

Assignment 2

Mahir Hambiralovic

February 21, 2022

1 Summary

A Computer Exercise on classification of a text dataset was completed. The exercises were to classify two corpuses of text, one in English and one in French of *Salammbô* using linear regression, perceptrons and logistic regression, all using gradient descent.

2 My implementations

2.1 Linear Regression

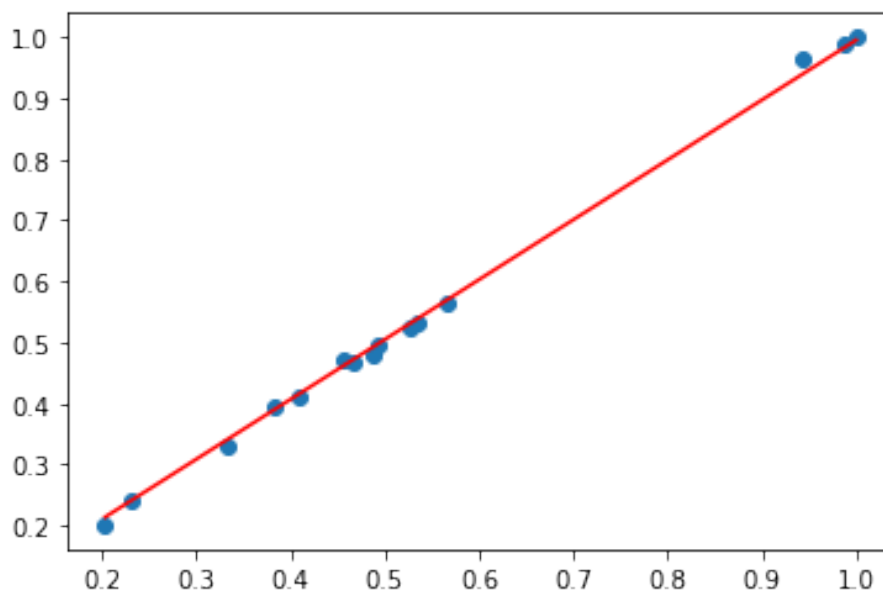
The task was to fit the number of A's in a text, given the total number of characters in the corpus using Linear Regression.

2.1.1 Batch descent

```
1 def fit_batch(X, y, alpha, w,
2               epochs=500,
3               epsilon=1.0e-5):
4     w = np.array(w)
5     epoch = 1
6     error = 10000
7
8     while error > epsilon and epoch <= epochs:
9         (b, m) = w
10        y_pred = m*X + b
11
12        loss_m = sum(X * (y - y_pred))
13        loss_b = sum(y - y_pred)
14
15        m = m + (alpha / len(y)) * loss_m
16        b = b + (alpha / len(y)) * loss_b
17
18        epoch += 1
19        error = np.linalg.norm(y - y_pred)
20        w = (b, m)
21
22    print('stopped at Error:', error, 'Epoch:', epoch)
23    return w
24
25 (b, m) = fit_batch(X, y, 0.01, [random.random(), random.random()])
```

2.1.2 Stochastic descent

```
1 # Write your code here
2 def fit_stoch(X, y, alpha, w,
3               epochs=500,
4               epsilon=1.0e-5):
5     epoch = 1
6     error = 10000
7     idx = list(range(len(X)))
8     while error > epsilon and epoch <= epochs:
9         errors = []
10        random.shuffle(idx)
11        for i in idx:
12            x = X[i]
13            y_i = y[i]
14
15            (b, m) = w
16            y_i_pred = m*x + b
17
18            loss_m = x * (y_i - y_i_pred)
19            loss_b = y_i - y_i_pred
20
21            m = m + alpha * loss_m
22            b = b + alpha * loss_b
23
24            y_pred = m*X + b
25            error = np.linalg.norm(y - y_pred)
26            w = (b, m)
27
28        # error = np.linalg.norm(errors)
29        epoch += 1
30    return w
31
32 (b, m) = fit_stoch(X, y, 0.01, [random.random(), random.random()])
```



2.2 Perceptron

The task was to classify whether a text was in French or in English based on the number of letter's and the number of A's in the text.

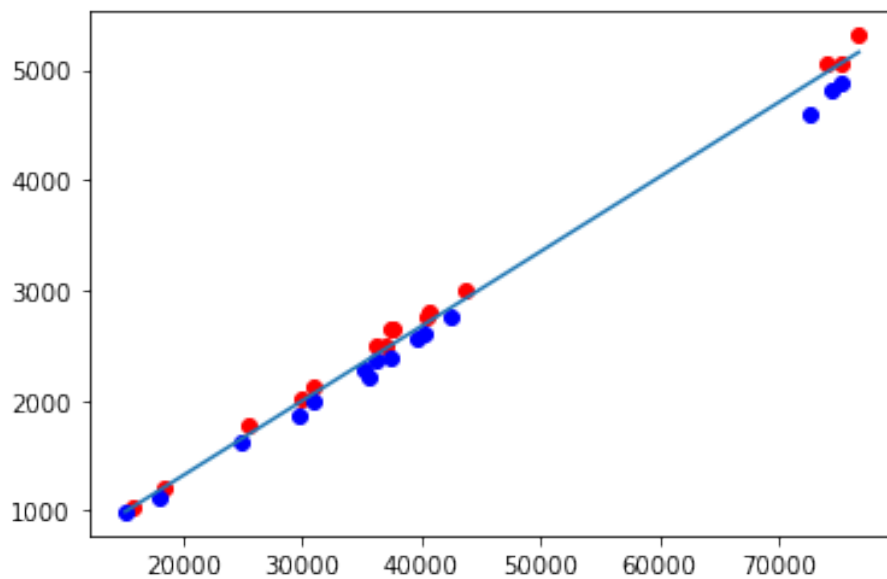
```
1 def predict(X, w):
2     y_pred = np.matmul(w, X)
3     if y_pred >= 0:
4         return 1
5     else:
6         return 0
7
8 def fit_stoch(X, y,
9               alpha=0.01,
10              epochs=1000,
11              max_misclassified=0,
12              verbose=True):
13     epoch = 1
14     misclassified = 0
15     correct_classified = 0
16
17     idx = list(range(len(X)))
18     w = [random.random() for _ in range(len(X[0]))]
19
20     while epoch <= epochs:
21         if max_misclassified != 0:
22             if misclassified > max_misclassified:
23                 break
24
25         random.shuffle(idx)
26
27         for i in idx:
28             x_i = X[i]
29             y_i = y[i]
30
31             y_i_pred = predict(x_i, w)
32             error = (y_i - y_i_pred)
33
34             for j in range(len(w)):
35                 w_j = w[j]
36                 x_j = x_i[j]
37                 w[j] += alpha*(error)*x_j
38
39             if abs(error):
40                 misclassified += 1
41             else:
42                 correct_classified += 1
43
44         epoch += 1
45
46     correct_classified, 'epoch', epoch)
47     return w
48
49 def leave_one_out_cross_val(X, y, fitting_function):
50     n = len(X)
51     n_correct = 0
52
53     for fold in range(n):
54         X_train = X[:fold] + X[fold + 1:]
55         X_test = X[fold]
56
57         y_train = y[:fold] + y[fold + 1:]
58         y_test = y[fold]
```

```

59     w = fitting_function(X_train, y_train)
60
61     y_pred = predict(X_test, w)
62     correct = not abs(y_test - y_pred)
63     n_correct += correct
64
65     print('Fold', fold+1, 'on', n, ': weights:', w, 'Correct.')
66     if correct else 'Wrong.')
67     # print('Fold', fold, 'on 30: epochs:', epochs, 'weights:',
68         w)
69     print('Correct', n_correct, '/', n)
70     return n_correct / n

```

With this implementation, a cross validation accuracy of 86.7% was achieved. Restored weights [7.083916039112965, -0.014127309008252075, 0.20878803577563473]. I would consider these good results for practically a linear regression implementation of the ANN (we can see in the decision boundary and the data points in the figure, that a straight line cannot fit the points perfectly). With one or two more perceptrons, we should be able to fit the curve perfectly, but risk overfitting.



2.2.1 Logistic Regression

Implementing the same task as above, but using Logistic Regression.

```

1 def logistic(x):
2     try:
3         return 1 / ( 1 + math.exp(-x) )
4     except OverflowError:
5         x_rounded = 500 * x/abs(x)
6         return 1 / ( 1 + math.exp(-x_rounded) )
7
8 def predict_proba(X, w):
9     return logistic(np.matmul(w, X))
10

```

```

11 def predict(X, w):
12     y_pred = predict_proba(X, w)
13     return 1 if y_pred >= 0.5 else 0
14
15 def fit_stoch(X, y, alpha=0.5,
16               epochs=4000,
17               epsilon=1.0e-4,
18               verbose=False):
19     epoch = 1
20
21     idx = list(range(len(X)))
22     w = [random.random() for _ in range(len(X[0]))]
23
24     while epoch <= epochs:
25         random.shuffle(idx)
26
27         for i in idx:
28             x_i = X[i]
29             y_i = y[i]
30
31             y_i_pred = predict_proba(x_i, w)
32             # error = (y_i - y_i_pred)
33
34             for j in range(len(w)):
35                 # w_j = w[j]
36                 x_j = x_i[j]
37                 w[j] += alpha * (y_i - y_i_pred) * y_i_pred * (1 -
y_i_pred) * x_j
38
39
40         epoch += 1
41
42     return w
43
44 # Write your code here
45 def leave_one_out_cross_val(X, y, fitting_function):
46     n = len(X)
47     n_correct = 0
48
49     for fold in range(n):
50         X_train = X[:fold] + X[fold + 1:]
51         X_test = X[fold]
52
53         y_train = y[:fold] + y[fold + 1:]
54         y_test = y[fold]
55
56         w = fitting_function(X_train, y_train)
57
58         y_pred = predict(X_test, w)
59         # print(y_pred)
60         correct = not abs(y_test - y_pred)
61         n_correct += correct
62
63         print('Fold', fold+1, 'on', n, ': weights:', w, 'Correct.'
if correct else 'Wrong.')
64         # print('Fold', fold, 'on 30: epochs:', epochs, 'weights:',
w)
65     print(n_correct, n)
66     return n_correct / n

```

With this implementation, a cross validation accuracy of 86.7% was achieved.
 Restored weights [-0.1618444091889996, -0.0005504545172716554, 0.008356119511984951]

3 Improvements

In the above tasks, I decided to use stochastic descent where possible. One improvement to this approach was to use larger batches than 1. Furthermore, as discussed in the perceptron section, the ANN approach could reach a perfect score on the data if it would have been allowed to be more complex. However, this would have let it overfit, especially seeing that we tested on the training data.

4 Dissertation of An overview of gradient descent optimization algorithms

Ruth writes, in his article "An overview of gradient descent optimization algorithms" an overview, summary and motivation of the various approaches to gradient descent, the main algorithm for every deep learning implementation. The main gradient descent variants to chose from are explained in the following sections.

4.1 Gradient descent variants

4.1.1 Batch gradient descent

The traditional and "vanilla" approach to gradient descent, calculating the gradient for the entire dataset and calculating the step to take in each step for a fixed number of epochs. Some draw-backs of this approach is that it is very time consuming, as it must compute the gradient of the entire dataset in each epoch. Also it does not support online-training.

4.1.2 Stochastic gradient descent

A "stochastic" approach for updating the weights by creating a gradient for each training example. This approach does however lead to many overfitted and biased steps, that are each perfect to only one training example.

4.1.3 Mini-batch gradient descent

A compromise of the above two, that predefines a batch-size of data to work with, which is trained upon. Thus online training is supported, faster training, and fewer biased "missteps" are taken.

4.1.4 Other optimization methods

Later in the paper, Ruth discusses a myriad of different methods used to optimize the descent method in various ways. The optimizations are used for e.g. avoiding local minimas, speeding up learning and improve learning on sparse data. The final evaluation of the methods show that there can be a large difference in the speed of convergence between the different methods. While not all optimizations were created equal, all of them seem to have some improvement to the vanilla approach, which Ruth criticizes as being overused in many papers. Ruths' paper concludes that Adam is generally the overall best choice.