

Auxiliary Feature Learning for Small Dataset Regularization (January 2019)

Jordà. Jaume, *Data Science and Decision Making, University of Essex*

Abstract – An autoencoder is a neuronal network able to generate new data.

The autoencoder can learn by a few data and then, with different machine learning methods like lineal regression, classification, etc the neural network can take some data of the same type and classify it, obviously, with some error.

With three different small datasets we will understand how an autoencoder works.

I.INTRODUCTION

The autoencoder consists of two parts, the encoder and the decoder, the encoder compress the input data in a latent variable space and the decoder try to rebuild the input based in previous information that it recollects.

The purpose of this paper is the selection of the datasets to train the autoencoder, it is done based on the number of attributes that every dataset has, because with fewer attributes we will understand it better. Also, we need to use small datasets (number of instances) or if it is not possible, we must choose a little sample of the dataset.

The datasets are provided by *UCI Machine Learning Repository* and the choose datasets are *flags of the world*, *wines*, and *zoo*.

II.BACKGROUND

To start working on these datasets we must know the previous works that the creators of the datasets did with it.

The *flags* and *zoo* datasets were used only in the document that Forsyth's, the creator of the datasets, used to explain what is BEAGLE and how he program it called "BEAGLE User Notes" and we can find it [here](#).

The *wines* dataset is created by Forina, M from the Institute of Pharmaceutical and Food Analysis and Technologies, and it was used in 1992 by S. Aeberhard, D. and O. de Vel of University of North Queensland to compare different classifiers in high dimensional settings and for a paper called "The Classification performance of RDA".

III.METHODOLOGY

Before start working it is necessary how are the datasets.

The *flags* dataset consists on 194 instances of 30 attributes. The attributes are the name, the landmass, the graphic quadrant zone, the area, the population, the language, the religion, the bars in the flag, the stripes in the flag, the number of different colours, (red, green, blue, gold/yellow, white, black, orange) as a boolean if is present or not in the flag, the predominant colour in the flag, the number of circles, the number of upright crosses, the number of diagonal crosses, the number of quartered sections, the number of sun or star symbols, is crescent moon present, is triangle present, if there are inanimate symbology (e.g., a boat), if there are animate image like eagle, if there are letters or writing, if there is colour in the top-left corner, and finally, if there is colour in the bottom-left. Is possible to see that this dataset is a little bit large but before start working we have to select a sample.

The *wines* dataset is distributed in three classes, the class 1 have 59 instances, the class 2 have 71 and the class 3 have 48, and every instance have 13 attributes. The attributes are the alcohol, the malic acid, the ash, the alkalinity of the ash, the magnesium, the total phenols, the flavonoids, the nonflavanoid phenols, the proanthocyanins, the color intensisty, the hue, the OD280/OD315 of diluted wines, and the proline.

The *zoo* dataset have 101 instances of 18 attributes, the animal name, 15 boolean attributes, and two numeric values. The boolean attributes make sense to if the animal is/have/do something like hair, feathers, eggs, milk, airborne, aquatic, predator, toothed, backbone, breathes, venomous, fins, tail, domestic and catsize, and the numeric values are the number of legs (set of values: {0, 2, 4, 5, 6, 8}) or the type (integer in range [1, 7]).

IV. EXPERIMENTS

The experiment I did with the datasets is simply charge the data to the program and print it in the screen. The aim to do that is to be sure that the choosed datasets will be useful to learn how an autoencoder works.

In the images below we can see the piece of code that allows us to do it and how is shown in the screen with the *flag* dataset to show the example.

```
dataReader.py
1 import numpy as np
2 #import panda as pd
3 import matplotlib.pyplot as plt
4
5 # Open the data file in reading mode
6 f = open("flag.data", "r+")
7
8 # Save the information in the data variable
9 data = f.read()
10
11 # Show de data on the screen
12 print(data)
13
14 f.close # Close the data file
```

```
Qatar,5,1,11,0,8,2,0,0,2,0,0,0,1,0,1,brown,0,0,0,0,0,0,0,0,white,brown
Romania,3,1,237,22,0,0,3,0,7,1,1,1,0,1,red,0,0,0,0,2,0,0,1,1,blue,red
Rwanda,4,2,26,5,10,5,3,0,4,1,1,0,1,0,red,0,0,0,0,0,0,0,1,red,green
San-Marino,3,1,0,0,0,0,2,2,0,0,1,0,1,0,white,0,0,0,0,0,0,0,0,white,blue
Sao-Tome,4,1,0,0,0,0,0,3,4,1,0,1,0,1,green,0,0,0,2,0,1,0,0,green,green
Saudi-Arabia,5,1,2150,9,8,2,0,0,2,0,1,0,0,1,0,green,0,0,0,0,0,0,1,0,green,green
Senegal,4,4,196,0,3,2,3,0,3,1,1,0,1,0,green,0,0,0,0,1,0,0,0,green,red
Seychelles,4,2,0,0,1,1,0,0,3,1,1,0,0,1,red,0,0,0,0,0,0,0,0,red,green
Sierra-Leone,4,4,72,3,1,5,0,3,3,0,1,1,0,1,0,green,0,0,0,0,0,0,0,0,green,blue
Singapore,5,1,1,3,7,3,0,2,2,1,0,0,0,1,0,white,0,0,0,0,5,1,0,0,red,white
Solomon-Islands,6,2,30,0,1,1,0,0,4,0,1,1,1,1,0,green,0,0,0,5,0,1,0,0,blue,green
Somalia,4,1,0,0,3,7,5,10,2,0,0,2,0,0,blue,0,0,0,0,1,0,0,0,blue,blue
South-Africa,4,2,1221,29,6,1,0,3,5,1,1,1,0,1,orange,0,1,1,0,0,0,0,0,orange,blue
South-Korea,5,1,99,39,10,7,0,0,4,1,0,1,0,1,white,1,0,0,0,0,0,1,0,white,white
South-Yemen,5,1,288,2,0,2,0,3,4,1,0,0,1,1,red,0,0,0,0,1,0,1,0,red,black
Spain,3,4,505,38,2,0,0,3,2,1,0,0,1,0,red,0,0,0,0,0,0,0,0,red,red
Sri-Lanka,5,1,66,15,10,3,2,0,4,0,1,0,1,0,gold,0,0,0,0,0,0,0,1,0,gold,gold
St-Helena,4,3,0,0,1,1,0,0,7,1,1,1,1,1,blue,0,1,1,0,0,0,1,0,white,blue
St-Kitts-Nevis,4,1,0,0,1,1,0,0,0,1,1,0,1,1,0,green,0,0,0,0,2,0,1,0,green,red
St-Lucia,1,4,0,0,1,1,0,0,4,0,0,1,1,1,0,blue,0,0,0,0,0,0,1,0,blue,blue
St-Vincent,1,4,0,0,1,1,5,0,4,0,1,1,1,1,0,green,0,0,0,0,0,0,1,1,blue,green
Sudan,4,1,2586,10,2,0,3,4,1,1,0,0,1,1,red,0,0,0,0,0,0,1,0,red,black
Surinam,2,4,03,0,0,1,0,5,4,1,1,0,1,0,red,0,0,0,0,1,0,0,0,green,green
Swaziland,4,2,17,1,10,1,0,5,7,1,0,1,1,1,1,blue,0,0,0,0,0,0,1,0,blue,blue
Sweden,3,1,450,0,0,1,0,0,2,0,0,1,1,0,0,blue,0,1,0,0,0,0,0,0,blue,blue
Switzerland,3,1,41,0,4,1,0,0,2,1,0,0,0,1,0,red,0,1,0,0,0,0,0,0,red,red
Syria,5,1,185,10,0,2,0,3,4,1,1,0,0,1,1,red,0,0,0,2,0,0,0,0,red,black
Taiwan,5,1,36,10,7,0,0,0,0,1,0,1,0,1,0,red,1,0,0,1,1,0,0,0,red,red
Tanzania,4,2,945,10,10,0,0,0,0,4,1,1,0,1,0,green,0,0,0,0,0,0,1,0,green,blue
Thailand,5,1,514,49,10,3,0,5,3,1,0,1,0,1,0,red,0,0,0,0,0,0,0,0,red,red
Togo,4,1,57,2,3,7,0,5,4,1,1,0,1,1,0,green,0,0,0,1,1,0,0,0,red,green
Tonga,6,2,1,0,10,1,0,0,0,1,0,0,0,1,0,red,0,1,1,0,0,0,0,white,red
Trinidad-Tobago,2,4,5,1,1,1,0,0,3,1,0,0,0,1,1,red,0,0,0,0,0,1,0,0,white,white
Tunisia,4,1,164,7,0,2,0,0,2,1,0,0,0,1,0,red,1,0,0,0,1,1,0,0,red,red
Turkey,5,1,781,45,0,2,0,0,2,1,0,0,0,1,0,red,0,0,0,0,1,1,0,0,red,red
Turks-Cocos-Islands,1,4,0,0,1,1,0,0,0,0,1,1,1,1,0,blue,0,1,1,0,0,0,1,1,white,blue
Tuvalu,0,2,0,0,1,1,0,0,5,1,0,1,1,1,0,blue,0,1,1,9,0,0,0,0,white,blue
UK,5,1,84,1,0,2,1,3,4,1,1,0,0,1,1,0,green,0,0,0,0,0,0,0,0,red,black
Uganda,4,1,230,13,10,5,0,0,5,1,0,0,1,1,0,gold,1,0,0,0,0,0,1,0,black,red
UK,3,4,245,56,1,1,0,0,3,1,0,0,1,0,red,0,1,1,0,0,0,0,0,white,red
Uruguay,2,3,178,3,2,0,0,0,3,0,0,1,1,0,0,white,0,0,0,1,1,0,0,0,white,white
US-Virgin-Isles,1,0,0,0,1,1,0,0,0,0,1,1,1,1,0,white,0,0,0,0,0,0,1,1,white,white
USA,1,4,9363,231,1,1,0,13,3,1,0,1,0,1,0,white,0,0,0,1,50,0,0,0,0,blue,red
USSR,5,1,22482,274,5,0,0,0,2,1,0,0,1,0,red,0,0,0,0,1,0,0,1,0,red,red
Vanuatu,6,2,15,0,0,0,1,0,0,4,1,0,0,1,0,red,0,0,0,0,0,0,1,0,black,green
Vatican-City,3,1,0,0,0,0,2,0,4,1,0,0,1,1,0,gold,0,0,0,0,0,0,1,0,gold,white
Venezuela,2,4,912,15,2,0,0,3,7,1,1,1,1,1,red,0,0,0,0,7,0,0,1,1,gold,red
Vietnam,5,1,333,60,10,0,0,0,2,1,0,0,1,0,0,red,0,0,0,0,1,0,0,0,red,red
Western-Samoa,6,3,0,1,1,0,0,3,1,0,1,0,1,0,red,0,0,0,1,5,0,0,0,blue,red
Yugoslavia,3,1,256,22,0,0,0,3,4,1,0,1,1,1,0,red,0,0,0,0,1,0,0,0,blue,red
Zaire,4,2,905,28,10,5,0,0,4,1,1,0,1,0,1,green,1,0,0,0,0,0,1,1,green,green
Zambia,4,2,753,0,10,5,3,0,4,1,1,0,0,0,1,1,green,0,0,0,0,0,0,1,0,green,brown
Zimbabwe,4,2,391,8,10,5,0,7,5,1,1,0,1,1,1,0,green,0,0,0,0,1,0,1,1,0,green,green
```

V. DISCUSION

Once the datasets are selected it is important to select how much data is necessary to train the autoencoder and how much is for testing.

In this case the 80% will be used as training data and the remaining 20% for tests, that is, see if the autoencoder is well trained because it is about unsupervised training.

VI. CONCLUSION

Finally we can see that the datasets are enough for our proposes, so in the next paper we can see if our proposal it is good or not.

VII. NOTES

To see the github repository it is necessary to use this link:

https://github.com/jj18367/ce888_assignment1.git