

Logistic Regression with Rust

Arpit Bhat
501808

Archit Bhonsle
501810

Shalmali Kulkarni
501835

1 Problem Statement

We aim to create a desktop application using GTK, a toolkit to create GUI apps that guides use through the a simplified version of the machine learning pipeline. The algorithm we'll be using is logistic regression which is used in various simple classification tasks.

2 Dataset

The dataset we've chosen is reduced version of a larger dataset with 76 attributes found here: <https://archive.ics.uci.edu/ml/datasets/Heart+Disease>. This reduced dataset has 14 attributes (described in detail in table 1) and one target variable which indicates whether heart disease is present or not. The reduced dataset has numeric, discrete, continuous and binary attributes all of which can be represented as numbers which in our case are 64-bit floating point values. This is an assumption that helps us optimize our model and skip certain steps like mapping nominal attributes to unique numbers.

Column	Attribute
age	Age in years
sex	Sex of person: 0 is female, 1 is male
cp	Chest pain type
trestbps	Resting blood pressure (in mm Hg)
chol	Serum cholesterol in mg/dl
fbs	1 if fasting blood sugar > 120 mg/dl else 0
restecg	Resting electrocardiograph results
thalch	Maximum heart rate achieved (in mm Hg)
exang	1 if exercise induced angina else 0
oldpeak	ST depression induced by exercise relative to rest
slope	The slope of the peak exercise ST segment
ca	Number of major vessels (0-3) colored by fluoroscope
thal	3 = normal; 6 = fixed defect; 7 = reversable defect
target	1 if heart disease is present 0 otherwise

Table 1: Attributes in the heart dataset

3 Model

The model we use here is logistic regression. It is a model used for binary classification problems like the heart disease problem we are tackling here. A fundamental part of this model is the sigmoid function which maps an arbitrary number of real values to a probability between 0 and 1:

$$\frac{1}{1 + e^{-x}}$$

The cost function we use is:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right)$$

And thus in gradient descent, which is basically applying this repeatedly:

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta)$$

After resolving the derivative using calculus the equation becomes:

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m x_j^{(i)} (h_{\theta}(x^{(i)}) - y^{(i)})$$

4 Creating the model and testing

The core code for the model is encapsulated inside the following `forward_backward` function:

```
fn forward_backward(
  weights: &Array2<f64>,
  bias: &f64,
  x_train: &Array2<f64>,
  y_train: &Array2<f64>,
) -> (f64, Array2<f64>, f64) {
  // forward
  let y_head: Array2<f64> = sigmoid(weights.t().dot(x_train)
    .mapv(|x| x + bias));
  let loss = (y_train * &(y_head.mapv(|z| z.ln()))
    + &((y_train.mapv(|z| 1. - z)) * &y_head.mapv(|z| (1. - z).ln()))
    .mapv(|z| -1. * z);
  let cost = loss.sum() / x_train.ncols() as f64;
```

```

// backward
let d_weights = (x_train.dot(&(&y_head - y_train).t()))
    .mapv(|z| z / x_train.ncols() as f64);
let d_bias = (&y_head - y_train).sum() / x_train.ncols() as f64;

(cost, d_weights, d_bias)
}

```

Here, `y_head` is $h_{\theta}(x^{(i)})$ or the hypothesis. `loss` is $J(\theta)$ and `cost` is the average loss. This comprises the “forward” step. In the “backward” step we calculate the change in the weights, `d_weights` and the change in bias as `d_bias`.

These values are then used by the `update` function which tweaks the weights and the biases for the specified `iterations`. Thus, the model is trained.

```

fn update(
    weights: Array2<f64>,
    bias: f64,
    x_train: Array2<f64>,
    y_train: Array2<f64>,
    learning_rate: f64,
    iterations: usize,
) -> (Vec<f64>, Array2<f64>, f64) {
    let mut costs: Vec<f64> = Vec::new();
    let mut weights = weights;
    let mut bias = bias;

    for _ in 0..iterations {
        let (cost, d_weight, d_bias) = forward_backward(
            &weights, &bias, &x_train, &y_train);
        weights -= &d_weight.mapv(|x| learning_rate * x);
        bias -= learning_rate * d_bias;

        costs.push(cost);
    }

    (costs, weights, bias)
}

```

5 Using the application

The application has three pages. Each page has represents one section of the machine learning pipeline.

1. In the first step, as shown in figure 1, we select our dataset. We can review it in the tabular view and inspect its various attributes.
2. After that we perform preprocessing on this data. As one may observe in the tabular view some attributes go in the high 100s while some are simple binary values, 0 and 1. If we give this data to the model it may cause problems. To counter this we normalize the data. This is demonstrated in figure 2.
3. In the third step, as shown in the figure 3, we actually train our model. We specify the learning rate and the number of iterations. Once trained a graph of the model's cost as it goes through various iterations is graphed. Based on this we can evaluate the model. Once trained we can test the model against the testing data. Then the various metrics like accuracy, precision, recall and f1-score are shown.

gtk-rs-experiment
Choose a file
Must be a csv file containing only numeric data

heart.csv

age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
63.000000	1.000000	3.000000	145.000000	233.000000	1.000000	0.000000	150.000000	0.000000	2.300000	0.000000	0.000000	1.000000	1.000000
37.000000	1.000000	2.000000	130.000000	250.000000	0.000000	1.000000	187.000000	0.000000	3.500000	0.000000	0.000000	2.000000	1.000000
41.000000	0.000000	1.000000	130.000000	204.000000	0.000000	0.000000	172.000000	0.000000	1.400000	2.000000	0.000000	2.000000	1.000000
56.000000	1.000000	1.000000	120.000000	236.000000	0.000000	1.000000	178.000000	0.000000	0.800000	2.000000	0.000000	2.000000	1.000000
57.000000	0.000000	0.000000	120.000000	354.000000	0.000000	1.000000	163.000000	1.000000	0.600000	2.000000	0.000000	2.000000	1.000000
57.000000	1.000000	0.000000	140.000000	192.000000	0.000000	1.000000	148.000000	0.000000	0.400000	1.000000	0.000000	1.000000	1.000000
56.000000	0.000000	1.000000	140.000000	294.000000	0.000000	0.000000	153.000000	0.000000	1.300000	1.000000	0.000000	2.000000	1.000000
44.000000	1.000000	1.000000	120.000000	263.000000	0.000000	1.000000	173.000000	0.000000	0.000000	2.000000	0.000000	3.000000	1.000000
52.000000	1.000000	2.000000	172.000000	199.000000	1.000000	1.000000	162.000000	0.000000	0.500000	2.000000	0.000000	3.000000	1.000000
57.000000	1.000000	2.000000	150.000000	168.000000	0.000000	1.000000	174.000000	0.000000	1.600000	2.000000	0.000000	2.000000	1.000000
54.000000	1.000000	0.000000	140.000000	239.000000	0.000000	1.000000	160.000000	0.000000	1.200000	2.000000	0.000000	2.000000	1.000000
48.000000	0.000000	2.000000	130.000000	275.000000	0.000000	1.000000	139.000000	0.000000	0.200000	2.000000	0.000000	2.000000	1.000000
49.000000	1.000000	1.000000	130.000000	266.000000	0.000000	1.000000	171.000000	0.000000	0.600000	2.000000	0.000000	2.000000	1.000000
64.000000	1.000000	3.000000	110.000000	211.000000	0.000000	0.000000	144.000000	1.000000	1.800000	1.000000	0.000000	2.000000	1.000000
58.000000	0.000000	3.000000	150.000000	283.000000	1.000000	0.000000	162.000000	0.000000	1.000000	2.000000	0.000000	2.000000	1.000000
50.000000	0.000000	2.000000	120.000000	219.000000	0.000000	1.000000	158.000000	0.000000	1.600000	1.000000	0.000000	2.000000	1.000000
58.000000	0.000000	2.000000	120.000000	340.000000	0.000000	1.000000	172.000000	0.000000	0.000000	2.000000	0.000000	2.000000	1.000000
66.000000	0.000000	3.000000	150.000000	226.000000	0.000000	1.000000	114.000000	0.000000	2.600000	0.000000	0.000000	2.000000	1.000000
43.000000	1.000000	0.000000	150.000000	247.000000	0.000000	1.000000	171.000000	0.000000	1.500000	2.000000	0.000000	2.000000	1.000000
69.000000	0.000000	3.000000	140.000000	239.000000	0.000000	1.000000	151.000000	0.000000	1.800000	2.000000	2.000000	2.000000	1.000000
59.000000	1.000000	0.000000	135.000000	234.000000	0.000000	1.000000	161.000000	0.000000	0.500000	1.000000	0.000000	3.000000	1.000000
44.000000	1.000000	2.000000	130.000000	233.000000	0.000000	1.000000	179.000000	1.000000	0.400000	2.000000	0.000000	2.000000	1.000000

Preprocessing

Figure 1: Choose Page

gtk-rs-experiment

Processing

Apply transformations on the data

Normalize

age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0.708333	1.000000	1.000000	0.481132	0.244292	1.000000	0.000000	0.603053	0.000000	0.370968	0.000000	0.000000	0.333333	1.000000
0.166667	1.000000	0.666667	0.339623	0.283105	0.000000	0.500000	0.885496	0.000000	0.564516	0.000000	0.000000	0.666667	1.000000
0.250000	0.000000	0.333333	0.339623	0.178082	0.000000	0.000000	0.770992	0.000000	0.225806	1.000000	0.000000	0.666667	1.000000
0.562500	1.000000	0.333333	0.245283	0.251142	0.000000	0.500000	0.816794	0.000000	0.129032	1.000000	0.000000	0.666667	1.000000
0.583333	0.000000	0.000000	0.245283	0.520548	0.000000	0.500000	0.702290	1.000000	0.096774	1.000000	0.000000	0.666667	1.000000
0.583333	1.000000	0.000000	0.433962	0.150685	0.000000	0.500000	0.587786	0.000000	0.064516	0.500000	0.000000	0.333333	1.000000
0.562500	0.000000	0.333333	0.433962	0.383562	0.000000	0.000000	0.625954	0.000000	0.209677	0.500000	0.000000	0.666667	1.000000
0.312500	1.000000	0.333333	0.245283	0.312785	0.000000	0.500000	0.778626	0.000000	0.000000	1.000000	0.000000	1.000000	1.000000
0.479167	1.000000	0.666667	0.735849	0.166667	1.000000	0.500000	0.694656	0.000000	0.080645	1.000000	0.000000	1.000000	1.000000
0.583333	1.000000	0.666667	0.528302	0.095890	0.000000	0.500000	0.786260	0.000000	0.258065	1.000000	0.000000	0.666667	1.000000
0.520833	1.000000	0.000000	0.433962	0.257991	0.000000	0.500000	0.679389	0.000000	0.193548	1.000000	0.000000	0.666667	1.000000
0.395833	0.000000	0.666667	0.339623	0.340183	0.000000	0.500000	0.519084	0.000000	0.032258	1.000000	0.000000	0.666667	1.000000
0.416667	1.000000	0.333333	0.339623	0.319635	0.000000	0.500000	0.763359	0.000000	0.096774	1.000000	0.000000	0.666667	1.000000
0.729167	1.000000	1.000000	0.150943	0.194064	0.000000	0.000000	0.557252	1.000000	0.290323	0.500000	0.000000	0.666667	1.000000
0.604167	0.000000	1.000000	0.528302	0.358447	1.000000	0.000000	0.694656	0.000000	0.161290	1.000000	0.000000	0.666667	1.000000
0.437500	0.000000	0.666667	0.245283	0.212329	0.000000	0.500000	0.664122	0.000000	0.258065	0.500000	0.000000	0.666667	1.000000
0.604167	0.000000	0.666667	0.245283	0.488584	0.000000	0.500000	0.770992	0.000000	0.000000	1.000000	0.000000	0.666667	1.000000
0.770833	0.000000	1.000000	0.528302	0.228311	0.000000	0.500000	0.328244	0.000000	0.419355	0.000000	0.000000	0.666667	1.000000
0.291667	1.000000	0.000000	0.528302	0.276256	0.000000	0.500000	0.763359	0.000000	0.241935	1.000000	0.000000	0.666667	1.000000
0.833333	0.000000	1.000000	0.433962	0.257991	0.000000	0.500000	0.610687	0.000000	0.290323	1.000000	0.500000	0.666667	1.000000
0.625000	1.000000	0.000000	0.386792	0.246575	0.000000	0.500000	0.687023	0.000000	0.080645	0.500000	0.000000	1.000000	1.000000
0.312500	1.000000	0.666667	0.339623	0.244292	0.000000	0.500000	0.824477	1.000000	0.064516	1.000000	0.000000	0.666667	1.000000

Model

Figure 2: Preprocessing Page



Figure 3: Model Page

6 Analyzing the results

The results of various values of the hyperparameters, learning rate and iterations are shown in table 2.

LR	Iterations	Accuracy	Precision	Recall	F1-Score
0.1	100	79.1%	90.2%	76.7%	82.9%
0.1	1000	80.2%	90.2%	78.0%	83.6%
1	100	80.1%	90.2%	78.0%	83.6%
1	1000	83.5%	92.2%	81.0%	86.2%
10	100	74.2%	68.6%	83.3%	75.3%
10	1000	76.9%	70.6%	85.7%	77.4%

Table 2: Metrics for various hyperparameters

As we can see, the model with the learning rate of 1 and 1000 iterations performs the best. The learning rate of 0.1 is too slow and it takes 1000 to match the performance of the (1, 1000) model. The learning rate of 10 is too fast which is also apparent from its cost vs iterations graph. Ideally, we'd like a mechanism that reduces the learning rate whenever the cost stops decreasing. This is known as an "adaptive learning rate".

7 Setup

Regardless of the operating system you're on you might need a dataset to test this on. The app only supports datasets with numeric columns so other datasets might crash the app. To download `heart.csv` used in this report, click on this link.

7.1 Windows

1. Go here.
2. In the downloaded folder double click on the `gtk-rs-experiment.exe`.

7.2 Linux

This may or may not work based on the libraries installed on your system.

1. Go to this link.
2. Open the location of the downloaded executable in the file explorer and right click and open a terminal there.
3. Enter the command `chmod +x ./linux-executable` followed by `./linux-executable`.