Rapidminer USA-Housing

1. Import  *.csv file,  Get csv file from local, click **import Configuration Wizard** , connect [readcsv]--res and [Run]





Step 2. Operators:  type and move [Select Attributes] to **design** pane.

Parameters: attribute filter type: **a subset** , click [select Attributes], the dialog window opens, choose every items beside **Address** [Apply]

Step 3. Operators:  type,  move [Set Role] to **design** pane. We set dependent value: **price => label**

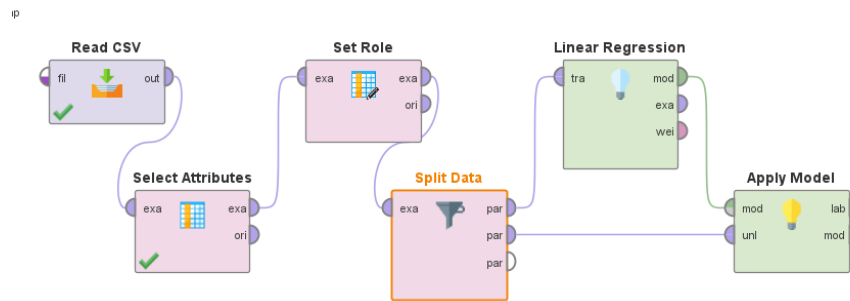Step 4. Operators:  type and move **[Split Data]** to **design** pane.
 Parameters: Partition, click Edit Enumertation,  click add entry 0.8, then add entry: 0.2  [ok]
Step 5. Operators:  type and move **[Linear Regression]** to **design** pane.
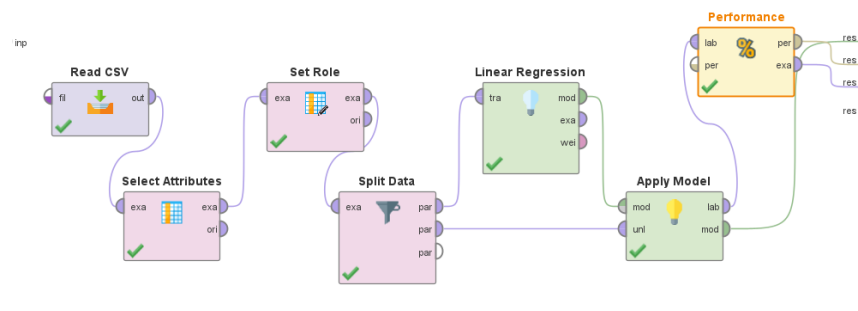  **[Split Data]** par ----- tra **[Linear Regression]**
Step 6: Operators:  type and move **[Apply Mode]** to **design** pane.
Connection as following



Step 7: Operators:  type and move **[ Performance (Regression)]** to **design** pane.
Parameters: select: root mean squared error( default), squared correlation



Data for LinearRegression(Linear Regression)

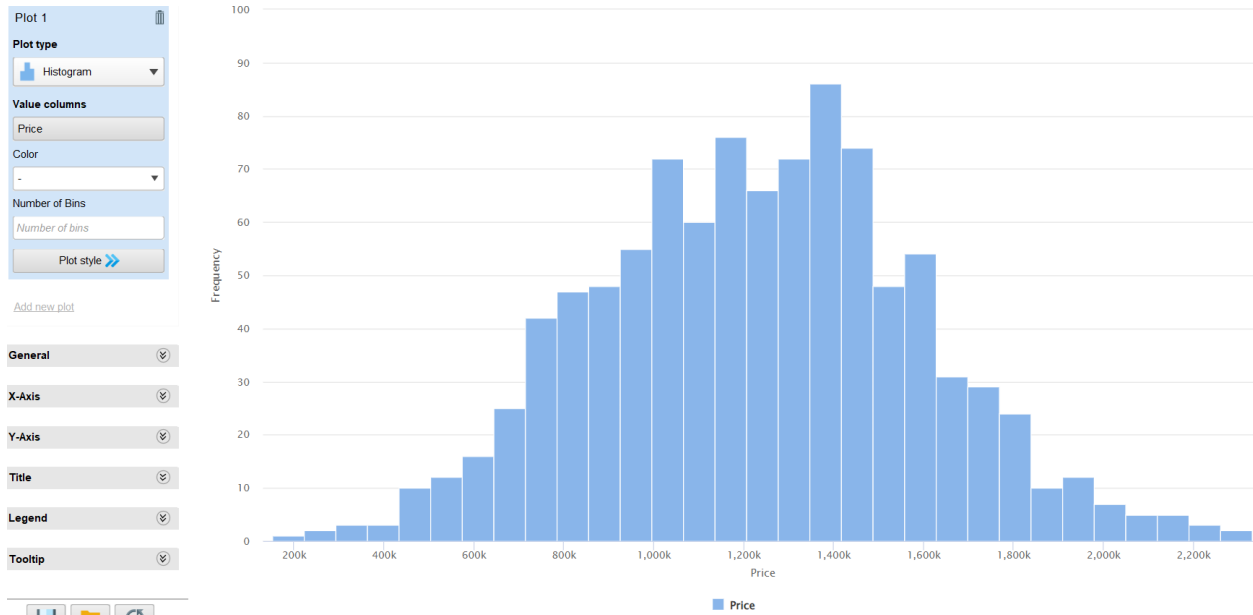| Attribute | Coefficient | Std. Error | Std. Coefficient | Tolerance | t-Stat | p-Value | Code |
|---|---|---|---|---|---|---|---|
| Avg. Area Income | 21.639 | 0.151 | 0.653 | 1.000 | 143.358 | 0 | **** |
| Avg. Area House Age | 165905.349 | 1617.971 | 0.467 | 0.999 | 102.539 | 0 | **** |
| Avg. Area Number of Rooms | 120333.459 | 1772.026 | 0.349 | 0.996 | 67.907 | 0 | **** |
| Avg. Area Number of Bedroo... | 2381.785 | 1464.263 | 0.008 | 0.970 | 1.627 | 0.104 | |
| Area Population | 15.222 | 0.162 | 0.428 | 1.000 | 94.004 | 0 | **** |
| (Intercept) | -2644341.426 | 19296.747 | ? | ? | -137.036 | 0 | **** |

## LinearRegression

```
  21.639 * Avg. Area Income
+ 165905.349 * Avg. Area House Age
+ 120333.459 * Avg. Area Number of Rooms
+ 2381.785 * Avg. Area Number of Bedrooms
+ 15.222 * Area Population
- 2644341.426
```

PerformanceVector

Root_mean_squared_error : 101010.879 +/- 0.000

Square correlation: 0.921
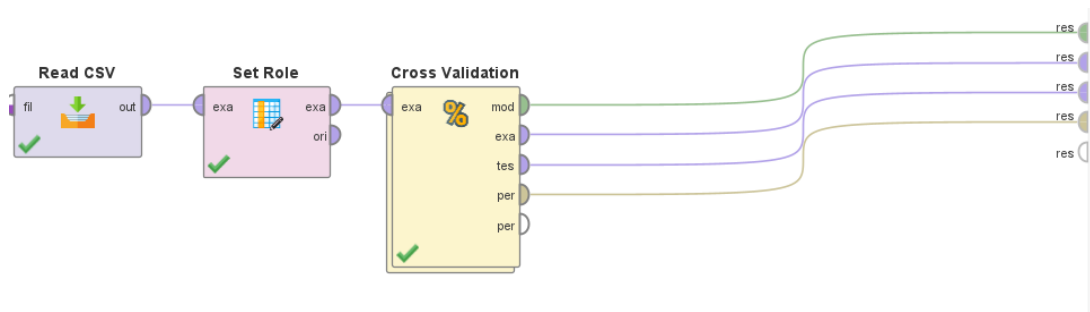
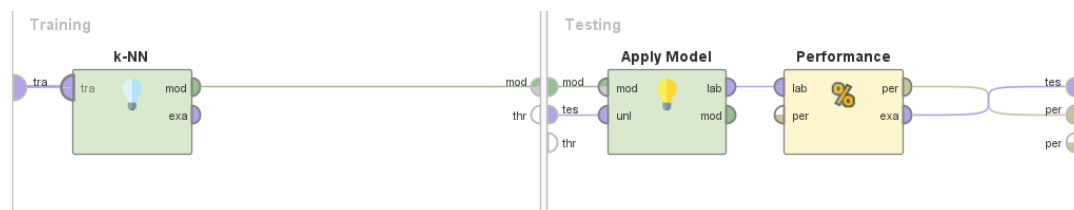ExampleSet(Apple Model)  ==Visualizations



- Data

| Name | | Type | Missing | Statistics | | |
|------|---|------|---------|-----------|---|---|
| | | | | Min | Max | Average |
| Label **Price** | ⌄ | Real | 0 | 152071.875 | 2332110.740 | 1233226.818 |
| Prediction **prediction(Price)** | ⌄ | Real | 0 | 220122.390 | 2476938.337 | 1233959.136 |
| **Avg. Area Income** | ⌄ | Real | 0 | 35454.715 | 107701.748 | 68755.094 |
| ⚠ **Avg. Area House Age** | ⌄ | Real | 0 | 2.683 | 8.973 | 5.983 |
| **Avg. Area Number of Rooms** | ⌄ | Real | 0 | 4.028 | 9.802 | 7.023 |
| **Avg. Area Number of Bedrooms** | ⌄ | Real | 0 | 2 | 6.500 | 4.024 |
| **Area Population** | ⌄ | Real | 0 | 6821.950 | 67727.229 | 35688.821 |

Comparing the train data price and the predicated price, the house price will grow, but it also has a lot of conditions to relate to the house price. It is not simple.
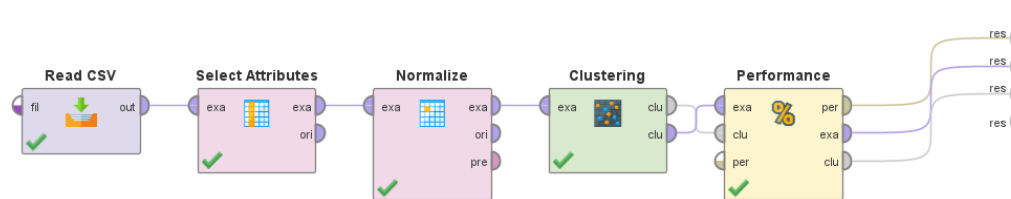
# Knn with Cross Validation ( USA-Housing)



## Inside Cross Validation



1. Set role ➜ price, Label

2. Apply do not do anything, Performance ➜ Performance

---

KMean ( USA-housing , This file is not good to use Kmean, just try for fun)



Select Attributes - Select every field, except [Address]
Normalize: The files need to be clustered
Clustering number 3
Performance: Performance( Cluster Distance Performance)