

Logistic regression

Logistic regression uses an equation as the representation, very much like linear regression.

Input values (x) are combined linearly using weights or coefficient values (referred to as W) to predict an output value (y). A key difference from linear regression is that the output value being modeled is a binary values (0 or 1) rather than a continuous value.

$$y^w(x) = \frac{1}{1 + \exp(-(w_0 + w_1x_1 + \dots + w_px_p))}$$

Multiclass Logistic regression using Softmax

The softmax function, also known as softargmax or normalized exponential function, is a generalization of the logistic function to multiple dimensions. It is used in multinomial logistic regression and is often used as the last activation function of a neural network.

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i=1, \dots, K \text{ and } z=(z_1, \dots, z_K) \in \mathbb{R}^K$$

Here K is the number of class and each z_i is calculated using

$$z_i = w_0 + w_1x_1 + \dots + w_px_p$$

Dataset

The dataset is available at

https://github.com/mishravipul/data/raw/main/obesity_data.csv

Original

Source: <https://archive.ics.uci.edu/ml/datasets/Estimation+of+obesity+levels+based+on+eating+habits+and+physical+condition+>

This dataset include data for the estimation of obesity levels in individuals from the countries of Mexico, Peru and Colombia, based on their eating habits and physical condition. The data contains 17 attributes and 2111 records, the records are labeled with the class variable NObesity (Obesity Level), that allows classification of the data using the values of Insufficient Weight, Normal Weight, Overweight Level I, Overweight Level II, Obesity Type I, Obesity Type II and Obesity Type III.

Features (X)

1. Gender {Female, Male}
2. Age {numeric}
3. Height {numeric}

4. Weight {numeric}
5. family_history_with_overweight: Has a family member suffered or suffers from overweight? {yes,no}
6. FAVC : Do you eat high caloric food frequently? {yes,no}
7. FCVC : Do you usually eat vegetables in your meals? {numeric}
8. NCP : How many main meals do you have daily? {numeric}
9. CAEC : Do you eat any food between meals? {no,Sometimes,Frequently,Always}
10. SMOKE : Do you smoke? {yes,no}
11. CH2O : How much water do you drink daily? {numeric}
12. SCC : Do you monitor the calories you eat daily? {yes,no}
13. FAF : How often do you have physical activity? {numeric}
14. TUE : How much time do you use technological devices such as cell phone, videogames, television, computer and others? {numeric}
15. CALC : how often do you drink alcohol? {no,Sometimes,Frequently,Always}
16. MTRANS : Daily Transportation {Automobile,Motorbike,Bike,Public_Transportation,Walking}

Take a look above at the source of the original dataset for more details.

Target (y)

17. NObeyesdad
{Insufficient_Weight,Normal_Weight,Overweight_Level_I,Overweight_Level_II,Obesity_Type_I,Obesity_Type_II,Obesity_Type_III}

Objective

To gain understanding of multiclass classification using logistic regression through implementing the model from scratch

Tasks

- Download and load the data
- Add intercept column with all values=1
- Feature transformation:
 - Convert 'Gender' column to numbers where 'Female' is 1 and 'Male' is 0

- Convert yes/no columns
['family_history_with_overweight','FAVC','SMOKE','SCC'] to 1/0
- One-Hot encode 'MTRANS', and 'NObeyesdad' columns. *Note:* One-hot encoding class/target variable is required for comparing binary predictions during training.
- Label encode 'CAEC', and 'CALC' columns
- Since the features have relatively different ranges, normalize the dataset
- Define X matrix (independent features) and y matrix (target features) as numpy arrays
- Print the shape and datatype of both X and y
- Split the dataset into 80% for training and rest 20% for testing (sklearn.model_selection.train_test_split function)
- Follow logistic regression class and fill code where highlighted:
 - Write softmax function to predict probabilities for all classes
 - Write cross entropy loss function
 - Write fit function where gradient descent is implemented
 - Write predict_proba function where we predict probabilities for input data
 - Write predict function to select single class for given input from probabilities
- Train the model
- Write function for calculating accuracy
- Compute accuracy on train and test data

keyboard_arrow_down

Further Fun

- Play with learning rate and max_iterations
- Preprocess data with different feature scaling methods (i.e. scaling, normalization, standardization, etc) and observe accuracies on both X_train and X_test
- Train model on different train-test splits such as 60-40, 50-50, 70-30, 80-20, 90-10, 95-5 etc. and observe accuracies on both X_train and X_test

- Shuffle training samples with different random seed values in the `train_test_split` function. Check the model error for the testing data for each setup.
- Print other classification metrics such as:
 - classification report (`sklearn.metrics.classification_report`),
 - confusion matrix (`sklearn.metrics.confusion_matrix`),
 - precision, recall and f1 scores (`sklearn.metrics.precision_recall_fscore_support`)

Helpful links

- Multiclass Logistic Regression from scratch: https://gluon.mxnet.io/chapter02_supervised-learning/softmax-regression-scratch.html
- Softmax tutorial : <http://deeplearning.stanford.edu/tutorial/supervised/SoftmaxRegression/>
- Softmax function detailed history: <https://medium.com/data-science-bootcamp/understand-the-softmax-function-in-minutes-f3a59641e86d>
- Understand gradients for cross entropy loss: https://rstudio-pubs-static.s3.amazonaws.com/337306_79a7966fad184532ab3ad66b322fe96e.html
- OnevsRest (OVR) strategy for multiclass classification: <https://medium.com/analytics-vidhya/logistic-regression-from-scratch-multi-classification-with-onevsall-d5c2acf0c37c>