

# Introduction to Natural Language Processing

BUS 243F: Summer 2023

Tuesday and Thursday 9:00 am – 11:00 am

Lemberg Acad Center 180

Instructor: Yeabin Moon ([yeabinmoon@brandeis.edu](mailto:yeabinmoon@brandeis.edu))

**Office Hours:** After class in my office (*Sachar International Center 209B*), or by appointment online through the following [link](#).

## Course Description

Natural language processing (NLP) is becoming increasingly widespread. Applications of NLP have become embedded in our everyday lives, and these applications are based somewhere between formal linguistics and statistical physics. Especially over the past decade, neural network approaches have become the de facto standard for many NLP tasks. This course aims to provide a survey of these foundations, but we will take NLP in a narrow sense to cover the text analysis tasks. The course assumes a background in multivariate calculus, college level linear algebra/statistics, and proficiency in Python. The goal of this course is to enable you to build your language applications using the Python framework.

Success in this course is based on the expectation that students would need to study for about three hours for every hour of in-class time. Hence, students will spend a *minimum of 12 hours* of study time per week in preparation for this class.

## Learning Goals

With this course, you will gain knowledge in the following areas:

1. Basic techniques for processing textual data
2. Familiarity with the domain-specific language commonly used in the field of Natural Language Processing (NLP) and its core applications
3. A range of text representations
4. Building sentiment analyzers and understanding evaluation metrics for assessing their performance
5. Utilizing a deep learning approach using popular frameworks such as Keras and TensorFlow

## Main Reference

We will use ***Speech and Language Processing (3<sup>rd</sup> edition)*** by Dan Jurafsky and James H. Martin as a main reference. You can access the book through the following [link](#). There is a partial list of useful books that will be touched during the course.

- Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, [\*Introduction to Information Retrieval\*](#)

- Hannes Hapke, Hobson Lane, Cole Howard, *Natural Language Processing in Action*
- Steven Bird, Ewan Klein, Edward Loper, *Natural Language Processing with Python – Analyzing Text with the Natural Language Toolkit*
- Sowmya Vajjala, Bodhisattwa Majumder, Anuj Gupta, Harshit Surana, *Practical Natural Language Processing*

You have an **online access** for all the references listed above through Brandeis Library. Other useful reference is *Introduction to Natural Language Processing* by Jacob Eisenstein for avid students of mathematical exposition.

## Prerequisites

1. Competency in Python (**Bus215f**)
  - All class exercises will be using Python. You should be familiar with NumPy, pandas and data structures in Python. Note that you should be fine if you have ample experience in coding with a different language.
2. Calculus, Linear Algebra, Probability, and Statistics (**Econ213a**)
  - You should know college-level calculus and the basics of probabilities
3. Machine Learning (**recommended**)
  - If you have basic machine learning or deep learning experience, the course would be much easier. *You can take it without knowing them.* If you need a top-bottom textbook treatment, I highly recommend: “*Hands-on machine learning with scikit-learn and TensorFlow*” by Geron Aurelien

## Class Participation

There is no such thing as a stupid question. Dialogue is not only strongly encouraged, it is critical to your understanding of the material. Vocalizing your questions often helps you solidify what you do and do not understand. It also provides me important feedback on the areas in which we need to spend more time. During lectures, I will encourage questions, and I will solicit input. If I call on you, please relax, I am NOT trying to intimidate you or embarrass you in any way. I am trying to encourage active listening and keep you engaged in the course. This will greatly assist you in learning the material. If you do not know the answer, I will move on to another student. **Hence, attendance is mandatory for this class.**

## Course Requirements

There are five assignments designed to enhance your theoretical understanding and practical skills. Each assignment consists of both written and programming components. It is essential to complete all assignments. Additionally, reading materials are mandatory for the course, as they provide comprehensive coverage beyond the lectures. Finally, there will be a final exam to evaluate your overall comprehension. The school has not yet announced the specific date for the final exam in each course for the summer semester. However, it is expected to be scheduled on either **August 10 or 11**.

Please be aware that the final exam will be conducted **in-person**, and **there will be no online option available for this class**.

You can submit your late work with a 10 percent penalty if you can make it within ONE day after the deadline.

The grade consists of

1. five assignments: 55 %
2. Exam: 40 %
3. Participation / Attendance: 5 %

## Course Plan

The class covers the basic building blocks used in NLP. We will mainly examine the practical use cases and delve into theories where necessary. Each lecture will be dedicated to one concept. However, some additional concepts would be introduced due to the compact class structure. The following outline provides a high-level overview of the course. The mandatory readings are indicated by \*, and you must be prepared before the class meeting. The additional readings will be posted on the latte one week in advance. The in-class quizzes primarily focus on the weekly readings.

### 1. Introduction to NLP (July 11)

(a) Course Logistics: Textual data

(b) Information Retrieval Primer

- Chapter 1 in [Introduction to Information Retrieval](#) \*
- Thought process approach
- Evaluation Metric

(c) Python and Math Reviews

- Google Colab
- Object Oriented Programming
- Regular Expressions
- Bayes in math

### 2. Vector Space Model (July 13)

(a) Search mechanism

- Chapter 6 in [Introduction to Information Retrieval](#) \*
- Boolean search and Ranked Retrieval

(b) Vector reasoning

- Count and Incidence
- Similarity

(c) Heap's Law

(d) Zipf's Law

**Assignment 1 deadline (11:59 AM, July 15)**

### **3. Text Representation primer (July 18)**

- (a) Chapter 2 (more important) and Chapter 3 (less important) in Speech and Language Processing \*
- (b) Regular Expression and Text search
- (c) Discussion: how to represent a large text data
  - Park, M., Leahey, E. & Funk, R.J. Papers and patents are becoming less disruptive over time. Nature 613, 138–144 (2023). <https://doi.org/10.1038/s41586-022-05543-x> \*

### **4. Introduction to Sentiment Analysis (July 20)**

- (a) Text categorization primer
  - Chapter 4 in Speech and Language Processing \*
- (b) Naïve Model
  - Theory vs. Practice
  - Evaluation metric
  - Math review
- (c) Discussion
  - COHEN, L., MALLOY, C. and NGUYEN, Q. (2020), Lazy Prices. The Journal of Finance, 75: 1371-1415. <https://doi.org/10.1111/jofi.12885> \*

### **Assignment 2 deadline (11:59 AM, April 22)**

### **5. Logistic Regression (July 25)**

- (a) Understand Discriminative classifiers
  - Chapter 5 in Speech and Language Processing \*
- (b) SGD algorithm review
- (c) Discussion on ML approach
  - Discriminative classifiers vs. Generative classifiers
  - Intuition of numerical optimization

### **6. Introduction to Word Embeddings (July 27)**

- (a) Introduction to Word2vec
  - Chapter 6 in Speech and Language Processing \*
  - Mikolov, T., Chen, K., Corrado, G. and Dean, J., 2013. Efficient estimation of word representations in vector space \*
- (b) Discussion on word embeddings
- (c) Gensim in practice

### **Assignment 3 deadline (11:59 AM, July 29)**

### **7. Deep learning Primer (August 1)**

- (a) Chapter 7 in Speech and Language Processing \*
  - XOR problems
- (b) Introduction to Deep learning software
  - Keras and Tensorflow

- Revisit: Sentiment analysis
- (c) Language model (if time permitted)
  - a. Discussion on input vector representations
  - b. Distributional semantics

## **8. Text Analysis with RNN (August 3)**

- (a) Introduction to Recurrent neural network
  - Chapter 9 in Speech and Language Processing \*
    - i. Keras and Tensorflow application
- (b) Sentiment analysis
- (c) Attention Mechanisms
- (d) Formal Language Theory (if time permitted)

## **Assignment 4 deadline (11:59 AM, August 5)**

## **9. NLP with Transformers (August 8)**

- (e) Encoder-Decoder problem
  - Readings will be provided on Latte \*
  - Bert Algorithm
- (f) Hugging Face's Transformers
- (g) Discussion on FinBert
  - Huang, A.H., Wang, H. and Yang, Y. (2023), FinBERT: A Large Language Model for Extracting Information from Financial Text\*. Contemp Account Res, 40: 806-841. <https://doi.org/10.1111/1911-3846.12832> \*

## **Assignment 5 deadline (11:59 AM, August 12)**

## **10. Final exam (TBA: August 10 or 11)**

The course plan is subject to change due to a weather condition/delayed start/early closing. If this situation happens, **the class will be held on zoom**. I will announce it accordingly.

## **Accommodations**

Brandeis seeks to create a learning environment that is welcoming and inclusive of all students, and I want to support you in your learning. Live auto transcription is available for all meetings or classes hosted on Zoom and you can turn it on or off to support your learning. Please [check for Zoom updates](#) to take advantage of this new feature. To learn more, visit the [Zoom Live Transcription webpage](#). For questions, contact [help@brandeis.edu](mailto:help@brandeis.edu)

If you think you may require disability accommodations, you will need to work with Student Accessibility Support (SAS) (781-736-3470, [access@brandeis.edu](mailto:access@brandeis.edu)). You can find helpful student FAQs and other resources on the [SAS website](#), including guidance on

how to know whether you might be eligible for support from SAS. If you already have an accommodation letter from SAS, please provide me with a copy as soon as you can so that I can ensure effective implementation of accommodations for this class.

## **Academic Integrity**

Every member of the University community is expected to maintain the highest standards of academic integrity. A student shall not submit work that is falsified or is not the result of the student's own effort. Infringement of academic integrity by a student subjects that student to serious penalties, which may include failure on the assignment, failure in the course, suspension from the University or other sanctions. Please consult [Brandeis University Rights and Responsibilities](#) for all policies and procedures related to academic integrity. Students may be required to submit work via TurnItIn.com or similar software to verify originality. A student who is in doubt regarding standards of academic integrity as they apply to a specific course or assignment should consult the faculty member responsible for that course or assignment before submitting the work. Allegations of alleged academic dishonesty will be forwarded to the Department of Student Rights and Community Standards. Citation and research assistance can be found at [Brandeis Library Guides - Citing Sources](#).

## **Classroom Health and Safety**

- Register for the [Brandeis Emergency Notification System](#). Students who receive an emergency notification while attending class should notify their instructor immediately. In the case of a life-threatening emergency, call 911. As a precaution, review [this active shooter information sheet](#).
- Brandeis provides [this shuttle service](#) for traveling across campus or to downtown Waltham, Cambridge and Boston.
- On the Brandeis campus, all students, faculty, staff and guests are required to observe the university's policies on physical distancing and mask-wearing to support the health and safety of all classroom participants. Review up to date [COVID-related health and safety policies](#) regularly

## **Student Support**

Brandeis University is committed to supporting all our students so they can thrive. If you want to learn more about support resources, the [Support at Brandeis](#) webpage offers a comprehensive list that includes these staff colleagues you can consult, along with other support resources:

- The [Care Team](#)
- [Academic Services](#) (undergraduate)
- [Graduate Student Affairs](#)
- Directors of Graduate Studies in each department, School of Arts & Sciences
- Program Administrators for the Heller School and International Business School
- [University Ombuds](#)
- [Office of Equal Opportunity](#)