



**FACULTY
OF INFORMATION
TECHNOLOGY
CTU IN PRAGUE**

ASSIGNMENT OF MASTER'S THESIS

Title: Comic2vec: Vector representation of comics
Student: Bc. Martin Piták
Supervisor: MSc. Juan Pablo Maldonado Lopez, Ph.D.
Study Programme: Informatics
Study Branch: Knowledge Engineering
Department: Department of Applied Mathematics
Validity: Until the end of summer semester 2019/20

Instructions

In this work we explore the feasibility of a vector-space embedding for unstructured data (comics). This can provide the basis for a content-based recommendation system. As online evaluation of recommender systems is costly, we propose an alternative criteria to measure the performance of such system in this setting.

- 1) Study different approaches for abstract representation of rich media.
- 2) Provide a comprehensive survey of the existing work.
- 3) Explore different alternative embeddings.
- 4) Propose a baseline system based on previous work. Consider possible optimizations and improvements.
- 5) Use a proxy to measure the performance of the achieved results.

References

Will be provided by the supervisor.

Ing. Karel Klouda, Ph.D.
Head of Department

doc. RNDr. Ing. Marcel Jiřina, Ph.D.
Dean

Prague October 12, 2018



**Czech
Technical
University
in Prague**

**Faculty of Information Technology
Department of Applied Mathematics**

Master Thesis

**Comic2vec: Vector representation of
comics**

Bc. Martin Piták

Supervisor: MSc. Juan Pablo Maldonado Lopez, Ph.D.

Reviewer: Ing. Karel Klouda, Ph.D.

May 2019



Acknowledgements

I would like show my gratitude to my supervisor MSc. Juan Pablo Maldonado Lopez Ph.D. for his excellent guidance. I would also like to express gratitude to Ing. Karel Klouda Ph.D., Ing Daniel Vařata Ph.D., Ing. Tomáš Bartoň, Ondřej Bíža, Evka Šimková, Štěpánka Jislová and Václav Roháč for their time in which they were kind enough to discuss and find solutions for any issues concerning this paper. I would also like to thank my family and my friends for supporting me along the journey. I would like to thank Lenka Kubičová for helping me create the dataset of comics. And on a final note I must sincerely thank Andrea Malá for the correction of this paper.



Declaration

I hereby declare that the presented thesis is my own work and that I have cited all sources of information in accordance with the Guideline for adhering to ethical principles when elaborating an academic final thesis.

I acknowledge that my thesis is subject to the rights and obligations stipulated by the Act No. 121/2000 Coll., the Copyright Act, as amended. In accordance with Article 46(6) of the Act, I hereby grant a nonexclusive authorization (license) to utilize this thesis, including any and all computer programs incorporated therein or attached thereto and all corresponding documentation (hereinafter collectively referred to as the “Work”), to any and all persons that wish to utilize the Work. Such persons are entitled to use the Work in any way (including for-profit purposes) that does not detract from its value. This authorization is not limited in terms of time, location and quantity. However, all persons that makes use of the above license shall be obliged to grant a license at least in the same scope as defined above with respect to each and every work that is created (wholly or in part) based on the Work, by modifying the Work, by combining the Work with another work, by including the Work in a collection of works or by adapting the Work (including translation), and at the same time make available the source code of such work at least in a way and scope that are comparable to the way and scope in which the source code of the Work is made available.

In Prague May 5, 2019



Abstract

The world of comics is receiving a lot of attention lately. Not only are thematic movies based on comics being released almost daily but on the top of that scientists are starting to study comics as a research field now as well.

This paper tries to describe comics so the reader can understand them. It explores the possibilities of embedding comics into vector space. It explains various methods and algorithms that will be used in the process of creating and evaluating the accuracy of the embeddings. There are two embeddings created: one for the style and the other one for the text. A special metric is used to measure the accuracy of these embeddings. The style embedding is created using Inception V3 which is a Convolutional neural network (CNN) re-trained on TPU and achieves accuracy of 98%. The text embedding is created using a method named Doc2Vec and achieves accuracy of over 83%. Two datasets are created in the process of making this work, unfortunately the one used for style embedding cannot be made public.

Keywords: Comic, TensorFlow, CNN, Recommendation system, Embedding, Clustering



Abstrakt

V dnešní době se komiksům dostává velké popularity. Jedním z faktorů toho nárůstu je fakt, že vycházejí filmy založené na komiksových světech, kterým se dostává velké popularity. Dalším faktorem je zájem vědců o komiksy a jejich následné studování.

Tato práce se snaží popsat komiksy tak, aby i komiksový začátečník pochopil, co to přesně komiksy jsou. Zkoumá různé možnosti vnoření komiksů do vektorového prostoru. Vysvětluje různé metody a algoritmy, které jsou použity při tvorbě těchto vnoření a jejich následného vyhodnocení přesnosti. Jsou vytvořena dvě vnoření, jedno pro styl a druhé pro text. Pro vyhodnocení přesnosti je použita speciální metrika. Vnoření stylu je vytvořeno pomocí Inception V3, což je konvoluční neuronová síť (CNN), která byla přetrénovaná pomocí TPU. Toto vnoření dosahuje přesnosti 98%. Vnoření textu je vytvořené pomocí Doc2Vec a dosahuje přesnosti více jak 70%. Při tvorbě této práce byly vytvořeny dva datasety, jeden obsahující panely komiksů a druhý texty. Bohužel dataset panelů komiksů nemůže být zveřejněn.

Klíčová slova: Komiks, TensorFlow, CNN, Doporučovací systémy, Vnoření, Shlukování



Contents

- 1 At first a bit of philosophy and history** **3**
- 1.1 Comic 3
- 1.2 History of comics 4
 - 1.2.1 Europe 5
 - 1.2.2 America 5
 - 1.2.3 Japan 6
 - 1.2.4 Modern 7
- 1.3 Basic categories of comics 7
 - 1.3.1 Region 8
 - 1.3.2 Format 8
 - 1.3.3 Layout 8
 - 1.3.4 Age 9
 - 1.3.5 Gender 9
- 1.4 Language of comics 9
 - 1.4.1 Panels and gutters 9
 - 1.4.2 Text 10
 - 1.4.3 Lettering/sounds 10
- 1.5 Style 10
 - 1.5.1 Line Style 11
 - 1.5.2 Movement Style 11
 - 1.5.3 Colouring and Shading 12
- 2 Related Work** **15**
- 2.1 Style embedding 15
 - 2.1.1 FaceNet 15
 - 2.1.2 Illustration2Vec 15
 - 2.1.3 Are Anime Cartoons? 15
 - 2.1.4 The Amazing Mysteries of the Gutter: Drawing Inferences Between Panels in
Comic Book Narratives 16
 - 2.1.5 A Neural Algorithm of Artistic Style 16
 - 2.1.6 Recognizing Image Style 16
 - 2.1.7 Item2Vec: Neural Item Embedding for Collaborative Filtering 16

2.1.8 A Visual Embedding for the Unsupervised Extraction of Abstract Semantics	17
2.1.9 CNN-based Classification of Illustrator Style in Graphic Novels: Which Features Contribute Most?	17
2.2 Text Embedding	17
2.2.1 Efficient Estimation of Word Representations in Vector Space	17
2.2.2 Distributed Representations of Sentences and Documents	17
3 Style embedding	19
3.1 Panel vs. Page	19
3.2 Dataset	19
3.2.1 Data augmentation	20
3.2.2 Data split	21
3.3 Distance Metric	21
3.4 Clustering algorithms	21
3.4.1 K-means	22
3.4.2 HDBScan	22
3.4.3 Markov clustering	24
3.5 Dimensionality reduction	24
3.5.1 PCA	24
3.5.2 MDS	25
3.5.3 Isomap	26
3.5.4 t-SNE	27
3.5.5 UMAP	28
3.5.6 Conclusion	30
3.6 Baseline	31
3.6.1 Colour histograms	31
3.6.2 Histogram of oriented gradients	32
3.7 CNN	34
3.7.1 Inception	35
3.7.2 RMSPROP	37
3.7.3 Loss	37
3.7.4 Metrics	38
4 Text embedding	41
4.1 Dataset	41
4.2 Doc2vec	42
4.2.1 PV-DM	42
4.2.2 PV-DBOW	43
5 Results	45
5.1 CNN	45
5.2 Doc2Vec	48

5.3 Combined	49
6 Future improvements	51
7 Conclusion	53
A Contents of CD	61



Figures

1.1 Examples of the Clear line and the Atomic style.	12
1.2 Examples of colouring and shading techniques.	12
3.1 Subset of MNIST dataset reduced by all previously described algorithms.	30
3.2 Accuracy of histogram baseline.	31
3.3 Accuracy of histogram of oriented gradients baseline.	32
5.2 Results of clustering reducec train data centroids.	47
5.1 2D representation of sample from test data.	48
5.3 2D representation of the textual evaluation data embeddings.	49
5.4 2D representation of the combined train data embeddings.	50



Tables

5.1 Results from Inception V3 experiment	46
5.2 Results from Inception V4 experiments	47
5.3 Results from Inception-ResNet V2 experiments	47
5.4 Results from Doc2Vec experiments	49



Introduction

The world of comics is receiving a lot of attention lately. Not only are thematic movies based on comics being released almost daily but on the top of that scientists are starting to study comics as a research field now as well. I am personally a huge fan of comics. Nowadays comics are progressively consumed online, these online sites often provide a recommendation for the user of what comics to read next. Typically, recommendations are done by collaborative filtering which gives recommendations based on similarities with other users. This recommendation system is not bad but it lacks the explanatory technique that content based recommendation systems have. If someone were to ask for recommendations normally, people would recommend comics based on their content similarity and not based on what others like. That is why I decided to create such embeddings which can be used in content based recommendation systems to give more personalized results.

The world of comics is relatively new and a scientific field called comic studies about them is also young. I am going to explain what comics are and what they are made of so we have an equal understanding of comics.

In this paper, I am going to explore possible solutions for comic embeddings. I am going to explore ways of embedding comics style and text. Many possible models can be used. For the style embedding I will be focusing on three models, them being colour histograms, histograms of oriented gradients and convolutional neural networks. As for the text embeddings I will describe and use state of the art method Doc2Vec.

I am also going to propose a metric to evaluate accuracy of these embeddings. This metric is only for the purpose of measuring the accuracy of these models. In practise a similarity between the features would be used.

I am also going to propose and create two datasets. Since the proposed ones would require deeper knowledge of comics and access to all kinds of comics, I will create different ones. One dataset will be used for style embedding and the other one for text embedding. These datasets cannot be public.

At the end, I am going to evaluate the results and assess the overall success of this paper. I will also talk about some future improvements that will be necessary to fully explore this topic.

Chapter 1

At first a bit of philosophy and history

1.1 Comic

First, we need to know what comic actually is. Since a comic is part of literature, there are many definitions and interpretations of that word. Here I am listing some major foreign definitions of the term comic.

“There must be a sequence of separated images. There must be a preponderance of image over text. The medium in which the strip appears and for which it is originally intended must be reproductive, that is in printed form, a mass medium. The sequence must tell a story which is both moral and topical”[1]

“Comics are juxtaposed pictorial and other images in deliberate sequence, intended to convey information and/or produce an aesthetic response in the viewer.”[2]

“x is a comic if x is a sequence of discrete, juxtaposed pictures that comprise a narrative, either in their own right or when combined with text.” [3]

I have not only relied on written sources I also did a couple of interviews with people from the comics industry and as I expected I received different answers.

“A comic is a medium. A comic may contain text. A single image can be considered as Comic, but it must contain some story or an idea, which is highly subjective. Some people might even consider hieroglyphs to be comics.” [4]

“A comic is a medium similar to literature or film. There must be image sequentiality, continuity between pictures. It should hold up without the text. Minimum of two images.” [5]

As you can clearly see the definitions vary in nature. Some of them ex-

regional ones are the most prominent and popular.

1.2.1 Europe

The earliest work that is known and could be classified as a comic are the wall paintings in the mortuary chapel of Menna in Egypt which is dated to around 1300 b.c.² Another example is the Bayeux Tapestry in France which is dated to around 1100 a.c. The Torture of St. Erasmus is the oldest recent³ example from around the year 1460. There are not many records about comics before the 1400s. Most of them were damaged or lost because of storing methods or wars. [2]

European comics between the years of 1400s and 1830s were mostly painted or drawn in a realistic manner and not as a cartoon style⁴ as we know them today. Their story was not intended as light entertainment as it is today but they mostly told important and serious stories. Since printing technology was not as advanced compared to today, mass production of books or even comics on the scale as we know today was impossible. As a consequence of this, the illustrations/images they contained were limited. Mass produced multi-coloured images were not available until the invention of Chromolithography in 1837. [1]

David Kunzle in his book “The early comic strip 1450 - 1825” shows about six hundred comics or “pictorial stories” as he calls them. But these are just mere examples from this time period, there were countless more comics made in this era. Many of which we might never be able to find.

Two most prominent European countries in the comics industry are France and Belgium. Some of the most popular and most influential comics are Adventures of Tintin, Asterix, The Smurfs, Lucky Luke.

1.2.2 America

American comics are currently the world’s most popular ones. Their popularity spiked up with the first issue of Action Comics in 1938, where Superman was introduced. This day also marks the beginning of the era called the Golden Age of comics. At this time these comics were still similar to newspaper comics with their layout, style and story. The Golden Age ended in 1950 after most of the soldiers returned from World War II.

The next era was called the Silver Age. Silver Age was much more experimental with its narrative and art style. They employed a more surrealistic art style. Comics also explored much deeper and not so simplistic plots, they

²Old cave paintings could also be classified as comics.

³That could be characterized as a comic without knowing the definition.

⁴Or a caricature style.

1.2.4 Modern

American cartoonist Coultou Waugh marks the beginning of the modern comic with the creation of *The Yellow Kid* by Richard Outcalt in 1896. On the other hand, European comic researchers mark the beginning with Rudolphe Töpffer in the early 19th century. But as we saw before there were many more comics even before that. [7]

The modern comic is heavily influenced by Belgium's early 1900s comic books. Mainly *Tintin* and *Spirou*⁶. They defined two most commonly used ways to show movement in comic. *Tintin* used "ligne claire", the Clear line and *Spirou* used "école de Charleroi", Atomic style. [9]

Early American comics were in shades of grey because coloured printing was not available and if it was, it was rather expensive. Later, when printing technology progressed, comics could be printed in colour but still, the colour pallet was quite limited as it was expensive to have a wide range of colours⁷. European comics, on the other hand, were printed with much better colour, because of a better printing technology. Japanese comics were and still are in the shades of grey, with the occasional full-colour pages. [2, 9]

Today, there are many more comics than there were ever before mainly thanks to of the increasing popularity of the Internet. The Internet makes it possible for anyone to create and publish comics online. And some previously printed comics migrate to being online only. As a result of this a new term emerged, Webcomic. This term describes comics published on the Web which are most often not in the usual grid or a horizontal strip form but in a vertical form. This way they are much more suitable for smartphones. Special platforms for these webcomics emerged like Korean Line Webtoon or Czech Nanits, which is trying to change the way electronic comics are shared and enjoyed.

1.3 Basic categories of comics

There are many different ways we can categorize comics and I am going to list a few of the basics and describe why they are not suitable methods of

⁶These are not the actual comics rather than comic books where one could find plethora of comics, but what they shared in common was the style in which movement was captured and what type of an audience they were designed for.

⁷They were so limited that they had a list of colors they may use and how big of an amount of one colour per page they may use. The colours they could use were the basic CMYK(Cyan, Magenta, Yellow, and Key. Key is black) colour in various opacities/shades with some basic mixing.

is one page with multiple lines of panels. This categorization also suffers from not distinguishing styles and storytelling. This categorization is helpful when viewing comics on a specific device, for example, it is much easier to view grid comics on paper than it is on a mobile phone or on a computer.[12]

■ 1.3.4 Age

Comics have been, from the very beginning, categorized for what demographic they are intended. There are commonly three categories but these can be split even further to focus on a specific niche demographic. The three categories are young, teenage and adult.

These categories can shed some light on the story of the comics but not so much on their art style. Commonly comics aimed for kids are drawn more pleasantly and simplistically so the children can understand them more easily. The later categories are generally drawn much more realistically.[12]

■ 1.3.5 Gender

There are hundreds of genres and subgenres, and they are one of the most commonly used descriptions of comics, movies, music and books. They provide a good description of the story and its elements but still do not provide a clear description of the art style.[12]

■ 1.4 Language of comics

From the definition of comics it is clear that in every comic there must contain at least one picture, panel. Besides that, the other parts of comics are the following: panel size, panel layout, lettering, text and empty space⁸.

■ 1.4.1 Panels and gutters

As stated before, panels are the centrepiece of every comic and they are separated by gutters, which makes them the next most important piece of language structure. The size of panel influences the flow of time in this panel. For example a wide panel would naturally take longer to “read”. This inherently makes the panel seem to be longer than just a moment.

A panel layout is also important; some can be neatly organised into a grid; others might not be organised into any specific layout. The grid layout increases the speed of reading and the level of comprehension. It is generally easier to read grid layouts than others. This is why grid layouts are mostly

⁸Are sometimes also called gutters.

vision.”[14] This description is not very helpful in our case. To simplify the issue, let us assume that style stands for the kind of lines the author used, how he draws movements, what kind of shading or colouring techniques he uses, what kind of characters he uses and how abstract or concrete the drawings are. There are other style elements that can be described. I will only describe types of lines, movement and coloring styles in more detail as they are the most prominent.

Most styles are grouped into what is called a movement, where very similar styles are grouped under one name, for example impressionism. But every artists has been, or is working towards his personal art style. But one artists can draw in multiple styles.

Overall, there is an infinite number of styles that comics author may use, but most often they use caricature style⁹ There are some comics that use a realistic style or abstract style but they are mostly considered to be an exception. Caricature style is mostly used to appeal to user’s imagination. This way readers can identify with the characters more easily. From the interviews I did, I learned that colours are also part of the style. Same applies for the character’s design. There are different kinds of styles of how the characters are portrayed. Another thing that influences the style is the difference in how the character and the background are drawn.[14]

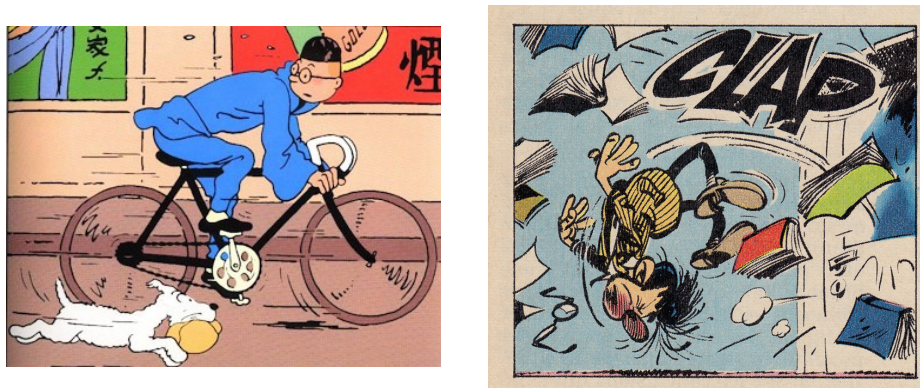
■ 1.5.1 Line Style

When we think about lines, we shortly conclude that there are thin lines, thick lines and lines in between. There are also several ways to change direction with lines; for example, we can change direction with a sharp turn or with a curve. Or we can even change the thickness of the line to emphasise the change of direction. There is also a possibility to combine all these techniques. This gives authors many possible options to draw comics.[15]

■ 1.5.2 Movement Style

There are two main possibilities of how to draw movement in comics. One is called the Atomic style, and the other one is named the Clear line. The Clear line uses fine lines to show movement; their description of movement is subtle and sometimes can be hard to notice. On the other hand, Atomic style uses different techniques to show movement. The author may use shadow images, mirages, white streaks or lots of lines behind the moving object. On the picture 1.1 we can see the difference between the Clear line, 1.1a, and the

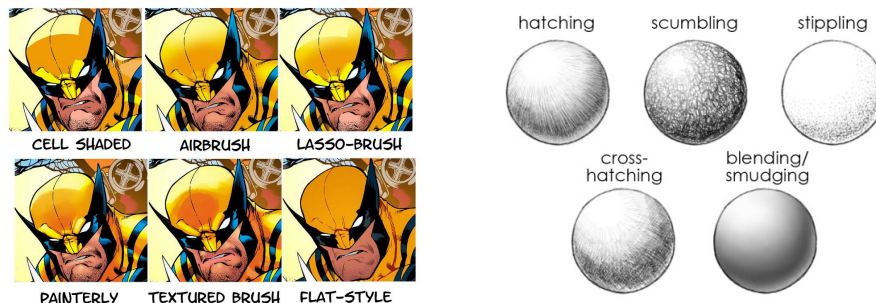
⁹They use what can be called an icon, a picture that represents a person, a place, a thing or an idea.[2]



(a) : Example of the Clear line. [16]

(b) : Example of the Atomic Style. [17]

Figure 1.1: Examples of the Clear line and the Atomic style.



(a) : Colouring techniques.[18]

(b) : Shading techniques.[19]

Figure 1.2: Examples of colouring and shading techniques.

Atomic style, 1.1b. Notice how the Clear line uses fine lines that are not too obstructive of the whole picture. On the other hand, Atomic style uses more dynamic lines that can sometimes obstruct parts of the picture. But there are more ways to indicate movement. Sometimes artists use more than one style in their comics, and depending on the scene use different movement style.[2]

1.5.3 Colouring and Shading

There are multiple ways of colouring and shading a comic. From the previous sections, we know that there are black and white comics and that there are also coloured ones.

On the picture above 1.2a we can see six basic styles of colouring for coloured comics. On the other picture 1.2b we can see five shading styles for black and white comics. But these are only some examples of techniques that can be used. The author may use different variations and combinations of these techniques or completely different techniques.

Colouring or shading are used to mimic light sources and surface textures.

Some artists can only use these techniques to mimic light sources but it is upto the reader to imagine what kind of material may be used and of what sort of texture it has.[2, 20]

Chapter 2

Related Work

Before we can commence our work we must know what was done previously on this or related topic. We could use parts of these works for comparison or as a way to avoid creating already created methods. Works listed and described in this section are in chronological order as I discovered them.

2.1 Style embedding

2.1.1 FaceNet

As stated in the introduction, FaceNet is what inspired me to write this paper. They describe a revolutionary method for face recognition in low dimensional space. They use Inception v1 from Google with triplet loss as their main model. They also describe efficient methods for triplet selection. Their method can be applied to general embedding creation. The current best FaceNet model achieves accuracy of 99,65%. [21]

2.1.2 Illustration2Vec

Illustration2vec is a similar paper that tries to create embeddings from illustrations. The only difference is that they are extracting a predefined set of features from the given image. They use 4096 different features and some of them describe the illustration features and others describe the style or a specific character. [22]

2.1.3 Are Anime Cartoons?

In a paper called “Are Anime Cartoons?” they are trying to see whether neural networks can differentiate anime and cartoons. They portray the fact that anime is different from a cartoon. But they only do a classification between these two classes. [23]

items. They achieve better results than classical SVD method.[27]

■ 2.1.8 A Visual Embedding for the Unsupervised Extraction of Abstract Semantics

In this paper, authors describe a method for image content embedding based on ImageNet classes and their connection to WorldNet. They produce one million dimensional embeddings. They created the embeddings by concatenating 27 layers of GoogLeNet, Inception v1. On these embeddings they show that the connection to WordNet can be used to make arithmetic searches in the resulted vector-space.[28]

■ 2.1.9 CNN-based Classification of Illustrator Style in Graphic Novels: Which Features Contribute Most?

This is a similar paper but with the difference that it is focused on classification and then on embedding. The resulted classification vector can be viewed as an embedding. They are working on their own dataset and are using Inception V3. They show that features from five mixed-layers have an accuracy of over 97%. According to them, more layers show signs of overfitting. They use their private dataset called Graphic Narrative Corpus and Manga109, where they classify the comics by illustrator and book. They use whole pages as an input into the network. They also discuss their finding with regards to embeddings and their surprise with how well it works like an embedding for Nearest Neighbors.[29]

■ 2.2 Text Embedding

■ 2.2.1 Efficient Estimation of Word Representations in Vector Space

Word2Vec is a key paper in the field of word embedding. They describe two methods of how to produce word embeddings efficiently. Their proposed methods are Continuous Bag-of-Words Model and a Continuous Skip-gram Model. They also show that simple vector arithmetic's can be done on these word vectors.[30]

■ 2.2.2 Distributed Representations of Sentences and Documents

This paper builds on top of Word2Vec 2.2.1 and is using a similar method. They show a way of creating embeddings of documents. These embeddings

are easily reusable. Same length embeddings can be compared to embeddings created from different document sets. They also provide comparison with current state of the art methods.[31]

Chapter 3

Style embedding

3.1 Panel vs. Page

In [29] they use sections of whole comic pages as training inputs. I decided to use panels because in my interviews [4, 5, 32] I learned that panel layout is mostly used as part of the story telling and not part of the art style 1.5. But other crucial factor is that webcomics are often a vertical column of panels and thus the pages are high and long or wide if they decide to use horizontal strip format. But there are also comics that play more with the space and the medium, they are given and the panel layout may stretch across several classical pages.

Another concern is that the model can learn the layout of the page and describe the comic based on the layout, or conversely treat comics with different art styles as similar comics only because of the similar layout.

This decision complicated the dataset creation process. But it was shown that in [24] there are panel extraction techniques with good results.

3.2 Dataset

Before we can do any work we must find or create a suitable dataset. Looking through related work, there can be found two good examples of a comic dataset. First is the Graphic Narrative Corpus[33] the second being the dataset used in [24], which will from now on be referenced to as the COMIC dataset.

Sadly the Graphic Narrative Corpus is not publicly available. So I decided to use the COMIC dataset. On the initial inspection, the dataset seemed suitable for this work but required to be cleaned from badly detected panels and advertisements that were present in the magazines the comics were scanned from. After a few weeks of cleaning the dataset, I started noticing that there are more issues with the dataset than initially thought. Mainly that the style of the drawing is really similar and also that there are multiple comics treated as one.

I started looking through the comics and tried to find their authors with the hopes of reorganizing the dataset. But the only information I could find was the company who produced them. From my interviews I learned that comics from Golden Age were made by teams where each member did only one part of the comic and most of the times there was no real artist involved.

With these discoveries in mind, again I started looking at possible datasets I could use. I found two datasets, one being the Japanese Manga109[34, 35] with 109 Mangas. Sadly this dataset is only limited to Japanese comics. The other dataset was eBDtheque[36] which is comprised of European comics, unfortunately this dataset is rather small and there are not many samples of comics.

Not being able to find any suitable dataset I decided to create my own. I started by acquiring the Manga109 dataset. Then I started collecting comics from Line Webtoon[37] and other internet sources. I also included some illustrations[38] since we allowed them in our definition. Some of these illustrations are computer generated. Since I am creating my own dataset lets propose what would the ideal dataset we would want to create.

Proposed dataset: Comics and Mangas, manhwas, webtoons from all over the world. Mainstream and underground work. Spanning the time period of at least 30 years. With at least two thousand of different comics where each would have at least 250 panels.

The proposed dataset would be representative enough to encompass all the art styles and narratives. Sadly it is unrealistic for me to create this dataset on my own in such a short time and with the limited resources.

As a result of this I will be using much smaller dataset but similarly structured. The dataset I will be using is composed of the Managa109 dataset, 85 comics from the COMIC dataset, 42 comics from various internet sources and 223 illustrations from various authors. Totalling in 163 thousand images. This dataset is not as representative as the proposed one but it much better than publicly available datasets.

■ 3.2.1 Data augmentation

Since our dataset has only one hundred sixty thousand samples which are all similar with only small differences to avoid overfitting and making the resulting solution more robust I am gonna apply some data augmentation methods. There are lots of techniques that can be used, for the sake of simplicity I will only use three described below.[39, 40]

- Gaussian or salt and pepper noise

- Features with high occurrence tend to lead to overfitting. To avoid that we will apply Gaussian noise with zero mean which distorts high-frequency features. It also affects lower frequency features but the benefit of overfitting is much greater.
- Flip
 - Most common comics are read right to left but Japanese comics are read left to right. To compensate for this difference we augment the data with flipped copies to avoid overfitting on these two regional categories.
- Colour intensities
 - As we learned before colours are also part of art style. Some comics could be badly scanned and the colour could be altered from the original. And also what looks like one shade of colour could be different on a different viewing device. But the most important reason to alter colour is to make the network more robust. Images have randomly altered brightness and contrast.

■ 3.2.2 Data split

I decided to split the dataset into training and evaluation parts in 80:20 ratio. I have decided to not to split the dataset into three parts for cross-validation because of the computational intensity of the training. I choose the 80:20 ratio because the dataset does not have an equal number of representatives for each comic and also because it's according to Pareto principle a good ratio.

■ 3.3 Distance Metric

Initially, I thought about testing different distance metrics and finding the one that works best in this use case, but since I am planning to create the embeddings in Euclidean space, there is no need to look for best metric as the best metric is given by the space I am building the embeddings in.

■ 3.4 Clustering algorithms

I will be using clustering algorithms mostly for analysis of the results and also for showing how well the embedding can be clustered. There exist lots of clustering algorithms, but for the sake of this work, I will be describing only

three. Them being K-means as representative of clustering where we need to specify the number of clusters we want to find. HDBSCAN an extension of DBSCAN. And Markov clustering as representative of graph clustering.

3.4.1 K-means

K-means is one of the most used and simplest unsupervised clustering algorithms. It tries to create K centroids of clusters where close neighbours of this centroid share some underlying pattern. The main disadvantage is that we need to have some prior knowledge about the data and specifically the number of clusters we expect.

K-means is an iterative algorithm that tries to minimize an inertia defined as follows:

$$\sum_{i=0}^n \min_{\mu \in C} (\|x_i - \mu_j\|^2).$$

Where n is number of samples. C is set of centroids, μ_j is centroid and x_i is data point.

The algorithm first creates K centroids, most often random points in data space or random points from data are used. Then it iterates between two steps until it reaches minimal changes in the clusters or after a defined number of steps. The first step is “Data assignment step” where each data point is assigned to its nearest centroid, any distance metric can be used. The next step is “Centroid update step” where the centroids are updated to the mean of the data points assigned to the specific clusters. Formally the new centroids are defined by:

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

Where c_i is the new cluster centroid S_i are the data points that were assigned to this cluster. And x_j is a specific data point.

The algorithm will converge to a minimum but it can be a local minimum, and thus often multiple runs are done with different initial centroids. The disadvantage of knowing the number of clusters can be overcome is by measuring the average distance to centroid within the cluster and finding an elbow point in these measures.[41, 42, 43]

3.4.2 HDBScan

HDBSCAN, as stated before, is an extension of DBSCAN. From a users point of view, the main difference is that DBSCAN uses ϵ which is radius in which neighbours have to be co cluster core to form. On the other hand, HDBSCAN uses a minimum number of data points to create a cluster. It also solves an issue DBSCAN has with data points that have varying densities.

■ How it works

First HDBSCAN transforms the data space so that close points are closer and distant points are farther away. This is achieved by defining a new distance called mutual reachability distance which uses core distance and our provided distance metric. The core distance is defined as

$$core_k(x) = d(x, y_k).$$

Where the x is some point and y_k k -th neighbour of this point. $d(x, y)$ is our provided distance metric. So the core distance is the distance for each point to its k -th neighbour. The mutual reachability distance is defined as:

$$d_{mreach-k}(a, b) = \max\{core_k(a), core_k(b), d(a, b)\}.$$
¹

Now that we have the new space as defined by the mutual reachability distance matrix the algorithm builds a minimum spanning tree. From this tree, we create cluster hierarchy by taking the edges sorted by the distance in increasing order and iterating over them and creating a new merged cluster for each edge. The merging of a cluster is done using a union-find data structure. Technically this is where DBSCAN stops and makes cut in this hierarchy based on the ϵ parameter. But HDBSCAN continues by condensing the cluster tree. For this, we need a hyperparameter called minimum cluster size. And the condensation is done by first assigning all points to one cluster and then going through the cluster hierarchy tree from its root. For each split, we ask if the new cluster created by the split has fewer points than the minimum cluster size then these points “fell out of cluster”. But if there are more points than the minimum then this split created two new clusters.

From this new condensed cluster hierarchy, we want to extract the clusters. For this, we need to measure the stability of persistence of the clusters. We can do that by defining measure of persistence as follows $\lambda = \frac{1}{distance}$ using this measure we define three values λ_{birth} , λ_{death} and λ_p . The λ_{birth} is value for when the cluster split off and became its own cluster. The λ_{death} is when the cluster split into smaller clusters or when it “died” meaning it was split into clusters of size less than the minimum. The λ_p is value for each point where this point “fell out of cluster”. Using these values, we define stability for each cluster as:

$$\sum_{p \in cluster} (\lambda_p - \lambda_{birth}).$$

To extract the clusters we first set all leaf nodes to be our selected clusters and go through the condensed tree in reverse order, from leaves to root. Now

¹Why this works is described in more detail in [44]

if the sum of stabilities of child clusters is greater than the stability of the current cluster, then we set the cluster stability to be the sum of the child stabilities. But if the current cluster stability is greater than the sum of its children stabilities then we declare this cluster as our selected cluster, and we unselect all of its descendants. This gives us a flat clustering.[45, 46]

■ 3.4.3 Markov clustering

Markov clustering is an unsupervised cluster algorithm for graphs, networks. The algorithm is based on simulation of stochastic flow in graphs. The flow is simulated by alternating two algebraic operations on matrices. The first operation is called expansion and the second one is inflation. This algorithm is fairly complex and yet explored by popular media. The description of how it works in detail is out of the scope for this work. More information can be found at [47] and in the original work [48].

■ 3.5 Dimensionality reduction

Dimensionality reduction can mean two things one is feature extraction and the other is feature selection. Feature selection is trying to select the best features from the original feature set that holds the most information so we can still accurately and efficiently do the task on hand. On the other hand feature extraction try to select new features based on the original features.

I am looking for dimensionality reduction that would allow me to easily visualise the resulting embeddings and also to decrease the embedding dimensionality. I might end up choosing two methods. There are many feature extraction methods but for the sake of this work, will be focusing only on selected few. Them being Principal Component Analysis (PCA), Isomap, Multi-Dimensional Scaling (NMDS), t-distributed Stochastic Neighbor Embedding (t-SNE) and Uniform Manifold Approximation and Projection (UMAP).

■ 3.5.1 PCA

PCA is an unsupervised method. It is one of the most used methods for dimensionality reduction. PCA tries to represent correlated features with new uncorrelated features that are linear combinations of the original features. PCA assumes that relationships between features are linear. If they are nonlinear then the new features created by PCA are not an informative or efficient way of reducing the dimension of the features.

■ In detail

PCA in its most straightforward explanation is just linear transformation that is defined by a simple equation

$$Y = UX : (k * m)(m * n) = (k * n)$$

where X is a matrix of source data, Y is a matrix of transformed data and U is some transformation matrix. With PCA the U is a matrix of eigenvectors of covariance or convergence matrix.

Let's assume we have matrix X with shape $m \times n$. X is the input data in matrix form, where rows give us values of particular feature for all samples. And columns give us features for a specific example. And we want to transform this data from m -dimensional space to lower k -dimensional space using PCA.

First, we calculate the covariance matrix using this formula

$$C_x = \frac{1}{n}XX^T$$

where n is the number of samples. We know that C_x is a square symmetric $m \times m$ matrix and that on the diagonal of this matrix we have the variance of particular features. The off-diagonal values are covariances between features. We then compute the eigenvectors and eigenvalues of this matrix. We sort these eigenvectors according to their eigenvalues in descending order. Then we select top k eigenvectors from this sorted set and form a matrix of them where each eigenvector is a row of the matrix. This matrix of eigenvectors is the U matrix in the equation above. So we can multiply this U matrix by the X matrix, and we will get the transformed data Y .

There are other ways of calculating PCA; for instance, we can instead of covariance matrix use correlation matrix, or we can use a different approach from the eigenvectors altogether and use Singular value decomposition, SVD.[49]

■ 3.5.2 MDS

MDS is an unsupervised method for dimensionality reduction. MDS tries to map the original N -dimensional features into M -dimensional space as such that distances in the original space are similar to distances in the new space. There exists another variant of MDS which is called non-metric MDS. Which aims to preserve ordinal relations of variables instead of the distances in the new space. MDS has a similar problem as PCA with nonlinear data.

■ In detail

The process to calculate MDS is almost the same as for PCA. The only difference is that we don't use covariance or correlation matrix but instead we use pairwise distance matrix. On the other hand, non-metric MDS is much more complex. I will be focusing on the non-metric variant as I already described PCA in the previous section.

Let's assume we have M-dimensional data and we want to reduce the dimension to N-dimensions using non-metric MDS. First, we have to calculate D which is a pairwise distance matrix of the original M-dimensional data. Alternatively, we can use a pairwise similarity matrix. Then we make random representations of the original data in the N-dimensional space. And we calculate D' which is pairwise distance matrix of the data in the new N-dimensional space. We only need the lower or upper triangle of these two matrices.

Now the algorithm starts iteratively moving the new data points in the new space while minimizing loss function called Stress which is defined as follows:

$$Stress = \sqrt{\frac{\sum (f(x) - d)^2}{\sum d^2}}$$

Where f is a monotonic approximation of the relationship between x and d . d is a distance from matrix D' and x is similarity or distance from the D matrix. The minimization is done by using gradient descent². After every step, we have to recalculate D' and change function f .

It's also important to note the value of the Stress in the final step as it describes how well the new N-dimensional space represents the original data. It's best to have stress lower than 0.1.[50]

■ 3.5.3 Isomap

Isomap is a nonlinear unsupervised method for dimensionality reduction. Isomap is based on MDS. The main difference between MDS and Isomap is that MDS uses straight distance, but Isomap uses geodesic distance. The geodesic distances are the shortest paths along the curved surface of the manifold measured as if the surface were flat. This distance is approximated with distances between neighbours on graph. This solves the issue MDS has with nonlinear data. Isomaps have a problem with nonconvex manifolds. It also has issues with manifolds that have "holes", areas with no data points.

²Alternative minimization algorithm exists and is called SMACOF algorithm.

The last weakness is that there can be an erroneous connection in the geodesic distances approximation and thus creating a connection where the distance should be big but because of this error shortcut appeared and the distance is short.

■ In detail

Let's assume we have M-dimensional data and we want to reduce the dimension to N-dimensions using Isomap. First, we are going to construct a weighted graph where each data point is connected to its neighbours in a fixed radius or alternatively we can use k-nearest neighbours. The weights of these connections, edges, are the distances between the data points. Then we calculate a pairwise distance for all the points but in this graph. Dijkstra's or Floyd-Warshall algorithm can be used here. These distances are then used to classical MDS, which I described above.

Isomap uses classical MDS but if we were to calculate the geodesic distances separately we could use these distances with non-metric MDS.[51]

■ 3.5.4 t-SNE

t-SNE is a nonlinear unsupervised method for dimensionality reduction. It tries to place the points in the lower dimension not based on distance, straight or geodesic, from other local points. But on a normalized probability of points in some distribution, where the probability represents the average pairwise similarity between the points in the original space. In the original space for the probability, a normal distribution is used and for the new space, student distribution is used. t-SNE mainly preserves the local structure of the data. t-SNE is mostly used for data visualisations.

■ In detail

Let's assume we have M-dimensional data and we want to reduce the dimension to N-dimensions using t-SNE. First, we calculate a pairwise probability of picking the other points as their neighbours if the neighbours were to be picked in proportion to their probability density under Gaussian distribution. First, we calculate conditional probabilities with this formula:

$$p_{i|j} = \frac{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}\right)}$$

Where x_i , x_j and x_k are points in our source M-dimensional space. And σ^2 is a variance of Normal distribution for given x_i . We are setting the σ^2 so that each point has a fixed number of neighbours.

And now we calculate the final pairwise probabilities with $p_{ij} = \frac{p_{i|j} + p_{j|i}}{2n}$. Where N is the number of samples. Then we create random N -dimensional representation of the source data. And calculate their probabilities, but using a different formula. This time we use Student t -distribution with one degree of freedom instead of Gaussian distribution. The formula is as follows:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

Now we want to make the p_{ij} be as similar to the q_{ij} as possible. By doing so the N -dimensional representation has a similar structure to that from M -dimensional space. We achieve this by minimizing Kullback–Leibler divergence. Defined as follow:

$$KL(P||Q) = \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

We are gonna minimize the loss by using gradient descent. This version of t -SNE is fairly slow $O(n^2)$. To speed it up can apply Barnes-Hut approximation. What this approximation does is it takes points which are similar, close to each other, and far enough from the source point³. We then calculate the mean of these points and then use this mean for the calculation of the source point movement. Points close to each other are found using a quadtree. With this, the complexity is now $O(n \log(n))$.

There exists a special variant of t -SNE which creates multiple N -dimensional embeddings. In these embeddings, we are also calculating the weights of the data points. And the algorithm optimizes all the maps at once.[52, 53, 54]

■ 3.5.5 UMAP

UMAP is a nonlinear unsupervised method for dimensionality reduction. It analyses the topology of the original high dimensional data. It works by representing the original topology with simplices, which are easy representations of topological space. The simplices are constructed by using fuzzy “cover”. The cover here describes radius or area in which neighbours for construction the simplices are found. Then it tries to represent the data in lower dimension iteratively similar to what t -SNE does but also with regards to the topology. UMAP assumes that the data is uniformly distributed on Riemannian manifold. That the Riemannian metric is locally constant (or can be approximated as such). And that the manifold is locally connected, every point is connected to at least another.

³The point which movement we are calculating.

In detail

Let's assume we have M-dimensional data and we want to reduce the dimension to N-dimensions using UMAP. First, we are gonna construct k nearest neighbours with the given metric d. Then for each point from our input dataset, we define ρ_i and σ_i . Where ρ_i is as follows:

$$\rho_i = \min\{d(x_i, x_{i_j}) \mid 1 \leq j \leq k, d(x_i, x_{i_j}) > 0\}$$

And we set σ_i to satisfy this:

$$\sum_{j=1}^k \exp\left(\frac{-\max(0, d(x_i, x_{i_j}) - \rho_i)}{\sigma_i}\right) = \log_2(k)$$

In these two formulas, x_{i_j} is a jth neighbour of x_i . Now we are going to construct a weighted graph. Where the vertices are our data points, edges are defined by $\{(x_i, x_{i_j}) \mid 1 \leq j \leq k, 1 \leq i \leq N\}$. The weights are defined by this function:

$$w((x_i, x_{i_j})) = \exp\left(\frac{-\max(0, d(x_i, x_{i_j}) - \rho_i)}{\sigma_i}\right)$$

Now let A be the weighted adjacency matrix of this weighted graph and consider the symmetric matrix $B = A + A^T - A \circ A^T$ where \circ is Hadamard⁴ product. Then B is an adjacency matrix of an undirected weighted graph that the high dimensional data.

Now we use spectral layout⁵ to initialize the low dimensional representation of the data. Then we use a force-directed graph layout algorithm to move the points in the low dimension space. We need to provide attractive and repulsive forces to the force directed graph layout algorithm. The attractive force between two vertices i and j at coordinates y_i and y_j is as follows:

$$\frac{-2ab \|y_i - y_j\|_2^{2(b-1)}}{1 + \|y_i - y_j\|_2^2} w((x_i, x_j))(y_i - y_j)$$

And the repulsive force by:

$$\frac{b}{(\varepsilon + \|y_i - y_j\|_2^2)(1 + \|y_i - y_j\|_2^2)} (1 - w((x_i, x_j)))(y_i - y_j)$$

Where a and b are hyper-parameters and ε is a really small number⁶ to avoid division by zero. The repulsive force is sampled meaning that when an attractive force is applied to an edge one of that edge's vertices is repulsed by a sampling of other vertices.[55, 56]

⁴Or pointwise product can be used here instead.

⁵Random initialization can also be used but has slower convergence.

⁶Value of 0.001 is used in practise.

3.5.6 Conclusion

For testing the different methods I used digit dataset which is a smaller version of MNIST dataset.

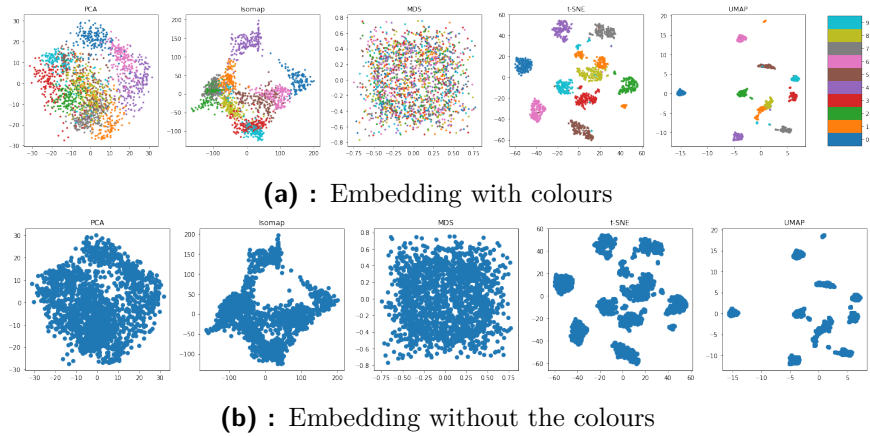


Figure 3.1: Subset of MNIST dataset reduced by all previously described algorithms.

First, let's focus on the visualisation aspect of the data. From the picture 3.1 it's clear that MDS provides the worst results, maybe if I would tune its hyperparameters it would provide better results. But since the others provide better results without tuning their hyperparameters I am gonna assume that it's not as robust as the other methods. PCA and Isomap give much better results than MDS. PCA clearly shows the global relationships between the classes. Isomap gives a much clearer picture of the data and even shows us some hints of clusters. But when we look at the pictures without the labels 3.1b it's hard to see any clusters. Some may be found in the Isomap features but still not the desired ten clusters.

On the other hand, t-SNE and UMAP provide the best results. They clearly show clusters and how each class is similar to others. UMAP clusters are more tightly packed. They both try to arrange the clusters more logically so that clusters of a similar class are closer together. The only major difference between them is that t-SNE puts zeros closer to nines and twos and far from twos where UMAP does the exact opposite. UMAP is faster than t-SNE, where t-SNE takes 22 minutes UMAP takes only 98 seconds to analyze the whole MNIST dataset.[57] This is why I choose the UMAP as a method to visualise clusters in the data.

Now let's look at the dimension reduction aspect of these methods. Since this is what is going to be calculated after every change in the underlying models it needs to be something that is noticeably fast. The dataset used in this work has more than 160000 images and according to benchmark that

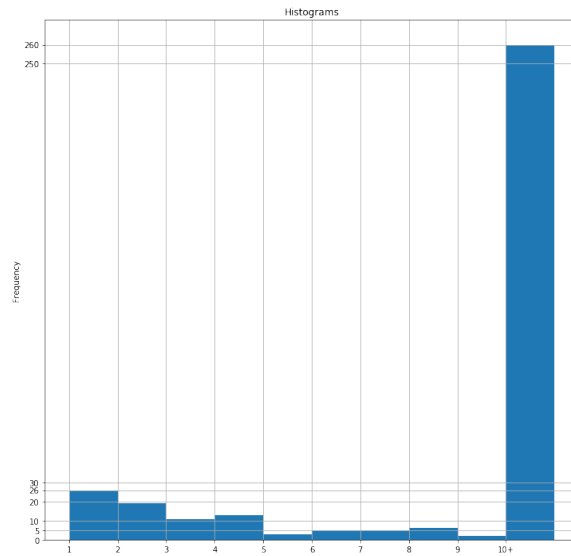


Figure 3.2: Accuracy of histogram baseline.

was run on MNIST dataset, the fastest two are PCA and UMAP. Since we are already using UMAP for visualisation I will also use it for dimension reduction.

3.6 Baseline

We need a baseline to compare our results to. The closest work that tries something with images scores over 90 % accuracy but they are doing classification and not embedding. FaceNet which is the inspiration for this work got 99 % accuracy but on faces. So there is no comparison for this work. So I decided to make two baselines one naive and one more robust. For the naive one, I choose colour histograms as they are one of the most basic features that can be extracted from an image. As the more robust baseline a selected Histograms of oriented gradients, HOG.

3.6.1 Colour histograms

I calculated histograms of the three channels and concatenated them together for all the images so as a result, I had 768-dimensional vectors. Then I made an average of the histograms for each comic, created centroids for each comic cluster. Then I made a basic kNN model and queried random panel from each of the comics, with up to ten neighbours. The results can be seen in picture 3.2.

When we look at the results on the picture 3.2 we can see that most of the results are the in the tenth or farther neighbour column and only some are

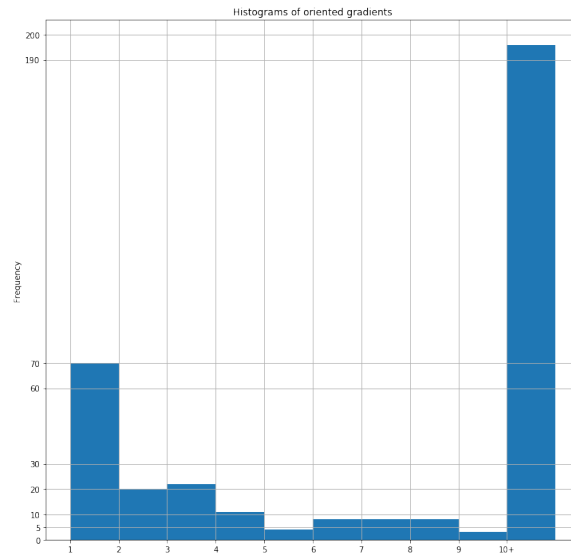


Figure 3.3: Accuracy of histogram of oriented gradients baseline.

in the top three. When looking at visual similarities, Appendix A, it only managed to partially separate black and white comics from full-colour comics.

This model does not detect style and only relies on the colours of the images. If we were to measure the accuracy, then it only managed to correctly detect twenty-six comics, which is only 7,4%. As the last success measure, I run HDBSCAN on the colour histograms, and it detected only one cluster and the rest were classified as noise.

■ 3.6.2 Histogram of oriented gradients

I calculated HOGs with 18x18 pixels per cell and nine cells per block. This resulted in only four thousand dimensional feature vector. As for a normalization method, I used L2-Hys. Similarly to the colour histograms, I averaged these vectors across all the comics. Afterwards, I again created the same basic kNN model with up to ten neighbours. Where I queried the same pictures i used for the colour histogram. The results can be seen in picture 3.3. When we look at it we can see that the model classified correctly more comics in the top three neighbours in comparison to the colour histograms and it even managed to classify seventy comics correctly. Nevertheless, a majority of the comics were in the section of ten or more

This model does a better job at detecting comics than colour histograms, but it still doesn't completely understand the style of comics. If we were to measure the accuracy then we could see it only managed to correctly classify seventy comics which is only 20 %. As the last measure of success, I ran HDBSCAN on the HOGs and it detected only three clusters and the rest was

classified as noise.

■ How it works

First, we have to calculate the gradients which is done by simply convolving the image channels with two vectors $\begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$ and $\begin{bmatrix} -1 & 0 & 1 \end{bmatrix}^T$, so we end up with two $m \times n \times u$ matrices; one for each convolution. With these two matrices, we have to calculate the matrix of magnitudes and the matrix of angles of the resulting gradients. To calculate the matrix of magnitudes we use the following formula

$$g = \sqrt{g_x^2 + g_y^2}$$

where g stands for the magnitude of the resulting gradient and g_x and g_y are the matrices we get from the convolution. As for the matrix of angles, we use

$$\theta = \arctan\left(\frac{g_y}{g_x}\right)$$

where again the g_x and g_y are the matrices we get from the convolution and θ is the resulting matrix of angles. For coloured pictures, the magnitude and angle of each pixel is set according to the channel where the pixel has the highest magnitude. Here we can decide whether we want to make the gradients signed or unsigned if we decide to use signed gradients we don't need to do anything. On the contrary, if we decide to use unsigned we need to transform the angle so that it is an element of $[0, 180]$.

The second step is to calculate a histogram over sections of the image, cells. The cell dimensions depend on the task at hand, and there is no general rule as to what size to use for them. Also, the number of bins depends on the task that we are dealing with. The histogram is created in a weighted way where the bins are the angles, and the values that are added to the bins are proportional to the distance from the bin centre. For example if we assume that we have bins "0", "20" and "80" and we have two gradients $(80, 2)^7$ and $(10, 4)^8$ the first gradients add its whole magnitude to the "80" bin but angle of the second gradient lies between 0 and 20 degrees so its magnitude is split between "0" and "20" bins. When calculating this histogram, we have to keep in mind that the angles wrap and so if our last bin is "160" and we have an angle of more than 160 then the magnitude is split between the "160" bin and the "0" bin.

The third step is block normalisation. The normalisation method used in this step also depends on the task, but L2-Hys⁹ was found to perform the

⁷80 is the angle and 2 is the magnitude.

⁸10 is the angle and 4 is the magnitude.

⁹L2 normalisation then clipping the maximum to 0,2 and again L2 normalisation.

best. Input for block normalisation is a block of cells, again the size of the block depends on the task. And we normalise histograms in this block as if they were one concatenated vector. We are saving the normalised vector to a new matrix. After the normalisation, we move the block one cell to the left and normalise again. When we normalise the leftmost block we start again from the right, but now we move one cell downwards, and we repeat this process until we reach the last block. This way each cell can contribute more than once in the resulting vector.

The final vector consists of all these normalised histograms concatenated.[58, 59]

3.7 CNN

Convolutional Neural Networks, CNN, is Neural network that uses convolutions as its main component. Every CNN has four main components:

■ Convolution

- It's a general purpose filter. Where mask, matrix, is moved across source data and mathematical calculation is done within this mask.
- The mask movement is defined by stride, which is how many pixels it has to move to make new convolution.
- Output of this operation is value that is weighted sum of all elements from the source data. The weights are taken from the mask.
- The mask always has odd dimensions, 3x3, 5x5, 7x1, 7x3, etc.
- The mathematical formula for convolution is

$$(f * g)(x, y) = \sum_{s=-a}^a \sum_{t=-b}^b h(s, t) f(x - s, y - t)$$

where f are the source data, g is the mask x, y are coordinates of the resulting new pixel. s and t are indexes of the mask and a and b are $\lceil \frac{w}{2} \rceil$ resp $\lceil \frac{h}{2} \rceil$ where w and h are width and height respectively. Functions h and f return pixel value for given coordinates.

- If we are working with coloured images we also have to sum the convolutions across the three dimensions.
 - A bias is added to the resulting value.
- ### ■ Non Linearity
- Non linearity is done by applying activation function to each of the element, pixel, that is a result of convolution.

- The most used functions are ReLU $f(x) = \max(0, x)$, Sigmoid $(x) = 1/(1 + e^{-x})$, Tanh $\tanh(x) = 2/(1 + e^{-2x}) - 1$ and special variant of ReLU called Leaky ReLU $f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
- Pooling or Sub Sampling
 - They are reducing the spacial size of the outputs from convolution layers and thus reducing number of parameters and computations in the network. They also help prevent overfitting.
 - Again similar to convolution a mask moves over the data and aggregate a value from within the mask.
 - The most common aggregations are maximal or average.
- Fully connected layer
 - This is classical fully connected layer, where each of its inputs is connected to all outputs from the previous layer.

The most common pattern for CNN is as follows:

INPUT \rightarrow [[CONV \rightarrow RELU]*N \rightarrow POOL?]*M \rightarrow [FC \rightarrow RELU]*K \rightarrow FC

What is actually learned are the convolution filters, biases and fully connected layer weights. Training is done with backpropagation.[60, 61]

■ 3.7.1 Inception

Inception is deep CNN network. There are multiple versions of Inception style CNN. They all share one common characteristic, and that is that they are made of blocks. These blocks perform different convolutions on the same input. This allows us to run different sized kernels on the same input.

The first version of inception had a single block that consisted of 1x1, 3x3 and 5x5 convolutions and also 3x3 max pooling. This block was chained nine times before these block there was a stem section. The whole network has 22 layers. They also added two auxiliary classifiers that were to be used to help train the network.

The second version improved the first by creating three types of inception blocks. In the first block, they replaced 5x5 convolutions with two 3x3 convolutions as apparently 5x5 convolution is 2.78 times more computationally expensive than single 3x3 convolution. In the second block type, they did the same replacement as in the first one, and they also replaced all 3x3

convolutions with 1x3 and 3x1 convolutions as apparently, this leads to a 33 % efficiency increase. In the last block type, they again did the same replacement as in the first block, and they replaced the topmost 3x3 convolutions with 1x3 and 3x1 convolutions but not chained. Instead, they share the same input. The whole network has a total of 50 layers.

The third version is only a slight upgrade of the second version. They factorised the initial 7x7 convolution to three 3x3 convolutions. They also added batch normalisation to the auxiliary classifiers.

In the fourth version, they changed the stem of the network and made it into an inception block. They changed convolutions in the second block so that they are now 7x1 and 1x7 convolutions instead of the previous 1x3 and 3x1 convolutions. And they added two new types of block types called “Reduction Blocks” which reduce the dimensionality of the input. They also changed how the blocks are chained. The total number of layers increased to 97.

There are two more variants of the inception CNN which is called Inception-ResNet. Its hybrid between Inception and ResNet where each block has a residual connection. The only difference between these two variants is that they have different stems. The residual connection replaced the pooling layer. The second and third block types got completely changed. The second block has one 1x1 convolution, and one 7x7 convolution factorised into 1x7 and 7x1 factorisation. The third block is similar but instead the 7x7 convolution there is 3x3 convolution factorised into 1x3 and 3x1 convolutions. They also changed the second reduction block. In total, the hybrids have around 144 layers[62, 63]

First, I started by running the third version of the Inception model. I choose the third version as it provides good accuracy and requires only about 12 G-ops to compute one pass[64]. On this version, I tried to tune hyperparameters to get best possible results and later use the same hyperparameters for the other models. I was not training the models from scratch, but I used pretrained weights. The pretrained models were trained as Imagenet classification. The downside of this approach is that the networks output 1001 dimensional vector but since the model is pretrained the training should be faster and the final model should be more robust than if it were trained from scratch. Another downside of this approach is that the chosen hyperparameters are based on Inception V3 and there could be better for the other models. As for the other models, I will be training Inception V4 and Inception-ResNet-V2.

I will be training the CNN on single TPU. I can use the TPU because I got 30 days trial period from TensorFlow Research Cloud, which is a cluster of more than 1000 TPUs. There are multiple version of TPUs I will have access to TPU v2.8 which has 8 cores and 12GB of RAM.

The different models I will be using have a different computational cost. With the Inception V3, the TPU can process 1400 images per second. Inception V4 is a much deeper network, almost double that of Inception V3, and as result, this TPU can do only do 700 images per second. Despite Inception-ResNet V2 being almost triple the size of Inception V3 the TPU can do 580 images per second.

■ 3.7.2 RMSPROP

RMSProp is mini-batch optimisation algorithm with similarities to rprop and Adagrad. It builds on top of rprop by normalising the learning rate by moving averages of squared gradients, instead of fixed factor. This helps to solve the issue that mini-batches have different gradients.

The moving average is defined as follows:

$$E[g^2]_t = \beta E[g^2]_{t-1} + (1 - \beta) \left(\frac{\delta C}{\delta W} \right)^2$$

Where β is moving average parameter, $\frac{\delta C}{\delta W}$ is gradient of the loss, cost, function with respect to the weight.

And the new weight is defined as follows:

$$w_t = w_{t-1} - \frac{\eta}{\sqrt{E[g^2]_t}} \frac{\delta C}{\delta w}$$

Where the η is a learning rate.

Similarity with adagrad is that adagrad uses a sum of squares of the gradients instead of moving averages which are used here. The sum of squares has effect where the steps get smaller and smaller as the training proceeds.

I decided to use RMSProp because of its fast convergence rate and also for its ever-increasing popularity.[65]

■ 3.7.3 Loss

There three main losses that I could use. Unfortunately, I am unable to use one of them because of TPU limitations. There are Triplet loss and Lifted structured loss. I will be using the former one as the main loss despite the fact that the later one was proven to be better. The reason behind this is that it was used in FaceNet where it achieves really great results.

■ Triplet loss

The triplet loss tries to enforce a margin between each pair of panels from one comic to all other panels. This allows the panels for one comic to live

on a manifold, while still enforcing the distance and thus discriminability to other comics.

The embedding is represented in M-dimensional euclidean space, additionally it's normalized thus the embeddings lives on M-dimensional hypersphere. The loss is defined as follows:

$$\sum_i^N \left[\|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]_+$$

Where x_i^a, x_i^p and x_i^n are anchor sample, positive sample and negative sample respectively. And the α is margin. What this means is that the loss is trying to separate classes from each other by some margin. And samples from same class to be together.[66]

■ Lifted structured loss

Is a loss that tries to increase efficiency of the previous method. By not only having one anchor points per pair but that each sample from the same class is anchor and positive sample and also there is not only one negative sample but all other class samples are negative. This change allows to utilize the whole information from a minibatch. The loss is defined as follows:

$$J_{i,j} = \log \left(\sum_{(i,k) \in N} \exp(\alpha - D_{i,k}) + \sum_{(j,l) \in N} \exp(\alpha - D_{j,l}) \right) + D_{i,j}$$

$$J = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, J_{i,j})^2$$

Where P is a set of positive pairs and N is a set of negative pairs. D is a matrix of distances. And the α is margin.[67]

■ 3.7.4 Metrics

I am using multiple metrics to show and see how “good” are the models. While training I am using four different metrics from which only one is almost without flaws. The other metric I am using is Accuracy its similar metric to that i used when comparing baselines.

■ Average Class similarity

Simple metric that calculates average of averaged euclidean similarity for each comic. Can be defined as follows:

$$\frac{1}{|B|} \sum_{b \in B} \left(\frac{1}{|C_b|} \sum_{c \in C_b} \left(\frac{1}{\bar{D}_c} \sum_{d \in D_c} \frac{1}{1+d} \right) \right)$$

Where B is set of batches in the current epoch, b is set of embeddings. C_b is a set of classes for a particular batch, c set of embeddings for a particular class. D_c is a set of pairwise distances for a particular class, d is particular pairwise distance. \tilde{D}_c is count of nonzero distances for a particular pairwise distances. And d is Euclidean distance between two embeddings from a particular distance matrix. This metric flawed because if one class is split between two batches then the metric is not accurate.

■ Accuracy

This is another simple metric that shows how well are the embeddings for each class separated. The calculation is fairly simple. For each class, a centroid is created, from train dataset and then five random images from each class, from eval dataset are compared with the centroids. If more than three images from the same class have their corresponding class centroid as closest that it's classified as a match. And this metric counts how many percents matched.

Chapter 4

Text embedding

4.1 Dataset

Again we have to look or create a dataset that we can use. There are many textual datasets, but in this case, we need a dataset that contains stories or dialogues of people. I found a dataset of a transcribed movie scripts [68]. This dataset contains 1093 scripts from movies of different genres. This dataset could be used but it has one issue, and that is I would need to pair the movie script with the comic. This pairing could heavily influence the results.

My other option is to manually create my dataset by extracting text from the comics I already have. But here arises another issue. That issue being that COMIC dataset comes with pre-extracted text, but the quality of this extraction is really poor. There are missing words or even letters. The sentences don't make sense because the extraction works line by line and not per speech bubbles. This issue would come up if I attempted to do the text extraction on my own. Solving this would need to create a method for extracting the speech bubbles first and then extracting the text.

For the text extraction, I would choose Tesseract OCR as it is an open-source state of the art method for text extraction. There are other OCR systems, but they are either proprietary or are not as good as Tesseract.

But since part of my panel dataset is Manga109. And for these comics, I have the text for each speech bubble. But the downside is that this text is in Japanese.

Since the results of text extractions would not be accurate representations of the text in the comics, the reason being bad word order, missing letters and ghost text¹. I decided to make a compromise and for the Japanese comics use their respective text. And for the other comics, I will use a random script from [68]. If the comic did not have text in the first place, it would not receive a random text.

¹Text that is part of the picture but not part of the speech bubble, or other telling methods.

4.2 Doc2vec

I could explore multiple methods for text representation but since there is such a good² state of the art method I choose to use this method. Doc2Vec 2.2.2 builds on top of Word2Vec 2.2.1. Doc2Vec has three different methods on how to construct the embeddings of the document. The methods being paragraph vector with distributed memory(PV-DM), paragraph vector with a distributed bag of words(PV-DBOW) and a combination of the previous two.

The main issue of Doc2Vec is that it does not support multiple languages at once. So i will have to do the embeddings on the English text first and then on the Japanese text. Comics with no text will have get empty vector.

4.2.1 PV-DM

First, we map every paragraph to a Vector represented as a column in matrix D . And every word is also mapped to a vector represented as a column in matrix W . Then the Word vectors and paragraph vector are concatenated or averaged so they can predict the next word in context.

Formally the model is trying to maximize the average log probability:

$$\frac{1}{T} \sum_{t=k}^{T-k} \log p(w_t | w_{t-k}, \dots, w_{t+k}).$$

Where T is number of words, k is index of current world. And the w_t is vector of world t . The prediction task is done by softmax where the probability is defined as:

$$p(w_t | w_{t-k}, \dots, w_{t+k}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}.$$

Where the y is un-normalized log-probability defined as:

$$y = b + Uh(w_{t-k}, \dots, w_{t+k}; W, d_i; D).$$

Where the U and b are softmax parameters. And h is constructed by a concatenation or average of word vectors from W and paragraph vector from D

The context if fixed-length and is sampled from a sliding window over the paragraph. The paragraph vector is shared across all contexts generated from the same paragraph but across the paragraphs. The word vector matrix is shared across all paragraphs so in every paragraph one word has the same representation.

²Error of only 7.42 % on the IMDB dataset.[31]

■ 4.2.2 PV-DBOW

This method is much simpler and is similar to the skip-gram model used in 2.2.1. It works such as in every iteration of the algorithm we sample random text window from a paragraph and then sample random word from this text windows and perform classification task given the paragraph vector from matrix D . So, in reality, we change the paragraph vector in such a way that it predicts the words in a small window.

Chapter 5

Results

5.1 CNN

I ran every variant of the three models for at least 100,000 iterations. I did not recognize any major decrease in loss. Most of the time the loss started at around two, with the exception of models using the lifted struct loss. Even the loss of the variant with the random weights started at around two. This might be because of the way I was selecting the batches. I was randomly selecting two samples of panels from each comic. I could not implement the superior offline semi-hard selection because of the limitation in Tensorflow TPUEstimators. It would be inefficient and time-consuming to have a pipeline that would train the CNN a certain amount of steps then predict the whole dataset and then select the new triplets based on these predictions and then train again. This inefficiency in time is caused by the whole TPU initialisation process and also that we can't predict the embeddings on the TPUs and thus we have to use CPUs for the predictions. Online semi-hard selection is also out of the question because it would be yet again very time inefficient with the current implementation of TPUEstimators to log the resulted embeddings out and then select the triplets.

In the three tables below 5.1 5.2 5.3 we can see the results from the various models. What is interesting is that the loss is almost the same across the various models but it is still quite high. But the average class similarity is also quite high. This might suggest that all the embeddings are in one big blob or that the similarity metric is not suitable for the batch evaluation.

But when we look at the accuracies of the models we notice something odd. With the exception of two models none of the other models have achieved accuracy of 1%. The two exceptions achieved accuracy of 34.42% and 98.47%. These results are quite unexpected. The reason for this may be that the other models had a learning rate of 0.165¹ which is quite high, but since I was using RMSProp which is scaling the learning rate and I thought this should not be

¹This learning rate is taken directly from official Tensorflow Inception V3 model.

an issue. This high learning rate caused the weights to change a lot and thus not allowing it to settle in a local minimum. The accuracies of models which used the lifted struct loss are surprisingly low, despite its superiority. The best model for “style representation” is thus Inception V3 with learning rate of 0.001 and triplet loss.

I also ran HDBSCAN on the test embeddings to see the number of clusters and the result was which I received 459 clusters which is equal to the number of comics. In addition this proves that the model works efficiently when separating the comics. This result was achieved with default HDBSCAN parameters, changing them decreased the number of cluster. When trying to find global clusters I achieved 129 clusters by changing the parameters. Achieving lower number of clusters would not be possible without most of the points being classified as noise.

Model variant	final loss	Average class similarity	accuracy
A	1.9962643	0.9991348	0.22
B	1.9981924	0.9951266	0.22
C	1.9966104	0.9993636	0.22
D	1.9999185	0.9945911	34.42%
E	1.9997084	0.9741333	98.47%
F	21.00456	0.99058074	0
G	1.997663	0.99813783	0.22
H	1.9983609	0.99759775	0
I	1.9967259	0.9957692	0.22
J	1.9985619	0.9985321	0
K	1.9947969	0.9964487	0.22

Table 5.1: Results from Inception V3 experiment. A is exactly the same as Tensorflow official. All of the other models are same as A but with small changes. B uses three instead of two samples per class. C uses four samples instead of the two per class. D has learning rate of 0,01. E has learning rate of 0,001. F uses lifted struct loss instead of the triplet loss. G uses average of the axillary lost and the final lost instead of sum of these losses. H uses only the final loss. I has margin of 0,2 instead of 1, this margin is used in triplet loss. J uses margin of 1,8 instead of 1. K is the same as A but it outputs 128 dimensional embedding and also is not using pretrained model.

On picture 5.1 we can see that each of the small blobs consist of multiple points. And it is clear that the embedding nicely divides the comics into separate groups. You can also notice some overlap among the blobs. There is one group of points which are slightly farther apart and the group is also far from other groups. This comic is the only comic that uses stick figures and

Model variant	final loss	Average class similarity	accuracy
Inception V4	1.9990538	0.9990527	0.22%

Table 5.2: Results from Inception V4 experiments. This model uses the same hyperparameters and settings as model A in Inception V3

Model variant	final loss	Average class similarity	accuracy
Triplet loss	1.9991846	0.99745023	0.22%
lifted struct loss	20.782486	0.9915535	0%

Table 5.3: Results from Inception-ResNet V2 experiments. These models use the same hyperparameters as model A in Inception V3 but with different loss functions.

has really unique art style. I cannot analyze to resulted embedding fully as I am not an art expert. To make analysis of the similarities of the embedding a full knowledge of style is necessary.

In the picture 5.2 we can see the results of clustering on the reduced centroids created from train data. It is clear that Markov clustering found more clusters which seems to be a more reasonable result then what HDBSCAN found, because there are more then two styles of comics incorporated. Running the clustering on the non-reduced centroids resulted with Markov clustering giving only one cluster. To the contrary, HDBSCAN finds one cluster and the other points are classified as noise.

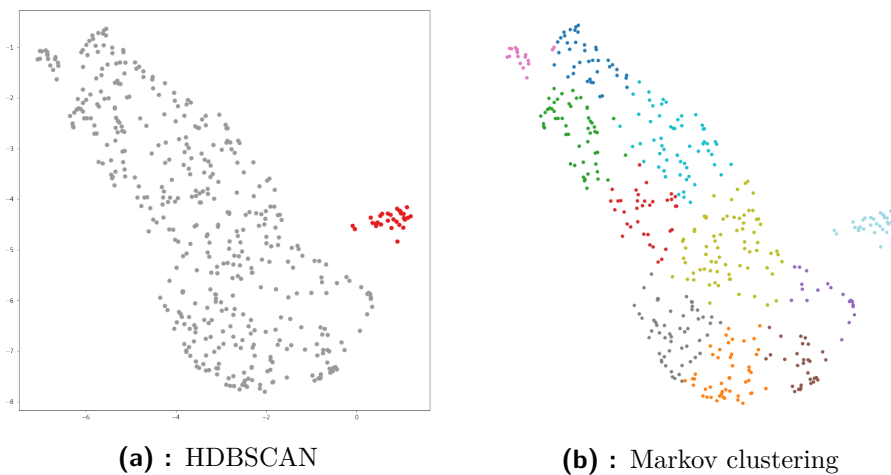


Figure 5.2: Results of clustering reducec train data centroids.

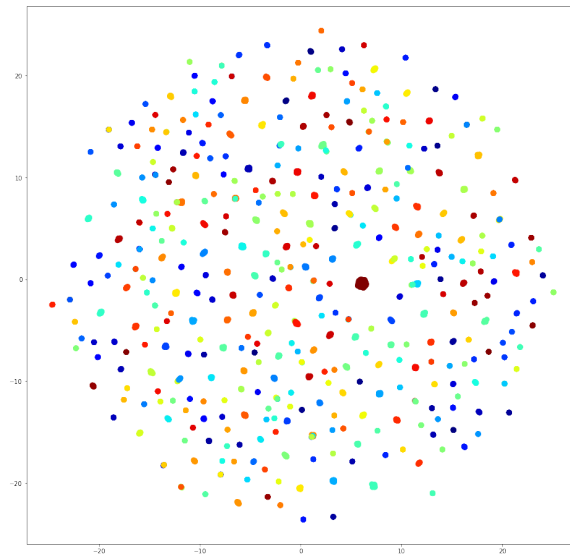


Figure 5.1: 2D representation of sample from test data.

5.2 Doc2Vec

The Doc2Vec experiments are quite straightforward. I trained five different variants of this model and evaluated their accuracy again with linear classifier. When we look at the results in table 5.4 we can see that PV-DBOW achieved the best accuracy with the English text and PV-DM was most successful with Japanese text with concatenation. I am going to assume that in production environment there are going to be multiple languages and it would be inconvenient to train different Doc2Vec model for each of the languages. Therefore I am going to select PV-DM with concatenation as the best model. It is important to note that this model also achieved the second best accuracy with the English text.

Since I am combining two separate models I cannot just merge the embeddings into one embedding space. To combat this, I added one more dimension to the embedding, this dimension represents the language. With this we can easily distinguish between the languages as the current model can not learn from multilingual data.

In the picture 5.3 you can see the two texts embeddings in 2-dimensions. The new feature added to the embeddings separated the embeddings. The texts are long and I do not have the knowledge to fully compare their similarities. Running HDBSCAN gave me two clusters.

Model	Eng Eval Time	Eng Accuracy	Jap Eval Time	Jap Accuracy
A	28.6 s	98.4%	4.87	66.06%
B	39.8 s	85.0%	7.91	56.90%
B'	88 s	92,1%	15.8	74.31%
A + B	73 s	81.9%	13.2	57.80%
A + B'	118 s	96.0%	20.5	63.30%

Table 5.4: Results from Doc2Vec experiments. A is PV-DBOW with 100 dimensional embedding. B is PV-DM with averaging and 100 dimensional embedding. B' is PV-DM with concatenation and 100 100 dimensional embedding

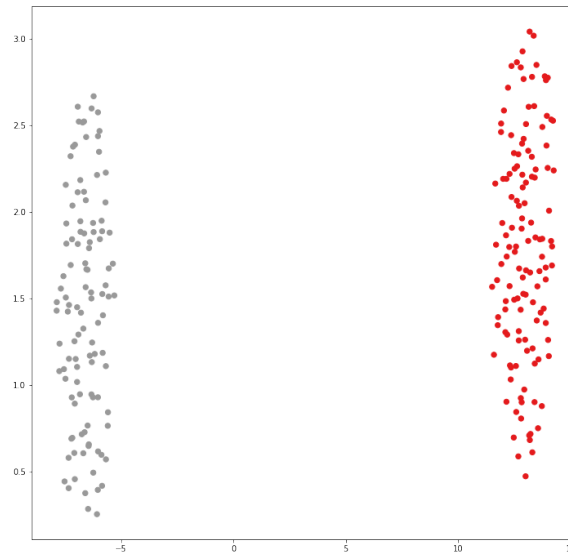
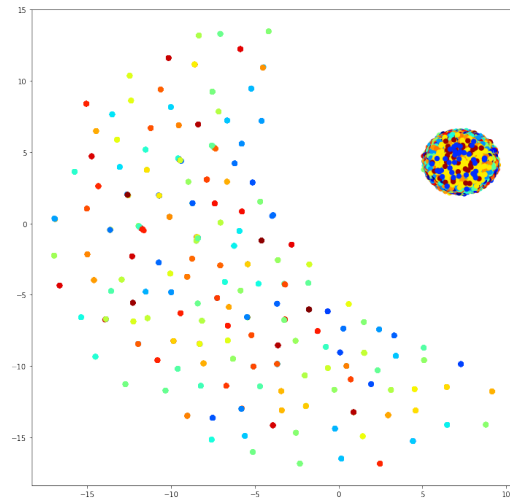


Figure 5.3: 2D representation of the textual evaluation data embeddings.

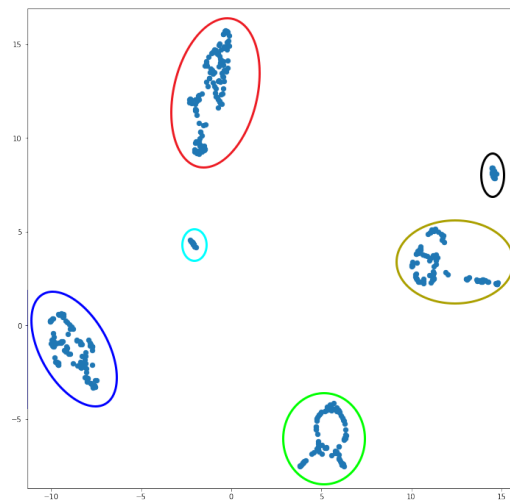
5.3 Combined

The final embedding is created from both the textual and the stylistic representation. First, I tried combining the embeddings as they are but that resulted in sparse embeddings. When I used UMAP on this sparse embedding the results of this operation can be seen on picture 5.4a. The big group (blob) are the comics with text and the surrounding mini-groups are the comics with no text. When I tried displaying only the comics with text the resulted reduced embedding looked like one big group (blob). This tells us that using these big sparse vectors is not the convenient option. Instead we should treat these two embeddings as two features that can be extracted from comics.

If we concatenate the reduced version of these two embeddings we are left with much better results. UMAP reduction to 2-dimensions can be seen on picture 5.4b where I marked clusters with coloured ellipses. The dark blue



(a) : Full Embeddings



(b) : Reduced Embeddings

Figure 5.4: 2D representation of the combined train data embeddings.

ellipse represents a cluster of Mangas. The red ellipse is a cluster of Comics with text. And the other clusters are comics or illustrations with no text. The black ellipse is a cluster of comics with no text that used very little or no colour. The green ellipse represents cluster of comics or illustrations where the illustrations are flat, which means that flat colouring technique was used. I was unable to see the splitting features of the other two clusters as I am no artist and I cannot pick up the small differences.

Chapter 6

Future improvements

One of the main issues of this work is the dataset of comic panels. Creation of the proposed dataset would be necessary to measure the real accuracy as proposed method seems to be sufficient in creating comic embeddings. Another way to improve this work would be to select the best model with its hyperparameters. Since i was limited by 30 days in which i was able to use the TPUs i could only train limited number of models and test limited number of hyperparameters. From the result section it looks like the learning rate of 0,001 is one of these hyperparameters. Different loss function should also be explored. I have explored only three models in this work, where there are plenty of other models to choose from.

Another possible improvement would be to create multilingual embeddings of the text since in the time of creation of this paper this functionality was not implemented in the Doc2Vec which was used for this task¹. To further improve the factual accuracy of the text embedding a better dataset would also be needed. For the dataset creation one would need to create text extraction algorithm which would extract text from speech bubbles in correct order.

To further improve the final embedding additional features could be added like colour histograms or the most dominant colour.

The accuracy measures used in this paper are not sufficient enough as they do not measure the stylistic similarity of the art or contextual similarity of the text. To measure these similarities experts in both fields would be needed to annotate the comic and text similarities.

¹This feature is planned with the next version of Doc2Vec.

Chapter 7

Conclusion

In this paper I provide a detailed explanation and description of comics. I also explore papers which are related to comics and to style. Then I describe different methods and algorithms that will be used in the process of creating or evaluating the resulted embeddings. I describe and propose two datasets that would be needed to create such embeddings. Simpler versions of these datasets were created but they are still sufficient enough to create such embeddings. I provide an efficient method for creating embeddings of comics. The method being Inception V3 with triplet loss for style representation and Doc2Vec for text representation. I also show the levels of accuracies of these methods and provide visualisations of the embedding spaces. The two embeddings, style and text, achieve accuracy of over 98% and 83% respectively. From the results it is rather clear that these embeddings should not be used concatenated as they make sparse vectors. Instead they should be reduced to fewer dimensions and then they can be used either separately as features or concatenated. I also show how the non-reduced and reduced embeddings look in 2-dimensional space. These reduced embeddings were made with a manifold learning algorithm called UMAP. I discuss possible improvements that would advance the results and fully explore this topic.

I explore two other possible methods for comic embeddings. Them being colour histograms and Histograms of oriented gradients. Both these prove to be worse methods then the one proposed.

The metrics that were used to measure the accuracies of these methods are not supposed to be perfect but should be sufficient enough that they showcase the strengths or weaknesses of these methods.

These embeddings solve the issue of expensive calculations of Collaborative recommender systems by providing efficient and compact feature vectors for comics. These embeddings can be used in Content-based recommender systems. The provided recommendations can be easily explained with similarities between these features.



Bibliography

1. KUNZLE, David. *The early comic strip: narrative strips and picture stories in the European broadsheet from c.1450 to 1825*. Berkeley: University of California Press, 1973. ISBN 05-200-1865-6.
2. MCCLOUD, Scott. *Jak rozumět komiksu*. Praha: BB/art, 2008. ISBN 978-80-7381-419-9.
3. HAYMAN, Greg; PRATT, Henry John. What Are Comics? In: *Aesthetics: a reader in philosophy of the arts*. Upper Saddle River, N.J.: Pearson Education Inc., 2005, pp. 419–424. ISBN 0134375912.
4. *Interview with Evka Šimková* [personal communication]. Prague, 2019. Five years of professional work with comics.
5. *Interview with Štěpánka Jislová* [personal communication]. Prague, 2019. Seven years of professional work with comics, leader of workshops, Ladies do Comics, author of 'Klášter nejsvětějšího srdce' and 'Češi 1938: Jak Beneš ustoupil Hitlerovi'.
6. MESKIN, AARON. Defining Comics? *Journal of Aesthetics and Art Criticism*. 2007, vol. 65, no. 4, pp. 369–379. ISSN 0021-8529. Available from DOI: 10.1111/j.1540-594X.2007.00270.x.
7. KUNZLE, David. *Comic strip* [online]. Encyclopædia Britannica, inc., 2017 [visited on 2019-03-24]. Available from: <https://www.britannica.com/art/comic-strip>.
8. MURRAY, Christopher. *Graphic novel* [online]. Encyclopædia Britannica, inc., 2007 [visited on 2019-03-24]. Available from: <https://www.britannica.com/art/graphic-novel>.
9. MAZUR, Dan; DANNER, Alexander. *Komiks: od roku 1968 do současnosti*. Praha: Knižní klub, 2015. ISBN 978-80-242-4856-1.
10. JOHNSON, Jamahl. *The amazing stylistic history of comic books* [online]. 2017 [visited on 2019-03-24]. Available from: <https://99designs.com/blog/design-history-movements/history-of-comic-book-styles/>.

11. THORN, Matt. *History Of Manga* [online]. 2005 [visited on 2019-03-24]. Available from: <http://www.tapanime.com/info/historymanga.php>.
12. KOŘÍNEK, Pavel; FORET, Martin; JAREŠ, Michal. *V panelech a bublinách: kapitoly z teorie komiksu*. Praha: Akropolis, 2015. ISBN 978-80-7470-113-9.
13. GROENSTEEN, Thierry. *Stavba komiksu*. Brno: Host, 2005. ISBN 80-729-4141-0.
14. MCARDLE, Thaneeya. *Explore art styles* [online] [visited on 2019-04-30]. Available from: <https://www.art-is-fun.com/art-styles>.
15. *Artistic Elements* [online] [visited on 2019-04-30]. Available from: http://www2.nkfust.edu.tw/~emchen/CLit/picturebook_artistic_elements.htm.
16. RÉMI, Georges Prosper. [Tintin on bicycle]. In: *The Adventures of Tintin: The Blue Lotus*. France: Casterman, 1946.
17. FRANQUIN, André; MESMAEKER, Jean De. [Gaston LaGaffe]. In: *Spirou: Gaston*. Belgium: Dupuis, 1968.
18. WICKLINE, Dan. *Breaking Down The Types Of Comic Coloring* [online] [visited on 2019-04-30]. Available from: <https://www.bleedingcool.com/2017/01/08/breaking-types-comic-coloring/>.
19. *Lesson four: Shading techniques* [online] [visited on 2019-04-30]. Available from: <https://drawing-with-lena.tumblr.com/post/165724936620/lesson-four-shading-techniques-in-lesson-three>.
20. *Shading* [online] [visited on 2019-04-30]. Available from: <http://www.creativecomicart.com/shading.html>.
21. SCHROFF, Florian; KALENICHENKO, Dmitry; PHILBIN, James. FaceNet: A Unified Embedding for Face Recognition and Clustering. 2015. Available from DOI: 10.1109/CVPR.2015.7298682.
22. SAITO, Masaki; MATSUI, Yusuke. Illustration2Vec: A Semantic Vector Representation of Illustrations. 2017.
23. XUN, Chai Jia; RAMESH, Haritha; YEO, Justin. Are Anime Cartoons? 2016.
24. IYYER, Mohit; MANJUNATHA, Varun; GUHA, Anupam; VYAS, Yogarshi; BOYD-GRABER, Jordan; III, Hal Daumé; DAVIS, Larry. *The Amazing Mysteries of the Gutter: Drawing Inferences Between Panels in Comic Book Narratives*. 2016. Available from arXiv: 1611.05118.
25. GATYS, Leon A.; ECKER, Alexander S.; BETHGE, Matthias. *A Neural Algorithm of Artistic Style*. 2015. Available from arXiv: 1508.06576.

26. KARAYEV, Sergey; TRENTACOSTE, Matthew; HAN, Helen; AGARWALA, Aseem; DARRELL, Trevor; HERTZMANN, Aaron; WINNEMOELLER, Holger. *Recognizing Image Style*. 2013. Available from arXiv: 1311.3715.
27. BARKAN, Oren; KOENIGSTEIN, Noam. *Item2Vec: Neural Item Embedding for Collaborative Filtering*. 2016. Available from arXiv: 1603.04259.
28. GARCIA-GASULLA, D.; BÉJAR, J.; CORTÉS, U.; AYGUADÉ, E.; LABARTA, J.; SUZUMURA, T.; CHEN, R. *A Visual Embedding for the Unsupervised Extraction of Abstract Semantics*. 2015. Available from arXiv: 1507.08818.
29. LAUBROCK, Jochen; DUBRAY, David. *CNN-based Classification of Illustrator Style in Graphic Novels: Which Features Contribute Most?* [EasyChair Preprint no. 557]. 2018. Available from DOI: 10.29007/z3f1.
30. MIKOLOV, Tomas; CHEN, Kai; CORRADO, Greg; DEAN, Jeffrey. *Efficient Estimation of Word Representations in Vector Space*. 2013. Available from arXiv: 1301.3781.
31. LE, Quoc V.; MIKOLOV, Tomas. *Distributed Representations of Sentences and Documents*. 2014. Available from arXiv: 1405.4053.
32. *Interview with Václav Roháč* [personal communication]. Prague, 2019. Illustrator of fifteen books. Author of many comics. Author of 'Pražáček Ten' a political comics which sold more than two million copies.
33. *The Graphic Narrative Corpus* [online] [visited on 2019-05-01]. Available from: <https://groups.uni-paderborn.de/graphic-literature/gncorpus/corpus.php>.
34. OGAWA, Toru; OTSUBO, Atsushi; NARITA, Rei; MATSUI, Yusuke; YAMASAKI, Toshihiko; AIZAWA, Kiyoharu. *Object Detection for Comics using Manga109 Annotations*. 2018. Available from arXiv: 1803.08670.
35. MATSUI, Yusuke; ITO, Kota; ARAMAKI, Yuji; YAMASAKI, Toshihiko; AIZAWA, Kiyoharu. *Sketch-based Manga Retrieval using Manga109 Dataset* [Multimedia Tools and Applications, Volume 76, Issue 20, 2017]. 2015. Available from DOI: 10.1007/s11042-016-4020-z.
36. GUÉRIN, Clément et al. eBDtheque: a representative database of comics. In: *Proceedings of the 12th International Conference on Document Analysis and Recognition (ICDAR)*. 2013.
37. *Line Webtoon* [online] [visited on 2019-05-01]. Available from: <https://www.webtoons.com/en/>.
38. *Illustration* [online] [visited on 2019-05-01]. Available from: <https://www.illustrationweb.com/>.

39. RAJ, Bharath. *Data Augmentation: How to use Deep Learning when you have Limited Data?—?Part 2* [online] [visited on 2019-04-30]. Available from: <https://medium.com/nanonets/how-to-use-deep-learning-when-you-have-limited-data-part-2-data-augmentation-c26971dc8ced>.
40. ZAWORSKI, Radek. *Data augmentation techniques and pitfalls for small datasets* [online] [visited on 2019-04-30]. Available from: <https://snow.dog/blog/data-augmentation-for-small-datasets>.
41. GARBADE, Michael. *Understanding K-means Clustering in Machine Learning* [online] [visited on 2019-04-30]. Available from: <https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1>.
42. TREVINO, Andrea. *Introduction to K-means Clustering* [online] [visited on 2019-04-30]. Available from: <https://www.datascience.com/blog/k-means-clustering>.
43. *Clustering: K-means* [online] [visited on 2019-04-30]. Available from: <https://scikit-learn.org/stable/modules/clustering.html%5C#k-means>.
44. ELDRIDGE, Justin; BELKIN, Mikhail; WANG, Yusu. *Beyond Hartigan Consistency: Merge Distortion Metric for Hierarchical Clustering*. 2015. Available from arXiv: 1506.06422.
45. MCINNIS, Leland; HEALY, John; ASTELS, Steve. hdbscan: Hierarchical density based clustering. *The Journal of Open Source Software*. 2017, vol. 2, no. 11, pp. 205.
46. MCINNIS, Leland; HEALY, John. Accelerated Hierarchical Density Based Clustering. In: *Data Mining Workshops (ICDMW), 2017 IEEE International Conference on*. 2017, pp. 33–42.
47. *MCL - a cluster algorithm for graphs* [online] [visited on 2019-04-30]. Available from: <https://micans.org/mcl/index.html>.
48. DONGEN, Stijn van. *Graph Clustering by Flow Simulation*. Netherlands, 2000. Dissertation. Utrecht University.
49. SHLENS, Jonathon. *A Tutorial on Principal Component Analysis*. 2014. Available from arXiv: 1404.1100.
50. KRUSKAL, Joseph B. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. In: 1964.

51. TENENBAUM, Joshua B.; SILVA, Vin de; LANGFORD, John C. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science*. 2000, vol. 290, no. 5500, pp. 2319–2323. ISSN 0036-8075. Available from DOI: [10.1126/science.290.5500.2319](https://doi.org/10.1126/science.290.5500.2319).
52. WATTENBERG, Martin; VIÉGAS, Fernanda; JOHNSON, Ian. How to Use t-SNE Effectively. *Distill*. 2016. Available from DOI: [10.23915/distill.00002](https://doi.org/10.23915/distill.00002).
53. MAATEN, Laurens van der; HINTON, Geoffrey. Visualizing Data using t-SNE. *Journal of Machine Learning Research*. 2008, vol. 9, pp. 2579–2605. Available also from: <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
54. *Visualizing Data Using t-SNE* [online] [visited on 2019-05-01]. Available from: <https://www.youtube.com/watch?v=RJVL80Gg31A>.
55. MCINNES, Leland; HEALY, John; MELVILLE, James. *UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction*. 2018. Available from arXiv: [1802.03426](https://arxiv.org/abs/1802.03426).
56. *UMAP Uniform Manifold Approximation and Projection for Dimension Reduction* [online] [visited on 2019-05-01]. Available from: <https://www.youtube.com/watch?v=nq6iPZVUxZU>.
57. *Performance Comparison of Dimension Reduction Implementations* [online] [visited on 2019-05-01]. Available from: <https://umap-learn.readthedocs.io/en/latest/benchmarking.html>.
58. DALAL, Navneet; TRIGGS, Bill. Histograms of oriented gradients for human detection. In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. 2005, vol. 1, pp. 886–893.
59. MALLICK, Satya. *Histogram of Oriented Gradients* [online] [visited on 2019-05-01]. Available from: <https://www.learnopencv.com/histogram-of-oriented-gradients/>.
60. KARPATY, Andrej. *CS231n: Convolutional Neural Networks for Visual Recognition* [online] [visited on 2019-05-01]. Available from: <https://cs231n.github.io/>.
61. KARN, Ujjwal. *An Intuitive Explanation of Convolutional Neural Networks* [online] [visited on 2019-05-01]. Available from: <https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/>.
62. SZEGEDY, Christian; VANHOUCKE, Vincent; IOFFE, Sergey; SHLENS, Jonathon; WOJNA, Zbigniew. *Rethinking the Inception Architecture for Computer Vision*. 2015. Available from arXiv: [1512.00567](https://arxiv.org/abs/1512.00567).

63. RAJ, Bharath. *A Simple Guide to the Versions of the Inception Network* [online]. 2018 [visited on 2019-05-01]. Available from: <https://towardsdatascience.com/a-simple-guide-to-the-versions-of-the-inception-network-7fc52b863202>.
64. CANZIANI, Alfredo; PASZKE, Adam; CULURCIELLO, Eugenio. *An Analysis of Deep Neural Network Models for Practical Applications*. 2016. Available from arXiv: 1605.07678.
65. BUSHAEV, Vitaly. *Understanding RMSprop?—?faster neural network learning* [online]. 2018 [visited on 2019-05-01]. Available from: <https://towardsdatascience.com/understanding-rmsprop-faster-neural-network-learning-62e116fcf29a>.
66. SCHROFF, Florian; KALENICHENKO, Dmitry; PHILBIN, James. *FaceNet: A Unified Embedding for Face Recognition and Clustering*. 2015. Available from DOI: 10.1109/CVPR.2015.7298682.
67. SONG, Hyun Oh; XIANG, Yu; JEGELKA, Stefanie; SAVARESE, Silvio. *Deep Metric Learning via Lifted Structured Feature Embedding*. 2015. Available from arXiv: 1511.06452.
68. ACERBI, Alberto. *imsdb movie scripts*. OSF, 2019. Available from DOI: 10.17605/OSF.IO/ZYTMP.

Appendix A

Contents of CD

Additional data could be requested from Ing. Karel Klouda Ph.D.

readme.txt	the file with CD contents description
data	the data files directory
models	the directory with the best models
inception_v3.pb	model for TF serving
inception_v3.ckpt	inception v3 checkpoint
doc2vec	Doc2Vec saved model
results	the directory with resulted embeddings
style_train.json	style embeddings of train data
style_eval.json	style embeddings of eval data
text_train.json	text embeddings of train data
text_eval.json	text embeddings of eval data
src	the directory of source codes
Comic2Vec	the directory of scripts
*.py	python scripts used in this paper
*.ipynb	python notebooks used in this paper
thesis	the directory of L ^A T _E X source codes of the thesis
figures	the thesis figures directory
*.tex	the L ^A T _E X source code files of the thesis
text	the thesis text directory
Comic2Vec.pdf	the Diploma thesis in PDF format
Comic2Vec.ps	the Diploma thesis in PS format