# IS 584 Course Term Project Proposal

Özkan Tuğberk Kartal

*Electrical-Electronics Engineering*
*METU*
Ankara, TURKEY
tugberk.kartal@metu.edu.tr

## I. INTRODUCTION

The purpose of this document is to give explanations regarding Proposal, EDA (Exploratory Data Analysis) and Quality Checks and Github Repository of IS 584 Course Term Project.

## II. PROPOSAL

### A. Research Questions

Two research questions are (i) "How to predict the acceptance/rejection decision of the articles by modeling the hierarchical relationships present in their abstract and (meta)reviews?" and (ii) "If we use Capsule Neural Networks, can we improve the semantic understanding of abstract texts and (meta)reviews compared to traditional models like CNNs (Convolutional Neural Networks) or LSTMs (Long Short Term Memories)?".

### B. Model Architecture and Framework

It is planned to use Capsule Neural Network architecture in PyTorch framework [1].

### C. Performance Metrics

Accuracy, F1-Score, Area Under the ROC (Receiver Operating Characteristic) Curve and Confusion Matrix are used to evaluate the performance of the Capsule Neural Network.

### D. Expected Outcomes

By using Capsule Neural Network, hierarchical and semantic relationships in the dataset can be captured in a much advanced manner and due to the dynamic routing mechanism of Capsule Neural Network, unseen review patterns can be handled much better.

### E. Potential Benefits

Capsule Neural Networks, in addition to its ability to capture hierarchical and semantic relationships in the dataset and its dynamic routing mechanism, can provide better robustness to the variations in the dataset, can guarantee less information loss and requires fewer training samples.

## III. EDA AND QUALITY CHECKS

The notebook file named "EDA_and_Quality_Checks.ipynb" was written in order to make EDA and Quality Checks and saved under the folder named "Notebooks" in [2]. Moreover, by using this notebook file various histograms were plotted and saved in the formats "eps", "jpg" and "svg" under the directory named "figures/histograms" in [2]. In this document, we can mention some of these analyses.

The dataset is about ICLR (International Conference on Learning Representations) papers from 2017-2020 and NIPS (Neural Information Processing Systems) papers from 2016-2019 and information related to these papers are located in their corresponding files in "json" format [3]. All of these files in "json" format were loaded in order to form pandas data frame variables "ICLR_data" (only ICLR data) having 5194 rows, "NIPS_data" (only NIPS data) having 3934 rows and "data" (ICLR and NIPS data together) having 9128 rows. The column names of the variable "data", the format of each column, unique values related to each column can be summarized in Table I and II. It must be emphasized that there exists NaN and None values used for missing data row-column values. Regarding decision column, the histogram can be seen in Figures 1.
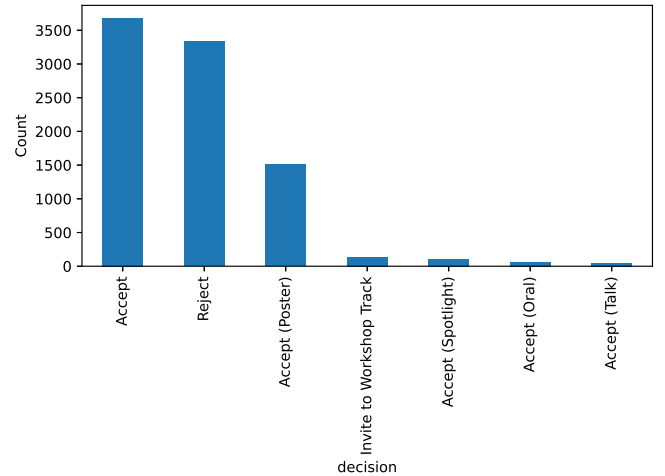


Fig. 1. Decision Histogram

TABLE I
DATA SUMMARY

| Column Name | Format | Unique Value |
|---|---|---|
| name | string<br>Some Examples<br>ICLR_2017_1.pdf,<br>NIPS_2016_1.pdf | 8856 unique values |
| id | string<br>Some Examples<br>ICLR_2017_1,<br>NIPS_2016_1 | each of them |
| metadata.source | string | CRF, META, NaN |
| metadata.title | string | 5459 unique values<br>(None and NaN<br>for missing<br>3652 time occurrences,<br>RECURRENT<br>NEURAL<br>NETWORKS<br>5 time occurrences,<br>DEEP<br>REINFORCEMENT<br>LEARNING<br>3 times occurrences,<br>some other titles<br>occur more than once,<br>most of the titles<br>occur only once) |
| metadata.authors | list<br>An Example<br>(Jonathon Cai,<br>Richard Shin,<br>Dawn Song) | not applicable |
| metadata.emails | list<br>An Example<br>(jonathon@<br>cs.berkeley.edu,<br>ricshin@<br>cs.berkeley.edu,<br>dawnsong@<br>cs.berkeley.edu) | not applicable |
| metadata.sections | list of dictionaries | not applicable |
| metadata.references | list of dictionaries | not applicable |
| metadata.referenceMentions | list of dictionaries | not applicable |
| metadata.year | int | 0, 1969, 2016, 2017,<br>2018, 2019, 2020, NaN |
| metadata.abstractText | string | 8844 unique values<br>(None and NaN<br>for missing<br>281 time occurrences,<br>some abstracts<br>occur twice<br>due to either<br>uploading to<br>the same conference<br>twice or to the two<br>distinct conferences,<br>most of the abstracts<br>occur only once) |
| metadata.creator | string,<br>An Example<br>Microsoft®<br>Word 2016 | 59 unique values |
| conference | string | ICLR, NIPS, NaN |
| decision | string | Accept<br>Accept (Oral)<br>Accept (Poster)<br>Accept (Spotlight)<br>Accept (Talk)<br>Invite to Workshop Track<br>Reject<br>NaN |

TABLE II
DATA SUMMARY (CONTINUATION)

| Column Name | Format | Unique Value |
|---|---|---|
| url | string | 8878 unique values |
| hasContent | boolean | true, false, NaN |
| hasReview | boolean | true, false, NaN |
| title | string | 8826 unique values<br>(None, NaN, NA and N\A<br>for missing<br>279 time occurrences,<br>some titles<br>occur twice<br>due to either<br>uploading to<br>the same conference<br>twice or to the two<br>distinct conferences,<br>most of the titles<br>occur only once) |
| authors | list<br>An Example<br>(Jonathon Cai,<br>Richard Shin,<br>Dawn Song) | not applicable |
| reviews | list of dictionaries | not applicable |
| metaReview | string | 5876 unique values |

## IV. GITHUB

The Github link of the repository is given by [2].

### REFERENCES

[1] Wikipedia contributors, Capsule neural network — Wikipedia, The Free Encyclopedia. 2024. [Online]. Available: https://en.wikipedia.org/wiki/Capsule_neural_network

[2] Ö. T. Kartal, IS584CourseTermProject. Online. Available: https://github.com/oztuka/IS584CourseTermProject/

[3] Volga Sezen, "IS 584 Course Term Project Dataset," 2025. Available: https://drive.google.com/file/d/1nJdljy468roUcKLbVwWUhMs7teirah75/view