

Neural Networks Final Project:

Comparing different data pre-processing techniques for convolution neural networks

Krystian Wojcicki, 101001444

Michael Kuang, 101000485

COMP 4107, Fall 2018

1 Introduction

Natural Images is a dataset currently featured on *kaggle.com*[1]. The set is comprised of 6899 images from 8 distinct classes of airplane, car, cat, dog, flower, fruit, motorbike and person. Each class has a varying number of examples. We examined different techniques in pre-processing images and implemented them to analyze their effects on performance in a convolution neural network.

2 Problem Statement

This project attempts to address the problems: "Do different image resizing techniques make a difference to classification accuracy?" and "Do balanced datasets improve classification accuracy?".

3 Background

One problem we looked at is the Class Imbalance problem, and this is the problem in machine learning where we have great discrepancy in the number of examples between class data. In other words, there are far more examples of one class than another in a set of data. This is a problem because as we know, machine learning works best when we have a more uniform distribution of data. In this project, we examine the performance between two up-sampling techniques called Synthetic Minority Over sampling Technique (SMOTE) and ADaptive SYNthetic (ADASYN) as solutions to class imbalance.

Another problem we investigated is how image resizing can affect performance. Given a set of non-uniform sizes of images, we must resize them to some constant shape in order to feed examples into the network. We examined two methods; resize the image by scaling it, and crop or pad the image about the center. The problem with resizing an image is that we inevitably lose information as images are scaled down or add noise as they are scaled up, but we are able to retain the information more uniformly. Cropping the image will directly lose the information as we resize the image by cropping out the outer most part of the image, while padding the image with black pixels evenly about center will retain the image aspect ratio.

3.1 The Dataset

The *Natural Images* dataset contains 6899 distinct RGB images of varying sizes for 8 classes as described below:

- airplane: 727
- car: 968
- cat: 885
- dog: 702
- flower: 843
- fruit: 1000
- motorbike: 788
- person: 986

3.2 Machine Learning Libraries

Four libraries were used to implement our code, and create our neural network model: Tensorflow, scikit-learn, imbalanced-learn and OpenCV.

3.3 Methodology