

Traditional Maltese Food and Snacks Classifier

Keane Schembri

Institute of Information & Communication Technology

Malta College of Arts, Science & Technology

Corradino Hill

Paola PLA 9032

keane.schembri.a106267@mcast.edu.mt

Abstract—With Malta’s tourism sector on continuous growth annually and being a very small island with rich culture, a tourist can engage and interest themselves into our cultivation. Using a classification tool such as this study proposes makes it easier for a foreigner to learn basic concepts in Malta’s traditional food. This has been achievable using YOLOv3 neural network to detect and classify different articles of food, from cheeselet to deep-fried date roll’s, labelled in Maltese as ‘gbejniet’ and ‘imqaret’ consecutively. On a dataset of 350 images, the model has achieved over 71% classification accuracy on images with multiple objects to classify.

Index Terms—Convolutional Neural Network (CNN), YOLOv3, Algorithm, Image Annotation

I. INTRODUCTION

Using computer vision as one of deep learning methods, it is possible to classify various traditional Maltese food. The idea behind this study is to have a more user-friendly introduction and an easier access for foreigners to traditional Maltese food, how to identify different articles of food and their label. Hence this study discusses and explores the field of computer vision and the reasoning behind a regression-based YOLOv3 algorithm high speed’s performance. This study also highlights the collection and construction of a dataset to train a YOLOv3 model and what such dataset consists. Finally, this study will question and answer whether or not computer vision is able to reach a high level of accuracy, which algorithm is best suited to classify food and how a reliable dataset is constructed. The significance of this research is to provide a tool that assists people who are not familiar with traditional Maltese food such as tourists to distinguish between various food easily.

II. LITERATURE REVIEW

Computer vision has become a very demanding technology that contributed to various sectors to identify various aspects; for instance, computer vision contributes a lot in the medical field and helps practitioners such as nurses and doctors in the process of identifying, study or treat diseases [1]. This technology is often found in agriculture, the tourism sector, sports, manufacturing amongst others. This technology works by taking an image, understand what the image is and interpret its content with a view to solve tasks and resolve queries.

Study relevant to pedestrian recognition [2] discusses how computer vision includes two operations; to filter the image

by reducing noise that could hinder the object segmentation algorithms performance, and the second operation is subtracting the background from the region of interest (ROI). The main concern about this method is that the captured image conditions may have different lighting, shadows and level of reflections that might influence the overall performance of subtracting the background. This step could be quite challenging and plays a very important role in object detection systems, a solution is to use a sliding window that moves across an image and collect window view to the histogram of oriented gradients (HOG) however, being an expensive operation this method is not very feasible. Therefore a less expensive solution are algorithms based on deep neural networks that are proven to have a better success rate.

Other relevant studies take into variable shape, size and the colour of the object to be detected. During a study regard vehicle intelligence and understanding of its surroundings [3], the researcher starts to question various aspects in relation to a typical driving scenario. Other variables discussed where the dynamic environment a vehicle is found in, thus the researcher takes the lighting of the environment and the background into consideration.

Throughout, various researchers have done intense study in the field of computer vision, some of which have solved major tasks. One vital study in the medical sector has researched the use of computer vision in medical imaging to localize and examine human long bones. Images also come in the form of sound waves which are emitted by a B-scan sensor and are translated into a digital image more understandable to the human eye. However, a B-scan sensor is not able to distinguish between long bones and other specimens such as muscles, soft tissue, veins and other internal organs under the skin tissue. Hence, this study implements a deep learning system to be able to recognize bone specimens. This study discusses various other algorithms but makes use of YOLOv3 convolution neural network algorithm to differentiate the said specimens in ultrasound images [1]. Similarly, my research makes use of the same concepts this study does which classifies long bones, or in my case traditional Maltese food, in digital images.

YOLOv3 Algorithm

YOLOv3 is a regression-based one-stage target detection algorithm just like SSD [4] which can detect and classify objects in an image at the same time [5]. YOLOv3 divides a digital image into $S \times S$ non-overlapped grids and determines whether each grid has a target to predict bounding boxes around [6]. Considering that YOLOv3 is a single end-to-end network that performs feature extraction, location and classification in a network, makes the algorithm extremely fast i.e. at 320×320 YOLOv3 runs at 22ms at 28.2mAP [7]. mAP is the mean average precision value which indicates the accuracy of detection in target. The precision is calculated using the following formula where true positive (TP) is a numerical representation of correct division, and false positive (FP) is a representation of wrong division. [1]

$$\text{Precision} = \frac{TP}{TP + FP}$$

Fig 1. mAP Formula

A convolution neural network (CNN) is a deep learning network that assigns importance to an image such as weights and can differentiate one aspect of an image from another. YOLOv3 is one of the fastest object detection algorithms that makes use of a CNN in real-time detection without loss of accuracy.

Different Algorithms

A study of vehicle detection [3] highlights and compares other object detection algorithms with a view of identifying which is the better algorithm.

Computer vision is categorized into two types; region-based such as R-CNN and Faster RCNN and regression-based such as SSD and YOLOv3. Region-based target detection algorithms extract region proposals from top to bottom of a given image as a candidate region for the model to analyse, extract features in each proposal, classifies them and performs border regression on the region proposals thus this process takes a lot of time since it is divided into several stages. R-CNN is a region-based algorithm that uses a selective search algorithm that makes the process run too slow therefore making the overall object detection run slow.

Faster RCNN was introduced to improve the quality of R-CNN. This algorithm makes use of region proposal network (RPN) to select the candidate regions instead of the selective search algorithm which results in faster detection and enables end to end detection by a neural network. RPN reduces the number and improves the quality of region proposals [7]. This method has a limited ability to extract other features and

generalization hence a neural network-based algorithm using the deep learning model solves several problems.

A one-stage algorithm, SSD, results that the algorithm achieves better detection on smaller objects since the SSD generates more anchor points to make the object position more accurate [3].

Dataset

Study in the medical field using computer vision [5] reseraches cancer tumour in mammogram images using a convolutional neural network. This study explains the use of a digital database for screening mammographs to train and tests the CNN system. The dataset used contains 2,620 cases of breast imaging, including four mammograms for each case made up from non-malignant to critical cases in order to have accurate results. Indeed, deep learning requires a complex dataset to train a model to be more accurate. During this study, researchers have discussed techniques of augmentation to the training dataset where such techniques require for the dataset to generate new instances using different transformation methods such as rotation of images, translation and scale [8]. Considering that very limited images may be available online to scrape which results in a small dataset, this technique could improve the training process of a model by expanding and generating more images. These researchers discuss how the original dataset was augmented by rotating the mammograms three times with angles 90, 180 and 270 degrees.

III. METHODOLOGY

To start the process of classifying different articles of Maltese food I had to first establish which articles I would like my model to classify hence I choose the traditional Maltese cheeselet commonly known as 'ġbejniet' and another Maltese snack deep-fried date roll's commonly known as 'imqaret'. Knowing which articles of food my model can classify I could start the training process which requires a large set of annotated images to train. However, since classifying these types of food is not a very popular concept it was difficult to find a ready annotated dataset that consists of hundreds of images of cheeselet and date roll's, thus I had to create my own. I have adopted to use the YOLOv3 algorithm to classify food due to being the better algorithm to classify such objects as explained previously.

Following is a list of software and hardware components used throughout this project;

Hardware

- Computer
- Smartphone

Software

- Google Colab
- Google Drive
- Google Images
- BBox-Label-Tool
- Fatkun Batch Download
- Batch Resize

To successfully result with a working model that is able to make use of computer vision, a pipeline was followed to visualize the workflow better. I have made use of the following pipeline;

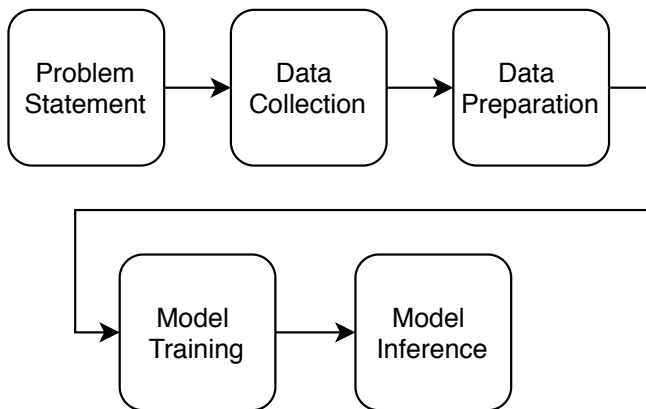


Fig 2. The Pipeline

Problem Statement

I had to establish which type of computer vision technique my model will make use of, thus in order to differentiate objects from each other my model will use a classification technique which is the ability to understand whats inside an image according to its visual content.

Data Collection and Preparation

To acquire several images I made use of an online tool called 'Fatkun Batch Download', this enabled me to search and scrape several images found on Google then filter through each image and decide their relevance under some circumstances i.e. angle of the article of food, lighting, resolution and if other food articles are overlapping each other. Filtering through the scraped images I was able to delete irrelevant images downloaded unintentionally or duplicates of the same image. This process causes the dataset to have a smaller number than desired to train the model which results in a less accurate prediction. To resolve I have also populated the dataset with images I took myself of both food articles to classify.

Annotation process enables the images to be labelled for the algorithm of choice to train the model. Before feeding an annotation tool with images to label, YOLOv3 requires images to all be the same size, this step was easily achieved by using 'Batch Resize' online tool to resize all images at 416 x 416.

Making use of 'BBox-Label-Tool' is a great tool that annotated images ready for YOLOv3 by providing a JPG, PNG or JPEG image as input and highlights a border around the article of food. In this process, I made sure to highlight each individual article of food alone as grouping several foods altogether will influence the training. Image annotation is another way of filtering through each image and deciding their relevance. After highlighting the regions of food, the tool will provide a text file populated with the image label which includes the index of class ('gbejniet', 'imqaret'), x-axis, y-axis, width and height.

Model Training

Training a YOLOv3 model requires a great amount of processing power hence I trained my model using Google Colab as this environment provides Googles own GPU's for use. I opted for this method as my personal computer doesn't have a graphics card thus training a model would have been almost impossible. Before training the model, colab required the setup of installing darknet and uploading the dataset to darknet's environment where two text files were generated for training and testing. Configuring darknet enabled me to provide the environment with all relevant files necessary including classes and the destination to dataset, destination of the backup folder to save weights and paths to train and test files. Finally, it was time to train in darknet using *darknet53.conv.74* according to the *yolov3_custom.cfg* file which consists of configurations and arithmetic solutions related to my model.

	Type	Filters	Size	Output
1x	Convolutional	32	3 × 3	256 × 256
	Convolutional	64	3 × 3 / 2	128 × 128
	Convolutional	32	1 × 1	
	Convolutional	64	3 × 3	
	Residual			128 × 128
2x	Convolutional	128	3 × 3 / 2	64 × 64
	Convolutional	64	1 × 1	
	Convolutional	128	3 × 3	
	Residual			64 × 64
	Convolutional	256	3 × 3 / 2	32 × 32
8x	Convolutional	128	1 × 1	
	Convolutional	256	3 × 3	
	Residual			32 × 32
	Convolutional	512	3 × 3 / 2	16 × 16
	Convolutional	256	1 × 1	
8x	Convolutional	512	3 × 3	
	Residual			16 × 16
	Convolutional	1024	3 × 3 / 2	8 × 8
	Convolutional	512	1 × 1	
	Convolutional	1024	3 × 3	
4x	Residual			8 × 8
	Avgpool		Global	
	Connected		1000	
	Softmax			

Fig 3. Darknet-53 Framework

Testing and Results

The testing dataset contains a set of both cheeselet and deep-fried date roll's images from different angles and some images made up of multiple objects to detect. Classifying images using *yolov3_custom_final.weights* which are the trained weights using the food dataset. After testing with test sample images, the food classifier model was able to reach over 71% mean average precision (mAP) rate when an image consists of multiple objects as seen below.

```
calculation mAP (mean average precision)...
80
detections_count = 542, unique_truth_count = 341
class_id = 0, name = gbejniet, ap = 81.90% (TP = 153, FP = 46)
class_id = 1, name = imqaret, ap = 60.18% (TP = 97, FP = 42)

for conf_thresh = 0.25, precision = 0.74, recall = 0.73, F1-score = 0.74
for conf_thresh = 0.25, TP = 250, FP = 88, FN = 91, average IoU = 56.19 %

IoU threshold = 50 %, used Area-Under-Curve for each unique Recall
mean average precision (mAP@0.50) = 0.710420, or 71.04 %
Total Detection Time: 2 Seconds
```

Fig 4. mean Average Precision

When the object to detect is relatively large the model detects good results, in fact the model was able to reach a 100% accuracy rate when classifying images with a single object as seen below. We can understand the high level of accuracy rate YOLOv3 has.

```
test_images/(6).jpg: Predicted in 90.348000 milli-seconds.
gbejniet: 89%
gbejniet: 98%
gbejniet: 100%
gbejniet: 92%
Unable to init server: Could not connect: Connection refused

(predictions:4826): Gtk-WARNING **: 16:38:59.584: cannot open display:
```



Fig 5. Multiple Object Classification

```
Done! Loaded 107 layers from weights-file
test_images/(2).jpg: Predicted in 90.143000 milli-seconds.
gbejniet: 100%
```

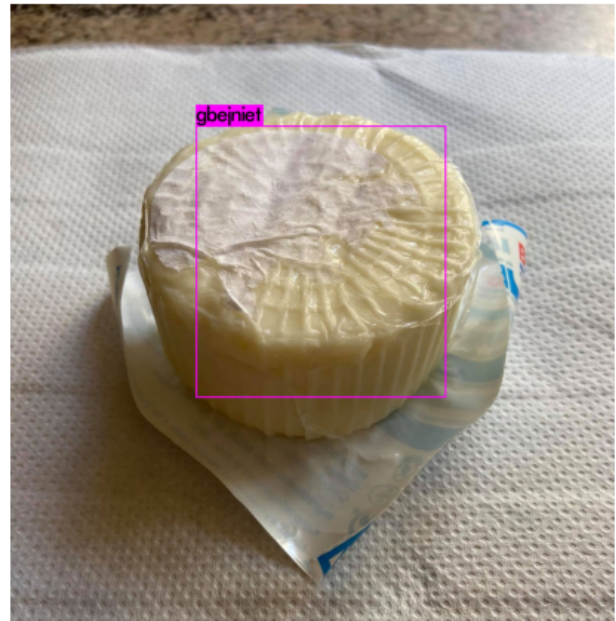


Fig 6. Single Object Classification

When testing with low-resolution images that consist of a too complicated detection scene the model doesn't detect anything.

```
Done! Loaded 107 layers from weights-file
test_images/(85).jpg: Predicted in 89.985000 milli-seconds.
Unable to init server: Could not connect: Connection refused

(predictions:4797): Gtk-WARNING **: 16:38:08.743: cannot open display:
```



Fig 7. Distorted Image

IV. EVALUATION

The model was trained using 350 images and has reached a level of 71% accuracy. This result reflects a reliable and good dataset constructed with various images capturing object from different angles, in different environments and making use of augmentation techniques. However, the dataset could improve in view of having a more accurate model. Improving a dataset would consist of having a larger number of images made up from different augmentation techniques such as image rotation and scale thus allowing a more diverse dataset. Improving the model would consist of having other relevant food to classify.

This result have clearly pointed that the bounding boxes around object being detected are classifying a smaller area of the object rather than the whole object in target, this is due to the method used to label object during the annotation process thus to resolve the dataset should be annotated again with more coverage around objects in target.

After gathering results and evaluating data in previous scenarios we are able to answer this paper's research questions;

1. How can a reliable dataset be obtained or constructed?

As highlighted in the methodology, this study describes how a dataset is constructed from scratch by making using of a bulk download tool, an annotation tool to label object for detection and a resize tool to resize all images in preparation for an algorithm. Whilst keeping in view the need of augmentation techniques such as rotation and angles of images.

2. Which algorithm is best used to classify food?

From various other studies and the result of this one as well, we can justify that YOLOv3 is the ideal algorithm of choice as it is regression-based rather than region-based. Results justifies the performance YOLOv3 has on images of food for classification.

3. Will image processing be able to reach a high level of accuracy?

Results using diverse test images result in high accuracy rates thus justifying that image processing, in conjunction with an annotated dataset and a good classification algorithm, is able to reach a high level of accuracy.

V. CONCLUSION

This research paper proposes to build a traditional Maltese food classifier based on the YOLOv3 algorithm combined with a neural network to be applied in scenarios to help the tourism sector in Malta. Results of the model reflect a very high accuracy level and can be improved by training of larger and more complex datasets including different food articles.

Hence, for future work, the algorithm could expand on a larger range of food articles and improve accuracy by training on a better dataset to be able to discriminate better between different objects in the same image.

REFERENCES

- [1] R. A. F. Lazuardi, T. Karlita, E. M. Yuniarno, I. K. E. Purnama, and M. H. Purnomo, "Human bone localization in ultrasound image using yolov3 cnn architecture," in *2019 International Conference on Computer Engineering, Network, and Intelligent Multimedia (CENIM)*. IEEE, 2019, pp. 1–6.
- [2] P. Tumas and A. Serackis, "Automated image annotation based on yolov3," in *2018 IEEE 6th Workshop on Advances in Information, Electronic and Electrical Engineering (AIEEE)*. IEEE, 2018, pp. 1–3.
- [3] H. Wang, Y. Yu, Y. Cai, X. Chen, L. Chen, and Q. Liu, "A comparative study of state-of-the-art deep learning algorithms for vehicle detection," *IEEE Intelligent Transportation Systems Magazine*, vol. 11, no. 2, pp. 82–95, 2019.
- [4] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [5] K. Djebbar, M. Mimi, K. Berradja, and A. Taleb-Ahmed, "Deep convolutional neural networks for detection and classification of tumors in mammograms," in *2019 6th International Conference on Image and Signal Processing and their Applications (ISPA)*. IEEE, 2019, pp. 1–7.
- [6] Y. Zheng, H. Bao, X. Xu, N. Ma, J. Zhao, and D. Luo, "A method of detect traffic police in complex scenes," in *2018 14th International Conference on Computational Intelligence and Security (CIS)*. IEEE, 2018, pp. 83–87.
- [7] F. Miao, Y. Tian, and L. Jin, "Vehicle direction detection based on yolov3," in *2019 11th International Conference on Intelligent Human-Machine Systems and Cybernetics (IHMSC)*, vol. 2. IEEE, 2019, pp. 268–271.
- [8] M. Heath, K. Bowyer, D. Kopans, P. Kegelmeyer, R. Moore, K. Chang, and S. Munishkumaran, "Current status of the digital database for screening mammography," in *Digital mammography*. Springer, 1998, pp. 457–460.