

3rd World Conference on Technology, Innovation and Entrepreneurship (WOCTINE)

## Data Science for personal use – practical examples

Erol Mrzic<sup>1</sup>, Tarik Zaimovic<sup>1</sup>

<sup>1</sup> *University of Sarajevo, School of Economics and Business,  
Department of Information Systems Development*

---

### Abstract

Due to the unprecedented rise of data content over the last decade an opportunity for data-based personalization and analysis has become a norm in the modern world. Implementing Machine Learning algorithms and Data Science methods no industry remained unchanged. This paper applied these methods and algorithms in personal, practical examples in order to see the benefits in our day-to-day lives. In our case study, we focus on three use cases: a personal movie recommender, messages analysis and real estate trends and predictions on the local market. For this research we used global and personal data, and applied a suitable machine learning model. The purpose of this paper is to establish how one individual, and in what measure, with the use of these new technological tools, can ease his decision making process and manage a more tailored lifestyle.

© 2019 The Author(s). Published by Elsevier B.V.

Peer-review under responsibility of the scientific committee of the 3rd World Conference on Technology, Innovation and Entrepreneurship

**Keywords:** Data Science; Machine Learning; Personal use; Message analysis; Movie recommender; Price prediction;

---

---

\* Corresponding author. E-mail address: [erol.mrzic@student.efsa.unsa.ba](mailto:erol.mrzic@student.efsa.unsa.ba)

### 1. Introduction

In the recent decade, we have witnessed the rise of social networks and growth of Internet-based media platforms resulting in vast amounts of data. All that data has functioned as a fuel for data exploration and data-based predictive models, giving rise to better and more complex models as well as new and innovative use cases [1, 2].

Such cases of Machine Learning algorithms together with data science and analysis methods has caused profound shifts across all industries, with the emergence of new business models, reshaping of transportation [3], healthcare [4], education [5], production and political campaigns, referendums and governments [1, 6].

This techno-analytical golden age is evolving at an exponential pace given that the world we live in now is more interconnected and multi-faced than has ever been before and everything we use is a data generator. This presented a new way to transform entire systems, across (and within) countries, institutions and society. Scaling down from society to individual, the goal of this paper is to establish could we use these methods and algorithms, these ground-breaking tools that large institutions and companies use for much more important purposes, to improve seemingly unimportant personal tasks, and make our lives more productive and effective.

There has been extensive research and projects using Machine Learning across all fields, including the topics that we are discussing in our paper such as: movie recommendations for users that have connection and similar ratings [7], social media message analysis as in alert messages during crises [8] and sentiment analysis [9] and apartment price analysis and prediction for certain areas [10, 11]. Previous research have taken a rather wide approach, while we are trying to use these methods for personal use.

In this paper, we will give an overview of these complex methods and algorithms for non-scientific use, for improving the quality of our private decisions and management of our time and resources. In our case study, we focus on a number of topics, including:

A personal recommender algorithm for movies based on private viewing habits and ratings of a single user. Given that now most of our interaction is text-based, we will make an analysis of personal texting habits giving us an insight of contact priorities and better management for important people in our lives. Finally, we take an example of an important and complex decision such as buying or selling an apartment and we try to make it easier by analysing apartment prices and trends in Sarajevo real estate market based on web collected data.

Our methodology involves collection of global and personal data, its per-processing followed up with a suitable machine learning model.

## **2. Methodology**

In this section, we will introduce our process of data gathering and extraction, some methods which are widely used, then describe the pre-processing for every use case. We will continue to explain our features and feature selection and offer our decision for Machine Learning algorithm. Our methodology contains following steps: Data Extraction, Data Pre-processing, Data Integration and Transformation, Feature Selection and ML algorithm

### **2.1. Data Extraction**

Data extraction is a process that involves retrieval of data from various sources. Extended pre-processing, transformation and integration of data is required in order to further analyse it.

#### ***2.1.1. Use case #1: Movie recommender system:***

In this use case we have used an online kaggle.com dataset containing over 45.000 movie metadata, and over 270.000 user reviews to be the base for our future movie selection. And for our recommendation base we have used a personal 12 month dataset which consisted of 250 movies watched in the period of (May, 2018 – May, 2019), on which we then applied our recommender algorithm.

#### ***2.1.2. Use case #2: Message analysis:***

Here we have extracted data from private social media accounts, including Viber, WhatsApp and Facebook messenger. Total amount of data exceeded 200.000 messages. Data was extracted using official app extraction methods (WhatsApp, and Facebook) and third-party app designed for data extraction of specific app (Viber). Data collected was in html format, and further processing was needed.

#### ***2.1.3. Use case #3: Apartment analysis:***

This use case had us scraping web content for data. Making a custom web scraper script combined with a third party app (Parse Hub) for complicated web page manoeuvring we selected a national web page (Olx.ba) for real estate

information and applied our system. Final result gave us over 80 web pages content with 40 apartment links per page, from which we gained useful data for over than 2500 apartments.

## 2.2. Data Pre-processing

Extracted data is often incomplete, inconsistent, and/or lacking in certain behaviours or trends, and is likely to contain errors. Data pre-processing is a method of transforming raw data into an understandable format and resolving issues. Data pre-processing prepares data for further processing.

### *2.2.1. Use case #1: Movie recommender system:*

In this case included removing unnecessary features and adequate formatting on other columns such as date-time formatting and rounding up values of ratings and votes.

### *2.2.2. Use case #2: Message analysis:*

This case requested more complicated approach. Our data collected was in html format for every contact in our messages. We used a custom script with Beautiful Soup library (bs4) for scraping every html document. Once done, we had a data frame for contact with selected features (Date, Time, Sender, and Content) info. A pipeline was made and applied to every single contact html document, which gave us adequate datasets with more suitable formatting (.csv). Also a special emoji library was imported for extracting emoji's from the content of the messages.

### *2.2.3. Use case #3: Apartment analysis:*

Pre-processing in case #3 consisted of extracting some of the data from the URLs in order for it to be useful (latitude and longitude) and making feature columns out of all raw data we scraped, again using Beautiful Soup library (bs4). Some of the features required type and formatting adjustments, from numerical/string to categorical and float to integer type adjustments. All missing data was given a special category to analyse the magnitude of the missing data, and then subsequently was removed if most of the items data was missing. Data was then formed into dataset with appropriate column names and exported as a single file.

## 2.3. Data Integration and Transformation

Data Integration stands for combining data from several separate sources, which is done using various methods, libraries and technologies to provide a unified view of our data. Data Transformation involves methods to transform or consolidate data into forms suitable for further data mining and analysis.

### *2.3.1. Use case #1: Movie recommender system:*

We grouped our personal dataset by custom users, in this case it was with family, alone or with girlfriend, was it on workdays or weekends. We reshaped our datasets so that our movies were now index rows instead of columns, so we could apply similarity algorithms on the movies that users watched instead of the usual finding similar users approaches. Once the movie ID number was our index rows we had to modify our columns to best describe the movie in question. Python's library Pandas method *get dummies* was used to sort keywords and genres in this situation.

### *2.3.2. Use case #2: Message analysis:*

Data integration in case #2 was done by grouping all contact messages documents into one single dataset that represented its social media app. All missing values were removed from the dataset and column names were modified so they were the same across all datasets. Finally, all datasets were concatenated and exported as a one single file (allMsgs.csv).

### *2.3.3. Use case #3: Apartment analysis:*

Our initial dataset contains 2196 apartment entries, each described with 9 features. We first grouped the apartments per location (municipalities) and analysed prices within each group (Fig. 1). As seen on Fig.1.a there are several outliers which could interfere in future analysis and model application. These outliers were removed from the dataset,

resulting in a more adequate price distribution (Fig.1.b). In the next step, we studied correlation between the features and the target price and removed features with insufficient percentile, i.e. features that do not impact the target price and therefore would not contribute to the model accuracy.

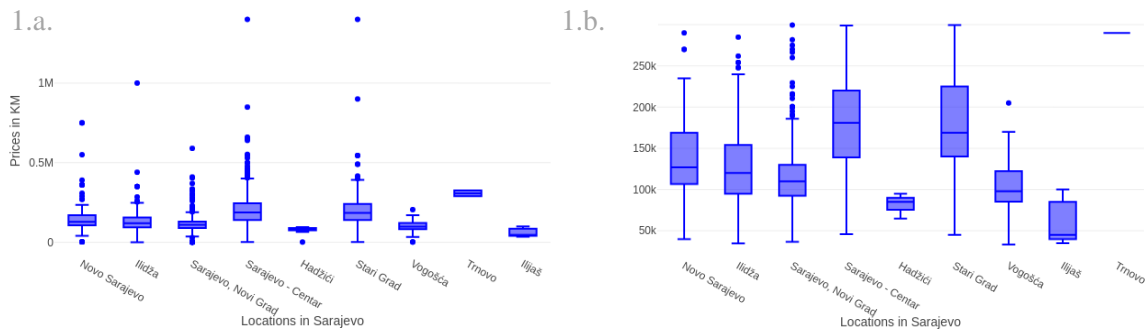


Fig. 1. Price distribution per municipality in Sarajevo; 1.a: Price distribution including outliers; 1.b: Price distribution after removing outliers.

## 2.4. Feature selection and ML algorithm

Feature selection, also known as attribute selection or subset selection, involves selecting specific data points and discarding redundant or irrelevant data, maximizing efficiency and making our Machine Learning model more precise in its predictions.

### 2.4.1. Use case #1: Movie recommender system:

A few features were removed because there was insufficient data and would only hurt the algorithm, such were the budget and revenue columns. Hopefully our taste was not based merely on those exact features. We chose cosine similarity algorithm, which was used to predict similar items in previous papers [12], but this time it was modelled on movies and not on similar users [7].

Cosine similarity is a metric used to determine how similar the documents based on counting the maximum number of common words between the documents irrespective of their size.

We have sorted our movie data in different datasets, based on the day and our company when the movie was watched.

### 2.4.2. Use case #2: Message analysis:

This case again requested more complex approach. For Feature selection we already had a small number of columns to select from. Our data consisted of Date and time of the message, Sender as in the name of the one who sent the message, and Content as in the content of the message. We added a Label column so we could distinguish groups in those contacts such as family, friends or girlfriend, and Messages length column for the number of characters in the content of the message. Also, during specific analysis, such as emoji analysis, data was divided into sent and received categories, which was accomplished by grouping the messages by the Label column.

Data was then subjected to further processing before the NMF topic modelling method was applied which was selected among other methods [13]. NMF (Non-negative Matrix Factorization) is a matrix factorization method where we constrain the matrices to be non-negative, essentially NMF decomposes each data point into an overlay of certain components and other further uses [14].

Our NMF model required use of documents, as in the messages we will choose to analyse for creating topics based on their token content. We established that the messages are too short to find clear patterns in the usage of words, and also the texting habit of sending messages in several lines. We solved this issue by combining the messages into groups of maximum 5 messages based on the time they were sent and the sender. After we made our message groups we dealt with removing punctuation, duplicate letters, conversion to lower case and finally tokenisation (slicing documents into words by using *nltk* package and its word tokenize function) and removing stop words. On that dataset TfidfVectorizer was fitted and NMF method was applied which got us our 30 items with 5 word each that was representing topics in our messages. A handmade classification was then made to label each subsequent topic.

### 2.4.3. Use case #3: Apartment analysis:

Our cleaned dataset included 1980 apartments and 9 feature columns. After using Python's Pandas *get dummies* method on some columns, we ended up with 994 feature columns. We decided on Random Forrest regressor in this use case because of its great performances so far [10, 11].

Random Forrest is an addition to the decision tree algorithm and comprises of a random collection of multiple decision trees across our data. A decision tree is a flowchart-like structure in which each internal node represents a test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label.

We removed the prices from the dataset and applied the model on the selected features. We used 80% of the data for training and 20% for testing, and the number of estimators in Random Forrest was 35.

We also used k-fold cross-validation for our results so that we could use all of our data for testing and training. In k-fold cross validation, the data is divided into k subsets. Now the holdout method is repeated k times, such that each time, one of the k subsets is used as the test set / validation set and the other k-1 subsets are put together to form a training set. The error estimation is averaged over all k trials to get total effectiveness of our model.

We used k-means unsupervised learning method on our error data to establish most likely clusters of data with substantial error margin. The goal of this algorithm is to find groups in the data, with the number of groups represented by the variable *k*. The algorithm works iteratively to assign each data point to one of *k* groups based on the features that are provided. Further use can be found here [15].

## 3. Results

### 3.1. Use case #1: Movie recommender system:

For every movie in that dataset our personal rating was selected, and based on our criteria (rating bigger than 3.0 - highest 5.0, lowest 1.0) those movies were then pushed into our cosine similarity algorithm, and from our movie database of 45.000 movies the movies that coincided with the ones we watched on specific day, with common keywords, similar average popular rating, and genre was then returned.

For weekend with SO:	For workdays with family:	For weekends alone:
Shanghai Triad	Casino	Ace Ventura: When Nature Calls
The City of Lost Children	Father of the Bride Part II	Shanghai Triad
Jumanji	Ace Ventura: When Nature Calls	Father of the Bride Part II
Babe	Othello	The City of Lost Children
Lamerica	Restoration	Dangerous Minds
Home for the Holidays	To Die For	Four Rooms
Don't Be a Menace to South Cen...	Se7en	Sense and Sensibility
Two If by Sea	Pocahontas	Heat
The Postman	Mighty Aphrodite	Get Shorty
Kids of the Round Table	Cutthroat Island	Copycat

Fig. 2. Results from our recommendation system based on day and company.

Results are now subjective (result sample shown on Fig. 2.), based on do we really like the movies recommended, but the algorithm did find similar movies based on small datasets to work with. Our personal data was 250 movies all together, given that we divided our data into categories based on day and company when the movie was watched, our datasets were very small in comparison to the movie data we were finding similarities to (45.000 movies).

### 3.2 Use case #2: Message analysis:

Our analysis showed our activity over 1-year period and we added the scale of sent and received messages. As shown in Fig. 3.

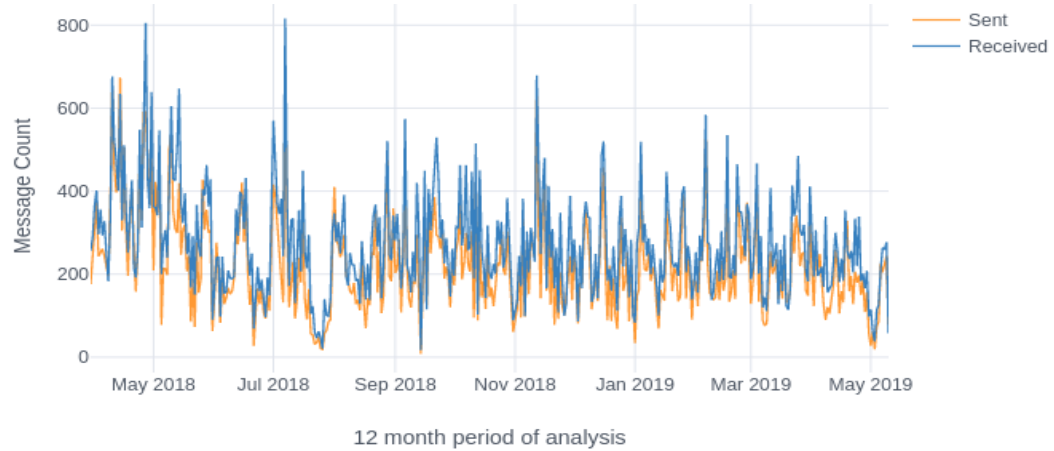


Fig. 3. Graph showing full data distribution over the period of 12 months.

We see that the number of received messages is higher than the number of sent messages, but the overall count of messages is quite high. The conclusion is that the subject relies on texting quite much but is less active than his contacts. Detailed view of our hourly message rate is offered in Fig. 4. and Fig. 5.

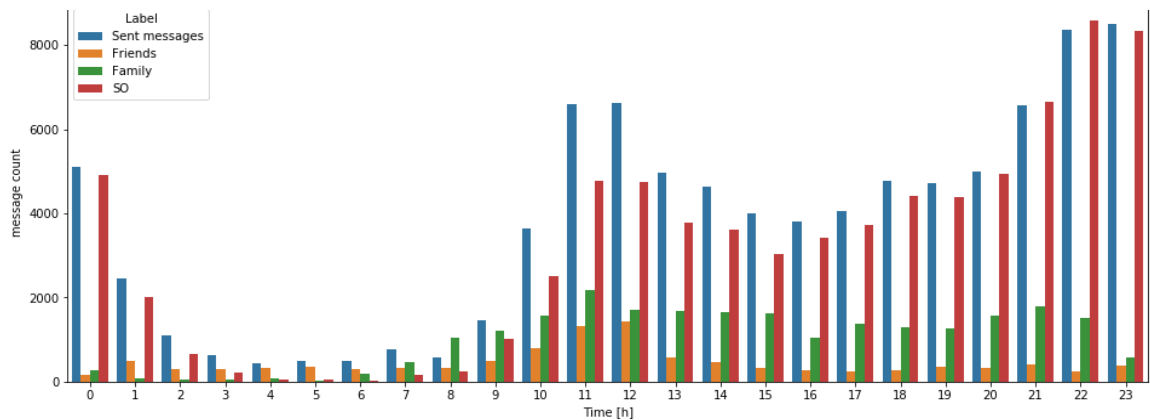


Fig. 4. Showing hourly rate of messages with labelled data corresponding to our grouped contacts

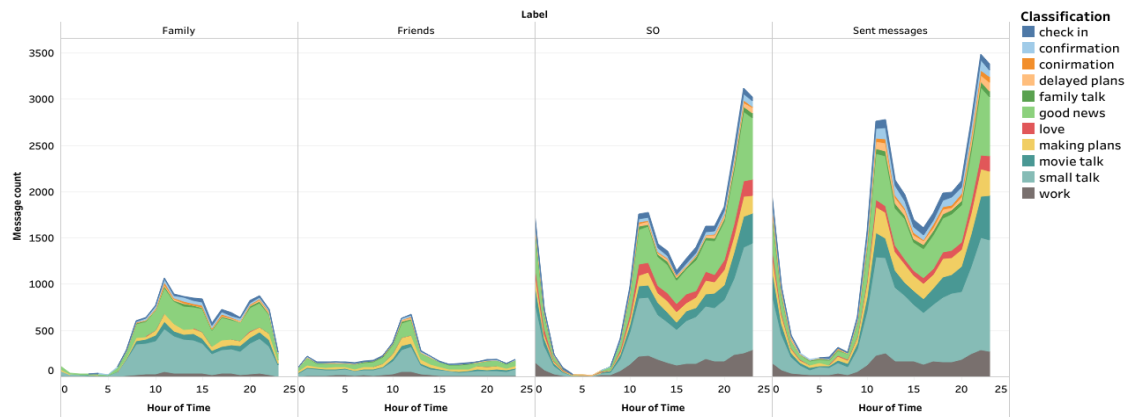


Fig. 5. The volume of conversation topics per hour.

From our hourly rate of messages analysis, we can conclude that group labelled *Friends* is the least active group, but spikes during noon, as do all other messages. Conversation with the *Significant Other* is present during the entire day and increasing towards late hours. *Family* labelled conversation are usually after work hours until bed time. We can also see that our *Sent messages* peak occurs in the 10-12 both PM and AM time of day.

Overall view of our conversation topics classification in general, as well as for each group can be seen on Fig. 6.

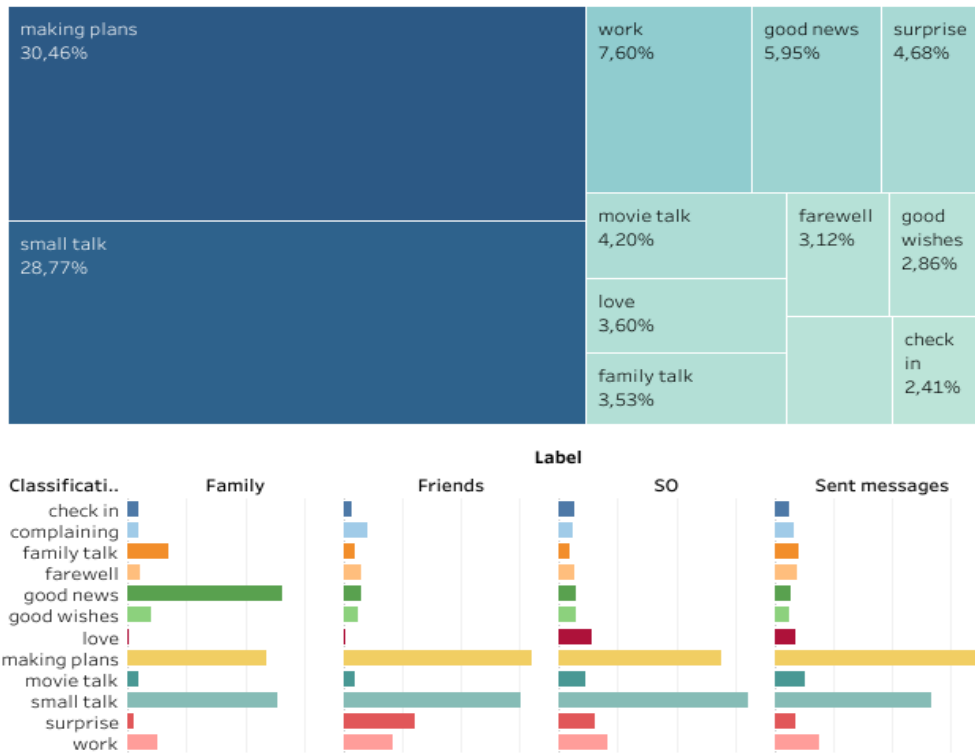


Fig. 6. Distribution of topics across all messages in labelled groups.

In our message topics analysis, we could establish that nearly 60% (59.23%) of our conversation data falls under the classification of ‘making plans’ and ‘small talk’. Those being among the most popular topics across all contact groups with small margin between them, except for *Sent Messages* where ‘making plans’ topic and *Family* messages where ‘good news’ topic have the highest percentage.

### 3.3 Use case #3: Apartment analysis:

Our model delivered a Cross-Validation result of 79% accuracy on price prediction (as seen on Fig. 7.), with an average percentage error in the ran ge of 2-5%, but some high outliers (as seen on Fig. 8.).

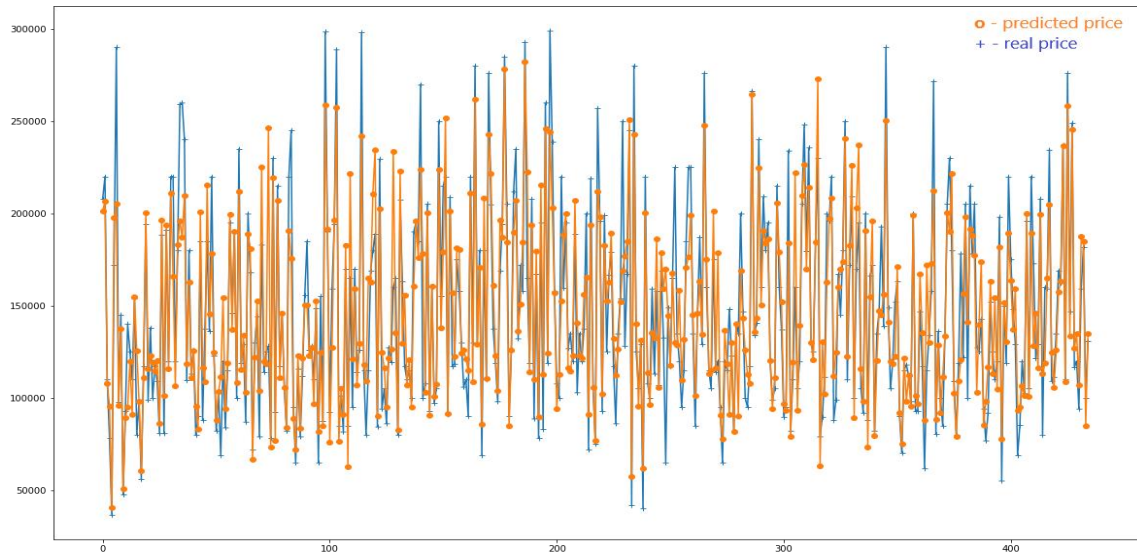


Fig. 7. Random Forreest prediction model

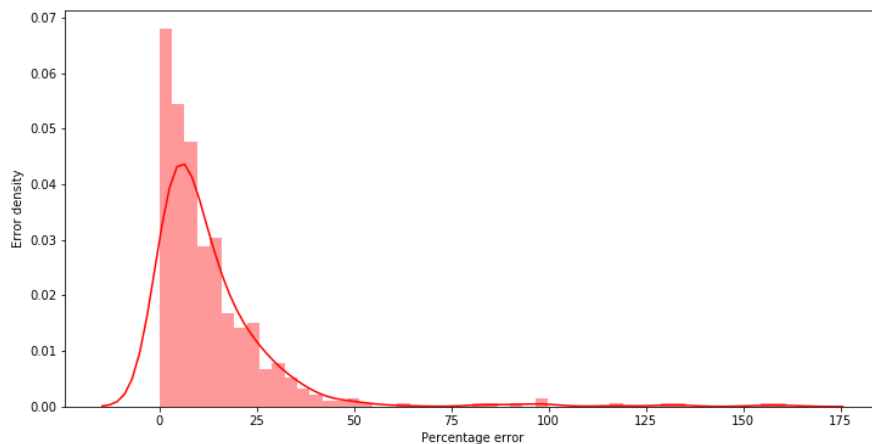


Fig. 8. Error distribution for Apartment prediction

We went on to further analyse our error data by grouping apartments with over 30% error in price prediction and applying k-means method on it.

As a result, we got clusters of apartments with high error percentage, and after looking into it we saw that the high error in prediction was mostly because of false description segment (e.g. ‘Damaged construction’, ‘leaking’) which we, avoiding text analysis in this occasion, did not include in our features.

#### 4. Discussion

This paper consisted of 3 projects that served personal usage, and the aim was to conclude if they could help manage our time as well as resource and help with our decision-making process. In Use case #1 we made a movie recommender system that could help our viewing habits and ease the movie selection. We achieved the recommender system efficiency that we wanted, movies are selected on the basis of our watch history and weekday activity. Also having labelled our users for these activities, we can expect over time to develop a full profile and tailored recommendations for each of them. In use case #2 we gained great insight through our message analysis in



terms of usage and conversation topics distributed through social media apps. We can now tell our communication style and timeline, and our preferred topics over different groups. We could see our messaging habits during the day and decide to act on them. Use case #3 gave us a predicting model for apartment prices which we can use to establish if the apartment had a fair value, as in was the price fair or not. If we came across an apartment with the price much lower than our predicted value, it could indicate a possible bargain and a good deal. We could then focus on specific areas as we had a visual geographical look at our data.

## 5. Conclusion

Our analysis showed how by rather simple use of Machine Learning and Data Science when regarding mundane problems with different degrees of complexity can help us make decisions or gain insight in our personal life. It can ease our research and simplify our selection, without any major setbacks. Each model can be further developed and more features can be added that would increase accuracy in our prediction. For use case #1 personal dataset will have to be bigger to further test our algorithm. And adding features like budget and revenue, as well as director and cast will certainly improve the prediction results on this model. It will be able to find more personalized patterns in our data and hence recommend more similar movies. Use case #2 can be tracked for a longer period than 1 year and data can be labelled to show which social media has more priority. Use case #3 can be improved by adding features like proximity to train stations and night clubs or university distance. Adding a sentiment analysis on the apartment description section would make a valuable feature as well. Ongoing process of digitalization will necessitate even further industrialization and monetization of ML applications across variety of services and areas.

## References

- [1] Palash Dey, Pravesh K. Kothari, and Swaprava Nath. 2019. The Social Network Effect on Surprise in Elections. In *Proceedings of the ACM India Joint International Conference on Data Science and Management of Data (CoDS-COMAD '19)*. ACM, New York, NY, USA, 1-9. DOI: <https://doi.org/10.1145/3297001.3297002>
- [2] J. Chen, M.Li (2019) “Chained Predictions of Flight Delay Using Machine Learning”, *AIAA Scitech 2019 Forum* 2019-1661
- [3] C. Jacob, A. Hadayeghi, B. Abdulhai, B. J. Malone (2006) “Highway Work Zone Dynamic Traffic Control Using Machine Learning” 2006 *IEEE Intelligent Transportation Systems Conference*, Toronto, Canada, September 17– 20, 2006
- [4] J. T. Polletini, R. Tinós, J. A. Baranauskas, A. A. Macedo, S. R. G. Panico, J. C. Daneluzzi (2012) “Using Machine Learning Classifiers to Assist Healthcare-Related Decisions: Classification of Electronic Patient Records” *Springer Science+Business Media, LLC* 2012
- [5] N. SUN, X. PEI 2, S. ZHOU (2008) “*FACIAL EMOTION RECOGNITION IN MODERN DISTANT EDUCATION SYSTEM USING SVM*” *Proceedings of the Seventh International Conference on Machine Learning and Cybernetics, Kunming*.
- [6] A. P. Patil, D. Doshi, D. Dalsaniya (2018) “Applying Machine Learning Techniques for Sentiment Analysis in the Case Study of Indian Politics”, *ISpringer International Publishing AG 2018 S.M. Thampi et al. (eds.), Advances in Signal Processing and Intelligent Recognition Systems, Advances in Intelligent Systems and Computing* 678, DOI 10.1007/978-3-319-67934-1\_31
- [7] P. Perny and J. Zucker (2001) “Preference-based Search and Machine Learning for Collaborative Filtering: the “Film-Conseil” Movie Recommender System”, *University Paris 6 Computer Science Laboratory*.
- [8] Brynielsson et al.: “Emotion classification of social media posts for estimating people’s reactions to communicated alert messages during crises” *Security Informatics* 2014 3:7.
- [9] G. Peterson and S. Shenoi (Eds.): *Advances in Digital Forensics X*, IFIP AICT 433, pp. 253–265, 2014. c IFIP International Federation for Information Processing 2014
- [10] M. Čeh, M. Kilibarda A. Lisec and B. Bajat (2018) “Estimating the Performance of Random Forest versus Multiple Regression for Predicting Prices of the Apartments”, *ISPRS Int. J. Geo-Inf.* 2018, 7, 168
- [11] Pow, N., Janulewicz, E., & Liu, L. (2014). *Applied Machine Learning Project 4 Prediction of real estate property prices in Montréal*.
- [12] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl. 2001. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th international conference on World Wide Web (WWW '01)*. ACM, New York, NY, USA, 285-295. DOI: <https://doi.org/10.1145/371920.372071>
- [13] Lee, M., Wang, W., & Yu, H. (2006). Exploring supervised and unsupervised methods to detect topics in biomedical text. *BMC bioinformatics*, 7, 140. doi:10.1186/1471-2105-7-140
- [14] Lee, D. D., & Seung, H. S. (2001).
- [15] Hartigan, J., & Wong, M. (1979). Algorithm AS 136: A K-Means Clustering Algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1), Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems* (pp. 556-562). 100-108. doi:10.2307/2346830