

COMPUTER SCIENCE TRIPOS - PART II PROJECT

Deep Learning Techniques for Credit Card Fraud Detection

May 18, 2018

Proforma

Name: **Harry Graham**
College: **Christ's College**
Project Title: **Deep Learning Techniques for Credit Card Fraud Detection**
Examination: **Computer Science Tripos – Part II, June 2018**
Word Count:
Project Originators: H. Graham & B. Dimanov
Supervisors: B. Dimanov & Dr M. Jamnik

Original Aims of the Project

The primary aim of the project was to implement and compare some deep learning techniques, alongside some baseline models, for credit card fraud detection (CCFD). More specifically, I aimed to experiment with two popular types of architecture, namely Convolutional Neural Networks (CNNs) [1] and Generative Adversarial Networks (GANs). These have been successful in the image classification space and the aim of this project was to shed light on their use in the credit card fraud space. This kind of experimentation of predominantly image-based models, on single dimensional, time-series data is a relatively novel approach for CCFD.

Work Completed

All of the core project aims set out in the proposal have been met, meaning results have been collated and evaluated across the three main components of the project: Baseline Models, CNN methods and GAN methods. I have also gone on to do some extension work relating to further investigation on the models I have experimented with. This is in the form of parameter tuning and further analysis not originally set out in the project proposal.

Special Difficulties

None.

Declaration

I, Harry Graham of Christ's College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

SIGNED

DATE

Contents

1	Introduction	5
2	Preparation	6
2.1	Software Engineering	6
2.1.1	Requirements	6
2.1.2	Tools and Technologies Used	7
2.1.3	Starting Point	7
2.2	Convolutional Neural Networks	8
2.3	Generative Adversarial Networks	9
2.4	Machine Learning Evaluation Practices	11
2.4.1	Train and Test Splits	11
2.4.2	Cross-validation	12
2.4.3	Resampling Methods	12
3	Implementation	16
3.1	Models Overview	16
3.2	Baseline Models	17
3.2.1	Data Preparation	17
3.2.2	Cross Validation Function	17
3.2.3	Cross Validation the Wrong Way	18
3.3	Convolutional Neural Network Models	19
3.4	Generative Adversarial Network Models	21
4	Evaluation	22
4.1	Evaluation Methodology	22
4.1.1	Metrics	22
4.1.2	Visual Inspections	23
4.2	Results Overview	23
5	Conclusions	24
	Bibliography	25

Chapter 1

Introduction

Credit card fraud is a globally significant and increasing problem. According to the Nilson Report [2], annual global fraud losses reached \$22.80 billion in 2016, up 4.4% over 2015. Machine learning has contributed a lot to this problem over the years, helping to automatically learn to classify fraudulent transactions. However, this is still somewhat tedious and clearly, the money lost due to fraud is not decreasing. Not to mention, we still have this difficult business decision of when to draw the cutoff points between classifying fraud but perhaps allowing more benign transactions to be blocked.

A lot of machine learning concepts have been around for decades but ongoing research into deep learning architectures and their applications, makes for an interesting experimentation space. In this project I explore the performance of some particular models, focusing on deep learning, applied to the particular problem of credit card fraud detection (CCFD).

In particular, the aim is to shed light on the use of architectures that have had success in the image classification/generation space, in the context of non-image data i.e transactional vectors and time series data. This is something that has recently seen some success [3] and is novel to credit card fraud data. I first explore a set of baseline classifiers, which are primarily a handful of out-of-the-box supervised learning classifiers such as Random Forest. The point of these is to set the scene for experimenting with the data and to see what can be achieved with what is easily available, in other words without any 'deep' learning components. Here, I also establish techniques and methods for processing and evaluating the data i.e cross-validation, datapoint scaling, and data visualisation.

Then the project shifts to experimenting with Convolutional Neural Networks (CNNs) [1] and Generative Adversarial Networks (GANs) [4]. The main aim of the project is to determine whether these can be applied to time-series data and in particular, the credit card fraud data. I experiment with two different variations of each, which I will outline in this dissertation. In the context of the GAN, I focus on whether this adversarial model is useful in learning the data before trying a semi-supervised, multitask learning approach in order to draw comparisons as a classification model.

Chapter 2

Preparation

2.1 Software Engineering

This section details the project requirements and early design decisions that were made.

2.1.1 Requirements

The main success criteria of the project is outlined as follows:

1. Baseline Models
 - Compare a handful of supervised learning classifiers, from SciKit-Learn.
 - Using metrics described in the Evaluation section of the project.
 - Experiment with resampling techniques.
 - Implement appropriate cross-validation.
2. CNN Models
 - Implement CNN version 1 - Single vector input approach.
 - Implement CNN version 2 - Time series, sliding image window approach.
3. GAN Models
 - Implement GAN version 1 - Dense generator network.
 - Implement GAN version 2 - Using a CNN model from previous work as the generator network.
 - Experiment with how GAN is used and how it performs on the data.

These were all done more or less in order. Some work overlapped, namely work on auxiliary functionality to allow appropriate cross validation or data preparation etc. More details on specifics is outlined in the implementation chapter.

2.1.2 Tools and Technologies Used

Below I describe and justify where necessary the tools and technologies that I used.

Machine Learning

I implemented work predominately making use of Keras¹(with TensorFlow² backend) for CNN and GAN work and SciKit-Learn³ for baseline models and some general data manipulation/metric functions.

The reasons for these choices were a mixture of good documentation, popularity & ease-of-use. Using a TensorFlow backend meant that I could use GPU acceleration if needed. I used Keras as a TensorFlow wrapper, so I could avoid writing models completely from scratch but still giving me the flexibility to develop around models and customise to a large extent. Similarly with SciKit-Learn, which has a lot of helpful utility functions for evaluating models and processing data.

Version Control and Project Tools

I hosted my project in a repository on GitHub⁴, used Git for version control, and used virtual environments with pip for project package management and requirements.

I made heavy use of Jupyter Notebooks for writing code in an experimental manner, with immediate execution and feedback.

Languages

My project was entirely written in Python, using the libraries and APIs described previously. This is mainly due to the large ecosystem and documentation surrounding these machine learning libraries in python and also for the ease of use of tools such as Jupyter Notebooks for experimentation.

2.1.3 Starting Point

My project codebase was written from scratch, with the assistance of the tools and libraries mentioned above. Apart from a basic knowledge of supervised learning covered by the part IB Artificial Intelligence course, I had to learn about most of the models and best practises myself, through thorough reading around the topics.

In terms of technologies, I had little prior experience with SKLearn and Keras/TensorFlow. I reviewed some documentation of these before-hand but I also took an agile approach to

¹<https://keras.io>

²<https://www.tensorflow.org>

³ <http://scikit-learn.org/stable/>

⁴<https://github.com/harrygraham/DeepLearning-CreditCardFraud>

the project, whereby I consulted documentation as and when I needed, during my various milestone sprints. I had significant experience in Python, however, from a summer internship in industry as well as experience with Git.

2.2 Convolutional Neural Networks

Convolutional Neural Networks are a class of deep artificial neural networks, that have seen success in image recognition and other computer vision areas. Unlike typical neural networks, CNNs exploit spatial locality by shared weights.

A CNN typically consists of an input layer, multiple hidden layers and an output layer. Hidden layers are usually convolutional layers, pooling layers and normalisation layers.

Convolutional Layer:

A Convolutional Layer can take many parameters, the most prominent though are:

1. Input shape
2. Kernel size
3. Number of filters
4. Stride size

The **kernel size** defines the size of the filter that will be moved over the input shape, shifting by an amount defined by **stride size**. The **number of filters** simply defines how many separate filters we initialise and convolve with the input, to give multiple outputs. Convoluting a filter with an input is essentially taking the dot product of all overlapping cells between the input and the filter, thus producing a single value in the output shape.

This can be visualised with an example as follows:

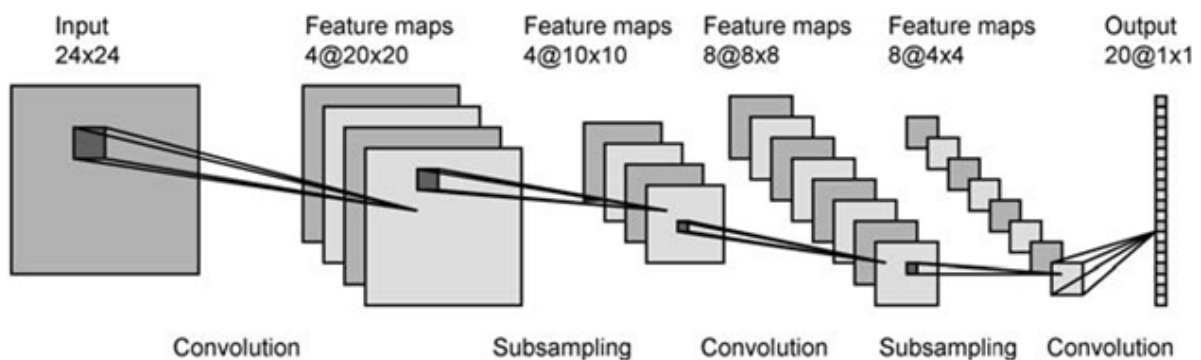


Figure 2.1: Example of a CNN network

We see here that in the first convolutional layer, the input shape is 24x24 and the kernel size is 5x5, which strides over the input shape to produce 4x(20x20) shapes. Here the number of filters was 4, which is why 4 outputs have been produced. Also, typically the

stride size is 1. In this 24x24 input, with a 5x5 filter we can have $(24 - 5 + 1)$ possible positions of the filter in one direction.

Subsampling / Pooling:

Subsampling is essentially taking the average across a group of cells to produce a smaller shape. This therefore produces the same number of 'feature maps' but smaller in size.

Fully Connected Layers:

Fully Connected Layers (FCs) take all the feature maps from a convolutional layer and stack them all together in a traditional neural network in which every node is connected to every other node (fully connected). This then allows the network to learn the relationships between small parts of the image (or input in general) on a fine grained level.

2.3 Generative Adversarial Networks

GAN is a framework proposed by Ian Goodfellow, Yoshua Bengio and others in 2014. The idea is that we have two networks: a generator and a discriminator. The generator network tries to produce, from random noise, data in the form of the training data we want. This could be an image, in the common use-case, or in this one, a fraudulent transaction vector. The discriminator, takes in real-life data (from the training set) and also the generated data from the generator network. The discriminator tries to determine whether the current input is real or generated. Effectively these two networks play a game with each other: the generator tries to fool the discriminator whilst the discriminator tries to catch out the generator. This continues until each network doesn't get any better and the GAN stabilises. Figure 2.1 represents an overview of a GAN network.

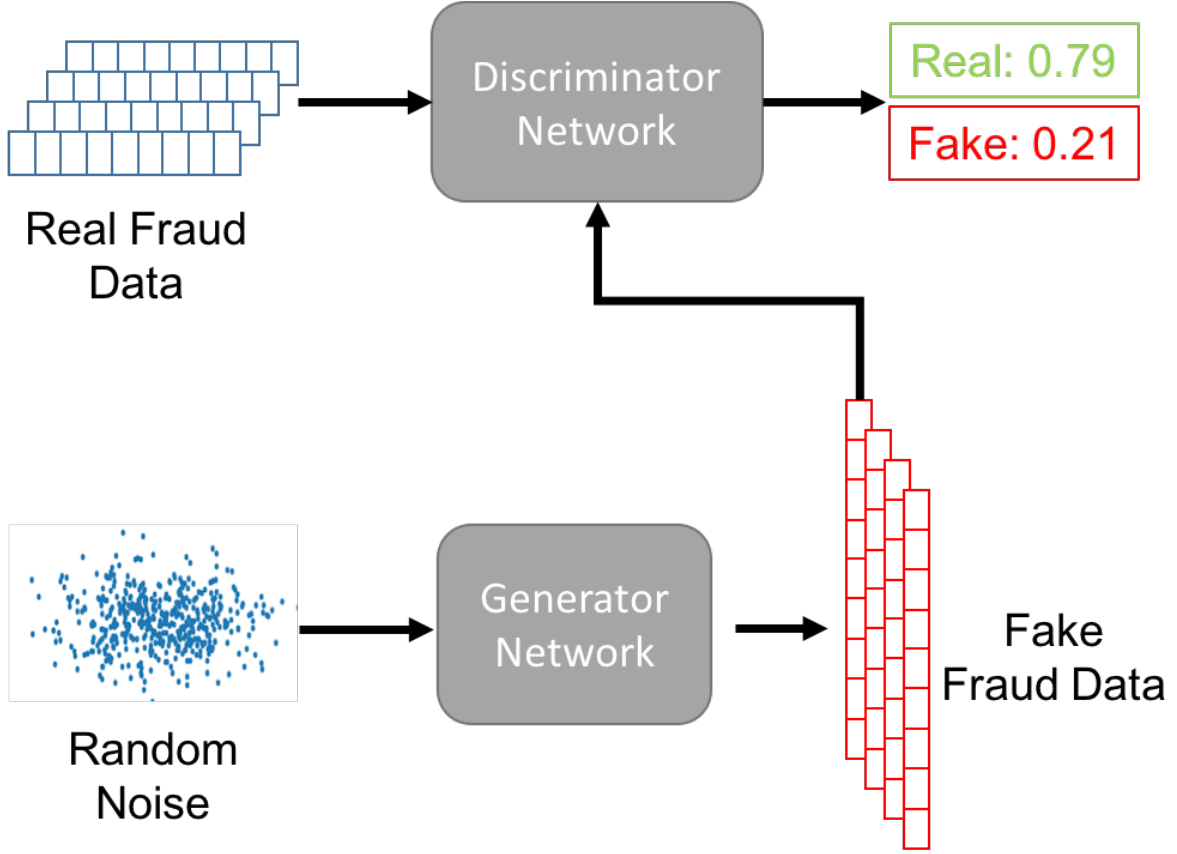


Figure 2.2: Overview of a GAN network

Mathematically, we can define the following quantities:

$X_{x \sim p_{\text{data}}(x)}$ = Sample from distribution of real data

$Z_{z \sim p_z(x)}$ = Sample from distribution of generated data

$G(z)$ = Generator Network

$D(x)$ = Discriminator Network

The process of training for a GAN is like a min-max game between the two networks, and can thus be represented by the following value cost function:

$$\min_G \max_D V(D, G)$$

where

$$V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

The first term in this equation represents the quantity of the real-distributed data passed through the discriminator network. The discriminator tries to maximise this such that $D(x) \rightarrow 1$. The second term represents generated data passed through the discriminator. The generator tries to minimise such that $D(G(z)) \rightarrow 1$ (i.e the discriminator is fooled by the generated sample).

The steps for training a GAN can be outlined as followed:

Step 1: (a) Take a batch of real data and train discriminator to correctly predict them as real

(b) Take a batch of generated data and train discriminator to correctly predict them as fake

Step 2: Freeze the training of the discriminator network

Step 3: Generate a batch of fake data and use the frozen discriminator to train the generator

Step 4: Repeat the above for n epochs until neither network makes any further improvements

In summary, we alternate between training of the discriminator to correctly determine real or fake data and training the generator on fooling the discriminator. The reason for freezing the weights of the discriminator while we train the generator is exactly so that we don't alter the weights during this process and the generator can use the current state of the discriminator to become better.

2.4 Machine Learning Evaluation Practices

Here I describe some of the main machine learning evaluation concepts that I had to have a clear idea of, prior to implementation of this project.

2.4.1 Train and Test Splits

Let's say we have m examples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ each in \mathbb{R}^n . In a supervised learning problem, we also have m labels $\{y_1, y_2, \dots, y_m\}$ in a set \mathbf{Y} .

In machine learning we wish to find a hypothesis $h : \mathbb{R}^n \rightarrow Y$ that is defined by a vector of weights \mathbf{w} . It is then common to write $h_{\mathbf{w}}(x)$.

We define a training set as $\mathbf{s} = [(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_m, \mathbf{y}_m)]$.

When training a model, to try and achieve our $h_{\mathbf{w}}(x)$, we use data from the training set. When we train a model, it will often report some metrics or loss function to reflect how the performance increased during the training process.

The problem with this however is that our model may not perform well on **unseen** data. In other words, it might not **generalise** well. To overcome this, it is good practise to preserve a testing portion of the data, called the test set, that is not used in the training process. The test set is then used to evaluate how the model performs after training.

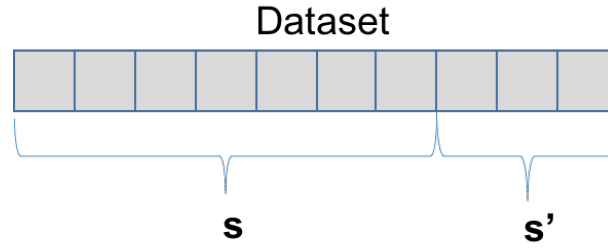


Figure 2.3: A 70:30 train test split of data.

2.4.2 Cross-validation

Cross-validation takes this a step further. The hypothesis of our model may in fact generalise and perform better on some portions of the data over others and thus performing just one split, may not give the most confident results. Also, when performing these train-test splits, the element of randomness in the split can work against us as even if we average this process say 3 times, it is likely that there is some overlap in the runs and we are not really using the data available to its full capacity.

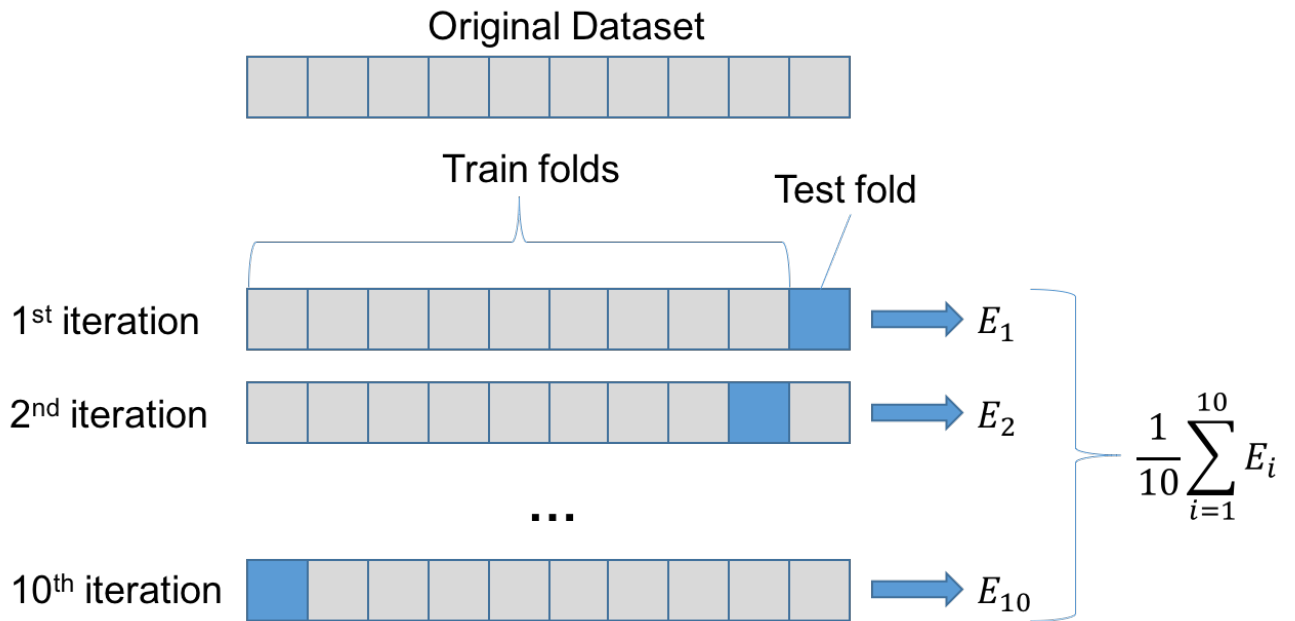


Figure 2.4: K-Fold Cross Validation.

2.4.3 Resampling Methods

During the baseline work of the project I also experiment with resampling methods. These are ways of balancing the dataset such that the ratio of positive class to negative class is brought to 50:50. I do this to give an example of how these methods affect our results before moving onto the deep learning architectures. The three main methods I explored were:

1. Undersampling
2. Duplicate Oversampling
3. Synthetic Minority Oversampling Technique (SMOTE)

Undersampling

Undersampling is the process of reducing the majority class down to a lower amount, to balance more with the minority class. This can be achieved by randomly removing samples until the ratio is 50:50 or another specified amount.

Duplicate Oversampling

Duplicate Oversampling is a naive method of increasing the amount of the minority data class, to match that of the majority. This works by taking existing minority data points and simply duplicating them.

SMOTE

Synthetic Minority Oversampling Technique (SMOTE) is another method of increasing the amount of the minority data but not by duplication as before. SMOTE takes the k nearest neighbours of a data point, randomly selects one of these k neighbours and creates a vector between the two. Then a new, synthetic data point is created some random distance along this vector line, by an amount in the range $[0, 1]$.

Cross-validating Correctly

When employing oversampling techniques, this completely affects the way we handle cross-validation. It no longer is valid to simply oversample the dataset and then perform cross-validation. The reason for this is the potential of overfitting. Figure 2.5 shows what can happen in this situation.

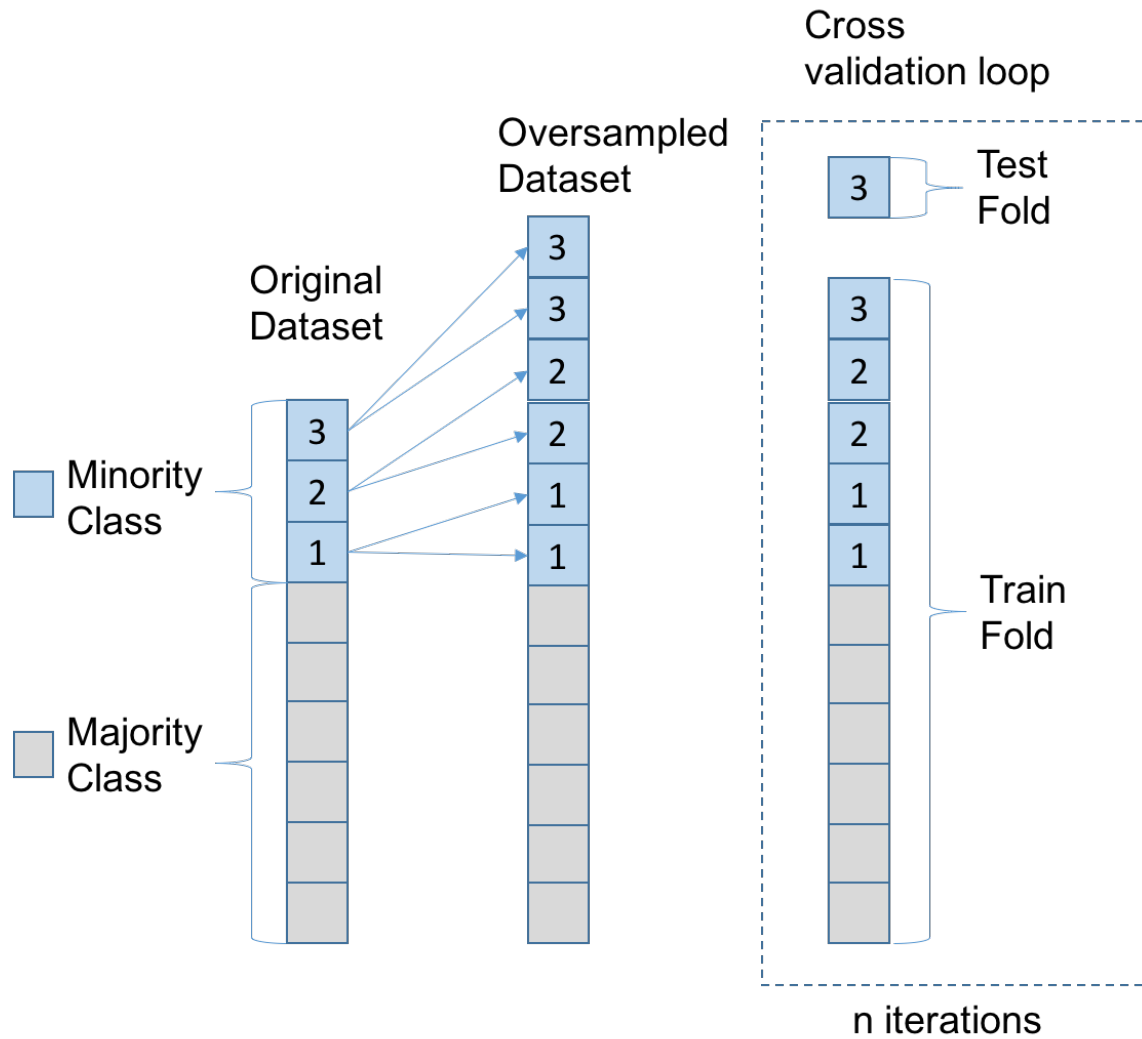
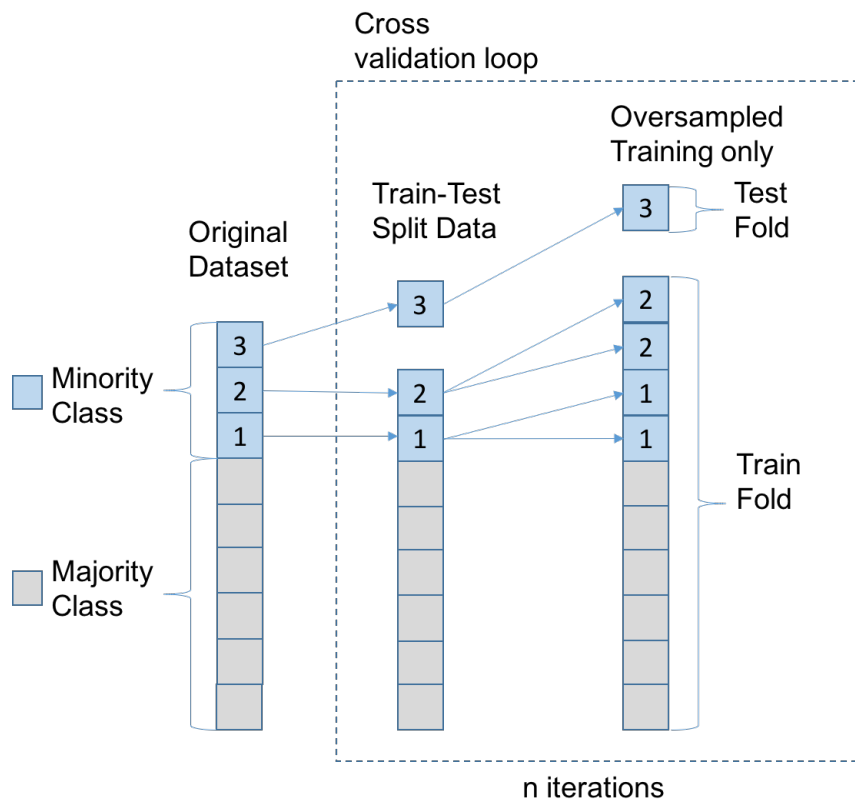
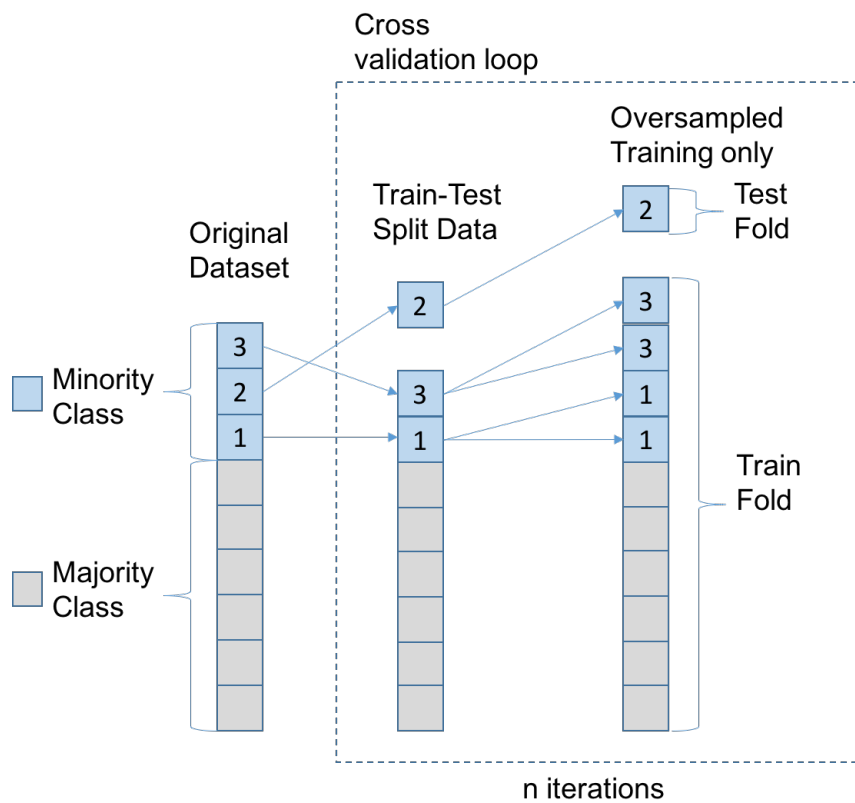


Figure 2.5: Oversampling OUTSIDE the cross-val loop.

Inside the cross-validation loop, the current train-test folds split means that some of the oversampled minority class (node 3) that was portioned into the test set is also in the training set. This means that our model would have seen this data during the training process and will therefore have a bias during the testing and evaluation phase as it already knows how to classify this data.

It is therefore crucial that this is adapted such that any oversampling occurs inside the cross-validation loop, and not before/outside it. Figure 2.6 and Figure 2.7 represent the correct way to handle this, for $n = 1$ and $n = 2$ respectively. For every iteration of the loop, we split the data into train-test first and then oversample the training data fold only. This preserves the testing fold, keeping a portion of the original dataset untouched, for evaluation. This is a lot more effective in testing the generalisability of a model.

Figure 2.6: Oversampling INSIDE the cross-val loop [$n=1$].Figure 2.7: Oversampling INSIDE the cross-val loop [$n=2$].

Chapter 3

Implementation

3.1 Models Overview

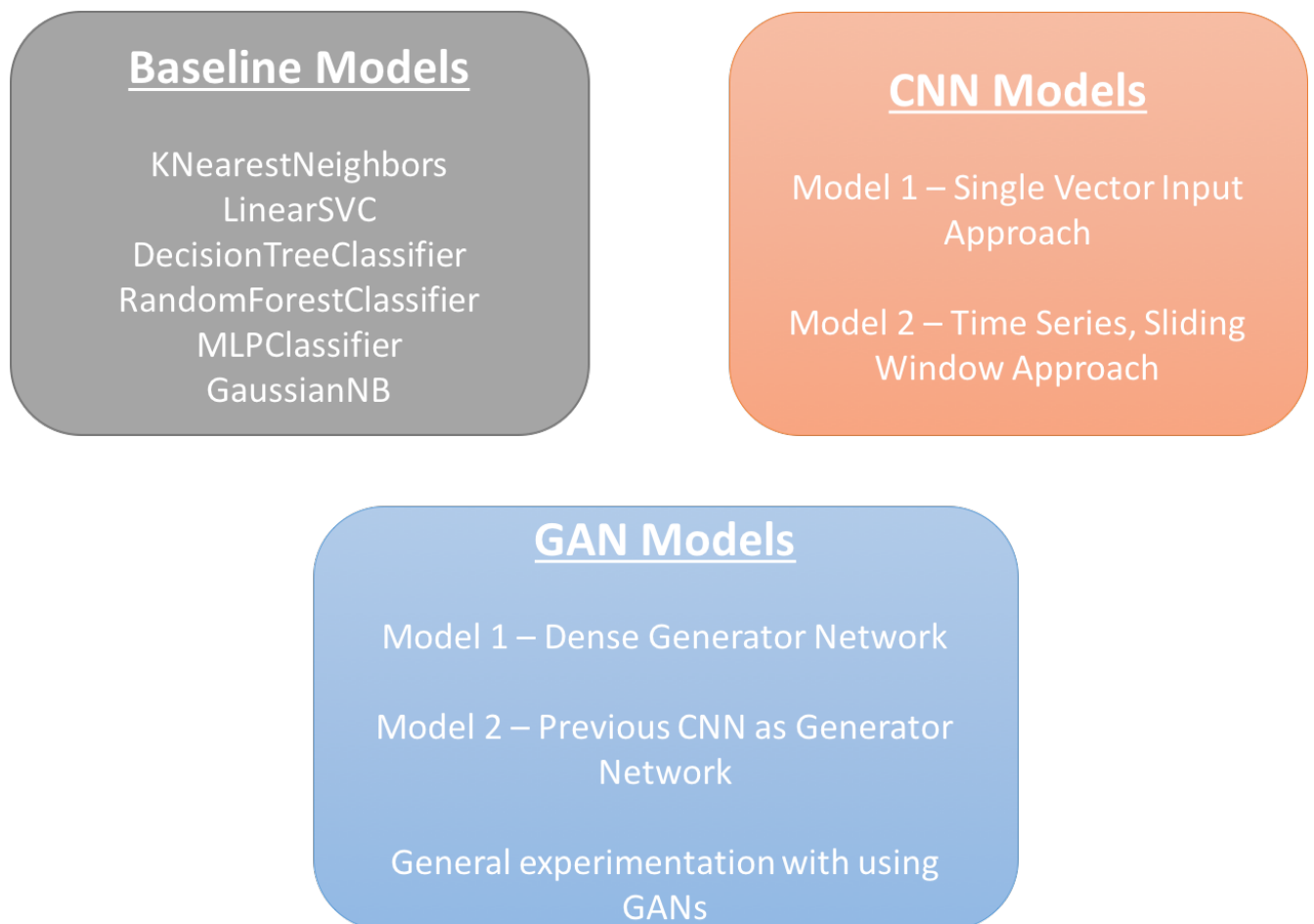


Figure 3.1: Overview of models experimented with.

3.2 Baseline Models

3.2.1 Data Preparation

Feature Scaling

Standardisation involves rescaling the features such that they have the properties of a standard normal distribution with a mean of zero and a standard deviation of one. This is very important in machine learning when features can be a lot larger in magnitude than others. For example, if one feature is age and another is salary in USD then these two quantities will be on different ranges and the salary feature may negatively impact classifier algorithms by deeming that feature to be more important and confusing the weights assignment. Scaling solves this problem.

Therefore one of the first steps was to scale the features, which I did using `StandardScaler` from `SKLearn`.

Feature Engineering

In many machine learning and data science applications, feature engineering is something used to add additional features based on calculations and domain knowledge of the data. In this case, though, as the data is already post-PCA and the time-series span of the transactions is only over 2 days. I decided not to do any feature engineering and to just employ algorithms on the original dataset features.

Sanity Checks

Also, I checked the data to see if any values were missing, or null. This was to ensure that the data was clean to be passed through machine learning networks. There were no missing data in this case.

3.2.2 Cross Validation Function

Defining my own cross-validation function early on was a crucial step in this project. As explained briefly in the preparation chapter, the way in which cross-validation is performed is very important to ensure no overfitting occurs and that the results are as confident as they can be. As far as implementation is concerned, It was not enough to use the built in `SKLearn` cross-val function as it did not provide enough granularity in it's output. I wanted very specific metrics and not to mention, I needed to resample data in the correct place which meant I needed to define my own function, based on the standard library which uses `Stratified KFold`. This essentially means that the splitting of the data in each fold of the `KFold` procedure, contains at least some portion of the underrepresented class. This is important as it would be very common that some folds don't contain any of the minority class.

The function uses a utility function from `SKLearn`, `precision_recall_f1`, which simply reports the precision, recall and f1-score for predicted values. You give it a set of predicted

values and the set of true values and it calculates the metrics. I use this inside my function to report the metrics we are interested in, something I could not do originally.

Algorithm 1 Cross-Validation function

```

1: procedure CUSTOM-CROSS-VAL( $X, Y, CLF, N$ )
2:    $skfolds \leftarrow \mathbf{StratifiedKFold}(n\_splits = N)$ 
3:
4:    $precisions \leftarrow []$ 
5:    $recalls \leftarrow []$ 
6:    $f1scores \leftarrow []$ 
7:    $elapsedtimes \leftarrow []$ 
8:
9:   for  $train\_indices, test\_indices$  in  $skfolds.split(X, Y)$  do
10:
11:      $X\_train\_folds \leftarrow X[train\_indices]$ 
12:      $Y\_train\_folds \leftarrow Y[train\_indices]$ 
13:      $X\_test\_folds \leftarrow X[test\_indices]$ 
14:      $Y\_test\_folds \leftarrow Y[test\_indices]$ 
15:
16:      $X\_res, Y\_res \leftarrow \mathbf{Resample}(X\_train\_folds, Y\_train\_folds)$ 
17:
18:      $start \leftarrow currenttime$ 
19:      $CLF.fit(X\_res, Y\_res)$ 
20:      $end \leftarrow currenttime$ 
21:      $elapsed \leftarrow end - start$ 
22:      $elapsedtimes.append(elapsed)$ 
23:      $y\_pred \leftarrow CLF.predict(X\_test\_folds)$ 
24:      $stats \leftarrow \mathbf{precision\_recall\_f1}(y\_pred, Y\_test\_folds)$ 
25:
26:      $precisions.append(stats[0])$ 
27:      $recalls.append(stats[1])$ 
28:      $f1scores.append(stats[2])$ 
29:   end for
30:   return [  $\text{mean}(precisions), \text{mean}(recalls), \text{mean}(f1scores), \text{mean}(elapsedtimes)$  ]
31: end procedure

```

This is something I could then use for varying **Resample** methods.

3.2.3 Cross Validation the Wrong Way

Test_Train_Split vs Custom cross_val_score using KFold

This concerns the underlying approach for the cross validation loop. We have already seen we can use sklearn's `test_train_split` function to split the data into train and test portions. We can then oversample our training data, fit the model and make predictions using the test set. Approach 1 would be to simply do this multiple times and average the results. This approach, however, will give better results than expected (which will

be shown in a later section, I go on to implement this version for comparison). The reason for this is due to the nature of `test_train_split` and its randomness. Problems with this method: `test_train_split` allows you to randomly split your data, by giving a parameter that specifies the ratio. However as the split is random, it is likely that there will be overlap in the CV iterations as to which data points are put in the test set. In other words, values selected during one iteration, could be selected again during another iteration. The consequences of this means that the model may not be exposed to particular portions of the data whereby it does not generalise well and we are not capturing that in our results. Also, It is not making maximal use of the data we have.

3.3 Convolutional Neural Network Models

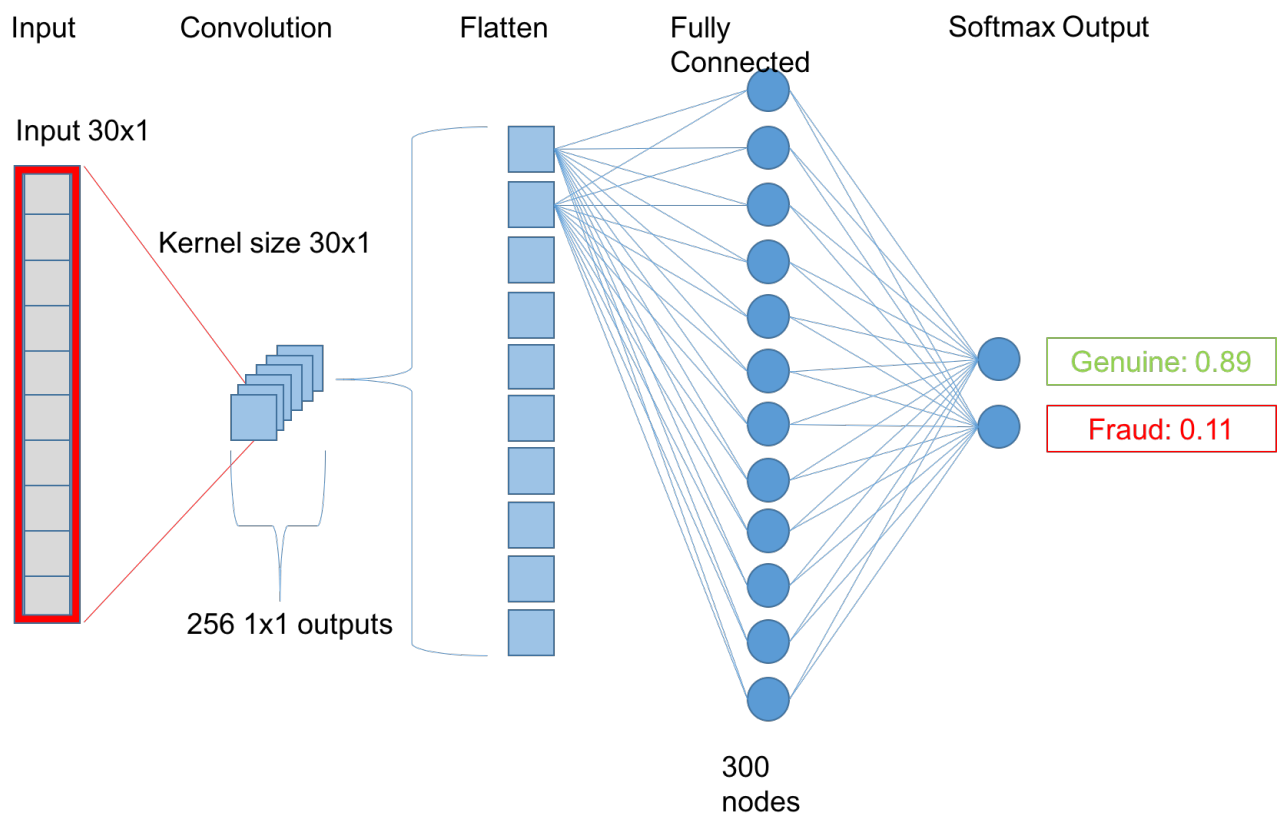


Figure 3.2: Overview of CNN Model 1.

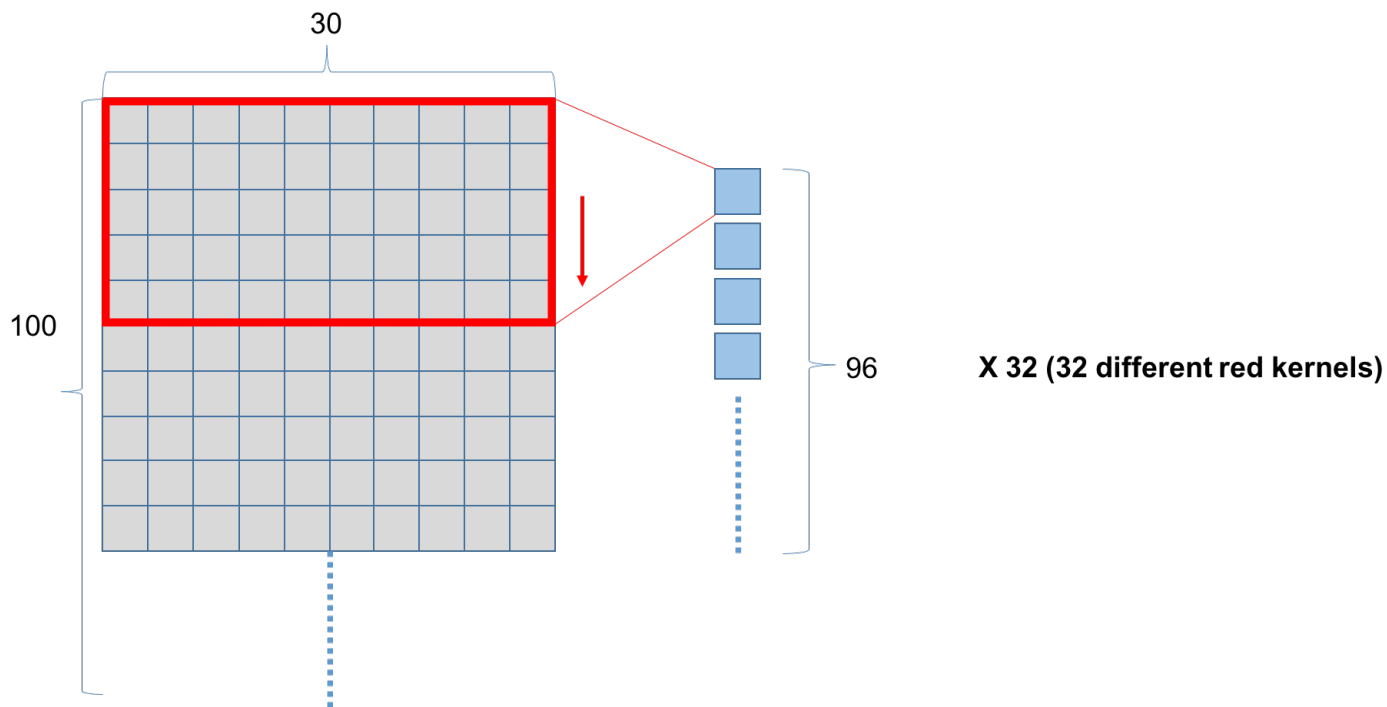


Figure 3.3: Overview of CNN Model 2 - Part 1.

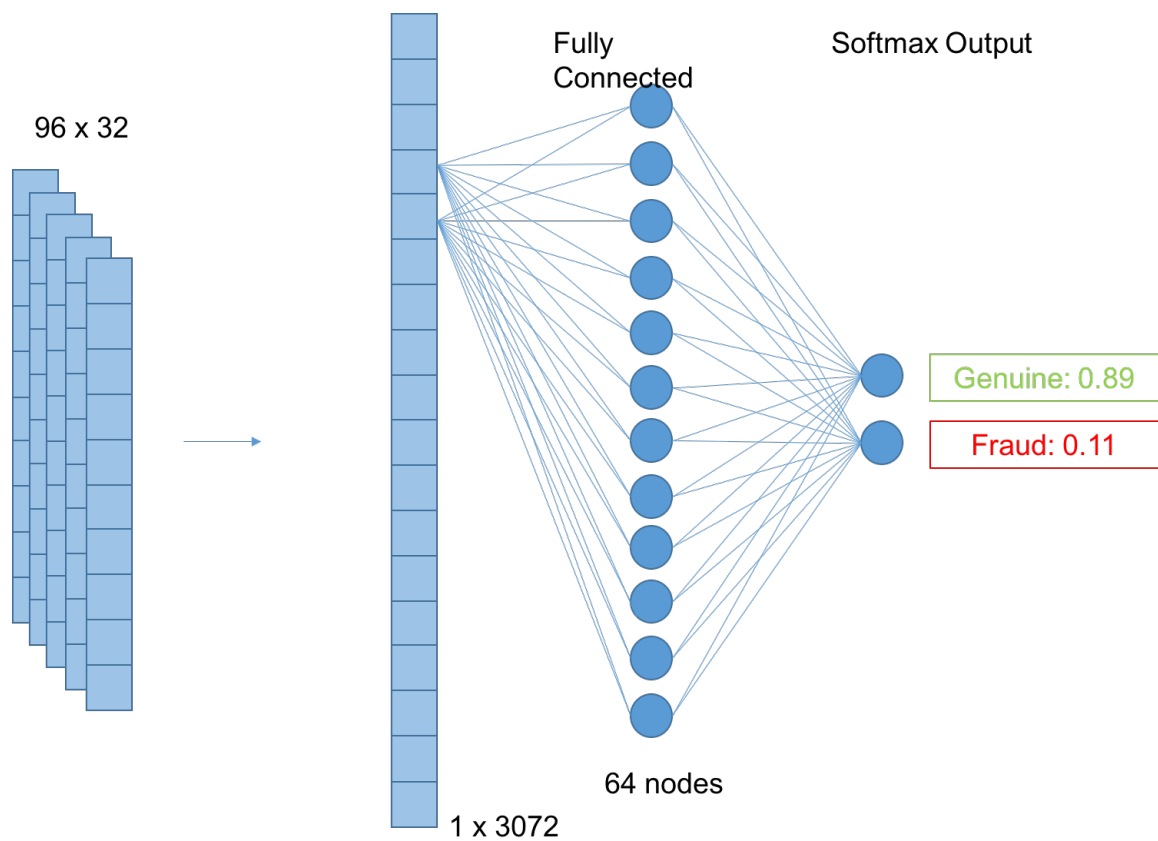


Figure 3.4: Overview of CNN Model 2 - Part 2.

3.4 Generative Adversarial Network Models

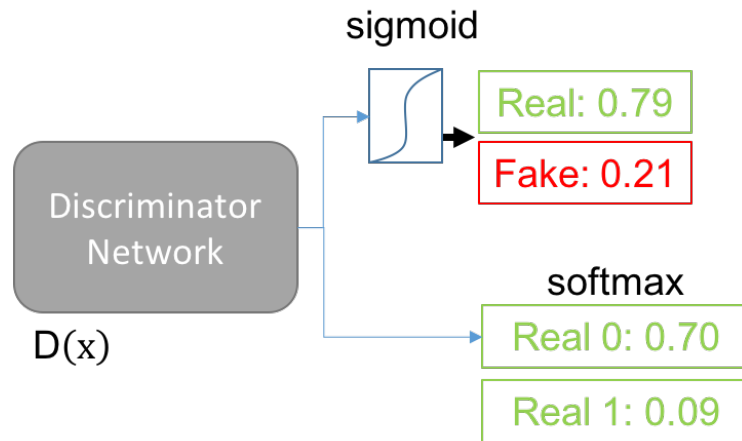


Figure 3.5: Overview of GAN-SSL discriminator output.

Chapter 4

Evaluation

4.1 Evaluation Methodology

4.1.1 Metrics

In this project, the metrics that were of interest were not just simply the accuracy, as with many machine learning applications. Accuracy simply measures how many correct classifications are made. This, in the context of credit card fraud, is useless. This is because, even if we classified all fraud as benign, the classifier would still be over 99% accurate, due to the high imbalance of the fraud/non-fraud classes. In the dataset being used, fraud accounts for only 0.17% of the total.

Precision and Recall

Therefore we need to look at metrics that tell us more about what we're interested in. Two main measures of interest are *Precision* and *Recall*.

$$Precision = \frac{T_p}{T_p + F_p} \quad , \quad Recall = \frac{T_p}{T_p + F_n} \quad T_p = TruePositive \quad , \quad F_n = FalseNegative$$

Precision is intuitively the ability of the classifier to not misclassify. Recall is intuitively how well to catch fraudulent examples.

Of course the threshold by which we allow benign transactions to be misclassified as fraud, is a business decision. In the context of a bank, of course the number of fraudulent transactions caught is very important but at the same time we care about precision as this amounts to freezing customers accounts, sending text messages even though they are not subject to fraud, which ultimately costs the bank time and money and potentially customers.

4.1.2 Visual Inspections

Precision-Recall Curves and ROC

When classifying, the true positive rate (TPR) and false positive rate (FPR) changes for increasing values of precision and recall. This can be visualised by looking at TPR-FPR curves. This gives an indication of how the trade-off varies for differing values.

Taking the area under this curve is known as AUC, therefore gives an overall statistic for comparing classifiers.

GAN loss function graphs and generated data visualising

With GANs, the way in which we can visualise performance is different to others due to it's specific adversarial nature. Two prominent use cases for GANs is to view the loss function graph for the generator and discriminator networks to see how the two networks fight with each other. In an ideal, theoretical world, we would expect to see the two loss curves converging to some stable value, which indicates the two networks can no longer improve or better themselves.

Also, seeing as a large part of the GAN architecture is generating data, we can inspect this to see if it looks appropriate. A common example with the MNIST dataset is to see if the generated images look like handwritten digits. In this case, we can inspect the shape and distribution of the generated fraud transactions and view them side-by-side with real fraud transactions to see if any similarities have been picked up by the network.

4.2 Results Overview

Chapter 5

Conclusions

1. Can you apply DL architectures to this problem? Yes, you can apply them, good/bad?
2. What can you compare it to, is it useful, how does it compare
3. What could the recommendation be, would you recommend to use DNN

Bibliography

- [1] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014.
- [2] T. N. R. I. . . O. 2017, “Card fraud losses reach \$22.80 billion.” <https://www.nilsonreport.com/>, 2017.
- [3] Z. Wang, W. Yan, and T. Oates, “Time series classification from scratch with deep neural networks: A strong baseline,” in *Neural Networks (IJCNN), 2017 International Joint Conference on*, pp. 1578–1585, IEEE, 2017.
- [4] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative Adversarial Networks,” *ArXiv e-prints*, June 2014.