Allison Molitor

M.L. Learning Journal

Start Date: 01/05/2022

Most Recent Edit: 08/24/2023

Machine Learning Regression Models

## *Linear Regression Modeling*

- Mathematical Representation
  - $Y = \beta 0 + \beta 1 * X + \varepsilon$
    - Y = dependent variable
    - B0 = Y int when X=0
    - $\beta 1$ = Slope coefficient (How much Y changes for 1 unit change in X)
    - X = independent variable
    - E = (Error) Variability not explained by model
- Real-world example
  - Required bandwidth and Network users
    - Required Bandwidth $= \beta 0 + \beta 1 *$ Network users $+ \varepsilon$
      - Ordinary Least Squares determine $\beta 0$ and $\beta 1$

- Mathematical computation

| Network users (X) | Required bandwidth in Mbps(Y) |
|---|---|
| 1 | 3 |
| 2 | 5 |
| 3 | 7 |
| 4 | 9 |
| 5 | 11 |

Equation setup
$\Sigma Y - n\beta 0 - \beta 1 \Sigma X = 0$

$\Sigma Y = 3 + 5 + 7 + 9 + 11 = 35$
$\Sigma X = 1 + 2 + 3 + 4 + 5 = 15$
n = 5 (sample size)

$35 - 5\beta_0 - \beta_1 15 = 0$
$5\beta_0 + 15\beta_1 = 35$

$\Sigma(XY) - \beta_0\Sigma X - \beta_1\Sigma(X^2) = 0$

$\Sigma(XY) = = (1 * 3) + (2 * 5) + (3 * 7) + (4 * 9) + (5 * 11) = 115$
$\Sigma(X^2) = = (1^2) + (2^2) + (3^2) + (4^2) + (5^2) = 55$
$115 - \beta_0 15 - \beta_1 55 = 0$
$15\beta_0 + 55\beta_1 = 115$

Solve for $\beta_1$ and $\beta_0$ given
$5\beta_0 + 15\beta_1 = 35$, $15\beta_0 + 55\beta_1 = 115$
$B_1 = 1$
$B_0 = 4$

$Y = \beta_0 + \beta_1*X + \varepsilon$

$Y = 4+1*X + e$

Meaning for every 1 unit increase in network users, approximately 4 more Mbps is required.

If a model was fitted to a dataset using the predict function, one could predict the required Mbps given the number of network users like so

(Predicted bandwidth requirement) = 4+1*(Network users) + e

RSS or Residual Sum of Squares can be used to determine a model's accuracy, it displays the overall discrepancy between predicted values and observed values

The lower the RSS, the better

The calculation for RSS in this model, first one must calculate Y_hat
X = 1
4 + 1*1 = 5
X = 2
4 + 1*2 = 6
X = 3
4 + 1*3 = 7
X = 4
4 + 1*4 = 8
X = 5

4 + 1*5 = <u>9</u>

Then subtract the Y_hat values from the actual values in the table
X = 1
3 - 5 = <u>-2</u>
X = 2
5 - 6 = <u>-1</u>
X = 3
7 - 7 = <u>0</u>
X = 4
9 - 8 = <u>1</u>
X = 5
11 - 9 = <u>2</u>

Then, square all values and sum them
$(-2)^2 = 4$
$(-1)^2 = 1$
$0^2 = 0$
$1^2 = 1$
$2^2 = 4$
RSS = 4 + 1 + 0 + 1 + 4
RSS = 10

Another measure to determine model accuracy is MSE or Mean Squared Error. This is calculated by dividing the RSS by the number of observations. For the scenario above, this would be

RSS/n = MSE
10/5 = MSE
MSE = 2

If one were to improve the model, one could use a multi-linear or non-linear regression model that considers multiple features like user demographics. One could confirm findings of a higher accuracy by comparing RSS values and MSE

- Primary Advantages
  - Simplicity
    - Straightforward
    - Simple mathematical approach
  - Easily Interpretable
    - Coefficients represent the variable relationship
  - Efficiency
    - Light computational load in relativity to more complex models
    - Quicker training on large datasets

- ○ Valuable insight on feature importance
  - ■ Given coefficient values, impact significance can be determined
- ○ Works better on smaller sample sizes than more complex models
- ● Primary Disadvantages
  - ○ Assumption of Linearity
    - ■ If a relationship is potentially non-linear, the model will be inaccurate
  - ○ Assumption of Independence
    - ■ Non-independence between variables can lead to inaccuracies and potentially invalidate conclusions
  - ○ Poor for complex relationships
    - ■ Complex patterns can be hard to discover through a Linear Regression.
  - ○ Susceptible to Overfitting or Underfitting
    - ■ Underfitting: Not enough data to capture the relationship
    - ■ Overfitting: too much data that causes noise/distortion to model accuracy
  - ○ Very susceptible to outliers
    - ■ Outliers can heavily distort coefficients
  - ○ Potentially requires feature engineering
    - ■ Fitting the data with the best features may need the implementation of other processes
- ● Common uses
  - ○ Predictive Analysis
    - ■ Forecasting for sales
    - ■ Forecasting demand
    - ■ Financial trend analysis
  - ○ Relationship Analysis
    - ■ Measures strength of correlation between two variables
  - ○ Risk assessment
    - ■ In insurance, it can predict expected claims given age, location, health, etc
  - ○ Price evaluation
    - ■ Given stats on competing products, one can determine the fair value of their own product
  - ○ Performance evaluation
    - ■ Can analyze correlations between performance and variables like hours of training
  - ○ Quality Control
    - ■ Helps find factors that are most tied to product quality
  - ○ Economic analysis
    - ■ Relationships between varying economic indicators