

LONDON'S GLOBAL UNIVERSITY



Anomaly detection in parcel X-rays using Siamese and triplet neural networks

Martynas Janonis¹

BSc Computer Science

Supervisor: Dr. Lewis Griffin

Submission date: 19th May 2020

¹**Disclaimer:** This report is submitted as part requirement for the BSc Computer Science at UCL. It is substantially the result of my own work except where explicitly indicated in the text. The report may be freely copied and distributed provided the source is explicitly acknowledged

Abstract

Detecting anomalies in X-ray security imaging is quite a challenging task. Previous solutions used features extracted from a CNN trained only on photographic images. This work aims to improve the performance of previous approaches by using a Siamese/triplet CNN that is pre-trained on ImageNet and fine tuned on an X-ray image dataset. In particular, this paper evaluates whether a representation extracted from a Siamese/triplet network, trained to distinguish if two X-ray images were different views of the same parcel or not, is effective for anomaly detection. Comparison of various network architectures and training methods is provided in detail. The best results, achieved by an isolation forest trained on embeddings generated by a triplet Densenet-201 model (pre-trained on ImageNet and fine tuned on the X-ray images dataset), were: $F_1 = 0.4048$, AUROC = 0.6922, $\kappa = 0.0908$. The results indicate that these embeddings show minor potential for use in anomaly detection.

Contents

1	Introduction	2
1.1	Problem outline	2
1.2	Aims and goals	2
1.3	Project development overview	3
1.4	Chapter overview	3
2	Context	4
2.1	Background research	4
2.1.1	Anomaly detection	4
2.1.2	Self-supervised representation learning	5
2.2	Siamese networks	5
2.3	Triplet networks	7
2.4	CAST dataset	9
2.5	Software & other tools used for the project	10
3	Design and Implementation	11
3.1	High-level design overview	11
3.2	Image pre-processing pipeline	11
3.3	Embedding networks	12
3.3.1	Siamese networks	12
3.3.2	Triplet networks	13
3.4	Anomaly detection using the generated embeddings	14
3.4.1	Using a linear SVM	14
3.4.2	Using isolation forests	14
4	Results evaluation	15
4.1	Overall performance	15
4.1.1	Linear SVM version	15
4.1.2	Isolation forest version	16
4.2	Computational performance	16
4.3	Generated embedding suitability for anomaly detection	16
5	Conclusions and evaluation	17
5.1	Summary	17
5.2	Future work	17
A	Quick guide	22
B	Source code	23
B.1	datasets.py	23

B.2	generate_vectors.py	29
B.3	losses.py	30
B.4	metrics.py	31
B.5	networks.py	33
B.6	train_isolation.py	34
B.7	train_siamese.py	35
B.8	train_svm.py	35
B.9	train_triplet.py	37
B.10	trainer.py	38
B.11	trim.py	40
B.12	utils.py	42
C	Project plan	46
D	Interim report	49

Chapter 1

Introduction

1.1 Problem outline

X-ray imaging is prevalent in modern society. It is used to detect various forms of cancers in humans [1] and inspect parcels, luggage and vehicles for illicit items [2].

While multiple advances in X-ray screening equipment have been made, the resulting images are still manually inspected by humans. As a result, the performance of the screening system ultimately depends on human factors and training [3]. Additionally, due to the massive volume of parcels, the repetitive nature of the job, as well as the visual challenge of discerning individual items in translucent and overlapping X-ray images, human operators are prone to making costly mistakes [4, 5, 6].

Dense and large objects, such as laptops and game consoles, are especially difficult to analyse. Their construction makes it hard for the X-rays to penetrate and image properly, thus making them an ideal place to store illegal goods. Threats placed inside or behind the laptop are especially difficult for human operators to spot [7]. Furthermore, in an airport setting, X-ray operators have about 3 seconds to make a decision on whether the bag contains illegal items [8]. This time pressure further diminishes the effectiveness of current day screening systems. An automated, ML based appearance anomaly detection solution, could potentially solve these issues by flagging abnormal looking packages, indicating that their contents might be malicious.

1.2 Aims and goals

Some previous works indicate that transfer learnt representations extracted from a CNN can be used for anomaly detection [9]. The question is whether these representations can be improved by training the CNN on another task, such as determining whether two X-ray images show the same parcel. Potentially, the representations generated after training a network on that task will have certain features which are more useful for detecting anomalies.

The aim of this project is to find out whether representations extracted from a Siamese or a triplet network, trained on discerning whether a pair of X-ray images are different views of the same parcel or not, can be used for effective appearance anomaly detection.

This project has the following goals:

- Train a Siamese and a triplet network to decide if two parcel X-rays are different views of the same parcel or not.
- Assess whether the resulting representations extracted from those networks are effective or not for anomaly detection.

1.3 Project development overview

This project was carried out in multiple stages. Firstly, a literature review was performed to find out the current approaches to appearance anomaly detection. Secondly, multiple network architectures were compared and the best performing one was picked for further training. Thirdly, a Siamese and a triplet network was trained and evaluated. Finally, the representations generated by the best performing neural network were passed through an isolation forest and a linear SVM to determine their usefulness in anomaly detection.

1.4 Chapter overview

The rest of the paper is structured as follows:

Chapter 2 gives a summary of the background research, network architectures and datasets used for training.

Chapter 3 provides a detailed explanation of the design and implementation of the anomaly detection system.

Chapter 4 evaluates the results achieved by the final models.

Chapter 5 concludes the paper by giving a summary of what was achieved and provides an evaluation of the entire project. It also briefly discusses potential future work that could be done if the project was developed for a longer period of time.

Chapter 2

Context

2.1 Background research

The background research is based on the works of Mouton *et al.* [10], Rogers *et al.* [11] and Akcay *et al.* [12].

2.1.1 Anomaly detection

Currently, there aren't many works published on the topic of unsupervised anomaly detection within X-ray imaging. Most authors instead focus on supervised learning for the purposes of object classification, detection and segmentation.

One of the earliest anomaly detection approaches [13] trains sparse feed-forward auto-encoder networks on normal data and uses the hidden layer activations of the network as representations. An SVM is then used to classify the images as normal or anomalous. Results on the MNIST dataset and a cargo container dataset showed some success as the normal class of images had a low variability, which was easily captured by the network. The work shows that representations extracted from the network are significant for detecting anomalies in images. Unfortunately, the data contained in X-ray images of parcels has a high variability, which means that the auto-encoder networks might not be suitable for the task of anomaly detection in X-ray imaging [14, 15].

The work by Andrews *et al.* [9] shows a promising representation learning framework for anomaly detection. The extracted representations from a CNN's pooling layer, that was trained to distinguish empty cargo containers from non-empty ones proved to be useful for anomaly detection. This method has the advantage of making no assumptions about the generating distributions of the normal or anomaly classes. The authors show that the transfer learnt representations, with some tweaking, show promise for use in anomaly detection.

Griffin *et al.* [14] describes another approach to anomaly detection in X-ray imaging that uses features extracted from a network similar to Inception V3, which is trained on the ImageNet dataset, to train a multivariate Gaussian model. The Gaussian model is then used to capture the distribution of the dataset. The resulting anomaly status of an image is based on the computed likelihoods relative to the Gaussian model. This method achieved an AUC (Area Under the Curve) of 92.5%. This paper also suggested multiple ways to increase the performance of the system. One method is to use a more advanced CNN architecture to generate embeddings and another is to use a CNN specifically trained on X-ray images. The authors also mention that while there is major potential to achieve better performance by using a CNN trained on X-ray images, sourcing the data might be a significant challenge.

This project seeks to implement and test both performance increasing suggestions of Griffin *et al.* [14].

2.1.2 Self-supervised representation learning

Self-supervised learning enables the ability to train a neural network using the labels that can be inferred from the data. Labelled datasets are hard and expensive to produce, but, in contrast, unlabelled data is abundant. One way to make use of that unlabelled data is to set the training task so that the supervision can be achieved from the data itself [16].

The network is trained to perform a self-supervised task (a pretext task). The task can be quite arbitrary (such as predicting image rotations [17]), but the idea is that by learning to do that task well, the network learns certain representations of the input that can be later used to complete a more complex task. In the context of this project, the pretext task is to classify whether two X-ray images show the same parcel from different views or not, and the more complex, downstream task, is to use those generated embeddings to classify images as anomalous or benign.

The work by Doersch *et al.* [18] explores a similar concept. The authors extract random pairs of patches from images of a large, unlabelled dataset, and train a CNN to predict the position of the second patch relative to the first. This, in theory, makes the network learn to recognize shapes, objects and spatial context of the images. This work shows that the resulting CNN trained on the aforementioned task outperformed a randomly initialized CNN, resulting in a state-of-the-art performance on the Pascal VOC 2011 dataset. This provides the motivation, that by training a CNN to determine if two image patches belong to the same parcel, it could be used to generate useful intermediate representations of X-ray images.

2.2 Siamese networks

A Siamese neural network consists of two identical networks, that share the same weights, joined at their output [19]. This architecture allows a model to tell how similar its inputs are, instead of making it classify one input. Siamese networks are commonly used in one-shot learning [20] (where one has very few training examples of each class) and facial recognition [21].

Siamese networks are used with a distance-based loss, in this case — contrastive loss. Each sister network takes one input each. The outputs of the last layers are used to compute the distance-based loss, which shows how similar the inputs are (how close are their embeddings in n-dimensional space). Training the network makes it generate embeddings which have low distances between positive pairs and distances greater than some margin between negative pairs.

Positive pairs consist of an anchor sample and a positive sample, which are similar in certain metrics. In contrast, negative pairs consist of an anchor sample and a negative sample, which are dissimilar in those metrics. In this project, a positive pair contains two images of the same parcel, but from different perspectives, and a negative pair contains two images of unrelated parcels.

In the context of this project, let Img_1 and Img_2 be a pair of images. Let Y be a binary label, $Y = 0$ if the images show different parcels and $Y = 1$ if the images show different views of the same parcel. Let \mathbf{W} be the shared weight vector that is learned. Then $\text{CNN}_{\mathbf{W}}(\text{Img})$ is an embedding of the image Img that is extracted from the network's last pooling layer after running the image through it.

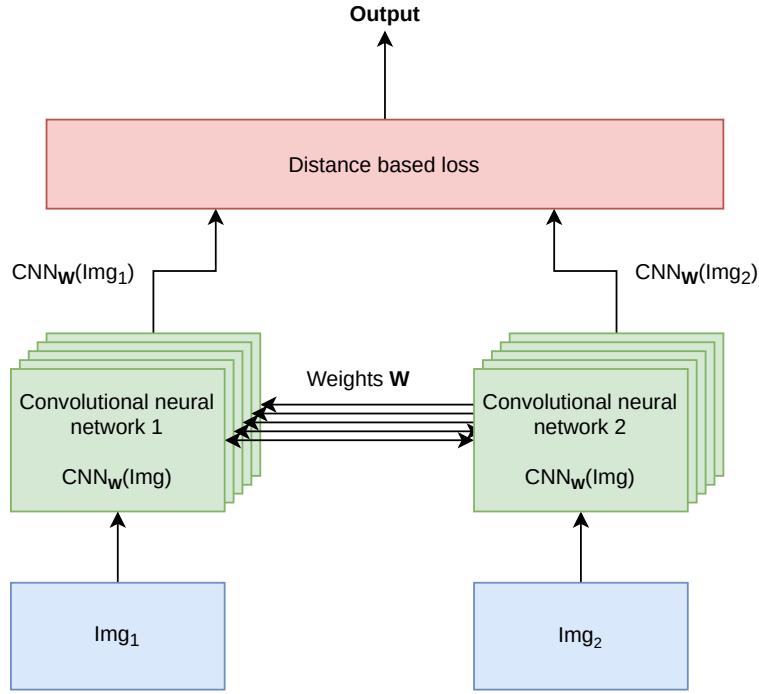


Figure 2.1: Siamese network architecture diagram [22]

Then the contrastive loss is defined as [23]:

$$\mathcal{L}_{\text{contrastive}}(\mathbf{W}, Y, \text{Img}_1, \text{Img}_2) = (1 - Y) \frac{1}{2} (D_{\mathbf{W}})^2 + (Y) \frac{1}{2} (\text{ReLU}(m - D_{\mathbf{W}}))^2$$

where $m > 0$ is a margin that defines a radius around $\text{CNN}_{\mathbf{W}}(\text{Img})$ and $D_{\mathbf{W}}$ is the euclidean distance between the generated embeddings [23]:

$$D_{\mathbf{W}} = \|\text{CNN}_{\mathbf{W}}(\text{Img}_1) - \text{CNN}_{\mathbf{W}}(\text{Img}_2)\|_2$$

Positive pairs will have $\mathcal{L}_{\text{contrastive}} = 0$ only if the distance between the two generated representations is 0. The bigger the distance between the positive pair, the bigger the resulting loss.

Negative pairs will have $\mathcal{L}_{\text{contrastive}} = 0$ only if the distance between the two representations is greater than margin m . Otherwise, the smaller the distance, the bigger the loss (up to m).

The margin is used to avoid wasting efforts on maximizing the distance between already distant negative pairs. This allows the network to focus on more difficult training examples [24].

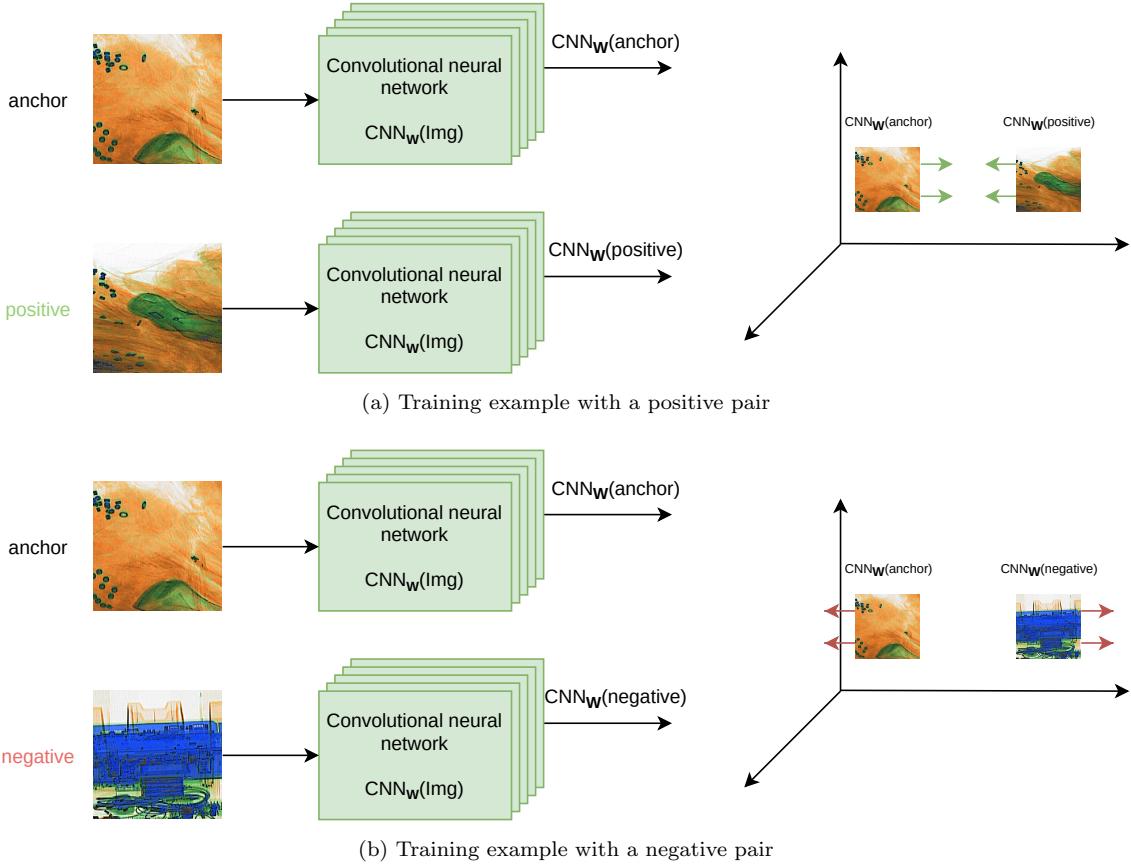


Figure 2.2: Training setup using a distance-based loss¹

Figure 2.2 shows an example of two training steps, one containing a positive pair and another containing a negative pair. In the positive case, the network “pulls” the representations together while in the negative case, the network “pushes” them apart.

The architecture shown in figure 2.1 takes quite a lot of memory in the GPU, as two big convolutional neural networks have to be loaded at the same time. To reduce the memory footprint, instead of having two identical networks loaded at the same time and passing Img_1 to one sister network and Img_2 to another, only one CNN is loaded and the images are passed to it sequentially. As the weights of the two CNNs are the same, the resulting network will perform the same.

2.3 Triplet networks

A triplet neural network consists of three identical networks that share the same weights. When the network gets fed three inputs, it outputs 2 values — the L_2 distance between the anchor and the positive sample as well as the L_2 distance between the anchor and the negative sample [25]. Those values are then fed into the triplet loss function.

¹ Adapted from [24]

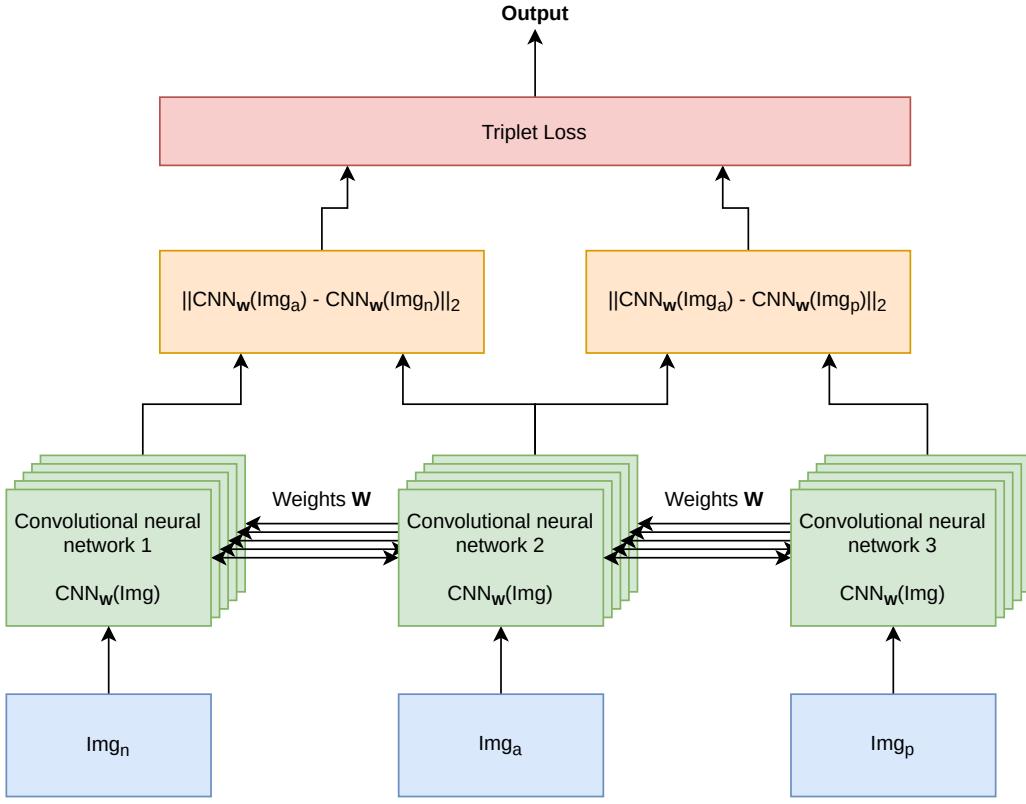


Figure 2.3: Triplet network architecture diagram

Let Img_a , Img_p and Img_n be a triplet of images. Let \mathbf{W} be the shared weight vector that is learned. Let $d(x, y)$ be the distance between embeddings x and y in some embedding space. Then $\text{CNN}_{\mathbf{W}}(\text{Img})$ is an embedding of the image Img that is extracted from the network's last pooling layer after running the image through it. In this case, the anchor and the positive samples are X-ray images of the same parcel, but taken from different views while the negative sample is an X-ray image of a completely different parcel.

Then the triplet loss is defined as:

$$\mathcal{L}_{\text{triplet}} (\mathbf{W}, \text{Img}_a, \text{Img}_p, \text{Img}_n) = \text{ReLU}(m + \|\text{CNN}_{\mathbf{W}}(\text{Img}_a) - \text{CNN}_{\mathbf{W}}(\text{Img}_p)\|_2^2 - \|\text{CNN}_{\mathbf{W}}(\text{Img}_a) - \text{CNN}_{\mathbf{W}}(\text{Img}_n)\|_2^2)$$

There are three categories of triplets [26]:

- Easy triplets: $d(a, p) + m < d(a, n)$ have a loss of 0 as the negative sample is further away from the anchor than the positive one in the embedding space.
- Hard triplets: $d(a, n) < d(a, p)$ have a loss greater than m as the negative sample is closer to the anchor than the positive one.
- Semi-hard triplets: $d(a, p) < d(a, n) < d(a, p) + m$ have a loss $0 < \mathcal{L} < m$ as the negative sample is further away from the anchor than the positive, but the negative sample is still within the margin m .

The learning objective of a triplet network is to generate embeddings such that the anchor and the positive sample are closer together than the anchor and the negative sample in n -dimensional embedding space by some margin $m > 0$.

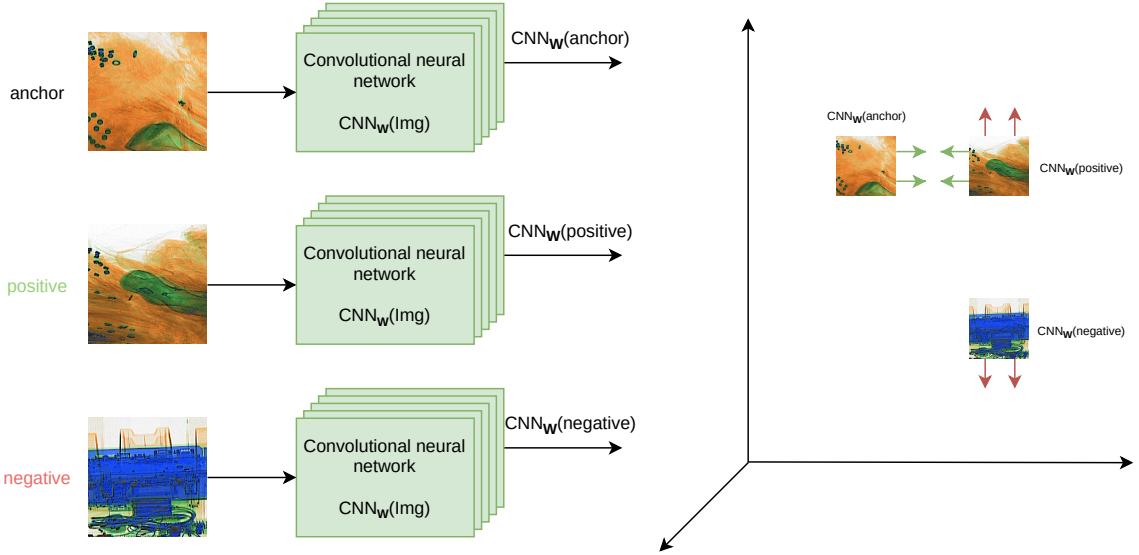


Figure 2.4: Training step of a triplet network. The network “pulls” the anchor and the positive representations together while, at the same time, “pushes” the anchor and the negative representations away from each other.²

2.4 CAST dataset

This project uses the OSCT Borders X-ray Image dataset, made by the UK Home Office Centre for Applied Science and Technology (CAST). The dataset contains X-ray images of parcels, divided into two sets.

The first, stream-of-commerce set, contains X-ray images of 5000 parcels. The images were taken in a UK postal distribution centre. Most parcel contents are unrecognizable, but various clothing, machine and electronic items can be seen in some.

The second, staged threat set, contains images of 160 parcels. These images show firearms or their parts hidden inside parcels containing normal items.

All parcel images are dual-view, meaning they have been taken from two different perspectives (i.e. one shows a top-down view and another — side-view). They are also pseudo-coloured using information from dual-energy imaging. An example is shown in figure 2.5

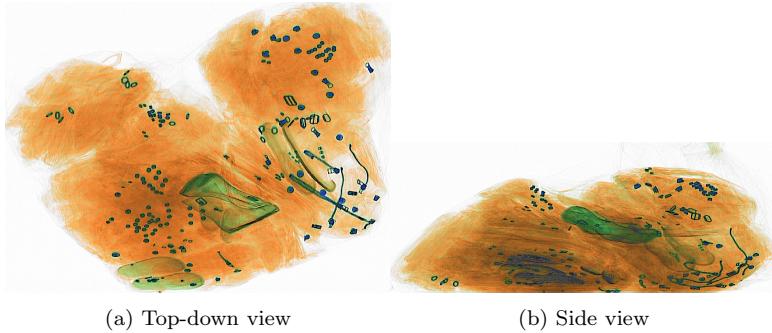


Figure 2.5: An example of a dual-view parcel image from the dataset, containing clothing and footwear

²Adapted from [24]

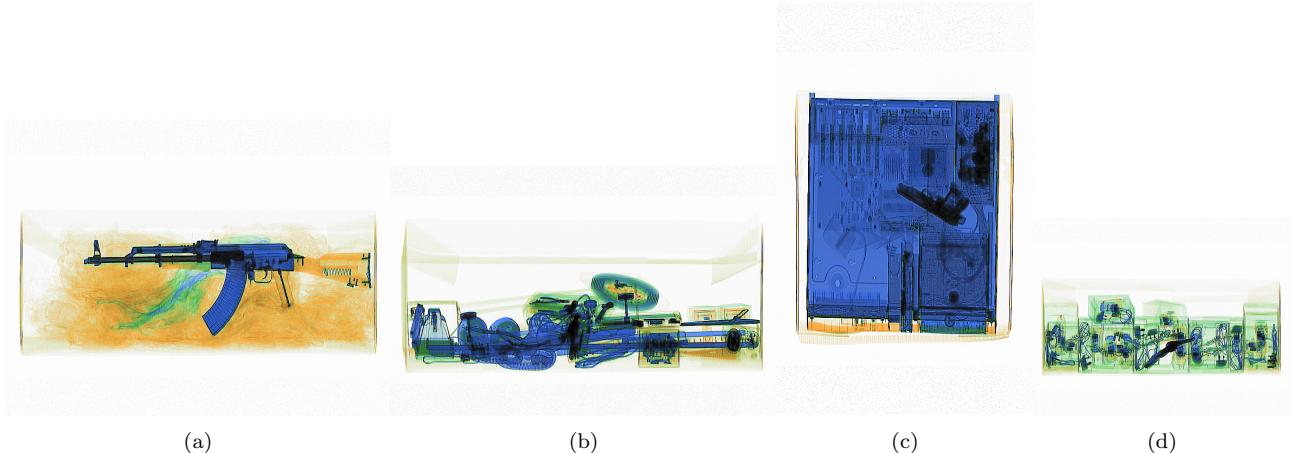


Figure 2.6: Examples from the staged threat set

2.5 Software & other tools used for the project

This project makes heavy use of the PyTorch deep-learning library [27]. All training scripts, network architecture and loss implementations, etc. were written using PyTorch. The TensorFlow library was a viable alternative. Both libraries are more or less functionally equivalent, so the choice ultimately boils down to personal preference.

Scikit-learn [28] was used to train SVMs and isolation forests for anomaly detection. It was chosen due to previous experience with it, its popularity in the field and great documentation.

OpenCV was used in various stages of the data pre-processing pipeline. It was used to trim, resize and pad the images. It was selected as it is easy to use, is well documented and performs well.

The starting codebase was taken from Adam Bielski's `siamese-triplet` repository on GitHub [29], which provides sample implementations of Siamese and triplet networks using PyTorch. It is licensed under the BSD 3-Clause license. The code was modified and extended during the development of this project. All modifications are licensed under the EUPL v1.2-or-later license.

Chapter 3

Design and Implementation

3.1 High-level design overview

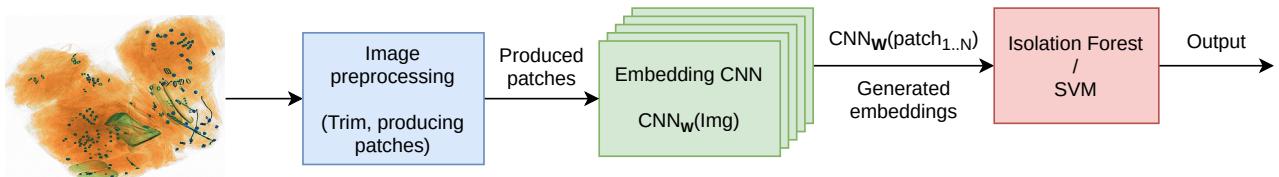


Figure 3.1: A high level architecture of the anomaly detection system

The complete anomaly detection system works in multiple stages. First, an image is loaded, trimmed to remove whitespace and reduced to a set of 224×224 patches. Those image patches are then fed into a Densenet-201 CNN that was pre-trained on ImageNet and then trained to determine whether two images show the same parcel or not. The network takes the patches and produces a set of embeddings, one embedding for each patch. The resulting set of embeddings is passed to an isolation forest that was trained using embeddings generated from the X-ray images dataset. The isolation forest then outputs a set of labels (one for each embedding), 1 if the embedding was normal and -1 if the embedding was anomalous. The resulting classification of the whole image is then taken as a minimum value of the whole set of labels. This means that if at least one embedding was classified as anomalous, the whole image is treated as such.

In total, 4 different configurations of this system, each with a different combination of components, were implemented and evaluated:

Embedding network	Anomaly detection model
Densenet-201 (pre-trained on ImageNet and trained on X-ray dataset)	Linear SVM
Densenet-201 (pre-trained on ImageNet and trained on X-ray dataset)	Isolation forest
Densenet-201 (trained only on ImageNet)	Linear SVM
Densenet-201 (trained only on ImageNet)	Isolation forest

Table 3.1: Different configurations of the anomaly detection system

The following sections describe each component of this system in detail.

3.2 Image pre-processing pipeline

First, all images were trimmed to remove the air around the parcel (which shows up as white in the image). To achieve that, a grayscale image together with thresholding is used to get a bounding box around the contents

of the image. That bounding box is then used to crop the RGB image. If two images are of the same parcel, then they are trimmed in a way that keeps their widths equal.

When training a Siamese or a triplet network, a 224×224 random crop is taken from the image. That specific resolution is used as the networks were pre-trained on images of the same resolution. If the image isn't large enough, it is padded with white pixels before taking the crop. If two images belong to the same parcel, the crops are taken from the same x coordinate. This makes it so the two patches correspond to roughly the same part of the parcel. When the networks are evaluated, a 224×224 center crop is taken instead to keep the validation between epochs consistent.

Each 3-channel RGB image ($3 \times H \times W$) is normalized using mean = [0.485, 0.456, 0.406] and std = [0.229, 0.224, 0.225] before feeding it to the network. That is required as the CNN models were pre-trained on the ImageNet dataset and they expect their inputs to be normalized using those values.

For anomaly detection, each image is reduced to a set of 224×224 patches using a stride of 112. Those patches are then preprocessed as explained above.

3.3 Embedding networks

The embedding networks are used to generate 1024-D vector representations of the images, which are later fed into an isolation forest model or a linear SVM. The networks are first trained in a Siamese or triplet architecture and then “decoupled” for use in the anomaly detection pipeline (by using the `get_embedding` function as seen in listing 1). The idea is that the CNN generated embeddings potentially have certain features that makes them useful for anomaly detection.

All networks were pre-trained on the ImageNet dataset. Then their final classification layer (a 1000-D dense layer, each dimension corresponding to a unique image class in the ImageNet dataset) was stripped and replaced by a 3-D average pooling layer and a flattening layer after it. That made them output 1024-D vector representations of their inputs.

```

class SiameseNet(nn.Module):
    def __init__(self, embedding_net):
        super(SiameseNet, self).__init__()
        self.embedding_net = embedding_net

    def forward(self, x1, x2):
        output1 = self.embedding_net(x1)
        output2 = self.embedding_net(x2)
        return output1, output2

    def get_embedding(self, x):
        return self.embedding_net(x)

```

```

class TripletNet(nn.Module):
    def __init__(self, embedding_net):
        super(TripletNet, self).__init__()
        self.embedding_net = embedding_net

    def forward(self, x1, x2, x3):
        output1 = self.embedding_net(x1)
        output2 = self.embedding_net(x2)
        output3 = self.embedding_net(x3)
        return output1, output2, output3

    def get_embedding(self, x):
        return self.embedding_net(x)

```

Listing 1: Implementation of Siamese and triplet architectures

3.3.1 Siamese networks

The task for the Siamese networks was to determine whether the two images it received as input were of the same parcel or not. If the two pictures show the same parcel, but from a different perspective, the network outputs two embeddings which are very close in 1024-D space, otherwise it outputs embeddings which are far apart.

Multiple Siamese networks were trained and evaluated, each using a different embedding CNN. The embedding networks used were: ResNeXt-101- $32 \times 8d$ [30], Densenet-201 [31], ResNet-50 [32] and Wide ResNet-50-2 [33].

These exact CNNs were selected because they represent the state-of-the-art in image classification performance on the ImageNet dataset. As the task of determining whether two images belong together is relatively similar to image classification, they should also perform quite well here.

All networks were trained for 50 epochs. The contrastive loss function had a margin set to 2. The batch size was set to 16. The Adam optimizer was used (with parameters $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$). The learning rate was initially set to 0.001, but with each epoch it was updated to be 99% of the previous value. All hyperparameters were tuned by experimentation. During training, the networks were fed randomly generated pairs of images, which were preprocessed as explained in the previous section. During validation, the networks were fed pre-determined pairs to maintain consistency between epochs.

After 50 epochs, the networks in order of validation loss (lowest to highest) were: Densenet-201 (0.1081), ResNeXt-101-32 \times 8d (0.11), ResNet-50 (0.1221) and Wide ResNet-50-2 (0.1226). The best model correctly identified whether the images belong to the same parcel 95% of the time.

A Densenet-201 model, which was not pre-trained on ImageNet, was also trained on this task as an experiment to check whether the pre-training makes any difference. Out of all tested models, this one performed the worst. At the end of 50 epochs, its validation loss was 0.1413, which is significantly higher than the loss of the same model pre-trained on ImageNet. This shows that feature discrimination, edge and shape detection learnt during training on ImageNet is useful for this task as well.

Another experiment was carried out by first freezing the pre-trained model's weights and then adding a couple of dense and convolutional layers at the end. This proved to be unsuccessful. The model failed to learn to correctly identify pairs of images and performed no better than random guessing. That is most likely because it didn't have enough trainable parameters and the features generated by the model trained only on ImageNet are not good enough to perform the aforementioned task.

3.3.2 Triplet networks

The task for the triplet networks was to generate embeddings which had a lower distance between the anchor and the positive sample and a longer distance between the anchor and the negative sample. The better it is at this task, the lower the triplet loss. The validation triplet loss was used to evaluate the performance of these networks.

Because the Densenet-201 and ResNeXt-101-32 \times 8d networks performed the best in the Siamese network, they were selected for further testing in a triplet network scenario.

The triplet networks were trained very similarly to the Siamese ones, except the batch size was set to 8 and the triplet margin loss was used with the margin set to 2. During training, the network was fed randomly generated triplets, while during testing it was fed pre-determined ones to keep consistency between epochs.

After 50 epochs, the Densenet-201 model outperformed the ResNeXt-101 one in terms of validation loss (0.0834 and 0.101 respectively). As a result of the Densenet-201 model being the best in the Siamese and triplet case, it was selected to be the embedding network for the final anomaly detection solution.

Various triplet selection methods were tried — random, online semi-hard negative and online hardest negative:

- Random triplet selection: for each anchor and positive sample a random negative is selected. Because no additional embeddings need to be calculated, this method is the fastest.
- Online semi-hard negative selection: for each anchor and positive sample, a negative is selected so that the resulting triplet loss is larger than 0 and smaller than the margin. That negative is selected from a mini-batch of 16 random negatives.
- Online hardest negative selection: for each anchor and a positive sample, a negative is selected so that the resulting triplet loss is larger than 0. In the end, the negative producing the largest loss is selected

out of a mini-bach of 16 random negatives.

In theory, the online hardest negative selection method should result in the best performing network [34]. Curiously, all methods produced very similar results, indicating that this task is not very sensitive to triplet selection or that the network wasn't trained long enough for the selection to matter (i.e. most triplets still produced a loss larger than 0).

3.4 Anomaly detection using the generated embeddings

In the end, the Densenet-201 model trained in a triplet setup was selected to be the embedding network. All tests indicated that this model performed better than the one trained in a Siamese setup.

The 1024-D representations this model generates are then used to train a linear SVM in a supervised manner and an isolation forest in an unsupervised manner. After the training is complete, a novel image can be fed through the network and the generated embedding can be analysed by an SVM/isolation forest.

This is the final stage of the whole pipeline. A set of patches produced by the data preprocessor was previously fed into the embedding CNN, which produced a set of embeddings. That set of embeddings is now passed to the SVM/isolation forest, which produces a set of labels, one for each embedding. The label 1 indicates a normal embedding and -1 indicates an anomalous one. The image is deemed anomalous if at least one label is -1 and normal if all labels are 1.

3.4.1 Using a linear SVM

To determine whether the generated embeddings contain any potential for anomaly detection, a linear SVM was trained to determine whether an embedding is anomalous or not (i.e. whether the image contained a gun or not). A Densenet-201 model trained only on ImageNet was selected as a baseline embedding network for comparison.

This is not the optimal method to do anomaly detection as the vectors generated by the CNN might not be linearly separable, which would render the linear SVM useless. Additionally, supervised learning was used which would not be possible in an anomaly detection setting in the real world, as it is very difficult to get enough labelled data, especially of the anomalous class. This method was only used to verify if the vectors are linearly separable, which would make the final solution much simpler. It was never intended to be the final solution on its own.

The SVMs were trained using the ASGD (Averaged Stochastic Gradient Descent) optimizer [35] and class weights of 1 for the normal class and 50 for the anomalous class. Number of training epochs was set to 50 and the batch size was 256.

3.4.2 Using isolation forests

For the final stage in the anomaly detection pipeline, an isolation forest was selected. It has many benefits over other solutions: low memory usage, great performance, ability to deal with high-dimensional data and, most importantly, it can be trained without many examples of anomalies [36].

To train the isolation forest, a train and validation set was created. The training set had 80041 generated embeddings and the test set had 14125. The embeddings were generated by the triplet Densenet-201 model. The isolation forest was initially configured to have 1024 base estimators in the ensemble. It was then trained using the training set.

Another isolation forest was also trained, but using the embeddings from a Densenet-201 model only trained on ImageNet for comparison purposes.

Chapter 4

Results evaluation

4.1 Overall performance

Embedding network	Anomaly detection model	F1 Macro score	AUROC	Cohen's kappa
Densenet-201 (ImageNet + X-ray)	Linear SVM	0.1849	0.8489	0.1454
Densenet-201 (ImageNet only)	Linear SVM	0.5143	0.9697	0.5009
Densenet-201 (ImageNet + X-ray)	Isolation forest	0.4048	0.6922	0.0908
Densenet-201 (ImageNet only)	Isolation forest	0.5271	0.5383	0.0581

Table 4.1: Results of all anomaly detection system configurations on a validation set of the X-ray dataset

Table 4.1 shows the results achieved by different anomaly detection systems on a validation set containing 14125 1024-D embeddings.

The following metrics were used:

- F1 Macro — the unweighted mean of F1 score of each label. The best value is 1, which indicates perfect precision and recall for every label.
- AUROC — area under the receiving operating characteristic curve. Best value is 1. A random classifier can achieve a value of 0.5.
- Cohen's kappa — a coefficient measuring the agreement between the classifier and the correct labels. The best value is 1, which indicates perfect agreement between the two classifiers. The worst value is -1, which indicates total disagreement. Negative values imply that the classifier is performing worse than random guessing.

Interestingly, the embedding network, trained on ImageNet only, generated very sparse and low variance embeddings, while the model fine-tuned on the X-ray image dataset generated embeddings with high variance. This might mean that the ImageNet model prioritizes more generic features (e.g. edges and shapes) while the X-ray trained model can detect features more relevant to X-ray imaging.

4.1.1 Linear SVM version

The linear SVM version was used as a way to gauge whether the embeddings generated by the model trained on the X-ray dataset were even remotely useful for anomaly detection. If an SVM trained in a supervised manner showed abysmal performance (e.g. a negative Cohen's kappa or an F1 close to 0), that would be a strong indication that those embeddings are completely useless.

The results from table 4.1 show that the SVM trained using embeddings from a CNN that was trained on

ImageNet + X-ray datasets performed significantly worse than the one trained using embeddings from a CNN trained on ImageNet only. This could be an indication that the embeddings generated by an X-ray trained model are not linearly separable in 1024-D space. However, that model didn't perform too badly, which means the embeddings have some potential to be useful.

The performance differential between these models wasn't unexpected. The supervised anomaly detection task is very similar to simple image classification, in which the Densenet-201 (trained on ImageNet only) performs reasonably well.

4.1.2 Isolation forest version

The isolation forest models are used to provide results in a "real-world" scenario. These models were taught in an unsupervised manner, so their performance is a relatively good indication of how useful the embeddings actually are.

The embedding CNN trained on ImageNet + X-ray datasets outperformed the baseline model trained only on ImageNet in all metrics but F1 score. This means that in the process of training the triplet model, the embedding network managed to generate embeddings that have certain features which are useful for anomaly detection.

Overall, the performance of the isolation forest using embeddings from Densenet-201 (ImageNet + X-ray) is mediocre. Even though the model achieved a decent AUROC value of ~ 0.7 , the Cohen's kappa coefficient indicates that it misclassified many examples. This would not be acceptable in a "real-world" environment, where even one false negative could cause a loss of human life.

The isolation forest's results are by no means state-of-the-art. Previous works using multivariate Gaussian models trained on similar embeddings achieved significantly better performance [14].

4.2 Computational performance

The whole process of loading the image into memory, pre-processing and producing patches, feeding the patches to the CNN and then feeding the generated embeddings to an isolation forest model takes ~ 1.5 seconds on a GeForce GTX Titan X GPU and a Intel Xeon E5-2620 v3 CPU. This does not include the time it takes to load the CNN or the isolation forest model into memory.

The biggest bottleneck is passing image patches to the CNN (which includes copying them from RAM to VRAM) and then copying the resulting tensors back to RAM from VRAM, as the isolation forest model isn't loaded to VRAM. The image processing pipeline is also not very well optimized. With careful optimizations, the whole anomaly detection system could analyse an image in ~ 1 second (depending on the resolution of the input image).

4.3 Generated embedding suitability for anomaly detection

The results show that the embeddings generated by a network, trained in a Siamese or a triplet architecture to determine whether two X-ray images show the same parcel or not, show minor potential for use in anomaly detection.

The isolation forest model, trained using embeddings from a CNN trained on ImageNet + X-ray dataset outperformed the baseline model. This means that the triplet model's embeddings are more useful for anomaly detection than the baseline ones. Other anomaly detection models, such as the multivariate Gaussian model mentioned in Griffin *et al.* [14], might see an increase in performance using the embeddings from the triplet network, but further testing is required to see if that's the case.

Chapter 5

Conclusions and evaluation

5.1 Summary

This work assessed whether embeddings generated from a CNN, trained to determine whether two X-ray images show the same parcel, but from a different perspective, are useful for anomaly detection.

Multiple CNN architectures were trained in a Siamese and triplet setup to distinguish pairs/triplets of images. Their performance was evaluated and the best performing one was selected as the final embedding network. Online triplet mining impact on the models' performance was also discussed.

The usefulness of the generated embeddings was assessed by training a linear SVM and an isolation forest using those embeddings. 4 configurations of embedding CNN/anomaly detection model were evaluated. The best version of the anomaly detection system achieved $F_1 = 0.4048$, AUROC = 0.6922, $\kappa = 0.0908$. This version used embeddings generated by a triplet Densenet-201 model, trained on the X-ray dataset, which shows that embeddings extracted from a CNN, trained to distinguish pairs of X-ray images, can be used for anomaly detection.

Thus, all goals set at the beginning of the project were successfully achieved.

5.2 Future work

Potential future work includes:

- Evaluating the usefulness of these embeddings in a more sophisticated anomaly detection model. Griffin *et al.* [14], for instance, used a multivariate Gaussian model trained on similar embeddings. It would be interesting to see if the performance of that model would increase if it was trained using the embeddings generated by the triplet Densenet-201 model. Other anomaly detection models, such as the Local Outlier Factor or Elliptic Envelope could also be tested and evaluated.
- Evaluating the usefulness of embeddings generated by other CNN architectures. The Densenet-201 model was selected only because it achieved the lowest validation loss in both Siamese and triplet training setups. There is a possibility, that even though the other CNNs finished with a higher loss, their embeddings might contain richer and more useful features for anomaly detection. That can only be determined via extensive testing.
- Deeper analysis of the generated embeddings and a more detailed comparison of embeddings before and after training the network in a Siamese or triplet setup. Currently, only variance and sparsity of the generated embeddings was measured and discussed. There might be other important changes in the

embeddings after the network was trained. A detailed analysis might give more insight into why certain embeddings are good/bad for anomaly detection.

In conclusion, an anomaly detection system, that detects anomalies in X-ray security imaging, can achieve a moderate degree of performance using the embeddings generated by an image classification CNN, that was fine-tuned on a pair finding task in X-ray images. Better performance could most likely be achieved by using a more sophisticated anomaly detection model.

Bibliography

- [1] S. Hisamichi and N. I. SUGAWARA, “Mass screening for gastric cancer by x-ray examination,” *Japanese journal of clinical oncology*, vol. 14, no. 2, pp. 211–223, 1984.
- [2] H. Vogel and D. Haller, “Luggage and shipped goods,” *European journal of radiology*, vol. 63, no. 2, pp. 242–253, 2007.
- [3] A. Bolfing, T. Halbherr, and A. Schwaninger, “How image based factors and human factors contribute to threat detection performance in x-ray aviation security screening,” in *symposium of the Austrian HCI and Usability Engineering Group*, pp. 419–438, Springer, 2008.
- [4] G. Chen, “Understanding x-ray cargo imaging,” *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms*, vol. 241, no. 1-4, pp. 810–815, 2005.
- [5] G. Tomei, M. Cinti, D. Cerratti, and M. Fioravanti, “Attention, repetitive works, fatigue and stress,” *Annali di igiene: medicina preventiva e di comunita*, vol. 18, no. 5, pp. 417–429, 2006.
- [6] “Passenger jets carried dubai bomb.” <https://web.archive.org/web/20101101203811/http://english.aljazeera.net/news/middleeast/2010/10/20101031144429122829.html>, 2010. Accessed: 2010-11-01.
- [7] M. Mendes, A. Schwaninger, N. Strelbel, and S. Michel, “Why laptops should be screened separately when conventional x-ray screening is used,” in *2012 IEEE International Carnahan Conference on Security Technology (ICCST)*, pp. 267–273, IEEE, 2012.
- [8] B. Burns, “Why do laptops have to be removed when tablets can stay in the bag?,” 2012.
- [9] J. Andrews, T. Tanay, E. J. Morton, and L. D. Griffin, “Transfer representation-learning for anomaly detection,” in *Proceedings of the 33rd International Conference on Machine Learning. JMLR: New York, NY, USA. (2016)*, JMLR, 2016.
- [10] A. Mouton and T. P. Breckon, “A review of automated image understanding within 3d baggage computed tomography security screening,” *Journal of X-ray science and technology*, vol. 23, no. 5, pp. 531–555, 2015.
- [11] T. W. Rogers, N. Jaccard, E. J. Morton, and L. D. Griffin, “Automated x-ray image analysis for cargo security: Critical review and future promise,” *Journal of X-ray science and technology*, vol. 25, no. 1, pp. 33–56, 2017.
- [12] S. Akcay and T. Breckon, “Towards automatic threat detection: A survey of advances of deep learning within x-ray security imaging,” *arXiv preprint arXiv:2001.01293*, 2020.
- [13] J. T. Andrews, E. J. Morton, and L. D. Griffin, “Detecting anomalous data using auto-encoders,” *International Journal of Machine Learning and Computing*, vol. 6, no. 1, p. 21, 2016.
- [14] L. D. Griffin, M. Caldwell, J. T. Andrews, and H. Bohler, ““unexpected item in the bagging area”: Anomaly detection in x-ray security images,” *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 6, pp. 1539–1553, 2018.

- [15] J. An and S. Cho, “Variational autoencoder based anomaly detection using reconstruction probability,” *Special Lecture on IE*, vol. 2, no. 1, 2015.
- [16] “Self-supervised representation learning.” <https://lilianweng.github.io/lil-log/2019/11/10/self-supervised-learning.html>, 2019. Accessed: 2020-04-25.
- [17] S. Gidaris, P. Singh, and N. Komodakis, “Unsupervised representation learning by predicting image rotations,” *arXiv preprint arXiv:1803.07728*, 2018.
- [18] C. Doersch, A. Gupta, and A. A. Efros, “Unsupervised visual representation learning by context prediction,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1422–1430, 2015.
- [19] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, “Signature verification using a” siamese” time delay neural network,” in *Advances in neural information processing systems*, pp. 737–744, 1994.
- [20] G. Koch, R. Zemel, and R. Salakhutdinov, “Siamese neural networks for one-shot image recognition,” in *ICML deep learning workshop*, vol. 2, Lille, 2015.
- [21] L. Zheng, S. Duffner, K. Idrissi, C. Garcia, and A. Baskurt, “Siamese multi-layer perceptrons for dimensionality reduction and face identification,” *Multimedia Tools and Applications*, vol. 75, no. 9, pp. 5055–5073, 2016.
- [22] S. Chopra, R. Hadsell, and Y. LeCun, “Learning a similarity metric discriminatively, with application to face verification,” in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, vol. 1, pp. 539–546, IEEE, 2005.
- [23] R. Hadsell, S. Chopra, and Y. LeCun, “Dimensionality reduction by learning an invariant mapping,” in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, vol. 2, pp. 1735–1742, IEEE, 2006.
- [24] R. Gómez, “Understanding ranking loss, contrastive loss, margin loss, triplet loss, hinge loss and all those confusing names,” 2019.
- [25] E. Hoffer and N. Ailon, “Deep metric learning using triplet network,” in *International Workshop on Similarity-Based Pattern Recognition*, pp. 84–92, Springer, 2015.
- [26] O. Moindrot, “Triplet loss and online triplet mining in tensorflow,” 2018.
- [27] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32* (H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, eds.), pp. 8024–8035, Curran Associates, Inc., 2019.
- [28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [29] A. Bielski, “Siamese and triplet networks with online pair/triplet mining in pytorch.” <https://github.com/adambielski/siamese-triplet>, 2019.
- [30] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1492–1500, 2017.
- [31] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.

- [32] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [33] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016.
- [34] A. Hermans, L. Beyer, and B. Leibe, “In defense of the triplet loss for person re-identification,” *arXiv preprint arXiv:1703.07737*, 2017.
- [35] W. Xu, “Towards optimal one pass large scale learning with averaged stochastic gradient descent,” *arXiv preprint arXiv:1107.2490*, 2011.
- [36] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly detection: A survey,” *ACM computing surveys (CSUR)*, vol. 41, no. 3, pp. 1–58, 2009.

Appendix A

Quick guide

The following code is a script to run a novel image through the anomaly detection system.

```
1 import torch
2 import torchvision
3 import cv2
4 import numpy as np
5
6 from networks import TripletNet, EmbeddingNet
7 from joblib import load
8 from generate_vectors import produce_patches
9
10 model = TripletNet(EmbeddingNet(torchvision.models.densenet201()))
11
12 # Map location should be changed depending on whether CUDA is available or not
13 model.load_state_dict(torch.load('triplet_densenet201_m2.pth', map_location=torch.device('cpu')))
14 model.eval()
15 embedding_net = model.embedding_net
16
17 isf = load('isolation_forest.joblib')
18
19 embeddings = np.empty((0, 1024))
20
21 patches = produce_patches(PATH_TO_IMAGE)
22
23 for patch in patches:
24     with torch.no_grad():
25         vec = embedding_net(patch).cpu().detach().numpy()
26         if vec.shape == (1024,):
27             vec = np.expand_dims(a, axis=0)
28         embeddings = np.append(embeddings, vec, axis=0)
29
30 y_pred = min(isf.predict(embeddings))
31 print(y_pred)
```

`triplet_densenet_m2.pth` is a file containing the saved weights of the model, which was trained on the X-ray image dataset.

`isolation_forest.joblib` is a file containing the saved isolation forest, which was trained using the embeddings generated by the model saved in `triplet_densenet_m2.pth`.

The latest source code is provided at <https://github.com/mjanonis/anomaly-detection>. Any questions can be asked in the issue section.

Appendix B

Source code

B.1 datasets.py

```
1 # SPDX-License-Identifier: EUPL-1.2
2 # Unmodified code written by Adam Bielski is licensed under the BSD-3-Clause license
3 # All further additions and modifications: Copyright (c) 2020, Martynas Janonis
4 # Licensed under the EUPL-1.2-or-later
5
6 import numpy as np
7 import cv2
8 import os
9 import re
10 import csv
11
12 from PIL import Image
13 from random import shuffle, choice
14
15 import torch
16 from torch.utils.data import Dataset
17 from torch.utils.data.sampler import BatchSampler
18 from torchvision.transforms.functional import to_tensor
19 from torchvision.transforms import Normalize
20
21 from trim import trim_xray_images
22
23 CROP_MARGIN = 32
24
25 # Returns the index of the top view image in a pair
26 def top_view(pair):
27     img1 = int(re.sub("\D", "", pair[0]))
28     img2 = int(re.sub("\D", "", pair[1]))
29     return int(img1 > img2)
30
31
32 # Returns the index of the side view image in a pair
33 def side_view(pair):
34     img1 = int(re.sub("\D", "", pair[0]))
35     img2 = int(re.sub("\D", "", pair[1]))
36     return int(img1 < img2)
37
38 """
39 Generates two .csv files from the root directory with the structure:
40
41 IMAGE, POSITIVE_PAIR
42 """
43
44 def siameses_train_test_csv(root, train_size=0.8):
45     # Get all the filepaths of the images
46     filepaths = []
47     pairs = []
48     for subdir, dirs, files in os.walk(root):
49         for file in files:
```

```

50         filepath = subdir + os.sep + file
51         filepaths.append(filepath)
52
53     # Sort the filenames
54     filepaths.sort(key=lambda f: int(re.sub("\D", "", f)))
55
56     for i in range(len(filepaths)):
57         if i % 2 == 0:
58             pairs.append([filepaths[i], filepaths[i + 1]])
59         else:
60             pairs.append([filepaths[i], filepaths[i - 1]])
61
62     train_samples = int(train_size * len(pairs))
63
64     # Generate the training .csv
65     with open("train.csv", "w") as csvfile:
66         train = pairs[:train_samples].copy()
67         shuffle(train)
68         writer = csv.writer(csvfile)
69         writer.writerow(train)
70
71     # Generate the testing .csv
72
73     test = pairs[train_samples:]
74     for pair in test:
75         target = np.random.randint(0, 2)
76         if target:
77             pair.append(target)
78         else:
79             pr = pair
80             while pr == pair or pr[0] == pair[1]:
81                 pr = choice(test)
82                 v1_t = top_view(pair)
83
84             if v1_t == 0:
85                 pair[1] = pr[side_view(pr)]
86             else:
87                 pair[1] = pr[top_view(pr)]
88
89         pair.append(target)
90
91     with open("test.csv", "w") as csvfile:
92         shuffle(test)
93         writer = csv.writer(csvfile)
94         writer.writerow(test)
95
96
97     # Outputs all pairs to a .csv file
98     def triplet_train_test_csv(root, train_size=0.8):
99         # Get all the filepaths of the images
100        filepaths = []
101        pairs = []
102        for subdir, dirs, files in os.walk(root):
103            for file in files:
104                filepath = subdir + os.sep + file
105                filepaths.append(filepath)
106
107        # Sort the filenames
108        filepaths.sort(key=lambda f: int(re.sub("\D", "", f)))
109
110        for i in range(len(filepaths)):
111            if i % 2 == 0:
112                pairs.append([filepaths[i], filepaths[i + 1]])
113
114        train_samples = int(train_size * len(pairs))
115        shuffle(pairs)
116
117        # Generate the training .csv
118        with open("triplet_train.csv", "w") as csvfile:
119            writer = csv.writer(csvfile)
120            writer.writerow(pairs[:train_samples])
121
122        # Generate anchor, positive, negative triplets

```

```

123     test = pairs[train_samples:]
124     test_set = []
125     for idx in range(0, (len(pairs) - train_samples) * 2):
126         anchor = test[idx // 2][idx % 2]
127         positive = test[idx // 2][(idx + 1) % 2]
128         n_idx = idx
129         while n_idx == idx:
130             n_idx = np.random.randint(0, (len(pairs) - train_samples) * 2)
131         negative = test[n_idx // 2][(idx + 1) % 2]
132         test_set.append((anchor, positive, negative))
133
134     with open("triplet_test.csv", "w") as csvfile:
135         writer = csv.writer(csvfile)
136         writer.writerows(test_set)
137
138
139 def svm_train_test_csv(root_neg, root_pos, train_size=0.8):
140     # Get all the filepaths of the images
141     filepaths = []
142     for subdir, dirs, files in os.walk(root_neg):
143         for file in files:
144             filepath = subdir + os.sep + file
145             filepaths.append((filepath, 0))
146
147     for subdir, dirs, files in os.walk(root_pos):
148         for file in files:
149             filepath = subdir + os.sep + file
150             filepaths.append((filepath, 1))
151
152     shuffle(filepaths)
153     train_samples = int(train_size * len(filepaths))
154
155     with open("svm_train.csv", "w") as csvfile:
156         writer = csv.writer(csvfile)
157         writer.writerows(filepaths[:train_samples])
158
159     with open("svm_test.csv", "w") as csvfile:
160         writer = csv.writer(csvfile)
161         writer.writerows(filepaths[train_samples:])
162
163
164 class SiameseXRayParcels(Dataset):
165     def __init__(self, xray_csv, image_size=224, train=True, transform=False):
166         self.train = train
167         self.transform = transform
168         self.pairs = []
169         self.image_size = image_size
170
171         # Read the .csv
172         with open(xray_csv, newline="") as csvfile:
173             reader = csv.reader(csvfile)
174             for row in reader:
175                 self.pairs.append(row)
176
177     def __getitem__(self, index):
178
179         # Final images have to be at least 224x224
180
181         if self.train:
182             target = np.random.randint(0, 2)
183             img1 = cv2.imread(self.pairs[index][0])
184             if target:
185                 img2 = cv2.imread(self.pairs[index][1])
186             else:
187                 pr = self.pairs[index]
188                 while pr == self.pairs[index] or pr[0] == self.pairs[index][1]:
189                     pr = choice(self.pairs)
190                 v1_t = top_view(self.pairs[index])
191                 if v1_t == 0:
192                     img2 = cv2.imread(pr[side_view(pr)])
193                 else:
194                     img2 = cv2.imread(pr[top_view(pr)])
195

```

```

196     else:
197         img1 = cv2.imread(self.pairs[index][0])
198         img2 = cv2.imread(self.pairs[index][1])
199         target = int(self.pairs[index][2])
200
201     img1, img2 = trim_xray_images(img1, img2, self.image_size + CROP_MARGIN, target)
202
203     CROP_TOP_MAX_1 = img1.shape[0] - self.image_size
204     CROP_TOP_MAX_2 = img2.shape[0] - self.image_size
205     CROP_LEFT_MAX_1 = img1.shape[1] - self.image_size
206     CROP_LEFT_MAX_2 = img2.shape[1] - self.image_size
207
208     if self.train:
209         if self.transform:
210             # Apply random cropping
211             top_1 = np.random.randint(0, CROP_TOP_MAX_1)
212             top_2 = np.random.randint(0, CROP_TOP_MAX_2)
213             left_1 = np.random.randint(0, CROP_LEFT_MAX_1)
214             left_2 = np.random.randint(0, CROP_LEFT_MAX_2)
215
216             img1 = img1[
217                 top_1 : self.image_size + top_1, left_1 : self.image_size + left_1
218             ]
219
220             if target:
221                 img2 = img2[
222                     top_2 : self.image_size + top_2,
223                     left_1 : self.image_size + left_1,
224                 ]
225             else:
226                 img2 = img2[
227                     top_2 : self.image_size + top_2,
228                     left_2 : self.image_size + left_2,
229                 ]
230
231         else:
232             # Apply a center crop
233             top_1 = CROP_TOP_MAX_1 // 2
234             top_2 = CROP_TOP_MAX_2 // 2
235             left_1 = CROP_LEFT_MAX_1 // 2
236             left_2 = CROP_LEFT_MAX_2 // 2
237
238             img1 = img1[
239                 top_1 : self.image_size + top_1, left_1 : self.image_size + left_1
240             ]
241             img2 = img2[
242                 top_2 : self.image_size + top_2, left_2 : self.image_size + left_2
243             ]
244
245     else:
246         # Apply a center crop
247         top_1 = CROP_TOP_MAX_1 // 2
248         top_2 = CROP_TOP_MAX_2 // 2
249         left_1 = CROP_LEFT_MAX_1 // 2
250         left_2 = CROP_LEFT_MAX_2 // 2
251
252         img1 = img1[
253             top_1 : self.image_size + top_1, left_1 : self.image_size + left_1
254         ]
255         img2 = img2[
256             top_2 : self.image_size + top_2, left_2 : self.image_size + left_2
257         ]
258
259     # Normalize using the mean and std of ImageNet
260     norm = Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
261
262     # Convert the images to tensors and return
263     return (norm(to_tensor(img1)), norm(to_tensor(img2))), target
264
265 def __len__(self):
266     return len(self.pairs)
267
268

```

```

269 class TripletXRayParcels(Dataset):
270     def __init__(self, pair_csv, image_size=224, train=True, transform=False):
271         self.train = train
272         self.transform = transform
273         self.pairs = []
274         self.image_size = image_size
275
276         # Read the .csv
277         with open(pair_csv, newline="") as csvfile:
278             reader = csv.reader(csvfile)
279             for row in reader:
280                 self.pairs.append(row)
281
282     def __getitem__(self, index):
283         # Prepare the triplet
284         if self.train:
285             anchor = cv2.imread(self.pairs[index // 2][index % 2])
286             positive = cv2.imread(self.pairs[index // 2][(index + 1) % 2])
287             n_ind = index
288             while n_ind == index:
289                 n_ind = np.random.randint(0, len(self.pairs) * 2)
290             negative = cv2.imread(self.pairs[n_ind // 2][(index + 1) % 2])
291         else:
292             anchor = cv2.imread(self.pairs[index][0])
293             positive = cv2.imread(self.pairs[index][1])
294             negative = cv2.imread(self.pairs[index][2])
295
296         # Trim the images
297         anchor, positive = trim_xray_images(
298             anchor, positive, self.image_size + CROP_MARGIN, 1
299         )
300         _, negative = trim_xray_images(
301             anchor, negative, self.image_size + CROP_MARGIN, 0
302         )
303
304         # Get maximum crop size
305         CROP_TOP_MAX_A = anchor.shape[0] - self.image_size
306         CROP_TOP_MAX_P = positive.shape[0] - self.image_size
307         CROP_TOP_MAX_N = negative.shape[0] - self.image_size
308
309         CROP_LEFT_MAX_A = anchor.shape[1] - self.image_size
310         CROP_LEFT_MAX_P = positive.shape[1] - self.image_size
311         CROP_LEFT_MAX_N = negative.shape[1] - self.image_size
312
313         if self.train:
314             if self.transform:
315                 # Apply random cropping
316                 top_a = np.random.randint(0, CROP_TOP_MAX_A)
317                 top_p = np.random.randint(0, CROP_TOP_MAX_P)
318                 top_n = np.random.randint(0, CROP_TOP_MAX_N)
319
320                 left_a = np.random.randint(0, CROP_LEFT_MAX_A)
321                 # Anchor and positive must be cropped from the same X coordinate
322                 left_p = left_a
323                 left_n = np.random.randint(0, CROP_LEFT_MAX_N)
324
325                 anchor = anchor[
326                     top_a : self.image_size + top_a, left_a : self.image_size + left_a
327                 ]
328                 positive = positive[
329                     top_p : self.image_size + top_p, left_p : self.image_size + left_p
330                 ]
331                 negative = negative[
332                     top_n : self.image_size + top_n, left_n : self.image_size + left_n
333                 ]
334
335             else:
336                 # Apply a center crop
337                 top_a = CROP_TOP_MAX_A // 2
338                 top_p = CROP_TOP_MAX_P // 2
339                 top_n = CROP_TOP_MAX_N // 2
340
341                 left_a = CROP_LEFT_MAX_A // 2

```

```

342         left_p = CROP_LEFT_MAX_P // 2
343         left_n = CROP_LEFT_MAX_N // 2
344
345         anchor = anchor[
346             top_a : self.image_size + top_a, left_a : self.image_size + left_a
347         ]
348         positive = positive[
349             top_p : self.image_size + top_p, left_p : self.image_size + left_p
350         ]
351         negative = negative[
352             top_n : self.image_size + top_n, left_n : self.image_size + left_n
353         ]
354
355     else:
356         # Apply a center crop
357         top_a = CROP_TOP_MAX_A // 2
358         top_p = CROP_TOP_MAX_P // 2
359         top_n = CROP_TOP_MAX_N // 2
360
361         left_a = CROP_LEFT_MAX_A // 2
362         left_p = CROP_LEFT_MAX_P // 2
363         left_n = CROP_LEFT_MAX_N // 2
364
365         anchor = anchor[
366             top_a : self.image_size + top_a, left_a : self.image_size + left_a
367         ]
368         positive = positive[
369             top_p : self.image_size + top_p, left_p : self.image_size + left_p
370         ]
371         negative = negative[
372             top_n : self.image_size + top_n, left_n : self.image_size + left_n
373         ]
374
375     # Normalize using the mean and std of ImageNet
376     norm = Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
377
378     # Convert the images to tensors and return
379     return (
380         (
381             norm(to_tensor(anchor)),
382             norm(to_tensor(positive)),
383             norm(to_tensor(negative)),
384         ),
385         [],
386     )
387
388     def __len__(self):
389         if self.train:
390             return len(self.pairs * 2)
391         else:
392             return len(self.pairs)
393
394
395     class XRayParcels(Dataset):
396         def __init__(self, csv_file, image_size=224, train=True, transform=False):
397             self.train = train
398             self.transform = transform
399             self.data = []
400             self.image_size = image_size
401
402             # Read the .csv
403             with open(csv_file, newline="") as csvfile:
404                 reader = csv.reader(csvfile)
405                 for row in reader:
406                     self.data.append(row)
407
408         def __getitem__(self, index):
409             img = cv2.imread(self.data[index][0])
410             label = np.int(self.data[index][1])
411
412             img, _ = trim_xray_images(img, img, self.image_size + CROP_MARGIN, 0)
413
414             CROP_TOP_MAX = img.shape[0] - self.image_size

```

```

415     CROP_LEFT_MAX = img.shape[1] - self.image_size
416
417     if self.train:
418         if self.transform:
419             # Apply random cropping
420             top = np.random.randint(0, CROP_TOP_MAX)
421             left = np.random.randint(0, CROP_LEFT_MAX)
422
423             img = img[top : self.image_size + top, left : self.image_size + left]
424
425         else:
426             # Apply a center crop
427             top = CROP_TOP_MAX // 2
428             left = CROP_LEFT_MAX // 2
429
430             img = img[top : self.image_size + top, left : self.image_size + left]
431
432     else:
433         # Apply a center crop
434         top = CROP_TOP_MAX // 2
435         left = CROP_LEFT_MAX // 2
436
437         img = img[top : self.image_size + top, left : self.image_size + left]
438
439     # Normalize using the mean and std of ImageNet
440     norm = Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
441
442     return norm(to_tensor(img)), label
443
444     def __len__(self):
445         return len(self.data)

```

B.2 generate_vectors.py

```

1  # SPDX-License-Identifier: EUPL-1.2
2  # Copyright (c) 2020, Martynas Janonis
3
4  # Licensed under the EUPL-1.2-or-later
5
6  import os
7  import re
8  import numpy as np
9  import cv2
10 import math
11 import torch
12
13 import trim
14
15 from skimage.util import view_as_windows
16 from torchvision.transforms import Normalize
17 from torchvision.transforms.functional import to_tensor
18
19
20 def generate_vectors(root, dest, model):
21     model.eval()
22
23     filepaths = []
24     for subdir, dirs, files in os.walk(root):
25         for file in files:
26             filepath = subdir + os.sep + file
27             filepaths.append(filepath)
28
29     # Sort the filenames
30     filepaths.sort(key=lambda f: int(re.sub("\D", "", f)))
31
32     for file in filepaths:
33         print(file)
34         patches = produce_patches(file)
35         print(len(patches))
36         f = open(dest + os.path.splitext(os.path.basename(file))[0] + ".vec", "ab")
37         for patch in patches:
38             with torch.no_grad():

```

```

39         vec = model(patch).cpu().detach().numpy()
40         np.savetxt(f, vec, delimiter=",")
41     f.close()
42
43
44 def produce_patches(path):
45     img = cv2.imread(path)
46     patches = []
47     # Trim the image
48     img = trim.trim_xray_images(img, img, 224, 0)[0]
49
50     # Add a white border around the image to make sure that stride of 112 covers it entirely
51     img = cv2.copyMakeBorder(
52         img,
53         (round_to(img.shape[0], 112) - img.shape[0]) // 2,
54         (round_to(img.shape[0], 112) - img.shape[0]) // 2,
55         (round_to(img.shape[1], 112) - img.shape[1]) // 2,
56         (round_to(img.shape[1], 112) - img.shape[1]) // 2,
57         borderType=cv2.BORDER_CONSTANT,
58         value=[255, 255, 255],
59     )
60
61     window = view_as_windows(img, (224, 224, 3))
62
63     # Normalize using the mean and std of ImageNet
64     norm = Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
65
66     for h in range(0, window.shape[0], 112):
67         for w in range(0, window.shape[1], 112):
68             patches.append(norm(to_tensor(window[h][w][0])).unsqueeze(0))
69
70     return patches
71
72
73 def round_to(num, rnd):
74     return int(math.ceil(num / rnd)) * rnd

```

B.3 losses.py

```

1  # SPDX-License-Identifier: BSD-3-Clause
2  # Copyright (c) 2019, Adam Bielski
3
4  import torch
5  import torch.nn as nn
6  import torch.nn.functional as F
7
8
9  class ContrastiveLoss(nn.Module):
10    """
11    Contrastive loss
12    Takes embeddings of two samples and a target label == 1 if samples are from the same class and label == 0 otherwise
13    """
14
15    def __init__(self, margin):
16        super(ContrastiveLoss, self).__init__()
17        self.margin = margin
18        self.eps = 1e-9
19
20    def forward(self, output1, output2, target, size_average=True):
21        distances = (output2 - output1).pow(2).sum(1) # squared distances
22        losses = 0.5 * (
23            target.float() * distances
24            + (1 + -1 * target).float()
25            * F.relu(self.margin - (distances + self.eps).sqrt()).pow(2)
26        )
27        return losses.mean() if size_average else losses.sum()
28
29
30 class OnlineContrastiveLoss(nn.Module):
31    """
32    Online Contrastive loss

```

```

33     Takes a batch of embeddings and corresponding labels.
34     Pairs are generated using pair_selector object that take embeddings and targets and return indices of positive
35     and negative pairs
36     """
37
38     def __init__(self, margin, pair_selector):
39         super(OnlineContrastiveLoss, self).__init__()
40         self.margin = margin
41         self.pair_selector = pair_selector
42
43     def forward(self, embeddings, target):
44         positive_pairs, negative_pairs = self.pair_selector.get_pairs(
45             embeddings, target
46         )
47         if embeddings.is_cuda:
48             positive_pairs = positive_pairs.cuda()
49             negative_pairs = negative_pairs.cuda()
50         positive_loss = (
51             (embeddings[positive_pairs[:, 0]] - embeddings[positive_pairs[:, 1]])
52             .pow(2)
53             .sum(1)
54         )
55         negative_loss = F.relu(
56             self.margin
57             - (embeddings[negative_pairs[:, 0]] - embeddings[negative_pairs[:, 1]])
58             .pow(2)
59             .sum(1)
60             .sqrt()
61         ).pow(2)
62         loss = torch.cat([positive_loss, negative_loss], dim=0)
63         return loss.mean()
64
65
66 class OnlineTripletLoss(nn.Module):
67     """
68     Online Triplets loss
69     Takes a batch of embeddings and corresponding labels.
70     Triplets are generated using triplet_selector object that take embeddings and targets and return indices of
71     triplets
72     """
73
74     def __init__(self, margin, triplet_selector):
75         super(OnlineTripletLoss, self).__init__()
76         self.margin = margin
77         self.triplet_selector = triplet_selector
78
79     def forward(self, embeddings, target):
80
81         triplets = self.triplet_selector.get_triplets(embeddings, target)
82
83         if embeddings.is_cuda:
84             triplets = triplets.cuda()
85
86         ap_distances = (
87             (embeddings[triplets[:, 0]] - embeddings[triplets[:, 1]]).pow(2).sum(1)
88         ) # .pow(.5)
89         an_distances = (
90             (embeddings[triplets[:, 0]] - embeddings[triplets[:, 2]]).pow(2).sum(1)
91         ) # .pow(.5)
92         losses = F.relu(ap_distances - an_distances + self.margin)
93
94         return losses.mean(), len(triplets)

```

B.4 metrics.py

```

1 # SPDX-License-Identifier: EUPL-1.2
2
3 # Unmodified code written by Adam Bielski is licensed under the BSD-3-Clause license
4
5 # All further additions and modifications: Copyright (c) 2020, Martynas Janonis
6

```

```

7 # Licensed under the EUPL-1.2-or-later
8
9 from torch.nn.modules import distance
10
11
12 class Metric:
13     def __init__(self):
14         pass
15
16     def __call__(self, outputs, target, loss):
17         raise NotImplementedError
18
19     def reset(self):
20         raise NotImplementedError
21
22     def value(self):
23         raise NotImplementedError
24
25     def name(self):
26         raise NotImplementedError
27
28
29 class AccumulatedAccuracyMetric(Metric):
30     """
31     Works with classification model
32     """
33
34     def __init__(self):
35         self.correct = 0
36         self.total = 0
37
38     def __call__(self, outputs, target, loss):
39         pred = outputs[0].data.max(1, keepdim=True)[1]
40         self.correct += pred.eq(target[0].data.view_as(pred)).cpu().sum()
41         self.total += target[0].size(0)
42
43     def reset(self):
44         self.correct = 0
45         self.total = 0
46
47
48     def value(self):
49         return 100 * float(self.correct) / self.total
50
51     def name(self):
52         return "Accuracy"
53
54
55 class AccumulatedDistanceAccuracyMetric(Metric):
56     """
57     If the distance between the two outputs is less than the margin,
58     classify as positive; else negative
59     """
60
61     def __init__(self, margin):
62         self.correct = 0
63         self.total = 0
64         self.margin = margin
65
66     def __call__(self, outputs, target, loss):
67         pred = distance.PairwiseDistance().forward(outputs[0], outputs[1])
68         pred = pred.flatten() < self.margin
69         self.correct += sum(pred == target[0])
70         self.total += target[0].size(0)
71
72     def reset(self):
73         self.correct = 0
74         self.total = 0
75
76
77     def value(self):
78         return 100 * float(self.correct) / self.total
79

```

```

80     def name(self):
81         return "Accuracy"
82
83
84 class TripletAccumulatedDistanceAccuracyMetric(Metric):
85     """
86     If the distance between the anchor and the positive is less than the distance between the anchor and the negative,
87     classify as positive; else negative
88     """
89
90     def __init__(self):
91         self.correct = 0
92         self.total = 0
93
94     def __call__(self, outputs, target, loss):
95         dist_pos = distance.PairwiseDistance().forward(outputs[0], outputs[1])
96         dist_neg = distance.PairwiseDistance().forward(outputs[0], outputs[2])
97         pred = dist_pos.flatten() < dist_neg.flatten()
98         self.correct += sum(pred == True)
99         self.total += pred.size(0)
100    return self.value()
101
102    def reset(self):
103        self.correct = 0
104        self.total = 0
105
106    def value(self):
107        return 100 * float(self.correct) / self.total
108
109    def name(self):
110        return "Accuracy"

```

B.5 networks.py

```

1 # SPDX-License-Identifier: EUPL-1.2
2
3 # Unmodified code written by Adam Bielski is licensed under the BSD-3-Clause license
4
5 # All further additions and modifications: Copyright (c) 2020, Martynas Janonis
6
7 # Licensed under the EUPL-1.2-or-later
8
9 import torch.nn as nn
10
11
12 class EmbeddingNet(nn.Module):
13     def __init__(self, network):
14         super().__init__()
15         self.convnet = network
16
17         # Strip the classification layer
18         self.convnet = nn.Sequential(*list(self.convnet.children())[:-1])
19         # Add an avg pooling layer (need output shape of [?, 1024, 1, 1])
20         self.aavgp3d = nn.AdaptiveAvgPool3d(output_size=(1024, 1, 1))
21         self.flat = nn.Flatten()
22
23     def forward(self, x):
24         output = self.convnet(x)
25         output = self.aavgp3d(output)
26         output = self.flat(output)
27         return output
28
29     def get_embedding(self, x):
30         return self.forward(x)
31
32
33 class SiameseNet(nn.Module):
34     def __init__(self, embedding_net):
35         super(SiameseNet, self).__init__()
36         self.embedding_net = embedding_net
37

```

```

38     def forward(self, x1, x2):
39         output1 = self.embedding_net(x1)
40         output2 = self.embedding_net(x2)
41         return output1, output2
42
43     def get_embedding(self, x):
44         return self.embedding_net(x)
45
46
47 class TripletNet(nn.Module):
48     def __init__(self, embedding_net):
49         super(TripletNet, self).__init__()
50         self.embedding_net = embedding_net
51
52     def forward(self, x1, x2, x3):
53         output1 = self.embedding_net(x1)
54         output2 = self.embedding_net(x2)
55         output3 = self.embedding_net(x3)
56         return output1, output2, output3
57
58     def get_embedding(self, x):
59         return self.embedding_net(x)

```

B.6 train_isolation.py

```

1  from sklearn.ensemble import IsolationForest
2  from random import shuffle
3  from joblib import dump, load
4
5  import numpy as np
6  import os
7  import re
8
9
10 def load_vectors(root):
11     X = np.empty((0, 1024))
12
13     filepaths = []
14     for subdir, dirs, files in os.walk(root):
15         for file in files:
16             filepath = subdir + os.sep + file
17             filepaths.append(filepath)
18
19     # Sort the filenames
20     filepaths.sort(key=lambda f: int(re.sub("\D", "", f)))
21
22     for file in filepaths:
23         a = np.loadtxt(file, delimiter=",")
24         if a.shape == (1024,):
25             a = np.expand_dims(a, axis=0)
26
27         X = np.append(X, a, axis=0)
28
29     print(X.shape)
30     return X
31
32
33 X_train = load_vectors("./vectors")
34 shuffle(X_train)
35
36 isf = IsolationForest(n_estimators=1024, verbose=1)
37
38 isf.fit(X_train)
39
40 dump(isf, "isolation_forest.joblib")

```

B.7 train_siamese.py

```
1 # SPDX-License-Identifier: EUPL-1.2
2
3 # Unmodified code written by Adam Bielski is licensed under the BSD-3-Clause license
4
5 # All further additions and modifications: Copyright (c) 2020, Martynas Janonis
6
7 # Licensed under the EUPL-1.2-or-later
8
9 import torch
10 from torch.optim import lr_scheduler
11 import torch.optim as optim
12 from torchvision.models import resnext101_32x8d, densenet201
13
14 from trainer import fit
15
16 cuda = torch.cuda.is_available()
17 from networks import EmbeddingNet, SiameseNet
18 from datasets import SiameseXRayParcels
19 from losses import ContrastiveLoss
20 from metrics import AccumulatedDistanceAccuracyMetric
21
22 siamese_train_dataset = SiameseXRayParcels("train.csv", train=True, transform=True)
23 siamese_test_dataset = SiameseXRayParcels("test.csv", train=False, transform=False)
24 batch_size = 1
25 kwargs = {"num_workers": 1, "pin_memory": True} if cuda else {}
26 siamese_train_loader = torch.utils.data.DataLoader(
27     siamese_train_dataset, batch_size=batch_size, shuffle=True, **kwargs
28 )
29 siamese_test_loader = torch.utils.data.DataLoader(
30     siamese_test_dataset, batch_size=batch_size, shuffle=False, **kwargs
31 )
32
33 margin = 2.0
34 embedding_net = EmbeddingNet(resnext101_32x8d(pretrained=True))
35 model = SiameseNet(embedding_net)
36 if cuda:
37     model.cuda()
38
39 loss_fn = ContrastiveLoss(margin)
40 lr = 1e-3
41 optimizer = optim.Adam(model.parameters(), lr=lr)
42 scheduler = lr_scheduler.StepLR(optimizer, 1, gamma=0.99, last_epoch=-1)
43 n_epochs = 50
44 log_interval = 10
45
46 fit(
47     siamese_train_loader,
48     siamese_test_loader,
49     model,
50     loss_fn,
51     optimizer,
52     scheduler,
53     n_epochs,
54     cuda,
55     log_interval,
56     [AccumulatedDistanceAccuracyMetric(margin // 2)],
57 )
```

B.8 train_svm.py

```
1 # SPDX-License-Identifier: EUPL-1.2
2 # Copyright (c) 2020, Martynas Janonis
3
4 # Licensed under the EUPL-1.2-or-later
5
6 import torch
7 import sys
8 import numpy as np
9
```

```

10 from torchvision.models import resnext101_32x8d, densenet201
11 from sklearn.linear_model import SGDClassifier
12 from sklearn.metrics import f1_score, roc_auc_score, accuracy_score, cohen_kappa_score
13 from networks import EmbeddingNet, TripletNet
14 from datasets import XRayParcels
15 from joblib import dump, load
16
17 cuda = torch.cuda.is_available()
18 device = torch.device("cuda") if cuda else torch.device("cpu")
19
20 xray_train_dataset = XRayParcels("svm_train.csv", train=True, transform=True)
21 xray_test_dataset = XRayParcels("svm_test.csv", train=False, transform=False)
22 batch_size = 64
23 kwargs = {"num_workers": 1, "pin_memory": True} if cuda else {}
24 xray_train_loader = torch.utils.data.DataLoader(
25     xray_train_dataset, batch_size=batch_size, shuffle=True, **kwargs
26 )
27 xray_test_loader = torch.utils.data.DataLoader(
28     xray_test_dataset, batch_size=batch_size, shuffle=False, **kwargs
29 )
30
31 # Initialize and load the embedding network
32 embedding_net = EmbeddingNet(densenet201())
33 model = TripletNet(embedding_net)
34 model.load_state_dict(torch.load("triplet_densenet201_m2.pth", map_location=device))
35 model.to(device)
36 model.eval()
37
38 # Initialize the SVM
39 svm = SGDClassifier(
40     loss="hinge", verbose=0, class_weight={0: 1, 1: 50}, warm_start=True, average=True
41 )
42
43 n_epochs = 50
44 highest_f1 = 0
45
46 for epoch in range(n_epochs):
47
48     print("Starting epoch {}".format(epoch))
49     # Train stage
50     # Generate #batch_size vectors to pass as a dataset to the SVM
51     for batch_idx, (data, target) in enumerate(xray_train_loader):
52         target = target if len(target) > 0 else None
53         if not type(data) in (tuple, list):
54             data = (data,)
55
56         data = tuple(d.to(device) for d in data)
57         if target is not None:
58             target = target.to(device)
59
60         with torch.no_grad():
61             vectors = model.embedding_net(*data)
62
63         # Convert from PyTorch tensors to NumPy arrays
64         vectors = vectors.detach().cpu().numpy()
65         target = target.detach().cpu().numpy()
66
67         # Do one epoch of SGD for the SVM
68         svm.partial_fit(vectors, target, classes=[0, 1])
69
70         message = "Train: [{}/{} ({:.0f}%)]".format(
71             batch_idx * len(data[0]),
72             len(xray_train_loader.dataset),
73             100.0 * batch_idx / len(xray_train_loader),
74         )
75
76         sys.stdout.write("\x1b[2K") # Clear to the end of line
77         sys.stdout.write("\r" + message)
78         sys.stdout.flush()
79
80     print()
81     print("Starting validation")
82     # Test stage

```

```

83     # Generate #batch_size vectors to pass as a dataset to the SVM
84     y_pred = []
85     y_true = []
86     for batch_idx, (data, target) in enumerate(xray_test_loader):
87         target = target if len(target) > 0 else None
88         if not type(data) in (tuple, list):
89             data = (data,)
90
91         data = tuple(d.to(device) for d in data)
92         if target is not None:
93             target = target.to(device)
94
95         with torch.no_grad():
96             vectors = model.embedding_net(*data)
97
98         # Convert from PyTorch tensors to NumPy arrays
99         vectors = vectors.detach().cpu().numpy()
100
101        y_true = np.append(y_true, target.detach().cpu().numpy())
102        y_pred = np.append(y_pred, svm.predict(vectors))
103
104    print(
105        "Epoch {}/. Validation set: Avg. accuracy: {:.4f}, avg. F1 score: {:.4f}, avg. AUC: {:.4f}, avg. Kappa {:.4f}" .format(
106            epoch,
107            n_epochs,
108            accuracy_score(y_true, y_pred),
109            f1_score(y_true, y_pred),
110            roc_auc_score(y_true, y_pred),
111            cohen_kappa_score(y_true, y_pred),
112        )
113    )
114
115    # Save the model if F1 is larger
116    if f1_score(y_true, y_pred) > highest_f1:
117        print("F1 score increased, saving model")
118        dump(svm, "svm.joblib")
119        highest_f1 = f1_score(y_true, y_pred)

```

B.9 train_triplet.py

```

1  # SPDX-License-Identifier: EUPL-1.2
2
3  # Unmodified code written by Adam Bielski is licensed under the BSD-3-Clause license
4
5  # All further additions and modifications: Copyright (c) 2020, Martynas Janonis
6
7  # Licensed under the EUPL-1.2-or-later
8
9  import torch
10 from torch.optim import lr_scheduler
11 import torch.optim as optim
12 from torchvision.models import resnext101_32x8d, densenet201
13
14 from trainer import fit
15
16 cuda = torch.cuda.is_available()
17 from networks import EmbeddingNet, TripletNet
18 from datasets import TripletXRayParcels
19 from torch.nn import TripletMarginLoss
20 from metrics import TripletAccumulatedDistanceAccuracyMetric
21
22 triplet_train_dataset = TripletXRayParcels(
23     "triplet_train.csv", train=True, transform=True
24 )
25 triplet_test_dataset = TripletXRayParcels(
26     "triplet_test.csv", train=False, transform=False
27 )
28 batch_size = 1
29 kwargs = {"num_workers": 1, "pin_memory": True} if cuda else {}
30 triplet_train_loader = torch.utils.data.DataLoader(
31     triplet_train_dataset, batch_size=batch_size, shuffle=True, **kwargs

```

```

32     )
33     triplet_test_loader = torch.utils.data.DataLoader(
34         triplet_test_dataset, batch_size=batch_size, shuffle=False, **kwargs
35     )
36
37     margin = 2.0
38     embedding_net = EmbeddingNet(resnext101_32x8d(pretrained=True))
39     model = TripletNet(embedding_net)
40     if cuda:
41         model.cuda()
42
43     loss_fn = TripletMarginLoss(margin=margin)
44     lr = 1e-3
45     optimizer = optim.Adam(model.parameters(), lr=lr)
46     scheduler = lr_scheduler.StepLR(optimizer, 1, gamma=0.99, last_epoch=-1)
47     n_epochs = 50
48     log_interval = 10
49
50     fit(
51         triplet_train_loader,
52         triplet_test_loader,
53         model,
54         loss_fn,
55         optimizer,
56         scheduler,
57         n_epochs,
58         cuda,
59         log_interval,
60         [TripletAccumulatedDistanceAccuracyMetric()],
61     )

```

B.10 trainer.py

```

1  # SPDX-License-Identifier: EUPL-1.2
2
3  # Unmodified code written by Adam Bielski is licensed under the BSD-3-Clause license
4
5  # All further additions and modifications: Copyright (c) 2020, Martynas Janonis
6
7  # Licensed under the EUPL-1.2-or-later
8
9  import torch
10 import numpy as np
11 from torch.utils.tensorboard import SummaryWriter
12 import sys
13
14
15 def fit(
16     train_loader,
17     val_loader,
18     model,
19     loss_fn,
20     optimizer,
21     scheduler,
22     n_epochs,
23     cuda,
24     log_interval,
25     metrics=[],
26     start_epoch=0,
27 ):
28     """
29     Loaders, model, loss function and metrics should work together for a given task,
30     i.e. The model should be able to process data output of loaders,
31     loss function should process target output of loaders and outputs from the model
32
33     Examples: Classification: batch loader, classification model, NLL loss, accuracy metric
34     Siamese network: Siamese loader, siamese model, contrastive loss
35     Online triplet learning: batch loader, embedding model, online triplet loss
36     """
37
38     writer = SummaryWriter()

```

```

39     lowest_val_loss = float("inf")
40     for epoch in range(start_epoch, n_epochs):
41
42         # Train stage
43         train_loss, metrics = train_epoch(
44             train_loader, model, loss_fn, optimizer, cuda, log_interval, metrics
45         )
46
47         message = "Epoch: {}/{}. Train set: Average loss: {:.4f}".format(
48             epoch + 1, n_epochs, train_loss
49         )
50         for metric in metrics:
51             message += "\t{}: {}".format(metric.name(), metric.value())
52             writer.add_scalar(metric.name() + "/train", metric.value(), epoch)
53
54         writer.add_scalar("Loss/train", train_loss, epoch)
55
56         val_loss, metrics = test_epoch(val_loader, model, loss_fn, cuda, metrics)
57         val_loss /= len(val_loader)
58
59         scheduler.step()
60
61         message += "\nEpoch: {}/{}. Validation set: Average loss: {:.4f}".format(
62             epoch + 1, n_epochs, val_loss
63         )
64
65         writer.add_scalar("Loss/test", val_loss, epoch)
66
67         for metric in metrics:
68             message += "\t{}: {}".format(metric.name(), metric.value())
69             writer.add_scalar(metric.name() + "/test", metric.value(), epoch)
70
71         print(message)
72
73         if val_loss < lowest_val_loss:
74             print(
75                 "{:.6f} < {:.6f}: saving model_state_dict to ./model.pth".format(
76                     val_loss, lowest_val_loss
77                 )
78             )
79             lowest_val_loss = val_loss
80             torch.save(model.state_dict(), "./model.pth")
81
82         writer.flush()
83
84
85     def train_epoch(train_loader, model, loss_fn, optimizer, cuda, log_interval, metrics):
86         for metric in metrics:
87             metric.reset()
88
89         model.train()
90         losses = []
91         total_loss = 0
92
93         for batch_idx, (data, target) in enumerate(train_loader):
94             target = target if len(target) > 0 else None
95             if not type(data) in (tuple, list):
96                 data = (data,)
97             if cuda:
98                 data = tuple(d.cuda() for d in data)
99                 if target is not None:
100                     target = target.cuda()
101
102             optimizer.zero_grad()
103             outputs = model(*data)
104
105             if type(outputs) not in (tuple, list):
106                 outputs = (outputs,)
107
108             loss_inputs = outputs
109             if target is not None:
110                 target = (target,)
111                 loss_inputs += target

```

```

112
113     loss_outputs = loss_fn(*loss_inputs)
114     loss = loss_outputs[0] if type(loss_outputs) in (tuple, list) else loss_outputs
115     losses.append(loss.item())
116     total_loss += loss.item()
117     loss.backward()
118     optimizer.step()
119     for metric in metrics:
120         metric(outputs, target, loss_outputs)
121
122     if batch_idx % log_interval == 0:
123         message = "Train: [{}/{} ({:.0f}%)]\tLoss: {:.6f}".format(
124             batch_idx * len(data[0]),
125             len(train_loader.dataset),
126             100.0 * batch_idx / len(train_loader),
127             np.mean(losses),
128         )
129         for metric in metrics:
130             message += "\t{}: {}".format(metric.name(), metric.value())
131
132         sys.stdout.write("\x1b[2K") # Clear to the end of line
133         sys.stdout.write("\r" + message)
134         sys.stdout.flush()
135     losses = []
136
137     # Add a new line
138     print()
139     total_loss /= batch_idx + 1
140     return total_loss, metrics
141
142
143 def test_epoch(val_loader, model, loss_fn, cuda, metrics):
144     with torch.no_grad():
145         for metric in metrics:
146             metric.reset()
147         model.eval()
148         val_loss = 0
149         for batch_idx, (data, target) in enumerate(val_loader):
150             target = target if len(target) > 0 else None
151             if not type(data) in (tuple, list):
152                 data = (data,)
153             if cuda:
154                 data = tuple(d.cuda() for d in data)
155                 if target is not None:
156                     target = target.cuda()
157
158             outputs = model(*data)
159
160             if type(outputs) not in (tuple, list):
161                 outputs = (outputs,)
162             loss_inputs = outputs
163             if target is not None:
164                 target = (target,)
165                 loss_inputs += target
166
167             loss_outputs = loss_fn(*loss_inputs)
168             loss = (
169                 loss_outputs[0] if type(loss_outputs) in (tuple, list) else loss_outputs
170             )
171             val_loss += loss.item()
172
173         for metric in metrics:
174             metric(outputs, target, loss_outputs)
175
176     return val_loss, metrics

```

B.11 trim.py

```

1 # SPDX-License-Identifier: EUPL-1.2
2 # Copyright (c) 2020, Martynas Janonis
3

```

```

4  # Licensed under the EUPL-1.2-or-later
5
6  import cv2
7  import numpy as np
8
9
10 def get_bounding_box(img):
11     # convert to grayscale
12     gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
13
14     # threshold to get just the contents
15     retval, thresh_gray = cv2.threshold(
16         gray, thresh=205, maxval=255, type=cv2.THRESH_BINARY
17     )
18
19     # find where the contents are and make a cropped region
20
21     # find where the black pixels are
22     points = np.argwhere(thresh_gray == 0)
23
24     # store them in x,y coordinates instead of row,col indices
25     points = np.fliplr(
26         points
27     )
28
29     # create a rectangle around those points
30     x, y, w, h = cv2.boundingRect(points)
31
32     return x, y, w, h
33
34
35 def trim_xray_images(img1, img2, image_size, label):
36     # add a white border around the images to make sure that there's enough space for cropping
37     img1 = cv2.copyMakeBorder(
38         img1,
39         img2.shape[0] // 2 + image_size,
40         img2.shape[0] // 2 + image_size,
41         img2.shape[1] // 2 + image_size,
42         img2.shape[1] // 2 + image_size,
43         borderType=cv2.BORDER_CONSTANT,
44         value=[255, 255, 255],
45     )
46     img2 = cv2.copyMakeBorder(
47         img2,
48         img1.shape[0] // 2 + image_size,
49         img1.shape[0] // 2 + image_size,
50         img1.shape[1] // 2 + image_size,
51         img1.shape[1] // 2 + image_size,
52         borderType=cv2.BORDER_CONSTANT,
53         value=[255, 255, 255],
54     )
55
56     # get bounding box for both images
57     x1, y1, w1, h1 = get_bounding_box(img1)
58     x2, y2, w2, h2 = get_bounding_box(img2)
59
60     # the bounding box has to be at least the size of the image_size
61     if w1 < image_size:
62         x1 -= (image_size - w1) // 2
63         w1 += image_size - w1
64     if h1 < image_size:
65         y1 -= (image_size - h1) // 2
66         h1 += image_size - h1
67
68     if w2 < image_size:
69         x2 -= (image_size - w2) // 2
70         w2 += image_size - w2
71     if h2 < image_size:
72         y2 -= (image_size - h2) // 2
73         h2 += image_size - h2
74
75     # if two images are of the same parcel
76     # make sure that the widths match

```

```

77     if label:
78         if w1 > w2:
79             x2 -= (w1 - w2) // 2
80             w2 = w1
81         elif w2 > w1:
82             x1 -= (w2 - w1) // 2
83             w1 = w2
84
85     # adjust the border if x or y is negative
86     img1 = cv2.copyMakeBorder(
87         img1,
88         abs(min(0, y1)),
89         0,
90         abs(min(0, x1)),
91         0,
92         borderType=cv2.BORDER_CONSTANT,
93         value=[255, 255, 255],
94     )
95     x1 = max(x1, 0)
96     y1 = max(y1, 0)
97     img2 = cv2.copyMakeBorder(
98         img2,
99         abs(min(0, y2)),
100        0,
101        abs(min(0, x2)),
102        0,
103        borderType=cv2.BORDER_CONSTANT,
104        value=[255, 255, 255],
105    )
106     x2 = max(x2, 0)
107     y2 = max(y2, 0)
108
109     return img1[y1 : y1 + h1, x1 : x1 + w1], img2[y2 : y2 + h2, x2 : x2 + w2]

```

B.12 utils.py

```

1  # SPDX-License-Identifier: BSD-3-Clause
2  # Copyright (c) 2019, Adam Bielski
3
4  from itertools import combinations
5
6  import numpy as np
7  import torch
8
9
10 def pdist(vectors):
11     distance_matrix = (
12         -2 * vectors.mm(torch.t(vectors))
13         + vectors.pow(2).sum(dim=1).view(1, -1)
14         + vectors.pow(2).sum(dim=1).view(-1, 1)
15     )
16     return distance_matrix
17
18
19 class PairSelector:
20     """
21     Implementation should return indices of positive pairs and negative pairs that will be passed to compute
22     Contrastive Loss
23     return positive_pairs, negative_pairs
24     """
25
26     def __init__(self):
27         pass
28
29     def get_pairs(self, embeddings, labels):
30         raise NotImplementedError
31
32
33 class AllPositivePairSelector(PairSelector):
34     """
35     Discards embeddings and generates all possible pairs given labels.

```

```

36     If balance is True, negative pairs are a random sample to match the number of positive samples
37     """
38
39     def __init__(self, balance=True):
40         super(AllPositivePairSelector, self).__init__()
41         self.balance = balance
42
43     def get_pairs(self, embeddings, labels):
44         labels = labels.cpu().data.numpy()
45         all_pairs = np.array(list(combinations(range(len(labels)), 2)))
46         all_pairs = torch.LongTensor(all_pairs)
47         positive_pairs = all_pairs[
48             (labels[all_pairs[:, 0]] == labels[all_pairs[:, 1]]).nonzero()
49         ]
50         negative_pairs = all_pairs[
51             (labels[all_pairs[:, 0]] != labels[all_pairs[:, 1]]).nonzero()
52         ]
53         if self.balance:
54             negative_pairs = negative_pairs[
55                 torch.randperm(len(negative_pairs))[: len(positive_pairs)]
56             ]
57
58         return positive_pairs, negative_pairs
59
60
61 class HardNegativePairSelector(PairSelector):
62     """
63     Creates all possible positive pairs. For negative pairs, pairs with smallest distance are taken into consideration,
64     matching the number of positive pairs.
65     """
66
67     def __init__(self, cpu=True):
68         super(HardNegativePairSelector, self).__init__()
69         self.cpu = cpu
70
71     def get_pairs(self, embeddings, labels):
72         if self.cpu:
73             embeddings = embeddings.cpu()
74             distance_matrix = pdist(embeddings)
75
76             labels = labels.cpu().data.numpy()
77             all_pairs = np.array(list(combinations(range(len(labels)), 2)))
78             all_pairs = torch.LongTensor(all_pairs)
79             positive_pairs = all_pairs[
80                 (labels[all_pairs[:, 0]] == labels[all_pairs[:, 1]]).nonzero()
81             ]
82             negative_pairs = all_pairs[
83                 (labels[all_pairs[:, 0]] != labels[all_pairs[:, 1]]).nonzero()
84             ]
85
86             negative_distances = distance_matrix[negative_pairs[:, 0], negative_pairs[:, 1]]
87             negative_distances = negative_distances.cpu().data.numpy()
88             top_negatives = np.argpartition(negative_distances, len(positive_pairs))[
89                 : len(positive_pairs)
90             ]
91             top_negative_pairs = negative_pairs[torch.LongTensor(top_negatives)]
92
93         return positive_pairs, top_negative_pairs
94
95
96 class TripletSelector:
97     """
98     Implementation should return indices of anchors, positive and negative samples
99     return np array of shape [N_triplets x 3]
100     """
101
102     def __init__(self):
103         pass
104
105     def get_triplets(self, embeddings, labels):
106         raise NotImplementedError
107
108

```

```

109 class AllTripletSelector(TripletSelector):
110     """
111     Returns all possible triplets
112     May be impractical in most cases
113     """
114
115     def __init__(self):
116         super(AllTripletSelector, self).__init__()
117
118     def get_triplets(self, embeddings, labels):
119         labels = labels.cpu().data.numpy()
120         triplets = []
121         for label in set(labels):
122             label_mask = labels == label
123             label_indices = np.where(label_mask)[0]
124             if len(label_indices) < 2:
125                 continue
126             negative_indices = np.where(np.logical_not(label_mask))[0]
127             anchor_positives = list(
128                 combinations(label_indices, 2)
129             ) # All anchor-positive pairs
130
131             # Add all negatives for all positive pairs
132             temp_triplets = [
133                 [anchor_positive[0], anchor_positive[1], neg_ind]
134                 for anchor_positive in anchor_positives
135                 for neg_ind in negative_indices
136             ]
137             triplets += temp_triplets
138
139         return torch.LongTensor(np.array(triplets))
140
141
142     def hardest_negative(loss_values):
143         hard_negative = np.argmax(loss_values)
144         return hard_negative if loss_values[hard_negative] > 0 else None
145
146
147     def random_hard_negative(loss_values):
148         hard_negatives = np.where(loss_values > 0)[0]
149         return np.random.choice(hard_negatives) if len(hard_negatives) > 0 else None
150
151
152     def semihard_negative(loss_values, margin):
153         semihard_negatives = np.where(
154             np.logical_and(loss_values < margin, loss_values > 0)
155         )[0]
156         return np.random.choice(semihard_negatives) if len(semihard_negatives) > 0 else None
157
158
159 class FunctionNegativeTripletSelector(TripletSelector):
160     """
161     For each positive pair, takes the hardest negative sample (with the greatest triplet loss value) to create a triplet
162     Margin should match the margin used in triplet loss.
163     negative_selection_fn should take array of loss_values for a given anchor-positive pair and all negative samples
164     and return a negative index for that pair
165     """
166
167     def __init__(self, margin, negative_selection_fn, cpu=True):
168         super(FunctionNegativeTripletSelector, self).__init__()
169         self.cpu = cpu
170         self.margin = margin
171         self.negative_selection_fn = negative_selection_fn
172
173     def get_triplets(self, embeddings, labels):
174         if self.cpu:
175             embeddings = embeddings.cpu()
176             distance_matrix = pdist(embeddings)
177             distance_matrix = distance_matrix.cpu()
178
179             labels = labels.cpu().data.numpy()
180             triplets = []
181

```

```

182     for label in set(labels):
183         label_mask = labels == label
184         label_indices = np.where(label_mask)[0]
185         if len(label_indices) < 2:
186             continue
187         negative_indices = np.where(np.logical_not(label_mask))[0]
188         anchor_positives = list(
189             combinations(label_indices, 2)
190         ) # All anchor-positive pairs
191         anchor_positives = np.array(anchor_positives)
192
193         ap_distances = distance_matrix[
194             anchor_positives[:, 0], anchor_positives[:, 1]
195         ]
196         for anchor_positive, ap_distance in zip(anchor_positives, ap_distances):
197             loss_values = (
198                 ap_distance
199                 - distance_matrix[
200                     torch.LongTensor(np.array([anchor_positive[0]])),
201                     torch.LongTensor(negative_indices),
202                 ]
203                 + self.margin
204             )
205             loss_values = loss_values.data.cpu().numpy()
206             hard_negative = self.negative_selection_fn(loss_values)
207             if hard_negative is not None:
208                 hard_negative = negative_indices[hard_negative]
209                 triplets.append(
210                     [anchor_positive[0], anchor_positive[1], hard_negative]
211                 )
212
213             if len(triplets) == 0:
214                 triplets.append(
215                     [anchor_positive[0], anchor_positive[1], negative_indices[0]]
216                 )
217
218             triplets = np.array(triplets)
219
220         return torch.LongTensor(triplets)
221
222
223     def HardestNegativeTripletSelector(margin, cpu=False):
224         return FunctionNegativeTripletSelector(
225             margin=margin, negative_selection_fn=hardest_negative, cpu=cpu
226         )
227
228
229     def RandomNegativeTripletSelector(margin, cpu=False):
230         return FunctionNegativeTripletSelector(
231             margin=margin, negative_selection_fn=random_hard_negative, cpu=cpu
232         )
233
234
235     def SemihardNegativeTripletSelector(margin, cpu=False):
236         return FunctionNegativeTripletSelector(
237             margin=margin,
238             negative_selection_fn=lambda x: semihard_negative(x, margin),
239             cpu=cpu,
240         )

```

Appendix C

Project plan

COMP0029: Project Plan

Anomaly detection in parcel X-ray images using a Siamese neural network

Martynas Janonis
Supervisor: Lewis Griffin

13th November 2019

Contents

1 Aims and objectives	1
1.1 Aims	1
1.2 Objectives	1
2 Deliverables	1
3 Work plan	1

1 Aims and objectives

1.1 Aims

The aim is to develop a working solution that is able to detect anomalies in parcel X-ray images with reasonable accuracy and performance.

1.2 Objectives

- Review the previous approaches to anomaly detection.
- Train a Siamese network to decide if two parcel X-ray images were different views of the same parcel or not.
- Assess the resulting representation as effective or not for anomaly detection.
- Measure how well humans find a pair for an image (i.e. given an image Z of one parcel, does image X or image Y show the same parcel but from a different view).

2 Deliverables

- A fully documented and functional anomaly detection solution.
- A performance and accuracy comparison between the different iterations of the anomaly detection model.
- A strategy for testing and evaluating the solution.

3 Work plan

- Project start to end of October (4 weeks): Literature search and review. Analysing previous anomaly detection solutions.

- Mid-October to mid-November (4 weeks): Refining objectives and starting work on the baseline solution.
- Mid-November to end of November (2 weeks): Completing and evaluating the baseline model.
- End of November to beginning of January (4 weeks): Completing increasingly complex iterations of the model and evaluating their performance.
- Mid-December to mid-February (8 weeks): Have all iterations of the model finished and tested. Fix bugs that were found during testing. Measure human performance on finding image pairs.
- Mid-February to end of March (6 weeks): Work on the Final Report.

Appendix D

Interim report

COMP0029: Interim Report

Anomaly detection in parcel X-ray images using a Siamese neural network

Martynas Janonis
Supervisor: Lewis Griffin

7th February 2020

Contents

1	Progress made and current status	1
1.1	Baseline solution	1
1.2	Model training	1
2	Remaining work to be done	1
2.1	Measuring human performance	1
2.2	Experimenting with triplet networks	1
2.3	Online pair/triplet mining	2

1 Progress made and current status

1.1 Baseline solution

Gathered results from the baseline models (DenseNet-201, ResNext-101). Evaluated other architectures but decided not to continue working on them due to their relatively large top-5 error on the ImageNet dataset.

1.2 Model training

Trained a ResNext-101 (32x8d) model on the parcel dataset. Unfortunately, the model stopped learning after \sim 20 epochs (i.e. the loss stopped decreasing).

The problem was that one of the parameters for the contrastive loss function (the margin) wasn't correct. I am currently training the same model but with the margin set to 2 instead of 1.

2 Remaining work to be done

2.1 Measuring human performance

I am now working on a script that would show three images: an anchor sample, a positive sample and a negative sample. A human would then have to find the correct pair for the anchor sample (i.e. given an image Z of one parcel, does image X or image Y show the same parcel but from a different view). I will then compare human and neural network performance.

2.2 Experimenting with triplet networks

My current approach uses a Siamese network. However, most modern literature describe triplet networks as being superior for similar tasks (e.g. facial recognition). I will train a triplet network and compare its performance to a similar Siamese one.

2.3 Online pair/triplet mining

Currently, during training the network is fed random pairs. After a while, the network can easily differentiate between the easy pairs, but struggles when given a more difficult example.

Online pair/triplet mining would allow the network to learn from only the useful pairs/triplets without having to waste time on pairs/triplets it already knows how to differentiate. This, in theory, should greatly improve the network's performance.