

A First Dataset for Film Age Appropriateness Investigation

Emad Mohamed, Le An Ha

RGCL, RILP, University of Wolverhampton
Wulfruna Street, Wolverhampton, WV1 1LY, UK
E.Mohamed2@wlv.ac.uk, L.A.Ha@wlv.ac.uk

Abstract

Film age appropriateness classification is an important problem with a significant societal impact that has so far been out of the interest of Natural Language Processing and Machine Learning researchers. To this end, we have collected a corpus of 17000 film transcripts along with their age ratings. We use the textual contents in an experiment to predict the correct age classification for the United States (G, PG, PG-13, R and NC-17) and the United Kingdom (U, PG, 12A, 15, 18 and R18). Our experiments indicate that gradient boosting machines beat FastText and various Deep Learning architectures. We reach an overall accuracy of 79.3% for the US ratings compared to a projected super human accuracy of 84%. For the UK ratings, we reach an overall accuracy of 65.3% (UK) compared to a projected super human accuracy of 80.0%.

Keywords: Machine Learning, Deep Learning, XGBoost, Age Appropriateness, Movies

1. Introduction

Interest in film from a computational linguistics perspective has been massive. Several studies in CL have examined the genre in terms of Sentiment Analysis (Phan and Matsumoto, 2018), turn-taking (Banchs, 2012), and many other things. Most of these studies have focused on film reviews (from Amazon.com and IMDB, but the actual film content (the script, audio, and video for example) has not received as much interest in spite of the potential availability of huge datasets as will be demonstrated below.

In this study, we investigate a completely new problem with a novel corpus. We investigate whether we can use Machine Learning and Artificial Neural Networks to predict the age classification accompanying the films based on their textual content. For this purpose, we collect a custom corpus from the internet and complement it with age rating certificates from the Internet Movie Database (IMDB). The data combination makes a dataset that can be used for many purposes part of which is the focus of this paper: *age appropriateness classification*.

Entities like the Motion Picture Association of America (MPAA)¹ and the British Board of Film Classification (BBFC)² issue film ratings that determine the age appropriateness of each film based on the film’s content. The latter define their classification as ‘the process of giving age ratings and content advice to films and other audiovisual content to help children and families choose what’s right for them and avoid what’s not.’ This is a human labour intensive endeavor that requires at least two Compliance Officers to watch each film and report on it. So far, there does not seem to be enough interest from the Computational Linguistics community in the problem of automatic age appropriateness classification, possibly due to the lack of resources.

We have collected a fairly large dataset for the purpose and

we have run experiments to test whether we can predict the age rating certificates based on the textual content of the film (i.e. scripts). We have used state-of-the-art classification methods including neural networks and gradient boosting machines (GBM’s). The best results were obtained using the XGBoost implementation of GBM’s. For the USA and the UK, the accuracies of the prediction reach 79.3%, and 65.3% respectively compared to two projected ceilings of 84% and 80%. The rest of this paper goes as follows: in Section 2, we describe the data and the methods. In Section 3, we present the results and perform some analysis. In Section 4, we discuss related datasets and experiments. Finally, the conclusion discusses our projected future research that seeks to combine textual and visual inputs to tackle the problem.

2. Data & Methods

2.1. Data Collection

We have collected a large number of film scripts from <https://www.springfieldspringfield.co.uk/>. We have matched these with their age certificates using IMDbPY³, a community-based API to access IMDB data. The films were identified mainly through combinations of IMDB film ID’s, directors, actors, and the production companies. Where ambiguity could not be automatically resolved, we simply dropped the film from our dataset, which left us with 17018 unambiguously age-rated film transcripts, where each transcript contains only the dialogue. Meta information, like scene settings and description, was only available for a small minority of the films and has thus not been utilised in the current paper although we believe it may be useful for our future research.

2.2. Data Cleaning

Movie dialogue is designed to match screens and is thus not made up of valid linguistic entities, i.e. sentences. The following extract is from “Terminator Salvation”:

¹<https://www.filmratings.com/>

²<https://bbfc.co.uk/>

³<https://imdbpy.readthedocs.io/en/latest/index.html>

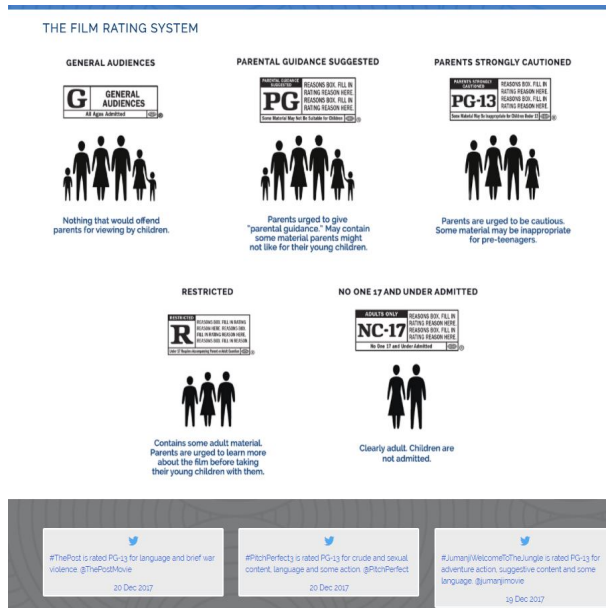


Figure 1: MPAA age classification (left), and BBFC age classification (right)

Our intel has located a hidden signal under the primary channel. It allows for direct control on the machines. Skynet's a machine, and like all machines, it has an off-switch.

For our modeling purposes, we have used the Spacy NLP library (Honribal and Montani, 2017) to tokenize the script into sentences. The Spacy sentence segmentation uses dependency parsing. We compared Spacy with the NLTK (Bird et al., 2009) in terms of sentence segmentation. We have tested both tools on a sample sample of film scenes, and we have found Spacy to be more accurate. While most of the models we use do not have the concept of a sentence, this conversion is necessary to allow for various linguistic contexts (in terms of ngrams) to be included. Following the sentence tokenization, the previous excerpt from the *Terminator Salvation* becomes:

Our intel has located a hidden signal under the primary channel. It allows for direct control on the machines. Skynet's a machine, and like all machines, it has an off-switch.

2.3. Annotation

We do not perform any manual annotation on the data. We, instead, use distant annotation, i.e. metadata that can be considered annotation though not originally intended as such. We have noticed, for example, that there is a striking similarity between the process of assigning a film to an age

category and that of linguistic annotation used in computational linguistics research. To quote the British censors (BBFC) ⁴:

Films for cinema release are usually seen by at least two of our Compliance Officers, and in most cases, their age rating recommendation is approved by the Compliance Manager or the Head of Compliance. If Compliance Officers are in any doubt, if a film is on the borderline between two categories, or if important policy issues are involved, it may be seen by other members of the BBFC, up to and including the Chief Executive, the President and Vice Presidents. Occasionally, we may also call for expert advice about the legal acceptability of film content or its potential for harm. DVDs and VoD films and series are normally seen by one Compliance Officer, but opinions from other Officers, the Compliance Manager, the Head of Compliance and Board of Classification may be required for more difficult content.

This matches our experience of linguistic annotation on various projects with which we were involved. Figure 1. contrasts the classification schemes used in the USA and the UK.

We use the Internet Movies Database (IMDB) ⁵ as the source of age ratings. For each film on IMDB, there is a section on *Certificates*, which lists age ratings from several countries in addition to the country of origin. For example, in the case of *Terminator Salvation*, Figure 2 shows the various classifications of the film.

⁴<https://bbfc.co.uk/what-classification>

⁵<https://imdb.com>

MPAA	Rated PG-13 for intense sequences of sci-fi violence and action, and language
Certification	Argentina:13 Australia:M Austria:14 Brazil:14 Canada:14A (Canadian Home Video rating) Canada:13+ (Quebec) Canada:14+ (TV rating) Denmark:15 Finland:K-15 France:Tous publics Germany:16 Germany:12 (TV: cut version) Hong Kong:IIB Iceland:12 India:UA Ireland:12A Italy:T Japan:G Malaysia:P13 Mexico:B Netherlands:12 New Zealand:M Norway:15 (2009, cinema rating) Peru:14 Philippines:PG-13 (MTRCB) Poland:12 Portugal:M/12 Russia:12+ Singapore:PG Singapore:PG13 (TV rating) Singapore:NC-16 (director's cut) South Africa:13 South Korea:15 Spain:13 Sweden:15 Switzerland:14 (canton of Geneva) Switzerland:14 (canton of Vaud) Taiwan:PG-12 United Kingdom:12A United Kingdom:12 (director's cut) United States:TV-14 United States:PG-13 (certificate #45308) United States:R (certificate #45600, director's cut)

Figure 2: Terminator Salvation age certificates from various countries, source: IMDB

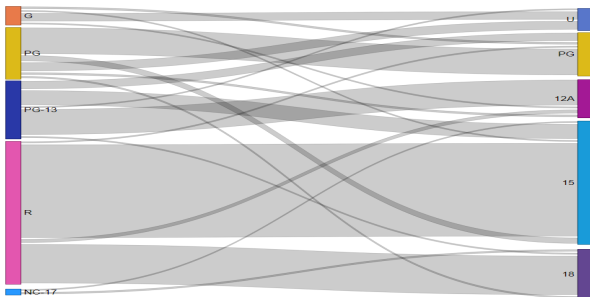


Figure 3: Sankey map between USA and UK rating systems

2.4. Dataset Description and Statistics

In total, we have 17018 titles, 181M words, with an average of 10651 words per film.

We focus on the theatre ratings rather than TV or video ratings. For the USA, they are G, PG, PG-13, R, and NC-17. For the UK ratings, they are U, PG, 12⁶, 15, 18, and R18. In total, we have these US ratings for 8923 titles, and the UK ratings for 10920 titles. These two sets are overlapped by 7068 titles. It is worth noting that there is a one to many mapping between films and certificates. A single film could have several certificates within and across countries. We use the main certificates provided by the IMDB based on the country of origin. For example, in Figure 2, while “Terminator Salvation” has several USA ratings for two different cuts, we use the PG-13 one provided on top. It is also of note that there is not a one to one mapping between the UK and the USA ratings. Figure 3 shows possible mapping between the two rating systems found in our dataset.

The statistics for these classes are shown in Table 1. We have divided the data into a train section (70%), a dev sec-

⁶In the UK, there are two 12 certificates: 12 and 12A, with the latter reserved for TV ratings. In our datasets, we have treated both as one and the same thing since they target the exact same group

tion (10%), and a test section (20%). The data was divided using random by-year stratification.

Table 1: Basic statistics about the dataset; TT = Total Number of Texts, NW = Number of Words, AWT = Average Number Words per Texts

		TT	NW	AWT
All		17018	181M	11K
US	G	294	3,107K	11K
	PG	1493	16,867K	11K
	PG-13	2150	25,062K	12K
	R	4965	52,863K	11K
	NC-17	21	234K	11K
	All	8923	98,133K	11K
UK	U	1095	12,799K	12K
	PG	1723	20,397K	12K
	12A	1268	15,719K	12K
	15	5093	53,768K	11K
	18	1740	15,473K	9K
	R18	1	6K	6K
	All	10920	118,432K	11K

2.5. Methods

We use three main methods/tools: (1) FastText, (2) Gradient Boosting Machines with the XGBoost implementation, and (3) Artificial Neural Networks, including Hierarchical Attention (HAtt), Character-based convolutional neural network (CharCNN), ELMo and BERT.

FastText (Joulin et al., 2016) is a document classifier that relies on a language model built on the principle that words can be represented as the sum of subword vectors, which could be useful in representing languages with large vocabularies, e.g. morphologically rich languages.

XGBoost (Chen and Guestrin, 2016) is an efficient implementation of gradient boosting machines that has been shown to be successful in many real life applications, for example (Volkovs et al., 2017; Sandulescu and Chiru,

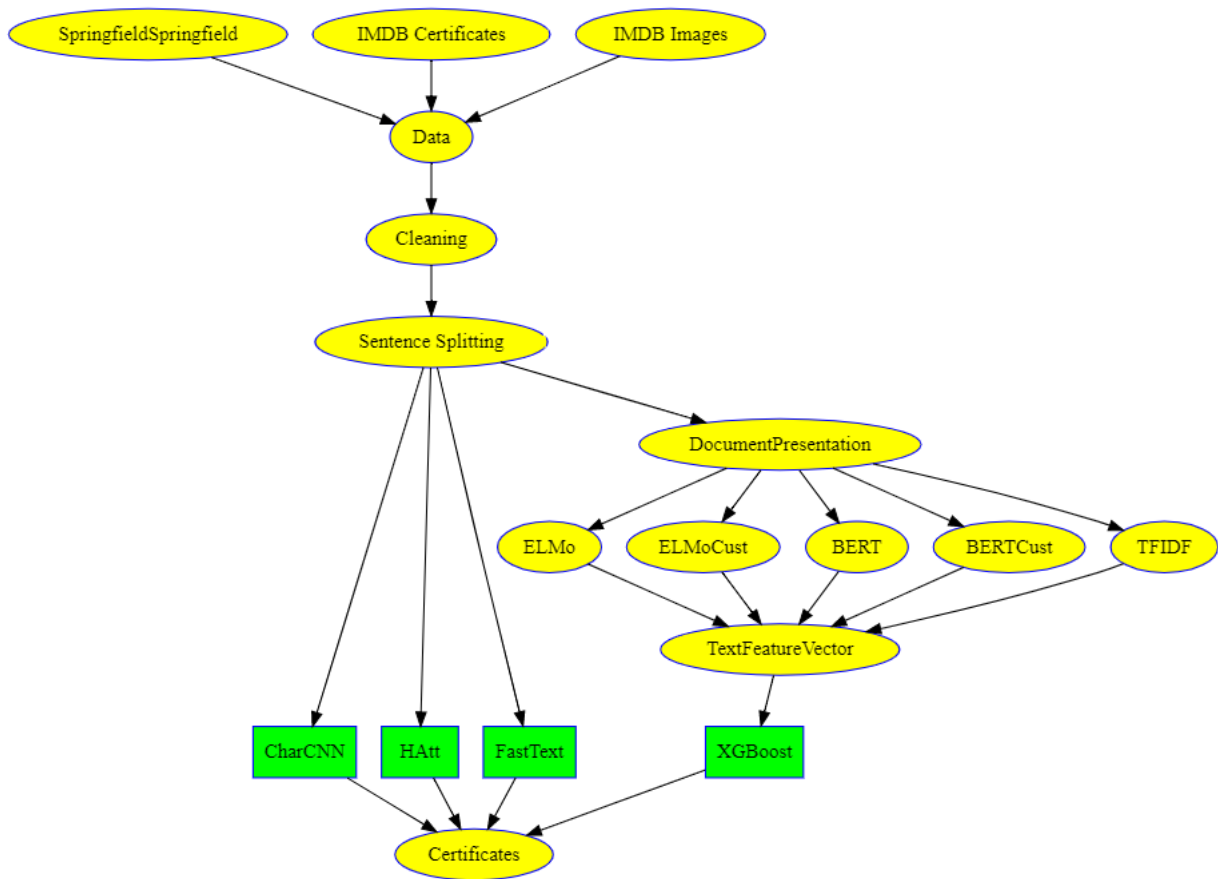


Figure 4: Dataflow of the experiment (green boxes represent machine learning algorithms)

2016).

Hierarchical Attention (Yang et al., 2016), HAtt, is a neural network architecture that represents a text as a sequence of sentences. Each sentence, in turn, is represented as a list of words. The attention mechanism is used first to take advantage of relations among words in a sentence, and then among sentences in a document.

BERT (Devlin et al., 2018) is a context sensitive embedding neural architecture that also uses the attention mechanism extensively, in which the model learns to guess the missing words and whether one sentence follows another. Using what it learns, the model can produce a vector representation of an arbitrary sentence, which can then be used for downstream tasks.

CharCNN (Zhang et al., 2015) is a character-based convolutional network that puts the input text through one or more 1d convolutional layers to construct the representations for the classifier.

ELMo (Peters et al., 2018) is another context-sensitive embedding architecture that makes use of character-based 1d convolutional layer and lstm layers to learn language models.

We use BERT, ELMo, both their generic models and their respective models that have been fine-tuned using our dataset (hereinafter referred to as BERTCust and ELMoCust), and char-based ngram tfidf, as alternatives for feature extractions. For BERT and ELMo, the document representation of the film is the mean pooling of representa-

tions of all sentences in the film transcripts, as calculated by the respective models. For char-based ngram tfidf (abbreviated as TFIDF), the document representation is the tfidf values of all the char-based ngrams presented in the film transcripts. The idf values are calculated using the training portion of the dataset. These document representation models have been shown to be successful in text classification tasks (Cavnar et al., 1994; Devlin et al., 2018; Peters et al., 2018). The document representations (in the form of vectors of various dimensions) will also be fed into XGBoost to learn classification models.

Alternatively, FastText, HAtt, and CharCNN, models that have also been found to be successful in document classification tasks, will be used directly to learn classification models from text. Figure 4 shows our workflow.

For evaluation, we use the standard metrics of accuracy and the Area Under the Curve of the Receiver Operating Characteristic (hereinafter referred to as AUC). As the AUC incorporates the trade-off between precision and recall, we will report them instead of recall, precision, and F1 score. We have two settings for evaluation of accuracy:

- *Strict evaluation* in which a prediction is considered correct only if it exactly matches the gold standard data, and
- *Relaxed Evaluation* where a prediction is considered correct if it is within one class of the correct gold standard class. For example, if the true class is PG-13,

then in a relaxed setting, both a prediction of PG and R would be considered correct.

Table 2 shows how some examples of whether or not the prediction is correct in the two settings.

Country	Prediction	Gold	Strict	Relaxed
US	G	PG	False	True
	G	G	True	True
	PG	G	False	True
	PG	R	False	False
UK	U	U	True	True
	U	PG	False	True
	15	18	False	True
	15	PG	False	False

Table 2: Strict and relaxed decisions for various prediction-gold pairs

While we believe that *strict evaluation* should be considered the only valid measure, we realize that the same film may have several ratings, but we have not seen a film whose rating variance goes beyond two classes, which gives some validity to relaxed evaluation.

For the experiments using XGBoost and FastText, we have run a limited version of grid search based on what we have found in the literature concerning the optimisation of both tools. Based on the dev set, we found the best parameters for XGBoost to be a learning rate of 0.3, a max depth of 3 and with the number of estimators being 300.

For FastText, the best performing experiment, as measured by performance on the dev set is where we use word unigrams, 100 epochs and a learning rate of 0.3. An analysis by regression to examine the effect of the factors in the 10 experiments of FastText shows the number of epochs to be the most important factor and the number of ngrams the least.

2.5.1. Setting a Baseline and projected human performance

We adopt the majority class as our baseline. For the US, the majority class is “R”, whose accuracy on the dev set is 55%. For the UK, the majority class is “15”, and the accuracy on the dev set is 44%. For relaxed evaluation, the majority classes are PG-13 and 15 with accuracies of 96% and 70% respectively.

For the performance ceiling, as there is no data within a single country to indicate the level of inter-person agreement of the ratings, we use a classification model that uses ratings from other countries to predict the ratings of the target country, and use its performance as an indicator of human performance. For example, if we use UK ratings to predict the USA ratings, it would yield a strict accuracy of 80.6% and a relaxed accuracy of 95.1% (SingleCt in Tables 3 and 4). On the other hand, if we use all the available ratings from the 69 remaining countries whose ratings are available on IMDB to predict the USA ratings (OtherCts in Tables 3 and 4), the strict and relaxed accuracies are 84.8% and 96.7% respectively. We consider this as our performance ceiling. The model used in this experiment is also XGBoost.

3. Results & Analysis

Tables 3 and 4 show the experiments we have carried out and the results we obtained. We observe that the deep learning models, with (BERT and ELMo) or without transfer learning (CharCNN and HAtt) do not perform as well as the character ngram based features. We hypothesise that for this task, each film transcript contains enough information for the task; such information may not be available in tasks in which deep learning models excel.

The performances of the classifiers are approaching or surpassing those using ratings from one country to predict those of another country. Around 95% of the predictions are within one rating of the correct ones. Figures 5 and 6 show the ROC curves for individual classes and micro averages for the United States and the United Kingdom certificates predictions.

Inspection of the confusion matrices (Tables 5 and 6) indicates that in the case of the US, the classifier only misclassifies one R title to be a G title, and only one 18 title to be an U in the UK case.

Algorithm	Input	Acc	RelaxAcc	AUC
FastText	text	74.7	95.0	0.945
HAtt	text	69.6	95.8	0.935
CharCNN	text	57.7	88.8	0.870
XGBoost	ELMo	62.8	87.5	0.896
XGBoost	ELMoCust	63.4	88.8	0.904
XGBoost	BERT	62.7	87.2	0.900
XGBoost	BERTCust	63.3	88.6	0.899
XGBoost	TFIDF	79.1	96.2	0.962
XGBoost	SingleCt	80.6	95.1	0.957
XGBoost	OtherCts	84.7	96.7	0.978

Table 3: Results on the USA categories

Algorithm	Input	Acc	RelaxAcc	AUC
FastText	text	61.5	94.6	0.915
HAtt	text	58.9	91.9	0.906
CharCNN	text	41.1	75.6	0.765
XGBoost	ELMo	54.2	86.8	0.872
XGBoost	ELMoCust	54.6	85.4	0.876
XGBoost	BERT	55.6	86.5	0.871
XGBoost	BERTCust	57.5	89.2	0.878
XGBoost	TFIDF	65.3	94.2	0.930
XGBoost	SingleCt	61.8	91.5	0.899
XGBoost	OtherCts	80.0	96.3	0.972

Table 4: Results on the UK categories

		Predicted			
		G	PG	PG-13	R
Gold	G	13	37	4	3
	PG	6	227	60	32
	PG-13	1	65	310	69
	R	1	26	65	863

Table 5: Confusion matrix of the predicted certificates for the USA

		Predicted				
		U	PG	12A	15	18
Gold	U	152	60	4	9	0
	PG	54	217	27	67	0
	12A	6	48	102	114	1
	15	2	34	42	868	46
	18	1	1	1	241	88

Table 6: Confusion matrix of the predicted certificates for the UK

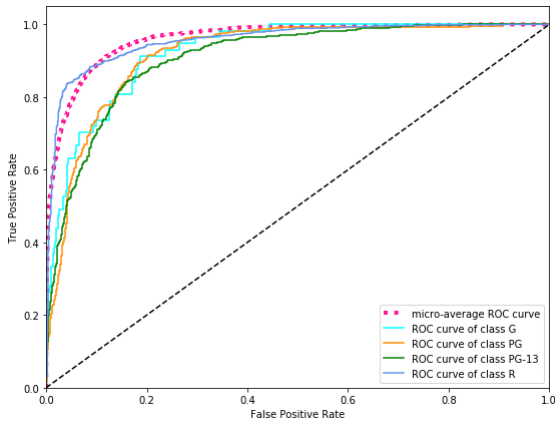


Figure 5: Receiver operating characteristics (ROC) of the XGBoost classifier using char-ngram tfidf for United States certificates predictions

3.1. Error Analysis

Table 7 lists some films whose classification is off by two ratings, for example if a G film is rated as R. While the error dataset is too small for us to make any generalisations, we have noticed that 6 out of the 8 films are from the seventies and the early eighties, which may be an indication that the time factor should be considered in the classification, which we will do in future research. We also hope that the inclusion of video in the classification will help improve the classification.

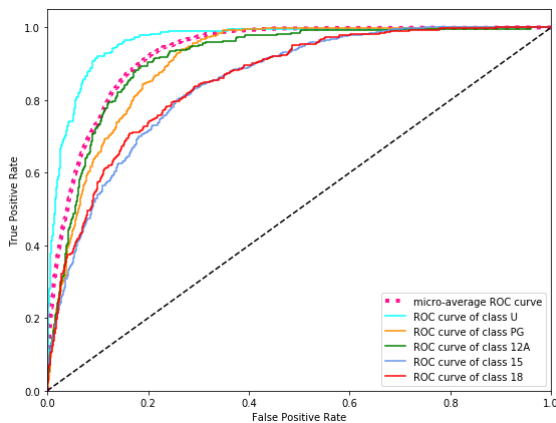


Figure 6: Receiver operating characteristics (ROC) of the XGBoost classifier using char-ngram tfidf for United Kingdom certificates predictions

Country	FilmID	Pred	Actual
UK	tt3078242	U	18
	tt0092048	PG	18
	tt0059578	U	15
	tt0083851	U	15
USA	tt0066434	G	R
	tt3421514	R	G
	tt0065462	R	G
	tt0067065	R	G

Table 7: Examples of films erroneously classified within a distance of more than two ratings

4. Related datasets

There are several datasets that contain movies’ scripts or dialogs. None of them are developed with the focus on age appropriateness. The OpenSubtitles collection of parallel corpora (Lison and Tiedemann, 2016), albeit a very large dataset (17.2G tokens), focuses on bitexts aspects of movies dialogs.

Gorinski and Lapata (Gorinski and Lapata, 2015) build a collection of 1,276 movie scripts, by automatically crawling web-sites which host or link entire movie scripts for the purpose of summarisation. 132,229 dialogues from 753 movies have been collected by (Banchs, 2012). The objective is to study the semantic and pragmatic aspects of human communication within a wide variety of contexts, scenarios, styles and socio-cultural settings. ramakrishna-et al-2015-quantitative ramakrishna-et al-2015-quantitative analyse differences in portrayal of characters in movies with respect to characters’ gender, race, age and other metadata using 945 screenplay files from two primary sources: *IMSDB* and *Daily Scripts*. (Phan and Matsumoto, 2018)’s corpus includes conversations from movie with more than 2.1 millions utterances which are partly annotated for emotions. (Kar et al., 2018a), (Kar et al., 2018b), (Battu et al., 2018) use plot synopses for various purposes. To the best of our knowledge, our dataset is the largest movie content dataset apart from the OpenSubtitles, and the only one with age appropriateness certificates.

5. Conclusion & Future Work

We have collected a dataset for film age appropriateness classification which we believe to be of importance to the NLP, ML and Film Studies communities. We have used traditional and DL algorithms to predict the various categories for the USA and the UK, and we have found XGBoost to be a clear winner for this task, beating more modern architectures.

In our future work, we will move in two directions: (1) examining the characteristics of these categories from a Digital Humanities perspective thus providing more interpretable models and results, and (2) we will integrate videos and images in the classification task, thus creating a more realistic setting for real-world applications. The dataset, including the ids of the titles used in training, developing, and testing, as well as codes can be obtained by contacting the authors.

6. Bibliographical References

- Banchs, R. E. (2012). Movie-DiC: a movie dialogue corpus for research and development. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 203–207, Jeju Island, Korea, July. Association for Computational Linguistics.
- Battu, V., Batchu, V., Gangula, R. R. R., Dakannagari, M. M. K. R., and Mamidi, R. (2018). Predicting the genre and rating of a movie based on its synopsis. In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation*, Hong Kong, 1–3 December. Association for Computational Linguistics.
- Bird, S., Klein, E., and Loper, E. (2009). *Natural Language Processing with Python*. O’Reilly Media.
- Cavnar, W. B., Trenkle, J. M., and Mi, A. A. (1994). N-Gram-Based Text Categorization. In *In Proceedings of SDAIR-94, 3rd Annual Symposium on Document Analysis and Information Retrieval*.
- Chen, T. and Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’16, pages 785–794, New York, NY, USA. ACM.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. (2018). BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.
- Gorinski, P. J. and Lapata, M. (2015). Movie script summarization as graph-based scene extraction. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1066–1076, Denver, Colorado, May–June. Association for Computational Linguistics.
- Honnibal, M. and Montani, I. (2017). spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.
- Joulin, A., Grave, E., Bojanowski, P., and Mikolov, T. (2016). Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.
- Kar, S., Maharjan, S., López-Monroy, A. P., and Solorio, T. (2018a). MPST: A corpus of movie plot synopses with tags. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May. European Language Resources Association (ELRA).
- Kar, S., Maharjan, S., and Solorio, T. (2018b). Folksonomication: Predicting tags for movies from plot synopses using emotion flow encoded neural network. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2879–2891, Santa Fe, New Mexico, USA, August. Association for Computational Linguistics.
- Lison, P. and Tiedemann, J. (2016). OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 923–929, Portorož, Slovenia, May. European Language Resources Association (ELRA).
- Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana, June. Association for Computational Linguistics.
- Phan, D.-A. and Matsumoto, Y. (2018). EMTC: Multilabel corpus in movie domain for emotion analysis in conversational text. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May. European Language Resources Association (ELRA).
- Sandulescu, V. and Chiru, M. (2016). Predicting the future relevance of research institutions - the winning solution of the kdd cup 2016. *ArXiv*, abs/1609.02728.
- Volkovs, M., Yu, G. W., and Poutanen, T. (2017). Content-based neighbor models for cold start in recommender systems. In *Proceedings of the Recommender Systems Challenge 2017, RecSys Challenge ’17*, pages 7:1–7:6, New York, NY, USA. ACM.
- Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., and Hovy, E. (2016). Hierarchical attention networks for document classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1480–1489, San Diego, California, June. Association for Computational Linguistics.
- Zhang, X., Zhao, J., and LeCun, Y. (2015). Character-level convolutional networks for text classification. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS’15*, pages 649–657, Cambridge, MA, USA. MIT Press.