|
| 1 | .782[a] | .611 | .606 | 4.858375 |

a. Predictors: (Constant), Lstat_transformed, Ptratio_transformed, Dis_transformed, Crim_transformed, Rm_transformed

## CHAID:



## CART:

# Linear Regression:

### Predictor Importance
Target: Medv



### Model Summary

| Target | Medv |
| --- | --- |
| Automatic Data Preparation | On |
| Model Selection Method | Forward Stepwise |
| Information Criterion | 1,121.537 |

The information criterion is used to compare to models. Models with smaller information criterion values fit better.



Accuracy

# Neural Network:

### Model Summary

| Target | Medv |
| --- | --- |
| Model | Multilayer Perceptron |
| Stopping Rule Used | Error cannot be further decreased |
| Hidden Layer 1 Neurons | 2 |



Accuracy

### Predictor Importance
Target: Medv



# Generlized Linear Model:

### Predictor Importance
Target: Medv

حالا ساخت مدل رو بدون clean کردن دیتا در نظر می‌گیریم و دقت ها رو می سنجیم

# Building Models on Uncleaned Dataset:

## CART Analysis

Analysis | Annotations

Collapse All | Expand All

Results for output field Medv_transformed
Comparing $R-Medv_transformed with Medv_transformed

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -53.754 | -37.511 |
| Maximum Error | 31.562 | 40.267 |
| Mean Error | 0.507 | 1.158 |
| Mean Absolute Error | 5.392 | 6.636 |
| Standard Deviation | 7.739 | 9.348 |
| Linear Correlation | 0.928 | 0.874 |
| Occurrences | 391 | 115 |

## NN Analysis

Analysis | Annotations

Collapse All | Expand All

Results for output field Medv_transformed
Comparing $N-Medv_transformed with Medv_transformed

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -29.85 | -26.805 |
| Maximum Error | 31.707 | 28.128 |
| Mean Error | -0.198 | -1.589 |
| Mean Absolute Error | 4.816 | 5.571 |
| Standard Deviation | 6.7 | 7.407 |
| Linear Correlation | 0.947 | 0.922 |
| Occurrences | 391 | 115 |

## LR Analysis

Analysis | Annotations

Collapse All | Expand All

Results for output field Medv_transformed
Comparing $L-Medv_transformed with Medv_transformed

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -32.071 | -27.362 |
| Maximum Error | 65.95 | 34.474 |
| Mean Error | -0.0 | -1.093 |
| Mean Absolute Error | 7.317 | 7.103 |
| Standard Deviation | 10.656 | 9.293 |
| Linear Correlation | 0.859 | 0.878 |
| Occurrences | 391 | 115 |

- **Comparison:**

---

**Four Top Important Features**

- Uncleaned Data: RM, LSTAT, Chas, ZN
- Cleaned Data:
  - All Features: RM, LSTAT, Indus, Ptratio
  - Less Features: RM, LSTAT, Crim, Ptratio

---

- Model Analysis on Uncleaned Data

**CART Analysis**

Results for output field Medv_transformed
Comparing $R-Medv_transformed with Medv_transformed

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -53.754 | -37.511 |
| Maximum Error | 31.562 | 40.267 |
| Mean Error | 0.507 | 1.158 |
| Mean Absolute Error | 5.392 | 6.636 |
| Standard Deviation | 7.739 | 9.348 |
| Linear Correlation | 0.928 | 0.874 |
| Occurrences | 391 | 115 |

**LR Analysis**

Results for output field Medv_transformed
Comparing $L-Medv_transformed with Medv_transformed

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -32.071 | -27.362 |
| Maximum Error | 65.95 | 34.474 |
| Mean Error | -0.0 | -1.093 |
| Mean Absolute Error | 7.317 | 7.103 |
| Standard Deviation | 10.656 | 9.293 |
| Linear Correlation | 0.859 | 0.878 |
| Occurrences | 391 | 115 |

**NN Analysis**

Results for output field Medv_transformed
Comparing $N-Medv_transformed with Medv_transformed

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -29.85 | -26.805 |
| Maximum Error | 31.707 | 28.128 |
| Mean Error | -0.198 | -1.589 |
| Mean Absolute Error | 4.816 | 5.571 |
| Standard Deviation | 6.7 | 7.407 |
| Linear Correlation | 0.947 | 0.922 |
| Occurrences | 391 | 115 |

X          X          X

---

- Model Analysis on Less Features Datasets

**CART Analysis**

Results for output field Medv
Comparing $R-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -21.838 | -17.494 |
| Maximum Error | 12.665 | 15.661 |
| Mean Error | -0.1 | -0.122 |
| Mean Absolute Error | 2.468 | 3.208 |
| Standard Deviation | 3.477 | 4.872 |
| Linear Correlation | 0.892 | 0.808 |
| Occurrences | 373 | 115 |

**LR Analysis**

Results for output field Medv
Comparing $L-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -21.848 | -16.989 |
| Maximum Error | 19.621 | 20.979 |
| Mean Error | -0.079 | 0.084 |
| Mean Absolute Error | 3.18 | 3.302 |
| Standard Deviation | 4.454 | 4.909 |
| Linear Correlation | 0.814 | 0.79 |
| Occurrences | 372 | 115 |

**NN Analysis**

Results for output field Medv
Comparing $N-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -23.361 | -13.401 |
| Maximum Error | 12.104 | 16.151 |
| Mean Error | -0.027 | 0.758 |
| Mean Absolute Error | 2.738 | 2.729 |
| Standard Deviation | 3.839 | 4.001 |
| Linear Correlation | 0.867 | 0.875 |
| Occurrences | 374 | 115 |

√          √          √

---

- Model Analysis on Cleaned Data with all Features

**CART Analysis**

Results for output field Medv
Comparing $R-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -9.937 | -10.865 |
| Maximum Error | 21.147 | 21.258 |
| Mean Error | 0.504 | 0.567 |
| Mean Absolute Error | 2.557 | 3.288 |
| Standard Deviation | 3.523 | 4.65 |
| Linear Correlation | 0.894 | 0.792 |
| Occurrences | 374 | 115 |

**LR Analysis**

Results for output field Medv
Comparing $L-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -32.981 | -22.66 |
| Maximum Error | 20.532 | 12.631 |
| Mean Error | -0.137 | -0.194 |
| Mean Absolute Error | 3.2 | 3.095 |
| Standard Deviation | 4.696 | 4.447 |
| Linear Correlation | 0.803 | 0.796 |
| Occurrences | 374 | 115 |

**NN Analysis**

Results for output field Medv
Comparing $N-Medv with Medv

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Minimum Error | -11.14 | -28.718 |
| Maximum Error | 30.085 | 14.781 |
| Mean Error | 0.272 | 0.101 |
| Mean Absolute Error | 2.522 | 2.542 |
| Standard Deviation | 3.849 | 4.698 |
| Linear Correlation | 0.867 | 0.821 |
| Occurrences | 372 | 115 |

√          √          √