

Improving product classification using generative recurrent networks

RODRIGUES, Marcos <<http://orcid.org/0000-0002-6083-1303>>

Available from Sheffield Hallam University Research Archive (SHURA) at:

<http://shura.shu.ac.uk/21693/>

This document is the author deposited version. You are advised to consult the publisher's version if you wish to cite from it.

Published version

RODRIGUES, Marcos (2018). Improving product classification using generative recurrent networks. In: PAPANIKOS, Gregory T., (ed.) Abstracts 2nd International Conference on Electrical Engineering, 23-26 July 2018, Athens, Greece. Athens Institute for Education and Research, 36-37.

Copyright and re-use policy

See <http://shura.shu.ac.uk/information.html>

[Athens Institute for Education and Research](#)

Abstract Submitting Form

Conference	A Stream on "Data Science", 23-26 July 2018, Athens, Greece
Title of Paper	Improving Product Classification using Generative Recurrent Networks
For more than one author, please copy and paste the following eight rows for each additional author.	
Title	<input checked="" type="checkbox"/> Dr <input type="checkbox"/> Mr <input type="checkbox"/> Ms <input type="checkbox"/> Other Specify: Professor
First Name	Marcos
Family Name	Rodrigues
Position	Professor of Computer Science
University/ Organization	Sheffield Hallam University
Country	UK
E-mail	m.rodrigues@shu.ac.uk
Telephone(s)	+44 (0)114 225 6911
Fax	+44 (0)114 225 6702

Abstract

The issue addressed in this paper is related to machine learning techniques for automatic classification of product descriptions. The problem arises when database entries do not perfectly match and so it is questionable whether a description is related or not to the same item, product, or service. A typical example is merging disparate databases that is required, for instance, when one business buys off a competitor. An obvious solution would be to train an AI system to perform classification. The problem is that AI deep learning networks require vast amounts of training data, normally in tens or hundreds of thousand samples and normally such data are not available. The specific classification problem we are addressing can be illustrated as follows.

Product	Category	Level 2	Level 3
Actimel Yogurt Drink 0.1% Fat Original 12x100g	Dairy	Yogurts	Actimel
Actimel Yogurt Drink 0.1% Fat Strawberry 8x100g	Dairy	Yogurts	Actimel
Actimel Yogurt Drink Blueberry 8x100g	Dairy	Yogurts	Actimel
Actimel Yogurt Drink Coconut 8x100g	Dairy	Yogurts	Actimel
Actimel Yogurt Drink Kids Strawberry and Raspberry 6x100g	Chilled	Easy Lunches	Lunchbox favourites

Note that while the first four records have been manually classified as 'Dairy', the last entry was classified as 'Chilled' (classification is accepted as correct for all entries). In order to learn the nuances of classification, an AI system needs a vast number of additional samples to be able to distinguish what characterizes Dairy and Chilled. Therefore, we have investigated network models to augment the training data set in a flexible but reliable way. The principle is to train a network with the objective of generating new data similar but not exactly the same as the input data. Validation of the newly generated data is performed by a second network which has been trained on the original data. A simple binary decision (yes/no) is output whether or not generated data has enough or acceptable similarity with the original data. Accepted data would eventually make part of an augmented training set, improving the network ability to classify unseen data. We designed and implemented a recurrent network with Keras, an open source neural network library written in Python. The network is based on the LSTM-Long-Short Term Memory model which has proved useful to a large number of problems with time dependencies. The encoding of product description is character-based so, once trained, the network outputs a character and tries to predict what the next character would be. With an appropriate training set to learn the structure of the data, such networks can output valid vectors. We set the network to train over 20 epochs outputting the description (with a limited number of characters) at the end of each epoch. At epoch 0 (before training) it can only output random characters:

R22QQQOOVVVV000000aa33aKTTTTTTTTTT**eLLLePPPCJl1mvao

At epoch 2, things start to get better as the net begins to learn to separate words properly:

X Crisps and snacks

	<p style="color: red;">Supermarket's Crisps and Crisps and Cream</p> <p>At epoch 3 the data now starts to resemble the training file with one description per line (ignoring the nonsense meaning of generated data such as chicken yogurt):</p> <p style="color: red;">Chilled > Fresh pasta and sauces > Fresh pasta</p> <p style="color: red;">British Chicken and Strawberry and Corner Yogurt 4x125g Dairy > Yogurts > Muller</p> <p style="color: red;">British Pork Sausages x8 200g Meat and fish > Fish and seafood > All fish and seafood</p> <p style="color: red;">Supermarket's British Pork Light and Coconut and Cheese</p> <p>Network outputs get increasingly better and, at the end of training, valid samples are generated for an augmented database. Note that the generated data are not the same as the original. The main outcome of such generative recurrent network is that it works for text generation, giving us the ability to generate valid data from a limited set of samples. In this paper, we provided a justification for using recurrent networks to solve a significant limitation of small data sets in deep learning. We also showed that LSTMs are a good solution to the problem together with character-based text encoding and these represent the state-of-the-art in recurrent neural networks. Future work involves improvements to the network design model and testing SimpleRNN or GRU-Gate Recurrent Unit in place of LSTMs and fine-tuning of network parameters.</p>
Keywords	AI, Deep Learning, Recurrent Networks

Please email to: atiner@atiner.gr as an attached file or fax it to +30 210 3634209