

# Bank of Belgium Project: Datasets from the National Bank of Belgium

Wesley Matheus

# Contents

<b>Financial Wealth of Households</b>	<b>2</b>
<b>Dataset Information</b>	<b>3</b>
Dataset Information, Cleaning and Transformation . . . . .	3
<b>Full Code</b>	<b>4</b>
<b>Appendix</b>	<b>5</b>
<b>Formulas</b>	<b>6</b>

# Financial Wealth of Households

# Dataset Information

## Dataset Information, Cleaning and Transformation

- Datasets: `stat.nbb.be`
- Articles: `archive.is/https://www.brusselstimes.com/belgium/1392510/belgian-households-are-richer-than-ever`

## Full Code

# Appendix

# Formulas

## Mean

$$m = \frac{1}{n} \sum_{i=1}^n v_i$$

## Variance

$$V = \frac{1}{n} \sum_{i=1}^n (v_i - m)^2$$

- Deviation from the mean:  $v_i - mean$

## Standard Deviation

$$SD = \sqrt{V}$$

## Root Mean Square Error

$$RMSE(\mathbf{X}, h) = \sqrt{\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2}$$

$$RMSE(\mathbf{Dataset}, MLAlgorithm) = \sqrt{\frac{1}{rows} \sum_{i=1}^{rows} (MLAlgorithm(predicted\ value^{(i)}) - label\ value^{(i)})^2}$$

- Euclidean distance: straight line  $d = \sqrt{\Delta x^2 + \Delta y^2}$
- The ML Algorithm takes into consideration all the column values of the dataset to form a column of predicted values.
- The RMSE measures the standard deviation of the predicted values from the label values.

### Mean Absolute Error

$$MAE(\mathbf{X}, h) = \frac{1}{m} \sum_{i=1}^m |h(x^{(i)}) - y^{(i)}|$$

$$MAE(Dataset, MLAlgorithm) = \frac{1}{rows} \sum_{i=1}^{rows} |MLAlgorithm(predictedvalue^{(i)}) - labelvalue^{(i)}|$$

- Manhattan distance: grid  $d = |\Delta x| + |\Delta y|$
- Both the RMSE and the MAE are ways to measure the distance between two vectors: the column of predicted values from the column of label values.
- The mean absolute error is preferred when the data has many outliers.

### Difference between RMSE and Standard Deviation:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n \left( y_{predicted}^{(i)} - y_{label}^{(i)} \right)^2}$$

$$STD = \sqrt{\frac{1}{n} \sum_{i=1}^n \left( y_{predicted}^{(i)} - mean \right)^2}$$

RMSE (Root Mean Squared Error) measures the average magnitude (value) of the differences (errors) between the predicted values and the true values (labels). In other words, it's the average "distance" between the predicted values and the label values. It is the deviation from the label.

Standard Deviation measures the average distance of the differences between the predicted values from their own mean. It measures how spread out the values (in a dataset) are from the mean value. When applied to predictions, it measures how spread out the predicted values are from their own mean. It is the deviation from the mean.



## Standardization of a Column

$$Column = V_0, V_1, V_2, V_3, \dots V_n \rightarrow Column' = Z_0, Z_1, Z_2, Z_3, \dots Z_n$$

$$Z_i = \frac{V_i - \text{mean}(Column)}{\text{Standard Deviation}(Column)}$$

$$\text{Mean}(Column') = \frac{1}{n} \sum_{i=1}^n (Z_i) \simeq 0$$

$$SD(Column') = \sqrt{\frac{1}{n} \sum_{i=1}^n (Z_i - \text{mean}')^2} \simeq 1$$

# Bibliography

- [1] Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems
- [2] Learning Deep Learning: Theory and Practice of Neural Networks, Computer Vision, Natural Language Processing, and Transformers Using Tensorflow
- [3] Mathematics for Machine Learning