

# Unstructured Data Analysis (Text Mining)

2017 Spring

School of Industrial Management Engineering

## 1. Overview

- ✓ This module aims to provide students with the theoretical and practical knowledge and skills to collect, modify, and analyze a large amount of unstructured data, especially texts, from various sources.
- ✓ Topics covered in this module include data collection methods from various sources, preprocessing methods including natural language processing, document representation & summarization, feature selection and extraction, document clustering, document classification, and topic models.
- ✓ The students are assessed by one final exam at the end of the semester, three presentations (proposal, interim, and final) and the final manuscript for their term projects.

## 2. Lecturer & Course homepage

- ✓ Pilsung Kang, Assistant professor at School of Industrial Management Engineering, Korea University
  - E-mail: pilsung\_kang@korea.ac.kr
  - Course homepage: <https://github.com/pilsung-kang/text-mining>

## 3. Textbook and additional resources (not mandatory)

- ✓ Weiss, S.M., Indurkha, N., and Zhang, T. (2010). Fundamentals of Predictive Text Mining. Springer.
- ✓ Feldman, R. and Sanger, J. (2007). The Text Mining Handbook. Cambridge University Press.
- ✓ Kao, A. and Poteet, S.R. (2007). Natural Language Processing and Text Mining. Springer.
- ✓ Manning, C.D., Raghavan, P., and Schütze, H. (2008). Introduction to Information Retrieval. Cambridge University Press.
- ✓ Jurafsky, D. and Martin, J.H. (2008). Speech and Language Processing, 2<sup>nd</sup> Ed. Prentice Hall. (Free online course available: <https://www.youtube.com/playlist?list=PL6397E4B26D00A269>)
- ✓ Socher, R. (2016). CS224d @Stanford: Deep learning for natural language processing (course homepage: <http://cs224d.stanford.edu/>, video lectures are available at Youtube)
- ✓ Blunsom, P. et al. (2017). Deep natural language processing @Oxford (course homepage: <https://github.com/oxford-cs-deepnlp-2017/lectures>)

## 4. Assessments

- ✓ Final exam (30%): Closed book
- ✓ Term project (70%): three presentations
  1. Group project: maximum 4 students in a group
  2. Proposal (10%): purpose of the project (task), data description, expected effects, etc.
  3. Interim presentation (10%): data collection/preprocessing, feature extraction, issues to be discussed.
  4. Final presentation (20%): employed/developed models, experimental results including interesting

patterns discovered, limitations and future research directions.

5. Research paper (20%): each team **must** write a research paper based on the results of the term project. The final manuscript should be completed by the end of the 16<sup>th</sup> week. The type of the manuscript determines the maximum grade.

1. A+: international journal paper (must be written in English)
2. A: domestic journal paper (written in either English or Korean)

## 5. Introduce yourself

- ✓ Submit your self-introduction slide (max. 5 pages) to the lecturer via E-mail by the end of the 2<sup>nd</sup> week.

## 6. Schedule & Topics

| Week  | Contents  | Exercises                                |
|-------|---|--|
| 1     | Orientation<br>Introduction to Text Mining<br>✓ The usefulness of large amount of text data and the challenges<br>Overview of text mining methods   |  |
| 2     | From Texts to Data<br>✓ From text files, databases, Facebook APIs, and web scraping   | R Exercise                               |
| 3     | Natural Language Processing<br>✓ Morphological analysis: tokenization, stemming/lemmatization, POS tagging, parsing, chunking, named entity recognition, language model, etc.   | R Exercise                               |
| 4     | Document Representation<br>✓ Bag-of-words, word weighting, N-gram, and distributed representation   | <b>Term Project Proposal</b>             |
| 5     | Document Summarization<br>✓ Summarize a large amount of texts to understand at a glance (Wordcloud)<br>✓ Interpret the relationship between features (Association rules, Graphs)  |  |
| 6     | Dimensionality Reduction: Feature Selection and Extraction<br>✓ Supervised feature selection: index term selection, information gain, cross entropy, etc.<br>✓ Unsupervised feature selection: latent semantic analysis (LSA) | R Exercise                               |
| 7     | Document Similarity & Clustering<br>✓ Document similarity measures: cosine similarity, Euclidean distances, etc.<br>Clustering algorithms: K-means clustering, hierarchical clustering  | R Exercise                               |
| 8     | Document Classification<br>✓ Naïve Bayesian classifier, k-nearest neighbor classifier<br>Classification performance evaluation  | <b>Term Project Interim Presentation</b> |
| 9     | Topic Modeling I: Probabilistic Latent Semantic Analysis (pLSA)   | R Exercise                               |
| 10    | Topic Modeling II: Latent Dirichlet Analysis (LDA)  | R Exercise                               |
| 11    | Sentiment Analysis<br>✓ Dictionary-based sentiment analysis<br>✓ Model-based sentiment analysis   | R Exercise                               |
| 12-14 | Prepare the final presentation and write the manuscript   |  |
| 15    | <b>Final Exam</b>   |  |
| 16    | <b>Term Project Final Presentation</b>  |  |